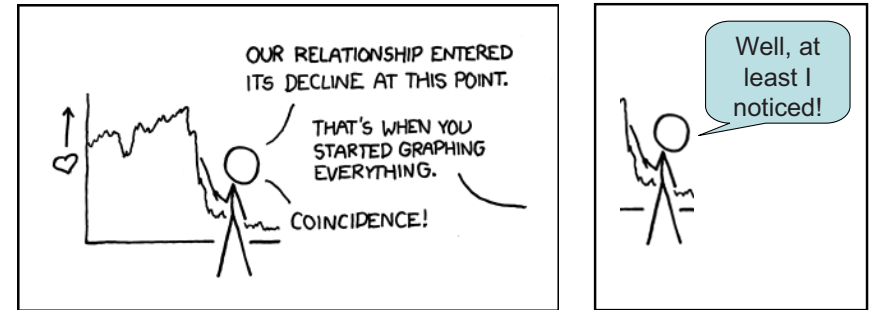


Graphical Methods for Data Analysis & Multivariate Statistics

Michael Friendly
Psychology 6140

Why plot your data?

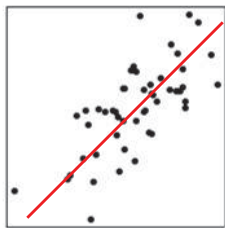
Graphs help us to see patterns, trends, and other features not otherwise easily apparent from numerical summaries.



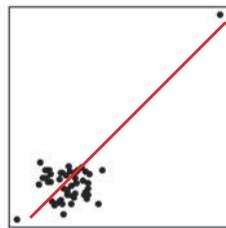
Why plot your data?

Three data sets with exactly the same bivariate summary statistics:

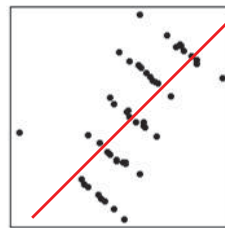
- Same correlations, linear regression lines, etc
- Indistinguishable from standard printed output



Standard data



$r=0$ with 2 outliers

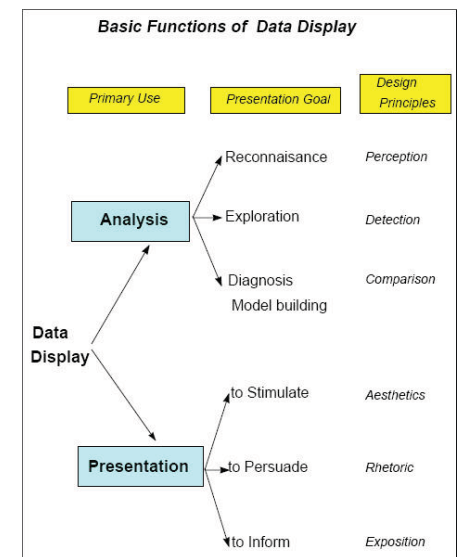


Lurking variable?

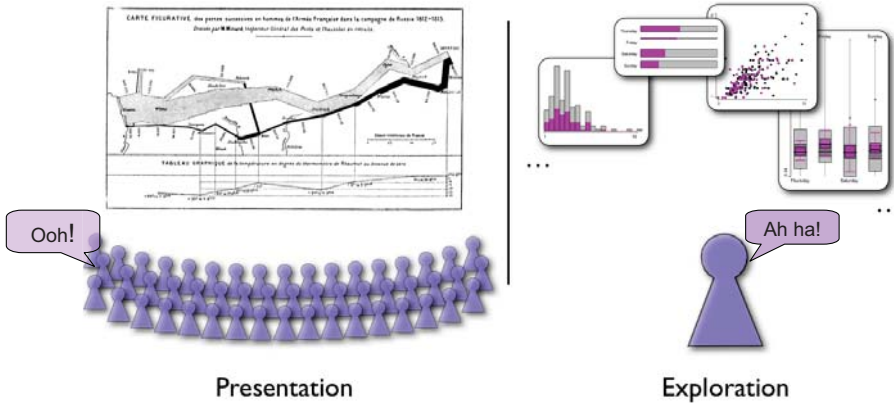
Different graphs for different purposes

Graphs (& tables) as communication:

- What audience?
- What message?
- **Analysis graphs:** design to see patterns, trends, aid the process of data description, interpretation
- **Presentation graphs:** design to make a point, illustrate a conclusion



Different graphs for different purposes

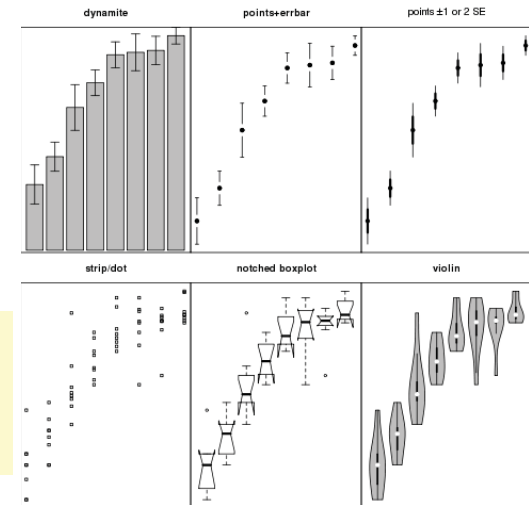


Presentation graphs: single image for a large audience
Exploratory graphs: many images for a narrow audience (you!)

Comparing groups

Six different graphs for comparing groups in a one-way design

- which group means differ?
- equal variability?
- distribution shape?
- what do error bars mean?



