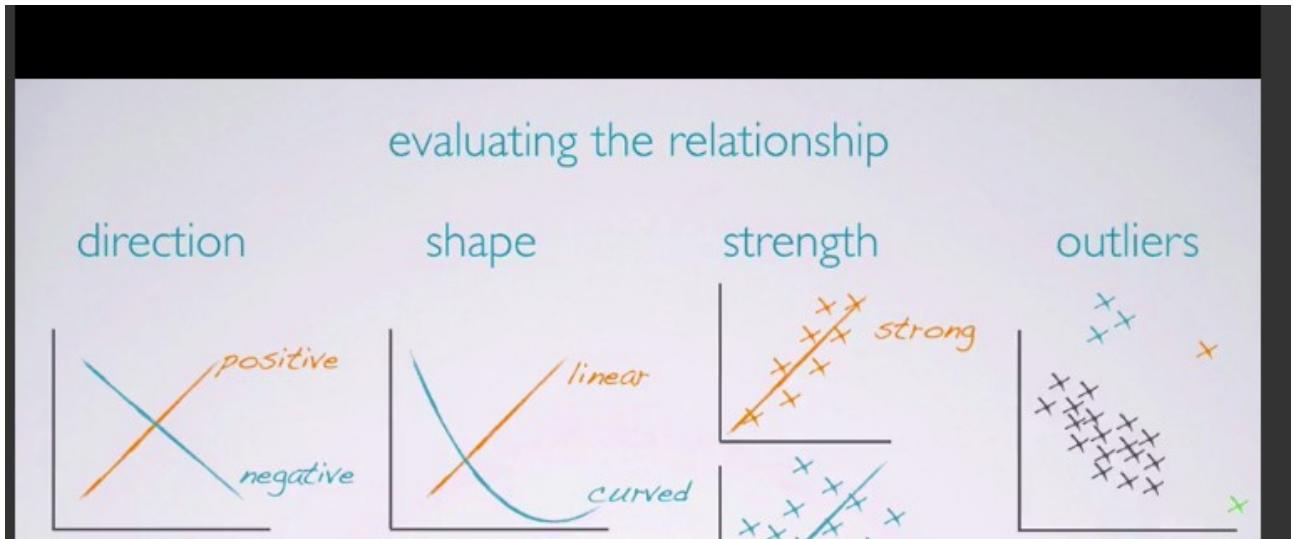


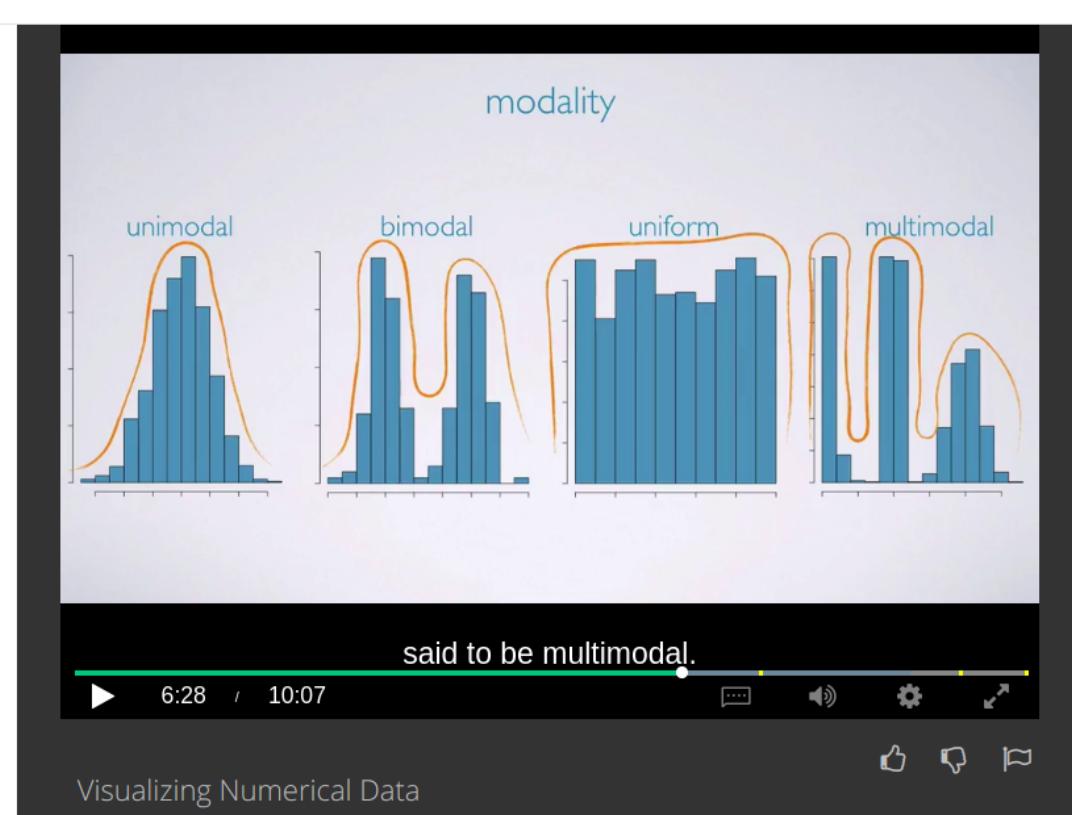
COURSE 1 → WEEK 2 NOTES.

Topic 1.6 of the book.
Visualizing the Numeric data:



histogram helps us to determine the dist of the data, it provides a view of data density.

Skewness : dists are skewed to the side of the long tail.



The bimodal shows that there are 2 groups in a data set while uniform shows that there's no trend every value in a data is equally likely to occur.

