
finding a story to tell: graphical research methods

Robert.Chung@berkeley.edu

October 2024

analytical graphics

- not presentation graphics, not exactly statistical graphics, not quite exploratory graphics

presentation graphics tell your story; analytical graphics help you figure out what your story is

analytical graphics are “work product”

- not always refined enough for publication or presentation

Tufte? A lot of his techniques focus on effectively communicating quantitative findings. We'll focus on steps before that: uncovering interesting stories and questions in data

Very valuable for grad students and early career researchers

analytical graphics are for analysis

- often analytical graphics are used not to prove hypotheses but to help generate them
 - we don't always answer questions; we use graphical techniques to help us ask new ones
- often, the audience is YOU
 - visualize differences and contrasts
 - across time
 - across places
 - across treatments or policies
 - across conditions
 - there are tips, tricks, and techniques that help you in visualization

what if you already have a question?

- no problem. sometimes analytical graphics can help you focus on where to look, or to refine your question
- doesn't replace theory, or your research question. You can do both: that's allowed

basic approach

- maximize insight
- uncover underlying structure
- extract important variables
- detect outliers and anomalies
- develop (very) simple models
- not much testing; that's for later

what we'll do in these lectures

- some examples
- some principles
- some basic tricks
- some slightly more advanced tricks
- you'll have a chance to try out some of these tricks before next week

the three things we're looking for

- look for
 - pattern
 - unexpected pattern
 - deviations from pattern
- these generate questions
- questions and how you address them are often the basis for papers or chapters or dissertations or careers
- “The data speak for themselves, but their voices are soft and sly”
so we're looking for ways to amplify their voices

demography

- demography is the study of populations, their characteristics, relationships among characteristics, and how they change
- you probably already know how to examine characteristics; we'll look at ways to highlight relationships among the characteristics and how they change
- in particular, we'll often look for models to help us understand the relationships among characteristics

the purpose of models

- “The purpose of models is not to fit the data but to sharpen the questions” – Sam Karlin
- We’ll use analytical graphics to help us sharpen questions

simple tools

- simple tools used intelligently (well, we can always dream) rather than complex tools used stupidly
rules of thumb, not hard rules and regulations
- a handful of plots and a handful of tricks
lots of specialty type graphs, but we try to avoid too many of them until we know our story
- xy plots are a hugely useful invention
with one or two exceptions we'll focus mostly on ways to enhance xy plots
- decoding the language of graphs can be complicated, so we build on familiar beginnings

how graphing helps

- we can only make sense of a handful of numbers at a single time
pages of dense tables are good for detail, evidence, and re-analysis but poor for understanding
- eye-brain is good at seeing patterns in large numbers of values
though it can be fooled—we'll present some problems that can mislead the eye
- therefore
 - use graphs when pattern is important
 - use tables when exact details are important
 - graphs and tables are complements, not replacements. (You can do both: that's allowed)

apophenia

- apophenia is “the experience of seeing patterns of connections in random or meaningless data”
- we'll occasionally accept a little “type I error” when we're looking for interesting questions – as long as we back it up later with real confirmatory analysis



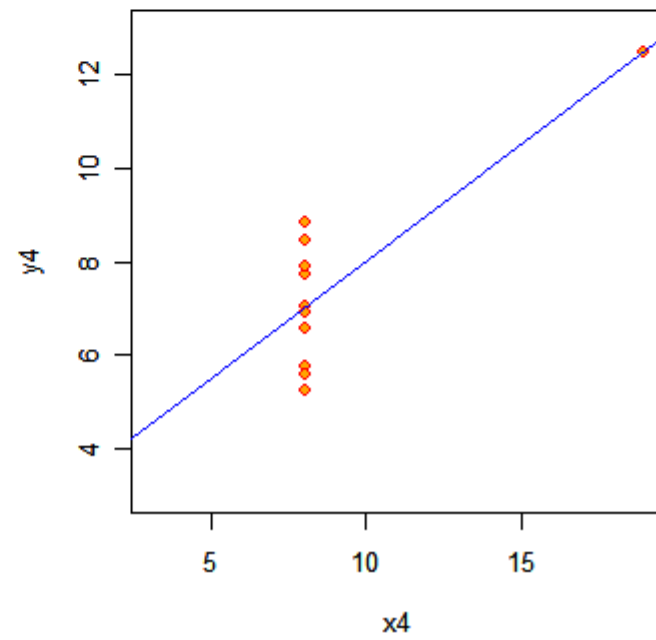
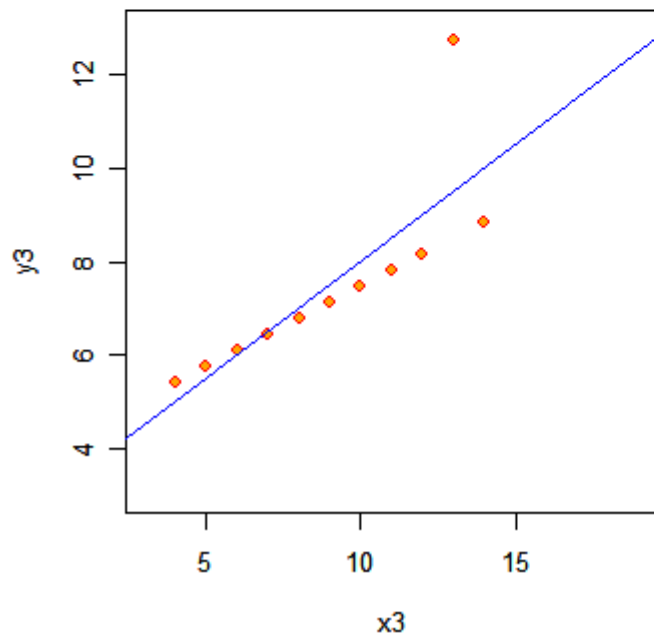
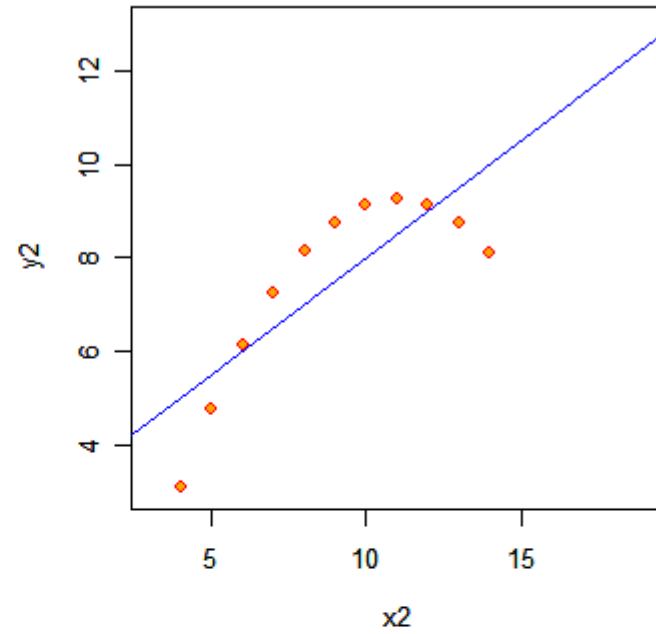
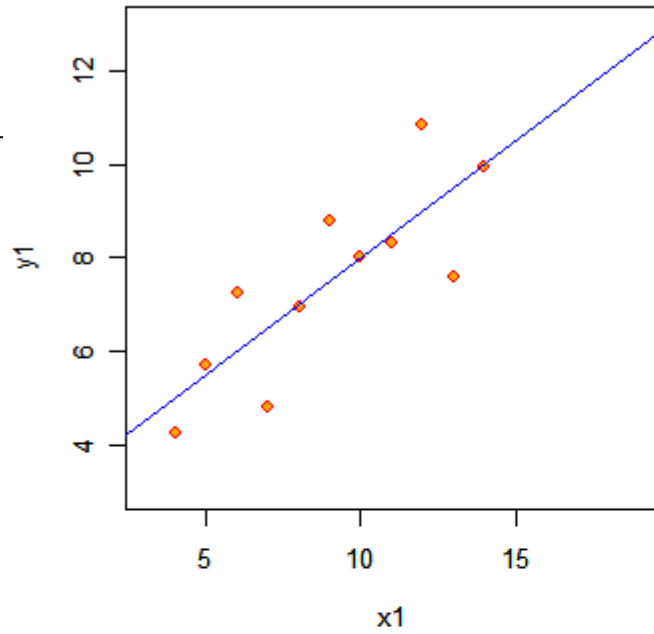
Anscombe's data

x1	y1	x2	y2	x3	y3	x4	y4
10	8.04	10	9.14	10	7.46	8	6.58
8	6.95	8	8.14	8	6.77	8	5.76
13	7.58	13	8.74	13	12.74	8	7.71
9	8.81	9	8.77	9	7.11	8	8.84
11	8.33	11	9.26	11	7.81	8	8.47
14	9.96	14	8.10	14	8.84	8	7.04
6	7.24	6	6.13	6	6.08	8	5.25
4	4.26	4	3.10	4	5.39	19	12.50
12	10.84	12	9.13	12	8.15	8	5.56
7	4.82	7	7.26	7	6.42	8	7.91
5	5.68	5	4.74	5	5.73	8	6.89

Anscombe's data

- same means, sd's, correlation, regression slope, fit
$$\begin{aligned}\text{mean}(x_1) &= \text{mean}(x_2) = \text{mean}(x_3) = \text{mean}(x_4) = 9 \\ \text{mean}(y_1) &= \text{mean}(y_2) = \text{mean}(y_3) = \text{mean}(y_4) = 7.5 \\ \text{sd}(x_1) &= \text{sd}(x_2) = \text{sd}(x_3) = \text{sd}(x_4) = 3.32 \\ \text{sd}(y_1) &= \text{sd}(y_2) = \text{sd}(y_3) = \text{sd}(y_4) = 2.03 \\ r(x_1, y_1) &= r(x_2, y_2) = r(x_3, y_3) = r(x_4, y_4) = 0.816 \\ y^* &= 3 + 0.5 x^* \text{ with } r^2 = 0.667\end{aligned}$$
- so, conventional linear models make them look alike
- what will you see if you graph the data?

Anscombe's 4 Regression data sets



the NJ Pick-It lottery

- each bettor selected a 3-digit number between 0 and 999
- each ticket cost 50 cents
- all bettors who held the winning number split the prize money. The size of the prize depended on selecting the winning number **and** on the number of players who chose that number
- what would you want to know?

winning numbers and prize amounts

(810, \$190.0)

(156, \$120.5)

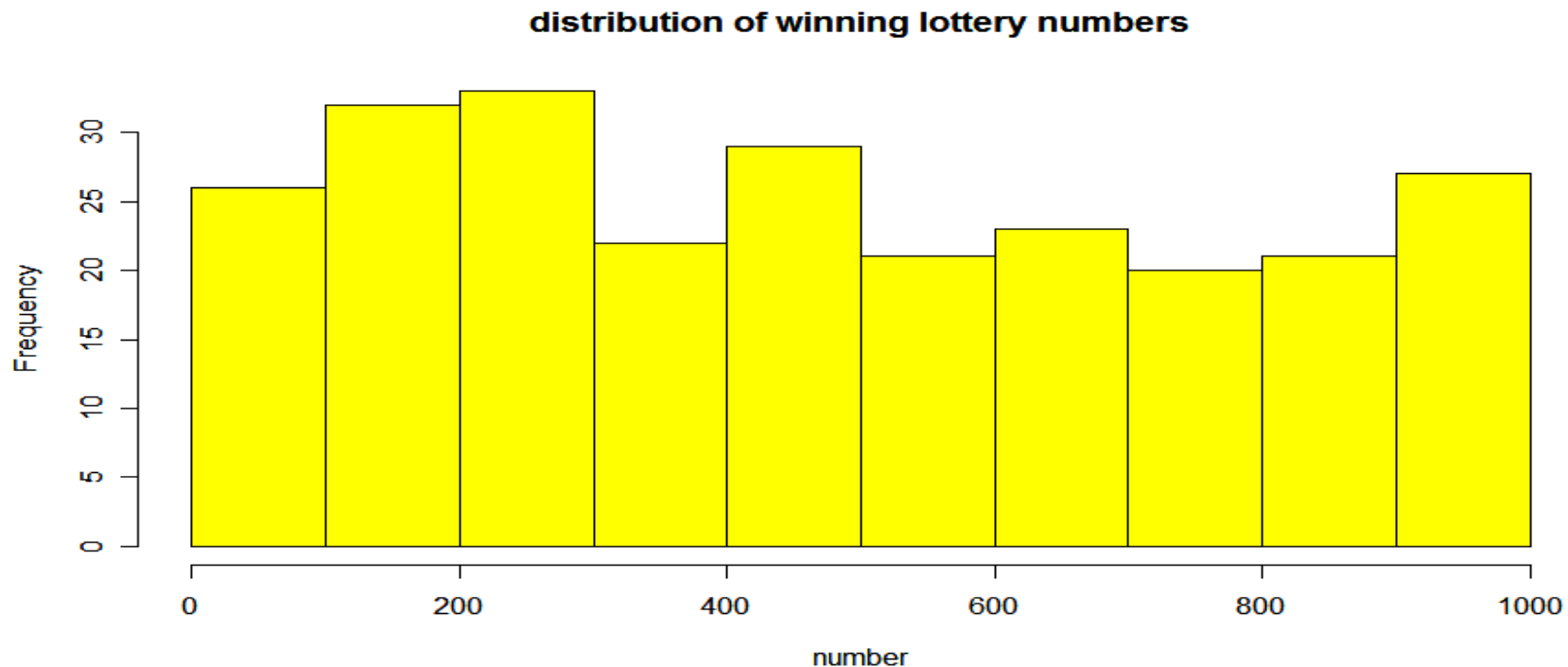
(140, \$285.5)

(542, \$184.0)

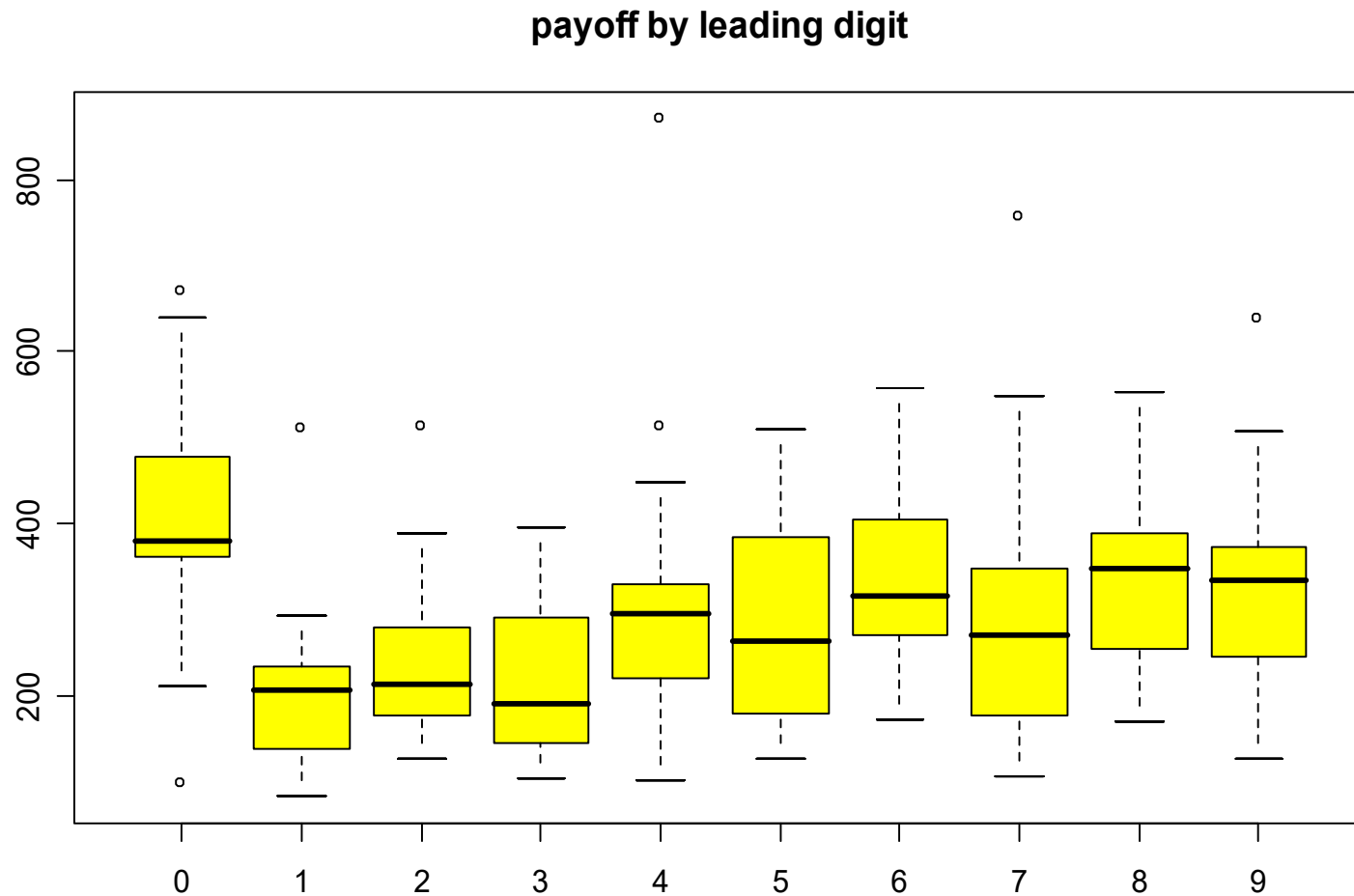
and so on for 254 consecutive days

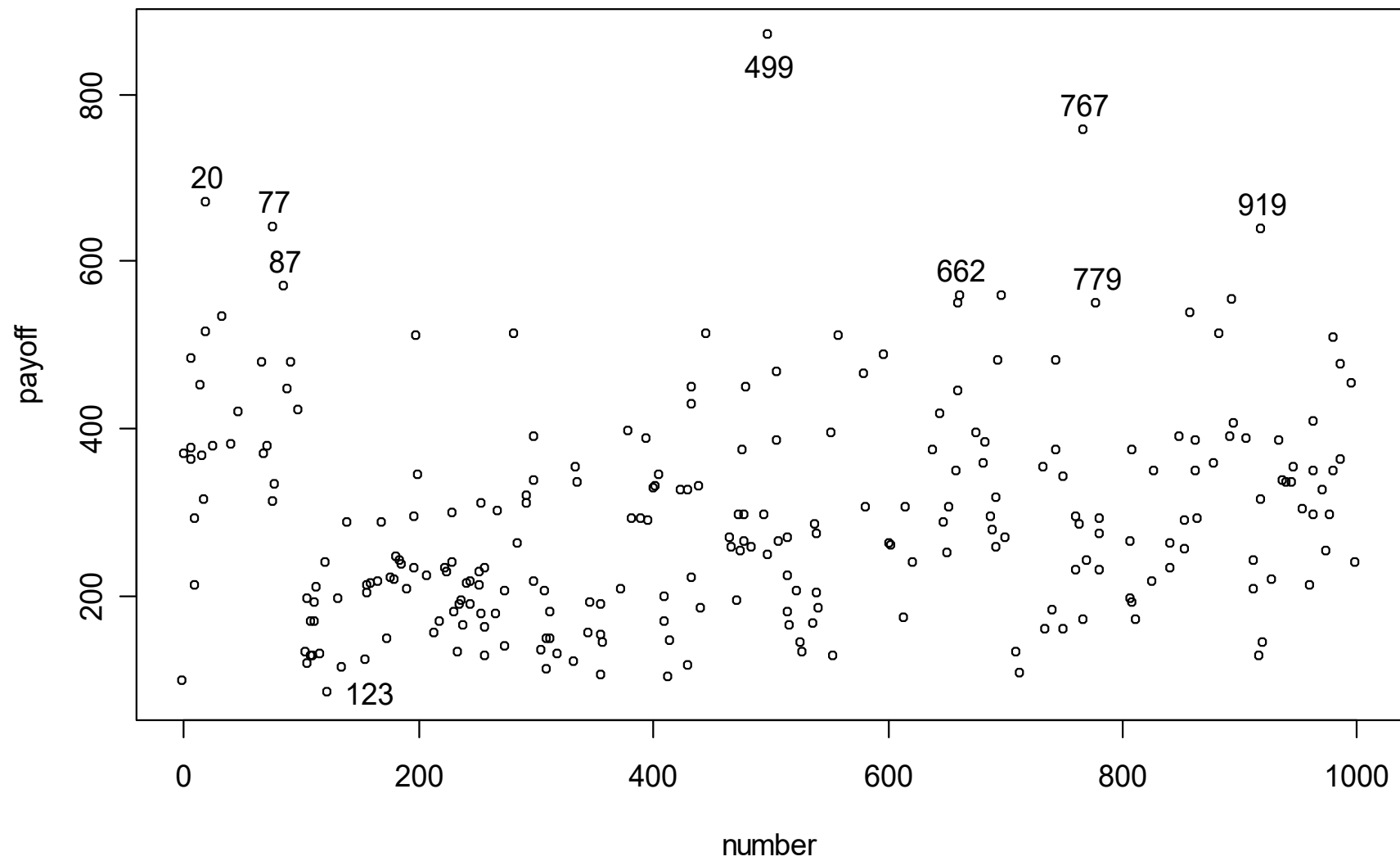
strategy 1: choose a winning number

- since we have data on the winning numbers, see if there's a pattern we can exploit to pick the winners
- examine the distribution of winning numbers using histograms or stem-and-leafs



strategy 2: choose a winning number that few others pick





why little circles?

- easier to distinguish overlapping points
- especially with jittering

five rules

- rule 1: graph lots
- rule 2: use what the eye is good at (and avoid what the eye is bad at)
- rule 3: find the right contrast and show it
- rule 4: make it easy to spot pattern, and deviations from pattern
- rule 5: plot models, not just the data

rule 1: graph lots

- only one out of 50 graphs will “work” so to get a handful of workable graphs, graph lots
- good graphing principles help raise your yield of workable graphs
- better if you can generate lots of simple graphs quickly even if they’re not perfect
- for you, not for presentation (at this stage) so don’t obsess on look (though I’m showing you the survivors of hundreds of graphs, so they’re cleaner and not quite representative of the messy graphs I usually produce: my working graphs usually don't have titles, clear axis labels, etc.)

rule 2: use what the eye is good at (and avoid what it's bad at)

- we need to know something about how the eye-brain perceives graphics
 - what it's good and bad at, and an ordering or hierarchy
 - “optical illusions” and traps to avoid
 - techniques to exploit strengths and minimize weaknesses

graphical perception

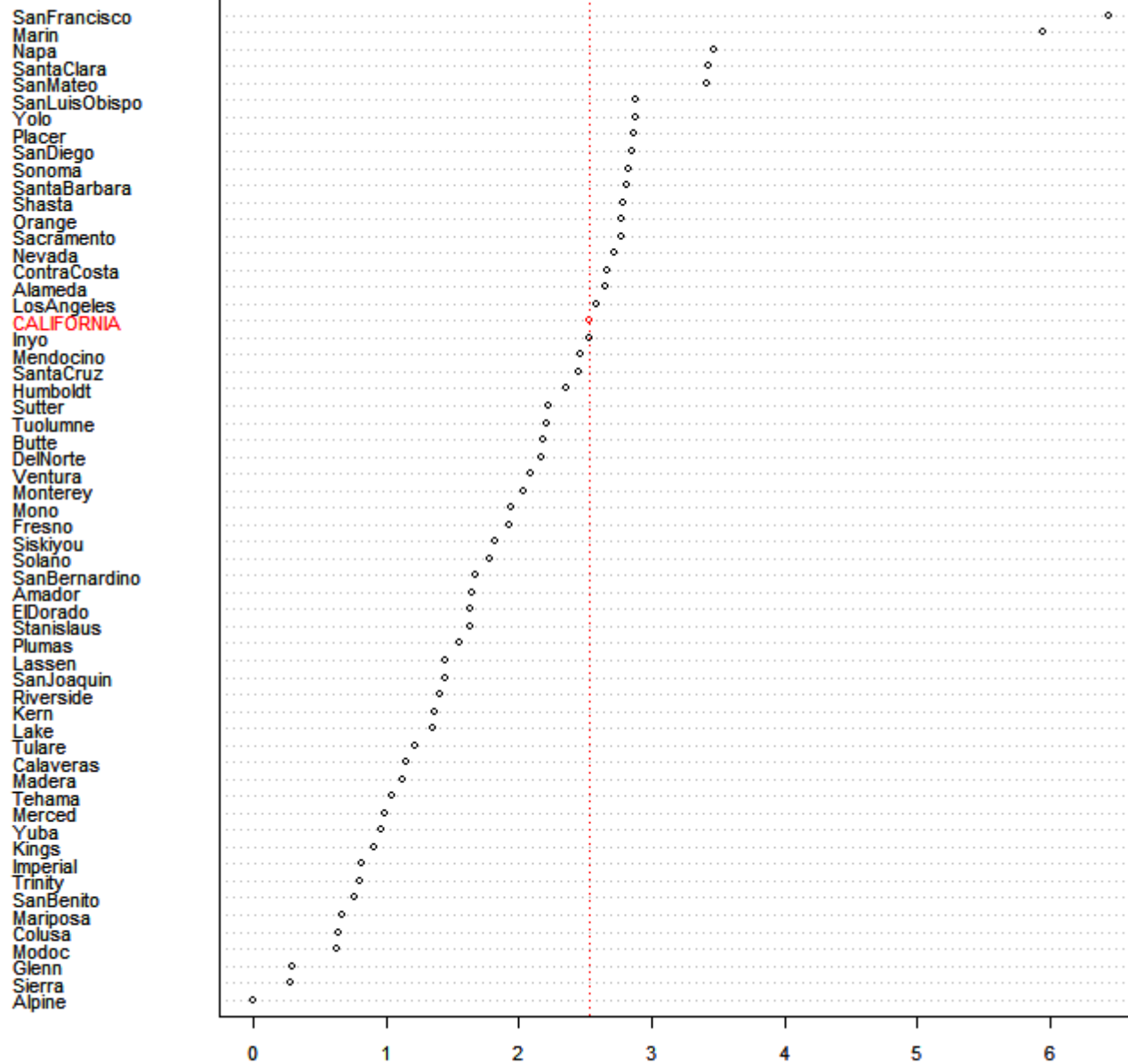
- quantitative pattern recognition by
 - detection: recognition of geometry
 - assembly: grouping of detected elements
 - estimation: assessment of relative magnitudes
- the human eye-brain can be fooled
 - optical illusions
- need to help it out
 - grouping, ordering, highlighting help to identify patterns

a hierarchy of graphical perception

- position along common scale
- position along identical non-aligned scales
- length
- angle, slope
- area
- volume
- shading, color (good discrimination but poor ordering)

dotplots

Licensed physicians per 1000 pop, California counties



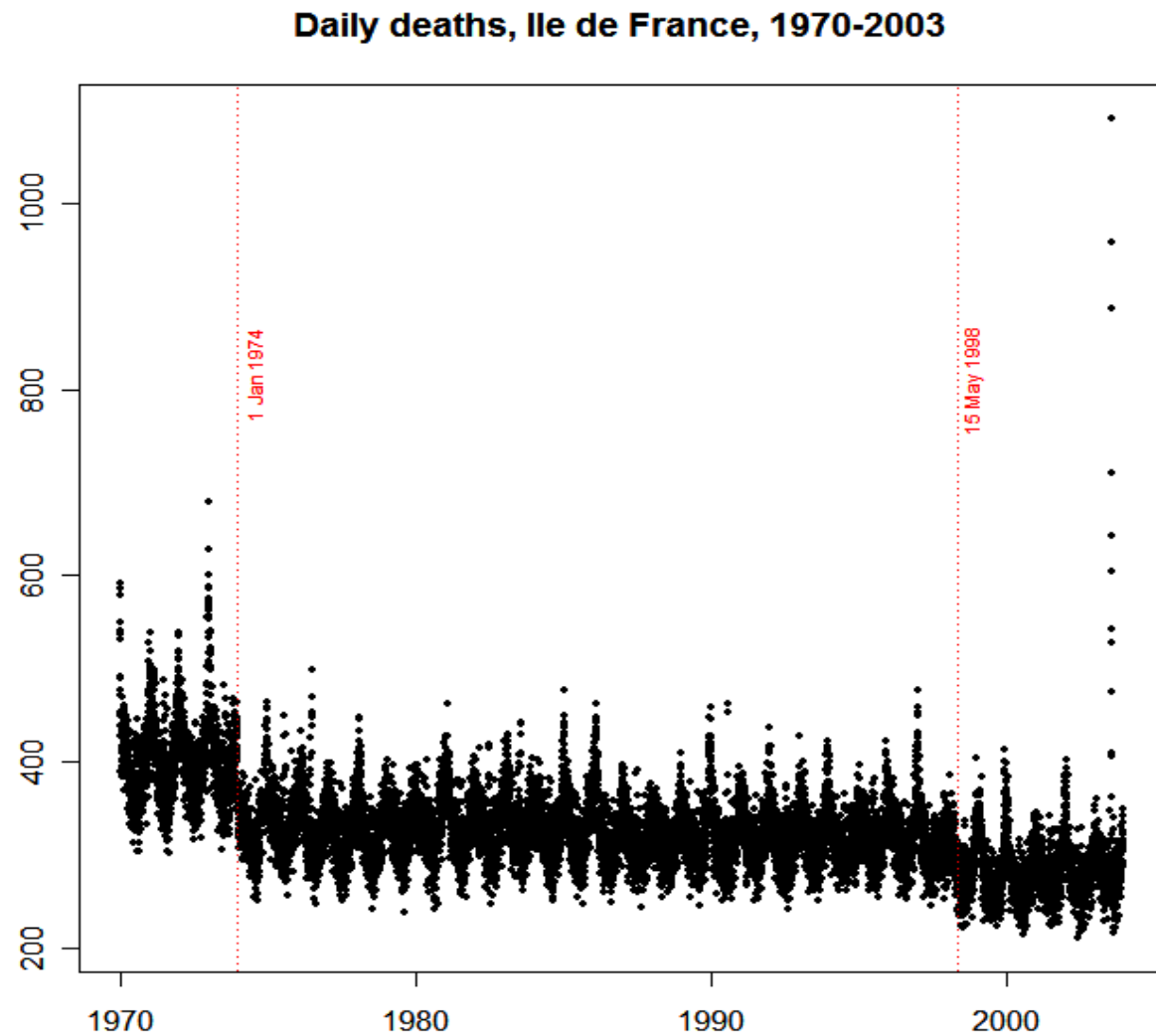
- pie charts require estimation of area
- human perception of relative areas is conservative, i.e., shrinkage toward 1.0
- shape affects estimation of area
 - concave shapes appear larger than convex
 - maps are good for context and clustering, not so good for comparisons of quantitative amounts
- color intensity affects estimation of area.
 - highly saturated colors appear larger

- human eye good at discrimination, poor at ordering
 - use for categories, not for quantitative coding
 - hues are not ordered
 - use for highlighting, patterning, especially in combination with small multiples
- more on color, later

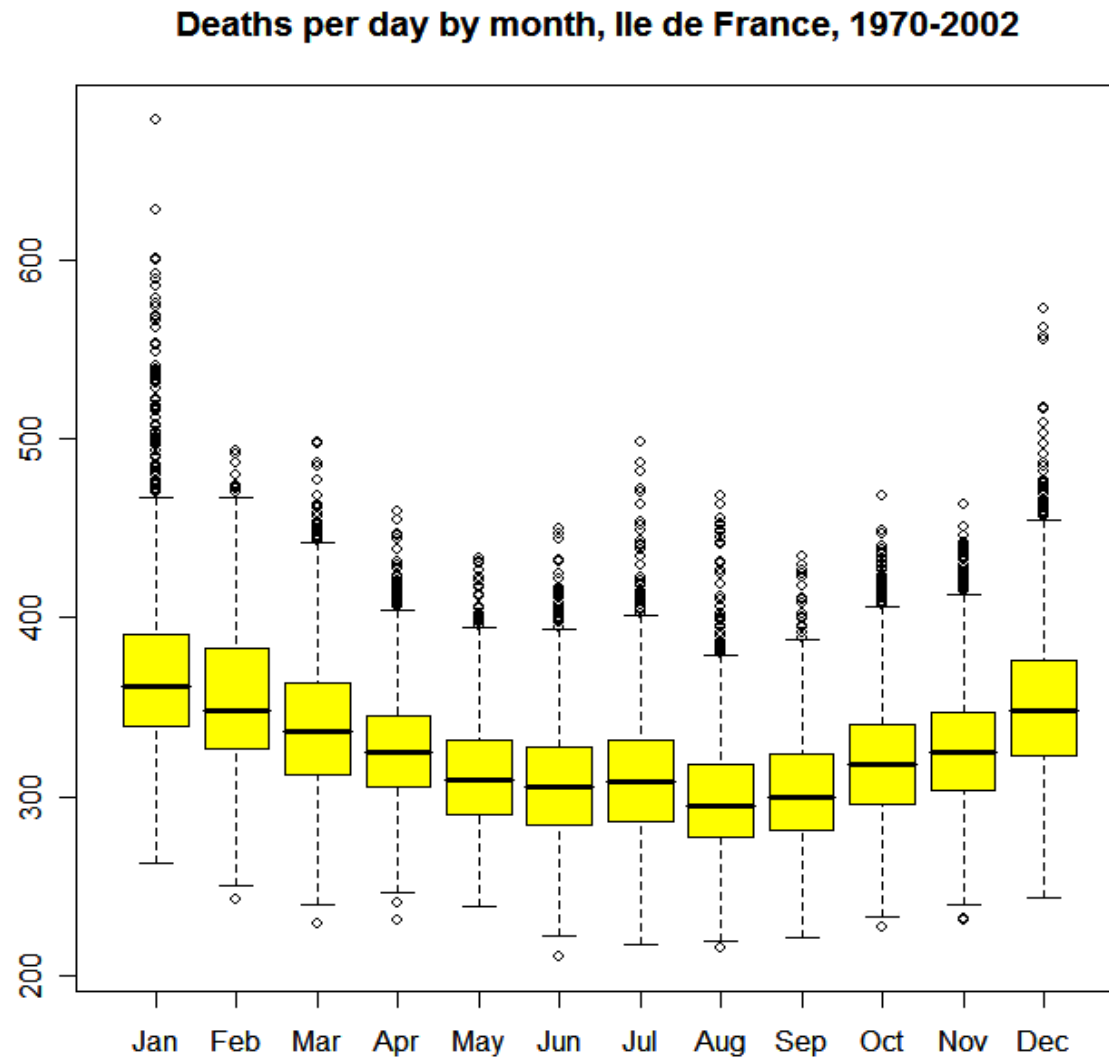
rule 3: find the right contrast and show it

- don't rely on the eye to do differencing
 - if you're interested in the difference between two lines, don't show the lines and rely on the eye to calculate the difference; calculate the difference itself and show it
- Tukey mean-difference (aka Bland-Altman) plots
 - levels vs. differences
- fits and residuals
- different contrasts can give different insights

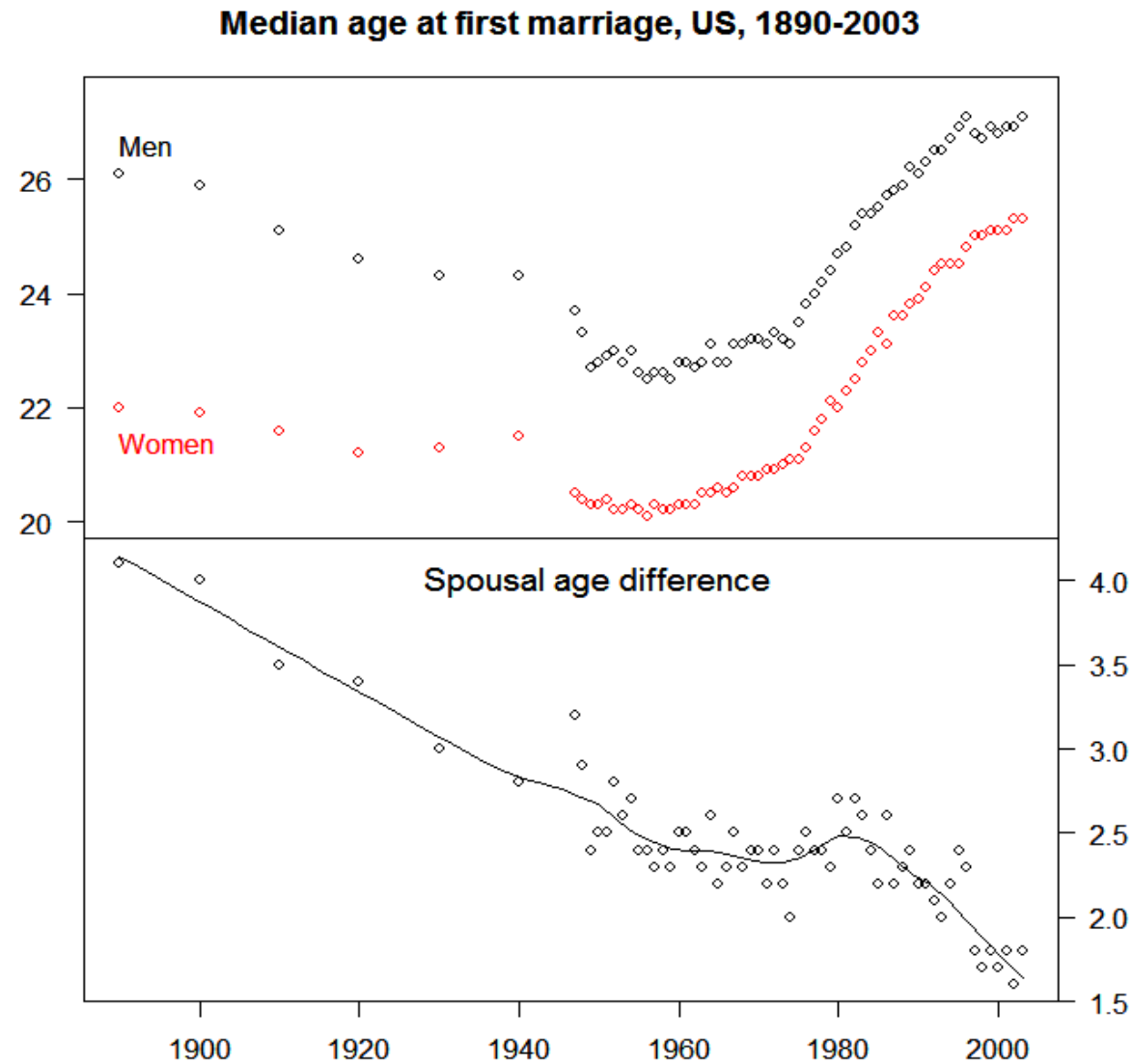
differences by time



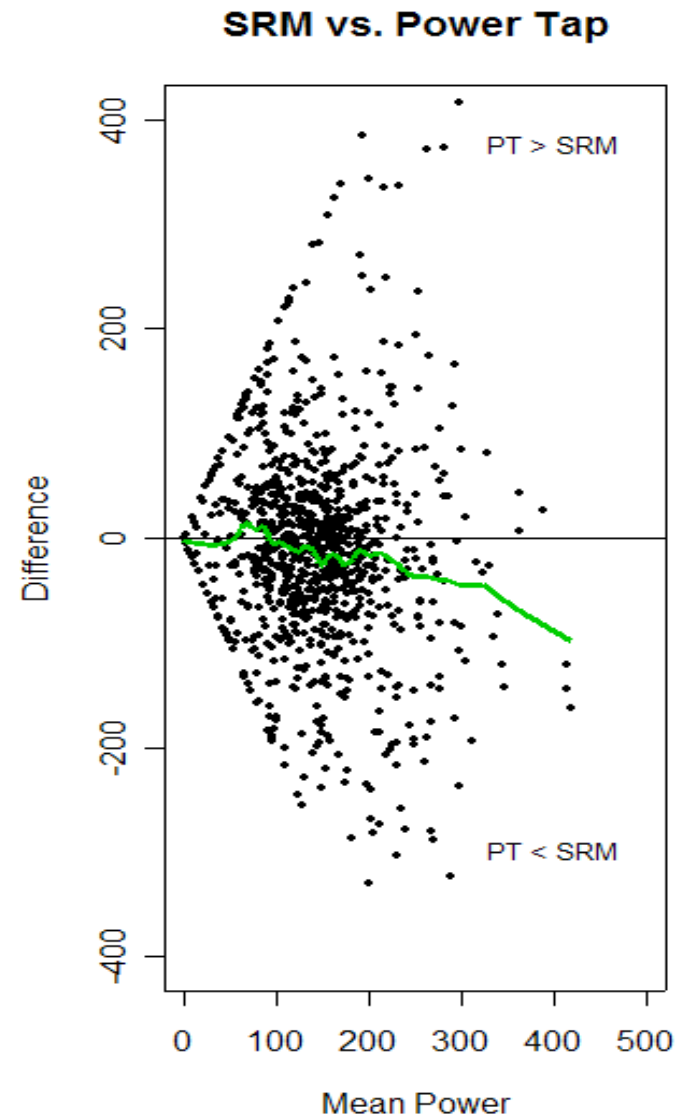
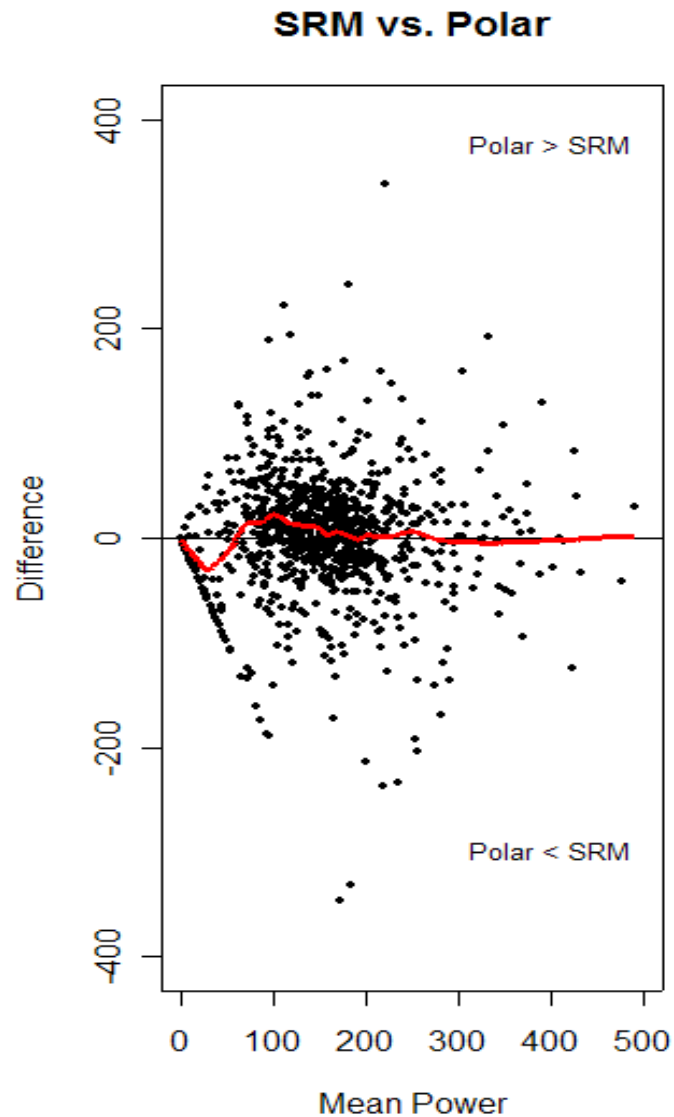
differences by category



differences between lines



differences and means

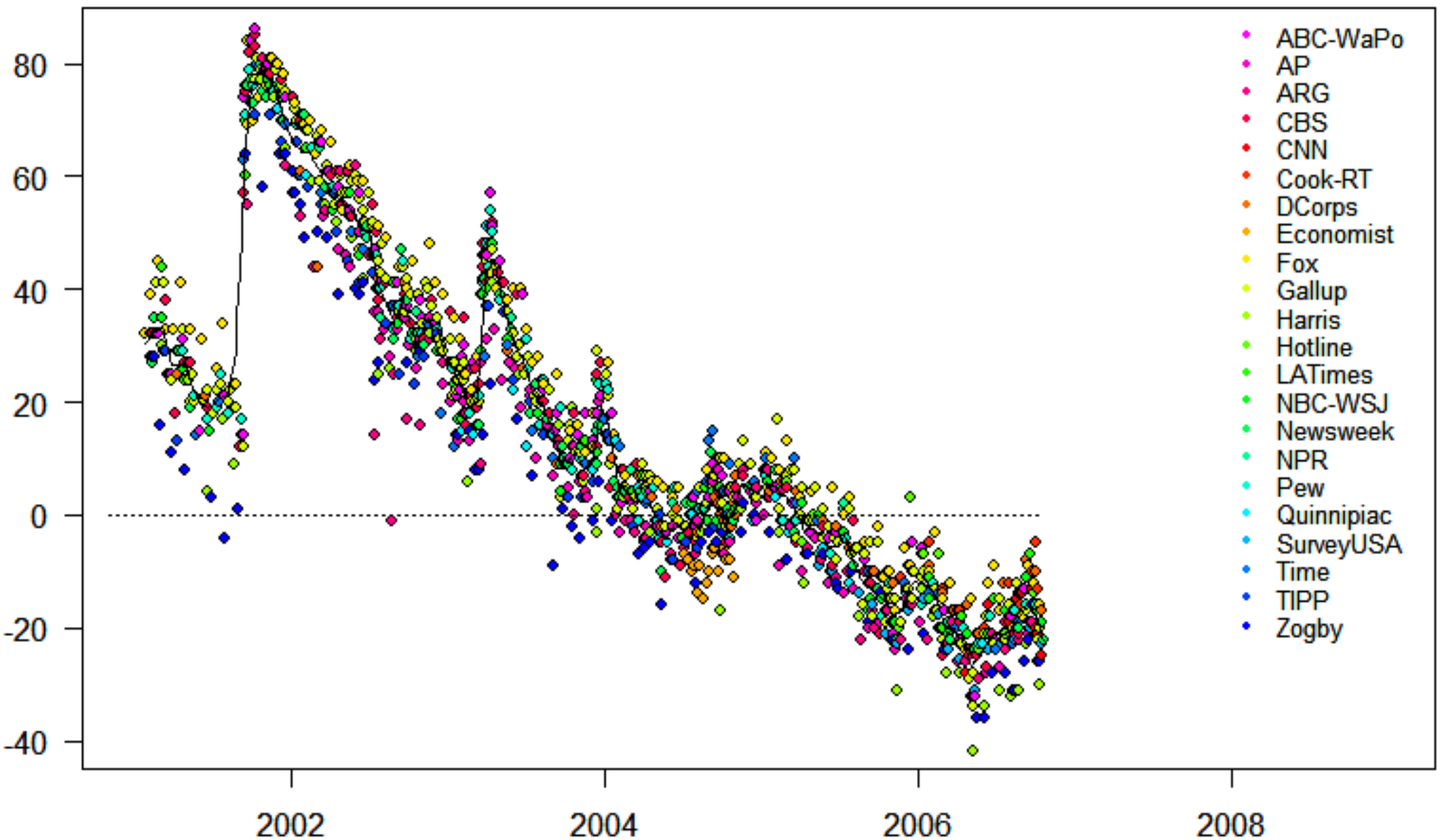


differences from fits

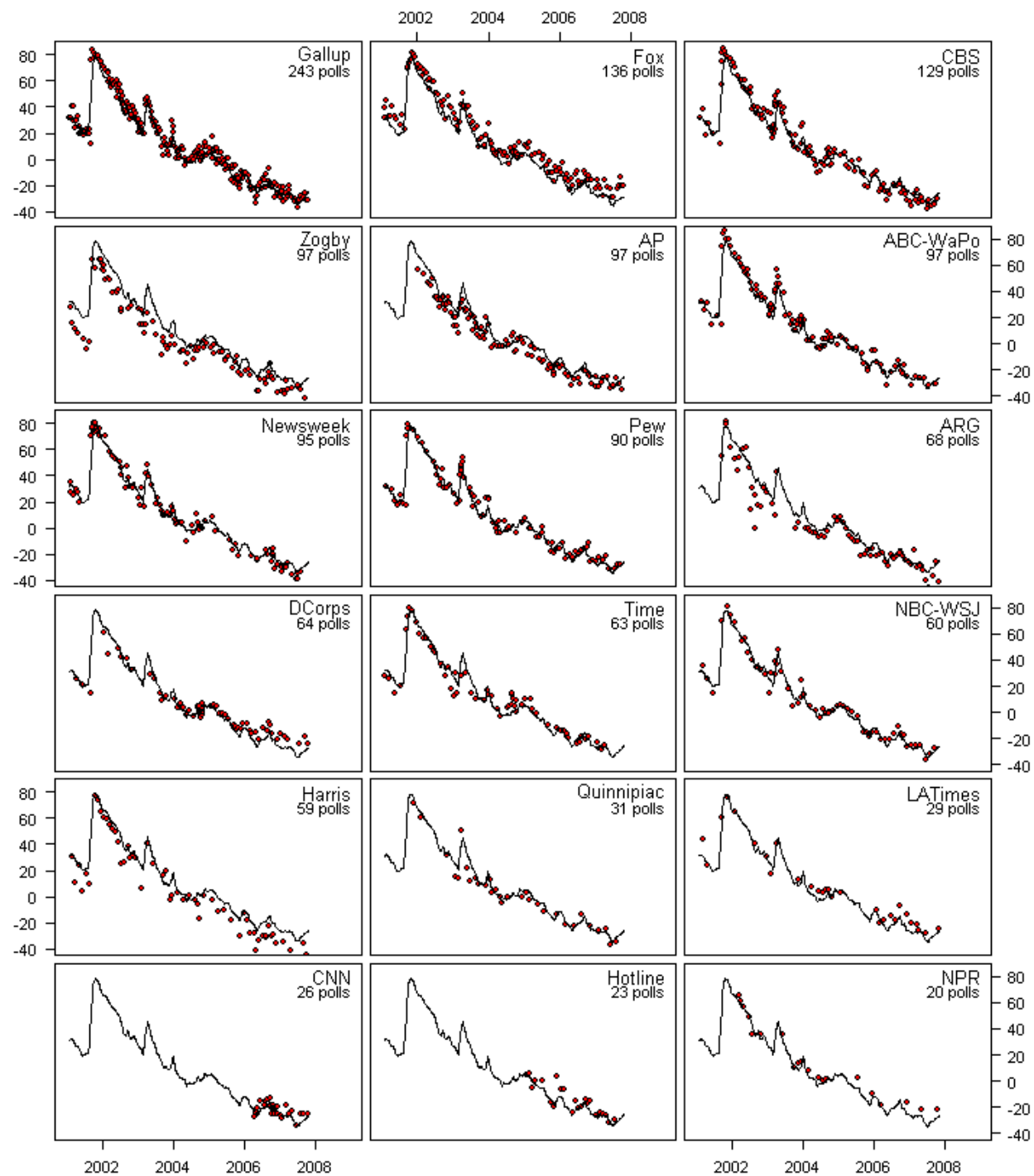
- $\text{data} = \text{fit} + \text{residual}$
- the classic residual plot
 - “fit” can be broadly defined: $\text{data} = \text{smooth} + \text{residual}$

