

Numerical (quantitative)		Categorical (qualitative)	
Continuous	Discrete	Regular Categorical	Ordinal
Take on any of an infinite number of values within a given range	Take on one of a specific set of numeric values	Categories that could be represented with numbers but difficult to do arithmetic operations on	Levels have inherent ordering
Height, Distance (Often rounded to look like a whole #)	Floors of a building, number of cars (can't have half a car)	Gender (0 for women, 1 for men)	Least Likely, Likely, Most Likely

Observational Study	Experimental Study
<p>Collect data in a way that does not directly interfere with how the data arises. Only establish an association</p> <p>Two types:</p> <ol style="list-style-type: none"> 1. Retrospective: uses past data 2. Prospective: data collected throughout the study 	<p>Randomly assign subjects to treatments</p> <p>Establish causal connections</p>

Sampling and Bias	
General Terms	<ul style="list-style-type: none"> • Census: entire population • Exploratory analysis • Representative sample • Inference
Sources of Bias	<ol style="list-style-type: none"> 1. Convenience sample - individuals who are easily accessible are more likely to be included in the sample 2. Non-response - only a non-random fraction of the random sampled people respond to the survey (initially a random sample) 3. Voluntary response - sample consist of people who respond because of strong opinion (no initial random sample)
Sampling Method	<ol style="list-style-type: none"> 1. Simple Random Sampling (SRS) - each case equally likely to succeed 2. Stratified Sampling - divided population into homogeneous strata, then randomly sample from each stratum 3. Cluster Sample - divided population into cluster, randomly

	<p>sample a few clusters, then sample all observations within these few clusters</p> <p>4. Multistage Sampling - divided population into cluster, randomly sample a few clusters, then random sample the observations within these few clusters</p>
--	--

Experimental Design	
4 Properties	<ol style="list-style-type: none"> 1. Control - control group 2. Random - equal probability 3. Replicate (collect a large enough sample or replicate the entire study) 4. Block - block variables known or suspected to affect the outcome
Explanatory vs Blocking Variables	Explanatory - conditions we can impose on the experimental units Blocking Variables - characteristics the samples come with that we would like to control for (gender)
More Terms	<ul style="list-style-type: none"> • Placebo/Placebo effect • Blinding (sample don't know which group they're in) • Double Blind (both sample and researcher don't know group assignment)

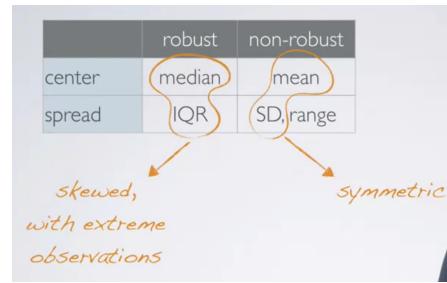
Random Sample Assignment		
<p>Causal - Effect on the experimental group is due to the treatment being applied Generalizable - what's happening to the sample can be generalized for the entire population</p>		
	Random Assignment	No Random Assignment
Random Sampling	Causal/Generalizable (ideal experiment)	Not Causal/Generalizable (most observational studies)
No Random Sampling	Causal/Not Generalizable (most experiments)	Not Causal/Not Generalizable (bad experiment)

Frequently used Visualizations	
Scatter Plots	Trends Identify outliers to observe
Histograms	<p>Data density Shape of distribution:</p> <ul style="list-style-type: none"> - Skewness: Left skewed, right skewed, symmetric - Modality: unimodal, bimodal, uniform, multimodal
Dot Plot	When individual values are of interest, may become busy
Box Plot	Useful for highlighting outliers, media, IQR
Intensity Map	For highlighting spatial distribution (heat map)

Measures of Data	
Measures of Center	<ul style="list-style-type: none"> - Mean, median, mode are measures of distribution. Each have different symbols for measures taken for a sample(x) vs taken for a population (μ) - In terms of skewness, remember that the mean chases the tail.
Measures of Spread (variability)	<ul style="list-style-type: none"> - Range - Variance (the average squared deviation from the mean). We square the variance to get rid of negatives so pos/neg doesn't cancel each other. Also so larger deviations are weighed more heavily than smaller numbers - Standard Deviation: (square root of variance) - Interquartile Range - 75th percentile - 25th percentile