While symmetric shows that most of the readings/values in a data set are around the mean.

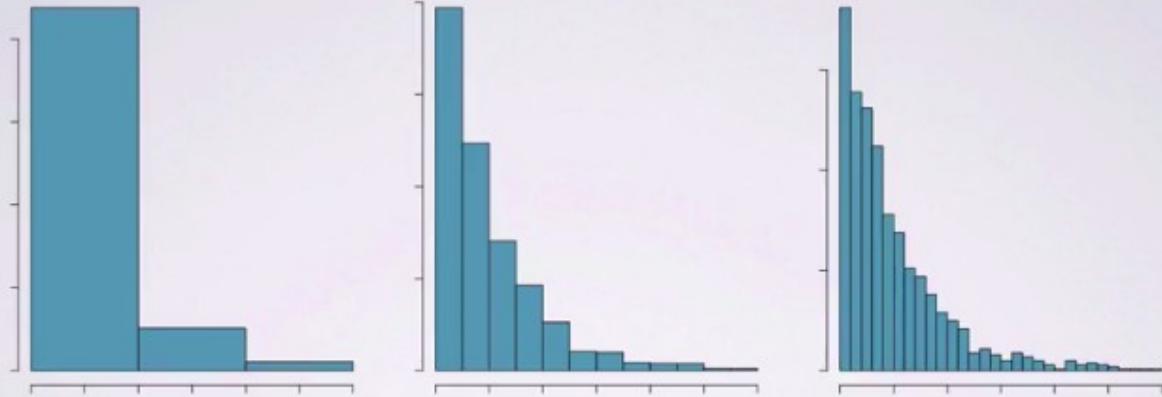
modality (cont.)

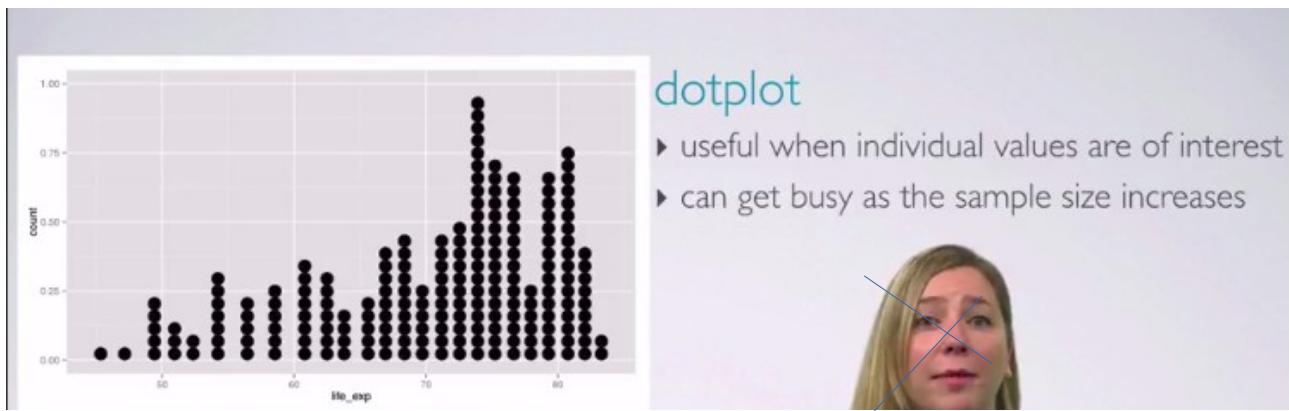


Here the last graph of the above pic is the plot of the last digit of Social security number of US citizen ,which is a Uniform dist as expected.

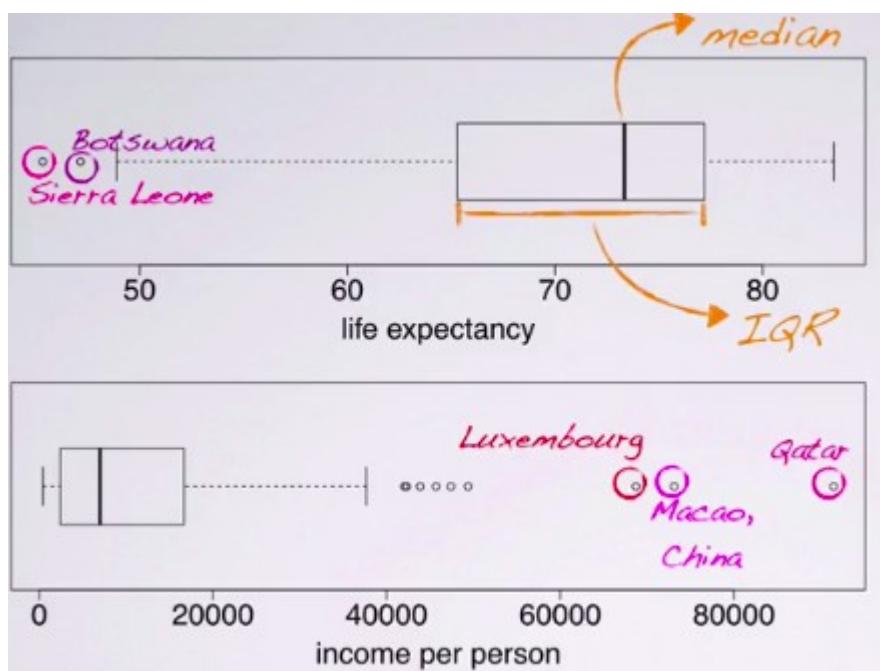
histogram & bin width

The chosen bin width can alter the story the histogram is telling.