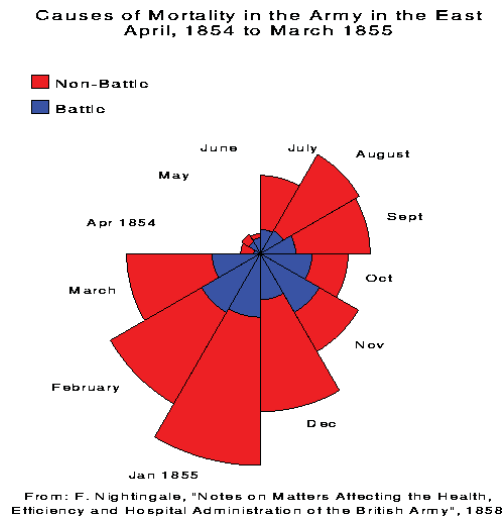
Never use dynamite plots
 Always explain what error bars mean
 Consider tradeoff between summarization & exposure

Presentation graph: Nightingale's coxcomb

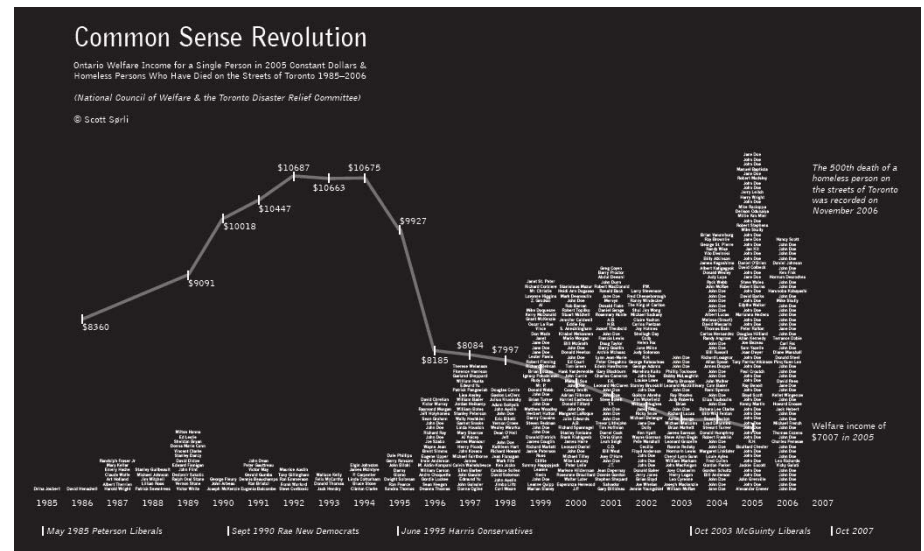
Florence Nightingale: Deaths in the Crimean war from battle vs. other causes (disease, wounds)

She used this to argue for better field hospitals (MASH units)

The best presentation graphs pass the **Interocular Traumatic Test:**
 The message hits you between the eyes!



Rhetorical graph: Common Sense Revolution

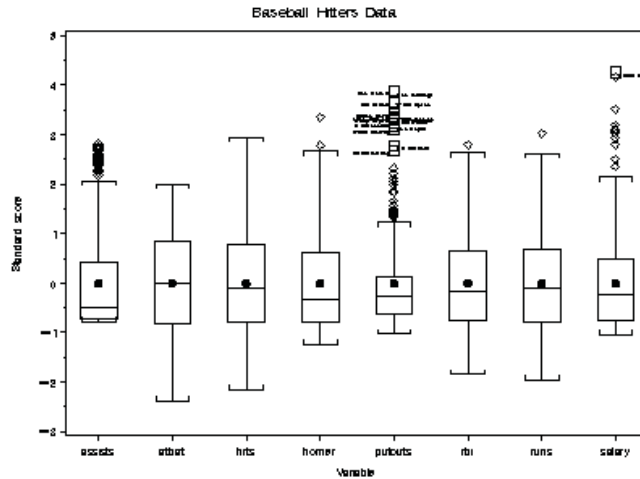


Analysis graph: Screening

Side-by-side boxplots of variables in the baseball data show the shapes of distributions --- aid to transformation

- Each variable is standardized to allow comparison.
- Plot is produced by **datachk** macro.

See: <http://datavis.ca/sas/mac/datachk.html>

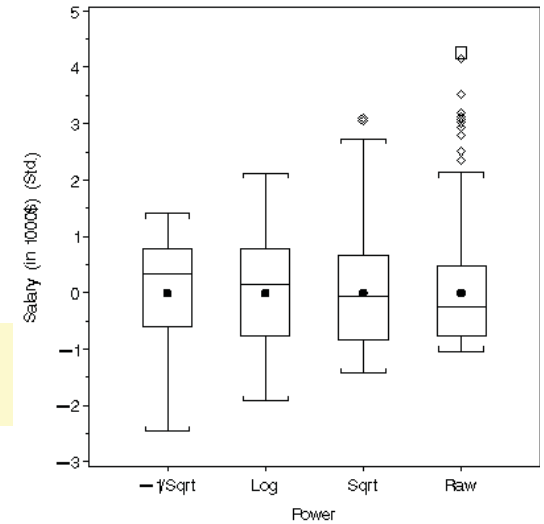


Exploratory graphs: Transformations

Data often needs to be transformed to meet analysis assumptions:

- Symmetry (~ Normality)
- Linear relations
- Constant variance

For symmetry, a **symbox** plot shows a variable transformed to various powers.



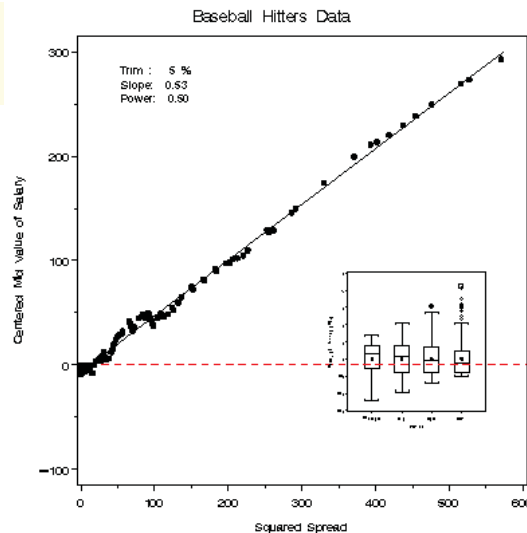
Diagnostic graphs: Transformations

Diagnostic plots can be used to suggest corrective action, often by a power transformation: $y \rightarrow y^p$

Symmetry transformation plot:

- Constructed so symmetric data plots as horizontal line
- Slope (b) of data line \rightarrow power: $p = 1 - b \rightarrow y^p = y^{(1-b)}$

