

Basics of UNIX

August 27, 2011

1 Connecting to a UNIX machine from {UNIX, Mac, Windows}

See the file on bspace on connecting remotely to SCF.

More generally, [the department computing FAQs](#) is the place to go for answers to questions about SCF.

For questions not answered there, the SCF help page requests: “please report any problems regarding equipment or system software to the SCF staff by sending mail to ‘trouble’ or by reporting the problem directly to room 498/499. For information/questions on the use of application packages (e.g., R, SAS, Matlab), programming languages and libraries send mail to ‘consult’. Questions/problems regarding accounts should be sent to ‘manager’.”

Note that for the purpose of this class, questions about application packages, languages, libraries, etc. can be directed to me (not *consult*).

2 Files and directories

1. Files are stored in directories (aka folders) that are in a (inverted) directory tree, with “/” as the *root* of the tree
2. Where am I?
> pwd
3. What’s in a directory?
> ls
> ls -a
> ls -al

4. Moving around

```
> cd /home/paciorek/teaching/243
> cd ~paciorek/teaching/243
> cd ~/teaching/243
> cd 243 # provided I am in 'teaching'
> cd ..
> cd -
```

5. Copying and removing

```
> cp
> cp -r
> cp -rp # preserves timestamps and other metainfo (VERY handy
for tracing your workflows if you move files between machines)
> mkdir
> rm
> rm -r
> rm -rf # CAREFUL!
```

To copy between machines, we can use *scp*, which has similar options to *cp*:

```
> scp file.txt paciorek@bilbo.berkeley.edu:~/research/.
> scp paciorek@bilbo.berkeley.edu:/data/file.txt ~/research/renamed.txt
```

6. File permissions: **ls -al** will show the permissions for the *'user'*, *'group'*, *'other'*

- to allow a file to be executed as a program:

```
> chmod ugo+x myProg # myProg should be compiled code or a
shell script
```

- to allow read and write access to all:

```
> chmod ugo+rw code.q
```

- to prevent write access:

```
> chmod go-w myThesisCode.q
```

7. Compressing files

- the *zip* utility compresses in a format compatible with zip files for Windows:


```
> zip files.zip a.txt b.txt c.txt
```
- *gzip* is the standard in UNIX:


```
> gzip a.txt b.txt c.txt # will create a.txt.gz, b.txt.gz,
c.txt.gz
```
- *tar* will nicely wrap up entire directories:


```
> tar cvf files.tar myDirectory
> tar cvfz files.tgz myDirectory
```
- To unwrap a tarball


```
> tar xvf files.tar
> tar xvfz files.tgz
```
- To peek into a zipped (text) file:


```
> gzip -cd file.gz | less
> zcat file.zip | less
```

3 A variety of UNIX tools/capabilities

Many UNIX programs are small programs (tools, utilities) that can be combined to do complicated things.

1. For help on a UNIX program, including command-line utilities like *ls*, *cp*, etc.

```
> man cp
```

2. What's the path of an executable?

```
> which R
```

3. Tools for remotely mounting the filesystem of a remote UNIX machine/filesystem

- *Samba* protocol - see "How can I mount my home directory" on SCF Help Desk FAQs

- *MacFUSE*

- Linux:

```
> cd; mkdir scf # create a directory as a mountpoint
```

```
> sshfs bilbo.berkeley.edu: /home/paciorek/scf # mount the
remote filesystem
```

```
> fusermount -u ~/scf # to unmount
```

4. Cloud storage: Dropbox and other services will mirror directories on multiple machines and on their servers

5. Doing something at the UNIX command line from within R:

```
> system("ls -al")
```

6. Editors

- traditional UNIX: *emacs*, *vi*
- Windows: *WinEdt*
- Mac: *Aquamacs Emacs*
- Be careful in Windows - file suffixes are often hidden
- For statistical computing, we need an editor, not a word processor, because we're going to be operating on files of code and data files, for which word processing formatting gets in the way.

7. Basic emacs:

- *emacs* has special modes for different types of files: R code files, C code files, Latex files – it's worth your time to figure out how to set this up on your machine for the kinds of files you often work on
 - For working with R, ESS (emacs speaks statistics) mode is helpful. This is built in to Aquamacs emacs.
- To open emacs in the terminal window rather than as a new window, which is handy when tunneling the graphical emacs window is slow or impossible:

```
> emacs -nw file.txt
```

8. Files that provide info about a UNIX machine:

- */proc/meminfo*
- */proc/cpuinfo*
- */etc/issue*
- Example: how do I find out how many processors a machine has:

```
> grep processor /proc/cpuinfo
```

9. There are (free) tools in UNIX to convert files between lots of formats (pdf, ps, html, latex, jpg). This is particularly handy when preparing figures for a publication. My [computing tips](#) page lists a number of these.

Table 1. Helpful emacs control sequences.

| Sequence | Result |
|--|---|
| C-x, C-c | Close the file |
| C-x, C-s | Save the file |
| C-x, C-w | Save with a new name |
| C-s | Search |
| ESC | Get out of command buffer at bottom of screen |
| C-a | Go to beginning of line |
| C-e | Go to end of line |
| C-k | Delete the rest of the line from cursor forward |
| C-space , then move to end of block | Highlight a block of text |
| C-w | Remove the highlighted block, putting it in the kill buffer |
| C-y (after using C-k or C-w) | Paste from kill buffer ('y' is for 'yank') |

The bash shell and UNIX utilities

September 7, 2011

Sections with more advanced material that are not critical for those of you just getting started with UNIX are denoted (***)).

Note that it can be difficult to distinguish what is shell-specific and what is just part of UNIX. Some of the material here is not bash-specific but general to UNIX.

Reference: Newham and Rosenblatt, Learning the bash Shell, 2nd ed.

1 Shell basics

The shell is the interface between you and the UNIX operating system. When you are working in a terminal window (i.e., a window with the command line interface), you're interacting with a shell.

There are multiple shells (*sh*, *bash*, *csh*, *tcsh*, *ksh*). We'll assume usage of *bash*, as this is the default for Mac OSX and on the SCF machines and is very common for Linux.

1. What shell am I using?

```
> echo $SHELL
```

2. To change to bash on a one-time basis:

```
> bash
```

3. To make it your default:

```
> chsh /bin/bash
```

/bin/bash should be whatever the path to the bash shell is, which you can figure out using **which bash**

Shell commands can be saved in a file (with extension *.sh*) and this file can be executed as if it were a program. To run a shell script called *file.sh*, you would type **./file.sh**. Note that if you just typed **file.sh**, the shell will generally have trouble finding the script and recognizing that it is executable.

2 Tab completion

When working in the shell, it is often unnecessary to type out an entire command or file name, because of a feature known as tab completion. When you are entering a command or filename in the shell, you can, at any time, hit the tab key, and the shell will try to figure out how to complete the name of the command or filename you are typing. If there is only one command in the search path and you're using tab completion with the first token of a line, then the shell will display its value and the cursor will be one space past the completed name. If there are multiple commands that match the partial name, the shell will display as much as it can. In this case, hitting tab twice will display a list of choices, and redisplay the partial command line for further editing. Similar behavior with regard to filenames occurs when tab completion is used on anything other than the first token of a command.

Note that R does tab completion for objects (including functions) and filenames.

3 Command history

By using the up and down arrows, you can scroll through commands that you have entered previously. So if you want to rerun the same command, or fix a typo in a command you entered, just scroll up to it and hit enter to run it or edit the line and then hit enter.

Note that you can use emacs-like control sequences (**C-a**, **C-e**, **C-k**) to navigate and delete characters, just as you can at the prompt in the shell usually.

You can also rerun previous commands as follows:

```
> !-n # runs the nth previous command
> !xt # runs the last command that started with 'xt'
```

If you're not sure what command you're going to recall, you can append **:p** at the end of the text you type to do the recall, and the result will be printed, but not executed. For example:

```
> !xt:p
```

You can then use the up arrow key to bring back that statement for editing or execution.

You can also search for commands by doing **C-r** and typing a string of characters to search for in the search history. You can hit return to submit, **C-c** to get out, or **ESC** to put the result on the regular command line for editing.

4 Basic UNIX utilities

Table 1 shows some basic UNIX programs, which are sometimes referred to as filters. The general syntax for a UNIX program is

Table 1. UNIX utilities.

| Name | What it does |
|-------------|---------------------------------------|
| <i>tail</i> | shows last few lines of a file |
| <i>less</i> | shows a file one screen at a time |
| <i>cat</i> | writes file to screen |
| <i>wc</i> | counts words and lines in a file |
| <i>grep</i> | finds patterns in files |
| <i>wget</i> | downloads files from the web |
| <i>sort</i> | sorts a file by line |
| <i>nl</i> | numbers lines in a file |
| <i>diff</i> | compares two files |
| <i>uniq</i> | removes repeated (sequential) rows |
| <i>cut</i> | extracts fields (columns) from a file |

> command -options argument1 argument2 ...

For example, > **grep -i graphics file.txt** looks for *graphics* (argument 1) in *file.txt* (argument2) with the option *-i*, which says to ignore the case of the letters. > **less file.txt** simply pages through a text file (you can navigate up and down) so you can get a feel for what's in it.

UNIX programs often take options that are identified with a minus followed by a letter, followed by the specific option (adding a space before the specific option is fine). Options may also involve two dashes, e.g., **R --no-save**. Here's another example that tells *tail* to keep refreshing as the file changes:

```
tail -f dat.txt
```

A few more tidbits about *grep*:

```
grep ^read code.r # returns lines that start with 'read'
```

```
grep dat$ code.r # returns lines that end with 'dat'
```

```
grep 7.7 dat.txt # returns lines with two sevens separated by a single character
```

```
grep 7.*7 dat.txt # returns lines with two sevens separated by any number of characters
```

If you have a big data file and need to subset it by line (if it's possible with *grep*) or by field (if possible with *cut*), then you can do it really fast from the UNIX command line, rather than reading it with R, SAS, Perl, etc.

Much of the power of these utilities comes in piping between them (see Section 5) and using wildcards (see Section 6) to operate on groups of files. The utilities can also be used in shell scripts to do more complicated things.

Table 2. Redirection.

| Syntax | What it does |
|---|--|
| <code>cmd > file</code> | sends stdout from <i>cmd</i> into <i>file</i> , overwriting <i>file</i> |
| <code>cmd >> file</code> | appends stdout from <i>cmd</i> to <i>file</i> |
| <code>cmd >& file</code> | sends stdout and stderr from <i>cmd</i> to <i>file</i> |
| <code>cmd < file</code> | execute <i>cmd</i> reading stdin from <i>file</i> |
| <code>cmd <infile >outfile 2>errors</code> | reads from <i>infile</i> , sending stdout to <i>file</i> and stderr to <i>errors</i> |
| <code>cmd1 cmd2</code> | sends stdout from <i>cmd1</i> as stdin to <i>cmd2</i> (a pipe) |

Note that *cmd* may include options and arguments as seen in the previous section.

5 Redirection

UNIX programs that involve input and/or output often operate by reading input from a stream known as standard input (*stdin*), and writing their results to a stream known as standard output (*stdout*). In addition, a third stream known as standard error (*stderr*) receives error messages, and other information that's not part of the program's results. In the usual interactive session, standard output and standard error default to your screen, and standard input defaults to your keyboard. You can change the place from which programs read and write through redirection. The shell provides this service, not the individual programs, so redirection will work for all programs. Table 2 shows some examples of redirection.

Operations where output from one command is used as input to another command (via the `|` operator) are known as pipes; they are made especially useful by the convention that many UNIX commands will accept their input through the standard input stream when no file name is provided to them.

Here's an example of finding out how many unique entries there are in the 2rd column of a data file whose fields are separated by commas:

```
cut -d',' -f2 mileage2009.csv | sort | uniq | wc
```

To see if there are any "M" values in certain fields (fixed width) of a set of files (note I did this on 22,000 files (5 Gb or so) in about 15 minutes on my old desktop; it would have taken hours to read the data into R):

```
> cut -b1,2,3,4,5,6,7,8,9,10,11,29,37,45,53,61,69,77,85,93,101,109,117,
    125,133,141,149,157,165,173,181,189,197,205,213,221,229,237,245,253,261,2
USC*.dly
| grep "M" | more
```

A closely related, but subtly different, capability is offered by the use of backticks (`). When the shell encounters a command surrounded by backticks, it runs the command and replaces the backticked expression with the output from the command; this allows something similar to a pipe, but is appropriate when a command reads its arguments directly from the command line instead of

Table 3. Wildcards.

| Syntax | What it matches |
|--|---|
| <code>?</code> | any single character |
| <code>*</code> | zero or more characters |
| <code>[<i>c₁c₂...</i>]</code> | any character in the set |
| <code>[!<i>c₁c₂...</i>]</code> | anything not in the set |
| <code>[<i>c₁ - c₂</i>]</code> | anything in the range from <i>c₁</i> to <i>c₂</i> |
| <code>{<i>string1, string2, ...</i>}</code> | anything in the set of strings |

through standard input. For example, suppose we are interested in searching for the text *pdf* in the last 4 R code files (those with suffix *.q*) that were modified in the current directory. We can find the names of the last 4 files which were modified using

```
> ls -t *q | head -4
```

and we can search for the required pattern using *grep*. Putting these together with the backtick operator we can solve the problem using

```
> grep pdf `ls -t *q | head -4`
```

Note that piping the output of the *ls* command into *grep* would not achieve the desired goal, since *grep* reads its filenames from the command line, not standard input.

You can also redirect output as the arguments to another program using the *xargs* utility. Here's an example:

```
> which bash | xargs chsh
```

And you can redirect output into a shell variable (see section 9) by putting the command that produces the output in parentheses and preceding with a *\$*. Here's an example:

```
> files=$(ls)
```

```
> echo $files
```

6 Wildcards in filenames

The shell will expand certain special characters to match patterns of file names, before passing those filenames on to a program. Note that the programs themselves don't know anything about wildcards; it is the shell that does the expansion, so that programs don't see the wildcards. Table 3 shows some of the special characters that the shell uses for expansion:

Here are some examples of using wildcards:

- List all files ending with a digit:

```
> ls *[0-9]
```

- Make a copy of *filename* in *filename.old*

```
> cp filename{,.old}
```

- Remove all files beginning with *a* or *z*:

```
> rm [az]*
```

- List all the R code files with a variety of suffixes:

```
> ls *{.r,.q,.R}
```

The *echo* command can be used to verify that a wildcard expansion will do what you think it will:

```
> echo cp filename{,.old} # returns cp filename filename.old
```

If you want to suppress the special meaning of a wildcard in a shell command, precede it with a backslash (`\`). Note that this is a general rule of thumb in many similar situations when a character has a special meaning but you just want to treat it as a character.

7 Job Control

Starting a job When you run a command in a shell by simply typing its name, you are said to be running in the foreground. When a job is running in the foreground, you can't type additional commands into that shell, but there are two signals that can be sent to the running job through the keyboard. To interrupt a program running in the foreground, use **C-c**; to quit a program, use **C-**. While modern windowed systems have lessened the inconvenience of tying up a shell with foreground processes, there are some situations where running in the foreground is not adequate.

The primary need for an alternative to foreground processing arises when you wish to have jobs continue to run after you log off the computer. In cases like this you can run a program in the background by simply terminating the command with an ampersand (`&`). However, before putting a job in the background, you should consider how you will access its results, since `stdout` is not preserved when you log off from the computer. Thus, redirection (including redirection of `stderr`) is essential when running jobs in the background. As a simple example, suppose that you wish to run an R script, and you don't want it to terminate when you log off. (Note that this can also be done using **R CMD BATCH**, so this is primarily an illustration).

```
> R --no-save < code.q >& code.Rout &
```

If you forget to put a job in the background when you first execute it, you can do it while it's running in the foreground in two steps. First, suspend the job using the **C-z** signal. After receiving the signal, the program will interrupt execution, but will still have access to all files and other resources. Next, issue the *bg* command, which will put the stopped job in the background.

Listing and killing jobs Since only foreground jobs will accept signals through the keyboard, if you want to terminate a background job you must first determine the unique process id (PID) for the process you wish to terminate through the use of the *ps* command. For example, to see all the R jobs running on a particular computer, you could use a command like:

```
> ps -aux | grep R # sometimes this needs to be "ps aux | grep R"
```

Among the output after the header (shown here) might appear a line that looks like this:

```
USER PID %CPU %MEM VSZ RSS TTY STAT START TIME COMMAND
```

```
paciorek 11998 97.0 39.1 1416644 1204824 pts/16 R+ Jul27 1330:01  
/usr/lib64/R/bin/exec/R
```

In this example, the *ps* output tells us that this R job has a PID of *11998*, that it has been running for 1330 minutes (!), is using 97% of CPU and 39% of memory, and that it started on July 27. You could then issue the command:

```
> kill 11998
```

or, if that doesn't work

```
> kill -9 11998
```

to terminate the job. Another useful command in this regard is *killall*, which accepts a program name instead of a process id, and will kill all instances of the named program. E.g.,

```
> killall R
```

Of course, it will only kill the jobs that belong to you, so it will not affect the jobs of other users. Note that the *ps* and *kill* commands only apply to the particular computer on which they are executed, not to the entire computer network. Thus, if you start a job on one machine, you must log back into that same machine in order to manage your job.

Monitoring jobs and memory use The *top* command also allows you to monitor the jobs on the system and in real-time. In particular, it's useful for seeing how much of the CPU and how much memory is being used, as well as figuring out a PID as an alternative to *ps*. You can also renice jobs (see below) and kill jobs from within *top*: just type *r* or *k*, respectively, and proceed from there.

One of the main things to watch out for is a job that is using close to 100% of memory and much less than 100% of CPU. What is generally happening is that your program has run out of memory and is using virtual memory on disk, spending most of its time writing to/from disk, sometimes called *paging*. If this happens, it can be a very long time, if ever, before your job finishes.

Nicing a job (IMPORTANT) The most important thing to remember when starting a job on a machine that is not your personal machine is how to be a good citizen. This often involves 'nicing' your jobs, and this is REQUIRED on the SCF machines (although the compute servers automatically nice your jobs). Nicing a job puts it at a lower priority so that a user working at the

keyboard has higher priority in using the CPU. Here's how to do it, giving the job a low priority of 19, as required by SCF:

```
> nice -19 R CMD BATCH --no-save in.R out.Rout
```

If you forget and just submit the job without nicing, you can reduce the priority by doing:

```
> renice +19 11998
```

where *11998* is the PID of your job. For short jobs (less than 5 minutes), it's ok to use a priority of 18 rather than 19.

On many larger UNIX cluster computers, all jobs are submitted via a job scheduler and enter a queue, which handles the issue of prioritization and jobs conflicting. Syntax varies by system and queueing software, but may look something like this for submitting an R job:

```
> bsub -q long R CMD BATCH --no-save in.R out.Rout
```

8 Aliases (***)

Aliases allow you to use an abbreviation for a command, to create new functionality or to insure that certain options are always used when you call an existing command. For example, I'm lazy and would rather type **q** instead of **exit** to terminate a shell window. You could create the alias as follow

```
> alias q="exit"
```

As another example, suppose you find the *-F* option of *ls* (which displays / after directories, * after executable files and @ after links) to be very useful. The command

```
> alias ls="ls -F"
```

will insure that the *-F* option will be used whenever you use *ls*. If you need to use the unaliased version of something for which you've created an alias, precede the name with a backslash (\). For example, to use the normal version of *ls* after you've created the alias described above, just type

```
> \ls
```

The real power of aliases is only achieved when they are automatically set up whenever you log in to the computer or open a new shell window. To achieve that goal with aliases (or any other bash shell commands), simply insert the commands in the file *.bashrc* in your home directory.

9 Shell Variables (***)

We can define shell variables that will help us when writing shell scripts. Here's an example of defining a variable:

```
> name="chris"
```

The shell may not like it if you leave any spaces around the = sign. To see the value of a variable

we need to precede it by \$:

```
> echo $chris
```

You can also enclose the variable name in curly brackets, which comes in handy when we're embedding a variable within a line of code to make sure the shell knows where the variable name ends:

```
> echo ${chris}
```

There are also special shell variables called environment variables that help to control the shell's behavior. These are generally named in all caps. Type **env** to see them. You can create your own environment variable as follows:

```
> export NAME="chris"
```

The *export* command ensures that other shells created by the current shell (for example, to run a program) will inherit the variable. Without the *export* command, any shell variables that are set will only be modified within the current shell. More generally, if one wants a variable to always be accessible, one would include the definition of a variable with an *export* command through your *.bashrc* file.

Here's an example of modifying an environment variable:

```
> export CDPATH=.:~/research:~/teaching
```

Now if you have a subdirectory *bootstrap* in your *research* directory, you can type **cd bootstrap** no matter what your *pwd* is and it will move you to *~/research/bootstrap*. Similarly for any subdirectory within the *teaching* directory.

Here's another example of an environment variable that puts the username, hostname, and *pwd* in your prompt. This is handy so you know what machine you're on and where in the filesystem you are.

```
> export PS1="\u@\h:\w> "
```

For me, this is one of the most important things to put in my *.bashrc* file. The ** syntax tells bash what to put in the prompt string: *u* for username, *h* for hostname, and *w* for working directory.

10 Functions (***)

You can define your own utilities by creating a shell function. This allows you to automate things that are more complicated than you can do with an alias. Here's an example from my *.bashrc* that uses *ssh* to locally mount (on my desktop machine) remote filesystems for several systems I have access to:

```
function mounts() {
    sshfs hsph.harvard.edu: /home/paciorek/hsph
    sshfs hpcc.sph.harvard.edu: /home/paciorek/hpcc
}
```

```
sshfs bilbo.berkeley.edu: /home/paciorek/scf
}
```

Now I just type `> mounts` and the shell will execute the code in the function. (Note: this could be done with an alias by just separating the three items with semicolons.)

One nice thing is that the shell automatically takes care of arguments for you. It places the arguments given by the user into local variables in the function called (in order): `$1 $2 $3` etc. It also fills `$#` with the number of arguments given by the user. Here's an example of using arguments in a function that saves me some typing when I want to copy a file to the SCF filesystem:

```
function putscf() {
    scp $1 paciorek@bilbo.berkeley.edu:~/$2
}
```

To use this function, I just do the following to copy `unit1.pdf` from the current directory on whatever non-SCF machine I'm on to the directory `~/teaching/243` on SCF:

```
> putscf unit1.pdf teaching/243/.
```

Of course you'd want to put such functions in your `.bashrc` file.

11 If/then/else (***)

We can use if-then-else type syntax to control the flow of a shell script. For an example, see `niceR()` in the demo code file for this unit.

For more details, look in Newham&Rosenblatt or search online.

12 For loops (***)

For loops in shell scripting are primarily designed for iterating through a set of files or directories. Here's an example:

```
for file in $(ls *txt)
do
    mv $file ${file/.txt/.q} # this syntax replaces .txt with .q
in $file
done
```

One use of for loops is automating file downloads: see the demo code file. And, in my experience, for loops are very useful for starting a series of jobs: see the demo code file.

13 How much shell scripting should I learn?

You can do a fair amount of what you need from within R using the *system()* function. This will enable you to avoid dealing with the shell programming syntax (but you'll still need to know how to use UNIX utilities, wildcards, and pipes to be effective). Example: a fellow student in grad school programmed a tool in R to extract concert information from the web for bands appearing in her iTunes library. Not the most elegant solution, but it got the job done (until Apple included that functionality in iTunes).

For more extensive shell programming, it's probably worth learning Python or Perl and doing it there rather than using a shell script.

Introduction to R

September 11, 2011

References:

- Adler
- Chambers
- Murrell, Introduction to Data Technologies.
- [R intro manual](#) on CRAN (R-intro).
- Venables and Ripley, Modern Applied Statistics with S

If you don't have a basic familiarity with R, you should work through the tutorial/sample session in either the Adler book, in Venables & Ripley, or in Appendix A of R-intro.

I'm going to try to refer to R syntax as statements, where a statement is any code that is a valid, complete R expression. I'll try not to use the term *expression*, as this actually means a specific type of object within the R language.

1 Some basic functionality

- To figure out what directory R is treating as the pwd, and to set it to something different:

```
> getwd()  
> setwd("~/teaching/243")
```

- To run the R code in a file, do

```
> source("code.r")
```

- To print graphical output to a file

```
> x = 1:3; y=1:3; plot(x, y)
```

```
> dev.copy2pdf(file = "myplot.pdf", width = 5, height = 4)
```

or

```
> pdf("myplot.pdf", width = 5, height = 4)
```

```
> x = 1:3; y=1:3; plot(x, y)
```

```
> dev.off()
```

- To open a new graphics window: *quartz()* [Mac], *windows()* [Windows], *X11()* [UNIX], or *dev.new()* [any OS]. To switch, use *dev.set()*, e.g., **dev.set(3)**, and to close the active device, use *dev.off()*.

- To save your current workspace (to the file *.RData* by default) (in a form that will transfer across OS)

```
> save.image()
```

To save only the objects you specify

```
> save(myMat, nIts, file = "ps1.RData")
```

I tend to run R as **R --no-save** and only save objects or the workspace manually, to avoid having large hidden files. Then to load stuff back in:

```
> load("ps1.RData")
```

- To run UNIX commands, use *system()*, e.g.,

```
> system("ls -al")
```

- There are also a bunch of functions that will do specific queries of the filesystem, including

```
> file.exists("file.txt")
```

```
> list.files("~/research")
```

- To get some info on the system you're running on:

```
> Sys.info()
```

- To see some of the options that control how R behaves, try

```
> options()
```

One useful option to change is the number of characters of width printed to the screen:

```
> options(width = 120)
```

- You can also prevent too much output from being printed to the screen when printing a large object


```
> options(max.print = 5000)
```
- Use **C-c** to interrupt execution. This will generally back out gracefully, returning you to a state as if the command had not been started. Note that if R is exceeding memory availability, there can be a long delay. This can be frustrating, particularly since a primary reason you would want to interrupt is when R runs out of memory.
- The [R mailing list archives](#) are very helpful for getting help - always search the archive before posting a question.
- Any code that you wanted executed automatically when starting R can be placed in `~/.Rprofile` (or in individual `.Rprofile` files in specific directories). This could include loading packages (see below), sourcing files with user-defined functions (or put the function code itself in `.Rprofile`), assigning variables, and specifying options via `options()`.

2 Types of objects in R

Everything in R is stored as an object, each of which has a class that describes what the object contains and standard functions that operate on objects in the class. The basic classes are

- *character* vectors: these are sets of character values (i.e., strings). Examples of character values are: “*Sam*”, “*0923*”, “*Sam9*”, “*Sam is*”, “*Sam\t9*”, “*Sam’s the man.\nNo doubt.\n*”. Each of these is a character vector of length 1.
- *numeric* vectors (i.e., double precision real numbers)
- *integer* vectors
- *logical* (T/F) vectors
- *complex* vectors
- *lists*: vectors of arbitrary objects
- *factors*: vector-like objects of categorical variables with a pre-defined set of labels

Notes:

- Scalars are actually vectors of length 1 in R.

- Unlike in compiled languages, objects can be created without explicitly initializing them and allocating memory.
- Factors can be ordered - see *ordered()*.

More complicated objects:

- Data frames can be thought of as a list of vectors of the same length, where the vectors may have different types.
- Matrices are different from data frames in that all the elements are of a single type. Furthermore, matrices and arrays (which allow dimensions of 1, 2, 3, ...) can be thought of as vectors with information on dimensions (layout). If you pass a matrix into a function expecting a vector, it will just treat it as a vector (of concatenated columns).

We can check on the class of an object using *class()* or with specific boolean queries: *is.numeric()*, *is.factor()*, *is.vector()*, etc.

We convert between vector classes using variants on *as()*: e.g.,

```
> as.character(c(1, 2, 3))
> as.numeric(c("1", "2.73"))
> as.factor(c("a", "b", "c"))
```

Some common conversions are converting numbers that are being interpreted as characters into actual numbers, converting between factors and characters, and converting between logical T/F vectors and numeric 1/0 vectors. In some cases R will automatically do conversions behind the scenes in a smart way. In other languages, converting between different classes is sometimes called *casting* a variable.

3 Data structures

Assignment We assign into an object using either '=' or '<='. One situation in which you want to use '<=' is if it is being used as part of an argument to a function, so that R realizes you're not indicating one of the function arguments, e.g.:

```
> apply(mat, 1, sum.isna <- function(vec) {return(sum(is.na(vec)))})
> apply(mat, 1, sum.isna = function(vec) {return(sum(is.na(vec)))})
# NOPE
> mean(x <- c(3,7)) # this is allowable, though perhaps not useful;
what does it do?
```

A rule of thumb is that for basic assignments where you have an object name, then the assignment operator, and then some code, '=' is fine, but otherwise use '<-'.

```
> x[4] = mean(rnorm(7)) # OK
> system.time(x = rnorm(10000)) # NOT OK
```

Type vs. class We've discussed vectors as the basic data structure in R, with character, integer, numeric, etc. classes. Objects also have a type, which relates to what type of values are in the objects and how objects are stored internally in R (i.e., in C).

Let's look at Adler's Table 7.1 to see some other types.

```
> a = data.frame(x = 1:2)
> class(a)
> typeof(a)
> m = matrix(1:4, nr = 2)
> class(m)
> typeof(m)
```

Type coercion occurs when you use the *as()* function as well as implicitly at times. We'll see implicit coercion when we read in characters into R using *read.table()* - there was an automatic coercion to a factor variable. Consider also this example:

```
> x = rnorm(5)
> x[3] = 'hat' # What do you think is going to happen?
```

We've said that everything in R is an object and all objects have a class. For simple objects class and type are often closely related, but this is not the case for more complicated objects. The class describes what the object contains and standard functions associated with it. In general, you mainly need to know what class an object is rather than its type. Classes can inherit from other classes; for example, the *glm* class inherits characteristics from the *lm* class. We'll see more on the details of object-oriented programming in the R programming unit.

We can create objects with our own defined class.

```
> me = list(firstname = 'Chris', surname = 'Paciorek')
> class(me) = 'person'
> typeof(me)
> typeof(me$firstname)
```

Information about objects Some functions that give information about objects are:

```
> is(me, 'person')
> is.list(me)
> str(me)
```

```

> attributes(me)
> mat = matrix(1:4, 2)
> class(mat)
> typeof(mat)
> length(mat) # recall that a matrix can be thought of as a vector
with dimensions
> attributes(mat)
> dim(mat)

```

Attributes are information about an object attached by name (in the form of a list) to the object. Attributes are often copied when operating on an object. This can lead to some weird-looking formatting:

```

> x = rnorm(10 * 365)
> qs = quantile(x, c(.025, .975))

```

If you now do subsequent operations with *qs*, the *names* attribute will often get carried along. We can get rid of it:

```

> names(qs) = NULL

```

A common use of attributes is that rows and columns may be named in matrices and data frames, and elements in vectors:

```

> row.names(mtcars)
> names(mtcars)
> mat = data.frame(x = 1:2, y = 3:4)
> row.names(mat) = c("first", "second")

```

Special objects There are also some special objects, which often begin with a period, like hidden files in UNIX. One is *.Last.value*, which stores the last result.

```

> rnorm(10)
> .Last.value

```

The workspace Objects exist in a workspace.

```

> objects() # what objects are in my workspace
> ls() # synonymous
> dat = 7; dat2 = 9; subdat = 3
> objects(pattern = "^dat")
> rm(dat2, subdat)
> rm(list = c("dat2", "subdat")) # a bit confusing - the 'list'
argument should be a character vector

```

```

> rm(list = ls(pattern = "^dat"))
> exists('dat') # can be helpful when programming
> is.null(dat) # can be helpful for checking if a function argument
was supplied by the user
> dat = rnorm(500000)
> object.size(dat)
> print(object.size(dat), units = "Mb") # this seems pretty clunky!

```

Information about functions To get help on functions:

```

> ?lm # or help(lm)
> ??lm # or help.search('lm')
> apropos(lm)
> ? '[' # operators are functions too - you don't need the space
here, but I'm having trouble with my Latex file...
> args(lm)

```

Strings and quotation Working with strings and quotes (see **?Quotes**). Generally one uses double quotes to denote text. If we want a quotation symbol in text, we can do something like the following, either combining single and double quotes or escaping the quotation:

```

> ch1 = "Chris's"
> ch2 = 'He said, "hello."'
> ch3 = "He said, \"hello.\""

```

4 Packages

One of the killer apps of R is the extensive collection of add-on packages on **CRAN** (www.cran.r-project.org) that provide much of R's functionality. To make use of a package it needs to be installed on your system (using *install.packages()* once only) and loaded into R (using *library()* every time you start R).

Here are some helpful commands related to packages .

To make a package available (loading it) and to see all the installed packages that you could load:

```

> library(help = gam)
> library()

```

Notice that some of the packages are in a system directory and some are in my home directory.

If a package is on CRAN but not on your system, you can install it easily (usually). You don't need root permission on a machine to install a package (though sometimes you run into hassles if you are installing it just as a user, so if you have administrative privileges it may help to use them).

```
> install.packages("fields", lib = "~/Rlibs") # ~/Rlibs needs to exist!
```

The objects in a package (primarily functions, but also data) are in their own workspaces, and are accessible after you load the package using *library()*, but are not directly visible when you use *ls()*. We'll talk more about this when we talk about scope and environments. If we want to see the objects in one of the other workspaces, we can do the following:

```
> search()
> ls(pos = 3) # for the graphics package
```

Some packages are *installed* by default with R and of these, some are *loaded* by default, while others require a call to *library()*. For packages I use a lot, I install them once and then load them automatically every time I start R using my *~/.Rprofile* file.

Data Manipulation and Management

September 16, 2011

While we focus on working with data in R, note that concepts of reading in data, database manipulations, missing values, and regular expressions are common to other languages, so familiarity with these should allow you to pick up other tools more easily. The main downside to working with datasets in R is that the entire dataset resides in memory, so R is not so good for dealing with very large datasets. More on alternatives in a bit. Another common frustration is controlling how the variables are interpreted (numeric, character, factor) when reading data into a data frame.

References:

- Murrell, Introduction to Data Technologies.
- [R Data Import/Export manual](#) on CRAN.
- Chambers, Software for Data Analysis
- Venables and Ripley, Modern Applied Statistics with S

1 Shell commands for pre-processing

We've already covered some helpful UNIX utilities - in particular *grep* and *cut* are particularly useful for selecting rows and columns. If you're doing a lot of pre-processing of complicated input data, *Perl* or *Python* are powerful scripting languages with extensive capabilities in manipulating text.

2 Data storage and formats (outside R)

R has the capability to read in a wide variety of file formats. Let's get a feel for some of the common ones.

1. Flat text files (ASCII files): data are often provided as simple text files. Often one has one record or observation per row and each column or field is a different variable or type of information about the record. Such files can either have a fixed number of characters in each field (fixed width format) or a special character (a delimiter) that separates the fields in each row. Common delimiters are tabs, commas, one or more spaces, and the pipe (`|`). Common file extensions are `.txt` and `.csv`. Metadata (information about the data) is often stored in a separate file. I like CSV files but if you have files where the data contain commas, tab-delimited can be good (or text can be put in quotes in CSV files). A couple more intricate details:
 - (a) Sometimes if you have a text file created in Windows, the line endings are coded differently than in UNIX (`\r\n` in Windows vs `\n` in UNIX); `dos2unix` is a UNIX utility that can do the necessary conversion. If you see `^M` at the end of the lines in a file, that's the tool you need. Alternatively, if you open a UNIX file in Windows, it expects a different coding for line endings and may treat all the lines as a single line. You can fix this with `unix2dos`.
 - (b) Another difficulty can be dealing with the text encoding (the mapping of individual characters (including tabs, returns, etc.) to a set of numeric codes) when the file is not simply an ASCII file (i.e., contains non-ASCII characters - stuff that's not on a US-style keyboard). There are a variety of different encodings for text files (ASCII is the most basic encoding, and allows for 2^7 different characters; others allow for 2^{16}), with different ones common on different operating systems. The UNIX utility `file`, e.g.

```
> file tmp.txt
```

 can help provide some information, `read.table()` in R takes an argument `fileEncoding` and the UNIX utility `iconv` can help with conversions.
2. In some contexts, such as textual data and bioinformatics data, the data may in a text file with one piece of information per row, but without meaningful columns/fields.
3. In scientific contexts, netCDF (`.nc`) (and the related HDF5) are popular format for gridded data that allows for highly-efficient storage and contains the metadata within the file. The basic structure of a netCDF file is that each variable is an array with multiple dimensions (e.g., latitude, longitude, and time), and one can also extract the values of and metadata about each dimension. The `ncdf` package in R nicely handles working with netCDF files. These are examples of a binary format, which is not (easily) human readable but can be more space-efficient and faster to work with (because they can allow random access into the data rather than requiring sequential reading).

4. Data may also be in the form of XML or html files. We won't deal with these in 243, but it may be the topic of one of the 244 modules in coming years.
5. Data may already be in a database or in the data storage of another statistical package (*Stata*, *SAS*, *SPSS*, etc.). The *foreign* package in R has excellent capabilities for importing *Stata* (*read.dta()*), *SPSS* (*read.spss()*), *SAS* (*read.ssd()*) and, for XPORT files, *read.xport()*, and *dbf* (a common database format) (*read.dbf()*), among others.
6. For Excel, there are capabilities to read an Excel file (see the *XLConnect* package among others), but you can also just go into Excel and export as a CSV file or the like and then read that into R. In general, it's best not to pass around data files as Excel or other spreadsheet format files because (1) someone may not have Excel, (2) Excel imposes limits on the number of rows, (3) one can easily manipulate text files such as CSV using UNIX tools, but this is not possible with an Excel file, (4) Excel files often have more than one sheet, graphs, macros, etc., so they're not a data storage format per se.

3 Reading data from text files into R

We can use *read.fwf()* to read from a fixed width text file into a data frame. *read.table()* is probably the most commonly-used function for reading in data, it reads in delimited files (*read.csv()* and *read.delim()* are special cases of *read.table()*). The key arguments are the delimiter (the *sep* argument) and whether the file contains a header, a line with the variable names.

The most difficult part of reading in such files can be dealing with how R determines the classes of the fields that are read in. There are a number of arguments to *read.table()* and *read.fwf()* that allow the user to control the classes. One difficulty is that character and numeric fields are sometimes read in as factors. Basically *read.table()* tries to read fields in as numeric and if it finds non-numeric and non-NA values, it reads in as a factor. This can be annoying.

Let's work through an example in the demo code file. First let's look at the arguments to *read.table()*. Note that *sep=""* separates on any amount of white space.

The basic function *scan()* simply reads everything in, ignoring lines, which works well and very quickly if you are reading in a numeric vector or matrix. *scan()* is also useful if your file is free format - i.e., if it's not one line per observation, but just all the data one value after another; in this case you can use *scan()* to read it in and then format the resulting character or numeric vector as a matrix with as many columns as fields in the dataset. Remember that the default is to fill the matrix by column.

If the file is not nicely arranged by field (e.g., if it has ragged lines), we'll need to do some more work. *readLines()* will read in each line into a separate character vector, after which we can

process the lines using text manipulation. Let's work through another example in the demo code file.

If a file is large, we may want to read it in in chunks (of lines), do some computations to reduce the size of things, and iterate. `read.table()`, `read.fwf()` and `readLines()` all have the arguments that let you read in a fixed number of lines. The problem then is how we can read in the subsequent blocks beyond the first one. The answer is to set up something called a *connection*.

Basically R allows you to read in not just from a file but from a more general construct. Here are some examples of connections:

```
> dat = readLines(pipe("ls -al"))
> dat = readLines(pipe("gunzip dat.gz"))
> dat = readLines(gzfile("dat.gz"))
> dat = read.table("http://www.berkeley.edu/users/dee/a.txt")
```

To read-on-the-fly in blocks, we need to first establish the connection and then read from it sequentially. See the demo code file for an example.

One cool trick that can come in handy is to create a text connection.

```
> con = textConnection(myCharVec, "r")
```

This lets you 'read' from an R character vector as if it were a text file and could be handy for processing text. For example, you could then use `read.fwf()` applied to `con`.

We can create connections for writing output too. Just make sure to open the connection first.

Be careful with the directory separator in Windows files: you can either do `"C:\\mydir\\file.txt"` or `"C:/mydir/file.txt"`, but not `"C:\mydir\file.txt"`.

4 Text manipulations and regular expressions

Text manipulations in R have a number of things in common with Perl, Python and UNIX, as many of these evolved from UNIX. When I use the term *character* or *string* here, I'll be referring to any string of characters that may include numbers, white space, and special characters, rather than to the character class of R objects.

A few of the basic R functions for manipulating strings are `paste()`, `strsplit()`, and `substr()`. `paste()` and `strsplit()` are basically inverses of each other: `paste()` concatenates together an arbitrary set (or a vector, if using the `collapse` argument) with a user-specified separator character, while `strsplit()` splits apart based on a delimiter/separator. `substr()` splits apart a string of characters based on fixed widths. Note that all of these operate in a vectorized fashion.

```
> paste("My", "name", "is", "Chris", ".", sep = " ")
> paste(c("My", "name", "is", "Chris", "."), collapse=" ")
```

`nchar()` tells the number of characters in a string.

Let's consider the different ways to identify substrings: (1) match based on position, using `substr()`, (2) match based on a particular text string, using R's `grep()` function, and (3) match based on regular expressions using `grep()` (covered next). `grep()` will look for a specified character string and report back indices identifying the elements of the vector the string was found in (using the `fixed=TRUE` argument ensures that regular expressions are NOT used). `gsub()` can be used to replace a specified character string with a replacement string (use `sub()` if you only want to replace only the first occurrence).

```
> vars = c("date98", "size98", "x98weights98", "sdfsd")
> grep("98", vars)
> gsub("98", "04", vars)
```

Regular expressions The `grep()` and `sub()` functions are more powerful when used with regular expressions. Regular expressions are a domain-specific language for finding patterns and are one of the key functionalities in scripting languages such as Perl and Python. We'll just cover their use in R, but once you know that, it would be easy to use them elsewhere. The basic idea of regular expressions is that they allow us to find matches of strings or patterns in strings, as well as do substitution. Regular expressions are good for tasks such as:

- extracting pieces of text - for example finding all the links in an html document;
- creating variables from information found in text;
- cleaning and transforming text into a uniform format;
- mining text by treating documents as data; and
- scraping the web for data.

Regular expression are constructed from three things:

Literal characters are matched only by the characters themselves,

Character classes are matched by any single member in the class, and

Modifiers operate on either of the above or combinations of them.

Note that the syntax is very concise, so it's helpful to break down individual regular expressions into the component parts to understand them. Here's the key syntax in regular expression (you'll notice some of this is quite similar to the wildcards we discussed for UNIX, but much more flexible).

This syntax can be put together in arbitrarily-complicated ways. For example, `"\\.[*\\]"` would search for any number of characters (.) separated by square brackets. What does this match: `^[^:lower:]]+ $` ? Let's practice the basic syntax for matching.

