# Graphics in R

There are three main types of functions that deal with graphics in R

High-level functions that produce complete plots

Low-level functions or graphical primitives that can be added to an existing plot or assembled to build a new plot type

A limited set of interactive features for working with graphical output

## Graphics in R

When making graphics in R, you typical issue a series of calls to graphics functions, each producing a complete plot or adding to an existing plot

Sometimes, people refer to this as the "Painter's model" meaning you add layers to a plot in steps, with later output obscuring what came before

## Standard plots

A range of standard plots can be made with R and these are typically produced by a single function call (but can be added to)

These functions have embedded in them a number of "good choices" (both in terms of layout and design as well as any parameters that might need setting) to facilitate rapid iterations to support analysis

These, however, can also just be the starting point for more elaborate graphics, adding annotations, overlaying other plotting elements

# The BRFSS

The Behavioral Risk Factor Surveillance System is the world's largest telephone survey and it is designed to track health risks in the United States; like many surveys, **the BRFSS works with only a *sample* of a larger *population***

With over 200 million adults in the United States, the CDC couldn't possibly contact their entire population*; instead, they selected around 400 thousand adults, calling roughly 30 thousand per month

## 1. Background

The Behavioral Risk Factor Surveillance System (BRFSS) is a collaborative project of the Centers for Disease Control and Prevention (CDC), and U.S. states and territories. The BRFSS, administered and supported by the Behavioral Surveillance Branch (BSB) of the CDC, is an ongoing data collection program designed to measure behavioral risk factors in the adult population 18 years of age or over living in households. The BRFSS was initiated in 1984, with 15 states collecting surveillance data on risk behaviors through monthly telephone interviews. The number of states participating in the survey increased, so that by 2000, 50 States, the District of Columbia, Puerto Rico, and the Virgin Islands were participating in the BRFSS.

The objective of the BRFSS is to collect uniform, state-specific data on preventive health practices and risk behaviors that are linked to chronic diseases, injuries, and preventable infectious diseases in the adult population. Factors assessed by the BRFSS include tobacco use, health care coverage, HIV/AIDS, physical activity, and fruit and vegetable consumption. Data are collected from a random sample of adults (one per household) through a telephone survey.

BRFSS
**Turning Information Into Health**

## Our data

The BRFSS in 2008 consists of responses from 400 thousand people; in this discussion, we will only look at a subset (a subsample, if you will) of **40 thousand people**\*

Each respondent receiving the survey is asked a series of questions and **the original BRFSS data set has 292 different fields**, most of which are questions; to make things easier, we've only pulled 34 variables

## Variables

**state**
   Where does the respondent live?

**imonth**
   Interview Month

**iday**
   Interview Day

**iyear**
   Interview Year

**nattempts**
   Number of Attempts

**numadults**
   Number of Adults in Household

**nummen**
   Number of Adult Women in Household

**numwomen**
   Number of Adult Women in Household

## Variables

**genhlth**

Respondents were asked to evaluate their general health values are excellent, very good, good, fair, poor

**physhlth**

The number of days out of the last 30 that the respondent was in poor health

**menthlth**

The number of days out of the last 30 that the respondent was in poor mental health

**hlthplan**

1 if the respondent has some form of health coverage and 2 else

**medcost**

Was there a time in the past 12 months when you needed to see a doctor but could not because of cost?

**checkup1**

Time since the respondent's last routine checkup

## Variables

**qlrest2**

During the past 30 days, for about how many days have you felt you did not get enough rest or sleep?

**cvdinfr4**

Has the respondent ever had a heart attack? (1 yes, 2 no)

**cvdcrhd4**

Has the respondent ever had angina or coronary heart disease?

**cvdstrk3**

Has the respondent ever had a stroke?

**asthma2**

Does the respondent have asthma?

**smoke100**

1 if the respondent has smoked at least 100 cigarettes in their entire life and 2 otherwise

## Variables

**age** in years

**marital** Is the respondent married?

**children** Number of children (< 18 years old) living at the household

**educ** The highest grade or year of school the respondent completed

**employ** Is the respondent currently employed?

**income2** range

**weight** in pounds

**height** in inches

**wtyrago** desired weight in pounds

**sex** of the respondent

**drnkany**
  Has the respondent had at least one alcoholic beverage in the last 30 days?

## Variables

**lsatisfy**
How satisfied is the respondent with their life overall?

# Preliminary examination

When faced with a new data set, we often have a look at a few cases; you should do this before and after the data are "loaded" into R or whatever analysis package you might end up using

What do we notice?

```
> brfss[1:20,1:13]
           state    imonth iday iyear nattempts numadults nummen numwomen   genhlth physhlth menthlth hlthplan medcost
1        Illinois   January   12  2008         1         1      1        0 Excellent        3        0        1       2
2         Florida     March   10  2008         7         2      1        1      Fair        7        0        2       1
3        Missouri  February    6  2008         5         2      1        1 Very good        2        1        1       2
4    South Dakota     March   18  2008         4         2      1        1 Very good        0        1        1       2
5     Connecticut   October   17  2008        10         3      2        1      Good        0        0        1       2
6    Pennsylvania      July   10  2008         6         4      1        3      Good        0       30        1       1
7       Tennessee      July   21  2008        12         2      1        1 Excellent        0        0        1       2
8           Texas    August    5  2008         3         2      1        1      Poor        0        0        1       2
9   New Hampshire     April    9  2008         3         2      1        1 Very good        0        0        1       2
10        Indiana      July   13  2008         5         2      1        1 Excellent        0        0        1       2
11        Florida September   14  2008        10         4      2        2      Fair       30        0        1       2
12       Illinois September   13  2008         1         3      1        2      Fair        0        0        1       2
13         Kansas    August   28  2008         8         1      0        1 Excellent        0        0        2       2
14    Puerto Rico  February    2  2008         1         4      1        3      Fair       25        0        1       2
15        Alabama    August   18  2008         5         4      2        2 Very good        0        0        1       2
16      Louisiana      June   28  2008         6         1      1        0 Very good        0        0        1       2
17          Texas     March    6  2008         5         2      1        1      Good        0        0        1       2
18       Missouri     March   28  2008         4         2      1        1 Excellent        0        0        1       2
19      Minnesota      June   10  2008         2         3      2        1 Very good        0        0        1       2
20       Illinois  December   11  2008        16         3      0        3 Very good        3        2        1       2
```

# A look

The survey responses are a mix of qualitative and quantitative data; let's start slow with a look at a couple of the categorical variables

What is the gender breakdown?