Other diagnostic plots use the same idea: slope (b) $\rightarrow y^{(1-b)}$

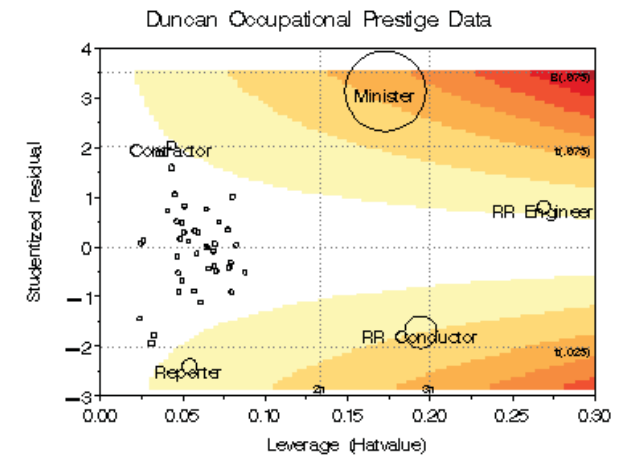


Model diagnosis: Influence in regression

Multiple regression model: prestige ~ income + education

Influence plots can show:

- model residual
- leverage (potential impact)
- influence ~ residual x leverage (Cook D statistic)
- contour map of influence



Model diagnosis: regression quartet

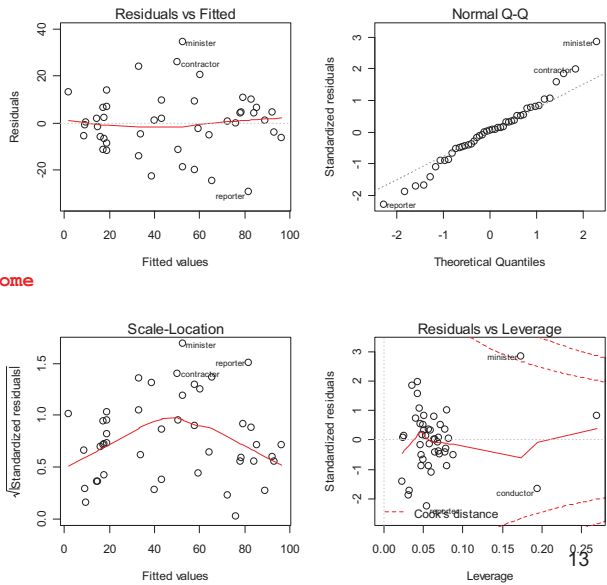
Statistical software should make it easy to get informative diagnostic plots

In R, plotting a `lm` model object → the “regression quartet” of plots

```
> model <- lm(prestige ~ income + education)
```

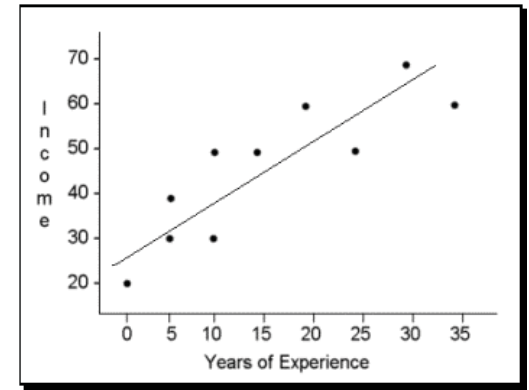
```
> plot(model)
```

(SAS has similar, using ODS graphics)



Scatterplots: A basic workhorse for quantitative data

- Show the relation between two Q variables (ignoring all others!)
- More useful when enhanced to show visual summaries
- Vary point color/shape to show strata/groups
- Combine in multi-panel displays to show more
 - Scatter plot matrix: all pairs
 - Conditional relations: Y vs. X stratified by Group

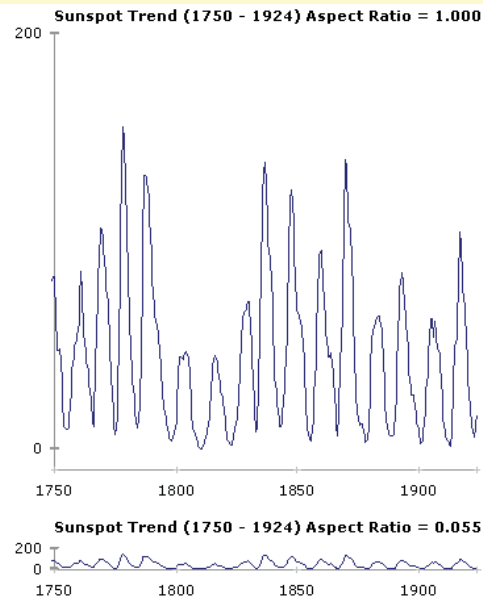


Scatterplots: Scales matter

Computer plots are usually generated with a given aspect ratio, to conform to the page or screen.

A better idea is to scale the plot so that slopes of lines or curves average ~ 45 degrees.

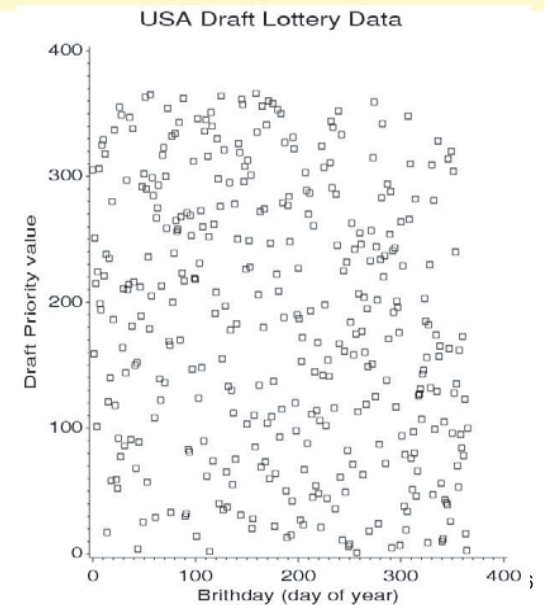
In the rescaled version, we can see that, within each cycle, sunspots tend to increase more quickly than they decline.



