

UNIVERSITY *of* WASHINGTON

# **Data Science UW**

# **Methods for Data**

# **Analysis**

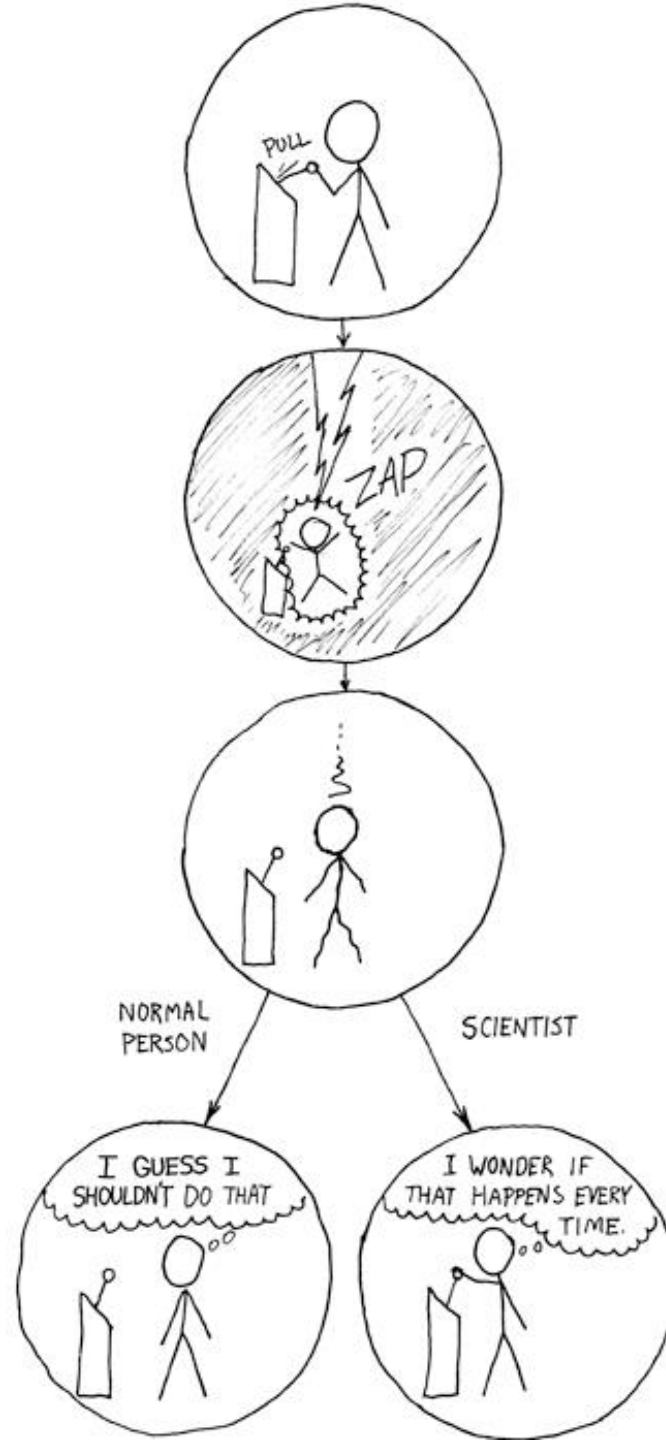
---

Hypothesis Testing and Outliers

Lecture 3

Nick McClure





W

# Topics

---

- > Review
- > Sampling Methods
- > Hypothesis Testing
- > Detecting outliers



# Review

## > Counting

- Factorials
- Permutations
- Combinations

## > Probability

- The 3 axioms
- Conditional probability
- Independent events

## > Missing Data

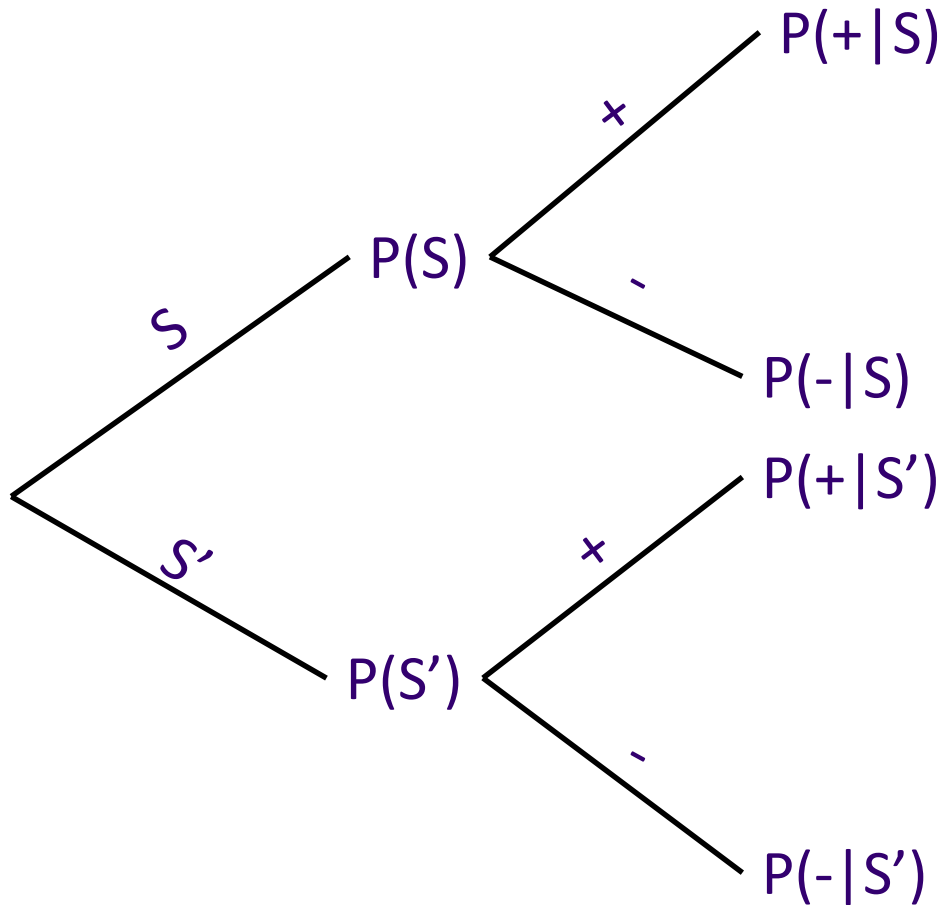


# Conditional Probability Trees

- > Let's consider a test for a Sickle Cell Anemia.
- > Events:
  - $S$  = patient has Sickle Cell Anemia
  - $S'$  = patient does not have Sickle Cell Anemia
  - $+$  = patient tests positive
  - $-$  = patient tests negative
- > Rate in US =  $1/3200$ .  $P(S) = 1/3200 = 0.0003125$ .
- > Medical company tells us that a test is 99% accurate.
  - $P(+ | S) = 0.99$
  - $P(- | S') = 0.99$

