

1/27/15

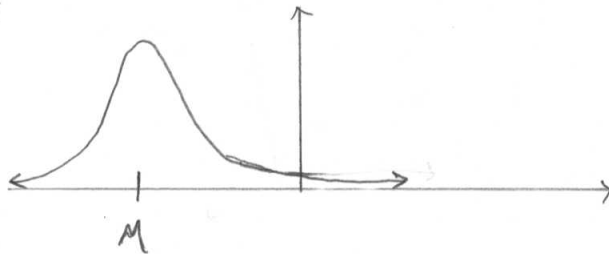
263 Lecture 4

Announcements:

- Chapter 2 webAssign DUE FRI 1/29/14
- EXCEL HW 2 Posted, DUE NEXT TUE (2/3)
TOPIC IS Scatter Plots & regressions
- EXAM 1 is 1 week from Thursday
- Tomorrow (1/28) is the last day to
DROP on VACCESS w/out a "W".

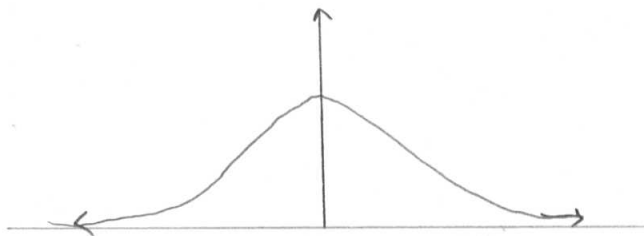
Last time

$N(\mu, \sigma)$:



"Normal Distribution
with mean μ ,
Std dev σ "

$N(0, 1)$:



"Standard Normal"

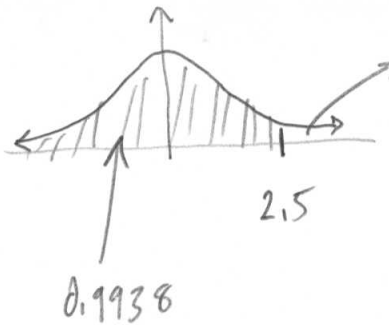
$z = \frac{x - \mu}{\sigma}$ * to Find Proportion of a Normally distributed Variable, convert to z-scores & use z-table

* z table gives Proportion of z Scores to the left of the given z-score.

Ex Suppose $X \sim N(10, 2)$. Find % of observations bigger than 15. $z = \frac{15 - 10}{2} = \frac{5}{2} = 2.5$

Go to Z Table:

	0.0
2.5	0.9938



Above $z=2.5$
 is $1 - 0.9938$
 $= 0.0062$

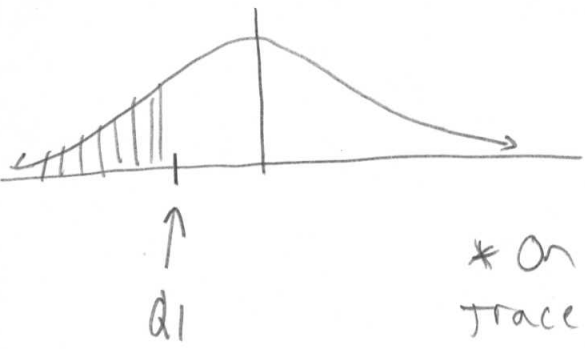
So 0.62%

ON Calculator:

2ND: VARS (DISTR) → normalcdf(a, b, μ , σ)

will give area between a & b under $N(\mu, \sigma)$.

Inverse: What is the IQR of $N(0, 1)$?



Q1 is where "CDF" (cumulative area) is 0.25

* On Z-table: look for # in table, trace back to row & col.

* On Calc: DISTR: $invNorm(0.25, 0, 1)$
 ↑ ↑ ↑
 Area μ σ

Chapter 2: Relationships

Two variables can be Associated

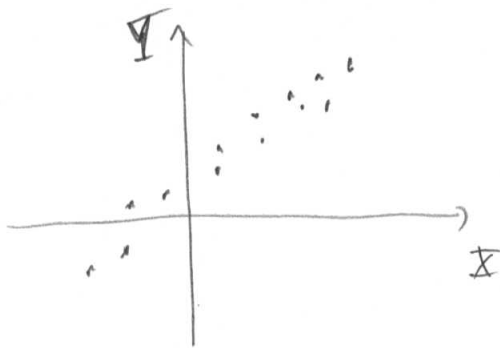
IE Their Scatter Plot Shows some Pattern:



ASSOCIATIONS can be Weak:



Strong:



Linear:



Nonlinear:



For linear associations, we can have either a Positive or Negative association

"An increase in the cost of gas will decrease the total miles driven"

"An increase in the # of calories consumed will increase your BMI"

Measuring amount of linear association:

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

$$z_{x_i} = \frac{x_i - \bar{x}}{s_x}$$

$$z_{y_i} = \frac{y_i - \bar{y}}{s_y}$$

- If $z_{x_i} z_{y_i} > 0$, There is a Positive association between X_i & Y_i (1st or 3rd Quad on re-centered Scatter Plot)