Table 1. Regular expression syntax.

| Syntax | What it matches |
|--------------------------|--|
| <code>^ab</code> | match 'ab' at the beginning of the string |
| <code>ab\$</code> | match 'ab' at the end of the string |
| <code>[abc]</code> | match a or b or c anywhere |
| <code>[\t]</code> | match a space or a tab |
| <code>(ab cd def)</code> | match any of the strings in the set |
| <code>(ab){2,9}</code> | match 'ab' repeated at least 2 and no more than 9 times |
| <code>(ab){2,}</code> | match 'ab' repeated 2 or more times |
| <code>[0-9a-z]</code> | match a digit or a lower-case alphabetical |
| <code>[^0-9]</code> | match anything except a digit |
| <code>a.b</code> | match a and b separated by a single character |
| <code>a.*b</code> | match a and b separated by any number of (or no) characters |
| <code>a.+b</code> | like a.*b but must have at least one character in between |
| <code>[:digit:]</code> | match <i>digit</i> class; other classes are <i>alpha</i> , <i>alnum</i> , <i>lower</i> , <i>upper</i> , <i>punct</i> , <i>blank</i> , <i>space</i> |
| <code>\\</code> | double backslashes are used if we want to search for a meta-character used in regexp syntax |

We can use the regexp syntax as a pattern to look for in `grep()` or a pattern to replace in `gsub()`.

```
> grep("[ \t]", c("john", "jennifer pierce", "juan carlos rey"))
> gsub("^j", "J", c("john", "jennifer pierce", "juan carlos rey"))
```

`regexpr()` is like `grep()`, but for each element of a vector, it reports back if the pattern was found and where in each string the first instance of the pattern was found (use `gregexpr()` to find them all).

```
> regexpr("[ \t]", c("john", "jennifer pierce", "juan carlos rey"))
```

`regexpr()` is particularly useful in seeing `_what_` was matched to help in understanding and learning regexp syntax. Note that if you know Perl regular expressions, you can add the argument `perl=TRUE` to ensure the regexp syntax is compatible with Perl regular expressions.

We can also refer back to any pattern in parentheses using the syntax `\\1` (or `\\2`, `\\3` for the second, third, etc.) patterns given in parentheses. So parentheses are used both for grouping and for extraction. Here's an example of removing the (single) space around arbitrary strings by substituting in whatever is in the brackets for the entire string including the brackets:

```
> gsub(" (.*) ", "\\1", myString)
```

Suppose `myString="safdsf df sdfds asdfd"`. It turns out the pattern matching is 'greedy' - it looks for the longest match possible. So instead, let's try

```
> gsub(" ([:alnum:])* ", "\\1", myString)
```

Regular expression can be used in a variety of places. E.g., to split by any number of white space characters

```
> line = "a dog\tjumped\nover \t"
> cat(line)
> strsplit(line, split = "[[:space:]]+")
> strsplit(line, split = "[[:blank:]]+")
```

Challenge: how would we proceed if we wanted to switch the order of two numeric substrings in a given string?

Example: let's work through some code to extract IP addresses from text (see demo code). IP addresses have four numbers (each 1-3 digits) separated by periods.

Challenge: How would I search for email titles that contain punctuation or numbers within words, like "Vi@gra" or "Viagra" or "Fancy repl!c@ted watches"? What might we have to be careful to exclude?

The meta-characters that require a double backslash when used literally are: `^ $ + ? () [] { }` `| \` We have to use two backslashes instead of a single backslash because R uses a single backslash to symbolize certain control characters, such as `\n` for newline. We call `\\` an escape sequence.

As Murrell notes, since regexp are their own language, it's a good idea to build up a regexp in pieces as a way of avoiding errors just as we would with any computer code.

Duncan Temple Lang (UC Davis Stats) has written a nice tutorial that I've posted on bspace (*regexpr-Lang.pdf*) or check out Sections 9.9 and 11 of Murrell.

Manipulating dates One common form of text information is dates and times. These can be a hassle because of the variety of ways of representing them, issues of time zones and leap years, the complicated mapping between dates and day of week, and the irregularity of the number of days in a month. We can use `as.Date()` to convert from text to an object of class `Date`, providing a character string such as `"%Y-%m-%d"` or `"%d/%m/%Y"` as the format for R to expect.

```
> date = as.Date("03-01-2011", format = "%m-%d-%Y")
> date = as.Date("03/01/11", format = "%m/%d/%Y")
> date = as.Date("03-Feb-11", format = "%d-%b-%Y")
```

With an object in the `Date` class, you can get information with functions such as `weekdays()`, `months()`, `julian()`. Note that the Julian date is the number of days since a fixed initial day (in R's case this defaults to Jan 1, 1970). One handy feature is to use `as.Date()` to convert from Julian date to a more informative format:

```
> dates = as.Date(julianValues, origin="1990-01-30")
```

One can also use the `POSIXlt` and `POSIXct` classes to work with dates and times – see **?DateTimeClasses**.

The *chron* package provides classes for dates and times and for manipulating them. Here's some *chron* syntax:

```
> d1 = chron("12/25/2004", "11:37:59") # default format of m/d/Y  
and h:m:s
```

You can use *as.chron()* to convert from a *Date* object.

5 Missing values

The basic missing value token in R is *NA*. Most functions (and any software you write) handle NAs as input gracefully. The user can generally provide input on how to handle NAs, with omitting them being a common desire (e.g., `> mean(x, na.rm = TRUE)`) One checks whether a value is an NA using *is.na()*. *NaN* indicates an operation has returned something that is not a number, such as 0/0, while *Inf* and *-Inf* indicate infinity and can be used in calculations.

```
> 3 / Inf  
> 3 + Inf
```

To check for NA or NaN, use *is.na()* and *is.nan()*. `==` won't work for `x == NA` and may cause errors for `x == NaN`.

6 Database manipulations

6.1 Databases in general

A relational database stores data as a database of tables (or relations), which are rather similar to R data frames, in that they are made up of columns or fields of one type (numeric, character, date, currency, ...) and rows or records containing the observations for one entity. One principle of databases is that if a category is repeated in a given variable, you can more efficiently store information about each level of the category in a separate table; consider information about people living in a state and information about each state - you don't want to include variables that only vary by state in the table containing information about individuals (at least until you're doing the actual analysis that needs the information in a single table. You can also interact with databases from a variety of database systems (DBMS=database management system) (some systems are SQLite, MySQL, Oracle, Access). We'll concentrate on accessing data in a database rather than management of a database.

The *DBI* package provides a front-end for manipulating databases from a variety of DBMS (MySQL, SQLite, Oracle, among others). Basically, you tell the package what's DBMS is being used on the backend, link to the actual database, and then you can use the syntax in the package.

Here's an example of using *DBI* to interact with a database called *myDatabase* (containing the table *MarsData*) using the SQLite system:

```
> drv = dbDriver("SQLite")
> con = dbConnect(drv, dbname = "myDatabase") # so we're using a
connection again
> mars = dbReadTable(con, "MarsData")
> dbDisconnect(con)
```

SQL is a common database query language, with variants on it used in the various database systems. It has statements like:

```
SELECT Id, City, State FROM Station WHERE Lat_N > 39.7 AND Lon_W
> 80 ORDER BY Lon_W
```

You can use DBI to send a database query from R rather than reading in the whole table:

```
> myQuery = "SELECT City, State FROM Station WHERE Lat_N > 39.7"
> marsSub = dbGetQuery(con, myQuery)
```

You can also use *dbSendQuery()* combined with *fetch()* to pull in a fixed number of records at a time, if you're working with a big database. Finally you can do database joins (i.e., merges) using SQL, but we won't go into the syntax here.

SAS can also be used for database management and manipulation and is quite fast for working with large datasets, storing them on disk rather than in memory. Some statisticians (including myself) often do database pre-processing in SAS and then output to R for analysis.

6.2 Database manipulations in R

Ok, let's assume you've gotten your data into an R data frame. It might be numeric, character, or some combination. R now allows you to do a variety of database manipulations.

Subsetting For subsetting, we generally use square bracket indexing and the power of logical subsetting:

```
> state = as.data.frame(state.x77)
> subdata = state[state$Population > 10000 & state$Frost > 100,
]
```

but one can also use the *subset()* function (this can be good for undergrad classes). To subset based on matching a field against multiple values, you can do:

```
> subdata = state[row.names(state) %in% smallStates, ] # where smallStates
is a vector of states you want included
```

This is called logical subscripting because **row.names(state) %in% smallStates** re-

turns a logical vector (of T/Fs). Also, you can find unions and intersections of elements of two vectors with *union()* and *intersect()*. This can be useful when merging datasets. One type of subsetting that is not as widely known is to subset with a matrix:

```
> mat = matrix(1:9, 3)
> ids = matrix(c(1, 1, 2, 2, 3, 1), nr = 3, byrow =TRUE)
> mat[ids]
```

You can also subset with a vector of row and column indices to extract blocks or arbitrary rows and columns:

```
> mat[1:2, 1:2]
> mat[c(1, 3), c(1, 3)]
```

Getting information about variable(s) One can get a variety of information about a variable:

```
> sum(mtcars$cyl == 6) # counts the number of TRUE values
> mean(mtcars$cyl == 6) # proportion of TRUE values
> unique(mtcars$cyl)
> length(unique(mtcars$cyl))
> duplicated(mtcars$cyl) # tells us which elements are repeated
> table(mtcars$gear) # tabulates # of records in each unique value
of mtcars$gear
> table(mtcars$cyl, mtcars$gear) # cross-tabulation by two variables
```

The output from *table()* can be hard to read, but *fable()* can help. Compare

```
> table(mtcars[c("cyl", "vs", "am", "gear")])
> ftable(mtcars[c("cyl", "vs", "am", "gear")])
```

Sorting Dataset sorting in R is a bit of a hassle, as the *sort()* function only sorts a single vector. Instead we rely on the *order()* function which gives the set of indices required to order a given vector (or multiple vectors, with ties within the first vector broken based on subsequent vectors) in either increasing or decreasing order. Here's a one-liner that does a sort based on giving the function a set of one or more vectors (usually columns of the dataset) to be used for sorting:

```
> f.sort = function(data, ...){
+ return(data[order(...), ])
```

Merging/joining Database joins are used to combine information from multiple database tables. In R, this is done with the *merge()* function. Let's look at an example in the demo code file. We'll also see that the output of *merge()* can have an arbitrary ordering, which may be unappealing.

Long and wide formats Finally, we may want to convert between so-called 'long' and 'wide' formats, which are motivated by working with longitudinal data (multiple observations per subject). The wide format has repeated measurements for a subject in separate columns, while the long format has repeated measurements in separate rows, with a column for differentiating the repeated measurements. *stack()* converts from wide to long while *unstack()* does the reverse. *reshape()* is similar but more flexible and it can go in either direction. Let's use the precipitation data as an example (see demo code file). The wide format is useful for doing separate analyses by group, while the long format is useful for doing a single analysis that makes use of the groups, such as ANOVA or mixed models. Check out *melt()* and *cast()* in the *reshape* package for easier argument formats than *reshape()*.

Working with factors You can create a factor with the *cut()* function:

```
> x = rnorm(100)
> f = cut(x, breaks = c(-Inf, -1, 1, Inf), labels = c('low', 'medium',
'high'))
> levels(f) # note that f is not ordered
> f = relevel(f, 'low') # puts low as first level
> f = relevel(f, 'medium') # puts medium as first and moves others
up
> f = relevel(f, 'high') # puts high as first and moves others up
```

Changing the order of the levels can be helpful for controlling the order of plotting of levels and for controlling the baseline category in an ANOVA setting, though in this case it's a bit strange to have *high* be the lowest level.

7 Output

Functions for text output are generally analogous to those for input. *write.table()*, *write.csv()*, and *writeLines()* are analogs of *read.table()*, *read.csv()*, and *readLines()*. *write()* can be used to write a matrix to a file, specifying the number of columns desired. *cat()* can be used when you want fine control of the format of what is written out and allows for outputting to a connection (e.g., a file).

One thing to be aware of when writing out numerical data is how many digits are included. For example, the default with *write()* and *cat()* is the number of digits displayed to the screen, controlled by **options()\$digits**. (to change this, do **> options(digits = 5)** or specify as an argument to *write()* or *cat()*) If you want finer control, use *sprintf()*, e.g. to print out print out temperatures as reals ("f"=floating points) with four decimal places and nine total character positions, followed by a C for Celsius:

```
> temps = c(12.5, 37.234324, 1342434324.79997234, 2.3456e-6, 1e15)
> cat(sprintf("%9.4f C", temps))
12.5000 C 37.2343 C 1342434324.8000 C 0.0000 C 1000000000000000.0000
C
```

And of course you can always save to an R data file using *save.image()* (to save all the objects in the workspace or *save()* to save only some objects. Happily this is platform-independent so can be used to transfer between different OS.

8 Tip

Note: to have your own utility functions available to you in any R session, you can include them in your *.Rprofile* file or source a separate file of R code within *.Rprofile*. E.g., I have *f.sort()* in a file called *utils.q* that is sourced by *.Rprofile*. Or if you're more organized, create an R package!

Good coding practices

September 20, 2011

References:

- Murrell, Introduction to Data Technologies, Ch. 2
- Journal of Statistical Software v42: 19 ways of looking at statistical software

Some of these tips apply more to software development and some more to analyses done for specific projects; hopefully it will be clear in most cases.

1 Editors

Use an editor that supports the language you are using (e.g., *Emacs*, *WinEdt*, *Tinn-R* on Windows). Some advantages of this can include: (1) helpful color coding of different types of syntax and of strings, (2) automatic indentation and spacing, (3) code can often be run or compiled from within the editor, (4) parenthesis matching, (5) line numbering (good for finding bugs).

2 Coding syntax

- Header information: put meta-info on the code into the first few lines of the file as comments. Include who, when, what, how the code fits within a larger program (if appropriate), possibly the versions of R and key packages that you wrote this for
- Indentation: do this systematically (your editor can help here). This helps you and others to read and understand the code and can help in detecting errors in your code because it can expose lack of symmetry.
- Whitespace: use a lot of it. Some places where it is good to have it are (1) around operators (assignment and arithmetic), (2) between function arguments and list elements, (3) between matrix/array indices, in particular for missing indices.

- Use blank lines to separate blocks of code and comments to say what the block does
- Split long lines at meaningful places.
- Use parentheses for clarity even if not needed for order of operations. For example, `a/y*x` will work but is not easy to read and you can easily induce a bug if you forget the order of ops.
- Documentation - add lots of comments (but don't belabor the obvious). Remember that in a few months, you may not follow your own code any better than a stranger. Some key things to document: (1) summarizing a block of code, (2) explaining a very complicated piece of code - recall our complicated regexp, (3) explaining arbitrary constant values.
- For software development, break code into separate files (<2000-3000 lines per file) with meaningful file names and related functions grouped within a file.
- Choose a consistent naming style for objects and functions: e.g. *nIts* vs. *n.its* vs *numberOfIts* vs. *n_its*
- Try to have the names be informative without being overly long.
- Don't overwrite names of objects/functions that already exist in R. E.g., don't use 'lm'.

```
> exists("lm")
```
- Use active names for functions (e.g., *calc.loglik*, *calcLogLik*)
- Learn from others' code

This semester, someone will be reading your code - me, when I look at your assignments. So to help me in understanding your code and develop good habits, put these ideas into practice in your assignments.

3 Coding style

This is particularly focused on software development, but some of the ideas are useful for data analysis as well.

- Break down tasks into core units
- Write reusable code for core functionality and keep a single copy of the code (w/ backups of course) so you only need to change it once

- Smaller functions are easier to debug, easier to understand, and can be combined in a modular fashion (like the UNIX utilities)
- Write functions that take data as an argument and not lines of code that operate on specific data objects. Why? Functions allow us to reuse blocks of code easily for later use and for recreating an analysis (reproducible research). It's more transparent than sourcing a file of code because the inputs and outputs are specified formally, so you don't have to read through the code to figure out what it does.
- Functions should:
 - be modular (having a single task);
 - have meaningful name; and
 - have a comment describing their purpose, inputs and outputs (see the help file for an R function for how this is done in that context).
- Object orientation is a nice way to go
- Don't hard code numbers - use variables (e.g., number of iterations, parameter values in simulations), even if you don't expect to change the value, as this makes the code more readable:


```
> speedOfLight = 3e8
```
- Use lists to keep disparate parts of related data together
- Practice defensive programming
 - check function inputs and warn users if the code will do something they might not expect or makes particular choices;
 - check inputs to *if* and the ranges in *for* loops;
 - provide reasonable default arguments;
 - document the range of valid inputs;
 - check that the output produced is valid; and
 - stop execution based on checks and give an informative error message.
- Try to avoid system-dependent code that only runs on a specific version of an OS or specific OS
- Learn from others' code

4 Tips for avoiding bugs

- Write an initial version of the code in the simplest way, without trying to be efficient (e.g., use for loops even if you're coding in R); then make a second version that employs efficiency tricks and check that both produce the same output.
- Make big changes in small steps, sequentially checking to see if the code has broken on test case(s).
- Plan out your code in advance, including all special cases/possibilities.
- Figure out if a functionality already exists in (or can be adapted from) an R package (or potentially in a C/Fortran library/package): code that is part of standard mathematical/numerical packages will probably be more efficient and bug-free than anything you would write.
- When doing software development, write tests for your code early in the process.

5 Running analyses

Save your output at intermediate steps (including the random seed state) so you can restart if an error occurs or a computer fails. Using *save()* and *save.image()* to write to *.RData* files work well for this.

Run your code on a small subset of the problem before setting off a job that runs for hours or days. Make sure that the code works on the small subset and saves what you need properly at the end.

6 Versioning and backup

For data analysis and your own project work:

- Basic version control ideas: whenever you make changes to the structure of your code/program, make it a new version, numbering sequentially (or with more sophisticated version numbering used by software developers).
- Keep a text file that documents what changes from version to version (or do this in the document header). I find this easier than trying to name my code files informatively.
- DO NOT get rid of old code versions (text is cheap to store) and DO NOT get rid of blocks of code that work.

- Keep a running file (or comments in the code) of what you'd like to change in the future.

In addition, for software development you might consider:

- Using version control software: *cvs* and *svn* are two tools.

7 Documenting an analysis

Provenance is becoming increasingly important in science. It basically means being able to trace the steps of an analysis back to its origins. *Replicability* is a related concept - the idea is that you or someone else could replicate the analysis that you've done. This can be surprisingly hard as time passes even if you're the one attempting the replication. Let's think through what is required for something to be replicable.

- Have a directory for each project with meaningful subdirectories: e.g., *code*, *data*, *paper*
- Keep a document describing your running analysis with dates in a text file (i.e., a lab book)
- Note where data were obtained (and when, which can be helpful when publishing) and pre-processing steps in the lab book. Have data version numbers with a file describing the changes and dates (or in lab book).
- Have a file of code for pre-processing, one or more for analysis, and one for figure/table preparation.
 - Have the code file for the figures produce the EXACT manuscript figures (or use Sweave), operating on an RData file that contains all the objects necessary to run the figure-producing code; the code producing the RData file should be in your analysis file (or somewhere else sensible).
- Note what code files do what in the lab book.

R programming

October 18, 2011

References:

- Adler
- Chambers
- Murrell, Introduction to Data Technologies.
- [R intro manual](#) (R-intro) and [R language manual](#) (R-lang), both on CRAN.
- Venables and Ripley, Modern Applied Statistics with S

1 Working with data structures

1.1 Lists and data frames

Extraction You extract from lists with “[[“ or with “[“

```
> x = list(a = 1:2, b = 3:4, c = rnorm(4))
> x[[2]] # extracts the indicated component, which can be anything,
in this case just an integer vector
> x[c(1, 3)] # extracts subvectors, which since it is a list, will
also be a list
```

When working with lists, it’s handy to be able to use the same function on each element of the list:

```
> lapply(x, length)
> sapply(x, length) # returns things in a user-friendly way
```

Note that to operate on a data frame, we’ll generally want to use *lapply()* or *sapply()*, as *apply()* is really designed for working with elements that are all of the same type:

```
> state = as.data.frame(state.x77); state$region = state.region
```

```
> apply(state, 2, class) # hmmm
> sapply(state, class)
```

Here's a nice trick to pull out a specific component from each element of a list:

```
> params = list( a = list(mn = 7, sd = 3), b = list(mn = 6, sd =1),
c = list(mn = 2, sd = 1))
> sapply(params, '[' , 1)
```

Note the use of the additional argument(s) to *sapply()* - this can also be done in the other *apply()* variants.

Calculations in the context of stratification We can also use an *apply()* variant to do calculations on subgroups, defined based on a factor or factors.

```
> tapply(mtcars$mpg, mtcars$cyl, mean)
> tapply(mtcars$mpg, list(mtcars$cyl, mtcars$gear), mean)
```

with() provides a clean way to use a function (or any R code, specified as R statements enclosed within `{}`) within the context of a data frame (or an environment).

```
> with(mtcars, tapply(mpg, list(cyl, gear), mean))
```

Check out *aggregate()* and *by()* for nice wrappers to *tapply()* when working with data frames. *aggregate()* returns a data frame and works when the output of the function is univariate, while *by()* returns a list, so can return multivariate output:

```
> aggregate(mtcars, list(cyl = mtcars$cyl), mean) # this uses the
function on each column of mtcars, split by the 'by' argument
> by(warpbreaks, warpbreaks$tension, function(x) {lm(breaks ~ wool,
data = x)})
```

In some cases you may actually need an object containing the subsets of the data, for which you can use *split()*:

```
> split(mtcars, mtcars$cyl)
```

The *do.call()* function will apply a function to the elements of a list. For example, we can *rbind()* together (if compatible) the elements of a list of vectors instead of having to loop over the elements or manually type them in:

```
> myList = list(a = 1:10, b = 11:20, c = 21:30)
> mat = do.call(rbind, myList)
```

Why couldn't we just use *rbind()* directly? Basically we're using *do.call()* to use functions that take "..." as input and to use a list as the input instead.

The various *apply()* functions for lists may be faster than a loop, but if the dominant part of the calculation lies in the time required by the function on each of the elements, then the main reason for using an *apply()* variant is code clarity.

```

> n = 500000; nr = 1000; nCalcs = n/nr
> mat = matrix(rnorm(n), nr = nr)
> times = 1:nCalcs
> system.time(out1 <- apply(mat, 1, function(vec) {mod = lm(vec
~ times); mod$coef[2]}))
> system.time({out2 = NA; length(out2) = nCalcs;
+   for(i in 1:nCalcs){out2[i] = lm(mat[i, ] ~ times)$coef[2]}})

```

1.2 Vectors and matrices

Column-major vs. row-major matrix storage Matrices in R are column-major ordered, which means they are stored by column.

```

> mat = matrix(rnorm(500), nr = 50)
> identical(mat[100], mat[50, 2])
> mat = rnorm(500, nr = 5)
> vec = c(mat)
> mat2 = matrix(vec, nr = 50)
> identical(mat, mat2)

```

If you want to fill a matrix row-wise:

```

> matrix(1:4, 2, byrow = TRUE)

```

Column-major ordering is also used in Matlab and Fortran, while row-major ordering is used in C.

Recycling R will recycle the right-hand side of an assignment if the left-hand side has more elements, and will give an error or warning if there is partial recycling. Recycling is powerful, but dangerous.

```

> mat = matrix(1:16, 4)
> mat[, 1] = c(1, 2)
> mat[, 1] = 1:3

```

Sequences The *rep()* and *seq()* functions are helpful for creating sequences with particular structure. Here are some examples of combining them:

```

> rep(seq(1, 9, by = 2), each = 2)
> rep(seq(1, 9, by = 2), times = 2) # repeats the whole vector 'times'
times

```

```
> rep(seq(1, 9, by = 2), times = 1:5) # repeats each element based
on 'times'
```

Sidenote: sometimes you can take a sequence generated with *rep()* and *seq()*, put it into a matrix, do some manipulations (such as transposing the matrix), and pull it out as a vector to get sequences in the order you want. For example, how could I produce “1 3 5 7 9 3 5 7 9 5 7 9 7 9 9” without using a loop? [Hint: you might need a function that you might not be familiar with - but you might be able to guess that it exists.] You can also get sequences using mathematical expressions and *round()* and modulo operations. We’ll play around with this more when we talk about efficiency.

You can create regular combinations of the values in two (or more) vectors as follows:

```
> expand.grid(x = 1:5, y = -1:1, theta = c("lo", "hi"))
```

This is useful for setting up simulations and for creating a grid of values for image and surface plots.

Identifying elements by index You can figure out the indices of elements having a given characteristic using *which()*:

```
> x = c(1, 10, 2, 9, 3, 8)
```

```
> which(x < 3)
```

which.max() and *which.min()* have similar sort of functionality.

Random sampling We can draw random samples with or without replacement using *sample()*.

Apply() The *apply()* function will apply a given function to either the rows or columns of a matrix or a set of dimensions of an array:

```
> x = matrix(1:6, nr = 2)
```

```
> apply(x, 1, min) # by row
```

```
> apply(x, 2, min) # by column
```

```
> x = array(1:24, c(2, 3, 4))
```

```
> apply(x, 2, min) # for(j in 1:3) print(min(x[ , j, ]))
```

```
> apply(x, c(2, 3), min) # for(j in 1:3) {for(k in 1:4)
print(min(x[ , j, k]))}
```

This can get confusing, but can be very powerful. Basically, I’m collapsing across only the rows in the second example and across rows and “3rd dimension” in the final example. Caution: if the result of each subcalculation is a vector, these will be *cbind*’ed together (recall column-major order), so if you used *apply()* on the row margin, you’ll need to transpose the result.

1.3 Linear algebra

We'll focus on matrices here. A few helpful functions are `nrow()` and `ncol()`, which tell the dimensions of the matrix. The `row()` and `col()` functions will return matrices of the same size as the original, but filled with the row or column number of each element. So to get the upper triangle of a matrix, X , we can do:

```
> X[col(X) >= row(X)]
```

See also the `upper.tri()` and `lower.tri()` functions, as well as the `diag()` function. `diag()` is quite handy - you can either extract the diagonals or you can assign into the diagonals:

```
> diag(X) = 1
```

or you can create a diagonal matrix:

```
> d = diag(c(rep(1, 3), rep(2, 2)))
```

To transpose a matrix:

```
> t(X)
```

Basic operations The basic matrix-vector operations are:

```
> X %*% Y # multiplication
```

```
> X * Y # direct product
```

```
> x %o% y # outer product:  $xy^T$ 
```

```
> outer(x, y) # same thing
```

```
> outer(x, y, function(x, y) cos(y)/(1 + x^2)) # evaluation of  $f(x, y)$ 
```

for all pairs of x, y values

```
> crossprod(X, Y) #  $X^T Y$ , and faster than t(X) %*% Y!
```

For inverses and solutions of systems of linear equations:

```
> solve(X) #  $X^{-1}$ 
```

```
> solve(X, y) #  $X^{-1}y$ 
```

Note that if all you need to do is solve the linear system, you should never explicitly find the inverse, UNLESS (1) you need the actual matrix, e.g., to get a covariance matrix of parameters, or (2) you are finding many solutions all with the same matrix (and at different points in time, otherwise, you would just column bind the various 'y' vectors as the second argument to `solve()`). Even if you're finding many solutions, it may be better to use a matrix decomposition.

We'll defer matrix decompositions (Cholesky, eigendecomposition, SVD, QR) until the numerical linear algebra unit.

2 Functions

Functions are at the heart of R. In general, you should try to have functions be self-contained - operating only on arguments provided to them, and producing no side effects, though in some cases there are good reasons for making an exception.

Functions that are not implemented internally in R (e.g., user-defined functions) are also referred to officially as *closures* (this is their type) - this terminology sometimes comes up in error messages.

2.1 Inputs

Arguments can be specified in the correct order, or given out of order by specifying *name = value*. In general the more important arguments are specified first. You can see the arguments and defaults for a function using `args()`:

```
> args(lm)
```

Functions may have unspecified arguments, which is designated using `'...'`. Unspecified arguments occurring at the beginning of the argument list are generally a collection of like objects that will be manipulated (consider `paste()`), while unspecified arguments occurring at the end are often optional arguments (consider `plot()`). These are sometimes passed along to a function within the function. For example, here's my own wrapper for plotting, where any additional arguments specified by the user will get passed along to `plot`:

```
> pplot = function(x, y, pch = 16, cex = 0.4, ...) plot(x, y, pch
= pch, cex = cex, ...)
```

If you want to manipulate what the user passed in as the `...` args, rather than just passing them along, you can extract them (the following code would be used within a function to which `...` is an argument:

```
> args = list(...)
> print(c(..1, ..3))
```

Arguments can also have default values, which may be `NULL`. The default values are assigned in the *frame* of the function evaluation (more on frames in a bit). The default values can also relate to other arguments. As an example, consider `dgamma()`:

```
> args(dgamma)
```

If you are writing a function and designate the default as `argname = NULL`, you can check whether the user provided anything using

```
> is.null(argname)
```

Functions can be passed in as arguments (e.g., see the variants of *apply()*). Note that one does not need to pass in a named function - you can create the function on the spot - this is called an *anonymous function*:

```
> mat = matrix(1:9, 3)
> apply(mat, 2, function(vec) vec - vec[1])
> apply(mat, 1, function(vec) vec - vec[1]) # explain why the result
is transposed
```

We can see the arguments using *args()* and extract the arguments using *formals()*. *formals()* can be helpful if you need to manipulate the arguments.

```
> f = function(x, y = 2, z = 3 / y) { x + y + z }
> a1 = args(f)
> a2 = formals(f)
```

match.call() will show the user-supplied arguments explicitly matched to named arguments.

```
> match.call(mean, quote(mean(y, na.rm = TRUE)))
```

Pass by value vs. pass by reference Note that R makes a copy of all objects that are arguments to a function, with the copy residing in the frame (the environment) of the function (we'll see more about frames and environments a bit later). This is a case of *pass by value*. In other languages it is also possible to *pass by reference*, in which case, changes to the object made within the function affect the value of the argument in the calling environment. R's designers chose not to allow pass by reference because they didn't like the idea that a function could have the side effect of changing an object.

An important exception is *par()*. If you change graphics parameters in a user-defined function (but not by passing specific arguments to a plotting function - e.g., `> plot(x, y, pch = '+')` doesn't change things except in that call), they are changed permanently. One trick is as follows:

```
> f = function() {
+   oldpar = par()
+   par(cex = 2)
+   # body of code
+   par() = oldpar
+ }
```


2.2 Variable scope

Variables in the parent environment (the environment in which a function is defined, also called the enclosing environment) are available within a function. This is the analog of *global variables* in other languages. Note that the parent/enclosing environment is NOT the environment in which the function was called.

Be careful when using variables from the parent environment as the value of that variable in the parent environment may well not be what you expect it to be. In general you should avoid this, but it is a useful way to reduce memory use since a copy is made of any variable passed into a function as an argument – in some cases we may not want to make a copy of a very large object.

```
> x = 3
> f = function() {x = x^2; print(x)}
> f(x)
> x # what do you expect?
> f = function() {assign('x', x^2, env = .GlobalEnv)} # careful,
this could be dangerous as a variable is changed as a side effect
```

Here's another example of scope:

```
> x = 3
> f = function() {
  f2 = function() { print(x) }
  f2()
}
> f() # what will happen?
> f = function() {
  f2 = function() { print(x) }
  x = 7
  f2()
}
> f() # what will happen?
> f2 = function() print(x)
> f = function() {
  x = 7
  f2()
}
> f() # what will happen?
```

Here's a somewhat tricky example:

```
> f = function() {
```

```
y = 10
g = function(x) x + y
return(g)
}
> h = f()
> h(3)
```

Let's work through this:

1. What is the enclosing environment of the function $g()$?
2. What does $g()$ use for y ?
3. When $f()$ finishes, does its environment disappear? What would happen if it did?
4. What is the enclosing environment of $h()$?

User-supplied arguments are evaluated in the calling frame, while default arguments are evaluated in the frame of the function:

```
> z = 3
> f = function(x, y = 3) {x+y}
> f(z*5)
```

Here, when $f()$ is called, z is evaluated in the calling frame and assigned to x in the frame of the function, while $y = 3$ is evaluated in the frame of the function.

2.3 Alternatives to pass by value

There are occasions we do not want to pass by value. These include when the inputs are large and when we want a function to modify a complicated object without having to return it and re-assign it in the parent environment. There are several work-arounds:

1. To avoid creating a second copy of an object, we can access the object in the parent environment as a 'global variable'. However we cannot change the value of the object in the parent environment without some additional tools. We can use the ' \llleftarrow ' operator to assign into an object in the parent environment (provided an object of that name exists in the parent environment). We can also access the object using information about the environment or about stack of frames created by the function calls (more on this later in the unit).
2. We can use replacement functions (we'll cover this in a bit), which hide the reassignment in the parent environment from the user. Note that a second copy is created in this case, but the original copy is quickly removed.

3. We can use a closure. This involves creating functions within a function call and returning the functions as a list. When one executes the enclosing function, the list is created and one can call the functions of that object. Those functions then can access objects in the enclosing environment (the environment of the original function) and can use ‘<<-‘ to assign into the enclosing environment, to which all the functions have access. Let’s see the example in Chambers Sec. 5.4.
4. We can have the objects we want to provide as inputs (and that we may want to change) in an environment and pass that environment as an argument to the function of interest. Then within the function, we use *get()* (with the *inherits* argument set to *FALSE*) and *assign()* to access the objects. I’ve used this approach on several occasions, though the R team frowns on this. See more in the section on environments.

2.4 Pointers

By way of contrast to R’s pass by value system, I want to briefly discuss the idea of a pointer, common in compiled languages such as C.

```
int x = 3;
int* ptr;
ptr = &x;
*ptr * 7; // returns 21
```

Here *ptr* is the address of the integer *x*.

Vectors in C are really pointers to a block of memory:

```
int x[10];
```

In this case *x* will be the address of the first element of the vector. We can access the first element as ***x*[0]** or ****x***.

Why have we gone into this? In C, you can pass a pointer as an argument to a function. The result is that only the scalar address is copied and not the entire vector, and inside the function, one can modify the original vector, with the new value persisting on exit from the function. For example:

```
int myCal(int *ptr){
  *ptr = *ptr + *ptr;
}
```

When calling C or C++ from R, one (implicitly) passes pointers to the vectors into C. Let’s see an example:

```
> out = rep(0, n)
```

```
> out = .C("logLik", out = as.double(out),
theta = as.double(theta))$out
```

In C, the function definition looks like this:

```
void logLik(double* out, double* theta)
```

2.5 Outputs

return(x) will specify x as the output of the function. By default, the output is the result of the last evaluated statement. *return()* can occur anywhere in the function, and allows the function to exit as soon as it is done. *invisible(x)* will return x and the result can be assigned in the calling environment but it will not be printed if not assigned:

```
> f = function(x){invisible(x^2)}
> f(3)
> a = f(3)
```

A function can only return a single object, but of course we can tack things together as a list and return that, as with `lm()` and many other functions.

```
> class(mod)
> is.list(mod)
```

2.6 Operators

Operators, such as `'+'`, `'/'` are just functions, but their arguments can occur both before and after the function call:

```
> a = 7; b = 3
> a + b # let's think of this as a mathematical function -- what's
the function call?
> `+`(a, b)
```

In general, you can use back-ticks to refer to the operators as operators instead of characters. In some cases single or double quotes also work.

Finally, since an operator is just a function, you can use it as an argument in various places:

```
> myList = list(list(a = 1:5, b = "sdf"), list(a = 6:10, b = "wer"))
> myMat = sapply(myList, `[`, 1) # note that the index "1" is the
additional argument to the `[` function
> outer(x, y, `+`)
```

We can look at the code of an operator as follows using back-ticks to escape out of the standard R parsing:

```
> `%*%`
```

You can define your own *binary* operator (an operator taking two arguments) using a string inside `%` symbols:

```
> `%2%` = function(a, b) {2 * (a + b)}
> 3 %2% 7
```

Since operators are just functions, there are cases in which there are optional arguments:

```
> mat = matrix(1:4, 2, 2)
> mat[ , 1]
> mat[ , 1, drop = FALSE] # what's the difference?
```

Take a look at help on “`[`”: `[` to figure out what the `drop` argument does and the syntax you would need.

2.7 Unexpected functions and replacement functions

All code in R can be viewed as a function call.

What do you think the functional version of the following code? What are the arguments?

```
> if(x > 27){
+ print(x)
+ } else{
+ print("too small")
+ }
```

Assignments that involve functions or operators on the left-hand side (LHS) are called *replacement expressions*. These can be quite handy. E.g.,

```
> diag(mat) = c(3, 2)
> is.na(vec) = 3
> names(df) = c('var1', 'var2')
```

Replacement expressions are actually function calls. The R interpreter calls the replacement function (which creates a new object that includes the replacement) and then assigns the result to the name of the original object.

```
> mat = matrix(rnorm(4), 2, 2)
> diag(mat) = c(3, 2)
> mat = `diag<-`(mat, c(3, 2))
> base::`diag<-`
> mat[1, 2] = 7
> mat <- `[<-`(mat, 1, 2, 7.5)
```

The old version of `mat` still exists until R's memory management cleans it up, but it's no longer referred to by the symbol `'mat'`. Occasionally this sort of thing might cause memory usage to

increase (for example it's possible if you're doing replacements on large objects within a loop), but in general things should be fine.

You can define your own replacement functions like this, with the requirements that the last argument be named *'value'* and that the function return the entire object:

```
> me = list(name = 'Chris', age = 25)
> `name<-` <- function(obj, value){
  obj$name <- value
  return(obj)
}
> name(me) = 'Christopher'
```

2.8 Functions as objects

Note that a function is just an object.

```
> x = 3
> x(2)
> x = function(z) z^2
> x(2)
```

We can call a function based on the text name of the function

```
> myFun = 'mean'; x = rnorm(10)
> eval(as.name(myFun))(x)
> f = function(fxn, x) match.fun(fxn)(x) # handy for extracting
a function when passed as an argument to a function
> f("mean", x)
```

This allows us to write functions in which the user passes in the name of an auxiliary function (as a character vector) they want to use (as an example, this works when using *outer()*). You can also pass around the actual function as an object. Caution: scoping issues may arise.

Function objects contain three components: an argument list, a body (a parsed R statement), and an environment.

```
> f = function(x) x
> f2= function(x) y = x^2
> f3 = function(x) {y = x^2; z = x^3; return(list(y, z))}
> class(f)
> typeof(body(f)); class(body(f))
> typeof(body(f2)); class(body(f2))
```

```
> typeof(body(f3)); class(body(f3))
```

We'll see more about objects relating to the R language and parsed code in a later section. For now, just realize that the parsed code itself is treated as an object(s) with certain types and certain classes.

We can extract the argument object as

```
> f4 = function(x, y = 2, z = 1/y) {x + y + z}
> args = formals(f4)
> class(args)
```

A *pairlist* is like a list, but with pairing that in this case pairs argument names with default values.

3 Flow control and logical operations

3.1 Logical operators

Everyone is probably familiar with the comparison operators, `<`, `<=`, `>`, `>=`, `==`, `!=`. Logical operators are slightly trickier:

Subsetting: `&` and `|` are the “AND” and “OR” operators used when subsetting - they act in a vectorized way:

```
> x = rnorm(10)
> x[x > 1 | x < -1]
> x = 1:10; y = c(rep(10, 9), NA)
> x > 5 | y > 5 # note that TRUE | NA evaluates to TRUE
```

Flow control: `&&` and `||` use only the first element of a vector and also proceed from left to right, returning the result as soon as possible and then ignoring the remaining comparisons (this can be handy because in some cases the second condition may give an error if the first condition is not passed). They are used in flow control (i.e., with *if* statements). Let's consider how the single and double operators differ:

```
> a = 7; b = NULL
> a < 8 | b > 3
> a < 8 || b > 3
```

You can use `!` to indicate negation:

```
> a = 7; b = 5
> !(a < 8 && b < 6)
```

3.2 If statements

If statements are at the core of programming. In R, the syntax is **if(condition) statement else other_statement**, e.g.,

```
> if (x > 7) {
+ x = x + 3
+ } else{
+ x = x - 3
+ }
```

When one of the statements is a single statement, you don't need the curly braces around that statement. An extension of *if* looks like:

```
> if (x > 7) {
+ x = x + 3
+ print(x)
+ } else if (x > 4) {
+ x = x + 1
+ print(x)
+ } else if (x > 0) {
+ x = x - 3
+ print(x)
+ } else{
+ x = x - 7
+ print(x)
+ }
```

Finally, be careful that you avoid the following syntax as R will just evaluate the first part since it is a valid statement and R doesn't know it should also parse the *else* piece as part of the *if* statement:

```
if(condition) {
statement1 } # R sees a complete, valid statement and executes
else{ # R gets here and tries to evaluate a new statement and fails
statement2
}
```

Instead, good syntax is:

```
if (condition) {
statement1 # R gets here and the statement is not yet complete
} else{ # so it keeps reading and parses the else as part of the
if
statement2
```



```
}
```

There's also the *ifelse()* function, which operates in a vectorized fashion:

```
> x = rnorm(100)
> truncx = ifelse(x > 0, x, 0)
```

Common bugs in the condition of an if statement include the following:

1. Only the first element of *condition* is evaluated. You should be careful that *condition* is a single logical value and does not evaluate to a vector as this would generally be a bug. [see p. 152 of Chambers]
2. If *condition* includes some R code, it can fail and produce something that is neither TRUE nor FALSE. Defensive programming practice is to check the condition for validity.

```
> if(min(vals) > eps) { } # not good practice
> minval = min(vals) # better practice
> if(!is.na(minval) && minval > eps) { }
```
3. Use *identical()* or *all.equal()* rather than “==” to ensure that you deal properly with vectors and always get a single logical value back. We'll talk more about issues that can arise when comparing decimal numbers on a computer in the next unit.

3.3 Switch

switch() is handy for choosing amongst multiple outcomes depending on an input, avoiding a long set of if-else syntax. The first argument is a statement that determines what choice is made and the second is a list of the outcomes, in order or by name:

```
> x = 3; y = 10
> switch(x, log(y), sqrt(y), y)
> center = function(x, type){
>   switch(type,
>     mean = mean(x),
>     median = median(x),
>     trimmed = mean(x, trim = .1)
>   )
> }
> x = rgamma(100, 1)
> center(x, 'median')
> center(x, 'mean')
```

3.4 Loops

Loops are at the core of programming in other functional languages, but in R, we generally try to avoid them as they're often (but not always) slow. In most cases looping can be avoided by using vectorized calculations, versions of *apply()*, and other tricks. But sometimes they're unavoidable and for quick and dirty coding and small problems, they're fine. And in some cases they may be faster than other alternatives. One case is that in working with lists they may be faster than their *lapply*-style counterpart, though often they will not be.

The workhorse is the *for* loop: **for(var in sequence) statement**, where, as with the *if* statement, we need curly braces around the body of the loop if it's more than one valid R statement:

```
> nIts = 1000; means = NA; length(means) = nIts
> for(it in 1:nIts){
+ means[it] = mean(rnorm(100))
+ if(identical(it %% 100, 0)) print(c(date(), it))
+ }
```

Challenge: how do I do this much faster?

