

CS 109: Data Science

Process, Data, and Visual Attributes

Hanspeter Pfister
pfister@seas.harvard.edu

Joe Blitzstein
blitzstein@stat.harvard.edu

This Week

- HW0 - due next Tuesday (not graded)
 - Install Anaconda & IPython frameworks
 - Sign up for Piazza and introduce yourself
 - Fill out survey
- Friday lab **10-11:30 am** in MD G115
 - *Intro to Python* with Ian Stokes-Rees
 - Make sure to have IPython installed and ready
- Readings - post comments on Piazza

Outline

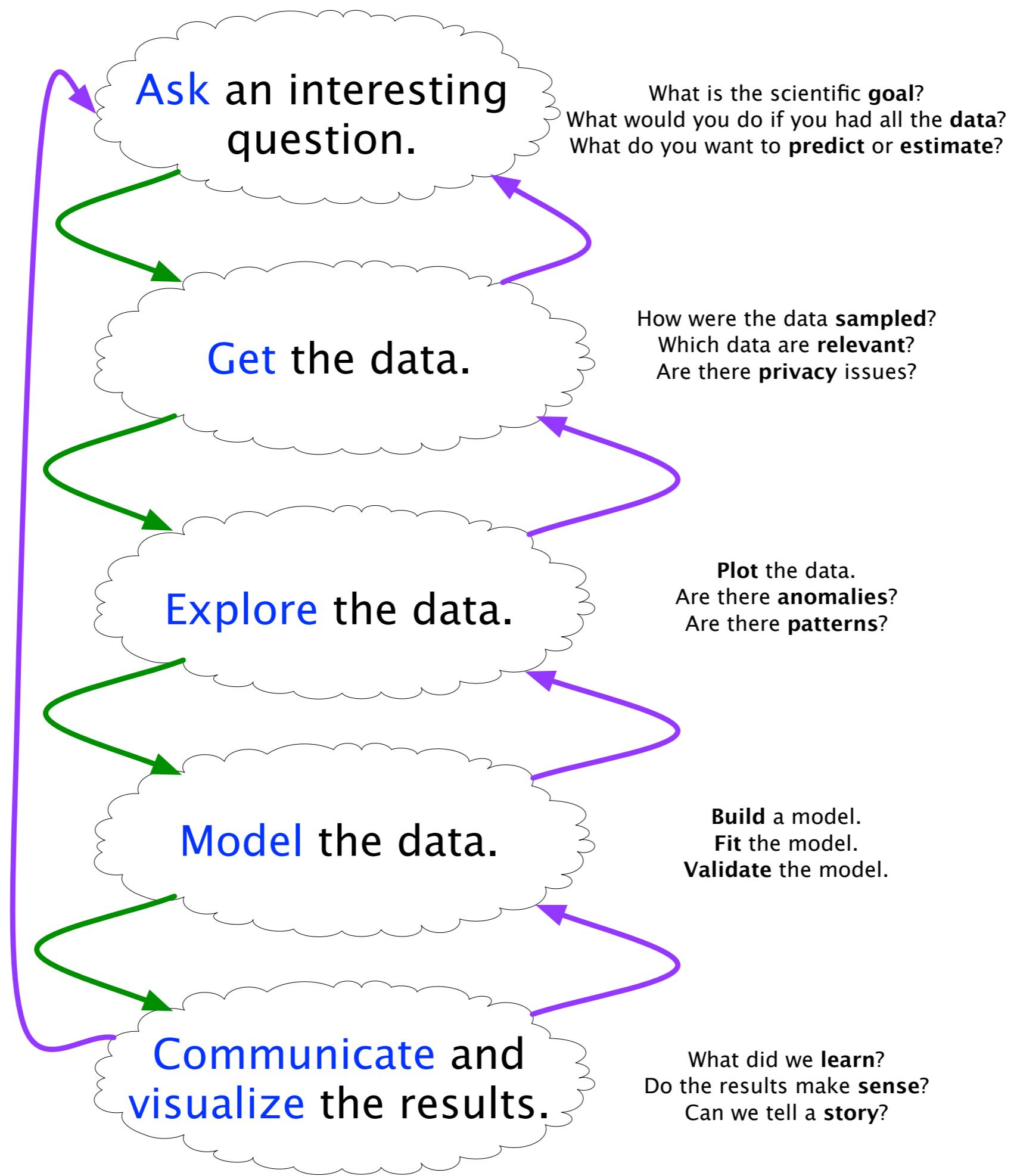
- Process & Process Books
- What makes visualizations effective?
- Data Sources & Data Cleanup

Process

Data Exploration

Not always sure what we are looking for
(until we find it)





What do analysts do?

Enterprise Data Analysis and Visualization: An Interview Study

Sean Kandel, Andreas Paepcke, Joseph M. Hellerstein, and Jeffrey Heer

Abstract—Organizations rely on data analysts to model customer engagement, streamline operations, improve production, inform business decisions, and combat fraud. Though numerous analysis and visualization tools have been built to improve the scale and efficiency at which analysts can work, there has been little research on how analysis takes place within the social and organizational context of companies. To better understand the enterprise analysts’ ecosystem, we conducted semi-structured interviews with 35 data analysts from 25 organizations across a variety of sectors, including healthcare, retail, marketing and finance. Based on our interview data, we characterize the process of industrial data analysis and document how organizational features of an enterprise impact it. We describe recurring pain points, outstanding challenges, and barriers to adoption for visual analytic tools. Finally, we discuss design implications and opportunities for visual analysis research.

Index Terms—Data, analysis, visualization, enterprise.

1 INTRODUCTION

Organizations gather increasingly large and complex data sets each year. These organizations rely on data analysis to model customer engagement, streamline operations, improve production, inform sales and business decisions, and combat fraud. Within organizations, an increasing number of individuals—with varied titles such as “business analyst”, “data analyst” and “data scientist”—perform such analyses. These analysts constitute an important and rapidly growing user population for analysis and visualization tools.

Enterprise analysts perform their work within the context of a larger organization. Analysts often work as a part of an analysis team or business unit. Little research has observed how existing infrastructure, available data and tools, and administrative and social conventions within an organization impact the analysis process within the enterprise. Understanding how these issues shape analytic workflows can inform the design of future tools.

ery and wrangling, often the most tedious and time-consuming aspects of an analysis, are underserved by existing visualization and analysis tools. We discuss recurring pain points within each task as well as difficulties in managing workflows across these tasks. Example pain points include integrating data from distributed data sources, visualizing data at scale and operationalizing workflows. These challenges are typically more acute within large organizations with a diverse and distributed set of data sources.

We conclude with a discussion of future trends and the implications of our interviews for future visualization and analysis tools. We argue that future visual analysis tools should leverage existing infrastructures for data processing to enable scale and limit data migration. One avenue for achieving better interoperability is through systems that specify analysis or data processing operations in a high-level language, enabling retargeting across tools or platforms. We also note

What do analysts do?

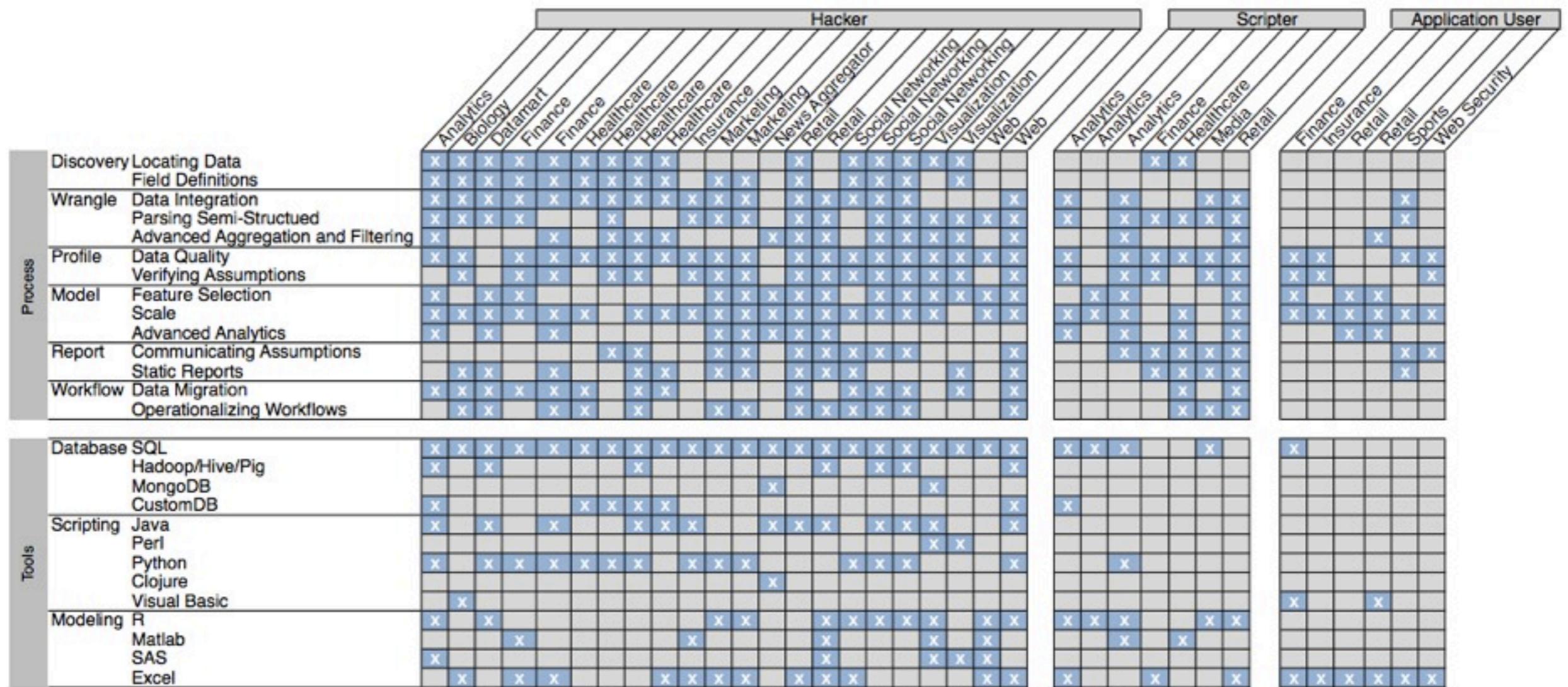


Fig. 1. Respondents, Challenges and Tools. The matrix displays interviewees (grouped by archetype and sector) and their corresponding challenges and tools. *Hackers* faced the most diverse set of challenges, corresponding to the diversity of their workflows and toolset. *Application users* and *scripters* typically relied on the IT team to perform certain tasks and therefore did not perceive them as challenges.

What do analysts do?

I spend more than half of my time integrating, cleansing and transforming data without doing any actual analysis. Most of the time I'm lucky if I get to do any analysis. Most of the time once you transform the data you just do an average... the insights can be scarily obvious. It's fun when you get to do something somewhat analytical.

Exploratory Data Analysis

“The greatest value of a picture is when it forces us to notice what we never expected to see.”



John Tukey

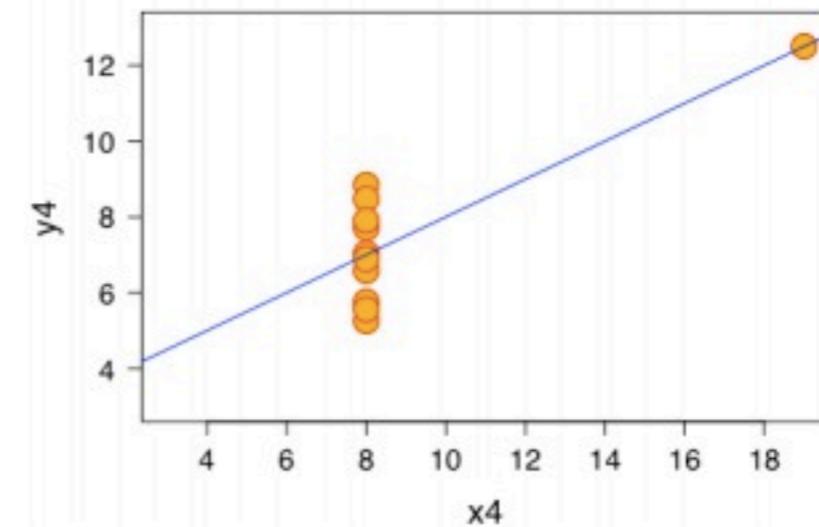
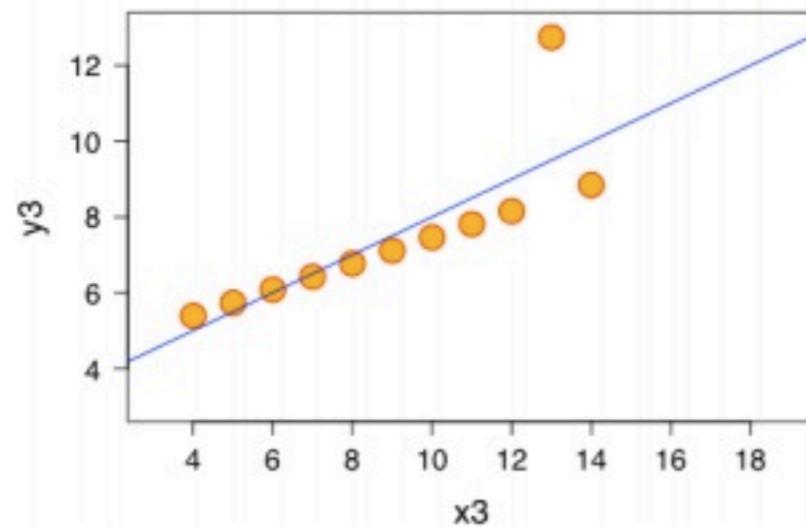
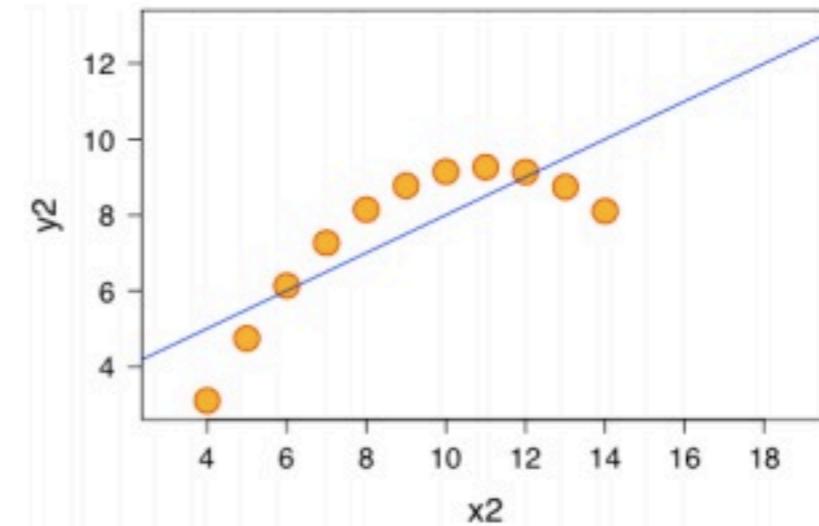
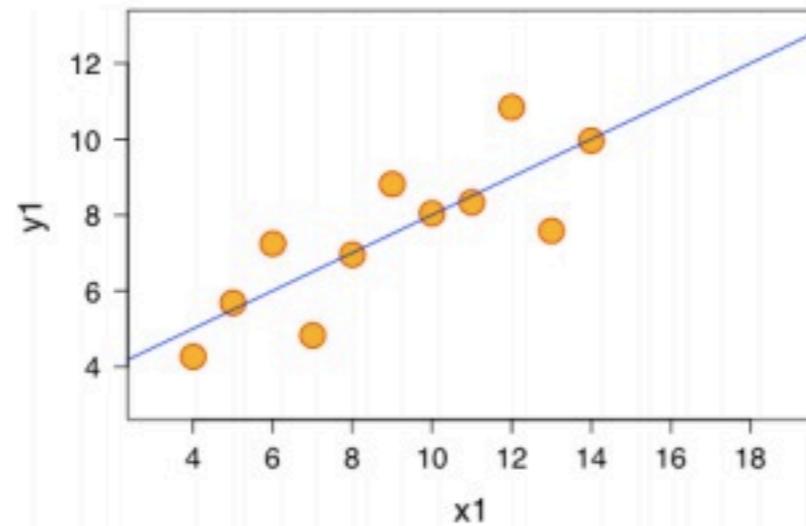
Anscombe's Quartet

Same mean, variance, correlation, and linear regression line

Anscombe's Quartet: Raw Data								
	I		II		III		IV	
	x	y	x	y	x	y	x	y
	10.0	8.04	10.0	9.14	10.0	7.46	8.0	6.58
	8.0	6.95	8.0	8.14	8.0	6.77	8.0	5.76
	13.0	7.58	13.0	8.74	13.0	12.74	8.0	7.71
	9.0	8.81	9.0	8.77	9.0	7.11	8.0	8.84
	11.0	8.33	11.0	9.26	11.0	7.81	8.0	8.47
	14.0	9.96	14.0	8.10	14.0	8.84	8.0	7.04
	6.0	7.24	6.0	6.13	6.0	6.08	8.0	5.25
	4.0	4.26	4.0	3.10	4.0	5.39	19.0	12.50
	12.0	10.84	12.0	9.13	12.0	8.15	8.0	5.56
	7.0	4.82	7.0	7.26	7.0	6.42	8.0	7.91
	5.0	5.68	5.0	4.74	5.0	5.73	8.0	6.89
mean	9.0	7.5	9.0	7.5	9.0	7.5	9.0	7.5
var.	10.0	3.75	10.0	3.75	10.0	3.75	10.0	3.75
corr.		0.816		0.816		0.816		0.816

Anscombe's Quartet

Same mean, variance, correlation, and linear regression line



Example: Antibiotics
Will Burtin, 1951

Effectiveness of Antibiotics

Table 1: Burtin's data.

Bacteria	Antibiotic			Gram Staining
	Penicillin	Streptomycin	Neomycin	
<i>Aerobacter aerogenes</i>	870	1	1.6	negative
<i>Brucella abortus</i>	1	2	0.02	negative
<i>Brucella anthracis</i>	0.001	0.01	0.007	positive
<i>Diplococcus pneumoniae</i>	0.005	11	10	positive
<i>Escherichia coli</i>	100	0.4	0.1	negative
<i>Klebsiella pneumoniae</i>	850	1.2	1	negative
<i>Mycobacterium tuberculosis</i>	800	5	2	negative
<i>Proteus vulgaris</i>	3	0.1	0.1	negative
<i>Pseudomonas aeruginosa</i>	850	2	0.4	negative
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	negative
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	negative
<i>Staphylococcus albus</i>	0.007	0.1	0.001	positive
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	positive
<i>Streptococcus faecalis</i>	1	1	0.1	positive
<i>Streptococcus hemolyticus</i>	0.001	14	10	positive
<i>Streptococcus viridans</i>	0.005	10	40	positive

Data & Questions

- What are the data types?
- What are possible questions?

Bacteria	Penicillin	Antibiotic Streptomycin	Neomycin	Gram stain
<i>Aerobacter aerogenes</i>	870	1	1.6	-
<i>Brucella abortus</i>	1	2	0.02	-
<i>Bacillus anthracis</i>	0.001	0.01	0.007	+
<i>Diplococcus pneumoniae</i>	0.005	11	10	+
<i>Escherichia coli</i>	100	0.4	0.1	-
<i>Klebsiella pneumoniae</i>	850	1.2	1	-
<i>Mycobacterium tuberculosis</i>	800	5	2	-
<i>Proteus vulgaris</i>	3	0.1	0.1	-
<i>Pseudomonas aeruginosa</i>	850	2	0.4	-
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	-
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	-
<i>Staphylococcus albus</i>	0.007	0.1	0.001	+
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	+
<i>Streptococcus fecalis</i>	1	1	0.1	+
<i>Streptococcus hemolyticus</i>	0.001	14	10	+
<i>Streptococcus viridans</i>	0.005	10	40	+

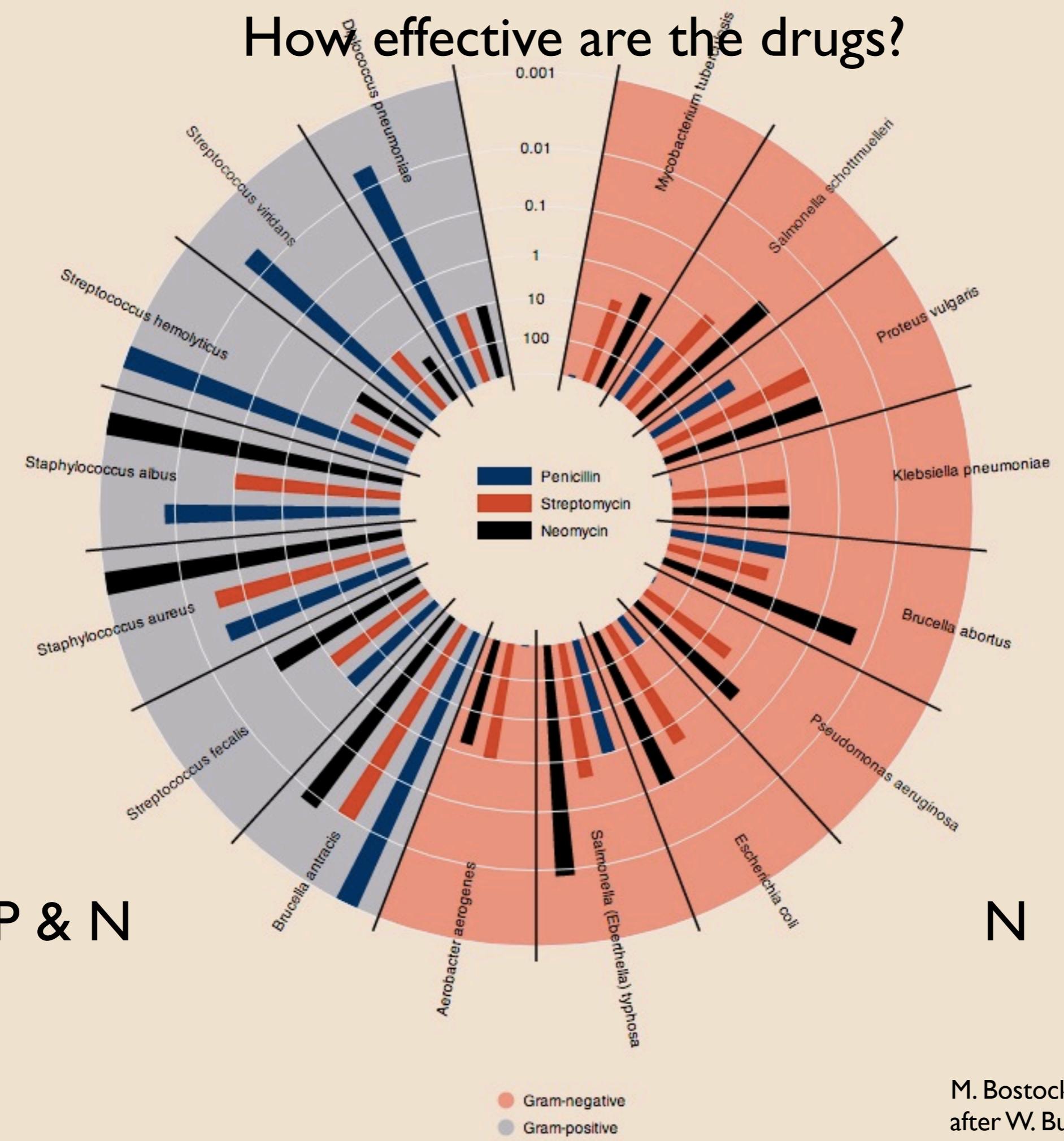
Data

- Genus & species of bacteria [string]
- Antibiotic name [string]
- Gram staining? [pos/neg]
- Minimum inhibitory concentration (mg/ml) [float]
(lower == more effective)

