

## 1.1 Intro & Basic Term

### QRUDAC

#### Scientific Method

- 1) Question
  - define population of interest
- 2) Research
  - what other research has been done?
  - what other variables do you need to measure

#### 5) Analyze

- Stats! {
  - hypo test
  - Conf Intervals

Population: All units of interest

- not all

- 'target' population

Sampling unit: smallest unit selected from pop.

obs. unit: unit actually measured

variable: characteristic of units we want to learn

Sample: subgroup of pop. from which we collect info

notation: 'n' = sample size

typically these two are the same

but if they differ and unintent.  
it means bias

Parameter: summary of a variable for entire pop.  
Statistic: summary of a variable for a sample

e.g. average of student height if we can measure every student  
~~Statistic~~ Statistic: summary of a variable for a sample

Stat. Inference: using sample info to make conclusions about the population

## 1.2 Types of Samples

Randomization: units selected for the sample randomly  
- avoid systematic bias

Representativeness: Sample is similar to the population  
- easier to say results in sample hold true to the pop

Precision: Stats from the sample are close to the true population  
Sampling plan can be reasonably implemented

Feasibility: Sampling plan can be reasonably implemented

Biased Sample: more likely to produce some outcomes over others

- Convenience Samples: too easy
  - Volunteer Response: tend to be very happy or very unhappy
  - Avoid Biased Samples by using Probability / Random Sample
- Avoidability / Random Sample: each member of pop. has a non-zero chance to be in the sample
- identically distr: all values follow same pattern, including mean & variance

Simple Random Sample:

every set of 'n' units has an equal chance of being picked

1) every set of all units: Sampling Frame

2) Compile list of all units: random # generator

3) Select units from list using random process e.g. random # generator

gold std in avoiding bias

Not easy

Stratified Sample: 1) Compile a Sampling Frame

Pop. divided into strata for each Strata

units in strata are similar

2) Take an SRS from each strata

random samples taken of each sampling frame

Cluster Sample:

Units in a cluster are diverse

• Units in a cluster are diverse

• Cheaper to collect; especially when units for a part

units for a part

• Sampling frame needed only from clusters

• Sampling frame needed only from clusters

if cluster units not variable,

sample may not represent population

## 1.3 Bias in Surveys & Samples

Selection Bias: only a subset of people selected to be in sample

Under-Coverage: when sampling frame does not include all of the population

- use a randomized/probability sampling method
- use a complete sampling frame

Non-Sampling Errors & Biases can occur during entry errors or formatting

Non-Response Bias: some part of pop refuses to participate aka not every sampled unit do not become an observational unit

- small response rate may indicate problem

Limit By:

- multiple contact attempts
- monetary reward
- ensure anonymity

Response Bias: responses given differ from the truth

- Results from question
  - ordering
  - working
  - illegal
  - face to face
  - anonymity

Results from question)

Limit:

- randomize order

- studies on wording

## 4 Study Design

Observational Study: ~~not~~ deliberate assignment; only witness current behavior

Lurking Variables: variables that influence response but not explicitly studied

\* evidence of relationship does not imply causation

Response Variable: measures outcome variable

- "Y" - ~~dependent~~ variable

Explanatory Variable: explains changes in response variable

- "X" - independent variable

Treatment: different levels of explanatory variable  
experiments: impose a difference in expl. variable to determine difference in

outcome variable

↳ C.g. 2 methods assigned to 2 different classes  
2 different treatments [But its individual student to measure the response

RP Units: units on which treatment assigned

US Units: units on which treatments applied

## 1.5 Designing Experiments

Randomization: any variable may be randomly assigned

- experimental units are randomly assigned to different levels by a random mechanism
- this allows for all variables equal except for explanatory variable or lurking variables

Replication: when two or more exp. units are assigned to the same unit

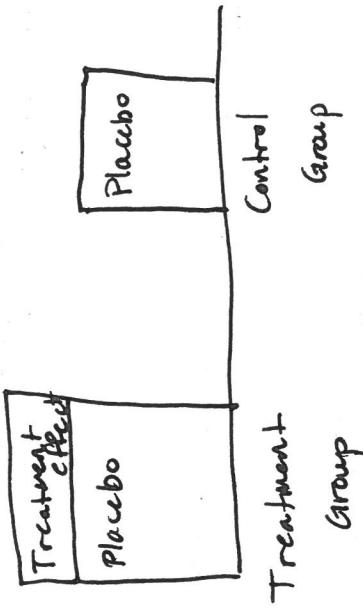
- this allows you to determine if any treatment effect is due to random chance or not; this is done by allowing one to calculate the variance of the experimental error to be compared to the size of the effect

Local Error Control: design techniques used to minimize experimental error; these errors can range from random error from natural variability to systematic bias  
"Control" can also mean to reference the 'control' group that allows one to determine if a treatment is more effective than if nothing had happened

"Blocking": when exp. units are grouped based on some characteristic and ~~treatments~~ different levels of treatments are assigned within each block; this allows for precise measurements by controlling for variation in the response due to blocking variables

Placebo Effect: likelihood of a response regardless of an active 'agent'

- \* the placebo effect is in place for all units that receive any treatment, even if it's the active agent  
because they expect a response ... \*



### Obs. Studies vs. Experiments

- . Obs. studies can still have treatment & control groups
- . In obs. studies, the subject chooses for themselves which group they want to be in
  - as opposed - to the researcher selecting the group

### Random Selection vs. Random Assignment

- RS: selects a random sample to represent population \* Both avoid bias & balances lurking variables between groups to attribute the differences to the treatment

