

Week 1

Monday, December 11, 2017 8:22 PM

Suggested Reading
Open Intro Statistics Ch1: 1.1 - 1.5

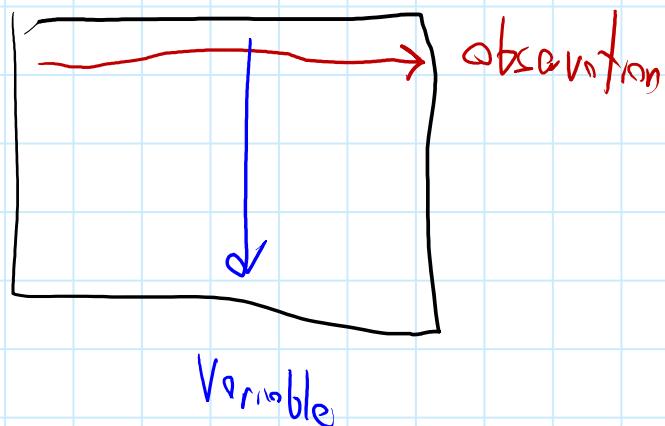
Exercises

1, 1, 3, 11, 13, 17, (9, 25, 27, 3)

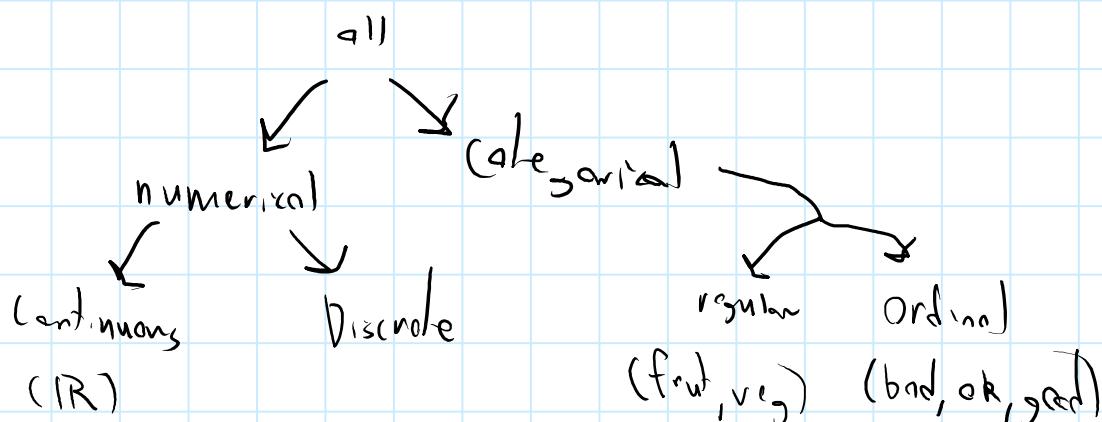
Designing Studies

Monday, December 11, 2017 9:21 PM

Data Matrix



Types of Variables



Observational Studies & Experiments

Observational: No interference, can only establish correlation

Experimental: Randomly assign treatments to subjects, can establish causation

Ex: Does exercise increase energy?

Poll people who do/don't exercise on energy level

Randomly assign a group to exercise or not

Confounding Variable: extraneous variable controlling correlated values
(Why correlation does not imply causation!)

Sampling & Sources of Bias

Census: poll entire population

issues: hard to locate or measure people's data won't be recorded
populations in constant flux

Ex. Cooking a pot of soup (census is eating the pot!)

Exploratory analysis: Tasting a spoonful

Representative sample: The spoonful

Inference: deciding to add salt

ISSUE: what if spices had settled at the bottom? (bad sample)

To prevent, stir (randomize)

Sampling Bias!