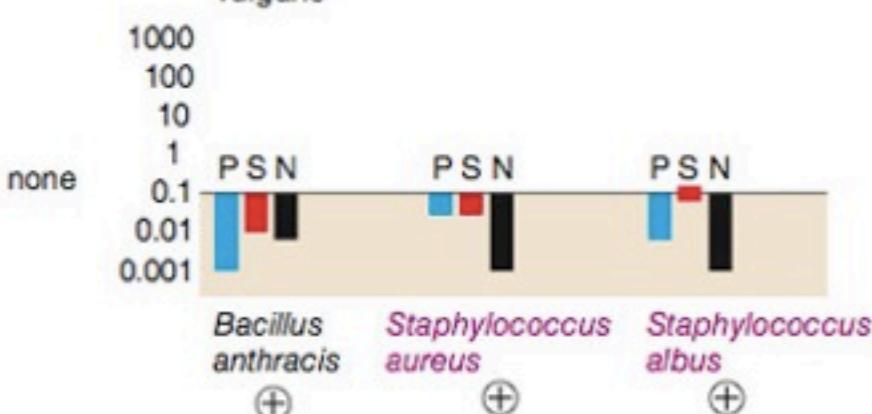
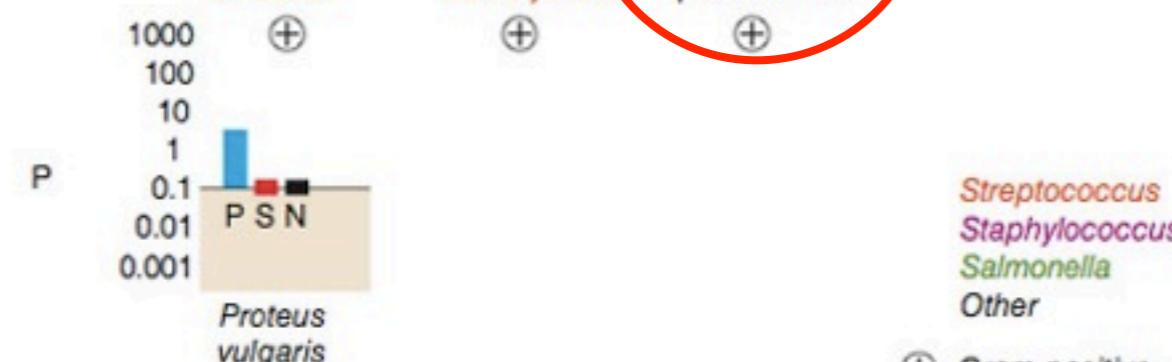
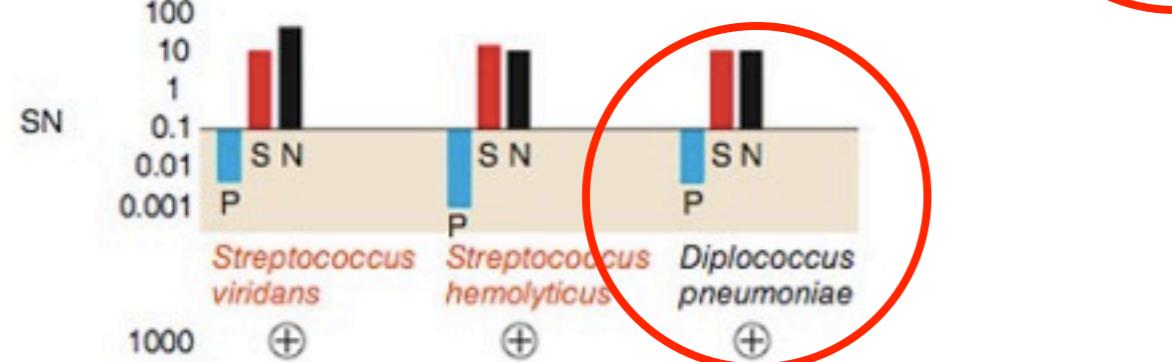
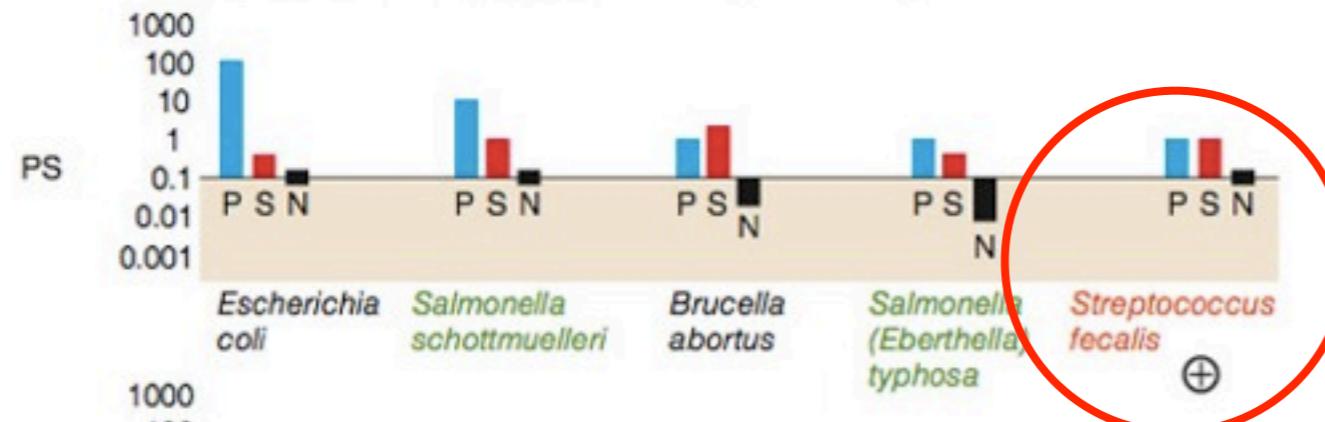
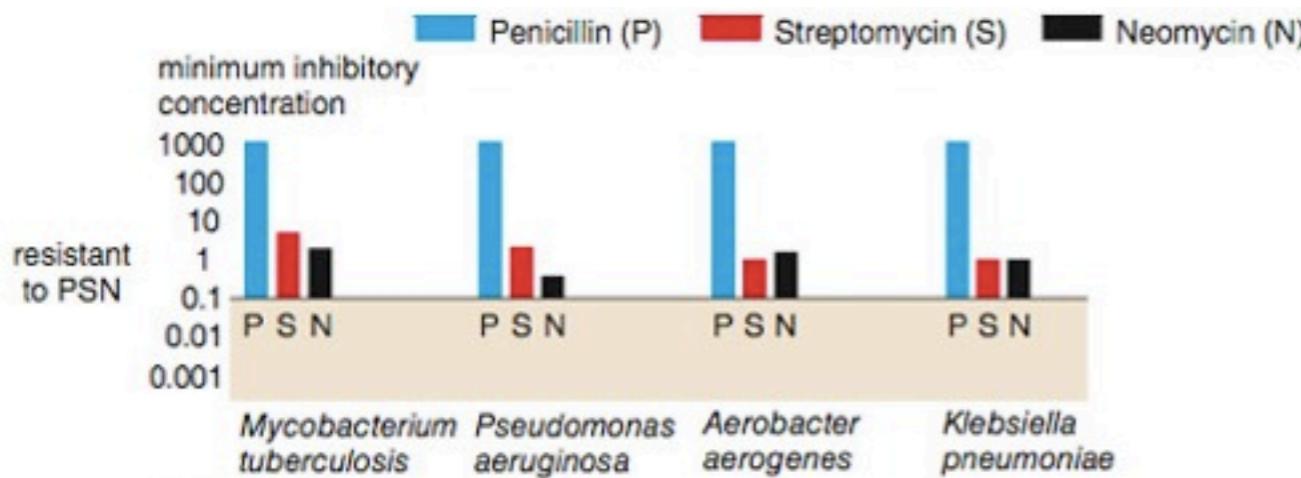
Bacteria	Penicillin	Antibiotic Streptomycin	Neomycin	Gram stain
<i>Aerobacter aerogenes</i>	870	1	1.6	-
<i>Brucella abortus</i>	1	2	0.02	-
<i>Bacillus anthracis</i>	0.001	0.01	0.007	+
<i>Diplococcus pneumoniae</i>	0.005	11	10	+
<i>Escherichia coli</i>	100	0.4	0.1	-
<i>Klebsiella pneumoniae</i>	850	1.2	1	-
<i>Mycobacterium tuberculosis</i>	800	5	2	-
<i>Proteus vulgaris</i>	3	0.1	0.1	-
<i>Pseudomonas aeruginosa</i>	850	2	0.4	-
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	-
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	-
<i>Staphylococcus albus</i>	0.007	0.1	0.001	+

What Questions?

How effective are the drugs?



M. Bostock, Protovis
after W. Burtin, 1951



How do the bacteria compare?

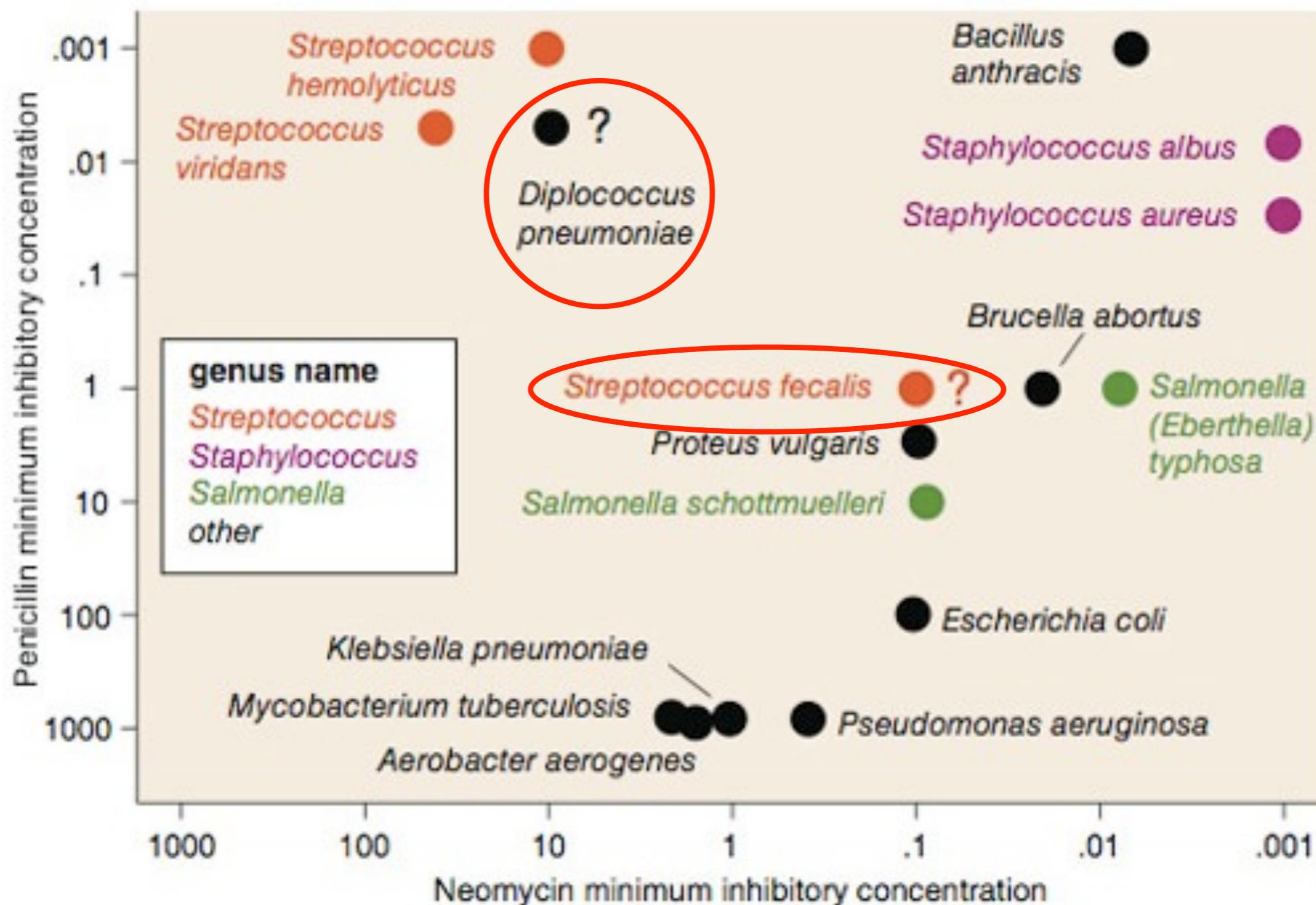
Not a streptococcus!
(realized ~30 years later)

Really a streptococcus!
(realized ~20 years later)

Streptococcus
Staphylococcus
Salmonella
Other
⊕ Gram positive

Wainer & Lysen, "That's funny..."
American Scientist, 2009
Adapted from Brian Schmotzer

How do the bacteria compare?



“The greatest value of a picture is when it forces us to notice what we never expected to see.”



John Tukey

Process Books

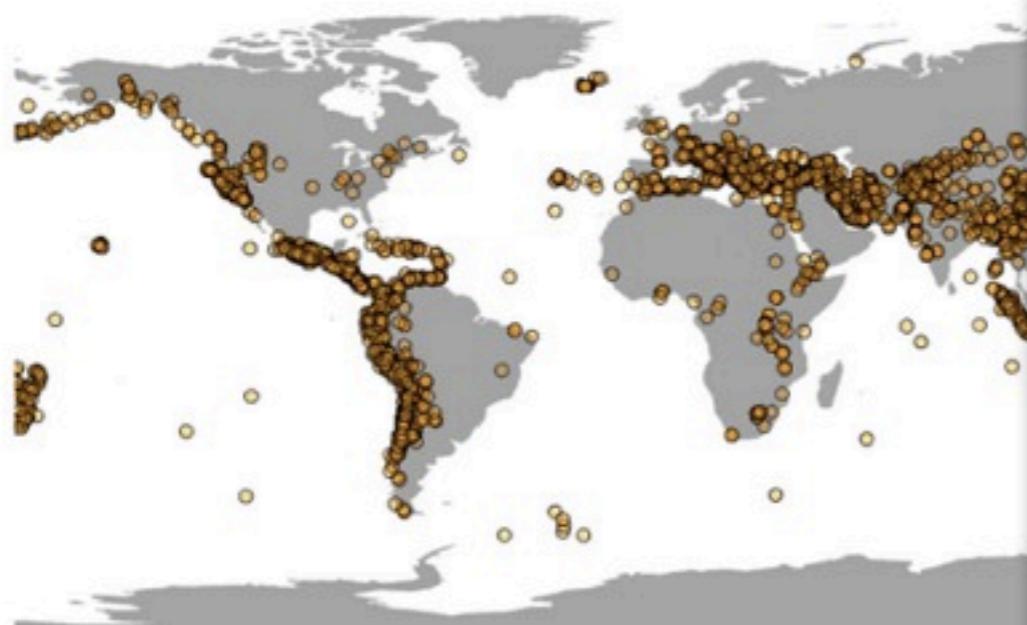
Process Books

[Blake Walsh, Gabriel Trevino, Antony Bett, CS171, 2013]

More Data Cleaning and Automation of Point Plotting in Prototype [entry by btwalsh, 04.01]

Today, now that I have confirmed that my implementation of a few sample points is functional, I took our XLS file of earthquake data and began to clean it up some more with Google Refine specifically for the purpose of easily integrating it with our existing Processing code. After the cleaning that was performed, I renamed the columns to read more sensibly, removed all entries that did not indicate a latitude and longitude for the location, and then removed all rows of data that similarly did not actually report a magnitude.

Now, we have implemented a way to automatically map our latitudinal/longitude data onto our grid of x-y pixels, as demonstrated below:



[note by gtrevino 04.01]

Thankfully we got this to work (after dealing with some null pointer exceptions) and we are currently working on our data density, i.e. how to deal with multiple points overlapping, a task which will be especially important as we implement user interaction. This will probably require a lot of time and effort on the part of the data, but we will be able to add it to our system.

[Varun Bansal, Cici Cao, Sofia Hou, CS171, 2013]

project process book

Users

The target audience is the general public, lovers of music, or simply those who want to learn about music. To appeal to the eclectic tastes of each individual, we have included nine main genres of music ranging from classical to country to Rap. This is not a visualization on the fundamentals of music theory so everyone has equal access to it - both for education and personal enjoyment!

Related Work

Everyone loves music, and with the physics of sound the data on music is limitless. Each individual piece of song can have a vast amount of frequency, amplitude, and beats data. However, we were inspired to analyze music after seeing a project named the Shape of a Song (<http://bewitched.com/song.html>) in which the author visualized, through arch diagrams, similar patterns within different types of songs. This

provided analysis on two levels: song and genre specific. On the song specific level, the arches showcased observations such as how the beginning of the song was similar to the end. On the genre specific level, comparing the arch diagrams of songs from different genres (i.e. modern techno or pop versus classical) shows characteristic patterns of each genre, with modern electronic and synthesized music displaying more repetition and similarity throughout the piece.



Chopin, Mazurka in F# Minor

The image illustrates the complex, nested structure of the piece.

IPython Notebooks

<http://nbviewer.ipython.org/>

Home FAQ IPython Bookmarklet

IPython Notebook Viewer

A Simple way to share your IP[y]thon Notebook as Gists.

Share your own notebook, or browse others'

Enter a gist number or url Go!

IP[y]: Notebook 01 Documenting your Research Journey

File Edit View Insert Cell Kernel Help

Documenting your Research Journey

The purpose of this code is to show how IPython notebooks can be used to document your GPU and the CPU. We compare the performance of each method using the system to document.

load image

```
In [1]: import PIL  
import PIL.Image  
  
image = PIL.Image.open("cinque_terre.jpg")  
image_array_rgb = numpy.array(image)  
  
r_original,g_original,b_original = numpy.split(image_array_rgb,  
a_original = numpy.ones_like(r_original)  
rgba_original = numpy.concatenate((r_original,
```

Figure: A grayscale image titled "cinque_terre.jpg".

```
figsize(6,4)  
  
matplotlib.pyplot.imshow(rgba_original);  
matplotlib.pyplot.title("rgba_original");
```

Figure: A grayscale image titled "rgba_original".

Probabilistic Programming

Why would I want samples from the posterior, anyways?

We will deal with this question for the remainder of the book, and it is an understatement to say we can perform amazingly useful things. For now, let's finish with using posterior samples to answer the follow question: what is the expected number of texts at day t , $0 \leq t \leq 30$? Recall that the expected value of a Poisson is equal to its parameter λ , then the question is equivalent to what is the expected value of λ at time t ?

In the code below, we are calculating the following: Let i index a particular sample from the posterior distributions. Given a day t , we average over all λ_i on that day t , using $\lambda_{t,i}$ if $t < i$; else we use $\lambda_{i,t}$.

```
import numpy as np  
import matplotlib.pyplot as plt  
  
# observed texts per day  
observed_texts_per_day = np.loadtxt('observed_texts_per_day.txt')  
  
# calculate the expected number of texts received  
expected_texts_per_day = np.zeros(len(observed_texts_per_day))  
for i in range(len(expected_texts_per_day)): # loop over days  
    for j in range(i+1): # loop over posterior samples  
        expected_texts_per_day[i] += observed_texts_per_day[j] / len(expected_texts_per_day)
```

Figure: A bar chart titled "Expected # text messages" showing the distribution of text messages per day. A red line shows the average expected value over time.

XKCD Plot With Matplotlib

XKCD plots in Matplotlib

Out [1]:

CHECK IT OUT!

DAMPED SINE DAMPED COSINE

Sometimes when showing schematic plots, this is the type of figure I want to display. But drawing it by hand is a bit of a pain in matplotlib. The problem is, matplotlib is a bit too precise. Attempting to duplicate this figure in matplotlib leads to:

Exploring R formula

```
In [4]: show_xkcd.random.normal(size=(10, 10), "My design")
```

Figure: A plot titled "XKCD plots in Matplotlib" showing two oscillating curves: a blue dashed line for "DAMPED SINE" and a red dashed line for "DAMPED COSINE". The plot is labeled "CHECK IT OUT!".

Non Parametric Regression

Covariance function

The behavior of individual realizations from the GP is governed by the covariance function. The Matern class of functions is a flexible choice.

```
In [4]: from gpyt.gp.cov_func import matern  
import numpy as np  
C = CovarianceMatrix_Burnhamers_matern42, diff_degree=1.4, amp=1.0, scale=1., rank_limit=1000  
  
subplot(1,2,1)  
contourf(x, y, C(x,y).view(ndarray), origin='lower', extent=(-1,1,-1,1), cmap=cm.bone)  
colorbar()  
  
subplot(1,2,2)  
plot(x, C(x,1).view(ndarray), 'k-')
```

Figure: A contour plot titled "Non Parametric Regression" showing a covariance function. Below it is a line plot showing the same function.

Exploring R formula

Let's test that with a fake design matrix

```
In [4]: show_xkcd.random.normal(size=(10, 10), "My design")
```

Figure: A 10x10 grid titled "My design" showing a checkerboard pattern of black and white squares.

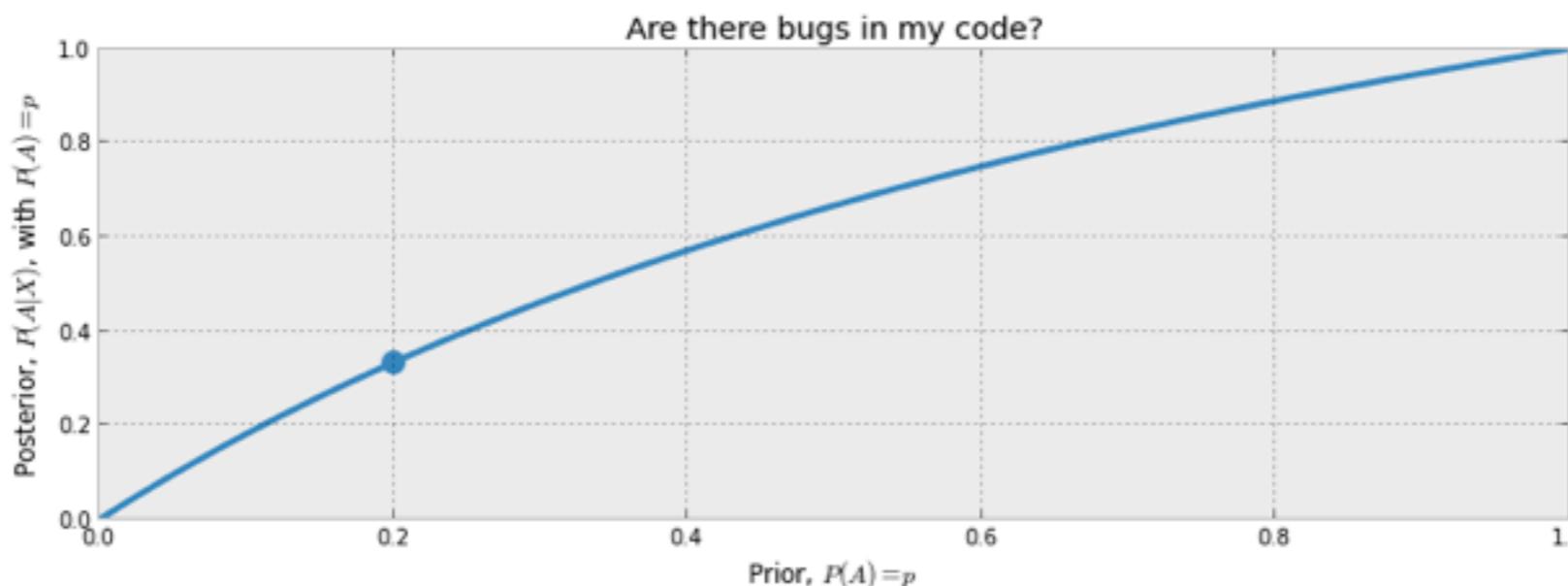
IPython is Great

(for large-scale computation, data exploration, and creating reproducible research artifacts)

$$\begin{aligned} P(A|X) &= \frac{1 \cdot p}{1 \cdot p + 0.5(1 - p)} \\ &= \frac{2p}{1 + p} \end{aligned}$$

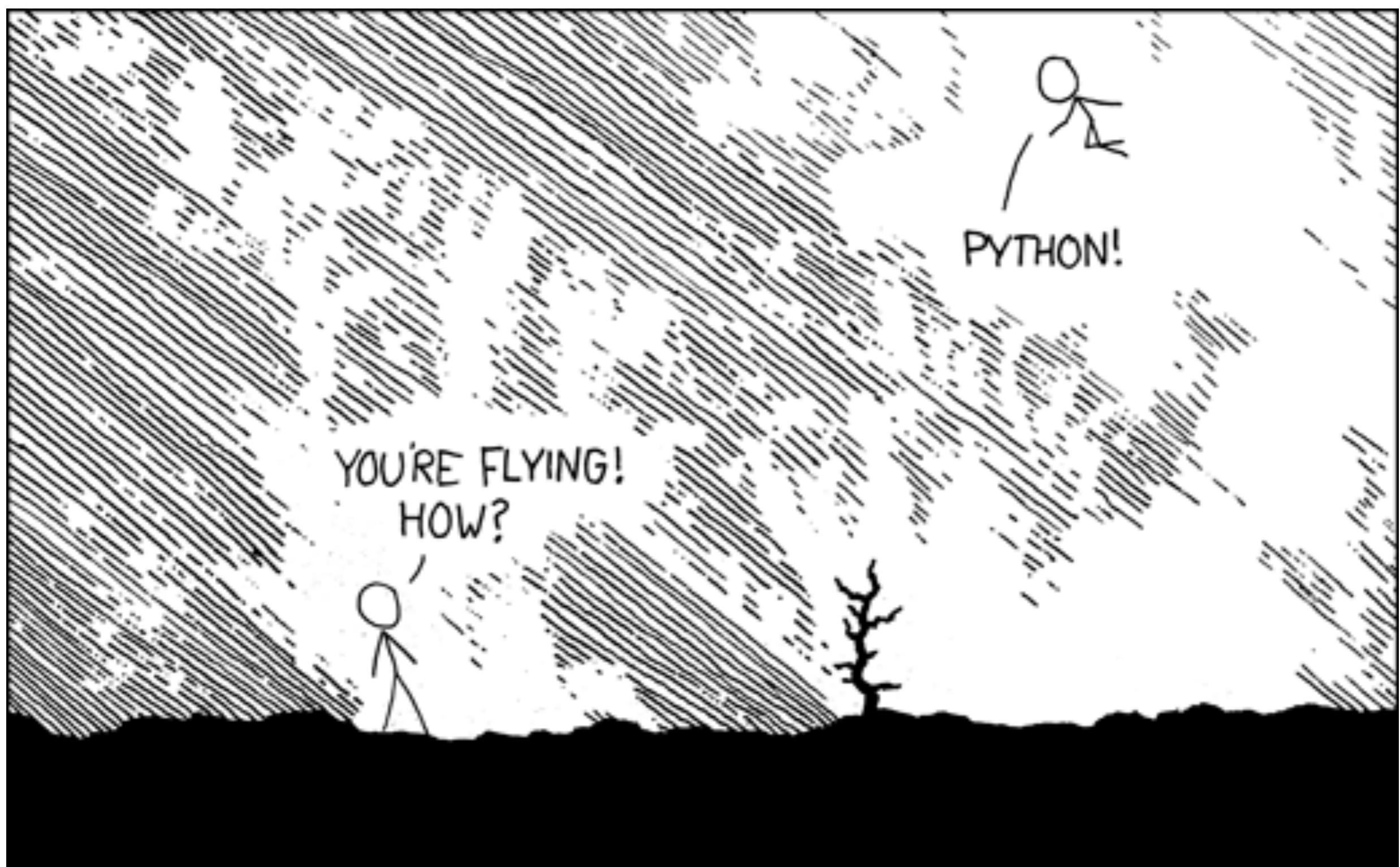
This is the posterior probability. What does it look like as a function of our prior, $p \in [0, 1]$?

```
figsize(12.5,4)
p = np.linspace( 0,1, 50)
plt.plot( p, 2*p/(1+p), color = "#348ABD", lw = 3 )
# plt.fill_between( p, 2*p/(1+p), alpha = .5, facecolor = ["#A60628"])
plt.scatter( 0.2, 2*(0.2)/1.2, s = 140, c ="#348ABD"  )
plt.xlim( 0, 1)
plt.ylim( 0, 1)
plt.xlabel( "Prior, $P(A) = p$")
plt.ylabel("Posterior, $P(A|X)$, with $P(A) = p$")
plt.title( "Are there bugs in my code?");
```



[http://nbviewer.ipython.org/urls/raw.github.com/CamDavidsonPilon/Probabilistic-Programming-and-Bayesian-Methods-for-Hackers/master/Chapter1_Introduction/Chapter1_Introduction.ipynb]

Python is great.