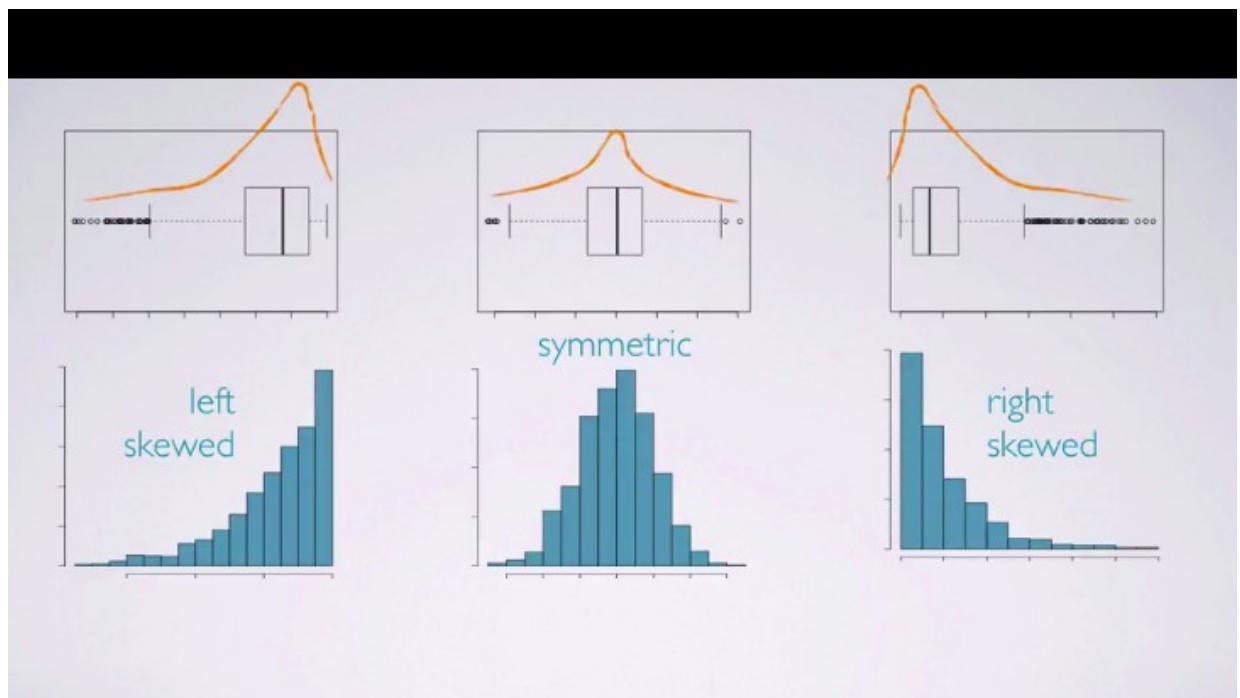




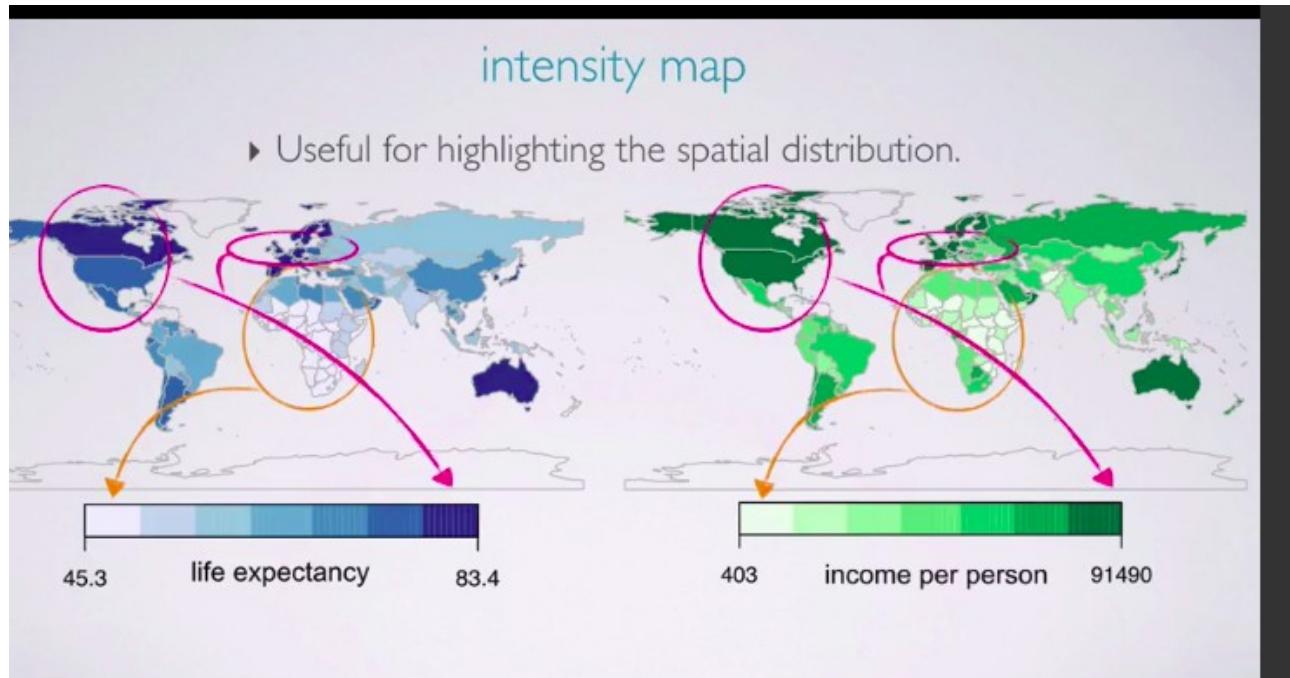
BOX PLOTS: useful for highlighting outliers and IQR.(inter quartile range):



