

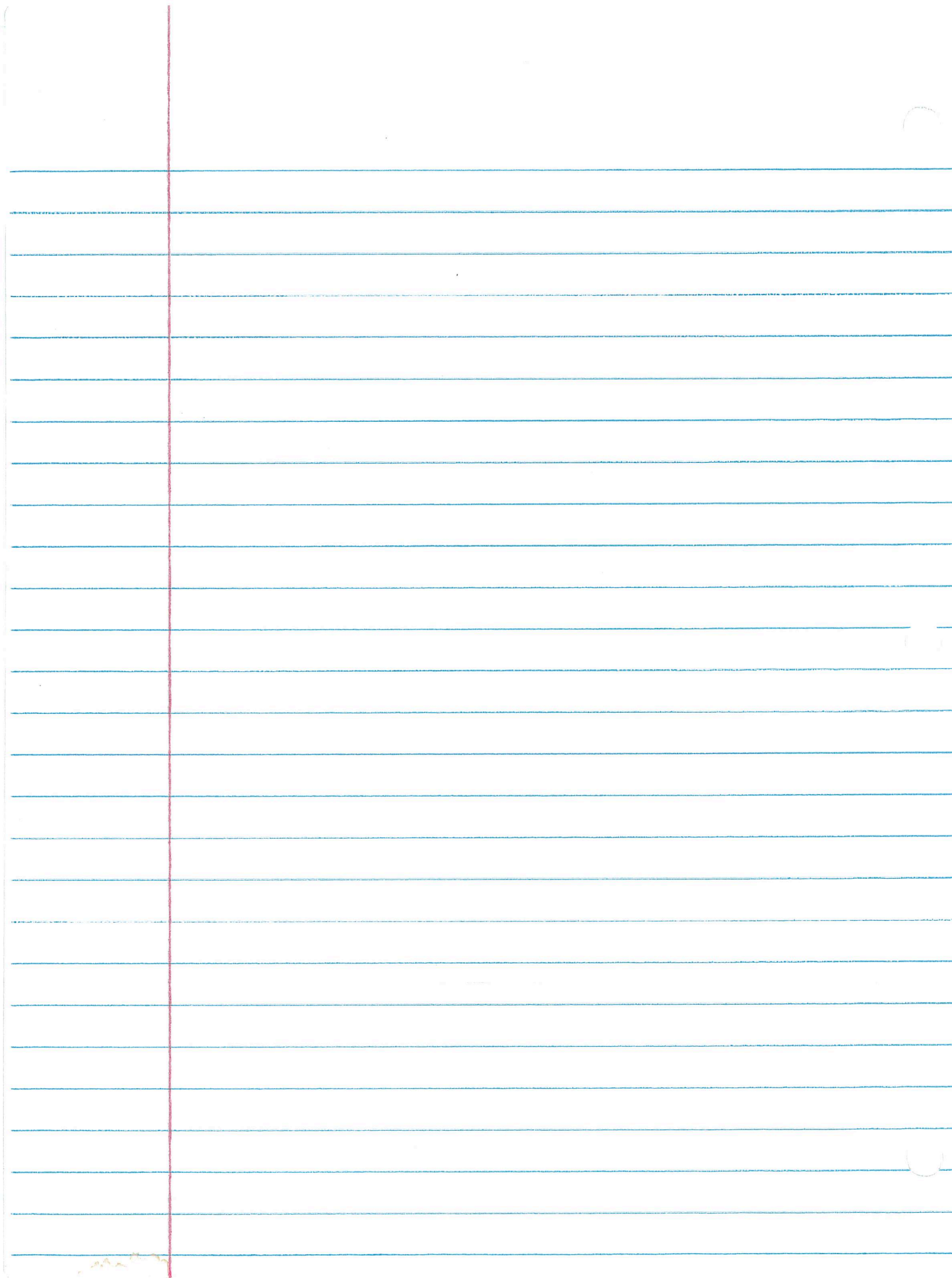
In A

S T A I S T I S T I S

nut-shell



By Annie ♡





# Unit 1 - Different distributions and how to describe them

## Populations and Samples

- Parameters are facts about population.  $\sigma$  and  $M_x$
- Statistics are drawn from samples; do tests:  $\bar{x}$  and  $s_x$  from samples to infer information about population.