E.g.: <https://github.com/jrjohansson/scientific-python-lectures>

Home FAQ IPython Bookmarklet

Download Notebook

Introduction to Python programming

J.R. Johansson (robert@riken.jp) <http://dml.riken.jp/~rob/>

The latest version of this IPython notebook lecture is available at <http://github.com/jrjohansson/scientific-python-lectures>.

The other notebooks in this lecture series are indexed at <http://jrjohansson.github.com>.

Python program files

- Python code is usually stored in text files with the file ending ".py":
myprogram.py
- Every line in a Python program file is assumed to be a Python statement
 - The only exception is comment lines, which start with the character '#'. Comment lines are usually ignored by the Python interpreter.
- To run our Python program from the command line we use:

```
$ python myprogram.py
```

- On UNIX systems it is common to define the path to the interpreter on the first line of the program (note that this is a comment line as far as the Python interpreter is concerned):

```
#!/usr/bin/env python
```

If we do, and if we additionally set the file script to be executable, we can run the program like this:

```
$ myprogram.py
```

Example:

```
In [1]: ls scripts/hello-world*.py
```

```
scripts/hello-world-in-swedish.py  scripts/hello-world.py
```

**Intro to Python Lab
this Friday!!
10-11:30 am, MD G115**

~~You probably all know the default Python interpreter.~~

Don't bother with

```
Last login: Mon May 20 17:53:32 on ttys000
dn0a2100e6:~ mike$ python
Enthought Python Distribution -- www.enthought.com
Version: 7.3-2 (32-bit)

Python 2.7.3 |EPD 7.3-2 (32-bit)| (default, Apr 12 2012, 11:28:34)
[GCC 4.0.1 (Apple Inc. build 5493)] on darwin
Type "credits", "demo" or "enthought" for more information.
>>> 2 + 2
4
>>> █
```

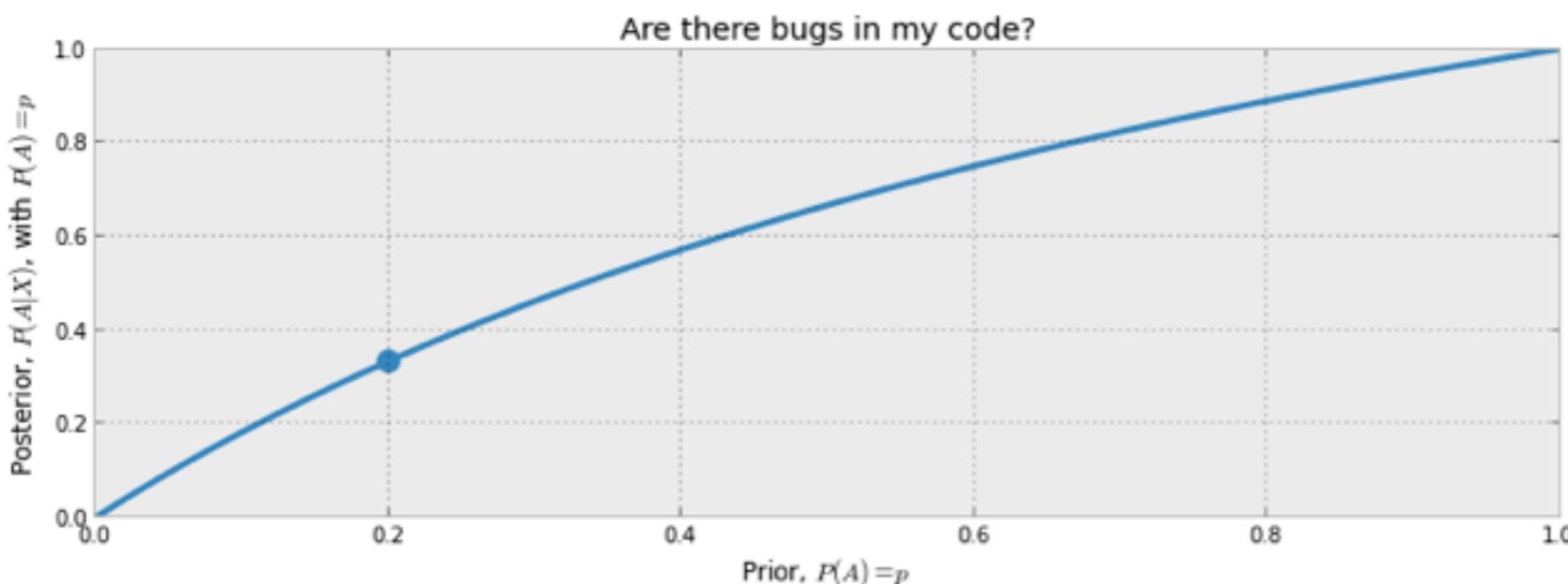
I Python is a more powerful interactive Python interpreter.

write and execute
Python code in
snippets

$$\begin{aligned} P(A|X) &= \frac{1 \cdot p}{1 \cdot p + 0.5(1 - p)} \\ &= \frac{2p}{1 + p} \end{aligned}$$

This is the posterior probability. What does it look like as a function of our prior, $p \in [0, 1]$?

```
figsize(12.5,4)
p = np.linspace( 0,1, 50)
plt.plot( p, 2*p/(1+p), color = "#348ABD", lw = 3 )
# plt.fill_between( p, 2*p/(1+p), alpha = .5, facecolor = ["#A60628"])
plt.scatter( 0.2, 2*(0.2)/1.2, s = 140, c ="#348ABD" )
plt.xlim( 0, 1)
plt.ylim( 0, 1)
plt.xlabel( "Prior, $P(A) = p$")
plt.ylabel("Posterior, $P(A|X)$, with $P(A) = p$")
plt.title( "Are there bugs in my code?");
```



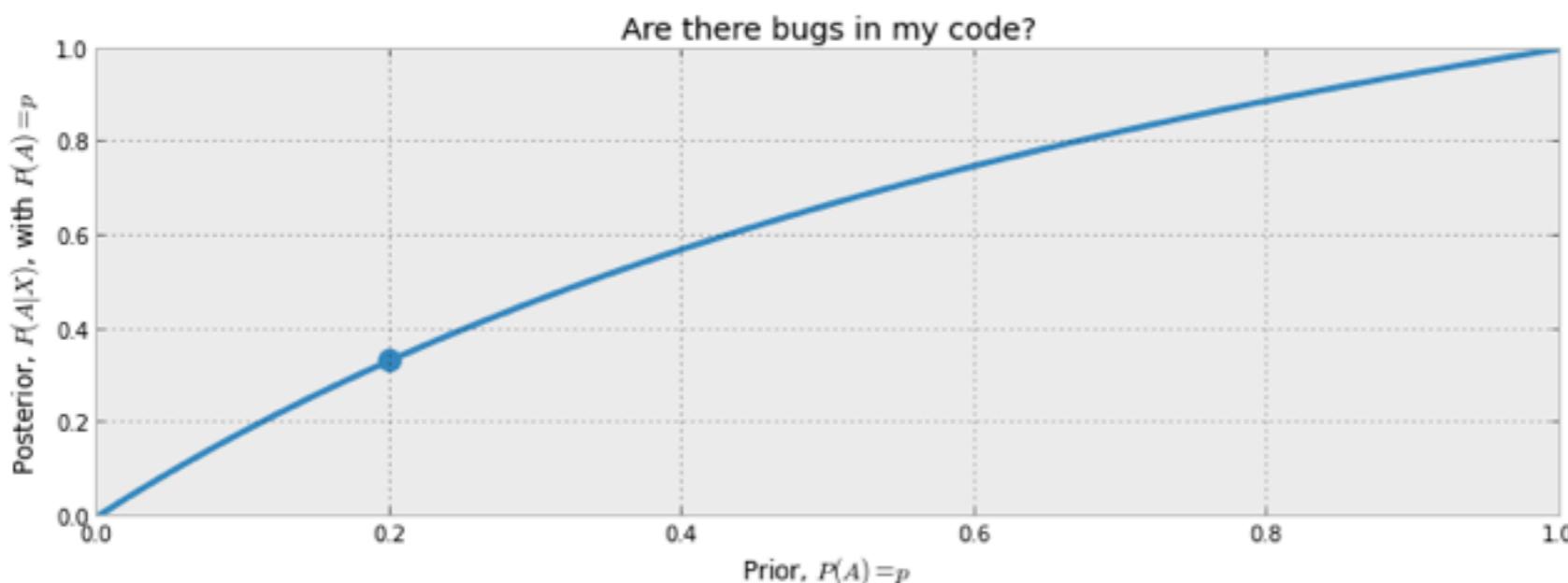
I Python is a more powerful interactive Python interpreter.

comments can include Latex math

$$\begin{aligned} P(A|X) &= \frac{1 \cdot p}{1 \cdot p + 0.5(1 - p)} \\ &= \frac{2p}{1 + p} \end{aligned}$$

This is the posterior probability. What does it look like as a function of our prior, $p \in [0, 1]$?

```
figsize(12.5,4)
p = np.linspace( 0,1, 50)
plt.plot( p, 2*p/(1+p), color = "#348ABD", lw = 3 )
# plt.fill_between( p, 2*p/(1+p), alpha = .5, facecolor = ["#A60628"])
plt.scatter( 0.2, 2*(0.2)/1.2, s = 140, c ="#348ABD" )
plt.xlim( 0, 1)
plt.ylim( 0, 1)
plt.xlabel( "Prior, $P(A) = p$")
plt.ylabel("Posterior, $P(A|X)$, with $P(A) = p$")
plt.title( "Are there bugs in my code?");
```



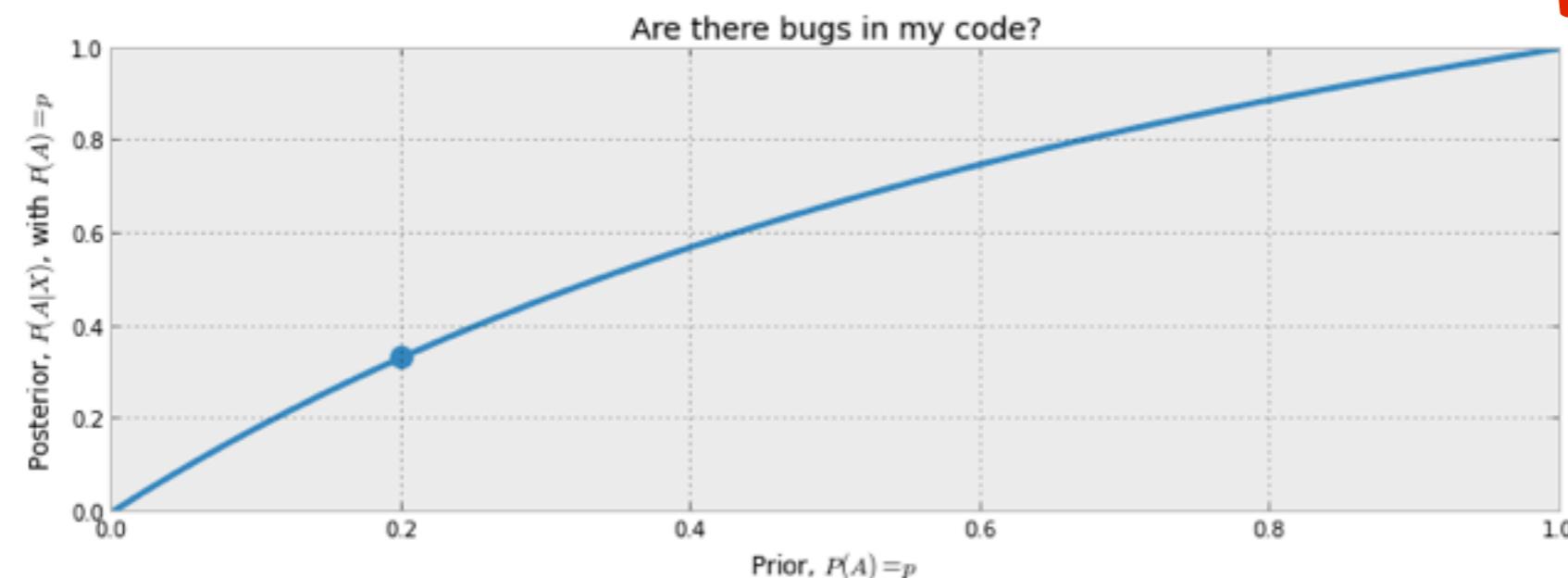
I Python is a more powerful interactive Python interpreter.

any plots generated by your code are displayed inline

$$\begin{aligned} P(A|X) &= \frac{1 \cdot p}{1 \cdot p + 0.5(1 - p)} \\ &= \frac{2p}{1 + p} \end{aligned}$$

This is the posterior probability. What does it look like as a function of our prior, $p \in [0, 1]$?

```
figsize(12.5,4)
p = np.linspace( 0,1, 50)
plt.plot( p, 2*p/(1+p), color = "#348ABD", lw = 3 )
#plt.fill_between( p, 2*p/(1+p), alpha = .5, facecolor = ["#A60628"])
plt.scatter( 0.2, 2*(0.2)/1.2, s = 140, c ="#348ABD" )
plt.xlim( 0, 1)
plt.ylim( 0, 1)
plt.xlabel( "Prior, $P(A) = p$")
plt.ylabel("Posterior, $P(A|X)$, with $P(A) = p$")
plt.title( "Are there bugs in my code?");
```



Write Python code interactively in a web browser instead of a terminal window.

The screenshot shows a web-based IPython Notebook interface. At the top, there's a toolbar with various icons and tabs like 'IPy IPython Dashboard' and 'IPy HW2 (Gaussian Blur)'. Below the toolbar, the address bar shows the URL `127.0.0.1:8888/d795d6fe-e6bb-44a4-a11c-f947e19a2560`. The main content area has tabs for 'Announcements', 'Matrix calculus', 'Differentiation rules', 'Lagrange multiplier', 'Likelihood function', 'List of logarithmic id', 'List of trigonometric', 'G3D Innovation Eng', and 'scikit-learn: machine'. The notebook title is 'IP[y]: Notebook HW2 (Gaussian Blur)' with a note 'Last saved: May 21 11:24 AM'. A menu bar includes 'File', 'Edit', 'View', 'Insert', 'Cell', 'Kernel', and 'Help'. Below the menu is a toolbar with icons for cell operations. The text area contains the following content:

common gaussian blur code

We begin by defining the non-normalized Gaussian Function $G(x,y)$ as follows:

$$G(x,y) = e^{-\frac{x^2+y^2}{2\sigma^2}}$$

We then define the normalized Gaussian Function as $N(x,y) = cG(x,y)$, where c is a normalization constant:

$$c = \frac{1}{\iint G(x,y) dx dy}$$

To perform a Gaussian Blur, we use $N(x,y)$ as a convolution kernel.

```
In [4]: import scipy
gaussian_blur_kernel_width      = numpy.int32(9)
gaussian_blur_kernel_half_width = numpy.int32(4)
gaussian_blur_sigma            = numpy.float32(2)

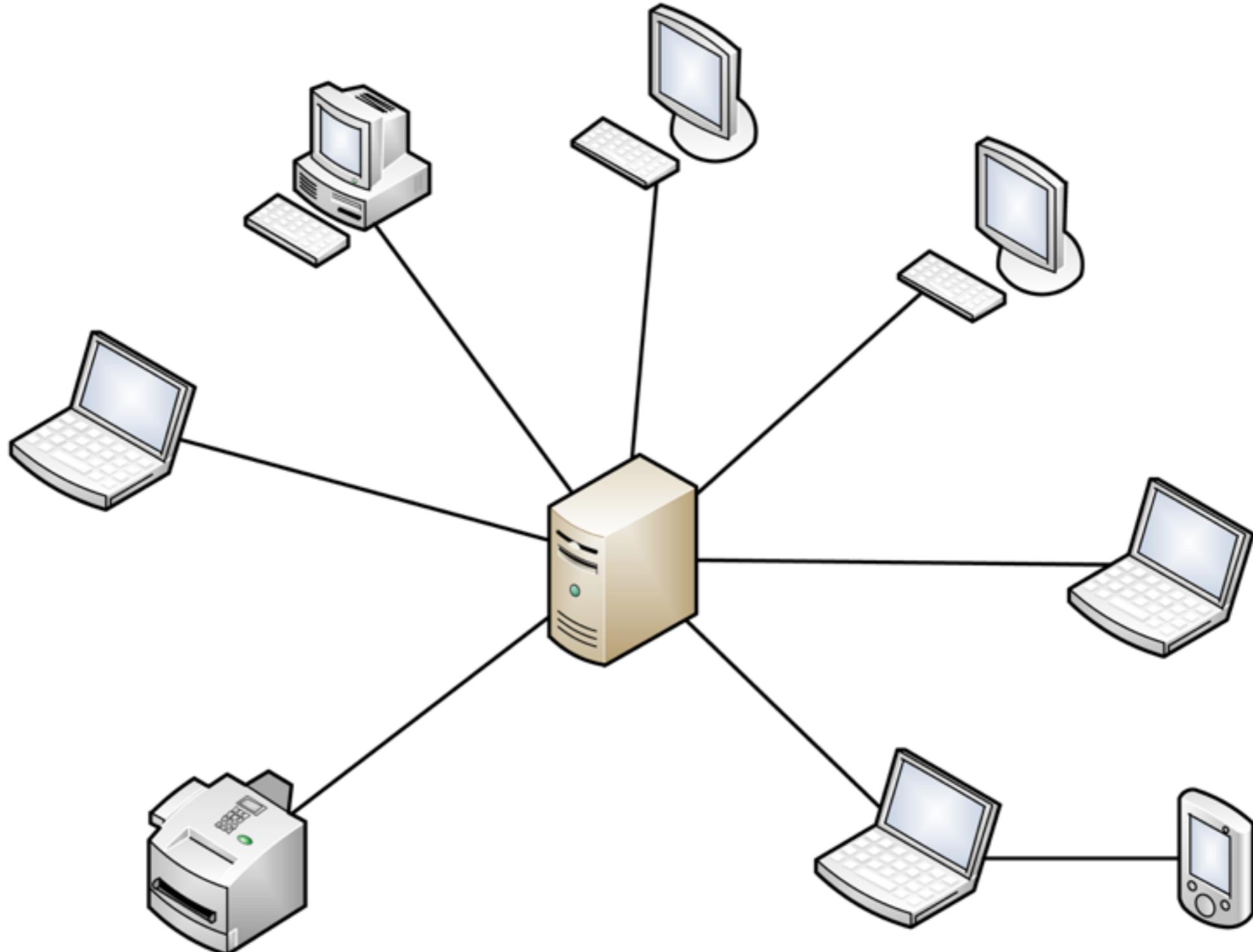
y, x = \
    scipy.mgrid[-gaussian_blur_kernel_half_width:gaussian_blur_kernel_half_width+1,-gaussian_blur_kernel_half_width:gaussian_blur_kernel_half_width+1]

gaussian_blur_kernel_not_normalized = numpy.exp( ( - ( x**2 + y**2 ) ) / ( 2 * gaussian_blur_sigma**2 ) )
normalization_constant           = numpy.float32(1) / numpy.sum(gaussian_blur_kernel_not_normalized)
gaussian_blur_kernel             = (normalization_constant * gaussian_blur_kernel_not_normalized).astype(numpy.float32)

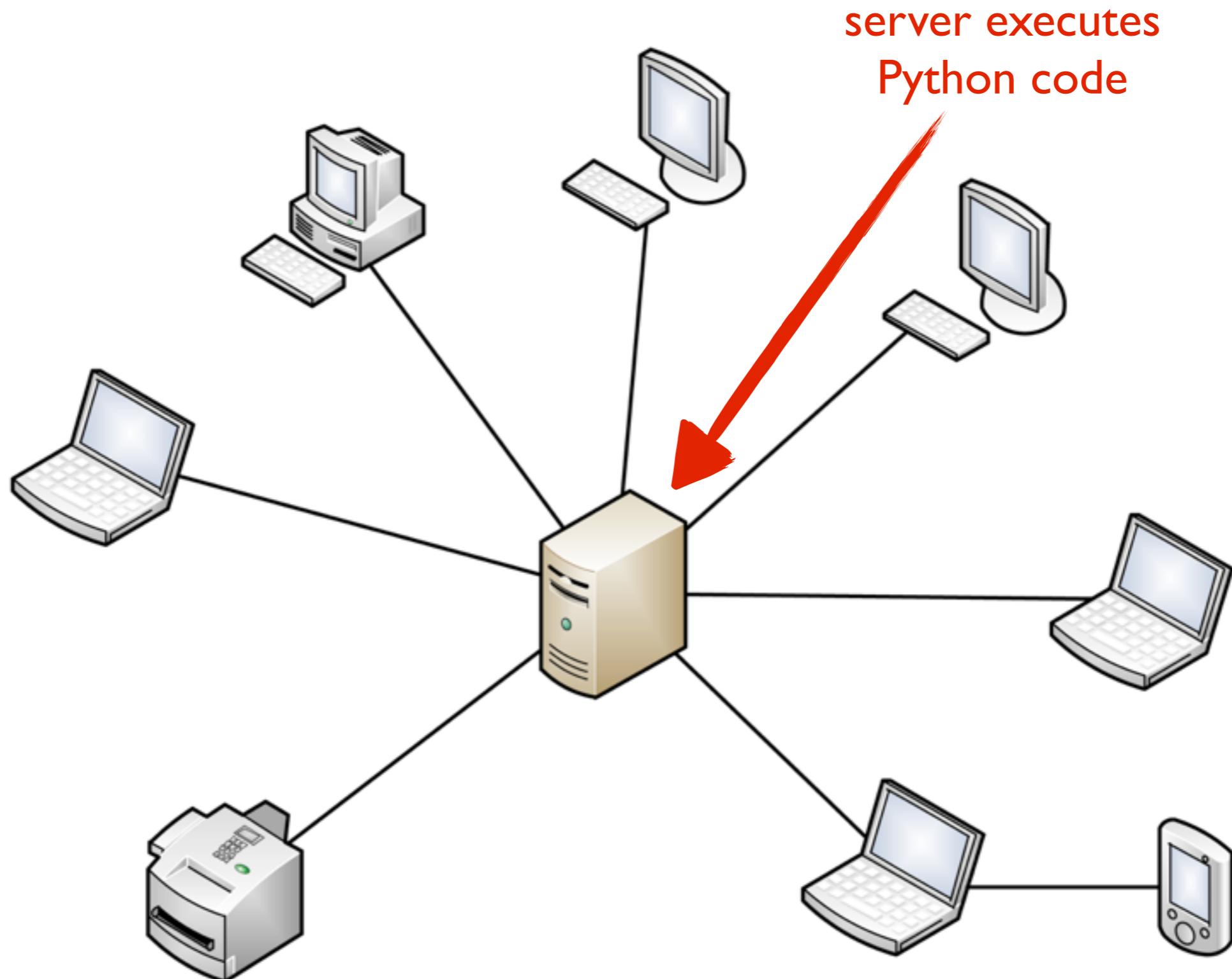
matplotlib.pyplot.imshow(gaussian_blur_kernel, cmap="gray", interpolation="nearest");
matplotlib.pyplot.title("gaussian_blur_kernel");
matplotlib.pyplot.colorbar();
```

The output of the code is a heatmap titled "gaussian_blur_kernel". The x-axis ranges from 0 to 8, and the y-axis ranges from 0 to 8. The color scale on the right ranges from 0.005 (dark) to 0.040 (light). The heatmap shows a symmetric bell-shaped distribution centered at (4, 4).