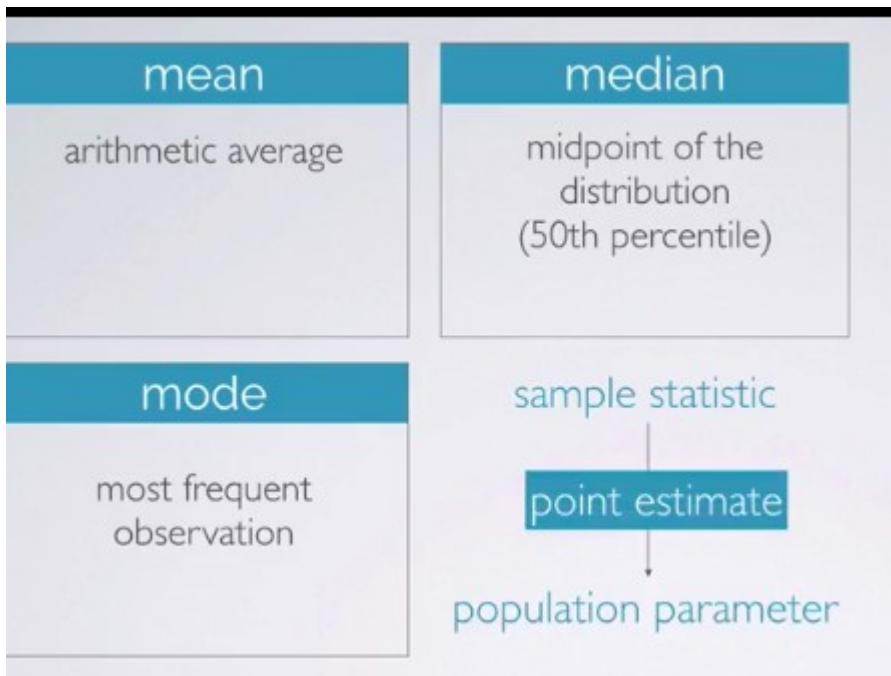
they do not show the modality, but they do show skewness.



INTENSITY MAP:



MEASURE OF CENTER:



Another key characteristic that is of interest is the center of the distribution, commonly used measures of center are the mean, which is simply the arithmetic average. The median, which is the mid point of the distribution or in other words the 50th percentile and the mode which is the most frequent observation. If these measurements are calculated from a sample, they're called sample statistics. Sample statistics are point estimates for the unknown population parameters.

We usually use letters from the Latin alphabet when denoting sample statistics, and letters from the Greek alphabet when denoting population parameters. For example, the sample mean is \bar{x} and the population mean is μ .

Which statement is true:

(Hint: Sketching the distributions might be useful.)

less than 50 is noticeable.

- In a symmetric distribution, more than 50% of the data are below and less than 50% are above the mean.
- In a left skewed distribution, roughly 50% of the data are below and 50% are above the mean.
- In a right skewed distribution, less than 50% of the data are below the mean.
- In a left skewed distribution, less than 50% of the data are below the mean.

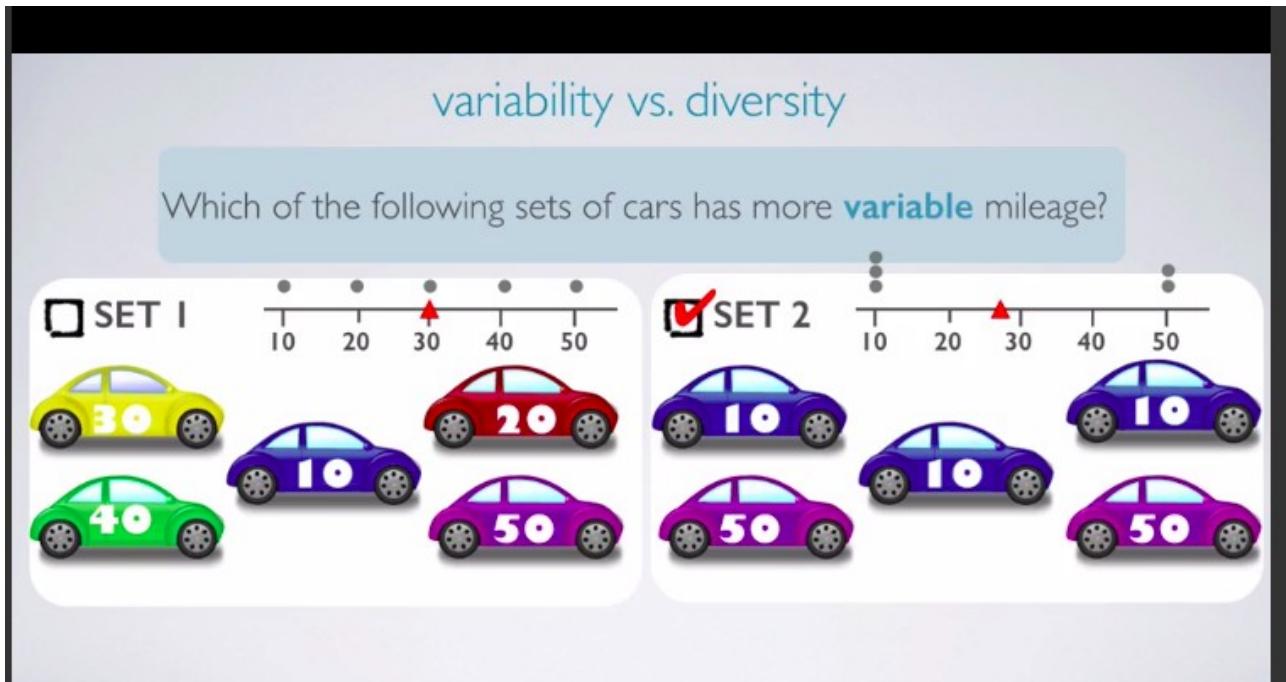
Correct

Since the median marks the 50th percentile, and in a left-skewed distribution the mean is smaller than the median, less than 50% of the data will be smaller than the mean.