You can also loop over a non-numeric vector of values.

```
> for(state in c('Ohio', 'Iowa', 'Georgia'))
print(state.x77[row.names(state.x77) == state, "Income"])
```

Challenge: how can I do this faster?

Note that to print to the screen in a loop you explicitly need to use *print()* or *cat()*; just writing the name of the object will not work. This is similar to *if* statements and functions.

```
> for(i in 1:10) i
```

You can use the commands *break* (to end the looping) and *next* (to go to the next iteration) to control the flow:

```
> for(i in 1:10){
+ if(i == 5) break
+ print(i)
+ }
> for(i in 1:10){
+ if(i == 5) next
+ print(i)
+ }
```

while loops are used less frequently, but can be handy: *while(condition) statement*, e.g. in optimization. See p. 59 of Venables and Ripley, 4th ed., whose code I've included in the demo code file.

A common cause of bugs in for loops is when the range ends at zero or a missing value:

```
> mat = matrix(1:4, 2)
> submat = mat[mat[, ] > 5]
> for(i in 1:nrow(submat)) print(i)
```

4 Efficiency

In general, make use of R's built-in functions, as these tend to be implemented internally (i.e., via compiled code in C or Fortran). In particular, if R is linked to optimized BLAS and Lapack code (e.g. Intel's MKL, GotoBLAS, and AMD's ACML), you should have good performance (potentially comparable to Matlab and to coding in C). Often you can figure out a trick to take your problem and transform it to make use of the built-in functions.

Note that I run a lot of MCMCs so pay attention to making sure my calculations are fast as they are done repeatedly. Similarly, one would want to pay attention to speed when doing simulations and bootstrapping, and in some cases for optimization. And if you're distributing code, it's good to have it be efficient. But in other contexts, it may not be worth your time. Also, it's good practice to code it transparently first to reduce bugs and then to use tricks to speed it up and make sure the fast version works correctly.

Results can vary with with your system setup and version of R, so the best thing to do is figure out where the bottlenecks are in your code, and then play around with alternative specifications.

Finally, as you gain more experience, you'll get some intuition for what approaches might improve speed, but even with experience I find myself often surprised by what matters and what doesn't. It's often worth trying out a bunch of different ideas; *system.time()* is your workhorse tool in this context.

4.1 Vectorized calculations

One key way to write efficient R code is to take advantage of R's vectorized operations.

```
> n = 100000
> x = rnorm(n)
> system.time(x2 <- x^2)
> x2 = NA
> system.time({length(x2) <- n; for(i in 1:n){ x2[i] <- x[i]^2}})
# how many orders of magnitude slower?
```

So what's different in how R handles the calculations that explains the disparity? The vectorized

calculation is being done natively in C in a for loop. The for loop above involves executing the for loop in R with repeated calls to C code at each iteration. You can usually get a sense for how quickly an R call will pass things along to C by looking at the body of the relevant function(s) being called and looking for *.Primitive*, *.Internal*, *.C*, *.Call*, or *.Fortran*. Let's take a look at the code for '+', *mean.default()*, and *chol.default()*.

Many R functions allow you to pass in vectors, and operate on those vectors in vectorized fashion. Here's another example from PS2:

```
> startIndices = seq(31, by = 12, length = nValues)
> substring(line, startIndices, startIndices + 11)
```

Challenge: Consider the chi-squared statistic involved in a test of independence in a contingency table:

$$\chi^2 = \sum_i \sum_j \frac{(y_{ij} - e_{ij})^2}{e_{ij}}, \quad e_{ij} = \frac{y_{i \cdot} y_{\cdot j}}{y_{\cdot \cdot}}$$

where $f_{i \cdot} = \sum_j f_{ij}$. Write this in a vectorized way without any loops. Note that 'vectorized' calculations also work with matrices and arrays.

Vectorized operations can also be faster than built-in functions, and clever vectorized calculations even better, though sometimes the code is uglier:

```
> x = rnorm(100000)
> system.time(truncx <- ifelse(x > 0, x, 0))
> system.time({truncx <- x; truncx[x < 0] <- 0})
> system.time(truncx <- x * (x > 0))
```

The demo code has a surprising example where combining vectorized calculations with a *for* loop is actually faster than using *apply()*. The goal is to remove rows of a large matrix that have any NAs in them.

4.2 Using *apply()* and specialized functions

Another core efficiency strategy is to use the *apply()* functionality. Even better than *apply()* for calculating sums or means of columns or rows (it also can be used for arrays) is *{row,column}{Sums,Means}*:

```
> n = 3000; x = matrix(rnorm(n * n), nr = n)
> system.time(out <- apply(x, 1, mean))
> system.time(out <- rowMeans(x))
```

We can 'sweep' out a summary statistic, such as subtracting off a mean from each column, using *sweep()*

```
> system.time(out <- sweep(x, 2, STATS = colMeans(x), FUN = "-"))
```

Here's a trick for doing it even faster based on vectorized calculations, recognizing that if we

subtract a vector from a matrix, it subtracts each element of the vector from all the elements in the corresponding ROW.

```
> system.time(out2 <- t(t(x) - colMeans(x)))
> identical(out, out2)
```

4.3 Matrix algebra efficiency

Often calculations that are not explicitly linear algebra calculations can be done as matrix algebra. The following can be done faster with *rowSums()*, so it's not a great example, but this sort of trick does come in handy in surprising places.

```
> mat = matrix(rnorm(500*500), 500)
> system.time(apply(mat, 1, sum))
> system.time(mat %*% rep(1, ncol(mat)))
> system.time(rowSums(mat))
```

On the other hand, big matrix operations can be slow. Suppose you want a new matrix that computes the differences between successive columns of a matrix of arbitrary size. How would you do this as matrix algebra operations? [see demo code] Here it turns out that the *for* loop is much faster than matrix multiplication. However, there is a way to do it faster as matrix direct subtraction.

As a sidenote, we can also be clever about creating patterned matrices, as seen in the demo code. Challenge: how can we fill the matrix of 1, -1 values? How can we do it quickly? Note that this level of optimization is only worth it if you're doing something over and over again, or writing code that you will distribute.

When doing matrix algebra, the order in which you do operations can be critical for efficiency. How should I order the following calculation?

```
> n = 5000
> A = matrix(rnorm(5000 * 5000), 5000)
> B = matrix(rnorm(5000 * 5000), 5000)
> x = rnorm(5000)
> system.time(res <- A %*% B %*% x)
```

We can use the matrix direct product to do some manipulations much more quickly than using matrix multiplication. Challenge: How can I use the direct product to find the trace of XY ?

You can generally get much faster results by being smart when using diagonal matrices. Here are some examples, where we avoid directly doing: $X + D$, DX , XD :

```
> diagvals = sample(1:5, 10000, replace = TRUE); D = diag(diagvals)
> diag(X) = diag(X) + diagvals
```

```
> diagvals * x # instead of D %*% x
> t(t(X) * diagvals) # instead of X %*% D
```

More generally, sparse matrices and structured matrices (such as block diagonal matrices) can generally be worked with MUCH more efficiently than treating them as arbitrary matrices. The *spam* (for arbitrary sparse matrices) and *bdsmatrix* (for block-diagonal matrices) packages in R can help, as can specialized code available in other languages, such as C and Fortran packages.

4.4 Fast initialization

It is very inefficient to iteratively add elements to a vector (using *c()*) or iteratively use *cbind()* or *rbind()* to add rows or columns to matrices or dataframes. Instead, create the full object in advance (this is equivalent to variable initialization in compiled languages) and then fill in the appropriate elements. The reason is that when R appends to an existing object, it creates a new copy and as the object gets big, this gets slow.

```
> n = 10000
> x = 1; system.time(for(i in 2:n) x = c(x, i))
> system.time({x <- rep(NA, n); for(i in 1:n) x[i] = i})
```

We can actually speed up the initialization (though in most practical circumstances, the second approach here would be overkill):

```
> n = 100000
> system.time(x <- rep(NA, n))
> system.time({x <- NA; length(x) <- n})
```

For matrices, start with the right length vector and then change the dimensions

```
> nr = 2000; nc = 2000
> system.time(x <- matrix(NA, nr, nc))
> system.time({x <- NA; length(x) <- nr * nc; dim(x) <- c(nr, nc)})
```

4.5 Fast mapping

Sometimes you need to map between two vectors. E.g., $y_{ij} \sim \mathcal{N}(\mu_j, \sigma^2)$ is a basic ANOVA type structure. Here are some efficient ways to aggregate to the cluster level and disaggregate to the observation level.

Disaggregate: Create a vector, *idVec*, that gives a numeric mapping of the observations to their cluster. Then you can access the μ value relevant for each observation as: **> mus[idVec]**.

Aggregate: To find the sample means by cluster,

```
> sapply(split(dat$obs, idVec), mean)
```

R provides something akin to a hash table when you have named objects. For example:

```

> vals = rnorm(10)
> names(vals) = c("a", "b", "c", "d", "e", "f", "g", "h", "i", "j")
> labs = c("a", "a", "h", "i")
> vals[labs]

```

I haven't looked into this much, but according to Adler if you need to do something like this where you are looking up from amongst very many items, the fast way to do it is by looking up items within an environment rather than within a named vector or list as environments are implemented using hash tables.

4.6 Challenges

One or more of these challenges may appear on a problem set.

Challenge 1: here's a calculation of the sort needed in mixture component modeling. I have a vector of n observations. I need to find the likelihood of each observation under each of p mixture components (i.e., what's the likelihood if it came from each of the components). So I should produce a matrix of n rows and p columns where the value in the i th row, j th column is the likelihood of the i th observation under the j th mixture component. The idea is that the likelihoods for a given observation are used in assigning observations to clusters. A naive implementation is:

```

> lik <- matrix(NA, nr = n, nc = p)
> for(j in 1:p) lik[, j] <- dnorm(y, mns[j], sds[j])

```

Note that `dnorm()` can handle matrices and vectors as the observations **and** as the means and sds, so there are multiple ways to do this.

Challenge 2: Suppose you have $y_i \sim \mathcal{N}(\sum_{k=1}^{m_i} w_{ik}\mu_k, \sigma^2)$ for a large number of observations, n . I give you a vector of μ values and a ragged list of weights and a ragged list of id's identifying the cluster for each weight (note m_i varies by observation); this is a mixed membership type model. Figure out how to calculate the vector of means, $\sum_k w_{ik}\mu_k$ as fast as possible. Suppose that m_i never gets too big (but μ might have many elements) - could this help you? Part of thinking this through involves thinking about how you want to store the information so that the calculations can be done quickly.

Challenge 3: Do you think there will be any difference in speed between extracting a set of rows from a matrix and extracting a set of columns? How about row and column means? Any implications for how to store data? Note this issue comes up in other languages. For example, it arises when doing matrix algebra in C when one can do operations sequentially by moving along the correct dimension.

Challenge 4: Write code that simulates a random walk in 2-dimensions for n steps. First write out a straightforward implementation that involves looping. Then try to speed it up. The `cumsum()`

function may be helpful.

Caution: sometimes relying on R's recycling, such as is done when direct-multiplying a matrix by a vector to scale the rows, can be dangerous as the code is not transparent and poses greater dangers of bugs. If it's needed to speed up a calculation, the best approach is to (1) first write the code transparently and then compare the efficient code to make sure the results are the same and (2) comment your code.

5 Environments and frames

Environments and frames are closely related. An environment is a collection of named objects (a frame), with a pointer to the 'enclosing environment', i.e., the next environment to look for something in, also called the parent.

5.1 Frames and the call stack

R keeps track of the call stack, which is the set of nested calls to functions. The stack operates like a stack of cafeteria trays - when a function is called, it is added to the stack (pushed) and when it finishes, it is removed (popped). There are a bunch of functions that let us query what frames are on the stack and access objects in particular frames of interest. This gives us the ability to work with objects in the environment(s) from which a function was called.

`sys.nframe()` returns the number of the current frame:

```
> sys.nframe()
> f = function() print(sys.nframe())
> f() # hmm, why did we get '2' instead of '1'?
> f = function() {out <- sys.nframe(); return(out)}
> f()
> f = function() cat(sys.nframe(), "\n")
> f()
> f = function() cat(sys.nframe(), "\n")
> f2 = function() f()
> f2()
```

Let's look at the demo code file for getting more information about the call stack and the frames involved using `sys.status()`, `sys.calls()`, `sys.parents()` and `sys.frames()`. Be careful in that *parent* here refers to the parent in terms of the call stack, not the enclosing environment.

If you're interested in parsing a somewhat complicated example of frames in action, Adler provides a user-defined timing function that evaluates statements in the calling frame.

5.2 Environments and the search path

When R goes looking for an object (in the form of a symbol), it starts in the current frame (e.g., the frame of a function) and then runs up through the enclosing environments, until it reaches the global environment, which is where R starts when you start it (it actually continues further up; see below). In general, these are not the frames on the stack.

By default objects are created in the global environment, *.GlobalEnv*. As we've seen, the environment within a function call has as its enclosing environment the environment where the function was defined (not the environment from which it was called), and this is next place that is searched if an object can't be found in the frame of the function call. This is called lexical scoping (and differs from *S-plus*). As an example, if an object couldn't be found within the environment of an *lm()* function call, R would first look in the environment (also called the namespace) of the stats package, then in packages imported by the stats package, then the base package, and then the global environment.

If R can't find the object when reaching the global environment, it runs through the search path, which you can see with *search()*. The search path is a set of additional environments. Generally packages are created with namespaces, i.e., each has its own environment, as we see based on *search()*. Data frames or list that you attach using *attach()* generally are placed just after the global environment.

```
> search()
> searchpaths()
```

We can also see the nestedness of environments using the following code, using the *environmentName()*, which prints out a nice-looking version of the environment name.

```
> x <- .GlobalEnv
> parent.env(x)
> while (environmentName(x) != environmentName(emptyenv())) {
+ print(environmentName(parent.env(x))); x <- parent.env(x)}
```

Note that eventually the global environment and the environments of the packages are nested within the base environment (of the base package) and the empty environment. Note that here *parent* is referring to the enclosing environment.

We can look at the objects of an environment as follows:

```
> objects(pos = 3) # what does this do?
> objects("package:graphics")
> environment(lm)
```

We can retrieve and assign objects in a particular environment as follows:

```
> lm = function() {return(NULL)}
> x = 1:3; y = rnorm(3); mod = lm(y ~ x)
```

```

> mod = get('lm', pos = 'package:stats')(y ~ x)
> mod = mgcv::gam(y ~ x) # an alternative; in this to distinguish
from gam::gam

```

Note that our (bogus) *lm()* function masks but does not overwrite the default function. If we remove ours, then the default one is still there.

5.3 Creating and working in an environment

Sometimes you may want to create an environment in which to store objects.

```

> e = new.env()
> assign('x', 3, envir = e) # same as e$x = 3
> e$x
> get("x", envir = e, inherits = FALSE) # the FALSE avoids looking
for x in the enclosing environments
> e$y = 5
> objects(e)
> rm('x', envir = e)
> parent.env(e)

```

You might do this as a way to pass objects by reference into a function and avoid having to pass back and re-assign the output. Here's an example where we iteratively update a random walk.

```

> myWalk = new.env(); myWalk$pos = 0
> nextStep = function(walk) walk$pos = walk$pos + sample(c(-1, 1),
size = 1)
> nextStep(myWalk)

```

It can be handy to evaluate some code in a given environment (e.g., an environment you create, or within the elements of a list or data frame):

```

> with(myWalk, pos <- pos + sample(c(-1, 1), size = 1))
> with(mtcars, mean(cyl) * disp)
> mtcars2 = within(mtcars, crazyvar <- mean(cyl) * disp)

```

A different approach uses *eval()* to evaluate some code within a given environment. Note again the use of the *eval()* function, which will evaluate code in a user-specified environment. By default, it evaluates in the result of *parent.frame()*, which amounts to evaluating in the frame from which *eval()* was called.

```

> eval(quote(pos <- pos + sample(c(-1, 1), 1)), envir = myWalk)
> evalq(pos <- pos + sample(c(-1, 1), 1), envir = myWalk) # shorthand
to avoid using quote()

```

6 Formulas

Formulas were introduced to specify linear models, but are now used more generally.

Here are some examples of formulas in R, used to specify a model structure:

Additive model:

```
y ~ x1 + x2 + x3
```

Additive model without the intercept:

```
y ~ x1 + x2 + x3 -1
```

All the other variables in the data frame are used as covariates:

```
y ~ .
```

All possible interactions:

```
y ~ x1 * x2 * x3
```

Only specified interactions (in this case $x1$ by $x2$) (of course, you'd rarely want to fit this without $x2$):

```
y ~ x1 + x3 + x1:x2
```

Protecting arithmetic expressions:

```
y ~ x1 + I(x1^2) + I(x1^3)
```

Using functions of variables

```
y ~ x1 + log(x2) + sin(x3)
```

In some contexts, such as *lattice* package graphics, the “|” indicates conditioning, so $\mathbf{y} \sim \mathbf{x} \mid \mathbf{z}$ would mean to plot y on x within groups of z . In the context of lme-related packages (e.g., *lme4*, *nlme*, etc.), variables after “|” are *grouping* variables (e.g., if you have a random effect for each hospital, hospital would be the grouping variable) and multiple grouping variables are separated by “/”.

We can manipulate formulae as objects, allowing automation. Consider how this sort of thing could be used to write code for automated model selection.

```
> resp = 'y ~'  
> cov = 'x1'  
> for(i in 2:5){ cov = paste(cov, '+ x', i, sep="") }  
> form = as.formula(paste(resp, cov, sep = ""))
```

That's pretty clunky - let's do better:

```
> form = as.formula(paste(resp, paste('x', 1:5, sep = "", collapse = ' + ')))  
> lm(form, data = dat)
```

Standard arguments in model fitting functions, in addition to the formula are *weights*, *data* (indicating the data frame in which to interpret the variable names), *subset* (for using a subset

of data), and *na.action*. Note that the default *na.action* in R is set in *options()**\$na.action* and is *na.omit*, so be wary in fitting models in that you are dropping cases with NAs and may not be aware of it.

There is some more specialized syntax given in *R-intro.pdf* on CRAN.

7 Errors

When writing functions, and software more generally, you'll want to warn the user or stop execution when there is an error and exit gracefully, giving the user some idea of what happened. The *warning()* and *stop()* functions allow you to do this; in general they would be called based on an if statement. *stopifnot()* can stop based on checking multiple conditions.

You can control what happens when a warning occurs with *options()**\$warning*. This can be helpful for debugging - e.g., you can force R to stop if a warning is issued rather than continuing so you can delve into what happened.

Also, sometimes a function you call will fail, but you want to continue execution. For example, suppose you are fitting a bunch of linear models and occasionally the design matrix is singular. You can wrap a function call within the *try()* function (or *tryCatch()*) and then your code won't stop. You can also evaluate whether a given function call executed properly or not. Here's an example of fitting a model for extreme values:

```
> library(ismev)
> n = 100; nDays = 365
> x = matrix(rnorm(nDays * n), nr = nDays)
> x = apply(x, 2, max)
> x = cbind(rep(0, 100), x)
> params = matrix(NA, nr = ncol(x), nc = 3)
> for(i in 1:ncol(x)) {
  fit = try(gev.fit(x[,i], show = FALSE))
  if(!is(fit, "try-error")) params[i, ] = fit$mle
}
```

Challenge: figure out how to use *tryCatch()* to deal with the error above. Note that I haven't used it and it seemed somewhat inscrutable on quick look.

8 Object-oriented programming (OOP)

Popular languages that use OOP are C++, Java, and Python. In fact C++ is the object-oriented version of C. Different languages implement OOP in different ways.

The idea of OOP is that all operations are built around objects, which have a class and methods that operate on objects in the class. Classes are constructed to build on (inherit from) each other, so that one class may be a specialized form of another class, extending the components and methods of the simpler class (e.g., *lm* and *glm* objects).

Note that in more formal OOP languages, all functions are associated with a class, while in R, only some are.

Often when you get to the point of developing OOP code in R, you're doing more serious programming, and you're going to be acting as a software engineer. It's a good idea to think carefully in advance about the design of the classes and methods.

8.1 S3 approach

S3 classes are widely-used, in particular for statistical models in the *stats* package.

Inheritance Let's look at the *lm* class, which builds on lists, and *glm* class, which builds on the *lm* class. Here '*mod*' is an object (an instance) of class '*lm*'. An analogy is the difference between a random variable and a realization of that random variable.

```
> yb = sample(c(0, 1), 10, replace = TRUE)
> yc = rnorm(10)
> x = rnorm(10)
> mod1 = lm(yc ~ x)
> mod2 = glm(yb ~ x, family = binomial)
> class(mod1)
> class(mod2)
> is.list(mod)
> names(mod)
> is(mod2, "lm")
> methods(class = 'lm')
```

Creating our own class We can create an object with a new class (often a special case of a list - "inheriting" from a list) as follows:

```
> me = list(firstname = 'Chris', surname = 'Paciorek', age = NA)
> class(me) = 'indiv' # it turns out there is already a 'person'
class in R
```

Actually, if we want to create a new class that we'll use again, we want to create a *constructor* function that initializes new individuals:

```

> indiv = function(firstname = NA, surname = NA, age = NA) {
+ obj = list(firstname = firstname, surname = surname, age = age)
+ class(obj) = 'indiv'
+ return(obj)
}
> me = indiv('Chris', 'Paciorek')

```

Methods The real power of OOP comes from defining methods. For example,

```

> mod = lm(y ~ x)
> summary(mod)
> gmod = glm(y ~ x)
> summary(gmod)

```

Here *summary()* is a generic method (or generic function) that, based on the type of object given to it (the first argument), dispatches a class-specific function (method) that operates on the object. This is convenient for working with objects using familiar functions. Consider the generic methods *plot()*, *print()*, *summary()*, *['*, and others. We can look at a function and easily see that it is a generic method:

```
> mean
```

We can also see what classes have methods for a given generic method:

```
> methods(mean)
```

In many cases there will be a default method (here, *mean.default()*), so if no method is defined for the class, R uses the default. Sidenote: arguments to a generic method are passed along to the selected method by passing along the calling environment.

To define a generic method:

```
> summarize = function(object, ...) UseMethod("summarize")
```

Once *UseMethod()* is called, R searches for the specific method associated with the class of '*object*' and calls that method, without ever returning to the generic method. Let's try this out on for our individual class. In reality, we'd write either *summary.indiv()* or *print.indiv()* (and of course the generics for *summary* and *print* already exist) but for illustration, I wanted to show how we would write both the generic and the specific method.

```

> summarize.indiv = function(object)
+ return(with(object, cat("Person of age ", age, " with name ",
firstname, " ", surname, ".\n", sep = "")))

```

Note that the *print()* function is what is called when you simply type the name of the object, so we can have the object information printed out in a structured way. Recall that the output when we type the name of an *lm* object is NOT simply a regurgitation of the elements of the list - rather

`print.lm()` is called.

The `summary()` method generally doesn't actually print out information; rather it computes things not stored in the original object and returns it as a new class (e.g., class `summary.lm`), which is then automatically printed, per my comment above, unless one assigns it to a new object:

```
> out = summary(mod)
> out
> print(out)
```

More on inheritance As noted with `lm` and `glm` objects, we can assign more than one class to an object.

```
> class(me) = c('BerkeleyIndiv', 'indiv')
> summary(me) # note that this still works,
though I have not defined summarize.BerkeleyIndiv()
```

The classes should nest within one another with the more specific classes to the left, e.g., here a `BerkeleyIndiv` would have some additional objects on top of those of an individual, perhaps `CalnetID`, and perhaps additional or modified methods. `BerkeleyIndiv` inherits from `indiv`, and R uses methods for the first class before methods for the next class(es), unless no such method is defined for the first class. If no methods are defined for any of the classes, R looks for `method.default()`, e.g., `print.default()`, `plot.default()`, etc..

Class-specific operators We can also use operators with our classes.

```
> methods(`+`)
```

This example will be a bit silly (it would make more sense with a class that is a mathematical object) but indicates the power of having methods:

```
> `+`.indiv` = function(object, incr) {
+ object$age = object$age + incr
+ return(object)
+ }
> old.me = me + 15
```

Why use class-specific methods? We could have implemented different functionality (e.g., for `summary()`) for different objects using a bunch of `if` statements (or `switch()`) to figure out what class of object is the input, but then we need to have all that checking. Furthermore, we don't control the `summary()` function. The OOP framework makes things *extensible*, so we can build our own new functionality on what is already in R.

Final thoughts Recall the *Date* class we discussed in Unit 4. This is an example of an S3 class, with methods such as *julian()*, *weekdays()*, etc.

Challenge: how would you get R to quit immediately, without asking for any more information, when you simply type 'q' (no parentheses!)?

What we've just discussed are the old-style R (and S) object orientation, called S3 methods. The new style is called S4 and we'll discuss it next. S3 is still commonly used, in part because S4 can be slow (or at least it was when I last looked into this a few years ago). S4 is more structured than S3.

8.2 S4 approach

S4 methods are used a lot in *bioconductor*. They're also used in *lme4*, among other packages. Tools for working with S4 classes are in the *methods* package.

Note that components of S4 objects are obtained as **object@component** so do not use the usual list syntax. The components are called *slots*, and there is careful checking that the slots are specified and valid when a new object of a class is created. You can use the prototype argument to *setClass()* to set default values for the slots. There are default constructor (the function is *initialize()*) and print methods, but you can override these. One can create methods for operators and for replacement functions too.

Let's re-consider our *indiv* class using S4 in the demo code file.

As with S3, classes can inherit from one or more other classes. Chambers calls the class that is being inherited from a *superclass*. For example suppose we had spatially-indexed time series. We could have a time series class, a spatial location class, and a "location time series" class that inherits from both. Be careful that there are not conflicts in the slots or methods from the multiple classes. For conflicting methods, you can define a method specific to the new class to deal with this. Also, if you define your own *initialize()* method, you'll need to be careful that you account for any initialization of the superclass(es) and for any classes that might inherit from your class (see Chambers, p. 360).

You can inherit from other S4 classes (which need to be defined or imported into the environment in which your class is created), but not S3 classes. You can inherit (at most one) of the basic R types, but not environments, symbols, or other non-standard types. You can use S3 classes in slots, but Chambers says this requires that the S3 class be declared as an S4 class. To create S4 versions of S3 classes use *setOldClass()* - this creates a virtual class. This has been done, for example, for the *data.frame* class:

```
> showClass("data.frame")
```

You can use *setClassUnion()* to create what Adler calls *superclass* and what Chambers calls a

virtual class that allows for methods that apply to multiple classes. So if you have a person class and a pet class, you could create a “named lifeform” virtual class that has methods for working with name and age slots, since both people and pets would have those slots. You can’t directly create an object in the virtual class.

9 The R interpreter

Parsing When you run R, the R interpreter takes the code you type or the lines of code that are read in a batch session and parses each statement, translating the text into functional form. It substitutes objects for the symbols (names) that represent those objects and evaluates the statement, returning the resulting object. For complicated R code, this may be recursive.

Since everything in R is an object, the result of parsing is an object that we’ll be able to investigate, and the result of evaluating the parsed statement is an object.

We’ll see more on parsing in the next section.

.Primitive and .Internal Some functionality is implemented internally within the C implementation that lies at the heart of R. If you see *.Internal* or *.Primitive*, in the code of a function, you know it’s implemented internally (and therefore generally very quickly). Unfortunately, it also means that you don’t get to see R code that implements the functionality, though Chambers p. 465 describes how you can look into the C source code.

```
> plot.xy
> ‘%*%’
```

Evaluating a function What happens when an R function is evaluated? The function arguments (the user inputs) are matched to the argument names in the function definition. A new environment is created, and assignment to those names is done in the environment, including any default arguments. The body of the function is evaluated in the environment.

One additional concept that it’s useful to be aware of is the idea of a promise object. When R matches user input arguments to formal argument names, it does not (usually) evaluate the arguments until they are needed, which is called *lazy evaluation*. Instead the formal arguments are of a special type called a *promise*. Let’s see lazy evaluation in action:

```
> f = function(x) print("hi")
> system.time(mean(rnorm(100000)))
> system.time(f(3))
> system.time(f(mean(rnorm(100000))))
```

10 Computing on the language

10.1 Parsing code and understanding language objects

R code is just text and we can actually write R code that will create or manipulate R code. We can then evaluate that R code using `eval()`.

`quote()` will parse R code, but not evaluate it. This allows you to work with the code rather than the results of evaluating that code.

```
> quote(if (x > 1) "orange" else "apple")
```

The `print()` method for language objects is not very helpful! But we can see the parsed code by treating the result as a list.

```
> obj = quote(if (x > 1) "orange" else "apple")
> as.list(obj)
> class(obj)
```

Officially, the name that you assign to an object (including functions) is a *symbol*.

```
> x = 3; typeof(quote(x))
```

To access symbols that involve special syntax (such as special characters), you use backquotes:

```
> ? `<-`
```

We can create an *expression* object that contains R code as

```
> myExpr = expression(x <- 3)
> eval(myExpr)
> typeof(myExpr)
```

The difference between `quote()` and `expression()` is basically that `quote()` works with a single statement, while `expression()` can deal with multiple statements, returning a list-like object of parsed statements. Both of them parse R code.

Table 1 shows the *language* objects in R; note that there are three classes of language objects: *expressions*, *calls*, and *names*.

| | Example syntax to create | Class | Type |
|----------------|--------------------------------------|------------|-----------------------|
| functions | <code>function() {}</code> | function | closure |
| object names | <code>quote(x)</code> | name | symbol (language) |
| expressions | <code>expression(x <- 3)</code> | expression | expression (language) |
| function calls | <code>quote(f())</code> | call | language |
| if statements | <code>quote(if(x < 3) y=5)</code> | if (call) | language |
| for statement | <code>quote(for(i in 1:5) {})</code> | for (call) | language |
| assignments | <code>quote(x <- 3)</code> | <- (call) | language |
| operators | <code>quote(3 + 7)</code> | call | language |

Basically any standard function, operator, if statement, for statement, assignment, etc. are function

calls and inherit from the *call* class.

Objects of type language are not officially lists, but they can be queried as such. You can convert between language objects and lists with *as.list()* and *as.call()*.

An official expression is one or more syntactically correct R statements. When we use *quote()*, we're working with a single statement, while *expression()* will create a list of separate statements (essentially separate call objects). As I've said, I'm trying to use the term *statement* to refer colloquially to R code, rather than using the term *expression*.

```
> e0 = quote(3)
> e1 = expression(x <- 3)
> elm = expression({x <- 3; y <- 5})
> e2 = quote(x <- 3)
> e3 = quote(rnorm(3))
> class(e0); typeof(e0)
> class(e1); typeof(e1)
> class(e1[[1]]); typeof(e1[[1]])
> class(elm); typeof(elm)
> class(e2); typeof(e2)
> identical(e1[[1]], e2)
> class(e3); typeof(e3)
> e4 = quote(-7)
> class(e4); typeof(e4) # huh? what does this imply?
> as.list(e4)
```

We can evaluate language types using *eval()*:

```
> eval(e1)
> eval(e2)
> e2list = as.list(e2); elm1list = as.list(elm)
> eval(as.call(e2list)) # here's how to do it if the language object
is actually a list
> eval(as.expression(elm1list))
```

Let's look at the demo code file to delve more into parsing a language object. We'll be able to get a sense from this of how R evaluates code. We see that when R evaluates a parse tree, the first element says what function to use and the remaining elements are the arguments. But in many cases one or more arguments will themselves be call objects, so there's recursion.

10.2 Manipulating the parse tree

Of course since the parsed code is just an object, we can manipulate it, i.e., compute on the language:

```
> out = quote(y <- 3)
> out[[3]] = 4
> eval(out)
```

Here's another example:

```
> e1 = quote(2 + 5)
> e2 = quote(plot(x, y))
> e2[[1]] = '+'
> eval(e2)
> e1[[3]] = e2
> e1
> class(e1[[3]]) # note the nesting
> eval(e1) # what should I get?
```

We can also turn it back into standard R code, as a character, using *deparse()*, which turns the parse tree back into R code as text:

```
> codeText = deparse(out)
> parsedCode = parse(text = codeText) # works like quote except
on the code in the form of a string
> eval(parsedCode)
> deparse(quote(if (x > 1) "orange" else "apple"))
```

Note that the quotes have been escaped since they're inside a string.

To convert between names of objects as text to names that R interprets as symbols referring to objects:

```
> x3 = 7
> i = 3
> as.name(paste('x', i, sep=""))
> eval(as.name(paste('x', i, sep="")))
> deparse(quote(x))
```

10.3 Parsing replacement expressions

Let's consider replacement expressions. See the demo code file for some examples of parsing a replacement expression.

Now let's consider what was going on. It turns out that when R evaluates a call to '`<-`', if the first argument is a name, then it does the assignment, but if the first argument is a call then it calls the appropriate replacement function. The second argument (the value being assigned) is evaluated first.

This explains how the different parse trees we saw in the demo code give the same result - ultimately, the replacement function is used.

10.4 `substitute()`

The `substitute` function acts like `quote()`:

```
> identical(quote(z <- x^2), substitute(z <- x^2))
```

but if you also pass `substitute()` an environment, it will replace symbols with their object values in that environment.

```
> e = new.env(); e$x = 3
> substitute(z <- x^2, e)
```

This can do non-sensical stuff:

```
> e$z = 5
> substitute(z <- x^2, e)
```

Let's see a practical example of substituting for variables in statements:

```
> plot(x = rnorm(5), y = rgamma(5, 1)) # how does plot get the axis
label names?
> xlabel <- if(!missing(x)) deparse(substitute(x))
```

So what's going on is that within `plot.default()`, it substitutes in for '`x`' with the statement that was passed in as the `x` argument, and then uses `deparse()` to convert to character. The fact that `x` still has '`rnorm(5)`' associated with it rather than the actual 5 normal deviates (based on already having evaluated `rnorm(5)`) has to do with lazy evaluation and promises.

More generally, we can substitute into `expression` and `call` objects by providing a named list (or an environment) - the substitution happens within the context of this list.

```
> substitute(a + b, list(a = 1, b = quote(x)))
```

Things can get intricate quickly:

```
> e1 = quote(x + y)
> e2 = substitute(e1, list(x = 3))
```

The problem is that `substitute()` doesn't evaluate its first argument, "`e1`", so it can't replace the parsed elements in `e1`. Instead, we'd need to do the following, where we force the evaluation of `e1`:

```
> e2 = substitute(substitute(e, list(x = 3)), list(e = e1))
```

```
> substitute(substitute(e, list(x = 3)), list(e = e1)) # so e1 is
substituted as an evaluated object, which then allows for substitution
for 'x'
```

```
> e2
```

```
> eval(e2)
```

If this subsection is confusing, let me assure you that it confused me too. The indirection going on here is very involved.

10.5 Final thoughts

Challenge: figure out how a *for* loop is parsed in R. See how a *for* loop with one statement within the loop differs from one with two or more statements.

We'll see *expression()* again when we talk about inserting mathematical notation in plots.

11 Programming concepts

A few thoughts on what we've learned that is germane beyond R. Among other things:

- variable types
- passing by reference and by value
- variable scope
- the call stack
- flow control
- object-oriented programming
- matrix storage concepts
- parsing

If you go to learn other languages, these and other concepts that we've covered should be helpful.

Graphics

October 17, 2011

We'll focus on implementing graphics in R, but will also discuss some general concepts relevant more broadly than R. I won't focus much on specific types of plots but rather on methods for formatting graphics that are useful regardless of the specific types of plots you create.

References:

- Adler
- Chambers
- Venables and Ripley, Modern Applied Statistics with S
- Murrell, R Graphics (available electronically through OskiCat: <http://uclibs.org/PID/110697>)
- [R intro manual](#) (R-intro) on CRAN

If you are not comfortable with basic R plots, your initial homework assignment is to look through the “Unit 8 - Overview of R graphics commands” pdf (in particular the first page) and through Chapter 14 of Adler and play around a bit with some of the core plotting functions. These include high-level functions such as *plot()*, *matplot()*, *pairs()*, *coplot()*, *hist()*, *density()*, and *boxplot()*, as well as low-level functions for adding information to a plot such as *lines()*, *points()*, *abline()*. Adler goes into great detail on many of these, both in the context of the *graphics* package and that of the *lattice* package.

1 Good practices for graphics

1. Have a high density of information to space
2. Avoid 3-d graphics unless they truly add information
3. Strategies for going beyond two dimensions

- (a) Use color, but avoid if unnecessary (saves printing costs in journal articles)
 - (b) Use varying symbol or line types
 - (c) Use multiple panel plots
4. Avoid stacked barplots (see demo code) as it's hard to assess anything except the total and the baseline category.
 5. Studies indicate that humans have a hard time comparing areas, volumes, or angles, so try to avoid plots that are based on such comparisons (such as pie charts).
 6. Label axes
 7. Keep the range of axes (and other features) the same for multiple panels, when possible
 8. Use a legend where appropriate

2 Controlling standard R graphics

2.1 Background

R has three general graphics packages: *graphics*, *grid*, and *lattice*. *graphics* is the standard package, and we'll focus on it, with some discussion of *lattice*. *grid* is more object-oriented, and is quite powerful and useful if you're involved in serious graphics programming. In fact, *lattice* functions are built based on *grid*. More details on *grid* can be found in Murrell's *R Graphics*.

The basic components of R graphics are

1. high-level functions (e.g., *plot()*, *boxplot()*, etc.) for producing an entire plot
2. low-level functions for adding components to existing plots. The ability to build up a complicated plot piecewise is one of the strengths of R's graphics. These low-level functions include *abline()*, *arrows()*, *axis()*, *legend()*, *lines()*, *points()*, *rug()*, *text()*, *mtext()*, *title()*, *symbols()*.
3. graphics parameters controlled through (1) *par()* or (2) as arguments to a graphics function call. These parameters change the appearance of a plot or plots (e.g., *mai*, *pch*, *col*)

R graphics work by sequentially plotting graphics elements, so subsequent items may paint over initial elements. One technique when this is an issue is making your colors transparent - see the color subsection below.

Note that as we saw when talking about OOP, the core *plot()* function does both scatterplots and is a generic method used for plotting a wide variety of objects, e.g., *plot.lm()*.

I won't say much about the high-level and low-level functions, but a wide range of these exist. In particular low-level functions allow you to add almost anything you might want to add: arbitrary lines, polygons, symbols, text, arrows, boundary lines in maps, etc.

2.2 Graphics devices

Graphics are plotted on a *device*. In the old days when computer monitors were not high resolution or in color, this referred to a physical device, but nowadays this is a general term that denotes the context in which the plot is being made: typically on screen or as a file in a particular file format. The standard device in a UNIX environment is X11, basically a graphics window set up in the X11 windowing system. On-screen plotting is generally done with R interacting with a window manager for the operating system, so R is not interacting directly with the physical display. Often one needs to iterate to get a plot to look good when printed to a file; in particular the aspect (width to height ratio) (e.g., you can specify width and height in *pdf()*), the margin sizes relative to the size of the core plot, and size of plotting symbols and text relative to the size of the plot. That is the relative sizes when seen in a graphics window on the screen may be very different when printed to a file.

You can have multiple graphics windows open at once; you'll need to explicitly call the function for opening a device to set up any additional ones. *dev.cur()* tells the number of the active one and *dev.set()* allows to change it.

2.3 Graphics parameters

We can set graphics parameters through various plotting functions (for temporary changes applying only to the current plotting command) or through the *par()* function (for a permanent change), allowing us to customize the layout or appearance of a plot. Let's take a look at some of these and what they mean.

Most of you know how to create a multi-panel plot:

```
> par(mfrow=c(4, 2)) # 4 rows, 2 columns of subplots
```

Changing the margin sizes (and axis information spacing) is one of the most common modifications one needs to do, particularly when producing a multi-panel plot and when creating an output file. One often wants to reduce the size of the inner margins. This helps maximize the information to white space and increase the resolution of your plot. Sometimes one needs extra space in the outer margin of a multi-panel plot. Positioning and lengths within the graphics window can

occur with relative units; i.e., the device domain is $(0, 1) \times (0, 1)$, or in physical units (commonly inches), or in “lines” of text. Some values I commonly use are

```
> par(mai = c(.5, .5, .1, .1)) # manipulate inner margins of subplots
> par(mgp = c(1.8, .7, 0)) # manipulate spacing of axis ticks, axis
labels, axis text
> par(omi = c(0, 0, .3, 0)) # manipulate outer margin of full plot
```

Here for the margins, the first number is the bottom, the next the left, the third the top, and the fourth the right margin, so with *omi*, I’ve made some space (0.3 inches) in the outer margin at the top.

Note that if you change graphics options using *par()* (including within a function call!), the values change permanently (so it’s like pass by reference). One way to be able to go back is

```
> oldpar = par(no.readonly = TRUE)
> par(cex = 3); plot(x, y)
> par(oldpar); plot(x, y)
```

So if you create a function that changes the graphics parameters, you should use the above strategy to reset so you don’t surprise your users.

In the last page of this pdf, you’ll see a graphical representation of the various layout parameters and their units: *din* (device size, inches), *pin* (plot size, inches), *fin* (figure size, inches), *mai* (margin size, inches) or *mar* (margin size, lines), *omi* (outer margin, inches) or *oma* (outer margin, lines), *mex* (# text lines per interline spacing), and *plt* (plot region as fraction of figure region). *cex* controls size of points and text in general when called through *par()*, but only controls the size of points when called within a plotting function, while *cex.lab*, *cex.axis*, and *cex.main* control the size of axis labels, axis values, and the title. Similarly for *col*, *col.lab*, *col.axis* and *col.main*.

Here are some other things you can control: text justification (*adj*), font size (*cex*, *csi*), font type (*font*), rotation of text (*srt*), color (*col*), line type (*lty*), line width (*lwd*), plotting character symbol (*pch*), type/presence of boundary box (*bty*), log-scale axes (*log*), axis labels (*{x,y}lab*), axis limits (*{x,y}lim*), a variety of details about the axis limits, labels, and ticks (*lab*, *las*, *tck*, *{x,y}axp*, *{x,y}axs*, *{x,y}axt*), and whether to plot axes (*axes*).

Some additional tidbits You can force subplots to have the same axis ranges by manipulating *xlim* and *ylim*.

It can be handy to put axis labels in the outer margin of a multi-panel plot:

```
> x = rnorm(10); y = rnorm(10); par(mfrow = c(2, 2))
> for(i in 1:4) plot(x, y, xlab = "", ylab = "")
> mtext("my x variable", 1, line = -1, outer = TRUE)
> mtext("my y variable", 2, line = -1, outer = TRUE)
```

Note if we wanted to put the label further towards the edge of the plot, we'd need to use *omi* or *oma* to create an outer margin, and then we could use a value of *line* that is greater than -1 to put the text further away from the plotting regions.

You can create multi-line text by just using “\n” in your character string.

To plot outside the plot region, set **xpd = TRUE** - otherwise anything outside the region is “clipped”.