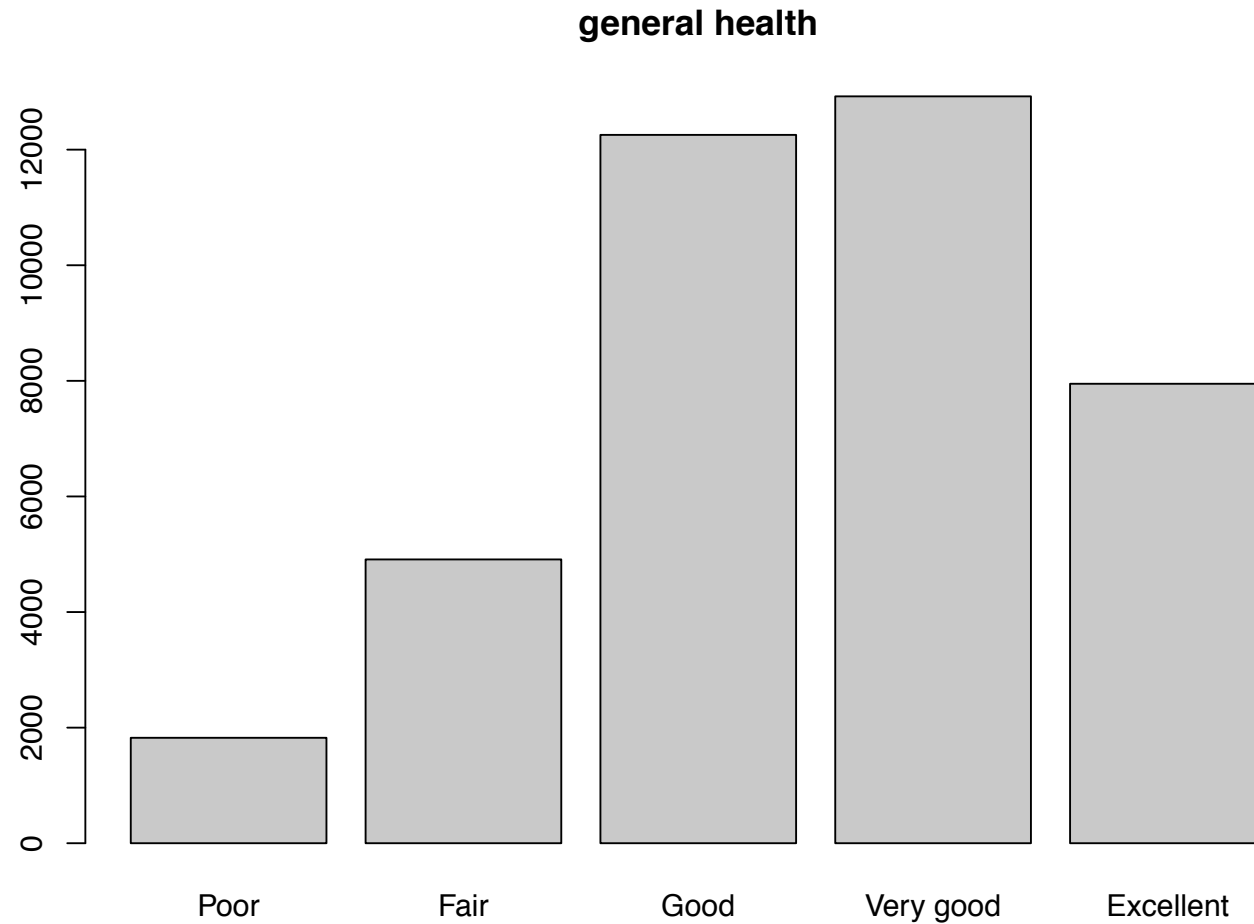What proportion of respondents have exercised in the last 30 days?

What about the respondents' general health?

Their overall satisfaction with their lives?

```
exerany

1    29726
2    10233


gender


male     19604
female   20396
```

```
genhlth

excellent    7949
very good   12922
good        12255
fair         4910
poor         1824
don't know     61
refused        79
```

```
lsatisfy

very satisfied       17590
satisfied            18806
dissatisfied         12255
very dissatisfied     4910
don't know             191
refused                364
```

# Graphical displays

A **barplot** can be formed to make comparisons easier

**general health**

```
# let's start with tabular output

table(brfss$genhlth)

# for the moment, drop the non-responders and the unsure :)

ta = table(brfss$genhlth)
ta = ta[5:1]

# a high-level plot routine

barplot(ta,main="general health")

# ... or flipped on its side (getting fancy!)

par(oma=c(0,2,0,0))
barplot(ta,main="general health",horiz=TRUE,las=1)

?barplot
```

**general health**

Graphical displays

Some have argued that comparisons are better made when the bars run horizontally*

The code on the previous slide skips ahead a fair bit, but don't worry, we'll cover all this

* Cleveland, W. S. (1993), *Visualizing Data,* Hobart Press

## Questions

While these one-dimensional summaries are interesting, they can't address certain questions we might bring to the data; for example, does exercise have any effect on what people feel about their general health?