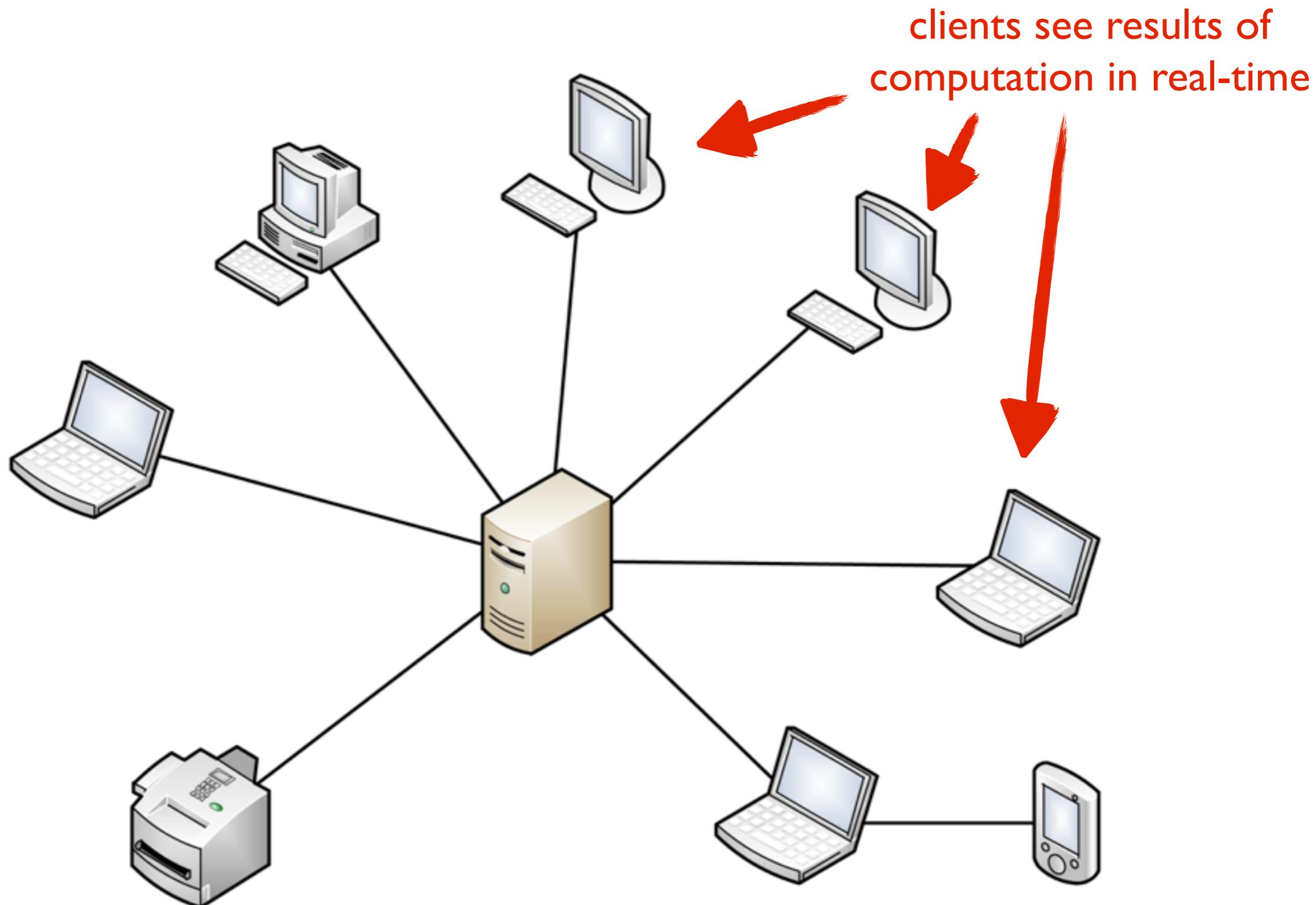
IPython has a decoupled client-server architecture.



I^{Python} has a decoupled client-server architecture.



IPython has a decoupled client-server architecture.

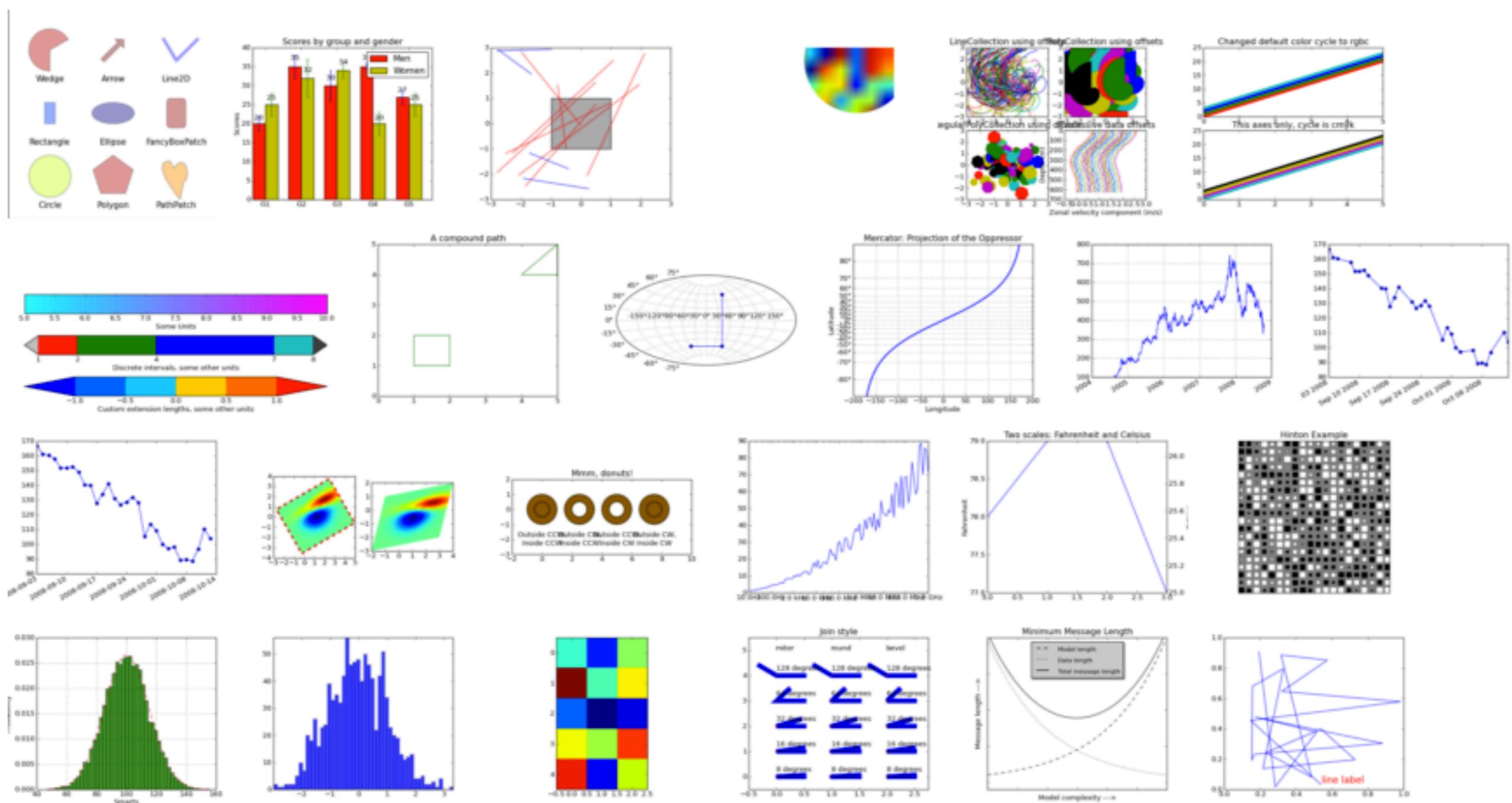


**You can stay productive on any computer with
an internet connection and a web browser.**



[image source unknown]

IPython integrates seamlessly with Matplotlib, making it well-suited for data exploration.



A woman in a vibrant green, flowing, futuristic outfit is captured in a dynamic, mid-motion pose. She is leaning forward with her left leg extended back and her right leg bent at the knee. Her arms are outstretched to the sides, creating a sense of balance and movement. The outfit features intricate, swirling patterns and metallic accents. The background is a soft, gradient blue, suggesting a sky or water environment.

Examples

Resources

- <http://nbviewer.ipython.org/4542975>
- [https://github.com/mroberts3000/
IPythonIsGreat](https://github.com/mroberts3000/IPythonIsGreat)
- <http://nbviewer.ipython.org/>
- [https://github.com/jrjohansson/scientific-python-
lectures](https://github.com/jrjohansson/scientific-python-lectures)
- <http://scipy-lectures.github.io/index.html>
- Many more online and on the course web site

Good Practices

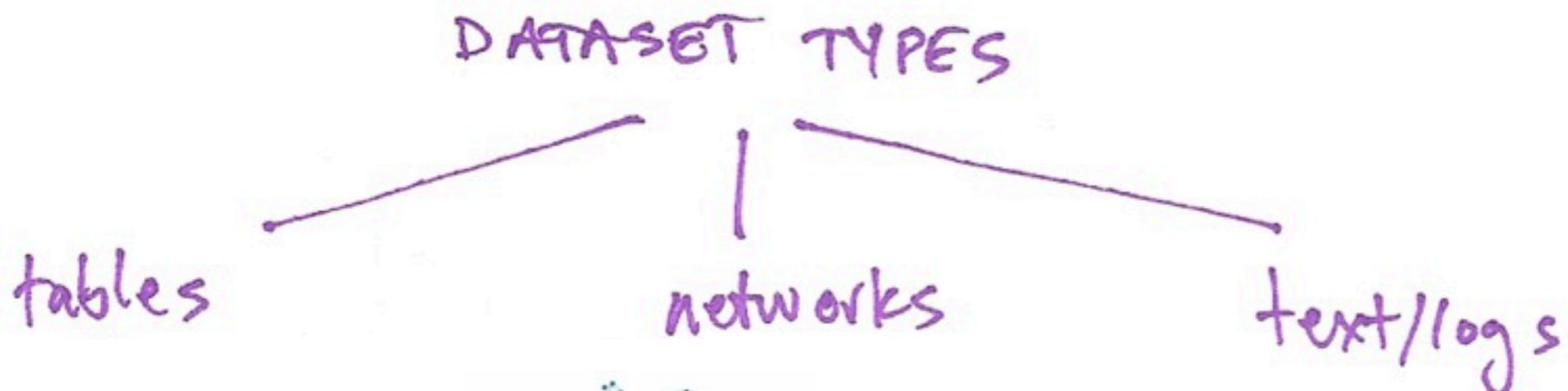
- IPython notebooks to document your process
- Visualizations for data exploration
- Comment your code!
- Modularity - breaking down code into small functional, composable pieces
- Array-oriented computing
- Using assert statements and tests
- Version control (svn, git, github)

Data Types

Ben Shneiderman, 1996

- 1D (sequences)
- Temporal
- 2D (maps)
- 3D (shaped)
- nD (relational)
- Trees (hierarchical)
- Networks (graphs)
- Others?

Tamara Munzner, 2013



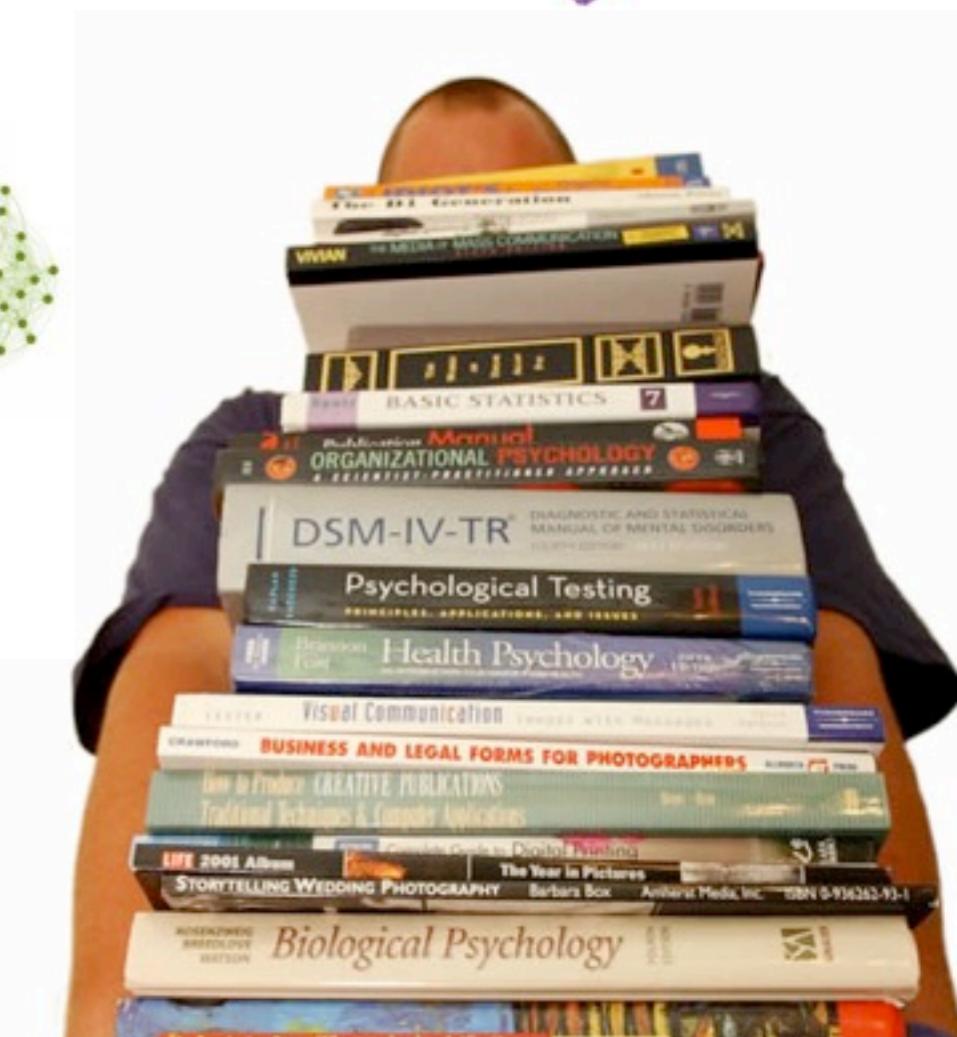
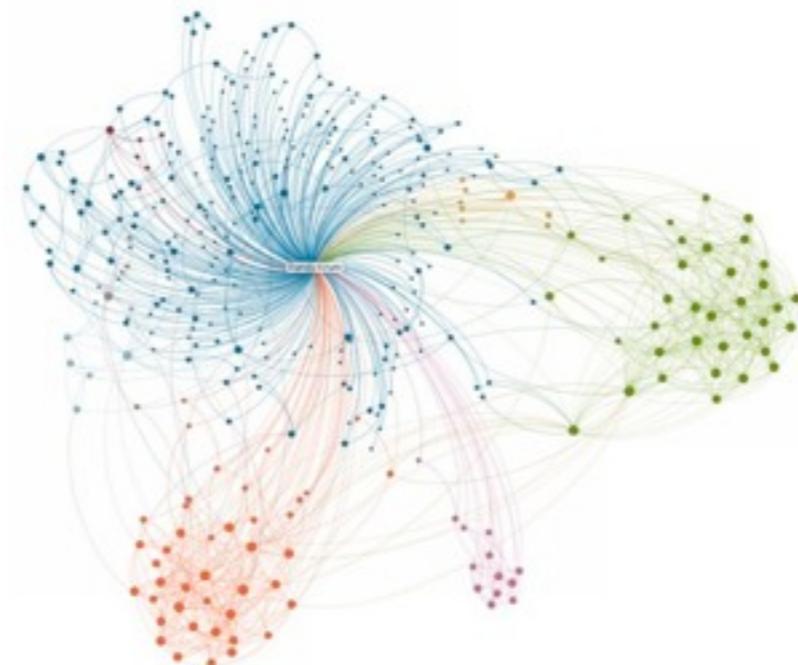
Google Docs

@gmail.com | New features | Docs Home

FriendFeed Audience

Share | Autosaved on 10:35 AM

	Site Name	Category	Compositor	Unique	Use	Country	Re	Page	Views	Google
1	friendfeed.com	/Online Cor	160000	150000	0.1	2000000				
2	twhirl.org	/Computer	47000	43000	0	74000				
3	tweetscan.com	/Online Cor	43000	18000	0	120000				
4	chrisbrogan.com	/Online Cor	39000	29000	0	74000				
5	brightkite.com	/Telecomm	29000	68000	0	910000				
6	twitpic.com	/Home & Gi	24000	71000	0	340000	TRUE			
7	web-strategist.c	/Online Cor	24000	32000	0	86000				
8	summize.com	/Arts & Hur	20000	54000	0	570000				



Semantics vs. Types

- Data Semantics: The real-world meaning
 - e.g., company name, day of the month, person height, etc.
- Data Type: Interpretation in terms of scales of measurements
 - e.g., quantity or category, sensible mathematical operations, data structure, etc.

SCIENCE

Vol. 103, No. 2684

Friday, June 7, 1946

On the Theory of Scales of Measurement

S. S. Stevens

Director, Psycho-Acoustic Laboratory, Harvard University

FOR SEVEN YEARS A COMMITTEE of the British Association for the Advancement of Science debated the problem of measurement. Appointed in 1932 to represent Section A (Mathematical and Physical Sciences) and Section J (Psychology), the committee was instructed to consider and report upon the possibility of "quantitative estimates of sensory events"—meaning simply: Is it possible to measure human sensation? Deliberation led only to disagreement, mainly about what is meant by the term measurement. An interim report in 1938 found one member complaining that his colleagues

by the formal (mathematical) properties of the scales. Furthermore—and this is of great concern to several of the sciences—the statistical manipulations that can legitimately be applied to empirical data depend upon the type of scale against which the data are ordered.

A CLASSIFICATION OF SCALES OF MEASUREMENT

Paraphrasing N. R. Campbell (Final Report, p. 340), we may say that measurement, in the broadest sense, is defined as the assignment of numerals to objects or events according to rules. The fact that numerals can be assigned under different rules leads

Scale	Basic Empirical Operations	Mathematical Group Structure	Permissible Statistics (invariantive)
Nominal Categorical Qualitative	Determination of equality	<i>Permutation group</i> $x' = f(x)$ $f(x)$ means any one-to-one substitution	Number of cases Mode Contingency correlation
Ordinal	Determination of greater or less	<i>Isotonic group</i> $x' = f(x)$ $f(x)$ means any monotonic increasing function	Median Percentiles
Interval	Determination of equality of intervals or differences	<i>General linear group</i> $x' = ax + b$	Mean Standard deviation Rank-order correlation Product-moment correlation
Ratio	Determination of equality of ratios	<i>Similarity group</i> $x' = ax$	Coefficient of variation

Data Types

- Nominal (Categorical) (N)

Are = or \neq to other values

Apples, Oranges, Bananas,...



- Ordinal (O)

Obey a $<$ relationship

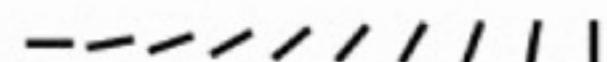
Small, medium, large



- Quantitative (Q)

Can do arithmetic on them

10 inches, 23 inches, etc.



Data Types

- Q - Interval (location of zero arbitrary)

Dates: Jan 19; Location: (Lat, Long)

Like a geometric point. Cannot compare directly.

Only differences (i.e., intervals) can be compared

- Q - Ratio (zero fixed)

Measurements: Length, Mass, Temp, ...

Origin is meaningful, can measure ratios & proportions

Like a geometric vector, origin is meaningful

Data Types

- N - Nominal (labels)
Operations: $=, \neq$
- O - Ordinal (ordered)
Operations: $=, \neq, >, <$
- Q - Interval (location of zero arbitrary)
Operations: $=, \neq, >, <, +, -$ (distance)
- Q - Ratio (zero fixed)
Operations: $=, \neq, >, <, +, -, \times, \div$ (proportions)

A	B	C	S	T	U
Order ID	Order Date	Order Priority	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low	Large Box	0.8	10/21/06
6	2/21/08	4-Not Specified	Small Pack		2/22/08
32	7/16/07	2-High	Small Pack		7/17/07
32	7/16/07	2-High	Jumbo Box	0.72	7/17/07
32	7/16/07	2-High	Medium Box	0.6	7/18/07
32	7/16/07	2-High	Medium Box	0.65	7/18/07
35	10/23/07	4-Not Specified	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Specified	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent	Small Box	0.55	11/3/07
65	3/18/07	1-Urgent	Small Pack	0.49	3/19/07
66	1/20/05	5-Low	Wrap Bag	0.56	1/20/05
69	6/4/05	4-Not Specified	Small Pack	0.44	6/6/05
69	6/4/05	4-Not Specified	Wrap Bag	0.6	6/6/05
70	12/18/06	5-Low	Small Box	0.59	12/23/06
70	12/18/06	5-Low	Wrap Bag	0.82	12/23/06
96	4/17/05	2-High	Small Box	0.55	4/19/05
97	1/29/06	3-Medium	Small Box	0.38	1/30/06
129	11/19/08	5-Low	Small Box	0.37	11/28/08
130	5/8/08	2-High	Small Box	0.37	5/9/08
130	5/8/08	2-High	Medium Box	0.38	5/10/08
130	5/8/08	2-High	Small Box	0.6	5/11/08
132	6/11/06	3-Medium	Medium Box	0.6	6/12/06
132	6/11/06	3-Medium	Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Specified	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified	Small Pack	0.64	10/23/07
166	9/12/07	2-High	Small Box	0.55	9/14/07
193	8/8/06	1-Urgent	Medium Box	0.57	8/10/06
194	4/5/08	3-Medium	Wrap Bag	0.42	4/7/08

A	B	C	S	T	U
Order ID	Order Date	Order Priority	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low	Large Box	0.8	10/21/06
6	2/21/08	4-Not Specified	Small Pack	0.55	2/22/08
32	7/16/07	2-High	Small Pack	0.79	7/17/07
32	7/16/07	2-High	Jumbo Box		7/17/07
32	7/16/07	2-High	Medium Box		7/18/07
32	7/16/07	2-High	Medium Box	0.65	7/18/07
35	10/23/07	4-Not Specified	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Specified	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent	Small Box	0.55	11/3/07
65	3/18/07	1-Urgent	Small Pack	0.49	3/19/07
66	1/20/05	5-Low	Wrap Bag	0.56	1/20/05
69	6/4/05	4-Not Specified	Small Pack	0.44	6/6/05
69	6/4/05	4-Not Specified	Wrap Bag	0.6	6/6/05
70	12/18/06	5-Low	Small Box	0.59	12/23/06
70	12/18/06	5-Low	Wrap Bag	0.82	12/23/06
96	4/17/05	2-High	Small Box	0.55	4/19/05
97	1/29/06	3-Medium	Small Box	0.38	1/30/06
129	11/19/08	5-Low	Small Box	0.37	11/28/08
130	5/8/08	2-High	Small Box	0.37	5/9/08
130	5/8/08	2-High	Medium Box	0.38	5/10/08
130	5/8/08	2-High	Small Box	0.6	5/11/08
132	6/11/06	3-Medium	Medium Box	0.6	6/12/06
132	6/11/06	3-Medium	Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Specified	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified	Small Pack	0.64	10/23/07
166	9/12/07	2-High	Small Box	0.55	9/14/07
193	8/8/06	1-Urgent	Medium Box	0.57	8/10/06
194	4/5/08	3-Medium	Wrap Bag	0.42	4/7/08

A	B	C	S	T	U
Order ID	Order Date	Order Priority	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low	Large Box	0.8	10/21/06
6	2/21/08	4-Not Specified	Small Pack		2/22/08
32	7/16/07	2-High	Small Pack		7/17/07
32	7/16/07	2-High	Jumbo Box		7/17/07
32	7/16/07	2-High	Medium Box		7/18/07
32	7/16/07	2-High	Medium Box		7/18/07
35	10/23/07	4-Not Specified	Wrap Bag		10/24/07
35	10/23/07	4-Not Specified	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent	Small Box	0.55	11/3/07
65	3/18/07	1-Urgent	Small Pack	0.49	3/19/07
66	1/20/05	5-Low	Wrap Bag	0.56	1/20/05
69	6/4/05	4-Not Specified	Small Pack	0.44	6/6/05
69	6/4/05	4-Not Specified	Wrap Bag	0.6	6/6/05
70	12/18/06	5-Low	Small Box	0.59	12/23/06
70	12/18/06	5-Low	Wrap Bag	0.82	12/23/06
96	4/17/05	2-High	Small Box	0.55	4/19/05
97	1/29/06	3-Medium	Small Box	0.38	1/30/06
129	11/19/08	5-Low	Small Box	0.37	11/28/08
130	5/8/08	2-High	Small Box	0.37	5/9/08
130	5/8/08	2-High	Medium Box	0.38	5/10/08
130	5/8/08	2-High	Small Box	0.6	5/11/08
132	6/11/06	3-Medium	Medium Box	0.6	6/12/06
132	6/11/06	3-Medium	Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Specified	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified	Small Pack	0.64	10/23/07
166	9/12/07	2-High	Small Box	0.55	9/14/07
193	8/8/06	1-Urgent	Medium Box	0.57	8/10/06
194	4/5/08	3-Medium	Wrap Bag	0.42	4/7/08