- 1.7 Types of Exp. Designs
- Good design is based in reducing variation  $\downarrow$
  - to detect a treatment effect  $\square$
  - All help
  - Local Err. Control
  - reduce variation
- Completely Randomized Design :
- each unit is randomly assigned to exactly one level of treatment
    - if a lot of variability it can be hard to measure a treatment effect
  - only 2 variables
  - simple
- Matched Pairs Design:
- Explanatory
  - Response
  - Matched Pair v. Block
  - Pair: on an individual level or often 2 units
  - Block: group level; multiple units per level of exp variable
- Matched Pairs Design:
- each unit receives each level of the explanatory variable
  - "unit" in this case could be either
    - a single unit that ~~is measured~~ is measured multiple times (order of treatments should be random!)
    - multiple individual units that have been matched by sex, age, or other
  - usually these pairs  $\uparrow$  will receive treatment and  $\uparrow$  will receive control and will always be compared to each other
- Block Design:
- Reduce variability between Exp Units due to similar characteristics
  - units divided into similar groups called blocks
  - each level of explan. variable is applied in each block
  - similar to Matched Pairs, just bigger scale
  - Matching can be difficult
  - Measurements are not independent between pairs and impacts analysis
  - Reduces variability between Exp units since units in same block are similar
- Assignment of treatments inside blocks
- is random
  - Assignment of treatment inside blocks is random and is by characteristic features of blocks
  - More representative sample
- Assignment of treatments for blocks
- Need to account for variability within blocks since Exp's differ between blocks

## 1.6 Biases in Experiments

Healthcare Effect: people work harder when they ~~know~~ know they're being monitored!

Subject Bias: Subject might want a specific outcome  
- might want to 'please researcher'

Researcher Bias: financial interest  
assign subjects or report in biased manner

Limit by

Single Blinding: 1 party (subject or researcher)  
is not aware of which group a subject is assigned to  
Double Blinding: Neither party knows which treatment group a subject is assigned

### Other issues

Non-response or dropout  
- subjects rarely drop out in a random fashion  
the control group is more likely to drop out than a treatment group if they're very ill and the drug actually helps (thus keeping treatment group and a bigger gap between the two)

### Non-adherence

- take other drugs
- not take all drugs prescribed

### Generalization

- tests college students? Can't apply to adults

## 2.1 Types of Data

- Categorical: places subject into a group aka qualitative
  - e.g. type of car

Subtypes of Categorical:

- Nominal: No ordering (e.g. blue, red ...)
- Ordinal: Ordered (small, med, large)

- Quantitative: numeric (e.g. # of years)
  - only specific values are possible (1, 2, 3, 4, 5...)
  - Discrete: only specific values are possible (.001, .002, .003, ...)
  - Continuous: any value possible (-infinity, infinity)
    - infinite # of possible values

Quant vs. Categ:

- Quant. data that are categorized after can be summarized with both methods ... e.g. ~~Age group 10-20, 20-30, etc.~~
- If data was collected in categories to begin with, it must be treated as categorical
  - e.g. Age = >18, 18-44, 45-64, 65+
    - this would be an ordinal variable
- \* if N/A or unknown is included as a value
  - it then becomes a nominal variable, unless that category is treated as missing

Likert Scale Data: Categorical but often assigned numeric values

- it can be invalid for: shorter scales (eg 5 vs 7 categories) small scale individual items vs. combined items (Likert Scales Data)
- We treat Likert as Categorical in this class

## Graphical Displays:

Categorical: Bar / Mosaic; Piecharts (Counts / %'s of category)

Quantitative: Histogram, Box Plot, Scatterplot, Time series, Heatmap  
Categorical summaries allow comparisons of individuals to the group  
or of a group to other groups

- \* Quantitative numeric summaries describe the main chunk of data w/ mean or median and variability/consistency w/ variation, Std. Dev., range, and IQR's

## 2.2 Summarizing Categorical Data

- Count in Category

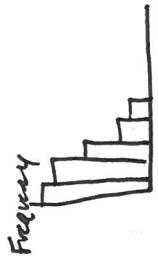
$$- \hat{p} = \frac{\text{Count}}{\text{sample size}} = \frac{y}{n} \quad \text{proportion}$$

- Bar Charts: distorts vol/scale; need baseline of 0
- Pie Charts:
- Display +1 or more category
  - only used if a unit can only be in 1 category
  - Grouped & Stacked & Mosaic stacked bar chart for represent a 2<sup>nd</sup> Categorical variable
  - Mosaic's scale the width of each bar by the sample size by each group
  - Flexibility in bar arrangement e.g. Ordinal / or in order

- Due to flexibility you can't talk about shape of a distr!

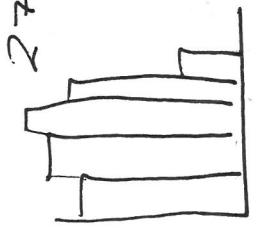
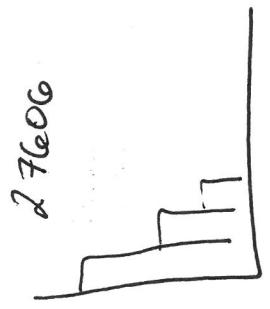
## 2.3 Histograms

- Bars: range of values



- Horizontal axis: values of variable
- Vertical axis: frequency or count of variable
- Big difference between Histograms & Bar charts...
- Can compare across levels of categorical distribution can be viewed in a Histogram

e.g. Sales Prices between Zip Codes



### Distributions:

- 3 major elements to distributions

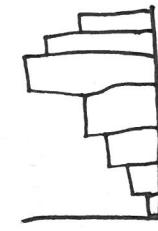
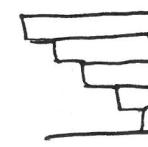
- 1) Shape
- 2) Center (Location)
- 3) Variability (spread)

### Shape

Right Skewed: long tail to the right  
(units stacked to lower limit)  
but unlimited higher values