For this, we might consider tabular displays

# Tables

Here is a two-by-two table (also referred to as **a contingency table**) describing how respondents answered both the question of how good they feel and whether or not they exercise; we have added **row and column sums** to this display also

What do we see?

```
> table(brfss$genhlth,brfss$exerany)
```

|                    | 1     | 2    |
|--------------------|-------|------|
| Excellent          | 6880  | 1065 |
| Very good          | 10561 | 2353 |
| Good               | 8649  | 3590 |
| Fair               | 2811  | 2089 |
| Poor               | 732   | 1091 |
| Don't know/Not Sure | 36    | 24   |
| Refused            | 57    | 21   |

## Mosaic plots

These displays represent the counts in a contingency table by tiles whose size (area) is proportional to the cell count

It is also possible to extend these displays to tabulations with more than two variables; how might this work?

Hartigan, J.A., and Kleiner, B. (1984) A mosaic of television  ratings. *The American Statistician*, **38**, 32-35

# general health by exerany

general health by exerany

```
# let's start with tabular output

table(brfss$genhlth,brfss$exerany)

# a high-level plot routine

ta = table(brfss$genhlth,brfss$exerany)
ta = ta[1:5,]

mosaicplot(ta,main="general health by exerany")

# colors!

mosaicplot(ta,color=c("red","blue"))
```

## Creating new variables

BMI (Body Mass Index) is defined to be

$$\text{BMI} = 703 \times \frac{\text{weight in pounds}}{(\text{height in inches})^2}$$

We can derive this from our data set and create a new quantitative variables

The CDC interprets these limits as follows

| BMI | Weight Status |
|---|---|
| Below 18.5 | Underweight |
| 18.5 – 24.9 | Normal |
| 25.0 – 29.9 | Overweight |
| 30.0 and Above | Obese |

```
bmi = 703*brfss$weight/brfss$height^2

bmicat = (bmi > 0) + (bmi >= 18.5) + (bmi >= 25) + (bmi >= 30)

levs = c("underweight","normal","overweight","obese")
bmicat = factor(levs[bmicat],levels=levs)

ta = table(brfss$genhlth,bmicat)
ta = ta[1:5,]

mosaicplot(ta,main="general health by exerany")

mosaicplot(ta,color=heat.colors(4),
           main="general health by exerany")
```

## Quantitative data

In the BMI example, we generated an ordinal variable from our computed BMI; in general, "seeing" the values of a quantitative variable (whether it be continuous or discrete with a large number of values) can be hard

But the idea of grouping comes to our rescue in the form of **grouped frequency displays or histograms**

# Histograms

A histogram groups or bins the data and, like a barplot, presents the number of data points that fall into each group

This display involves a "tuning parameter"; that is, we are free to choose how many bins we want to make the display -- this is what I meant before by there being both tuning in terms of the aesthetics (colors, fonts) as well as the underlying methodology that generates the display

In situations like this, **it is always good to vary the number of bins and examine the plot for any structure that emerge**s; in so doing, we want to get a sense of the "shape" of the data

What do we see?

**weight of respondents**

**weight of respondents**

**weight of respondents**

**weight of respondents**

Frequency

brfss$weight

# weight of respondents



Frequency vs. brfss$weight

# Varying bin sizes

By changing the bin size, we can uncover features in the data; in this case we uncover a basic fact about how people report their weights

**weight of respondents**

```
hist(brfss$weight,breaks=20,main="weight of respondents")

hist(brfss$weight,breaks=50,main="weight of respondents")

hist(brfss$weight,breaks=100,main="weight of respondents")

hist(brfss$weight,breaks=500,main="weight of respondents")

hist(brfss$weight,breaks=500,main="weight of respondents",xlim=c(100,150))
```

## Default bin size

It is often the case that we don't want to think very hard about how many bins or groups to use when drawing a histogram; the hist() function in R uses a rule of thumb for setting the number of bins based on our sample size

$$\text{number of bins} \approx \log_2(n) + 1$$

Where might a rule like this come from?

# Comparing distributions (I)

We can use these displays to compare distributions

At the left we have separate histograms of the heights of males and females in the sample

What do you see?



**height of male respondents**



**height of female respondents**

```
# compare two histograms...

# first, a common range...

ra = range(brfss$height,na.rm=T)
ra

# and now, the high-level command...

hist(brfss$height[brfss$sex=="male"],
     breaks=50,main="height of male respondents",xlim=ra,xlab="heights")

hist(brfss$height[brfss$sex=="female"],
     breaks=50,main="height of female respondents",xlim=ra,xlab="heights")
```

# Comparing distributions (I)

A more effective strategy would be to simply overlay one histogram over the other, perhaps adding a snappy color

At this point it should be clear how helpful it is to have a good rule of thumb for picking the number of bins



heights of respondents (cyan/male, white/female)

```
# compare two histograms...

# first, a common range...

ra = range(brfss$height,na.rm=T)
ra

# and now, the high-level command...

hist(brfss$height[brfss$sex=="male"],
     breaks=50,main="height of male respondents",xlim=ra,xlab="heights")

hist(brfss$height[brfss$sex=="female"],
     breaks=50,main="height of female respondents",xlim=ra,xlab="heights")

# overlay!

hist(brfss$height[brfss$sex=="female"],breaks=50,
 main="heights of respondents (cyan/male, white/female)",xlim=ra,xlab="heights")

hist(brfss$height[brfss$sex=="male"],breaks=50,add=T,col="cyan")
```

# Smoothed histograms

Overlaying histograms can get tricky if we aren't careful; one can obscure the features of the other -- here we're lucky in that both distributions are essentially unimodal

Another approach, however, is to create a simpler view; a smoothed histogram (technically, a kernel density estimate) is one such device

