

Computing, technology & Data Analysis in the Graduate Curriculum

Duncan Temple Lang
UC Davis
Dept. of Statistics

- ④ Statistics is much broader than we represent in our educational programs
 - ④ Context of Scientific Discovery, not statistical methods!
 - ④ Data analysis and “problems”
- ④ Technology has significantly altered science
And hence statistics. We must respond!
- ④ Computing & technology are essential elements of our practice, research and education

- ④ Extend stat. curricula with
 - ④ computing & technology
- ④ Mix with introduction to modern statistical methods and “real” applications of data analysis
- ④ This combination makes our students into more valuable contributors to scientific inquiry.

Calls for Change in view of statistics



1962 Annals

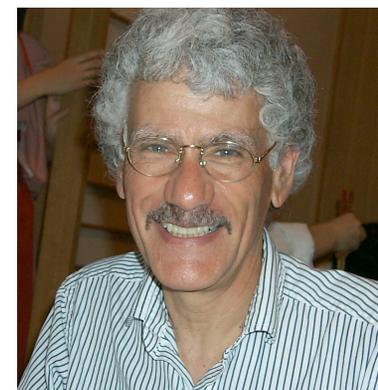


1977

Analysis of Large Complex Data

[http://www.stat.fi/isi99/
proceedings/arkisto/varasto/
frie0060.pdf](http://www.stat.fi/isi99/proceedings/arkisto/varasto/frie0060.pdf)

1997



Resistance

- ④ Nay-sayers often prioritize stat. topics over computing
- ④ defend the “mathematical” foundations of our discipline based on conservatism,
- ④ Frequently people who don't understand technology and computation and its role in practice and in reshaping opportunities for statistical thinking and research.
- ④ But not competition between computing or mathematics
both are tools for statistical concepts & practice.

- ④ After 46, 31, 11 years, it is time for action & change, not just talk.
- ④ We must do the best we can and strive to get computing, technology and data analysis adequately into our curricula.
- ④ Need to attract/retain researchers with modern, different perspective on data analysis & its impact.

Outline

- ④ Why?
- ④ What?
- ④ For whom?
- ④ By whom?
- ④ How to achieve this?

- ④ Need our students to be computationally literate to be able to
 - ④ Do computations for their own research (simulations, implement methodology)
 - ④ to help with our research
 - ④ interact with scientists using complex data from diverse sources
 - ④ disseminate new statistical methods as software
 - ④ understand, critically appreciate and exploit new technologies as they emerge to allow us to do new types of data analysis.

- ④ Opportunity to teach statistical methodology that the students wouldn't necessarily see in a heuristic manner
- ④ Improve their (exploratory) data analysis skills and intuition.
- ④ Expose them to research by implementing computations within a paper.

- ④ Omitting computing & technology from our curricula means we are “playing with one hand tied behind our back”
- ④ We cannot provide complete solutions to scientific problems, but merely prescriptions (not functioning/tangible tools) for how others can approach these problems.
- ④ Software for “doing” statistics in the analytic pipeline

What – Broad Topics

- ④ Fundamentals of scientific programming -
- ④ Computing for Research - profiling, C, basic parallel computing, object-oriented computing, “R packages”
- ④ Computational Statistics - Lin. Alg, Numeric optimization, RNG, MCMC, EM, resampling, numerical integration,
- ④ Data Technologies - Databases (SQL), Regular Expressions, XML, Web services.
- ④ Visualization technologies - graphical techniques & software; dynamic, interactive & Web displays

Intro to Stat. Computing

- Operating system concepts - commonalities & differences
 - file system (files, folders, binary/text), editors, ...
- Types of languages - compiled/interpreted, vectorized/scalar, task-specific languages, Perl, Python, R, MATLAB, SAS
- Language elements - data types, subsetting, function calls, vectorized looping (apply()), control flow
- Input and Output (I/O)
- Writing functions - mechanics, design,
 - *Debugging - tools, technique and philosophy/art,
- Efficiency - alg. complexity, idioms, profiling, interface to C/FORTRAN/...
- Batch computing & remote "shells"

- ③ Vital to avoid teaching just the syntax of a particular language, or how to cut-and-paste & modify templates
- ③ Need to teach concepts of computing, how to understand other languages, approach a computational task & abstract the ideas.