Left Skewed: long tail to the left  
(units stacked to higher limit)  
but unlimited lower values

Symmetric: tails are approx. even  
major cluster far from limits



### modes: peaks

uniform = 0 peaks  
unimodal = 1 peak  
bimodal = 2 peaks

Outliers: unusual values  
Data Entry: correct if poss  
remove if not

Invalid Point: remove  
Actual/Desired Values: what info do  
they give context to

\* Careful about changing bin width  
too wide: may hide important features  
too narrow: too many features, not enough summary

## 2.4 Boxplots

Visual Display of the 5 number summary

- 50% of data is inside the box

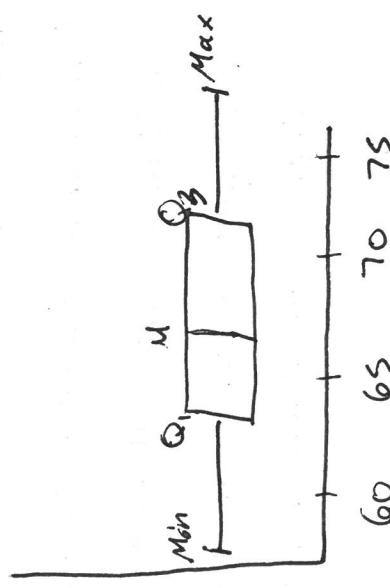
1) Minimum

2) 1<sup>st</sup> Quartile Q1: 25% of values below it

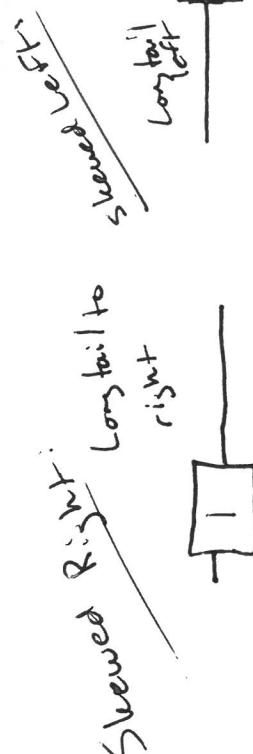
3) Median Q2 : 50% of values below it

4) 3<sup>rd</sup> Quartile : 75% of values below it

5) Max

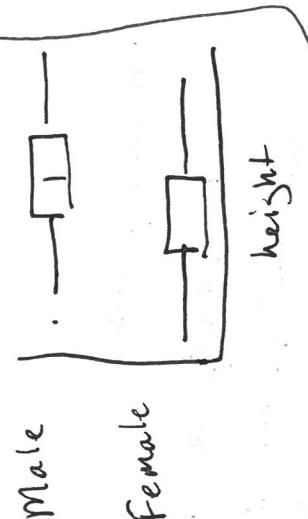


\* Can't determine if multi-modal



Side by Side box plots are a good way to summarize a quantitative variable within levels of a categorical variable

Time plots:



• Scatterplot w/ Horizontal or Time

• Time Series Data: Measurements taken at regular intervals over time

• Watch out for truncated y-axis!

Heat Maps:

• Represent w/ Spatial variable  
• Spatial / Geospatial Data

• Watch for colors that go against convention!

Two-Y-Axes Graphs:

• Putting two variables measured on two different scales

## 2.5 Numeric Summaries of Quantitative Data

### Measures of Central Tendency:

Mean, Median

focused on main chunk of data

### Measures of Variability:

Variation, Std. Dev., Range, IQR

focused on consistency

Mean - Average: symmetric data / or median is skewed

sensitive to outliers

"mu" . sensitive to skewed data, will pull towards whichever side the data is skewed to

Pop. Mean:  $\mu$

$$\text{Sample Mean: } \bar{y} = \frac{\sum y_i}{n} = \frac{1}{n} \sum y_i$$

"the expected value"

- "typical value" - but not always the most common.
- commonly used

Median: good for skewed or symmetric

- good for skewed or symmetric
- 50% of — are above the median
- middle value in a data set
- middle values? take their average
- 2 middle values?
- resistant to skew & outliers

. Variance / Std Dev sensitive to skew

. Std. Dev more often reported

. Variance / Std Dev sensitive to outliers

. Std. Dev more often reported

### Measures of Variability:

sensitive to outliers/skew

Range: Max - Min (spread of entire data set) quick looks

### Interquartile Range (IQR):

$Q_3 - Q_1$ : spread of middle 50% (length of box of IQR)

good for skewed or symmetric

$Q_1$ : median of lower half

$Q_3$ : median of upper half

· quick look

· summarizes distance between all points & the mean

Variance: summarizes distance between points and their mean

Pop. Variance:  $\sigma^2$

$$\text{Sample Variation: } s^2 = \frac{\sum (y_i - \bar{y})^2}{n-1}$$

- "Avg distance between points and their mean"
- or how far, on avg, values are from the mean"



## 2.6 Data Transformations

Adding or Subtracting: changes only measures of center

- measures of center move by same amounts (mean, median)
- measures of variability do not change ( $\sigma$ ,  $\sigma^2$ )

Multiplying or Dividing: changes both

- . mean, median,  $\sigma$  are multiplied & divided by the same amount



### 3.1 Random Variables & Probability

Random Variable: any characteristic whose value can change or vary between individuals

- are always quantitative

- for categorical, we assign numbers to categories  
(e.g. 1 for tails, 0 for heads or a coin flip)
- denoted by capital letters ( $X, Y, Z$ )
- specific values by lowercase ( $x, y, z$ )

- $P(X=x)$ : "probability that random variable  $X$  equals the value  $x$ "

$$\text{Probability: } P(A) = \frac{\# \text{ of ways } A \text{ can occur}}{\text{total # of possible outcomes}}$$

- must always be between 0 and 1

Rules of Prob:

- Disjoint Addition Rule: if two events are disjoint/mutually exclusive  
or do not occur at the same time

$$P(A \text{ or } B) = P(A) + P(B)$$

- Multiplication for Independent: if two events do not impact the other's occurrence

$$P(A \text{ and } B) = P(A) P(B)$$

Compliment of A: A and compliment of A are disjoint and means ...

$$P(A) = 1 - P(A^c)$$

Types of Variables:

Discrete : outcomes are finite (0, 1, 2)

Continuous : any value in an interval (height, weight, etc.)