Scatterplots: Annotations enhance perception

Data from the US draft lottery, 1970

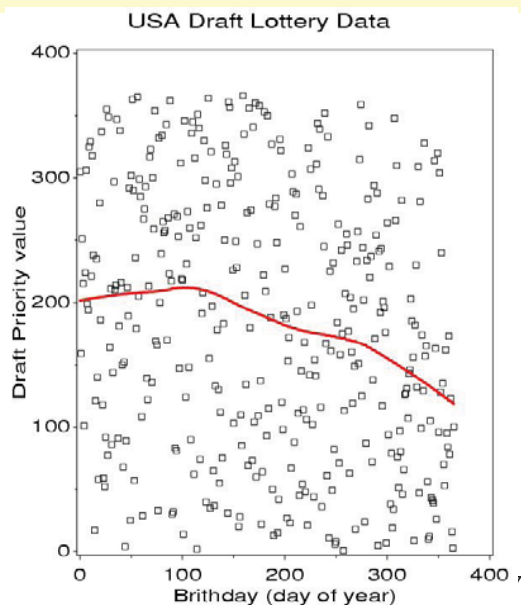
- Birth dates were drawn at random to assign a “draft priority value” (1=bad)
- Can you see any pattern or trend?



Scatterplots: Annotations enhance perception

Drawing a smooth curve shows a systematic decrease toward the end of the year.

- The smooth curve is fit by **loess**, a form of non-parametric regression.

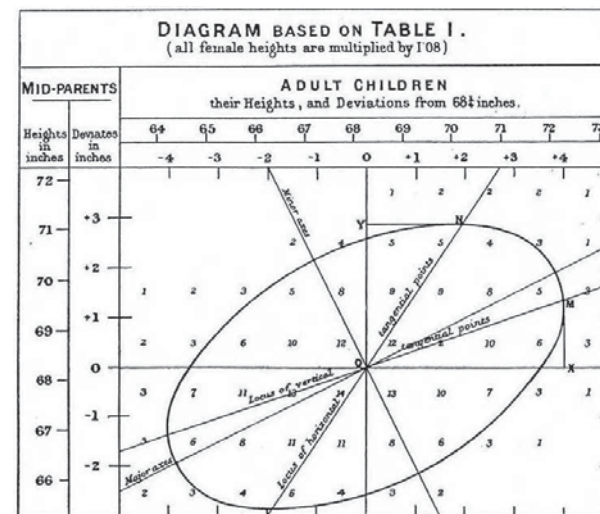


Scatterplots: Data ellipses

Galton's (1886) semi-graphic table, showing relation of mid-parent's height to children's height.

As shown:

- Contours of equal frequency formed ellipses
- Regression lines of Y on X and X on Y are the loci of vertical and horizontal tangents
- Major/minor axes are the principal components

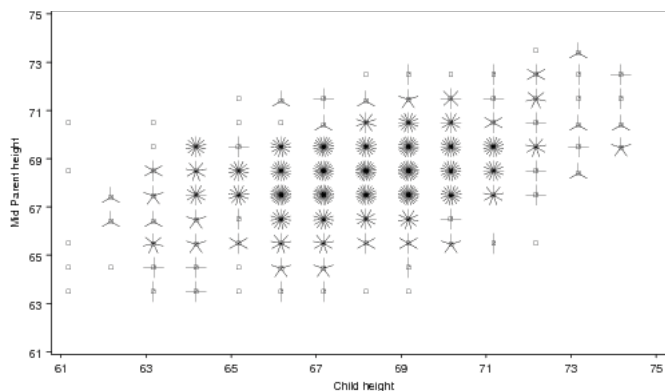
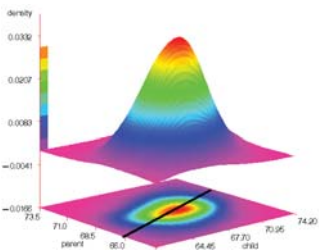


18

Scatterplots: Data ellipses

Galton's data on child & mid-parent heights, shown as a sunflower plot: each sunflower symbol shows the number of observations in the (x, y) cell.

2D density estimate of bivariate surface



19

Scatterplots: Data ellipses

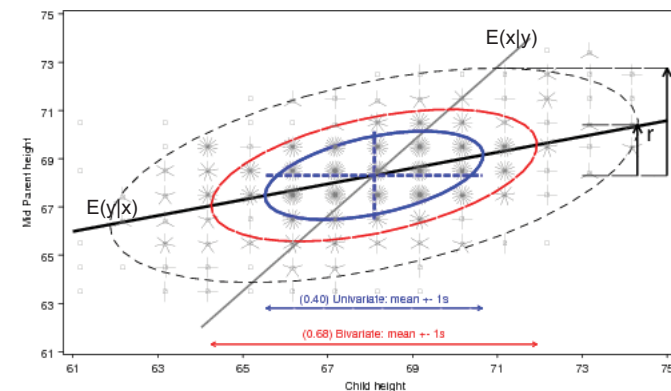
Any scatterplot can be summarized by data ellipses (assuming normality). These show: means, standard deviations, and allow correlations & regression lines to be visually estimated.

Data ellipse:

$$D^2(y) \approx \chi_p^2(1 - \alpha)$$

Galton data, 40%, 68% & 95% data ellipses. Sizes are:

- $X^2(0.40) = 1.0$
- $X^2(0.68) = 2.28$
- $X^2(0.95) = 6.0$



20

Visualizing multivariate data

Showing relations among 3 or more variables:

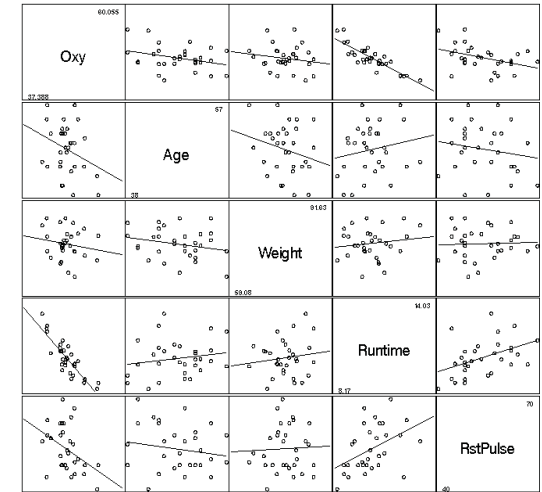
- Scatter plot matrices (enhance with visual summaries, thin for many variables)
- Conditional plots: $Y \sim X \mid (Z, \text{Group})$
- Seeing multivariate profiles, clusters:
 - Star plots, face plots, parallel coordinates
- Biplots: project data into low-D view