- If $z_{x_i} z_{y_i} < 0$, There is a negative association between X_i & Y_i (2nd or 4th Quad, on Scatter Plot)

r is the Average of all these, so

If "MOST" $z_{x_i} z_{y_i}$'s are > 0 , $r > 0$

If "MOST" $z_{x_i} z_{y_i}$'s are < 0 , $r < 0$

~~⊗~~ $-1 \leq r \leq 1$

* $r = -1$ is Perfect negative Correlation

So $\hat{Y} = mX + b$ w/ $m < 0$

* $r = 1$ is Perfect Positive Correlation

$\hat{Y} = mX + b$ w/ $m > 0$

* $r = 0$ No Correlation

- Most of the time, $-1 < r < -0.5$
 $-0.5 < r < 0$ or $0 < r < 0.5$, which means

That a linear model might be a good fit

IE $Y \approx mX + b$ is a good guess

" If I increase X , Y will increase (decrease)

by about $m \Delta X$ "
↑ change in X .

Notes about r :

- only measures linear association
- Both X & Y should be quantitative
- r is not robust

Linear Regression

- We have a data set with 2 variables X & Y ,
and we think Y might be adequately described

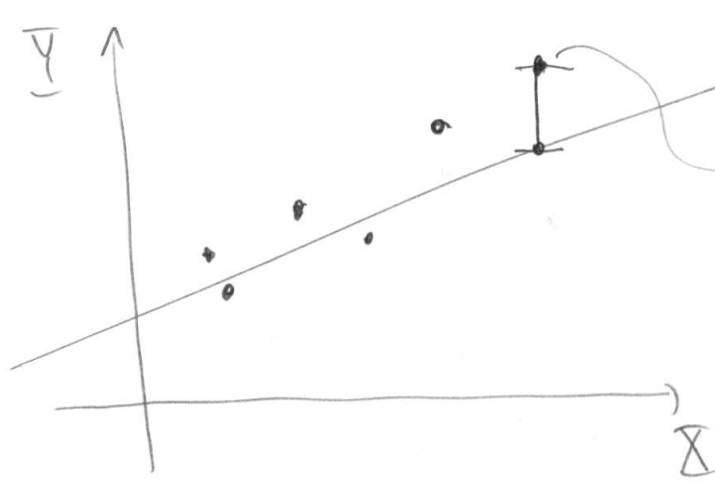
by $\boxed{Y = mX + b}$ (really, we think $Y \approx mX + b$
with small error!)

- m is the slope - it is interpreted as

"When X increases by 1 unit, Y increases or
decreases by m units"

IE $F = 1.8C + 32$ ($F \sim$ °Fahrenheit, $C \sim$ °Celsius)

- We would proceed w/ a linear regression if we had a correlation close to ± 1
- How do we find the "Best fit line"?
- Want to "minimize distance from line to data set"
- Standard (aka old) method: least squares



$Y = mX + b$

* Data: (x_i, y_i)

* linear model: $(x_i, mx_i + b)$

* squared error: $(y_i - (mx_i + b))^2$

least squares minimizes $\sum \text{error}^2$

$\hat{Y} = mX + b$	where $m = r \frac{S_Y}{S_X}$	* Book uses $Y = b_1X + b_0$
	$b = \bar{Y} - m\bar{X}$	

$r \sim$ Correlation Coeff

$S_X, S_Y \sim$ Std deviations of X & Y

$\bar{X}, \bar{Y} \sim$ means of X, Y

How Can We Use a Least Squares Lines?

* Interpolation (Don't have data point for $X=7$, what should it be?)

* NOT Extrapolation

Example using The Fruits/Veggies vs Smoking data, we have

$$\bar{X} = 23.9 \quad S_x = 4.37$$

$$\bar{Y} = 15.33 \quad S_y = 3.58$$

$$r = -0.745$$

So, $Y = mX + b$ w/ $m = r \frac{S_y}{S_x} = -0.745 \frac{3.58}{4.37} \approx -0.61$

$$b = 15.33 - (-0.61)23.9$$

$$= 29.91$$

So $Y = -0.61X + 29.91$ (to 2 Dec Places)

Interpretation:

If a state increases its Average fruit & veggie intake by 1 per day, 0.61% Fewer People will smoke everyday (Causation??)

FACTS About linear regression

① Alternate Interpretation of Slope:

$m = r \frac{S_y}{S_x}$ If X increases by 1 std dev.,

Y will increase (decrease) by r std deviations

IE

Student	X midterm	Y final
1	65	77
2	43	82
3	86	84

$$\bar{X} = 70 \quad S_x = 10$$

$$\bar{Y} = 72 \quad S_y = 8$$

$$r = 0.5$$

If a student's midterm score was 10 pts higher, we would expect their final exam score to be about 4 points higher.

② How good is the fit?

r^2 is the fraction of variation in Y that is explained by the least-squares fit.