## 3.2 Distributions

### Major Elements:

- Skew: Skew v. Symm, Modes, Outliers
- Center: Mean, Median
- Variability: Measured by Variance,  $\sigma$ , IQR

Probability Mass Function: how likely a value of a discrete RV is to occur

- plug in specific values to calculate probabilities
- Some PMF's are indexed by parameters
  - 1) Probabilities can never be negative
  - 2) sum of  $P(X=x)$  must = 1
- each possible value of a parameter defines a different pmf

$X \sim \text{Bernoulli}(\alpha)$

"the rv  $X$  follows the Bernoulli Dist w/ parameter  $\alpha$ "

$$P(X=x) = \begin{cases} \alpha & \text{if } x=1 \\ 1-\alpha & \text{if } x=0 \end{cases} \quad \text{for } 0 \leq \alpha \leq 1$$

e.g.  $\alpha = .5$   $P(X=x) = \begin{cases} .5 & \text{if } x=1 \\ .5 & \text{if } x=0 \end{cases}$

## Distributions for Continuous RV's

---

Probability Density Functions (pdf):

- $f(x)$  is just a function that defines shape
- $P(a \leq X \leq b) = \int_a^b f(x)dx$  needs to be integrated to find probabilities

Properties:  $f(x)$  is non-negative

= valid pdf

$$\int f(x)dx = 1$$

$$P(0 \leq x \leq 1) = \int_0^1 3x^2 dx$$

otherwise

$$\int_{.2}^{.5} 3x^2 dx = x^3 \Big|_{.2}^{.5} = .117 = 11.7\%$$

#### 4.1 Famous Discrete Distr. - ~~Bernoulli~~: Binomial

Ex Multiple choice test has 20 questions

a. probability of answering all correctly?

Multiplicative/Independent Rule from 3.1

$$\begin{bmatrix} 20 \\ p \end{bmatrix}^* \text{ where } p \text{ is the probability a student answers a question correctly}$$

b. probability of answering all incorrectly?

Compliment Rule of 3.1

$$(1-p)^{20}$$

c. 18 out of 20?

$$\text{To capture all the different combinations of answers we need the following:}$$

$$\binom{n}{x} : \text{reads "n choose } x" = \frac{n!}{x!(n-x)!}$$

$$\binom{20}{18} : \frac{20!}{18!2!} = \boxed{190}$$

So now we know there are 190 different ways of completing the test w 18 out of 20 correct answers... but what's the probability of getting 18/20 correct for any one way?

Probability of getting 18 correct:  $p^{18}$

Probability of getting 2 incorrect:  $(1-p)^2$

These events are independent, thus we multiply by the number of different ways you can do it:

Now we multiply the probability of answering the test to get a combination of 18/20 corr by the number of different ways you can do it:  $(190)(p^{18})(1-p)^2$

thus, if  $p$  represents the probability of getting 18/20 correct ...

$$P = .80 \text{ and we apply that to our formula } (190) (.8)^{18} (.2)^2 = \boxed{.1374}$$

thus, we have the Binomial Distn

$$PMF: P(X=x) = \binom{n}{x} p^x (1-p)^{n-x} \quad \text{Mean: } np \quad \mathbb{E}[X] = x(1)p^x(1-p)^{n-x}$$

$$\text{Variance: } np(1-p) \leq (x-np)^2 (1-p)^x (1-p)^{n-x}$$

$x$  can take integer values from  $n$  to  $0$

- $n = \#$  of trials
- $p$  = probability of success (must be same for each trial or PMF fails)
- $P$  = probability of failure

Conditions

- 1) Fixed # of trials ( $n$ )
- 2) Only 2 outcomes
- 3) Same prob. of success for each trial
- 4) Outcomes of trials must be independent

Ex1

$$p = \frac{\text{prob of getting heads}}{\text{prob of getting tails}} = \frac{1}{3}$$

$$P = \frac{70}{100} \text{ of getting heads 70 times} \quad \binom{100}{70} \cdot \frac{1}{3}^7 \cdot \left(\frac{2}{3}\right)^{100-70} = ?$$

$$(1-p)^{30} \text{ of getting tails 30 times}$$

$$P(X=70) = \binom{100}{70} \cdot \left(\frac{1}{3}\right)^7 \cdot \left(\frac{2}{3}\right)^{100-70}$$

$$n=100 \text{ thus, } P(X=70) = \binom{100}{70} \cdot \left(\frac{1}{3}\right)^7 \cdot \left(\frac{2}{3}\right)^{100-70}$$

=

$$= \frac{.0000237}{\sqrt{.0000237}}$$

$$1) 7 \text{ heads out of 10} \quad \binom{10}{7} \cdot \left(\frac{1}{3}\right)^7 \cdot \left(\frac{2}{3}\right)^{10-7}$$

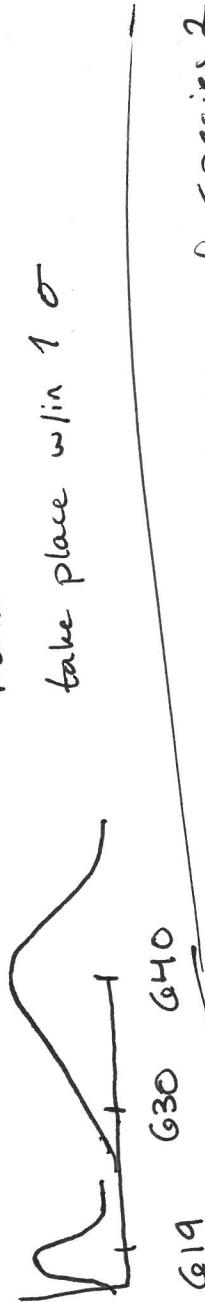
$$n=10 \quad \binom{10}{7} \cdot \left(\frac{1}{3}\right)^7 \cdot \left(\frac{2}{3}\right)^{10-7} = .0009765625$$

ex1 Which species of dinosaur does the skull belong to given the following info?