Convenience - Easily accessible individuals only, ones included

Non-response: Non-random fraction of sample responds, (only)

Voluntary response: Opinionated people respond more likely

Sampling Methods

Simple Random Sample - every case equally likely to be selected

Stratified Sample - Divide pop into homogeneous strata, randomly sample from these

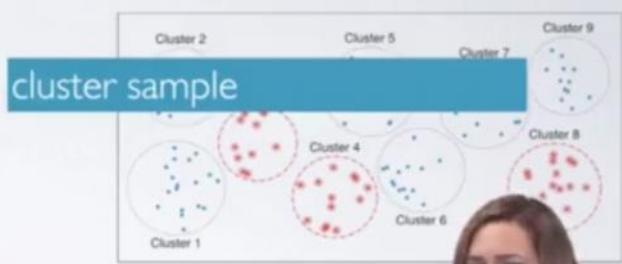
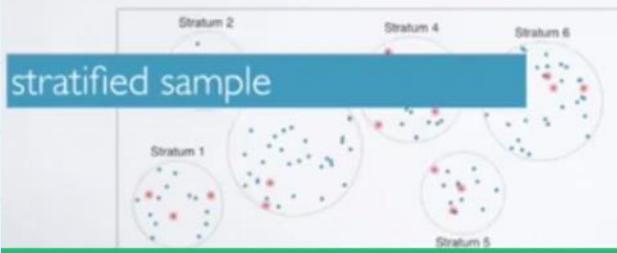
Cluster - Divide pop into clusters, randomly select clusters, sample all w/in selected

↓ used in econ

(heterogeneous)

Multi-stage - Cluster, but randomly sample within clusters

sampling methods



Experimental Design

Blocking ex. Do energy gels make you run faster?

Issue prob. might be affected differently

ISSUE: program might be affected differently
Solut'n stratify into programs

Blocking vs Explanatory

Explanatory can be imposed on an experiment (e.g., the energy gel)
Blocking, intrinsic (e.g., program or M/F)

Week 2

Wednesday, December 13, 2017

2:25 PM

Ch. 1 sec 6 - 8

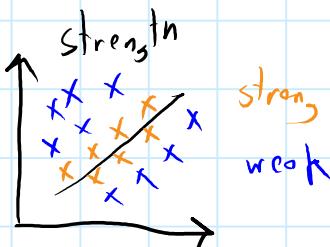
Exercises: Ch 1 39, 41, 45, 49, 51, 55, 59, 63, 65, 67

Exploring Numerical and Categorical Data

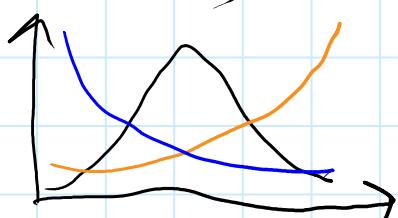
Wednesday, December 13, 2017 2:25 PM

Visualizing Numerical Data

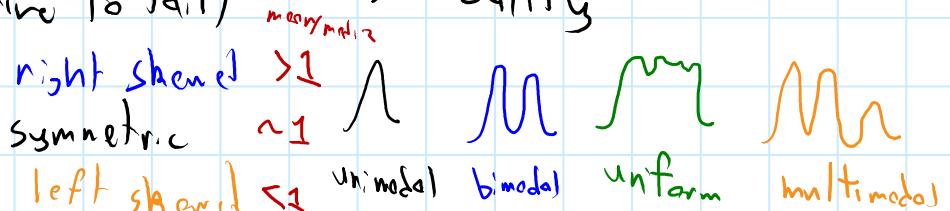
Correlations



Skewness (relative to tail)



Modality



Measures of Spread

Variance: roughly avg square deviation from mean

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

why? Gets rid of negatives to avoid cancellations
Amplifies large deviations

Standard deviation:

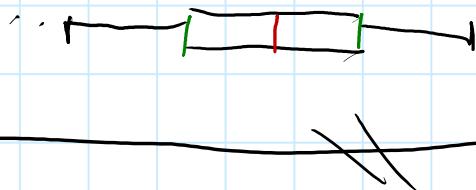
$$s = \sqrt{s^2}$$

Venability \neq Diversity

$$\{1, 1, 1, 5, 5\}$$

$$\{1, 2, 3, 4, 5\}$$

IQR = $Q_3 - Q_1$



Robust Statistics

- measures on which extreme observations have little effect

	robust	non-robust
center	median	mean
spread	IQR	SD

Good for skewed distributions
Good for symmetric distributions

Transformation

Used to make heavily skewed data more symmetric,
more useful for modelling,

e.g. log



Categorical Variables

Bar plot > pie chart

Trend in dependent categorical variables?