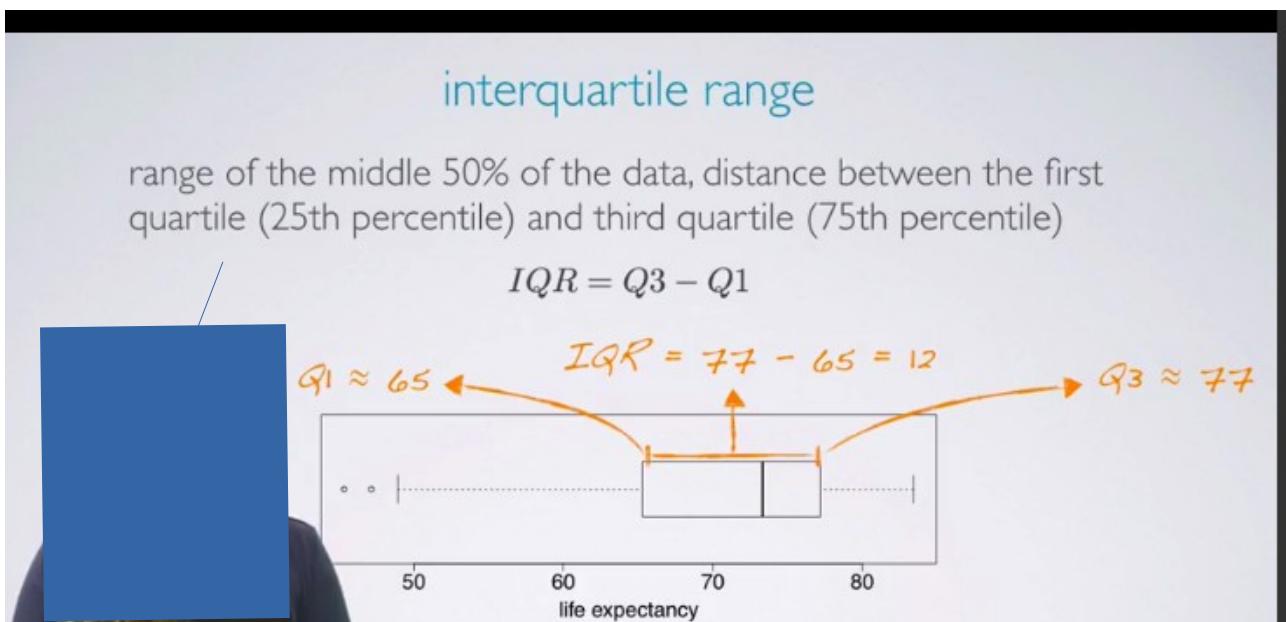
Measure of spreads:

range, variance, std and IQR.

The set 1 is more diverse but the set 2 is more variable bcz the **means** are 30 and 26 which is



roughly equal but the values aren't centered in the set 2 they are at end points.



Robust statistics:

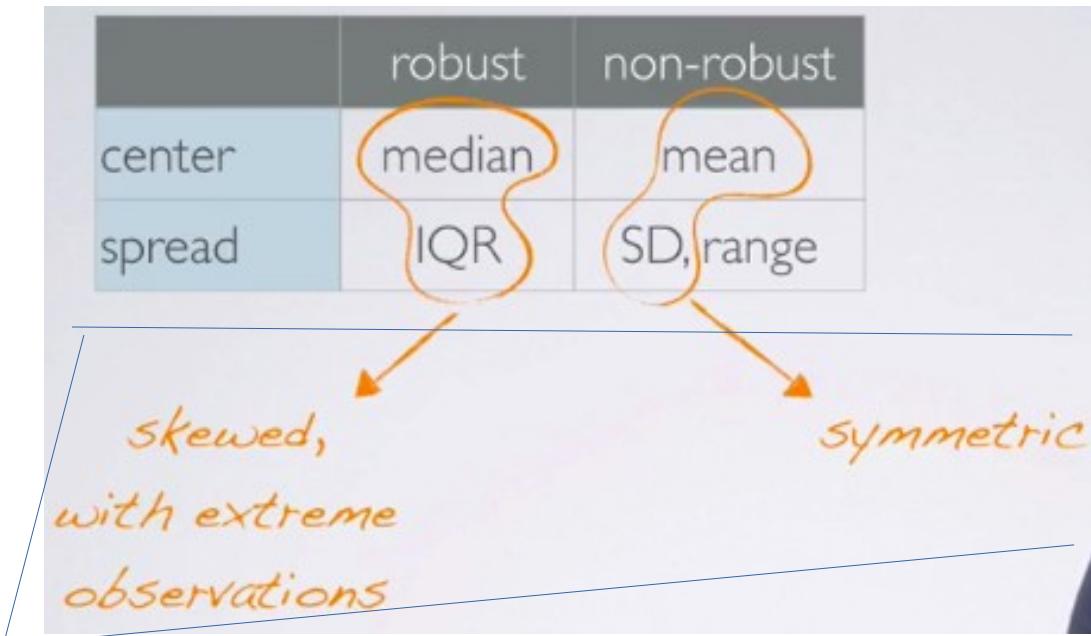
robust statistics

we define robust statistics as measures on which extreme observations have little effect

example

data	mean	median
1, 2, 3, 4, 5, 6	3.5	3.5
1, 2, 3, 4, 5, 1000	169	3.5

median is more robust than the mean because it doesn't rely on the extreme values.



