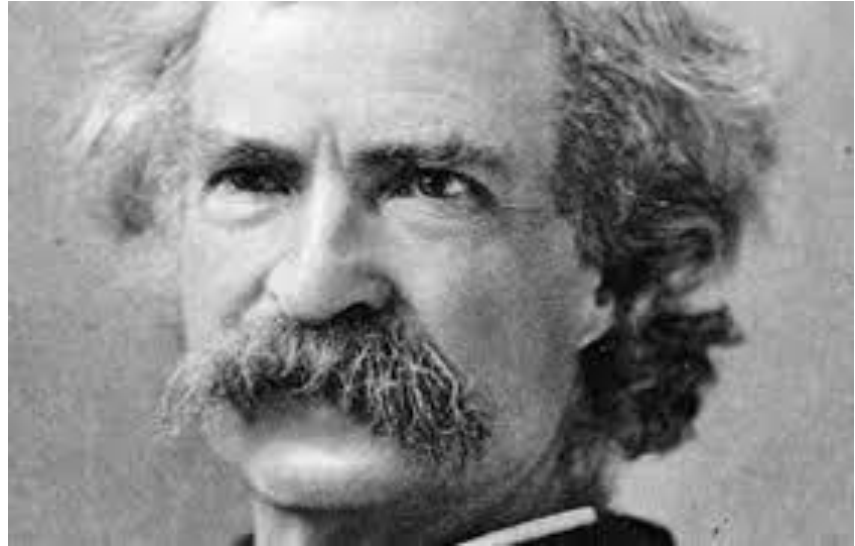


# Random Quote

“There are three kinds of lies: lies, damned lies, and statistics.”

- Mark Twain



CSE 312

# Foundations of Computing II

Lecture 29: How to lie/be misled/detect lies with statistics



PAUL G. ALLEN SCHOOL  
OF COMPUTER SCIENCE & ENGINEERING

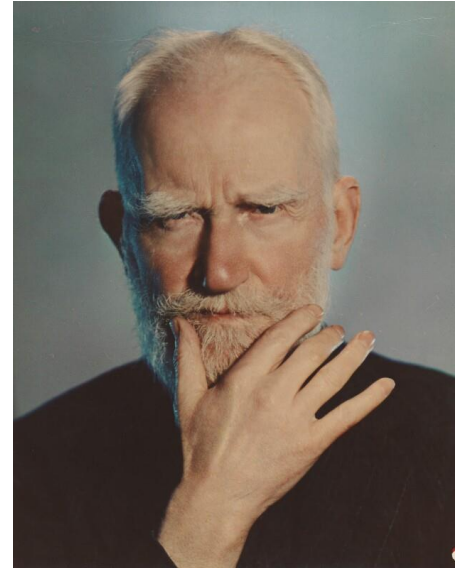
**Anna R. Karlin**

Slide Credit: Based on Stefano Tessaro's slides for 312 19au incorporating ideas from Alex Tsun, Maya Bar-Hillel & myself ☺

# Random Quote

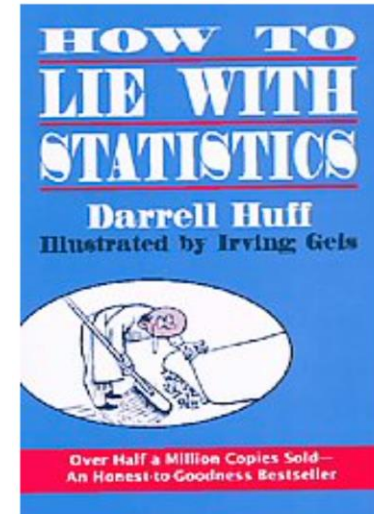
“It is the mark of a truly intelligent person to be moved by statistics”

- George Bernard Shaw



# The Book

- Published in 1954, over 500,000 copies sold
- “A great introduction to the use of statistics, and a great refresher for anyone who’s already well versed in it” - Bill Gates.



# The Book

- Published in 1954, over 500,000 copies sold
- “A great introduction to the use of statistics, and a great refresher for anyone who’s already well versed in it” - Bill Gates.
- Doesn’t teach how to lie with statistics, but how we are/can be lied to using statistics



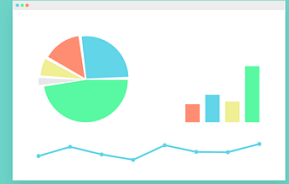


# To be clear...

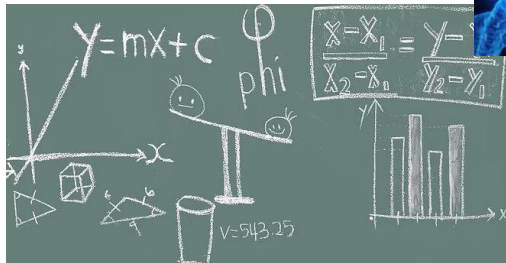
- Many lies are unintentional
- People passing on misinformation/bad information that they don't even know is bad.
- People using bad data to make inferences
- People not understanding statistics well enough



# What is “Statistics”?



- A way to make sense of information from data
- Framework for thinking, for reaching insights, and solving problems.
- Numbers alone mean very little without context
- Statistics is a marriage of:
  - Math
  - Science
  - Art





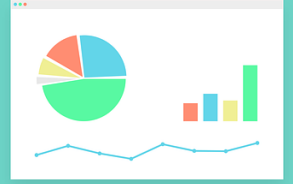
# Random Quote

“Statistical Thinking will one day be as necessary for efficient citizenship as the ability to read and write”

- H.G. Wells

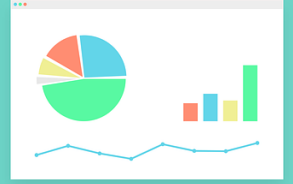


# Statistical Inference



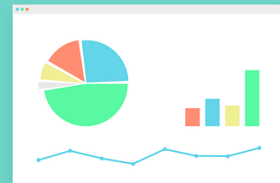
- Making an estimate or prediction about a population based on a sample.

# Statistical Inference



- Making an estimate or prediction about a **population** based on a **sample**.
  - Often very expensive/impossible to survey an entire population (all students at UW, all residents in the U.S)

# Statistical Inference



- Making an estimate or prediction about a **population** based on a **sample**.
  - Often very expensive/impossible to survey an entire population (all students at UW, all residents in the U.S)
  - Need to use a **random unbiased** *sample of population* to draw conclusions (with some chance/margin of error)

# Sampling Gone Wrong (Bias)

“The Literary Digest” Magazine wanted to predict 1936 election:

- Alfred Landon vs Franklin D Roosevelt
- Sent 10 million surveys and received 2.4 million responses
- From a “List” containing: their subscribers, owners of cars and telephones

Electoral Votes	Prediction	Actual
<b>Landon</b>		
<b>Roosevelt</b>		



# Sampling Gone Wrong (Bias)

“The Literary Digest” Magazine wanted to predict 1936 election:

- Alfred Landon vs Franklin D Roosevelt
- Sent 10 million surveys and received 2.4 million responses
- From a “List” containing: their subscribers, owners of cars and telephones

Electoral Votes	Prediction	Actual
<b>Landon</b>	370	
<b>Roosevelt</b>	161	



# Sampling Gone Wrong (Bias)

“The Literary Digest” Magazine wanted to predict 1936 election:

- Alfred Landon vs Franklin D Roosevelt
- Sent 10 million surveys and received 2.4 million responses
- From a “List” containing: their subscribers, owners of cars and telephones

Electoral Votes	Prediction	Actual
<b>Landon</b>	370	8
<b>Roosevelt</b>	161	523



# Sampling Gone Wrong (Bias)

“The Literary Digest” Magazine wanted to predict 1936 election:

- Alfred Landon vs Franklin D Roosevelt
- Sent 10 million surveys and received 2.4 million responses
- From a “List” containing: their subscribers, owners of cars and telephones

Electoral Votes	Prediction	Actual
<b>Landon</b>	370	8
<b>Roosevelt</b>	161	523



What went wrong?



# Sampling Gone Wrong (Bias)

Let  $x_1, x_2, \dots, x_n$  be iid samples...



# Sampling Gone Wrong (Bias)

Let  $x_1, x_2, \dots, x_n$  be iid samples...

- Not Representative
  - Voluntary Response Bias
    - Only 24% of respondents answered the poll.



# Sampling Gone Wrong (Bias)

Let  $x_1, x_2, \dots, x_n$  be iid samples...

- Not Representative
  - Voluntary Response Bias
    - Only 24% of respondents answered the poll.
  - Not the Right Population
    - Was biased toward people with more money, education, information, alertness than average American



# Sampling Gone Wrong (Bias)

Let  $x_1, x_2, \dots, x_n$  be iid samples...

- Not Representative
  - Voluntary Response Bias
    - Only 24% of respondents answered the poll.
  - Not the Right Population
    - Was biased toward people with more money, education, information, alertness than average American
- Not Random
  - Convenience Sampling
    - Only to people whose contact information they have.
    - Like standing outside a church and asking “Do you believe in God?”, using those samples to represent the US population.



# Sampling Gone Wrong (Bias)

Let  $x_1, x_2, \dots, x_n$  be iid samples...

- Not Representative
  - Voluntary Response Bias
    - Only 24% of respondents answered the poll.
  - Not the Right Population
    - Was biased toward people with more money, education, information, alertness than average American
- Not Random
  - Convenience Sampling
    - Only to people whose contact information they have.
    - Like standing outside a church and asking “Do you believe in God?”, using those samples to represent the US population.