Attribute
aka
Feature

A	B	C	S	T	U
Order ID	Order Date	Order Priority	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low	Large Box	0.8	10/21/06
6	2/21/08	4-Not Specified	Small Pack	0.55	2/22/08
32	7/16/07	2-High	Small Pack	0.79	7/17/07
32	7/16/07	2-High	Jumbo Box	0.72	7/17/07
32	7/16/07	2-High	Medium Box	0.6	7/18/07
32	7/16/07	2-High	Medium Box	0.65	7/18/07
35	10/23/07	4-Not Specified	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Specified	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent	Small Box	0.55	11/3/07
65	3/18/07	1-Urgent	Small Pack	0.49	3/19/07
66	1/20/05	5-Low	Wrap Bag	0.56	1/20/05
69	6/4/05	4-Not Specified	Small Pack	0.44	6/6/05
69	6/4/05	4-Not Specified		0.6	6/6/05
70	12/18/06	5-Low		0.59	12/23/06
70	12/18/06	5-Low		0.82	12/23/06
96	4/17/05	2-High		0.55	4/19/05
97	1/29/06	3-Medium		0.38	1/30/06
129	11/19/08	5-Low	Small Box	0.37	11/28/08
130	5/8/08	2-High	Small Box	0.37	5/9/08
130	5/8/08	2-High	Medium Box	0.38	5/10/08
130	5/8/08	2-High	Small Box	0.6	5/11/08
132	6/11/06	3-Medium	Medium Box	0.6	6/12/06
132	6/11/06	3-Medium	Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Specified	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified	Small Pack	0.64	10/23/07
166	9/12/07	2-High	Small Box	0.55	9/14/07
193	8/8/06	1-Urgent	Medium Box	0.57	8/10/06
194	4/5/08	3-Medium	Wrap Bag	0.42	4/7/08

I = Quantitative
2 = Nominal
3 = Ordinal

A	B	C	S	T	U
Order ID	Order Date	Order Priority	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low	Large Box	0.8	10/21/06
6	2/21/08	4-Not Specified	Small Pack	0.55	2/22/08
32	7/16/07	2-High	Small Pack	0.79	7/17/07
32	7/16/07	2-High	Jumbo Box	0.72	7/17/07
32	7/16/07	2-High	Medium Box	0.6	7/18/07
32	7/16/07	2-High	Medium Box	0.65	7/18/07
35	10/23/07	4-Not Specified	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Specified	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent	Small Box	0.55	11/3/07
65	3/18/07	1-Urgent	Small Pack	0.49	3/19/07
66	1/20/05	5-Low	Wrap Bag	0.56	1/20/05
69	6/4/05	4-Not Specified	Small Pack	0.44	6/6/05
69	6/4/05	4-Not Specified		0.6	6/6/05
70	12/18/06	5-Low		0.59	12/23/06
70	12/18/06	5-Low		0.82	12/23/06
96	4/17/05	2-High		0.55	4/19/05
97	1/29/06	3-Medium		0.38	1/30/06
129	11/19/08	5-Low	Small Box	0.37	11/28/08
130	5/8/08	2-High	Small Box	0.37	5/9/08
130	5/8/08	2-High	Medium Box	0.38	5/10/08
130	5/8/08	2-High	Small Box	0.6	5/11/08
132	6/11/06	3-Medium	Medium Box	0.6	6/12/06
132	6/11/06	3-Medium	Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Specified	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified	Small Pack	0.64	10/23/07
166	9/12/07	2-High	Small Box	0.55	9/14/07
193	8/8/06	1-Urgent	Medium Box	0.57	8/10/06
194	4/5/08	3-Medium	Wrap Bag	0.42	4/7/08
194	4/5/08	3-Medium	Wrap Bag	0.84	4/7/08

I = Quantitative
2 = Nominal
3 = Ordinal

A	B	C	S	T	U
Order ID	Order Date	Order Priority	Product Container	Product Base Margin	Ship Date
3	10/14/06	5-Low	Large Box	0.8	10/21/06
6	2/21/08	4-Not Specified	Small Pack	0.55	2/22/08
32	7/16/07	2-High	Small Pack	0.79	7/17/07
32	7/16/07	2-High	Jumbo Box	0.72	7/17/07
32	7/16/07	2-High	Medium Box	0.6	7/18/07
32	7/16/07	2-High	Medium Box	0.65	7/18/07
35	10/23/07	4-Not Specified	Wrap Bag	0.52	10/24/07
35	10/23/07	4-Not Specified	Small Box	0.58	10/25/07
36	11/3/07	1-Urgent	Small Box	0.55	11/3/07
65	3/18/07	1-Urgent	Small Pack	0.49	3/19/07
66	1/20/05	5-Low	Wrap Bag	0.56	1/20/05
69	6/4/05			44	6/6/05
69	6/4/05			0.6	6/6/05
70	12/18/06			59	12/23/06
70	12/18/06			82	12/23/06
96	4/17/05			55	4/19/05
97	1/29/06			38	1/30/06
129	11/19/08			37	11/28/08
130	5/8/08	2-High	Small Box	0.37	5/9/08
130	5/8/08	2-High	Medium Box	0.38	5/10/08
130	5/8/08	2-High	Small Box	0.6	5/11/08
132	6/11/06	3-Medium	Medium Box	0.6	6/12/06
132	6/11/06	3-Medium	Jumbo Box	0.69	6/14/06
134	5/1/08	4-Not Specified	Large Box	0.82	5/3/08
135	10/21/07	4-Not Specified	Small Pack	0.64	10/23/07
166	9/12/07	2-High	Small Box	0.55	9/14/07
193	8/8/06	1-Urgent	Medium Box	0.57	8/10/06
194	4/5/08	3-Medium	Wrap Bag	0.42	4/7/08
194	4/5/08	3-Medium	Wrap Bag	0.84	4/7/08

Nominal /Ordinal = Dimensions

Describe the data, independent variables

Quantitative = Measures

Numbers to be analyzed, dependent variables

Data vs. Conceptual Model

- Data Model: Low-level description of the data
Set with operations, e.g., floats with +, -, /, *
- Conceptual Model: Mental construction
Includes semantics, supports reasoning

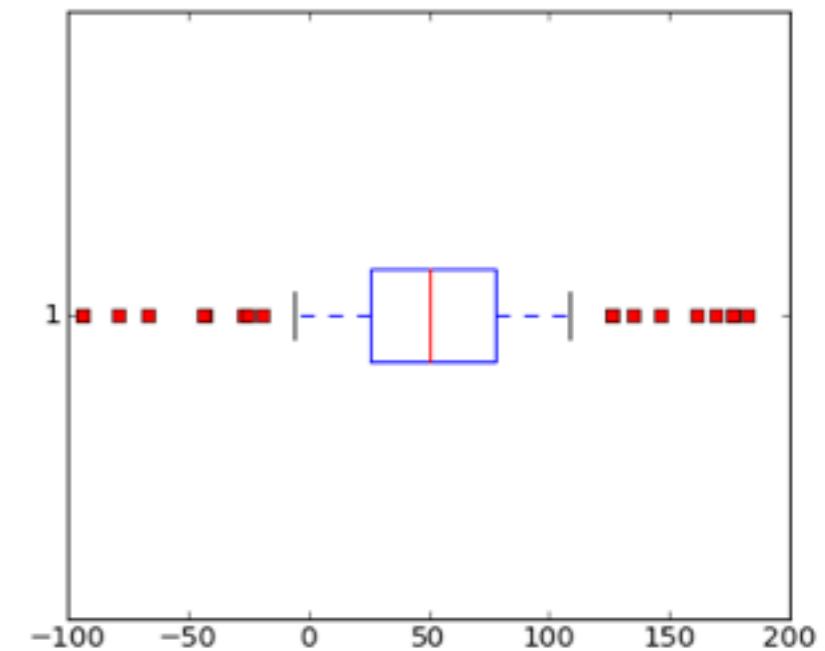
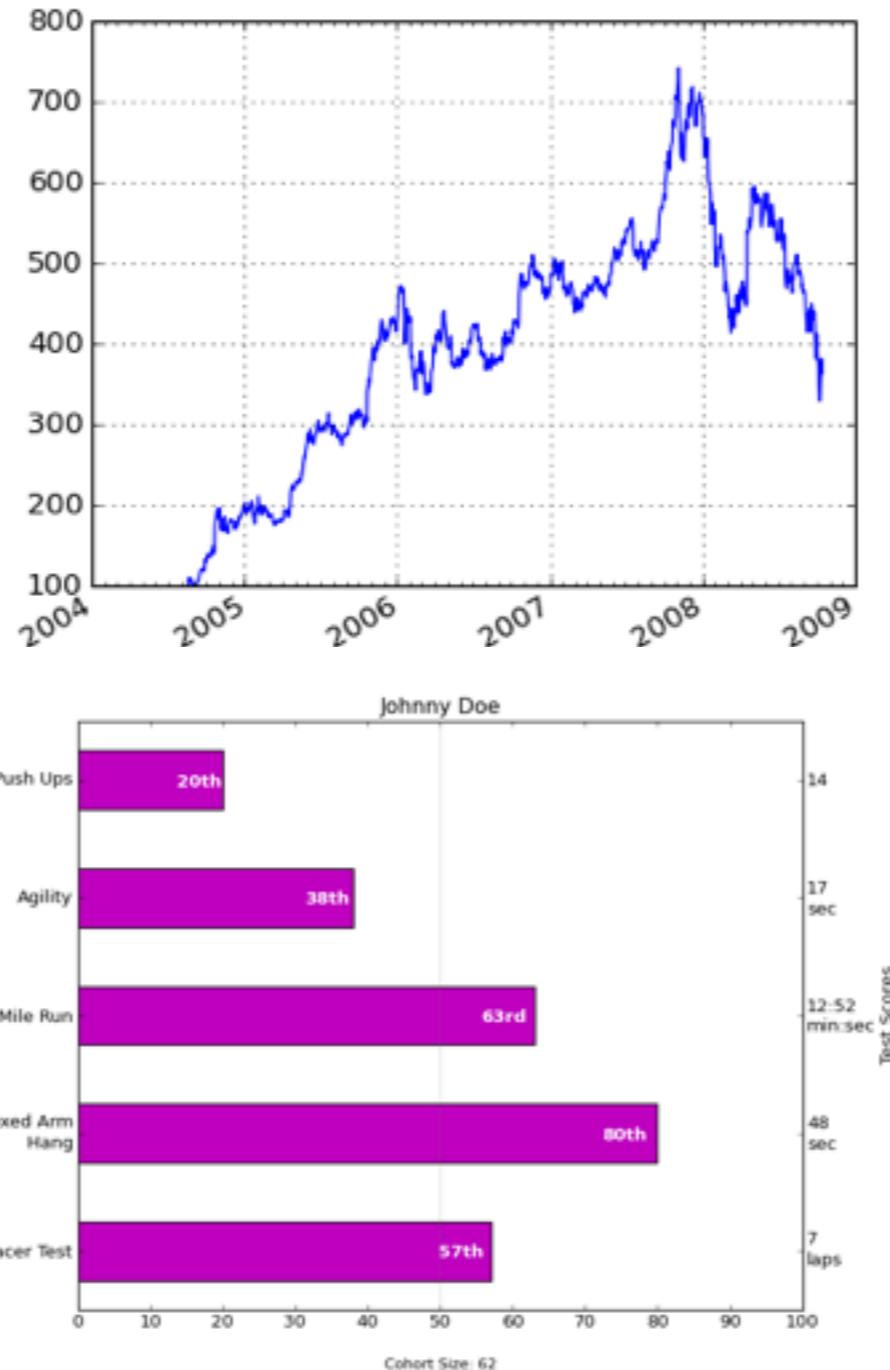
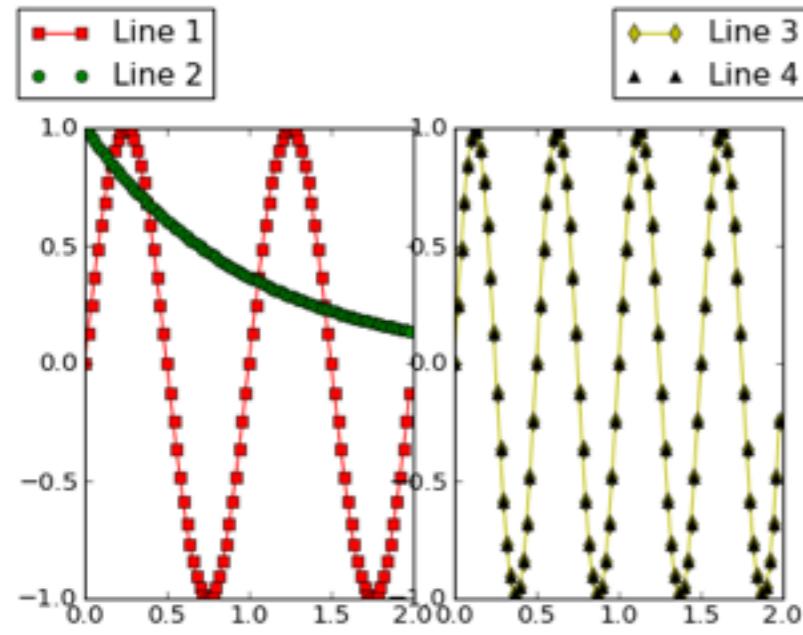
Data	Conceptual
1D floats	temperature
3D vector of floats	space

Data vs. Conceptual Model

- From data model...
32.5, 54.0, -17.3, ... (floats)
- using conceptual model...
Temperature
- to data type
Continuous to 4 significant figures (Q)
Hot, warm, cold (O)
Burned vs. Not burned (N)

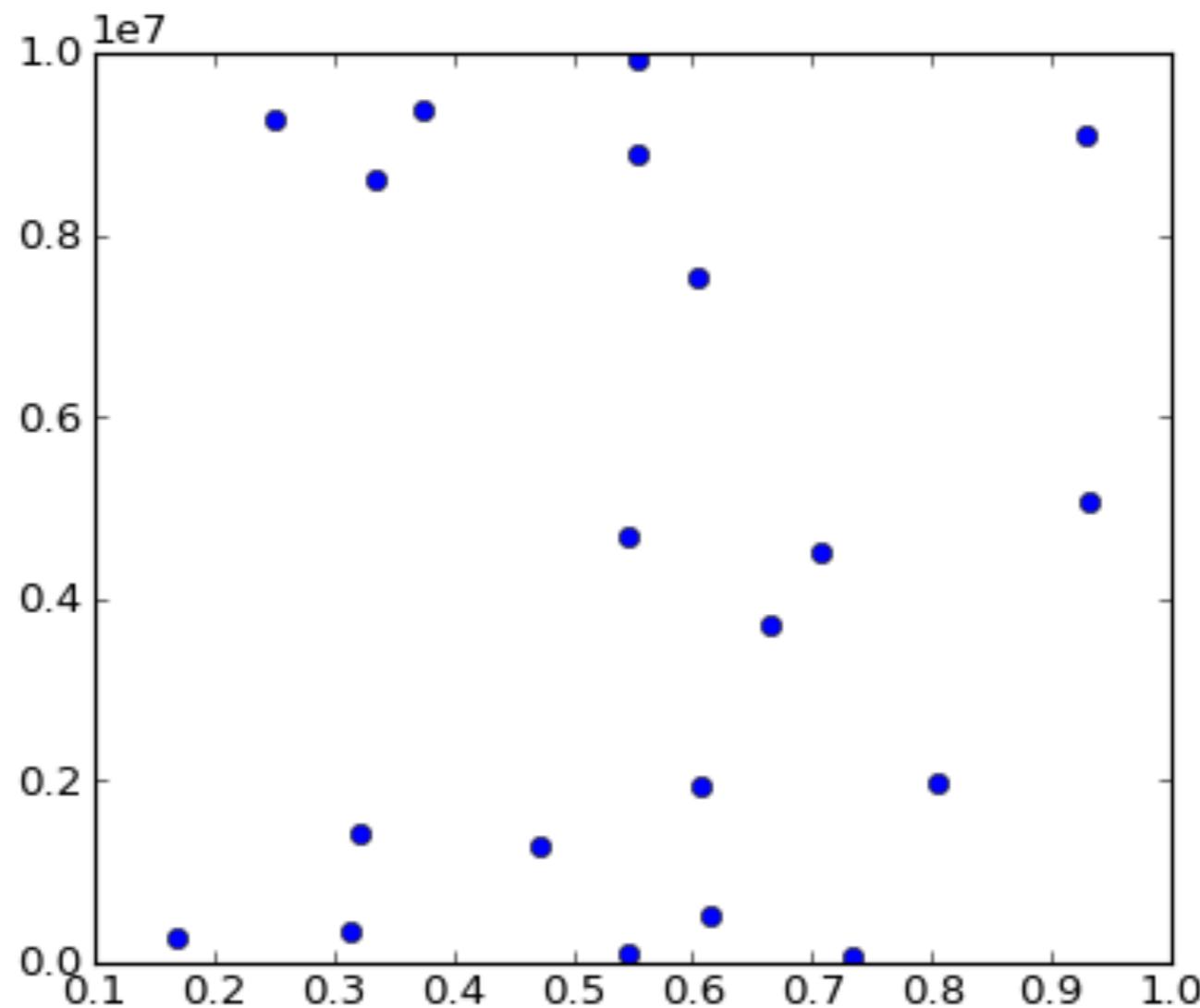
Data Dimensions

Univariate Data



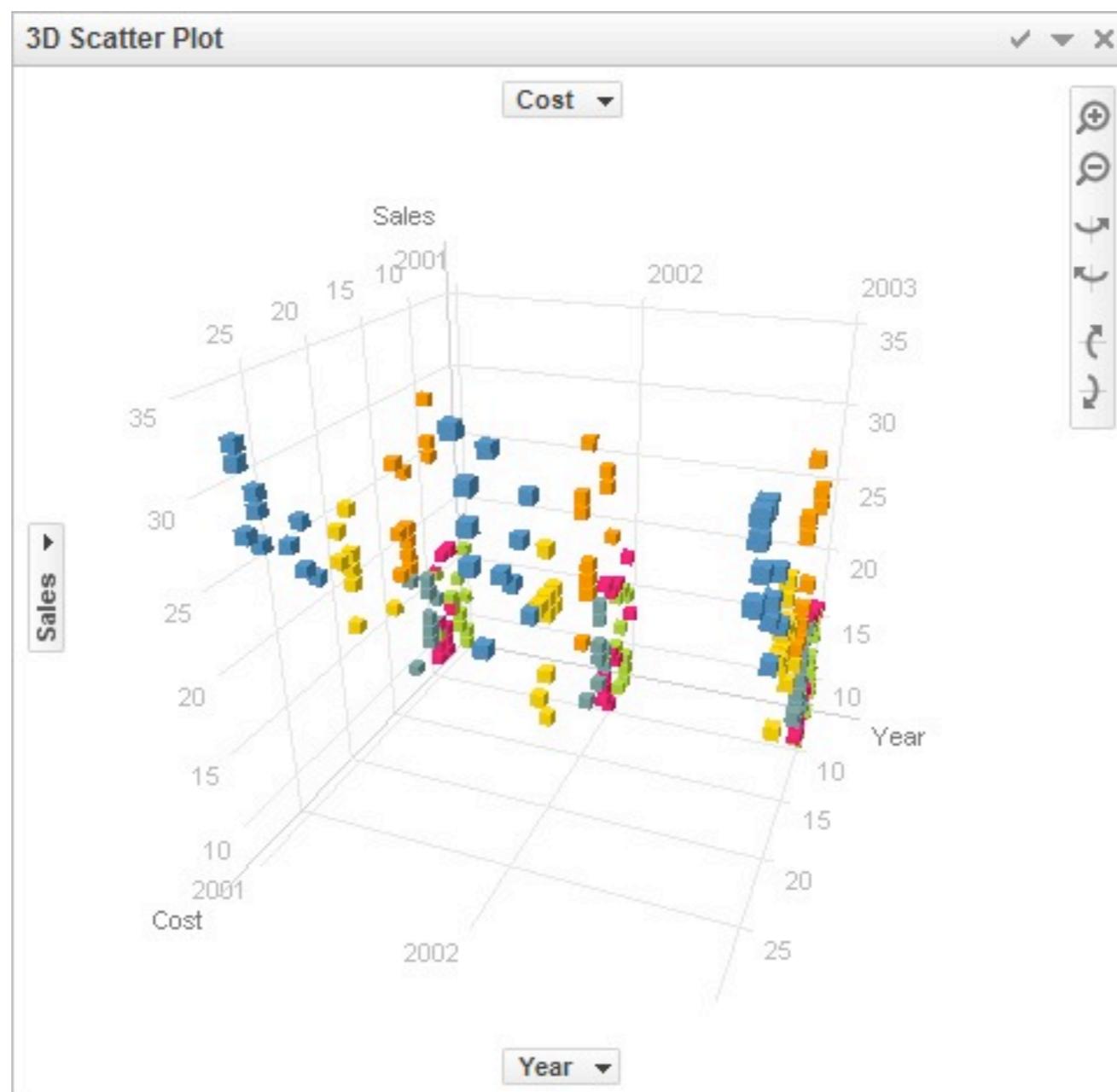
Bivariate Data

Scatterplot is common



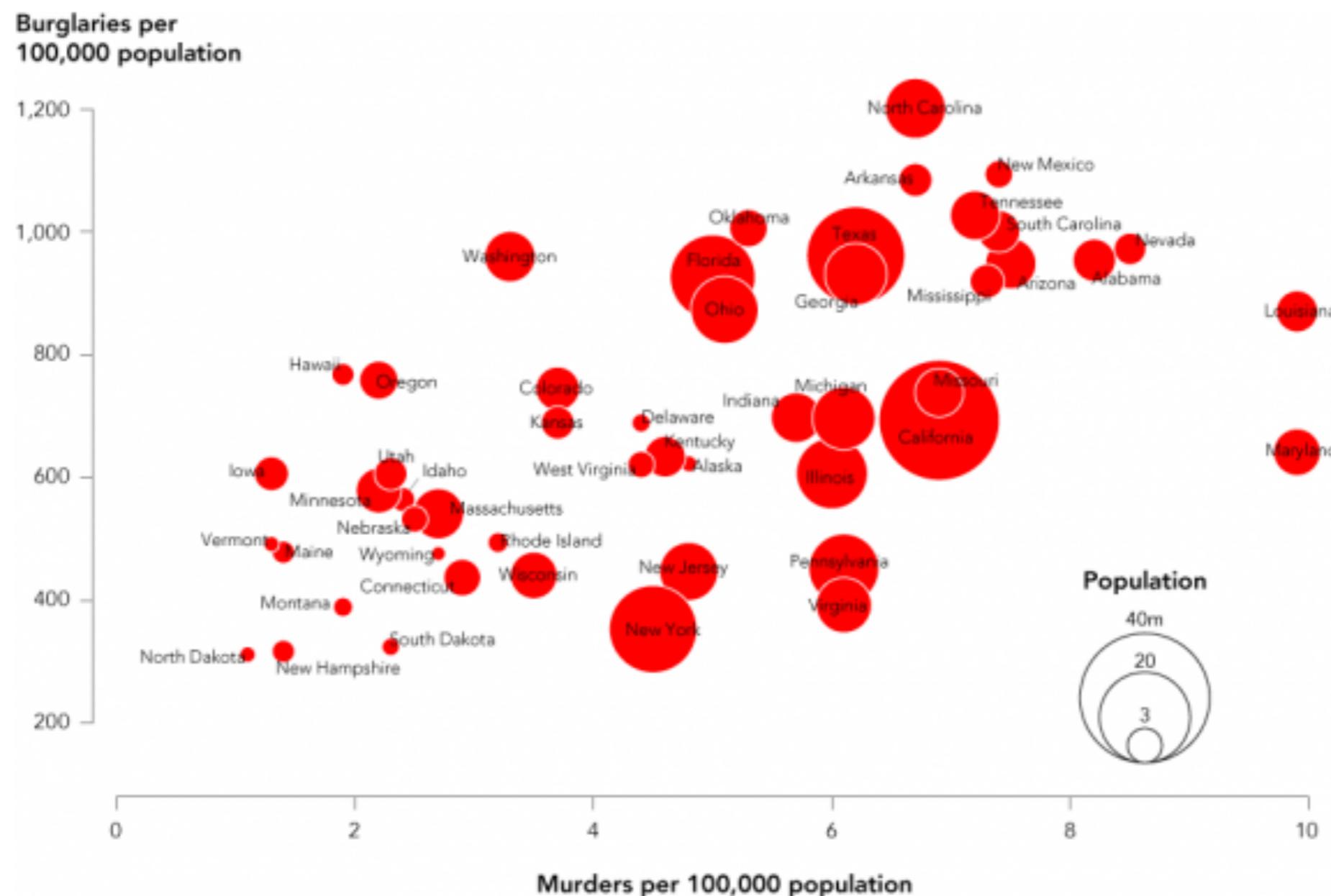
Trivariate Data

Do NOT use 3D scatterplots!



