

Lesson #19
6/27/11

Admin

- Cont r.v.'s not on final (p268)

Plan

- review

- rest of ch 6

- ch 19.1-19.3

- parts of ch 2 will also be covered

We talked about dependent r.v.'s and non-zero covariance:

$$\text{Cov}(X, Y) \triangleq$$

$$E[(X - \mu_X)(Y - \mu_Y)] \neq 0 \Rightarrow X, Y \text{ dependent}$$

Covariance was a hard metric to conceptualize

$$\Rightarrow \rho_{XY} \triangleq \text{Corr}(X, Y) \triangleq \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \in [-1, 1]$$

Where -1 is perfect ^{neg.} linear association

and $+1$ is perfect pos. linear association

Bayesian analysis of interest on interest rates.

So (X, Y) is viewed in pairs

(x_1, y_1)
 (x_2, y_2)
 \vdots
 (x_n, y_n)

As usual, we do not know the parameters $\mu_x, \mu_y, \sigma_x, \sigma_y$ nor ρ . so we use statistics to estimate them (inference).

$$\bar{x} \rightarrow \mu_x, \bar{y} \rightarrow \mu_y, s_x \rightarrow \sigma_x, s_y \rightarrow \sigma_y$$

$$s_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1} \rightarrow \text{Cov}[X, Y]$$

sample correlation coefficient
"corr"

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}} = \frac{s_{xy}}{s_x s_y} \rightarrow \rho_{xy}$$

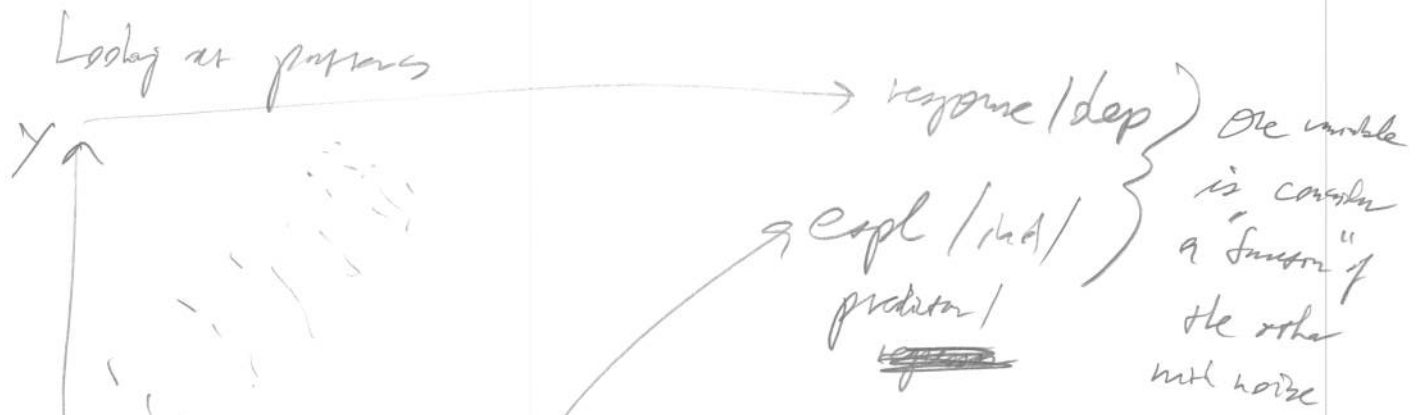
Just like there's hypothesis tests to test values of μ , there's hypothesis tests to test values of ρ which we won't cover. Since we're going to do something more painful.

If there are many ^{internal} variables (imagine a table with lots of cols) then a corr matrix lists the sample corr for each pair

pl/4

| | Age | Price | Size |
|-------|-----|-------|------|
| Age | | 0.58 | 0.67 |
| Price | | | 0.13 |
| Size | | | |

talk about diagonal and symmetry



Looks for and linear down common outline, take one step

- Things to look for
- ① direction
 - ② linearity
 - ③ variance / spread
 - ④ outliers

Derive of
Simpson's paradox
~~keep~~ best prof

ch 6.4: association in line (also all of ch 19)

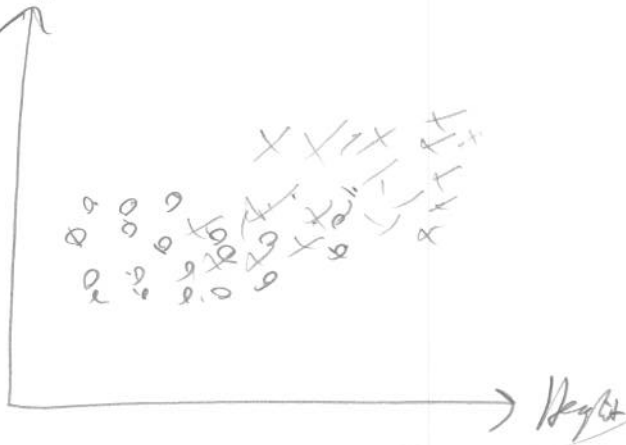
We will now derive all these formulas

The method is "regression" - a strange word which will talk about tomorrow... it has ~~two~~ main goals:

- ① To best predict a response given the ~~pre~~ explaining value
- ② To measure how significant a predictor is in affecting change in the response

Goal #2 you will cover ad-hoc in Stat 102. For now, our goal is #1.