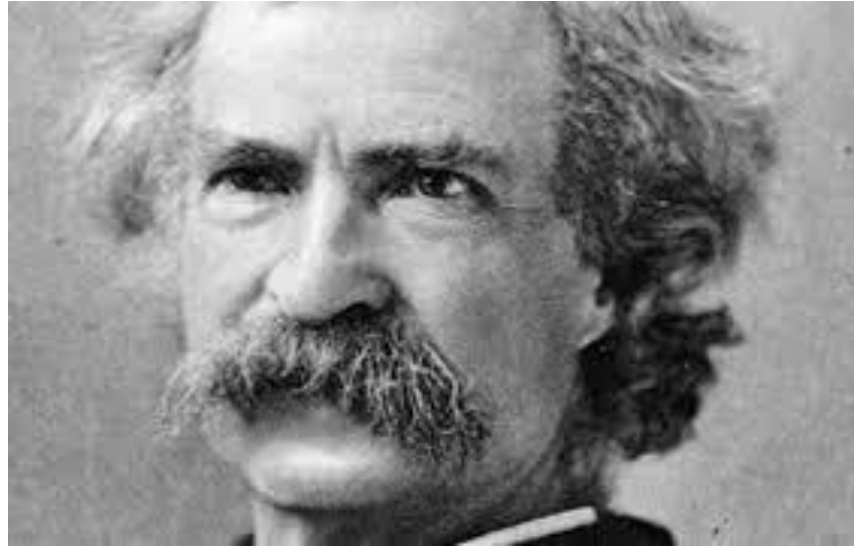


**More samples is NOT a solution for bad sampling technique...**

# Random Quote

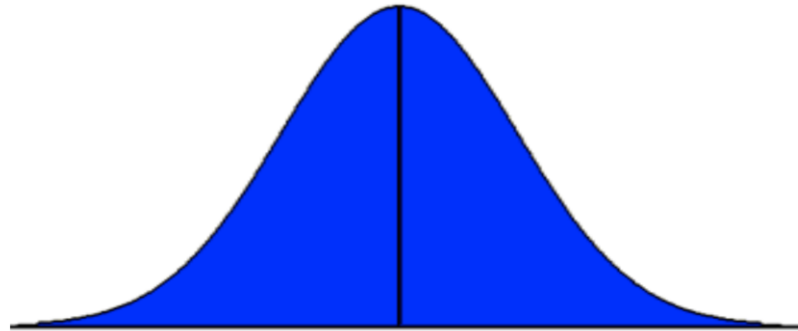
“Facts are stubborn, but statistics are more pliable.”

- Mark Twain



# Detecting lies with statistics

A story about the famous French mathematician Henri Poincare

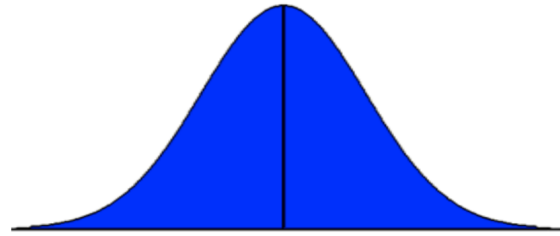


950 grams

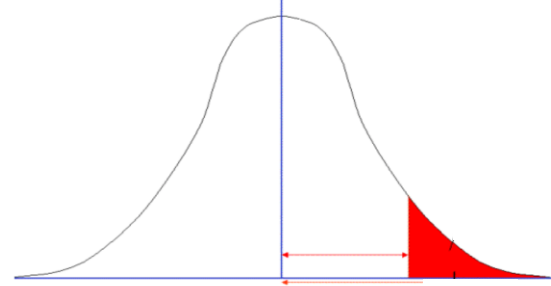


# Detecting lies with statistics

A story about the famous French mathematician Henri Poincaré



950 grams



950 grams





# To fake a distribution...

You'd better know what it looks like....

People that are untrained in statistics often don't.

For example, people are really bad at faking a sequence of fair coin tosses.

# Random Quote

“It’s easy to lie with statistics. It’s hard to tell the truth without statistics.”

- Andrejs Dunkels



# First digit phenomenon

Suppose that I pick a random integer in the range 1..999

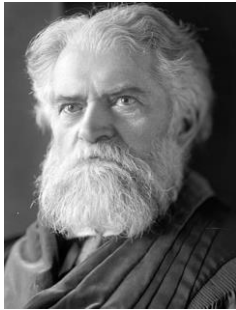
What's the chance that the first digit of the number I pick is a 1?

- a). About 1/9
- b). About 11%
- c) 30%
- d) I don't know.

# Benford's Law

How about in real life? Do certain digits in numbers collected randomly from the front pages of the newspaper or census statistics or from stock-market prices occur more often than others?

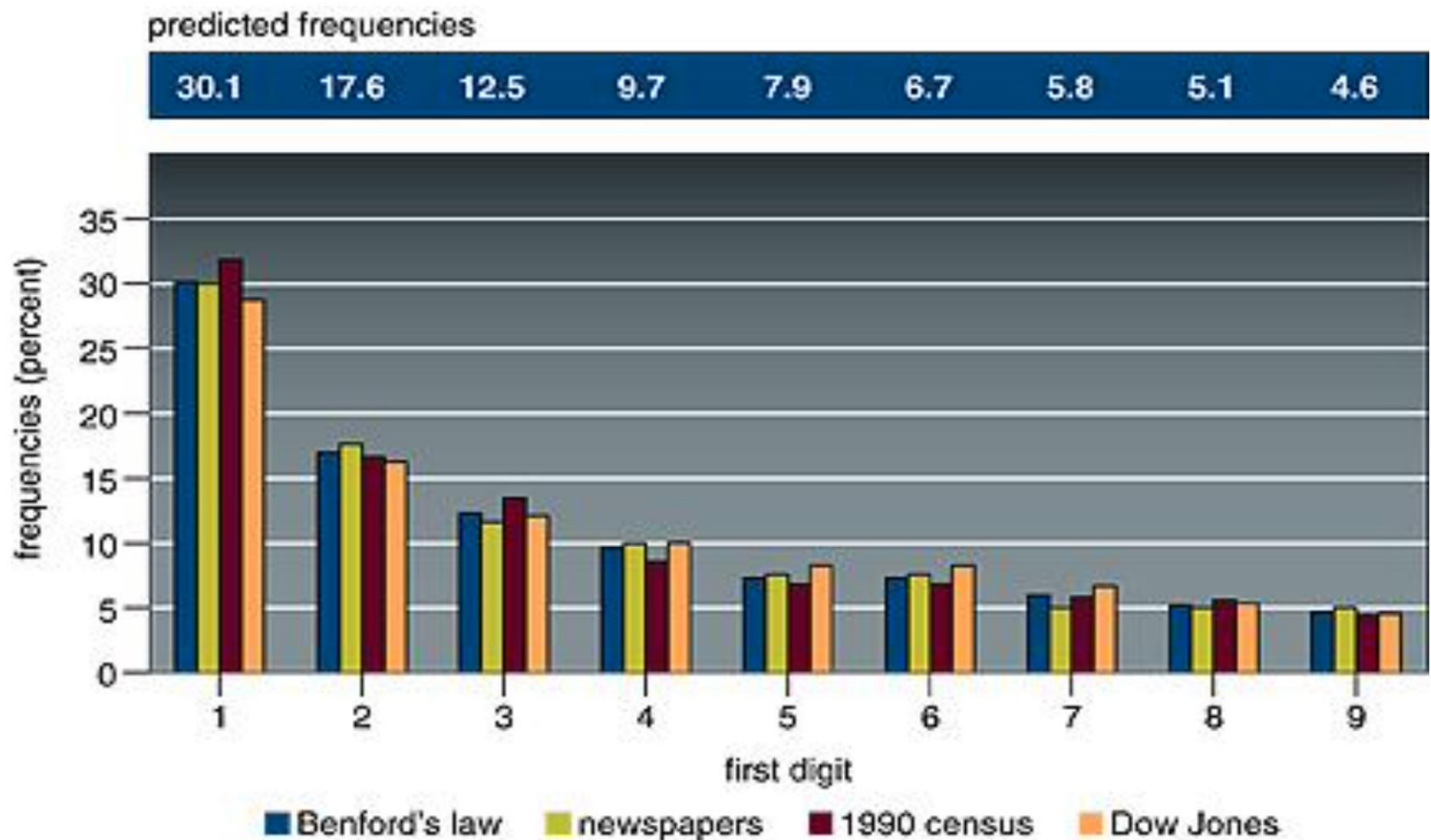
Frequency with which first significant digit is  $d = \log(1 + 1/d)$



Gr.	o	+	-	Sum
o	Sum	Logarithm	Difference	Sum
0	0	1.000000	0.000000	0.000000
1	1000	0.301030	0.698970	0.301030
2	10000	0.477121	0.522879	0.477121
3	100000	0.504130	0.495870	0.504130
4	1000000	0.530783	0.469217	0.530783
5	10000000	0.556469	0.443531	0.556469
6	100000000	0.581159	0.417841	0.581159
7	1000000000	0.605051	0.392949	0.605051
8	10000000000	0.628932	0.368068	0.628932
9	100000000000	0.652789	0.343211	0.652789
10	1000000000000	0.676626	0.318374	0.676626
11	10000000000000	0.700443	0.293557	0.700443
12	100000000000000	0.724242	0.268760	0.724242
13	1000000000000000	0.748023	0.244001	0.748023
14	10000000000000000	0.771786	0.219284	0.771786
15	100000000000000000	0.795532	0.194619	0.795532
16	1000000000000000000	0.819262	0.169997	0.819262
17	10000000000000000000	0.842976	0.145421	0.842976
18	100000000000000000000	0.866674	0.120891	0.866674
19	1000000000000000000000	0.890357	0.096407	0.890357
20	10000000000000000000000	0.914025	0.071969	0.914025
21	100000000000000000000000	0.937678	0.047578	0.937678
22	1000000000000000000000000	0.961317	0.023234	0.961317
23	10000000000000000000000000	0.984942	0.008937	0.984942
24	100000000000000000000000000	1.008554	0.004688	1.008554
25	1000000000000000000000000000	1.032153	0.002487	1.032153
26	10000000000000000000000000000	1.055740	0.001347	1.055740
27	100000000000000000000000000000	1.079315	0.000765	1.079315
28	1000000000000000000000000000000	1.102879	0.000441	1.102879
29	10000000000000000000000000000000	1.126433	0.000271	1.126433
30	100000000000000000000000000000000	1.149978	0.000164	1.149978