Bush job performance ratings

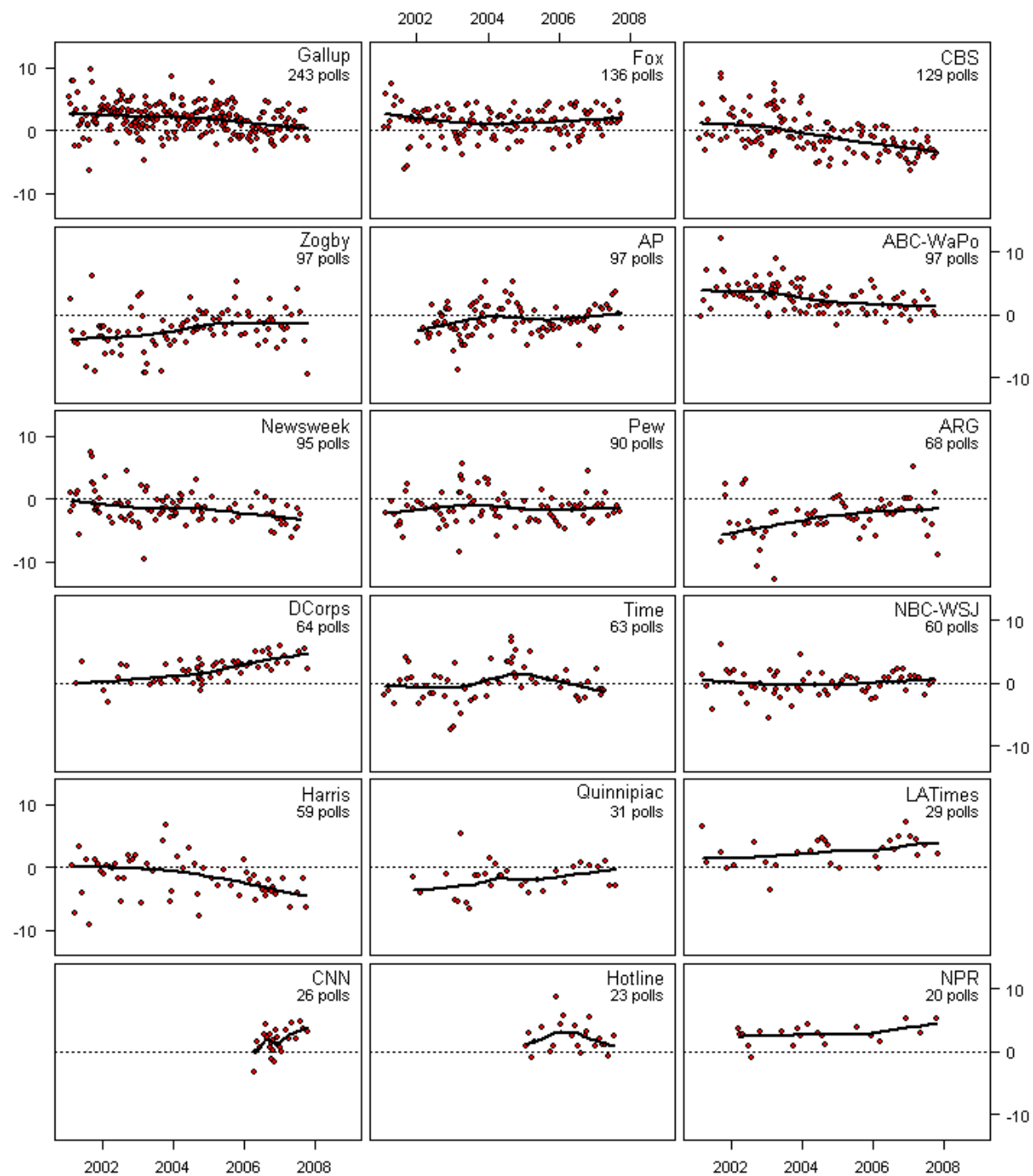
Approval - disapproval spread through 20 Oct 2006



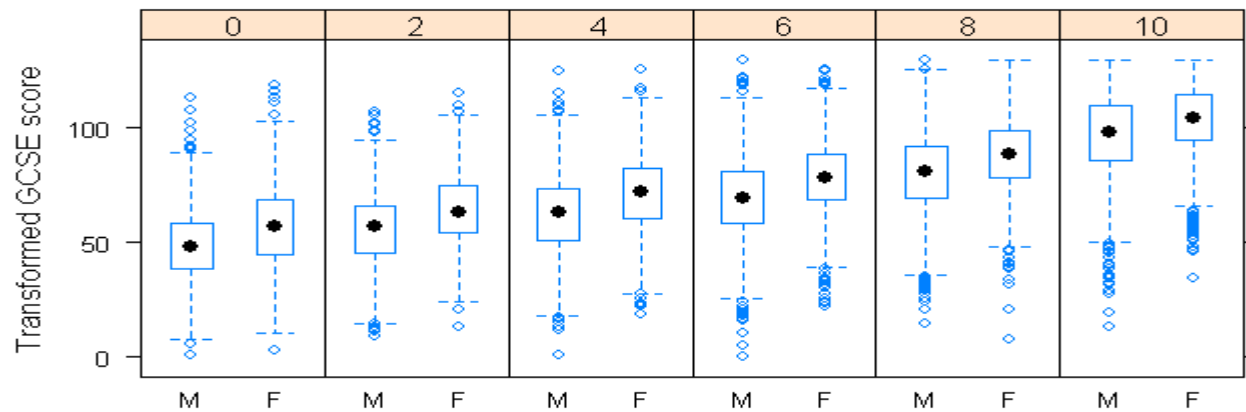
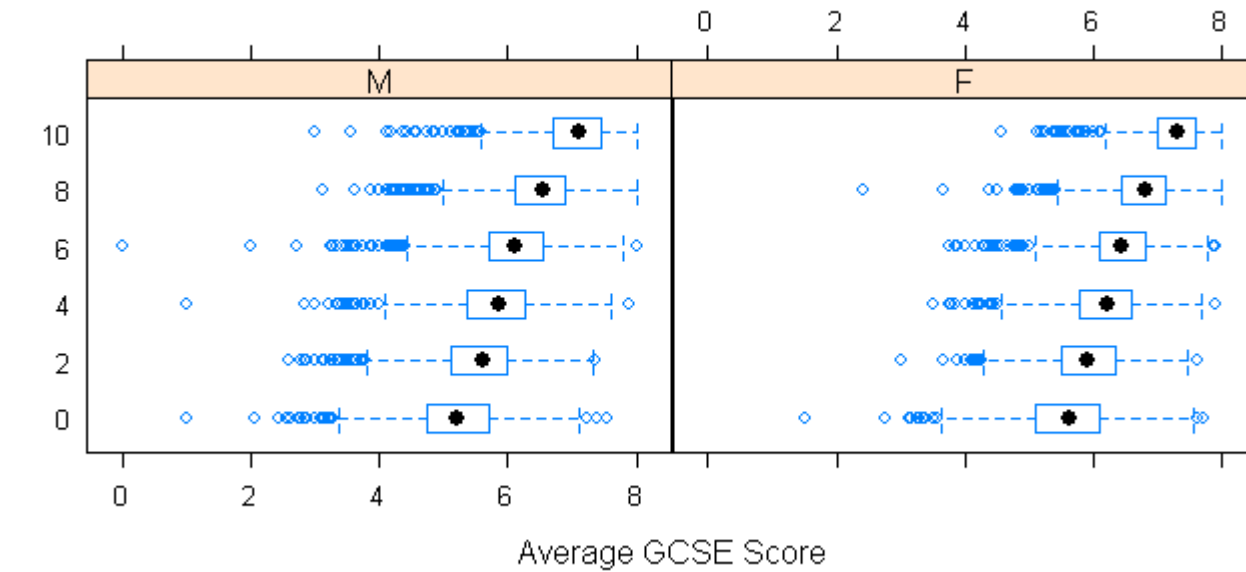
Net job approval for President Bush



Some polls consistently below trend; others above



same data, different contrast

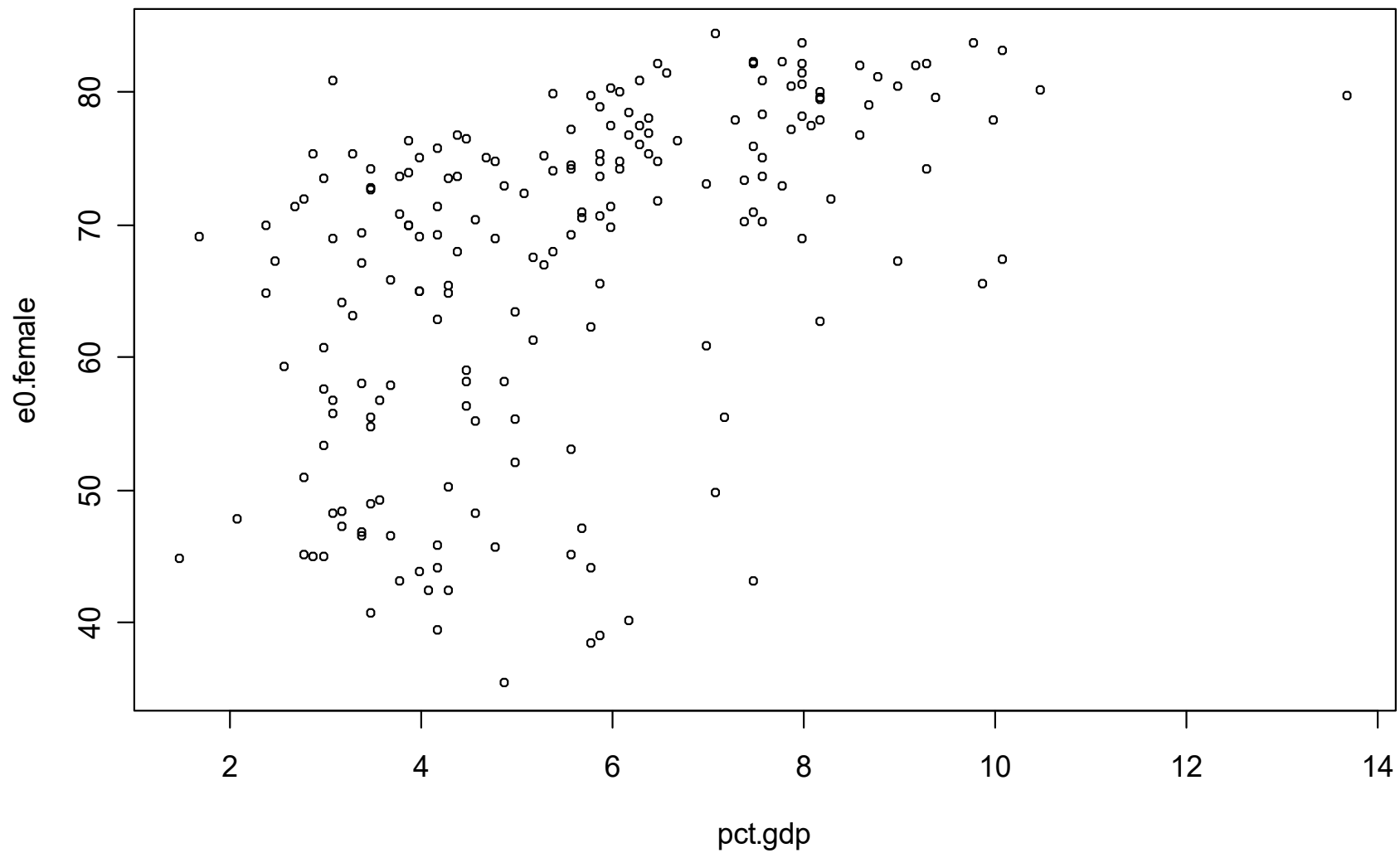


rule 4: make it easy to spot pattern

- add information depth, not (unnecessary) complexity
sometimes two plots are better than one complex plot (and sometimes it isn't)

direct labeling

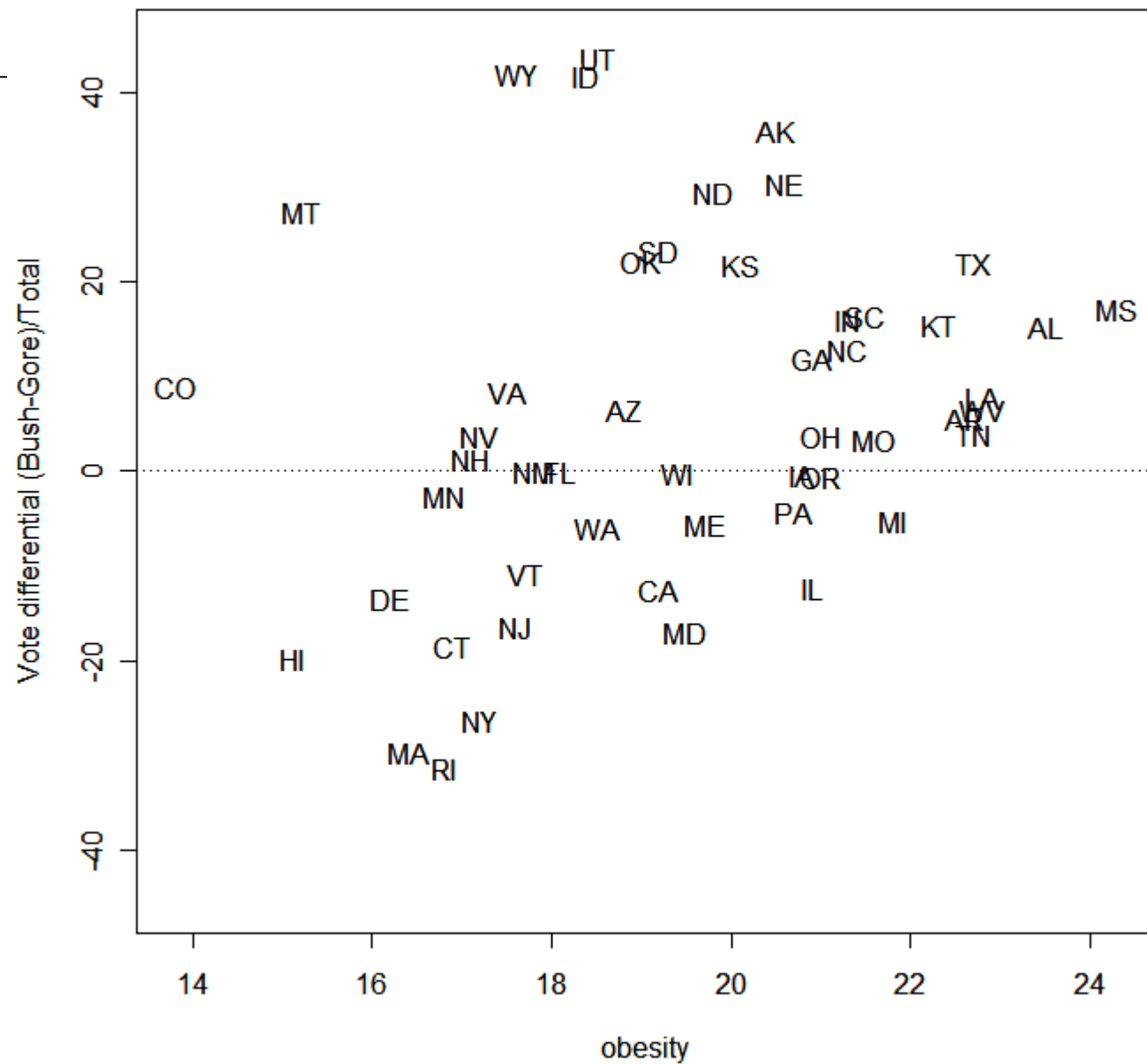
life expectancy and pct of gdp spent on healthcare



A scatter plot showing the relationship between the percentage of GDP (pct.gdp) on the x-axis and the expected life expectancy at birth for females (e0.female) on the y-axis. The x-axis ranges from 0 to 14, and the y-axis ranges from 40 to 80. The plot shows a positive correlation, with countries like Japan, Monaco, and France having high values for both variables, and countries like Sierra Leone and Malawi having low values for both.

Country	pct.gdp (approx.)	e0.female (approx.)
Japan	6.5	82
Monaco	7.5	82
France	9.5	82
Switzerland	10.5	81
Germany	10.5	79
USA	13.5	79
Singapore	2.5	80
Norway	6.0	80
Canada	7.5	79
Italy	8.5	79
Spain	7.5	78
Finland	6.5	78
Denmark	6.5	77
Sweden	7.0	77
Belgium	7.0	77
Netherlands	7.5	77
Austria	8.0	77
Slovenia	9.0	77
Uruguay	9.5	77
Colombia	8.5	70
Venezuela	3.5	75
Ghana	3.5	55
Sierra Leone	4.0	35
Malawi	5.5	38
Zimbabwe	6.0	38
Botswana	4.5	38
Niger	3.5	40
Uganda	4.0	42
Kenya	4.5	45
Guinea-Bissau	5.5	45
Lesotho	5.5	45
Mozambique	5.5	45
Namibia	7.5	45
SAfrica	6.5	50
Cambodia	7.5	55
Bhutan	6.5	60
India	5.5	60
Bolivia	5.5	62
Tuvalu	5.5	65
Marshall Is	8.5	65
Kiribati	9.5	65
Lebanon	10.5	65
Nicaragua	7.5	68
Honduras	7.5	68
El Salvador	7.5	68
Costa Rica	8.5	68
Panama	7.5	68
Bosnia	7.5	68
Suriname	8.5	68
Moldova	7.5	68
Peru	5.5	68
Chad	4.5	50
Angola	3.5	48
Madagascar	2.5	48
Somalia	1.5	45
Djibouti	2.5	45
Algeria	3.5	65
Syria	2.5	65
Libya	3.5	65
Egypt	3.5	65
Saudi Arabia	3.5	65
Yemen	3.5	58
Laos	3.5	55
Guinea	3.5	55
Sudan	3.5	55
Benin	3.5	52
Togo	3.5	48
Myanmar	2.5	58
NKorea	3.5	60
Viet Nam	4.5	68
Iran	4.5	68
Morocco	4.5	68
Israel	4.5	68
Malta	5.5	78
Cyprus	5.5	78
Brunei	5.5	78
Maldives	8.5	62

2000 Vote vs. self-reported obesity percentage



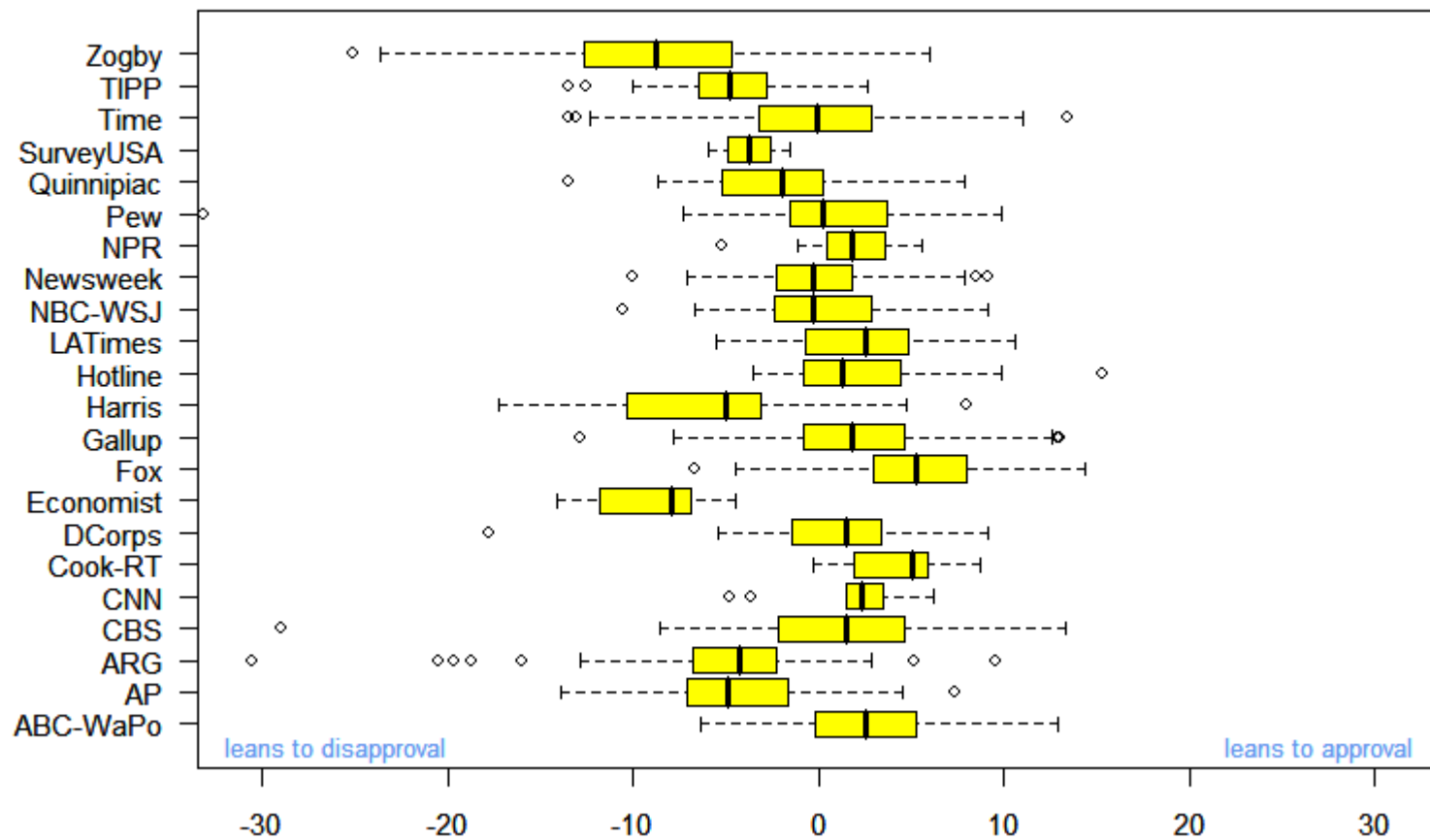
more on labeling

- direct labeling of lines often better than legends
 - particularly good when combined with line color
 - symbol plus line type often too busy to decode
 - looking back-and-forth at a legend is distracting

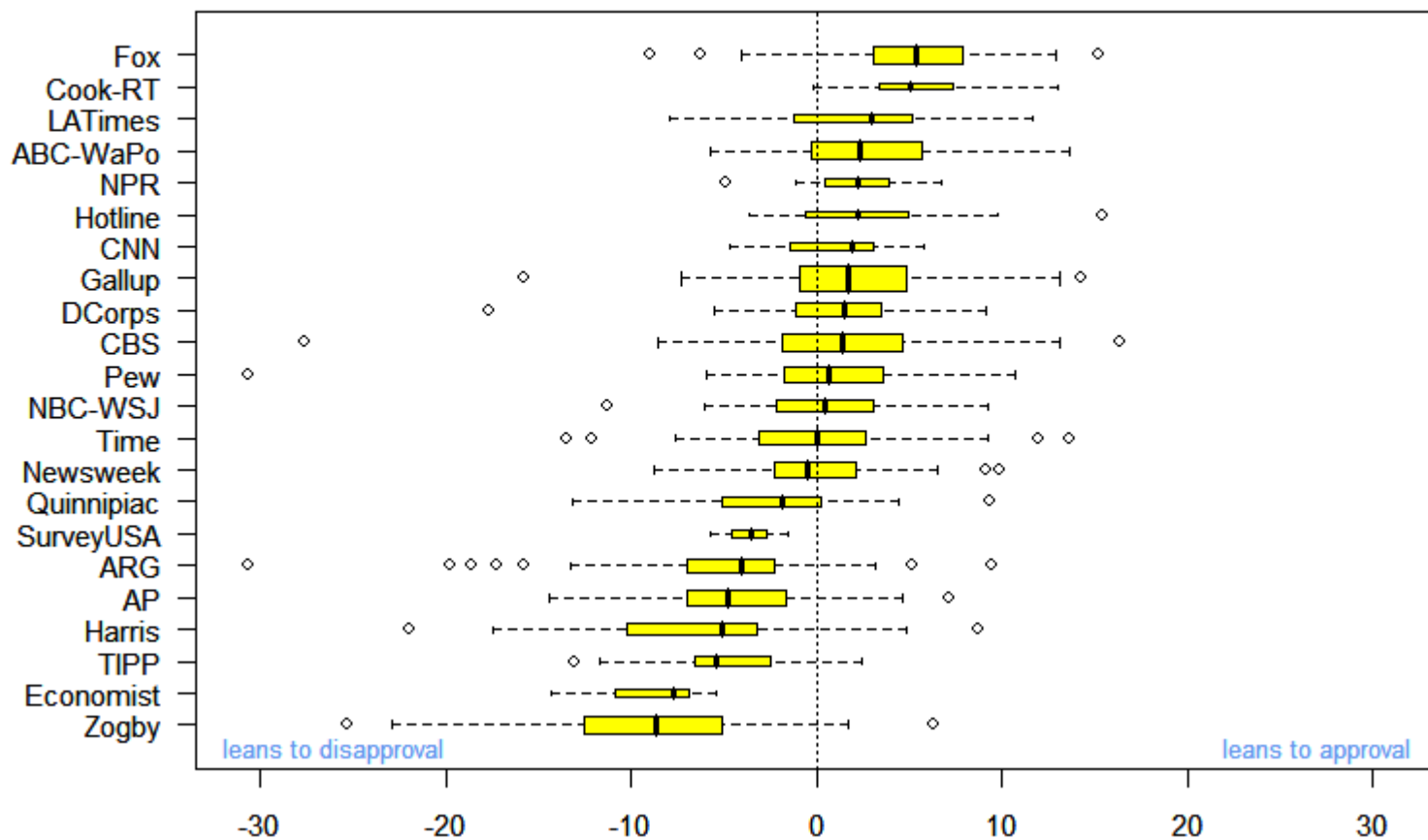
ordering

- default ordering for categorical variables is often alphabetical
that makes categories easy to find, but hard to compare
example: country data are often ordered by name of country
rather than by the variable you're interested in
- find an ordering that makes sense and use it
if you are interested in mortality differences among countries,
order by mortality not country name
this helps you spot and evaluate small differences between
countries

Differences from trend, net job performance



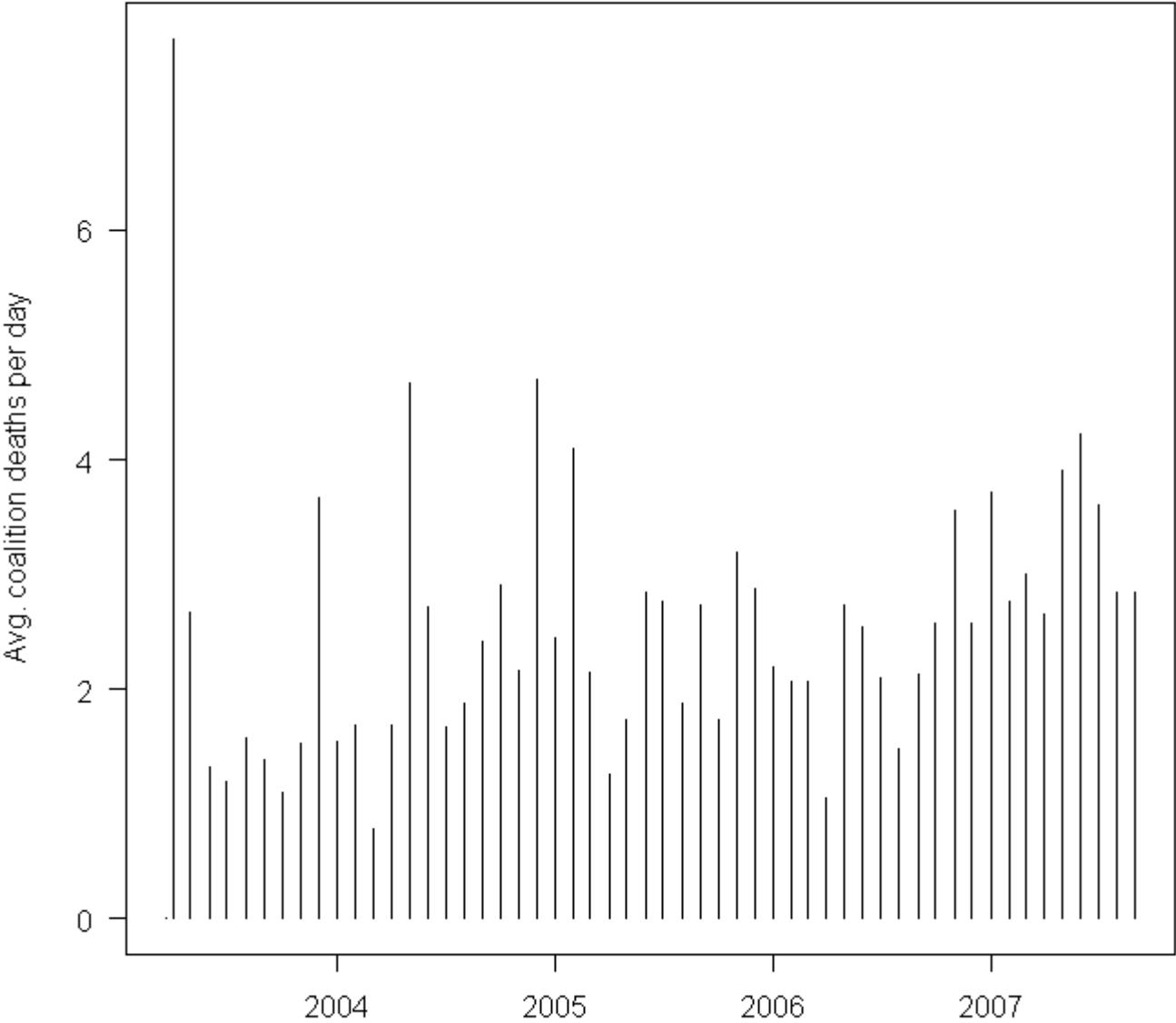
Differences from trend, net job performance



grouping

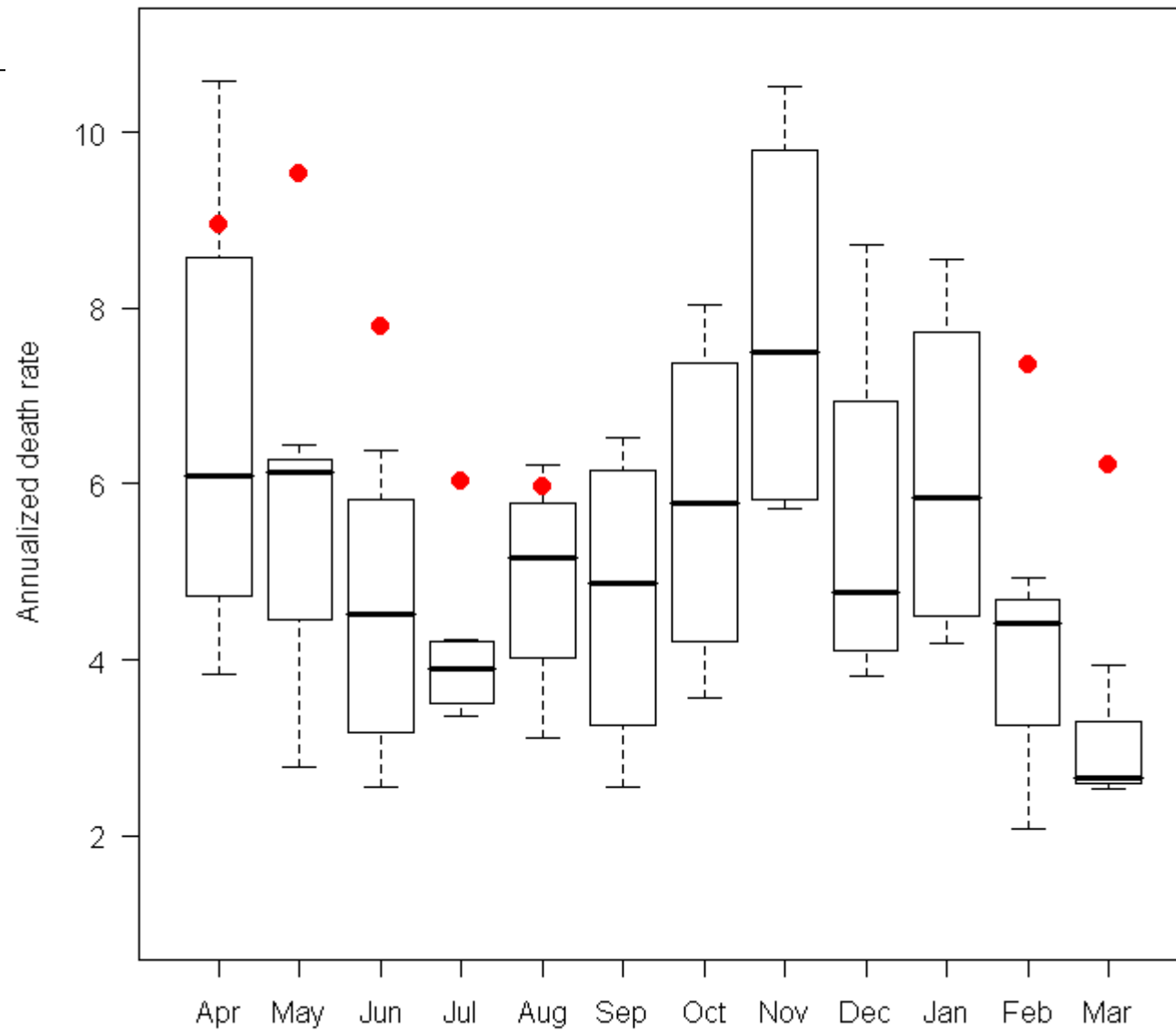
- grouping (done well) helps with pattern recognition
boxplots are a familiar way to group
- grouping (done poorly) obscures pattern
not all “obvious” groupings are informative
- next two slides show (almost) same data

Coalition deaths in Iraq, by month



Death rate in Iraq, coalition forces

excluding Mar 2003



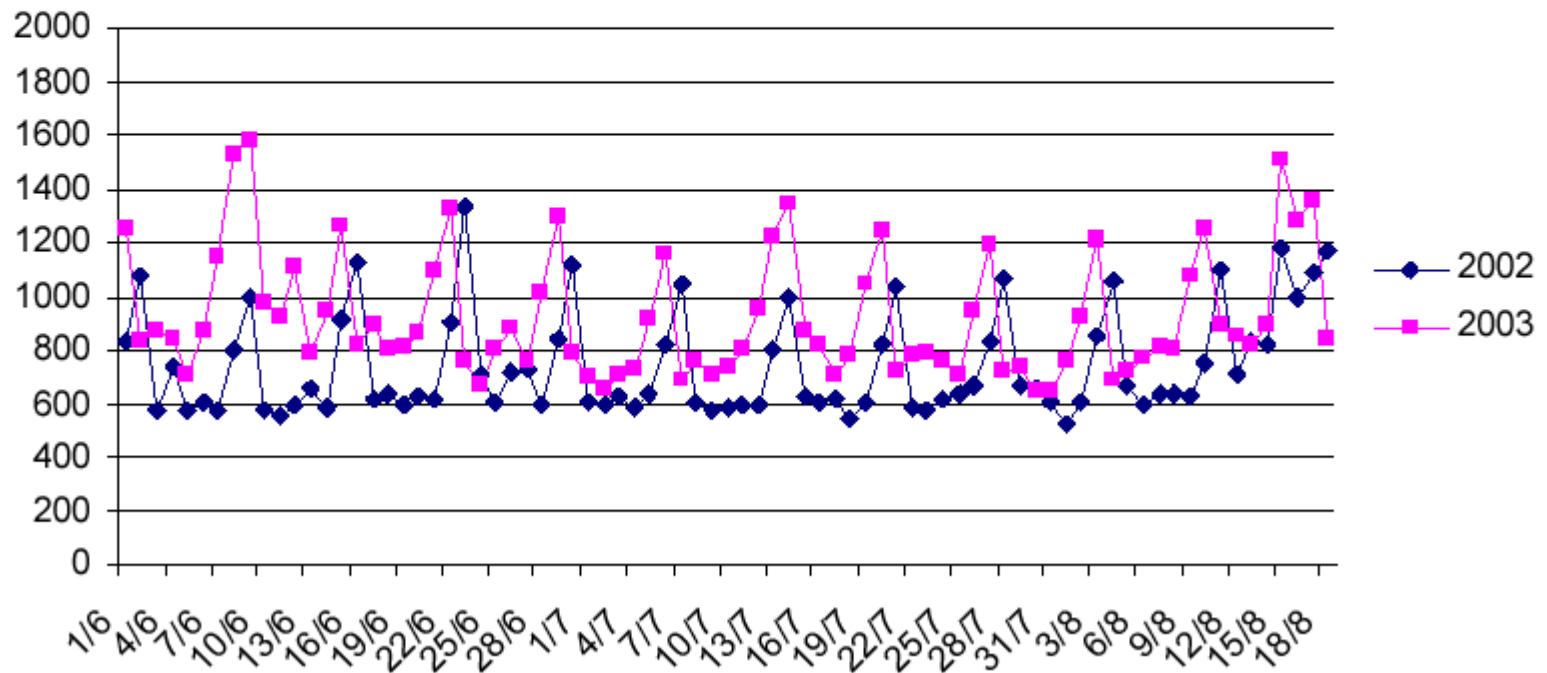
multivariate comparisons

- show relationships
more importantly, give you ideas
- time series plots show you what happened when, but rarely why they happened
we'll want to dig deeper into the data to generate new questions about the 'why?'

patterns in multivariate data

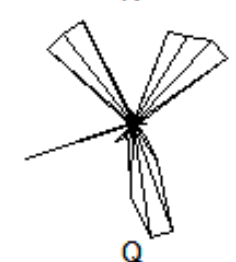
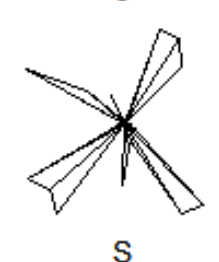
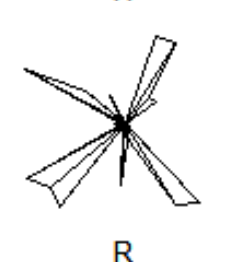
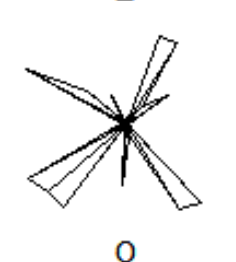
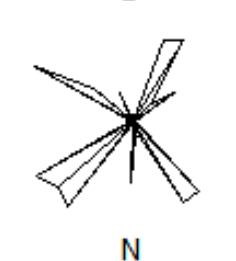
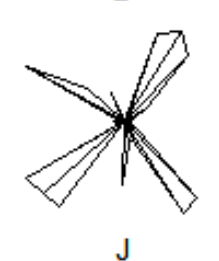
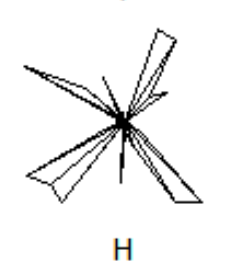
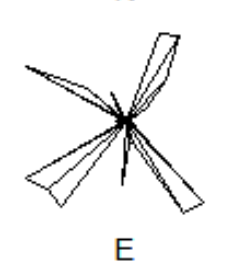
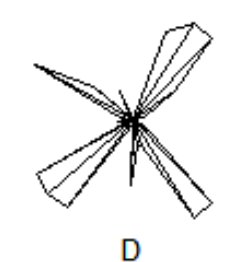
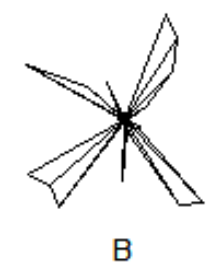
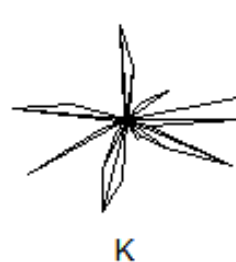
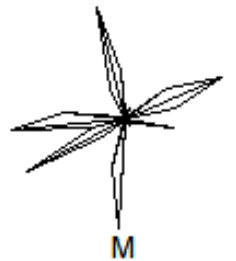
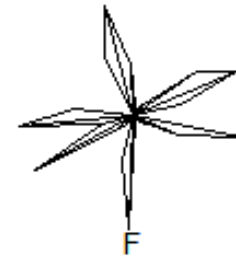
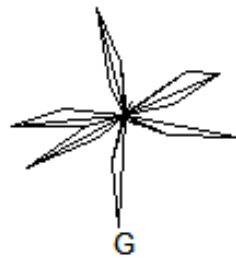
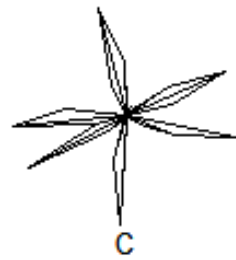
- twenty students read numbers of ambulance calls for July 2003 off a graph. How can we summarize the results?

Graphique 5 : nombre d'interventions du SAMU 13 en 2003 par rapport à l'année précédente (2002)





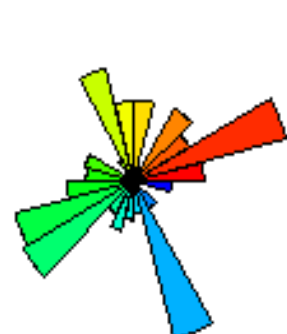
Rough grouping of PS1



dc hospitals

- thirteen hospitals
- twenty-four service lines
- do different hospitals specialize in different areas?