Relative Frequency Segmented Bar Plot

Mosaic Plot ~ aesthetic ~

Intro to Inference

Wednesday, December 13, 2017 8:48 PM

Case Study on Gender Discrimination

Identical file given to promotion managers, gender randomly assigned

Promoted					
♂	1				
♂	21	3	24	88%	H_0 Difference was due to chance
♀	14	10	24	58%	H_A Discrimination
	35	13	48		

Probability of observing an outcome as extreme as what is represented in original data is the p-value

Week 3

Tuesday, December 19, 2017 8:43 PM

Reading: (h2: 1, 2)

Defining Probability

Tuesday, December 19, 2017 8:44 PM

Intro

Random process: know what could happen, but not what will

$P(A)$ = probability of event A

Frequentist interpretation of Probability:

$P(A)$ = proportion of times A would occur if the random process were observed ∞ times

Bayesian interpretation:

Probability is a subjective degree of belief
Allows integration of prior info

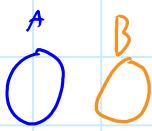
Law of large #'s:

As more observations are made, $\frac{A}{\text{all possibilities}}$ converges on $P(A)$



Disjoint Events & General Addition Rule

Disjoint = mutually exclusive



Non-disjoint = can, but don't have to, occur together



Venn diagram - mutually exclusive  Non-disjoint = can, but don't have to, occur together 

For disjoint events, $P(A \text{ or } B)$

$$= P(A \cup B) = P(A) + P(B)$$

Non-disjoint: $P(A \text{ or } B)$

$$\text{General Addn} = P(A) + P(B) - P(A \cap B)$$

Sample Space:

Collection of all possible outcomes of a trial

Probability distribution!

rules:

All outcomes in sample space
AND the probabilities they will occur

Events listed must be disjoint

$$0 < P_i < 1$$

$$\sum P_i = 1$$

Complementary events:

2 mutually exclusive events whose probability $\Sigma = 1$

e.g. 1 coin toss

$$A = \text{head}, B = \text{tail}$$

2 coin tosses

$$A = HH, B = \Sigma \text{ of the rest (HT or TH or TT)}$$



Independence

2 processes are independent if knowing the outcome of one provides no useful info on the other

e.g. 2 coin flips

Checking for independence.

$$\underbrace{P(A|B) = P(A)}_{\text{Probability of } A \text{ given } B}, A \text{ and } B \text{ are independent}$$

If data suggests dependence, do a hypothesis test!
or, use inference first:

Observe diff is large, stronger evidence diff is real

if sample is large, even small differences suggest it's real
(due to law of large numbers)

If A and B are independent, $P(A \text{ and } B) = P(A) \times P(B)$

$$\text{e.g. } P(T, T) = P(T) \times P(T)$$
$$= \frac{1}{2} \times \frac{1}{2}$$

$$P(\text{!}A) = 1 - \sum \text{all permutations of } A$$

∴

$$P(\text{at least 1 } A) = 1 - P(\text{!}A)$$

Conditional Probability

Wednesday, December 20, 2017 3:19 PM

Conditional Probability

$$P(A|B)$$

Bayes theorem, $P(A|B) = \frac{P(A \cap B)}{P(B)}$

General Product rule:

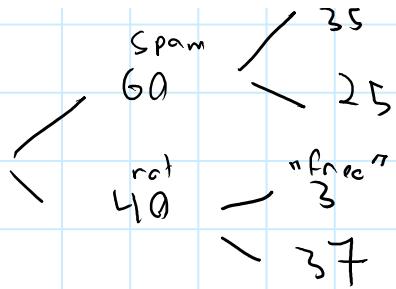
$$P(A \cap B) = P(B) \times P(A|B)$$

Probability Trees

Ex. Have 100 emails, 60 are spam, 40 aren't
of the 60 spam emails, 35 contain the word "free"
of the 40 non-spam, 3 contain the word "free"

$$P(\text{spam} | \text{"free"})$$





$$P(\text{spam} | \text{"freq"}) = \frac{35}{35+5} = 0.92$$

Ex. Swaziland has 25.9% HIV infection

For those who have HIV, ELISA is 99.7% Accurate

For those who don't, it's 92.6% accurate

$P(HIV+ | ELISA+)$

$$P(HIV) = 0.254, P(ELISA+(HIV+)) = 0.997, P(ELISA-(HIV-)) = 0.926$$

A handwritten probability tree diagram illustrating the detection of H⁺ and H⁻ particles. The tree starts with two main branches: "H⁺ V" leading to 0.259, and "H⁻ V" leading to 1 - 0.259 = 0.741.