From "The First-Digit Phenomenon" by T. P. Hill, American Scientist, July-August 1998)



# Long-term efforts to “prove” Benford’s Law

Properties of a random sample that result in such a distribution? E.g. not true for Unif {1,...999}

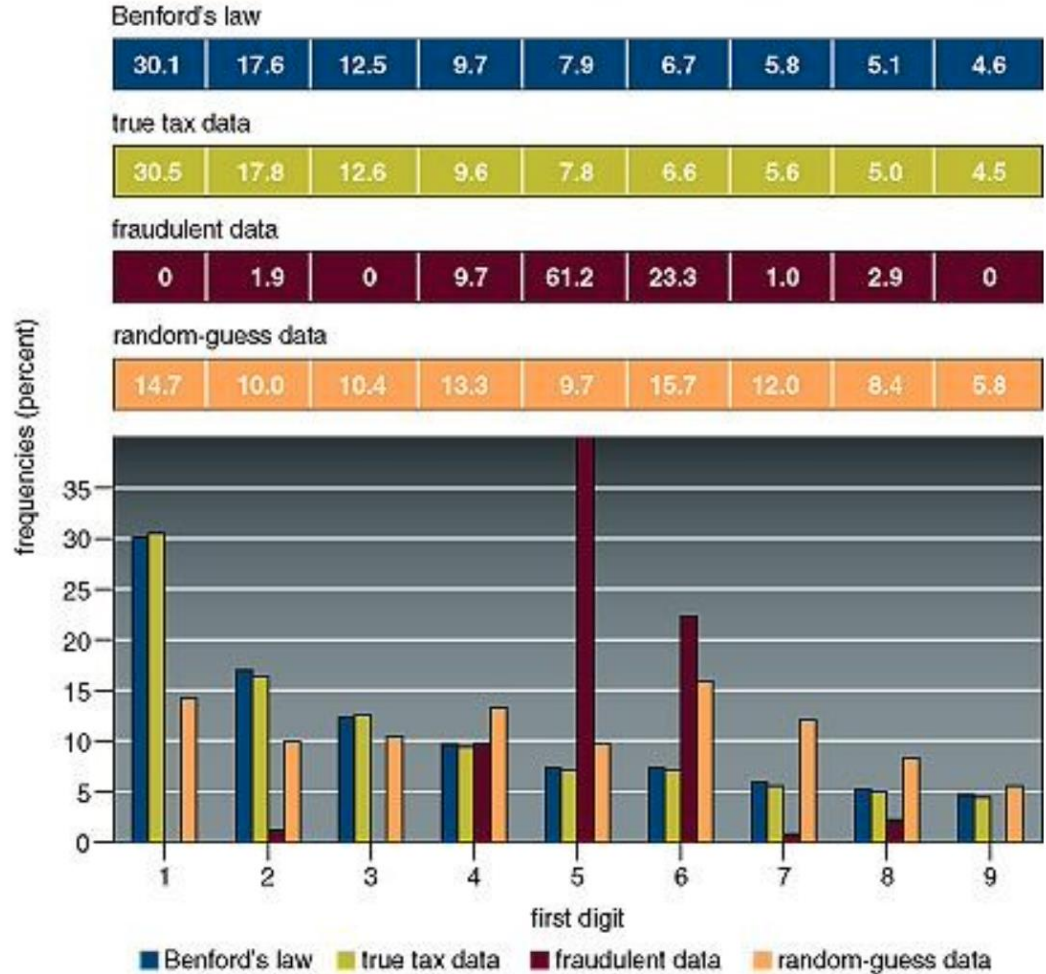
- **Scale invariance:** e.g. convert from dollars to pesos shouldn’t change the first digit frequencies much
- **Independent of base:** Equally valid when numbers expressed in base 10, base 100, or others

**The only distributions on numbers that satisfies these conditions satisfy**

$$\Pr(\text{first significant digit} = d) = \log(1 + 1/d)$$

# Modern Application

- Using Benford's law to detect fraud or fabrication of data in financial documents.



# Random Quote

“It is easy to lie with statistics, but easier to lie without them”.

Fred Mosteller



# “Too good to be true”

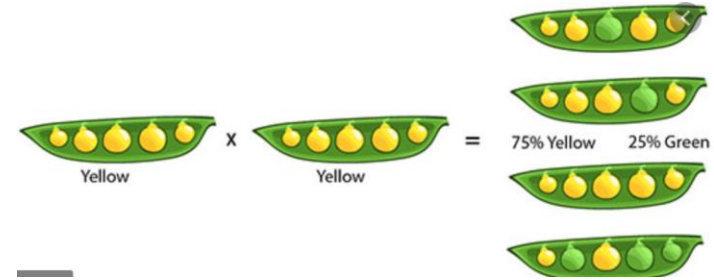
- The special case of not appreciated the expected magnitude of sampling error.
- Data comes out “too good to be true”, a telltale sign of having been tampered with, if not generated out of whole cloth.

# Gregor Mendel's Sweet Peas



Postulated that self fertilization Of hybrid yellow-seeded sweet peas would yield offspring with

- 0.75 chance yellow-seeded
- 0.25 chance green seeded.



1865, reported results of 8023 experiments:

- 0.7505 yellow-seeded
- 0.2495 green-seeded.

Probability of observations as close to expected value as he reported is minute.

# Some telltale signs of fakery....

- Wrong shape
- Too close to expected value (especially replicated)
- Too far from expected value
- Replications too good to be true.



Another famous example: Sir Cyril Burt's Twins

3 data sets: same to 3 decimal points.

# Random Quote

"82.123456789% of statistics are made up."

- Alex Tsun



# p-Hacking

Manipulating data or statistical analyses to get **“significant p-values”**

First, a brief primer on hypothesis testing and p-values.

Suppose that I believe that jelly beans cause acne. How might I provide evidence of this?

Approach – “probabilistic proof by contradiction”



# Hypothesis Testing

Want to provide evidence that the null hypothesis can be rejected!

Average teenager has amount of acne with mean  $\mu$  and variance  $\sigma^2$

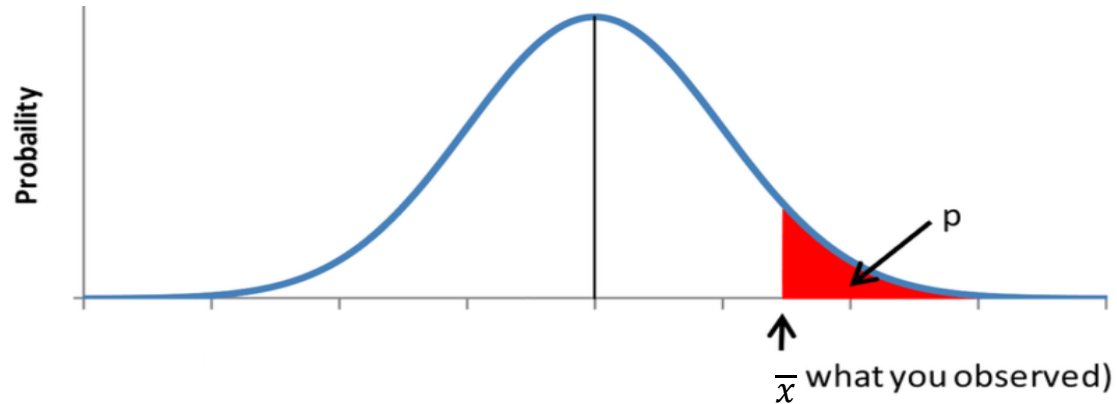
**$H_0$  - null hypothesis (baseline):** the mean amount of acne someone who eats jelly beans has is  $\mu$ , i. e., **jelly beans have no effect on acne**

**$H_A$  - Alternative hypothesis:** the mean amount of acne someone who eats jelly beans has is  $> \mu$

Choose **significance level**, say 0.05

Observe 100 jelly-bean-eating teenagers and measure their acne levels.

Suppose sample mean observed  $\bar{x}$

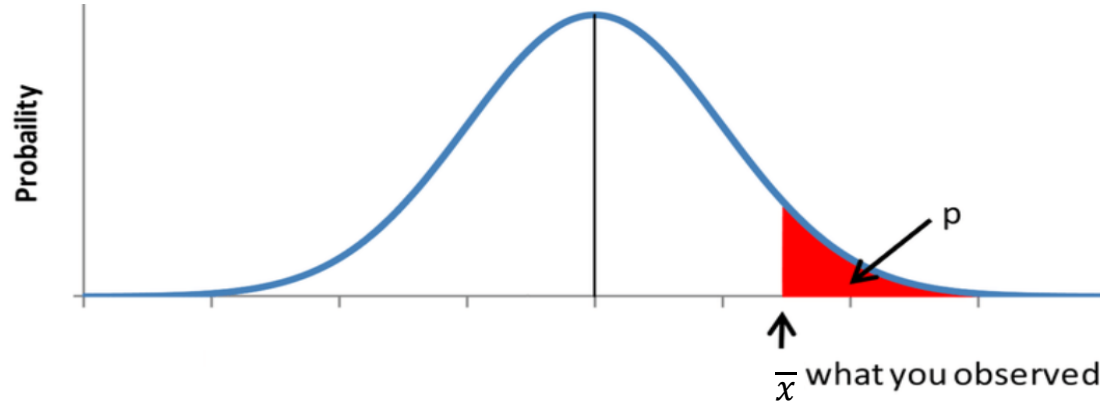


# Hypothesis Testing

**$H_0$  - null hypothesis (baseline):**  
jelly beans have no effect on acne

**$H_A$  - Alternative hypothesis:**  
Jelly beans increase acne

Suppose find that for measured  $\bar{x}$



$\Pr(\text{observing amount of acne this high if } H_0 \text{ true}) = \Pr(\bar{X} \geq \bar{x}) = 0.0162.$  This is our **p-value**.

