

1/29/15

263 Lecture 5

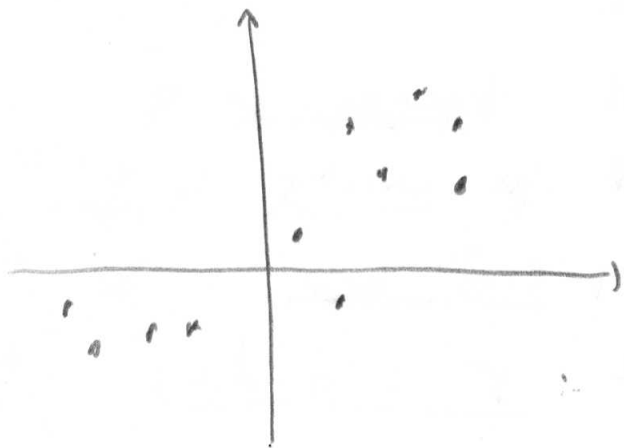
Announcements

- Chap 2 webassign due tomorrow
- Exam I next Thurs; Chapters 1, 2
- Quiz today
- Excel HW 2 DUE TUE 2/3

Last time

* Data Set w/ 2 Quant Vars X & Y

* Scatter Plot: Plot $\{(x_i, y_i)\}$.



* Correlation: Compute \bar{X} , \bar{Y} , S_X , S_Y , z-scores for all data pts

$$z_{x_i} = \frac{x_i - \bar{X}}{S_X} \quad z_{y_i} = \frac{y_i - \bar{Y}}{S_Y}$$

Then
$$r = \frac{1}{n-1} \sum_{i=1}^n z_{x_i} z_{y_i}$$

* $-1 \leq r \leq 1$

* Least Squares Line:

$$\hat{Y} = mX + b$$

$$m = r \frac{S_y}{S_x}$$

$$b = \bar{Y} - m\bar{X}$$

* Least Squares Line minimizes

$$\sum_{i=1}^n r_i^2$$

"Residuals"

$$r_i = Y_i - \hat{Y}(x_i)$$

↑ Observed ↑ Predicted

* Slope: an increase of 1 unit in X results in a change of m units in Y.

-OR- an increase of 1 std dev in X results in a change of r std devs in Y

* r^2 is the fraction of variation in Y "explained" by variation in X.

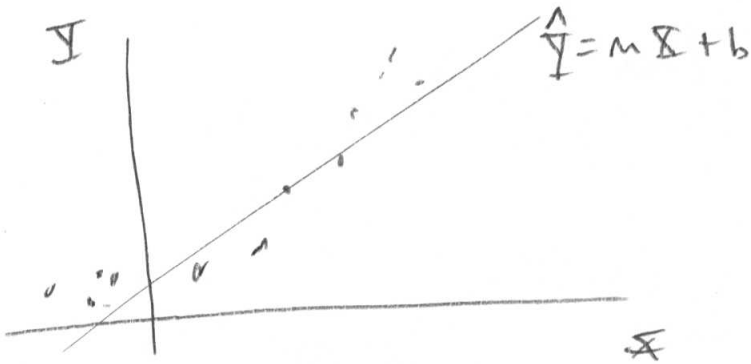
$$r^2 = \frac{\text{Variance of Predicted vals}}{\text{Variance of observed vals}}$$

Residual Plots

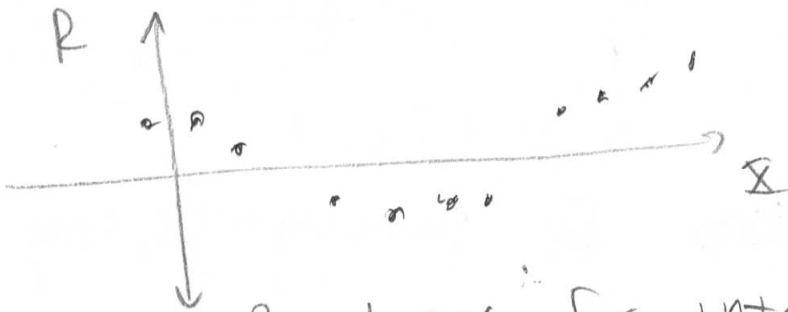
* Residual is Observed - Predicted

$$= Y_i - \hat{Y}_i$$

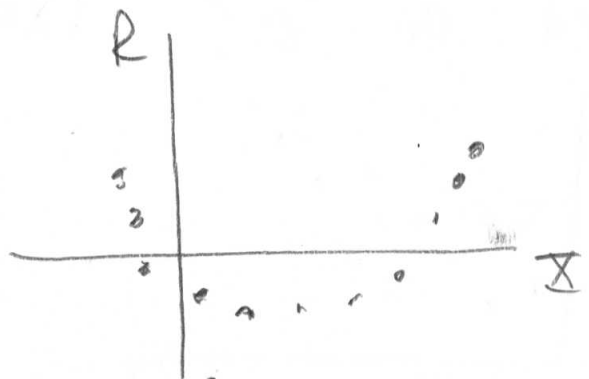
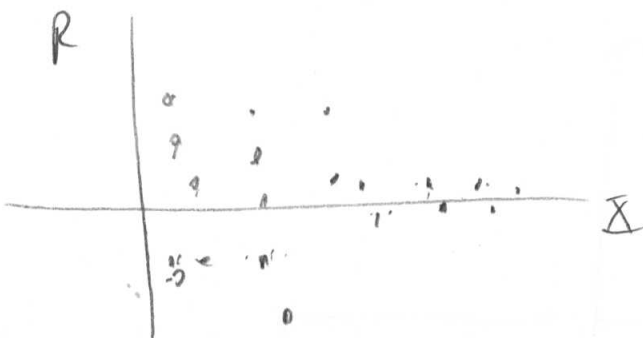
* Patterns in residuals are bad



Residual:



* Look @ Residuals for internet vs. Birth rate data
* From the Residual Plot, what do you think the data is doing?