Hospitals and their service lines, 2000-1



Washington



Georgetown



Providence



Sibley



GWU



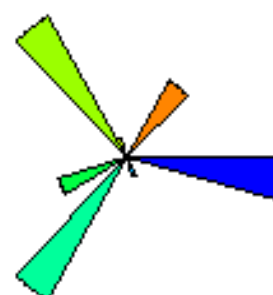
Howard



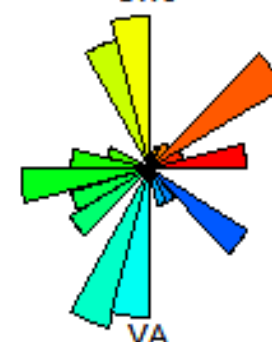
Children's



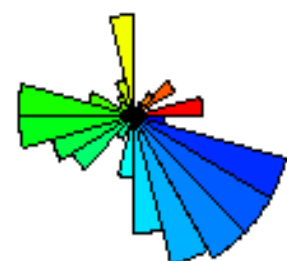
Greater SE



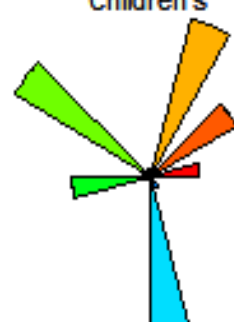
ColumbiaWomen



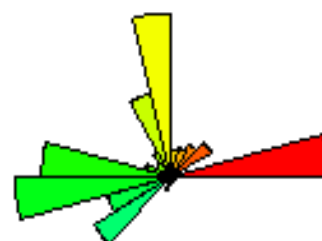
VA



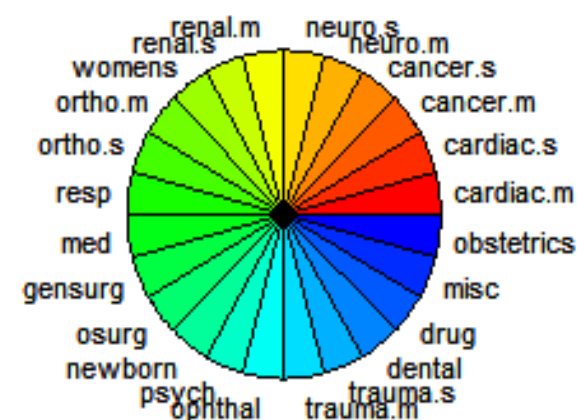
DC General



Natl Rehab



Hadley



stars are like pies

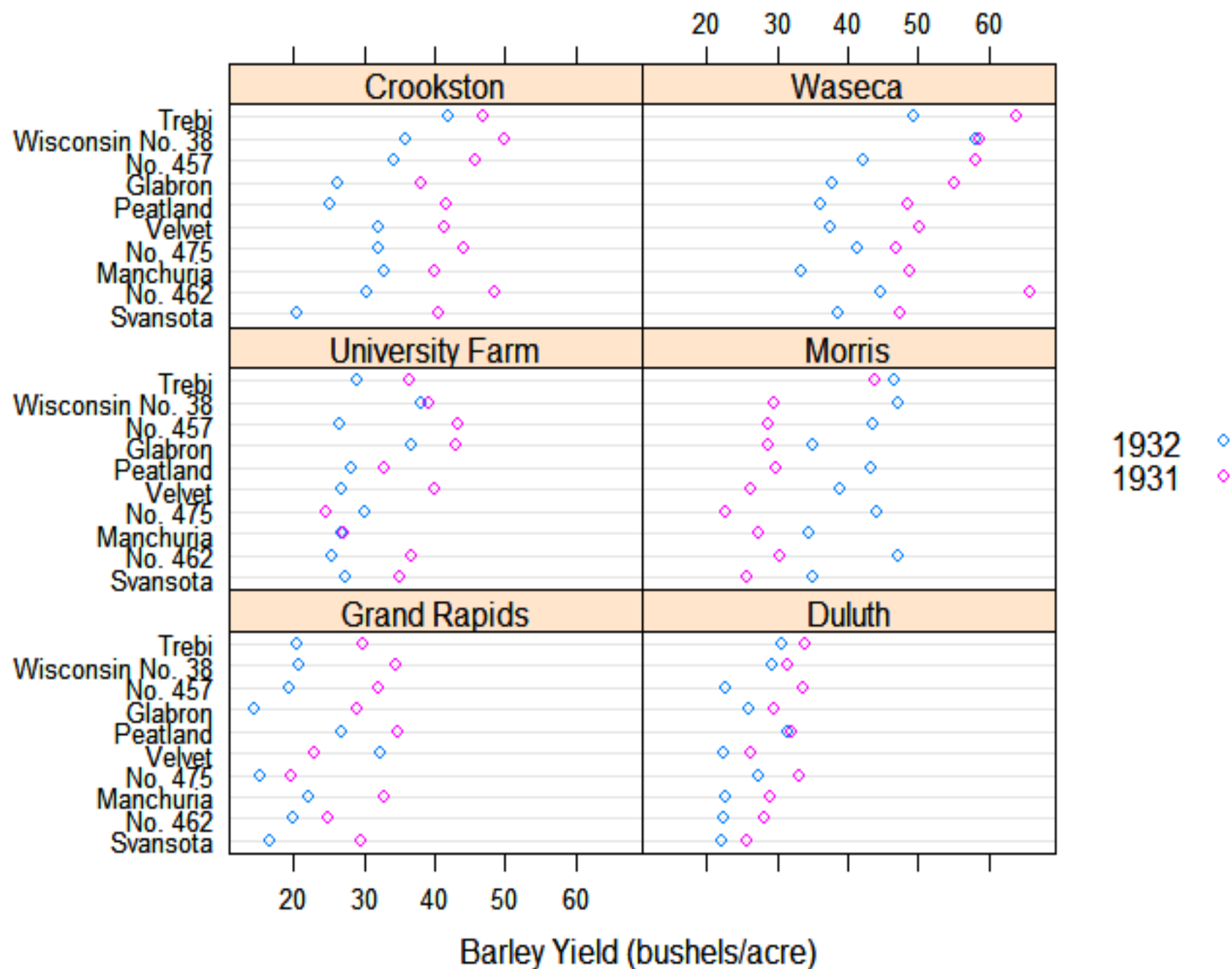
- except that angle is constant and radius varies
in pies, radius is constant and angle varies
that's why pie segments need labels
- watch out for colors

scatterplots

- can (sometimes) show more than two variables
 - can code categorical variables with color
 - can code some interval variables with size
- small multiples can show varying conditions
 - lattice (i.e., trellis) plots
 - use same scale and ranges, if possible, to enhance comparison

dotplots (and trellis)

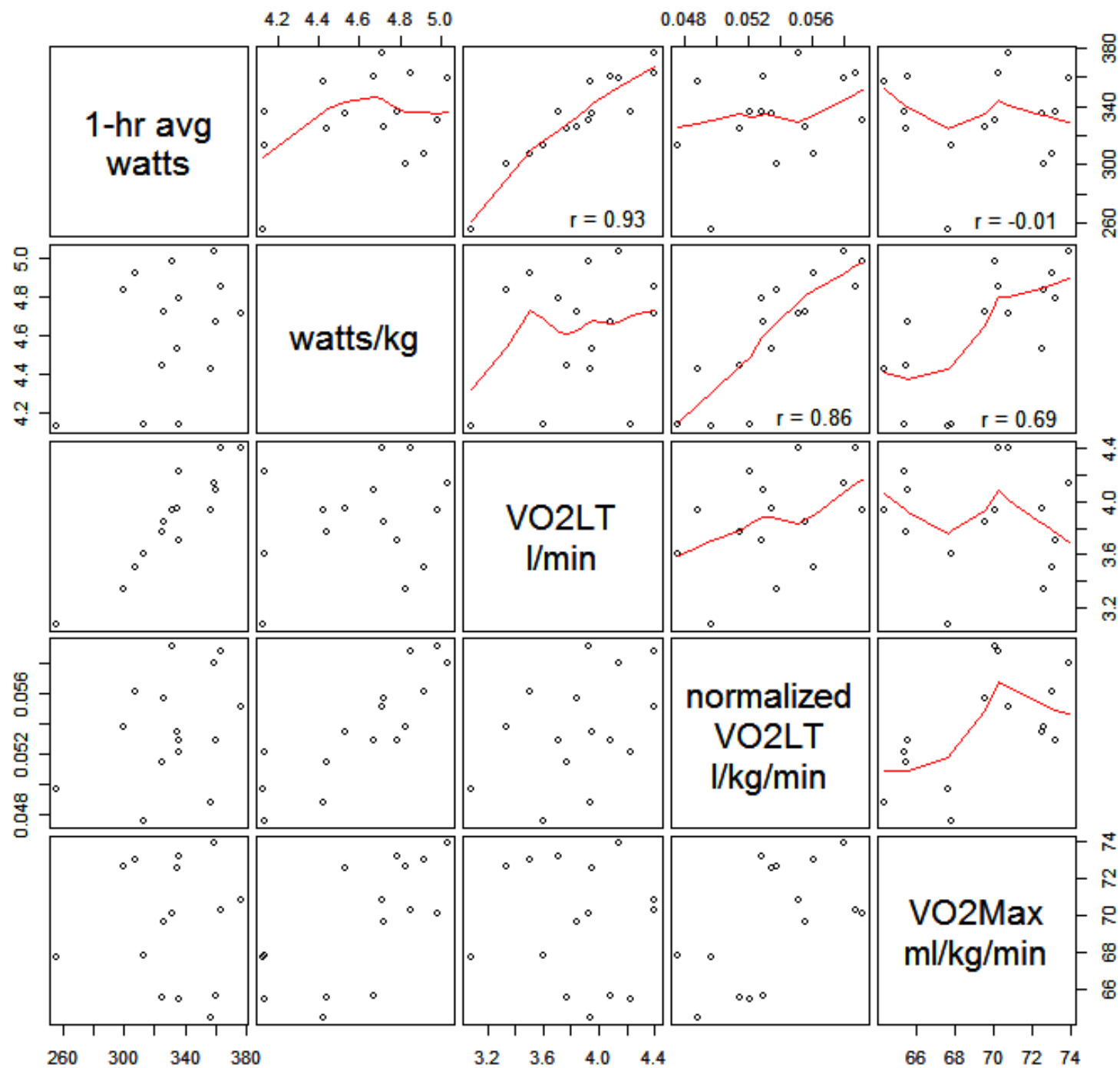
- conditioning plots
- barley yield
 - ten varieties
 - six plots
 - two years



scatterplot matrices

- scatterplot matrices compress a lot of information on bivariate relationships into a small space
 - useful for winnowing out uninteresting variables and deciding which variables might be worth further investigation

Coyle, 1991

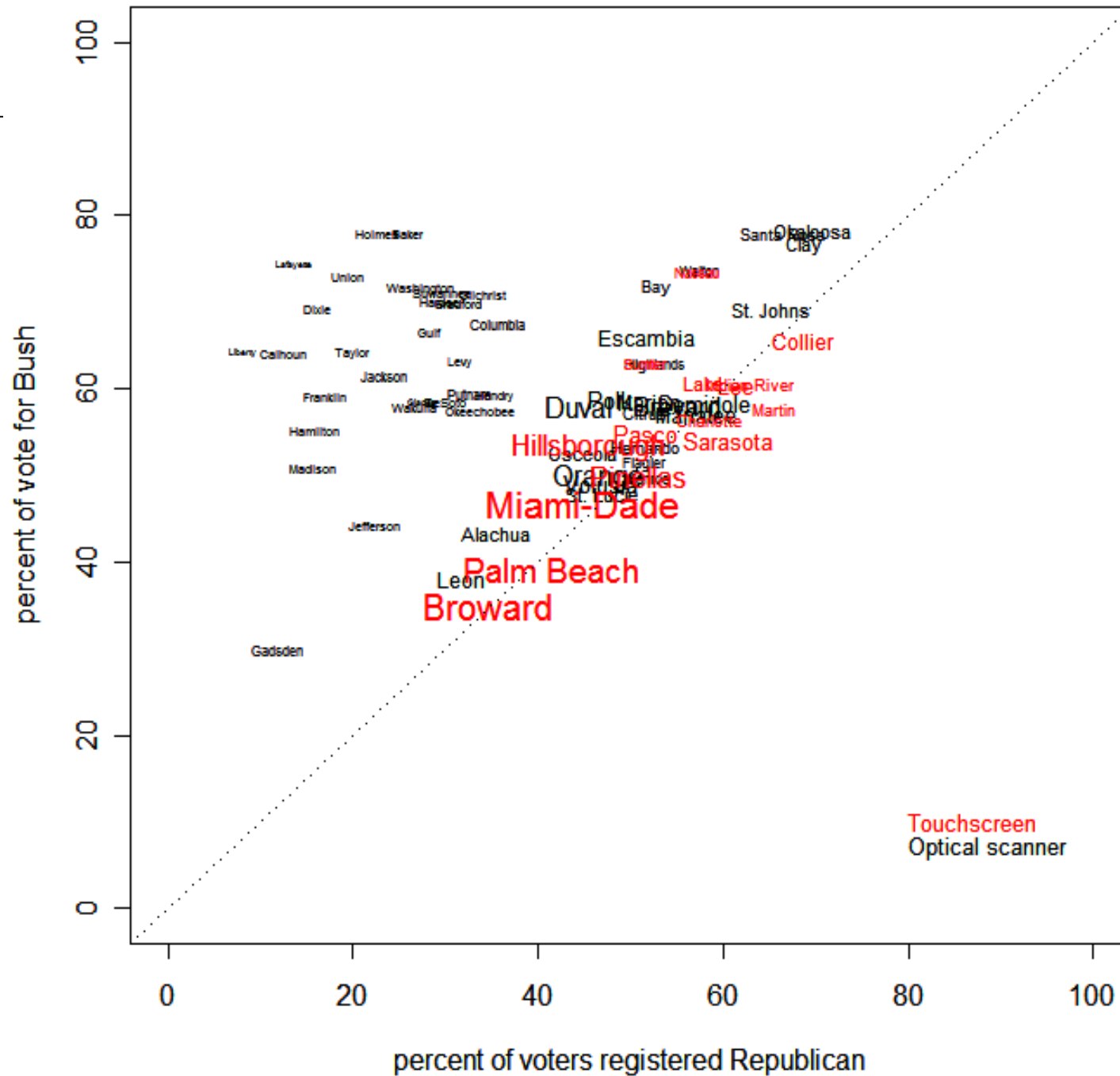


coding plotting symbols

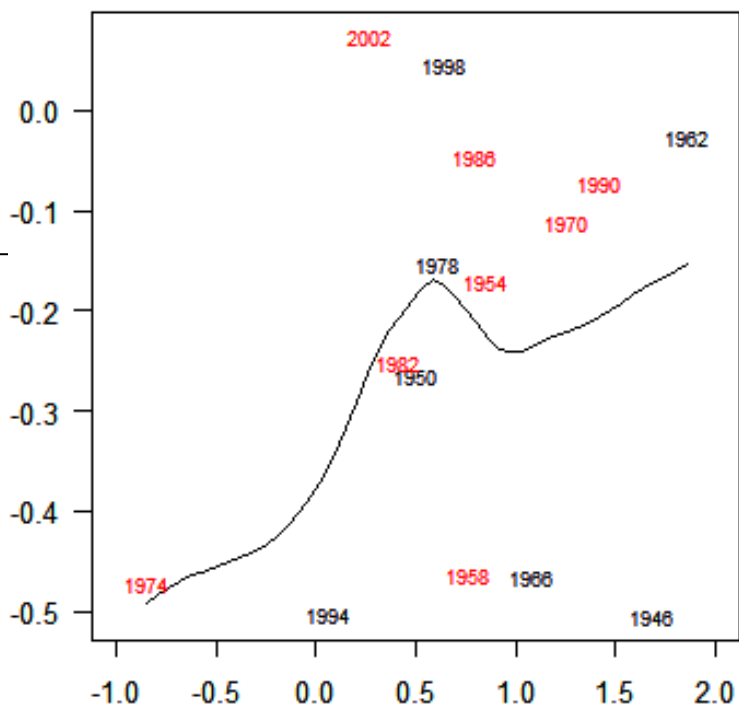
- improves information density by tagging plotting symbols with attributes

you've seen this before using color or shape; can often combine with direct labeling

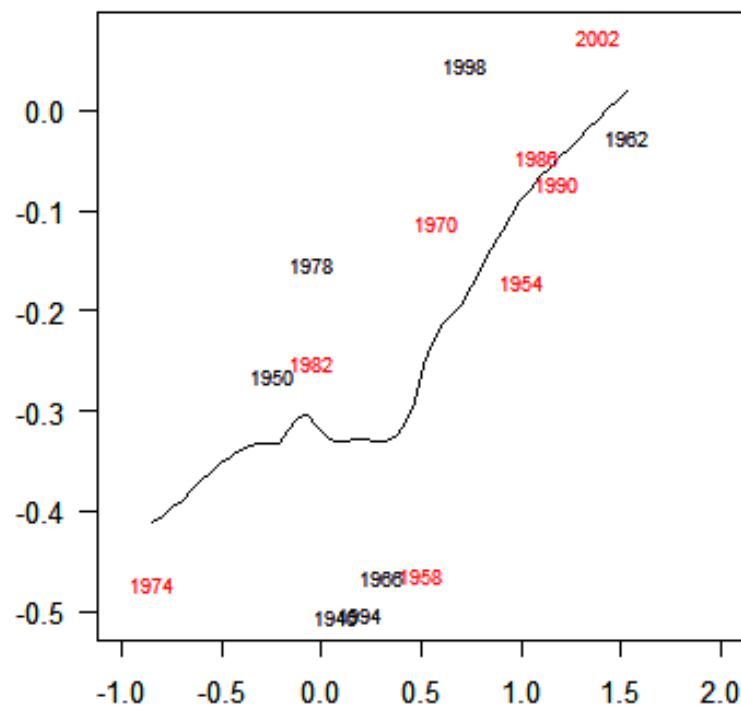
Florida vote by county



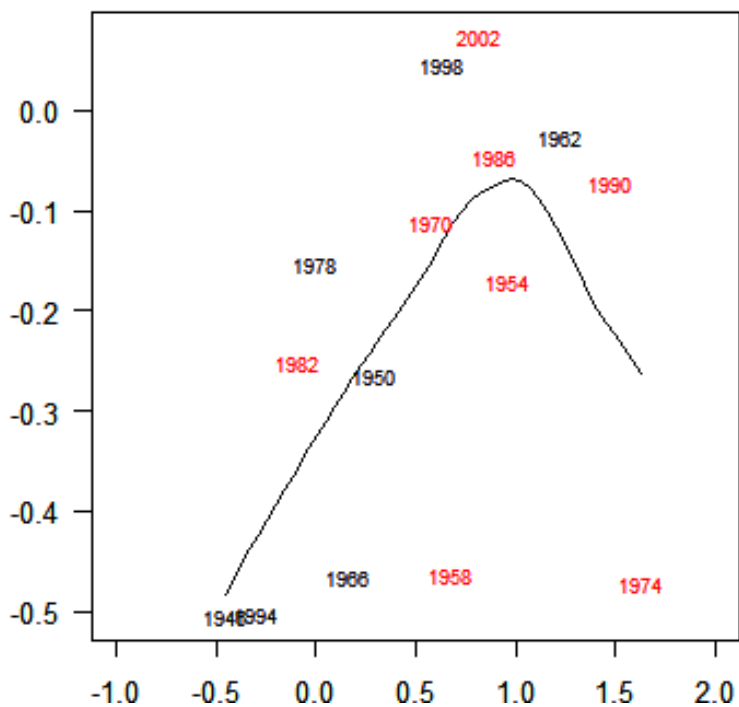
12 months before mid-term



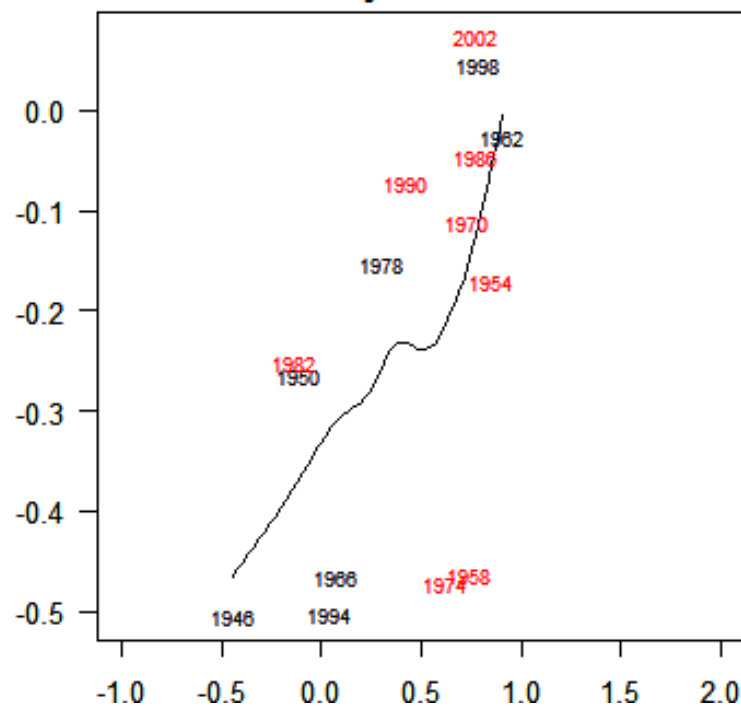
6 months before mid-term



2 months before mid-term



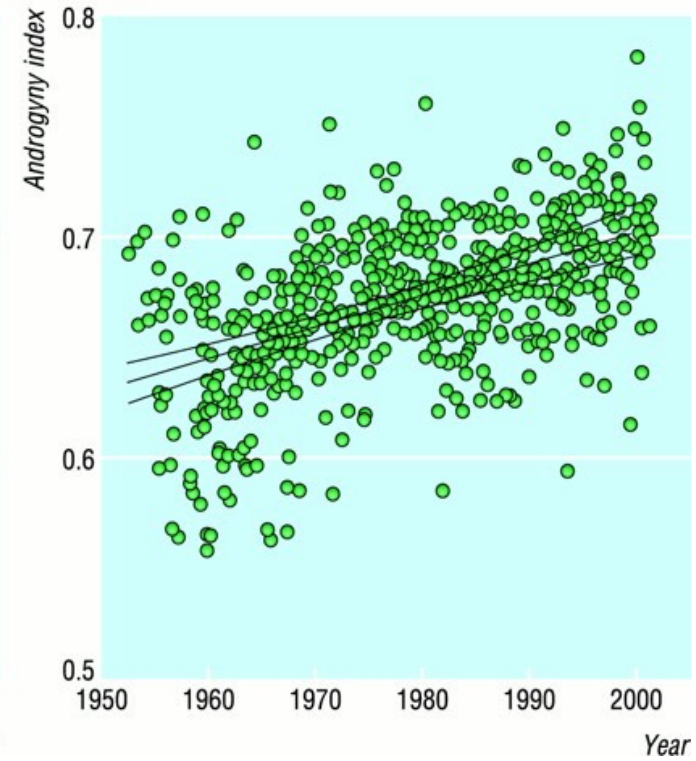
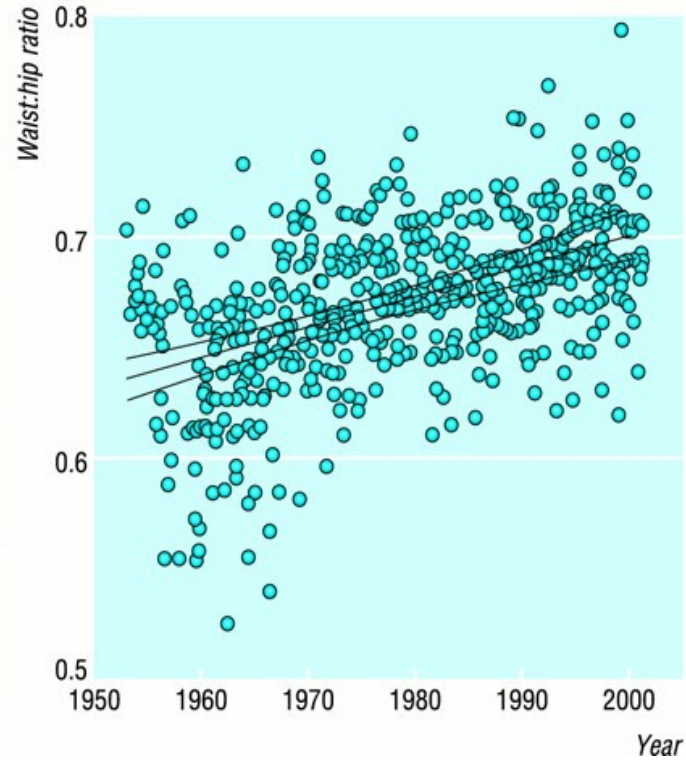
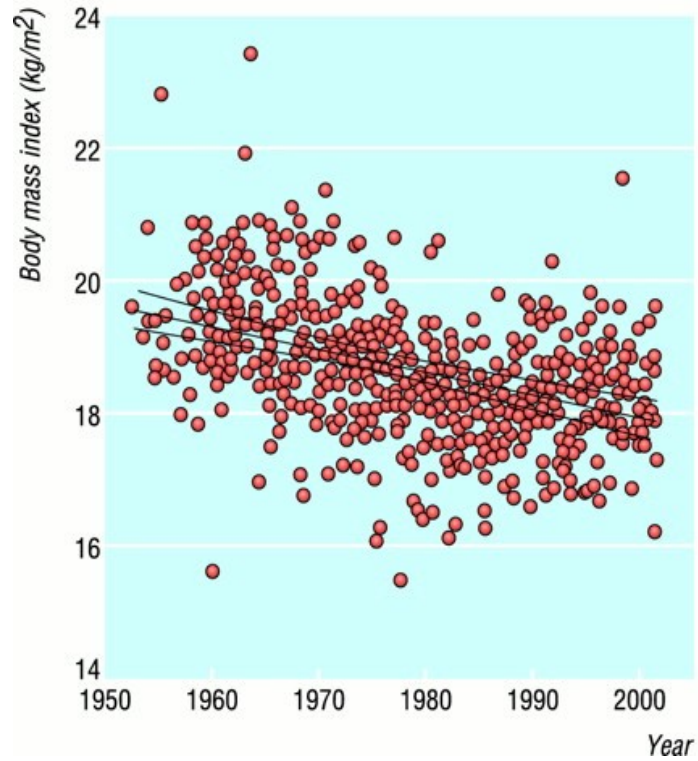
Immediately before mid-term



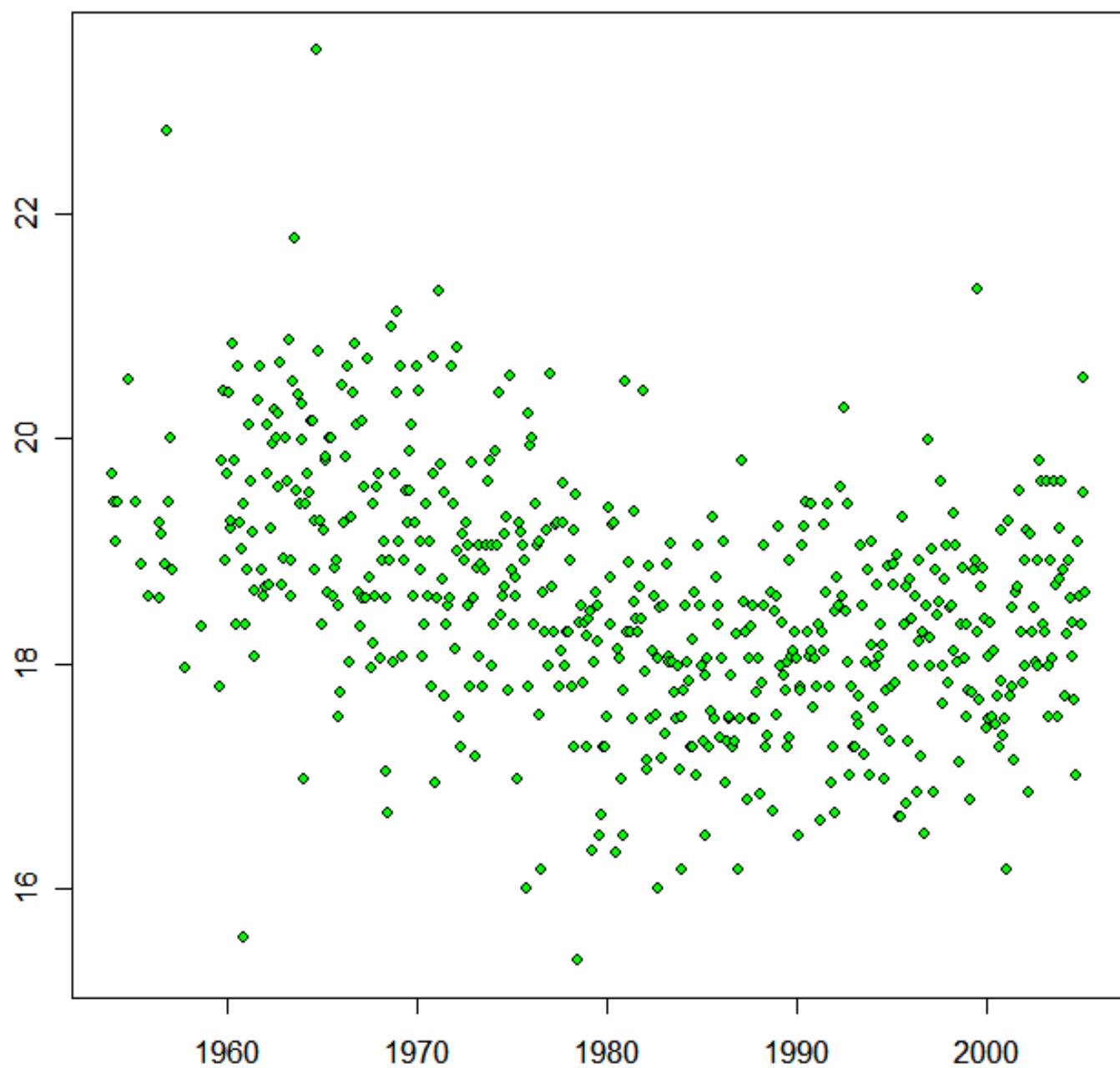
smoothing and straightening

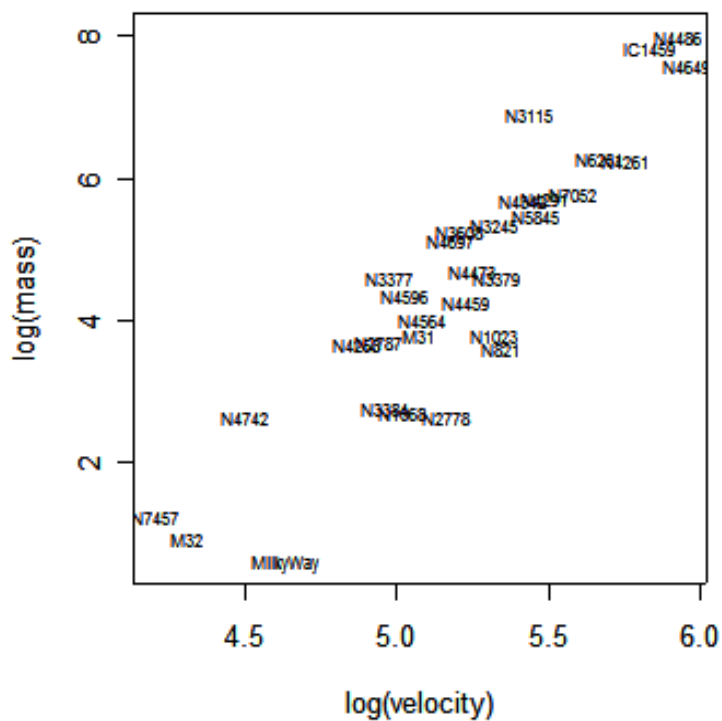
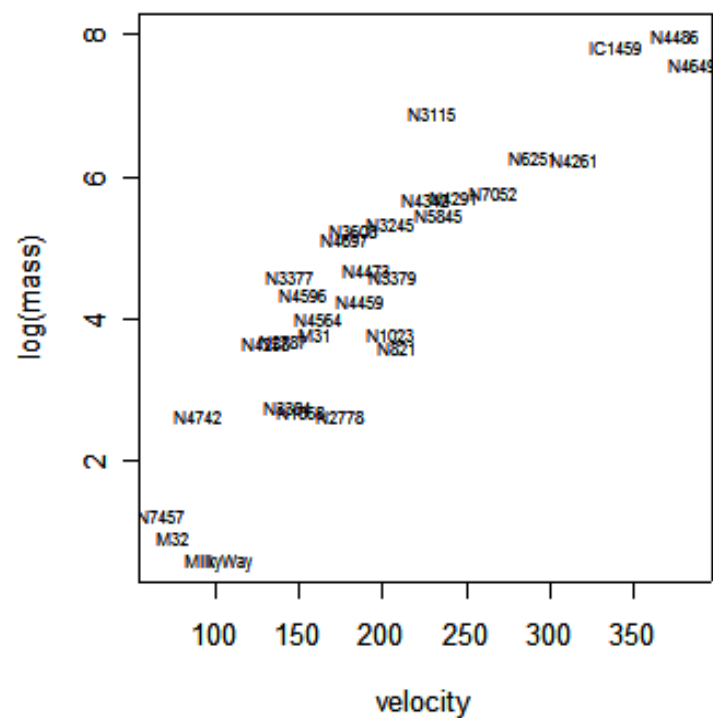
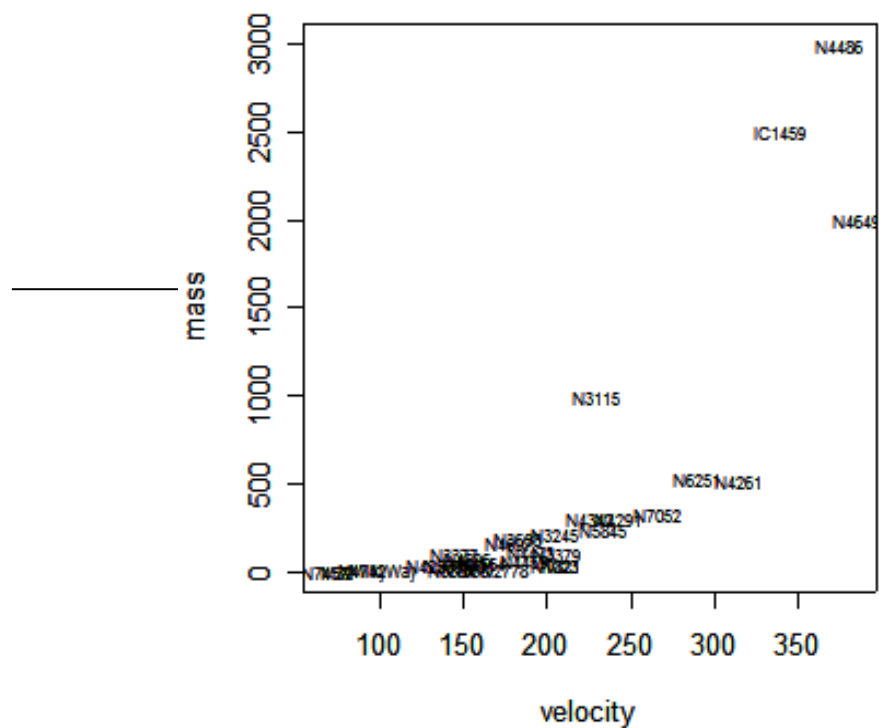
- smooth lines
 - piecewise linearity
 - splines and lo(w)ess
- a ladder of re-expression
- the re-expression rule

shapely centrefolds?



BMI for Playboy centerfolds, 1953-2005





a ladder of re-expressions...

3

2

1

$\frac{1}{2}$

#

$-\frac{1}{2}$

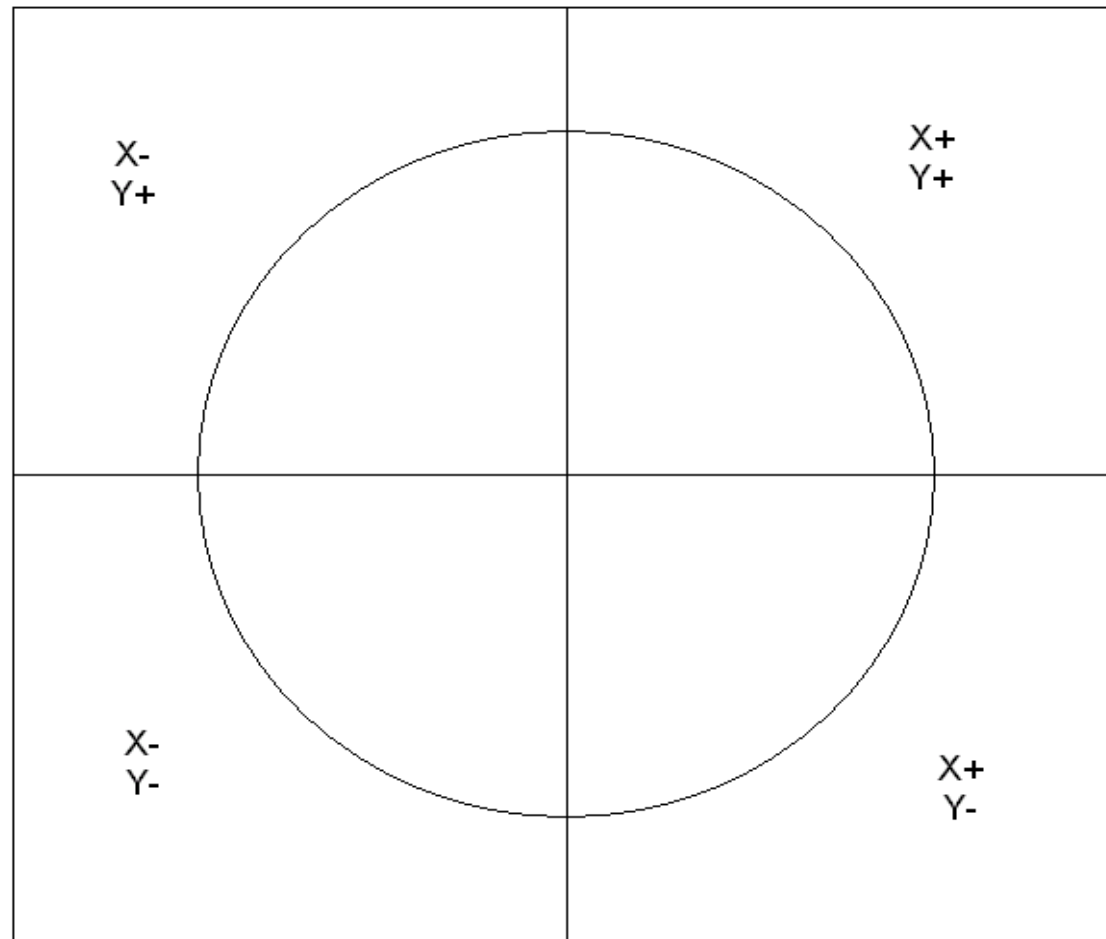
-1

-2

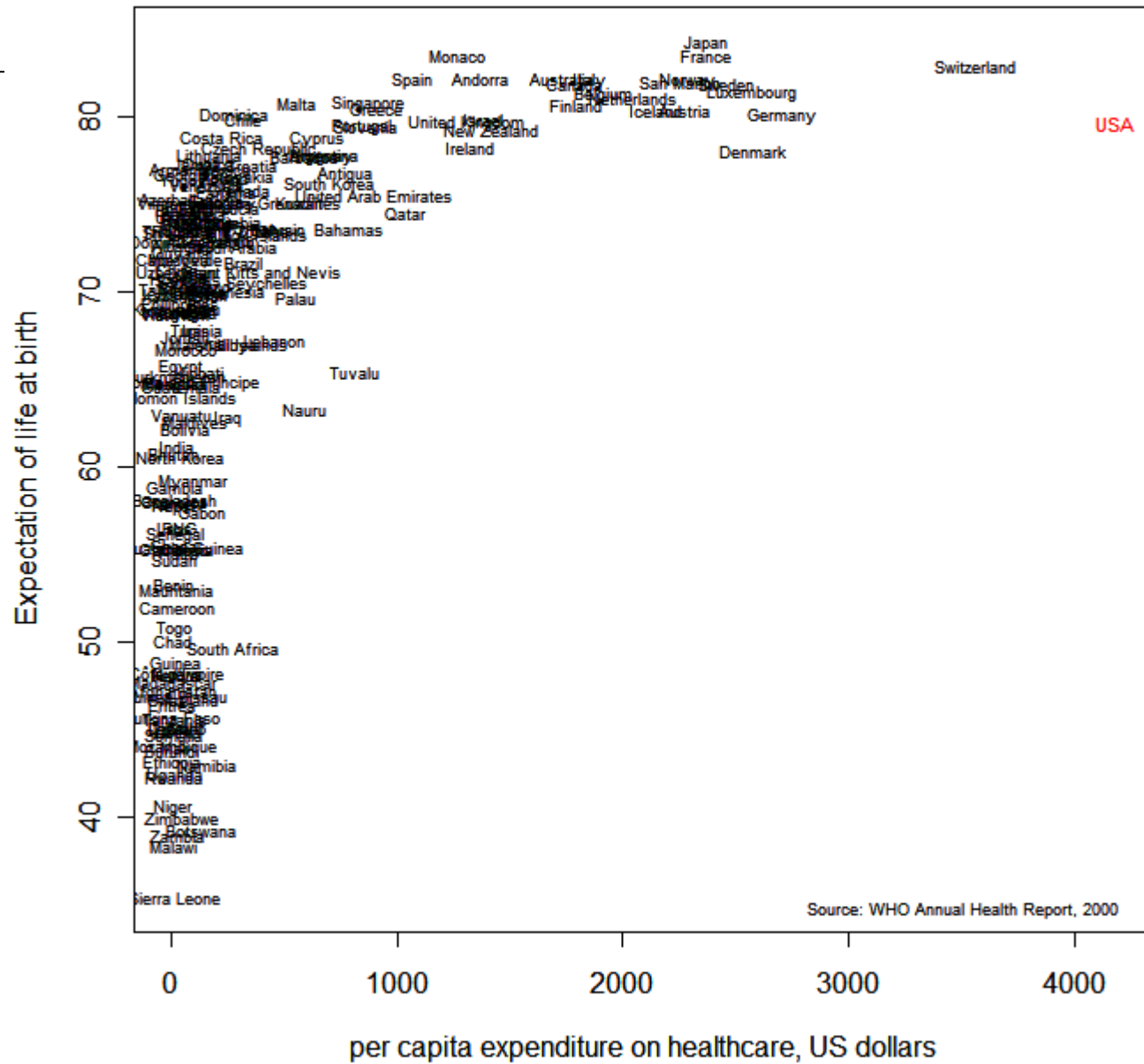
-3

...and a rule for using them

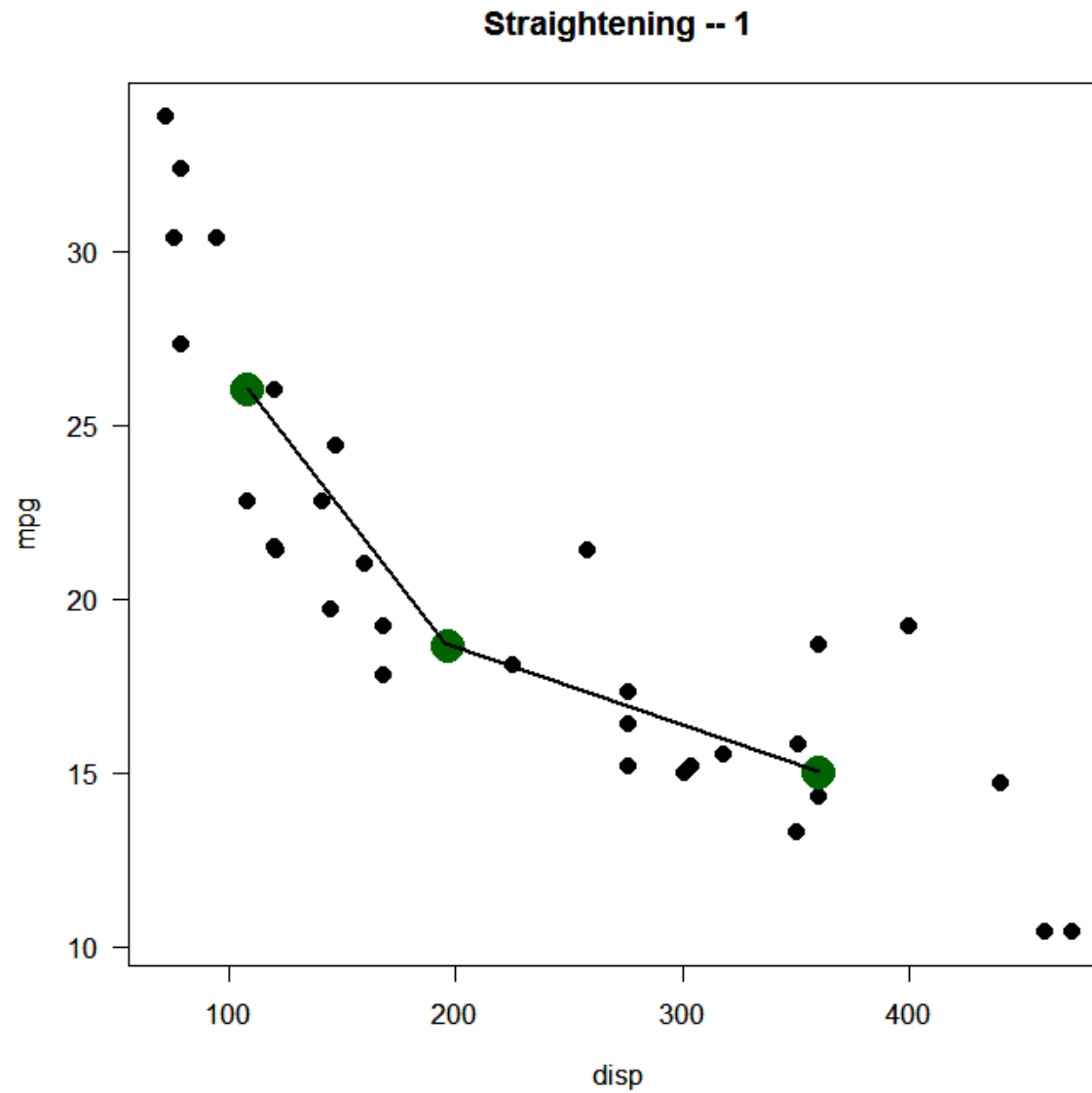
Straightening by re-expression



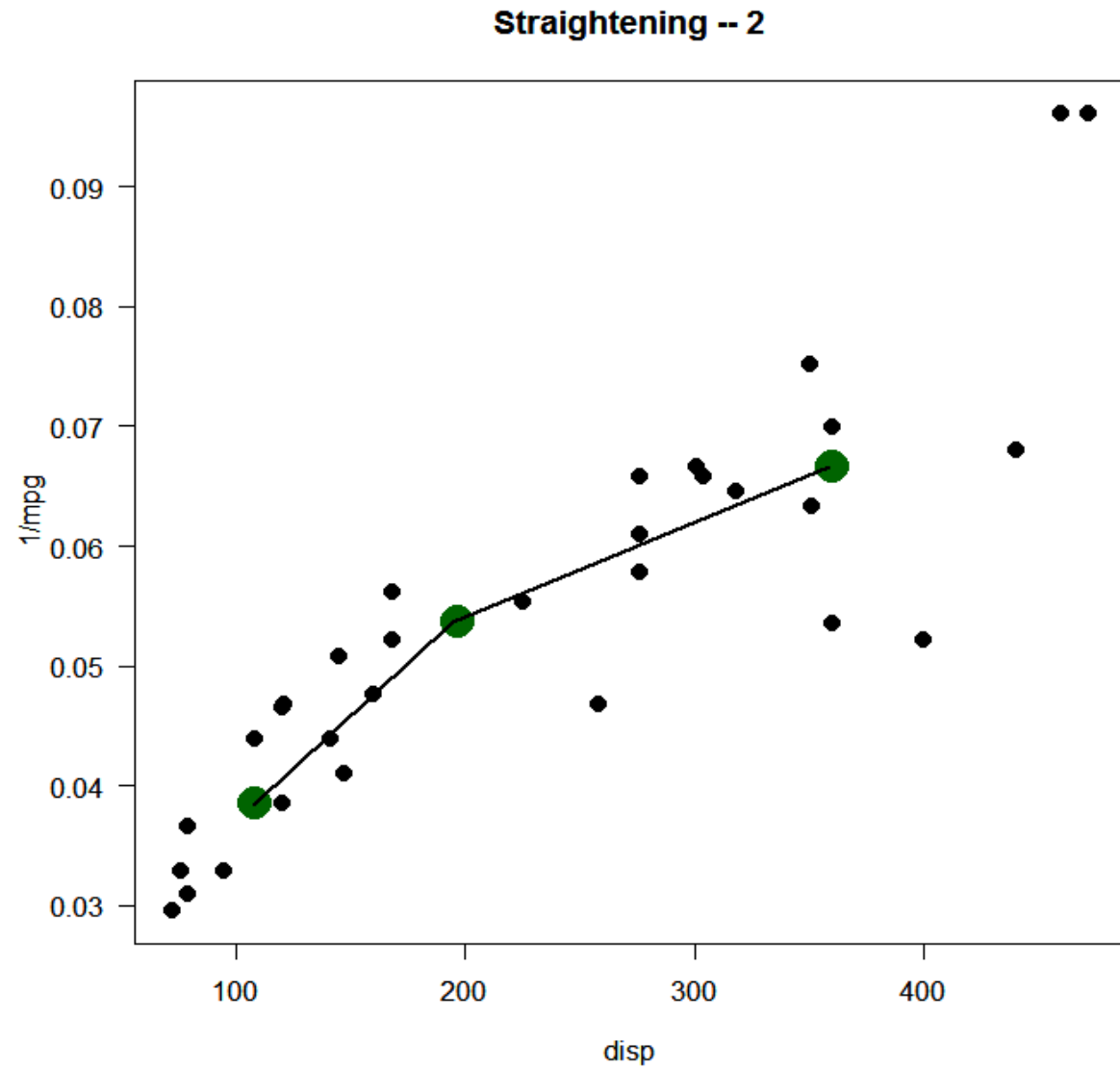
Life expectancy and national expenditures