If  $p < 0.05$  reject  $H_0$  at the 0.05 significance level, i.e., **strong statistical evidence that jelly beans cause an increase in acne.** (If  $H_0$  was true, this would be a very unlikely outcome).

If  $p > 0.05$ , fail to reject  $H_0$ ;  
**Not enough evidence to suggest the jelly bean effect on acne was significant.**

# Hacking

SIGNIFICANT

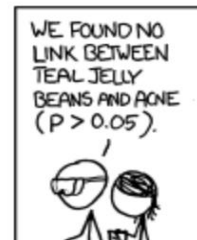
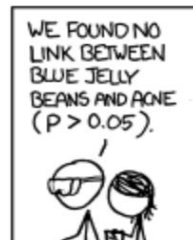
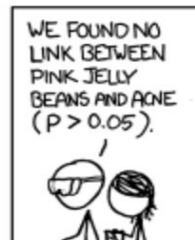
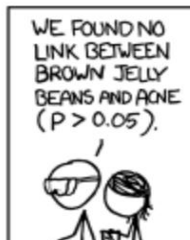
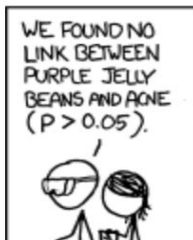
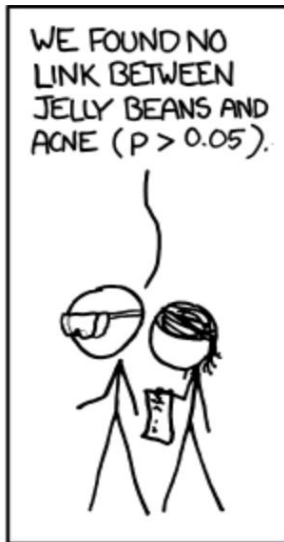
|<

< PREV

RANDOM

NEXT >

>|

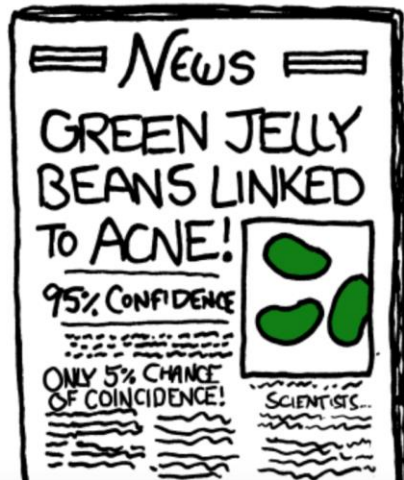
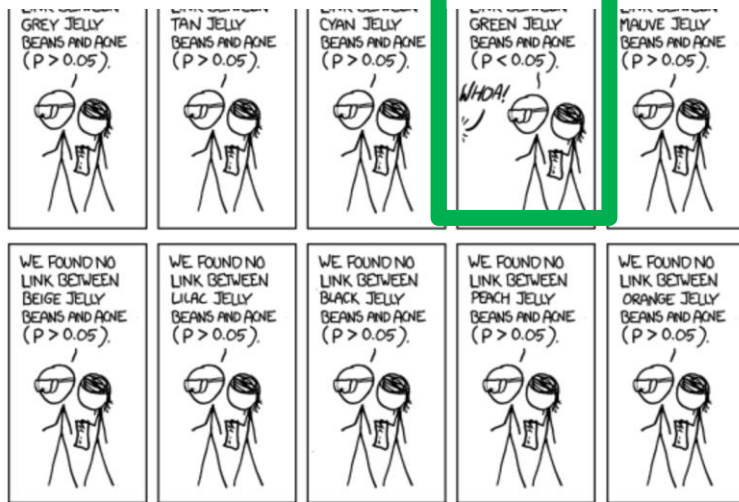
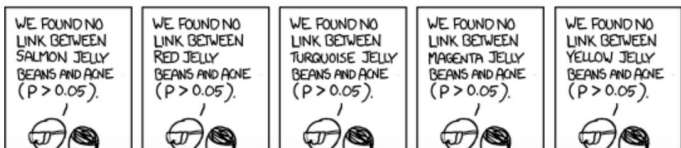
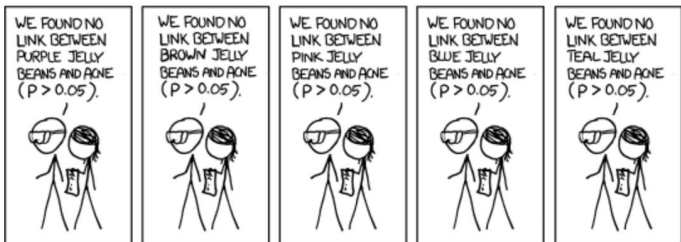
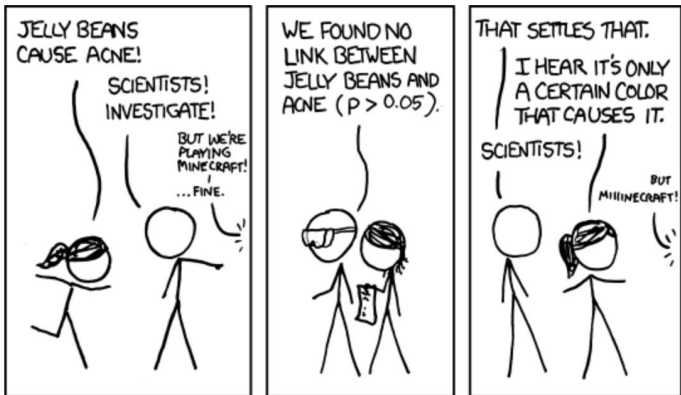




# p-Hacking

SIGNIFICANT

< < PREV RANDOM NEXT > >



# p-Hacking



- Scientists concluded that “Eating green jelly beans gives you more acne” after testing that teenagers who ate green jelly beans have more acne than those who don’t, with a p-value of 0.05”.
  - **The p-value means:** if the null hypothesis is true (teens who eat green jelly beans and those who don’t have the same amount of acne), the probability of observing at least as extreme an outcome as we did is p.
  - **Putting it another way, a p-value of 0.05 means:** only a 5% chance of seeing this much acne if green jelly beans don’t cause acne
  - But what if I repeat the experiment 20 times?
  - The chance that in 20 trials I will never get a p value < 0.05 is  $0.95^{20} \approx 0.358$

In other words 64% of the time one of these tests will be significant. This result has no significance! Happened by random chance!

# p-Hacking

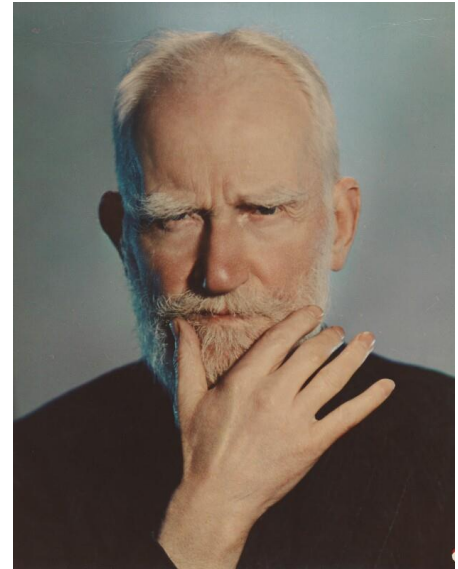
- **Definition**: Performing the same hypothesis test multiple times in order to get a statistically significant result.
- The particularly evil thing: reporting only the significant tests, but not reporting the other 19 tests.....



# Random Quote

“If at first you don’t succeed, try two more times so your failure is statistically significant”.

- George Bernard Shaw



# Random Quote

“Torture numbers, and they’ll confess to anything”

- George Easterbrook



# Another interesting misuse of statistics

Attali/Bar-Hillel noticed that SAT answer keys are not randomized.

Keys are balanced rather than randomized.

Was easy for statisticians to detect by examining published tests.

This is a case of thinking “**randomization is too important to be left to chance**”!

# Suggests a strategy for test-takers

- Answer all the questions you can.
- When guessing the rest, pick an answer position that occurs least frequently in your answers.

Simulations shows this adds 10-16 points over random guessing.

Claimed to be more gain than some very expensive SAT prep courses!

# Conclusions

1. Determine if the samples are **random** and **representative**.
2. Ask for a confidence interval.
3. Be dubious. Be extremely dubious.
4. Don't make up data or statistics. You'll get caught.
5. Be wary of p-hacking (and don't do it yourself)!
6. Be careful about seeing patterns where there are none.
7. Correlation does not imply causation.



# Random Quote



# Random Quote

“Data is the sword of the 21st century, those who wield it well, the Samurai.”

- Jonathan Rosenberg (ex-Google SVP)



# Random Quote

“Do not trust any statistics you did not fake yourself”

- Winston Churchill



# Staring Down a Statistic

1. Who says so?
2. How do they know this is true?
3. What's missing?
4. Did somebody change the subject?
5. Does it make sense?

# Correlation → Causation?



- “People who use Senserdime generally have less cavities than those who use generic brands”.