* Sometimes, a log can help:

$$Y = ab^X$$

$$\ln(Y) = X \ln(b) + \ln(a)$$

So a plot of $(X, \ln(Y))$ should be a straight line

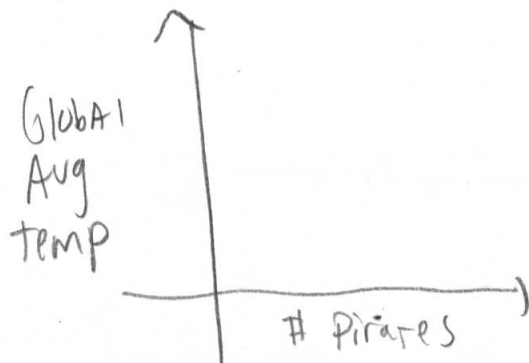
$$Y = aX^P$$

$$\ln(Y) = P \ln(X) + \ln(a)$$

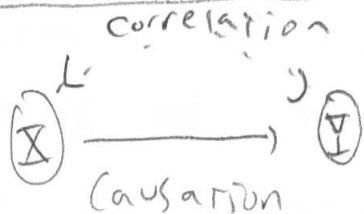
So a plot of $(\ln(X), \ln(Y))$ will be straight.

* Outliers: look @ data for Smoking vs, lung cancer. Outliers can be influential.

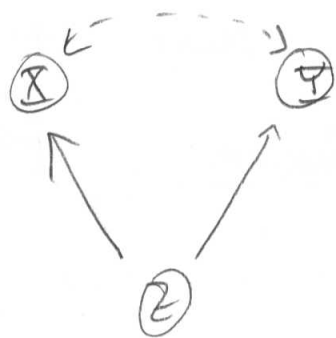
* Lurking Variables: maybe both X & Y depend on Z ? Examples



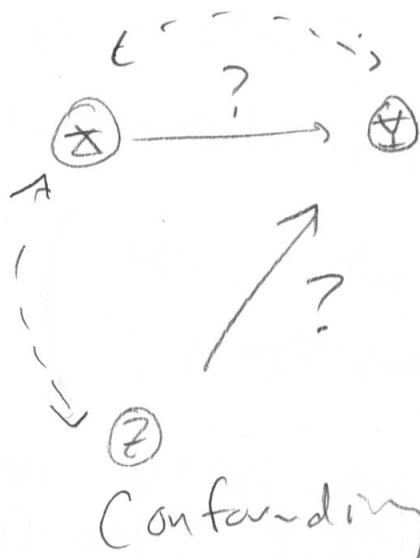
Correlation VS. Causation



* Correlation is Necessary for causation, but
but not Sufficient.



Common Response



Confounding

Example

Aspirin Prevents Heart attacks

(See link on website)

* How do we establish causation?

- Perform EXPERIMENTS
- Eliminate lurking vars / confounding
- ("Control the influence of all vars")

Criteria for Causation

Sometimes an experiment is NOT practical or ethical. Example EXPOSURE TO X-RAYS increases chance of cancer in human subjects. Can't do the experiment we want, which would be to expose many people to different levels of X-rays (controlling for other factors) So, we use existing data: Patients in hospitals;

X-ray exp Level (mSv)	# deaths from Cancer

- Strong Association between dose level & # Cancer deaths
- Consistent across many hospitals / region (controlling for environ. factors)

- Higher dose, stronger response
- Causality is logical in time
 - Can we eliminate those patients who already had cancer?
- The cause is plausible
 - Supporting evidence from e.g. bomb survivors, experiments on animals, etc.

Chapter 3: Producing data.

How do we collect data to answer statistical questions?

Conjectural / anecdotal data:

- I think that an increase in Vitamin C intake will decrease rate of common cold because I have heard it is TRUE.
- I have heard that individuals on welfare are more likely to do drugs
- I am told that vaccines cause autism.

Nullius in Verba

Motto of the Royal Society - "Take nobody's word for it"

Where do we get data from?

* "Available" data eg. US Census, NAEP, data.gov

* "Gathered" data: either observational

or EXPERIMENTAL

WE observe & measure
E.g. Sunlight vs. Size of Plants

↓
We impose a treatment on Cases / individuals
to modify a variable

E.g. Artificial sunlight vs. size of plants, Aspirin dose, etc.

FOR TUESDAY

- Read section 3.2
- Read MMR/AUTISM Paper posted on website
- Bring Q's for Exam (Half of class)