Species 1: normal distr.  $\mu = 619 \text{ mm}$   $\sigma = 3.4 \text{ mm}$

Species 2: normal distr.  $\mu = 641 \text{ mm}$   $\sigma = 10 \text{ mm}$

\* Remember 68% of all values



Species 2 since the distribution of Species 2 should have a higher frequency than Species 1 due to its larger std. dev & std score is smaller

$$\text{sp}_1: z_1 = \frac{629 - 619}{3.4} = +2.94$$

smaller absolute value

$$\text{sp}_2: z_2 = \frac{629 - 641}{10} = -1.2$$

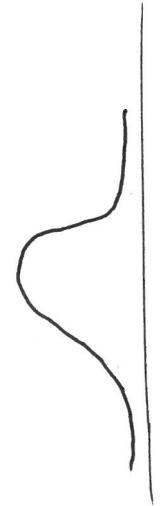
## 4.2 Normal Distr.

aka Gaussian

$$\text{PDF: } f(X; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(X-\mu)^2/2\sigma^2}$$

Notation:  $X \sim N(\mu, \sigma^2)$   
 - Bell shaped, unimodal, symmetric

- $X$  can take any values  $-\infty < X < \infty$
- $\mu$  must be  $-\infty < \mu < \infty$
- $\sigma$  must be  $\geq 0$



- 68% of possible values within 1  $\sigma$  of the  $\mu$
- 95% of possible values w/in 2  $\sigma$  of the  $\mu$
- 99.7% of possible values w/in 3  $\sigma$  of the  $\mu$
- Dist has about 60's overall
- Good guess for min value:  $\mu - 3\sigma$
- Good guess for max value:  $\mu + 3\sigma$

Standard Scores: values from different scales, a uniform measurement is

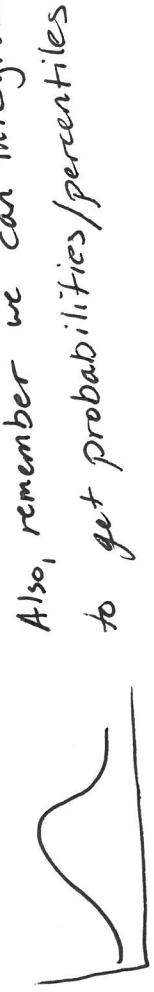
$$\text{To compare values from different scales, a uniform measurement is used in Standard Scores: } Z = \frac{X - \mu}{\sigma}$$

It is also found that if:  
 $X$  is a normal distr. w/  $\mu = 0$  &  $\sigma = 1$   
 then the Std. Score follows a normal distr. w/  $\mu = 0$  &  $\sigma = 1$

called Std. Normal distr.

### 4.3 Probabilities & Percentiles for a Normal Distr.

Recall:  $f(x) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-(x-\mu)^2/2\sigma^2}$  represents a normal distr. curve



Also, remember we can integrate this to get probabilities/percentiles

Recall: People used to standardize each normal distr. using the Std. Sc.

$$z = \frac{x - \mu}{\sigma} \text{ since we didn't have computers}$$

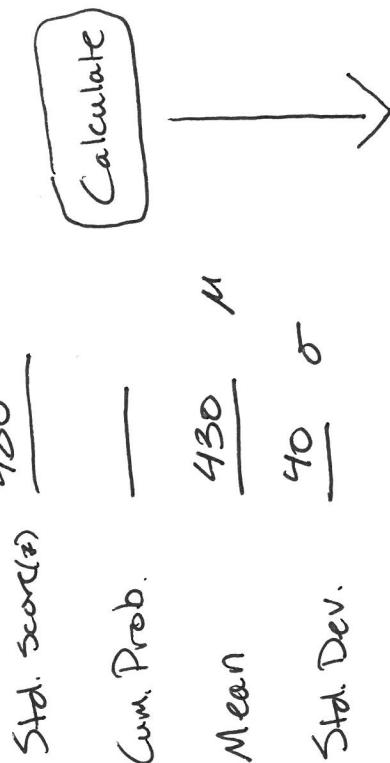
- Probability represents area under a curve

Methods to calc:

SAS: DATA temp; prob\_x = cdf('normal',  $\mu_B$ , mean, std-dev);  
 PROC PRINT; var prob\_x; run;

Online Calc:stattrek.com/online-calculator/normal.aspx

# of interest      E.g. Math test follows a normal dist.  
 Std. score(z)      480  
 Cum. Prob.      —  
 Mean      430  
 Std. Dev.      40       $\sigma$



Cum. Prob = .894
------------------

## Finding Percentiles:

value of a variable such that a specified % is below that value e.g. 75<sup>th</sup> percentile:  $X$  is a value such that 75% of area under the curve is less than  $X$

SAS: \* returns a z-score\* will have to plug value into  $z = \frac{X - \mu}{\sigma}$  for  $X$  to find value

```
DATA temp; z = probit(proportion); x = (z * std-dev) + mean;
```

```
PROC PRINT; var x; run;
```

Eg. what score is needed to be placed in the 90<sup>th</sup> percentile?