## Numerical data

- numbers

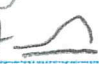

Discrete - whole #'s

continuous - #'s with fractions

Organize: dot plots, box + whisker plot, stem + leaf, histograms, cumulative.

## Describing Distributions

- shape: symmetric, normal,

skew left  skew right 

- Center: median - more robust

mean, IQR:  $Q_3 - Q_1$

- Variability: range, standard deviation

- Unusual features: outliers, clusters, gaps.

- Context: title, x and y axis, answer with it.

## Categorical Data

- nouns or proportions  
ex) 60% brown hair

Open - unlimited options

closed - limited set of options

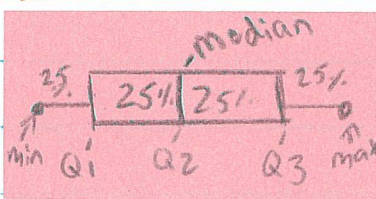
Organize: bar graphs, frequency tables.

## Calculator

For mean,  $S_x$ ,  $Q_3$ ,  $Q_1$ , etc:

Stat, edit, fill list, Calc, Vars<sup>1</sup>

## Box Plots ♥



• 25% each section

• Median Center

## Outliers:

$IQR \cdot 1.5 = \text{fence}$

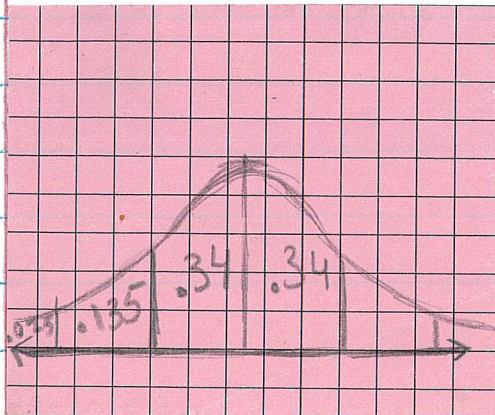
$Q_3 - Q_1$

$Q_1 - \text{fence}$

and  $Q_3 + \text{fence}$

If # is outside that, it's an outlier.





## Normal Graphs

- Bell Curve
- Mean in middle
- Symmetric
- unimodal
- Empirical rule: .34, .135, .025

## Z-Scores and Percentiles

Z-Score: How many standard deviations a value is away from the mean.

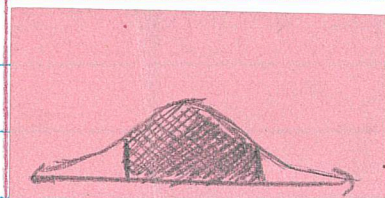
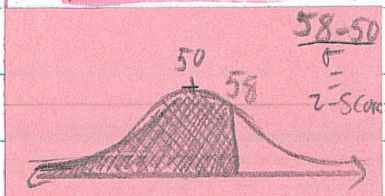
$$z = \frac{\text{obs} - \text{exp}}{\sigma}$$

• Use to find the percent of the data below or above the value. (or in between values)

Convert to % to find how much data is below/above it.

Calc: 2nd vars, normcdf, enter -10000 as lower if trying to find everything below point, 10000 as upper if trying to find all data above.

\* Draw Pictures



or... When finding % data between two values do the percent for both through z-scores, then subtract to get what you need.

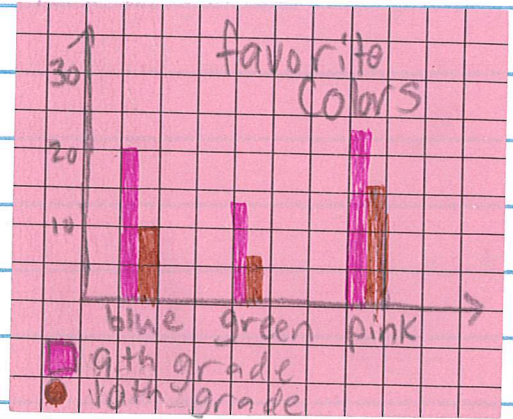
or... Do 1 - combined value percents to find everything less than one value but greater than other value.