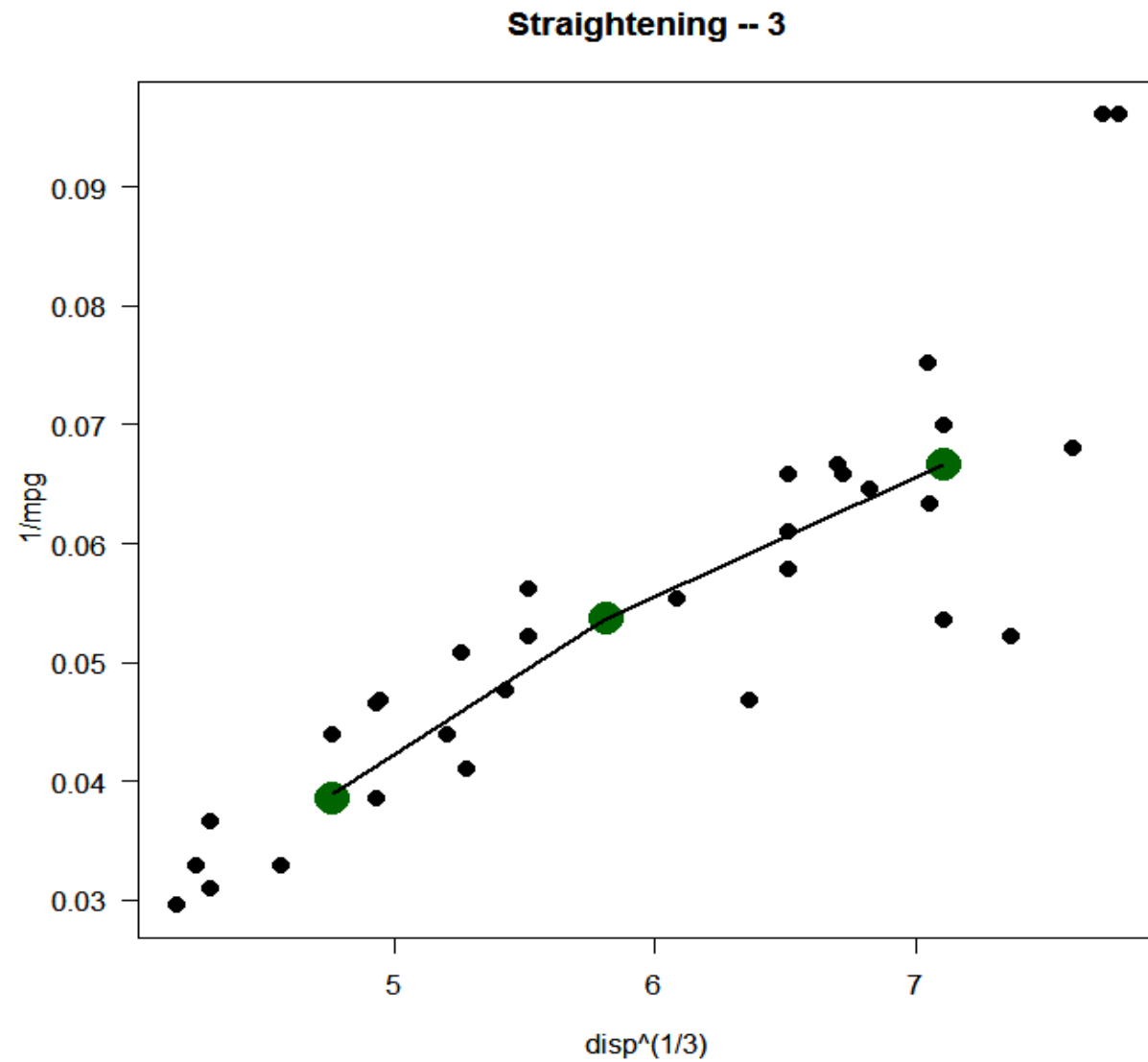
straightening



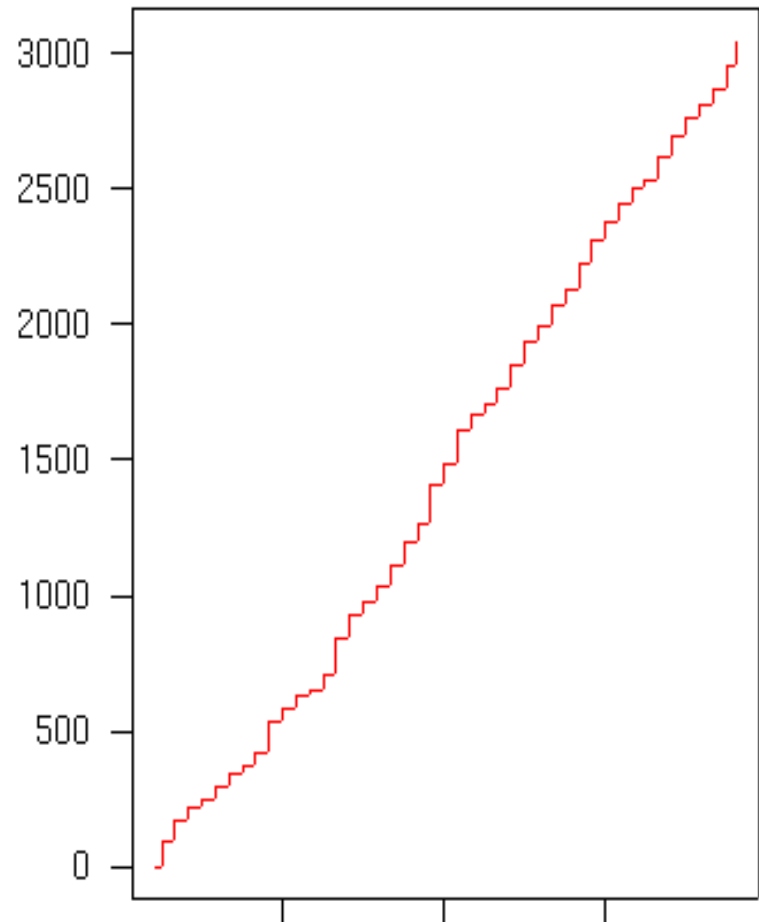
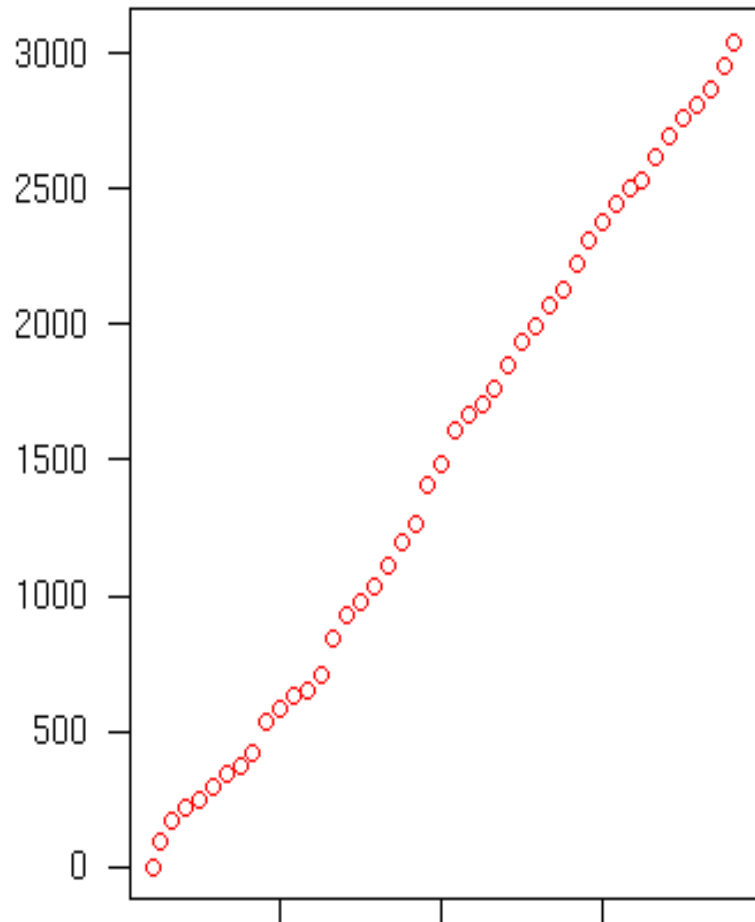
straightening 2



straightening 3



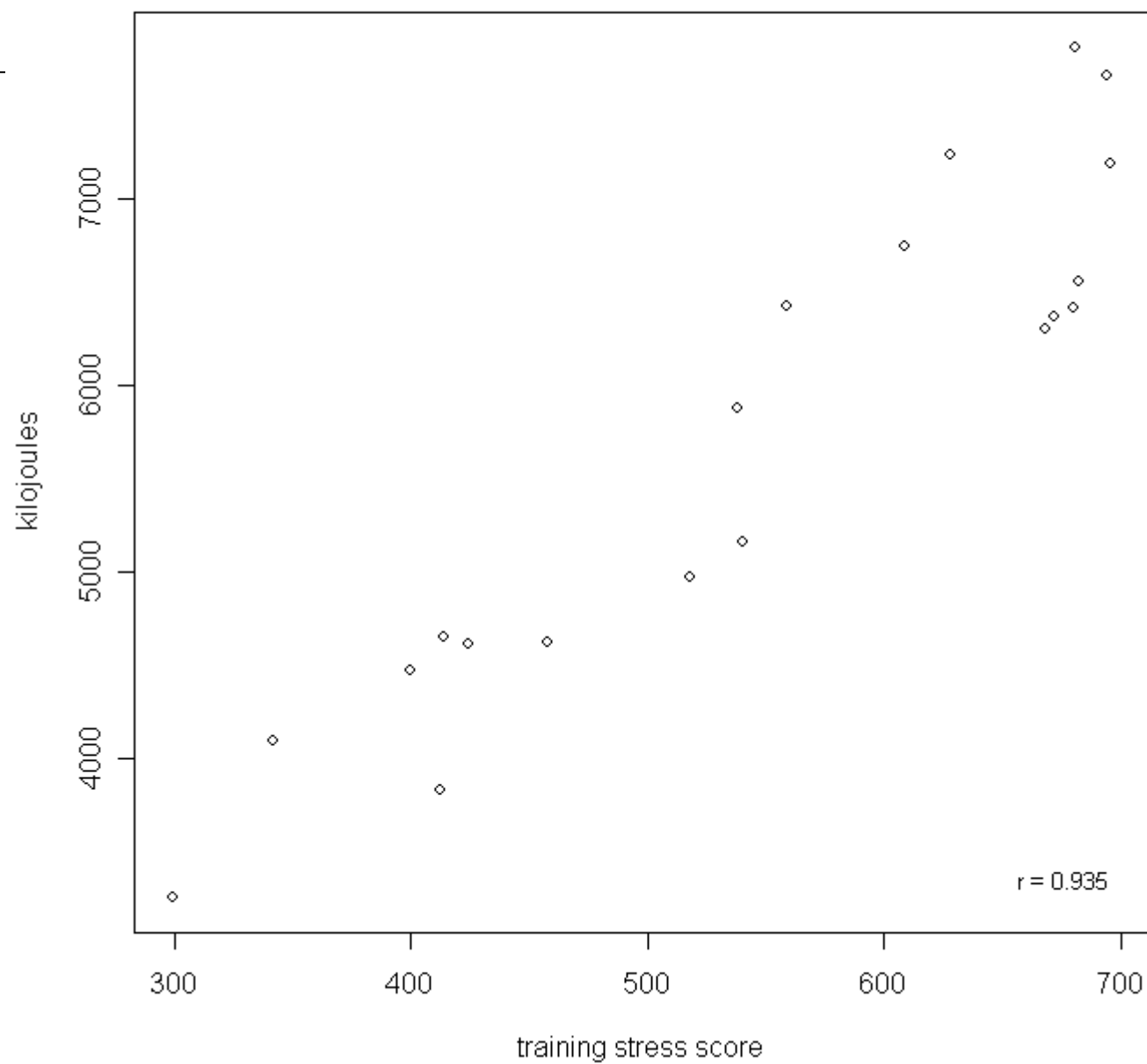
when is smooth too smooth?



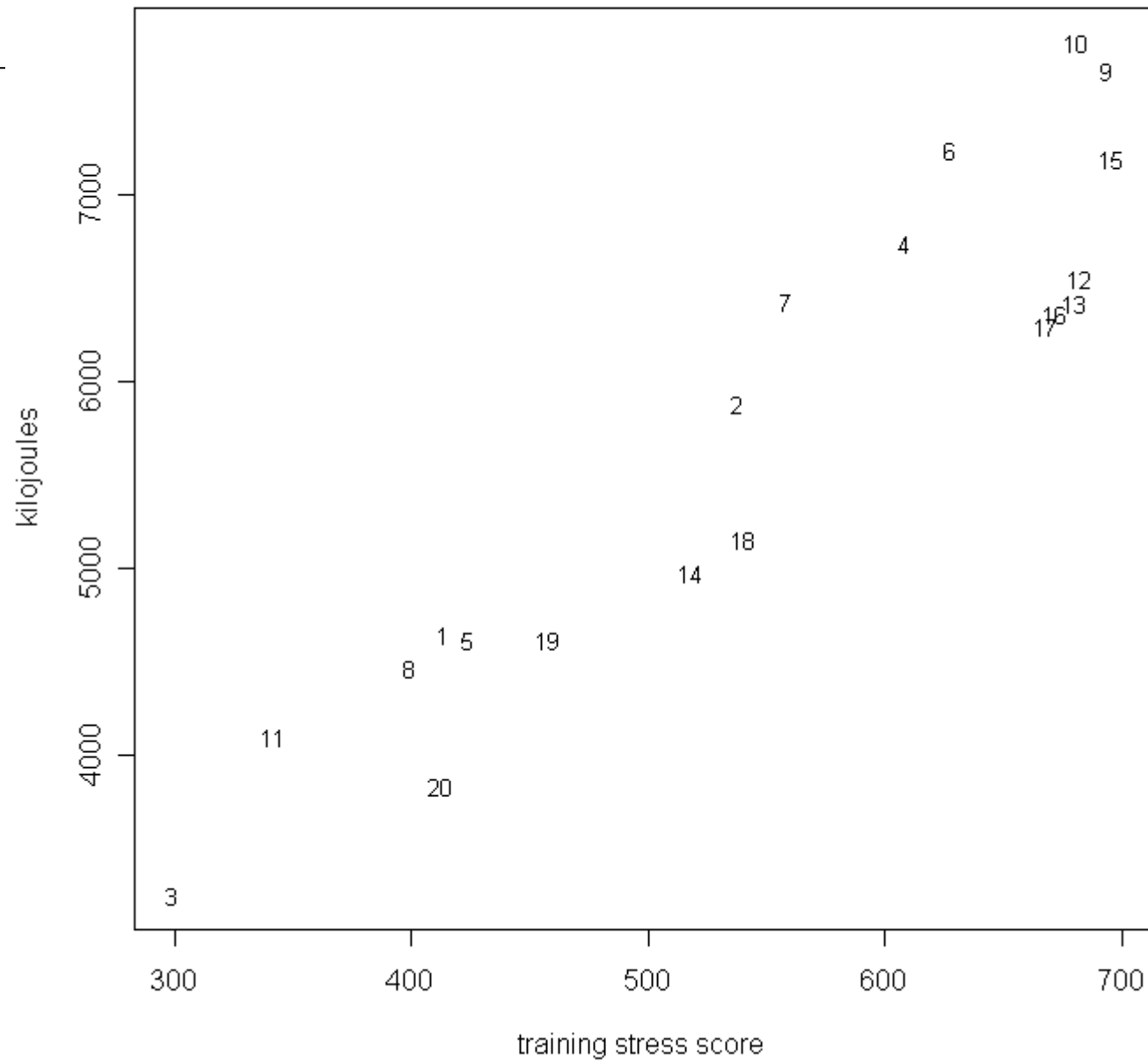
when is straight too straight?

- can straight be too straight?

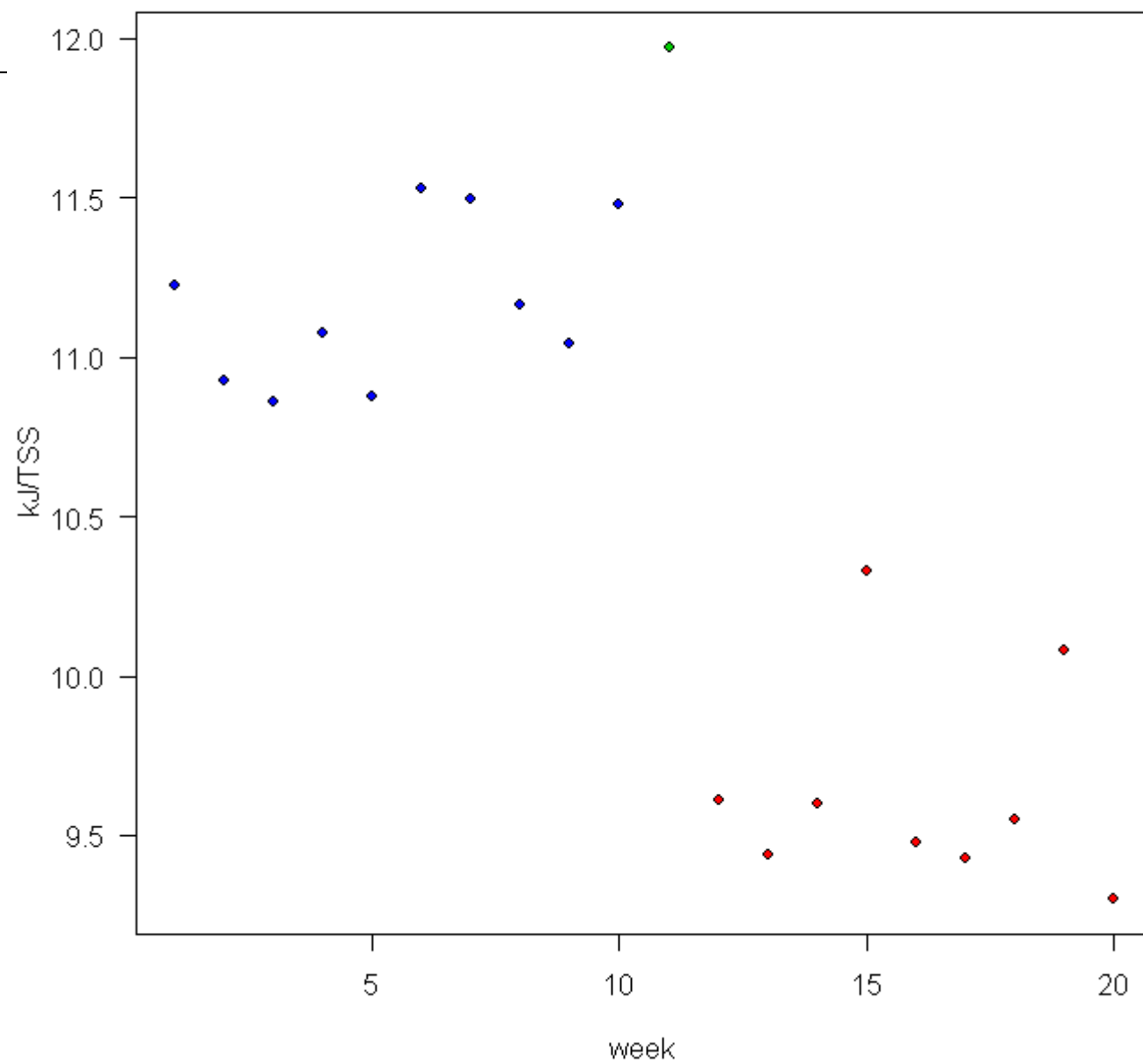
Andy's weekly training load



Andy's weekly training load



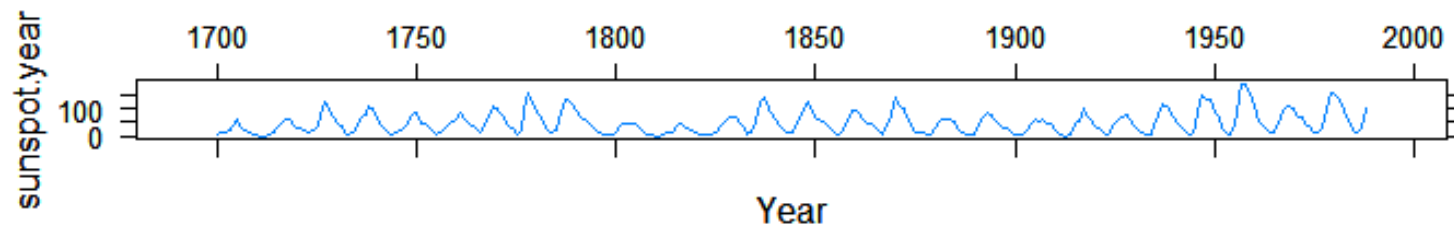
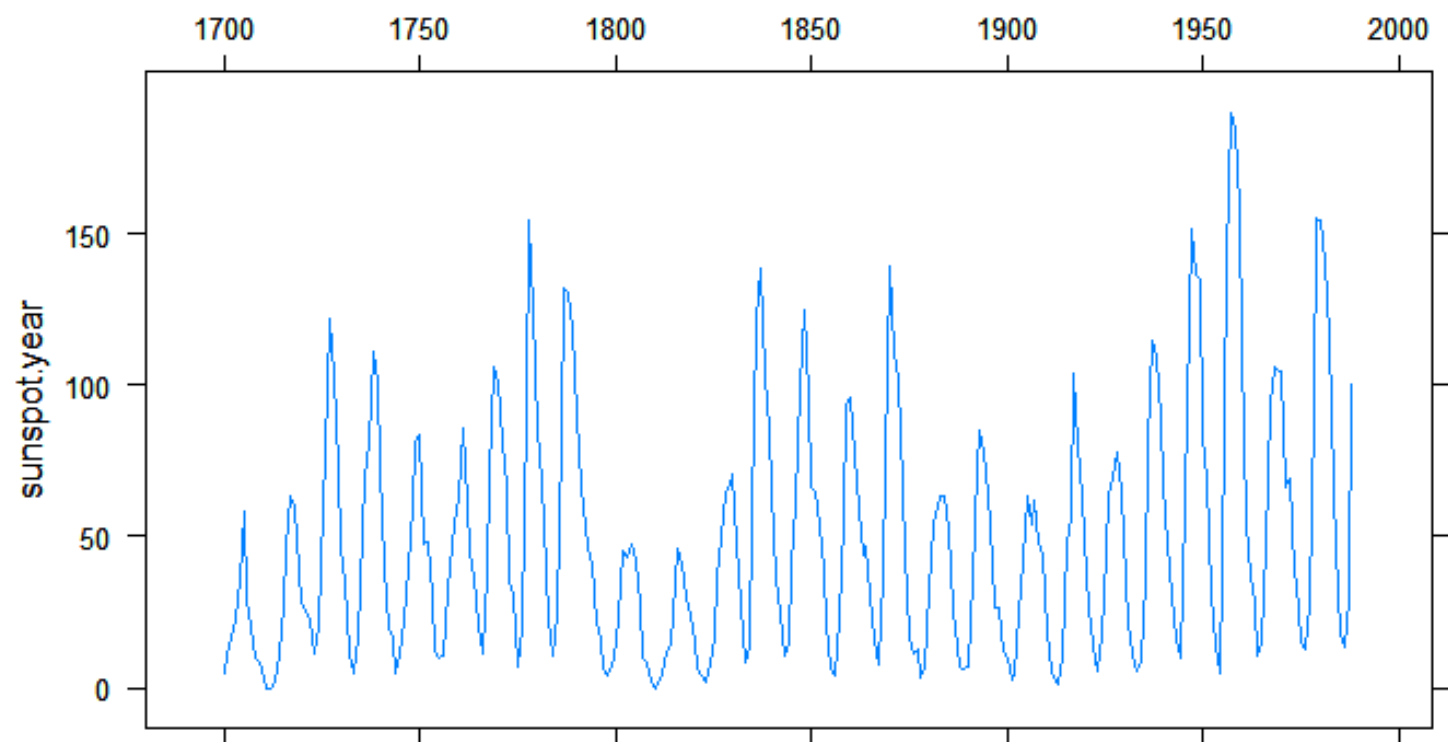
Andy's weekly training load



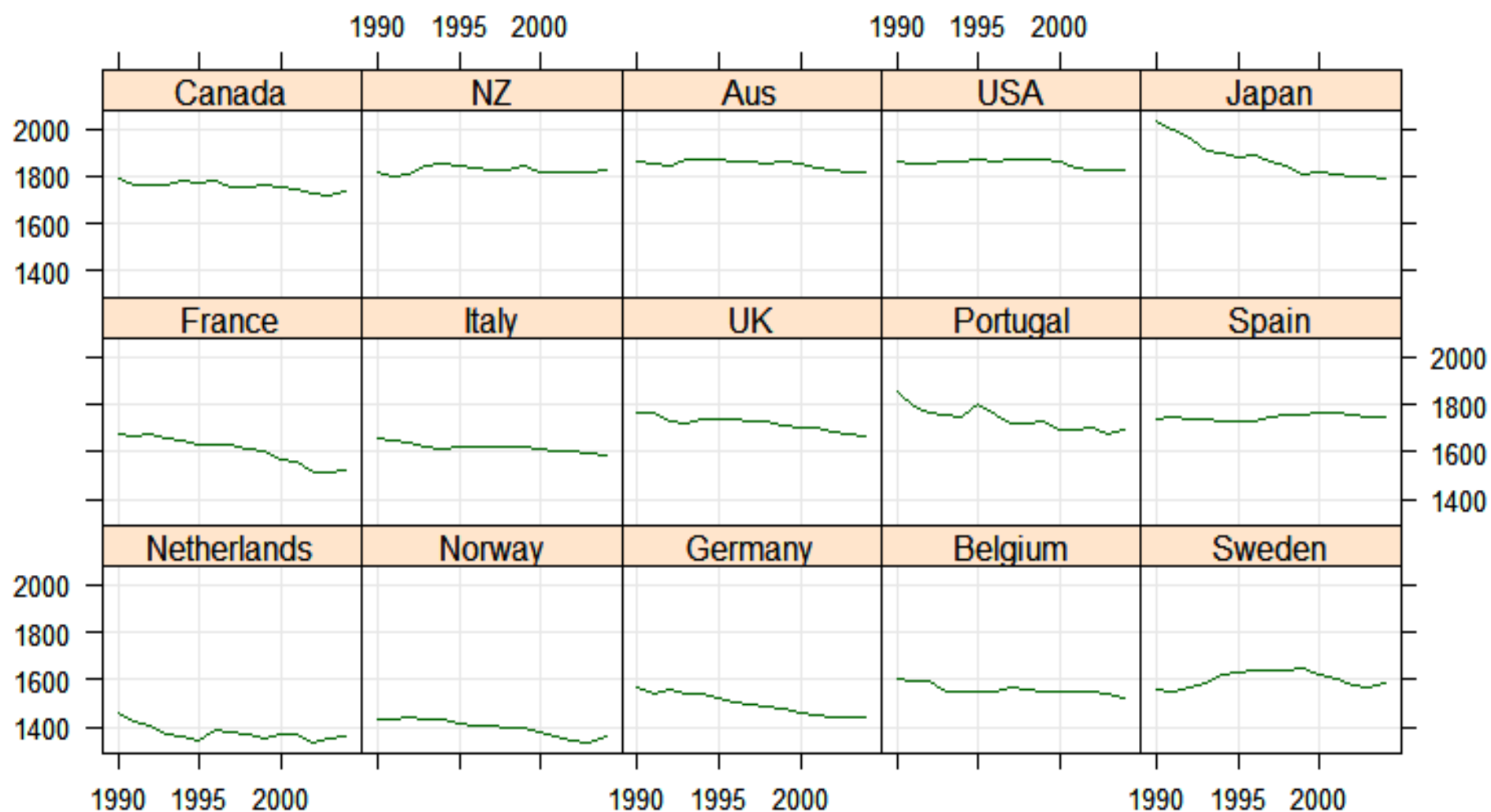
aspect ratios and banking

- human eye isn't great at decoding angles
banking helps the eye to decipher angles

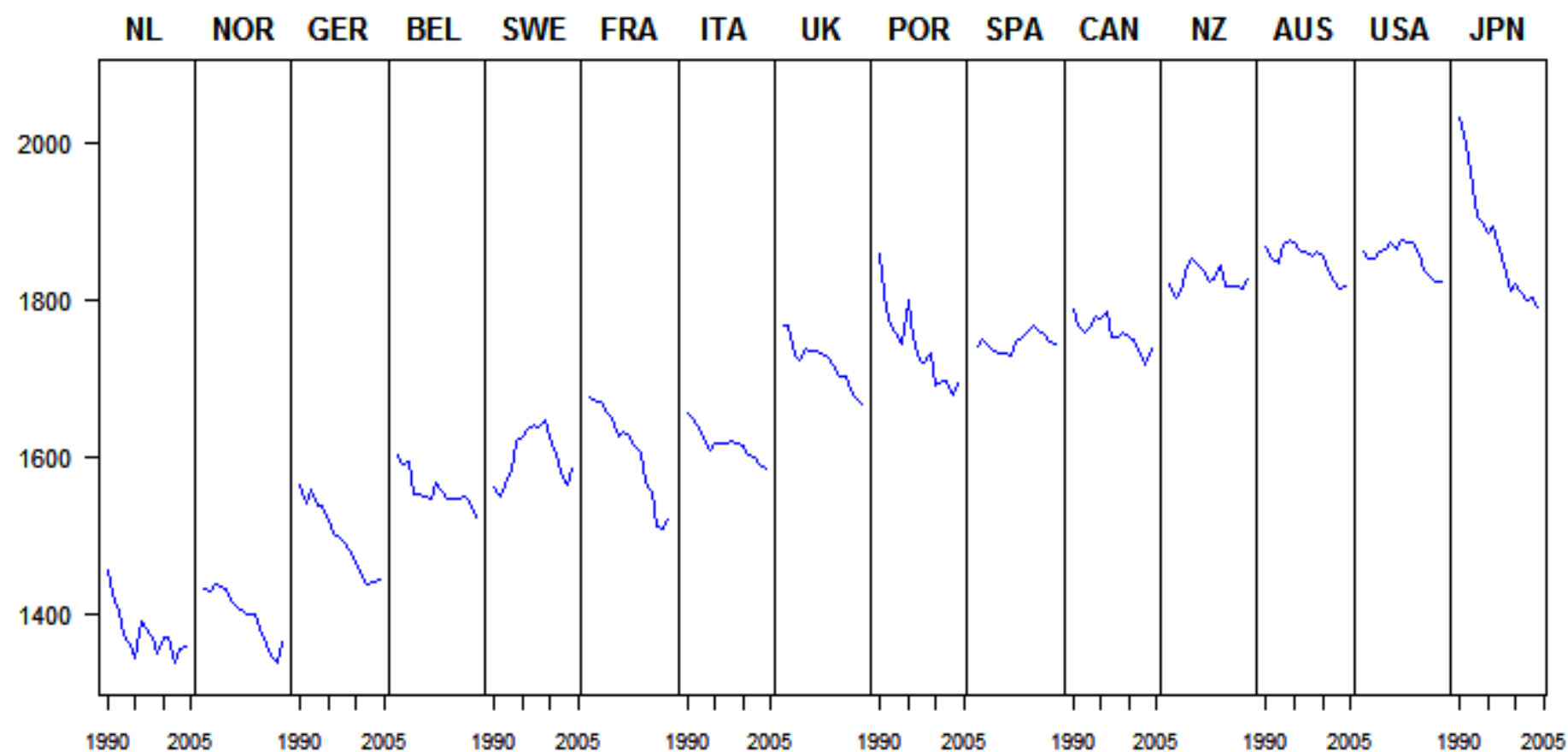
Yearly Sunspots



Average annual hours worked per worker, 1990-2004



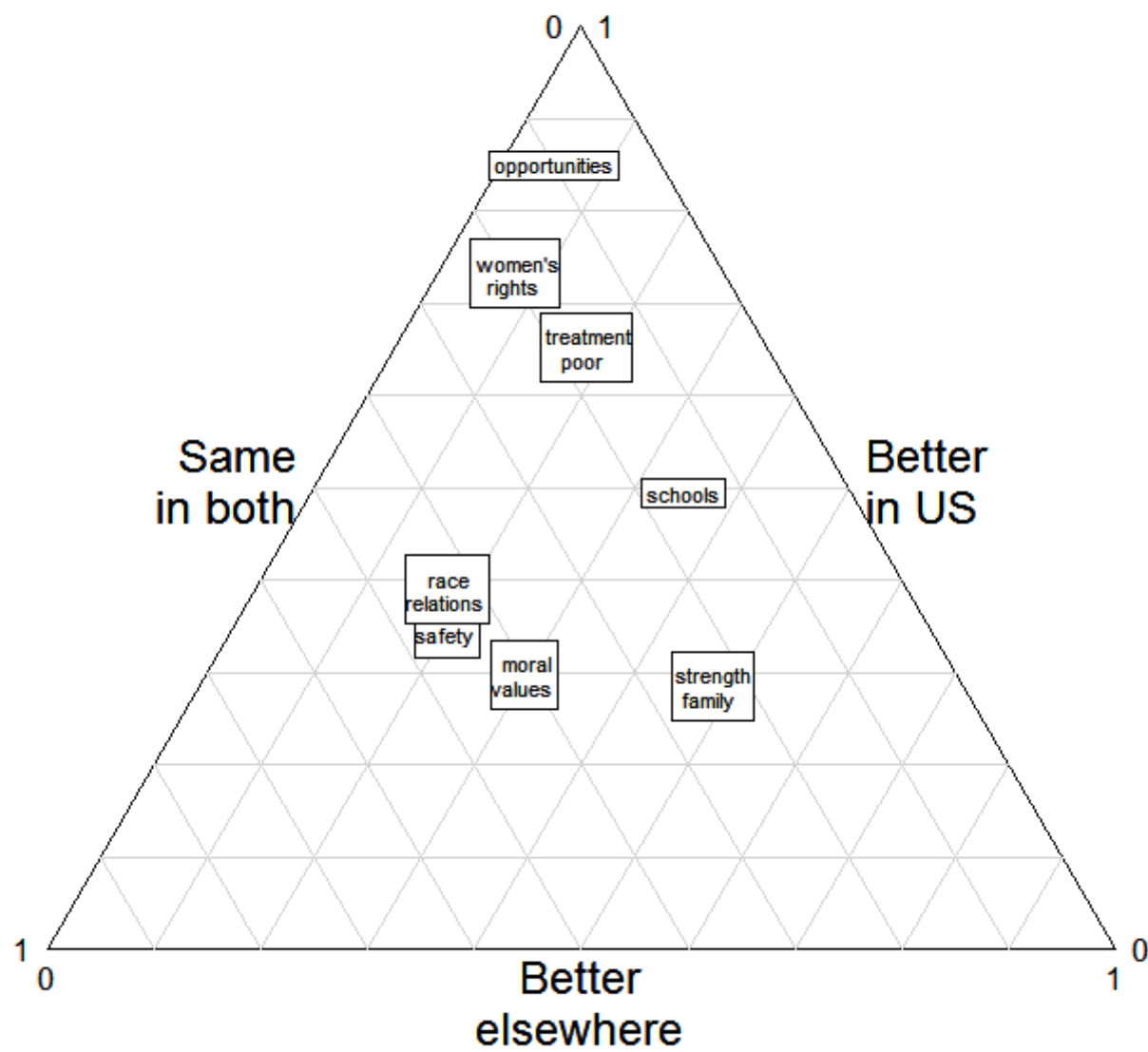
Average annual hours worked per worker, 1990-2004



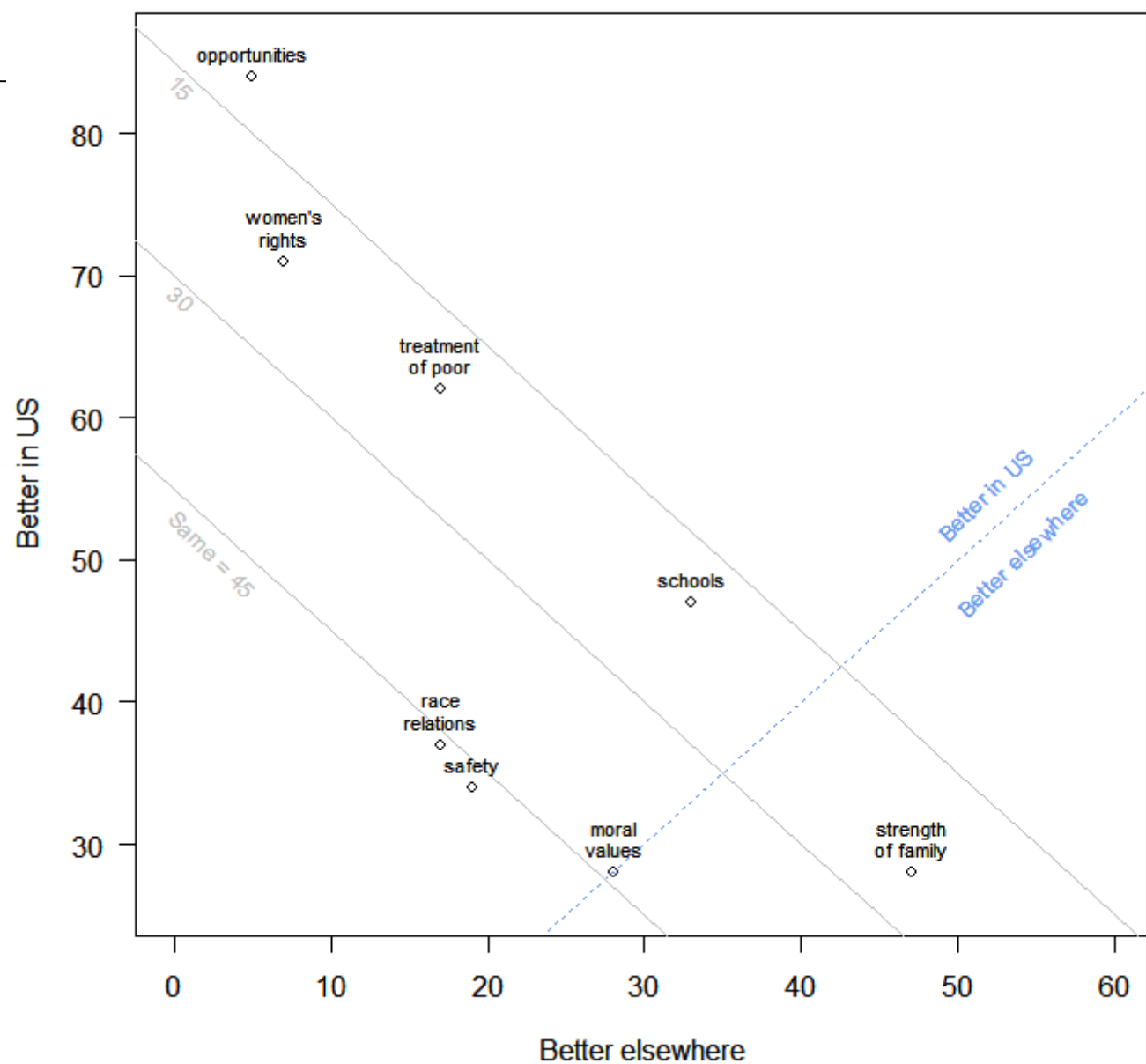
contours and bases

- triangle plots (= ternary plots)
 - soil texture plots
 - two degrees of freedom
- function contours can add context

Immigrant's assessments of US vs. country of origin

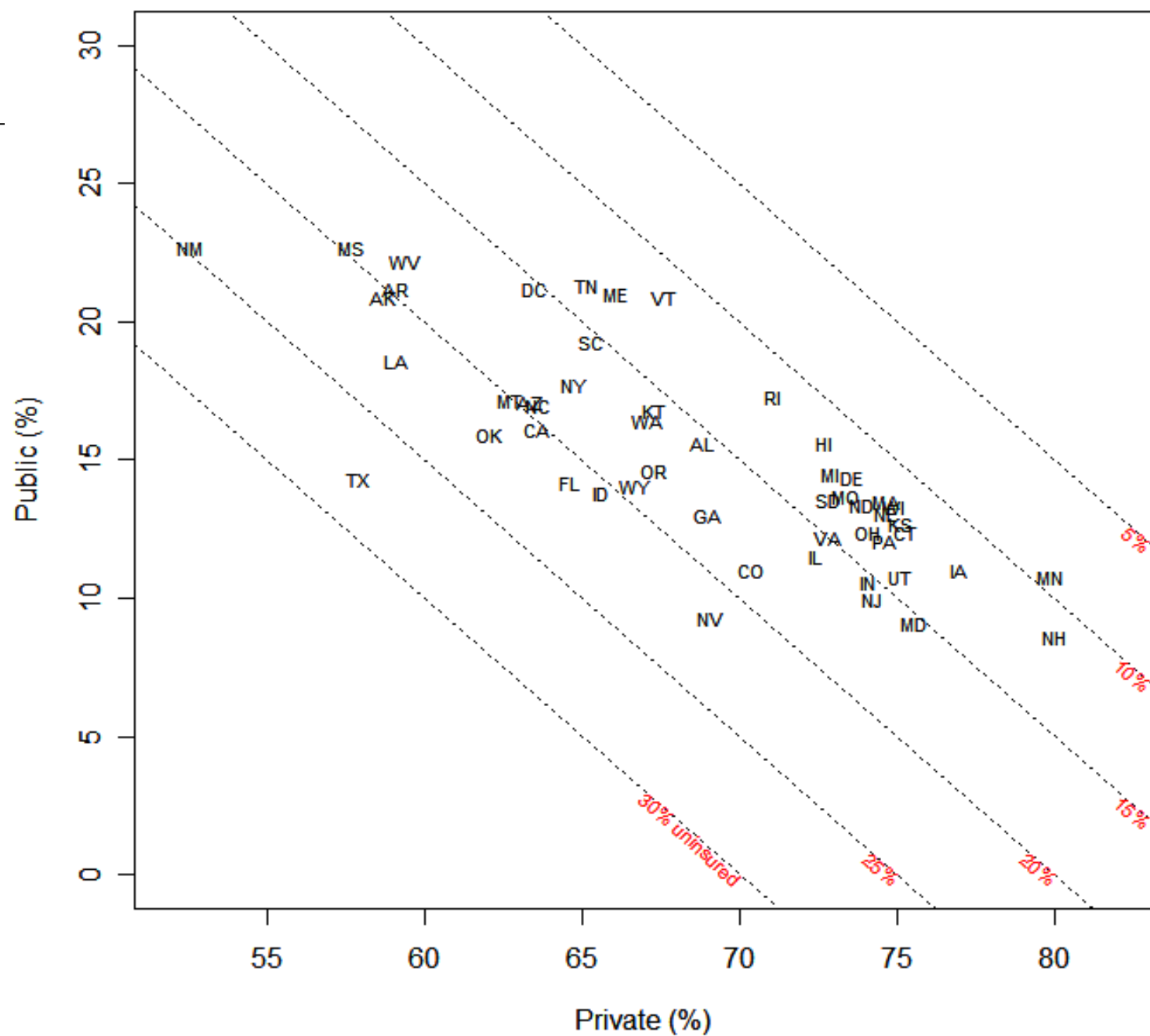


Immigrant's assessments of US vs. country of origin



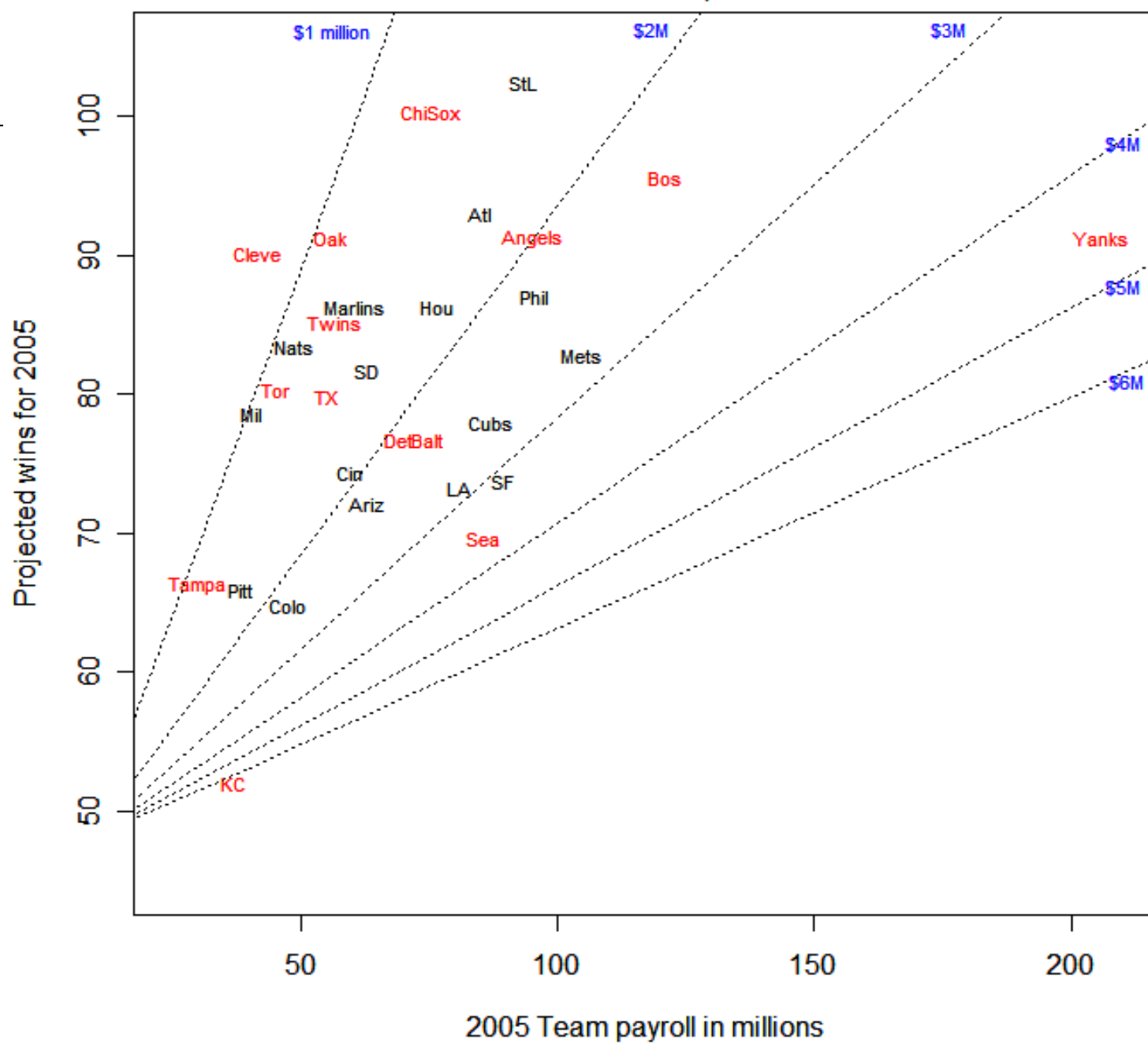
Type of health insurance for non-elderly