Robusts are used for skewed distributions while we use mean and std for symmetric.

You have collected annual salary data from a large company with many employees who make below \$100,000 per year, a fewer number of managers with salaries between \$100,000 - \$150,000, and a few high-level executives whose salaries can go beyond \$1 million per year. Determine what shape the distribution of these salaries would be expected to follow, and accordingly decide whether the median or the mean would best represent a typical salary for an employee working at this company.

- right skewed, mean is a better measure of typical salary
- right skewed, median is a better measure of typical salary

Correct

Majority of the distribution is below \$100,000 and as the salary increases the number of employees who make as high as salary decreases, hence giving the distribution a long tail on the right (right skewed). For skewed distributions the median is a better measure of the typical observation.

- symmetric, mean is a better measure of typical salary
- symmetric, median is a better measure of typical salary
- left skewed, mean is a better measure of typical salary
- left skewed, median is a better measure of typical salary

TRANSFORMING DATA:

Goals of transforming the data:

We might want to see the data structure a little differently.

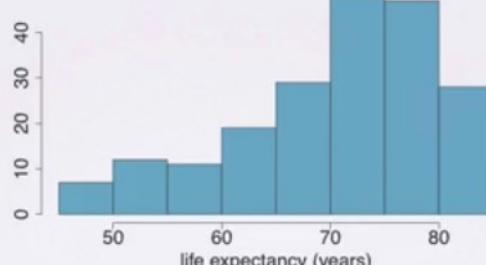
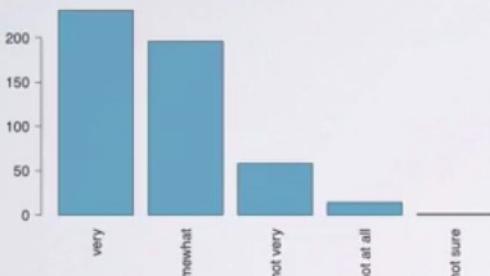
We might want to reduce skew to assist in modeling

We might want to straighten a nonlinear relationship in a scatterplot, so that we can model the relationship with simpler methods.

CATEGORICAL VARIABLES:

How are bar plots different than histograms?

- ▶ barplots for categorical variables, histograms for numerical variables
- ▶ x-axis on a histogram is a number line, and the ordering of the bars are not interchangeable



This [video](#) is must for exploring categorical variables.

For single categorical variable we use frequency distribution table while for more than 2 categorical variable we use contingency table.

We use segmented barplot for contingency table.

STATISTICAL INFERENCE:

We saw a difference of roughly 30% between the proportion of male and female files that are promoted. Based on this information, which of the below might be true? Select all that apply.

- If we were to repeat the experiment we will definitely see that more female files get promoted. This was a fluke.

Un-selected is correct

- Males are more likely to be promoted, and hence there is gender discrimination against women in promotion decisions.

Correct

Both of these options are possible – the difference might be indicative of discrimination against women in promotion decisions, or it might just be due to chance. We need further analysis to make a decision between these two competing claims.

- The difference in the proportions of promoted male and female files is due to chance, this is not evidence of gender discrimination.

Correct

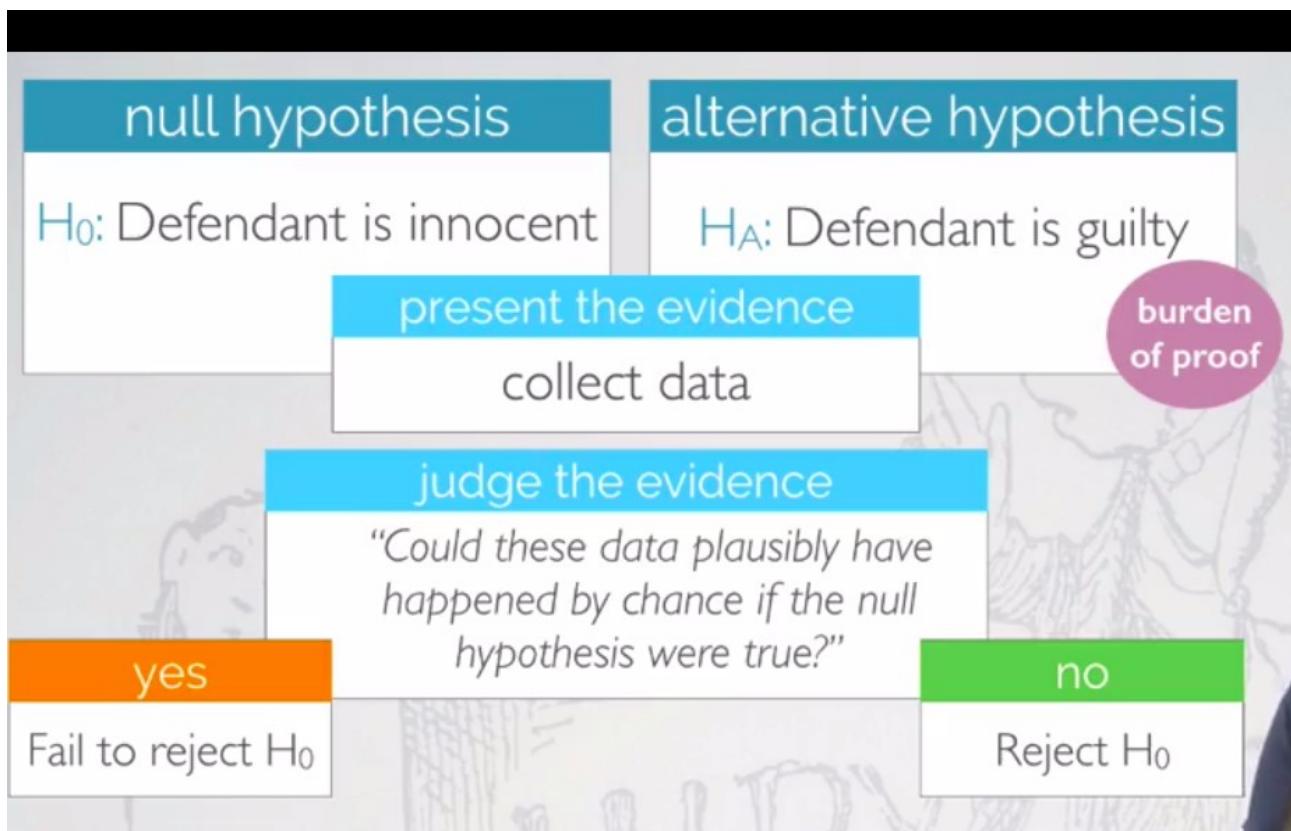
Both of these options are possible – the difference might be indicative of discrimination against women in promotion decisions, or it might just be due to chance. We need further analysis to make a decision between these two competing claims.