Salary



remember where the O's, X's are...

describe this association: linear, pos

Does this mean being taller causes you to earn more money? NO

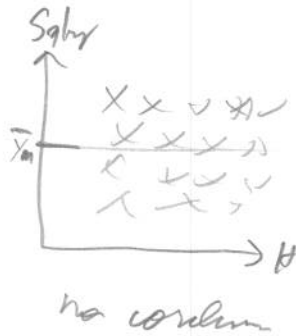
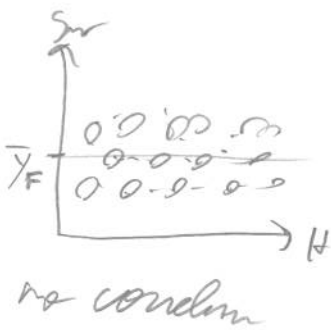
Correlation \Rightarrow Causation

"Spurious correlation" ch 6-8

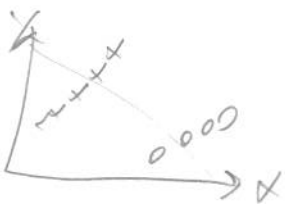
Diff between correlation and association?

We learn before there may be a lurking variable and we can get Simpson's Paradox (when controlling for a lurking variable, association can change)

Let's say O's are women, X's are men



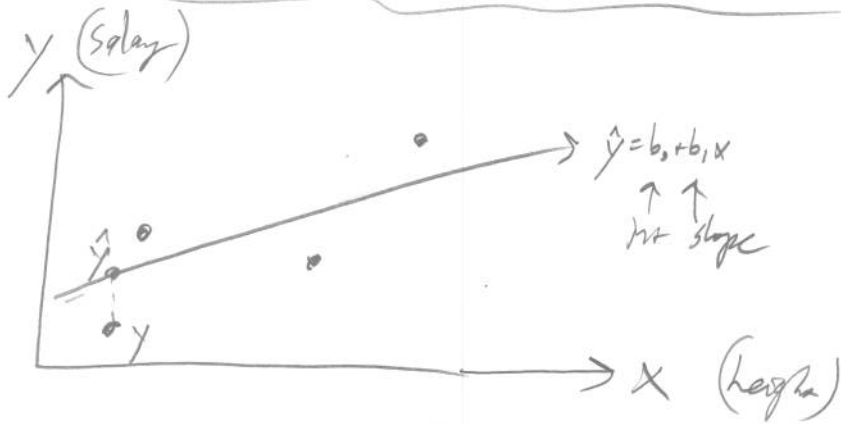
What happened to the association?
GONE! Why was there an association to begin with? Height tells you something about gender and ~~gender~~ men make more than women...



\Rightarrow



(Wikipedia)



Step I Assume
relationship is linear
which is why we call it
"linear regression"

We want to draw the "best fit line" or "fitted line" or
"least square regression line". How do we do this?

Maybe we want to put some sort of "error" function. Call it err .
 err should take in our x value and our y value and
spit out an error.

→ Our predicted value we call \hat{y} .

$$err(\hat{y}, y)$$

What should this be?

Possible choice:

the squared function

$$Var(X) = E[(X-\mu)^2] \text{ AND } X \sim N(\mu, \sigma^2) \triangleq f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

Something is very special about the squared function. It just seems
natural. In advanced courses you'll see the stats initial
assumption gets you a lot of benefits.

Now, we want to minimize the total amount of this error and the
 pos... so we want to minimize:

$$\sum_{i=1}^n \text{err}(y_i, \hat{y}_i) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (b_0 + b_1 x_i))^2$$

$$= \sum (y_i - b_0 + b_1 x_i)^2$$

with respect to all possible vals of b_0, b_1 . How do we do that?
 Penetration calculator? Let's do b_0 first!

$$\frac{\partial}{\partial b_0} \left[\sum_{i=1}^n (y_i^2 + b_0^2 + b_1^2 x_i^2 - 2y_i b_0 - 2y_i b_1 x_i + 2b_0 b_1 x_i) \right] = 0$$

$$\Rightarrow \frac{\partial}{\partial b_0} \left[\sum y_i^2 + n b_0^2 + b_1^2 \sum x_i^2 - 2 \sum y_i b_0 - 2 \sum y_i b_1 x_i + 2 b_0 b_1 \sum x_i \right] = 0$$

$$\Rightarrow 2n b_0 - 2 \sum y_i + 2 b_1 \sum x_i = 0$$

$$b_0 - \bar{y} + b_1 \bar{x} = 0$$

$$\Rightarrow \boxed{b_0 = \bar{y} - b_1 \bar{x}} \quad \checkmark$$

$$\frac{\partial}{\partial b_1} [\dots] = 0$$

$$\Rightarrow \boxed{b_1 = r \frac{s_y}{s_x}}$$

From bivariate and bivariate statistics, you can compute slope + intercept
 of best fit line.

17

$$\hat{y}_i = \bar{y} - b_1 \bar{x} + b_1 x_i = \bar{y} + b_1 (x_i - \bar{x}) = \bar{y} + r \frac{s_y}{s_x} (x_i - \bar{x})$$

$$\hat{y}_i = \bar{y} + r \frac{s_y}{s_x} (x_i - \bar{x})$$

Look on slides

$$\hat{y}_i - \bar{y} = r \frac{s_y}{s_x} (x_i - \bar{x})$$

$$\Rightarrow \frac{\hat{y}_i - \bar{y}}{s_y} = r \left(\frac{x_i - \bar{x}}{s_x} \right)$$

$$\Rightarrow \boxed{z_y = r z_x} \text{ // done w/ ch. 6}$$

- talk about residuals
- talk about $\bar{e} = 0$
- talk about \bar{x}, \bar{y}
- talk about RMSE