based on March 2004 CPS



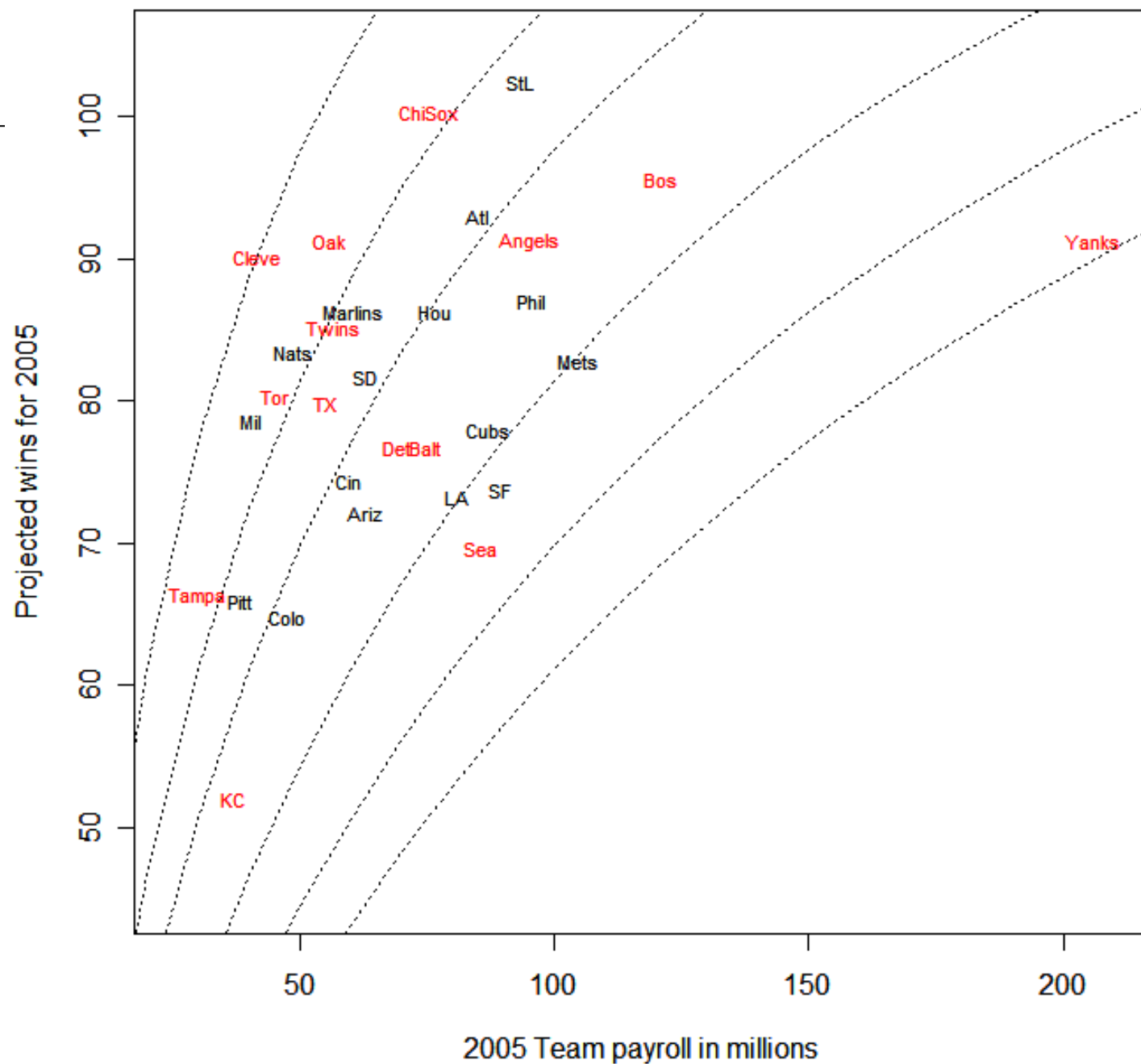
Pappas' efficiency metric

records as of 3 Sep 2005

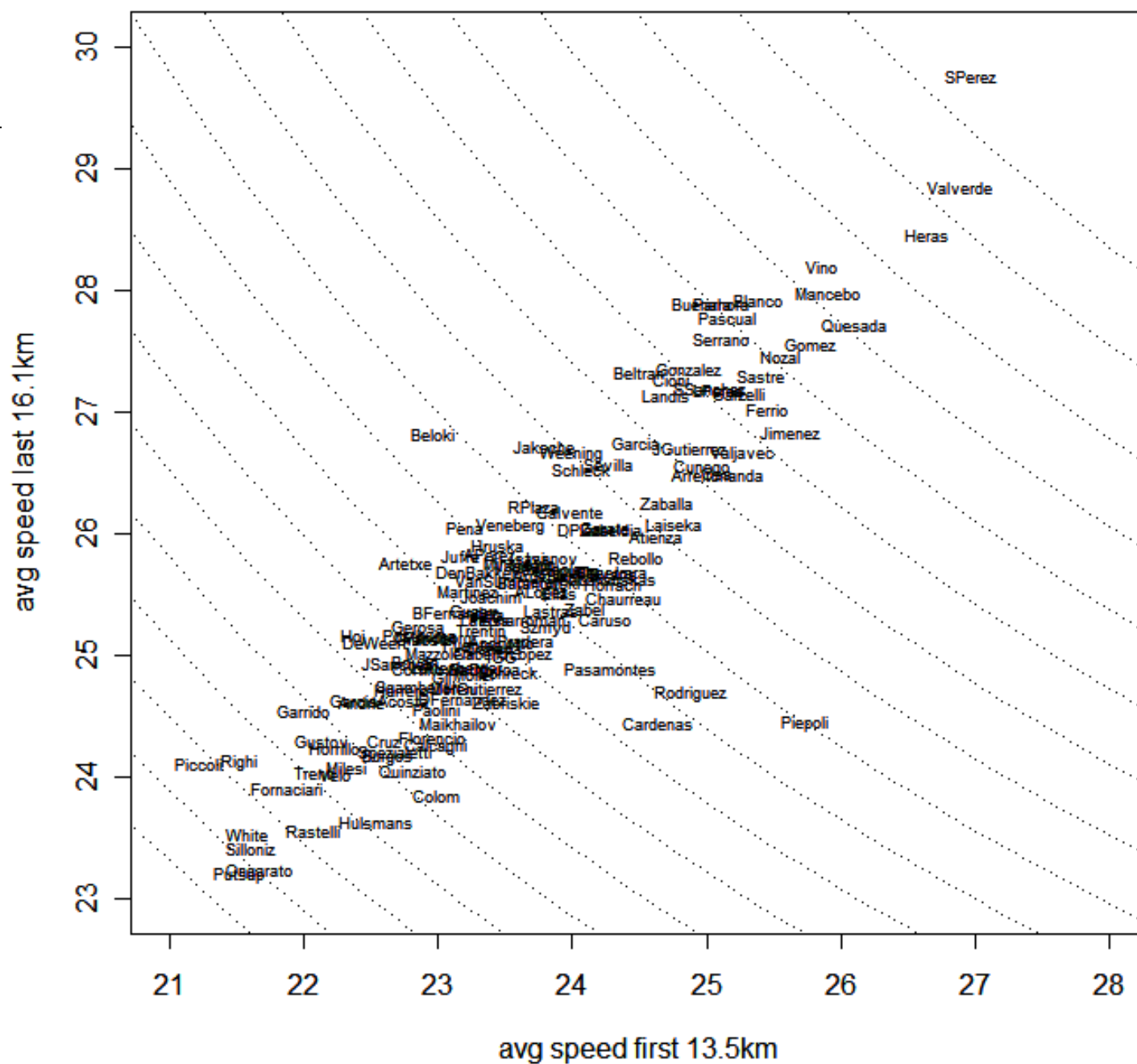


Payroll/median payroll ratio to win/loss ratio

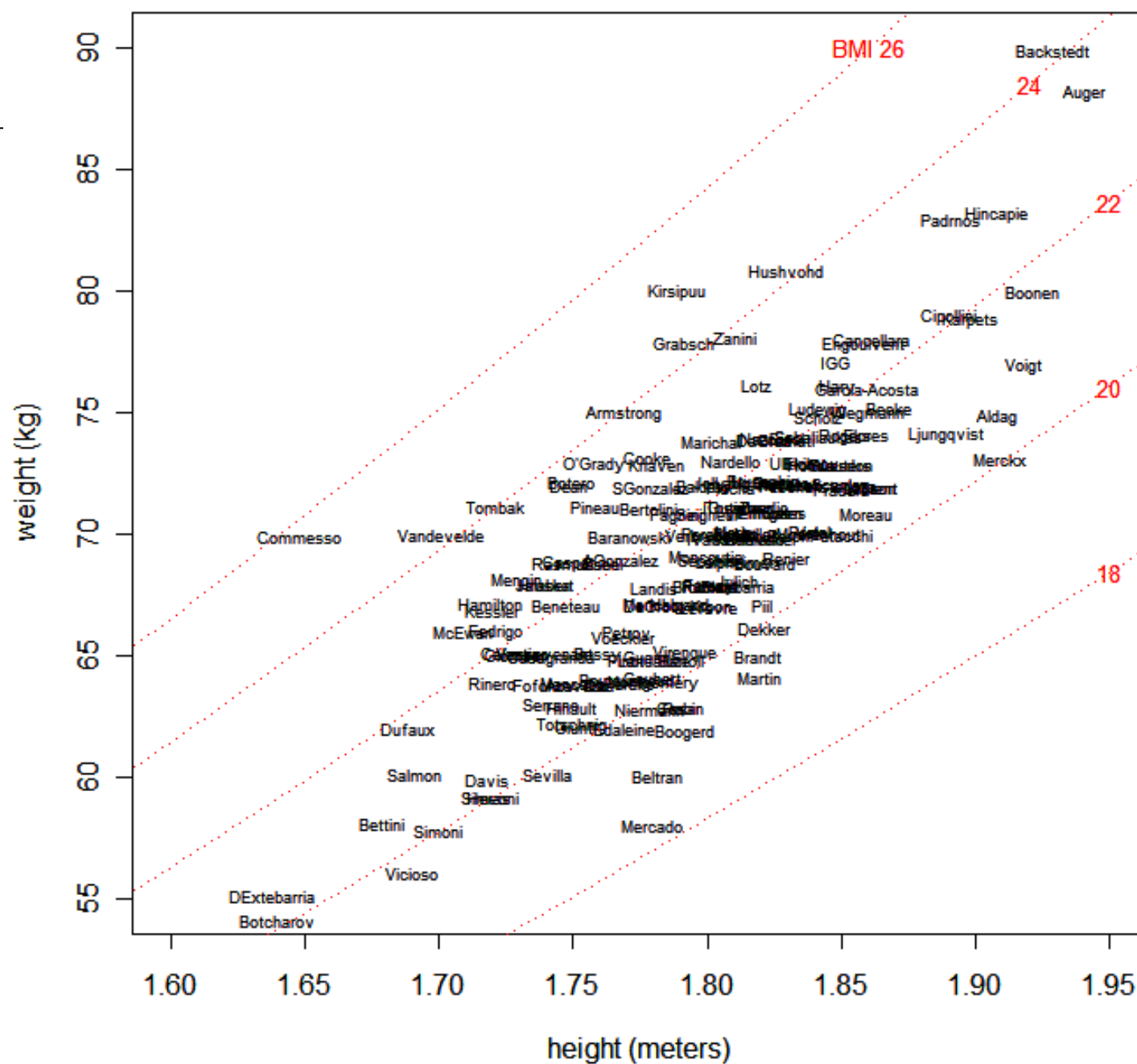
records as of 3 Sep 2005



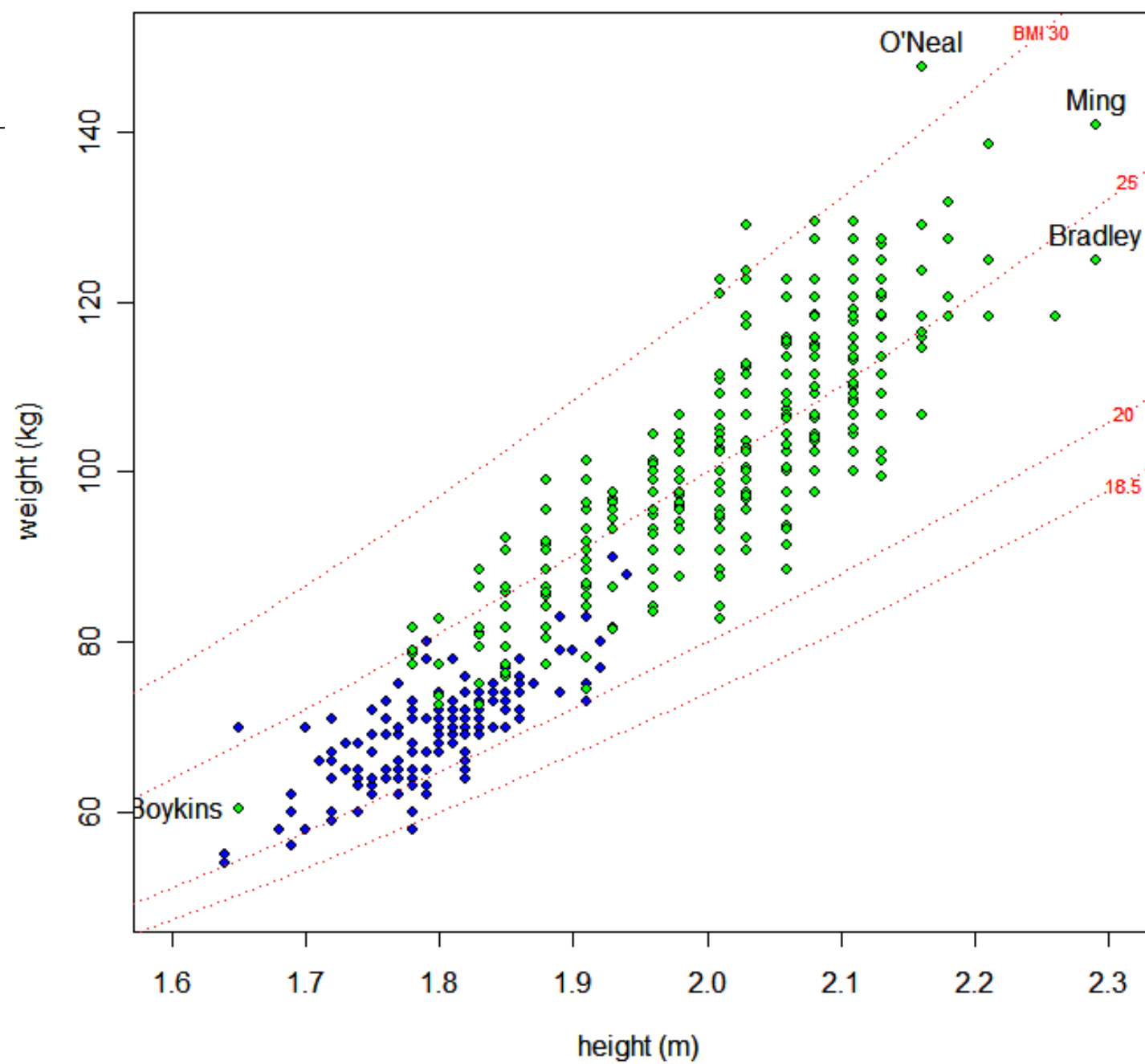
With one minute time contours



2004 TdF riders



NBA players and TdF riders, 2004

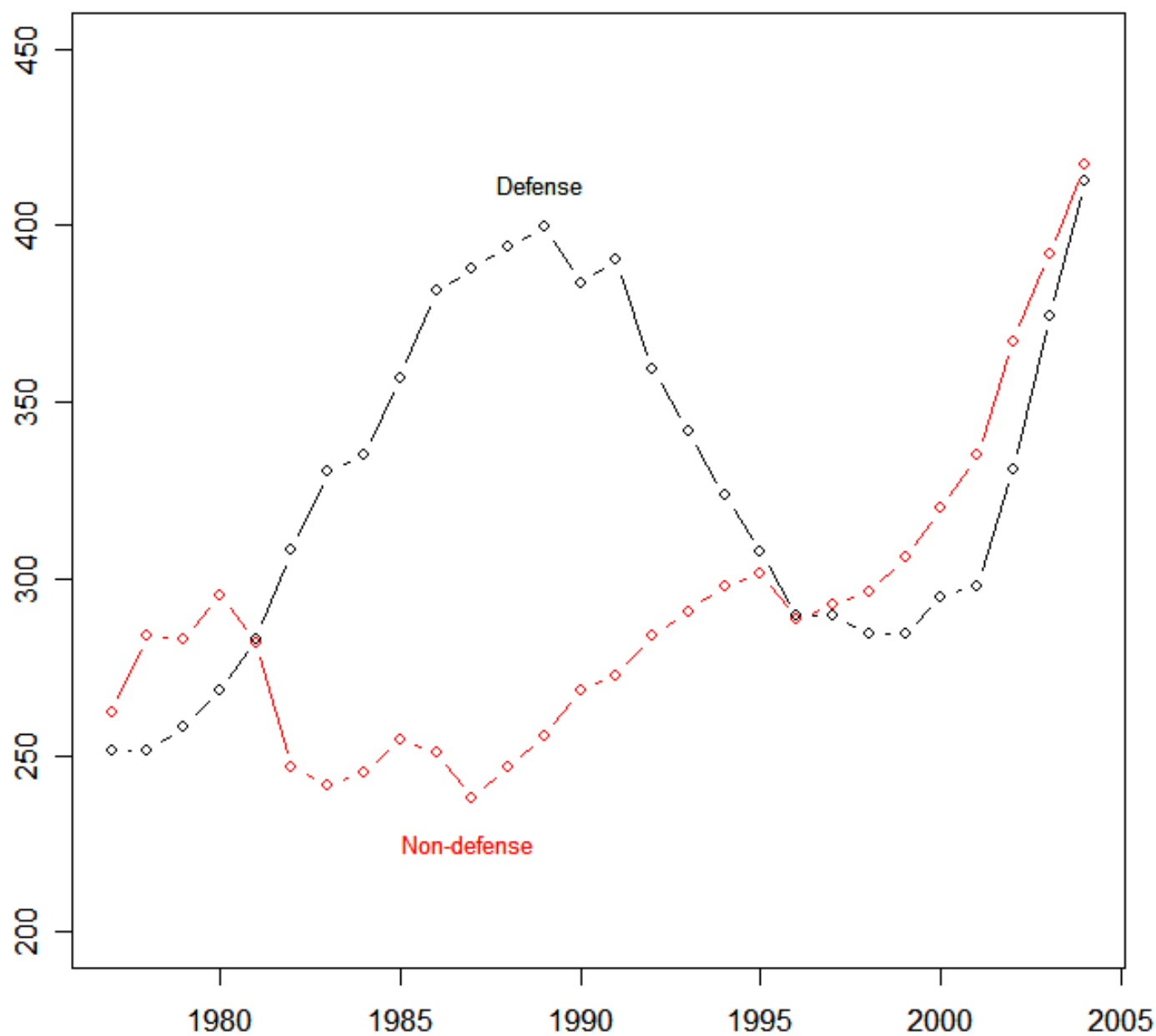


decomposing series into phase plots

- another version of “show the difference”

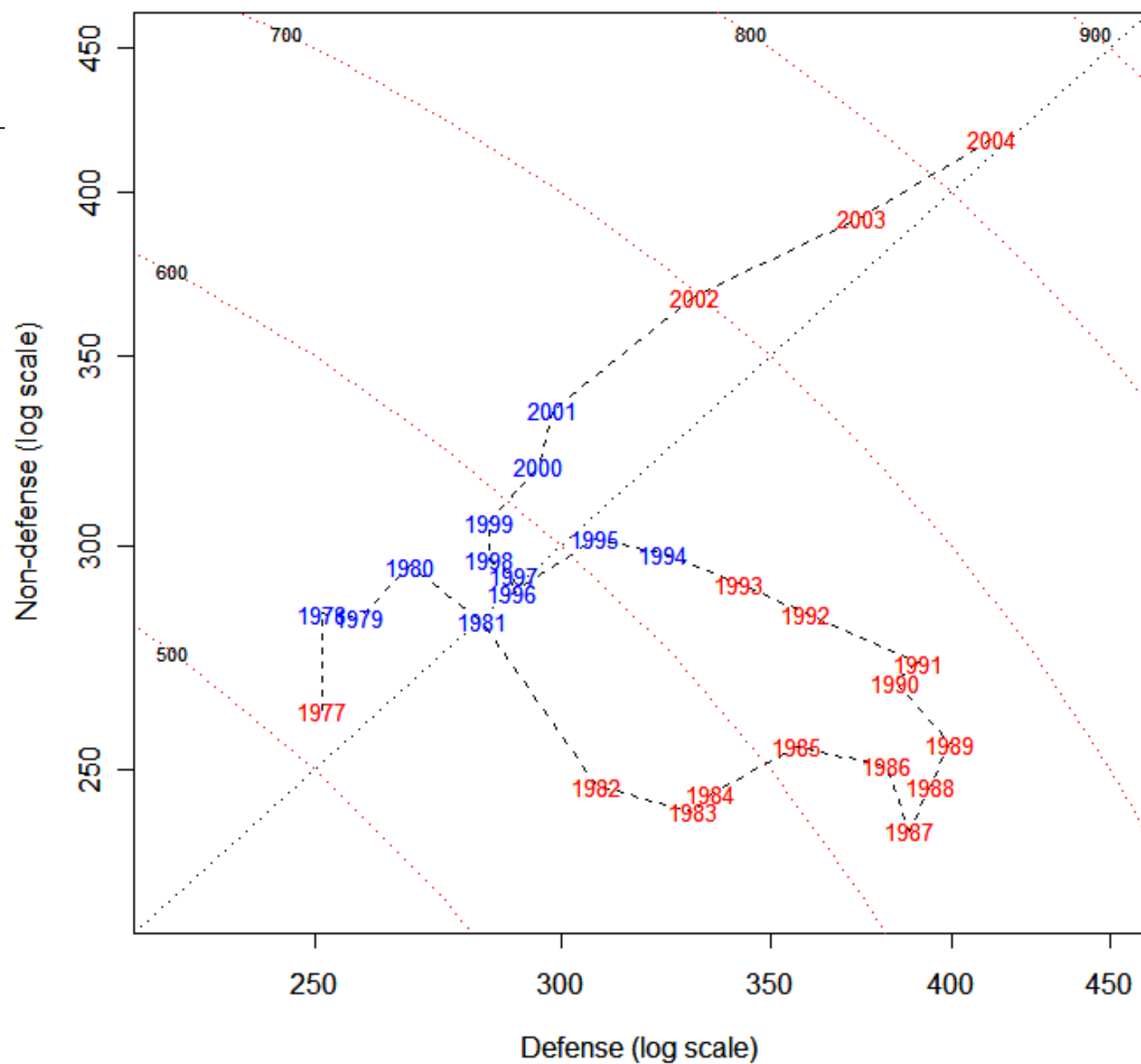
Federal buget: components of discretionary spending

in billions of 2000 constant dollars

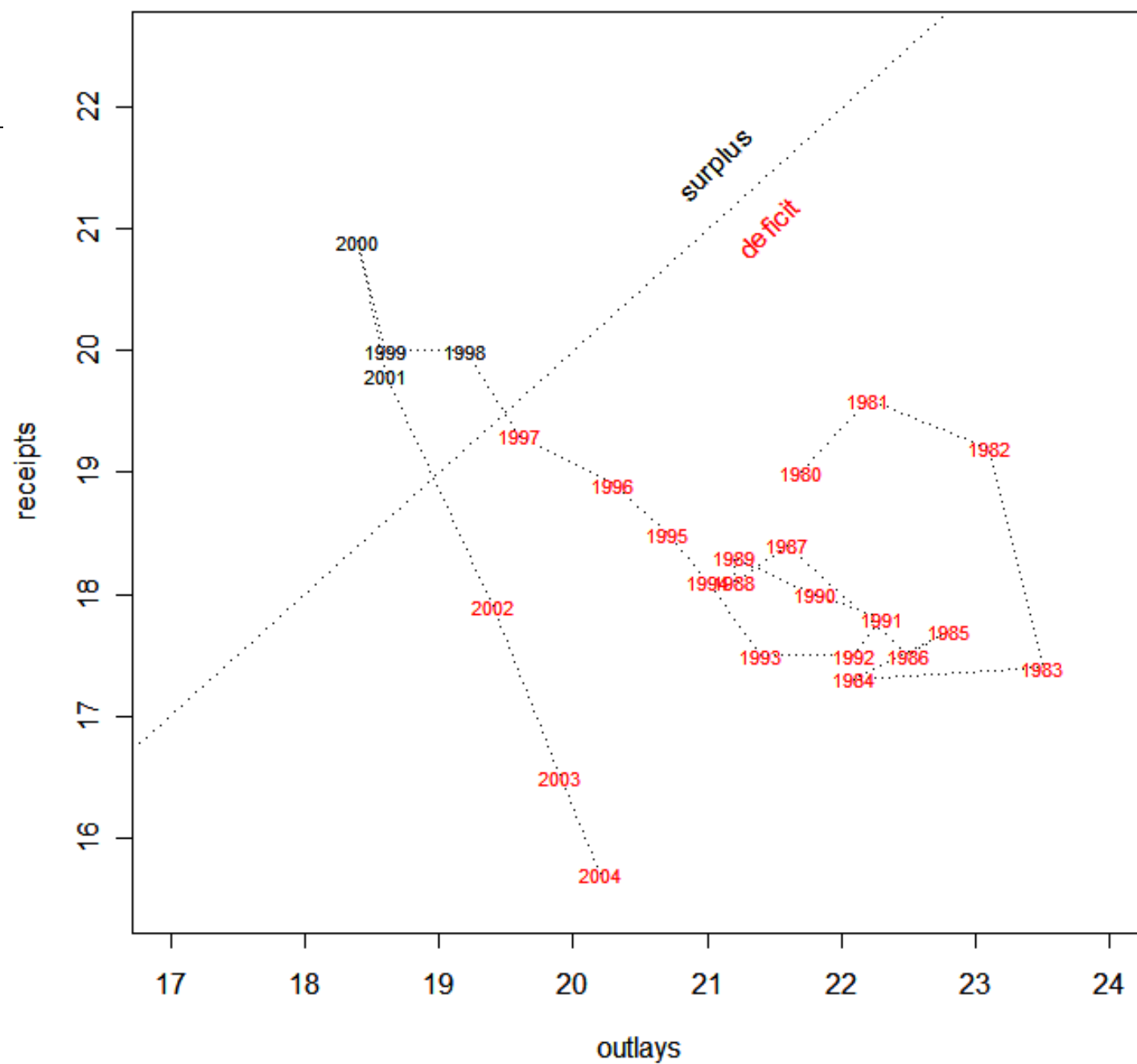


Federal budget: components of discretionary spending

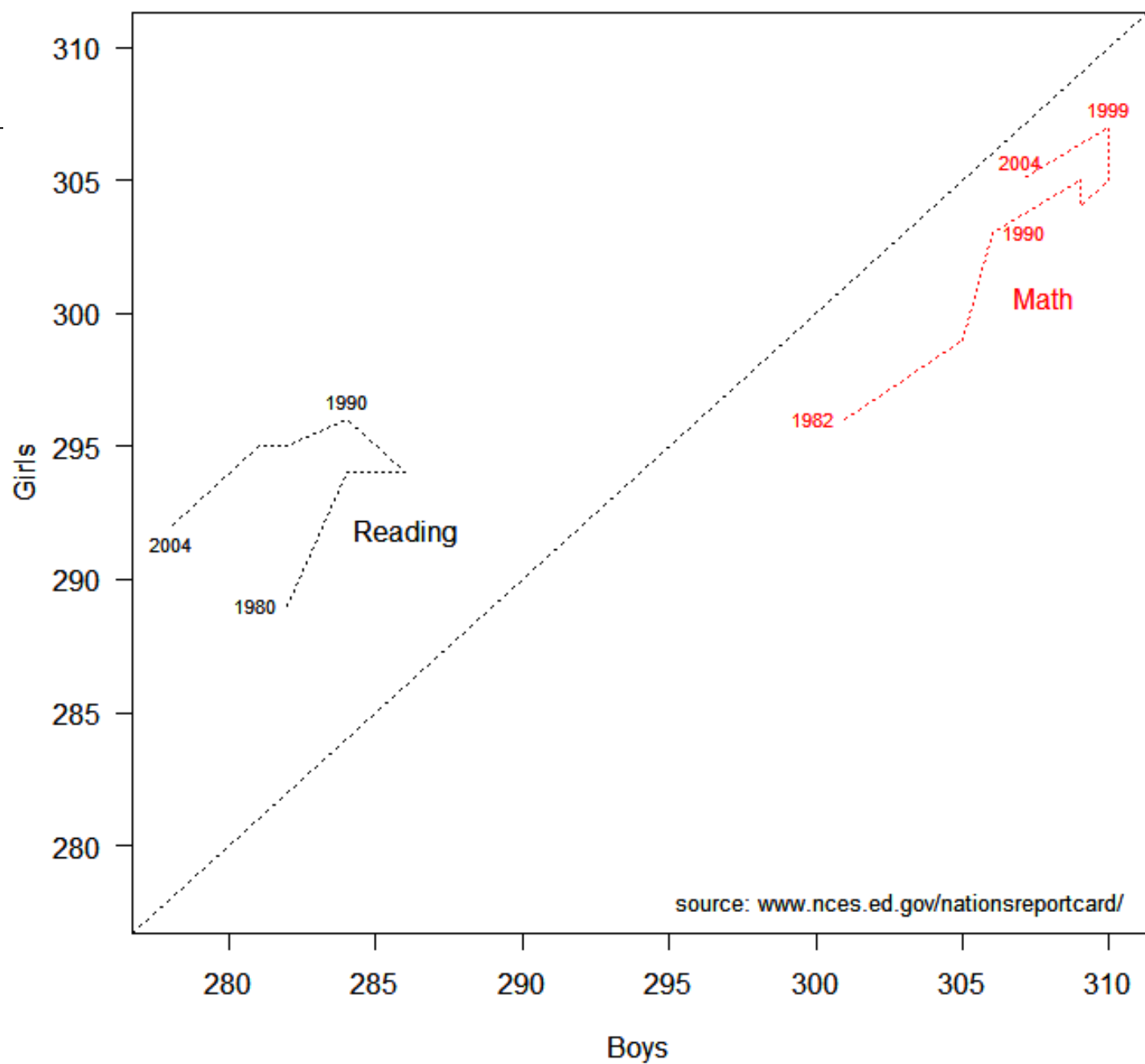
in billions of 2000 constant dollars



Federal outlays and receipts as percentage of GDP



NAEP reading and math scores, 17 year olds

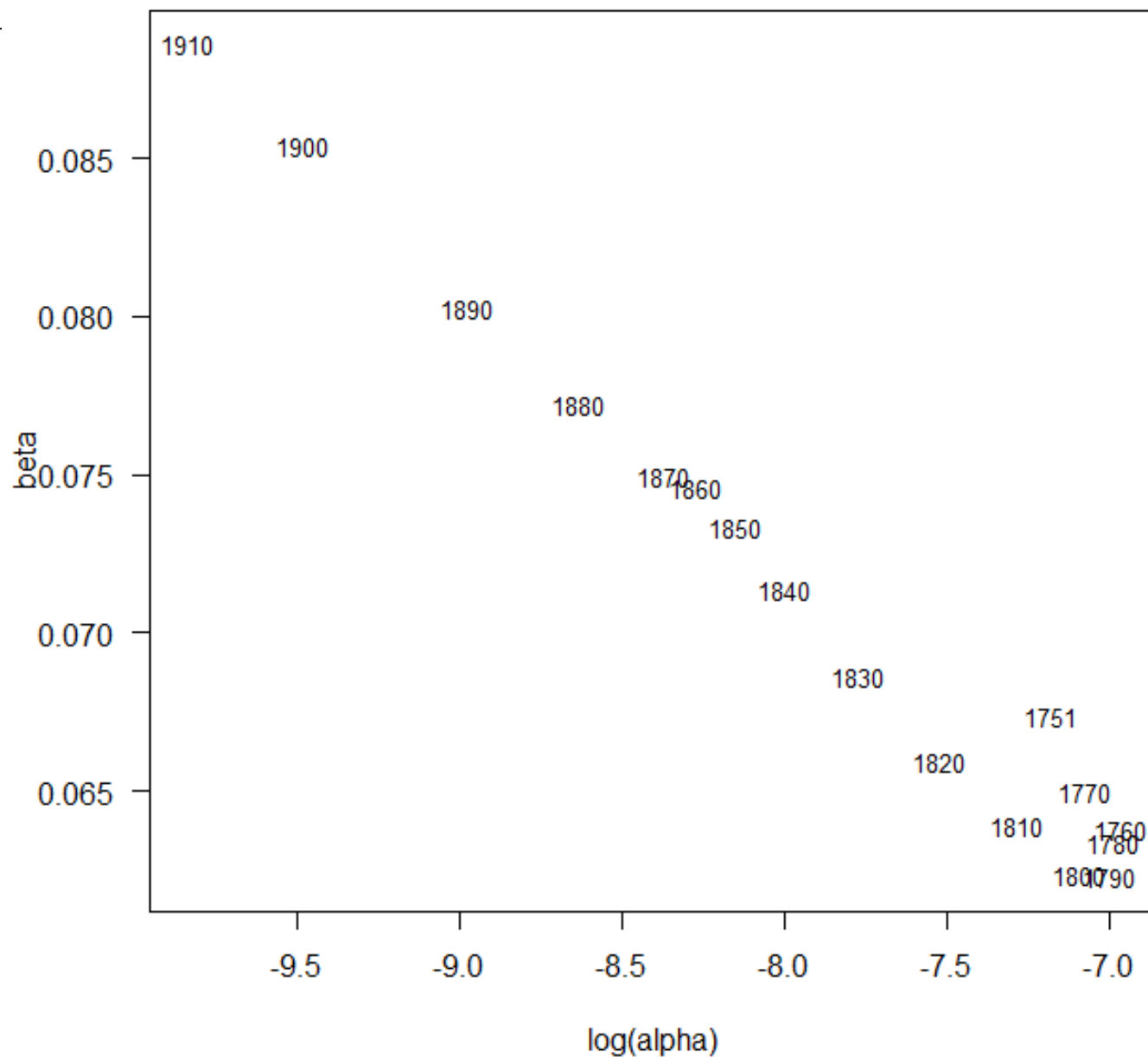


plot summaries for simplification

- when all subset have same contrasts, plot subset summaries
sometimes can get away with it even if not all subsets have all same contrasts—but then must be doubly careful
helps to identify patterns
plot and identify extremes, leave middle alone
this is the idea underlying “10 plus 10” plots
or, split into n groups (n small, like 3), and plot subsamples from each

Gompertz parameters, Sweden, 1751 to 1910 birth cohorts

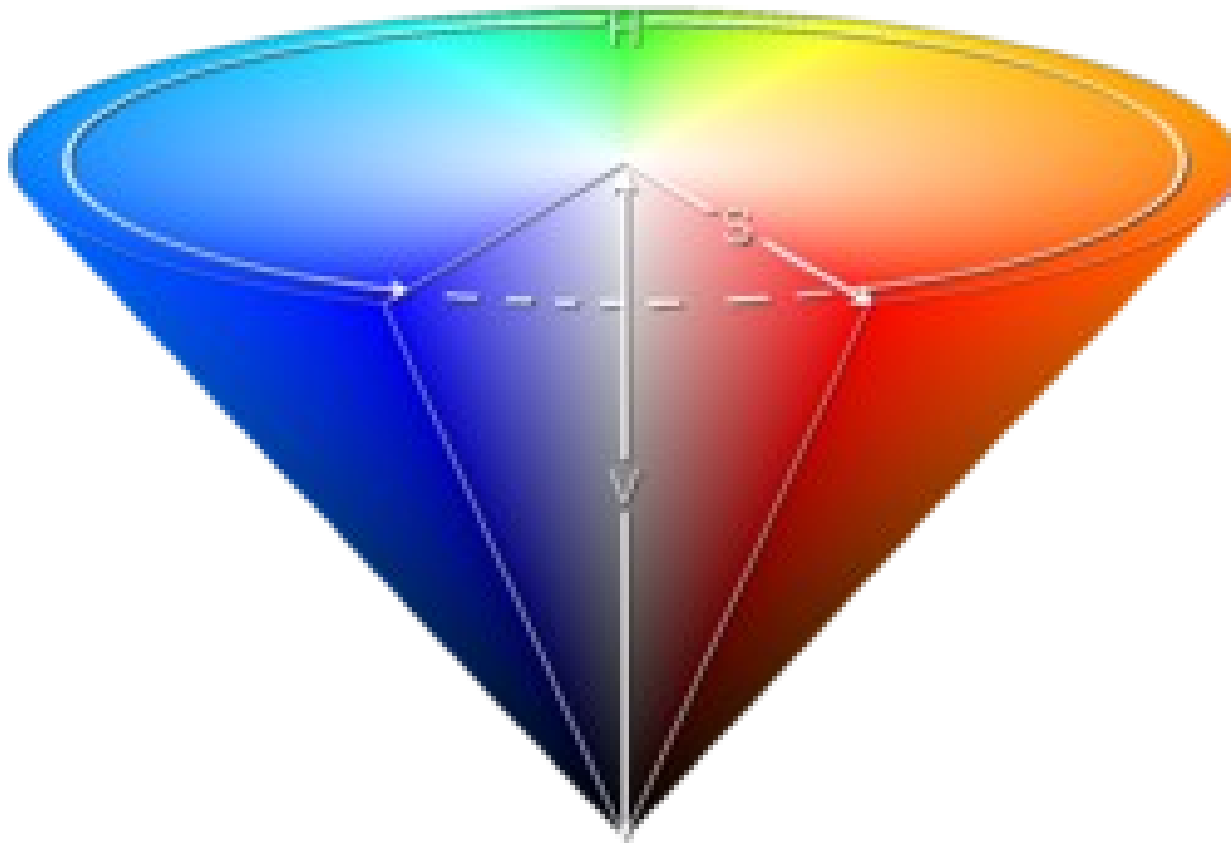
both sexes combined, ages 40 to 85



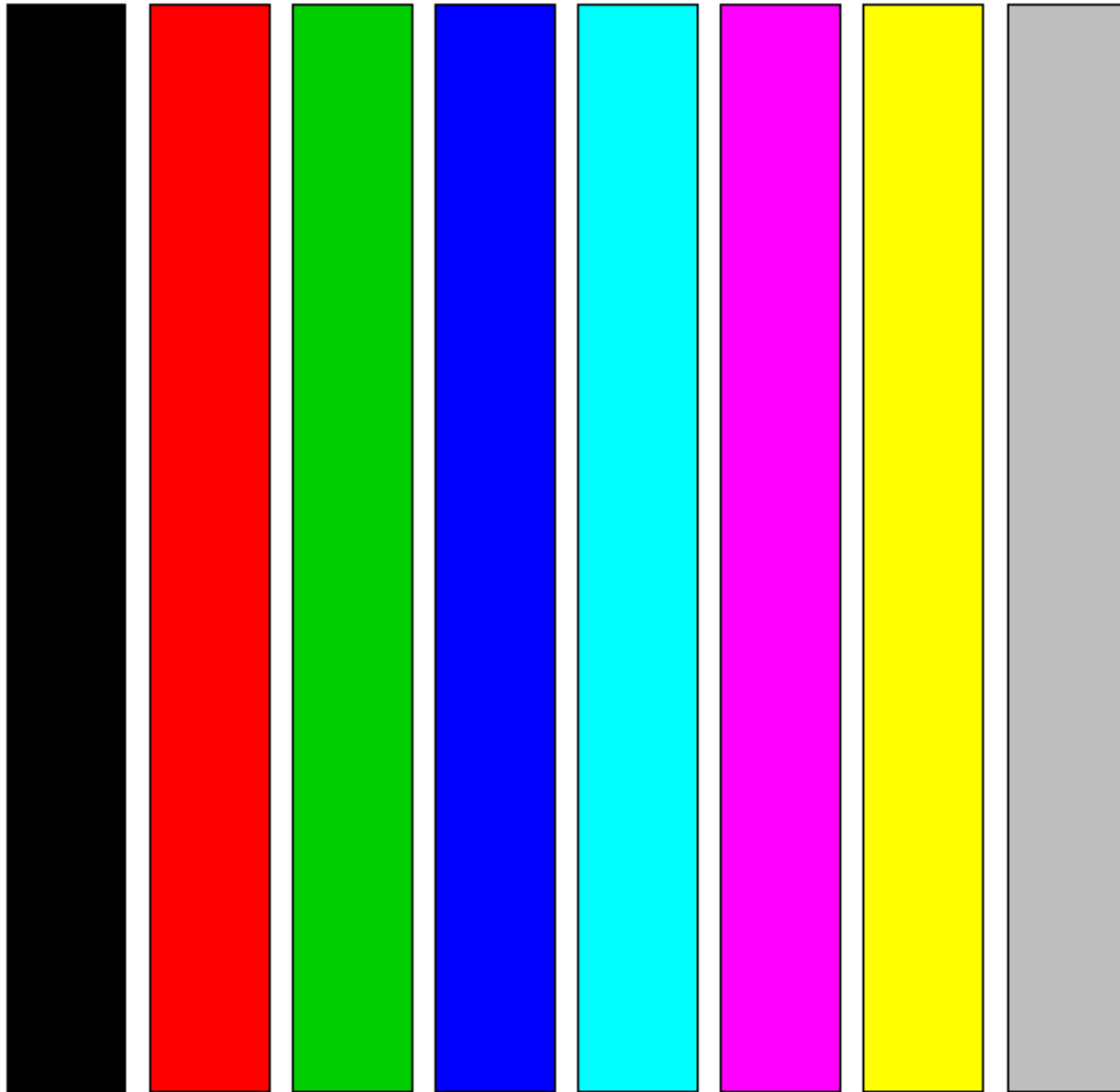
more on color

- HSV
 - h=hue, s=saturation, v=value
 - sometimes called HSL for hue, saturation, luminance
- equal impact colors
 - CIELUV and Munsell are systems of color perception
 - medium saturation, kind of pastel-like

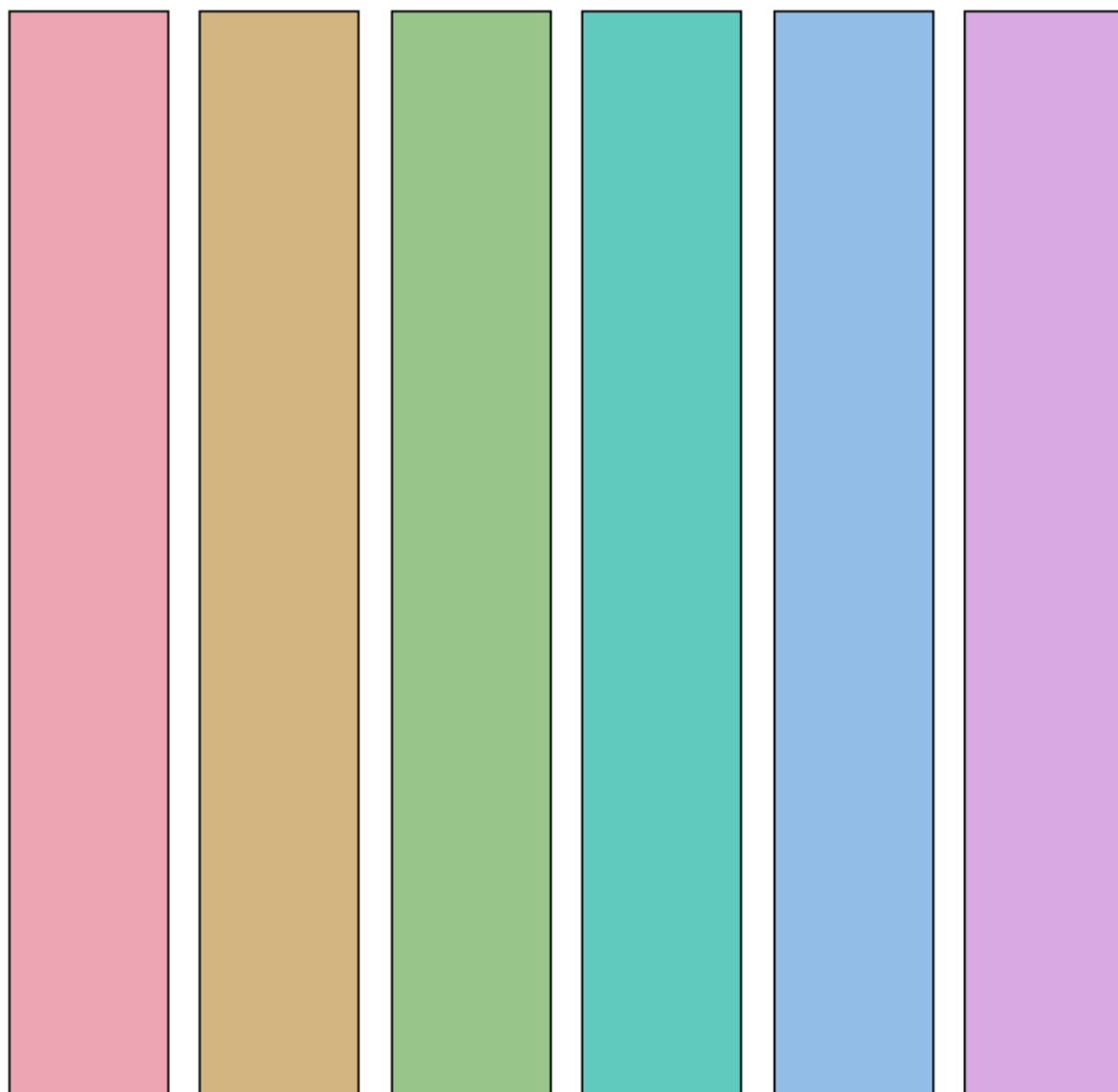
HSV cone



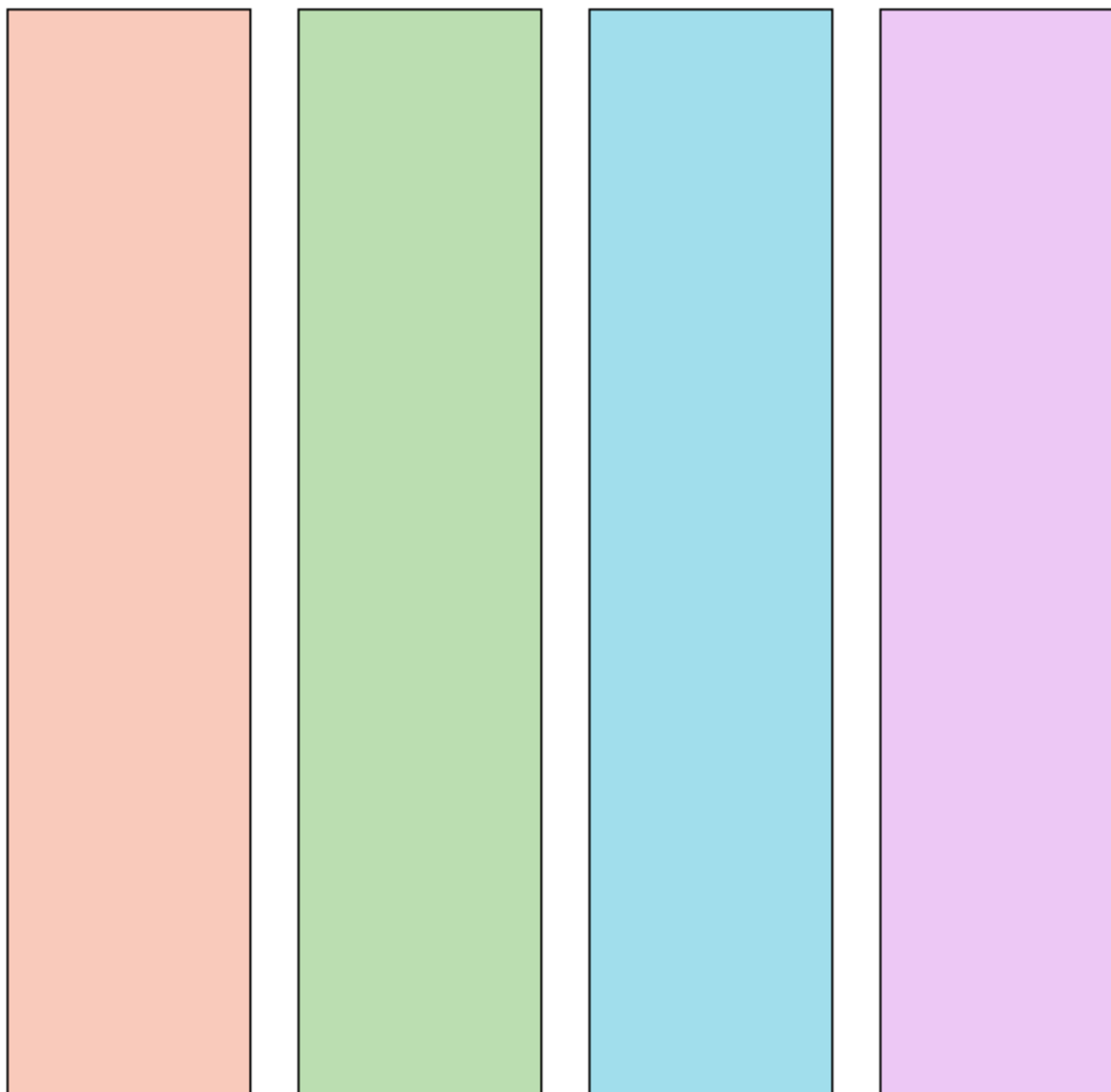
R's default colors



Example of equal impact colors



tetrad with maximal color differences



basic techniques

- show the difference
- identify outliers (or, label directly)
- group and order
- plot extremes
- multiple comparisons

slightly more advanced techniques

- smoothing
- straightening
- phase plots
- contours
- banking
- coloring

stuff I wanted to hide until the end

- friends don't let friends graph with Excel
but let's be realistic: sometimes you have no choice
dates in Excel are particularly a problem

nature

Explore content ▾

About the journal ▾

Publish with us ▾

Subscribe

[nature](#) > [news](#) > article

NEWS | 13 August 2021 | Correction [25 August 2021](#)

Autocorrect errors in Excel still creating genomics headache

Despite geneticists being warned about spreadsheet problems, 30% of published papers contain mangled gene names in supplementary data.