For whom?

- ④ Different types of students - different courses
- ④ Each student should take ≥ 2 computing classes
- ④ Required class - "Scientific" Programming
 - ④ teach how to think & reason about computing and express stat. tasks as computations.
 - ④ ideally also cover R/MATLAB fundamentals, interface to C, efficiency, parallel computing (in context of data analysis).
- ④ And one class in either Data Technologies, Computational Statistics, Advanced Computing.

Types of students & second course

- ④ Student studying methodology research (theory)
 - ④ simulation, software development (e.g. R packages), efficiency, algorithmic complexity, numerical algorithms, parallel computing, streaming data, visualization
- ④ Probability – simulation, RNG, efficiency, visualization.
- ④ Applications – data technologies for accessing data, additional languages, efficiency, parallel programming, visualization.

- ④ For me, programming and the basics of data technologies
 - ④ I/O for complex data
 - ④ text manipulation & regular expressions
 - ④ databases
 - ④ XML
- ④ are vital for all students working with data.

Masters Students

- ④ What do they end up doing?
 - ④ Data manipulation and processing – data technologies
 - ④ Exploratory Data Analysis & Reports – visualization
 - ④ Simulations – programming
 - ④ Modeling – R, MATLAB, SAS.
- ④ First class in stat. computing & then mix of visualization, data technologies, SAS
- ④ Data, data, data....

Can we weave topics into existing classes?

- ④ Not the programming class!
Starting from nothing
- ④ We need programming to be a fundamental class to
 - ④ emphasize its importance & establish culture of computing.
 - ④ provide solid, rich foundation for other topics,
 - ④ allow the students to absorb the concepts/reasoning over a quarter/semester,
 - ④ put in the context of data analysis/math. stat.

By whom?

- ④ Rarely in our graduate programs
- ④ More senior faculty haven't been exposed to this, so can't teach it, so students aren't exposed to it, so ...
- ④ Students left to learn it on their own with little encouragement or priority
 - ④ empirically results are poor with major misconceptions
- ④ So very few instructors who can teach computing and technology

Computer Science Classes

- ④ Can we send our students to Computer Science classes on programming? databases? text manipulation?
- ④ to Applied Math for numerical analysis? optimization? ...

⊖ No

- ⊖ we do a different type of programming (vectorized, interpreted languages versus compiled scalar languages)
- ⊖ a class in databases teaches internal details of database not how to use it.
- ⊖ importantly, don't put these methods in the context of statistical data analysis.

How?

- Have to train instructors or train themselves?
- NSF grant (Nolan, Hansen, Temple Lang) to
 - develop potential syllabi & topics for computing
 - create resources for teaching - lecture notes, exercises/homeworks/projects/case-studies, text book
 - teach instructors how to teach computing
 - evangelize computing, technologies & data analysis within the community via papers, talks, etc.

- ④ May 2007 - Workshop for syllabi
- ④ July 2008 - Workshop for teaching instructors
- ④ 2009 - final workshop. What form?
 - ④ teach additional instructors (same as 2008)?
 - ④ summer school on technology, computing & data analysis for recent graduates, like New Researchers Conf.?
 - ④ summer school for PhD students starting research?
 - ④ small working group to complete materials for others to pick up?

Materials at

- Workshop 1

<http://www.stat.berkeley.edu/twiki/Workshop/CompCurric>

- Workshop 2

<http://www.stat.berkeley.edu/~statcur/>

Summary

- ④ It is time to step up and do something about computing & technology & data analysis in our curricula.
- ④ Must have dedicated statistical/scientific programming course
- ④ Data technologies, advanced/research computing, numerical algorithms, visualization classes or individual topics

Summary ctd.

- ④ Grow pool of potential instructors by teaching these classes now
- ④ and teaching existing instructors via workshops & developing class materials
- ④ What form for next workshop?
Seek funding for additional workshops?
- ④ Strategic Initiative from ASA, ISI, ...

Actions

- ④ Time for action on technolog & computing.
- ④ Departments should introduce computing into graduate & undergraduate classes.
 - ④ Explicit “computing” classes
 - ④ Introductory programming, data technologies for scientific computing with data
 - ④ Second class or integrate topics into classes from: Data technologies, advanced/research computing, numerical algorithms, visualization classes or individual topics