# Correlation → Causation?



- “People who use Senserdime generally have less cavities than those who use generic brands”.
  - Even if we had a stat-sig p-value (and rejected  $H_0$ ), correlation does not imply causation.
  - Cannot say “Senserdime prevents cavities”.

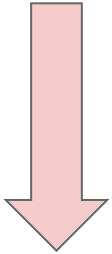
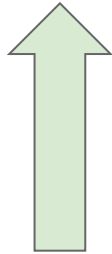
# Correlation → Causation?



- “People who use Senserdime generally have less cavities than those who use generic brands”.
  - Even if we had a stat-sig p-value (and rejected  $H_0$ ), correlation does not imply causation.
  - Cannot say “Senserdime prevents cavities”.
  - Turns out, Senserdime costs \$120,000 per tube. This means only wealthy people can afford it. Wealthy people often have access to good healthcare (e.g., dentists). Senserdime didn’t do anything!

# Correlation $\rightarrow$ Causation?

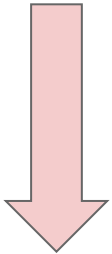
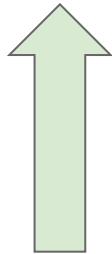
- “When ice cream sales go up, umbrella sales go down.”



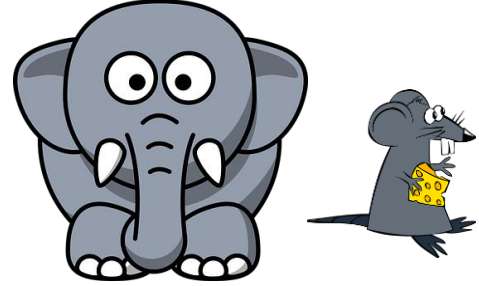


# Correlation → Causation?

- “When ice cream sales go up, umbrella sales go down.”
  - Both generally happen when the weather is sunny.
  - Ice cream sales rise did not CAUSE umbrella sales to go down. The weather CAUSED both of these things to happen.
  - Again, correlation does not imply causation!

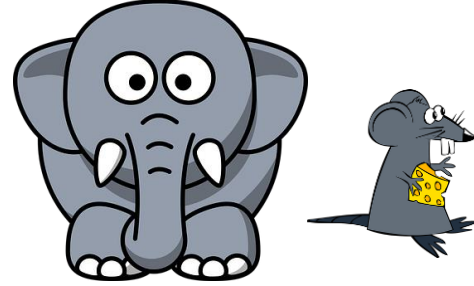


# Size-Based Sampling



- Let's say there are 100 families. 50 families have five children each, and 50 families only have a single child.

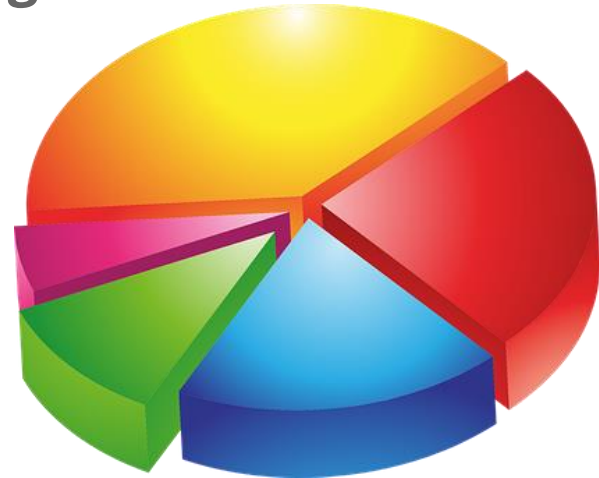
# Size-Based Sampling



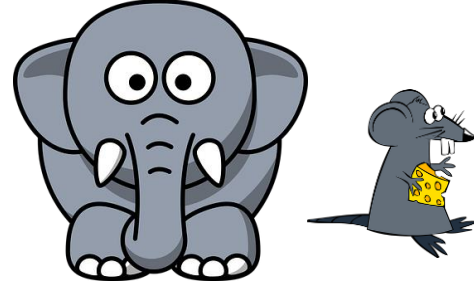
- Let's say there are 100 families. 50 families have five children each, and 50 families only have a single child.
  - What is the expected number of *siblings* a random child has?
    - **Choices:** 0, 1, 2, 2.5, 3, 3.33, 4

# Poll 5

- Let's say there are 100 families. 50 families have five children each, and 50 families only have a single child.
  - What is the expected number of *siblings* a random child has?
    - **Choices:** 0, 1, 2, 2.5, 3, 3.33, 4

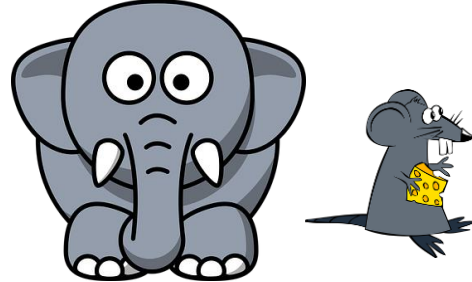


# Size-Based Sampling



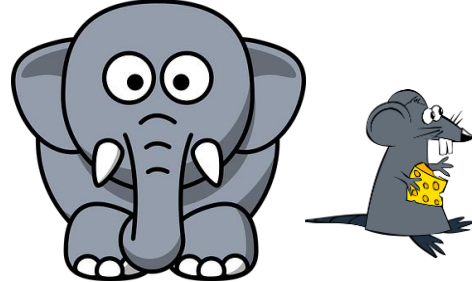
- Let's say there are 100 families. 50 families have five children each, and 50 families only have a single child.
  - What is the expected number of **siblings** a random child has?
    - **Choices:** 0, 1, 2, 2.5, 3, 3.33, 4 (you might guess 2?)

# Size-Based Sampling



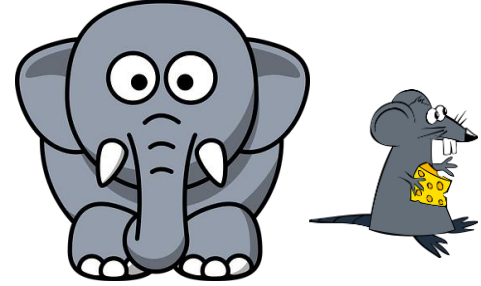
- Let's say there are 100 families. 50 families have five children each, and 50 families only have a single child.
  - What is the expected number of **siblings** a random child has?
    - **Choices:** 0, 1, 2, 2.5, 3, 3.33, 4 (you might guess 2?)
  - There are  $50 \times 5 = 250$  children with 4 siblings, and  $50 \times 1 = 50$  children with 0 siblings.
    - $250/300 * 4 + 50/300 * 0 = 3.3333$

# Size-Based Sampling



- Let's say there are 100 families. 50 families have five children each, and 50 families only have a single child.
  - What is the expected number of **siblings** a random child has?
    - **Choices:** 0, 1, 2, 2.5, 3, 3.33, 4 (you might guess 2?)
  - There are  $50 \times 5 = 250$  children with 4 siblings, and  $50 \times 1 = 50$  children with 0 siblings.
    - $250/300 * 4 + 50/300 * 0 = 3.3333$
  - Actually, it was ambiguous what "random child" meant:

# Size-Based Sampling

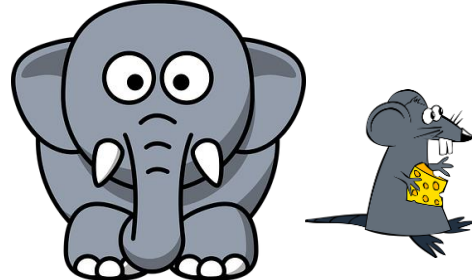


- UW says the average class size is 28. Do you think that is true, or does it feel that way?
- To simplify, let's say there are 300 students, and each student takes exactly one of three classes.

Class	# Students
1	278
2	10
3	12



# Size-Based Sampling

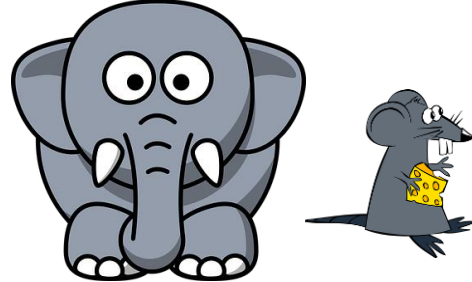


- UW says the average class size is 28. Do you think that is true, or does it feel that way?
- To simplify, let's say there are 300 students, and each student takes exactly one of three classes.

Class	# Students
1	278
2	10
3	12

$$\text{Average (over each class): } 278 \cdot \frac{1}{3} + 10 \cdot \frac{1}{3} + 12 \cdot \frac{1}{3} = 300 \cdot \frac{1}{3} = 100$$

# Size-Based Sampling



- UW says the average class size is 28. Do you think that is true, or does it feel that way?
- To simplify, let's say there are 300 students, and each student takes exactly one of three classes.