Scatter plot matrix

- **Fitness data:**

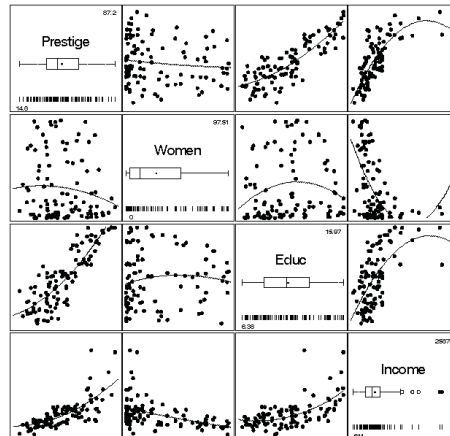
$\text{Oxy} \sim \text{Age} + \text{Weight} + \text{Runtime} + \text{Rstpulse}$

- Each panel shows row var vs. col var
- Reg line shows *linear* relation



Scatter plot matrix

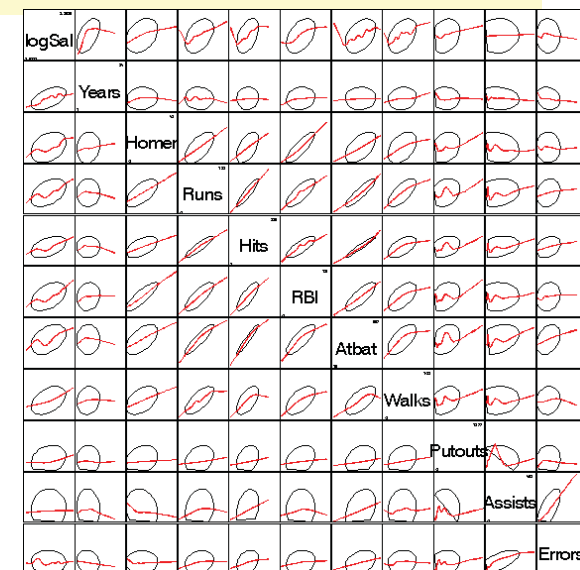
- **Occ. prestige:**
 $\text{Prestige} \sim \% \text{women} + \text{Educ} + \text{Income}$
- Box, rug plots show univar. distributions
- Quadratic regressions show linear/non-linear relations (loess would be better)



Larger data sets: Visual thinning

Baseball data: $\log(\text{Salary}) \sim$ performance variables

- Too much data to show individual points
- Each scatterplot is summarized by a loess smoothed curve and a data ellipse

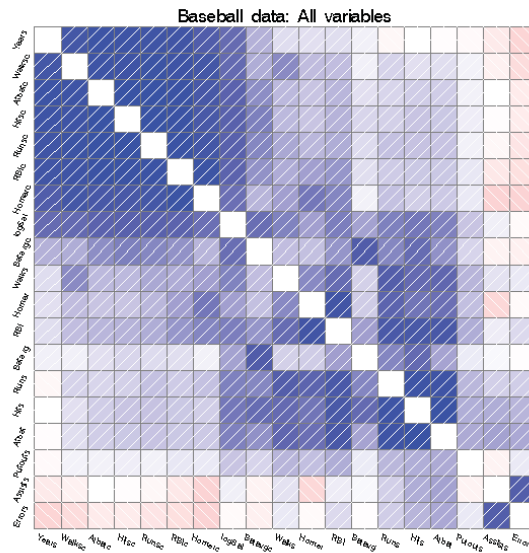


Larger data sets: Corrgrams

Correlation diagram shows **pattern** of correlations for many variables.

Variables are re-ordered to make the groupings most visually apparent.

This graphic assumes that all relations are linear, not necessarily always true

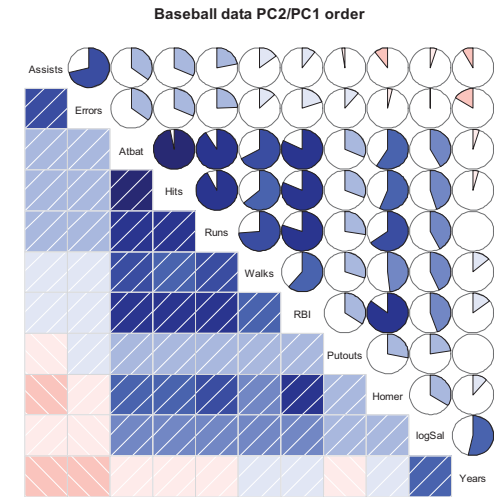


Graph using SAS **corrgram** macro, <http://datavis.ca/sasmacro/corrgram.html>

Corrgrams: Different renderings

The value of a correlation may be rendered in different ways, with different visual impact.

- Shading levels: help detect similar values
- Pie symbols: make it easier to compare for larger/smaller



Graph using R **corrgram** package

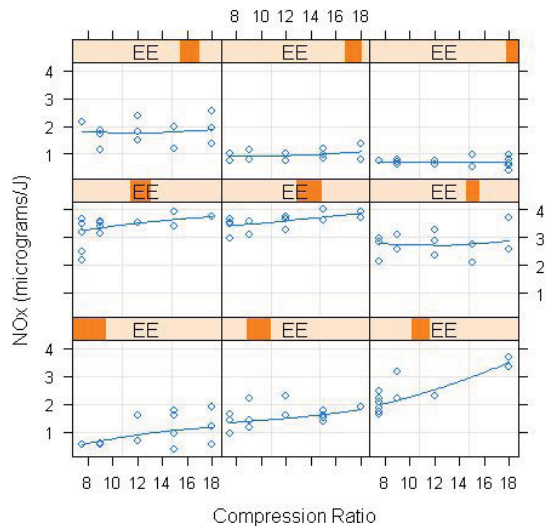
Conditional plots: $Y \sim X | Z$

Often want to explore how the relation between Y and X depends on/ varies with some other variable(s) Z.