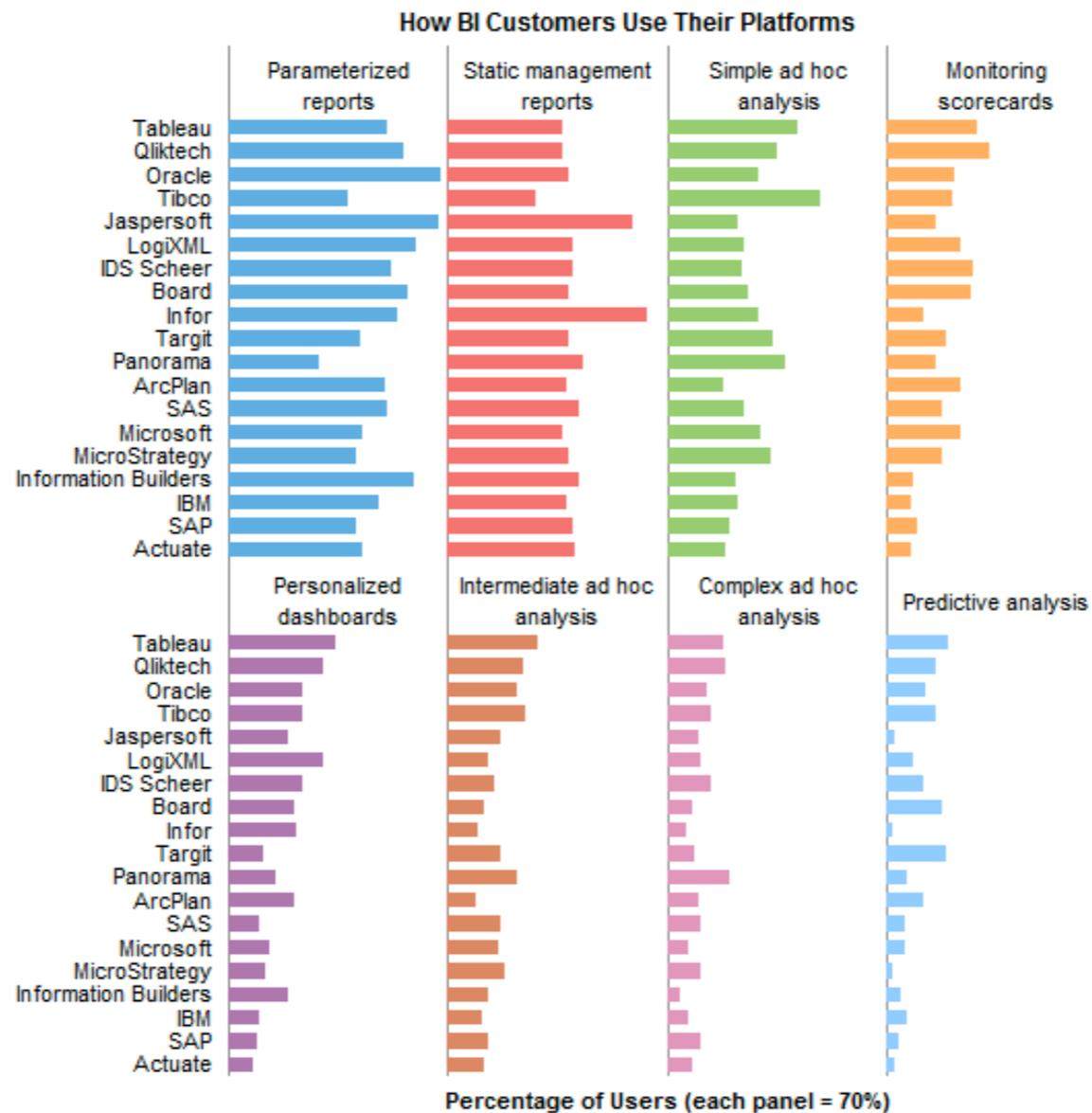
Trivariate Data

Map the third dimension to some other visual attribute

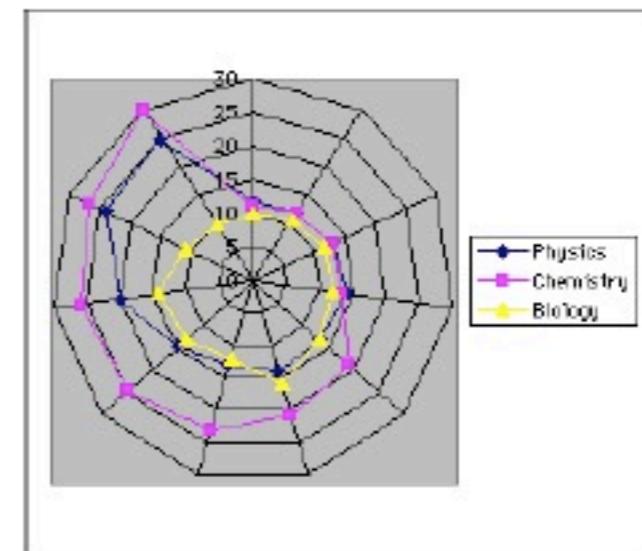
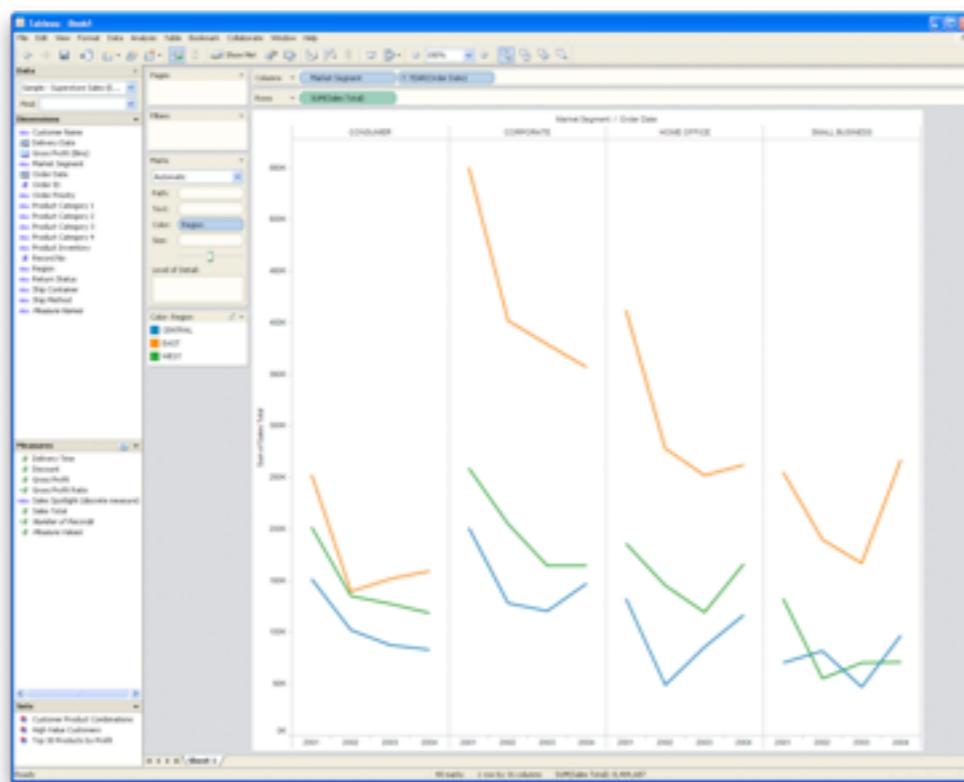
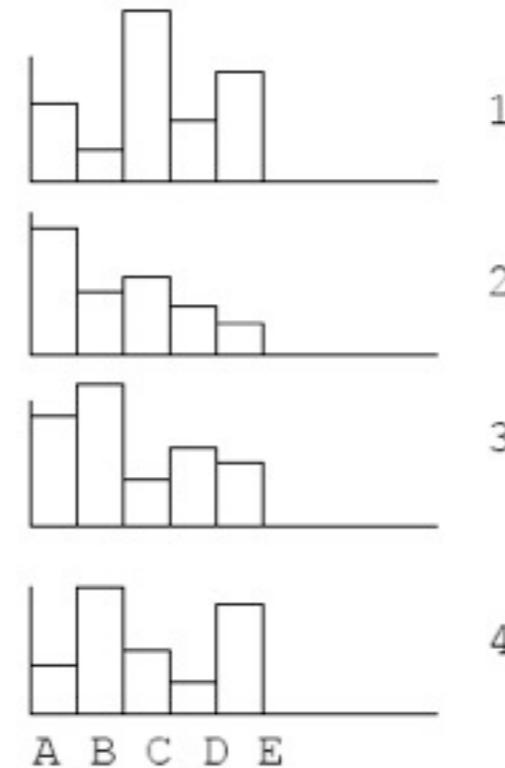
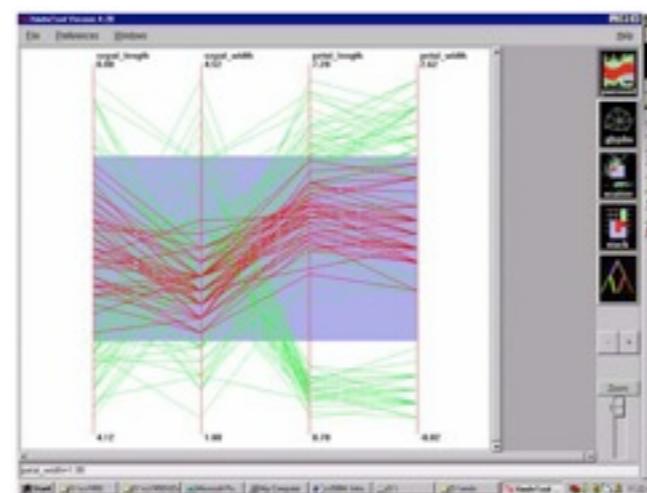
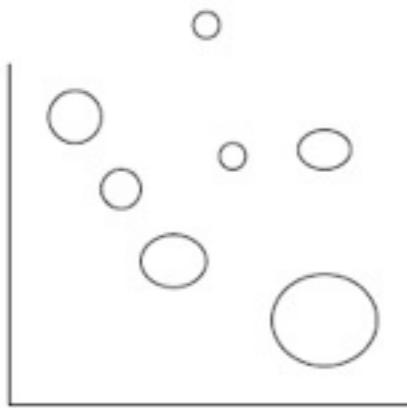


Multivariate Data

Give each attribute its own display (small multiples)



Multivariate Data Representations



Data Reduction

- **Filtering:** Eliminate some items or attributes
 - e.g., select range of interest, zoom in, remove outliers, etc.
- **Aggregation:** Represent a group of elements by a new derived element
 - e.g., take average, min, max, count, sum
 - Attribute aggregation a.k.a. dimensionality reduction

Mapping Data to Images

Jacques Bertin

- French cartographer
[1918-2010]
- Semiology of Graphics
[1967]
- Theoretical principles
for visual encodings



Bertin's Visual Attributes

Channels

Position

Size

(Grey)Value

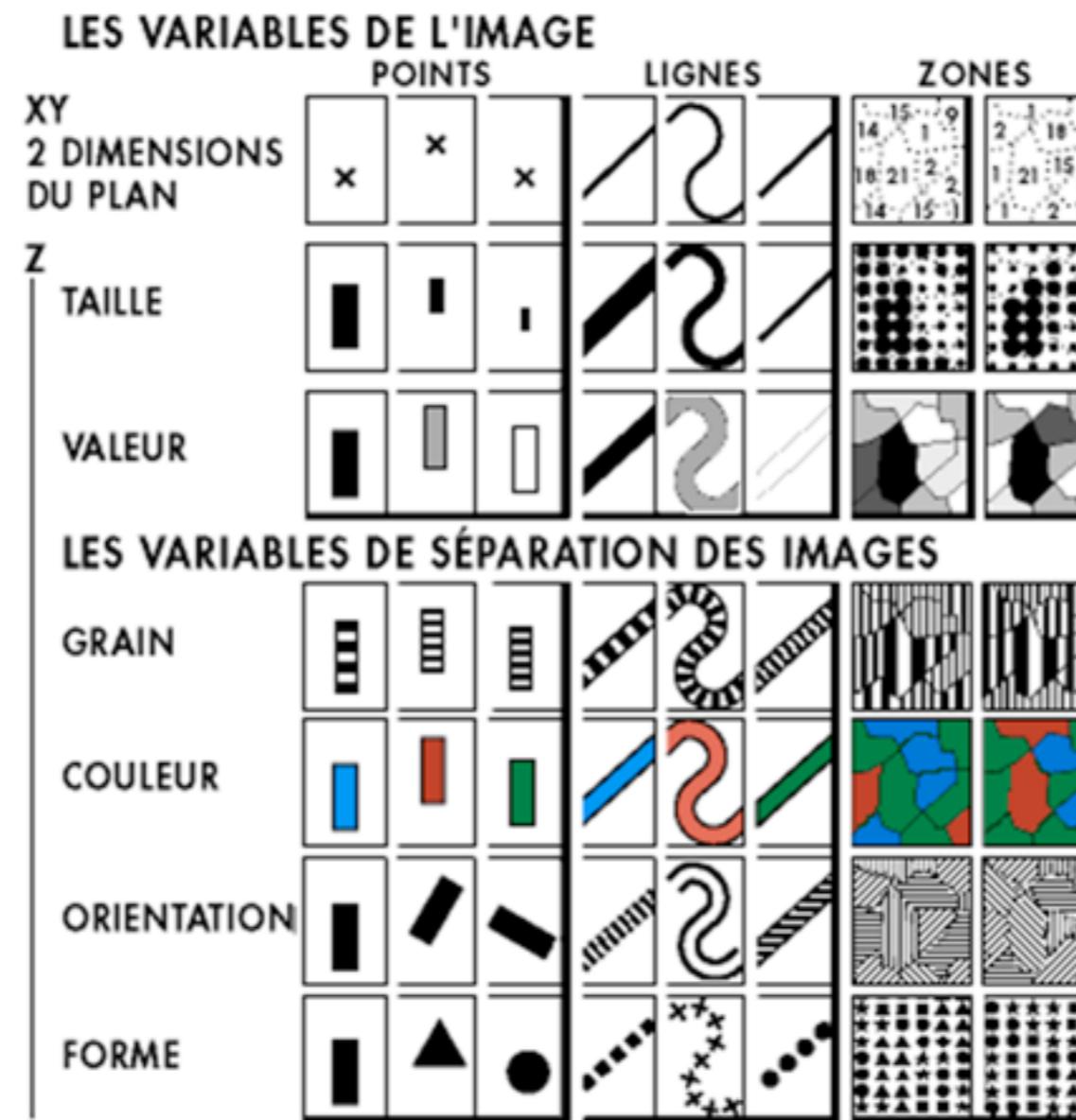
Texture

Color

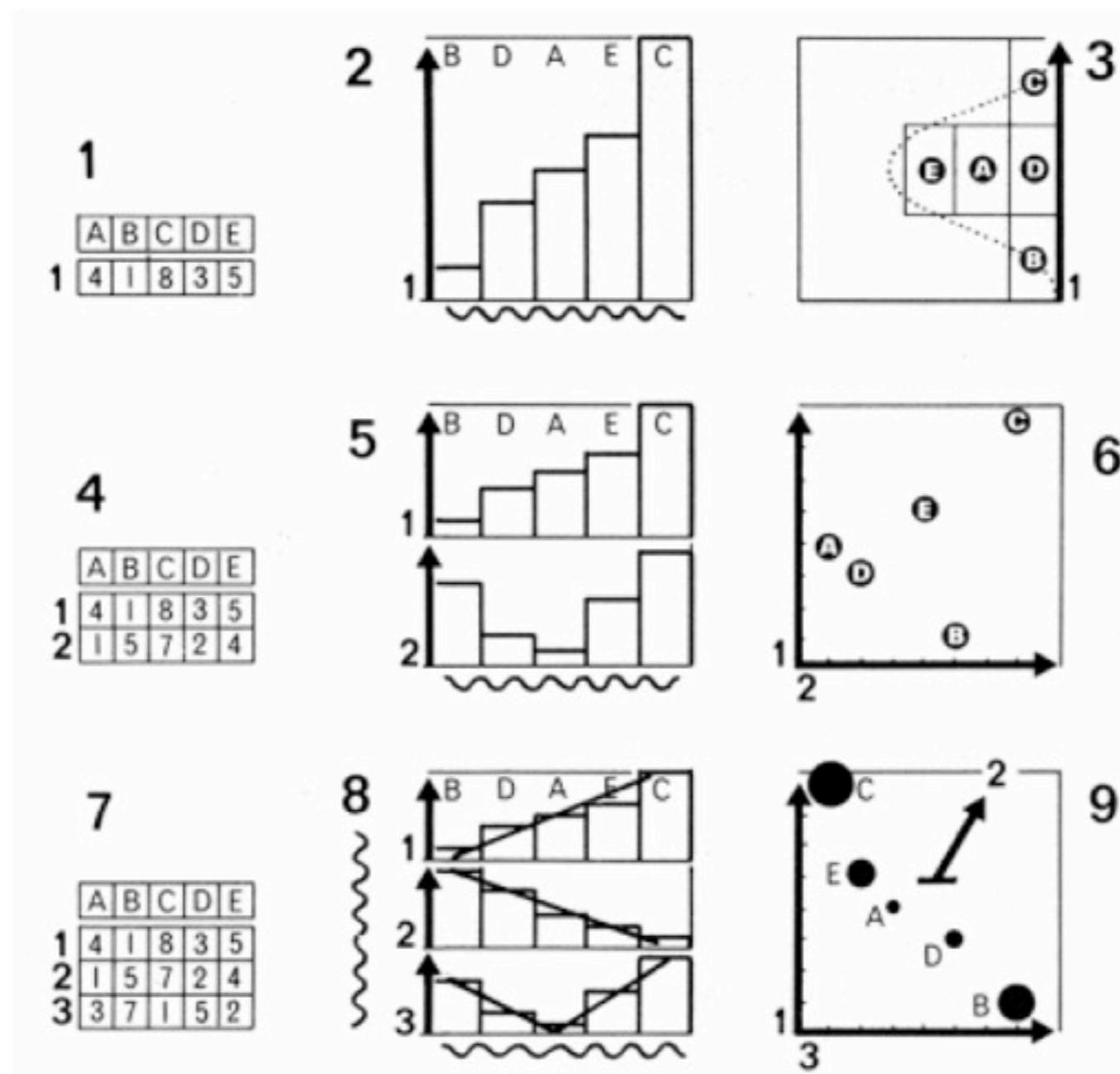
Orientation

Shape

Marks Points Lines Areas

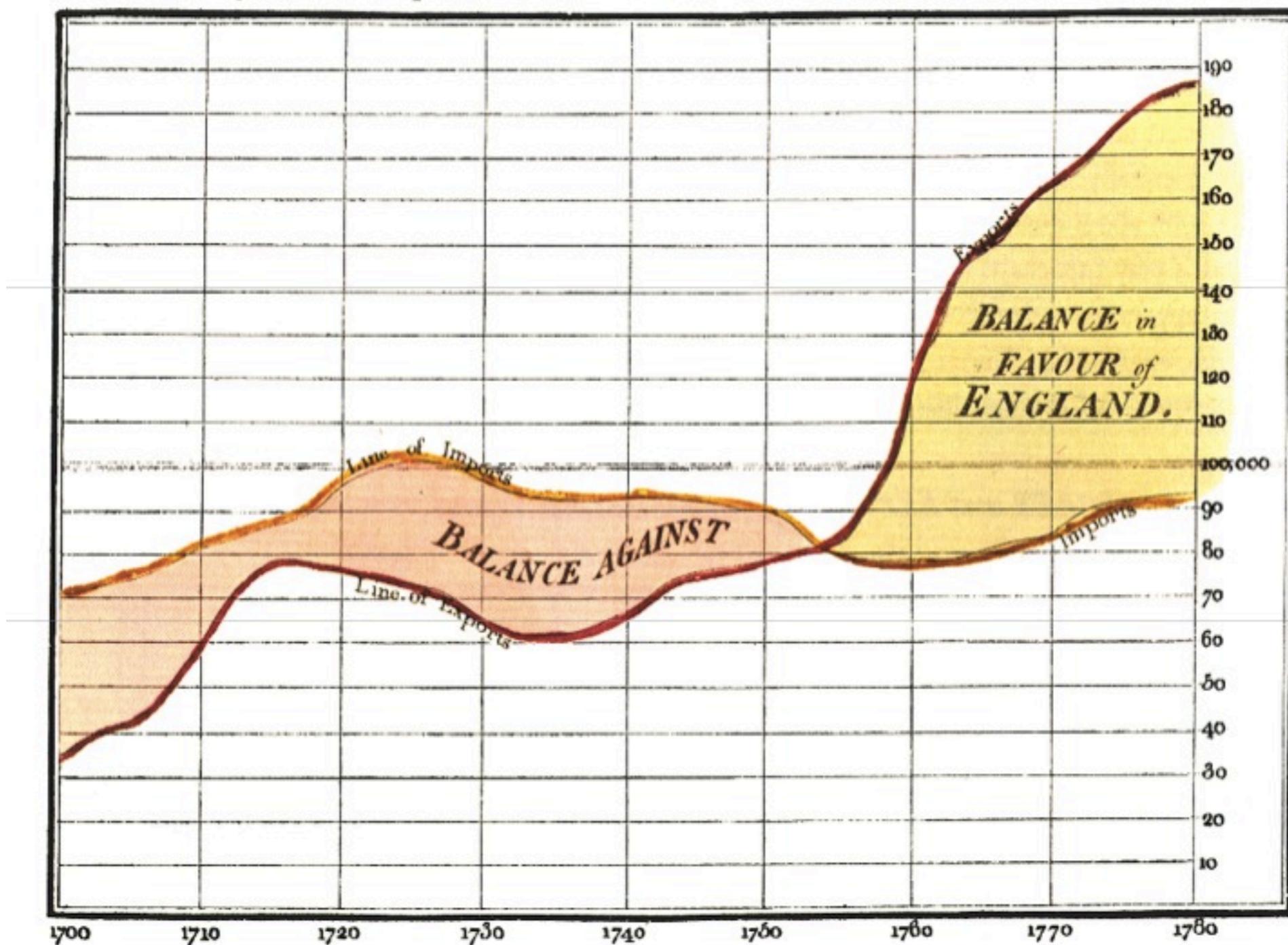


Large Design Space (Visual Metaphors)



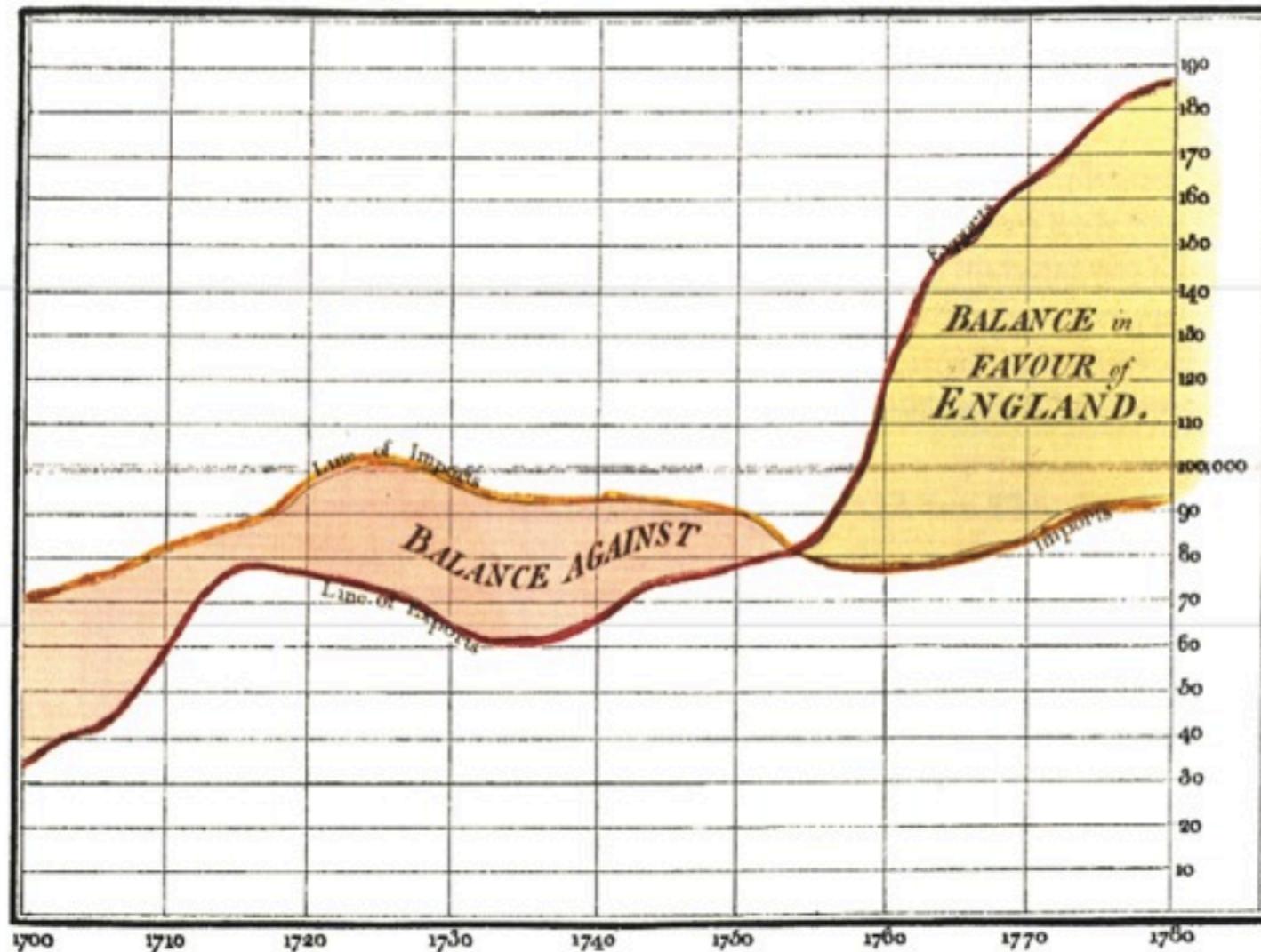
Example

Exports and Imports to and from DENMARK & NORWAY from 1700 to 1780.