## Unit 2 - Scatter plots, LSRL, Relative frequency, data graphs

### Side by Side Graphs

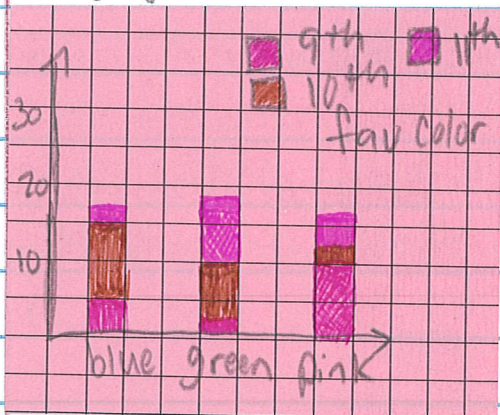
- Categorical, # values
- labels on sides
- next to each other
- swappable variables



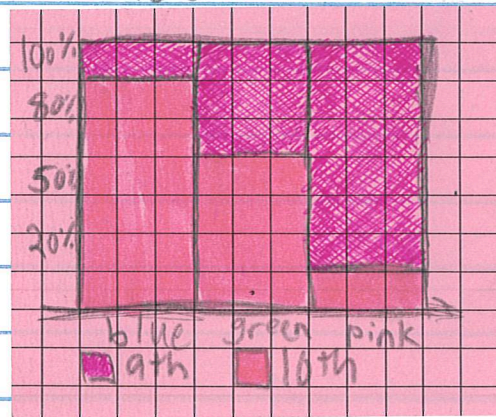
### Segmented

- Swappable Variables
- 1st var on bottom, second as color on side

Segmented:



Mosaic:



Mosaic:

- Uses relative frequencies to better compare variables.

- totals are not the same

beware! The

### Relative Frequency

- Also called conditional frequency

proportion of values may be higher but the real value can be different due to different totals.

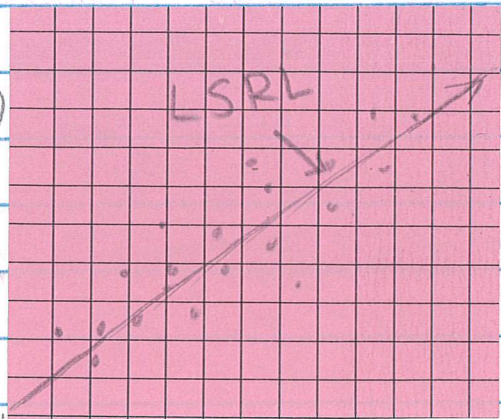
- Probability of one event given another event must occur
- Value / column or row total
- In rel fre table, values must add to the totals for rows and columns.



## Scatterplots

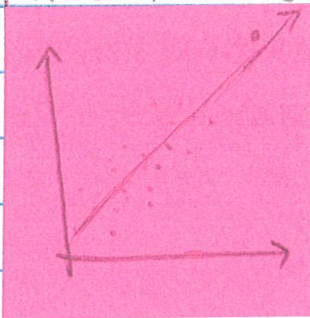
### Describe:

- Direction (pos or neg)
- unusual features: outliers, clusters, gaps
- linear, non linear
- strength: weak, med, strong. + Context



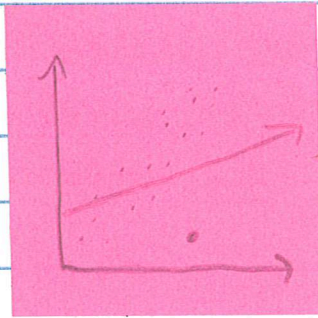
## Outliers

### High leverage point



- far from domain
- doesn't affect residuals

### Influential point



- In domain
- outside range
- very influential on residuals
- far from LSRL

- close to LSRL
- increases variance

• A point can be both influential and high leverage

## Correlation Coefficient ( $r$ ) -

How closely the scatterplot fits the least<sup>2</sup> Regline. From -1 to 1. -1 means perfect negative slope, 0 is no correlation, and 1 is perfect positive slope.

Coefficient of Determination ( $r^2$ ) - What % of the response variable is controlled by explanatory variable. Higher means X-var has more control over Y-var. Smaller than  $r$ .