- Moderator variables
- Interactions

Emission of NOx from ethanol in relation to engine compression ratio and richness of air/ethanol mixture (EE)

Graph using R **lattice** package



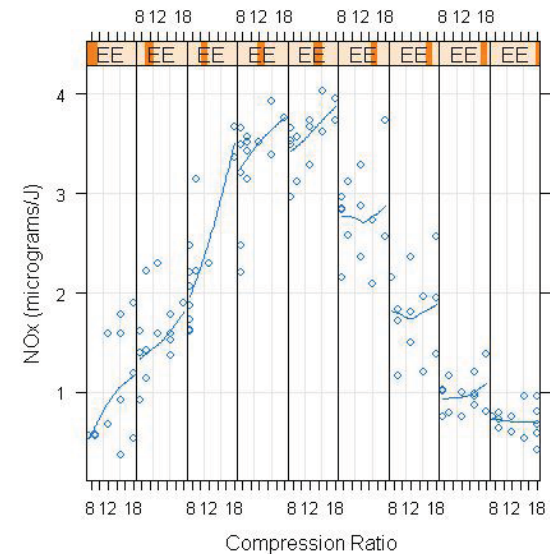
Conditional plots: $Y \sim X | Z$

The same data is shown in a different format, with

- loess smooth curves
- curves banked to ~ 45°

The joint dependence on CR and EE is now much clearer

(These are examples of **lattice plots**, produced using R software.)

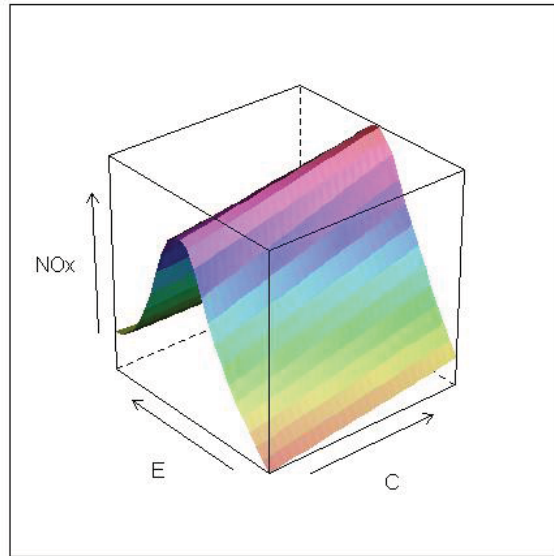


3D plots

Often not useful, unless done with great care.

This plot shows the loess **smoothed** predicted values of NOx in relation to EE and CR. (But, raw data not shown.)

Color is used to show the predicted NOx, using a "heatmap" color scale.

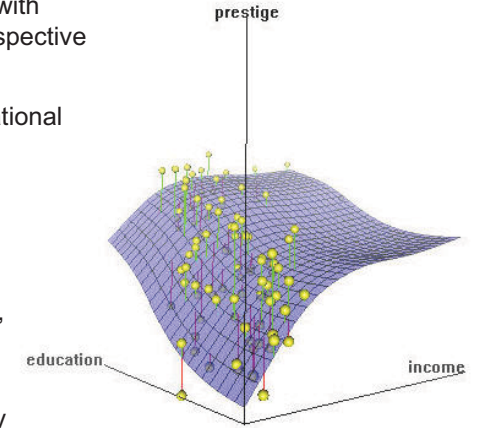


3D plots

3D plots can be enormously useful with dynamic, interactive software & perspective

This plot shows a relation of occupational prestige to income & education.

- points are shown in perspective, connected to the fitted surface
- the fitted surface (linear, quadratic, smoothed) can be changed interactively
- the plot can be rotated dynamically to see other views



30

Seeing multivariate clusters: face plot

A faces plot assigns variables to facial features, to show **configural patterns** of many variables.

Pros: Easy to see similar patterns in large data sets.

Cons:

- Hard to connect features to variables for interpretation
- No good rules/ideas for assigning variables to features.

Graph using SAS **faces** macro,
<http://datavis.ca/sasmac/faces.html>

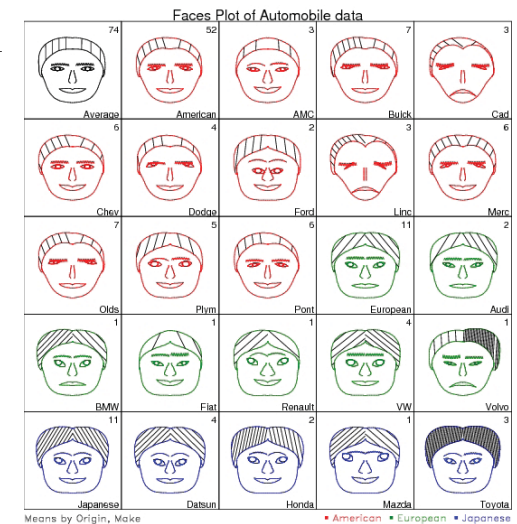


31

Seeing multivariate clusters: face plot

Parameter	Variable assignment key	
	Left Side Variable	Right Side Variable
Eye size	mpg	mpg
Pupil size	mpg	mpg
Pupil position	turn	turn
Eye slant	turn	turn
Eye X position	hroom	hroom
Eye Y position	hroom	hroom
Eyebrow curvature	rseat	rseat
Density of eyebrow	rseat	rseat
Eyebrow X position	displa	displa
Eyebrow Y position	length	length
Upper hair line	rep77	rep78
Lower hair line	weight	weight
Face line	weight	weight
Hair darkness	rep77	rep78
Hair shading slant	gratio	gratio
Nose line	length	length
Mouth size	price	price
Mouth curvature	price	price

Means, by make & origin

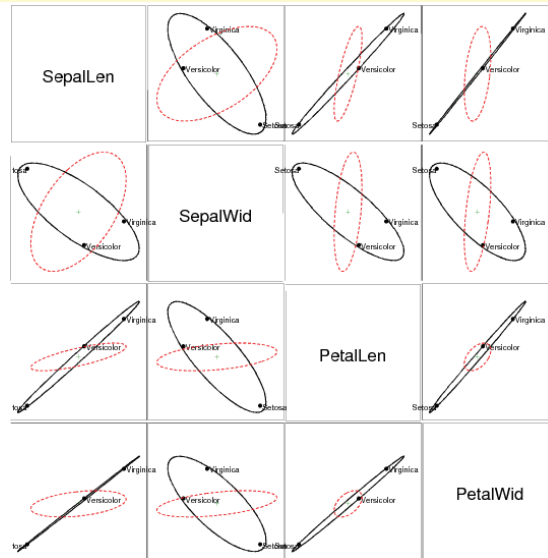


32

HE plot matrices

HE plots in a scatterplot matrix show effects for all pairs of responses.