Example

Exports and Imports to and from DENMARK & NORWAY from 1700 to 1780.



x-axis: Year (Q)

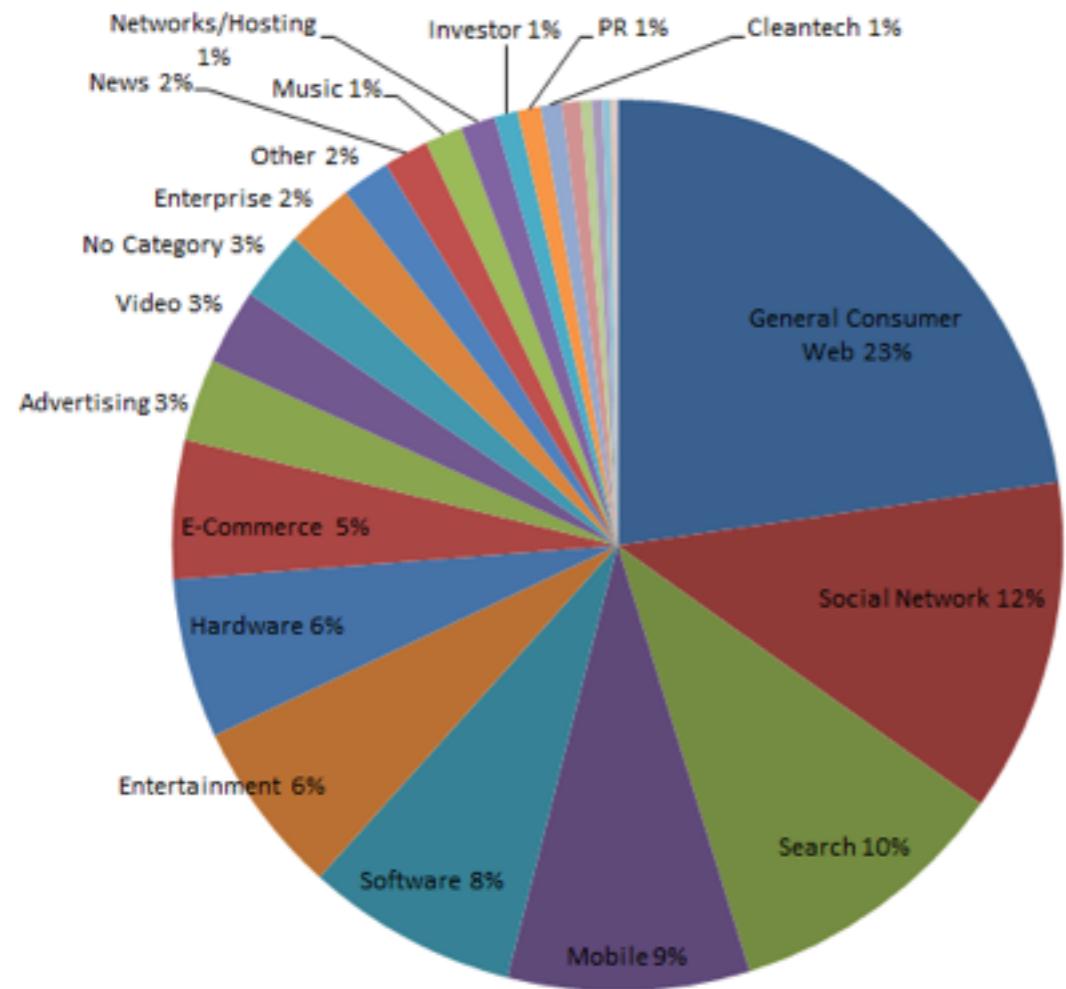
y-axis: Currency (Q)

Color: Imports / Exports (N, O)

W. Playfair, 1786

Effective Visual Attributes

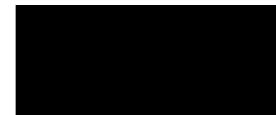
Compare These Values



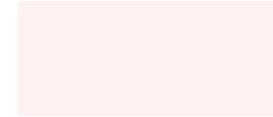
Order These Colors



Order These Colors

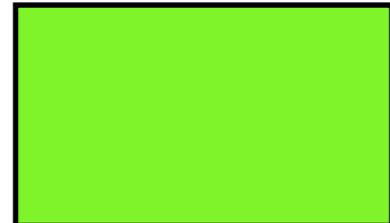


Order These Colors

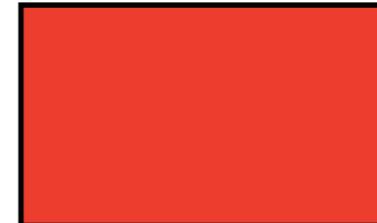


Perceived as Ordered

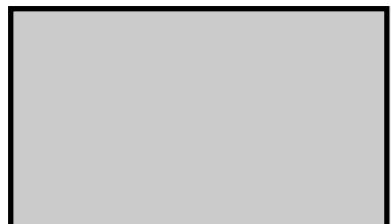
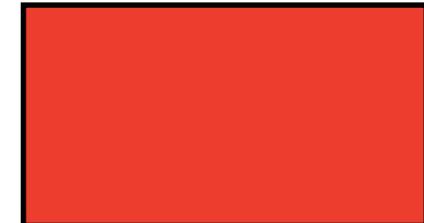
Brightness



Saturation



Hue: not as much



Visual Attributes per Data Type

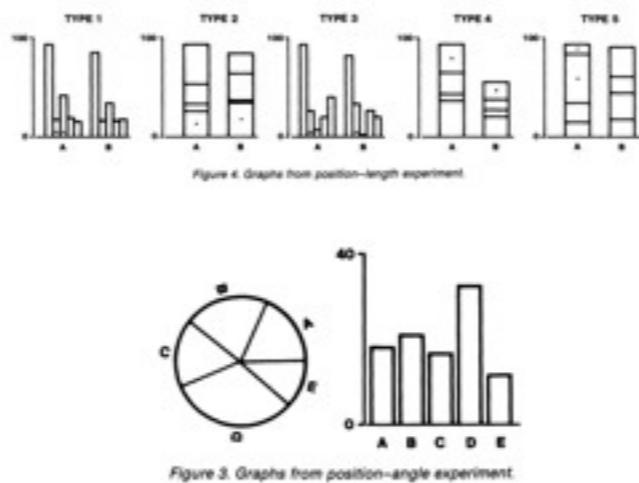
Bertin, 1967

	Categories	Ordinal	Quantitative
Position	✓	✓	✓
Length	✓	✓	✓
Brightness	✓	✓	~
Texture	✓	~	✗
Color	✓	~	✗
Angle	✓	✗	✗
Shape	✓	✗	✗

✓ = Good
~ = OK
✗ = Bad

Bertin, Semiology of Graphics, 1967

Cleveland / McGill, 1984

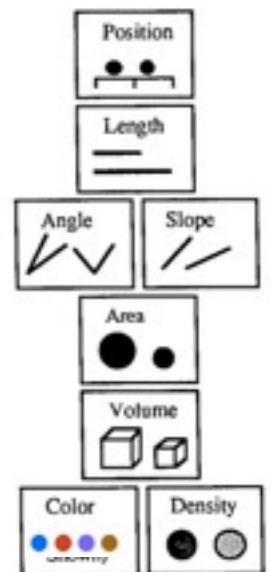


William S. Cleveland; Robert McGill ,
“Graphical Perception: Theory,
Experimentation, and Application to
the Development of Graphical Methods.” 1984

Mackinlay, 1986

More accurate

Less accurate

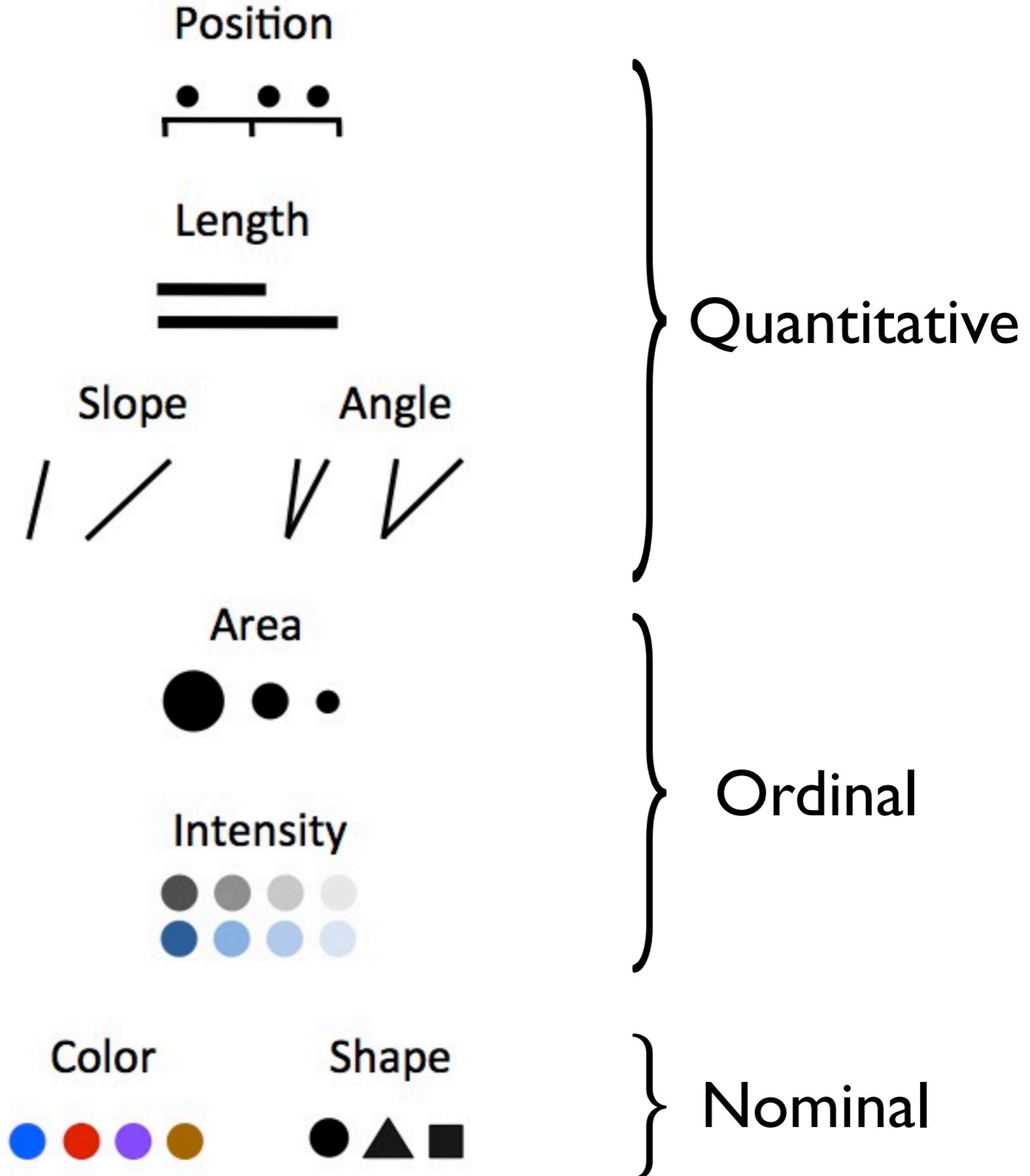


Jock Mackinlay “Automating The Design of
Graphical Presentations.” 1986

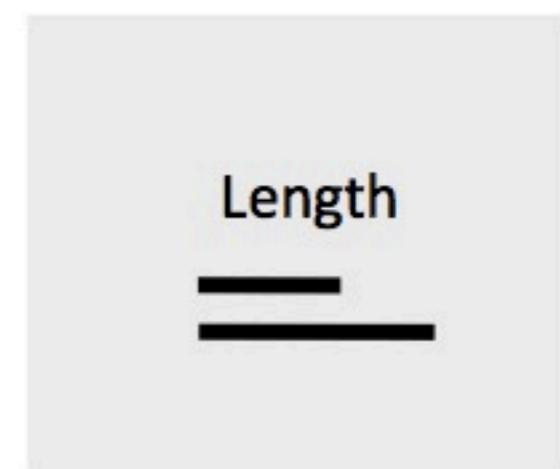
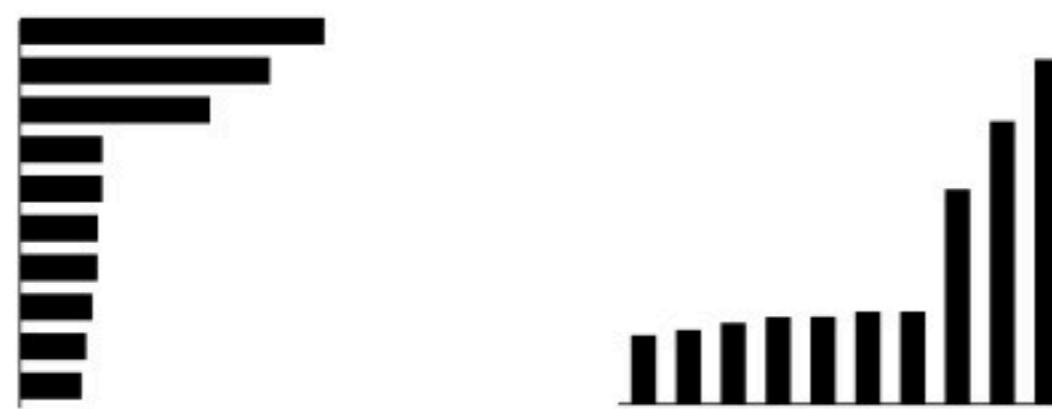
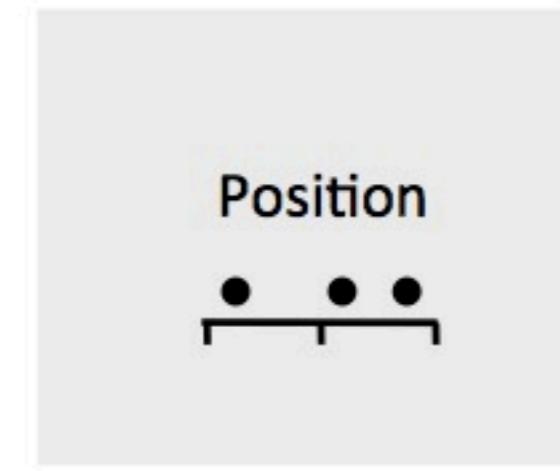
Most
Efficient



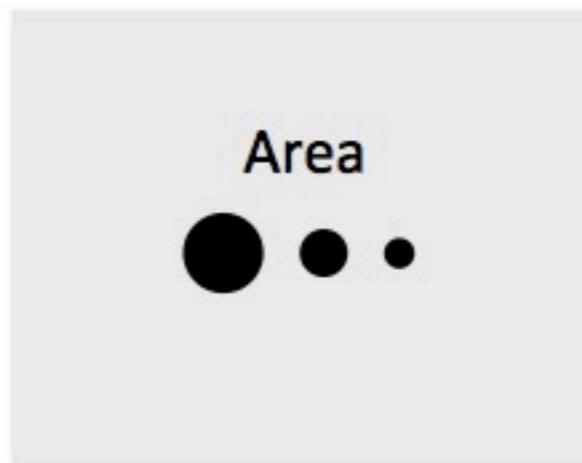
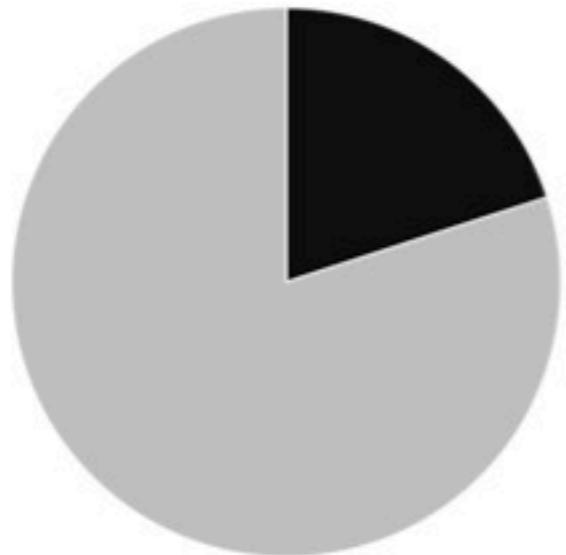
Least
Efficient



Most Effective



Less Effective



Least Effective

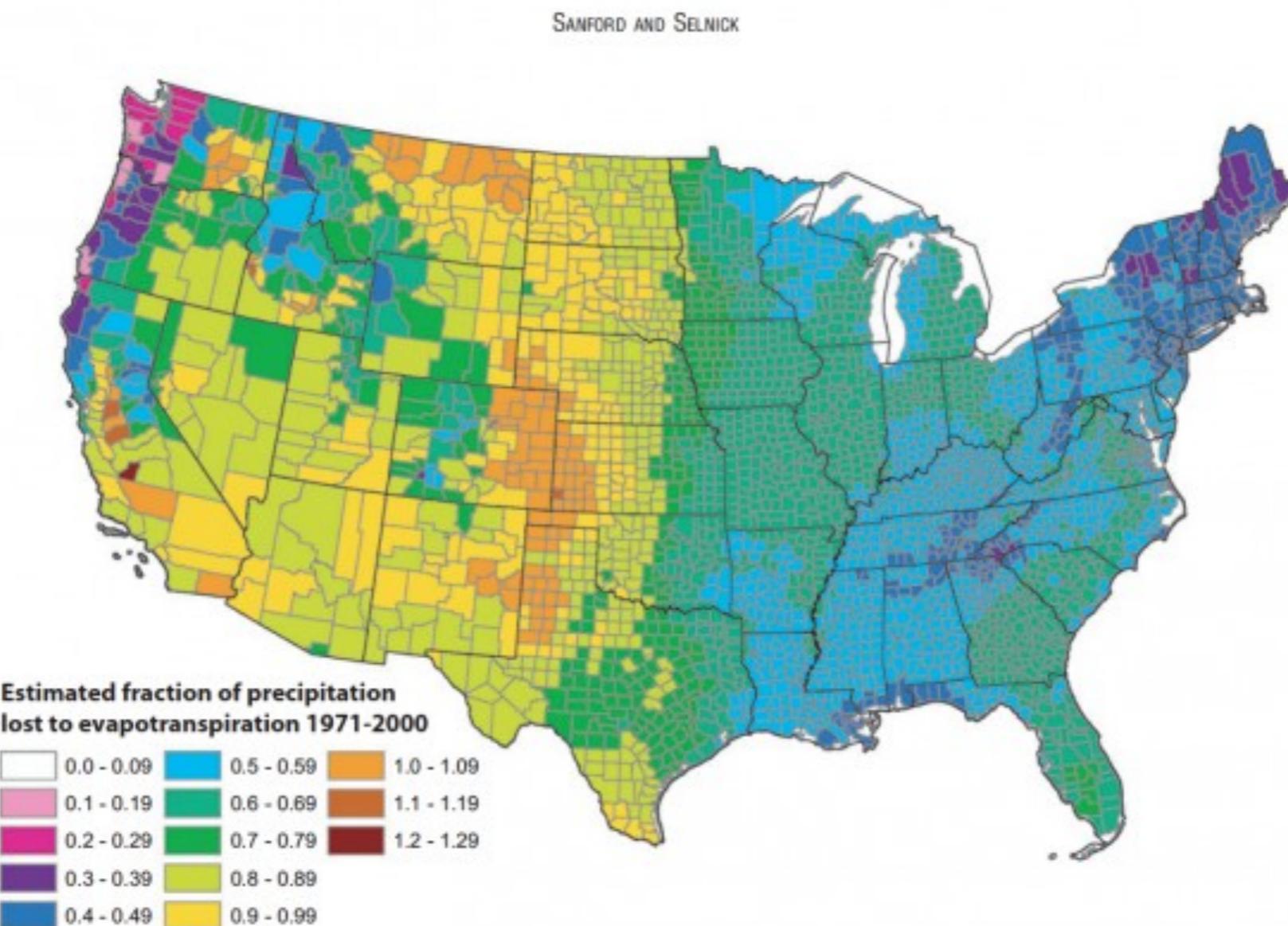


FIGURE 13. Estimated Mean Annual Ratio of Actual Evapotranspiration (ET) to Precipitation (P) for the Conterminous U.S. for the Period 1971-2000. Estimates are based on the regression equation in Table 1 that includes land cover. Calculations of ET/P were made first at the 800-m resolution of the PRISM climate data. The mean values for the counties (shown) were then calculated by averaging the 800-m values within each county. Areas with fractions >1 are agricultural counties that either import surface water or mine deep groundwater.

Effective Visualizations

Not Effective...