Class	# Students
1	278
2	10
3	12

$$\text{Average (over each class): } 278 \cdot \frac{1}{3} + 10 \cdot \frac{1}{3} + 12 \cdot \frac{1}{3} = 300 \cdot \frac{1}{3} = 100$$

$$\text{Average (over each student): } 278 \cdot \frac{278}{300} + 10 \cdot \frac{10}{300} + 12 \cdot \frac{12}{300} \approx 258.43$$

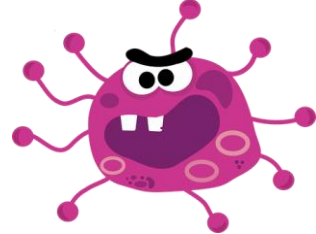
# Random Quote

“Statistics is the grammar of science”

- Karl Pearson

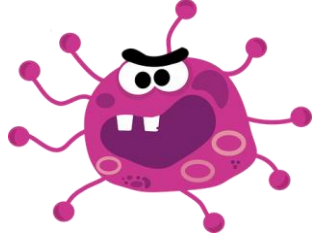


# Conditional Probability



- A disease test is 99% accurate, and 0.005% of the population has the disease. If you test positive: the probability you have the disease is:

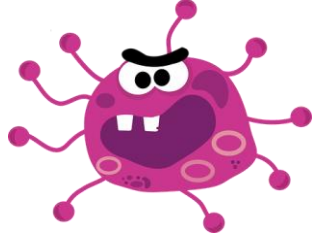
# Conditional Probability



- A disease test is 99% accurate, and 0.005% of the population has the disease. If you test positive: the probability you have the disease is only:

$$P(D|+) = \frac{P(+|D)P(D)}{P(+|D)P(D) + P(+|D^c)P(D^c)} =$$

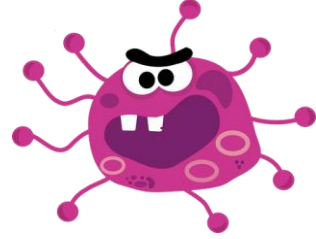
# Conditional Probability



- A disease test is 99% accurate, and 0.005% of the population has the disease. If you test positive: the probability you have the disease is only:

$$P(D|+) = \frac{P(+|D)P(D)}{P(+|D)P(D) + P(+|D^c)P(D^c)} = \frac{0.99 \cdot 0.00005}{0.99 \cdot 0.00005 + 0.01 \cdot 0.9995} \approx 0.49\%$$

# Conditional Probability



- A disease test is 99% accurate, and 0.005% of the population has the disease. If you test positive: the probability you have the disease is only:

$$P(D|+) = \frac{P(+|D)P(D)}{P(+|D)P(D) + P(+|D^c)P(D^c)} = \frac{0.99 \cdot 0.00005}{0.99 \cdot 0.00005 + 0.01 \cdot 0.99995} \approx 0.49\%$$

Much lower than we initially thought! Sometimes non-intuitive...

# Conditional Probability



$P(\text{Attacked by Alien}) = 0.10\%$

$P(\text{Attacked by Alien} \mid \text{AlienShield}) = 0.01\%$



# Conditional Probability



$P(\text{Attacked by Alien}) = 0.10\%$

$P(\text{Attacked by Alien} \mid \text{AlienShield}) = 0.01\%$

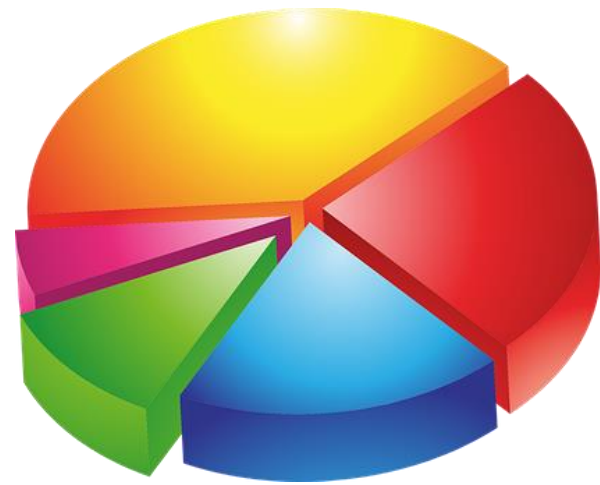
If you are AlienShield, which advertisement do you prefer?

1. (**Relative** Improvement) "AlienShield reduces your chance of getting attacked by an alien 10-fold!"
2. (**Absolute** Improvement) "AlienShield reduces your chance of getting attacked by an alien by 0.09%."

# Poll 7

$P(\text{Attacked by Alien}) = 0.10\%$

$P(\text{Attacked by Alien} \mid \text{AlienShield}) = 0.01\%$



If you are AlienShield, which advertisement do you prefer?

1. (**Relative** Improvement) "AlienShield reduces your chance of getting attacked by an alien 10-fold!"
2. (**Absolute** Improvement) "AlienShield reduces your chance of getting attacked by an alien by 0.09%."

# Conditional Probability



$P(\text{Attacked by Alien}) = 0.10\%$

$P(\text{Attacked by Alien} \mid \text{AlienShield}) = 0.01\%$

If you are AlienShield, which advertisement do you prefer?

1. (**Relative** Improvement) "AlienShield reduces your chance of getting attacked by an alien 10-fold!"
2. (**Absolute** Improvement) "AlienShield reduces your chance of getting attacked by an alien by 0.09%."

Watch for which type of improvement is cited, and consider if the original probability was already low or high.

# Conditional Probability



- Suppose there is a carnival game which gives out prizes, and three types of players: children, teenagers, and adults.
- Bob thinks the carnival unfairly gives more prizes to children over the other types of players. Is this true?

# Conditional Probability

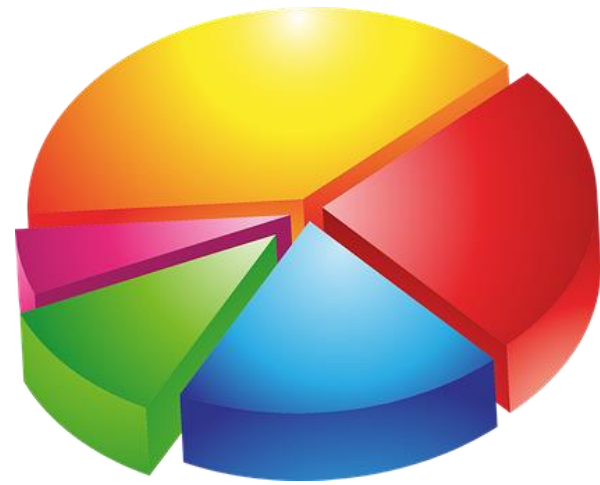


- Suppose there is a carnival game which gives out prizes, and three types of players: children, teenagers, and adults.
- Bob thinks the carnival unfairly gives more prizes to children over the other types of players. Is this true?

Player Type	% Prizes Won
Children	70%
Teenagers	5%
Adults	25%

# Poll 8a

Is this unfair?



Player Type	% Prizes Won
Children	70%
Teenagers	5%
Adults	25%

# Conditional Probability

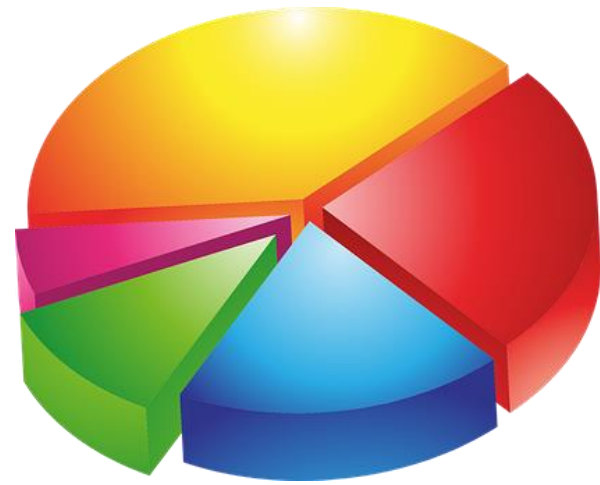


- Suppose there is a carnival game which gives out prizes, and three types of players: children, teenagers, and adults.
- Bob thinks the carnival unfairly gives more prizes to children over the other types of players. Is this true?

<b>Player Type</b>	<b>% Prizes Won</b>	<b>% Global Population</b>
Children	70%	25%
Teenagers	5%	15%
Adults	25%	60%

# Poll 8b

Is this unfair?



<b>Player Type</b>	<b>% Prizes Won</b>	<b>% Global Population</b>
Children	70%	25%
Teenagers	5%	15%
Adults	25%	60%



# Conditional Probability

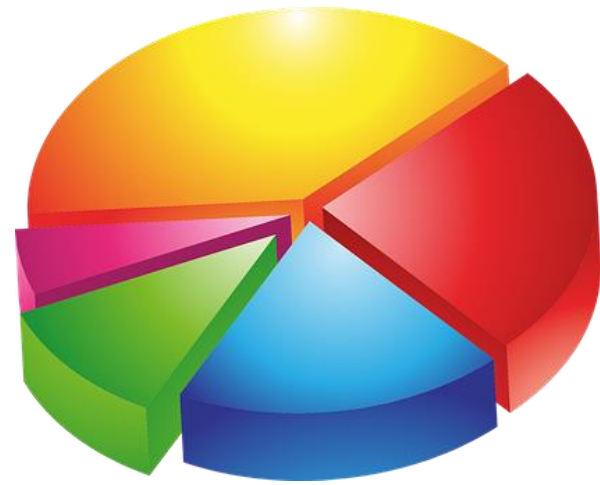


- Suppose there is a carnival game which gives out prizes, and three types of players: children, teenagers, and adults.
- Bob thinks the carnival unfairly gives more prizes to children over the other types of players. Is this true?

Player Type	% Prizes Won	% Global Population	% Carnival Population
Children	70%	25%	71%
Teenagers	5%	15%	4.5%
Adults	25%	60%	24.5%

# Poll 8c

Is this unfair?



<b>Player Type</b>	<b>% Prizes Won</b>	<b>% Global Population</b>	<b>% Carnival Population</b>
Children	70%	25%	71%
Teenagers	5%	15%	4.5%
Adults	25%	60%	24.5%

# Conditional Probability



- Suppose there is a carnival game which gives out prizes, and three types of players: children, teenagers, and adults.
- Bob thinks the carnival unfairly gives more prizes to children over the other types of players. Is this true?