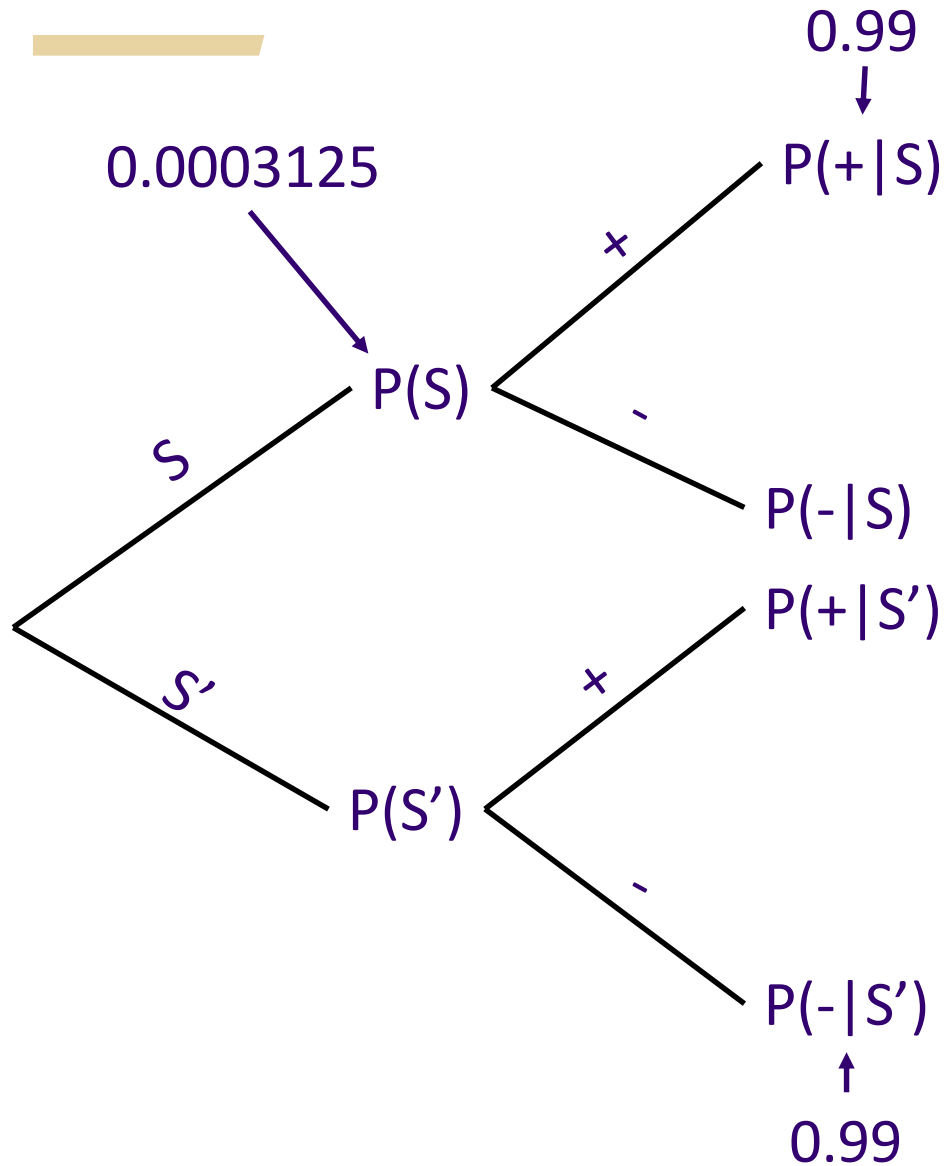


# Conditional Probability Trees



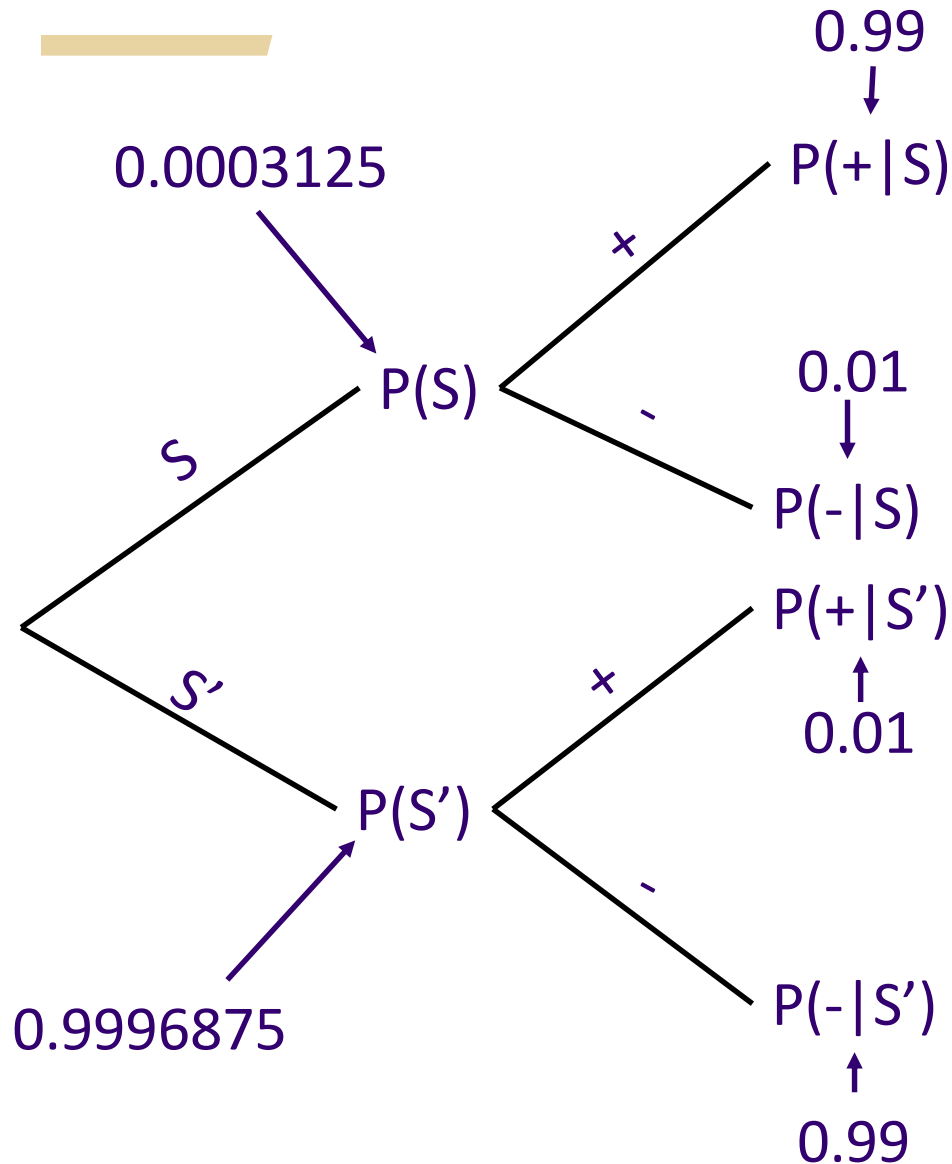
W

# Conditional Probability Trees



**W**

# Conditional Probability Trees



**W**



# Conditional Probability Trees

- > What we really want to know is:
  - What is the  $P(S|\oplus)$ ?
  - Also important to know:  $P(S| -)$ ?
- > From conditional probability definition:

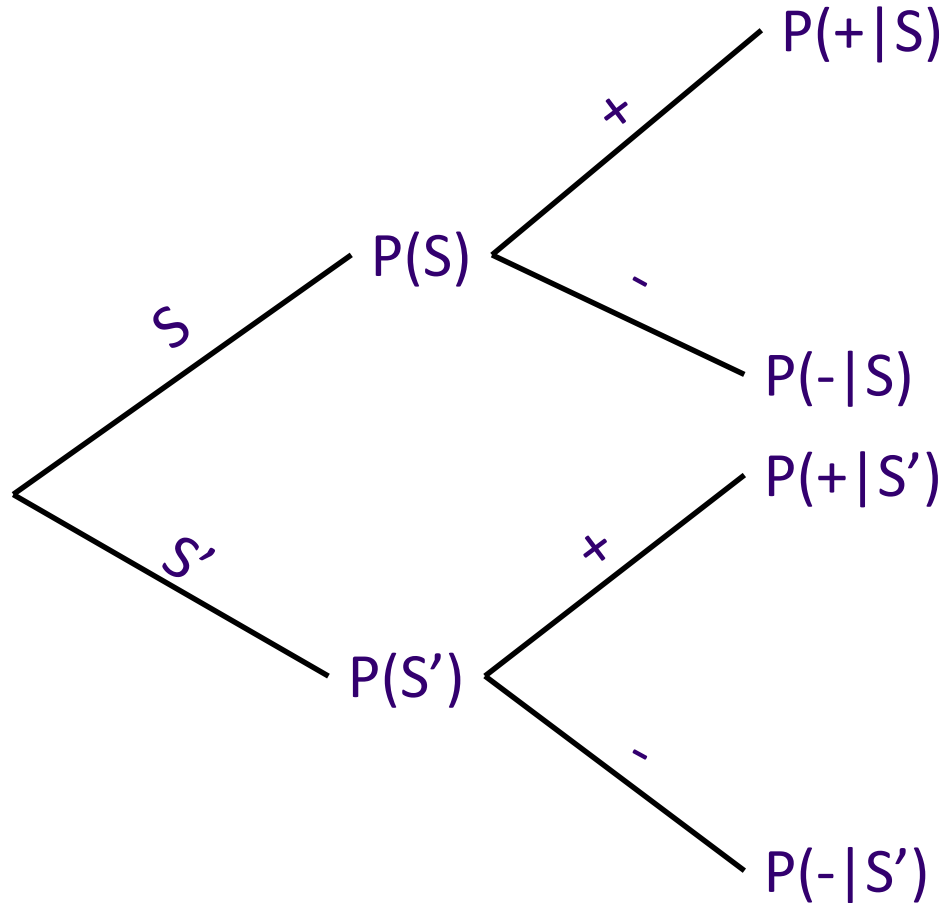
$$P(S|\oplus) = \frac{P(S \cap \oplus)}{P(\oplus)}$$