Online Calc: [Statattack](#)

Std. Score ( $z$ )

Cum. Prob.

.9

$$\text{Calc } P(X < 481.262) = .90$$

430

Mean

40

Std. Dev

$$[481.262]$$

Cum. Prob :

### 4.3 Examples:

School Scores    a) Prob of a school scoring below 70?

$$\mu = 75$$

$$\sigma = 5$$

$$\text{Online Calc} \quad \text{Std. Score} \quad 70$$

$$\text{Cum. Prob} \quad \underline{\hspace{2cm}} \quad = .159$$

$$\begin{array}{c} \text{Mean} \\ 75 \\ \hline P(x < 70) = 15.9\% \end{array} \quad \begin{array}{c} \text{Std. Dev} \\ 5 \end{array}$$

8

b) Prob of school scoring above 83?

$$\text{Std. Score} \quad \underline{\hspace{2cm}} \quad 83$$

$$\text{Cum. Prob} \quad \underline{\hspace{2cm}} \quad = .945$$

$$\mu \quad \quad \quad 75$$

$$5 \quad P(x > 83) = (1 - .945)$$

$$\quad \quad \quad \underline{\hspace{2cm}} \quad = .055$$

or

$$\underline{\hspace{2cm}} \quad \quad \quad 5.5\%$$

c) What score is required to be in top 20%?

$$\text{Std. Score} \quad \underline{\hspace{2cm}} \quad = 79.208$$

Cum. Prob .80 \* set to 80% since to be in top 20% we need to be ahead of 80% of area under curve

$$\mu \quad \quad \quad 75$$

$$\sigma \quad \quad \quad 5$$

## 4.4 Distributions for Sample Statistics

Sampling Variability: variation in statistics that results from selecting different random samples

Sampling Distr. of that statistic: consists of every value a sample stat could take based on every different possible sample taken

\* As sample size increases, std. dev of sample means decreases

\* As pop. std dev increases, variability in sample means increases

As pop. std dev.

$$\text{Std. Dev of Sample Mean: } \sigma_{\bar{Y}} = \frac{\sigma}{\sqrt{n}}$$

"sigma of  
y-bar"

$$\sigma_{\bar{Y}} = \frac{2.8}{\sqrt{100}} = .28$$

$$\mu = 65 \quad \sigma = 2.8 \text{ in} \quad n = 100$$

↑  
pop. std dev      sample size

Interpretation: the possible values of the sample mean are about .28 in, on average, from the true value of the pop. mean

Right Skewed Pop:

Means are normal when n is large

Bell Shaped Pop:

Means are normal even when n is small

\* Bi-Modal Pop:

Means normal when n > 30

## 4.5 Sampling Distr. of the Mean

### Central Limit Theorem:

of sample means

a distribution can be expected to be well modeled by a normal distr. if the sample size is large ( $n > 30$ ) and

$$\text{have: mean: } \mu = \frac{\sigma}{\sqrt{n}} \cdot \text{SampleMean}$$

- \* note that a distribution of samples will still look like the pop dist and the CLT does not apply

$$\mu = 68 \text{ in} \quad \sigma = 3 \text{ in}$$

Example: height follows normal distr.

a) What is prob they're over 72 in tall?



$$X : 72 \quad 1 - .909 = .091$$

CumProb: .909

Mean: 68

is taller than 72 in

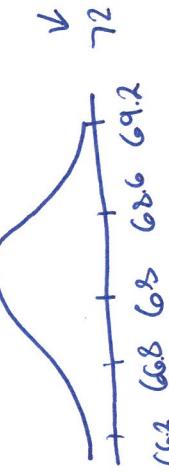
StdDev: 3

b) Select a random sample of 25 from pop<sup>3</sup>. What is the prob sample mean

is > than 72 in?

$$\sigma_{\bar{Y}} = \frac{3}{\sqrt{25}} = .6$$

$$N(68, .6)$$



Almost no chance the average from a sample of 25 people from this pop would

## 1.6 Sampling Dist. of the Sample Prop

$$P = \frac{\text{# of successes}}{\text{total #}} = \frac{Y}{n}$$

↑  
Sample prop

. as sample size ↑, variability ↓

- Std. Dev of the Proportion depends on
  - the population proportion
  - largest at  $p=.5$  because outcome is most uncertain

Std. Dev. of Sample Prop. :

$$\sigma_{\hat{P}} = \sqrt{\frac{p(1-p)}{n}}$$

format  
"z-hat"

Ex] 20% of city water customers are commercial  
n=100 ; What is the std error of the prop in this situation?

$$\sigma_{\hat{P}} = \sqrt{\frac{(0.2)(0.8)}{100}} = .04$$

Interpretation: Possible values of sample proportion are about 4% away from the true value of population proportion, on average

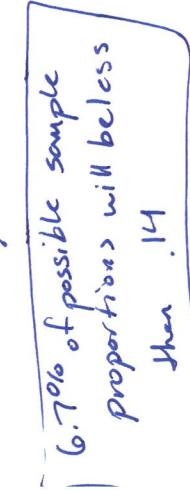
• Normal is good model if  $np \geq 10$  and  $n(1-p) \geq 10$

• Sample must be random and large Ex] What is the probability is the sample proportion is  $< 14\%$ ?

• Shape: approximately normal

• Center:  $\mu_{\hat{P}}$  on p (probability)

6.7% of possible sample proportions will be less than .14



$$\mu_{\hat{P}} = \frac{14}{2} = 0.7$$

$$\text{Mean} = \frac{14}{2} = 0.7$$

$$\text{Std.Dev.} = \sqrt{\frac{0.7(1-0.7)}{2}} = 0.14$$

$$\text{.14} \quad .2 \quad .26 \quad .28$$

## 5.1 Confidence Intervals

### 5.1 Point Estimation

We have gone over several pieces of the scientific process

- ... Step 4: Gathering Data & Step 5: Evaluating Data (sample Stat.)