**female heights, smoothed histogram**

Density

N = 19963   Bandwidth = 0.355

```
fsmooth = density(brfss$height[brfss$sex=="female"],na.rm=T)

plot(fsmooth,main="female heights, smoothed histogram")

fsmooth = density(brfss$height[brfss$sex=="female"],na.rm=T,adjust=2)
lines(fsmooth,col="cyan")

fsmooth = density(brfss$height[brfss$sex=="female"],na.rm=T,adjust=5)
lines(fsmooth,col="magenta")

# overlay!

fsmooth = density(brfss$height[brfss$sex=="female"],na.rm=T,adjust=1.5)
msmooth = density(brfss$height[brfss$sex=="male"],na.rm=T,adjust=1.5)

plot(fsmooth,main="respondent heights, smoothed histograms",
     xlab="respondent heights")

lines(msmooth,col="cyan")

# add a legend!

legend(40,0.04,col=c("black","cyan"),legend=c("female","male"),lty=c(1,1))
```

**respondent heights, smoothed histograms**

## A graphical measure

Often, we find ourselves asking if the distribution of data in question is normal or not; statisticians in the late 1800s were obsessed with finding normal curves in groups of body measurements

Histograms are one way to assess normality (does it look bell-shaped or not?) but a qqplot is a more refined measuring device...

**normal qq plot, male heights**

Sample Quantiles

Theoretical Quantiles

```
# normal quantile-quantile plot

qqnorm(brfss$height[brfss$sex=="male"],main="normal qq plot, male heights")

# add a guide line

qqline(brfss$height[brfss$sex=="male"])
```

## Boxplots

In many cases, we don't need to examine the complete distribution, but we can instead look at just a thumbnail sketch -- we're sort of creeping up on that idea as we simplify the smooth histograms

Boxplots are one form of thumbnail, focusing on the so-called five number summary; they were developed by one of the big thinkers in EDA, John Tukey

These plots let us relate a categorical and a continuous variable...

**respondent heights**

```
boxplot(brfss$height[brfss$sex=="female"],brfss$height[brfss$sex=="male"])

# another way to generate the same plot...

boxplot(height~sex,data=brfss, main="respondent heights")

# ... and yet another way

plot(height~sex,data=brfss,main="respondent heights")

# and another...

plot(bmi~genhlth,data=brfss,main="bmi by general health")
```

**bmi by general health**

# Graphics in R

We have seen some high-level plots; box plots, histograms, smoothed histograms, bar plots and mosaic plots

These were all called by special functions that execute one kind of graphical display; we started to see, however, that the function plot() itself, was a fairly flexible character -- we'll come back to that shortly

Now, these are by no means the only high-level functions out there; people are actively contributing all manner of interesting high-level, specialty plots

# Violin plots

The so-called violin plot might be more artistry than data analysis; but it uses the smoothed histogram tipped on its side and mirrored left and right in place of a box

Compare this plot to the boxplot three slides back; what do you think?



violin plot bmi given genhlth

# Weighty matters

To relate two continuous variables, we could consider a scatterplot (by rights, in any sane graphics introduction, this would come first)

Let's look at people's weights this year to their weights last year...

**weights then and now**

**weights then and now**

```r
# the many faces of plot!

plot(brfss$weight,brfss$wtyrago,main="weights then and now")

# or...

plot(weight~wtyrago,data=brfss,main="weights then and now")

# add another guide line...

abline(0,1)

abline(0,2)    # people who are twice as heavy

abline(0,0.5)    # people who are twice as heavy
```

```
library(hexbin)

# create a hexagonal grid over the data and count the points falling
# in each cell

h = hexbin(brfss$weight,brfss$wtyrago)

plot(h,main="weights then and now")

# look again what plot's doing!!
```
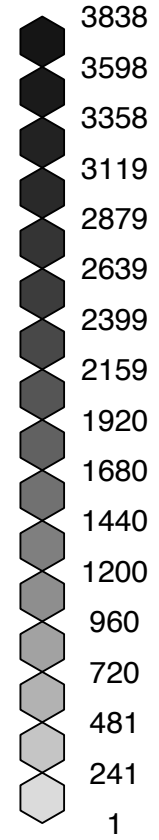
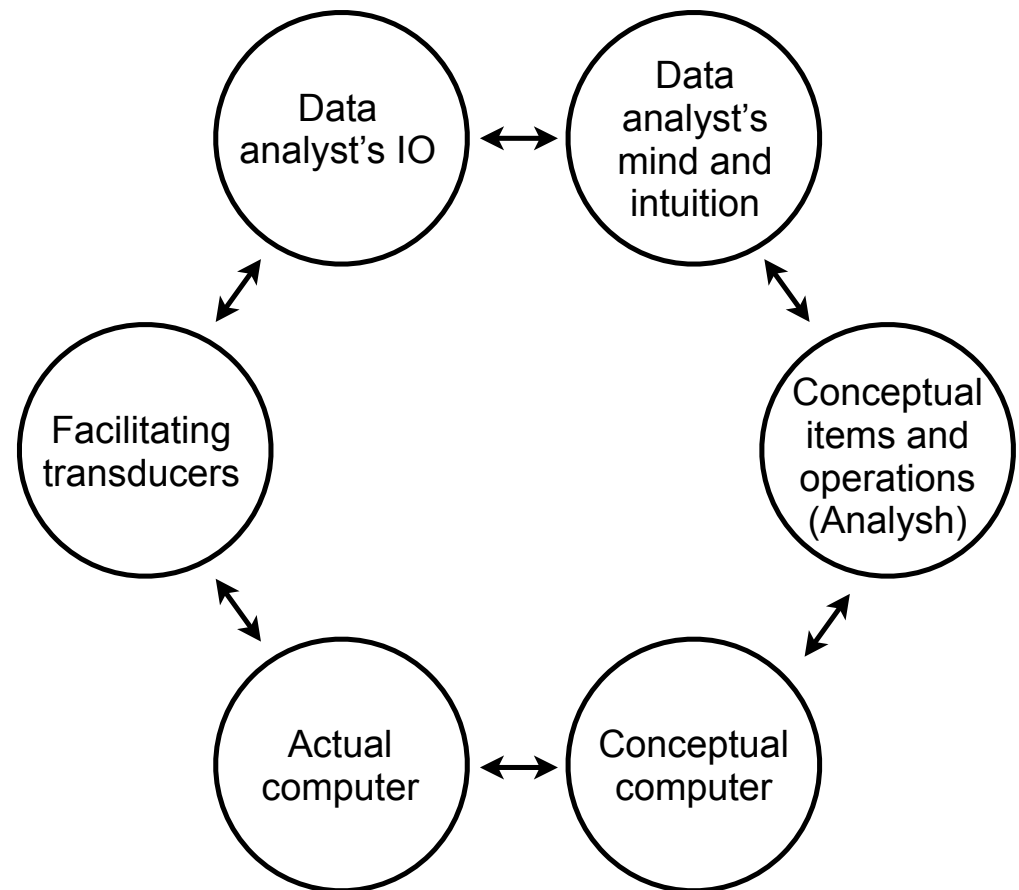weights then and now