- > We also know that

$$P(\oplus) = P(\oplus \cap S) + P(\oplus \cap S')$$



# Conditional Probability Trees



$$P(+ \cap S) = P(S)P(+|S)$$

$$0.0003125 * 0.99 = 0.000309375$$

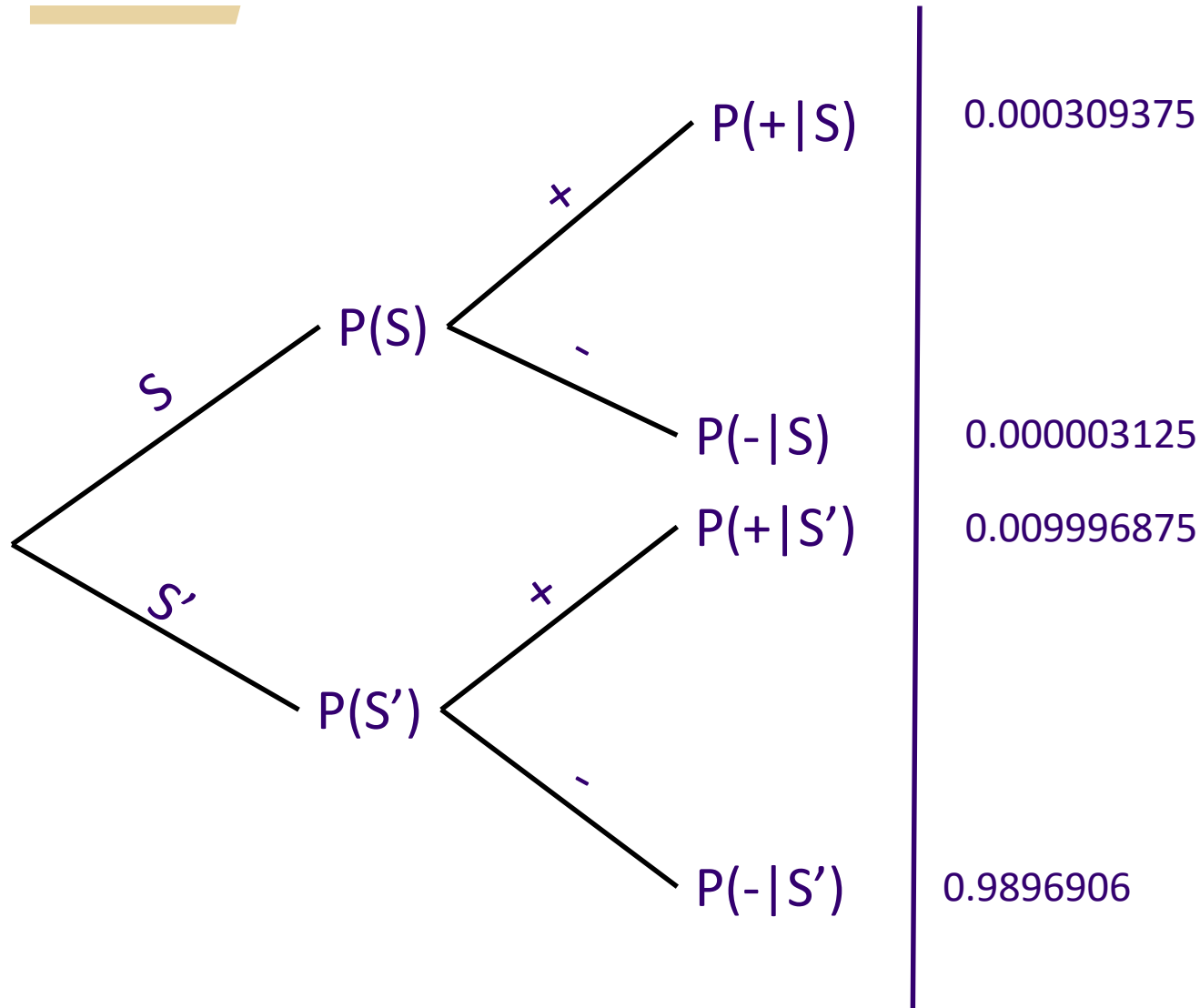
$$P(- \cap S) = P(S)P(-|S)$$

$$P(+ \cap S') = P(S')P(+|S')$$

$$P(- \cap S') = P(S')P(-|S')$$

W

# Conditional Probability Trees



**W**

# Conditional Probability Trees

$$P(\oplus) = P(\oplus \cap S) + P(\oplus \cap S')$$

$$P(\oplus) = 0.01030625$$

$$P(S|\oplus) = \frac{P(S \cap \oplus)}{P(\oplus)}$$

$$P(S|\oplus) = \frac{0.000309375}{0.01030625}$$

$$P(S|\oplus) = 0.03001819$$

Similarly,

$$P(S|-) = 0.000003157543$$

$$0.000309375 = P(+ \cap S)$$

$$0.000003125 = P(- \cap S)$$

$$0.009996875 = P(+ \cap S')$$

$$0.9896906 = P(- \cap S')$$

# W

\*Now see the probability interview question

# Sample vs. Population

- > Sampling is important because we can almost never look at the whole population.
- > We use inferences on the sample to say something about the population.
- > We need estimates of variances on the sample calculations to say something about the population.

	Sample	Population
AB Testing	The users we show A and B versions of the website.	All users that visit our site. (Past, present and future)
World Cup Soccer	Only 32 teams post qualification in one season.	All national teams in the world for four years.
Average height of Data Science Students	UW Methods for DS Class	All DS students.



# Sample vs. Population

- > If we sampled 4 beers and the ABV was [4%,5%,5%,6%], then the sample mean would be 5%.
- > There is NO variance in that value. But if we want to say something about the population, we provide the mean with a variance statistic.
  - This allows us to say something to the effect of ‘There is a 90% chance that the mean of all beers lies between 4.5% and 5.5%’.
  - In order to say something about the population we have to know how the sample was generated.



# Sampling

- > Convenience or Accidental Sampling (This is bad).
  - Grabbing whatever is easier.
- > Bernoulli Sampling
  - Every point subjected to a probability of being selected.
- > Cluster Sampling
  - Sampling in representative groups- very important later.
- > Simple Random Sample (most common)
  - Fixed size Bernoulli sampling.
- > Stratified Sampling
  - Sampling subpopulations in a representative fashion.
- > Systematic Sampling
  - Sampling every  $k$ -th element of a population.



# Sampling

- > Note that random sampling, if done properly, controls for database effects, like indexing.
- > R demo





# Large Samples and Law of Large Numbers

- > If we roll a die 60 times and 600 times, which of the dice will more likely have exactly  $1/6^{\text{th}}$  of the rolls equal to 6 appearing?
  - $P(x=10|60\text{trials})=?$
  - $P(x=100|600\text{trials})=?$
  
- > Which die will be more likely to be within 5%?
  - $P((1/6-1/20) < x < (1/6+1/20))=P(7/60 < x < 13/60)?$ 
    - >  $P(7 < x < 13 \mid 60 \text{ trials})=?$
    - >  $P(70 < x < 130 \mid 600 \text{ trials})=?$



# Law of Large Numbers!!!

- > Sample statistics converge to the population statistics as more unbiased experiments are performed.
  - E.g. The mean of 50 coin flips  $(0,1)=(T,H)$  is usually farther away from the true mean of 0.5 than 5,000 coin flips.
- > R demo of coin flips



# Standard Deviation vs. Standard Error

- > Standard Deviation: Measure of variability in a sample or population.
- > Standard Error: Measure of variability in the statistics of the sample.
- > For example:
  - Standard deviation of a sample.
  - Standard deviation of a set of means calculated from multiple samples.
    - > You can imagine that the larger my sample, the more confident we can be about the mean.
  - Standard error of a statistics decreases by a rate of  $1/\sqrt{n}$  , where  $n$  is your sample size.
- > Proof by R demo



# Hypothesis Testing

- > Identify a hypothesis that can be tested.
  - “Changing our web-site logo to be bigger on the front page will drive more than 100,000 customers to our site per day.”
- > Select a criteria to evaluate the hypothesis.
  - If our sample has a probability of  $\geq 90\%$  chance that there are more than 100,000 customers per day, we accept the hypothesis.
- > Select a random sample from the population.
  - Randomly assign a cookie to new site users that tells the server to show A or B website.
- > Compare observations to what we expect to observe and calculate statistic.



# Hypothesis Testing

- > We first state our population assumptions in the *null hypothesis*. ( $H_0$ )
- > We state our new *alternative hypothesis* as an alternative to the null. ( $H_a$ )
- > The null + alternative should make up all possible outcomes and be mutually exclusive

$H_0$ : The old website drives equal amount of traffic or more.

$H_a$ : The old website drives less traffic than the new one.