Robust Statistics

- Measures on which extreme observations have little effect



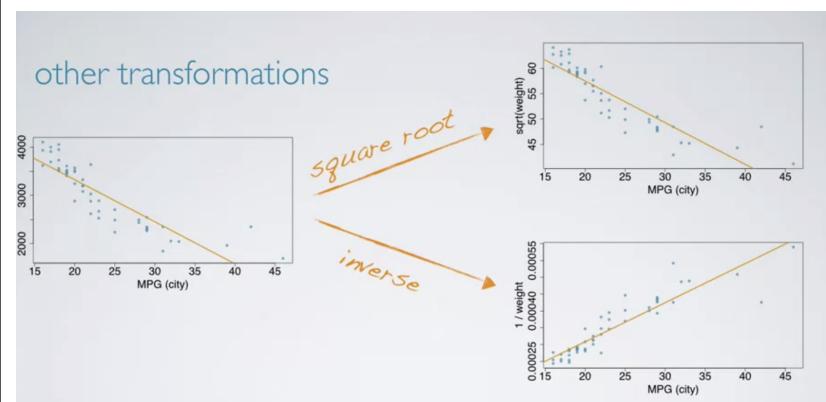
Data Transformation (rescale the data, especially when it is strong skewed)

- See the data structure differently
- Reduce skew in modeling
- Straighten a non-linear relationship

Natural Log Transformation

- Applied when much of the data cluster near zero relative to the larger values in the data set and all observations are positive
- To make relationship between variables more linear

Square root transformation Inverse transformation



Intro to Inference

Null Hypothesis (H_0)

There is nothing going on, no relation

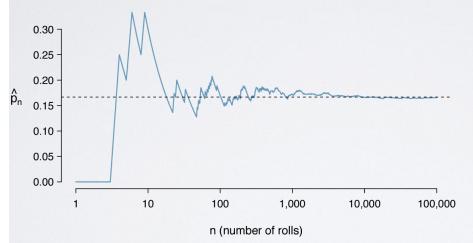
Alternative Hypothesis (H_A)

There is something going on, a relationship. Has the burden of proof to "fail to reject H ")

Process

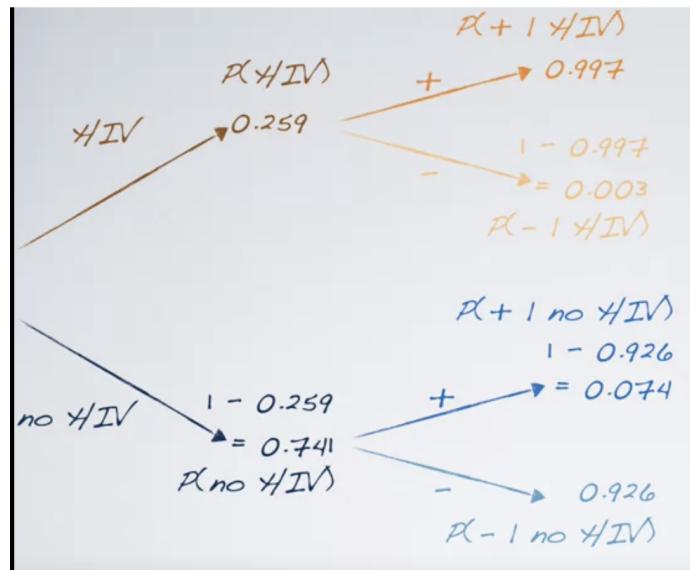
- Set a null and alternative hypothesis.
- Simulate the experiment assuming that the null hypothesis is true
- Evaluate the probability of observing an outcome at

	<p>least as extreme as the one observed in the original data (p-value)</p> <ul style="list-style-type: none"> - If probability is low, reject the null hypothesis
--	--

Probability	
Frequentist Interpretation	<ul style="list-style-type: none"> - Probability of an outcome is the proportion of times the outcome would occur if we observe the random process an infinite number of times
Bayesian Interpretation	<ul style="list-style-type: none"> - Interprets probability as a subjective degree of belief - When presented with new facts, the probability changes (posterior probability)
Law of Large Numbers	<ul style="list-style-type: none"> - As more observations are collected, the proportion of occurrences with a particular outcome converges to the probability of that outcome - ex) if you toss a coin 1000 times, expect heads 500 times than only 5 times 
Gambler's Fallacy	<ul style="list-style-type: none"> - Common misunderstanding of law of large numbers - When an individual erroneously believes that a certain random event is less/more likely to happen based on outcome of previous events - ex) roll heads 3 times, higher chance of head a 4th time.
Disjoint (Mutually Exclusive)	<ul style="list-style-type: none"> - Cannot happen at the same time - ex) cannot have both a head and a tail <div style="border: 1px solid black; padding: 10px; width: fit-content; margin-left: auto; margin-right: 0;"> <p>For disjoint events A and B, $P(A \text{ or } B) = P(A) + P(B)$</p> </div>
Non-disjoint	<ul style="list-style-type: none"> - Events can happen at the same time <div style="border: 1px solid black; padding: 10px; width: fit-content; margin-left: auto; margin-right: 0;"> <p>For non-disjoint events A and B, $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$</p> </div>
Sample Space	<ul style="list-style-type: none"> - Collection of all possible outcomes of a trial