For the iris data, the Species means are highly correlated on all variables except Sepal length.



37

HE plots: 2-way MANOVA

Plastic film data: 2x2 MANOVA

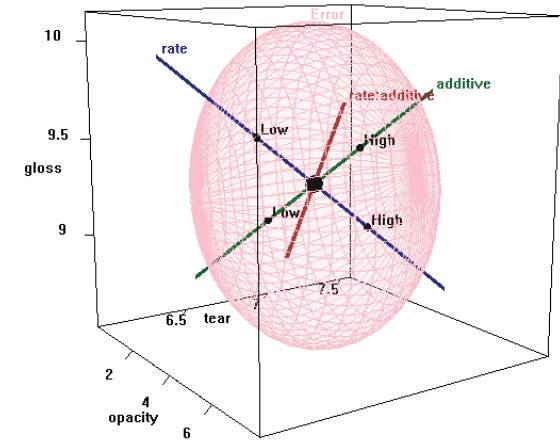
(gloss, opacity, tear) ~ rate*additive

MANOVA tests show that both main effects are significant:

```
Type II MANOVA Tests: Roy test statistic
                        Df approx F Pr(>F)
rate                    1 7.5543 0.003034 **
additive                1 4.2556 0.024745 *
rate:additive           1 1.3385 0.301782
```

HE plot shows the nature of these effects, e.g.,

high rate: ↑tear, ↑opacity, ↓gloss



1 df tests: H ellipsoid collapses to a line

38

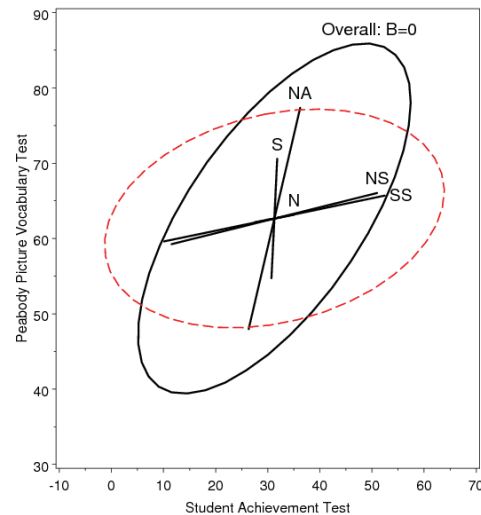
HE plots: MMRA

Rohwer data: Cognitive ability and PA tests: n=37, Low SES group

(SAT, PPVT, Raven) ~ n + s + ns + na + ss

- Only one predictor, NA, is (barely) significant

- Yet, overall multivariate test: $H_0: \mathbf{B} = \mathbf{0}$ is highly so!

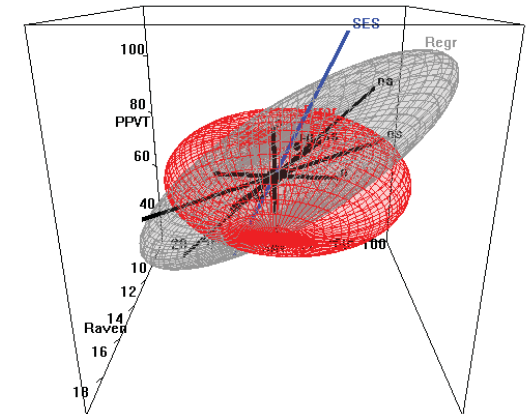


39

HE plots: MMRA & MANCOVA

Rohwer data: Low SES & Hi SES groups

(SAT, PPVT, Raven) ~ SES + n + s + ns + na + ss

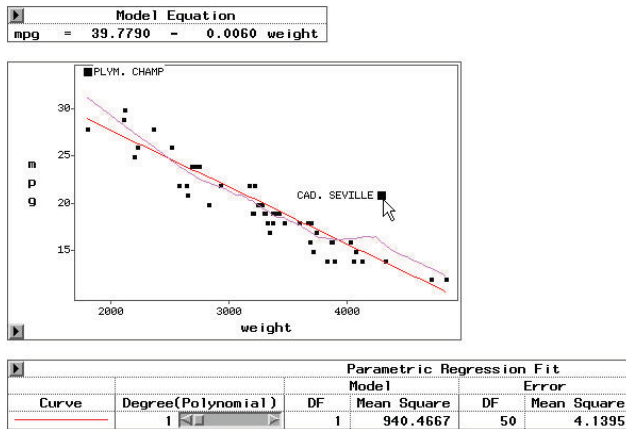


40

Dynamic, interactive graphics

Interactive graphics & data analysis provides:

- Identifying points
- Model & display controls

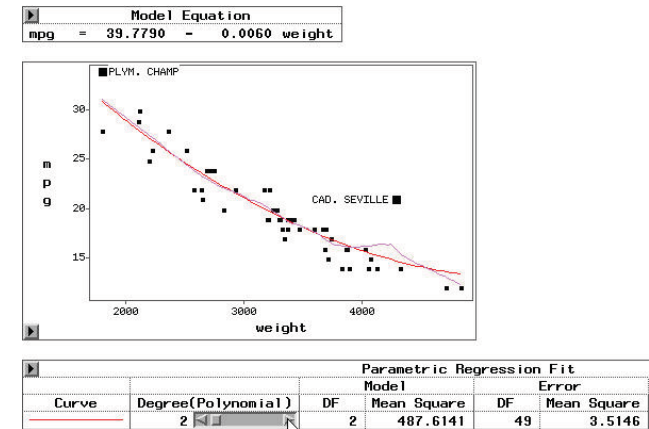


SAS/Insight: mpg ~ weight, linear fit

Dynamic, interactive graphics

Interactive graphics & data analysis provides:

- Identifying points
- Model & display controls



SAS/Insight: mpg ~ weight, quadratic fit

Dynamic, interactive graphics

Dynamic graphics provide multiple, linked views of a data set

Selecting points, regions in one plot (“brushing”) selects the same observations in all other plots