From each main branch, there are two further branches labeled "EL7SA".

- H⁺ V branch:**
 - EL7SA + leads to 0.997. Below it, "(false)" is written in parentheses. To its right, $P(H^+ \rightarrow E^+) = P(H^+ \wedge E^+)$ is shown, followed by the calculation $0.259 \times 0.997 = 0.2582$, with the result underlined.
 - EL7SA - leads to 0.003. Below it, "(true)" is written in parentheses. To its right, $P(H^+ \rightarrow E^-)$ is shown, followed by the calculation $0.259 \times 0.003 = 0.0008$.
- H⁻ V branch:**
 - EL7SA + (false) leads to 0.074. Below it, "(true)" is written in parentheses. To its right, $P(H^- \rightarrow E^+)$ is shown, followed by the calculation $0.741 \times 0.074 = 0.0548$, with the result underlined.
 - EL7SA - leads to 0.926. Below it, "(false)" is written in parentheses. To its right, $P(H^- \rightarrow E^-)$ is shown, followed by the calculation $0.741 \times 0.926 = 0.6862$.

$$P(H+ | E+) = \frac{P(H+E+)}{P(E+)} = \frac{0.2582}{(0.2582 + 0.0548)} = \underline{\underline{0.825}}$$

Bayesian Inference

Setup: 2 dice, 1 6 sided, 1 12 sided

Can ask to roll dice & will be told if ≥ 4

$$P(\geq 4 | 6s) = \frac{1}{2}, \quad P(\geq 4 | 12s) = \frac{3}{4}$$

Assume can't roll enough times to invoke law of Large Numbers

$$\begin{array}{c}
 \text{12 is} \\
 \diagup \quad \diagdown \\
 \begin{array}{cc}
 \begin{array}{c}
 R \\
 \text{roll} \\
 \geq 4
 \end{array}
 &
 \begin{array}{c}
 P(12_R \geq 4) \\
 \rightarrow 0.375
 \end{array}
 \\[-1ex]
 \begin{array}{c}
 R \\
 0.5
 \end{array}
 &
 \begin{array}{c}
 P(R \leq 4) \\
 \rightarrow 0.125
 \end{array}
 \\[-1ex]
 \begin{array}{c}
 L \\
 \text{roll} \\
 \geq 4
 \end{array}
 &
 \begin{array}{c}
 P(12_L \geq 4) \\
 \rightarrow 0.25
 \end{array}
 \\[-1ex]
 \begin{array}{c}
 L \\
 0.5
 \end{array}
 &
 \begin{array}{c}
 P(12_L \leq 4) \\
 \rightarrow 0.25
 \end{array}
 \end{array}
 \end{array}$$

$$P(12_R | R \geq 4) = \frac{P(12_R \geq 4)}{P(R \geq 4)} = \frac{0.375}{0.375 + 0.25} = 0.6$$

"Posterior probability" $P(H_{\text{posterior}} | \text{data})$

depends on prior probability AND current observation

[Note]: this is different from the probability of observed data being true i.e. $P(D|H)$
Updating the prior:

The prior is updated with posterior probability from previous iteration

Bayesian method allows:

taking advantage of previous studies or physical models

וְנִזְמָן מִלְּפָנֶיךָ וְנִזְמָן מִלְּפָנֶיךָ

taking advantage of previous studies or physical models
Integrate data as it's collected
Avoid using p-values, use posterior probability

Prior is important, but matters less w/ more data

Week 4

Wednesday, January 10, 2018 1:07 PM

68 - 95 - 99, 7 rule

(h. 3, 1, 2, 4)