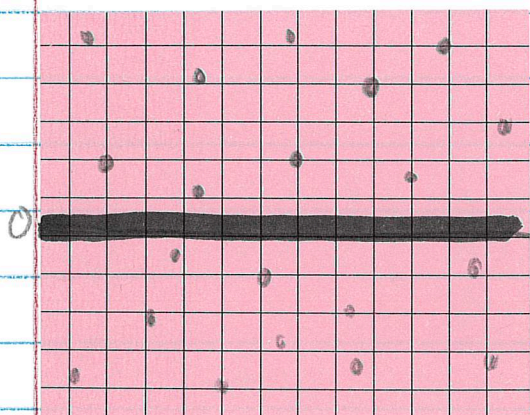
Residuals - How far observed value is from expected.

Expected = plug in X value to LSRL. Ex)  $y = 5.4(5) - 2$ . Positive means underestimated, negative is overestimated. 0 = perfect.



## Residual Graph

← linear model



- On x-axis is the observed value
- On y-axis is the residual value.
- A random pattern means linear is a good fit.
- $S_{res}$  means the residuals are spread out by standard deviation of observed and residual values.

### Calculator functions

How to get residuals, expected,  $r^2$ ,  $r$ , LSRL,  $S_{res}$ :

1. Enter x-values into  $L_1$ , y-values into  $L_2$
2. To get LSRL,  $r^2$ ,  $r$ : Stat, Calc, #4 Lin Reg
3. For expected values: go to  $L_3$ , enter  $y = m(L_1) + b$  or the LSRL.
4. For residual values: go to  $L_4$ , enter  $L_2 - L_3$ , or do obs-exp
5. For  $S_{res}$ : Calc, 1-var stats on  $L_4$ ,  $S_x$  = Standard deviation of Res.

♥ Use this data to find out if the expected values are accurate or not, and if linear is a good fit.



# Unit 3

♥ Essentially what type of study do you use in scenario.

## Population vs. Sample

♥ A sample or survey can measure part of the population to conclude something about the population. When looking for answer about what a study proves ... **taking a sample of population infers something about that specific population.** When a survey covers whole population, it is fact and not something we are inferring. **Samples! Only! generalize to population they came from.**

## Causation vs. Correlation

♥ Causation: is when one event IS the result of another event. Proved that one var controls another by **experiment.**

♥ Correlation: there is an observed correlation between variables through samples of **observational study.**

### Experiments Design

**Experimental units** - Subjects given the treatments. Manipulated variable

**factors** - What experiment manipulates

**Explanatory var** - Var manipulated

**Response var** - outcome of treatments

**Confounding Vars** - Vars related to

exp var that can affect resp var. (Sometimes falsely)

### Inferences through obs. study

♥ **Inferences** are made on data collected, and deemed to be **Significant** or **unSignificant**.

**Significant** - able to infer some fact about data.

**IDEA:** Random + bias free samples = inference for larger population.



## types of studies

• Experimental - assign treatment different ways:

• Randomly - randomly assign treatments to subjects. Good = repeatable, a control group, <sup>limited</sup> Confounding Vars.

• Block Designs - matches homogenous features to eliminate Conf Vars

• Matched Pairs - matches close characteristics in pairs to eliminate possible Confounding Vars. Each subject in pair randomly assigned to either the experiment or control group.

ex) two arms, one with sunscreen, one without.

• Randomized Block Design - Researcher divides subjects into homogeneous blocks, treatments randomly assigned to subjects in blocks  
ex) Divide subjects into blocks with same skin color, randomly assign sunscreen vs. none.

• Use experiments to infer causation

• Observational study

• Random - HAT METHOD! Put names in hat, shake, draw, collect samples.

• Use the data to do different tests:

• Z-test uses proportions of categorical data, t-test uses #'s of numerical data. One sample for 1 var, 2 sample for different

• Use observational studies to infer correlation. } in variables.  
Based on if data is significant from expected.

## Unit 4

### Different types of probability

- **Unions** - Either event happening, add probs.

$P(A \cup B)$  → probability A or B happens

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) \leftarrow \text{to make sure you don't double count prob of A and B occurring simultaneously.}$$

- **Intersections (and)** - Both events happening at same time. Multiply. Only works if Independent.