Probability Distribution	<ul style="list-style-type: none"> - Lists all possible outcomes in the sample space and the probabilities with which they occur - Rules: <ol style="list-style-type: none"> 1. Events listed must be disjoint 2. Each probability must be between 0 and 1 3. The probabilities must total to 1
Complementary Events	<ul style="list-style-type: none"> - Two mutually exclusive events whose probabilities add up to 1  - Disjoint vs Complementary - Disjoint are not necessarily complementary but complementary are always disjoint
Independence	<p>Checking for independence:</p> $P(A B) = P(A)$, then A and B are independent. <i>given</i> <p>determining dependence based on sample data</p> <p>observed difference between conditional probabilities \rightarrow dependence \rightarrow hypothesis test</p> <p>if sample size is large, even a small difference can provide strong evidence of a real difference</p>
Baye's Theorem (Conditional Probability)	<p>Probability of A given B</p> <p>Product rule for independent events:</p> $\text{If } A \text{ and } B \text{ are independent, } P(A \text{ and } B) = P(A) \times P(B)$ <p>Bayes' theorem:</p> $P(A B) = \frac{P(A \text{ and } B)}{P(B)}$ <p>General product rule:</p> $P(A \text{ and } B) = P(A B) \times P(B)$

Probability Trees

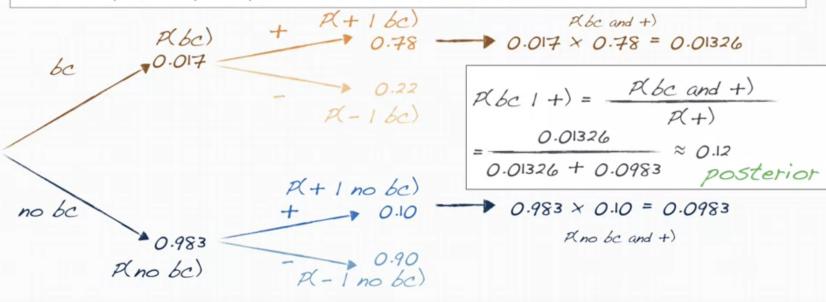


Bayesian Inference

Steps:

1. Setting a prior
2. Collecting Data
3. Obtaining a Posterior
4. Updating the prior with the previous posterior

When a patient goes through breast cancer screening there are two competing claims: patient has cancer and patient doesn't have cancer. If a mammogram yields a positive result, what is the probability that patient has cancer? $P(bc | +)$ = ?



Normal Distribution

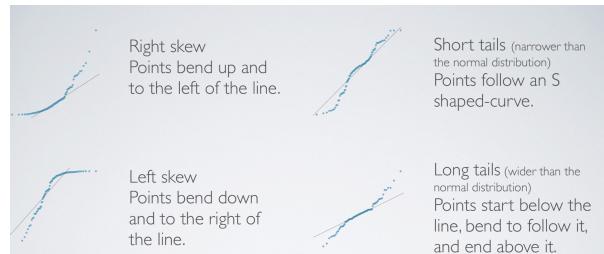
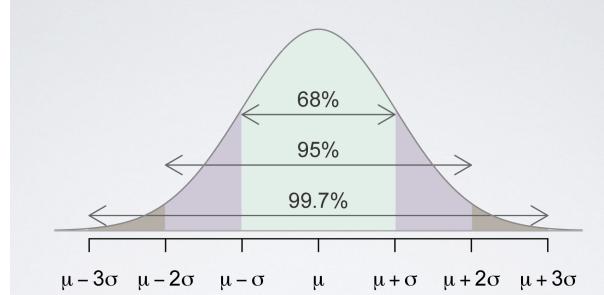
z -score (only under normal curves)