Now we need to start applying the evaluations and drawing conclusions with statistical inference between the sample and the population. Our first tool to help us is ... (while respecting sampling variability)

Confidence Intervals: a range of reasonable estimates of the true population parameter based on the Sampling Distr.

Hypothesis Tests: provides evidence of a statistically significant effect in a population

Statistical Modeling: (regression analysis) explores the relationship between variables, used for estimation & prediction

the Sampling Dist. & C.I. work together to connect the sample statistic to the best guess on the population parameter which we refer to as the point estimate

recall sample prop:

- random sample • center: pop. mean  $\mu$
- large  $n \geq 30$  • Std. Dev:  $\sigma_Y = \sqrt{\frac{\sigma^2}{n}}$
- Shape: normal

recall sample prop:

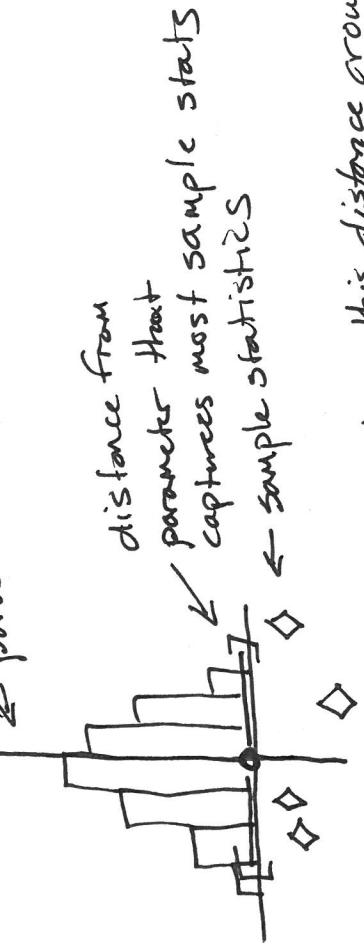
- random
- $n\bar{p} \geq 10$  and  $n(1-\bar{p}) \geq 10$
- center: Pop. Prop  $P$
- Std. Dev:  $\sigma_{\hat{p}} = \sqrt{\frac{P(1-P)}{n}}$



## 5.2 MoE

Margin of Error: numeric indicator of how far a statistic could be from a true parameter

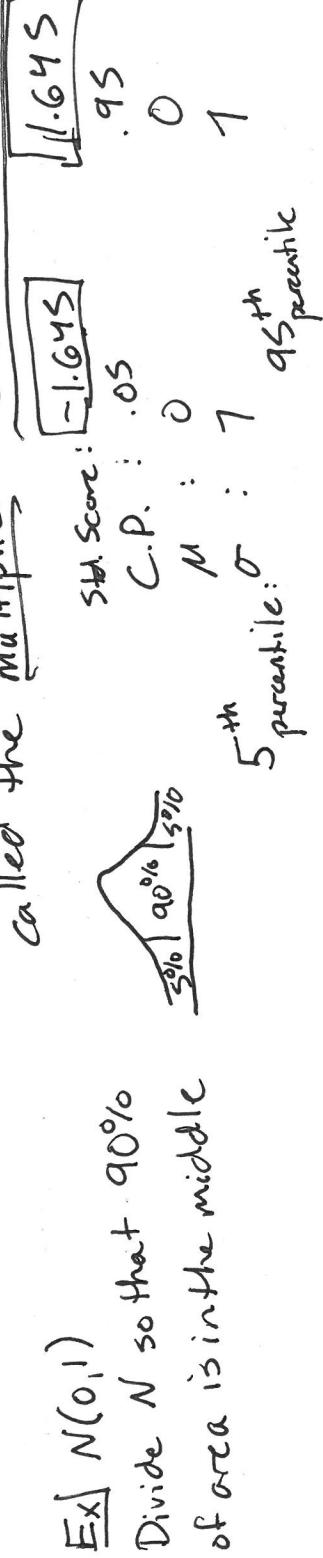
we determine a distance that would capture most values in a sample's distribution, then we place that distance on all the sample distributions, that still capture most of the potential values of a true parameter



\* we place this distance around each sample stat and find it captures most sample stats: this is the MoE  
depends on the std deviations of the statistic

for means:  $\sigma_{\bar{y}} = \frac{\sigma}{\sqrt{n}}$  □ How large should MoE be?

for prop:  $\sigma_p = \sqrt{\frac{p(1-p)}{n}}$ . Depends on the % of the sampling distribution we would like to capture... called the multiplier or confidence coefficient