## To sum up

So far, we've executed some simple high-level commands to make basic plots in R; there are many such functions that we will encounter over the week

Notice that for the most part the basic high-level commands do the right thing; they were designed to be part of a data analysis pipeline that Tukey and others envisioned where you go back and forth with plots and computation

We have started customizing these plots by adding annotations and graphical elements, and we'll talk more about that after the break -- we will also spend more time with the basic anatomy of an R plot

Follow the arrows clockwise from the Mind and Intuition block. Tukey's notion is that data analysts have an arsenal of operations applicable to data, which they describe to themselves and to each other in a combination of mathematics and (English) words, for which he coins the term Analysh. These descriptions can be made into algorithms (my term, not his) -- specific computational methods, but not yet realized for an actual computer (hence the conceptual computer). Then a further mapping implements the algorithm, and running it produces output for the data analyst. The output, of course, stimulates further ideas and the cycle continues. (The facilitating transducers I interpret to mean software that allows information to be translated back and forth between internal machine form and forms that humans can write or look at -- a transducer, in general, converts energy from one form to another. So parsers and formatting software would be examples.)

Taken from Chambers (2000)

Adapted from Chambers (2000)

◀ ▶ ⟳ ✕ ⌂ | http://earthquake.usgs.gov/eqcenter/dyfi/events/ci/10410337/us/ | ☆▾ · G▾ Google 🔍

# ⚡USGS
## science for a changing world

**USGS Home**
**Contact USGS**
**Search USGS**

## Earthquake Hazards Program

Home **Earthquake Center** Regional Information About Earthquakes Research & Monitoring Other Resources

You are here: Home » Earthquake Center » M4.7 – Greater Los Angeles Area, California

Latest Earthquakes
 USA
 World
 EQ Notification Service
 📶 Feeds & Data
 Animations
Recent Earthquakes:Last 8-30 Days
Earthquake Archives
 Lists & Maps
 Search EQ Database
 EQ Summary Posters
 Scientific Data
About EQ Maps
**Did You Feel It?**
 Archives
 About the Maps
 Website Feedback
 Disclaimer
 FAQ
Fast Moment Tensors
Media Info
PAGER
Seismogram Displays

## M4.7 – Greater Los Angeles Area, California
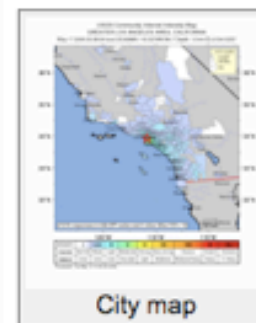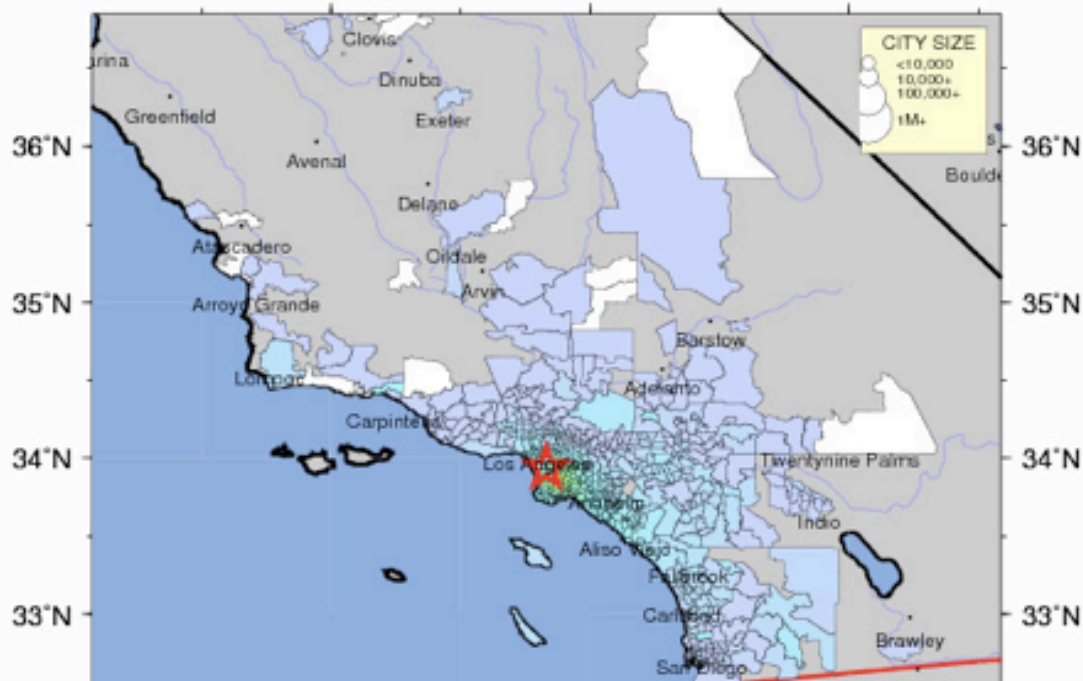Monday, May 18, 2009 at 03:39:36 UTC
Sunday, May 17, 2009 at 20:39:36 Local