- > Decide on a significance level (probability cutoff)
  - 0.9, 0.95, and 0.99 are common (problem specific)



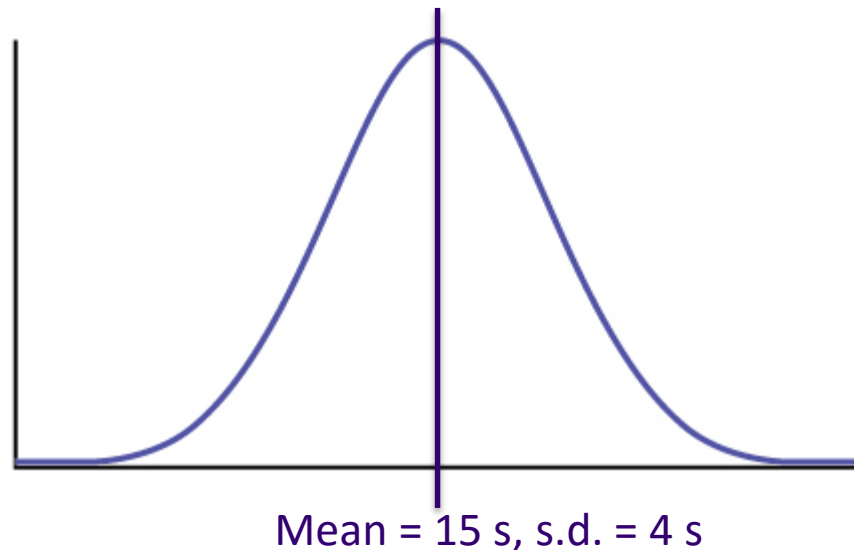
# Hypothesis Testing

- > Based on our findings we can only do two things
  - We reject the null-hypothesis.
    - > Since the alternative covers all other possibilities, we can say we accept the alternative hypothesis.
  - We *fail* to reject the null hypothesis.
    - > We do not accept the null hypothesis because we have already believed our null hypothesis from the start.
    - > We could have failed for two reasons:
      - The alternative hypothesis was false to begin with.
      - We did not collect enough evidence for the alternative hypothesis.



# Hypothesis Testing

- > We know that the average time a user spends on a page has a mean of 15 seconds and a s.d. of 4 seconds.
- > If we assume normality, how do we test if a change to the page has a higher view time?



**W**

# Hypothesis Testing

- > We know that the average time a user spends on a page has a mean of 15 seconds and a s.d. of 4 seconds.
- > If we assume normality, how do we test if a change to the page has a higher view time?

$H_0$ : The old website has the same or more viewership than the new website.

$H_a$ : The old website has less viewership than the original.

Or...

$H_0$ : The new website has the same or less viewership than the original.

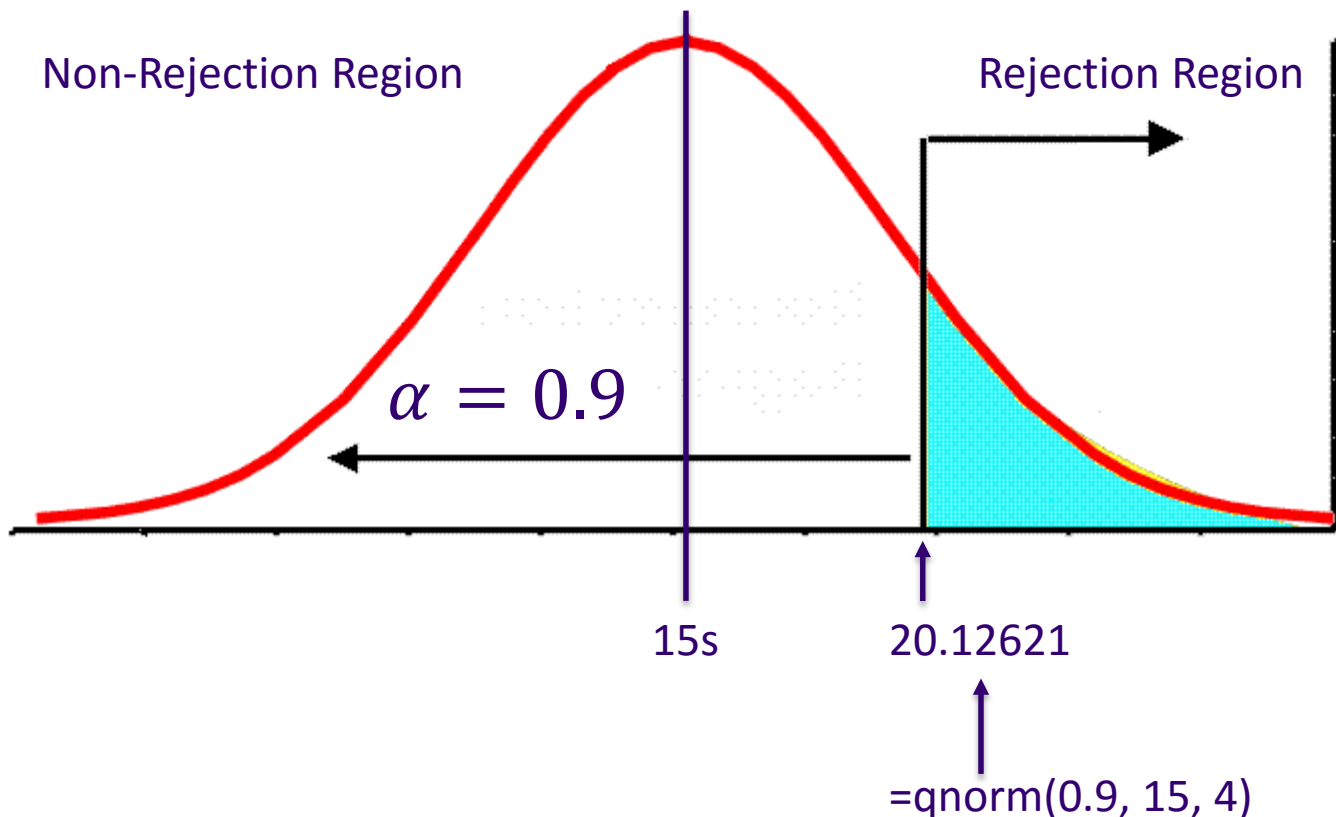
$H_a$ : The new website has more than the original website.