By [Dyani Lewis](#)

how a demographer changed
bicycle racing and design

finding a story to tell: graphical research methods

Robert.Chung@berkeley.edu

October 2024

analytical graphics

- not presentation graphics, not exactly statistical graphics, not quite exploratory graphics

presentation graphics tell your story; analytical graphics help you figure out what your story is

analytical graphics are “work product”

- not always refined enough for publication or presentation

Tufte? A lot of his techniques focus on effectively communicating quantitative findings. We'll focus on steps before that: uncovering interesting stories and questions in data

Very valuable for grad students and early career researchers

You may be familiar with the work of Edward Tufte. Excellent stuff, but his focus is on how to display quantitative data in a way that tells a clear and compelling story. What we're focusing on is at an earlier step: how to figure out what your story is. Tufte focuses on presentation graphics that clearly render the information that they're supposed to convey. What we're doing is sort of mid-way between exploratory graphics and presentation graphics. We're still trying to figure out what that information might be, or if there is any interesting information there at all. These are *graphical research methods* and, like all research methods, they reveal the most if you approach it in a (semi-) organized way. In these lectures we'll present some principles, some rules of thumb, some shortcuts, and some tricks that can help you identify questions that are deep enough to hold your attention long enough to finish a dissertation.

analytical graphics are for analysis

- often analytical graphics are used not to prove hypotheses but to help generate them
 - we don't always answer questions; we use graphical techniques to help us ask new ones
- often, the audience is YOU
 - visualize differences and contrasts
 - across time
 - across places
 - across treatments or policies
 - across conditions
 - there are tips, tricks, and techniques that help you in visualization

Many of the (very beautiful) graphics in Tufte's books are “one-offs” that have been hand-tuned in Adobe Illustrator. We don't have that luxury – we're looking for techniques we can use over and over in a way that makes production relatively easy. We're optimizing for analysis, not for presentation.

what if you already have a question?

- no problem. sometimes analytical graphics can help you focus on where to look, or to refine your question
- doesn't replace theory, or your research question. You can do both: that's allowed

Many of the (very beautiful) graphics in Tufte's books are “one-offs” that have been hand-tuned in Adobe Illustrator. We don't have the luxury – we're looking for techniques we can use over and over in a way that makes production relatively easy. We're optimizing for analysis, not for presentation.

basic approach

- maximize insight
- uncover underlying structure
- extract important variables
- detect outliers and anomalies
- develop (very) simple models
- not much testing; that's for later

Testing hypotheses is important, but this isn't (so much) about that. We're looking for stories to tell, and just like writing, you try to separate the writing from the editing.

what we'll do in these lectures

- some examples
- some principles
- some basic tricks
- some slightly more advanced tricks
- you'll have a chance to try out some of these tricks before next week

In past years, we've used a prepared data set for this, but that requires more time commitment than I had this year.

Warning: a lot of these examples are pretty old, because I've been working on versions of this for about 20 years. So if the examples seem anachronistic, that's cuz they are. The principles and tricks are (hopefully) still useful.

I sometimes wonder how AI will affect these lectures.

the three things we're looking for

- look for
 - pattern
 - unexpected pattern
 - deviations from pattern
- these generate questions
- questions and how you address them are often the basis for papers or chapters or dissertations or careers
- “The data speak for themselves, but their voices are soft and sly”
so we're looking for ways to amplify their voices

demography

- demography is the study of populations, their characteristics, relationships among characteristics, and how they change
- you probably already know how to examine characteristics; we'll look at ways to highlight relationships among the characteristics and how they change
- in particular, we'll often look for models to help us understand the relationships among characteristics

We have, in demography, a pretty sweet situation: we have relatively strong models that give us clues about how variables might (ought?) to be related. So sometimes we'll want to look at potential relationships. That's an advantage that not all fields have.

the purpose of models

- “The purpose of models is not to fit the data but to sharpen the questions” – Sam Karlin
- We’ll use analytical graphics to help us sharpen questions

simple tools

- simple tools used intelligently (well, we can always dream) rather than complex tools used stupidly
rules of thumb, not hard rules and regulations
- a handful of plots and a handful of tricks
lots of specialty type graphs, but we try to avoid too many of them until we know our story
- xy plots are a hugely useful invention
with one or two exceptions we'll focus mostly on ways to enhance xy plots
- decoding the language of graphs can be complicated, so we build on familiar beginnings

Simple tools don't always mean that the doing is simple, or that the results are simple. Here we mean “simple” in the sense that these are basic foundations on which we'll build. We'll learn about strengths and weaknesses of different tools.

Excel is a simple (and widely available) graphics tool but, as we will see, in certain uses it's too simple and in others it's not complex enough. However, because it is ubiquitous, later on if you insist we can go over some things you can do to neutralize some of the bad things that Excel does.

how graphing helps

- we can only make sense of a handful of numbers at a single time
pages of dense tables are good for detail, evidence, and re-analysis but poor for understanding
- eye-brain is good at seeing patterns in large numbers of values
though it can be fooled—we'll present some problems that can mislead the eye
- therefore
 - use graphs when pattern is important
 - use tables when exact details are important
 - graphs and tables are complements, not replacements. (You can do both: that's allowed)

People who focus on presentation (like, for example, Tufte) often appear to disdain tables in favor of graphics. However, because we tend to do both analysis and presentation, I think of them as complements. Use each (judiciously) when appropriate.

apophenia

- apophenia is “the experience of seeing patterns of connections in random or meaningless data”
- we'll occasionally accept a little “type I error” when we're looking for interesting questions – as long as we back it up later with real confirmatory analysis



This is the famous (?) “Martian face” as photographed by one of the Viking spacecraft missions to Mars.

Anscombe's data

x1 y1		x2 y2		x3 y3		x4 y4	
10	8.04	10	9.14	10	7.46	8	6.58
8	6.95	8	8.14	8	6.77	8	5.76
13	7.58	13	8.74	13	12.74	8	7.71
9	8.81	9	8.77	9	7.11	8	8.84
11	8.33	11	9.26	11	7.81	8	8.47
14	9.96	14	8.10	14	8.84	8	7.04
6	7.24	6	6.13	6	6.08	8	5.25
4	4.26	4	3.10	4	5.39	19	12.50
12	10.84	12	9.13	12	8.15	8	5.56
7	4.82	7	7.26	7	6.42	8	7.91
5	5.68	5	4.74	5	5.73	8	6.89

built into R.

`help(anscombe)`

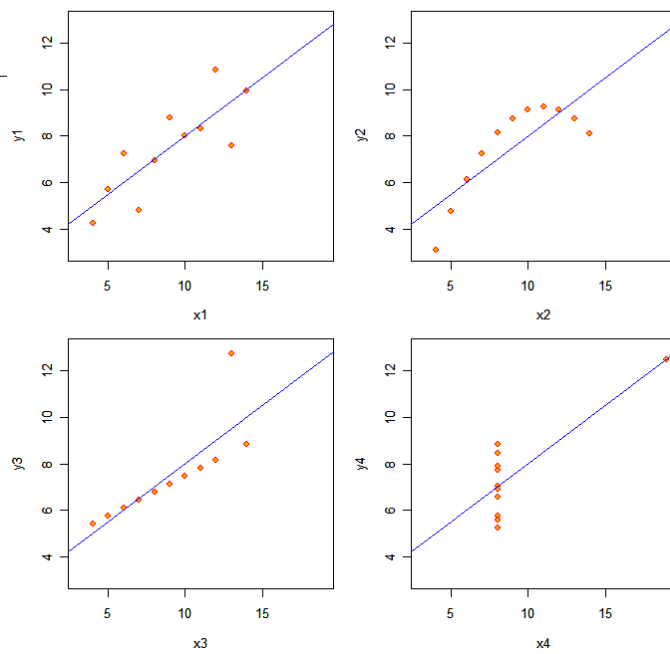
This is a pretty famous set of manufactured data. We'll see why in a moment.

Anscombe's data

- same means, sd's, correlation, regression slope, fit
$$\begin{aligned}\text{mean}(x_1) &= \text{mean}(x_2) = \text{mean}(x_3) = \text{mean}(x_4) = 9 \\ \text{mean}(y_1) &= \text{mean}(y_2) = \text{mean}(y_3) = \text{mean}(y_4) = 7.5 \\ \text{sd}(x_1) &= \text{sd}(x_2) = \text{sd}(x_3) = \text{sd}(x_4) = 3.32 \\ \text{sd}(y_1) &= \text{sd}(y_2) = \text{sd}(y_3) = \text{sd}(y_4) = 2.03 \\ r(x_1, y_1) &= r(x_2, y_2) = r(x_3, y_3) = r(x_4, y_4) = 0.816 \\ y^* &= 3 + 0.5 x^* \text{ with } r^2 = 0.667\end{aligned}$$
- so, conventional linear models make them look alike
- what will you see if you graph the data?

What's worse, a depressingly large proportion of analysts will stop right there.

Anscombe's 4 Regression data sets



`example(anscombe)`

the NJ Pick-It lottery

- each bettor selected a 3-digit number between 0 and 999
- each ticket cost 50 cents
- all bettors who held the winning number split the prize money. The size of the prize depended on selecting the winning number **and** on the number of players who chose that number
- what would you want to know?

built into S-plus; in fact, I stole this example from “The Blue Book”

```
dat= read.table("http://anonymous.coward.free.fr/misc/lottery.txt",header=T)
```

winning numbers and prize amounts

(810, \$190.0)

(156, \$120.5)

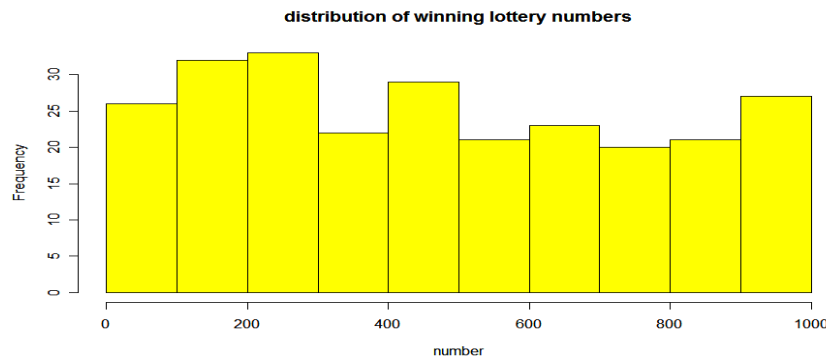
(140, \$285.5)

(542, \$184.0)

and so on for 254 consecutive days

strategy 1: choose a winning number

- since we have data on the winning numbers, see if there's a pattern we can exploit to pick the winners
- examine the distribution of winning numbers using histograms or stem-and-leafs



```
dat = read.table("lottery.txt",header=T)
attach(dat)
hist(number,col="yellow")
```

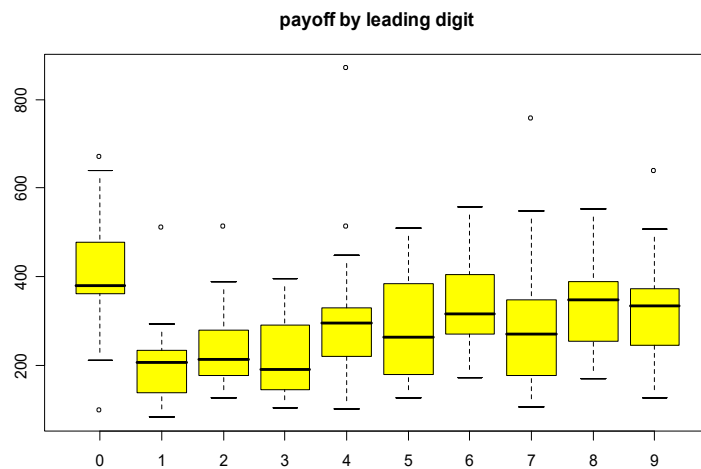
Are there more winning numbers between 100 and 299 than we'd expect by chance?

There are $n=254$ total "picks". We'd expect $n \cdot p$ in each of the 10 bars, or about 25.4, with $sd = \sqrt{n \cdot p \cdot q} = 4.8$.

So we could draw horizontal lines at $25.4 \pm 2 \cdot 4.8$ to mark 2SD CI's. None of bars exceed CI.

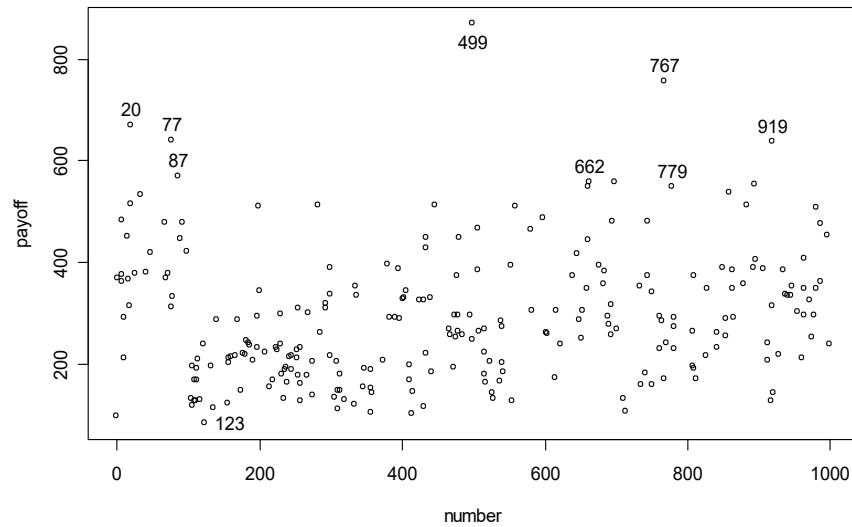
So since we can't pick a winning number, maybe we can see a pattern in how much each winning number won.

strategy 2: choose a winning number that few others pick



```
boxplot(payoff~number %/% 100, col="yellow")
```

Aha. It seems as if when asked to pick a number, not many people picked numbers below 100.



```
plot(number,payoff)
```

```
identify(number,payoff,number)
```

What do you see? Identifying the points helps you see that numbers with high returns either have a zero for the first digit, or double digits. Note that the lowest payoff for a winning number was “123”.

Identifying high and low payoffs tells you something interesting that you wouldn't have known from more standard analysis. We'll come back to this later when we talk about labeling.

why little circles?

- easier to distinguish overlapping points
- especially with jittering

little circles show up better when you have points close to another. This is part of the investigation into graphical perception that was done at Bell Labs.

five rules

- rule 1: graph lots
- rule 2: use what the eye is good at (and avoid what the eye is bad at)
- rule 3: find the right contrast and show it
- rule 4: make it easy to spot pattern, and deviations from pattern
- rule 5: plot models, not just the data

rule 1: graph lots

- only one out of 50 graphs will “work” so to get a handful of workable graphs, graph lots
- good graphing principles help raise your yield of workable graphs
- better if you can generate lots of simple graphs quickly even if they’re not perfect
- for you, not for presentation (at this stage) so don’t obsess on look (though I’m showing you the survivors of hundreds of graphs, so they’re cleaner and not quite representative of the messy graphs I usually produce: my working graphs usually don’t have titles, clear axis labels, etc.)

Here's a tip that often helps me raise the yield: sketch what you think your graph will look like on a piece of paper. This sounds archaic but, for me at least, it seems to help. This may or may not work for you—you have to pay attention to what works for you.

rule 2: use what the eye is good at (and avoid what it's bad at)

- we need to know something about how the eye-brain perceives graphics
 - what it's good and bad at, and an ordering or hierarchy
 - “optical illusions” and traps to avoid
 - techniques to exploit strengths and minimize weaknesses

graphical perception

- quantitative pattern recognition by
 - detection: recognition of geometry
 - assembly: grouping of detected elements
 - estimation: assessment of relative magnitudes
- the human eye-brain can be fooled
 - optical illusions
- need to help it out
 - grouping, ordering, highlighting help to identify patterns

a hierarchy of graphical perception

- position along common scale
- position along identical non-aligned scales
- length
- angle, slope
- area
- volume
- shading, color (good discrimination but poor ordering)

Cleveland and McGill, 1984

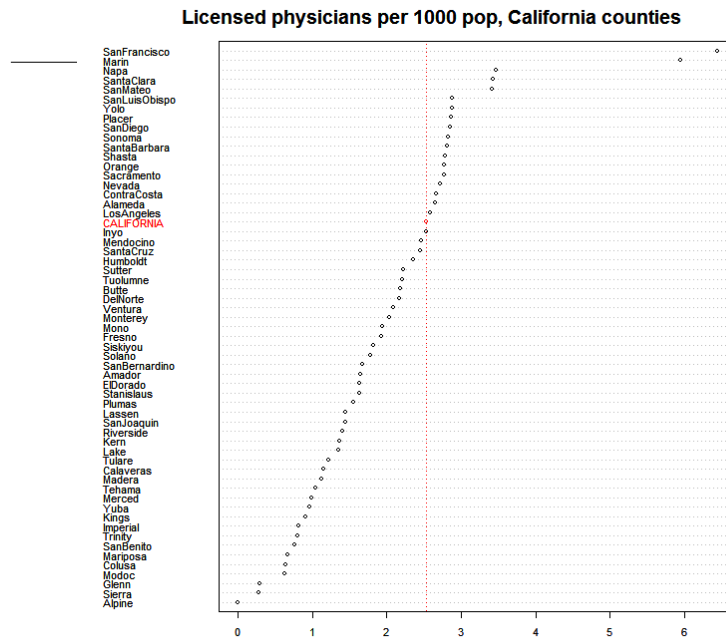
dotplots are plots where quantitative values are displayed as positions along a common scale.

Pie charts, which are ubiquitous and rely on decoding angles, are bad. (Stacked pie charts, aka “spie” charts, can be useful in certain circumscribed situations, but don't go hog wild).

Notice that the familiar age pyramid with males on one side and females on the other uses non-aligned scales, and lengths and areas when it could as easily use position on the same scale.

Color-coded maps are often good for detecting pattern *when the pattern is geographical*, but poor when the pattern is just about anything else. Because colors are poorly ordered by the human perception, it's often hard to decode quantitative variables by color – no matter how nice the final product looks. Use color to highlight, or to encode categorical variables, but avoid using it for quantities unless you know what you're doing. For example, if a pattern changes over time, sometimes maps just confuse things.

dotplots



dotplots are preferred to barcharts because bars have width and, thus, area and area is less reliably interpreted than position.

```
dat = read.csv("docs-percap.csv")
with(dat,dotchart(ratio,label=County))
```

area

- pie charts require estimation of area
- human perception of relative areas is conservative, i.e., shrinkage toward 1.0
- shape affects estimation of area
 - concave shapes appear larger than convex
 - maps are good for context and clustering, not so good for comparisons of quantitative amounts
- color intensity affects estimation of area.
 - highly saturated colors appear larger

shape affects perception of area: if you look at a map of the southeastern states, it appears that Florida is larger than Georgia – Florida is an odd shape, GA is more compact.

This is another potential problem with using maps to encode quantitative information: Denmark is about one-eighth the size of Norway in land area, but it's population is slightly larger, so encoding a mortality rate on top of a map adds one more layer of information for the eye-brain to decode (just so it can reject it as irrelevant).

One of the problems with “spie” aka “stacked pie” charts: are you comparing areas or lengths of radii?

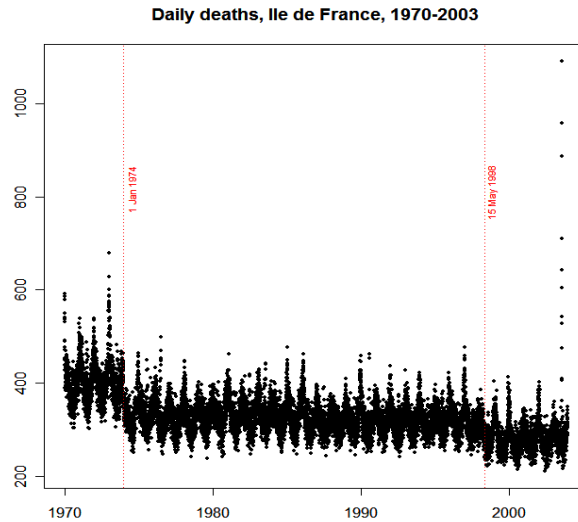
color

- human eye good at discrimination, poor at ordering
 - use for categories, not for quantitative coding
 - hues are not ordered
 - use for highlighting, patterning, especially in combination with small multiples
- more on color, later

rule 3: find the right contrast and show it

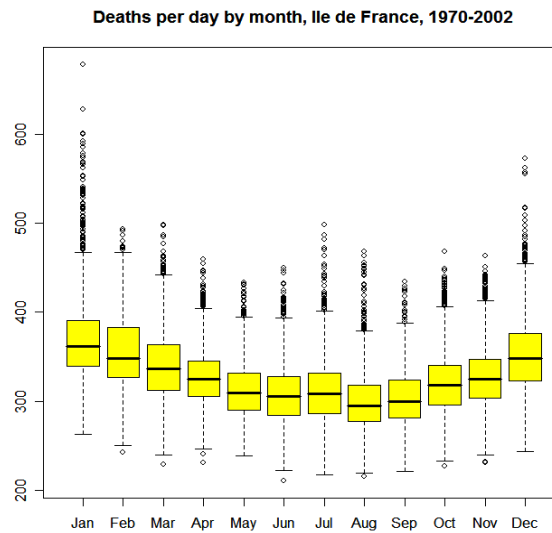
- don't rely on the eye to do differencing
 - if you're interested in the difference between two lines, don't show the lines and rely on the eye to calculate the difference; calculate the difference itself and show it
- Tukey mean-difference (aka Bland-Altman) plots
 - levels vs. differences
- fits and residuals
- different contrasts can give different insights

differences by time



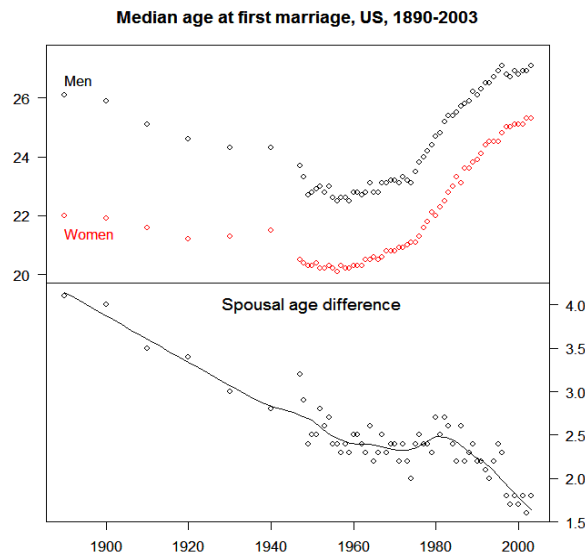
Time series plots are, next to pie charts, perhaps the most common of all charts. They do show how something changed over time; however, while they tell you when something happened, they don't tell you why it happened. You have to fill in the context. Most of what we're trying to do here is figure out ways to generate ideas and hypotheses, so showing when something happened isn't always very helpful.

differences by category



This gets a little more interesting. Simply by grouping by category, we see a pattern emerging that may be clearer than the previous slide.

differences between lines



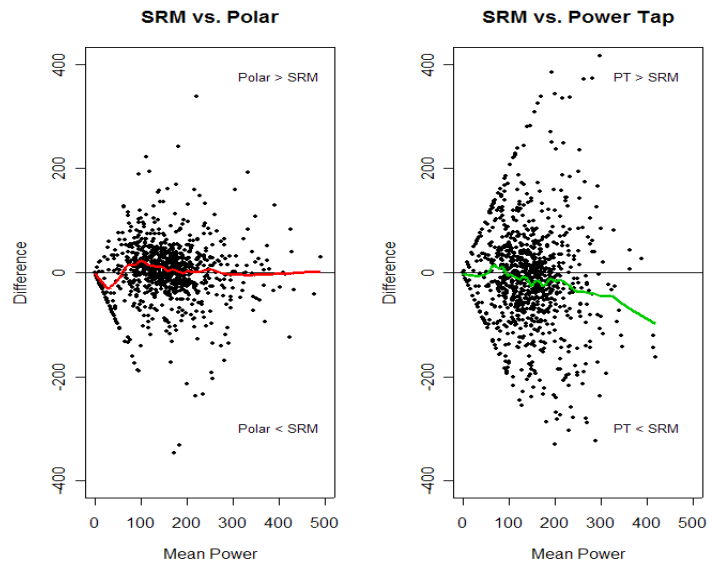
Just by looking at top panel, eye might have thought spousal age diff was constant after mid-1980's. Eye would never have picked up bump in late 1980's. Humans have trouble judging differences between two curves, *especially* when those two curves are steep.

If your story is about the difference, calculate and show the difference.

```
dat = read.csv("marriage-age.csv",comment="#")
head(dat)
attach(dat)
yl = range(Men,Women)
plot(Year,Men,ylim=yl)
points(Year,Women,col="red")

plot(Year,Men-Women)
lines(lowess(Year,Men-Women,f=0.5))
```

differences and means

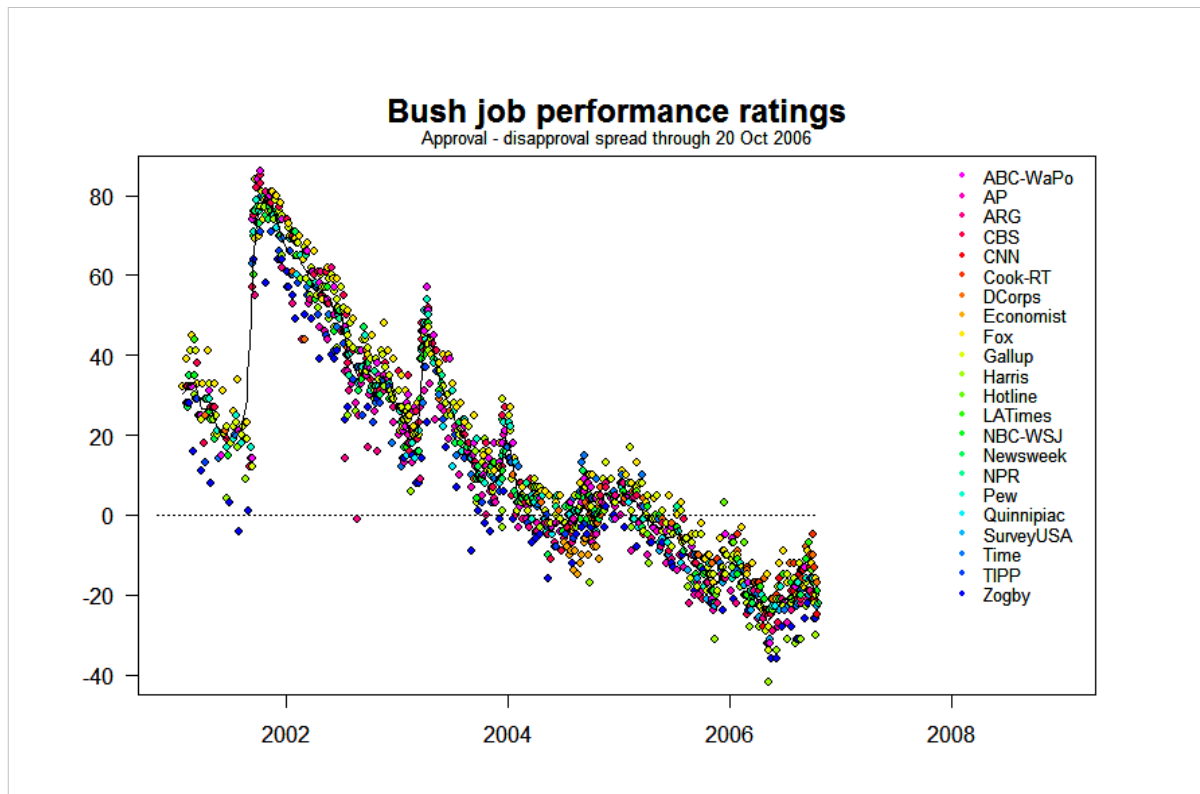


Tukey mean-difference plot. Shows whether difference grows with level.
x-axis is $(\text{SRM} + \text{PT})/2$, y-axis is $\text{PT} - \text{SRM}$

$\text{plot}((x+y)/2, (x-y))$

differences from fits

- $\text{data} = \text{fit} + \text{residual}$
- the classic residual plot
 - “fit” can be broadly defined: $\text{data} = \text{smooth} + \text{residual}$



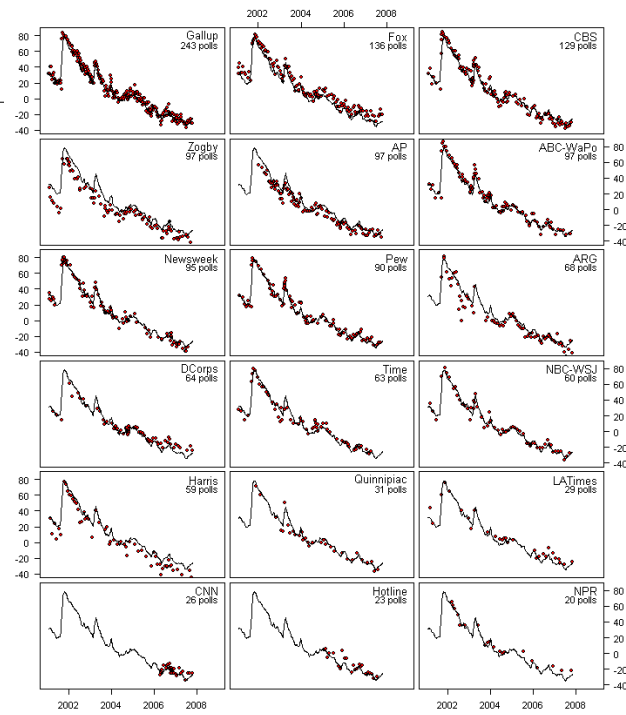
lowess

Too many colors

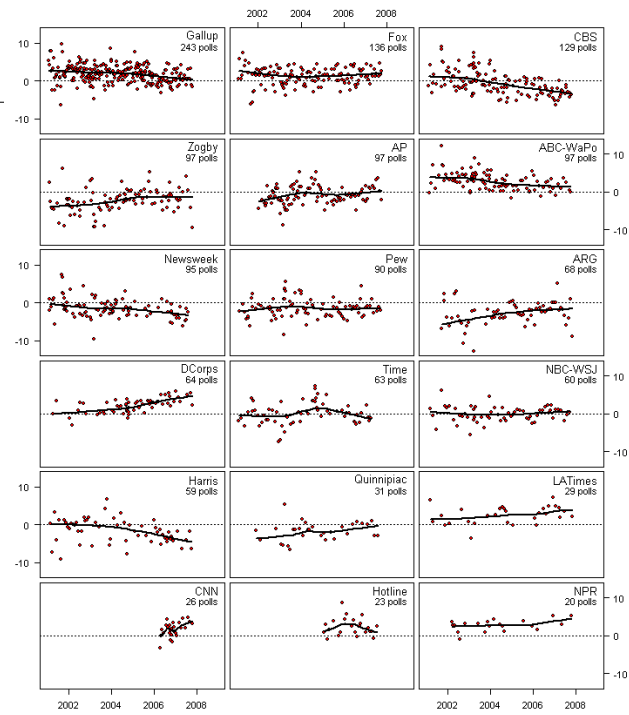
Legend hard to use

```
dat = read.csv("polldata.csv",comment="#")
head(dat)
dat$date = as.Date(dat$date,"%m/%d/%Y")
with(dat,plot(date,approval-disapproval))
```

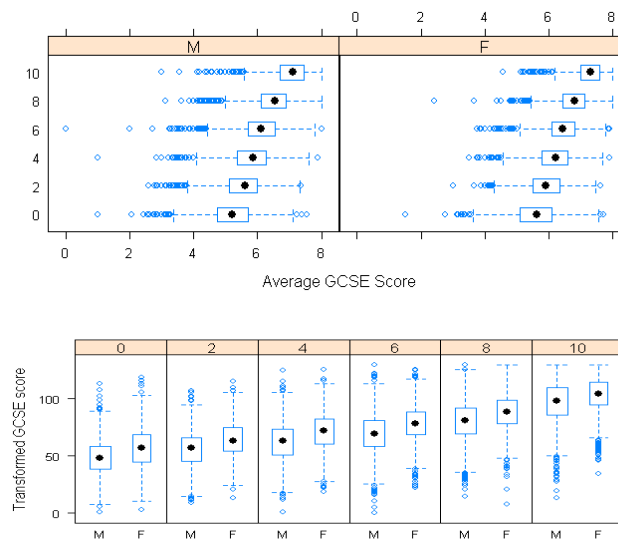
Net job approval for President Bush



Some polls consistently below trend: others above



same data, different contrast



These are Figs. 3.11 and 3.12 from Sarkar's Lattice book. The bottom figure contrasts male and female scores on a test and makes it easier to see that transformed male scores seem to improve more than female scores.

To be clear, it's **still** hard to see the difference. If you were doing a presentation you might want to show this in a different way. (Q: how might you do it?) The point isn't that this contrast is the best way to see this effect – it's that if you hadn't done this contrast you might not have seen it at all.

```
data(Chem97, package = "mlmRev")
```

```
## Figure 3.11
```

```
bwplot(factor(score) ~ gcsescore | gender, data = Chem97, xlab = "Average GCSE Score")
```

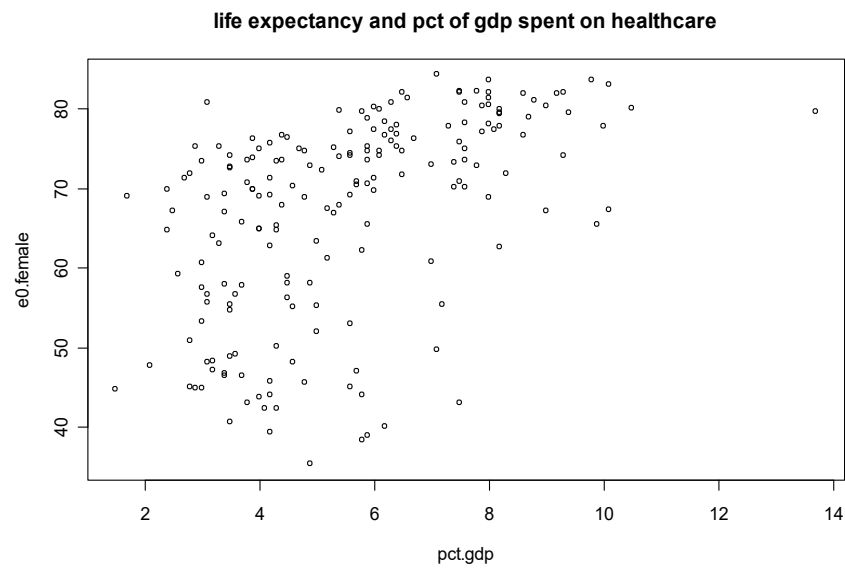
```
## Figure 3.12
```

```
bwplot(gcsescore^2.34 ~ gender | factor(score), Chem97, varwidth = TRUE, layout = c(6, 1),  
ylab = "Transformed GCSE score")
```

rule 4: make it easy to spot pattern

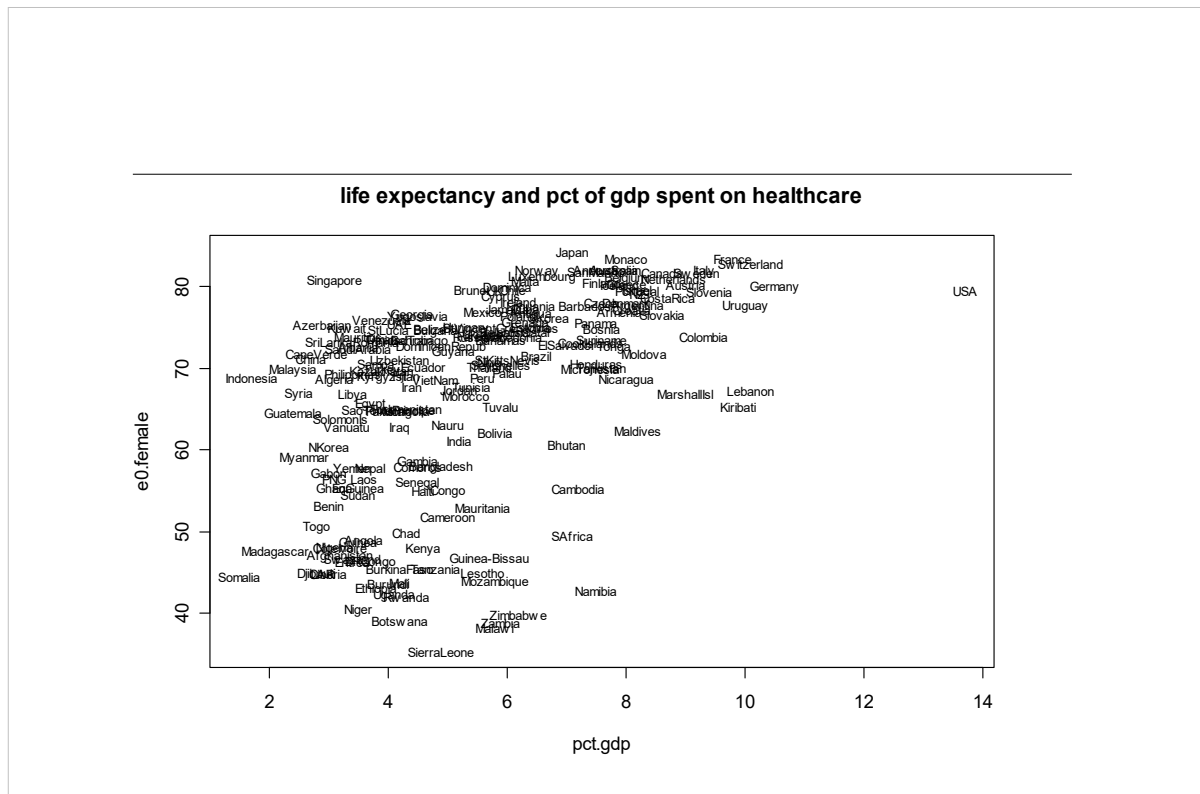
- add information depth, not (unnecessary) complexity
sometimes two plots are better than one complex plot (and
sometimes it isn't)

direct labeling



```
dat = read.csv("http://anonymous.coward.free.fr/misc/who.csv")  
with(dat, plot(pct.gdp, e0.female))
```

This tells you something about the pattern (which is good) but it still kind of sits there on the screen. Can we do better?



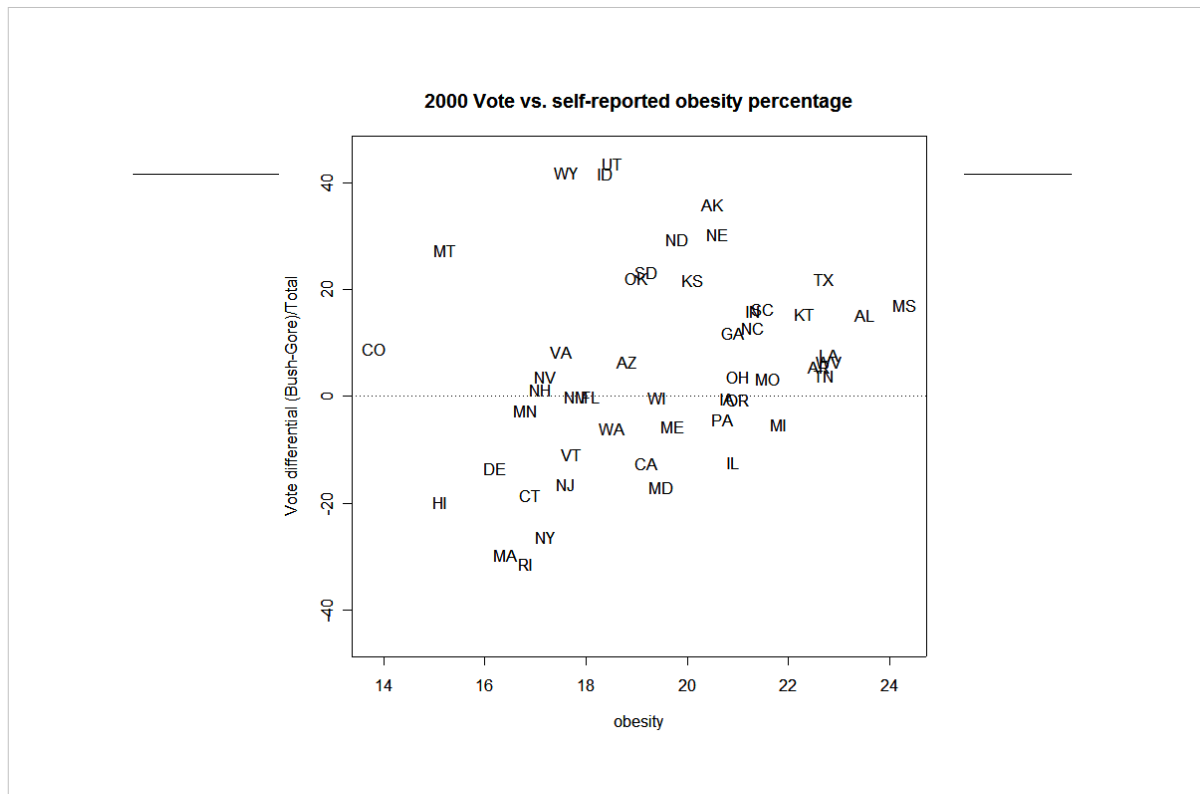
```
plot(pct.gdp,e0.female,type="n")
text(pct.gdp,e0.female,country,cex=.7)
```

Labeling helps a lot, even though you sacrifice exact placement of the point (since the names are of various lengths—in this example, the labels are centered on the data value).

You see the US, off to the right. You see Singapore way off to the left. You can spot sub-Saharan Africa. The labels have added context that helps you to think about new questions.

Note, your eye probably went to the outliers: we all tend to pick up on the outliers and to think about their story, and the mass in the middle looked squeezed and muddled anyway. This is a good technique: label the outliers and leave the middle as points. This is, in spirit, related to Tukey's 10-plus-10 plots: he recommended plotting the top 10 and bottom 10 in order to see if he could pick out patterns. In our xy-plots, you might want to play with identifying the "envelope" and leaving the inside as dots.

If you have lots of rows of data, you can try splitting them into thirds by one of the variables and then plotting a subsample of each third.



1. The president has broad-based support.
2. Notice mountain states. We could have colored them separately if we wanted. Basically, the slope for the mountain west is about the same as for the rest of the states but offset with a different intercept. You could do a regression with a dummy variable with mountain west states.
3. This isn't causation. We're looking for interesting stories, not proving that Bush supporters are fat-assed.

```
dat = read.csv("fat-vote.csv",comment="#")
```

```
head(dat)
```

```
with(dat,plot(obesity,(bush-gore)/total))
```

```
with(dat,plot(obesity,(bush-gore)/total,type="n"))
```

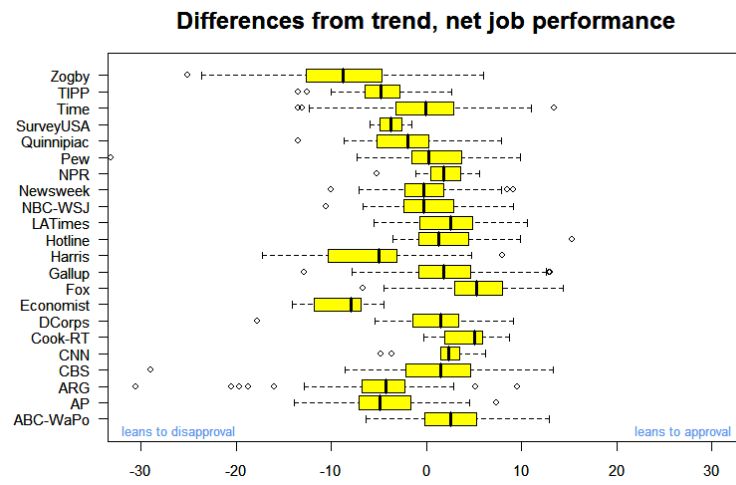
```
with(dat,text(obesity,(bush-gore)/total,as.character(state)))
```

more on labeling

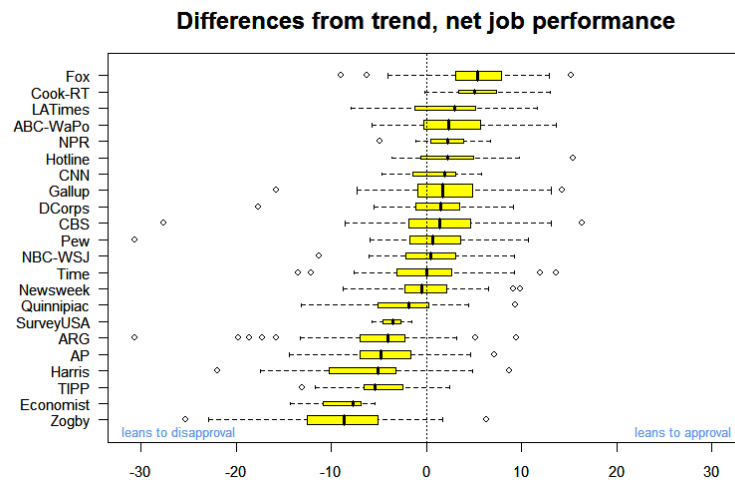
- direct labeling of lines often better than legends
 - particularly good when combined with line color
 - symbol plus line type often too busy to decode
 - looking back-and-forth at a legend is distracting

ordering

- default ordering for categorical variables is often alphabetical
 - that makes categories easy to find, but hard to compare
 - example: country data are often ordered by name of country rather than by the variable you're interested in
- find an ordering that makes sense and use it
 - if you are interested in mortality differences among countries, order by mortality not country name
 - this helps you spot and evaluate small differences between countries



Different polling houses exhibit different amounts of “house effects” in their polls. These boxplots show, for each national polling firm, the distribution of residuals from a model of presidential job approval. The ordering is alphabetical from bottom to top.

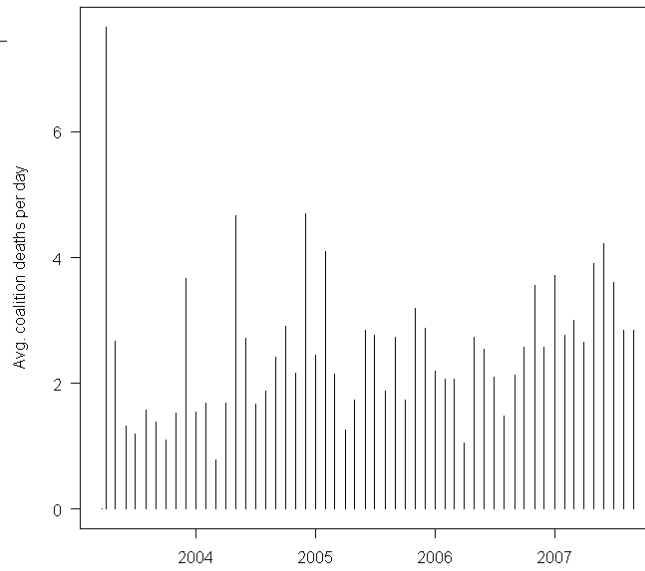


Similar plot to previous, but ordered by median residual, with boxwidth proportional to square root of sample size, and a reference line for zero average house effect.