$P(A \cap B)$  → prob A and B both happen

$$P(A \cap B) = P(A) \cdot P(B)$$

- **With Replacement** -

When drawing from group keep denominator same, and keep prob of certain event same:

$$\frac{7}{12} \cdot \frac{7}{12}$$

blue ball drawn → 2nd blue ball drawn

- **Without Replacement** -

Change denominator and prob of certain event if drawn.

$$\text{ex) } \frac{4}{12} \cdot \frac{3}{11} \cdot \frac{2}{10}$$

- **Conditional probability** - probability one event occurs given another event occurred.

$P(A|B)$  → prob A given B

divide prob  $A \cap B$  by prob of B.

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

- **Independent vs. Dependent**

Independent $P(A) = P(A B)$ dependent $P(A) \neq P(A B)$	} Disjoint Events: $P(A \cap B) = 0$ cause mutually exclusive.

- **Mutually exclusive** - 2 events can't occur at same time inclusive means can happen together.



## Geometric Distribution

- Solves for # of trials needed to obtain success.
- $(1-p)^{x-1} (p)$   $p$  = probability of success.
- Expected values =  $\frac{1}{p}$       $\sigma = \frac{\sqrt{1-p}}{p}$
- to easily see all values do  $y = (1-p)^{x-1} (p)$

## Binomial Distribution

- To find how many successes in a set # of trials.
- $nCr (p)^r (1-p)^{n-r}$      • Also set equal to  $y$  to view all trials

# trials  $\downarrow$   $x$      Prob Success  $\downarrow$

•  $M_x = n \cdot p$

• Variance:  $\sigma_x = \sqrt{n \cdot p(1-p)}$

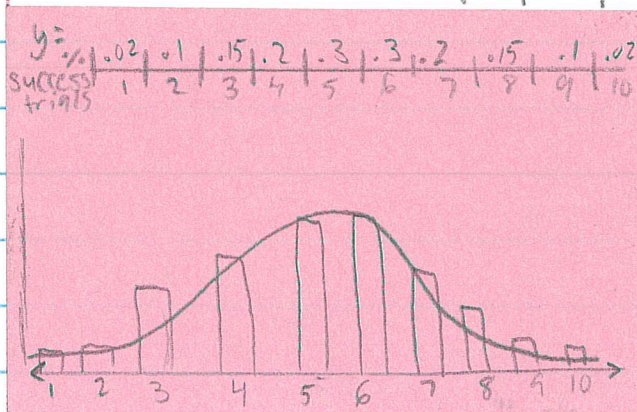
After finding Mean and  $S_x$

you can set up graph and table to answer questions.

If finding values greater/less than certain value, just find which values apply in the chart and add them up.

$M_x$  = how we solve expected value

$\sigma_x = \sqrt{n \cdot p(1-p)}$  so  $\sqrt{10 \cdot p(1-p)}$   
trials



## Circles Combinations

- total # of circle arrangements without restrictions =  $\frac{n!}{n}$
- ex) 7 people, how many ways to arrange:  $\frac{7!}{7} = 720$
- If people must sit together, take the # that must sit together! times the normal equation but take away one from  $n$  for every new person that must be together. Ex) 3 people together:  $3! \times \frac{5!}{5}$

Expected Value - doesn't add to 1 →

$p(x)$	.2	.05	.01	.01
#	10	20	50	100

- Multiply probability by value of events. Often this involves game and how much won from game; Expected value = mean of data

ex)  $.2(10) + .05(20) + .01(50) + .01(100) = .73(0)$

- Make sure the probability adds to 1; if not, see why and adjust. Possibly another factor like % 0, or divide prob by total probability there. ex)  $\frac{.2}{.27} = .74$

Fundamental Counting Principal - how many possible outcomes

flip coin, roll dice →  $2 \cdot 6 = 12$  outcomes

nothing repeats: SWAIN →  $5 \cdot 4 \cdot 3 \cdot 2 \cdot 1$  • Decreases each time

repeats: POOP:  $\frac{4 \cdot 3 \cdot 2 \cdot 1}{2 \cdot 1}$  ← cause 0's are same cause spots fill