Eg. For fruits/smoking, $r^2 = 0.55$. So about 55% of the variation in Y is explained by variation in X , via the least sq. line.

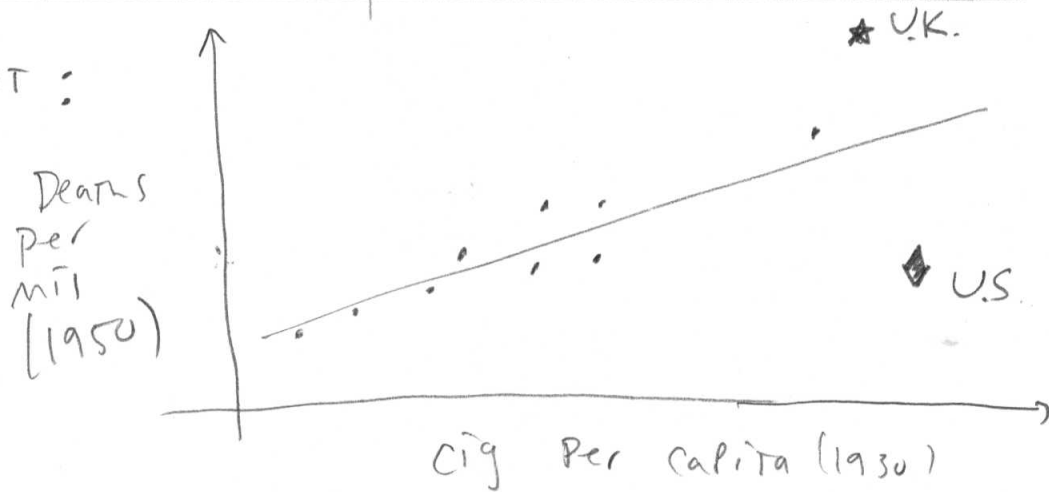
③ LS line always passes through (\bar{X}, \bar{Y})

④ OUTLIERS: Consider The following data.

Cig Per capita (1930) Lung Cancer death Per Million (1950)

	Cig Per capita (1930)	Lung Cancer death Per Million (1950)
Iceland	220	58
Norway	250	90
Sweden	310	115
Denmark	380	165
Australia	455	170
Holland	460	245
Canada	510	150
Switzerland	530	250
Finland	1115	350
UK	1145	465
US	1280	190

- Scatter Plot :



$$Y = 0.2291X + 65.749$$

$$R^2 = 0.549$$

Take out U.S. ;

$$Y = 0.3577X + 13.553$$

$$R^2 = 0.8855$$

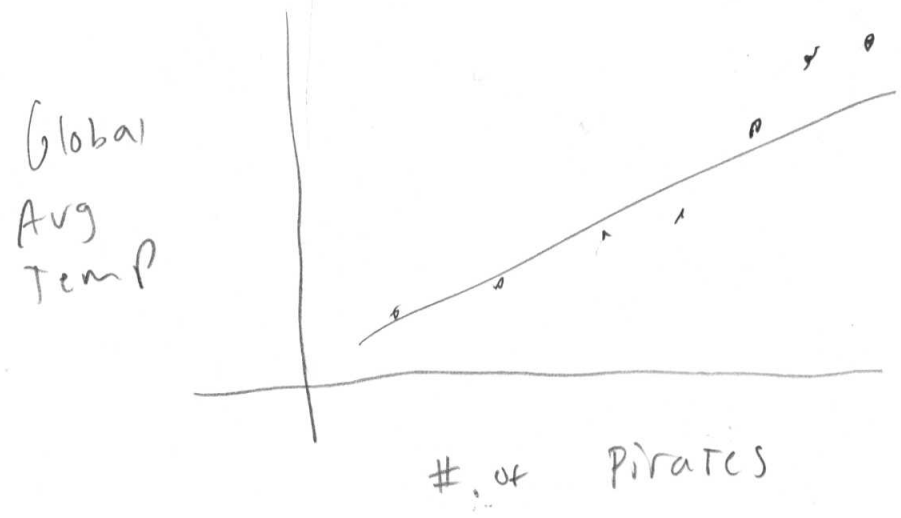
Warnings about Regression / Causation:

- "Everything" is a line if you get "localized" enough! ("Restricted range")
- Inspect Residuals to see Patterns:

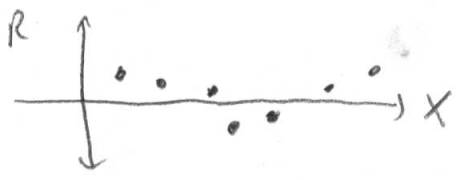
Residual = Observed y - Predicted \hat{y}

Lurking Vars

Pirates Cause Global Warming:



EX (Residuals) Sketch a graph of the scatter plot given the residuals:



Section 2.7: Read on your own. Correlation is Necessary for causation, but NOT sufficient. Experiment

- OR-5 criteria:
- 1) Strong Assoc.
 - 2) Consistent across many studies
 - 3) Strong response
 - 4) Time
 - 5) Plausibility