# Hypothesis Testing

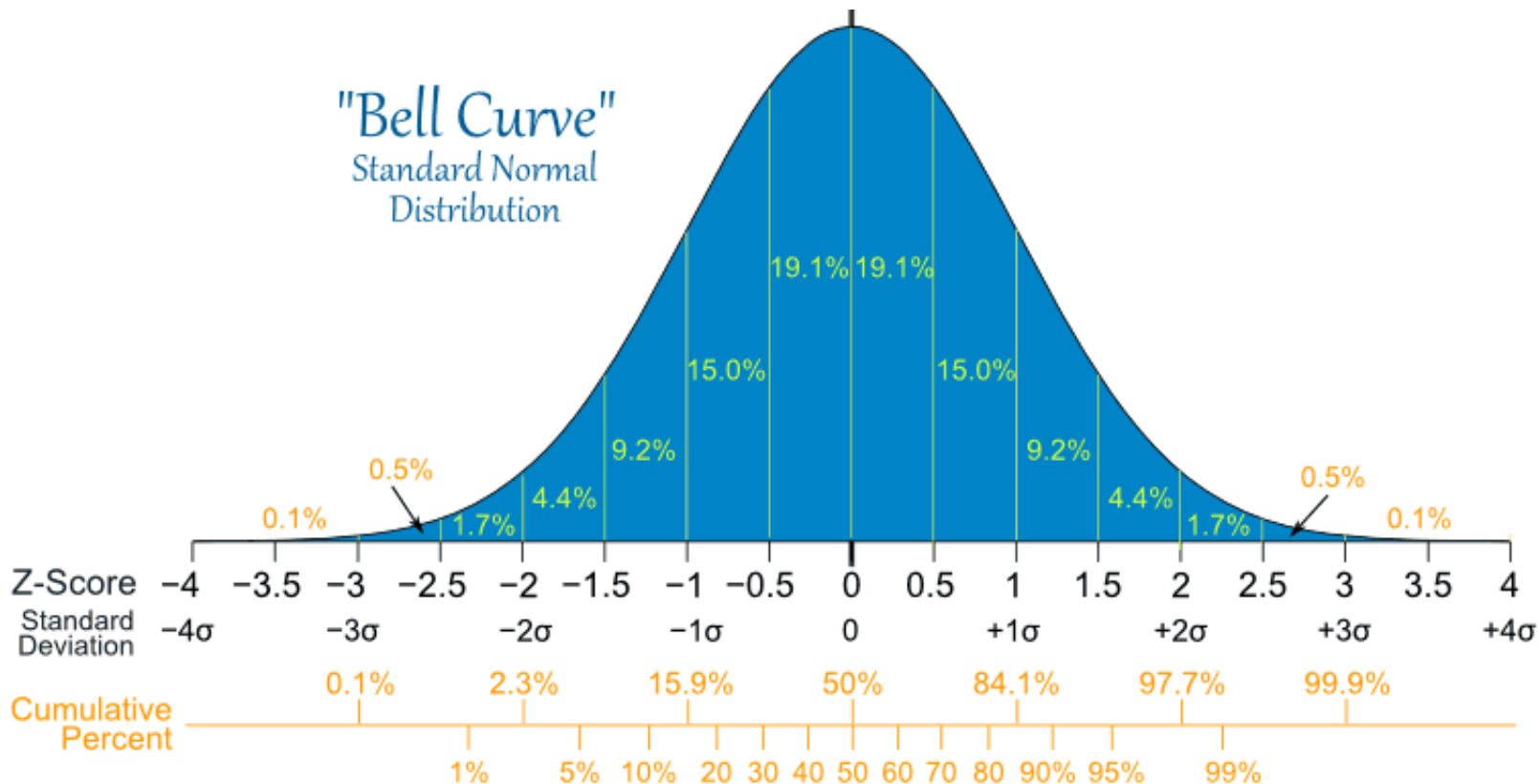
- > We now select a confidence value.
- > An event in the **blue** region will have a 10% chance or less of occurring.



W

# Hypothesis Testing

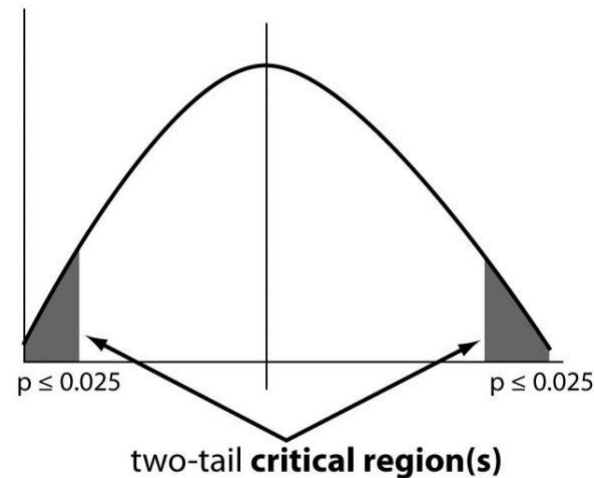
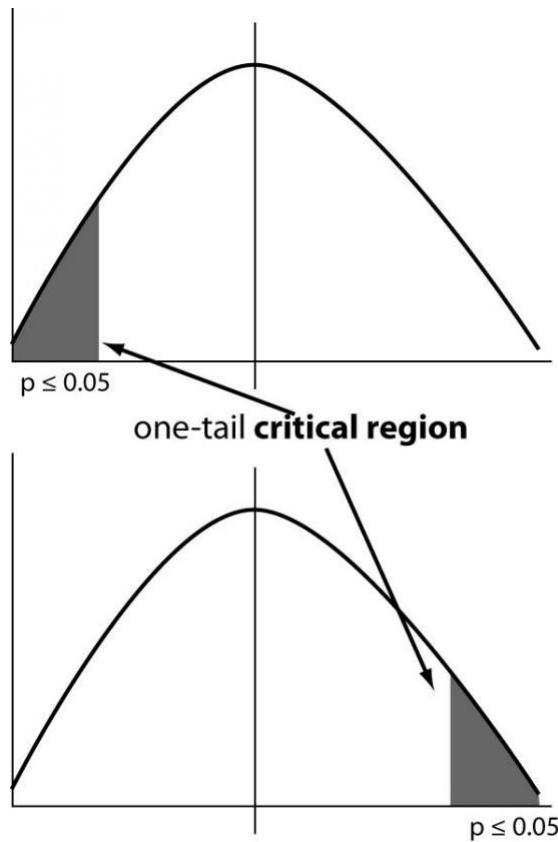
- > Probability areas on the normal curve are directly related to the distance to the mean.