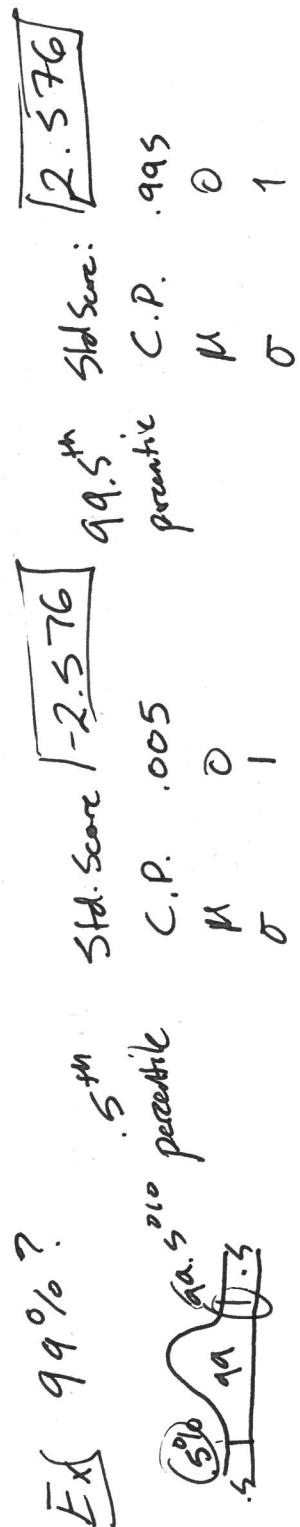


**Ex** Divide  $N(0, 1)$  so that 95% is in the middle  $\boxed{1.960}$

2. 5 <sup>th</sup> std. dev. $\boxed{-1.960}$	99.5%	C.P. 2.5	99.5%	C.P. 97.5
percentile $\mu$ 0	percentile $\mu$ 0	percentile $\mu$ 0	percentile $\mu$ 0	percentile $\mu$ 1
$\sigma$ 1	$\sigma$ 1	$\sigma$ 1	$\sigma$ 1	$\sigma$ 1

2.5 97.5

**Ex** 99%?



95% in Middle      2. 5<sup>th</sup> / 97.5<sup>th</sup> percentile

$$N(0, 1) = 1.96$$

99% in middle

99% in middle

.5<sup>th</sup> / 99.5<sup>th</sup> percentile

$$N(0, 1) = 2.574$$

$$\begin{aligned} & \text{MoE : Multiplier} \times \text{std. Deviation} \quad (1.96 \times \sigma_{\bar{Y}}) \\ & \text{MoE : Multiplier} \times \text{std. Deviation} \quad (2.576 \times \sigma_{\bar{Y}}) \end{aligned}$$

\* Std. Dev. of the statistic aka std. error

sample std.

$S = \text{std deviation}$

Mean:

$\frac{\alpha}{2} \%$



$\sigma$  known:  $Z^* \frac{S}{\sqrt{n}}$

$\sigma$  unknown:  $t^* \frac{S}{\sqrt{n}}$

- found using

-  $\frac{\alpha}{2}$  found by calculating tails

$$\text{multiplier} = Z^* \frac{\sqrt{p(1-p)}}{\sqrt{n}}$$

### 5.3 The t-distribution

The t-dist. is used because we rarely know the true pop. std dev. as is used in the MoE formula

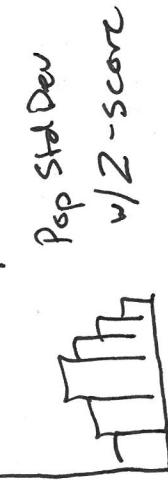
for means

$$2 * \frac{\sigma}{\sqrt{n}} \leftarrow \text{pop std dev}$$

Our solution was to use ' $\sigma'$ , the sample std dev ... but that

is a 2nd estimation to accommodate this 2nd estimation, that isn't accounted for let's see how big the difference is ~~is~~ in the first MoE formula let's see how big the difference is in the std. scores using Pop Std Dev & in the std. scores using Sample Std. Dev

$$y - \mu / \sigma$$



w/ Z-score

$$y - \mu / s$$

sample std Dev

A normal distribution curve with a vertical line drawn through the center labeled "sample std Dev".

w/ Z-score

Most values are between -4 and +4 std deviations  
Most values are between -3 and +3 std. deviations

thus we can see there is more variability in the distribution of std. scores of the sample std. Dev... which means the normal distribution doesn't capture it well .. we need a new distribution ... the t-distr.

## t-distr:

- commonly used for conducting inference for a pop mean
- models additional uncertainty

- bell shaped

uni-modal

$$\text{PDF: } f(x; v) = \frac{\Gamma\left(\frac{v+1}{2}\right)}{\sqrt{v\pi} \Gamma\left(\frac{v}{2}\right)} \left(1 + \frac{x^2}{v}\right)^{-\frac{v+1}{2}}$$

- has "gamma" coefficient:  $\Gamma$  we use tech to help us find the percentiles since
- depends on degrees of freedom (df)  $\Gamma$  is hard to solve/integrate
- $v = df (n-1)$
- means the # of values in calculation of a statistic that means the n-1 values allowed to vary are free to vary if you know the first n-1 data values logic:

Ex  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$  and the sample mean, then the last data value can be calculated by...

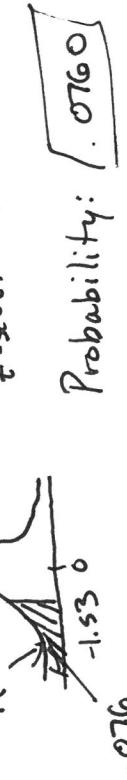
$$y_n = n\bar{y} - \sum_{i=1}^{n-1} y_i \quad \dots \quad \text{thus } n-1 \text{ values allowed to vary}$$

- \* as df grows larger the t-distr becomes more like the normal distr  $[N(0, 1)]$
- as df  $\rightarrow \infty$ , t-distr converges to  $N(0, 1)$
- Ex For t-dist w/ df = 12, find  $P(t < -1.53)$

- mean of t-distr: 0

Variance:  $\frac{v}{v+2}$

df: 12  
t-score: -1.53



Probability: 0.076

### Finding percentiles w/ t-distr

rv t-score

Ex]  $t(\alpha)$ , multiplier for 90% confidence  
t( $\alpha$ ) 5th P( $t < ?$ ) = .05  
q5th  
90th  
 $-1.895 \quad 1.895$  \* report positive value so

df: 7  
P( $t \leq t$ ): .05  
 $t(\alpha) = -1.895$

Ex]  $t(\alpha)$ , multiplier for 90% confidence  
df: 7  
P( $t \leq t$ ): .05  
 $t(\alpha) = -1.895$