- Women are less qualified than men, and this is why fewer females get promoted.

Un-selected is correct

if the null hypothesis is true than the Z statistic has a std normal distribution

Z-scores can be calculated for distributions that are not normal.



Hypothesis testing is very much like a court trial in the US. The null hypothesis says that the defendant is innocent and the alternative hypothesis says that the defendant is guilty. We then present evidence or, or in other words, collect data. Then, we judge this evidence and ask ourselves the question, could these data plausibly have happened by chance if the null hypothesis were true?

If the data were likely to have occurred under the assumption that the null hypothesis were true, then we would fail to reject the null hypothesis, and state that the evidence is not sufficient to suggest that the defendant is guilty. Note that when this happens, the jury returns with a verdict of not guilty. The jury does not say the defendant is innocent, just that there is not enough evidence to convict. The defendant may in fact be innocent but the jury has no way of being sure. Said statistically, we fail to reject the null hypothesis. We never declare the null hypothesis to be true. Because we do not know and cannot prove whether it's true or not. Therefore, we also never say that we would accept the null hypothesis. If the data were very unlikely to have occurred, then the evidence raises more than a reasonable doubt in our minds about the null hypothesis, and hence we reject the null hypothesis

recap: hypothesis testing framework

- ▶ start with a **null hypothesis (H_0)** that represents the status quo
- ▶ set an **alternative hypothesis (H_A)** that represents the research question, i.e. what we're testing for
- ▶ conduct a hypothesis test under the assumption that the null hypothesis is true, either via simulation or theoretical methods
 - ▶ if the test results suggest that the data do not provide convincing evidence for the alternative hypothesis, stick with the null hypothesis
 - ▶ if they do, then reject the null hypothesis in favor of the alternative

how to make decision?

Case study is the gender discrimination in promotion.

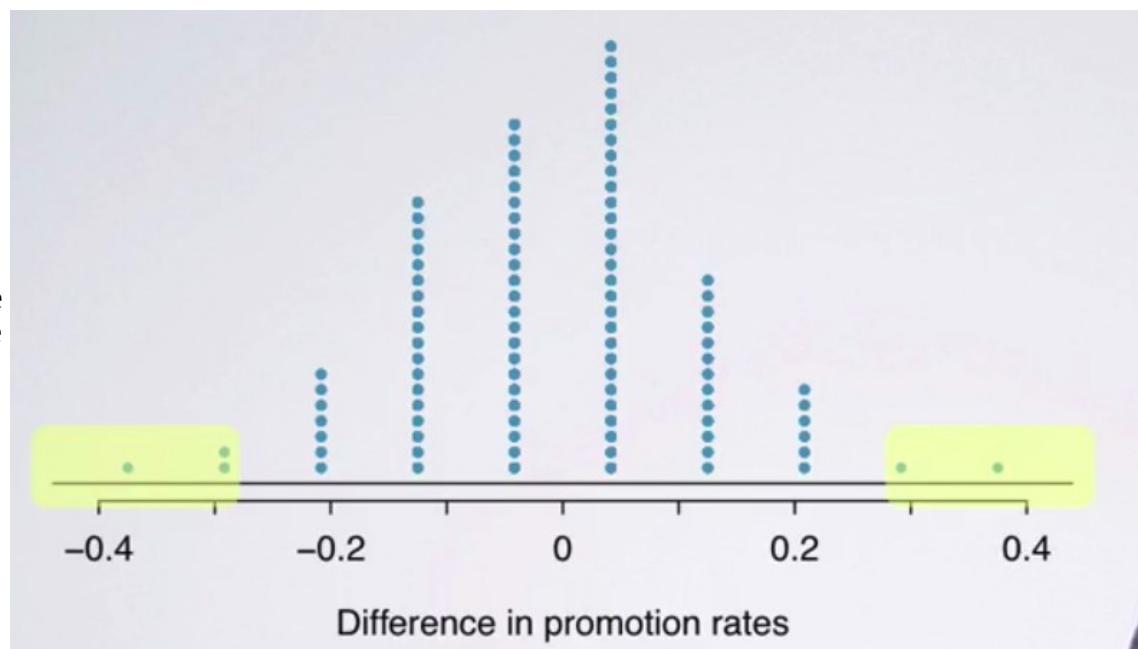
making a decision

- ▶ results from the simulations look like the data → the difference between the proportions of promoted files between males and females was **due to chance** (promotion and gender are **independent**)
- ▶ results from the simulations do not look like the data → the difference between the proportions of promoted files between males and females was **not** due to chance, but **due to an actual effect of gender** (promotion and gender are **dependent**)

We simulated the experiment and obtained 100 samples I.e randomly assigning the files to the officers for promotion of males and females and we observed the diff between the proportions of them being promoted or not we then took the difference between their proportions and plot it using dot plot.

