

- Discuss Quizzes/Answer Questions
- Statistical Graphics (Not Maps)
- Box Plots
- Linear Relationships
- Correlation

Lesson Plan

Why story does the graph tell?

change over time.

William Playfair, the man who invented the histogram, made many plots for Parliament. The one shown here is a longitudinal plot, showing

now begin with some other kinds of graphics.

In the last lecture, we examined statistical graphics related to maps. We

5.1 Statistical Graphics (Not Maps)

A Pareto chart is used in industry to identify patterns of problems. The third graph is a Pareto chart, based on responses students gave in their course evaluations a few years ago.

- Where is the baby boom?
- Why are there fewer men in certain age groups?
- Why are there „notches“ in the curve?

1968. The upper plot has unnecessary detail (chartjunk) which many graphics experts dislike.

The sixth graph is another set of side-by-side boxplots, showing draft rank against birth month for the 1970 lottery to go to Vietnam. What does this show?

The fifth graph is a set of side-by-side boxplots that show how the number of hours that students spent on my course is related to their overall evaluation of the course. Why is there an irregularity for those giving the lowest rating?

The fourth graph is a histogram of the number of hours students reported working on an introductory statistics course. Why is there a spike at 9?

The seventh graph is a scatterplot of the scores that my students got on the final against the number of hours of sleep that they had.

There were two sections for the class, denoted by open circles and dots. Is there any difference between the sections?

There is a weak trend to increasing scores with increasing hours of sleep. But the trend is not strong, and the correlation is probably only about .3 or so.

What are two possible reasons for this increasing trend?

The last two are the medians of the numbers less than the median and the numbers greater than the median, respectively (according to the convention for this class).

- the 75th percentile
- the 25th percentile
- the median
- the largest
- the least

numbers:

The boxplot is a graphic that summarizes a dataset by plotting 5

5.2 Boxplots

A boxplot looks like this:

Boxplots are often made for different groups and plotted side by side to highlight differences.

Suppose one made a side-by-side boxplot of the IQ scores for Duke students and the IQ scores for students at the University of Idaho. What do you think they would look like?

The use of coordinate axes to show functional relationships was invented by René Descartes (1596-1650). He was an artillery officer, and probably got the idea from pictures that showed the trajectories of cannibalists.

If you don't know this stuff, read the book carefully.

$$q + X^m = Y$$

The algebraic equation for a line is

5.3 Lines

Sir Francis Galton explored Africa, invented eugenics, studied whether ships that carried missionaries were less likely to be lost at sea, pioneered birth-and-death models and meteorology, and was Charles Darwin's cousin.

He also was the first to conceive of linear regression (although he did not have the mathematical skill to develop the formulae, and got a friend of his at Cambridge to do the derivations).

5.4 Correlation

are taller. Thus the sons' height tends to "regress towards the mean". That are shorter; an exceptionally short father tends to have sons that the father-son example. An exceptionally tall father tends to have sons Regression fits a line to the points in a scatterplot. The term comes from

scatterplot.
son's height, then all father-son pairs would lie on a straight line in a have short sons. If the father's height were a perfect predictor of the Clearly, tall fathers tend to have tall sons, and short fathers tend to

the height of fathers and the height of sons.
two continuous variables. The book discusses the relationship between Correlation is a measure of the strength of the linear association between

X . If this is 0, then X provides no insight about Y (assuming linearity).
is the proportion of the variation in Y that is explained by knowledge of
The square of the correlation is called the **coefficient of determination**. It

All correlations are between -1 and 1, inclusive.

If the correlation is near zero, then knowing one variable gives essentially
no information about the other.

The correlation tells you how much information one of the variables
provides about the other. If the absolute value of the correlation is near
1, then knowing one tells the other almost perfectly (if the relationship is
linear).

To calculate the correlation coefficient, just take the average of the products of the z -transforms of the X values with the z -transforms of the Y values.

If the assumption does not hold, then one can get strange and even misleading behavior.

Where the Y_i and X_i are the observed values for the i th case (e.g., a father-son pair) and the ϵ_i is random error (due to genetics, measurement error, etc.).

$$Y_i = \mu + q + \epsilon_i$$

The assumption of linearity is important. Specifically, it asserts that:

3	$(3-2)/.8165 = 1.2245$	5	$(5 - 2)/2.1602 = 1.3888$	
2	$(2-2)/.8165 = 0$	1	$(1-2)/2.1602 = -.4629$	
1	$(1-2)/.8165 = -1.2245$	0	$(0-2)/2.1602 = -.9258$	
	X	Y	XZ	X

The z -transformation subtracts the mean and divides by the standard deviation for each of the X and Y values, giving:

$$sdy = \sqrt{\frac{1}{3}(0^2 + 1^2 + 5^2) - 2^2} = 2.1602$$

$$sdx = \sqrt{\frac{1}{3}(1^2 + 2^2 + 3^2) - 2^2} = .8165$$

Example: Suppose your data are $(1,0)$, $(2,1)$, and $(3,5)$. Then the mean of the X values is 2, the mean of the Y values is 2, and the sds are:

$$r^2 = (.9446)^2 = .89.$$

How much of the variation in Y is explained by knowing X ? This is

Do we really believe this?

$$r = \frac{1}{3}[(-1.224)(-.9258) + (0)(-.4629) + (1.2245)(1.3888)] = .9446.$$

Then the correlation coefficient is