33.94°N 118.34°W
Depth: 13km

**Maps** | Graphs | Responses | Downloads | **Did You Feel It? — Tell Us !**

USGS Community Internet Intensity Map
GREATER LOS ANGELES AREA, CALIFORNIA
May 17 2009 20:39:36 local 33.9396N 118.3378W M4.7 Depth: 13 km ID:ci10410337

CITY SIZE
<10,000
10,000+
100,000+
1M+

City map

City map

Done

# U.S. Census Bureau

What's New | Map Products | Boundary Files | On-Line Mapping | Related Sites | Contact | Site Map

Cartographic Products   Geography

## Cartographic Boundary Files

**Download Boundary Files**

**Descriptions and Metadata**
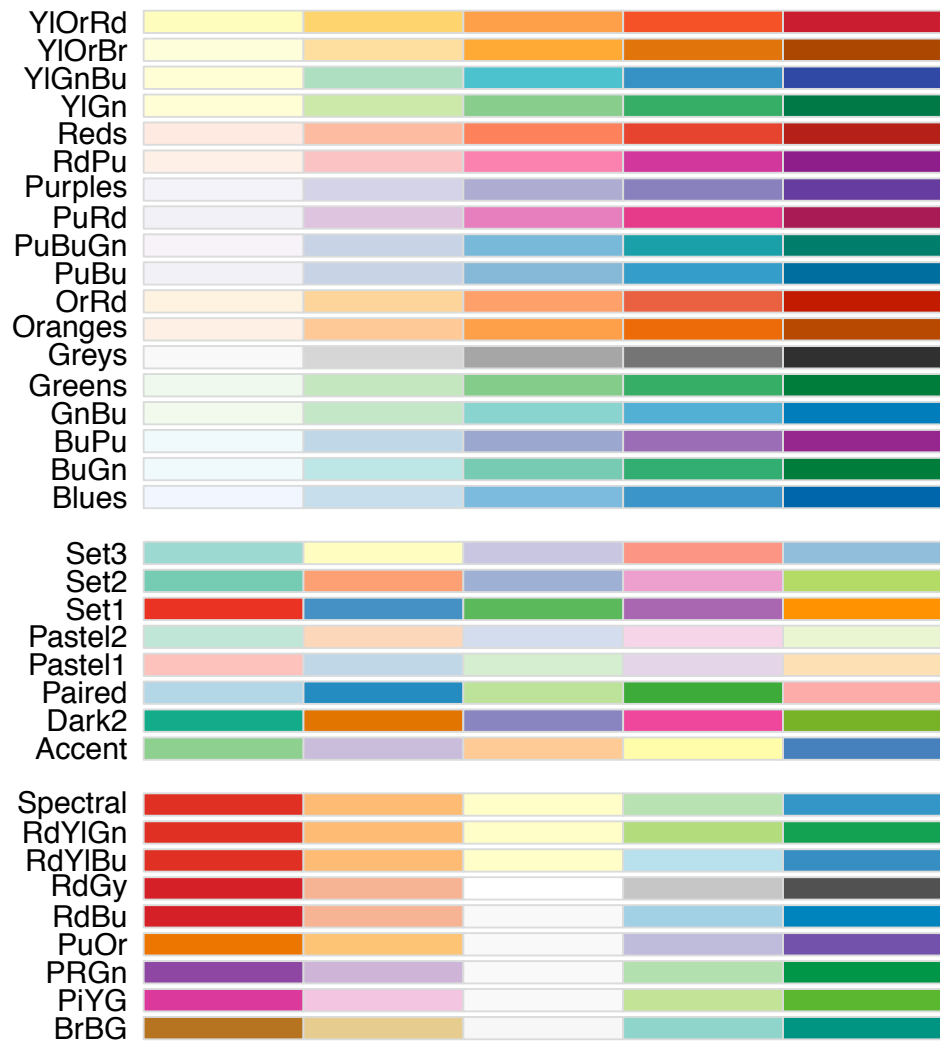
**Technical Information**

**Boundary Files Site News**

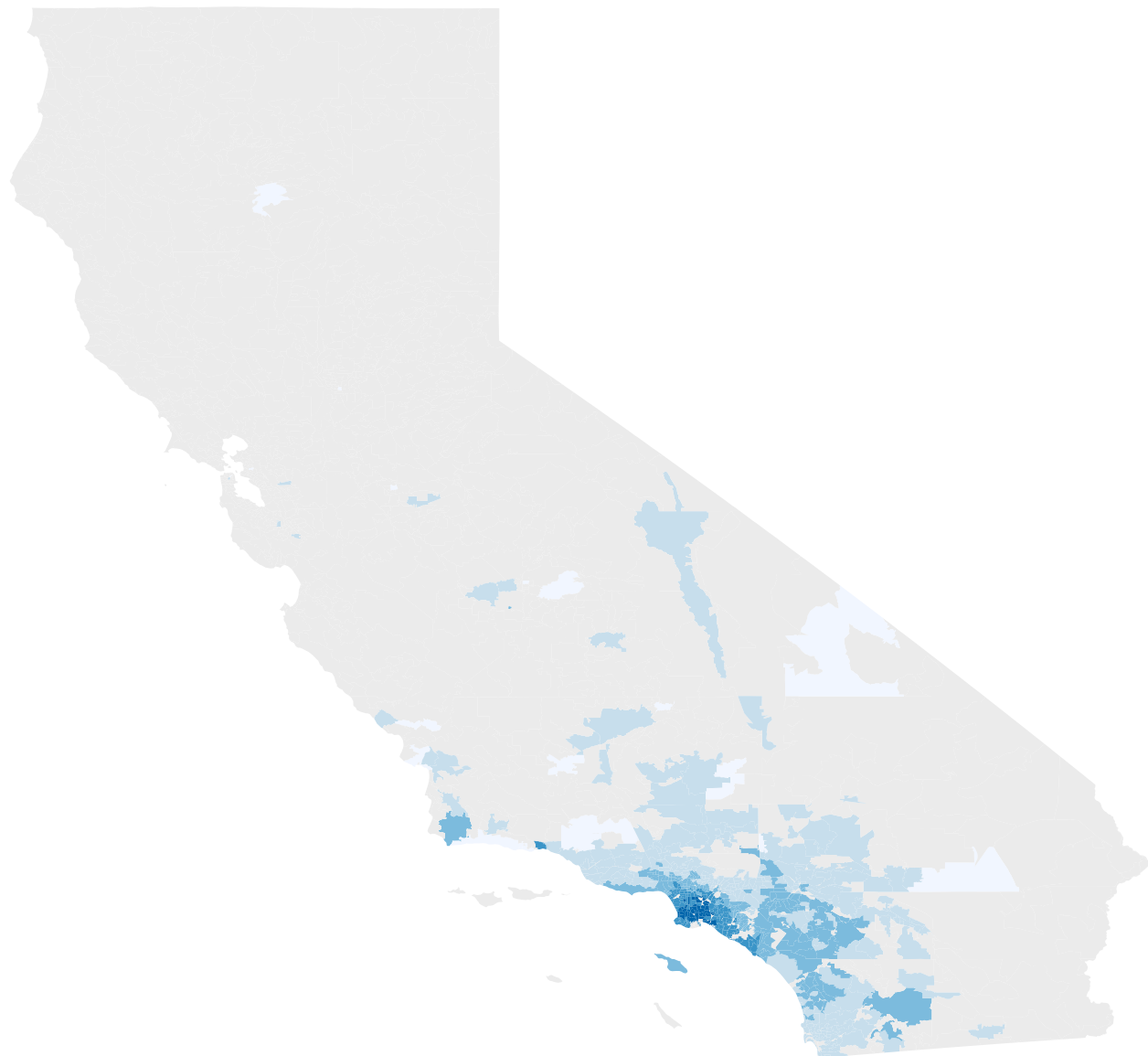**Census 2000 5-Digit ZIP Code Tabulation Areas (ZCTAs)
Cartographic Boundary Files**

ARC/INFO Export (.e00) | ArcView Shapefile (.shp) | ARC/INFO Ungenerate (ASCII)
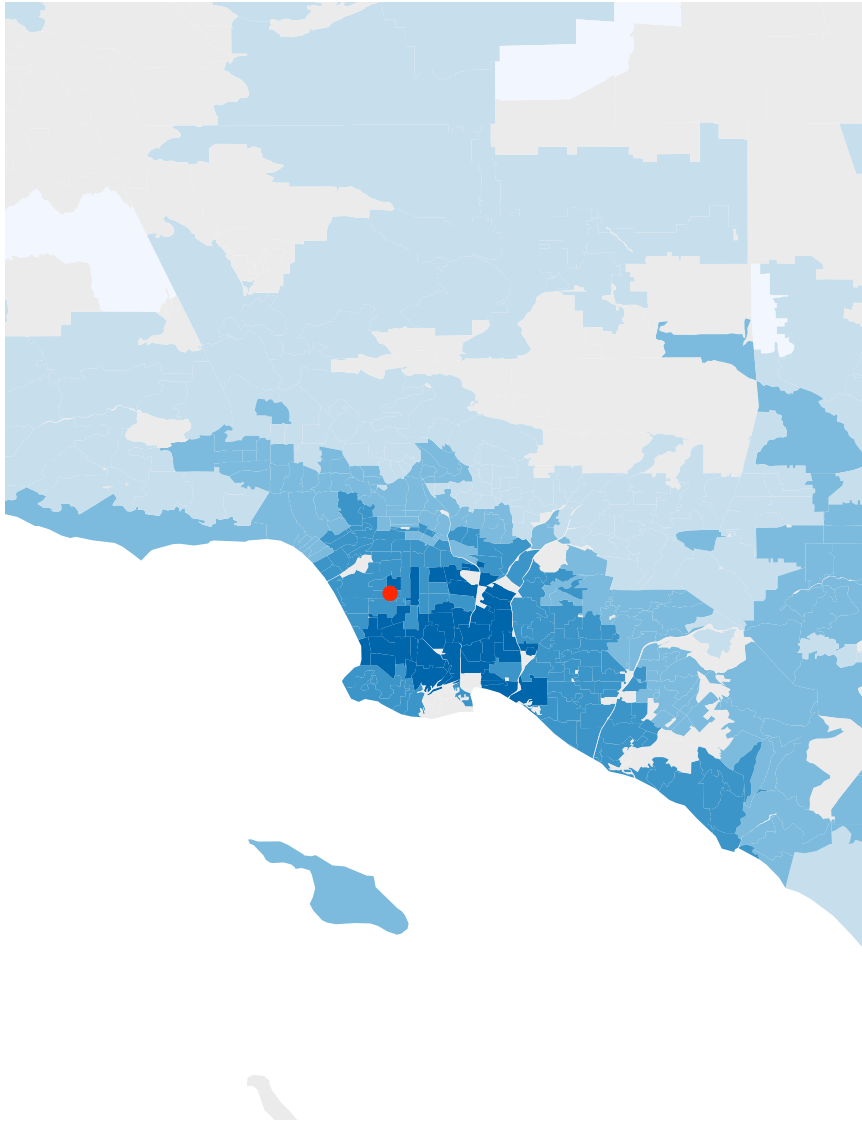Click on a file below to begin downloading...