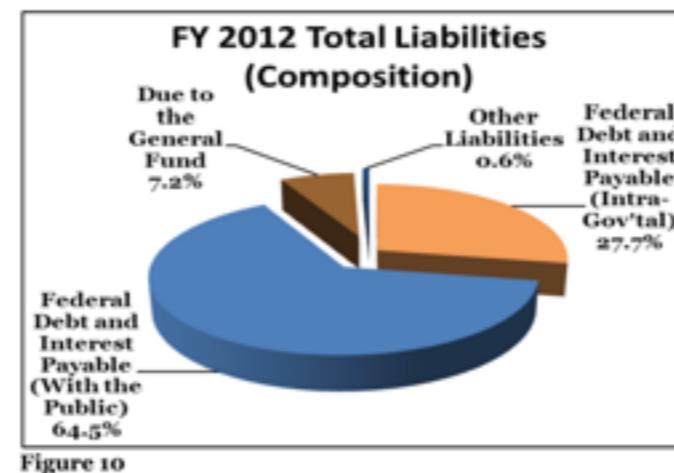
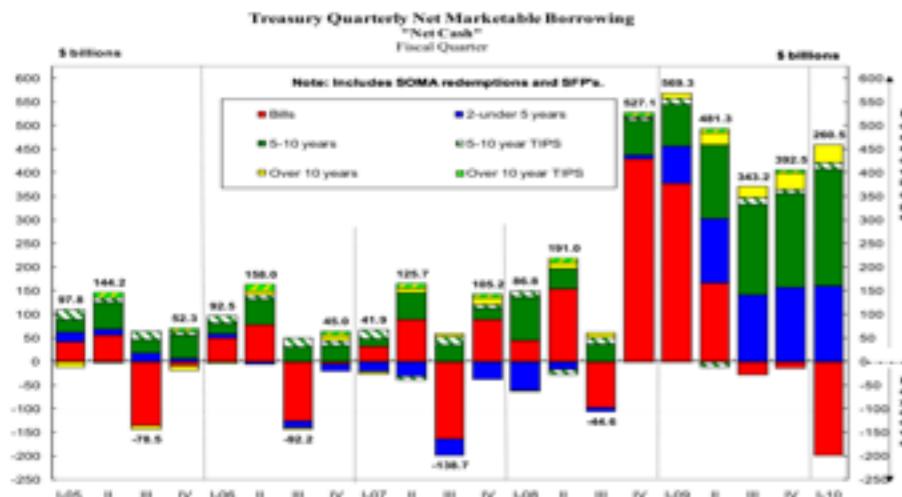


Figure 10

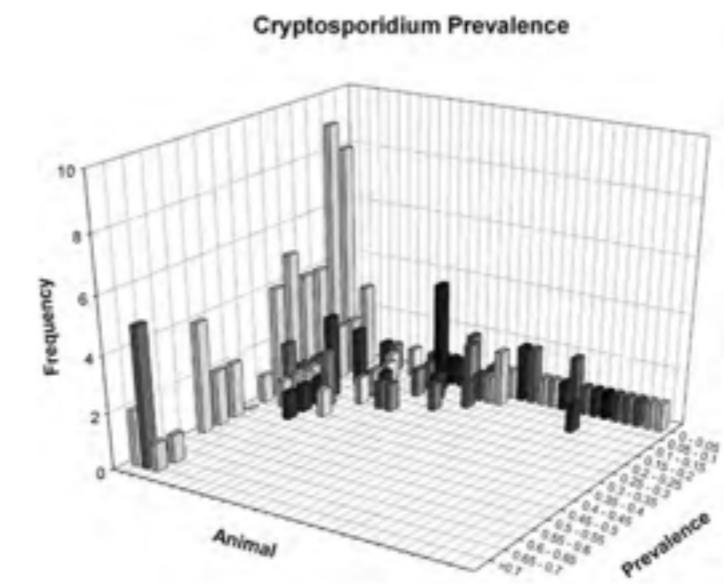
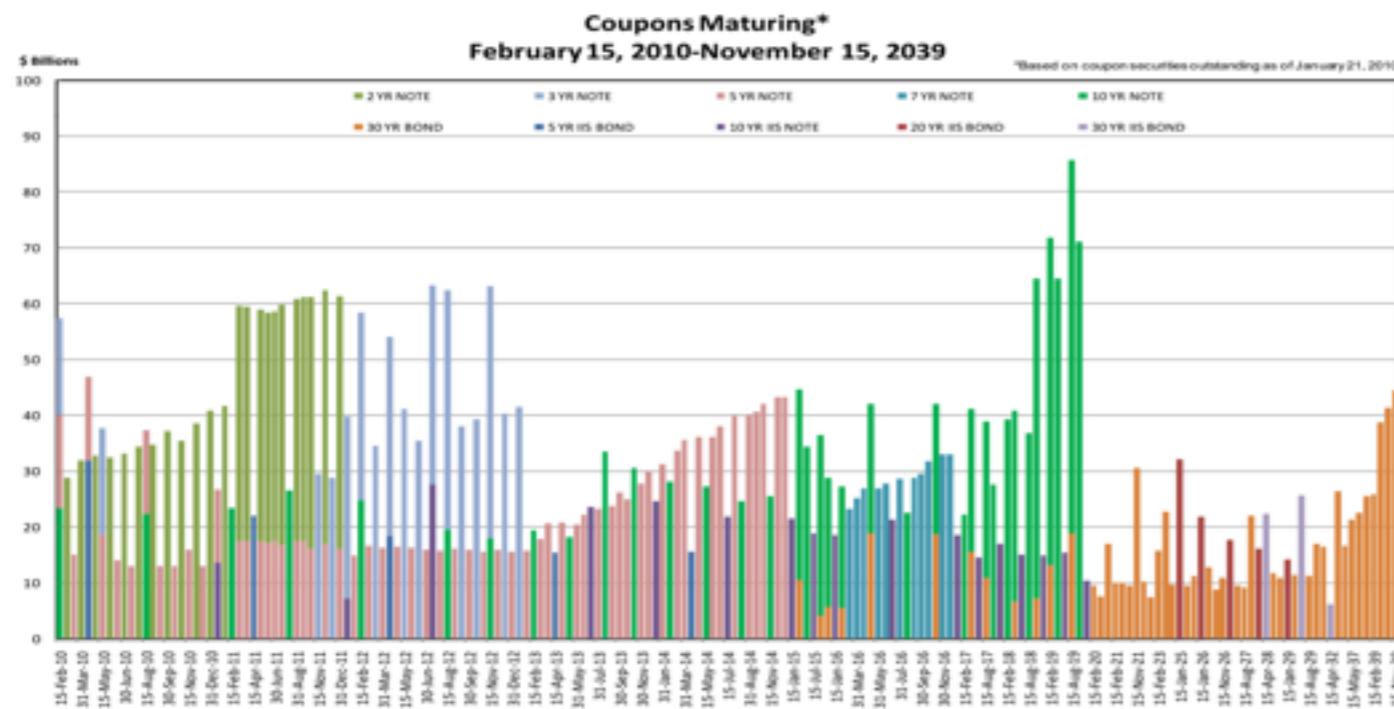
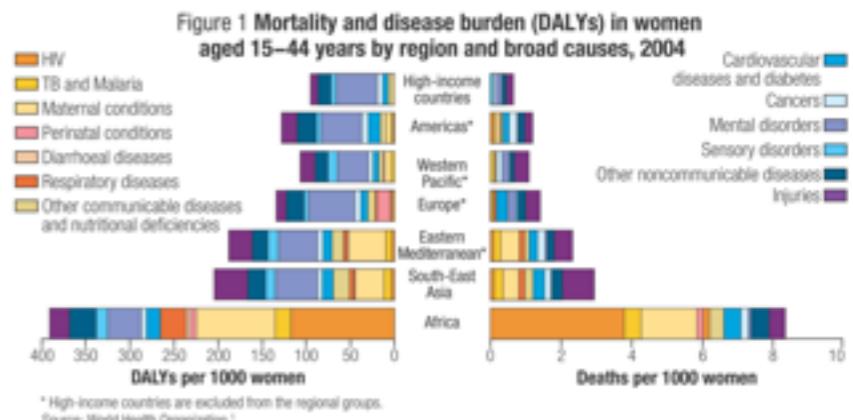
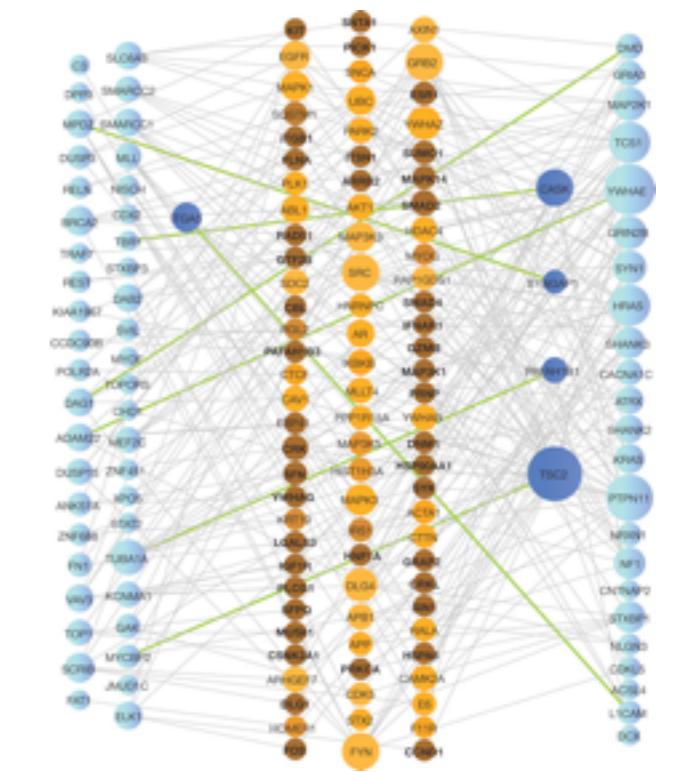
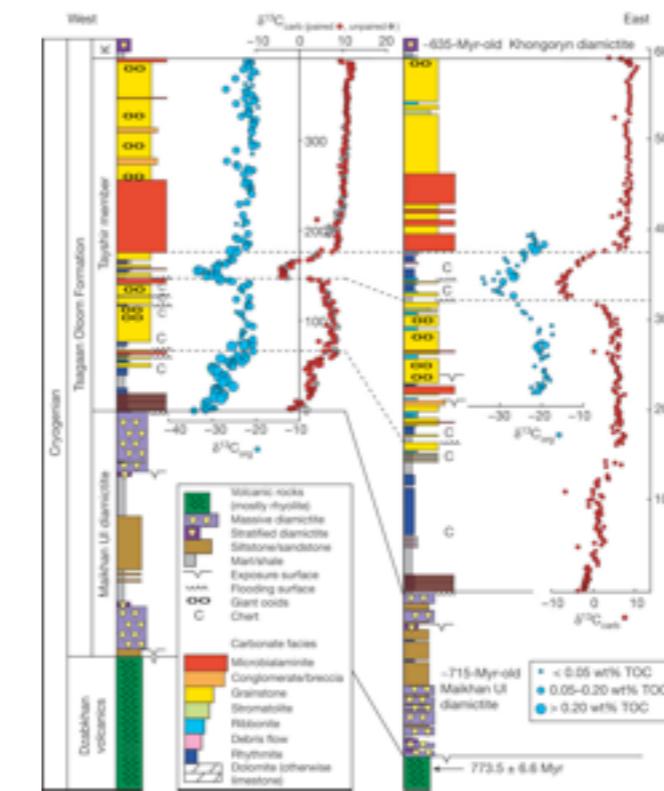
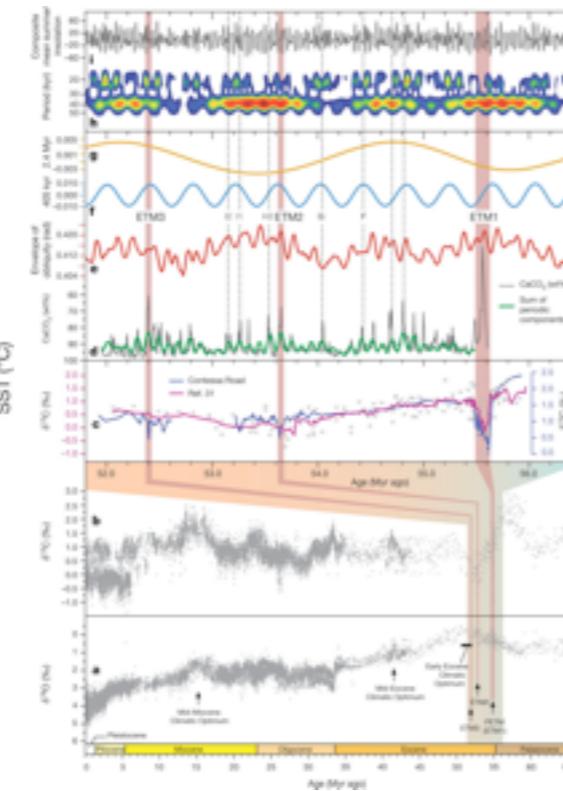
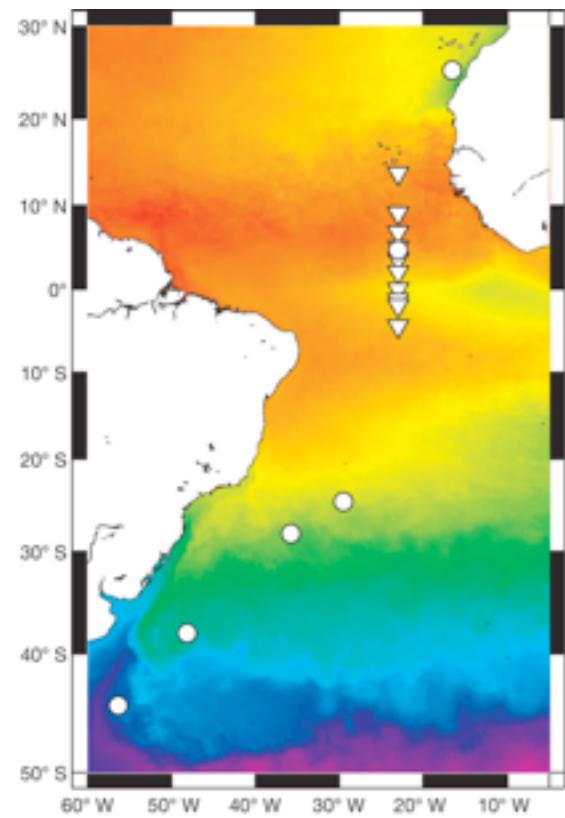
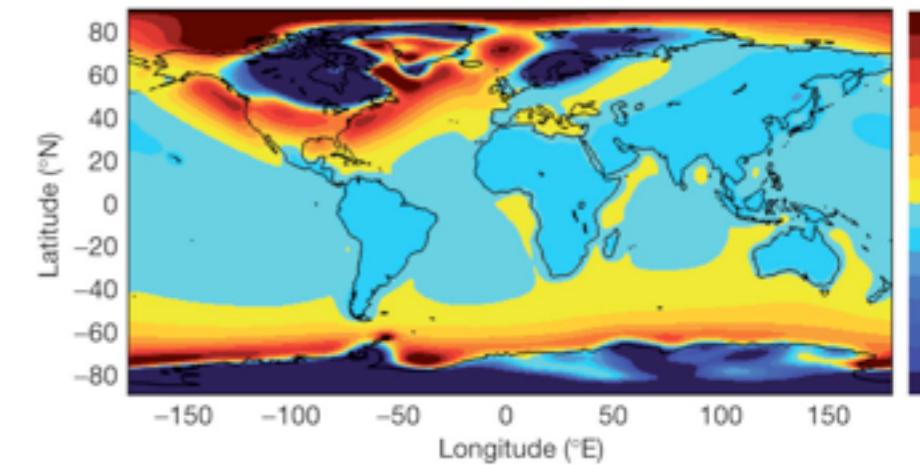
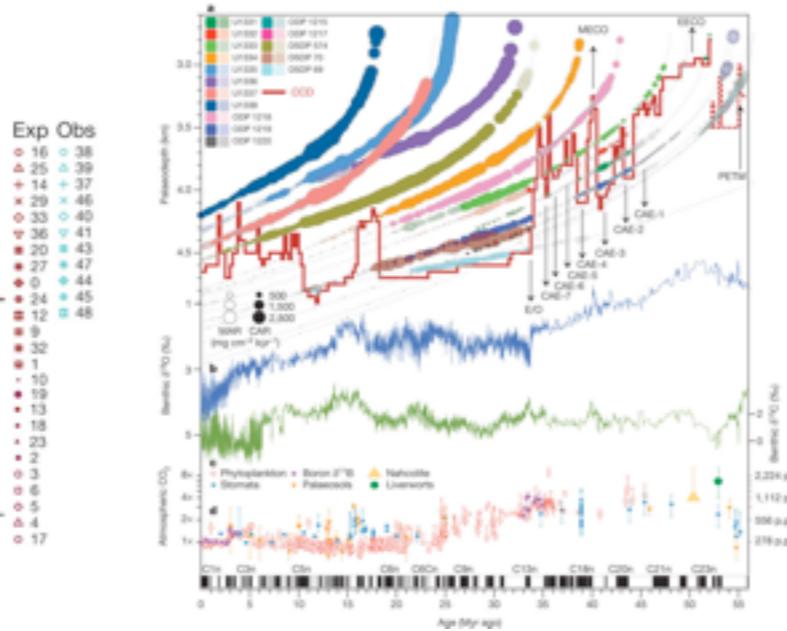
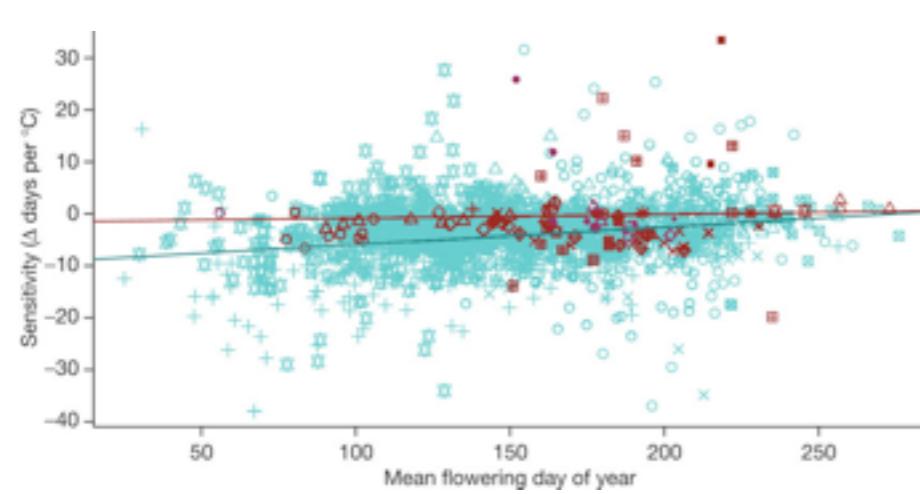


Figure 5.2 Mean prevalence rates of *Cryptosporidium* oocysts by animal species.

Sources: US Treasury and WHO reports

Also not effective...



Source: Nature

Much better...

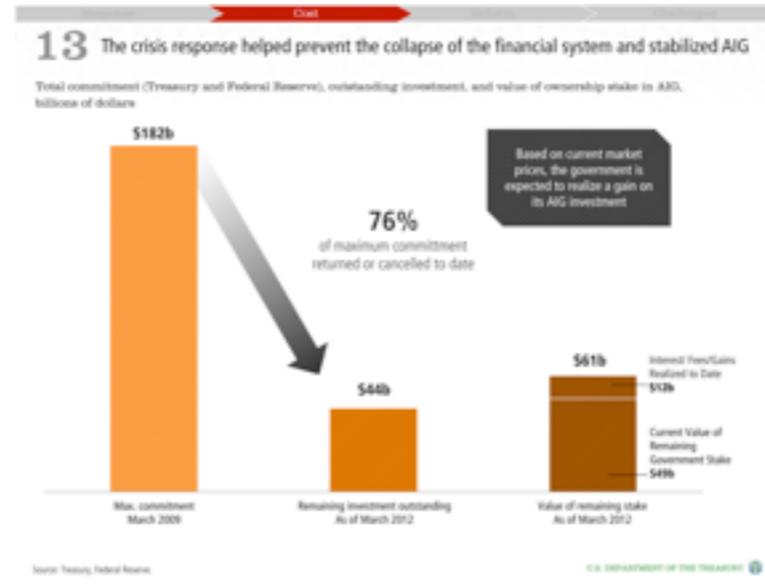
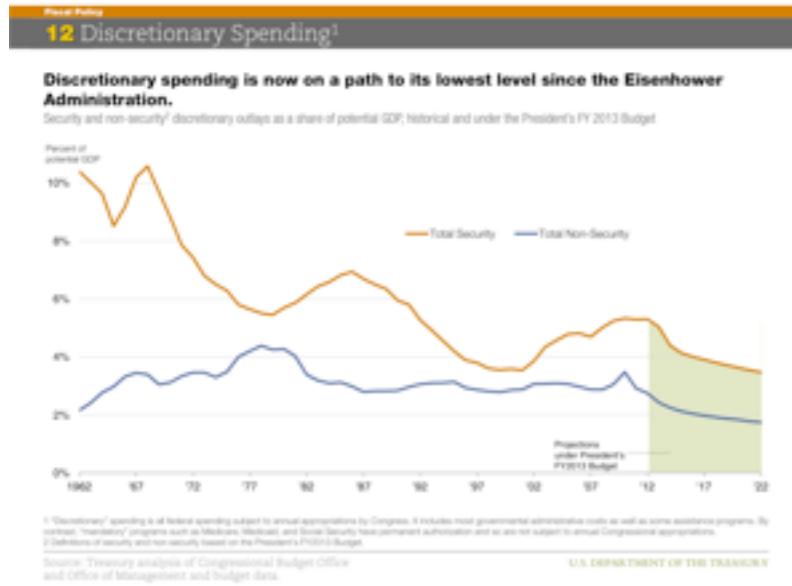
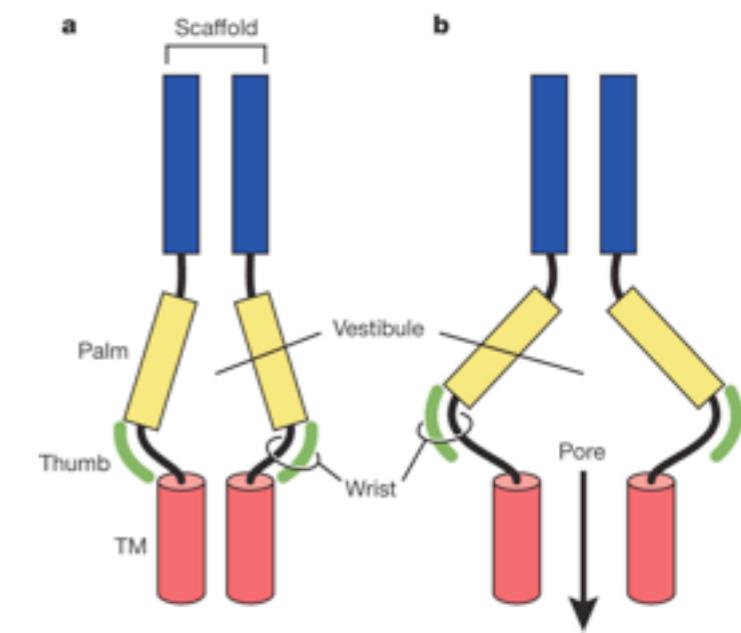
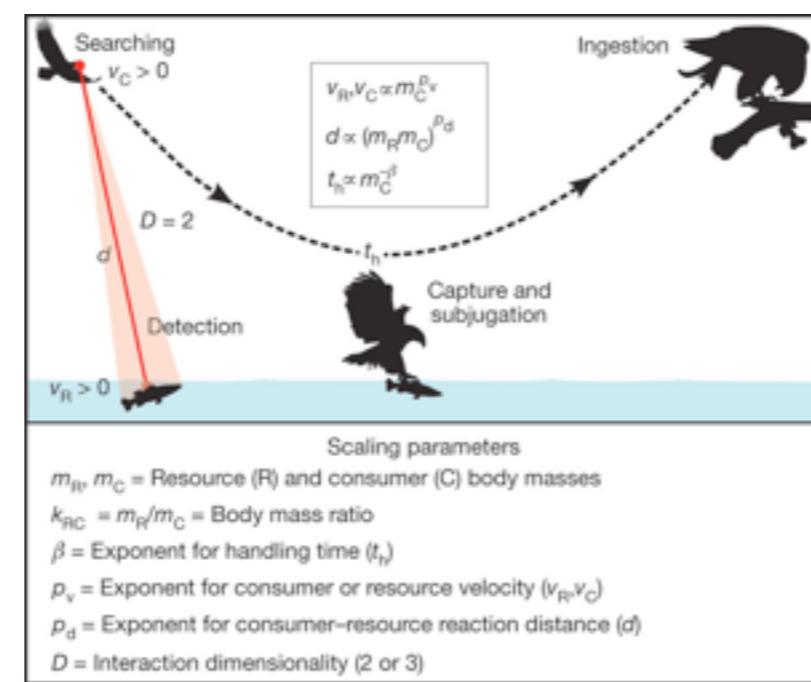
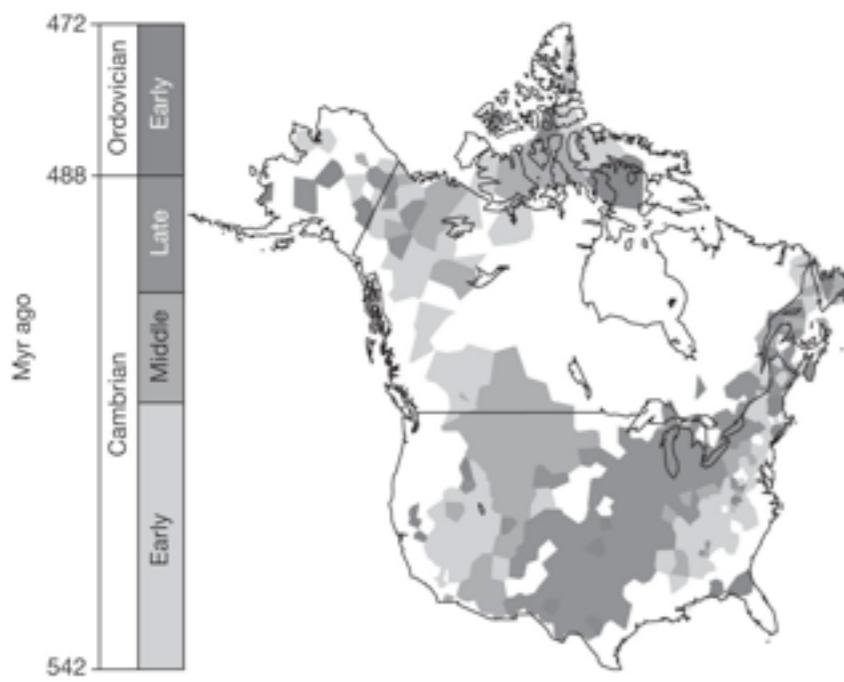
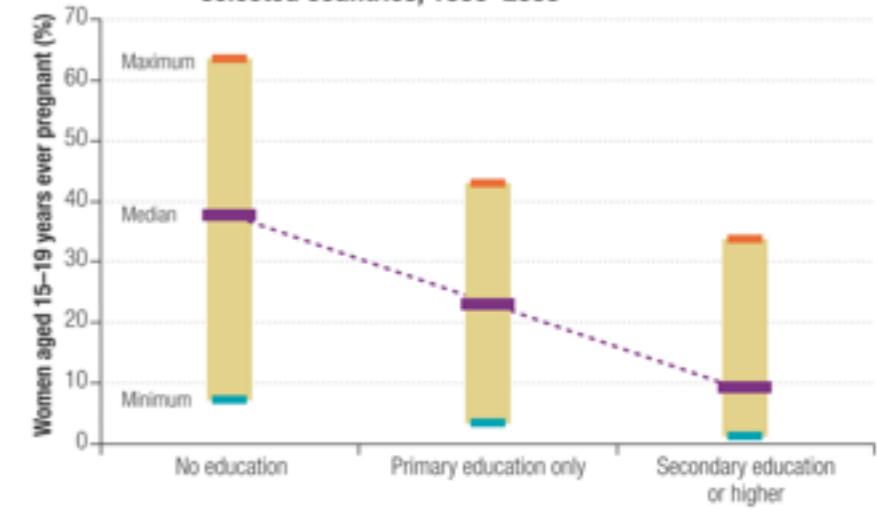


Figure 2 Adolescence pregnancy rates by educational level, selected countries, 1990–2005



Sources: US Treasury, WHO, Nature