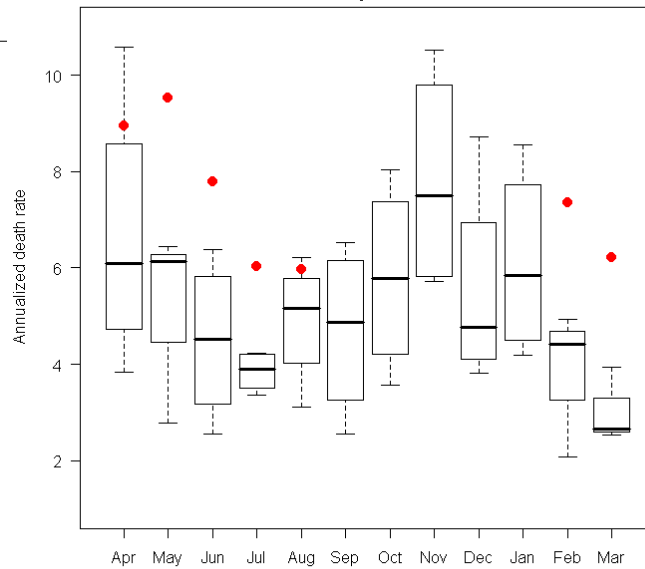
grouping

- grouping (done well) helps with pattern recognition
boxplots are a familiar way to group
- grouping (done poorly) obscures pattern
not all “obvious” groupings are informative
- next two slides show (almost) same data

Coalition deaths in Iraq, by month



Death rate in Iraq, coalition forces
excluding Mar 2003



multivariate comparisons

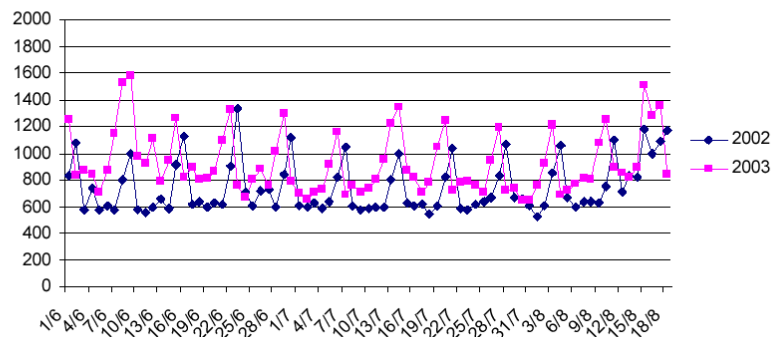
- show relationships
more importantly, give you ideas
- time series plots show you what happened when, but rarely why they happened
we'll want to dig deeper into the data to generate new questions about the 'why?'

Time-series plots are useful as descriptive summaries, but from an analytical point of view they don't usually provide enough motivation unless there are specific events that trigger sudden changes. They also tend to have relatively low information density, especially since one of the axes gets eaten up with time – time is a equi-interval variable, so it will often be more useful to use that axis for a variable whose observations are more complex. This is much more in the spirit of Descartes original insight.

patterns in multivariate data

- twenty students read numbers of ambulance calls for July 2003 off a graph. How can we summarize the results?

Graphique 5 : nombre d'interventions du SAMU 13 en 2003 par rapport à l'année précédente (2002)



```
dat = read.csv("http://anonymous.coward.free.fr/mpa/ps2/ps1results.csv")
```

```
head(dat)
```

```
attach(dat)
```

```
stars(dat)
```

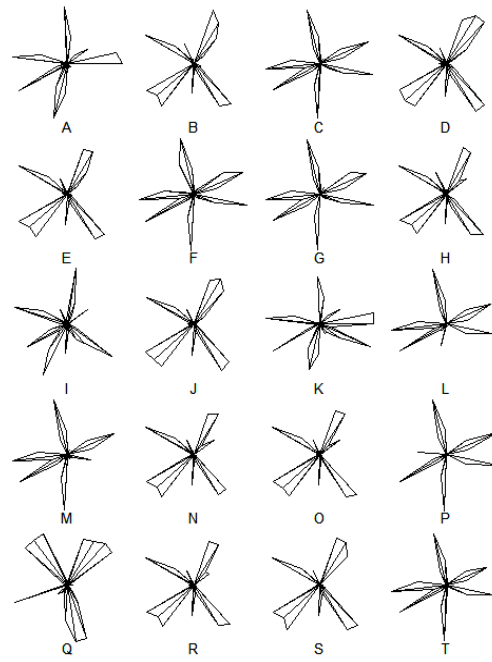
```
x <- as.matrix(dat)
```

```
stars(t(x))
```

```
stars(t(x),draw.seg=T)
```

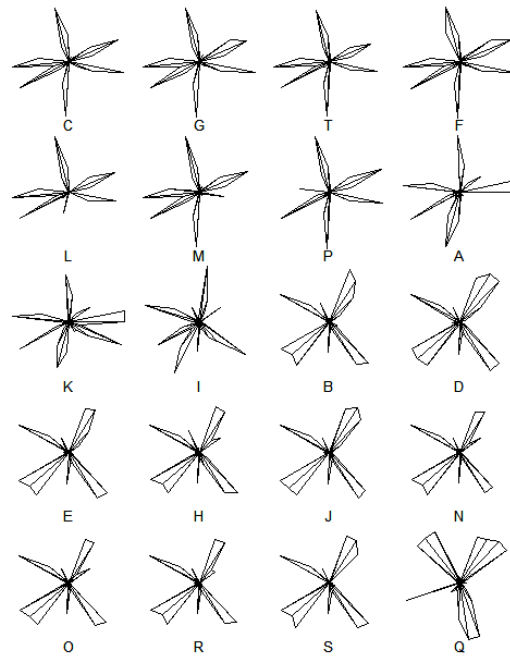
Notice, by the way, that the original graph appears to have been done using Excel-like graphics. Grid lines are good, but dark grid lines compete for your attention. If you must use Excel, dim down the grid lines by making them a light gray.

Starplot for PS1



stars(...)

Rough grouping of PS1



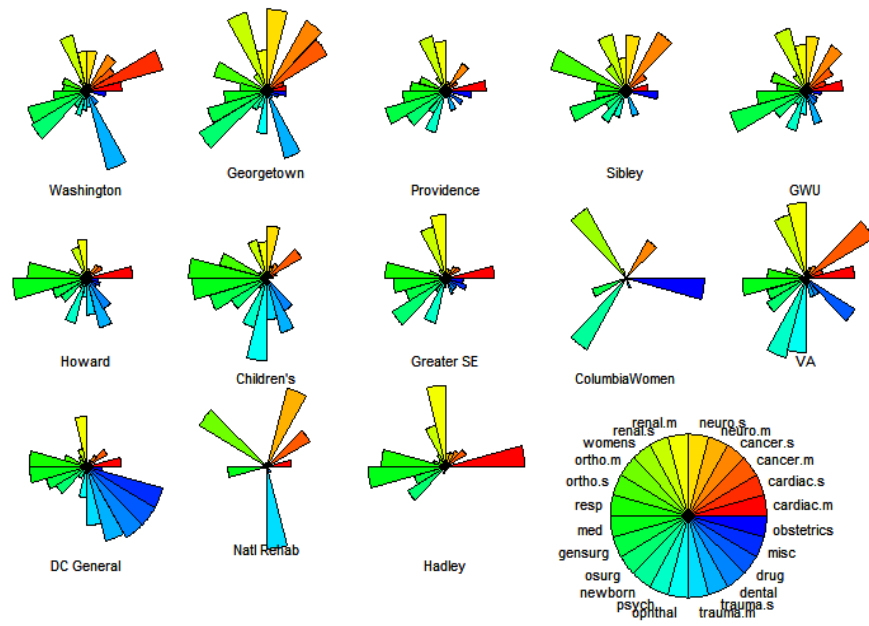
Stars can sometimes make it easy to spot patterns over many variables across individuals.

You can see another example using multivariate data on cars in R with `example(stars)`

dc hospitals

- thirteen hospitals
- twenty-four service lines
- do different hospitals specialize in different areas?

Hospitals and their service lines, 2000-1



stars(..., draw.segments=T)

Notice how DC General looks different from all other hospitals. Georgetown and GWU look vaguely similar. Washington Hospital Medical Center has a large cardiac surgery program.

stars are like pies

- except that angle is constant and radius varies
in pies, radius is constant and angle varies
that's why pie segments need labels
- watch out for colors

Though pie segments add to 100% and there's no such condition in stars.

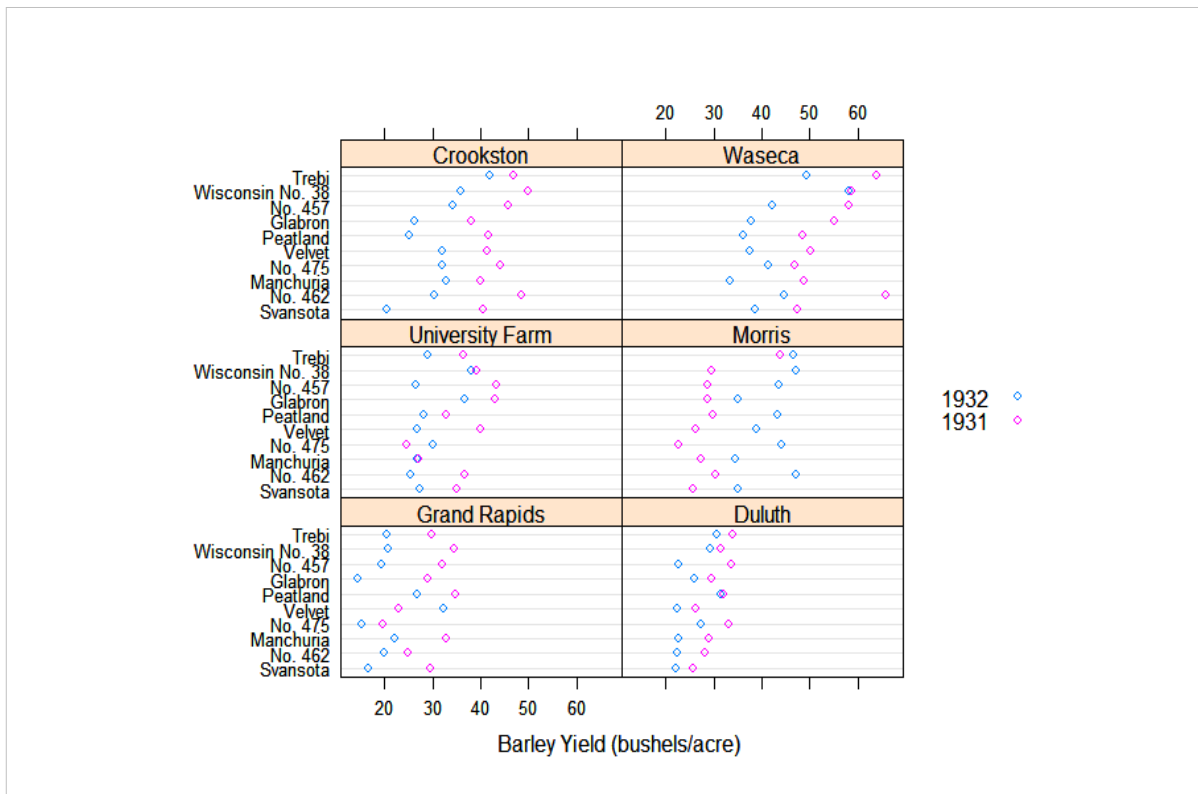
scatterplots

- can (sometimes) show more than two variables
 - can code categorical variables with color
 - can code some interval variables with size
- small multiples can show varying conditions
 - lattice (i.e., trellis) plots
 - use same scale and ranges, if possible, to enhance comparison

small multiples can be one of the best ways to exploit many comparisons at once. Try to use the same axis scales if at all possible.

dotplots (and trellis)

- conditioning plots
- barley yield
 - ten varieties
 - six plots
 - two years



```
library(lattice)

dotplot(variety ~ yield | site, data = barley, groups = year,
        key = simpleKey(levels(barley$year), space = "right"),
        xlab = "Barley Yield (bushels/acre) ",
        aspect = 0.5, layout = c(2,3), ylab = NULL)
```

Though many statisticians worked on this data set, it took nearly half a century for someone to notice that the data for Morris appear to have been switched (Cleveland, 1984). That was realized from a plot like this.

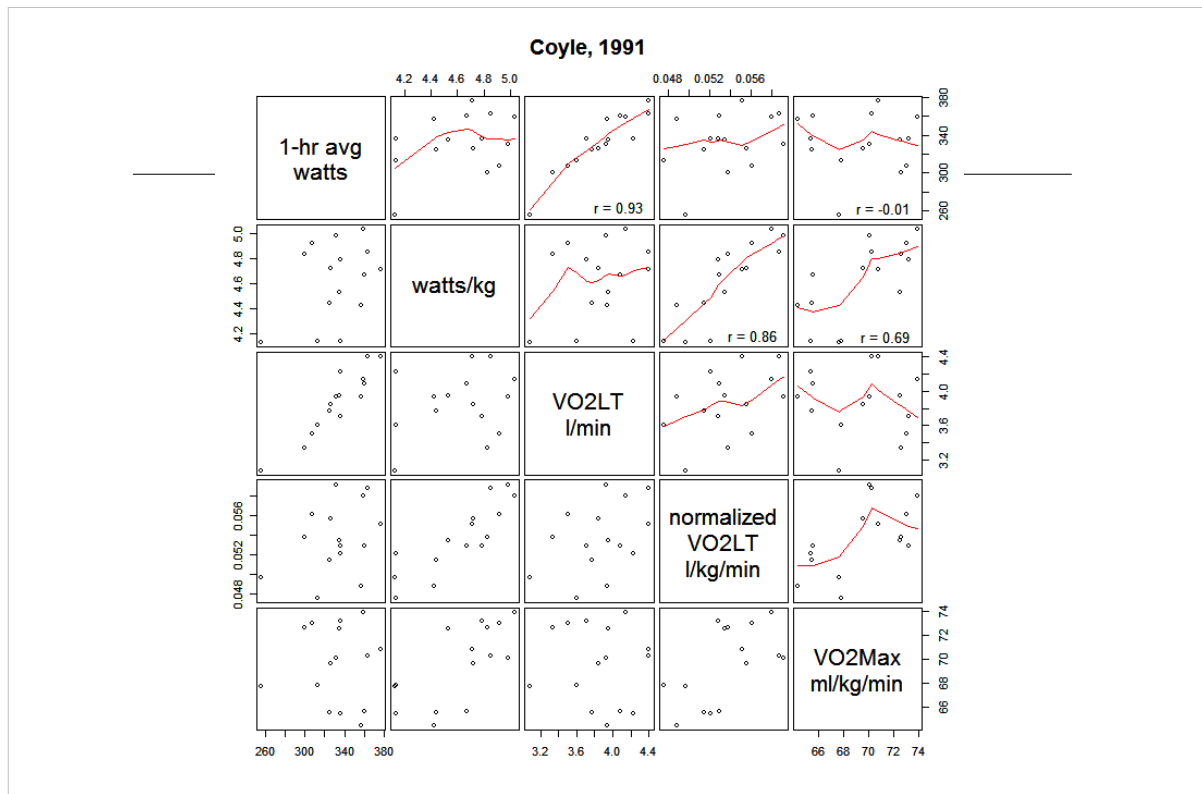
You can see the “Morris” problem even more clearly if all six of the plots had been arranged in one long column: in R,

```
library(lattice)

example(barley)
```

scatterplot matrices

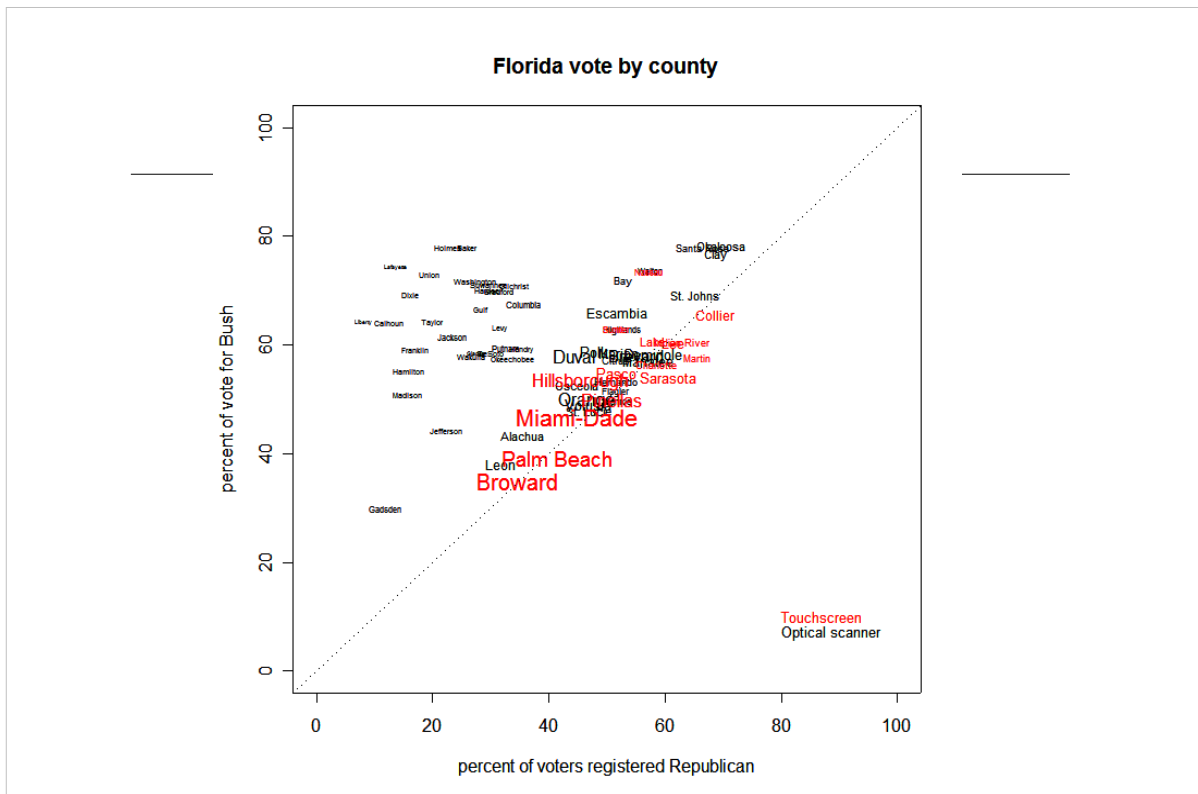
- scatterplot matrices compress a lot of information on bivariate relationships into a small space
 - useful for winnowing out uninteresting variables and deciding which variables might be worth further investigation



physiological measurements on 15 riders, elite and “state” class.

coding plotting symbols

- improves information density by tagging plotting symbols with attributes
 - you've seen this before using color or shape; can often combine with direct labeling



Here we show:

% repub

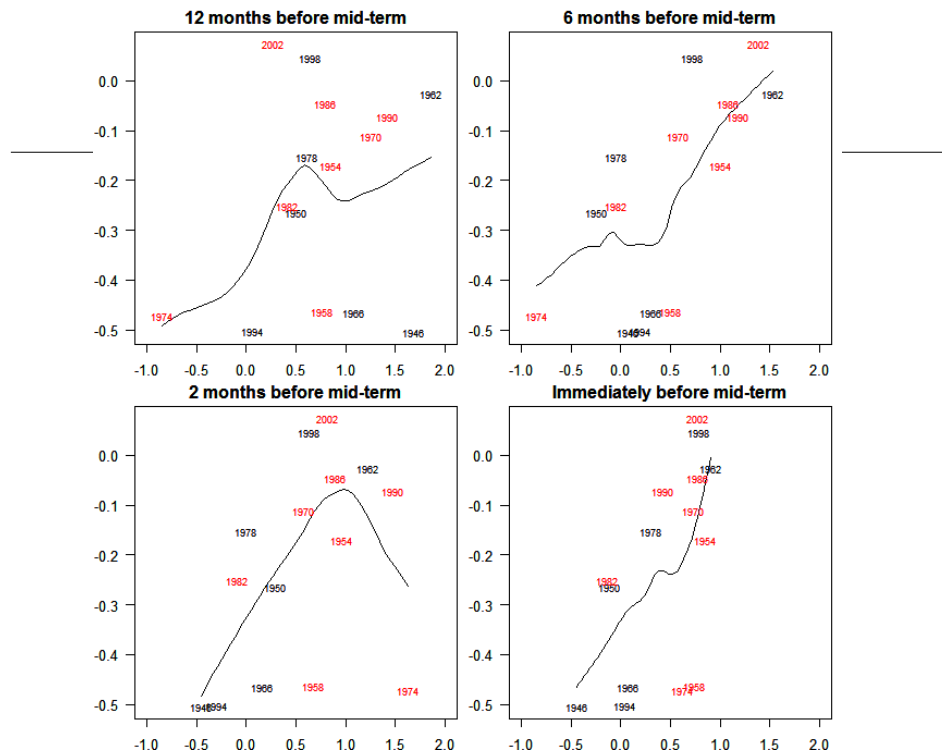
% Bush vote

name of county

type of voting machine

rough idea of county population

45-degree line to show that not many counties went more dem than their registration.



log(odds ratios) on y-axis, log(approval/disapproval) on x-axis

red means rep president, black is dem president

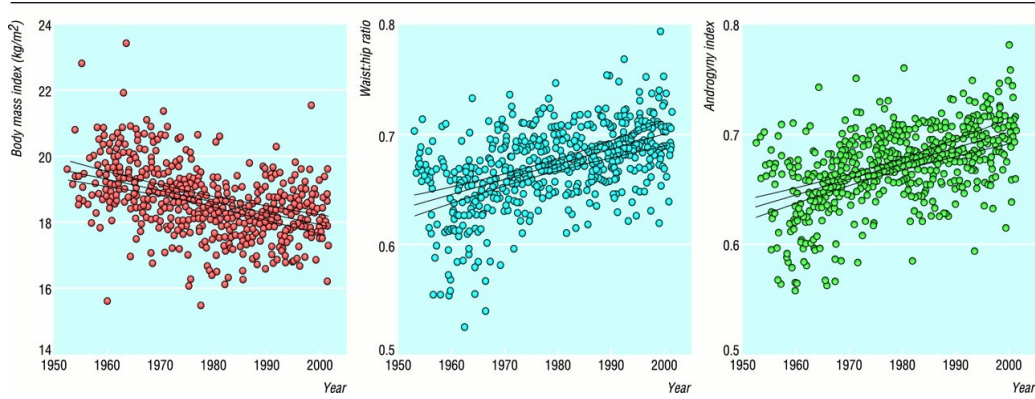
1974 and 1946 were strange years.

smoothing and straightening

- smooth lines
 - piecewise linearity
 - splines and lo(w)ess
- a ladder of re-expression
- the re-expression rule

CABG EF relationship with in-hospital mortality: piecewise linear.

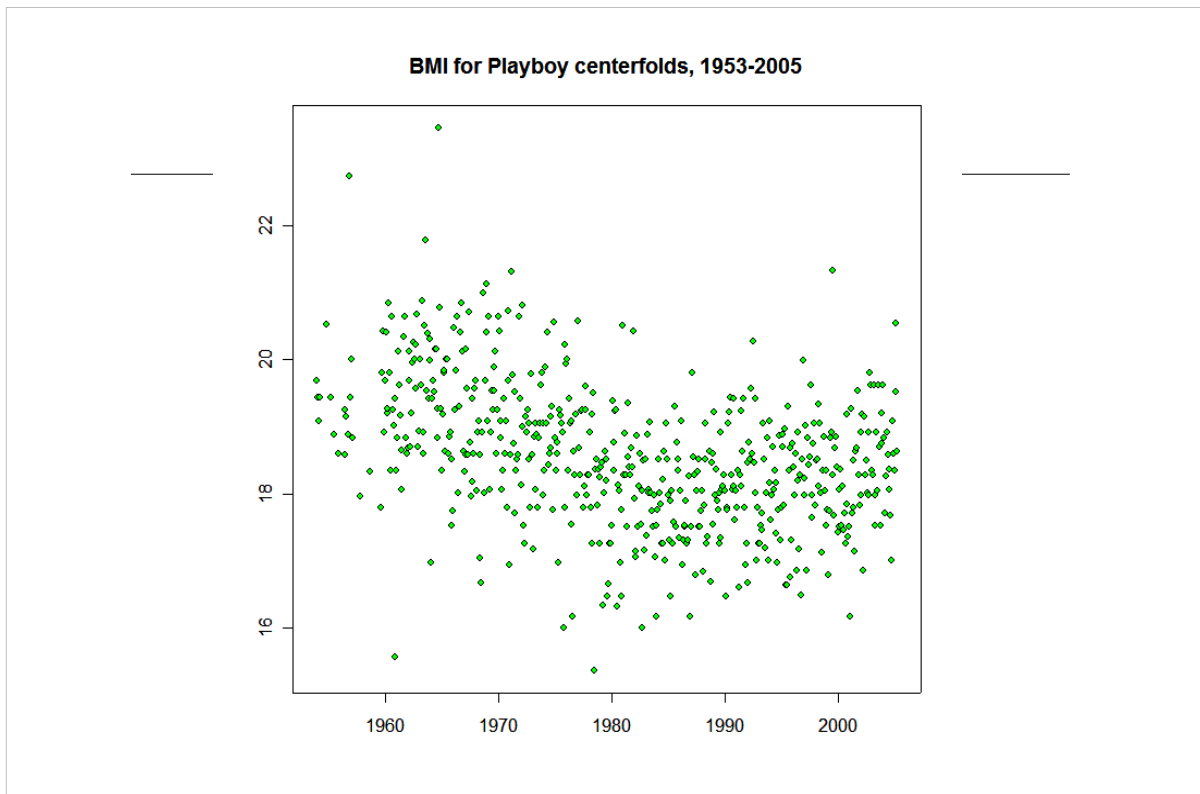
shapely centrefolds?



Voracek, M. et al. BMJ 2002;325:1447-1448

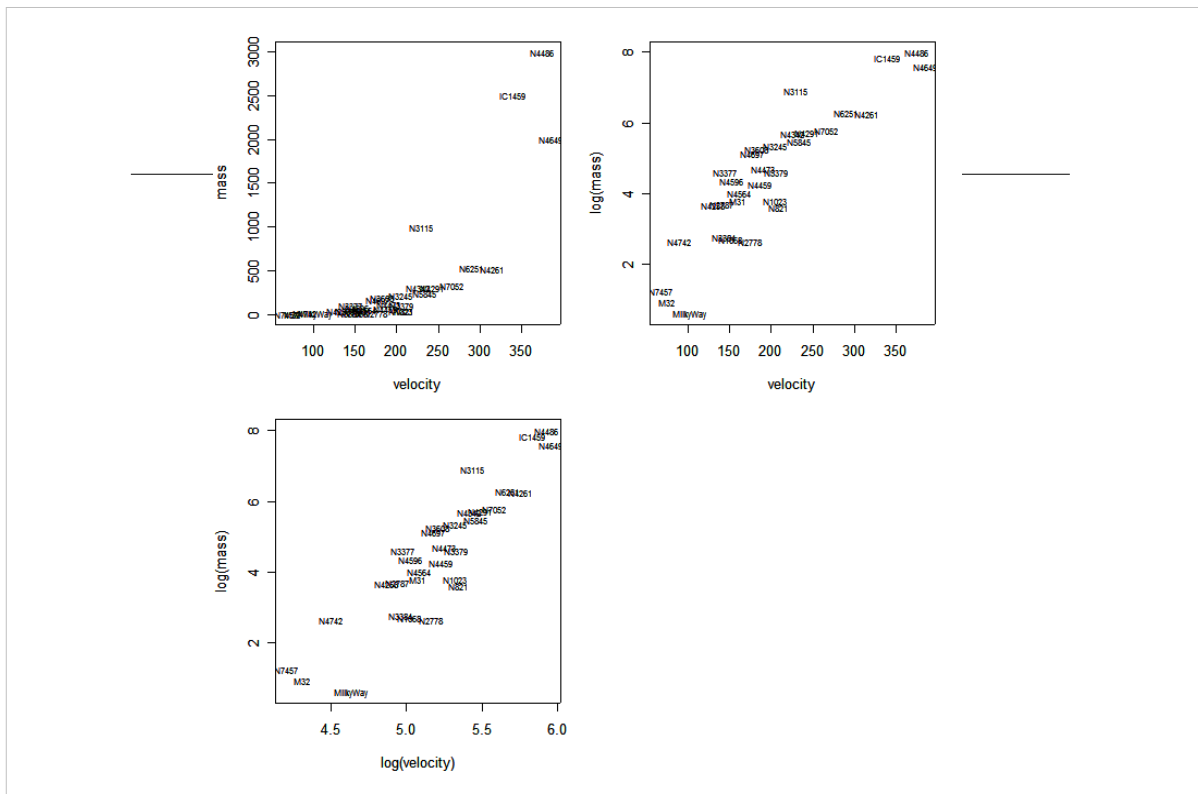
<http://bmj.bmjjournals.com/cgi/content/full/325/7378/1447>

Straight line or not?



Does that look straight to you?

I think the authors were looking for a straight line so they fit a straight line. We fit not-straight lines not only to help us guess relationships, but as a check to make sure our prejudices don't restrict our analyses.



transformation of axes in order to produce linear relationship

Black hole mass and dispersion of stellar velocities near galactic centers

Data from:

http://www.physics.ucsb.edu/~jatila/astro/astro2/b_hole_dispersion.html

a ladder of re-expressions...

3

2

1

$\frac{1}{2}$

#

$-\frac{1}{2}$

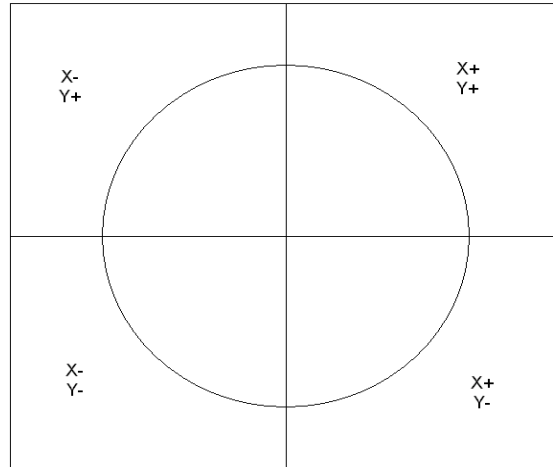
-1

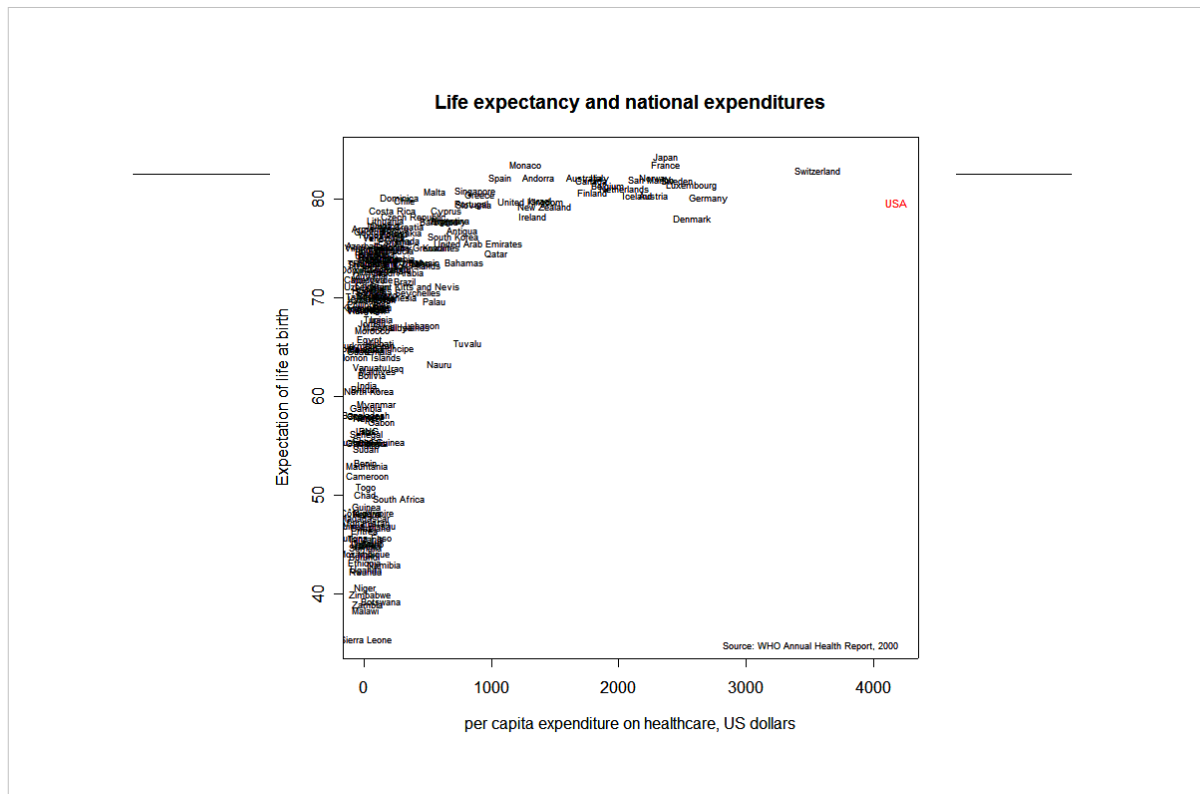
-2

-3

...and a rule for using them

Straightening by re-expression



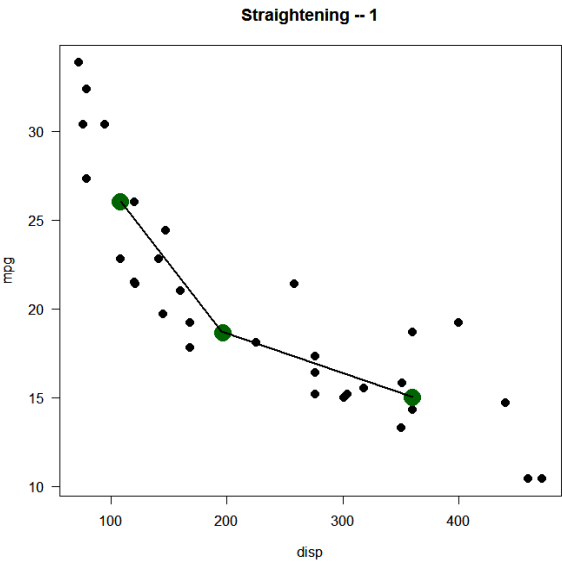


```

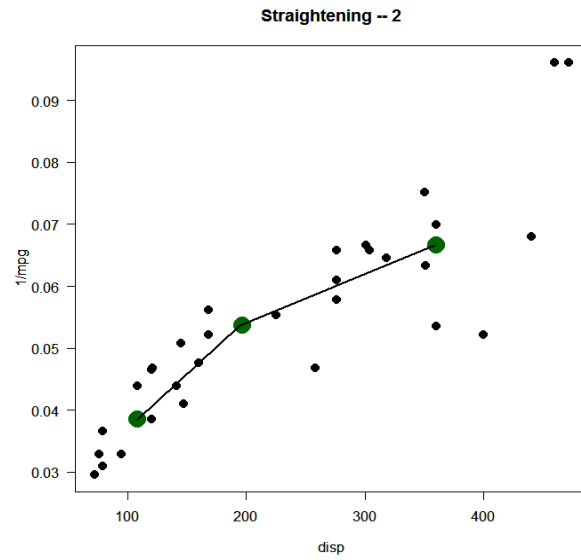
dat = read.csv("who.csv",comment="#")
head(dat)
summary(dat)
with(dat,plot(exp.percap,e0.female))
with(dat,plot(exp.percap,e0.female,log="x"))
with(dat,identify(exp.percap,e0.female,country))

```

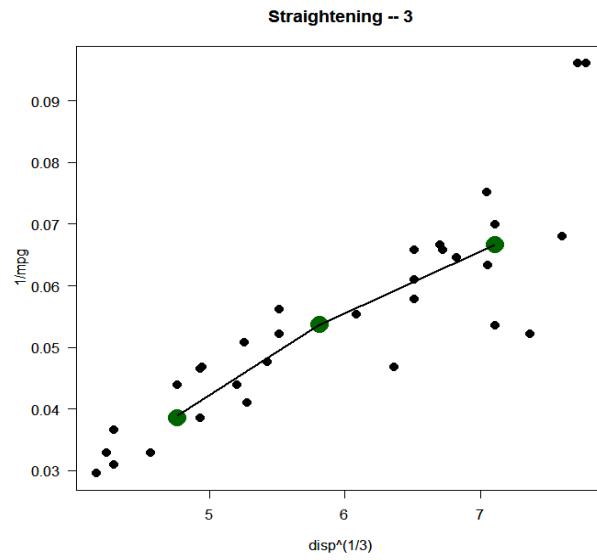
straightening



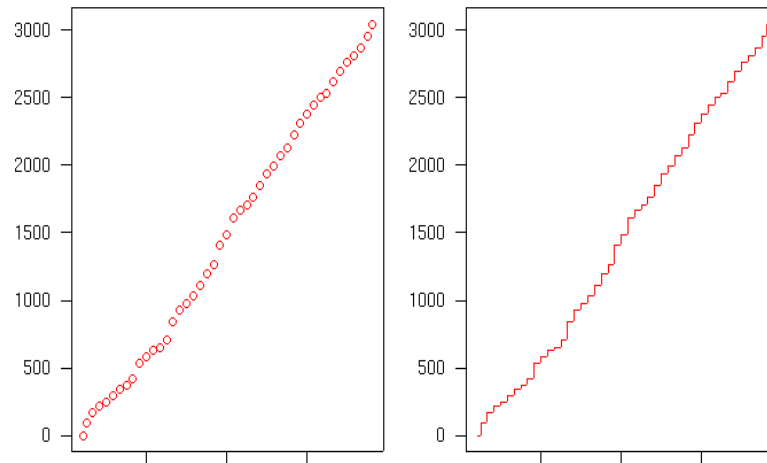
straightening 2



straightening 3



when is smooth too smooth?



Coalition deaths in Iraq, 30 Mar 2003 - 23 Oct 2006

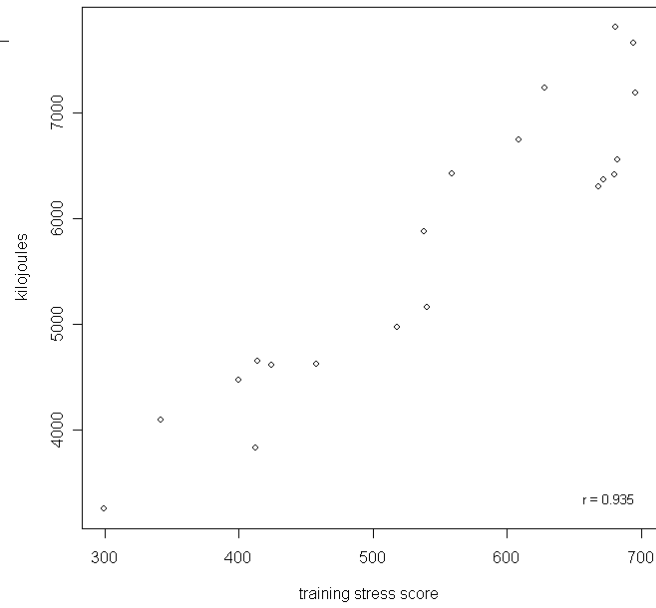
Data from <http://icasualties.org/oif/>

The eye wants to see “lines” connecting the dots, particularly when the dots are almost straight. In the right hand panel, it’s easier to force the eye not to connect the dots by using horizontal and vertical segments. It’s easier to see that a handful of months had large numbers of deaths.

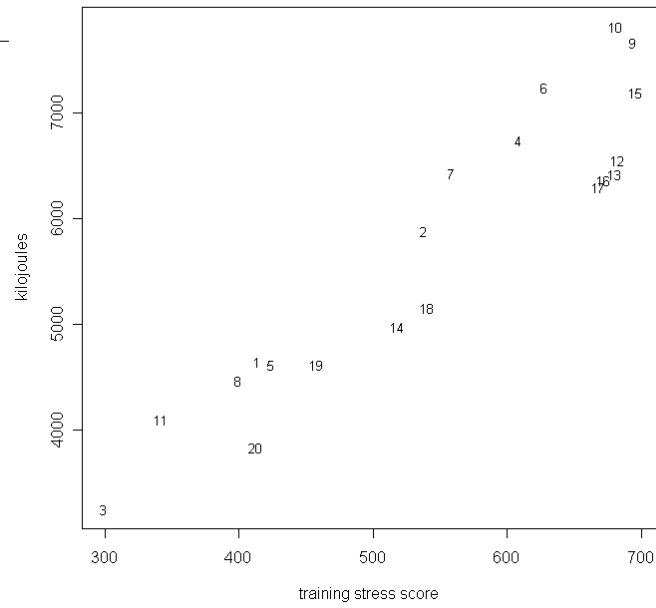
when is straight too straight?

- can straight be too straight?

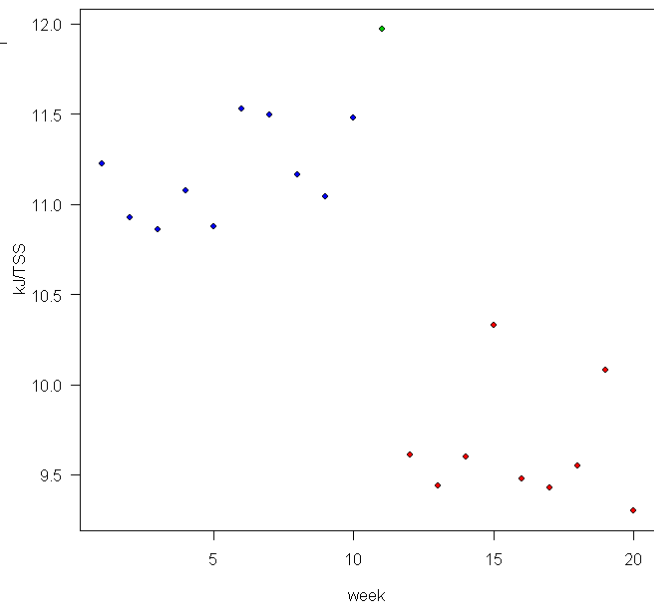
Andy's weekly training load



Andy's weekly training load



Andy's weekly training load

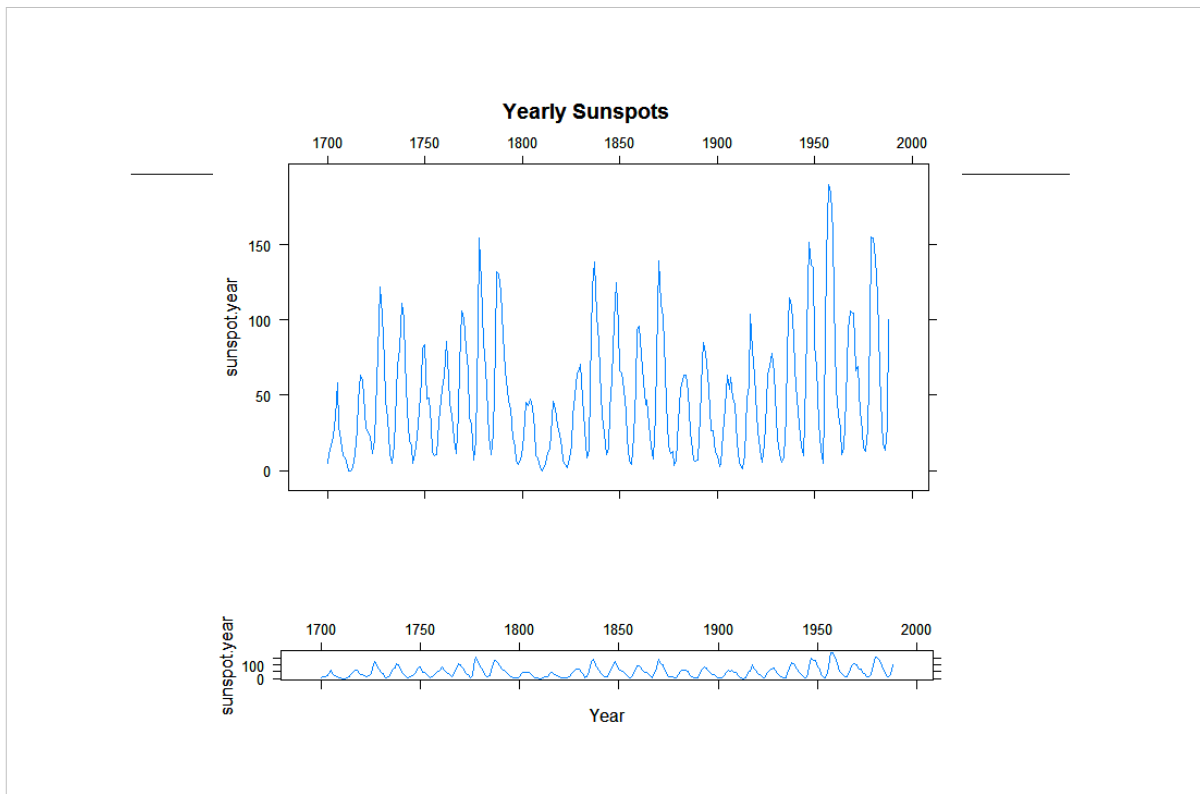


aspect ratios and banking

- human eye isn't great at decoding angles
banking helps the eye to decipher angles

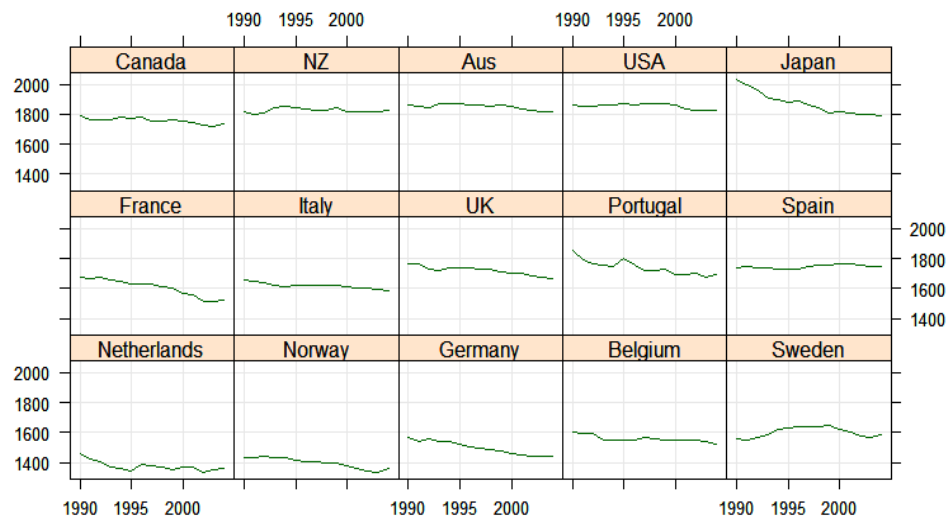
If you use base R, and you have separate graph window you can just grab the corner to change the apparent aspect ratio.

If you use Rstudio, you can write graphs to `dev.new()` and grab the corner to change the aspect ratio. When you want to close the window, use `dev.off()`



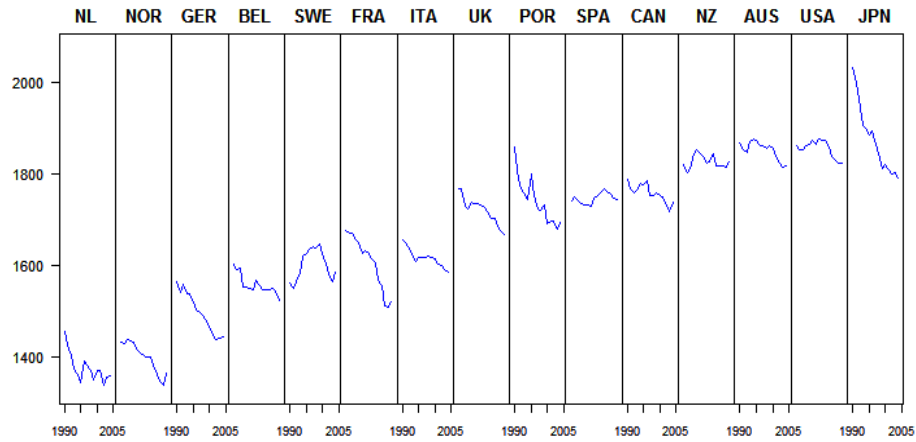
Lower panel is low-aspect ratio. Note that for the 18th and 19th C. that sunspots ramped up faster than they declined, but that in the 20th they were more symmetric. Can't see that in the top panel.