2.4 Adding information to a plot sequentially

There are lots of low-level functions you can use to add components to a plot. These include *abline()*, *arrows()*, *axis()*, *legend()*, *lines()*, *points()*, *rug()*, *text()*, *mtext()*, *title()*, *symbols()*, *hline()*, *vline()*. A basic strategy for customizing components of a plot is to use the high-level function to plot the basics, specifying parameter arguments that leave out some components (e.g., **xaxt = "n"**, **bty = "n"**, **xlab = ""**). At its most extreme, you can create the structure of a plot without plotting any data for customizing everything:

```
> plot(x, y, type= 'n')
```

With *title()*, you can add main and subtitles and axis labels, as we saw in the demo code for the previous subsection. This gives you control over positioning using the *line* and *outer* arguments. *rug()* will add tick marks for each observation along a chosen axis. *locator()* allows you to choose a location for certain items (e.g., a legend) with your mouse.

You can create custom axes with the *axis()* function, adding the axis, ticks, and axis value labels. Here's an example of plotting two sets of data with different axis values. It also shows the use of the new argument to overplot on an existing plot. Often one would initially exclude (some) axes and a boundary box from the plot:

```
> x = rnorm(10); y = rnorm(10); z = rnorm(10, 0, 5)
> plot(x, y, yaxt = 'n', bty = 'n', col = 'blue')
> box() # add box in black
> axis(side = 2, col = 'blue') # add y-axis, but in blue
> par(new = TRUE) # this tells R to overplot the next plot on the
current
> plot(x, z, col = 'red', bty = 'n', ann = FALSE, yaxt = 'n', yaxt
= 'n')
> axis(side = 4, col = 'red') # add right-side y axis
```

Adding an additional dimension Let's consider other ways to pack information on another dimension into a plot:

```

> groups = sample(1:3, 10, replace = TRUE) # fake strata
> plot(x, y, col = groups)
> plot(x, y, pch = groups)

```

You can also use this approach with different line widths (*lwd*) and line styles/types (*lty*) when that is relevant.

2.5 Interacting with graphics

R is not well set up for interactive work, but one handy function is *identify()*, which is handy for both identifying points (such as outliers) and marking a specific point.

```

> plot(x, y)
> identify(x, y) # now left-click on individual points and R will
tell you the index of the point

```

If you want to do more with interactive graphics (which can be very helpful for high dimensional data), check out the *GGobi* software.

2.6 Using mathematical expressions in plots

We can use Greek letters and mathematical symbols in labels using *expression()* to enclose the mathematical expression and *paste()* to combine elements.

```

> plot(x, y, xlab = expression(lambda),
+ ylab = expression(paste(E[lambda]^(2-Lambda), "(deviance) - ", Var(phi)
)))

```

Here's an example of using a variable within an expression. This is quite powerful if you are producing multiple plots and the text changes for each plot:

```

> plot(1:10, type="n", xlab="", ylab="", main = "plot math & numbers");
theta <- 1.23

```

```

> mtext(bquote(hat(theta) == .(theta)), side = 3, line = 1.5)

```

bquote() treats its argument as an R language object but substitutes in for terms wrapped in *.(.)*, behaving somewhat like *substitute()*. An alternative approach uses *substitute()*:

```

> mtext(substitute(paste("Estimate: ", hat(theta) == thetaVal),

```

```

+ list(thetaVal = thetaVal)), side = 1, line = 2)

```

For a full set of mathematical syntax and examples of such substitutions, see **?plotmath** as well as **example(plotmath)** and **demo(plotmath)**.

For plotting text, it's helpful to keep in mind that we can choose a display of text based on the encoding of the characters and a font in which to display the characters. So if you need characters

not in the standard Latin character scheme, it should be doable by choosing a different encoding.

2.7 Laying out panels (subplots)

We can create more complicated arrangements of panels than the rectangular layout of *mfrow* and *mfc*.

One approach is use *layout()*. Here we give it a matrix with as many rows and columns as we want in the plot and assign values, 1, 2, ... to the appropriate rows and columns of the matrix, with the values serving to identify each panel.

```
> layout(matrix(c(1, 1, 0, 2), nr = 2, byrow = TRUE))
> layout.show(n = 2)
> layout(matrix(c(1, 1, 1, 1, 4, 3, 2, 2), nr = 2, byrow = TRUE))
> layout.show(n = 4)
```

Let's see the example from **?layout**, which slickly plots a scatterplot with marginal histograms.

We can use *split.screen()* as an alternative to *layout()*. Here's a basic example that allows us to choose which panel to plot in:

```
> split.screen(figs = c(2, 3))
> screen(3)
> plot(x, y)
> hist(x) # so we need to change screens manually
```

split.screen() can be used with a 4-column matrix where each row indicates the xrange and yrange (xl, xu, yl, yu) of the given subplot where xl and xu are the lower and upper x-axis values and yl and yu for the y-axis values. The values should all be in between 0 and 1. Note that overplotting is possible.

```
> split.screen(figs = matrix(c(0, .4, 0, .5, .7, 1, .7, 1, .7, 1,
.4, .7, .2, 1, 0, .4, 0, .4, .7, 1), nc = 4, byrow = TRUE))
> screen(1); plot(x, y); screen(2); hist(x); screen(3); hist(y);
screen(4); plot(1:length(x), x, type = "l")
> screen(5)
> split.screen(figs = matrix(c(0, .5, 0, 1, .5, 1, 0, 1), nr = 2,
byrow = TRUE)) # you can continue subdividing, but need to be careful
that the margins don't get too big for the remaining space (and that
text is not too large)
> screen(6); par(mai = c(.05, .05, .05, .05), cex = .7); hist(x);
> screen(7); par(mai = c(.05, .05, .05, .05), cex = .7); hist(y)
```

```
> close.screen(all = TRUE)
```

In the *grid* package, positioning subpanels is done via *viewports*.

3 Lattice graphics

The *lattice* package provides an alternative set of functions based on the Trellis graphics system. The style is prescriptive, which in some cases is good (think *Latex* vs. *Microsoft Word*). Many people like the design of lattice graphics, and they're also well set up for dealing with conditioning either through sets of panels or through grouping by color or symbol. In the case of multi-panel plots, the style exploits the inter-relatedness of the information, with 'strips' giving labelling information, and avoidance of repetition of axis information.

Fig. 1 on the next page shows example from a paper of mine that (hopefully) illustrates the power of lattice graphics. I'm showing how the bias of an estimator varies with two parameters ranging in (0,1) (θ_c and θ_u on the x and y axes) and with three different levels of parameters p_c and p_z distinguishing the 9 subplots. Notice that labels get re-used and space is used efficiently because of this and avoiding space between the subplots.

Getting the strip text right took some doing, but the core of the plot is generated by:

```
> levelplot(vals ~ theta.c * theta.u | pc * pz, data = myData,
+ xlim=c(-.05,1.05), col.regions=tim.colors(), cuts=10,
+ at=seq(-.05,1.05,by=.1))
```

Notice the conditioning bar, which says to plot *vals* as a function of *theta.c* and *theta.u* within levels of *pc* and *pz*.

Modularity and the `panel()` function Lattice functions are highly modular and have consistent argument names (see p. 308 of Adler), unlike in *graphics*. The high-level functions take care of the overall logistics while low-level panel functions do the actual plotting within each panel. To add features to the panels, you need to do it through the *panel()* function argument:

```
> library(MASS)
> out = xyplot(time ~ dist, data = hills,
+ panel = function(x, y, ...){
+ panel.xyplot(x, y, ...)
+ panel.lmline(x, y, type = "l")
+ panel.loess(x, y, ...)
+ })
> splom(~ hills,
```

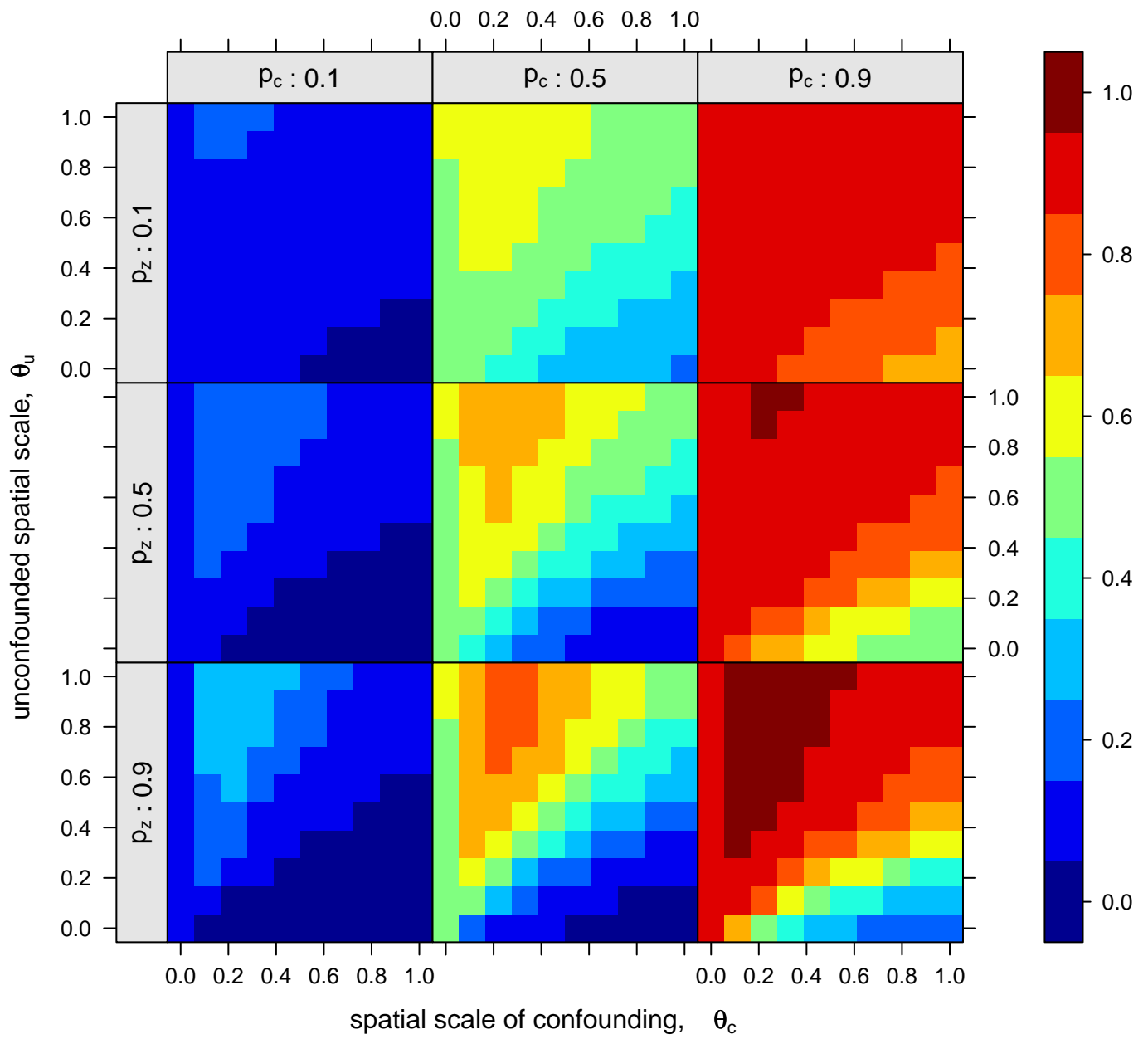


Figure 1. Example of lattice graphics.

```

+ panel = function(x, y, ...){
+ panel.xyplot(x, y, ...)
+ panel.loess(x, y, ...)
+ })

```

Conditioning plots (via panels) Lattice does nice job with conditioning plots.

```

> xyplot(Petal.Length ~ Sepal.Length | Species, data = iris, layout
= c(2, 2)) # conditioning on a factor

```

When conditioning on a continuous variable, we need to work with shingles:

```

> state$popShingle = equal.count(state$Population, number = 3, overlap
= 0)
> state$FrostShingle = equal.count(state$Frost, number = 3, overlap
= .25)
> xyplot(Income ~ LifeExp | popShingle * FrostShingle, data = state)

```

Note the use of a formula a la statistical models in R.

We can do something similar in graphics using *coplot()*; Note that *coplot()* takes a panel argument as well, so you can customize what is done in each of the panels.

Side note: we can use also formulas in the *graphics* high-level functions:

```

> par(mfrow=c(1,2)); plot(mpg ~ disp + hp , data = mtcars)

```

Conditioning plots (via color or symbol type) Alternatives for conditioning use color or symbol type: first a clunky version using the standard graphics and then the lattice version:

```

> mapping = data.frame(Species = unique(iris$Species), pch = c(1,
2, 4), col = c('red', 'purple', 'green'))
> mapping$col = as.character(mapping$col) # arggh, automatic coercion
bites us again
> iris2 = merge(iris, mapping, all.x = TRUE)
> with(iris2, plot(Sepal.Length, Petal.Length, pch = pch))
> with(iris2, plot(Sepal.Length, Petal.Length, col = iris2$col))
> legend(x = 4.5, y = 7, legend = unique(iris$Species), col = mapping$col,
pch = 1)

```

The lattice version is easier:

```

> xyplot(Petal.Length ~ Sepal.Length, groups = Species, data = iris,
auto.key = TRUE)

```

We can manipulate subpanels in a lattice graphic:


```
> xyplot(Petal.Length ~ Sepal.Length | Species, iris, layout = c(2,
2))
> trellis.focus("panel", 1, 2)
> do.call("panel.lmline", trellis.panelArgs())
```

Lattice graphics are optimized for the current plotting device, so using *dev.copy2pdf()* and the like is a bad idea since the graphic will have been optimized for the computer screen window. Instead use *pdf()* and the like to make the plot within the device.

The analog of *par()* for lattice is *trellis.par.set()* and *trellis.par.get()*. Often you can get a nice lattice graphic without much work, but modifying aspects of a lattice graphic may take some work. Chambers Chapter 7 discusses this to some degree. Note that *lattice* graphics are based on the *grid* package, so serious monkeying may require learning something about *grid*.

Final note: formally, lattice functions produce objects and the print method applied to the object produces the plot.

4 Graphics file formats

We've already seen *pdf()* and *dev.copy2pdf()* in the R intro unit. Analogues of *pdf()* include *postscript()*, *png()*, *tiff()*, *bmp()*, and *jpeg()*.

4.1 Vectorized vs. rasterized file formats

pdf and postscript images are vectorized - in general elements of the image are symbolic objects (such as points, text, symbols, line segments, etc.) and when an image is resized, the items rescale appropriately for the new size without losing resolution. In contrast, with a rasterized format such as JPEG, individual pixels are plotted, and when an image is rescaled, in particular enlarged, one is stuck with the resolution that one used in plotting the figure (i.e., one has the original pixels but if you zoom, you only show some of them, losing resolution). One downside to vectorized images is that with a lot of points or line segments, they can be very large. For 2-d images such as created by *image()*, rasterized formats do make some sense inherently, though the other features in the file (such as any text) is also rasterized.

As a sidenote, if you have a scatterplot with many points that creates a huge file, consider the *scatterSmooth()* function as well as the *hexbin* package. The former creates a two-d density plots with outlying individual points included, while the latter creates an empirical two-d density by binning into hexagonal areas:

```
> x = rnorm(10000); y = rnorm(10000)
> bin = hexbin(x, y)
```

```
> plot(bin)
> smoothScatter(x, y)
```

4.2 File formats for journal articles

For inserting into Latex files for journal articles, I typically use *pdf()* in combination with *pdflatex* to compile the Latex file. In some cases for large images I will use *jpeg* or *png* files for figures in a paper.

Journals sometimes request either *postscript* or *encapsulated postscript* (EPS) format images. If you create a single R plot (as opposed to multiple pages), using *postscript()*, the result should be EPS compatible. You can do this with the arguments **horizontal = FALSE**, **onefile = FALSE**, **paper = "special"**, or by calling *setEPS()*. If you have a postscript file that is not EPS, *ps2epsi* is a Ghostscript tool that can do the conversion from *ps* to *eps* (i.e., *epsi*). EPS is a standardized form of postscript and contains a bounding box describing the bounds of the image. EPSI files contain a bitmapped preview image in the “Interchange” format that systems can use to show an approximation of the image even if they can’t render (i.e., display) postscript.

Discussing much about postscript is beyond our scope, but note that postscript files are just text in a particular language that the printer uses for generating a printout. So you can use *less* or *emacs* to view the file in UNIX. You’ll see that it says that it’s EPS and has a bounding box. It’s occasionally handy to search for and replace text in a postscript file manually rather than regenerating the file. For example I could search for ‘noise’ in *example.eps* (created in the demo code) and replace with something else. I could also monkey with the positioning, which should involve the numbers just before the text string. You can also do this sort of thing in pdf files as well.

4.3 Conversion utilities

UNIX has a lot of utilities for converting between image formats. Windows and Mac also have GUI-style programs for doing this.

In UNIX, *pdftops* will convert pdf to postscript and with the optional argument *-eps* to encapsulated postscript, while *ps2epsi* will create encapsulated postscript. *gs* (Ghostscript) will do a lot of different manipulations of ps and pdf files, including converting to jpeg and other formats and merging and splitting pages of pdf files. Here are some examples of command-line calls from within a UNIX shell:

```
> gs -dNOPAUSE -r[xres]x[yres] -sDEVICE=jpeg -sOutputFile=file.jpg
file.ps # convert from ps to jpeg
> gs -sDEVICE=pdfwrite -dNOPAUSE -dQUIET -dBATCH -dFirstPage=m
```

```

+      -dLastPage=n -sOutputFile=out.pdf in.pdf # extract pages from
a pdf
> gs -dNOPAUSE -sDEVICE=pdfwrite -sOUTPUTFILE=out.pdf -dBATCH in1.pdf

+      in2.pdf in3.pdf # merge pdf files into one pdf

```

5 Colors

The default colors can be seen with `palette()`. Using `col=i` in a plot uses the *i*th element of the output of `palette()`. You can change the palette:

```
> palette(c("black", "yellowgreen", "purple"))
```

See `colors()` for the colors available by name. You can also use RGB levels, discussed next.

5.1 Colorspaces

Colors live in a 3-dimensional space that can be parameterized in several ways. One standard parameterization is RGB, which is a set of three numbers indicating the intensity of red, green and blue. We can use RGB levels to specify colors in R.

```

> rgb(0.5, 0.5, 0) # specified on scales of [0, 1]
> plot(x, y, col = rgb(0.5, 0.5, 0))
> col2rgb("yellowgreen") # on scale of {0, ..., 255}

```

Notice `rgb()` gives us back the color as a hexadecimal number (`#RRGGBB`), where each of *RR*, *GG*, and *BB* is 2-digit hexadecimal number (base 16) in the range 0 (00) to 255 (FF), so red is `#FF0000`. A string in this format can be used to specify colors and you'll run across this in R if you work with colors.

Another parameterization is HSV: *hue*, *saturation* (colorfulness metric), and *value* (brightness). Let's see the demo code to see how colors vary as we change HSV values using `rainbow()`.

A parameterization that uses a more absolute measure of colorfulness than saturation is HCL (hue, chroma, luminance). In the example in the demo code, none of the colors stands out more than the others.

5.2 Color sequences

If we're using color to illustrate a continuous range of values, we need a meaningful color sequence. To construct a continuous color set giving a sequence of colors you can use a variety of color schemes: `rainbow()`, `heat.colors()`, `terrain.colors()`, `topo.colors()`, `temp.colors()`, and (in the `fields` package), `tim.colors()`. I know Tim! He likes to fish.

The main thing to avoid is a sequence in which the colors do not appear to vary smoothly or in some cases may not even appear monotonic. Let's examine a variety of the sequences (see the demo code).

`temp.colors()` is a good blue to red “diverging” color scheme that emphasizes magnitudes around a central point, with two hues - one for each direction.

The *RColorBrewer* package is good for choosing colors for unordered levels, sequential ordering, and two-way diverging color ordering and the *ColorBrewer* website provides recommendations. We'll see an example in the section on mapping.

5.3 Overplotting and transparency

One nice trick if you have a lot of overplotting is to have your color be partly transparent, so that overplotting results in darker colors. Note that this may not work on all devices. We can specify transparency level as either the 4th number in `rgb()` on a scale of 0 (transparent) to 1 (opaque - the default), or as a fourth hexadecimal number on the scale of 00 to FF (0 to 255); e.g., `#FF000080` would be half-transparent red, since 80 is one-half of FF in base 16.

```
> x = rnorm(10000); y = rnorm(10000)
> plot(x, y, col = rgb(0, 0, 0, .2) # fairly transparent black
```

5.4 Colorblindness

One thing to be aware of is that 7-8% of men are color blind. As we see in the demo code, the standard result of this is to make it difficult to distinguish red and green, so one may want to avoid color schemes that have both of these in them. We can use `dichromat()` from the *dichromat* package to assess the effect of colorblindness on viewing of one's images; see more in the demo code file.

```
> showpal = function(colors){ n = length(colors); plot(1:n, rep(1,
n), col = colors, pch = 16, cex = 4) }
> showpal(palette())
> showpal(dichromat(palette))
```

6 Plotting in three dimensions and mapping

In general for plotting $z \sim x + y$, either in an actual spatial setting or as a function of two variables, I prefer the `image()` function (`levelplot()` in lattice), which uses color to represent different levels. `image.plot()` from the *fields* package nicely adds a color bar legend. `contour()` and `persp()` are other options but I personally find them less informative most of the time: `contour()` because it requires

the viewer to read the contour line values and *persp()* because some of the features often hide other features.

```
> n = 20; xs = ys = 1:n; gr = expand.grid(xs, ys);
> U = chol(exp(-rdist(gr)/6))
> image.plot(1:n, 1:n, matrix(crossprod(U, rnorm(n^2)), n, n), col
= tim.colors(32))
```

If you do want to make a 3-d plot of a surface, *persp()* will do it, with *x*, *y*, *z* arguments similar to *image()* and also having arguments for the angle of view (the rotation around the *z*-axis – *theta*) and the viewing angle (up or down – *phi*). You may be able to choose a set of angles that show the surface with the least amount of hidden material. See the demo code for using *image()*, *contour()*, and *persp()* to show topography in a mapping context.

R has a lot of tools for mapping in a variety of packages. In particular you can import standard GIS shape files and overlay boundaries on a map using the *spdep* package. R also has state and national boundaries as part of the *maps* package (see the demo code for more examples).

```
> plot(x, y)
> map('state', add = TRUE) # adds state boundaries
```

If you're using color in maps, *RColorBrewer* package is good for choosing colors for unordered levels, sequential ordering, and two-way diverging color ordering (see **?display.brewer.all**) and the *ColorBrewer* website provides recommendations.

Debugging and Profiling

October 8, 2011

Sources:

- Chambers
- Roger Peng's [notes](#) on debugging in R

1 Common syntax errors and bugs

Tips for avoiding bugs

1. Build up code in pieces, testing along the way.
2. Use core R functionality and algorithms already coded.
3. Remove objects you don't need, to avoid accidentally using values from an old object via the scoping rules.
4. Be careful that the conditions of *if* statements and the sequences of *for* loops are robust when they involve evaluating R code.
5. Write code for clarity and accuracy first; then worry about efficiency.
6. Code in a modular fashion, making good use of functions, so that you don't need to debug the same code multiple times.

Common syntax errors and bugs:

1. Parenthesis mis-matches
2. `[[. . .]]` vs. `[. . .]`
3. `==` vs. `=`

4. `||` vs. `|` and `&&` vs. `&`
5. Comparing real numbers exactly using `'=='` is dangerous (more in a later Unit). Suppose you generate $x = 0.333333$ in some fashion with some code and then check:

```
> x == 1/3 # FALSE is the result
```
6. Vectors vs. single values:
 - (a) You expect a single value but your code gives you a vector
 - (b) You want to compare an entire vector but your code just compares the first value (e.g., in an *if* statement) – consider using *identical()* or *all.equal()*
7. Silent type conversion when you don't want it, or lack of coercion where you're expecting it
8. Using the wrong function or variable name
9. Giving arguments to a function in the wrong order
10. In an if-then-else statement, the *else* cannot be on its own line because R will see the if-then part of the statement, which is a valid R statement, will execute that, and then will encounter the *else* and return an error. See the example in the R programming unit for syntax.
11. Forgetting to define a variable in the environment of a function and having the function, via lexical scoping, get that variable as a global variable from one of the enclosing environments. At best the types are not compatible and you get an error; at worst, you use a garbage value and the bug is hard to trace. In some cases your code may work fine when you develop the code (if the variable exists in the enclosing environment), but then may not work when you restart R if the variable no longer exists or is different.

2 Debugging Strategies

Debugging is about figuring out what went wrong and where it went wrong.

In compiled languages, one of the difficulties is figuring out what is going on at any given place in the program. This is a lot easier in R by virtue of the fact that R is interpreted and we can step through code line by line at the command line. However, beyond this, there are a variety of helpful tools for debugging R code. In particular these tools can help you step through functions and work inside of functions from packages.

2.1 Basic strategies

Read and think about the error message. Sometimes it's inscrutable, but often it just needs a bit of deciphering. Looking up a given error message in the [R mailing list archive](#) can be a good strategy.

Fix errors from the top down - fix the first error that is reported, because later errors are often caused by the initial error. It's common to have a string of many errors, which looks daunting, caused by a single initial error.

Is the bug reproducible - does it always happen in the same way at at the same point? It can help to restart R and see if the bug persists - this can sometimes help in figuring out if there is a scoping issue and we are using a global variable that we did not mean to.

Another basic strategy is to build up code in pieces (or tear it back in pieces to a simpler version). This allows you to isolate where the error is occurring.

To stop code if a condition is not satisfied, you can use `stopifnot()`, e.g.,

```
> stopifnot(is(x, "matrix"))
```

This allows you to catch errors that can be anticipated.

The `codetools` library has some useful tools for checking code, including a function, `findGlobals()`, that let's you look for the use of global variables

```
> findGlobals(lm)
> f = function() {y = 3; print(x + y)}
> findGlobals(f)
```

If you've written your code modularly with lots of functions, you can test individual functions. Often the error will be in what gets passed into and out of each function.

You can have warnings printed as they occurred, rather than saved, using `> options(warn = 1)`. This can help figure out where in a loop a warning is being generated. You can also have R convert warnings to error using `> options(warn = 2)`.

At the beginning of time (the 1970s?), the standard debugging strategy was to insert print statements in one's code to see the value of a variable and thereby decipher what could be going wrong. We have better tools nowadays.

2.2 Interactive debugging via the browser

The core strategy for interactive debugging is to use `browser()`, which pauses the current execution, and provides an interpreter, allowing you to view the current state of R. You can invoke `browser()` in four ways

- by inserting a call to `browser()` in your code if you suspect where things are going wrong
- by invoking the browser after every step of a function using `debug()`

- by using `options(error = recover)` to invoke the browser when `error()` is called
- by temporarily modifying a function to allow browsing using `trace()`

Once in the browser, you can execute any R commands you want. In particular, using `ls()` to look at the objects residing in the current function environment, looking at the values of objects, and examining the classes of objects is often helpful.

2.3 Using `debug()` to step through code

To step through a function, use `debug(nameOfFunction)`. Then run your code. When the function is executed, R will pause execution just before the first line of the function. You are now using the browser and can examine the state of R and execute R statements.

In addition, you can use “n” or return to step to the next line, “c” to execute the entire current function or current loop, and “Q” to stop debugging. We’ll see an example in the demo code.

To unflag the function so that calling it doesn’t invoke debug, use `undebug(nameOfFunction)`. In addition to working with functions you write you can use debug with standard R functions and functions from packages. For example you could do `debug(glm)`.

2.4 Tracing errors in the call stack

`traceback()` and `recover()` allow you to see the call stack (the sequence of nested function calls) at the time of an error. This helps pinpoint where in a series of function calls the error may be occurring.

If you’ve run the code and gotten an error, you can invoke `traceback()`. R will show you the stack, which can help pinpoint where an error is occurring.

More helpful is to be able to browse within the call stack. To do this invoke `options(error = recover)` (potentially in your `.Rprofile` if you do a lot of programming). Then when an error occurs, `recover()` gets called, usually from the function in which the error occurred. The call to `recover()` allows you to navigate the stack of active function calls at the time of the error and browse within the desired call. You just enter the number of the call you’d like to enter (or 0 to exit). You can then look around in the frame of a given function, entering an empty line when you want to return to the list of calls again.

You can also combine this with `options(warn = 2)`, which turns warnings into errors.

2.5 Using `trace()` to temporarily insert code

`trace()` lets you temporarily insert code into a function (including standard R functions and functions in packages!) that can then be easily removed. You can use `trace` in a few ways - here's how you would do it most simply, where by default the second argument is invoked at the start of the function given as the first argument, but it is also possible to invoke just before exiting a function:

```
> trace(lm, recover) # invoke recover() when the function starts
> trace(lm, exit = browser) # invoke browser() when the function
ends
> trace(lm, browser, exit = browser) # invoke browser() at start
and end
```

Then in this example, once the browser activates I can poke around within the `lm()` function and see what is going on.

The most flexible way to use `trace()` is to use the argument `edit = TRUE` and then insert whatever code you want wherever you want. If I want to ensure I use a particular editor, such as `emacs`, I can use the argument `edit = "emacs"`. A standard approach would be to add a line with `recover()` (or `browser()` if you don't care about the call stack). Alternatively, you can manually change the code in a function without using `trace()`, but it's very easy to forget to change things back, so `trace()` is a nice way to do things.

You call `untrace()`, e.g., `untrace(lm)`, to remove the temporarily inserted code; otherwise it's removed when the session ends.

Alternatively you can do `trace(warning, recover)` which will insert a call to `recover()` whenever `warning()` is called.

3 Memory management

3.1 Allocating and freeing memory

Unlike compiled languages like C, in R we do not need to explicitly allocate storage for objects. However, we have seen that there are times that we do want to allocate storage in advance, rather than successively concatenating onto a larger object.

R automatically manages memory, releasing memory back to the operating system when it's not needed via garbage collection. Occasionally you will want to remove large objects as soon as they are not needed. `rm()` does not actually free up memory, it just disassociates the name from the memory used to store the object. In general R will clean up such objects without a reference (i.e., a name) but you may need to call `gc()` to force the garbage collection. This uses some computation so it's generally not recommended.

In a language like C in which the user allocates and frees up memory, memory leaks are a major cause of bugs. Basically if you are looping and you allocate memory at each iteration and forget to free it, the memory use builds up inexorably and eventually the machine runs out of memory. In R, with automatic garbage collection, this is generally not an issue, but occasionally memory leaks do occur.

3.2 Monitoring memory use

There are a number of ways to see how much memory is being used. When R is actively executing statements, you can use *top* from the UNIX shell. In R, *gc()* reports memory use and free memory as *Ncells* and *Vcells*. As far as I know, *Ncells* concerns the overhead of running R and *Vcells* relates to objects created by the user, so you'll want to focus on *Vcells*. You can see the number of Mb currently used (the “*used*” column of the output) and the maximum used in the session (the “*max used*” column)”

```
> gc()
> x = rnorm(1e8) # should use about 800 Mb
> gc()
> rm(x)
> gc() # note the “max used” column
```

In Windows only, *memory.size()* tells how much memory is being used.

You can check the amount of memory used by individual objects:

```
> object.size(x)
```

The demo code file has a function, *ls.sizes()*, that wraps *object.size()* to report the largest *n* objects in a given environment.

One frustration with memory management is that if your code bumps up against the memory limits of the machine, it can be very slow to respond even when you're trying to cancel the statement with *Ctrl-C*. Leonard (one of the students in the class) pointed out that you can impose memory limits in Linux by starting R (from the UNIX prompt) in a fashion such as this

```
> R --max-vsize=1000M
```

Then if you try to create an object that will push you over that limit or execute code that involves going over the limit, it will simply fail with the message “*Error: vector memory exhausted (limit reached?)*”. So this approach may be a nice way to avoid paging by setting the maximum in relation to the physical memory of the machine. It might also help in debugging memory leaks because the program would fail at the point that memory use was increasing. I haven't played around with this much, so offer this with a note of caution.

4 Profiling

As we've seen, `system.time()` is very handy for comparing the speed of different implementations.

```
> x = matrix(rnorm(n^2), n)
> system.time({mns = rep(NA, n), for(i in 1:n) mns[i] = mean(x[i
, ])}))
> system.time(rowMeans(x))
```

The `Rprof()` function will show you how much time is spent in different functions, which can help you pinpoint bottlenecks in your code.

```
> Rprof("makeTS.prof")
> system.time(makeTS(0.1, 1000))
> Rprof(NULL)
> summaryRprof("makeTS.prof")
```

Here's the result for the `makeTS()` function from the demo code file:

```
$by.self   self.time self.pct total.time total.pct
".Call"    0.38     48.72     0.38     48.72
".Fortran" 0.22     28.21     0.22     28.21
"matrix"   0.08     10.26     0.30     38.46
"exp"      0.08     10.26     0.08     10.26
"/"        0.02      2.56     0.02      2.56
$by.total          total.time total.pct self.time self.pct
"makeTS"           0.78    100.00    0.00     0.00
".Call"            0.38     48.72    0.38    48.72
"chol.default"    0.38     48.72    0.00     0.00
"chol"            0.38     48.72    0.00     0.00
"standardGeneric" 0.38     48.72    0.00     0.00
"matrix"           0.30     38.46    0.08    10.26
"rdist"            0.30     38.46    0.00     0.00
".Fortran"         0.22     28.21    0.22    28.21
"exp"              0.08     10.26    0.08    10.26
"/"                0.02      2.56    0.02     2.56
$sample.interval [1] 0.02
$sampling.time   [1] 0.78
```

`Rprof()` tells how much time was spent in each function alone (the *by.self* bit) and aggregating the time spent in a function and all of the functions that it calls (the *by.total* bit). Usually the former

is going to be more useful, but in some cases we need to decipher what is going on based on the latter.

Let's figure out what is going on here. The self time tells us that `.Call` (a call to C code), `.Fortran` (a call to Fortran code) and `matrix()` take up most of the time. Looking at the total time and seeing in `chol.default()` that `.Call` is used and in `rdist()` that `.Fortran()` and `matrix()` are used we can infer that about 49% of the time is being spent in the Cholesky and 38% in the `rdist()` calculation, with 10% in `exp()`. As we increase the number of time points, the time taken up by the Cholesky would increase since that calculation is order of n^3 while the others are order n^2 (more in the linear algebra unit).

As another example you can try out `Rprof()` on different implementations of the random walk from PS3.

Apparently there is a memory profiler in R, `Rprofmem()`, but it needs to be enabled when R is compiled (i.e., installed on the machine), because it slows R down even when not used. So I've never gotten to the point of playing around with it.

5 Getting help online

Just searching the [R mailing list archive](#) often gives you a hint of how to fix things. An example occurred when I was trying to figure out how fix the problem with the “*multibyte string*” error in the emails in the `SpamAssassinTraining` object from PS2. I knew it had something to do with the encoding and R not interpreting the codes for non-ASCII characters correctly but I wasn't sure how to fix it. So I searched for “*invalid multibyte string*”. Around the 8th hit or so there was a comment about using `iconv()` to convert to UTF-8, which solved the problem.

If you've searched the archive and haven't found an answer to your problem, you can often get help by posting to the R-help mailing list. A few guidelines (generally relevant when posting to mailing lists beyond just the R lists):

1. Search the archives and look through relevant R books or manuals first.
2. Boil your problem down to the essence of the problem, giving an example, including the output and error message
3. Say what version of R, what operating system and what operating system version you're using.
4. Read the [posting guide](#).

The mailing list is a way to get free advice from the experts, who include some of the world's most knowledgeable R experts - seriously - members of the R core development team contribute

frequently. The cost is that you should do your homework and that sometimes the responses you get may be blunt, along the lines of “read the manual”. I think it’s a pretty good tradeoff - where else do you get the foremost experts in a domain actually helping you?

Note: of course the mailing list archive is also helpful for figuring out how to do things, not just for fixing bugs.

Computer numbers

October 26, 2011

References:

- Gentle, Computational Statistics, Chapter 2.
- <http://www.lahey.com/float.htm>
- And for more gory detail, see Monahan, Chapter 2.

A quick note that R's version of scientific notation is XeY , which means $X \cdot 10^Y$.

A second note is that the concepts developed here apply outside of R, but we'll illustrate the principles of computer numbers using R.

1 Basic representations

Everything in computer memory or on disk is stored in terms of bits. A *bit* is essentially a switch that can be either on or off. Thus everything is encoded as numbers in base 2, i.e., 0s and 1s. 8 bits make up a *byte*. For information stored as plain text (ASCII), each byte is used to encode a single character. One way to represent a byte is to write it in hexadecimal, rather than as 8 0/1 bits. Since there are $2^8 = 256$ possible values in a byte, we can represent it more compactly as 2 base-16 numbers, such as “3e” or “a0” or “ba”. A file format is nothing more than a way of interpreting the bytes in a file.

We can think about how we'd store an integer in terms of bytes. With two bytes, we could encode any value from $0, \dots, 2^{16} - 1 = 65535$. This is an unsigned integer representation. To store negative numbers as well, we can use one bit for the sign, giving us the ability to encode $-32767 - 32767 (\pm 2^{15} - 1)$. Note that in general, rather than be stored simply as the sign and then a number in base 2, integers are actually stored in a different binary encoding to facilitate arithmetic. Finally note that the set of computer integers is not closed under arithmetic:

```
> a = as.integer(3423333)
```

```
> a * a
```

and R reports an overflow.

Real numbers (or *floating points*) use a minimum of 4 bytes, for single precision floating points. Note that if you give me 4 bytes, depending on the format we are using to interpret it, it could represent 4 characters, 2 integers, or 1 single precision floating point. In general 8 bytes are used to represent real numbers and these are called *double precision floating points* or *doubles*. Let's see some examples in R of how much space different types of variables take up.

Let's see how this plays out in terms of memory use in R.

```
> x = rnorm(100000)
```

```
> y = 1:100000
```

```
> z = rep('a', 100000)
```

```
> object.size(x) # 800040 bytes
```

```
> object.size(y) # 400040 bytes - so how many bytes per integer  
in R?
```

```
> object.size(z) # 800088 bytes - hmm, what's going on here?
```

We can easily calculate the number of megabytes (Mb) a vector of floating points (in double precision) will use as the number of elements times 8 (bytes/double) divided by 10^6 to convert from bytes to megabytes. (In some cases when considering computer memory, a megabyte is $1,048,576 = 2^{20}$ bytes so slightly different than 10^6). Finally, R has a special object that tells us about the characteristics of computer numbers on the machine that R is running on called *.Machine*. For example, *.Machine\$integer.max* is $2147483647 = 2^{31} - 1$, which confirms how many bytes R is using for each integer.

2 Floating point basics

Reals (also called floating points) are stored on the computer as an approximation, albeit a very precise approximation. As an example, with a double, the error in the distance from the earth to the sun is around a millimeter. However, we need to be very careful if we're trying to do a calculation that produces a very small (or very large number) and particularly when we want to see if numbers are equal to each other.

```
> 0.3 - 0.2 == 0.1 # FALSE
```

```
> 0.3
```

```
> 0.2
```

```
> 0.1 # Hmm...
```

```
> options(digits=22)
```

```
> a = 0.3
```



```
> b = 0.2
> a - b
> 1/3
```

We notice that we can represent the result accurately only up to the 16th decimal place. This suggests no need to show more than 16 decimal places and no need to print out any more when writing to a file. And of course, often we don't need anywhere near that many. *Machine epsilon* is the term used for indicating the accuracy of real numbers and it is defined as the smallest float, x , such that $1 + x \neq 1$.

```
> 1 + 1e-16
> 1 + 1e-15
> 1 + 2e-16
> .Machine$double.eps # 2.220446049250313080847e-16
```

Floating point refers to the decimal point (or radix point since we'll be working with base 2 and *decimal* relates to 10). Consider Avogadro's number in terms of scientific notation: $+6.023 \times 10^{23}$. A real number on a computer is stored in what is basically scientific notation:

$$\pm 0.d_1d_2 \dots d_p \times b^e \quad (1)$$

where b is the base, e is an integer and $d_i \in \{0, \dots, b-1\}$. First, we need to choose the number of bits to represent e so that we can represent sufficiently large and small numbers. Second we need to choose the number of bits, p , to allocate to $d = d_1d_2 \dots d_p$, which determines the accuracy of any computer representation of a real.

Overflow and underflow The largest and smallest numbers are $b^{e_{\max}}$ and $b^{e_{\min}}$ where e_{\max} and e_{\min} are the smallest and largest possible values of e . ($b^{e_{\max}}$ is the max because we can have $0.d_1 \dots d_p \approx 1$, while $b^{e_{\min}}$ is essentially the min because we generally require $d_1 = 1$; otherwise we could just move the radix point over and change e). Let's consider the exponent and what we can infer about the range of possible numbers. Commonly a computer uses base 2, $b = 2$, because base-2 arithmetic is faster than base-10 arithmetic. For a double, we have 8 bytes=64 bits. Consider our representation as (S, d, e) where S is the sign. Commonly 1 bit is used for S and $p = 52$ bits for d . This leaves us with 11 bits for e , of which we need one for the sign of the exponent, so we can represent $\pm 2^{10} = \pm 1024$ different exponent values (see *.Machine\$double.max.exp*) (why is *.Machine\$double.min.exp* only -1022?). So the largest number we could represent is 2^{1024} since $b = 2$. What is this in base 10?

```
> log10(2^1024) # whoops ... we've actually just barely overflowed
Let's be smarter. log10(2^1024) = log2(2^1024)/log2(10) = 1024/3.32 ≈ 308. Analogously for the small-
```

est number, so we have that floating points can range between 1×10^{-308} and 1×10^{308} . Take a look at `.Machine$double.xmax` and `.Machine$double.xmin`. Producing something larger or smaller in magnitude than these values is called overflow and underflow respectively. When we overflow, R gives back an Inf or -Inf (and in other cases we might get an error message). When we underflow, we get back 0, which in particular can be a problem if we try to divide by the value. Values stored as integers should overflow if they exceed `.Machine$integer.max`.