Player Type	% Prizes Won	<del>% Global Population</del>	% Carnival Population
Children	70%	25%	71%
Teenagers	5%	15%	4.5%
Adults	25%	60%	24.5%

Looks very fair now!

# Conditional Probability



$$P(\text{child} \mid \text{prize}) = 70\%$$

$$P(\text{child}) = 71\%$$

$$P(\text{teen} \mid \text{prize}) = 5\%$$

$$P(\text{teen}) = 4.5\%$$

$$P(\text{adult} \mid \text{prize}) = 25\%$$

$$P(\text{adult}) = 24.5\%$$

Player Type	% Prizes Won	<del>% Global Population</del>	% Carnival Population
Children	70%	25%	71%
Teenagers	5%	15%	4.5%
Adults	25%	60%	24.5%

Player Type and Prize are (almost) independent!

# Conditional Probability



$$P(\text{child} \mid \text{prize}) = 70\%$$

$$P(\text{child}) = 71\%$$

$$P(\text{teen} \mid \text{prize}) = 5\%$$

$$P(\text{teen}) = 4.5\%$$

$$P(\text{adult} \mid \text{prize}) = 25\%$$

$$P(\text{adult}) = 24.5\%$$

Hypothesis Test: “chi-squared test of independence”

Player Type and Prize are (almost) independent!

# Conditional Probability

Statement: “Most people who win a nobel prize went to college.”

- $P(\text{college} \mid \text{nobel prize}) \approx 1$

# Conditional Probability

**Statement:** “Most people who win a nobel prize went to college.”

- $P(\text{college} \mid \text{nobel prize}) \approx 1$

**Misinterpretation:** “If you go to college, you’ll win a nobel prize!”

- $P(\text{nobel prize} \mid \text{college}) \approx 0$

# Conditional Probability

**Statement:** “Most people who win a nobel prize went to college.”

- $P(\text{college} \mid \text{nobel prize}) \approx 1$

**Misinterpretation:** “If you go to college, you’ll win a nobel prize!”

- $P(\text{nobel prize} \mid \text{college}) \approx 0$

There is a big difference between  $P(A \mid B)$  and  $P(B \mid A)$ !!!



# Gambler's Fallacy

- “Play another round of blackjack - you have to win soon! You’ve been losing too much.”



# Gambler's Fallacy



- “Play another round of blackjack - you have to win soon! You’ve been losing too much.”
  - Each game/trial is **independent**, and so even if you already lost 10 times, the probability you win the next game is the same as any other.

# Gambler's Fallacy



- “Play another round of blackjack - you have to win soon! You’ve been losing too much.”
  - Each game/trial is **independent**, and so even if you already lost 10 times, the probability you win the next game is the same as any other.
  - “Memorylessness” for Geometric RV.

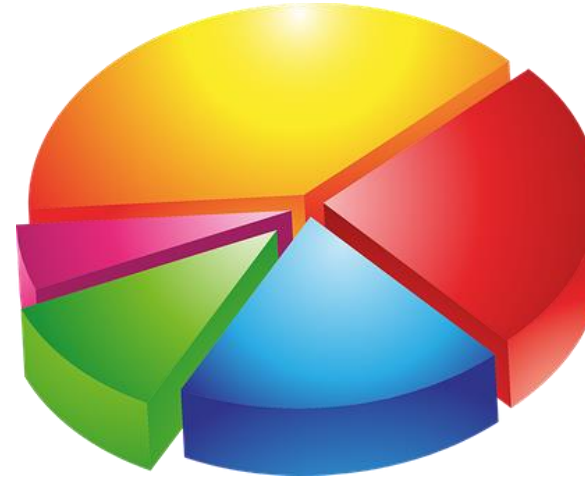
# Gambler's Fallacy



- “Play another round of blackjack - you have to win soon! You’ve been losing too much.”
  - Each game/trial is **independent**, and so even if you already lost 10 times, the probability you win the next game is the same as any other.
  - “Memorylessness” for Geometric RV.
  - $P(\text{win} \mid 100 \text{ losses}) = P(\text{win} \mid 10 \text{ losses}) = P(\text{win})$

# Poll 9

What advice would you give to your friend who has lost 10 consecutive hands of HoldEm and nearly \$1000?



# Gambler's Fallacy



**Terrible Advice:** “Play another round of blackjack - you have to win soon! You’ve been losing too much.”

# Gambler's Fallacy



**Terrible Advice:** “Play another round of blackjack - you have to win soon! You’ve been losing too much.”

**Good Advice:** “Cut your losses and go home. Quit while you’re ahead”.

# Gambler's Fallacy



**Terrible Advice:** “Play another round of blackjack - you have to win soon! You’ve been losing too much.”

**Good Advice:** “Cut your losses and go home. Quit while you’re ahead”.

**Best Advice:** “Stop gambling, you idiot.” - A Caring Friend (who understands statistics).



# Random Quote

“I guess I think of lotteries as a tax on the mathematically challenged.” - Roger Jones

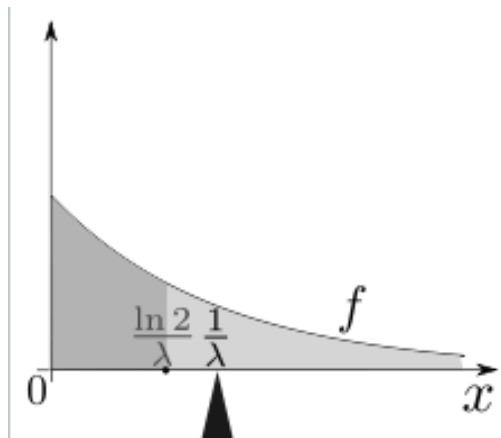


# The Well-Chosen Average

- **Mean** (average of all values weighted by probability or density)

Let  $X \sim \text{Exp}(\lambda)$ .

$$E[X] = \frac{1}{\lambda}$$



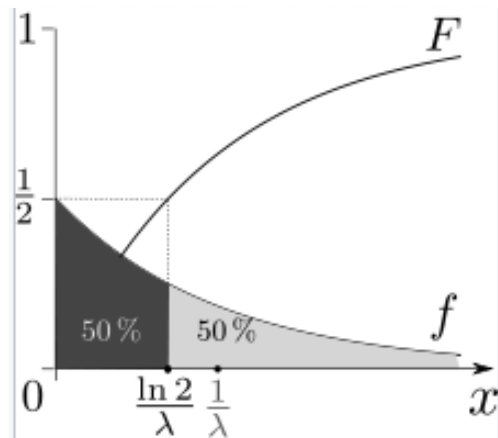
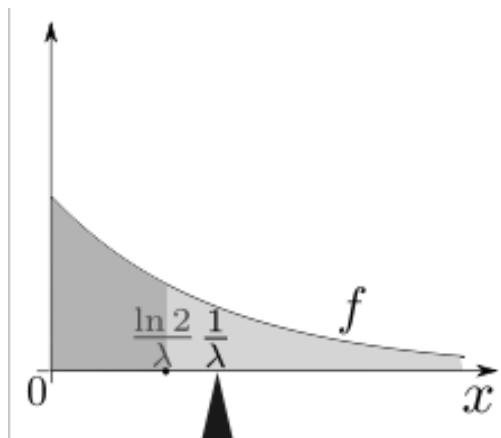
# The Well-Chosen Average

- **Mean** (average of all values weighted by probability or density)
- **Median** (the point  $m$  where 1/2 values are larger, and 1/2 are smaller)

Let  $X \sim \text{Exp}(\lambda)$ .

$$E[X] = \frac{1}{\lambda}$$

$$\text{median}[X] = \frac{\ln 2}{\lambda}$$



# The Well-Chosen Average

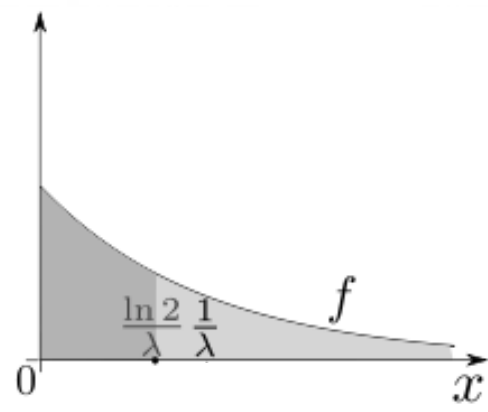
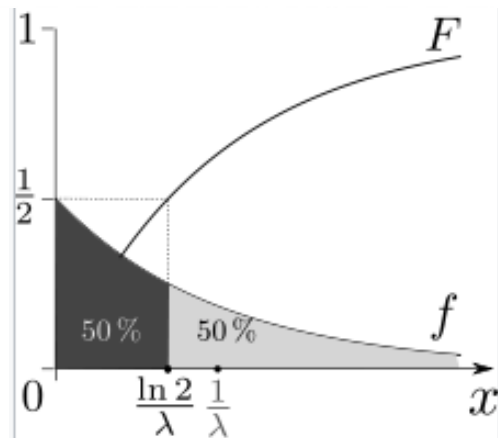
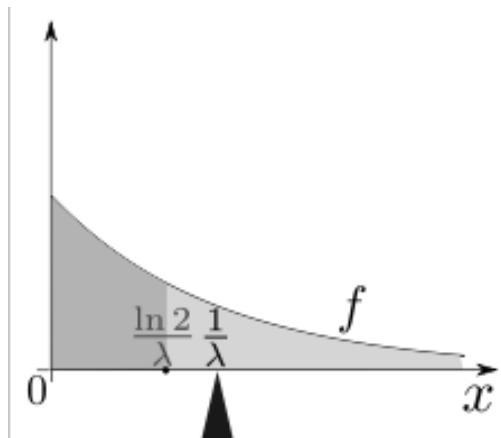
- **Mean** (average of all values weighted by probability or density)
- **Median** (the point  $m$  where 1/2 values are larger, and 1/2 are smaller)
- **Mode** (the point with highest probability or density)