W

# Hypothesis Testing

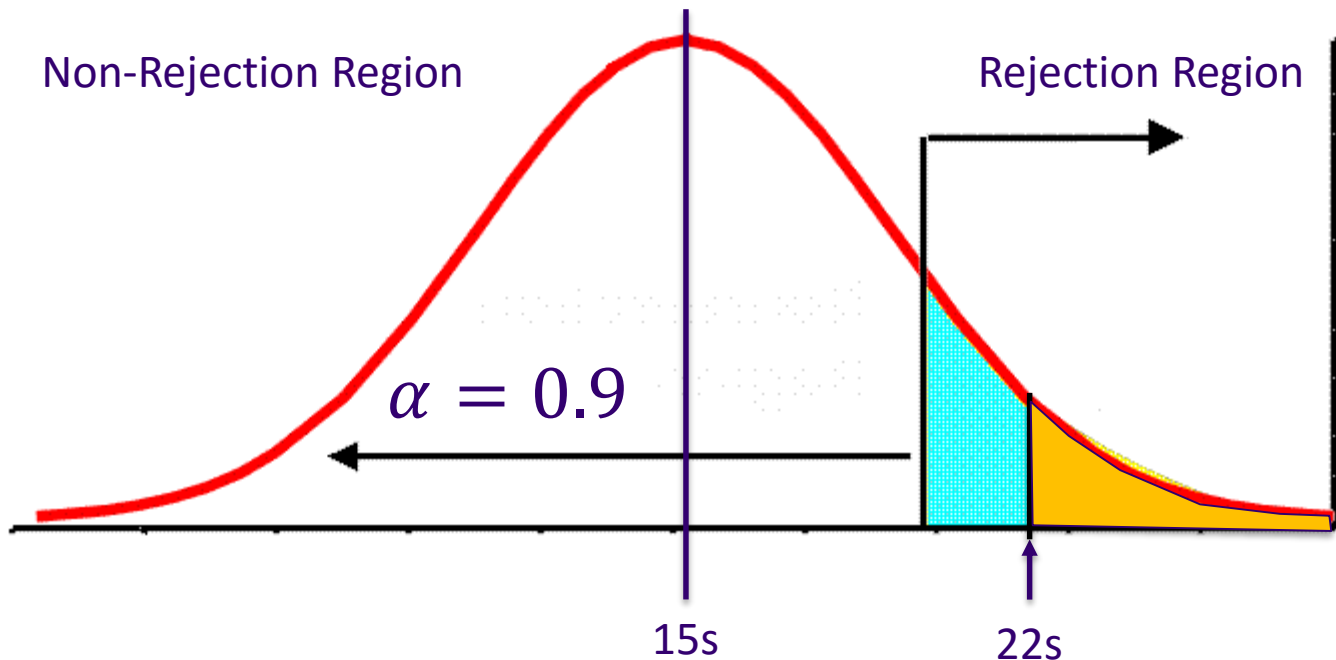
- > Asking if a new value is “greater than” or “less than” the null creates a **one-tailed hypothesis test**.
- > Asking if a new value is “not equal to” the null creates a **two tailed hypothesis test**.



W

# Hypothesis Testing

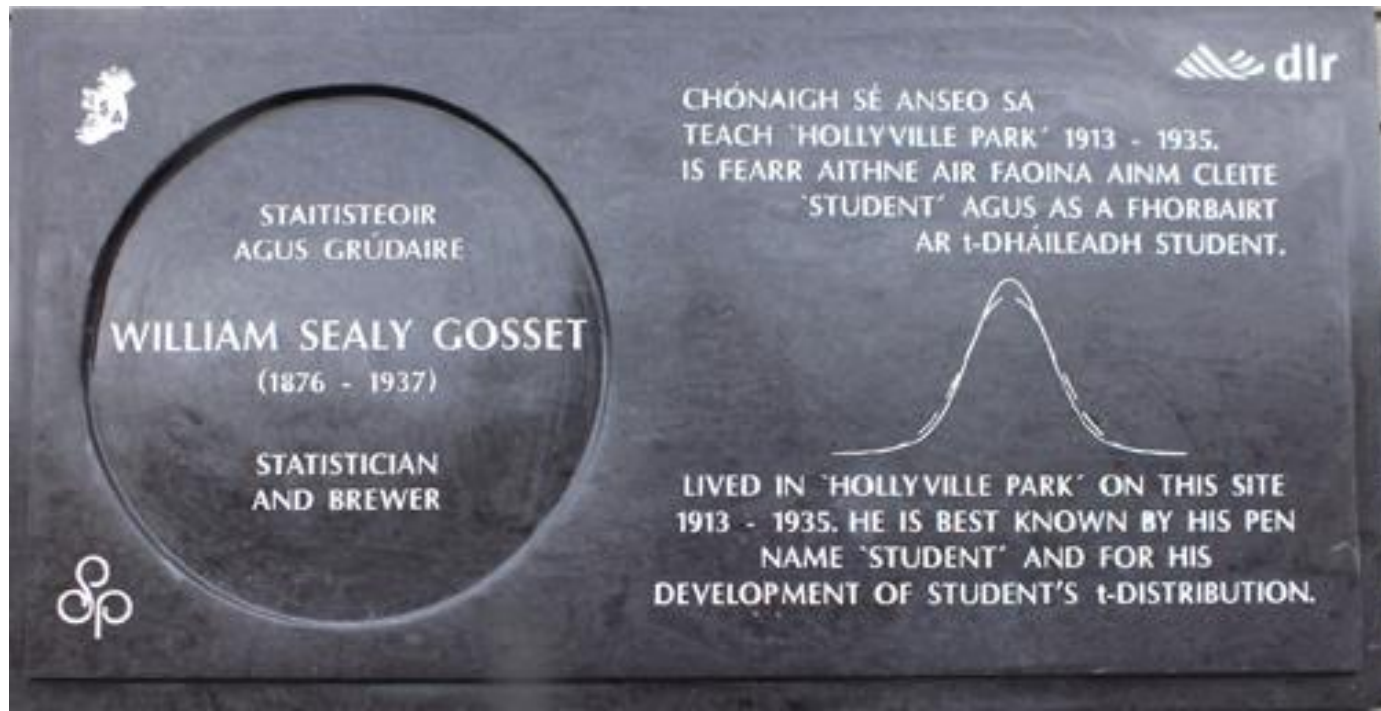
- > The p-value is the probability of obtaining the sample results or worse, assuming the null hypothesis is true.
- > What is the p-value of a sample mean of 22 seconds?



W

# Hypothesis Testing

- > What if we don't know our population's standard deviation?
- > There are involved probability rules that relax the normality assumption in favor of "heavier tails".
- > This distribution is known as the Student's T-distribution.



# t-test

- > Student's T-test: tests a hypothesis about the difference of two sets of data:
  - Test whether a population mean has a specified value.
  - Test the difference between two means (equal, unknown variances).
  - Test a paired-response difference from zero.
    - > E.g. a before/after drug treatment on patients.
  - Test whether the slope of a line is not zero.
    - > Important for testing the importance of variables (later in class).
- > Use 'Welch's T-test' for testing the difference between two means (unknown variances, potentially different).
- > Picking the right test changes test's results.