• Decreasing #'s can be written as  $5! \rightarrow 5, 4, 3, 2, 1$

Permutations - Order matters, line things up. Calc: Math → prob →  $nPr$

$nPr$   $n$  = Starting # to choose from  $r$  = # you need

ex) line up 6 people from class of 20 →  ${}_{20}P_6$

Combinations - When Order doesn't matter Calc: Math → prob →  $nCr$

$nCr$   $n$  = total to choose from  $r$  = # needed

Combining Variables

- Expected Value / Combined means:  $(P_x \cdot M_x) + (P_y \cdot M_y)$  probability  $\times$  mean  $x$ , + prob  $y$   $\cdot$  mean  $y$
- for Combined Standard deviation:  $\sqrt{S_x^2 + S_y^2}$  only add variation

Scaling Variables

• Mean: Just multiply Standard deviation:  $\sqrt{x\sigma^2 + x\sigma^2}$

- for adjusting means and standard deviation scale factor

for Sample, mean stays same, standard dev =  $\frac{\sigma}{\sqrt{n}}$



# Unit 5 ♡♡

## Normal Distribution

• normal: symmetric, mounded in middle, bell curve, from a normal population or sample size  $> 30$ .

- ★ Keep discrete graphs and continuous graphs separated.  
If discrete, add up probability values to find what's needed.  
If continuous with fractions, use z-score. Real world usually continuous.

## Central limiting theorem

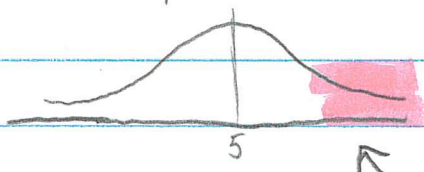
• Law of large #'s - More samples run = closer to actual value. This means that eventually  $\bar{M}_x = \mu$   
mean samples      true mean

Essentially:  $\bar{X}_x = \mu$  over repeated samples  
center samp. dist      pop center

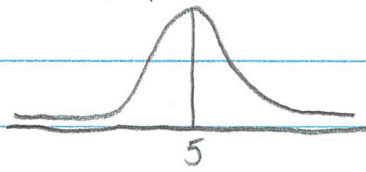
**Numerical**  $\sigma_x = \frac{\sigma}{\sqrt{n}}$  over repeated samples because  
as you increase the sample size, the variance will decrease due to central limiting theorem.  
n = sample size

Therefore... repeated samples = same center, less variation so higher mounded data. (less spread out)

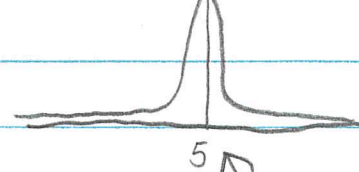
ex) 20 samples:



100 samples:



1000 samples:



Also note...  
less samples = larger tail ends  
more samples = mounded in middle  
more samples = closer to true mean

# Categorical

$M\hat{p} = M_p$  Center population = center samp dist.

$$\sigma_s = \sqrt{\frac{p(1-p)}{n}}$$

← probability  $p$   
← sample size  $n$

## Conditions

numerical	Categorical
random	random
Independent - small enough? $n < .1N$	Independent $n < .1N$
Normal - big enough? $n \geq 30$	Normal $n \cdot p \geq 10$ and $n(1-p) \geq 10$

## • Bias check •

- $M\hat{p} \pm M_p$ . If more than 2 standard deviations away from parameter, the samples are likely biased.  
(If z-score of  $M\hat{p}$  is greater than  $\pm 2$ , it is likely biased.)

**POINT ESTIMATOR** - Statistic gathered from sample that's used to estimate parameter of population. Check for bias and conditions to see if statistic is good point estimator.

## • Sampling Distribution Differences •

### Numerical

$$M\hat{p}_1 - M\hat{p}_2 = p_1 - p_2$$

$$\sigma_{p_1} - \sigma_{p_2} = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$$

- Use to find differences in means or SD from different groups.

### Categorical

$$M\bar{x}_1 - M\bar{x}_2 = M_1 - M_2$$

$$\sigma_{x_1} - \sigma_{x_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

- Use to find differences in proportions from different groups.