**Census 2000 5-Digit ZIP Code Tabulation Areas (ZCTAs)
in ARC/INFO Export (.e00) format**

Alabama - zt01_d00_e00.zip (939,702 bytes)
Alaska - zt02_d00_e00.zip (1,634,983 bytes)
Arizona - zt04_d00_e00.zip (482,399 bytes)
Arkansas - zt05_d00_e00.zip (805,069 bytes)
California - zt06_d00_e00.zip (1,868,987 bytes)
Colorado - zt08_d00_e00.zip (525,215 bytes)
Connecticut - zt09_d00_e00.zip (178,620 bytes)
Delaware - zt10_d00_e00.zip (68,765 bytes)
District of Columbia - zt11_d00_e00.zip (12,289 bytes)
Florida - zt12_d00_e00.zip (1,182,978 bytes)
Georgia - zt13_d00_e00.zip (943,514 bytes)
Hawaii - zt15_d00_e00.zip (101,248 bytes)
Idaho - zt16_d00_e00.zip (629,742 bytes)
Illinois - zt17_d00_e00.zip (1,036,995 bytes)
Indiana - zt18_d00_e00.zip (582,147 bytes)
Iowa - zt19_d00_e00.zip (701,864 bytes)
Kansas - zt20_d00_e00.zip (526,709 bytes)
Kentucky - zt21_d00_e00.zip (854,491 bytes)
Louisiana - zt22_d00_e00.zip (1,474,449 bytes)
Maine - zt23_d00_e00.zip (670,663 bytes)
Maryland - zt24_d00_e00.zip (420,863 bytes)
Massachusetts - zt25_d00_e00.zip (326,474 bytes)
Michigan - zt26_d00_e00.zip (948,937 bytes)
Minnesota - zt27_d00_e00.zip (1,141,294 bytes)
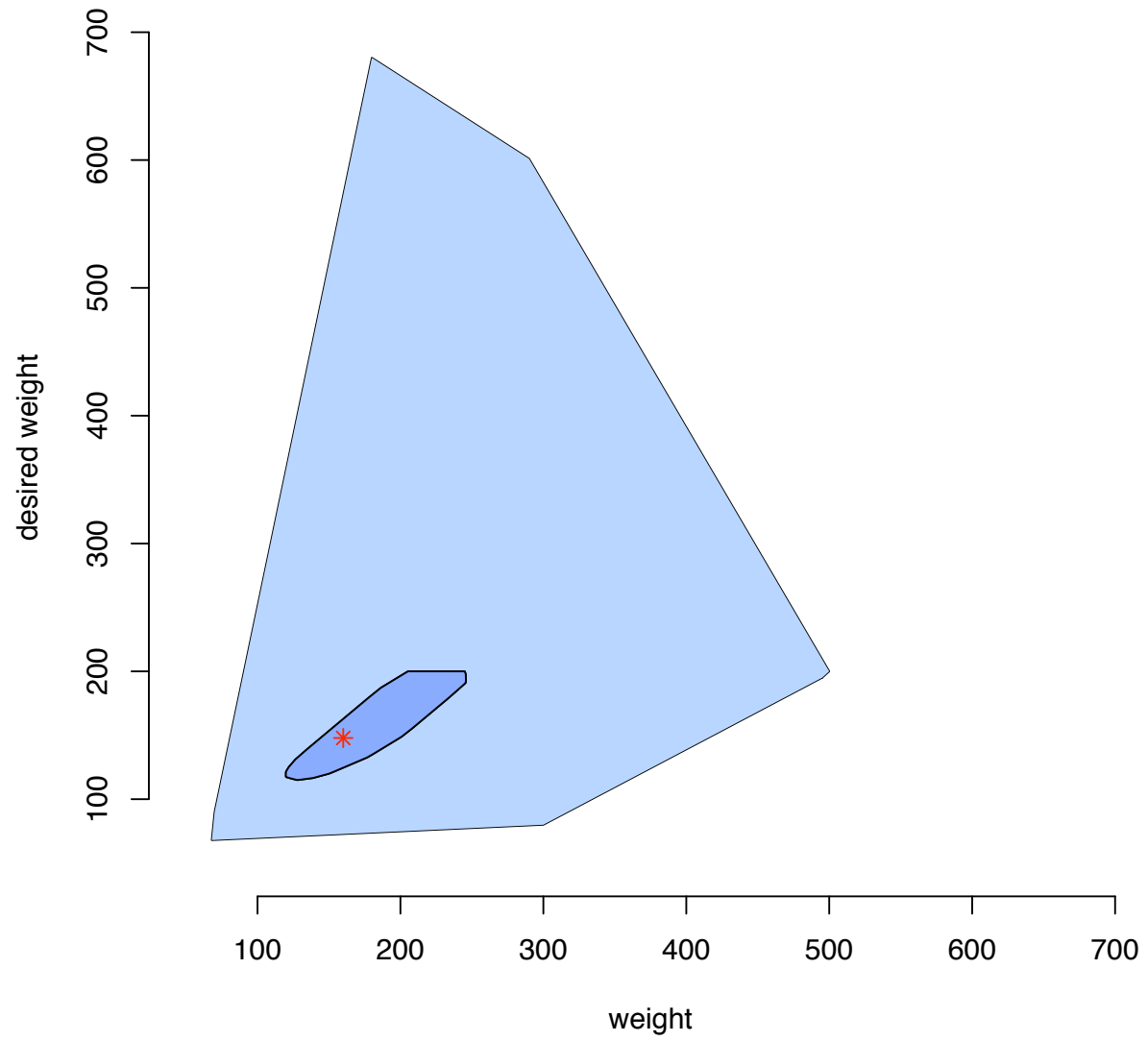
MMI Responses, Sunday's EQ in Los Angeles

## Scatterplots

We've added a line with unit slope to the plot; Why? What do you notice? What strikes you as expected? Unexpected?

Now, suppose we want to create something like a boxplot for these data; what concepts do we have to extend?

bagplot of weight and desired weight

## Bagplots

Bagplots are another Tukey innovation (along with the boxplot), but somehow they haven't caught on; why?

Can you see this being useful? Under what circumstances? How might they be interesting for our data?