Average annual hours worked per worker, 1990-2004



Because of the previous graph, Tufte generally recommends low-aspect ratios for time series. Here and on the next page are a counter-example.

Average annual hours worked per worker, 1990-2004

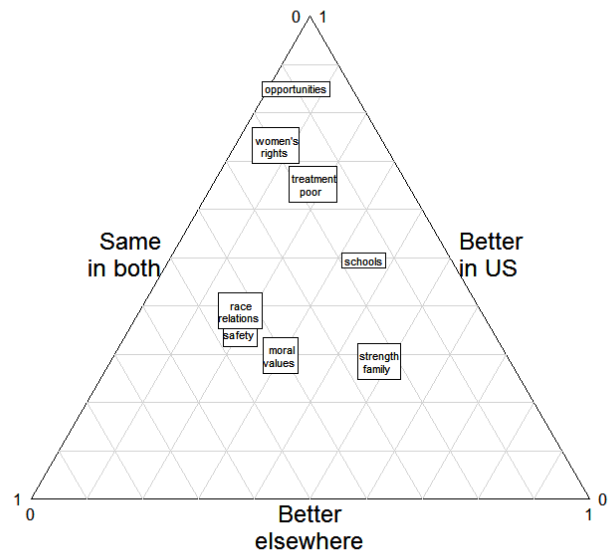


The aspect ratio is so high here that I had to abbreviate the country names, but it shows both the levels among the countries and also that Sweden's work hours went up and then came back down (sort of the same with the US).

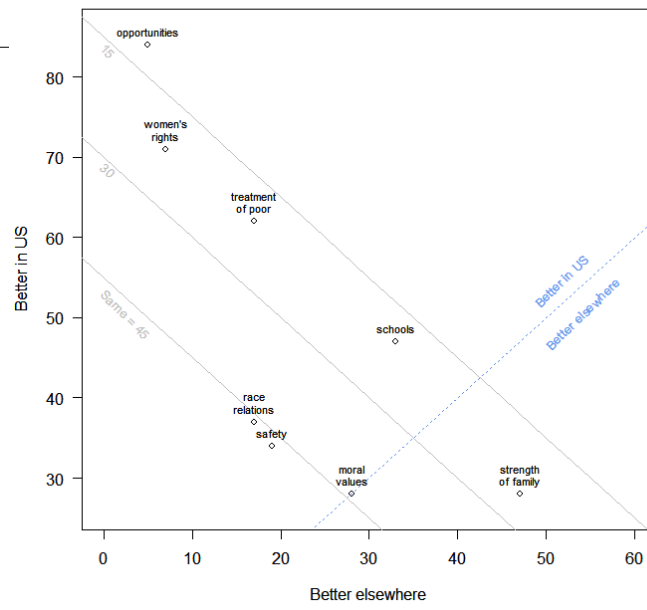
contours and bases

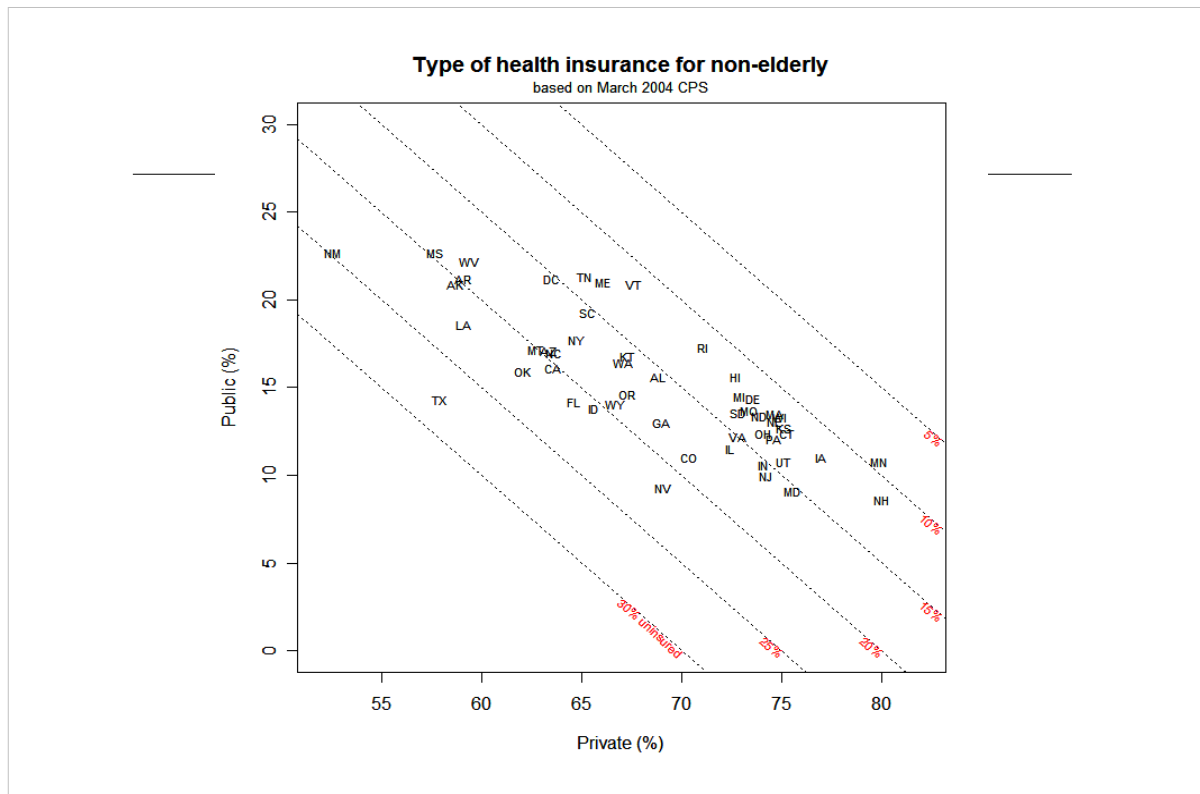
- triangle plots (= ternary plots)
soil texture plots
two degrees of freedom
- function contours can add context

Immigrant's assessments of US vs. country of origin



Immigrant's assessments of US vs. country of origin





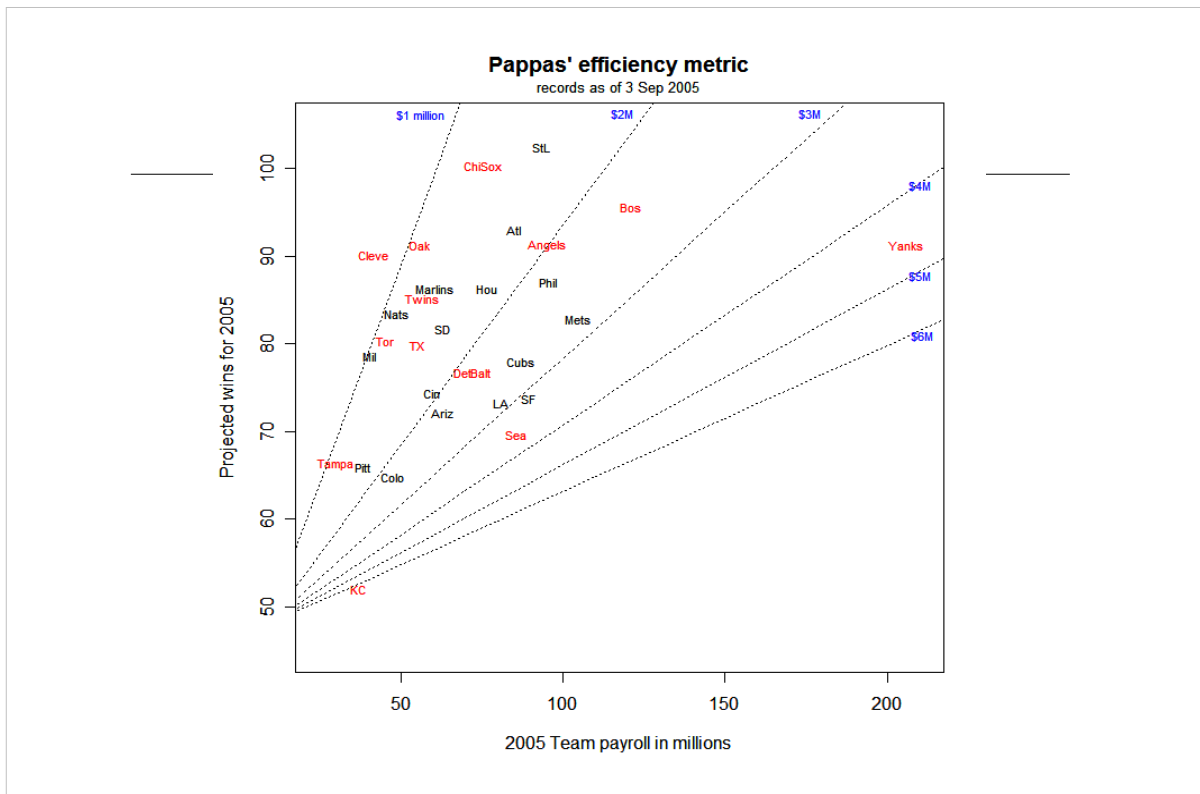
curve(...)

This is a ternary plot on normal x-y axes. Some people have problems “visualizing” a ternary plot.

Private insurance + public insurance + uninsured add to 100%, so plot any two and the third can be shown by contours counting down from upper right.

Texas has highest percentage of uninsured.

Note that across states, as private% decreases, public% increases – but not at the same rate. Cross-state, safety net isn’t uniform.

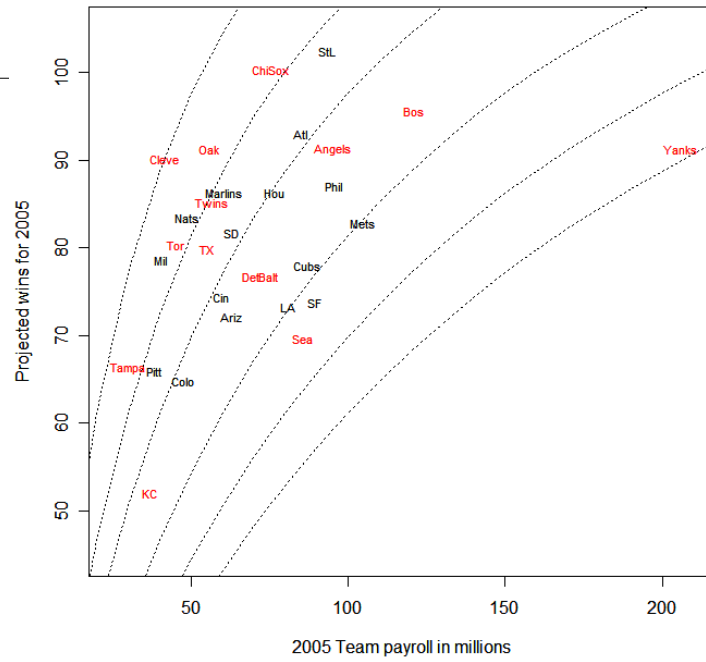


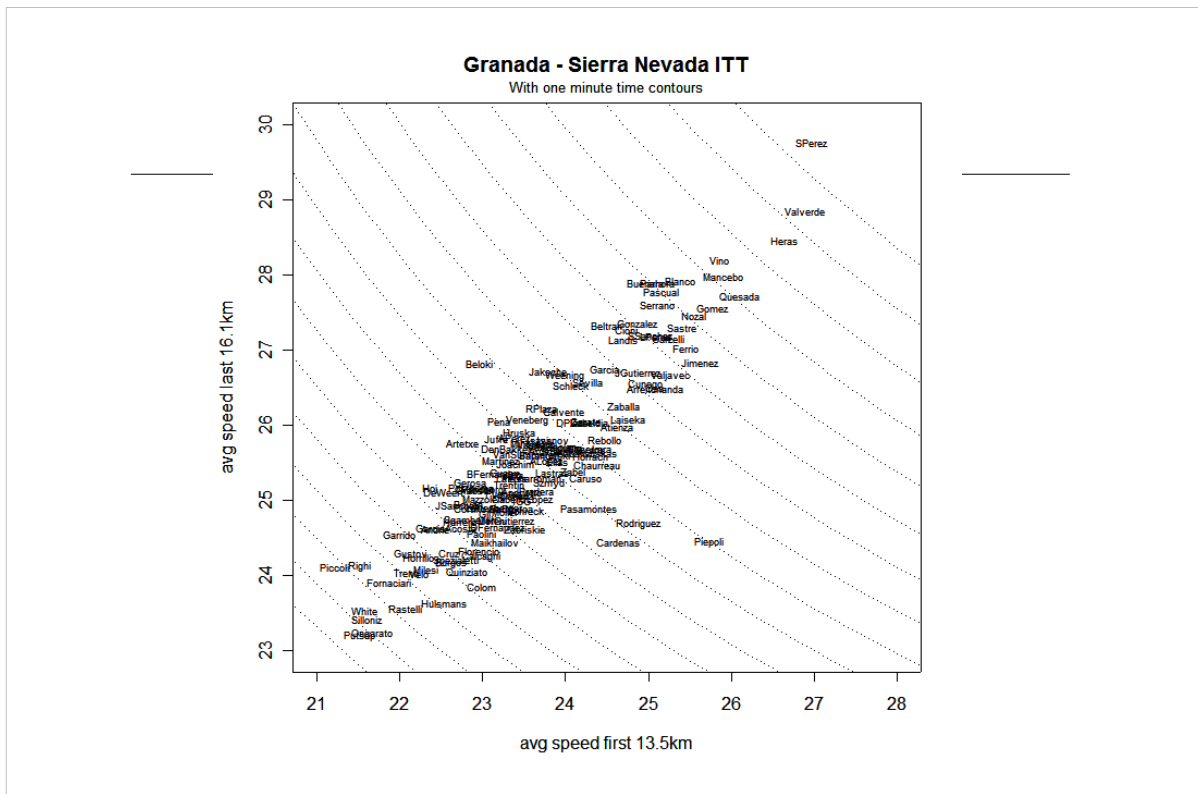
For context, read:

http://junkcharts.typepad.com/junk_charts/2005/08/baseball_roi_3_.html

Payroll/median payroll ratio to win/loss ratio

records as of 3 Sep 2005



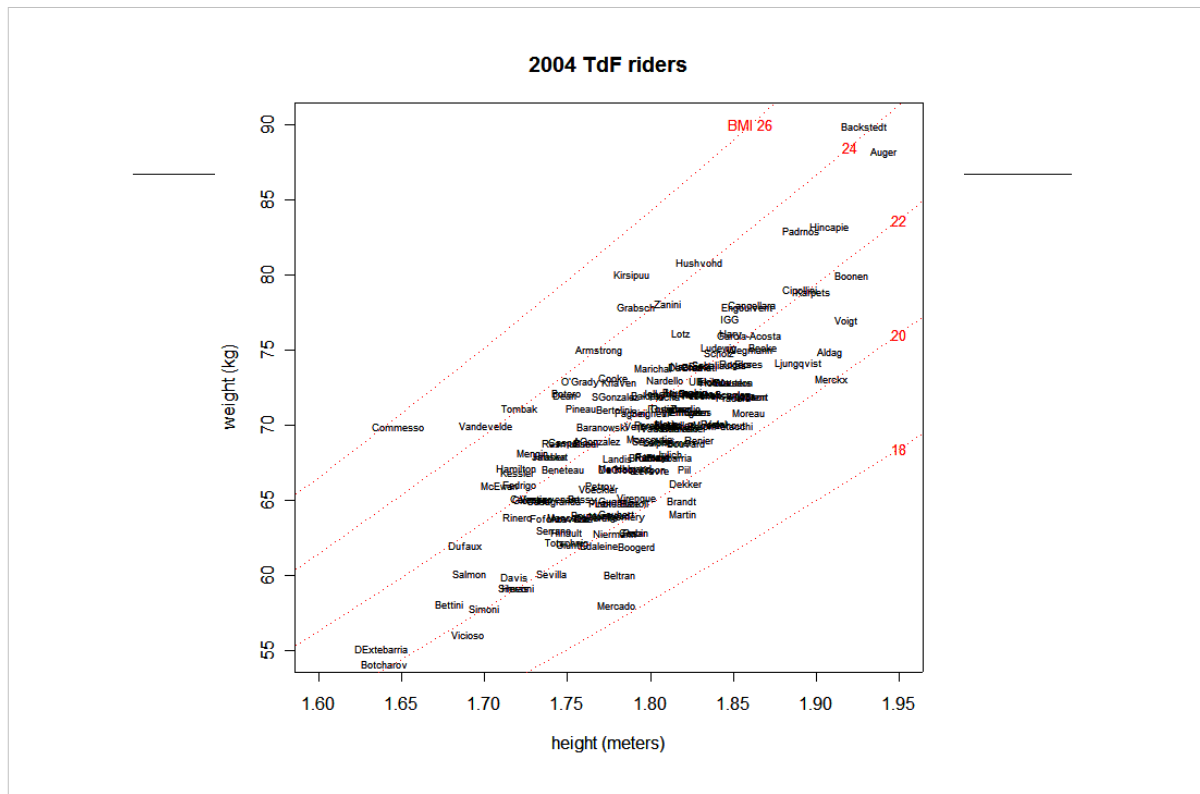


Similar to previous plot, but with one-minute contours added (I think the contours are too much info -- if I were to do this again, I'd restrict them to a smallish area around the data.)

2004 Vuelta a Espana. Stage 15

Perez was tossed out of Vuelta for using EPO

The contours extending across the page are distracting. If I re-made this thing I'd cut down on size or number.



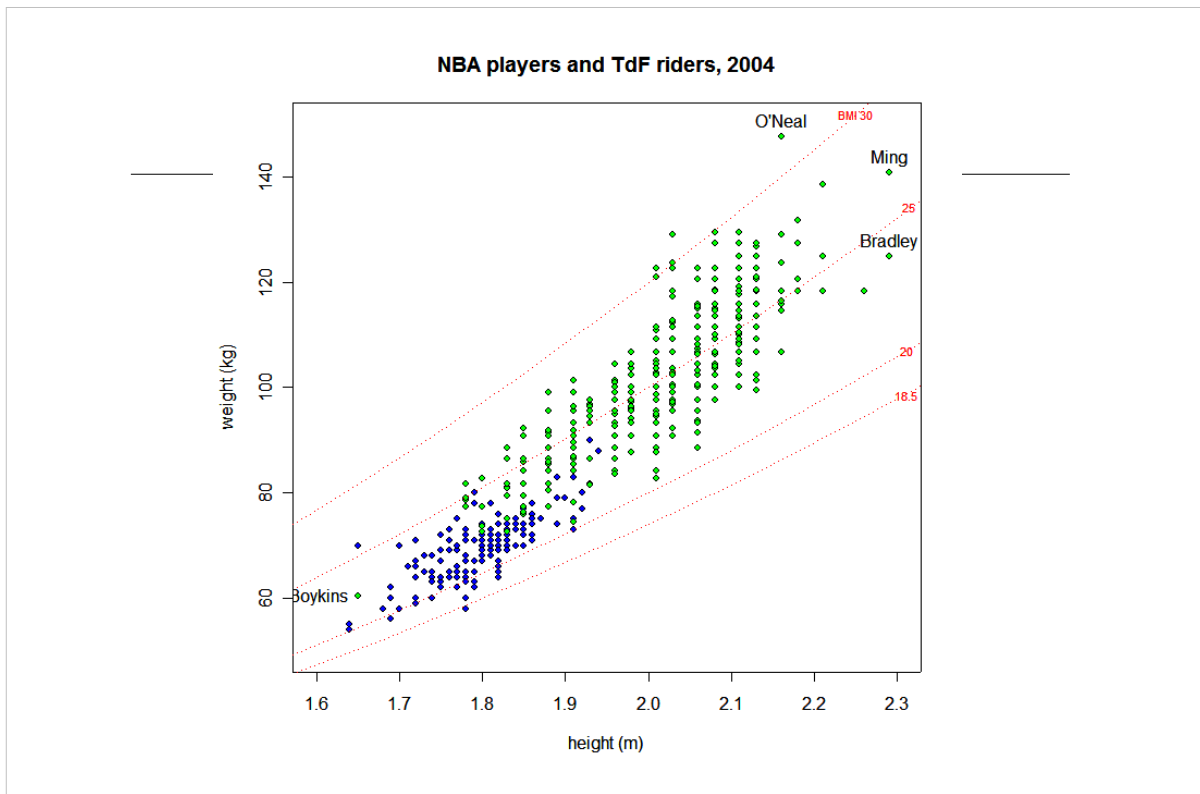
BMI = kg/m^2 . In the US, the current standard is

BMI < 18 = “underweight”

18 < BMI < 25 = “normal”

25 < BMI < 30 = “overweight”

BMI > 30 = “obese”

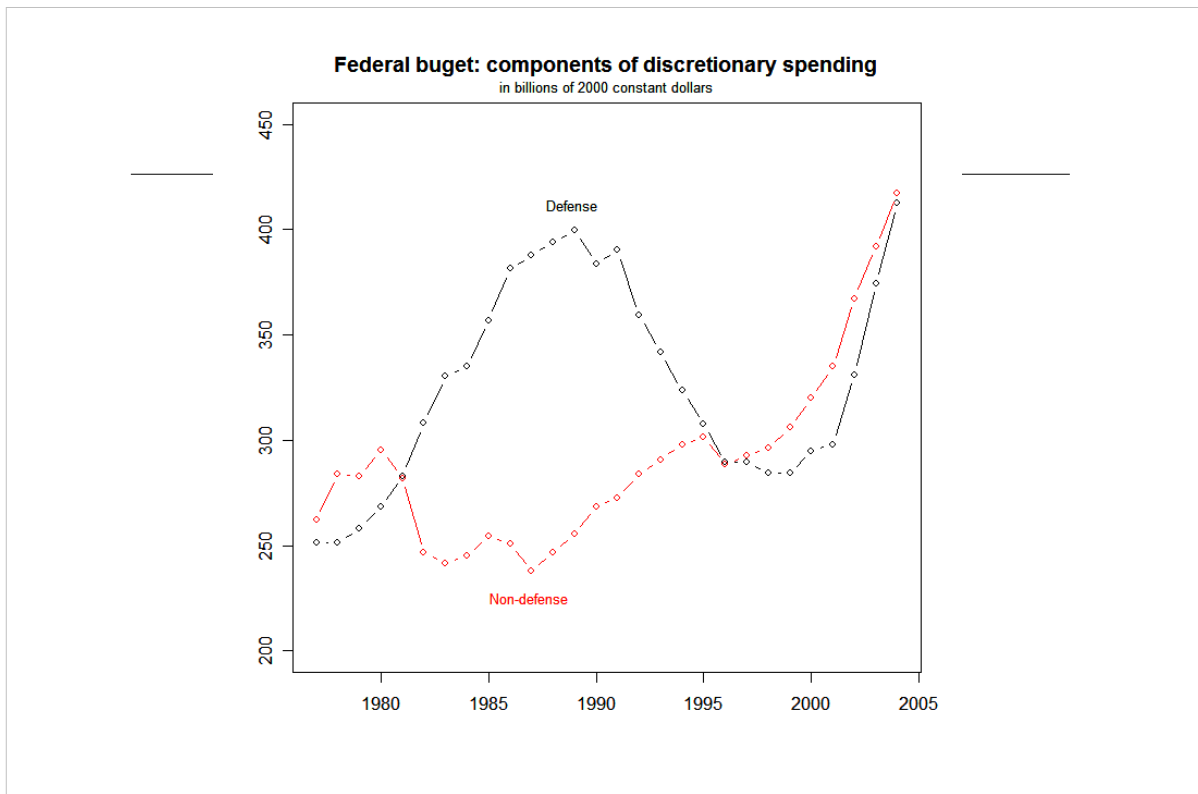


$$\text{BMI} = \text{kg}/(\text{m}^2)$$

This graph shows that BMI doesn't scale well with height -- Yao Ming is very, very tall but he doesn't appear overweight. Shaquille O'Neal isn't obese.

decomposing series into phase plots

- another version of “show the difference”



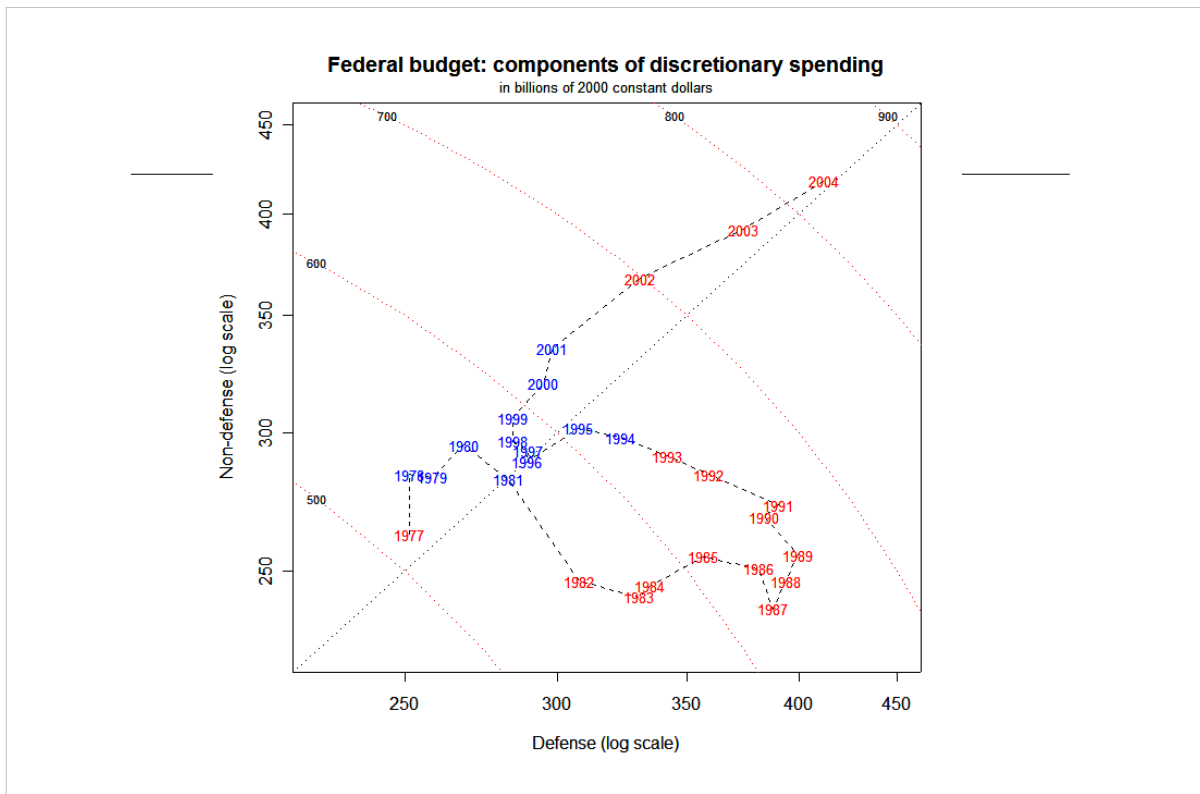
This is the conventional way to show the data – plotted against time. However, as we've already seen, it's hard for the eye to estimate distances between lines.

note: lines labeled, so no need for legend

Alternatively, you could look at percent of GDP

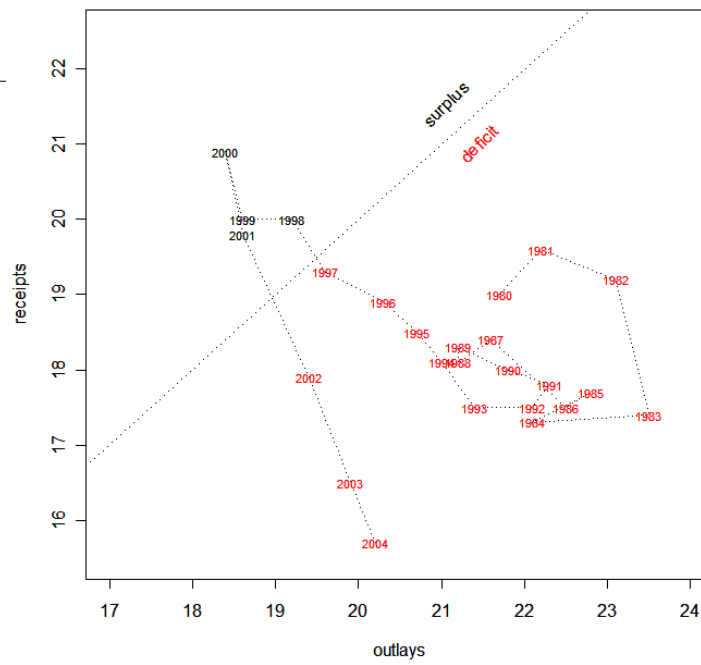
```
dat = read.table("discretionary.txt",comment="#",header=T)
head(dat)
with(dat,plot(year,defense))
with(dat,lines(year,domestic,col=2))
```

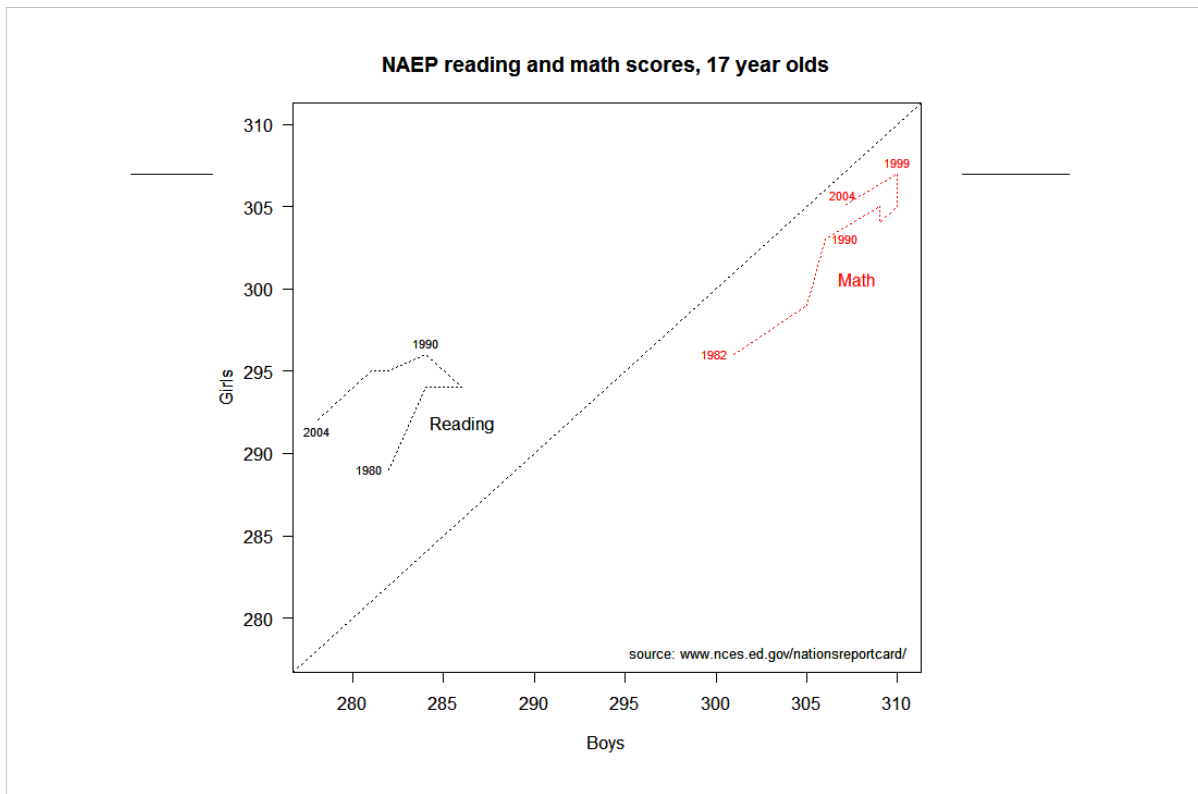
```
with(dat,plot(defense,domestic,type="n"))
with(dat,text(defense,domestic,year,cex=.7))
```



Of discretionary spending, this graph shows amount for defense, amount for non-defense, total, growth rate, year, presidential party (Red=budget proposed by a Republican president, blue=budget proposed by a Democratic president)

Federal outlays and receipts as percentage of GDP



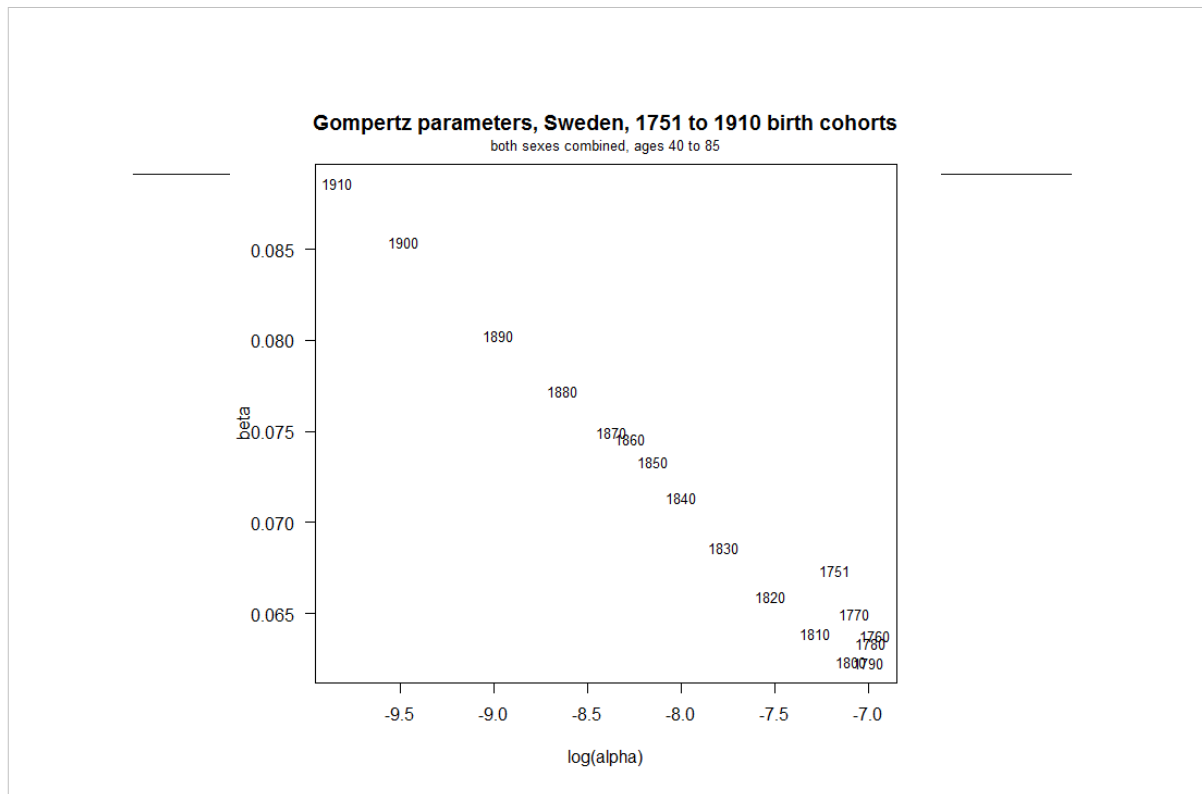


national assessment of educational progress.

You can see that math scores for boys and girls seems to move up and down together, but not so for reading scores. Also, it's easy to pick out and compare differences between boys and girls.

plot summaries for simplification

- when all subset have same contrasts, plot subset summaries
 - sometimes can get away with it even if not all subsets have all same contrasts—but then must be doubly careful
 - helps to identify patterns
 - plot and identify extremes, leave middle alone
 - this is the idea underlying “10 plus 10” plots
 - or, split into n groups (n small, like 3), and plot subsamples from each

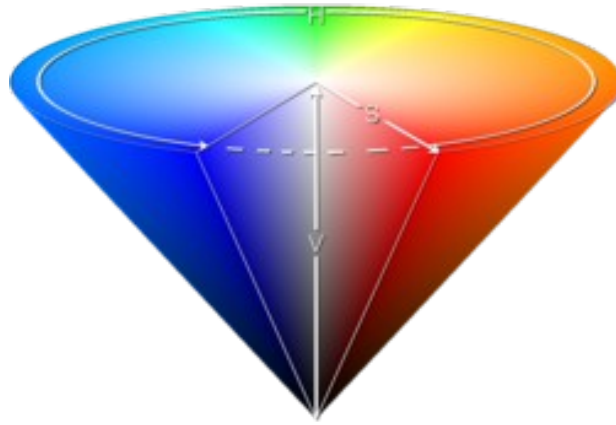


The Gompertz α and β parameters summarize a particular kind of “fit” to the hazard rates. Here’s a tip: if you’re estimating model parameters over time, or over space, or for some other kind of contrast, plot the parameters.

more on color

- HSV
 - h=hue, s=saturation, v=value
 - sometimes called HSL for hue, saturation, luminance
- equal impact colors
 - CIELUV and Munsell are systems of color perception
 - medium saturation, kind of pastel-like

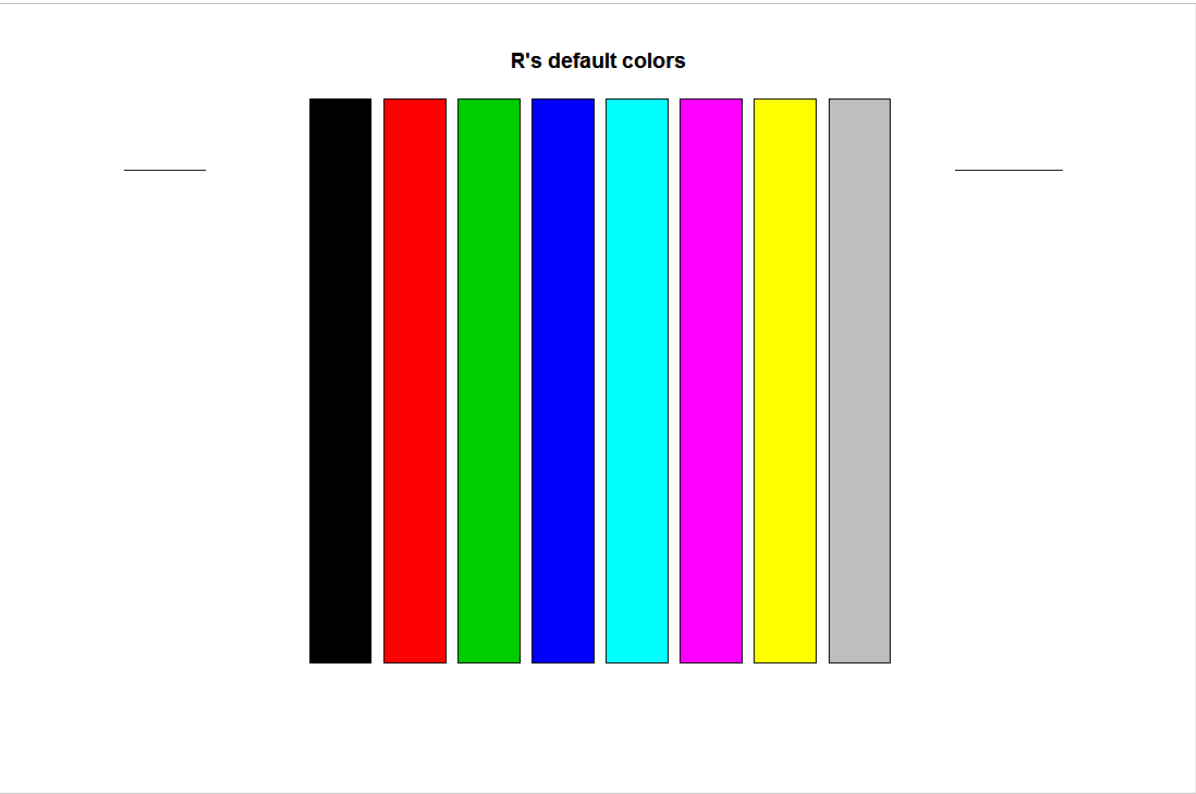
HSV cone



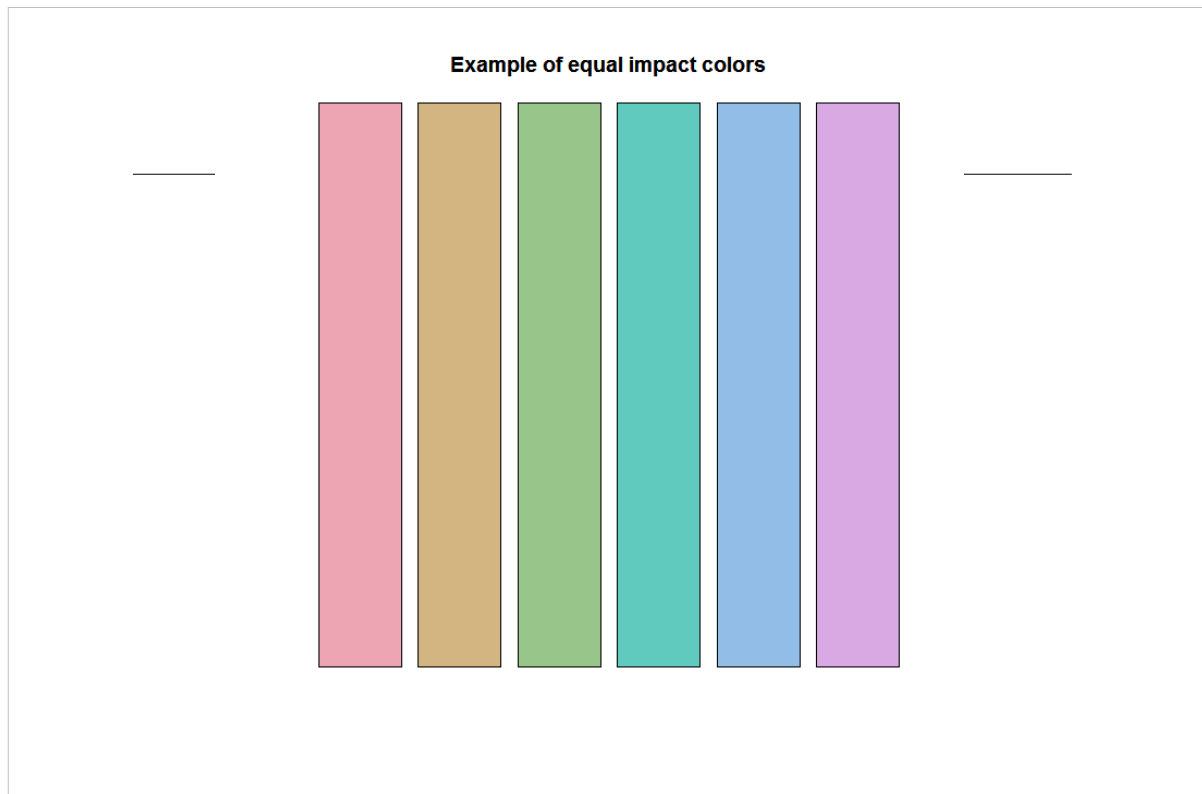
Think of the outside edge of the top of the cone as an old-fashioned color wheel. That's why hues aren't ordered – they form a circle.

As saturation goes up, you go from center to edge, so highly saturated colors appear brighter.

The human eye's ability to distinguish color depends on how much light there is. As value goes up, the luminance goes up, so the bottom of the cone is black. The top center is white because it has high value and low saturation for any particular hue.

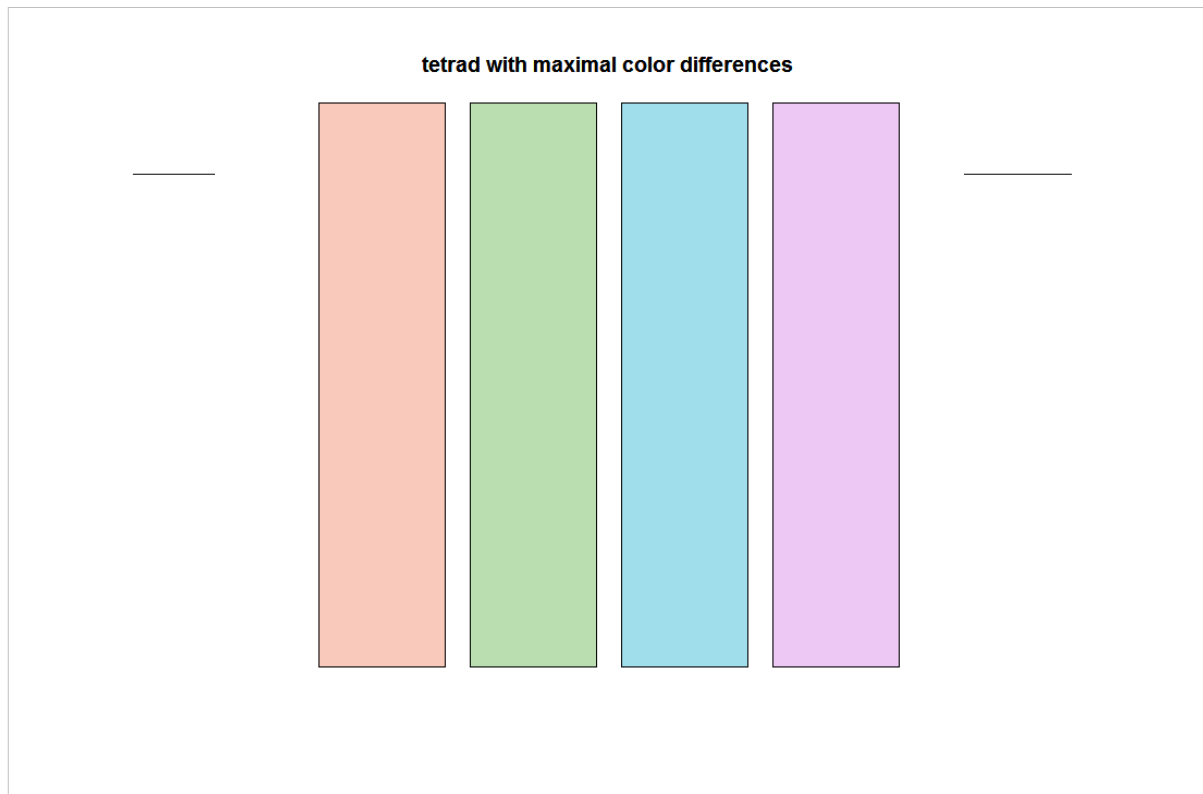


Yikes.



`hcl()` is one way to get equal impact colors. These bars were produced with:

```
barplot(rep(1,6),col=hcl(h=seq(0,300,len=6),l=75,c=45),axes=F,main="Example of equal  
impact colors")
```



```
barplot(rep(1,4),col=hcl(h=seq(30,300,len=4)),axes=F,main="tetrad with maximal color  
differences")
```


basic techniques

- show the difference
- identify outliers (or, label directly)
- group and order
- plot extremes
- multiple comparisons

slightly more advanced techniques

- smoothing
- straightening
- phase plots
- contours
- banking
- coloring

stuff I wanted to hide until the end

- friends don't let friends graph with Excel
but let's be realistic: sometimes you have no choice
dates in Excel are particularly a problem

nature

[Explore content](#) ▾ [About the journal](#) ▾ [Publish with us](#) ▾ [Subscribe](#)

[nature](#) > [news](#) > article

NEWS | 13 August 2021 | Correction [25 August 2021](#)

Autocorrect errors in Excel still creating genomics headache

Despite geneticists being warned about spreadsheet problems, 30% of published papers contain mangled gene names in supplementary data.

By [Dyani Lewis](#)



Science

One in five genetics papers contains errors thanks to Microsoft Excel

29 AUG 2016 • BY [JESSICA BODDY](#)

how a demographer changed
bicycle racing and design