Should 2^{45} overflow?

```
> x = 2^45
> z = 25
> class(x)
> class(z)
> as.integer(x)
> as.integer(z)
> 1e308
> as.integer(1e308)
> 1e309
```

Precision Recall our representation as (S, d, e) where S is the sign and commonly 1 bit is used for S and $p = 52$ bits for d . Since we have $2^{52} \approx 0.5 \times 10^{16}$, we can represent about that many discrete values, which means we can accurately represent about 16 digits (in base 10). The result is that floats on a computer are actually discrete (we have a finite number of bits), and if we get a number that is in one of the gaps (there are uncountably many reals), it's approximated by the nearest discrete value. The accuracy of our representation is to within $1/2$ of the gap between the two discrete values bracketing the true number. Let's consider the implications for accuracy in working with large and small numbers. By changing e we can change the magnitude of a number. So regardless of whether we have a very large or small number, we have about 16 digits of accuracy, since the absolute spacing depends on what value is represented by the least significant digit (the *ulp*, or unit in the last place) in d , i.e., the $p = 52$ nd one, or in terms of base 10, the 16th digit. Let's explore this:

```
> options(digits = 22)
> .1234123412341234
> 1234.1234123412341234 # not accurate to 16 places
> 123412341234.123412341234 # only accurate to 4 places
> 1234123412341234.123412341234 # no places!
> 12341234123412341234 # fewer than no places!
```

We can see the implications of this in the context of calculations:

```

> 1234567812345678 - 1234567812345677
> 12345678123456788888 - 12345678123456788887
> .1234567812345678 - .1234567812345677
> .12345678123456788888 - .12345678123456788887
> .00001234567812345678 - .00001234567812345677 # not as close as
we'd expect, should be 1e-20
> .000012345678123456788888 - .000012345678123456788887

```

Will the following calculation work? How many decimal places do we expect to be accurate?

```
> 123456781234 - .0000123456781234
```

The spacing of possible computer numbers that have a magnitude of about 1 is known as *machine epsilon* (an alternative, but essentially equivalent definition to that given above). This tells us also about the relative spacing of numbers at any point by virtue of numbers being represented as (S, d, e) . For numbers just above 1, if we use (1) with base 10, we have $1 = .100\dots00 \times 10^1$ while the next largest number is $.100\dots01 \times 10^1$, a difference of 10^{1-p} . Generalizing this, and dividing by 2 to represent the error if results are rounded (rather than truncated), we have that the spacing is $\frac{1}{2}2^{1-p} = 2^{-p}$. The difference is $2^{-p} = 2^{-52} \approx 2.2 \times 10^{-16} \equiv \epsilon$, which is what is given as *Machine\$double.eps*. Machine epsilon gives the *absolute spacing* for numbers near 1 and the *relative spacing* for numbers with a different order of magnitude and therefore a different absolute magnitude of the error in representing a real. The relative spacing at x is

$$\frac{(1 + \epsilon)x - x}{x} = \epsilon$$

since the next largest number from x is given by $(1 + \epsilon)x$. Suppose $x = 1 \times 10^6$. Then the absolute error in representing a number of this magnitude is $x\epsilon \approx 2 \times 10^{-10}$. We can see by looking at the numbers in decimal form, where we are accurate to the order 10^{-10} but not 10^{-11} .

```
> 1000000.1
```

3 Implications for calculations and comparisons

As mentioned for integers, computer number arithmetic is not closed, unlike real arithmetic. For example, if we multiply two computer floating points, we can overflow and not get back another computer floating point. One term that is used, which might pop up in an error message (though probably not in R) is that an “exception” is “thrown”. Another mathematical concept we should consider here is that computer arithmetic does not obey the associative and distribute laws, i.e., $(a + b) + c$ may not equal $a + (b + c)$ on a computer and $a(b + c)$ may not be the same as $ab + ac$.

3.1 Comparisons

As we saw, we should never test `a==b` unless a and b are represented as integers in R (they'll actually generally be doubles, but are stored exactly if not too large). Similarly we should never test `a==0`. One nice approach to checking for approximate equality is to make use of *Machine epsilon*. If the relative spacing of two numbers is less than Machine epsilon then for our computer approximation, we say they are the same. Here's an implementation that relies on the absolute error being $x\epsilon$ (see above):

```
> if(abs(a - b) < .Machine$double.eps * abs(a + b))
+   print("approximately equal")
```

Actually, we probably want to use a number slightly larger than `.Machine$double.eps` to be safe. You can also take a look at the R function `all.equal()`.

3.2 Calculations

Given the limited *precision* of computer numbers, we need to be careful when:

- Subtracting large numbers that are nearly equal (or adding negative and positive numbers of the same magnitude). You won't have the precision in the answer that you would like.

```
> 123456781234.56 - 123456781234.00
```

The absolute error here, based on the larger value (which has the fewest error-free decimal places) is of the order $\epsilon x = 2.2 \times 10^{-16} \cdot 1 \times 10^{12} \approx 1 \times 10^{-4} = .0001$, so while we might think that the result is close to the value 1 and should have error of about machine epsilon, we actually only have about four significant digits.

This is called *catastrophic cancellation*, because most of the digits that are left represent rounding error - all the significant digits have cancelled with each other.

It's best to do any subtraction on numbers that are not too large. For example, we can get catastrophic cancellation when computing a sum of squares in a naive way:

$$s^2 = \sum x_i^2 - n\bar{x}^2$$

```
> x = c(-1, 0, 1) + 1e8
> n = length(x)
> sum(x^2) - n*mean(x)^2 # that's not good!
> sum((x - mean(x))^2)
```

A good principle to take away is to subtract off a number similar in magnitude to the values

(in this case \bar{x} is obviously ideal) and adjust your calculation accordingly.

In general, you can sometimes rearrange your calculation to avoid catastrophic cancellation. Another example involves the quadratic formula for finding a root (p. 101 of Gentle).

- Adding or subtracting numbers that are very different in magnitude. The precision will be that of the large magnitude number, since we can only represent that number to a certain absolute accuracy, which is much less than the absolute accuracy of the smaller number:

> 123456781234 - 0.00001

The absolute error in representing the larger number is around 1×10^{-4} and the smaller number is smaller than this.

A work-around is to add a set of numbers in increasing order. However, if the numbers are all of similar magnitude, then by the time you add ones later in the summation, the partial sum will be much larger than the new term. A work-around is to add the numbers in a tree-like fashion, so that each addition involves a summation of numbers of similar size.

Given the limited *range* of computer numbers, be careful when you are:

- Multiplying or dividing many numbers, particularly large or small ones. Never take the product of many large or small numbers as this can cause over- or under-flow. Rather compute on the log scale and only at the end of your computations exponentiate. E.g.,

$$\prod_i x_i / \prod_j y_j = \exp(\sum_i \log x_i - \sum_j \log y_j)$$

- Challenge: Let's think about how we can handle the following calculation. Suppose I want to calculate a predictive density (e.g., in a model comparison in a Bayesian context):

$$\begin{aligned} f(y^*|y, x) &= \int f(y^*|y, x, \theta)\pi(\theta|y, x)d\theta \\ &\approx \frac{1}{M} \sum_{j=1}^m \prod_{i=1}^n f(y_i^*|x, \theta_j) \\ &= \frac{1}{M} \sum_{j=1}^m \exp \sum_{i=1}^n \log f(y_i^*|x, \theta_j) \\ &\equiv \frac{1}{M} \sum_{j=1}^m \exp(v_j) \end{aligned}$$

First, why do I use the log conditional predictive density? Second, let's work with an estimate of the unconditional predictive density on the log scale, $\log f(y^*|y, x) \approx \log \frac{1}{M} \sum_{j=1}^m \exp(v_j)$. Now note that e^{v_j} may be quite small as v_j is the sum of log likelihoods. So what hap-

pens if we have terms something like e^{-1000} ? So we can't exponentiate each individual v_j . Thoughts? I have one solution in mind, but there might be other approaches.

Numerical issues come up frequently in linear algebra. For example, they come up in working with positive definite and semi-positive-definite matrices, such as covariance matrices. You can easily get negative numerical eigenvalues even if all the eigenvalues are positive or non-negative. Here's an example where we use an squared exponential correlation as a function of time (or distance in 1-d), which is *mathematically* positive definite but not numerically positive definite:

```
> xs = 1:100
> dists = rdist(xs)
> corMat = exp(-(dists/10)^2)
> eigen(corMat)$values
```

How the computer actually does arithmetic with the floating point representation in base 2 gets pretty complicated, and we won't go into the details. These rules of thumb should be enough for our practical purposes. Monahan and the URL reference have many of the gory details.

4 Some additional details

It's important to note that operations with integers are fast and exact (but can easily overflow) while operations with floating points are slower and approximate. Because of this slowness, floating point operations (*flops*) dominate calculation intensity and are used as the metric for the amount of work being done - a multiplication (or division) combined with an addition (or subtraction) is one flop.

The actual storage of a number on a computer these days is generally as a double in the form:

$$(-1)^S \times 1.d \times 2^{e-1023}$$

where the leading 1 normalizes the number; i.e., ensures there is a unique representation for a given computer number. This avoids representing any number in multiple ways, e.g., either $1 = 1.0 \times 2^0 = 0.1 \times 2^1 = 0.01 \times 2^2$. The leading 1 gives an extra bit, the hidden bit. e is represented using 11 bits, and the subtraction takes the place of having a sign bit for the exponent.

Let's consider what can be represented exactly since we're working in base 2. First, let's do this on my Linux machine:

```
> .1
> .5
```

So why is 0.5 stored exactly and 0.1 not stored exactly? Consider the difficulty with 1/3 in base 10. Strangely, if we do this on my laptop or one of the department's Mac or Linux machines, we get

```
> 0.1 # 0.1
```

Perhaps the computer is doing something more sophisticated that we're not understanding. However, we see that $1/3$ only prints out 16 3's, so I think it's just ignoring my request for 22 digits and only giving 16.

In R, numbers are generally stored as doubles unless generated based on `seq()` or explicitly coerced:

```
> x = 3; typeof(x)
```

```
> x = as.integer(3); typeof(x)
```

This isn't generally a problem in part because integers will be stored exactly in base two provided the absolute value is less than 2^{53} . Why? [Consider the base 2 representation of an integer in terms of the form $-1^S \times 1.d \times 2^{e-127}$.]

We often encounter a certain integer as a symbol for missing values. This should generally be ok:

```
> x [ x == -9999 ] = NA
```

but only because integers of this magnitude are stored exactly. To be really careful, you can read in as character type and do the assessment before converting to numeric.

Finally, be wary of

```
> x[ x < 0 ] = NaN
```

if what you are looking for is values that might be *mathematically* less than zero, rather than whatever is *numerically* less than zero.

Numerical linear algebra

November 7, 2011

References:

- Gentle: Numerical Linear Algebra for Applications in Statistics (my notes here are based primarily on this source) [Gentle-NLA]
 - Unfortunately, this is not in the UCB library system and the copy I have is borrowed from another faculty member, so I can't lend it out.
- Gentle: Computational Statistics [Gentle-CS]
- Lange: Numerical Analysis for Statisticians
- Monahan: Numerical Methods of Statistics

In working through how to compute something or understanding an algorithm, it can be very helpful to depict the matrices and vectors graphically. We'll see this on the board in class.

1 Preliminaries

1.1 Goals

Here's what I'd like you to get out of this unit:

1. How to think about the computational order of a problem
2. How to choose a computational approach to a given linear algebra calculation you need to do.
3. An understanding of how issues with computer numbers play out in terms of linear algebra.

1.2 Key principle

The form of a mathematical expression and how it should be evaluated on a computer may be very different. Better computational approaches can increase speed and improve the numerical properties of the calculation.

Example 1: We do not compute $(X^T X)^{-1} X^T Y$ by computing $X^T X$ and finding its inverse. In fact, perhaps more surprisingly, we may never actually form $X^T X$ in some implementations.

Example 2: Suppose I have a matrix A , and I want to permute (switch) two rows. I can do this with a permutation matrix, P , which is mostly zeroes. On a computer, in general I wouldn't need to even change the values of A in memory in some cases. Why not?

1.3 Computational complexity

We can assess the computational complexity of a linear algebra calculation by counting the number multiplies/divides and the number of adds/subtracts. Sidenote: addition is a bit faster than multiplication, so some algorithms attempt to trade multiplication for addition.

In general we do not try to count the actual number of calculations, but just their order, though in some cases in this unit we'll actually get a more exact count. In general, we denote this as $O(f(n))$ which means that the number of calculations approaches $cf(n)$ as $n \rightarrow \infty$ (i.e., we know the calculation is approximately proportional to $f(n)$). Consider matrix multiplication, AB , with matrices of size $a \times b$ and $b \times c$. Each column of the second matrix is multiplied by all the rows of the first. For any given inner product of a row by a column, we have b multiplies. We repeat these operations for each column and then for each row, so we have abc multiplies so $O(abc)$ operations. We could count the additions as well, but there's usually an addition for each multiply, so we can usually just count the multiplies and then say there are such and such {multiply and add}s. This is Monahan's approach, but you may see other counting approaches where one counts the multiplies and the adds separately.

For two symmetric, $n \times n$ matrices, this is $O(n^3)$. Similarly, matrix factorization (e.g., the Cholesky decomposition) is $O(n^3)$ unless the matrix has special structure, such as being sparse. As matrices get large, we have speed decrease drastically and memory use increase drastically. In terms of memory use, to hold the result of the multiply indicated above, we need to hold $ab+bc+ac$ total elements, which for symmetric matrices sums to $3n^2$. So for a matrix with $n = 10000$, we have $3 \cdot 10000^2 \cdot 8/1e6 = 2.4\text{Gb}$.

When we have $O(n^q)$ this is known as polynomial time. Much worse is $O(b^n)$ (exponential time), while much better is $O(\log n)$ (log time). Computer scientists talk about NP-complete problems; these are essentially problems for which there is not a polynomial time algorithm - it turns out all such problems can be rewritten such that they are equivalent to one another.

In real calculations, it's possible to have the actual time ordering of two approaches differ from what the order approximations tell us. E.g., something that involves n^2 operations may be faster than one that involves $1000(n \log n + n)$ even though the former is $O(n^2)$ and the latter $O(n \log n)$. The problem is that the constant, $c = 1000$, can matter, as can the extra calculations from the lower order term(s), in this case $1000n$.

A note on terminology: *flops* stands for both floating point operations (the number of operations required) and floating point operations per second, the speed of calculation.

1.4 Notation and dimensions

I'll try to use capital letters for matrices, A , and lower-case for vectors, x . Then x_i is the i th element of x , A_{ij} is the i th row, j th column element, and $A_{.j}$ is the j th column and A_i the i th row. By default, we'll consider a vector to be a one-column matrix, and x^\top to be a one-row matrix. Some of the textbook resources also use a_{ij} for A_{ij} and a_j for the j th column.

Throughout, we'll need to be careful that the matrices involved in an operation are conformable: for $A + B$ both matrices need to be of the same dimension, while for AB the number of columns of A must match the number of rows of B . Note that this allows for B to be a column vector, with only one column, Ab . Just checking dimensions is a good way to catch many errors. Example: is $\text{Cov}(Ax) = A\text{Cov}(x)A^\top$ or $\text{Cov}(Ax) = A^\top\text{Cov}(x)A$. Well, if A is $n \times m$, it must be the former, as the latter is not conformable.

The inner product of two vectors is $\sum_i x_i y_i = x^\top y \equiv \langle x, y \rangle \equiv x \cdot y$. The outer product is xy^\top , which comes from all pairwise products of the elements.

When the indices of summation should be obvious, I'll sometimes leave them implicit. Ask me if it's not clear.

1.5 Norms

$\|x\|_p = (\sum_i |x_i|^p)^{1/p}$ and the standard (Euclidean) norm is $\|x\|_2 = \sqrt{\sum x_i^2} = \sqrt{x^\top x}$, just the length of the vector in Euclidean space, which we'll refer to as $\|x\|$, unless noted otherwise. The standard norm for a matrix is the Frobenius norm, $\|A\|_F = (\sum_{i,j} a_{ij}^2)^{1/2}$. There is also the induced matrix norm, corresponding to any chosen vector norm,

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|}$$

So we have

$$\|A\|_2 = \sup_{\|x\|_2=1} \|Ax\|_2$$

A property of any legitimate matrix norm (including the induced norm) is that $\|AB\| \leq \|A\|\|B\|$. Recall that norms must obey the triangle inequality, $\|A + B\| \leq \|A\| + \|B\|$.

A normalized vector is one with “length”, i.e., Euclidean norm, of one. We can easily normalize a vector: $\tilde{x} = x/\|x\|$

The angle between two vectors is

$$\theta = \cos^{-1} \left(\frac{\langle x, y \rangle}{\sqrt{\langle x, x \rangle \langle y, y \rangle}} \right)$$

1.6 Orthogonality

Two vectors are orthogonal if $x^\top y = 0$, in which case we say $x \perp y$. An orthogonal matrix is a matrix in which all of the columns are orthogonal to each other and normalized. Orthogonal matrices can be shown to have full rank. Furthermore if A is orthogonal, $A^\top A = I$, so $A^{-1} = A^\top$. Given all this, the determinant of orthogonal A is either 1 or -1. Finally the product of two orthogonal matrices, A and B , is also orthogonal since $(AB)^\top AB = B^\top A^\top AB = I$.

Permutations Sometimes we make use of matrices that permute two rows (or two columns) of another matrix when multiplied. Such a matrix is known as an elementary permutation matrix and is an orthogonal matrix with a determinant of -1. You can multiply such matrices to get more general permutation matrices that are also orthogonal. If you premultiply by P , you permute rows, and if you postmultiply by P you permute columns. Note that on a computer, you wouldn’t need to actually do the multiply (and if you did, you should use a sparse matrix routine), but rather one can often just rework index values that indicate where relevant pieces of the matrix are stored (more in the next section).

1.7 Some vector and matrix properties

$AB \neq BA$ but $A + B = B + A$ and $A(BC) = (AB)C$.

In R,

> A + B

> A %*% B

You don’t need the spaces, but they’re nice for code readability.

1.8 Trace and determinant of square matrices

For square matrices, $\text{tr}(A + B) = \text{tr}(A) + \text{tr}(B)$, $\text{tr}(A) = \text{tr}(A^\top)$.

We also have $\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$ - basically you can move a matrix from the beginning to the end or end to beginning. How is this helpful?

1. We can find the ordering that reduces computation the most if the individual matrices are not square.
2. $x^\top Ax = \text{tr}(x^\top Ax)$ since the quadratic form is a scalar, and this is equal to $\text{tr}(xx^\top A)$. It can be helpful to be able to go back and forth between a scalar and a trace in some statistical calculations.

For square matrices, the determinant exists and we have $|AB| = |A||B|$ and therefore, $|A^{-1}| = 1/|A|$ since $|I| = |AA^{-1}| = 1$. Also $|A| = |A^\top|$.

Other matrix multiplications The Hadamard or direct product is simply multiplication of the corresponding elements of two matrices by each other. In R this is simply

```
> A * B
```

Challenge: How can I find $\text{tr}(AB)$ without using `A %*% B`?

The Kronecker product is the product of each element of one matrix with the entire other matrix”

$$A \otimes B = \begin{pmatrix} A_{11}B & \cdots & A_{1m}B \\ \vdots & \ddots & \vdots \\ A_{n1}B & \cdots & A_{nm}B \end{pmatrix}$$

The inverse of a Kronecker product is the Kronecker product of the inverses,

$$B^{-1} \otimes A^{-1}$$

which is obviously quite a bit faster because the inverse (i.e., solving a system of equations) is $O(n^3 + m^3)$ rather than the naive approach being $O((nm)^3)$.

1.9 Linear independence, rank, and basis vectors

A set of vectors, v_1, \dots, v_n , is linearly independent (LIN) when none of the vectors can be represented as a linear combination, $\sum c_i v_i$, of the others. If we have vectors of length n , we can have at most n linearly independent vectors. The rank of a matrix is the number of linearly independent rows (or columns - it's the same), and is at most the minimum of the number of rows and columns. We'll generally think about it in terms of the dimension of the column space - so we can just think about the number of linearly independent columns.

Any set of linearly independent vectors (say v_1, \dots, v_n) span a space made up of all linear combinations of those vectors ($\sum_{i=1}^n c_i v_i$). The spanning vectors are known as basis vectors. We can express a vector x that is in the space with respect to (as a linear combination of) normalized basis vectors using the inner product: $x = \sum_i c_i v_i$ where we can find the weights as $c_k = \langle x, v_k \rangle$.

Consider a regression context. We have p covariates (p columns in the design matrix, X), of which q are linearly independent covariates. This means that $p - q$ of the vectors can be written as linear combos of the q vectors. The space spanned by the covariate vectors is of dimension q , rather than p , and $X^T X$ has $p - q$ eigenvalues that are zero. The q LIN vectors are basis vectors for the space - we can represent any point in the space as a linear combination of the basis vectors. You can think of the basis vectors as being like the axes of the space, except that the basis vectors are not orthogonal. So it's like denoting a point in \mathbb{R}^q as a set of q numbers telling us where on each of the axes we are - this is the same as a linear combination of axis-oriented vectors. When we have $n \leq q$, any observation vector can be represented exactly as a linear combination of the q basis vectors, so there is no residual. If $n = q$, then we have a single unique solution, while if $n < q$ we have multiple possible solutions and the system is ill-determined. Of course we usually have $n > p$, so the system is overdetermined - there is no exact solution, but regression is all about finding solutions that minimize some criterion about the differences between the observations and linear combinations of the columns of the X matrix (such as least squares or penalized least squares).

1.10 Invertibility, singularity, rank, and positive definiteness

For square matrices, let's consider how invertibility, singularity, rank and positive (or non-negative) definiteness relate.

Any square matrix has an eigendecomposition, $A = \Gamma \Lambda \Gamma^{-1}$ where Γ is a matrix with the eigenvectors as the columns and Λ is a diagonal matrix of eigenvalues. The number of non-zero eigenvalues is the same as the rank of the matrix. Square matrices that have an inverse are also called nonsingular, and this is equivalent to having full rank. If the matrix is symmetric, the eigenvectors and eigenvalues are real and Γ is orthogonal, so we have $A = \Gamma \Lambda \Gamma^T$. The determinant of the matrix is the product of the eigenvalues (why?), which is zero if it is less than full rank. Note that if none of the eigenvalues are zero then $A^{-1} = \Gamma \Lambda^{-1} \Gamma^T$.

Let's focus on symmetric matrices. The symmetric matrices that tend to arise in statistics are either positive definite (p.d.) or non-negative definite (n.n.d.). If a matrix is positive definite, then $x^T A x > 0$ for any x . Note that $x^T A x = \text{Cov}(x^T y)$ if $y \sim \mathcal{N}(0, A)$, so positive definiteness amounts to having linear combinations of random variables having positive variance. So we must have that positive definite matrices are equivalent to variance-covariance matrices (I'll just refer to

this as a variance matrix or as a covariance matrix). If A is p.d. then it has all positive eigenvalues and it must have an inverse, though as we'll see, from a numerical perspective, we may not be able to compute it if some of the eigenvalues are very close to zero.

If $x^T Ax \geq 0$ then A is nonnegative definite (also called positive semi-definite). In this case one or more eigenvalues can be zero. What does it mean statistically if a covariance matrix has one or more zero eigenvalues?

Any crossproduct or sum of squares matrix, such as $X^T X$ is positive definite (non-negative definite if $p > n$). This makes sense as it's just a scaling of an empirical covariance matrix.

Generalized inverses These come up in solving equations when A is not full rank. A generalized inverse is a matrix, A^- s.t. $AA^-A = A$. The Moore-Penrose inverse (the pseudo-inverse), A^+ , is a (unique) generalized inverse that also satisfies some additional properties. $x = A^+b$ gives the solution to the linear system, $Ax = b$, that has the shortest length

1.11 Matrices arising in regression

In regression, we work with $X^T X$. Some properties of this matrix are that it is symmetric and non-negative definite (hence our use of $(X^T X)^{-1}$ in the OLS estimator). When is it not positive definite?

Fitted values are $X\hat{\beta} = X(X^T X)^{-1}X^T Y = HY$. The "hat" matrix, H , projects Y into the column space of X . H is idempotent: $HH = H$, which makes sense - once you've projected into the space, any subsequent projection just gives you the same thing back. H is singular. Why? Also, under what special circumstance would it not be singular?

2 Computational issues

2.1 Storing matrices

We've discussed column-major and row-major storage of matrices. First, retrieval of matrix elements from memory is quickest when multiple elements are contiguous in memory. So in a column-major language (e.g., R, Fortran), it is best to work with values in a common column (or entire columns) while in a row-major language (e.g., C) for values in a common row.

In some cases, one can save space (and potentially speed) by overwriting the output into the space occupied by an input. This occurs in some clever implementations of matrix factorizations.

2.2 Algorithms

Good algorithms can change the efficiency of an algorithm by one or more orders of magnitude, and many of the improvements in computational speed over recent decades have been in algorithms rather than in computer speed.

Most matrix algebra calculations can be done in multiple ways. For example, we could compute $b = Ax$ in either of the following ways, denoted here in pseudocode.

1. Stack the inner products of the rows of A with x .
for(i=1:n){
 $b_i = 0$
 for(j=1:m) $b_i = b_i + a_{ij}x_j$ }
2. Take the linear combination (based on x) of the columns of A
for(i=1:n) $b_i = 0$
for(j=1:m)
 for(i=1:n) $b_i = b_i + a_{ij}x_j$

In this case the two approaches involve the same number of operations but the first might be better for row-major matrices and the second for column-major. Challenge - see which is faster in R and figure out what that implies about whether '%*%' calls a Fortran or C routine. (Write the code just doing the outer loop and doing the inner loop using vectorized calculation.)

General computational issues The same caveats we discussed in terms of computer arithmetic hold naturally for linear algebra, since this involves sums of (many) elements. Good implementations of algorithms are aware of the danger of catastrophic cancellation and of the possibility of dividing by zero or by values that are near zero.

2.3 Ill-conditioned problems

Basics A problem is ill-conditioned if small changes to values in the computation result in large changes in the result. This is quantified by something called the *condition number* of a calculation. For different operations there are different condition numbers.

Ill-conditionedness arises most often in terms of matrix inversion, so the standard condition number is the “condition number with respect to inversion”, which when using the L_2 norm is the

ratio of the absolute values of the largest to smallest eigenvalue. Here's an example:

$$A = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix}.$$

The solution of $Ax = b$ for $b = (32, 23, 33, 31)$ is $x = (1, 1, 1, 1)$, while the solution for $b + \delta b = (32.1, 22.9, 33.1, 30.9)$ is $x + \delta x = (9.2, -12.6, 4.5, -1.1)$, where δ is notation for a perturbation to the vector or matrix. What's going on?

Some manipulations with inequalities involving the induced matrix norm (for any chosen vector norm, but we might as well just think about the Euclidean norm) (see Gentle-CS Sec. 5.1) gives

$$\frac{\|\delta x\|}{\|x\|} \leq \|A\| \|A^{-1}\| \frac{\|\delta b\|}{\|b\|}$$

where we define the condition number w.r.t. inversion as $\text{cond}(A) \equiv \|A\| \|A^{-1}\|$. We'll generally work with the L_2 norm, and for a nonsingular square matrix the result is that the condition number is the ratio of the absolute values of the largest and smallest magnitude eigenvalues. This makes sense since $\|A\|_2$ is the absolute value of the largest magnitude eigenvalue of A and $\|A^{-1}\|_2$ that of the inverse of the absolute value of the smallest magnitude eigenvalue of A . We see in the demo code that the large disparity in eigenvalues of A leads to an effect predictable from our inequality above, with the condition number helping us find an upper bound.

The main use of these ideas for our purposes is in thinking about the numerical accuracy of a linear system solution (Gentle-NLA Sec 3.4). On a computer we have the system

$$(A + \delta A)(x + \delta x) = b + \delta b$$

where the 'perturbation' is from the inaccuracy of computer numbers. Our exploration of computer numbers tells us that

$$\frac{\|\delta b\|}{\|b\|} \approx 10^{-p}; \quad \frac{\|\delta A\|}{\|A\|} \approx 10^{-p}$$

where $p = 16$ for standard double precision floating points. Following Gentle, one gets the approximation

$$\frac{\|\delta x\|}{\|x\|} \approx \text{cond}(A) 10^{-p},$$

so if $\text{cond}(A) \approx 10^t$, we have accuracy of order 10^{t-p} instead of 10^{-p} . (Gentle cautions that this holds only if $10^{t-p} \ll 1$). So we can think of the condition number as giving us the number

of digits of accuracy lost relative to the precision of numbers on the computer. E.g., a condition number of 10^8 means we lose 8 digits of accuracy relative to our original 16 on standard systems. One issue is that estimating the condition number is itself subject to numerical error and requires computation of A^{-1} (albeit not in the case of L_2 norm with square, nonsingular A) but see Golub and van Loan (1996; p. 76-78) for an algorithm.

Improving conditioning Ill-conditioned problems in statistics often arise from collinearity of regressors. Often the best solution is not a numerical one, but re-thinking the modeling approach, as this generally indicates statistical issues beyond just the numerical difficulties.

A general comment on improving conditioning is that we want to avoid large differences in the magnitudes of numbers involved in a calculation. In some contexts such as regression, we can center and scale the columns to avoid such differences - this will improve the condition of the problem. E.g., in simple quadratic regression with $x = \{1990, \dots, 2010\}$, Monahan gives an example where the condition number is 5×10^{11} for a naive implementation, 5×10^5 after scaling the columns by their means, and 74 after simply centering each column.

The basic story is that simple strategies often solve the problem, and that you should be cognizant of the absolute and relative magnitudes involved in your calculations.

3 Matrix factorizations (decompositions) and solving systems of linear equations

Suppose we want to solve the following linear system:

$$\begin{aligned} Ax &= b \\ x &= A^{-1}b \end{aligned}$$

Numerically, this is never done by finding the inverse and multiplying. Rather we solve the system using a matrix decomposition (or equivalent set of steps). One approach uses Gaussian elimination (equivalent to the LU decomposition), while another uses the Cholesky decomposition. There are also iterative methods that generate a sequence of approximations to the solution but reduce computation (provided they are stopped before the exact solution is found).

Gentle-CS has a nice table overviewing the various factorizations (Table 5.1, page 219).

3.1 Triangular systems

As a preface, let's figure out how to solve $Ax = b$ if A is upper triangular. The basic algorithm proceeds from the bottom up. We solve for x_n trivially, and then move upwards plugging in the known values of x and solving for the remaining unknown in each row (each equation).

1. $x_n = b_n/A_{nn}$
2. Now for $k < n$, use the already computed $\{x_n, x_{n-1}, \dots, x_{k+1}\}$ to calculate $x_k = \frac{b_k - \sum_{j=k+1}^n b_j A_{kj}}{A_{kk}}$.
3. Repeat for all rows.

How many multiplies and adds are done?

In R, *backsolve()* solves upper triangular systems and *forwardsolve()* solves lower triangular systems:

```
> backsolve(A, b) # A is upper-triangular
> forwardsolve(A, b) # A is lower-triangular
```

We can also solve $(A^T)^{-1}b$ as

```
> backsolve(A, b, transpose = TRUE) # A is upper-triangular
> forwardsolve(A, b, transpose = TRUE) # A is lower-triangular
```

Let $U = A$ be upper-triangular. To reiterate the distinction between matrix inversion and solving a system of equations, when we write $U^{-1}b$, what we mean on a computer is to carry out the above algorithm, not to find the inverse and then multiply.

3.2 Gaussian elimination (LU decomposition)

Gaussian elimination is a standard way of directly computing a solution for $Ax = b$. It is equivalent to the LU decomposition. LU is primarily done with square matrices, but not always. Also LU decompositions do exist for some singular matrices.

The idea of Gaussian elimination is to convert the problem to a triangular system. In class, we'll walk through Gaussian elimination in detail and see how it relates to the LU decomposition. I'll describe it more briefly here. Following what we learned in algebra when we have multiple equations, we preserve the solution, x , when we add multiples of rows (i.e., add multiples of equations) together. This amounts to doing $L_1Ax = L_1b$ for a lower-triangular matrix L_1 that produces all zeroes in the first column of L_1A except for the first row. We proceed to zero out values below the diagonal for the other columns of A . The result is $L_{n-1} \cdots L_1A \equiv U = Bb \equiv b^*$ where U is upper triangular. This is the forward reduction step of Gaussian elimination. Then the backward elimination step solves $Ux = b^*$.

If we're just looking for the solution of the system, we don't need $L = (L_{n-1} \cdots L_1)^{-1}$, but it turns out to have a simple form that is computed as we go along, it is unit lower triangular and the values below the diagonal are the negative of the values below the diagonals in L_1, \dots, L_{n-1} (note that each L_j has non-zeroes below the diagonal only in the j th column). As a side note related to storage, it turns out that as we proceed, we can store the elements of L and U in the original A matrix, except for the implicit 1s on the diagonal of L .

If we look at `solve.default()` in R, we see that it uses `dgesv`. A Google search indicates that this is a Lapack routine that does the LU decomposition with partial pivoting and row interchanges, so R is using the algorithm we've just discussed.

Let's work out the computational complexity. The backward elimination step takes $n^2/2$ operations. What about the forward reduction? It turns out this dominates, and the full LU decomposition requires $n^3/3 + O(n^2)$ flops

Would it ever make sense to explicitly compute the inverse, such as to solve multiple right-hand sides? It turns out if we actually computed the inverse of A from the factorization, this would require $\frac{2}{3}n^3 + O(n^2)$ additional flops. Furthermore, for every right-hand side, we have n^2 flops. This is the same number of operations as the backsolve and forwardsolve involved in $U^{-1}(L^{-1}b)$.

One additional complexity is that we want to avoid dividing by very small values to avoid introducing numerical inaccuracy (we would get large values that might overwhelm whatever they are being added to, and small errors in the divisor will have large effects on the result). This can be done on the fly by interchanging equations to use the equation (row) that produces the largest value to divide by. For example in the first step, we would switch the first equation (first row) for whichever of the remaining equations has the largest value in the first column. This is called partial pivoting. The divisors are called pivots. Complete pivoting also considers interchanging columns, and while theoretically better, partial pivoting is generally sufficient and requires fewer computations. Note that partial pivoting can be expressed as multiplying along the way by permutation matrices, P_1, \dots, P_{n-1} that switch rows. Based on pivoting, we have $PA = LU$, where $P = P_{n-1} \cdots P_1$.

Finally $|PA| = |P||A| = |L||U| = |U|$ (why?) so $|A| = |U|/|P|$ and since the determinant of each permutation matrix, P_j is -1 (except when $P_j = I$ because we don't need to switch rows), we just need to multiply by minus one if there is an odd number of permutations. Or if we know the matrix is non-negative definite, we just take the absolute value of $|U|$. So Gaussian elimination provides a fast stable way to find the determinant.

3.3 Cholesky decomposition

When A is p.d., we can use the Cholesky decomposition to solve a system of equations. Positive definite matrices can be decomposed as $U^T U = A$ where U is upper triangular. U is called a square root matrix and is unique (apart from the sign, which we fix by requiring the diagonals to be positive). One algorithm for computing U is:

1. $U_{11} = \sqrt{A_{11}}$
2. For $j = 2, \dots, n$, $U_{1j} = A_{1j}/U_{11}$
3. For $i = 2, \dots, n$,
 - $U_{ii} = \sqrt{A_{ii} - \sum_{k=1}^{i-1} U_{ki}^2}$
 - for $j = i + 1, \dots, n$: $U_{ij} = (A_{ij} - \sum_{k=1}^{i-1} U_{ki}U_{kj})/U_{ii}$

We can then solve a system of equations as: $U^{-1}(U^T)^{-1}b$, which in R is done as follows:

```
> backsolve(U, backsolve(U, b, transpose = TRUE))
```

This is the same as, but more efficient than,

```
> backsolve(U, forwardsolve(t(U), b))
```

The Cholesky has some nice advantages over the LU: (1) while both are $O(n^3)$, the Cholesky involves only half as many computations, $n^3/6 + O(n^2)$ (2) the Cholesky factorization has only $(n^2 + n)/2$ unique values compared to $n^2 + n$ for the LU. Of course the LU is more broadly applicable. The Cholesky does require computation of square roots, but it turns out this is not too intensive. There is also a method for finding the Cholesky without square roots.

Uses of the Cholesky The standard algorithm for generating $y \sim \mathcal{N}(0, A)$ is:

```
> U = chol(A)
```

```
> y = crossprod(U, rnorm(n)) # i.e., t(U) %*% rnorm(n), but much faster
```

Question: where will most of the time in this two-step calculation be spent?

If a regression design matrix, X , is full rank, then $X^T X$ is positive definite, so we could find $\hat{\beta} = (X^T X)^{-1} X^T Y$ using either the Cholesky or Gaussian elimination. Question: what order would we want to carry out these operations in? However, it turns out that the standard approach is to work with X using the QR decomposition rather than working with $X^T X$; working with X is more numerically stable, though in most situations without extreme collinearity, either of the approaches will be fine.

Numerical issues with eigendecompositions and Cholesky decompositions for positive definite matrices Monahan comments that in general Gaussian elimination and the Cholesky decomposition are very stable. However if the matrix is very ill-conditioned we can get $A_{ii} - \sum_k U_{ki}^2$ being negative and then the algorithm stops when we try to take the square root. In this case, the Cholesky decomposition does not exist numerically although it exists mathematically. It's not all that hard to produce such a matrix, particularly when working with high-dimensional covariance matrices with large correlations - let's see this in the demo code.

We can think about the accuracy here as follows. When the condition number approaches 1×10^{16} , we know that the relative accuracy of results deteriorates such that the error approaches the actual values in magnitude. So it shouldn't be too surprising that we end up with matrices that are not numerically positive definite. These matrices are close to not being positive definite; we can see this in their having eigenvalues near zero. As the conditioning gets worse, we eventually get the situation where we have $A_{ii} - \sum_k U_{ki}^2$ being numerically negative, though mathematically it is positive (but very close to zero). Note that when the Cholesky fails, we can still compute an eigendecomposition, but we have negative numeric eigenvalues. Even if all the eigenvalues are numerically positive (or equivalently, we're able to get the Cholesky), errors in small eigenvalues near machine precision could have large effects when we work with the inverse of the matrix. This is what happens when we have columns of the X matrix nearly collinear. We cannot statistically distinguish the effect of two (or more) covariates, and this plays out numerically in terms of unstable results.

A strategy when working with mathematically but not numerically positive definite A is to set eigenvalues or singular values to zero when they get very small, which amounts to using a pseudo-inverse and setting to zero any linear combinations with very small variance. We can also use pivoting with the Cholesky and accumulate zeroes in the last $n - q$ rows (for cases where we try to take the square root of a negative number), corresponding to the columns of A that are numerically linearly dependent.

3.4 QR decomposition

3.4.1 Introduction

The QR decomposition is available for any matrix, $A = QR$, with Q orthogonal and R upper triangular. If A is non-square, $n \times p$ with $n > p$ then the leading p rows of R provide an upper triangular matrix (R_1) and the remaining rows are 0. (I'm using p because the QR is generally applied to design matrices in regression). In this case we really only need the first p columns of Q , and we have $A = Q_1 R_1$, the 'skinny' QR (this is what R's QR provides). For uniqueness, we can require the diagonals of R to be nonnegative, and then R will be the same as the upper-triangular

Cholesky factor of $A^T A$:

$$\begin{aligned} A^T A &= R^T Q^T Q R \\ &= R^T R \end{aligned}$$

There are three standard approaches for computing the QR, using (1) reflections (Householder transformations), rotations (Givens transformations), or Gram-Schmidt orthogonalization (see below).

For $n \times n$ A , the QR (for the Householder approach) requires $2n^3/3$ flops, so QR is less efficient than LU or Cholesky.

We can also obtain the pseudo-inverse of A from the QR: $A^+ = [R_1^{-1} \ 0] Q^T$. In the case that A is not full-rank, there is a version of the QR that will work (involving pivoting) and we end up with some additional zeroes on the diagonal of R_1 .

3.4.2 Regression and the QR

Often QR is used to fit linear models, including in R. Consider the linear model in the form $Y = A\beta + \epsilon$, finding $\hat{\beta} = (A^T A)^{-1} A^T Y$. Let's consider the skinny QR and note that R^T is invertible. Therefore, we can express the normal equations as

$$\begin{aligned} A^T A \beta &= A^T Y \\ R^T Q^T Q R \beta &= R^T Q^T Y \\ R \beta &= Q^T Y \end{aligned}$$

and solving for β is just a backsolve since R is upper-triangular. Furthermore the standard regression quantities, such as the hat matrix, the SSE, the residuals, etc. can be easily expressed in terms of Q and R .

Why use the QR instead of the Cholesky on $A^T A$? The condition number of A is the square root of that of $A^T A$, and the QR factorizes A . Monahan has a discussion of the condition of the regression problem, but from a larger perspective, the situations where numerical accuracy is a concern are generally cases where the OLS estimators are not particularly helpful anyway (e.g., highly collinear predictors).

What about computational order of the different approaches to least squares? The Cholesky is $np^2 + \frac{1}{3}p^3$, an algorithm called sweeping is $np^2 + p^3$, the Householder method for QR is $2np^2 - \frac{2}{3}p^3$, and the modified Gram-Schmidt approach for QR is $2np^2$. So if $n \gg p$ then Cholesky (and sweeping) are faster than the QR approaches. According to Monahan, modified Gram-Schmidt is most numerically stable and sweeping least. In general, regression is pretty quick unless p is large

since it is linear in n , so it may not be worth worrying too much about computational differences of the sort noted here.

3.4.3 Regression and the QR in R

Regression in R uses the QR decomposition via $qr()$, which calls a Fortran function. $qr()$ (and the Fortran functions that are called) is specifically designed to output quantities useful in fitting linear models. Note that by default you get the skinny QR, namely only the first p rows of R and the first p columns of Q , where the latter form an orthonormal basis for the column space of A . The remaining columns form an orthonormal basis for the null space of A . The analogy in regression is that we get the basis vectors for the regression, while adding the remaining columns gives us the full n -dimensional space of the observations.

$qr()$ returns the result as a list meant for use by other tools. R stores the R matrix in the upper triangle of $\$qr$, while the lower triangle of $\$qr$ and $\$aux$ store the information for constructing Q (this relates to the Householder-related vectors u below). One can multiply by Q using $qr.qy()$ and by Q^\top using $qr.qty()$. If you want to extract R and Q , the following will work:

```
> A.qr = qr(A)
> Q = qr.Q(A.qr)
> R = qr.R(A.qr)
```

As a side note, there are QR-based functions that provide regression-related quantities, such as $qr.resid()$, $qr.fitted()$ and $qr.coef()$. These functions (and their Fortran counterparts) exist because one can work through the various regression quantities of interest and find their expressions in terms of Q and R , with nice properties resulting from Q being orthogonal and R triangular.

3.4.4 Computing the QR decomposition

We'll work through some of the details of the different approaches to the QR, in part because they involve some concepts that may be useful in other contexts.

One approach involves reflections of vectors and a second rotations of vectors. Reflections and rotations are transformations that are performed by orthogonal matrices. The determinant of a reflection matrix is -1 and the determinant of a rotation matrix is 1.

Reflections If u and v are orthonormal vectors and x is in the space spanned by u and v , $x = c_1u + c_2v$, then $\tilde{x} = -c_1u + c_2v$ is a reflection (a *Householder* reflection) along the u dimension (since we are using the negative of that basis vector). We can think of this as reflecting across the plane perpendicular to u . This extends simply to higher dimensions with orthonormal vectors, u, v_1, v_2, \dots

Suppose we want to formulate the reflection in terms of a “Householder” matrix, Q . It turns out that

$$Qx = \tilde{x}$$

if $Q = I - 2uu^\top$. Q has the following properties: (1) $Qu = -u$, (2) $Qv = v$ for $u^\top v = 0$, (3) Q is orthogonal and symmetric.

One way to create the QR decomposition is by a series of Householder transformations that create an upper triangular R from A :

$$\begin{aligned} R &= Q_n \cdots Q_1 A \\ Q &= (Q_n \cdots Q_1)^\top \end{aligned}$$

where we make use of the symmetry in defining Q .