# t-test in R

---

> t.test() demo in R



# Summary

- > The normal test and t-test are used for testing values from a continuous distribution (or approximately so).
- > What if we wanted to test occurrences or count data?





# Chi-squared Test (Pearson's)

- > Unpaired test for counts in different categories.
- > These categories must be mutually exclusive.
  - Does the patient have cancer? (yes/no)
  - Rolling a die. (1,2,3,4,5,6)
  - Does a tweet contain a specific word? (yes/no)
- > This tests whether the different categories differ in some specific value.
- > In order to do this test, we have to specify the 'degrees of freedom' in the Chi-squared test.
  - This is equal to  $n-1$ . Where  $n$  equals the number of different categories.
- > The test looks at the sum of the outcome differences from expectations.



# Chi-squared Test (Pearson's)

- > Example: A-B test with three different outcomes.

	Occurrence	Expectation %	Expectation Counts	Difference	Squared Difference	(Squared Difference)/Expected
Leave Page	55	0.6	$=0.6*120=72$	$=55-72=-17$	289	$=289/72=4.014$
Continue Purchase	43	0.3	$=0.3*120=36$	$=43-36=7$	49	$=49/36=1.361$
Add More to Purchase	22	0.1	$=0.1*120=12$	$=22-12=10$	100	$=100/12=8.333$
Totals	120					13.708

- > Test statistic is 13.708 on a chi-squared distribution with  $(3-1)=2$  degrees of freedom.
- > Degree of freedom is (# of options minus 1).
- > R Demo



# Chi-squared Test (Pearson's)

- > Chi-squared is also used for a 'goodness of fit' test.
- > Test if sample is representative of population.
  - Test if your sample has expected make up of categories.
  - E.g. If our population is 50-50 men-women, then we test if our sample is different from those expected probabilities.
  - Or, in the case of the homework, we can test if the population is representative for all 50 states.
- > If our total sample size is small, we see a breakdown of the Chi-squared test. (a subgroup size  $\sim < 10$ )
  - We switch to a Fisher's Exact test in these cases.



# Fisher's Exact Test

- > Tests for difference between two groups based on ratios.
- > Exact test, because it calculates the probability of observing the sample under the null or worse in all possible cases.
  - Not as much statistical 'power' as Chi-Squared.
  - If you have larger sample sizes, and the two categories are sufficiently different, both tests should give similar p-values.
- > Probability of observing a specific outcome:

	Cat. 1	Cat. 2
Successes	A	B
Failures	C	D

$$prob = \frac{\binom{A+C}{A} \binom{B+D}{B}}{\binom{N}{A+B}}$$



# Fisher's Exact Test

> For example, consider:

2	3
3	4

> We sum up the probability of that outcome occurring or worse:

0	5	1	4	2	3	3	2	4	1	5	0
5	2	4	3	3	4	2	5	1	6	0	7



Same outcome or worse, with the same marginals (row sums).

> R-demo

# W

# Fisher's Exact Test

- > For paired data, look into McNemar's test (based on the binomial distribution).



# Outliers

## > Outlier causes:

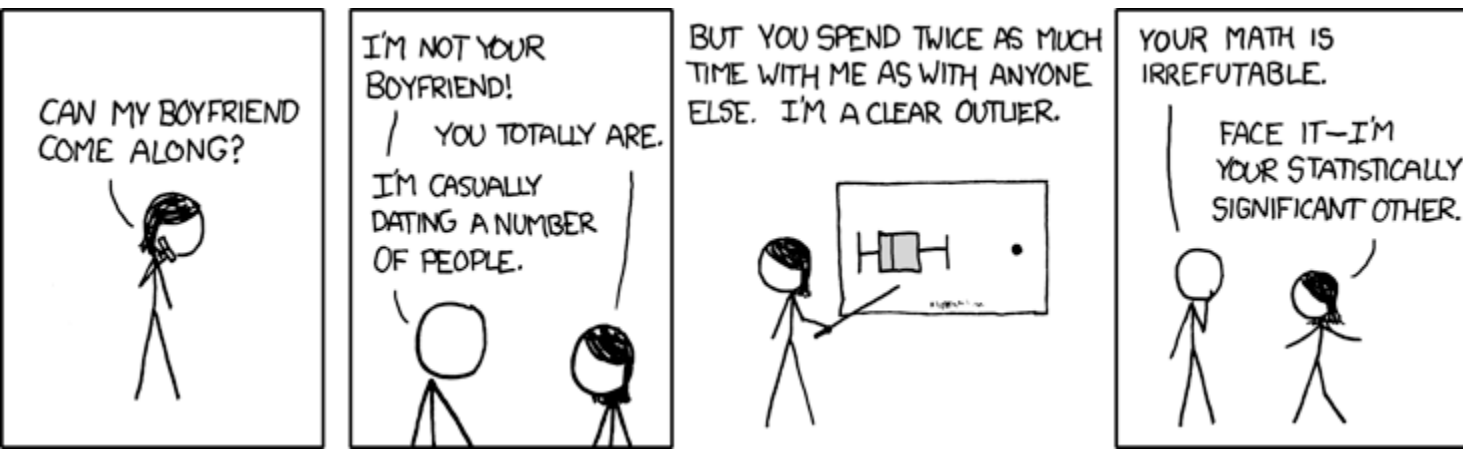
- Bad data
  - > Sensor misread, human error, software error
- Non-representative data
  - > Real data that can be argued to be out of our interest. E.g. a sample of annual salaries that includes Warren Buffet.
  - > Must provide a legitimate argument to consider as outlier.



# Outliers

## > Outlier Issues

- Identification
  - > Test whether or not an observation(s) is an outlier
- Accomodation
  - > Using robust statistical techniques that can deal with outliers. E.g. Using the median instead of the mean.
- Dealing/fixing
  - > Correcting data to not have outliers influence statistical conclusions.



W



# Assignment

---

## > Complete Homework 3:

- Test hypotheses for farm subsidy data:
  - > Does the sample represent all 50 states equally?
  - > Does the sample represent all 50 states, weighted by the number of farms per state equally?
  - > <http://www.fsa.usda.gov/FSA/webapp?area=newsroom&subject=land&topic=foi-er-fri-pfi>
- You should submit:
  - > **R-script** written in production level style.
  - > A text document summarizing your findings.
- Read Intro to Data Science Chapter 6.