Let  $X \sim \text{Exp}(\lambda)$ .

$$E[X] = \frac{1}{\lambda}$$

$$\text{median}[X] = \frac{\ln 2}{\lambda}$$

$$\text{mode}[X] = 0$$



# The Well-Chosen Average

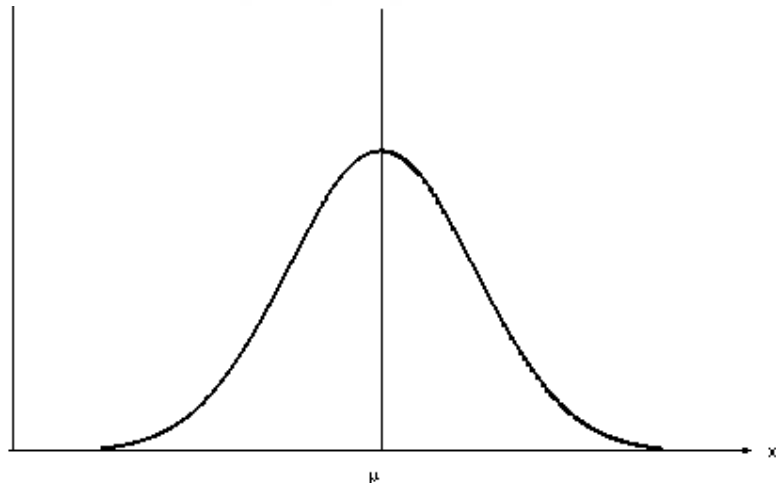
- **Mean** (average of all values weighted by probability or density)
- **Median** (the point  $m$  where 1/2 values are larger, and 1/2 are smaller)
- **Mode** (the point with highest probability or density)

Let  $X \sim N(\mu, \sigma^2)$ .

$$E[X] = \mu$$

$$\text{median}[X] = \mu$$

$$\text{mode}[X] = \mu$$



# The Well-Chosen Average

- Are haircuts more expensive in Toronto or Vancouver?

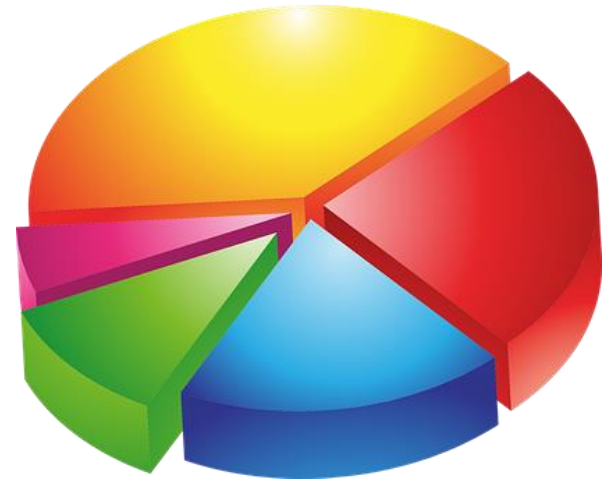


Haircut Prices	Vancouver	Toronto
$x_1$	\$20	\$15
$x_2$	\$20	\$25
$x_3$	\$22	\$25
$x_4$	\$24	\$29
$x_5$	\$25	\$35
$x_6$	\$28	\$45
$x_7$	\$400	\$65

# Poll 2

Are haircuts more expensive in Toronto or Vancouver?

Haircut Prices	Vancouver	Toronto
$x_1$	\$20	\$15
$x_2$	\$20	\$25
$x_3$	\$22	\$25
$x_4$	\$24	\$29
$x_5$	\$25	\$35
$x_6$	\$28	\$45
$x_7$	\$400	\$65



# The Well-Chosen Average

- Are haircuts more expensive in Toronto or Vancouver?

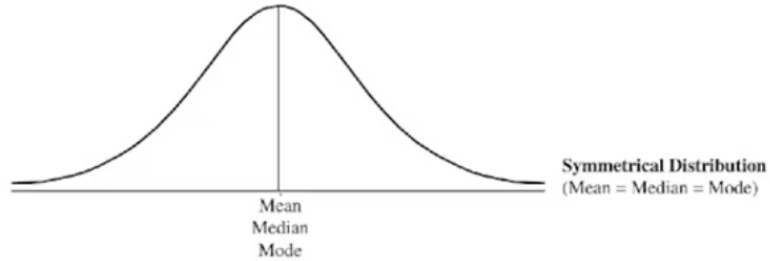


Haircut Prices	Vancouver	Toronto
$x_1$	\$20	\$15
$x_2$	\$20	\$25
$x_3$	\$22	\$25
$x_4$	\$24	\$29
$x_5$	\$25	\$35
$x_6$	\$28	\$45
$x_7$	\$400	\$65
Mean	\$77	\$36
Median	\$24	\$29
Mode	\$20	\$25



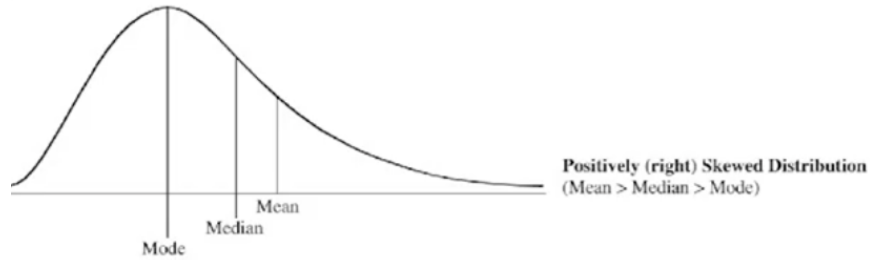
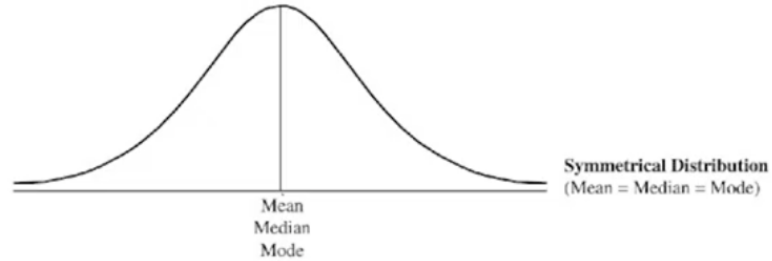
# The Well-Chosen Average

---



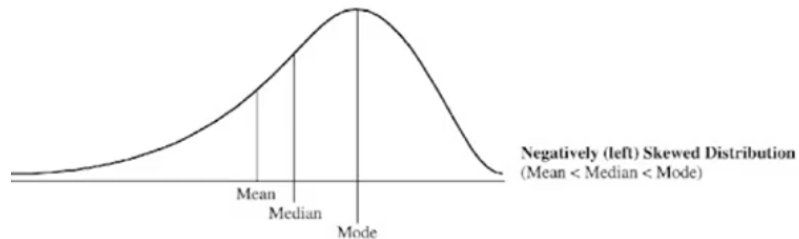
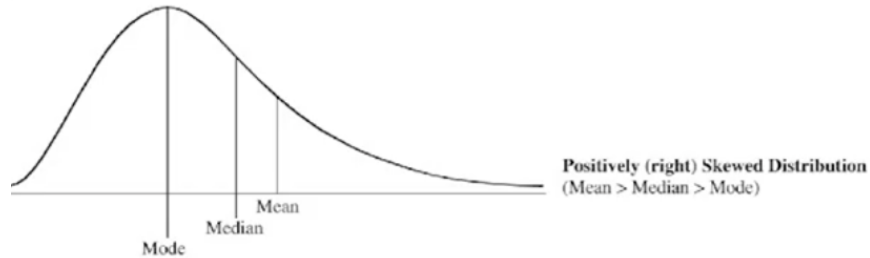
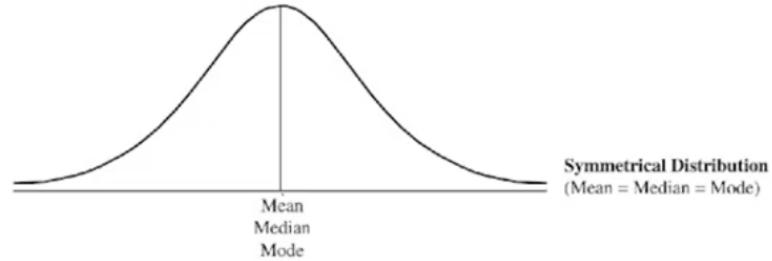
# The Well-Chosen Average

---



# The Well-Chosen Average

---



# The Well-Chosen Average

- **Mean:** Heavily affected/influenced by outliers. Any extreme value(s) may make this measure terrible.
- **Median:** About half the values are higher, and half are lower than this.
- **Mode:** The most frequently occurring value.

Which is “best”?

# The Well-Chosen Average

- **Mean:** Heavily affected/influenced by outliers. Any extreme value(s) may make this measure terrible.
- **Median:** About half the values are higher, and half are lower than this.
- **Mode:** The most frequently occurring value.

Which is “best”?

It depends, and it's good to know all of them for a better idea of the distribution!

# Conclusions

1. Determine if the samples are **random** and **representative**.
2. Ask for a confidence interval.
3. Be dubious. Be extremely dubious.
4. Don't make up statistics. You'll get caught.
5. Be wary of p-hacking (and don't do it yourself)!
6. Be careful about seeing patterns where there are none.
7. Correlation does not imply causation.
8. Be careful with interpreting conditional probabilities.  
Intuition sometimes doesn't work here!
9. Be wary of assuming things are independent that aren't independent.