- How many standard deviations above/below the mean

standardizing with Z scores

- ▶ standardized (Z) score of an observation is the number of standard deviations it falls above or below the mean
- ▶ Z score of mean = 0
- ▶ unusual observation: $|Z| > 2$
- ▶ defined for distributions of any shape

$$Z = \frac{\text{observation} - \text{mean}}{\text{SD}}$$

percentiles	<ul style="list-style-type: none"> - Area below the probability distribution curve <p>Three ways to find it.</p> <ol style="list-style-type: none"> 1. Can be used with z-score if normally distributed. Find the z-score and look up on a table. otherwise need calculus 2. Can be calculated with R <div style="border: 1px solid #ccc; padding: 5px; margin-top: 10px;"> <pre>R > pnorm(-1, mean = 0, sd = 1) [1] 0.1586553</pre>  </div> <ol style="list-style-type: none"> 3. Using a web applet
Graphs	<p>Percentiles are symmetric and unimodal Scatter plots will be a straight line</p> <div style="border: 1px solid #ccc; padding: 10px; margin-bottom: 10px;">  <p>Right skew: Points bend up and to the left of the line. Left skew: Points bend down and to the right of the line. Short tails (narrower than the normal distribution): Points follow an S shaped-curve. Long tails (wider than the normal distribution): Points start below the line, bend to follow it, and end above it.</p> </div> <div style="border: 1px solid #ccc; padding: 10px;">  <p>68% 95% 99.7%</p> <p>$\mu - 3\sigma, \mu - 2\sigma, \mu - \sigma, \mu, \mu + \sigma, \mu + 2\sigma, \mu + 3\sigma$</p> </div>

Binomial Distribution	
Bernoulli Random Variable	<p>When an individual trial has only two possible outcomes.</p> <p>Ex) Milgam's experiment where success is if a person doesn't administer a shock and failure if she does administer a shock</p>
Binomial Distribution	<p>The probability of having exactly k successes in n independent Bernoulli trials with probability of success p.</p> <p>Number of combinations * probability of success and failure</p>

	<p>The diagram illustrates the components of the binomial probability formula:</p> $\binom{n}{k} = \frac{n!}{k!(n-k)!}$ <p>is labeled "n choose k".</p> <p>$\# \text{ of scenarios} \times P(\text{single scenario})$</p> <p>$p^k(1-p)^{(n-k)}$</p> <p>is labeled "probability of success to the power of number of successes" and "probability of failure to the power of number of failures".</p>
	<p>Note that the combination formula is the “choose(,)” function in r</p> <p>Note that binomial calculation is “dbinom(k,size=sample size, p=%success)”</p>
Binomial Conditions	<ol style="list-style-type: none"> 1. Trials must be independent 2. Number of trials (n) must be fixed 3. Each trial outcome must be classified as a success or failure 4. The probability of success, p, must be the same for each trial
Example Problem	<p>According to a 2013 Gallup poll, worldwide only 13% of employees are engaged at work (psychologically committed to their jobs and likely to be making positive contributions to their organizations). Among a random sample of 10 employees, what is the probability that 8 of them are engaged at work?</p> $ \begin{aligned} n &= 10 & P(K=8) &= \binom{10}{8} 0.13^8 \times 0.87^2 \\ p &= 0.13 & &= \frac{10!}{8! \times 2!} \times 0.13^8 \times 0.87^2 \\ 1-p &= 0.87 & &= \frac{10 \times 9 \times 8!}{8! \times 2 \times 1} \times 0.13^8 \times 0.87^2 \\ k &= 8 & &= 45 \times 0.13^8 \times 0.87^2 \\ & & &= 0.00000278 \end{aligned} $
Expected Value and Standard Deviation of BinomD	<p>Among a random sample of 100 employees, how many would you expect to be engaged at work? Remember: $p = 0.13$.</p> $\mu = 100 \times 0.13 = 13$ <p>Expected value (mean) of binomial distribution: $\mu = np$</p> <p>Standard deviation of binomial distribution: $\sigma = \sqrt{np(1-p)}$</p> $\sigma = \sqrt{100 \times 0.13 \times 0.87} = 3.36$
Binomial vs Normal Distribution	<p>If a sample size is large enough, then the binomial distribution will closely resemble a normal distribution</p>

Success-failure rule: A binomial distribution with at least 10 expected successes and 10 expected failures closely follows a normal distribution.

$$\begin{aligned} np &\geq 10 \\ n(1-p) &\geq 10 \end{aligned}$$

Normal approximation to the binomial: If the success-failure condition holds,

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

$$\text{where } \mu = np \text{ and } \sigma = \sqrt{np(1-p)}$$