Basically Q_1 reflects the first column of A with respect to a carefully chosen u , so that the result is all zeroes except for the first element. We want $Q_1 x = \tilde{x} = (\|x\|, 0, \dots, 0)$. This can be achieved with $u = \frac{x - \tilde{x}}{\|x - \tilde{x}\|}$. Then Q_2 makes the last $n - 2$ rows of the second column equal to zero. We'll work through this a bit in class.

In the regression context, as we work through the individual transformations, $Q_j = I - 2u_j u_j^\top$, we apply them to A and Y to create R (note this would not involve doing the full matrix multiplication - think about what calculations are actually needed) and $QY = Q^\top Y$, and then solve $R\beta = Q^\top Y$. To find $\text{Cov}(\hat{\beta}) = (A^\top A)^{-1} = (R^\top R)^{-1} = R^{-1} R^{-\top}$ we do need to invert R , but it's upper-triangular and of dimension $p \times p$. It turns out that $Q^\top Y$ can be partitioned into the first p and the last $n - p$ elements, $z^{(1)}$ and $z^{(2)}$. The SSR is $\|z^{(1)}\|^2$ and SSE is $\|z^{(2)}\|^2$.

Rotations A Givens rotation matrix rotates a vector in a two-dimensional subspace to be axis oriented with respect to one of the two dimensions by changing the value of the other dimension. E.g. we can create $\tilde{x} = (x_1, \dots, \tilde{x}_p, \dots, 0, \dots, x_n)$ from $x = (x_1, \dots, x_p, \dots, x_q, \dots, x_n)$ using a matrix multiplication: $\tilde{x} = Qx$. Q is orthogonal but not symmetric.

We can use a series of Givens rotations to do the QR but unless it is done carefully, more computations are needed than with Householder reflections. The basic story is that we apply a series of Givens rotations to A such that we zero out the lower triangular elements.

$$\begin{aligned} R &= Q_{pn} \cdots Q_{22} Q_{1n} \cdots Q_{13} Q_{12} A \\ Q &= (Q_{pn} \cdots Q_{11})^\top \end{aligned}$$

Note that we create the $n - p$ zero rows in R (because the calculations affect the upper triangle of R), but we can then ignore those rows and the corresponding columns of Q .

Gram-Schmidt Orthogonalization Gram-Schmidt involves finding a set of orthonormal vectors to span the same space as a set of LIN vectors, x_1, \dots, x_p . If we take the LIN vectors to be the columns of A , so that we are discussing the column space of A , then G-S yields the QR decomposition. Here's the algorithm:

1. $\tilde{x}_1 = \frac{x_1}{\|x_1\|}$ (normalize the first vector)
2. Orthogonalize the remaining vectors with respect to \tilde{x}_1 :
 - (a) $\tilde{x}_2 = \frac{x_2 - \tilde{x}_1^\top x_2 \tilde{x}_1}{\|x_2 - \tilde{x}_1^\top x_2 \tilde{x}_1\|}$, which orthogonalizes with respect to \tilde{x}_1 and normalizes. Note that $\tilde{x}_1^\top x_2 \tilde{x}_1 = \langle \tilde{x}_1, x_2 \rangle \tilde{x}_1$. So we are finding a scaling, $c\tilde{x}_1$, where c is based on the inner product, to remove the variation in the x_1 direction from x_2 .
 - (b) For $k > 2$, find interim vectors, $x_k^{(2)}$, by orthogonalizing with respect to \tilde{x}_1
3. Proceed for $k = 3, \dots$, in turn orthogonalizing and normalizing the first of the remaining vectors w.r.t. \tilde{x}_{k-1} and orthogonalizing the remaining vectors w.r.t. \tilde{x}_{k-1} to get new interim vectors

Mathematically, we could instead orthogonalize x_2 w.r.t. \tilde{x}_1 , then orthogonalize x_3 w.r.t. $\{\tilde{x}_1, \tilde{x}_2\}$, etc. The algorithm above is the *modified* G-S, and is known to be more numerically stable if the columns of X are close to collinear, giving vectors that are closer to orthogonal. The resulting \tilde{x} vectors are the columns of Q . The elements of R are obtained as we proceed: the diagonal values are the the normalization values in the denominators, while the off-diagonals are the inner products with the already-computed columns of Q that are computed as part of the numerators.

Another way to think about this is that $R = Q^\top A$, which is the same as regressing the columns of A on Q , $(Q^\top Q)^{-1} Q^\top A = Q^\top A$. By construction, the first column of A is a scaling of the first column of Q , the second column of A is a linear combination of the first two columns of Q , etc., so R being upper triangular makes sense.

3.5 Determinants

The absolute value of the determinant of a square matrix can be found from the product of the diagonals of the triangular matrix in any factorization that gives a triangular (including diagonal) matrix times an orthogonal matrix (or matrices) since the determinant of an orthogonal matrix is either one or minus one.

$$|A| = |QR| = |Q||R| = \pm|R|$$

$$|A^\top A| = |(QR)^\top QR| = |R^\top R| = |R_1^\top R_1| = |R_1|^2$$

In R, the following will do it (on the log scale), since R is stored in the upper triangle of the qr element.

```

> myqr = qr(A)
> magn = sum(log(abs(diag(myqr$qr))))
> sign = prod(sign(diag(myqr$qr)))

```

An alternative is the product of the diagonal elements of D (the singular values) in the SVD factorization, $A = UDV^T$.

For non-negative definite matrices, we know the determinant is non-negative, so the uncertainty about the sign is not an issue. For positive definite matrices, a good approach is to use the product of the diagonal elements of the Cholesky decomposition.

One can also use the product of the eigenvalues: $|A| = |\Gamma\Lambda\Gamma^{-1}| = |\Gamma||\Gamma^{-1}||\Lambda| = |\Lambda|$

Computation Computing from any of these diagonal or triangular matrices as the product of the diagonals is prone to overflow and underflow, so we always work on the log scale as the sum of the log of the values. When some of these may be negative, we can always keep track of the number of negative values and take the log of the absolute values.

Often we will have the factorization as a result of other parts of the computation, so we get the determinant for free.

R's `determinant()` uses the LU decomposition. Supposedly `det()` just wraps determinant, but I can't seem to pass the *logarithm* argument into `det()`, so `determinant()` seems more useful.

4 Eigendecomposition and SVD

4.1 Eigendecomposition

The eigendecomposition (spectral decomposition) is useful in considering convergence of algorithms and of course for statistical decompositions such as PCA. We think of decomposing the components of variation into orthogonal patterns (the eigenvectors) with variances (eigenvalues) associated with each pattern.

Square symmetric matrices have real eigenvectors and eigenvalues, with the factorization into orthogonal Γ and diagonal Λ , $A = \Gamma\Lambda\Gamma^T$, where the eigenvalues on the diagonal of Λ are ordered in decreasing value. Of course this is equivalent to the definition of an eigenvalue/eigenvector pair as a pair such that $Ax = \lambda x$ where x is the eigenvector and λ is a scalar, the eigenvalue. The inverse of the eigendecomposition is simply $\Gamma\Lambda^{-1}\Gamma^T$. On a similar note, we can create a square root matrix based on the eigendecomposition. What would it be?

In R, `eigen(A)$vectors` is Γ , with each column a vector, and `eigen(A)$values` is the ordered eigenvalues.

Positive definite matrices have positive eigenvalues, while non-negative definite (positive semi-definite) matrices have non-negative values, with each zero eigenvalue corresponding to a rank deficiency. Statistically, we can think of this as constraints induced by the covariance structure - i.e., linear combinations that are fixed and therefore known with certainty. The pseudo-inverse is $\Gamma\Lambda^+\Gamma^\top$ where $\lambda_i^+ = 1/\lambda_i$ except when $\lambda_i = 0$ in which case $\lambda_i^+ = 0$. If we think of A as a precision matrix with zero precision (infinite variance) for certain linear combinations, then A^+ is the covariance matrix that forces those combinations to have zero variance.

The spectral radius of A , denoted $\rho(A)$, is the maximum of the absolute values of the eigenvalues. As we saw when talking about ill-conditionedness, for symmetric matrices, this maximum is the induced norm, so we have $\rho(A) = \|A\|_2$. It turns out that $\rho(A) \leq \|A\|$ for any induced matrix norm. The spectral radius comes up in determining the rate of convergence of some iterative algorithms.

Computation There are several methods for eigenvalues; a common one for doing the full eigen-decomposition is the *QR algorithm*. The first step is to reduce A to upper Hessenberg form, which is an upper triangular matrix except that the first subdiagonal in the lower triangular part can be non-zero. For symmetric matrices, the result is actually tridiagonal. We can do the reduction using Householder reflections or Givens rotations. At this point the QR decomposition (using Givens rotations) is applied iteratively (to a version of the matrix in which the diagonals are shifted), and the result converges to a diagonal matrix, which provides the eigenvalues. It's more work to get the eigenvectors, but they are obtained as a product of Householder matrices (required for the initial reduction) multiplied by the product of the Q matrices from the successive QR decompositions.

We won't go into the algorithm in detail, but note that it involves manipulations and ideas we've seen already.

If only the largest (or the first few largest) eigenvalues and their eigenvectors are needed, which can come up in time series and Markov chain contexts, the problem is easier and can be solved by the *power method*. E.g., in a Markov chain context, steady state is reached through $x_t = A^t x_0$. One can find the largest eigenvector by multiplying by A many times, normalizing at each step. $v^{(k)} = Az^{(k-1)}$ and $z^{(k)} = v^{(k)} / \|v^{(k)}\|$. There is an extension to find the p largest eigenvalues and their vectors.

4.2 Singular value decomposition

Let's consider an $n \times m$ matrix, A , with $n \geq m$ (if $m > n$, we can always work with A^\top). This often is a matrix representing m features of n observations. We could have n documents and m words, or n gene expression levels and m experimental conditions, etc. A can always be

decomposed as

$$A = UDV^T$$

where U and V are matrices with orthonormal columns (left and right eigenvectors) and D is diagonal with non-negative values (which correspond to eigenvalues in the case of square A and to squared eigenvalues of $A^T A$).

The SVD can be represented in more than one way. One representation is

$$A_{n \times m} = U_{n \times k} D_{k \times k} V_{k \times m}^T = \sum_{j=1}^k D_{jj} u_j v_j^T$$

where u_j and v_j are the columns of U and V and where k is the rank of A (which is at most the minimum of n and m of course). The diagonal elements of D are the singular values.

If A is positive semi-definite, the eigendecomposition is an SVD. Furthermore, $A^T A = V D^2 V^T$ and $AA^T = U D^2 U^T$, so we can find the eigendecomposition of such matrices using the SVD of A . Note that the squares of the singular values of A are the eigenvalues of $A^T A$ and AA^T .

We can also fill out the matrices to get

$$A = U_{n \times n} D_{n \times m} V_{m \times m}^T$$

where the added rows and columns of D are zero with the upper left block the $D_{k \times k}$ from above.

Uses The SVD is an excellent way to determine a matrix rank and to construct a pseudo-inverse ($A^+ = V D^+ U^T$).

We can use the SVD to approximate A by taking $A \approx \tilde{A} = \sum_{j=1}^p D_{jj} u_j v_j^T$ for $p < m$. This approximation holds in terms of the Frobenius norm for $A - \tilde{A}$. As an example if we have a large image of dimension $n \times m$, we could hold a compressed version by a rank- p approximation using the SVD. The SVD is used a lot in clustering problems. For example, the Netflix prize was won based on a variant of SVD (in fact all of the top methods used variants on SVD, I believe).

Computation The basic algorithm is similar to the QR method for the eigendecomposition. We use a series of Householder transformations on the left and right to reduce A to a bidiagonal matrix, $A^{(0)}$. The post-multiplications generate the zeros in the upper triangle. (A bidiagonal matrix is one with non-zeroes only on the diagonal and first subdiagonal above the diagonal). Then the algorithm produces a series of upper bidiagonal matrices, $A^{(0)}$, $A^{(1)}$, etc. that converge to a diagonal matrix, D . Each step is carried out by a sequence of Givens transformations:

$$A^{(j+1)} = R_{m-2}^T R_{m-3}^T \cdots R_0^T A^{(j)} T_0 T_1 \cdots T_{m-2}$$

$$= RA^{(j)}T$$

Note that the multiplication on the right is required to produce zeroes in the k th column and the k th row. This eventually gives $A^{(\dots)} = D$ and by construction, U (the product of the pre-multiplied Householder matrices and the R matrices) and V (the product of the post-multiplied Householder matrices and the T matrices) are orthogonal. The result is then transformed by a diagonal matrix to make the elements of D non-negative and by permutation matrices to order the elements of D in nonincreasing order.

5 Computation

5.1 Linear algebra in R

Speedups and storage savings can be obtained by working with matrices stored in special formats when the matrices have special structure. E.g., we might store a symmetric matrix as a full matrix but only use the upper or lower triangle. Banded matrices and block diagonal matrices are other common formats. Banded matrices are all zero except for $A_{i,i+c_k}$ for some small number of integers, c_k . Viewed as an image, these have bands. The bands are known as co-diagonals.

Note that for many matrix decompositions, you can change whether all of the aspects of the decomposition are returned, or just some, which may speed calculations.

Some useful packages in R for matrices are *Matrix*, *spam*, and *bdsmatrix*. *Matrix* can represent a variety of rectangular matrices, including triangular, orthogonal, diagonal, etc. and provides methods for various matrix calculations that are specific to the matrix type. *spam* handles general sparse matrices with fast matrix calculations, in particular a fast Cholesky decomposition. *bdsmatrix* focuses on block-diagonal matrices, which arise frequently in contexts where there is clustering that induces within-cluster correlation and cross-cluster independence.

In general, matrix operations in R go to compiled C or Fortran code very soon, so they can actually be pretty efficient and are based on the best algorithms developed by numerical experts. The core libraries that are used are LAPACK and BLAS (the Linear Algebra PACKage and the Basic Linear Algebra Subroutines). One way to speed up code that relies heavily on linear algebra is to make sure you have a BLAS library tuned to your machine. These include GotoBLAS (free), Intel's MKL, and AMD's ACML. These can be installed and R can be linked to the shared object library file (.so file) for the fast BLAS. These BLAS libraries are also available in threaded versions that farm out the calculations across multiple cores or processors that share memory. Let's see a demo on the Mac that illustrates that this is enabled automatically.

BLAS routines do vector operations (level 1), matrix-vector operations (level 2), and dense

matrix-matrix operations (level 3). Often the name of the routine has as its first letter “d”, “s”, “c” to indicate the routine is double precision, single precision, or complex. LAPACK builds on BLAS to implement standard linear algebra routines such as eigendecomposition, solutions of linear systems, a variety of factorizations, etc.

5.2 Sparse matrices

As an example of exploiting sparsity, here’s how the *spam* package in R stores a sparse matrix. Consider the matrix to be row-major and store the non-zero elements in order in a vector called *value*. Then create a vector called *rowptr* that stores the position of the first element of each row. Finally, have a vector, *colindices* that tells the column identity of each element.

We can do a fast matrix multiply, Ab , as follows in pseudo-code:

```
for(i in 1:nrows(A))
  x[i] = 0
  for(j in (rowptr[i]:rowptr[i+1]-1)) x[i]=x[i]+value[j]*b[colindices[j]]
```

How many computations have we done? Only k multiplies and $O(k)$ additions where k is the number of non-zero elements of A . Compare this to the usual $O(n^2)$ for dense multiplication.

Note that for the Cholesky of a sparse matrix, if the sparsity pattern is fixed, but the entries change, one can precompute an optimal re-ordering that retains as much sparsity in U as possible. Then multiple Cholesky decompositions can be done more quickly as the entries change.

Banded matrices Suppose we have a banded matrix A where the lower bandwidth is p , namely $A_{ij} = 0$ for $i > j + p$ and the upper bandwidth is q ($A_{ij} = 0$ for $j > i + q$). An alternative to reducing to $Ux = b^*$ is to compute $A = LU$ and then do two solutions, $U^{-1}(L^{-1}b)$. One can show that the computational complexity of the LU factorization is $O(npq)$ for banded matrices, while solving the two triangular systems is $O(np + nq)$, so for small p and q , the speedup can be dramatic.

Banded matrices come up in time series analysis. E.g., MA models produce banded covariance structures because the covariance is zero after a certain number of lags.

5.3 Low rank updates

A transformation of the form $A - uv^T$ is a rank-one update because uv^T is of rank one.

More generally a low rank update of A is $\tilde{A} = A - UV^T$ where U and V are $n \times m$ with $n \geq m$. The Sherman-Morrison-Woodbury formula tells us that

$$\tilde{A}^{-1} = A^{-1} + A^{-1}U(I_m - V^T A^{-1}U)^{-1}V^T A^{-1}$$

so if we know $x_0 = A^{-1}b$, then the solution to $\tilde{A}x = b$ is $x + A^{-1}U(I_m - V^T A^{-1}U)^{-1}V^T x$. Provided m is not too large, and particularly if we already have a factorization of A , then $A^{-1}U$ is not too bad computationally, and $I_m - V^T A^{-1}U$ is $m \times m$. As a result $A^{-1}(U(\dots)^{-1}V^T x)$ isn't too bad.

This also comes up in working with precision matrices in Bayesian problems where we may have A^{-1} but not A (we often add precision matrices to find conditional normal distributions). An alternative expression for the formula is $\tilde{A} = A + UCV^T$, and the identity tells us

$$\tilde{A}^{-1} = A^{-1} - A^{-1}U(C^{-1} + V^T A^{-1}U)^{-1}V^T A^{-1}$$

Basically Sherman-Morrison-Woodbury gives us matrix identities that we can use in combination with our knowledge of smart ways of solving systems of equations.

6 Iterative solutions of linear systems

Gauss-Seidel Suppose we want to iteratively solve $Ax = b$. Here's the algorithm, which sequentially updates each element of x in turn.

- Start with an initial approximation, $x^{(0)}$.
- Hold all but $x_1^{(0)}$ constant and solve to find $x_1^{(1)} = \frac{1}{a_{11}}(b_1 - \sum_{j=2}^n a_{1j}x_j^{(0)})$.
- Repeat for the other rows of A (i.e., the other elements of x), finding $x^{(1)}$.
- Now iterate to get $x^{(2)}$, $x^{(3)}$, etc. until a convergence criterion is achieved, such as $\|x^{(k)} - x^{(k-1)}\| \leq \epsilon$ or $\|r^{(k)} - r^{(k-1)}\| \leq \epsilon$ for $r^{(k)} = b - Ax^{(k)}$.

Let's consider how many operations are involved in a single update: $O(n)$ for each element, so $O(n^2)$ for each update. Thus if we can stop well before n iterations, we've saved computation relative to exact methods.

If we decompose $A = L + D + U$ where L is strictly lower triangular, U is strictly upper triangular, then Gauss-Seidel is equivalent to solving

$$(L + D)x^{(k+1)} = b - Ux^{(k)}$$

and we know that solving the lower triangular system is $O(n^2)$.

It turns out that the rate of convergence depends on the spectral radius of $(L + D)^{-1}U$.

Gauss-Seidel amounts to optimizing by moving in axis-oriented directions, so it can be slow in some cases.

Conjugate gradient For positive definite A , conjugate gradient (CG) reexpresses the solution to $Ax = b$ as an optimization problem, minimizing

$$f(x) = \frac{1}{2}x^\top Ax - x^\top b,$$

since the derivative of $f(x)$ is $Ax - b$ and at the minimum this gives $Ax - b = 0$.

Instead of finding the minimum by following the gradient at each step (so-called steepest descent, which can give slow convergence - we'll see a demonstration of this in the optimization unit), CG chooses directions that are mutually conjugate w.r.t. A , $d_i^\top Ad_j = 0$ for $i \neq j$. The method successively chooses vectors giving the direction, d_k , in which to move down towards the minimum and a scaling of how much to move, α_k . If we start at $x_{(0)}$, the k th point we move to is $x_{(k)} = x_{(k-1)} + \alpha_k d_k$ so we have

$$x_{(k)} = x_{(0)} + \sum_{j \leq k} \alpha_j d_j$$

and we use a convergence criterion such as given above for Gauss-Seidel. The directions are chosen to be the residuals, $b - Ax_{(k)}$. Here's the basic algorithm:

- Choose $x_{(0)}$ and define the residual, $r_{(0)} = b - Ax_{(0)}$ (the error on the scale of b) and the direction, $d_0 = r_{(0)}$ and set $k = 0$.
- Then iterate
 - $\alpha_k = \frac{r_{(k)}^\top r_{(k)}}{d_k^\top A d_k}$ (choose step size so next error will be orthogonal to current direction - which we can express in terms of the residual, which is easily computable)
 - $x_{(k+1)} = x_{(k)} + \alpha_k d_k$ (update current value)
 - $r_{(k+1)} = r_{(k)} - \alpha_k A d_k$ (update current residual)
 - $d_{k+1} = r_{(k+1)} + \frac{r_{(k+1)}^\top r_{(k+1)}}{r_{(k)}^\top r_{(k)}} d_k$ (choose next direction by conjugate Gram-Schmidt, starting with $r_{(k+1)}$ and removing components that are not A -orthogonal to previous directions, but it turns out that $r_{(k+1)}$ is already A -orthogonal to all but d_k).
- Stop when $\|r^{(k+1)}\|$ is sufficiently small.

The convergence of the algorithm depends in a complicated way on the eigenvalues, but in general convergence is faster when the condition number is smaller (the eigenvalues are not too spread out). CG will in principle give the exact answer in n steps (where A is $n \times n$). However, computationally we lose accuracy and interest in the algorithm is really as an iterative approximation where we stop

before n steps. The approach basically amounts to moving in axis-oriented directions in a space stretched by A .

In general, CG is used for large sparse systems.

See the [extensive description from Shewchuk](#) for more details and for the figures shown in class, as well as the use of CG when A is not positive definite.

Updating a solution Sometimes we have solved a system, $Ax = b$ and then need to solve $Ax = c$. If we have solved the initial system using a factorization, we can reuse that factorization and solve the new system in $O(n^2)$. Iterative approaches can do a nice job if $c = b + \delta b$. Start with the solution x for $Ax = b$ as $x^{(0)}$ and use one of the methods above.

Simulation

November 14, 2011

References:

- Gentle: Computational Statistics
- Monahan: Numerical Methods of Statistics

Many (most?) statistical papers include a simulation (i.e., Monte Carlo) study. The basic idea is that closed form analysis of the properties of a statistical method/model is often hard to do. Even if possible, it usually involves approximations or simplifications. A canonical situation is that we have an asymptotic result and we want to know what happens in finite samples, but often we do not even have the asymptotic result. Instead, we can estimate mathematical expressions using random numbers. So we design a simulation study to evaluate the method/model or compare multiple methods. The result is that the statistician carries out an experiment, generally varying different factors to see what has an effect on the outcome of interest.

The basic strategy generally involves simulating data and then using the method(s) on the simulated data, summarizing the results in a way that accounts for the randomness.

Most simulation studies aim to approximate an integral, generally an expected value (mean, bias, variance, MSE, probability, etc.). In low dimensions, methods such as Gaussian quadrature are best for estimating an integral but these methods don't scale well, so in higher dimensions we often use Monte Carlo techniques.

1 Monte Carlo considerations

1.1 Monte Carlo basics

The basic idea is that we often want to estimate $\mu \equiv E_f(h(X))$ for $X \sim f$. Note that if h is an indicator function, this includes estimation of probabilities, e.g., $p = P(X \leq x) = F(x) =$

$\int_{-\infty}^x f(t)dt = \int I(t \leq x)f(t)dt = E_f(I(X \leq x))$. We would estimate variances or MSEs by having h involve squared terms.

We get an MC estimate of μ based on an iid sample of a large number of values of X from f :

$$\hat{\mu} = \frac{1}{m} \sum_{i=1}^m h(X_i),$$

which is justified by the Law of Large Numbers. The simulation variance of $\hat{\mu}$ is $\text{Var}(\hat{\mu}) = \sigma^2/m$, with $\sigma^2 = \text{Var}(h(X))$. An estimator of $\sigma^2 = E_f((h(X) - \mu)^2)$ is $\hat{\sigma}^2 = \frac{1}{m-1} \sum (h(X_i) - \hat{\mu})^2$. So our MC simulation error is based on

$$\widehat{\text{Var}}(\hat{\mu}) = \frac{1}{m(m-1)} \sum_{i=1}^m (h(X_i) - \hat{\mu})^2.$$

Note

Sometimes the X_i are generated in a dependent fashion (e.g., sequential MC or MCMC), in which case this variance estimator does not hold, but the estimator $\hat{\mu}$ is still correct. [As a sidenote, a common misperception with MCMC is that you should thin your chains because of dependence of the samples. This is not correct - the only reason to thin a chain is if you want to save on computer storage or processing.]

Note that in this case the randomness in the system is very well-defined (as it is in survey sampling, but unlike in most other applications of statistics), because it comes from the RNG that we perform as part of our attempt to estimate μ .

Happily, we are in control of m , so in principle we can reduce the simulation error to as little as we desire. Unhappily, as usual, the standard error goes down with the square root of m .

1.2 Simple example

Suppose we've come up with a fabulous new estimator for the mean of a distribution. The estimator is to take the middle value of the sorted observations as our estimate of the mean of the entire distribution. We work out some theory to show that this estimator is robust to outlying observations and we come up with a snazzy new name for our estimator. We call it the 'median'. Let's denote it as \tilde{X} .

Unfortunately, we have no good way of estimating $\text{Var}(\tilde{X}) = E((\tilde{X} - E(\tilde{X}))^2)$ analytically. We decide to use a Monte Carlo estimate

$$\frac{1}{m} \sum_{i=1}^m (\tilde{X}_i - \tilde{X})^2$$

where $\tilde{X} = \frac{1}{m} \sum \tilde{X}_i$. Each \tilde{X}_i in this case is generated by generating a dataset and calculating the median. In evaluating the variance of the median and comparing it to our standard estimator, the sample mean, what decisions do we have to make in our Monte Carlo procedure?

1.3 Variance reduction

There are some tools for variance reduction in MC settings. One is importance sampling (see Section 3). Others are the use of control variates and antithetic sampling. I haven't personally run across these latter in practice, so I'm not sure how widely used they are and won't go into them here.

In some cases we can set up natural strata, for which we know the probability of being in each stratum. Then we would estimate μ for each stratum and combine the estimates based on the probabilities. The intuition is that we remove the variability in sampling amongst the strata from our simulation.

Another strategy that comes up in MCMC contexts is *Rao-Blackwellization*. Suppose we want to know $E(h(X))$ where $X = \{X_1, X_2\}$. Iterated expectation tells us that $E(h(X)) = E(E(h(X)|X_2))$. If we can compute $E(h(X)|X_2) = \int h(x_1, x_2)f(x_1, x_2|x_2)dx_1$ then we should avoid introducing stochasticity related to the X_1 draw (since we can analytically integrate over that) and only average over stochasticity from the X_2 draw by estimating $E_{X_2}(E(h(X)|X_2))$. The estimator is

$$\hat{\mu}_{RB} = \frac{1}{m} \sum_{i=1}^m E(h(X)|X_{2,i})$$

where we draw X , but only use X_2 , averaging over the simulated values of X_2 . This is called Rao-Blackwellization because it relates to the idea of conditioning on a sufficient statistic. It has lower variance because the variance of each term in the sum of the standard MC estimator is $\text{Var}(h(X))$, which is greater than $\text{Var}(E(h(X)|X_2))$ based on the usual iterated variance formula: $V(X) = E(V(X|Y)) + V(E(X|Y)) \Rightarrow V(E(X|Y)) < V(X)$.

2 Random number generation (RNG)

At the core of simulations is the ability to generate random numbers, and based on that, random variables. On a computer, our goal is to generate sequences of pseudo-random numbers that behave like random numbers but are replicable. The reason that replicability is important is so that we can reproduce the simulation.

2.1 Generating random uniforms on a computer

Generating a sequence of random standard uniforms is the basis for all generation of random variables, since random uniforms (either a single one or more than one) can be used to generate values from other distributions. Most random numbers on a computer are pseudo-random. The numbers are chosen from a deterministic stream of numbers that behave like random numbers but are actually a finite sequence (recall that both integers and real numbers on a computer are actually discrete and there are finitely many distinct values), so it's actually possible to get repeats. The seed of a RNG is the place within that sequence where you start to use the pseudo-random numbers.

Many RNG methods are sequential congruential methods. The basic idea is that the next value is

$$u_k = f(u_{k-1}, \dots, u_{k-j}) \bmod m$$

for some function, f , and some positive number m . Often $j = 1$. \bmod just means to take the remainder after dividing by m . One then generates the random standard uniform value as u_k/m , which by construction is in $[0, 1]$.

Given the construction, such sequences are periodic if the subsequence ever reappears, which is of course guaranteed because there is a finite number of possible values. One key to a good random number generator (RNG) is to have a very long period.

An example of a sequential congruential method is a basic linear congruential generator:

$$u_k = (au_{k-1}) \bmod m$$

with integer a , m , and u_i values. Here the periodicity can't exceed $m - 1$ (the method is set up so that we never get $u_k = 0$ as this causes the algorithm to break), so we only have $m - 1$ possible values. The seed is the initial state, u_0 - i.e., the point in the sequence at which we start. By setting the seed you guarantee reproducibility since given a starting value, the sequence is deterministic. In general a and m are chosen to be large, but of course they can't be too large if they are to be represented as a computer integer. The standard values of m are Mersenne primes, which have the form $2^p - 1$ (but these are not prime for all p), with $m = 2^{31} - 1$ common.

A wide variety of different RNG have been proposed. Many have turned out to have substantial defects based on tests designed to assess if the behavior of the RNG mimics true randomness. Some of the behavior we want to ensure is uniformity of each individual random deviate, independence of sequences of deviates, and multivariate uniformity of subsequences. One test of a RNG that many RNGs don't perform well on is to assess the properties of k -tuples - subsequences of length k , which should be independently distributed in the k -dimensional unit hypercube. Unfortunately, linear congruential methods produce values that lie on a simple lattice in k -space, i.e., the points

are not selected from q^k uniformly spaced points, where q is the the number of unique values. Instead, points often lie on parallel lines in the hypercube.

Combining generators can yield better generators. The Wichmann-Hill is an option in R and is a combination of three linear congruential generators with $a = \{171, 172, 170\}$, $m = \{30269, 30307, 30323\}$, and $u_i = (x_i/30269 + y_i/30307 + z_i/30323) \bmod 1$ where x , y , and z are generated from the three individual generators.

R uses something called the Mersenne twister, which is in the class of generalized feedback shift registers (GFSR). The basic idea of a GFSR is to come up with a deterministic generator of bits (i.e., a way to generate sequences of 0s and 1s), B_i , $i = 1, 2, 3, \dots$. The pseudo-random numbers are then determined as sequential subsequences of length L from $\{B_i\}$, considered as a base-2 number and dividing by 2^L to get a number in $(0, 1)$. In general the sequence of bits is generated by taking B_i to be the *exclusive or* [i.e., $0+0 = 0$; $0 + 1 = 1$; $1 + 0 = 1$; $1 + 1 = 0$] summation of two previous bits further back in the sequence where the lengths of the lags are carefully chosen.

Additional notes Generators should give you the same sequence of random numbers, starting at a given seed, whether you ask for a bunch of numbers at once, or sequentially ask for individual numbers.

When one invokes a RNG without a seed, they generally have a method for choosing a seed, often based on the system clock.

There have been some attempts to generate truly random numbers based on physical randomness. One that is based on quantum physics is <http://www.idquantique.com/true-random-number-generator/quantis-usb-pcie-pci.html>. Another approach is based on lava lamps!

2.2 RNG in R

We can change the RNG in R using `RNGkind()`. We can set the seed with `set.seed()`. The seed is stored in `.Random.seed`. The first element indicates the type of RNG (and the type of normal RV generator). The remaining values are specific to the RNG. In the demo code, we've seen that for Wichmann-Hill, the remaining three numbers are the current values of $\{x, y, z\}$.

In R the default RNG is the Mersenne twister (`?RNGkind`), which is considered to be state-of-the-art – it has some theoretical support, has performed reasonably on standard tests of pseudorandom numbers and has been used without evidence of serious failure. Plus it's fast (because bitwise operations are fast). In fact this points out one of the nice features of R, which is that for something as important as this, the default is generally carefully chosen by R's developers. The particular Mersenne twister used has a periodicity of $2^{19937} - 1 \approx 10^{6000}$. Practically speaking

this means that if we generated one random uniform per nanosecond for 10 billion years, then we would generate 10^{25} numbers, well short of the period. So we don't need to worry about the periodicity! The seed for the Mersenne twister is a set of 624 32-bit integers plus a position in the set, where the position is `.Random.seed[2]`.

We can set the seed by passing an integer to `set.seed()`, which then sets as many actual seeds as required for a given generator. Here I'll refer to the integer passed to `set.seed()` as *the* seed. Ideally, nearby seeds generally should not correspond to getting sequences from the stream that are closer to each other than far away seeds. According to Gentle (CS, p. 327) the input to `set.seed()` should be an integer, $i \in \{0, \dots, 1023\}$, and each of these 1024 values produces positions in the RNG sequence that are "far away" from each other. I don't see any mention of this in the R documentation for `set.seed()` and furthermore, you can pass integers larger than 1023 to `set.seed()`, so I'm not sure how much to trust Gentle's claim. More on generating parallel streams of random numbers below.

So we get replicability by setting the seed to a specific value at the beginning of our simulation. We can then set the seed to that same value when we want to replicate the simulation. We can also save the state of the RNG and pick up where we left off:

```
> savedSeed = .Random.seed
> # execute a bunch of code
> .Random.seed=savedSeed
```

At this point any generation of random numbers will pick where you had left off, ignoring what happened in between saving to `savedSeed` and resetting. In some cases you might want to reset the seed upon exit from a function so that a user's random number stream is unaffected:

```
> f = function(args) {
+   oldseed = .Random.seed
+   # other code
+   .Random.seed <- oldseed}
```

Note the need to reassign to the global variable `.Random.seed`.

`RNGversion()` allows you to revert to RNG from previous versions of R, which is very helpful for reproducibility.

The RNGs in R generally return 32-bit (4-byte) integers converted to doubles, so there are at most 2^{32} distinct values. This means you could get duplicated values in long runs, but this does not violate the comment about the periodicity because the two values after the two duplicated numbers will not be duplicates of each other.

One way to proceed if you're using both R and C is to have C use the R RNG, using the `unif_rand` (corresponding to `runif()`) and `norm_rand` (corresponding to `rnorm()`) functions in the `Rmath` library. This way you have a consistent source of random numbers and don't need to worry

about issues with RNG in C. If you call C from R, this should approach should also work; you could also generate all the random numbers you need in R and pass them to the C function.

Note that whenever a random number is generated, the software needs to retain information about what has been generated, so this is an example where a function must have a side effect not observed by the user. R frowns upon this sort of thing, but it's necessary in this case.

2.3 Random slippage

If the exact floating point representations of a random number sequence differ, even in the 14th, 15th, 16th decimal places, if you run a simulation long enough, such a difference can be enough to change the result of some conditional calculation. Suppose your code involves:

```
> if(x>0) { } else{ }
```

As soon as a small difference changes the result of testing $x>0$, the remainder of the simulation can change entirely. This happened to me in my thesis as a result of the difference of an AMD and Intel processor, and took a while to figure out.

2.4 RNG in parallel

We can generally rely on the RNG in R to give a reasonable set of values. One time when we want to think harder is when doing work with RNG in parallel on multiple processors. We would of course not want to set the same seed on each, or we'd get no replication at all. Furthermore, if we're splitting a simulation of IID replicates amongst multiple processors, we want to avoid having the RNG sequence overlap for any of the processors. So we need a methodology that gives us distinct subsequences from a given RNG sequence.

Several packages in R will do this for us. *Rsprng* is one that's easy to implement. The basic syntax on the i th processor is:

```
> init.sprng (nstream = 10, streamno = i,  
+           seed = 0, kindprng = "default")  
> runif(100) # i.e., insert your code here  
> free.sprng() # when you're done
```

Another approach is in the *rlecuyer* package.

3 Generating random variables

There are a variety of methods for generating from common distributions (normal, gamma, beta, Poisson, t, etc.). Since these tend to be built into R and presumably use good algorithms, we

won't go into them. A variety of statistical computing and Monte Carlo books describe the various methods. Many are built on the relationships between different distributions - e.g., a beta random variable (RV) can be generated from two gamma RVs.

Also note that you can call the C functions that implement the R distribution functions as a library (*Rmath*), so if you're coding in C or another language, you should be able to make use of the standard functions: $\{r,p,q,d\}\{norm,t,gamma,binom,pois,etc.\}$ (as well as a variety of other R math functions, which can be seen in *Rmath.h*). Phil Spector has a writeup on this (<http://www.stat.berkeley.edu/classes/s243/rmath.html>) and material can also be found in the *Writing R Extensions* manual on CRAN (section 6.16).

3.1 Multivariate distributions

The *mvtnorm* package supplies code for working with the density and CDF of multivariate normal and t distributions.

To generate a multivariate normal, we've seen the standard method based on the Cholesky decomposition:

```
> U = chol(covMat)
> crossprod(U, rnorm(n))
```

For a singular covariance matrix we can use the Cholesky with pivoting, setting as many rows to zero as the rank deficiency. Then when we generate the multivariate normals, they respect the constraints implicit in the rank deficiency.

3.2 Inverse CDF

Most of you know the inverse CDF method. To generate $X \sim F$ where F is a CDF and is an invertible function, first generate $Z \sim \mathcal{U}(0, 1)$, then $x = F^{-1}(z)$. For discrete CDFs, one can work with a discretized version. For multivariate distributions, one can work with a univariate marginal and then a sequence of univariate conditionals: $f(x_1)f(x_2|x_1) \cdots f(x_k|x_{k-1}, \dots, x_1)$, when the distribution allows this analytic decomposition.

3.3 Rejection sampling

The basic idea of rejection sampling (RS) relies on the introduction of an auxiliary variable, u . Suppose $X \sim F$. Then we can write $f(x) = \int_0^{f(x)} du$. Thus X is the marginal density of X in the joint density, $(X, U) \sim \mathcal{U}\{(x, u) : 0 < u < f(x)\}$. Now we'd like to use this in a way that relies only on evaluating $f(x)$ without having to draw from f .

To implement this we draw from a larger set and then only keep draws for which $u < f(x)$. We choose a density, g , that is easy to draw from and that can *majorize* f , which means there exists a constant c s.t. , $cg(x) \geq f(x) \forall x$. In other words we have that $cg(x)$ is an upper envelope for $f(x)$. The algorithm is

1. generate $y \sim g$
2. generate $u \sim \mathcal{U}(0, 1)$
3. if $u \leq f(y)/cg(y)$ then use y ; otherwise go back to step 1

The intuition here is graphical: we generate from under a curve that is always above $f(x)$ and accept only when u puts us under $f(x)$ relative to the majorizing density. A key here is that the majorizing density have fatter tails than the density of interest, so that the constant c can exist. So we could use a t to generate from a normal but not the reverse. We'd like c to be small to reduce the number of rejections because it turns out that $\frac{1}{c} = \frac{\int f(x)dx}{\int cg(x)dx}$ is the acceptance probability. This approach works in principle for multivariate densities but as the dimension increases, the proportion of rejections grows, because more of the volume under $cg(x)$ is above $f(x)$.

If f is costly to evaluate, we can sometimes reduce calculation using a lower bound on f . In this case we accept if $u \leq f_{\text{low}}(y)/cg_Y(y)$. If it is not, then we need to evaluate the ratio in the usual rejection sampling algorithm. This is called squeezing.

One example of RS is to sample from a truncated normal. Of course we can just sample from the normal and then reject, but this can be inefficient, particularly if the truncation is far in the tail (a case in which inverse CDF suffers from numerical difficulties). Suppose the truncation point is greater than zero. Working with the standardized version of the normal, you can use an translated exponential with lower end point equal to the truncation point as the majorizing density (Robert 1995; Statistics and Computing, and see calculations in the demo code). For truncation less than zero, just make the values negative.

3.4 Adaptive rejection sampling

The difficulty of RS is finding a good enveloping function. Adaptive rejection sampling refines the envelope as the draws occur, in the case of a continuous, differentiable, log-concave density. The basic idea considers the log of the density and involves using tangents or secants to define an upper envelope and secants to define a lower envelope for a set of points in the support of the distribution. The result is that we have piecewise exponentials (since we are exponentiating from straight lines on the log scale) as the bounds. We can sample from the upper envelope based on sampling from a discrete distribution and then the appropriate exponential. The lower envelope is

used for squeezing. We add points to the set that defines the envelopes whenever we accept a point that requires us to evaluate $f(x)$ (the points that are accepted based on squeezing are not added to the set). We'll talk this through some in class.

3.5 Importance sampling

Importance sampling (IS) allows us to estimate expected values, with some commonalities with rejection sampling.

$$E_f(h(X)) = \int h(x) \frac{f(x)}{g(x)} g(x) dx$$

so $\hat{\mu} = \frac{1}{m} \sum_i h(x_i) \frac{f(x_i)}{g(x_i)}$ where $w_i^* = f(x_i)/g(x_i)$ act as weights. Often in Bayesian contexts, we know $f(x)$ only up to a normalizing constant. In this case we need to use $w_i = w_i^* / \sum_i w_i^*$.

Here we don't require the majorizing property, just that the densities have common support, but things can be badly behaved if we sample from a density with lighter tails than the density of interest. So in general we want g to have heavier tails. More specifically for a low variance estimator of μ , we would want that $f(x_i)/g(x_i)$ is large only when $h(x_i)$ is very small, to avoid having overly influential points.

This suggests we can reduce variance in an IS context by oversampling x for which $h(x)$ is large and undersampling when it is small, since $\text{Var}(\hat{\mu}) = \frac{1}{m} \text{Var}(h(X) \frac{f(X)}{g(X)})$. An example is that if h is an indicator function that is 1 only for rare events, we should oversample rare events and then the IS estimator corrects for the oversampling.

What if we actually want a sample from f as opposed to estimating the expected value above? We can draw x from the unweighted sample, $\{x_i\}$, with weights $\{w_i\}$. This is called sampling importance resampling (SIR).

3.6 Ratio of uniforms

If U and V are uniform in $C = \{(u, v) : 0 \leq u \leq \sqrt{f(v/u)}\}$ then $X = V/U$ has density proportion to f . The basic algorithm is to choose a rectangle that encloses C and sample until we find $u \leq f(v/u)$. Then we use $x = v/u$ as our RV. The larger region enclosing C is the majorizing region and a simple approach (if $f(x)$ and $x^2 f(x)$ are bounded in C) is to choose the rectangle, $0 \leq u \leq \sup_x \sqrt{f(x)}$, $\inf_x x \sqrt{f(x)} \leq v \leq \sup_x x \sqrt{f(x)}$.

One can also consider truncating the rectangular region, depending on the features of f .

Monahan recommends the ratio of uniforms, particularly a version for discrete distributions (p. 323 of the 2nd edition).