Exercises

3. 3, 5, 9, 11, 17, 25, 27, 29, 33

Practice quiz

Q 6: $P(\text{identical} | \text{2 rolls})$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A) P(A)}{P(B)}$$

$$P(\text{2 rolls identical}) = 0.5$$

$$P(\text{2 rolls}) = 0.5 \times 0.3 + 0.25 \times 0.7$$

$$P(\text{identical}) = 0.3$$

Qn-2

Q6

$$P(\text{taken} | \text{complain}) = \frac{P(\text{complain} | \text{taken}) P(\text{taken})}{P(\text{complain})}$$

$$P(\text{complain} | \text{taken}) = 1 - 0.39$$

$$P(\text{complain}) = 0.23(1 - 0.39) + (1 - 0.23) \times (1 - 0.27)$$

$$P(\text{taken}) = 0.23$$

Q7 $n = 100, p = 0.28$

$$\mu = 44.8, \sigma = \sqrt{np(1-p)}$$

Q6

$$P(\text{pregnant} | +)$$

$$\frac{P(+|p) p(p)}{p(+)}$$

$$P(p) = 0.055$$

$$P(+|p) = 0.99$$

$$P(\bar{+}|\bar{p}) = 0.995$$

$$\frac{0.99 \times 0.055}{0.99 \times 0.055 + 0.005 \times (1 - 0.055)}$$

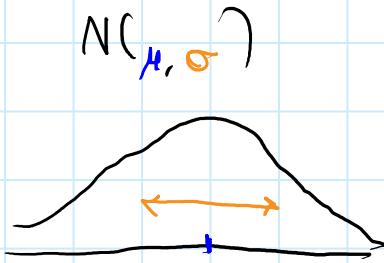
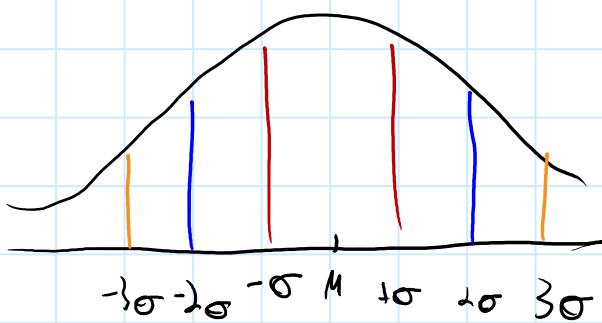
Normal Distribution

Wednesday, January 10, 2018 1:08 PM

Normal distributions

- Unimodal & symmetric

$$\mu = \text{avg}, \sigma = \text{sd}$$



data encompassed

68%

95%

99.7%

Z-score: number of std deviations away from mean

$$Z_x = \frac{x - \mu}{\sigma}$$

Unusual: $|Z| > 2$

Z-scores are not only used with normal distribution, but when the data is normal, a Z-score can be used to calculate Z-scores

I, R

$$\text{pnorm}(val, \mu, \sigma) = \%$$
$$\text{qnorm}(\%, \mu, \sigma) = val$$

~~Ever notice the Normal Distribution~~

Binomial Distribution

Saturday, January 13, 2018 9:55 PM

Milgram Experiment

- measured willingness to obey evil acts (shocking someone)

2 outcomes: Don't shock (success), shock (fail)



Trials w/ binary outcomes are **Bernoulli variables**

Binomial distribution describes P of having exactly k successes in n independent trials w/ p (success) p

$$\frac{\binom{n}{k} * p^k (1-p)^{n-k}}{k! (n-k)!}$$

$\underbrace{\binom{n}{k}}_{\# \text{ scenarios}} * \underbrace{p^k (1-p)^{n-k}}_{\text{choose}}$

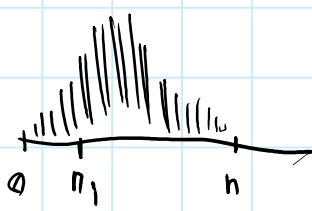
Expected value (mean) of binomial dist $\mu = np$
Std. Dev $\sigma = \sqrt{np(1-p)}$



Normal Approximation

Can treat many samples which resemble normal distribution as such

calculations



$$P(K \geq n_1) = \sum_{i=n_1}^n P(i) \quad \leftarrow \text{approximation}$$

R:
`sum(dbinom(n_1:n, size, p))`

Use $\mu + \sigma$ to find Z-score, get probability

Normal approx may not approx binomial dist well, accounted for by adding or subtracting 0.5 to x , depending, on if its above or below normal curve

When to use normal approx

$$np \geq 10$$

$$n(1-p) \geq 10$$

$$\text{Binomial}(n, p) \sim \text{Normal}(n, p)$$