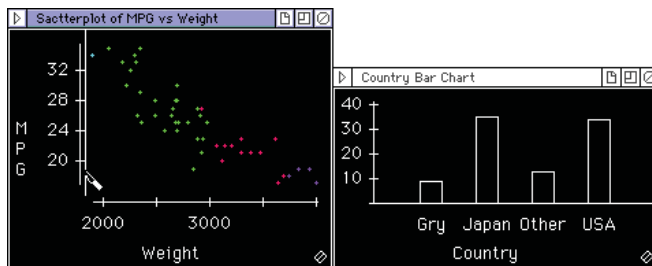
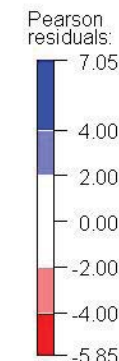
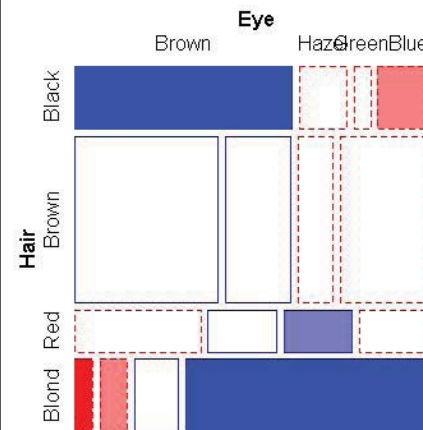


Image source: Data Desk (Paul Velleman)

See: <http://www.activstats.com/products/mediadx/custom/lessonbook/nyheart.shtml>

Multivariate frequency data: mosaic plots

Two-way table: [Hair][Eye]



A contingency table can be visualized by tiles whose area ~ cell frequency.

Shading: ~ Pearson residual,

$$d_{ij} = (O_{ij} - E_{ij}) / \sqrt{E_{ij}}$$

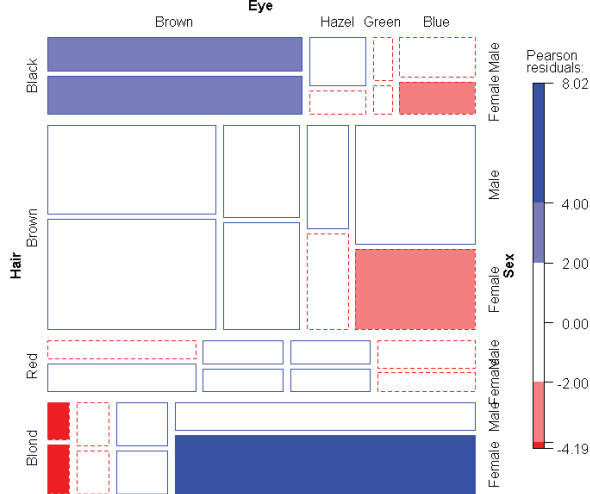
Color:

- blue: $O_{ij} > E_{ij}$; red: $O_{ij} < E_{ij}$

Interp: + association (dark hair, dark eyes), (light hair, light eyes)

N-way tables

Independence model: $[Hair][Eye][Sex]$



3+ way tables: split each tile ~ conditional proportions of the next variable.

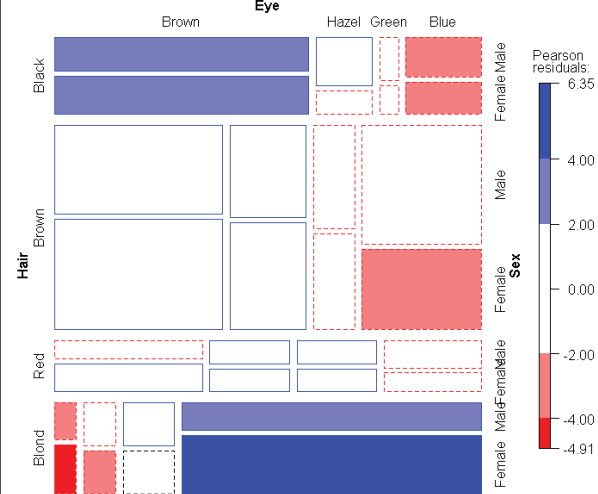
Now, there are several different models that can be fit.

- Mutual independence: $[H][E][S]$ → all vars unassociated
- Residuals: show associations not acct'd for by the model

45

N-way tables

Conditional independence: $[Hair, Sex][Eye, Sex]$



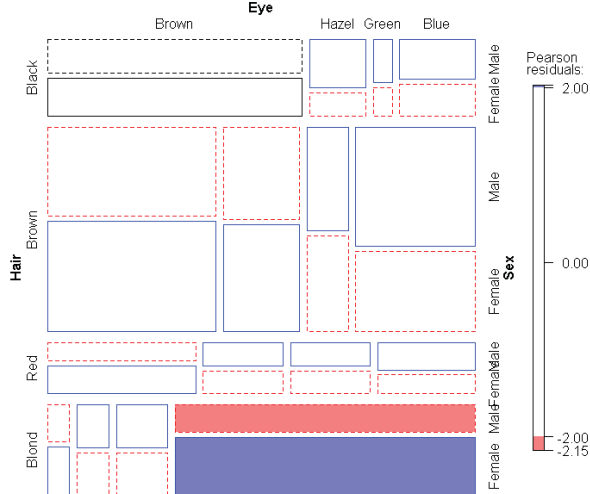
All models fit to the same table have **same-sized** tiles (O_{ijk}), but **different** residuals.

This model of conditional independence, $[HS][ES] \rightarrow H, E$ independent *given* Sex.

46

N-way tables

Joint independence: $[Hair, Eye][Sex]$



The model of joint independence, $[HE][S]$ allows Hair, Eye color association, but → $[HE]$ assoc. is independent of Sex.

This model obviously fits much better, except for blue-eyed blonds, where females are more prevalent than the model allows.

47

Summary

- Goal of statistical analysis: summarization
- Goals of graphical analysis: exposure!
 - Often more useful when enhanced with visual summaries (fitted curve, data ellipse)
- Different graphs for different purposes:
 - Reconnaissance (overview)
 - Exploration (detecting patterns, trends)
 - Model diagnosis (assumptions, outliers)

48

Summary

- Multivariate data requires novel graphs to display increasing # of variables
 - Enhanced scatterplot matrices
 - Visual thinning: less is often more
 - Low-D views (biplots)
 - HE plots to visualize multivariate tests
 - Mosaic plots to visualize n-way frequency tables.