4 Design of simulation studies

Let's pose a (relatively) concrete example. We have developed a statistical method that predicts whether someone will get cancer over the next 10 years based on gene expression data. A number of other methods for such classification also exist and we want to compare performance in a simulation study. What do we need to decide in terms of carrying out the study?

4.1 Basic steps of a simulation study

1. Specify what makes up an individual experiment: sample size, distributions, parameters, statistic of interest, etc.
2. Determine what inputs, if any, to vary; e.g., sample sizes, parameters, data generating mechanisms
3. Write code to carry out the individual experiment and return the quantity of interest
4. For each combination of inputs, repeat the experiment m times.
5. Summarize the results for each combination of interest, quantifying simulation uncertainty
6. Report the results in graphical or tabular form

4.2 Overview

Since a simulation study is an experiment, we should use the same principles of design and analysis we would recommend when advising a practitioner on setting up an experiment.

These include efficiency, reporting of uncertainty, reproducibility and documentation.

In generating the data for a simulation study, we want to think about what structure real data would have that we want to mimic in the simulation study: distributional assumptions, parameter values, dependence structure, outliers, random effects, sample size (n), etc.

All of these may become input variables in a simulation study. Often we compare two or more statistical methods conditioning on the data context and then assess whether the differences between methods vary with the data context choices. E.g., if we compare an MLE to a robust estimator, which is better under a given set of choices about the data generating mechanism and how sensitive is the comparison to changing the features of the data generating mechanism? So the “treatment variable” is the choice of statistical method. We're then interested in sensitivity to the conditions.

Often we can have a large number of replicates (m) because the simulation is fast on a computer, so we can sometimes reduce the simulation error to essentially zero and thereby avoid reporting uncertainty. To do this, we need to calculate the simulation standard error, generally, s/\sqrt{m} and see how it compares to the effect sizes. This is particularly important when reporting on the bias of a statistical method.

We might denote the data, which could be the statistical estimator under each of two methods as Y_{ijklq} , where i indexes treatment, j, k, l index different additional input variables, and $q \in \{1, \dots, m\}$ indexes the replicate. E.g., j might index whether the data are from a t or normal, k the value of a parameter, and l the dataset sample size (i.e., different levels of n).

One can think about choosing m based on a basic power calculation, though since we can always generate more replicates, one might just proceed sequentially and stop when the precision of the results is sufficient.

When comparing methods, it's best to use the same simulated datasets for each level of the treatment variable and to do an analysis that controls for the dataset (i.e., for the random numbers used), thereby removing some variability from the error term. A simple example is to do a paired analysis, where we look at differences between the outcome for two statistical methods, pairing based on the simulated dataset.

One can even use the “same” random number generation for the replicates under different conditions. E.g., in assessing sensitivity to a t vs. normal data generating mechanism, we might generate the normal RVs and then for the t use the same random numbers, in the sense of using the same quantiles of the t as were generated for the normal - this is pretty easy, along the lines of: `qt(pnorm(x))` - this helps to control for random differences between the datasets.

4.3 Experimental Design

A typical context is that one wants to know the effect of multiple input variables on some outcome. Often, scientists, and even statisticians doing simulation studies will vary one input variable at a time. As we know from standard experimental design, this is inefficient.

The standard strategy is to discretize the inputs, each into a small number of levels. If we have a small enough number of inputs and of levels, we can do a full factorial design (potentially with replication). For example if we have three inputs and three levels each, we have 3^3 different treatment combinations. Choosing the levels in a reasonable way is obviously important.

As the number of inputs and/or levels increases to the point that we can't carry out the full factorial, a fractional factorial is an option. This carefully chooses which treatment combinations to omit. The goal is to achieve balance across the levels in a way that allows us to estimate lower level effects (in particular main effects) but not all high-order interactions. What happens

is that high-order interactions are aliased to (confounded with) lower-order effects. For example you might choose a fractional factorial design so that you can estimate main effects and two-way interactions but not higher-order interactions.

In interpreting the results, I suggest focusing on the decomposition of sums of squares and not on statistical significance. In most cases, we expect the inputs to have at least some effect on the outcome, so the null hypothesis is a straw man. Better to assess the magnitude of the impacts of the different inputs.

When one has a very large number of inputs, one can use the Latin hypercube approach to sample in the input space in a uniform way, spreading the points out so that each input is sampled uniformly. Assume that each input is $\mathcal{U}(0, 1)$ (one can easily transform to whatever marginal distributions you want). Suppose that you can run m samples. Then for each input variable, we divide the unit interval into m bins and randomly choose the order of bins and the position within each bin. This is done independently for each variable and then combined to give m samples from the input space. We would then analyze main effects and perhaps two-way interactions to assess which inputs seem to be most important.

5 Implementation of simulation studies

5.1 Computational efficiency

Since we often need to do some sort of looping, writing code in C/C++ and compiling and linking to the code from R may be a good strategy, albeit one not covered in this course. Parallel processing is often helpful for simulation studies. The reason is that simulation studies are embarrassingly parallel - we can send each replicate to a different computer processor and then collect the results back, and the speedup should scale directly with the number of processors we used. Hopefully we'll get to see the use of *foreach()* to farm out replicates to different processors in the last unit of the course on parallel processing.

Handy functions in R include *expand.grid()* to get all combinations of a set of vectors and the *replicate()* function in R, which will carry out the same R expression (often a function call) repeated times. This can replace the use of a *for* loop with some gains in cleanliness of your code. Storing results in an array is a natural approach.

5.2 Analysis and reporting

Often results are reported simply in tables, but it can be helpful to think through whether a graphical representation is more informative (sometimes it's not or it's worse, but in some cases it may be

much better).

You should set the seed when you start the experiment, so that it's possible to replicate it. It's also a good idea to save the current value of the seed whenever you save interim results, so that you can restart simulations (this is particularly helpful for MCMC) at the exact point you left off, including the random number sequence.

To enhance reproducibility, it's good practice to post your simulation code (and potentially data) on your website or as supplementary material with the journal. One should report sample sizes and information about the random number generator.

Here are JASA's requirements on documenting computations:

“Results Based on Computation - Papers reporting results based on computation should provide enough information so that readers can evaluate the quality of the results. Such information includes estimated accuracy of results, as well as descriptions of pseudorandom-number generators, numerical algorithms, computers, programming languages, and major software components that were used.”

Optimization

November 30, 2011

References:

- Gentle: *Computational Statistics*
- Lange: *Optimization*
- Monahan: *Numerical Methods of Statistics*
- Givens and Hoeting: *Computational Statistics*
- Materials online from Stanford's [EE364a course](#) on convex optimization, including Boyd and Vandenberghe's (online) book *Convex Optimization*, which is linked to from the course webpage.

1 Notation

We'll make use of the first derivative (the gradient) and second derivative (the Hessian) of functions. We'll generally denote univariate and multivariate functions without distinguishing them as $f(x)$ with $x = (x_1, \dots, x_p)$. The (column) vector of first partial derivatives (the gradient) is $f'(x) = \nabla f(x) = (\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_p})^\top$ and the matrix of second partial derivatives (the Hessian) is

$$f''(x) = \nabla_f^2(x) = H_f(x) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_p} \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_p} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_p} & \frac{\partial^2 f}{\partial x_2 \partial x_p} & \cdots & \frac{\partial^2 f}{\partial x_p^2} \end{pmatrix}.$$

In considering iterative algorithms, I'll use $x_0, x_1, \dots, x_t, x_{t+1}$ to indicate the sequence of values as we search for the optimum, denoted x^* . x_0 is the starting point, which we must choose (sometimes carefully). If it's unclear at any point whether I mean a value of x in the sequence or

a subelement of the x vector, let me know, but hopefully it will be clear from context most of the time.

I'll try to use x (or if we're talking explicitly about a likelihood, θ) to indicate the argument with respect to which we're optimizing and Y to indicate data involved in a likelihood. I'll try to use z to indicate covariates/regressors so there's no confusion with x .

2 Overview

The basic goal here is to optimize a function numerically when we cannot find the maximum (or minimum) analytically. Some examples:

1. Finding the MLE for a GLM
2. Finding least squares estimates for a nonlinear regression model,

$$Y_i \sim \mathcal{N}(g(z_i; \beta), \sigma^2)$$

where $g(\cdot)$ is nonlinear and we seek to find the value of $\theta = (\beta, \sigma^2)$ that best fits the data.

3. Maximizing a likelihood under constraints

Maximum likelihood estimation and variants thereof is a standard situation in which optimization comes up.

We'll focus on **minimization**, since any maximization of f can be treated as minimization of $-f$. The basic setup is to find

$$\arg \min_{x \in D} f(x)$$

where D is the domain. Sometimes $D = \mathbb{R}^p$ but other times it imposes constraints on x . When there are no constraints, this is unconstrained optimization, where any x for which $f(x)$ is defined is a possible solution. We'll assume that f is continuous as there's little that can be done systematically if we're dealing with a discontinuous function.

In one dimension, minimization is the same as root-finding with the derivative function, since the minimum of a differentiable function can only occur at a point at which the derivative is zero. So with differentiable functions we'll seek to find x s.t. $f'(x) = \nabla f(x) = 0$. To ensure a minimum, we want that for all y in a neighborhood of x^* , $f(y) \geq f(x^*)$, or (for twice differentiable functions) $f''(x^*) = \nabla^2 f(x^*) = H_f(x^*) \geq 0$, i.e., that the Hessian is positive semi-definite.

Different strategies are used depending on whether D is discrete and countable, or continuous, dense and uncountable. We'll concentrate on the continuous case but the discrete case can arise in statistics, such as in doing variable selection.

In general we rely on the fact that we can evaluate f . Often we make use of analytic or numerical derivatives of f as well.

To some degree, optimization is a solved problem, with good software implementations, so it raises the question of how much to discuss in this class. The basic motivation for going into some of the basic classes of optimization strategies is that the function being optimized changes with each problem and can be tricky to optimize, and I want you to know something about how to choose a good approach when you find yourself with a problem requiring optimization. Finding global, as opposed to local, minima can also be an issue.

Note that I'm not going to cover MCMC (Markov chain Monte Carlo) methods, which are used for approximating integrals and sampling from posterior distributions in a Bayesian context and in a variety of ways for optimization. If you take a Bayesian course you'll cover this in detail, and if you don't do Bayesian work, you probably won't have much need for MCMC, though it comes up in MCEM (Monte Carlo EM) and simulated annealing, among other places.

Goals for the unit Optimization is a big topic. Here's what I would like you to get out of this:

1. an understanding of line searches,
2. an understanding of multivariate derivative-based optimization and how line searches are useful within this,
3. an understanding of derivative-free methods,
4. an understanding of the methods used in R's optimization routines, their strengths and weaknesses, and various tricks for doing better optimization in R, and
5. a basic idea of what convex optimization is and when you might want to go learn more about it.

3 Univariate function optimization

We'll start with some strategies for univariate functions. These can be useful later on in dealing with multivariate functions.

3.1 Golden section search

This strategy requires only that the function be unimodal.

Assume we have a single minimum, in $[a, b]$. We choose two points in the interval and evaluate them, $f(x_1)$ and $f(x_2)$. If $f(x_1) < f(x_2)$ then the minimum must be in $[a, x_2]$, and if the converse

in $[x_1, b]$. We proceed by choosing a new point in the new, smaller interval and iterate. At each step we reduce the length of the interval in which the minimum must lie. The primary question involves what is an efficient rule to use to choose the new point at each iteration.

Suppose we start with x_1 and x_2 s.t. they divide $[a, b]$ into three equal segments. Then we use $f(x_1)$ and $f(x_2)$ to rule out either the leftmost or rightmost segment based on whether $f(x_1) < f(x_2)$. If we have divided equally, we cannot place the next point very efficiently because either x_1 or x_2 equally divides the remaining space, so we are forced to divide the remaining space into relative lengths of 0.25, 0.25, and 0.5. The next time around, we may only rule out the shorter segment, which leads to inefficiency.

The efficient strategy is to maintain the *golden ratio* between the distances between the points using $\phi = (\sqrt{5} - 1)/2 \approx .618$, the golden ratio. We start with $x_1 = a + (1 - \phi)(b - a)$ and $x_2 = a + \phi(b - a)$. Then suppose $f(x_1) < f(x_2)$. We now choose to place x_3 s.t. it uses the golden ratio in the interval $[a, x_1]$: $x_3 = a + (1 - \phi)(x_2 - a)$. Because of the way we've set it up, we once again have the third subinterval, $[x_1, x_2]$, of equal length as the first subinterval, $[a, x_3]$. The careful choice allows us to narrow the search interval by an equal proportion, $1 - \phi$, in each iteration. Eventually we have narrowed the minimum to between x_{t-1} and x_t , where the difference $|x_t - x_{t-1}|$ is sufficiently small (within some tolerance - see Section 4 for details), and we report $(x_t + x_{t-1})/2$. We'll see an example of this on the board in class.

3.2 Bisection method

The bisection method requires the existence of the first derivative but has the advantage over the golden section search of halving the interval at each step.

We start with an initial interval (a_0, b_0) and proceed to shrink the interval. Let's choose a_0 and b_0 , and set x_0 to be the mean of these endpoints. Now we update according to the following algorithm, assuming our current interval is $[a_t, b_t]$:

$$[a_{t+1}, b_{t+1}] = \begin{cases} [a_t, x_t] & \text{if } f'(a_t)f'(x_t) \leq 0 \\ [x_t, b_t] & \text{if } f'(a_t)f'(x_t) > 0 \end{cases}$$

and set x_{t+1} to be the mean of a_{t+1} and b_{t+1} . The basic idea is that if the derivative at both a_t and x_t is negative, then the minimum must be between x_t and b_t , based on the intermediate value theorem. If the derivatives at a_t and x_t are of different signs, then the minimum must be between a_t and x_t .

Since the bisection method reduces the size of the search space by one-half at each iteration, one can work out that each decimal place of precision requires 3-4 iterations. Obviously bisection is more efficient than the golden section search because we reduce by $0.5 > 0.382 = 1 - \phi$, so we've gained information by using the derivative. It requires an evaluation of the derivative

however, while golden section just requires an evaluation of the original function.

Bisection is an example of a *bracketing* method, in which we trap the minimum within a nested sequence of intervals of decreasing length. These tend to be slow, but if the first derivative is continuous, they are robust and don't require that a second derivative exist.

3.3 Newton-Raphson (Newton's method)

3.3.1 Overview

We'll talk about Newton-Raphson (N-R) as an optimization method rather than a root-finding method, but they're just different perspectives on the same algorithm.

For N-R, we need two continuous derivatives that we can evaluate. The benefit is speed, relative to bracketing methods. We again assume the function is unimodal. The minimum must occur at x^* s.t. $f'(x^*) = 0$, provided the second derivative is non-negative at x_0 . So we aim to find a zero (a root) of the first derivative function. Assuming that we have an initial value x_0 that is close to x^* , we have the Taylor series approximation

$$f'(x) \approx f'(x_0) + (x - x_0)f''(x_0).$$

Now set $f'(x) = 0$, since that is the condition we desire (the condition that holds when we are at x^*), and solve for x to get

$$x_1 = x_0 - \frac{f'(x_0)}{f''(x_0)},$$

and iterate, giving us updates of the form $x_{t+1} = x_t - \frac{f'(x_t)}{f''(x_t)}$. What are we doing intuitively? Basically we are taking the tangent to $f(x)$ at x_0 and extrapolating along that line to where it crosses the x-axis to find x_1 . We then reevaluate $f(x_1)$ and continue to travel along the tangents.

One can prove that if $f'(x)$ is twice continuously differentiable, is convex, and has a root, then N-R converges from any starting point.

Note that we can also interpret the N-R update as finding the analytic minimum of the quadratic Taylor series approximation to $f(x)$.

Newton's method converges very quickly (as we'll discuss in Section 4), but if you start too far from the minimum, you can run into serious problems.

3.3.2 Secant method variation on N-R

Suppose we don't want to calculate the second derivative required in the divisor of N-R. We might replace the analytic derivative with a discrete difference approximation based on the secant line

joining $(x_t, f'(x_t))$ and $(x_{t-1}, f'(x_{t-1}))$, giving an approximate second derivative:

$$f''(x_t) \approx \frac{f'(x_t) - f'(x_{t-1})}{x_t - x_{t-1}}.$$

For this variant on N-R, we need two starting points, x_0 and x_1 .

An alternative to the secant-based approximation is to use a standard discrete approximation of the derivative such as

$$f''(x_t) \approx \frac{f'(x_t + h) - f'(x_t - h)}{2h}.$$

A sidenote on how computer precision affects numerical derivatives that we'll come back to in Unit 13 (Numerical Integration/Differentiation). Ideally we would take h very small and get a highly accurate estimate of the derivative. However, the limits of machine precision mean that the difference estimator can behave badly for very small h . Therefore we accept a bias in the estimate, often by taking h to be square root of machine epsilon.

3.3.3 How can Newton's method go wrong?

Let's think about what can go wrong - namely when we could have $f(x_{t+1}) > f(x_t)$? Basically, if $f'(x_t)$ is relatively flat, we can get that $|x_{t+1} - x^*| > |x_t - x^*|$. We'll see an example on the board and the demo code. Newton's method can also go uphill when the second derivative is negative, with the method searching for a maximum.

One nice, general idea is to use a fast method such as Newton's method *safeguarded* by a robust, but slower method. Here's how one can do this for N-R, safeguarding with a bracketing method such as bisection. Basically, we check the N-R proposed move to see if N-R is proposing a step outside of where the root is known to lie based on the previous steps and the gradient values for those steps. If so, we could choose the next step based on bisection.

Another approach is backtracking. If a new value is proposed that yields a larger value of the function, backtrack to find a value that reduces the function. One possibility is a line search but given that we're trying to reduce computation, a full line search is unwise computationally (also in the multivariate Newton's method, we are in the middle of an iterative algorithm for which we will just be going off in another direction anyway at the next iteration). A basic approach is to keep backtracking in halves. A nice alternative is to fit a polynomial to the known information about that slice of the function, namely $f(x_{t+1})$, $f(x_t)$, $f'(x_t)$ and $f''(x_t)$ and find the minimum of the polynomial approximation.

4 Convergence ideas

4.1 Convergence metrics

We might choose to assess whether $f'(x_t)$ is near zero, which should assure that we have reached the critical point. However, in parts of the domain where $f(x)$ is fairly flat, we may find the derivative is near zero even though we are far from the optimum. Instead, we generally monitor $|x_{t+1} - x_t|$ (for the moment, assume x is scalar). We might consider absolute convergence: $|x_{t+1} - x_t| < \epsilon$ or relative convergence, $\frac{|x_{t+1} - x_t|}{|x_t|} < \epsilon$. Relative convergence is appealing because it accounts for the scale of x , but it can run into problems when x_t is near zero, in which case one can use $\frac{|x_{t+1} - x_t|}{|x_t| + \epsilon} < \epsilon$. We would want to account for machine precision in thinking about setting ϵ . For relative convergence a reasonable choice of ϵ would be to use the square root of machine epsilon or about 1×10^{-8} . This is the *reltol* argument in *optim()* in R.

Problems with the optimization may show up in a convergence measure that fails to decrease or cycles (oscillates). Software generally has a stopping rule that stops the algorithm after a fixed number of iterations; these can generally be changed by the user. When an algorithm stops because of the stopping rule before the convergence criterion is met, we say the algorithm has failed to converge. Sometimes we just need to run it longer, but often it indicates a problem with the function being optimized or with your starting value.

For multivariate optimization, we use a distance metric between x_{t+1} and x_t , such as $\|x_{t+1} - x_t\|_p$, often with $p = 1$ or $p = 2$.

4.2 Starting values

Good starting values are important because they can improve the speed of optimization, prevent divergence or cycling, and prevent finding local optima.

Using random or selected multiple starting values can help with multiple optima.

4.3 Convergence rates

Let $\epsilon_t = |x_t - x^*|$. If the limit

$$\lim_{t \rightarrow \infty} \frac{|\epsilon_{t+1}|}{|\epsilon_t|^\beta} = c$$

exists for $\beta > 0$ and $c \neq 0$, then a method is said to have order of convergence β . This basically measures how big the error at the $t + 1$ th iteration is relative to that at the t th iteration.

Bisection doesn't formally satisfy the criterion needed to make use of this definition, but

roughly speaking it has linear convergence ($\beta = 1$). Next we'll see that N-R has quadratic convergence ($\beta = 2$), which is fast.

To analyze convergence of N-R, consider a Taylor expansion of the gradient at the minimum, x^* , around the current value, x_t :

$$f'(x^*) = f'(x_t) + (x^* - x_t)f''(x_t) + \frac{1}{2}(x^* - x_t)^2 f'''(\xi) = 0,$$

for some $\xi \in [x_t, x^*]$. Making use of the N-R update equation: $x_{t+1} = x_t - \frac{f'(x_t)}{f''(x_t)}$, and some algebra, we have

$$\frac{|x^* - x_{t+1}|}{(x^* - x_t)^2} = \frac{1}{2} \frac{f'''(\xi)}{f''(x_t)}.$$

If the limit of the ratio on the right hand side exists and is equal to c , then we see that $\beta = 2$. If c were one, then we see that if we have k digits of accuracy at t , we'd have $2k$ digits at $t + 1$, which justifies the characterization of quadratic convergence being fast. In practice c will moderate the rate of convergence. The smaller c the better, so we'd like to have the second derivative be large and the third derivative be small. The expression also indicates we'll have a problem if $f''(x_t) = 0$ at any point [think about what this corresponds to graphically - what is our next step when $f''(x_t) = 0$?]. The characteristics of the derivatives determine the domain of attraction (the region in which we'll converge rather than diverge) of the minimum.

Givens and Hoeting show that using the secant-based approximation to the second derivative in N-R has order of convergence, $\beta \approx 1.62$.

5 Multivariate optimization

Optimizing as the dimension of the space gets larger becomes increasingly difficult. First we'll discuss the idea of profiling to reduce dimensionality and then we'll talk about various numerical techniques, many of which build off of Newton's method.

5.1 Profiling

A core technique for likelihood optimization is to analytically maximize over any parameters for which this is possible. Suppose we have two sets of parameters, θ_1 and θ_2 , and we can analytically maximize w.r.t θ_2 . This will give us $\hat{\theta}_2(\theta_1)$, a function of the remaining parameters over which analytic maximization is not possible. Plugging in $\hat{\theta}_2(\theta_1)$ into the objective function (in this case generally the likelihood or log likelihood) gives us the profile (log) likelihood solely in terms of the obstinant parameters. For example, suppose we have the regression likelihood with correlated

errors:

$$Y \sim \mathcal{N}(X\beta, \sigma^2\Sigma(\rho)),$$

where $\Sigma(\rho)$ is a correlation matrix that is a function of a parameter, ρ . The maximum w.r.t. β is easily seen to be the GLS estimator $\hat{\beta}(\rho) = (X^\top \Sigma(\rho)^{-1} X)^{-1} X^\top \Sigma(\rho)^{-1} Y$. In general such a maximum is a function of all of the other parameters, but conveniently it's only a function of ρ here. This gives us the initial profile likelihood

$$\frac{1}{(\sigma^2)^{n/2} |\Sigma(\rho)|^{1/2}} \exp\left(-\frac{(Y - X\hat{\beta}(\rho))^\top \Sigma(\rho)^{-1} (Y - X\hat{\beta}(\rho))}{2\sigma^2}\right).$$

We then notice that the likelihood is maximized w.r.t. σ^2 at $\hat{\sigma}^2(\rho) = \frac{(Y - X\hat{\beta}(\rho))^\top \Sigma(\rho)^{-1} (Y - X\hat{\beta}(\rho))}{n}$. This gives us the final profile likelihood,

$$\frac{1}{|\Sigma(\rho)|^{1/2}} \frac{1}{(\hat{\sigma}^2(\rho))^{n/2}} \exp\left(-\frac{1}{2}n\right),$$

a function of ρ only, for which numerical optimization is much simpler.

5.2 Newton-Raphson (Newton's method)

For multivariate x we have the Newton-Raphson update $x_{t+1} = x_t - f''(x_t)^{-1} f'(x_t)$, or in our other notation,

$$x_{t+1} = x_t - H_f(x_t)^{-1} \nabla f(x_t)$$

In class we'll see an example of finding the nonlinear least squares fit to some weight loss data to fit the model

$$Y_i = \beta_0 + \beta_1 2^{-t_i/\beta_2} + \epsilon_i.$$

Some of the things we need to worry with Newton's method in general about are (1) good starting values, (2) positive definiteness of the Hessian, and (3) avoiding errors in finding the derivatives.

A note on the positive definiteness: since the Hessian may not be positive definite (although it may well be, provided the function is approximately locally quadratic), one can consider modifying the Cholesky decomposition of the Hessian to enforce positive definiteness by adding diagonal elements to H_f as necessary.

5.3 Fisher scoring variant on NR

The Fisher information (FI) is the expected value of the outer product of the gradient of the log-likelihood with itself

$$I(\theta) = E_f(\nabla f(y)\nabla f(y)^\top),$$

where the expected value is with respect to the data distribution. Under regularity conditions (true for exponential families), the expectation of the Hessian of the log-likelihood is minus the Fisher information. We get the observed Fisher information by plugging the data values into either expression instead of taking the expected value.

Thus, standard NR can be thought of as using the observed Fisher information to find the updates. Instead, if we can compute the expectation, we can use minus the FI in place of the Hessian. The result is the Fisher scoring (FS) algorithm. Basically instead of using the Hessian for a given set of data, we are using the FI, which we can think of as the average Hessian over repeated samples of data from the data distribution. FS and NR have the same convergence properties (i.e., quadratic convergence) but in a given problem, one may be computationally or analytically easier. Givens and Hoeting comment that FS works better for rapid improvements at the beginning of iterations and NR better for refinement at the end.

In the demo code, we try out Fisher scoring in the weight loss example.

The Gauss-Newton algorithm for nonlinear least squares involves using the FI in place of the Hessian in determining a Newton-like step. `nls()` in R uses this approach. Note that this is not exactly the same updating as our manual coding of FS for the weight loss example.

Connections between statistical uncertainty and ill-conditionedness When either the observed or expected FI matrix is nearly singular this means we have a small eigenvalue in the inverse covariance (the precision), which means a large eigenvalue in the covariance matrix. This indicates some linear combination of the parameters has low precision (high variance), and that in that direction the likelihood is nearly flat. As we've seen with N-R, convergence slows with shallow gradients, and we may have numerical problems in determining good optimization steps when the likelihood is sufficiently flat. So convergence problems and statistical uncertainty go hand in hand. One, but not the only, example of this is with nearly collinear regressors.

5.4 IRLS (IWLS) for GLMs

As most of you know, iterative reweighted least squares (also called iterative weighted least squares) is the standard method for estimation with GLMs. It involves linearizing the model and using working weights and working variances and solving a weighted least squares (WLS) problem (the generic WLS solution is $\hat{\beta} = (X^\top W X)^{-1} X^\top W Y$).

Exponential families can be expressed as

$$f(y; \theta, \phi) = \exp((y\theta - b(\theta))/a(\phi) + c(y, \phi)),$$

with $E(Y) = b'(\theta)$ and $\text{Var}(Y) = b''(\theta)$. If we have a GLM in the canonical parameterization (log link for Poisson data, logit for binomial), we have the natural parameter θ equal to the linear predictor, $\theta = \eta$. A standard linear predictor would simply be $\eta = X\beta$.

Considering NR for a GLM in the canonical parameterization (and ignoring $a(\phi)$, which is one for logistic and Poisson regression), we find that the gradient is the inner product of the covariates and a residual vector, $\nabla f(\beta) = (Y - E(Y))^\top X$, and the Hessian is $\nabla^2 f(\beta) = -X^\top W X$ where W is a diagonal matrix with $\{\text{Var}(Y_i)\}$ on the diagonal (the working weights). Note that both $E(Y)$ and the variances in W depend on β , so these will change as we iteratively update β . Therefore, the NR update is

$$\beta_{t+1} = \beta_t + (X^\top W_t X)^{-1} X^\top (Y - E(Y)_t)$$

where $E(Y)_t$ and W_t are the values at the current parameter estimate, β_t . For example, for logistic regression, $W_{t,ii} = p_{ti}(1 - p_{ti})$ and $E(Y)_{ti} = p_{ti}$ where $p_{ti} = \frac{\exp(X_i^\top \beta_t)}{1 + \exp(X_i^\top \beta_t)}$. In the canonical parameterization of a GLM, the Hessian does not depend on the data, so the observed and expected FI are the same, and therefore NR and FS are the same.

The update above can be rewritten in the standard form of IRLS as a WLS problem,

$$\beta_{t+1} = (X^\top W_t X)^{-1} X^\top W_t \tilde{Y}_t,$$

where the so-called working observations are $\tilde{Y}_t = X\beta_t + W_t^{-1}(Y - E(Y)_t)$. Note that these are on the scale of the linear predictor.

While Fisher scoring is standard for GLMs, you can also use general purpose optimization routines.

IRLS is a special case of the general Gauss-Newton method for nonlinear least squares.

5.5 Descent methods and Newton-like methods

More generally a Newton-like method has updates of the form

$$x_{t+1} = x_t - \alpha_t M_t^{-1} f'(x_t).$$

We can choose M_t in various ways, including as an approximation to the second derivative.

This opens up several possibilities:

1. using more computationally efficient approximations to the second derivative,

2. avoiding steps that do not go in the correct direction (i.e., go uphill when minimizing), and
3. scaling by α_t so as not to step too far.

Let's consider a variety of strategies.

5.5.1 Descent methods

The basic strategy is to choose a good direction and then choose the longest step for which the function continues to decrease. Suppose we have a direction, p_t . Then we need to move $x_{t+1} = x_t + \alpha_t p_t$, where α_t is a scalar, choosing a good α_t . We might use a line search (e.g., bisection or golden section search) to find the local minimum of $f(x_t + \alpha_t p_t)$ with respect to α_t . However, we often would not want to run to convergence, since we'll be taking additional steps anyway.

Steepest descent chooses the direction as the steepest direction downhill, setting $M_t = I$, since the gradient gives the steepest direction uphill (the negative sign in the equation below has us move directly downhill rather than directly uphill). A better approach is to scale the step

$$x_{t+1} = x_t - \alpha_t f'(x_t)$$

where the contraction, or step length, parameter α_t is chosen sufficiently small to ensure that we descend, via some sort of line search. The critical downside to steepest ascent is that when the contours are elliptical, it tends to zigzag; we'll see an example in the demo code. If the contours are circular, steepest descent works well. Newton's method deforms elliptical contours based on the Hessian. Another way to think about this is that steepest descent does not take account of the rate of change in the gradient, while Newton's method does.

The general descent algorithm is

$$x_{t+1} = x_t - \alpha_t M_t^{-1} f'(x_t),$$

where M_t is generally chose to approximate the Hessian and α_t allows us to adjust the step in a smart way. Basically, since the negative gradient tells us the direction that descends (at least within a small neighborhood), if we don't go too far, we should be fine and should work our way downhill. One can work this out formally using a Taylor approximation to $f(x_{t+1}) - f(x_t)$ and see that we make use of M_t being positive definite. (Unfortunately backtracking with positive definite M_t does not give a theoretical guarantee that the method will converge. We also need to make sure that the steps descend sufficiently quickly and that the algorithm does not step along a level contour of f .)

The conjugate gradient algorithm for iteratively solving large systems of equations is all about choosing the direction and the step size in a smart way given the optimization problem at hand.

5.5.2 Quasi-Newton methods such as BFGS

Other replacements for the Hessian matrix include estimates that do not vary with t and finite difference approximations. When calculating the Hessian is expensive, it can be very helpful to substitute an approximation.

A basic finite difference approximation requires us to compute finite differences in each dimension, but this could be computationally burdensome. A more efficient strategy for choosing M_{t+1} is to (1) make use of M_t and (2) make use of the most recent step to learn about the curvature of $f'(x)$ in the direction of travel. One approach is to use a rank one update to M_t .

A basic strategy is to choose M_{t+1} such that the secant condition is satisfied:

$$M_{t+1}(x_{t+1} - x_t) = \nabla f(x_{t+1}) - \nabla f(x_t),$$

which is motivated by the fact that the secant approximates the gradient in the direction of travel. Basically this says to modify M_t in such a way that we incorporate what we've learned about the gradient from the most recent step. M_{t+1} is not fully determined based on this, and we generally impose other conditions, in particular that M_{t+1} is symmetric and positive definite. Defining $s_t = x_{t+1} - x_t$ and $y_t = \nabla f(x_{t+1}) - \nabla f(x_t)$, the unique, symmetric rank one update (why is the following a rank one update?) that satisfies the secant condition is

$$M_{t+1} = M_t + \frac{(y_t - M_t s_t)(y_t - M_t s_t)^\top}{(y_t - M_t s_t)^\top s_t}.$$

If the denominator is positive, M_{t+1} may not be positive definite, but this is guaranteed for non-positive values of the denominator. One can also show that one can achieve positive definiteness by shrinking the denominator toward zero sufficiently.

A standard approach to updating M_t is a commonly-used rank two update that generally results in M_{t+1} being positive definite is

$$M_{t+1} = M_t - \frac{M_t s_t (M_t s_t)^\top}{s_t^\top M_t s_t} + \frac{y_t y_t^\top}{s_t^\top y_t},$$

which is known as the Broyden-Fletcher-Goldfarb-Shanno (BFGS) update. This is one of the methods used in R in `optim()`.

Question: how can we update M_t^{-1} to M_{t+1}^{-1} efficiently? It turns out there is a way to update the Cholesky of M_t efficiently and this is a better approach than updating the inverse.

The order of convergence of quasi-Newton methods is generally slower than the quadratic convergence of N-R because of the approximations but still faster than linear. In general, quasi-Newton methods will do much better if the scales of the elements of x are similar. Lange suggests

using a starting point for which one can compute the expected information, to provide a good starting value M_0 .

Note that for estimating a covariance based on the numerical information matrix, we would not want to rely on M_t from the final iteration, as the approximation may be poor. Rather we would spend the effort to better estimate the Hessian directly at x^* .

5.6 Gauss-Seidel

Gauss-Seidel is also known as back-fitting or cyclic coordinate descent. The basic idea is to work element by element rather than having to choose a direction for each step. For example backfitting used to be used to fit generalized additive models of the form $E(Y) = f_1(z_1) + f_2(z_2) + \dots + f_p(z_p)$.

The basic strategy is to consider the j th component of $f'(x)$ as a univariate function of x_j only and find the root, $x_{j,t+1}$ that gives $f'_j(x_{j,t+1}) = 0$. One cycles through each element of x to complete a single cycle and then iterates. The appeal is that univariate root-finding/minimization is easy, often more stable than multivariate, and quick.

However, Gauss-Seidel can zigzag, since you only take steps in one dimension at a time, as we'll see in the demo code.

In the notes for Unit 10 on linear algebra, I discussed the use of Gauss-Seidel to iteratively solve $Ax = b$ in situations where factorizing A (which of course is $O(n^3)$) is too computationally expensive.

The lasso The lasso uses an L1 penalty in regression and related contexts. A standard formulation for the lasso in regression is to minimize

$$\|Y - X\beta\|_2^2 + \lambda \sum_j |\beta_j|$$

to find $\hat{\beta}$ for a given value of the penalty parameter, λ . A standard strategy to solve this problem is to use coordinate descent, either cyclically, or by using directional derivatives to choose the coordinate likely to decrease the objective function the most (a greedy strategy). We need to use directional derivatives because the penalty function is not differentiable, but does have directional derivatives in each direction. The directional derivative of the objective function for β_j is

$$-2x_j \sum_i (Y_i - X\beta) \pm \lambda$$

where we add λ if $\beta_j \geq 0$ and you subtract λ if $\beta_j < 0$. If $\beta_{j,t}$ is 0, then a step in either direction contributes $+\lambda$ to the derivative as the contribution of the penalty.

Once we have chosen a coordinate, we set the directional derivative to zero and solve for β_j to obtain $\beta_{j,t+1}$.

The LARS (least angle regression) algorithm uses a similar strategy that allows one to compute $\hat{\beta}_\lambda$ for all values of λ at once.

The lasso can also be formulated as the constrained minimization of $\|Y - X\beta\|_2^2$ s.t. $\sum_j |\beta_j| \leq c$, with c now playing the role of the penalty parameter. Solving this minimization problem would take us in the direction of quadratic programming, a special case of convex programming, discussed in Section 9.

5.7 Nelder-Mead

This approach avoids using derivatives or approximations to derivatives. This makes it robust, but also slower than Newton-like methods. The basic strategy is to use a simplex, a polytope of $p + 1$ points in p dimensions (e.g., a triangle when searching in two dimensions, tetrahedron in three dimensions...) to explore the space, choosing to shift, expand, or contract the polytope based on the evaluation of f at the points.

The algorithm relies on four tuning factors: a reflection factor, $\alpha > 0$; an expansion factor, $\gamma > 1$; a contraction factor, $0 < \beta < 1$; and a shrinkage factor, $0 < \delta < 1$. First one chooses an initial simplex: $p + 1$ points that serve as the vertices of a convex hull.

1. Evaluate and order the points, x_1, \dots, x_{p+1} based on $f(x_1) \leq \dots \leq f(x_{p+1})$. Let \bar{x} be the average of the first p x 's.
2. (Reflection) Reflect x_{p+1} across the hyperplane (a line when $p + 1 = 3$) formed by the other points to get x_r , based on α .

$$\bullet x_r = (1 + \alpha)\bar{x} - \alpha x_{p+1}$$

3. If $f(x_r)$ is between the best and worst of the other points, the iteration is done, with x_r replacing x_{p+1} . We've found a good direction to move.
4. (Expansion) If $f(x_r)$ is better than all of the other points, expand by extending x_r to x_e based on γ , because this indicates the optimum may be further in the direction of reflection. If $f(x_e)$ is better than $f(x_r)$, use x_e in place of x_{p+1} . If not, use x_r . The iteration is done.

$$\bullet x_e = \gamma x_r + (1 - \gamma)\bar{x}$$

5. If $f(x_r)$ is worse than all the other points, but better than $f(x_{p+1})$, let $x_h = x_r$. Otherwise $f(x_r)$ is worse than $f(x_{p+1})$ so let $x_h = x_{p+1}$. In either case, we want to concentrate our polytope toward the other points.

(a) (Contraction) Contract x_h toward the hyperplane formed by the other points, based on β , to get x_c . If the result improves upon $f(x_h)$ replace x_{p+1} with x_c . Basically, we haven't found a new point that is better than the other points, so we want to contract the simplex away from the bad point.

- $x_c = \beta x_h + (1 - \beta)\bar{x}$

(b) (Shrinkage) Otherwise (if x_c is not better than x_{p+1}) shrink the simplex toward x_1 . Basically this suggests our step sizes are too large and we should shrink the simplex, shrinking towards the best point.

- $x_i = \delta x_i + (1 - \delta)x_1$ for $i = 2, \dots, p + 1$

Convergence is assessed based on the sample variance of the function values at the points, the total of the norms of the differences between the points in the new and old simplexes, or the size of the simplex.

This is the default in `optim()` in R, however it is relatively slow, so you may want to try one of the alternatives, such as BFGS.

5.8 Simulated annealing (SA)

Simulated annealing is a *stochastic* descent algorithm, unlike the deterministic algorithms we've already discussed. It has a couple critical features that set it aside from other approaches. First, uphill moves are allowed; second whether a move is accepted is stochastic, and finally, as the iterations proceed the algorithm becomes less likely to accept uphill moves.

Assume we are minimizing a negative log likelihood as a function of θ , $f(\theta)$.

The basic idea of simulated annealing is that one modifies the objective function, f in this case, to make it less peaked at the beginning, using a "temperature" variable that changes over time. This helps to allow moves away from local minima, when combined with the ability to move uphill. The name comes from an analogy to heating up a solid to its melting temperature and cooling it slowly - as it cools the atoms go through rearrangements and slowly freeze into the crystal configuration that is at the lowest energy level.

Here's the algorithm. We divide up iterations into stages, $j = 1, 2, \dots$ in which the temperature variable, τ_j , is constant. Like MCMC, we require a proposal distribution to propose new values of θ .

1. Propose to move from θ_t to $\tilde{\theta}$ from a proposal density, $g_t(\cdot|\theta_t)$, such as a normal distribution centered at θ_t .

2. Accept $\tilde{\theta}$ as θ_{t+1} according to the probability $\min(1, \exp((f(\theta_t) - f(\tilde{\theta}))/\tau_j))$ - i.e., accept if a uniform random deviate is less than that probability. Otherwise set $\theta_{t+1} = \theta_t$. Notice that for larger values of τ_j the differences between the function values at the two locations are reduced (just like a large standard deviation spreads out a distribution). So the exponentiation smooths out the objective function when τ_j is large.
3. Repeat steps 1 and 2 m_j times.
4. Increment the temperature and cooling schedule: $\tau_j = \alpha(\tau_{j-1})$ and $m_j = \beta(m_{j-1})$. Back to step 1.

The temperature should slowly decrease to 0 while the number of iterations, m_j , should be large. Choosing these 'schedules' is at the core of implementing SA. Note that we always accept downhill moves in step 2 but we sometimes accept uphill moves as well.

For each temperature, SA produces an MCMC based on the Metropolis algorithm. So if m_j is long enough, we should sample from the stationary distribution of the Markov chain, $\exp(-f(\theta)/\tau_j)$. Provided we can move between local minima, the chain should gravitate toward the global minima because these are increasingly deep (low values) relative to the local minima as the temperature drops. Then as the temperature cools, θ_t should get trapped in an increasingly deep well centered on the global minimum. There is a danger that we will get trapped in a local minimum and not be able to get out as the temperature drops, so the temperature schedule is quite important in trying to avoid this.

A wide variety of schedules have been tried. One approach is to set $m_j = 1 \forall j$ and $\alpha(\tau_{j-1}) = \frac{\tau_{j-1}}{1+a\tau_{j-1}}$ for a small a . For a given problem it can take a lot of experimentation to choose τ_0 and m_0 and the values for the scheduling functions. For the initial temperature it's a good idea to choose it large enough that $\exp((f(\theta_i) - f(\theta_j))/\tau_0) \approx 1$ for any pair $\{\theta_i, \theta_j\}$ in the domain, so that the algorithm can visit the entire space initially.

Simulated annealing can converge slowly. Multiple random starting points or stratified starting points can be helpful for finding a global minimum. However, given the slow convergence, these can also be computationally burdensome.

6 Basic optimization in R

6.1 Core optimization functions

R has several optimization functions.

- *optimize()* is good for 1-d optimization: “The method used is a combination of golden section search and successive parabolic interpolation, and was designed for use with continuous functions.”
- Another option is *uniroot()* for finding the zero of a function, which you can use to minimize a function if you can compute the derivative.
- For more than one variable, *optim()* uses a variety of optimization methods including the robust Nelder-Mead method, the BFGS quasi-Newton method and simulated annealing. You can choose which method you prefer and can try multiple methods. You can supply a gradient function to *optim()* for use with the Newton-related methods but it can also calculate numerical derivatives on the fly. You can have *optim()* return the Hessian at the optimum (based on a numerical estimate), which then allows straightforward calculation of asymptotic variances based on the information matrix.
- Also for multivariate optimization, *nlm()* uses a Newton-style method, for which you can supply analytic gradient and Hessian, or it will estimate these numerically. *nlm()* can also return the Hessian at the optimum.

6.2 Various considerations in using the R functions

As we’ve seen, initial values are important both for avoiding divergence (e.g., in NR), for increasing speed of convergence, and for helping to avoid local optima. So it is well worth the time to try to figure out a good starting value or multiple starting values for a given problem.

Scaling can be important. One useful step is to make sure the problem is well-scaled, namely that a unit step in any parameter has a comparable change in the objective function, preferably approximately a unit change at the optimum. *optim()* allows you to supply scaling information through the *parscale* sub-argument to the *control* argument. Basically if x_j is varying at p orders of magnitude smaller than the other x s, we want to reparameterize to $x_j^* = x_j \cdot 10^p$ and then convert back to the original scale after finding the answer. Or we may want to work on the log scale for some variables, reparameterizing as $x_j^* = \log(x_j)$. We could make such changes manually in our expression for the objective function or make use of arguments such as *parscale*.

If the function itself gives very large or small values near the solution, you may want to rescale the entire function to avoid calculations with very large or small numbers. This can avoid problems such as having apparent convergence because a gradient is near zero, simply because the scale of the function is small. In *optim()* this can be controlled with the *fnscale* sub-argument to *control*.

Always consider your answer and make sure it makes sense, in particular that you haven’t ‘converged’ to an extreme value on the boundary of the space.

Venables and Ripley suggest that it is often worth supplying analytic first derivatives rather than having a routine calculate numerical derivatives but not worth supplying analytic second derivatives. R can do symbolic (i.e., analytic) differentiation to find first and second derivatives using *deriv()*. This is one area where Maple and Mathematica are useful. We'll discuss this in detail in Unit 13 (Numerical Integration/Differentiation)

In general for software development it's obviously worth putting more time into figuring out the best optimization approach and supplying derivatives. For a one-off analysis, you can try a few different approaches and assess sensitivity.

The nice thing about likelihood optimization is that the asymptotic theory tells us that with large samples, the likelihood is approximately quadratic (i.e., the asymptotic normality of MLEs), which makes for a nice surface over which to do optimization.

7 Combinatorial optimization over discrete spaces

Many statistical optimization problems involve continuous domains, but sometimes there are problems in which the domain is discrete. Variable selection is an example of this.

Simulated annealing can be used for optimizing in a discrete space. Another approach uses *genetic algorithms*, in which one sets up the dimensions as loci grouped on a chromosome and has mutation and crossover steps in which two potential solutions reproduce. An example would be in high-dimensional variable selection.

Stochastic search variable selection is a popular Bayesian technique for variable selection that involves MCMC.

8 Convexity

Many optimization problems involve (or can be transformed into) convex functions. Convex optimization (also called convex programming) is a big topic and one that we'll only brush the surface of in Sections 8 and 9. The goal here is to give you enough of a sense of the topic that you know when you're working on a problem that might involve convex optimization, in which case you'll need to go learn more.

Optimization for convex functions is simpler than for ordinary functions because we don't have to worry about local optima - any stationary point (point where the gradient is zero) is a global minimum. A set S in \mathbb{R}^p is convex if any line segment between two points in S lies entirely within S . More generally, S is convex if any convex combination is itself in S , i.e., $\sum_{i=1}^m \alpha_i x_i \in S$ for non-negative weights, α_i , that sum to 1. Convex functions are defined on convex sets - f is convex

if for points in a convex set, $x_i \in S$, we have $f(\sum_{i=1}^m \alpha_i x_i) \leq \sum_{i=1}^m \alpha_i f(x_i)$. Strict convexity is when the inequality is strict (no equality).

The first-order convexity condition relates a convex function to its first derivative: f is convex if and only if $f(x) \geq f(y) + \nabla f(y)^\top (x - y)$ for y and x in the domain of f . We can interpret this as saying that the first order Taylor approximation to f is tangent to and below (or touching) the function at all points.

The second-order convexity condition is that a function is convex if (provided its first derivative exists), the derivative is non-decreasing, in which case we have $f''(x) \geq 0 \forall x$ (for univariate functions). If we have $f''(x) \leq 0 \forall x$ (a concave, or convex down function) we can always consider $-f(x)$, which is convex. Convexity in multiple dimensions means that the gradient is nondecreasing in all dimensions. If f is twice differentiable, then if the Hessian is positive semi-definite, f is convex.

There are a variety of results that allow us to recognize and construct convex functions based on knowing what operations create and preserve convexity. The Boyd book is a good source for material on such operations. Note that norms are convex functions (based on the triangle inequality), $\|\sum_{i=1}^n \alpha_i x_i\| \leq \sum_{i=1}^n \alpha_i \|x_i\|$.

We'll talk about a general algorithm that works for convex functions (the MM algorithm) and about the EM algorithm that is well-known in statistics, and is a special case of MM.

8.1 MM algorithm

The MM algorithm is really more of a principle for constructing problem specific algorithms. MM stands for majorize-minorize. We'll use the majorize part of it to minimize functions - the minorize part is the counterpart for maximizing functions.

Suppose we want to minimize a convex function, $f(x)$. The idea is to construct a majorizing function, at x_t , which we'll call g . g majorizes f at x_t if $f(x_t) = g(x_t)$ and $f(x) \leq g(x) \forall x$.

The iterative algorithm is as follows. Given x_t , construct a majorizing function $g(x_t)$. Then minimize g w.r.t. x (or at least move downhill, such as with a modified Newton step) to find x_{t+1} . Then we iterate, finding the next majorizing function. The algorithm is obviously guaranteed to go downhill, and ideally we use a function g that is easy to work with (i.e., to minimize or go downhill with respect to). Note that we haven't done any matrix inversions or computed any derivatives of f . Furthermore, the algorithm is numerically stable - it does not over- or undershoot the optimum. The downside is that convergence can be quite slow.

The tricky part is finding a good majorizing function. Basically one needs to gain some skill in working with inequalities. The Lange book has some discussion of this.

An example is for estimating regression coefficients for median regression (aka least absolute

deviation regression), which minimizes $f(\theta) = \sum_{i=1}^n |y_i - z_i^\top \theta| = \sum_{i=1}^n |r_i(\theta)|$. Note that $f(\theta)$ is convex because affine functions (in this case $y_i - z_i^\top \theta$) are convex, convex functions of affine functions are convex, and the summation preserves the convexity. We'll work through this example in class. We'll make use the following (commonly-used) inequality, which holds for any concave function, f :

$$f(x) \leq f(y) + f'(y)(x - y).$$

8.2 Expectation-Maximization (EM)

It turns out the EM algorithm that many of you have heard about is a special case of MM. For our purpose here, we'll consider maximization.

The EM algorithm is most readily motivated from a missing data perspective. Suppose you want to maximize $L(\theta|X = x) = f(x|\theta)$ based on available data in a missing data context. Denote the complete data as $Y = (X, Z)$ with Z is missing. As we'll see, in many cases, Z is actually a set of latent variables that we introduce into the problem to formulate it so we can use EM. The canonical example is when Z are membership indicators in a mixture modeling context.

In general, $L(\theta|x)$ may be hard to optimize because it involves an integral over the missing data, Z :

$$f(x|\theta) = \int f(x, z|\theta) dz,$$

but $L(\theta|y) = f(x, z|\theta)$ may be straightforward to optimize.

The algorithm is as follows. Let θ_t be the current value of θ . Then define

$$Q(\theta|\theta_t) = E(\log L(\theta|Y)|x, \theta_t)$$

The algorithm is

1. E step: Compute $Q(\theta|\theta_t)$, ideally calculating the expectation over the missing data in closed form
2. M step: Maximize $Q(\theta|\theta_t)$ with respect to θ , finding θ_{t+1} .
3. Continue until convergence.

Ideally both the E and M steps can be done analytically. When the M step cannot be done analytically, one can employ some of the numerical optimization tools we've already seen. When the E step cannot be done analytically, one standard approach is to estimate the expectation by Monte Carlo, which produces Monte Carlo EM (MCEM). The strategy is to draw from z_j from $f(z|x, \theta_t)$

and approximate Q as a Monte Carlo average of $\log f(x, z_j | \theta)$, and then optimize over this approximation to the expectation. If one can't draw in closed form from the conditional density, one strategy is to do a short MCMC to draw a (correlated) sample.

EM can be shown to increase the value of the function at each step using Jensen's inequality (equivalent to the information inequality that holds with regard to the Kullback-Leibler divergence between two distributions) (Givens and Hoeting, p. 95, go through the details). Furthermore, one can show that it amounts, at each step, to maximizing a minorizing function for $\log L(\theta)$ - the minorizing function (effectively Q) is tangent to $\log L(\theta)$ at θ_t and lies below $\log L(\theta)$.

A standard example is a mixture model. Suppose we have

$$f(x) = \sum_{k=1}^K \pi_k f_k(x; \phi_k)$$

where we have K mixture components and π_k are the (marginal) probabilities of being in each component. The complete parameter vector is $\theta = \{\{\pi_k\}, \{\phi_k\}\}$. Note that the likelihood is a complicated product (over observations) over the sum (over components), so maximization may be difficult. Furthermore, such likelihoods are well-known to be multimodal because of label switching.

To use EM, we take the group membership indicators for each observation as the missing data. For the i th observation, we have $z_i \in \{1, 2, \dots, K\}$. Introducing these indicators "breaks the mixture". If we know the memberships for all the observations, it's often easy to estimate the parameters for each group based on the observations from that group. For example if the $\{f_k\}$'s were normal densities, then we can estimate the mean and variance of each normal density using the sample mean and sample variance of the x_i 's that belong to each mixture component. EM will give us a variation on this that uses "soft" (i.e., probabilistic) weighting.

The complete log likelihood given z and x is

$$\log \prod_i f(x_i | z_i, \theta) \Pr(Z_i = z_i | \pi)$$

which can be expressed as

$$\log L(x, z | \theta) = \sum_i \sum_k I(z_i = k) (\log f_k(x_i | \phi_k) + \log \pi_k)$$

with Q equal to

$$Q(\theta | \theta_t) = \sum_i \sum_k E(I(z_i = k) | x_i, \theta_t) (\log f_k(y_i | \phi_k) + \log \pi_k)$$

where $E(I(z_i = k)|x_i, \theta_t)$ is equal to the probability that the i th observation is in the k th group given x_i and θ_t , which is calculated from Bayes theorem as

$$p_{ikt} = \frac{\pi_{k,t} f_k(x_i|\theta_t)}{\sum_j \pi_{j,t} f_j(x_i|\theta_t)}$$

We can now separately maximize $Q(\theta|\theta_t)$ with respect to π_k and ϕ_k to find $\pi_{k,t+1}$ and $\phi_{k,t+1}$, since the expression is the sum of a term involving the parameters of the distributions and a term involving the mixture probabilities. In the latter case, if the f_k are normal distributions, you end up with a weighted sum of normal distributions, for which the estimators of the mean and variance parameters are the weighted mean of the observations and the weighted variance.

9 Optimization under constraints

Constrained optimization is harder than unconstrained, and inequality constraints harder to deal with than equality constraints.

Constrained optimization can sometimes be avoided by reparameterizing. E.g., to optimize w.r.t. a variance component or other non-negative parameter, you can work on the log scale.

Optimization under constraints often goes under the name of 'programming', with different types of programming for different types of objective functions combined with different types of constraints.

9.1 Convex optimization (convex programming)

Convex programming minimizes $f(x)$ s.t. $h_i(x) \leq 0$, $i = 1, \dots, m$ and $a_i^\top x = b_i$, $i = 1, \dots, q$, where both f and the constraint functions are convex. Note that this includes more general equality constraints, as we can write $g(x) = b$ as two inequalities $g(x) \leq b$ and $g(x) \geq b$. It also includes $h_i(x) \geq b_i$ by taking $-h_i(x)$. Also note that we can always have $h_i(x) \leq b_i$ and convert to the above form by subtracting b_i from each side (note that this preserves convexity). A vector x is said to be feasible, or in the feasible set if all the constraints are satisfied for x .

There are good algorithms for convex programming, and it's possible to find solutions when we have hundreds or thousands of variables and constraints. It is often difficult to recognize if one has a convex program (i.e., if f and the constraint functions are convex), but there are many tricks to transform a problem into a convex program and many problems can be solved through convex programming. So the basic challenge is in recognizing or transforming a problem to one of convex optimization; once you've done that, you can rely on existing methods to find the solution.

Linear programming, quadratic programming, second order cone programming and semidefinite programming are all special cases of convex programming. In general, these types of optimization are progressively more computationally complex.

First let's see some of the special cases and then discuss the more general problem.

9.2 Linear programming: Linear system, linear constraints

Linear programming seeks to minimize

$$f(x) = c^\top x$$

subject to a system of m inequality constraints, $a_i^\top x \leq b_i$ for $i = 1, \dots, m$, where A is of full row rank. This can also be written in terms of generalized inequality notation, $Ax \preceq b$. There are standard algorithms for solving linear programs, including the simplex method and interior point methods.

Note that each equation in the set of equations $Ax = b$ defines a hyperplane, so each inequality in $Ax \preceq b$ defines a half-space. Minimizing a linear function (presuming that the minimum exists) must mean that we push in the correct direction towards the boundaries formed by the hyperplanes, with the solution occurring at a corner (vertex) of the solid formed by the hyperplanes. The simplex algorithm starts with a feasible solution at a corner and moves along edges in directions that improve the objective function.

9.3 General system, equality constraints

Suppose we have an objective function $f(x)$ and we have equality constraints, $Ax = b$. We can manipulate this into an unconstrained problem. The null space of A is the set of x s.t. $Ax = 0$. So if we start with a candidate x_c s.t. $Ax_c = b$ (e.g., by using the pseudo inverse, A^+b), we can form all other candidates (a candidate is an x s.t. $Ax = b$) as $x = x_c + Bz$ where B is a set of column basis functions for the null space of A and $z \in \mathbb{R}^{p-m}$. Consider $h(z) = f(x_c + Bz)$ and note that h is a function of $p - m$ rather than p inputs. Namely, we are working in a reduced dimension space with no constraints. If we assume differentiability of f , we can express $\nabla h(z) = B^\top \nabla f(x_c + Bz)$ and $H_h(z) = B^\top H_f(x_c + Bz)B$. Then we can use unconstrained methods to find the point at which $\nabla h(z) = 0$.

How do we find B ? One option is to use the $p - m$ columns of V in the SVD of A that correspond to singular values that are zero. A second option is to take the QR decomposition of A^\top . Then B is the columns of Q_2 , where these are the columns of the (non-skinny) Q matrix corresponding to the rows of R that are zero.

For more general (nonlinear) equality constraints, $g_i(x) = b_i$, $i = 1, \dots, q$, we can use the Lagrange multiplier approach to define a new objective function,

$$L(x, \lambda) = f(x) + \lambda^\top (g(x) - b)$$

for which, if we set the derivative (with respect to both x and the Lagrange multiplier vector, λ) equal to zero, we have a critical point of the original function and we respect the constraints.

An example occurs with quadratic programming, under the simplification of linear equality constraints (quadratic programming in general optimizes a quadratic function under affine inequality constraints - i.e., constraints of the form $Ax - b \preceq 0$). For example we might solve a least squares problem subject to linear equality constraints, $f(x) = \frac{1}{2}x^\top Qx + m^\top x + c$ s.t. $Ax = b$, where Q is positive semi-definite. The Lagrange multiplier approach gives the objective function

$$L(x, \lambda) = \frac{1}{2}x^\top Qx + m^\top x + c + \lambda^\top (Ax - b)$$

and differentiating gives the equations

$$\begin{aligned} \frac{\partial L(x, \lambda)}{\partial x} &= m + Qx + A^\top \lambda = 0 \\ \frac{\partial L(x, \lambda)}{\partial \lambda} &= Ax = b, \end{aligned}$$

which leads to the solution

$$\begin{pmatrix} x \\ \lambda \end{pmatrix} = \begin{pmatrix} Q & A^\top \\ A & 0 \end{pmatrix}^{-1} \begin{pmatrix} -m \\ b \end{pmatrix} \quad (1)$$

which gives us $x^* = -Q^{-1}m + Q^{-1}A^\top(AQ^{-1}A^\top)^{-1}(AQ^{-1}m + b)$.

Under inequality constraints there are a variety of methods but we won't go into them.

9.4 KKT conditions

Karush-Kuhn-Tucker (KKT) theory provides sufficient conditions under which a constrained optimization problem has a minimum, generalizing the Lagrange multiplier approach. The Lange and Boyd books have whole sections on this topic.

Let $f(x)$ be the function we want to minimize, under constraints $g_i(x) = 0$; $i = 1, \dots, q$ and $h_j(x) \leq 0$; $j = 1, \dots, m$. Here I've explicitly written out the equality constraints to follow the notation in Lange.

First we need a definition. A *tangent direction*, w , is a vector for which $\nabla g_i(x)^\top w = 0$. If we

are at a point, x^* , at which the constraint is satisfied, $g_i(x^*) = 0$, then we can move in the tangent direction (orthogonal to the gradient of the constraint function) (i.e., along the level curve) and still satisfy the constraint. This is the only kind of movement that is legitimate (gives us a feasible solution).

KKT theory tells us when a candidate x^* is a local minimum. Suppose that the function and the constraint functions are continuously differentiable near x^* and that we have the Lagrangian

$$L(x, \lambda, \mu) = f(x) + \sum_i \lambda_i g_i(x) + \sum_j \mu_j h_j(x)$$

If the gradient of the Lagrangian is equal to 0,

$$\nabla f(x^*) + \sum_i \lambda_i \nabla g_i(x^*) + \sum_j \mu_j \nabla h_j(x^*) = 0,$$

and if $w^\top H_L(x^*, \lambda, \mu)w > 0$ for all vectors w s.t. $\nabla g(x^*)^\top w = 0$ and, for all active constraints, $\nabla h(x^*)^\top w = 0$, then x^* is a local minimum. An active constraint is an inequality for which $h_j(x^*) = 0$ (rather than $h_j(x^*) < 0$, in which case it is inactive). Basically we only need to worry about the inequality constraints when we are on the boundary, so the goal is to keep the constraints inactive. Note that the KKT theory doesn't require convexity.

Some basic intuition is that we need positive definiteness only for directions that stay in the feasible region. That is, our only possible directions of movement (the tangent directions) keep us in the feasible region, and for these directions, we need the objective function to be increasing to have a minimum. If we were to move in a direction that goes outside the feasible region, it's ok for the Hessian to be negative.

One can also state the KKT conditions without requiring two derivatives of the Lagrangian (see the Boyd reference, Section 5.5).

9.5 The dual problem

Sometimes a reformulation of the problem eases the optimization. There are different kinds of dual problems, but we'll just deal with the Lagrangian dual. Basically one can convert the original optimization

$$L(x, \lambda, \mu) = f(x) + \sum_i \lambda_i g_i(x) + \sum_j \mu_j h_j(x)$$

which corresponds to constrained optimization with complicated constraints $g_i(x) = 0$ and $h_j(x) \leq 0$ into the dual problem of maximizing

$$d(\lambda, \mu) = \inf_x L(x, \lambda, \mu),$$

s.t. $\mu \succeq 0$. Provided we can find the infimum over x in closed form we then maximize $d(\lambda, \mu)$ w.r.t. the Lagrangian multipliers in a new constrained problem that is sometimes easier to solve.

Here's a simple example: suppose we want to minimize $x^\top x$ s.t. $Ax = b$. The Lagrangian is $L(x, \lambda) = x^\top x + \lambda^\top (Ax - b)$. Since $L(x, \lambda)$ is quadratic in x , the infimum is found by setting $\nabla_x L(x, \lambda) = 2x + A^\top \lambda = 0$, yielding $x = -\frac{1}{2}A^\top \lambda$. So the dual function is obtained by plugging this value of x into $L(x, \lambda)$, which gives

$$g(\lambda) = -\frac{1}{4}\lambda^\top AA^\top \lambda - b^\top \lambda,$$

which is concave quadratic. In this case we can solve the original constrained problem in terms of this unconstrained dual problem.

If the optima of the primal (original) problem and that of the dual do not coincide, there is said to be a “duality gap”. For convex programming, if certain conditions are satisfied (called *constraint qualifications*), then there is no duality gap, and one can solve the dual problem to solve the primal problem. Usually with the standard form of convex programming, there is no duality gap.

9.6 Interior-point methods

We'll briefly discuss one of the standard methods for solving a convex optimization problem. The barrier method is one type of interior-point algorithm. It turns out that Newton's method can be used to solve a constrained optimization problem, with twice-differentiable f and linear equality constraints. So the basic strategy of the barrier method is to turn the more complicated constraint problem into one with only linear equality constraints.

Recall our previous notation, in which convex programming minimizes $f(x)$ s.t. $h_i(x) \leq 0$, $i = 1, \dots, m$ and $a_i^\top x = b_i$, $i = 1, \dots, q$, where both f and the constraint functions are convex. The strategy begins with moving the inequality constraints into the objective function:

$$f(x) + \sum_{i=1}^m I_-(h_i(x))$$

where $I_-(u) = 0$ if $u \leq 0$ and $I_-(u) = \infty$ if $u > 0$.

This is fine, but the new objective function is not differentiable so we can't use a Newton-like approach. Instead, we approximate the indicator function with a logarithmic function, giving the

new objective function

$$\tilde{f}(x) = f(x) + \sum_{i=1}^m -(1/t) \log(-h_i(x)),$$

which is convex and differentiable. The new term pushes down the value of the overall objective function when x approaches the boundary, nearing points for which the inequality constraints are not met. The $-\sum(1/t) \log(-h_i(x))$ term is called the log barrier, since it keeps the solution in the feasible set (i.e., the set where the inequality constraints are satisfied), provided we start at a point in the feasible set. Newton's method with equality constraints ($Ax = b$) is then applied. The key thing is then to have t get larger as the iterations proceed, which allows the solution to get closer to the boundary if that is indeed where the minimum lies.

The basic ideas behind Newton's method with equality constraints are (1) start at a feasible point, x_0 , such that $Ax_0 = b$, and (2) make sure that each step is in a feasible direction, $A(x_{t+1} - x_t) = 0$. To make sure the step is in a feasible direction we have to solve a linear system similar to that in the simplified quadratic programming problem (1):

$$\begin{pmatrix} x_{t+1} - x_t \\ \lambda \end{pmatrix} = \begin{pmatrix} H_{\tilde{f}}(x) & A^\top \\ A & 0 \end{pmatrix}^{-1} \begin{pmatrix} -\nabla \tilde{f}(x) \\ 0 \end{pmatrix},$$

which shouldn't be surprising since the whole idea of Newton's method is to substitute a quadratic approximation for the actual objective function.

10 Summary

The different methods of optimization have different advantages and disadvantages.

According to Lange, MM and EM are numerically stable and computationally simple but can converge very slowly. Newton's method shows very fast convergence but has the downsides we've discussed. Quasi-Newton methods fall in between. Convex optimization generally comes up when optimizing under constraints.

One caution about optimizing under constraints is that you just get a point estimate; quantifying uncertainty in your estimator is more difficult. One strategy is to ignore the inactive inequality constraints and reparameterize (based on the active equality constraints) to get an unconstrained problem in a lower-dimensional space.

Numerical Integration and Differentiation

November 28, 2011

References:

- Gentle: Computational Statistics
- Monahan: Numerical Methods of Statistics
- Givens and Hoeting: Computational Statistics

Our goal here is to understand the basics of numerical (and symbolic) approaches to approximating derivatives and integrals on a computer. Derivatives are useful primarily for optimization. Integrals arise in approximating expected values and in various places where we need to integrate over an unknown random variable (e.g., Bayesian contexts, random effects models, missing data contexts).

1 Differentiation

1.1 Numerical differentiation

There's not much to this topic. The basic idea is to approximate the derivative of interest using finite differences.

A standard discrete approximation of the derivative is the forward difference

$$f'(x) \approx \frac{f(x+h) - f(x)}{h}$$

A more accurate approach is the central difference

$$f'(x) \approx \frac{f(x+h) - f(x-h)}{2h}$$

Provided we already have computed $f(x)$, the forward difference takes half as much computing as the central difference. However, the central difference has an error of $O(h^2)$ while the forward difference has error of $O(h)$.

For second derivatives, if we apply the above approximations to $f'(x)$ and $f'(x+h)$, we get an approximation of the second derivative based on second differences:

$$f''(x) \approx \frac{f'(x+h) - f'(x)}{h} \approx \frac{f(x+2h) - 2f(x+h) + f(x)}{h^2}.$$

The corresponding central difference approximation is

$$f''(x) \approx \frac{f(x+h) - 2f(x) + f(x-h)}{h^2}.$$

For multivariate x , we need to compute directional derivatives. In general these will be in axis-oriented directions (e.g., for the Hessian), but they can be in other directions. The basic idea is to find $f(x+he)$ in expressions such as those above where e is a unit length vector giving the direction. For axis oriented directions, we have e_i being a vector with a one in the i th position and zeroes in the other positions,

$$\frac{\partial f}{\partial x_i} \approx \frac{f(x+he_i) - f(x)}{h}.$$

Note that for mixed partial derivatives, we need to use e_i and e_j , so the second difference approximation gets a bit more complicated,

$$\frac{\partial^2 f}{\partial x_i \partial x_j} \approx \frac{f(x+he_j+he_i) - f(x+he_j) - f(x+he_i) + f(x)}{h^2}.$$

We would have analogous quantities for central difference approximations.

Numerical issues Ideally we would take h very small and get a highly accurate estimate of the derivative. However, the limits of machine precision mean that the difference estimator can behave badly for very small h , since we lose accuracy in computing $x+h$ and $x-h$ and from dividing by small h . Therefore we accept a bias in the estimate by not using h so small, often by taking h to be square root of machine epsilon (i.e., about 1×10^{-8} on most systems). Actually, we need to account for the order of magnitude of x , so what we really want is $h = \sqrt{\epsilon}|x|$ - i.e., we want it to be in terms relative to the magnitude of x . As an example, recall that if $x = 1 \times 10^9$ and we did $x+h = 1 \times 10^9 + 1 \times 10^{-8}$, we would get $x+h = 1 \times 10^9 = x$ because we can only represent 7 decimal places with precision.

Givens and Hoeting and Monahan point out that some sources recommend the cube root of machine epsilon (about 5×10^{-6} on most systems), in particular when approximating second

derivatives.

Let's assess these recommendations empirically in R using the demo code.

1.2 Numerical differentiation in R

There are multiple numerical derivative functions in R. *numericDeriv()* will do the first derivative. It requires an expression rather than a function as the form in which the function is input, which in some cases might be inconvenient. The functions in the *numDeriv* package will compute the gradient and Hessian, either in the standard way (using the argument `method = 'simple'`) or with a more accurate approximation (using the argument `method = 'Richardson'`). For optimization, one might use the simple option, assuming that is faster, while the more accurate approximation might be good for computing the Hessian to approximate the information matrix for getting an asymptotic covariance.

Note that by default, if you rely on numerical derivatives in *optim()*, it uses $h = 0.001$ (the *ndeps* sub-argument to *control*), which might not be appropriate if the parameters vary on a small scale. This relatively large value of h is probably chosen based on *optim()* assuming that you've scaled the parameters as described in the text describing the *parscale* argument.

1.3 Symbolic differentiation

We've seen that we often need the first and second derivatives for optimization. Numerical differentiation is fine, but if we can readily compute the derivatives in closed form, that can improve our optimization. (Venables and Ripley comment that this is particularly the case for the first derivative, but not as much for the second.)

In general, using a computer program to do the analytic differentiation is recommended as it's easy to make errors in doing differentiation by hand. Monahan points out that one of the main causes of error in optimization is human error in coding analytic derivatives, so it's good practice to avoid this. R has a simple differentiation ability in the *deriv()* function (which handles the gradient and the Hessian). However it can only handle a limited number of functions.

For more complicated functions, both Maple and Mathematica do symbolic differentiation. In the demo code, we'll consider an example in Mathematica. At the moment, the Mathematica license on the SCF machines is not available, but this should be resolved shortly.

2 Integration

We've actually already discussed numerical integration extensively in the simulation unit, where we considered Monte Carlo approximation of high-dimensional integrals. In the case where we

have an integral in just one or two dimensions, MC is fine, but we can get highly-accurate, very fast approximations by numerical integration methods known as quadrature. Unfortunately such approximations scale very badly as the dimension grows, while MC methods scale well, so MC is recommended in higher dimensions. More on this below.

2.1 Numerical integration methods

The basic idea is to break the domain into pieces and approximate the integral within each piece:

$$\int_a^b f(x)dx = \sum_{i=0}^{n-1} \int_{x_i}^{x_{i+1}} f(x)dx,$$

where we then approximate $\int_{x_i}^{x_{i+1}} f(x)dx \approx \sum_{j=0}^m A_{ij}f(x_{ij}^*)$ where x_{ij}^* are the *nodes*.

2.1.1 Newton-Cotes quadrature

Newton-Cotes quadrature has equal length intervals of length $h = (b-a)/n$, with the same number of nodes in each interval. $f(x)$ is replaced with a polynomial approximation in each interval and A_{ij} are chosen so that the sum equals the integral of the polynomial approximation on the interval.

A basic example is the *Riemann rule*, which takes a single node, $x_i^* = x_i$ and the “polynomial” is a constant, $f(x_i^*)$, so we have

$$\int_{x_i}^{x_{i+1}} f(x)dx \approx (x_{i+1} - x_i)f(x_i).$$

Of course using a piecewise constant to approximate $f(x)$ is not likely to give us high accuracy.

The *trapezoidal rule* takes $x_{i0}^* = x_i$, $x_{i1}^* = x_{i+1}$ and uses a linear interpolation between $f(x_{i0}^*)$ and $f(x_{i1}^*)$ to give

$$\int_{x_i}^{x_{i+1}} f(x)dx \approx \left(\frac{x_{i+1} - x_i}{2}\right) (f(x_i) + f(x_{i+1})).$$

Simpson’s rule uses a quadratic interpolation at the points $x_{i0}^* = x_i$, $x_{i1}^* = (x_i + x_{i+1})/2$, $x_{i2}^* = x_{i+1}$ to give

$$\int_{x_i}^{x_{i+1}} f(x)dx \approx \left(\frac{x_{i+1} - x_i}{6}\right) \left(f(x_i) + 4f\left(\frac{x_i + x_{i+1}}{2}\right) + f(x_{i+1})\right).$$

The error of various rules is often quantified as a power of $h = x_{i+1} - x_i$. The trapezoid rule gives $O(h^2)$ while Simpson’s rule gives $O(h^4)$.

Romberg quadrature There is an extension of Newton-Cotes quadrature that takes combinations of estimates based on different numbers of intervals. This is called Richardson extrapolation.

tion and when used with the trapezoidal rule is called *Romberg quadrature*. The result is greatly increased accuracy. A simple example of this is as follows. Let $\hat{T}(h)$ be the trapezoidal rule approximation of the integral when the length of each interval is h . Then $\frac{4\hat{T}(h/2) - \hat{T}(h)}{3}$ results in an approximation with error of $O(h^4)$ because the differencing is cleverly chosen to kill off the error term that is $O(h^2)$. In fact this approximation is Simpson's rule with intervals of length $h/2$, with the advantage that we don't have to do as many function evaluations ($2n$ vs. $4n$). Even better, one can iterate this approach for more accuracy as described in detail in Givens and Hoeting.

Note that at some point, simply making intervals smaller in quadrature will not improve accuracy because of errors introduced by the imprecision of computer numbers.

2.1.2 Gaussian quadrature

Here the idea is to relax the constraints of equally-spaced intervals and nodes within intervals. We want to put more nodes where the function is larger in magnitude.

Gaussian quadrature approximates integrals that are in the form of an expected value as

$$\int_a^b f(x)\mu(x)dx \approx \sum_{i=0}^m w_i f(x_i)$$

where $\mu(x)$ is a probability density, with the requirement that $\int x^k \mu(x)dx = E_\mu X^k < \infty$ for $k \geq 0$. Note that it can also deal with indefinite integrals where $a = -\infty$ and/or $b = \infty$. Typically μ is non-uniform, so the nodes (the quadrature points) cluster in areas of high density. The choice of node locations depends on understanding orthogonal polynomials, which we won't go into here.

It turns out this approach can exactly integrate polynomials of degree $2m + 1$ (or lower). The advantage is that for smooth functions that can be approximated well by a single polynomial, we get highly accurate results. The downside is that if the function is not well approximated by such a polynomial, the result may not be so good. The Romberg approach is more robust.

Note that if the problem is not in the form $\int_a^b f(x)\mu(x)dx$, but rather $\int_a^b f(x)dx$, we can reexpress as $\int_a^b \frac{f(x)}{\mu(x)}\mu(x)dx$.

Note that the trapezoidal rule amounts to μ being the uniform distribution with the points equally spaced.

2.1.3 Adaptive quadrature

Adaptive quadrature chooses interval lengths based on the behavior of the integrand. The goal is to have shorter intervals where the function varies more and longer intervals where it varies less. The reason for avoiding short intervals everywhere involves the extra computation and greater opportunity for rounding error.

2.1.4 Higher dimensions

For rectangular regions, one can use the techniques described above over squares instead of intervals, but things become more difficult with more complicated regions of integration.

The basic result for Monte Carlo integration (i.e., Unit 11 on simulation) is that the error of the MC estimator scales as $O(m^{-1/2})$, where m is the number of MC samples, regardless of dimensionality. Let's consider how the error of quadrature scales. We've seen that the error is often quantified as $O(h^q)$. In d dimensions, the error is the same as a function of h , but if in one dimension we need n function evaluations to get intervals of length h , in d dimensions, we need n^d function evaluations to get hypercubes with sides of length h . Let's re-express the error in terms of n rather than h based on $h = c/n$ for a constant c (such as $c = b - a$), which gives us error of $O(n^{-q})$ for one-dimensional integration. In d dimensions we have $n^{1/d}$ function evaluations per dimension, so the error for fixed n is $O((n^{1/d})^{-q}) = O(n^{-q/d})$ which scales as $n^{-1/d}$. As an example, suppose $d = 10$ and we have $n = 1000$ function evaluations. This gives us an accuracy comparable to one-dimensional integration with $n = 1000^{1/10} \approx 2$, which is awful. Even with only $d = 4$, we get $n = 1000^{1/4} \approx 6$, which is pretty bad. This is one version of the curse of dimensionality.

2.2 Numerical integration in R

R implements an adaptive version of Gaussian quadrature in `integrate()`. The '...' argument allows you to pass additional arguments to the function that is being integrated. The function must be vectorized (i.e., accept a vector of inputs and evaluate and return the function value for each input as a vector of outputs).

Note that the domain of integration can be unbounded and if either the upper or lower limit is unbounded, you should enter `Inf` or `-Inf` respectively.

2.3 Singularities and infinite ranges

A singularity occurs when the function is unbounded, which can cause difficulties with numerical integration. For example, $\int_0^1 \frac{1}{\sqrt{x}} dx = 2$, but $f(0) = \infty$. One strategy is a change of variables. For example, to find $\int_0^1 \frac{\exp(x)}{\sqrt{x}} dx$, let $u = \sqrt{x}$, which gives the integral, $2 \int_0^1 \exp(u^2) du$.

Another strategy is to subtract off the singularity. E.g., in the example above, reexpress as

$$\int_0^1 \frac{\exp(x) - 1}{\sqrt{x}} dx + \int_0^1 \frac{1}{\sqrt{x}} dx = \int_0^1 \frac{\exp(x) - 1}{\sqrt{x}} dx + 2$$

where we do the second integral analytically. It turns out that the first integral is well-behaved at

0.

It turns out that R's *integrate()* function can handle $\int_0^1 \frac{\exp(x)}{\sqrt{x}} dx$ directly without us changing the problem statement analytically. Perhaps this has something to do with the use of adaptive quadrature, but I'm not sure.

Infinite ranges Gaussian quadrature deals with the case that $a = -\infty$ and/or $b = \infty$. Another possibility is change of variables using transformations such as $1/x$, $\exp(x)/(1 + \exp(x))$, $\exp(-x)$, and $x/(1 + x)$.

2.4 Symbolic integration

Mathematica and Maple are able to do symbolic integration for many problems that are very hard to do manually. So this may be worth a try.

Parallel Processing

December 5, 2011

Reference:

Schmidberger M, Morgan M, Eddelbuettel D, Yu H, Tierney L, Mansmann U (2009). "State of the Art in Parallel Computing with R." *Journal of Statistical Software*, 31(1), 127.

<http://www.jstatsoft.org/v31/i01>

1 Computer architecture

Computers now come with multiple processors for doing computation. Basically, physical constraints have made it harder to keep increasing the speed of individual processors, so the industry is now putting multiple processing units in a given computer and trying/hoping to rely on implementing computations in a way that takes advantage of the multiple processors.

Everyday personal computers often have more than one processor (more than one chip) and on a given processor, often have more than one core (multi-core). A multi-core processor has multiple processors on a single computer chip. On personal computers, all the processors and cores share the same memory.

Supercomputers and computer clusters generally have tens, hundreds, or thousands of 'nodes', linked by a fast local network. Each node is essentially a computer with its own processor(s) and memory. Memory is local to each node (distributed memory). One basic principle is that communication between a processor and its memory is much faster than communication between processors with different memory. An example of a modern supercomputer is the Jaguar supercomputer at Oak Ridge National Lab, which has 18,688 nodes, each with two processors and each processor with 6 cores, giving 224,256 total processing cores. Each node has 16 Gb of memory for a total of 300 Tb.

There is little practical distinction between multi-processor and multi-core situations. The main issue is whether processes share memory or not. In general, I won't distinguish between cores and processors, but rather between shared memory and distributed memory.

1.1 Graphics processing units (GPUs)

GPUs were formerly a special piece of hardware used by gamers and the like for quickly rendering (displaying) graphics on a computer. They do this by having hundreds of processing units and breaking up the computations involved in an embarrassingly parallel fashion (i.e., without inter-processor communication). For particular tasks, GPUs are very fast. Nowadays GPUs are generally built onto PC motherboards whereas previously they were on a video card.

Researchers (including some statisticians) are increasingly looking into using GPUs to do massively parallel computations with the inexpensive hardware available on standard computers. This has promise. However, some drawbacks in current implementations include the need to learn a specific programming language (similar to C) and limitations in terms of transferring data to the GPU and holding information in the memory of the GPU.

1.2 Cloud computing

Amazon (through its EC2 service) and other companies offer computing through the cloud. The basic idea is that they rent out their servers on a pay-as-you-go basis. You get access to a virtual machine that can run various versions of Linux or Microsoft Windows server and where you choose the number of processing cores you want. You configure the virtual machine with the applications, libraries, and data you need and then treat the virtual machine as if it were a physical machine that you log into as usual.

2 Parallelization

2.1 Overview

A lot of parallel processing follows a master-slave paradigm. There is one master process that controls one or more slave processes. The master process sends out tasks and data to the slave processes and collects results back.

One comment about parallelized code is that it can be difficult to debug because communication problems can occur in addition to standard bugs. Debugging may require some understanding of the communication that goes on between processors. If you can first debug your code on one processor, that can be helpful.

2.2 Threading

One form of parallel processing is threading. Here an algorithm is implemented across multiple processes in a shared memory situation.

BLAS libraries come in threaded versions so that if you call a function in a multi-processor environment, they do the calculation in parallel behind the scenes. We saw a demo in the simulation unit where on my Mac laptop, a task used more than 100% of CPU because it was a threaded calculation distributed across two cores. A standard protocol for parallelization in a shared memory environment is called *openMP*.

To make use of threading, you need to have your job have control of multiple processors. In UNIX, the `OMP_NUM_THREADS` environment variable controls how many threads are being used (provided multiple processors are available). E.g., in bash

```
> export OMP_NUM_THREADS=2
```

On Macs, the threaded BLAS is generally implemented by default, so you don't need to do anything.

In R, the basic strategy is to make sure that the R installation uses a threaded BLAS so that standard linear algebra computations are done in a threaded fashion. We can do this by linking R to a threaded BLAS shared object library (i.e., a `.so` file). Details can be found in Section A.3.1.5 of the [R administration manual](#) or talk to your system administrator.

Note that in R, the threading only helps with linear algebra operations. In contrast, I believe that Matlab using threading more generally to speed up other types of calculations.

2.3 Embarrassingly parallel (EP)

An EP problem is one that can be solved by doing independent computations as separate processes without communication between the processes. You can get the answer by doing separate tasks and then collecting the results. Examples in statistics include

1. simulations with many independent replicates
2. bootstrapping
3. stratified analyses

The standard setup is that we have the same code running on different datasets. (Note that different processes may need different random number streams, as we discussed in the Simulation Unit.)

To do parallel processing in this context, you need to have control of multiple processes. In a system with a queueing setup, this will generally mean requesting access to a certain number of processors and then running your job in such a way that you use multiple processors.

In the demo code, we'll see the use of *foreach()* (from the *foreach* package,) which allows us to implement a *for* loop spread across multiple processors. *foreach()* can distribute the computations across multiple cores or multiple processors using several different back-end parallel processing implementations, including use of the *doMC* package and the *Rmpi* package. *doMC* handles multiple cores or processors with shared memory, while *Rmpi* handles distributed (or shared) memory situations via message passing (see Section 2.4). Let's see an example in the demo code. Note that *foreach()* has some nice functionality for collecting results back in a clean format, but our demonstration just illustrates the bare bones usage.

The *snow* and *snowfall* packages provide a way to set up an ad hoc network of computers for EP calculations, using the *Rmpi* functionality described in Section 2.4.

In general, except for some modest overhead, an EP problem can be solved with $1/p$ the amount of time for the non-parallel implementation, given p processors. This gives us a speedup of p , which is called linear speedup (basically anytime the speedup is of the form cp for some constant p).

One difficulty is load balancing. We'd like to make sure each slave process finishes at the same time. Often we can give each process the same amount of work, but if we have a mix of faster and slower processors, things become more difficult. To the extent it is possible to break up a job into many small tasks and have processors start new tasks as they finish off old tasks, this can be effective, but may involve some parallel programming.

2.4 Parallelization with communication

If we do not have an EP problem, we have one that involves some sort of serial calculation. As a result, different processes need to communicate with each other. There are standard protocols for such communication, with some common ones being *MPI* and *PVM*. You can use C libraries that implement these protocols. While *MPI* has many functions, a core of 6-10 functions (basic functions for functionality such as sending and receiving data between processes - either master-slave or slave-slave) are what we mostly need.

R provides the *Rmpi* library, which allows you to do message passing in R. It has some drawbacks, but may be worth exploring if you have a non-EP problem and don't want to learn C. Installing *Rmpi* may be tricky and on institutional machines will require you talk to your systems administrator. *Rmpi* is a basic building block for other parallel processing functionality such as *foreach* and *SNOW*.

For non-EP problems, the primary question is how the speed of the computation scales with p . This will generally be much worse than $1/p$ and as p increases, if communication must increase as well, then the speedup can be much worse.