The plot shows that the promotion is independent of gender because the difference is highly centered around the zero. And there's only 5% of the time when the difference between them is 30%.

therefore we reject our null hypothesis that the promotion is dependent on gender discrimination.



summary

p-value

- ▶ set a null and an alternative hypothesis
- ▶ simulate the experiment assuming that the null hypothesis is true
- ▶ evaluated the probability of observing an outcome at least as extreme as the one observed in the original data
- ▶ and if this probability is low, reject the null hypothesis in favor of the alternative

P-VALUE

The p -value is a measure of the strength of the evidence against the null hypothesis.

if the null hypothesis is true than the Z statistic has a std normal distribution

if p value is $<$ alpha then we will reject the null hypothesis.

Week 3

random process

In a **random process** we know what outcomes could happen, but we don't know which particular outcome will happen.



Image sources:

<https://61368956@N00/3100369536.tango-icon-theme-0.7.1.tar.gz>

probability

$P(A) =$
Probability
of event A

There are several possible interpretations of probability but they (almost) completely agree on the mathematical rules probability must follow:

$$0 \leq P(A) \leq 1$$

frequentist interpretation

The probability of an outcome is the proportion of times the outcome would occur if we observed the random process an infinite number of times.

bayesian interpretation

A Bayesian interprets probability as a subjective degree of belief.

Largely popularized by revolutionary advance in computational technology and methods during the last twenty years.

For the same event, two separate people could have different viewpoints and

2:22 / 5:49

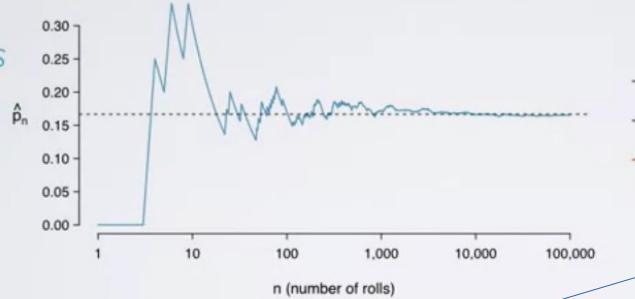
ckr.com/photos/dahlstroms/527634847

in bayesian interpretation 2 people can assign diff probability to the same event.

law of large numbers

law of large numbers states that as more observations are collected, the proportion of occurrences with a particular outcome converges to the probability of that outcome.

examples



- exactly 3 heads in 10 coin flips
- exactly 3 heads in 100 coin flips
- exactly 3 heads in 1000 coin flips

This is more unusual

Say you toss a coin 10 times, and it lands on Heads each time. What do you think the chance is that another head will come up on the next toss? 0.5, less than 0.5, or more than 0.5?

H H H H H H H H H H ?

The probability is still 50%:

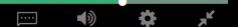
$$\begin{aligned} P(H \text{ on the 11th toss}) \\ = P(H \text{ on the 10th toss}) \\ = 0.50 \end{aligned}$$

The coin is **not** due for a tail.

Common misunderstanding of law of large numbers:
gambler's fallacy
(law of averages)

Photo by thegrid-ch on Flickr (<http://www.flickr.com/photos/thegrid-ch/5087085391>)

▶ 5:20 / 5:49



disjoint (mutually exclusive)

disjoint (mutually exclusive) events cannot happen at the same time.

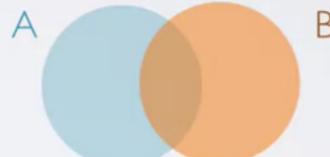
- the outcome of a single coin toss cannot be a head and a tail.
- a student can't both fail and pass a class.
- a single card drawn from a deck cannot be an ace and a queen.



$$P(A \text{ and } B) = 0$$

non-disjoint events can happen at the same time.

- a student can get an A in Stats and A in Econ in the same semester.



$$P(A \text{ and } B) \neq 0$$

of event A and B happening at the same time is non 0.

▶ 1:38 / 9:28



Can an event be disjoint and independent?

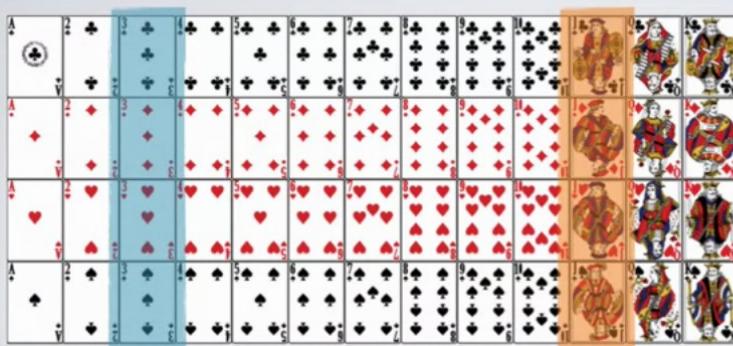


3 Answers. Two disjoint events can never be independent, except in the case that one of the events is null. ... Events are considered disjoint if they never occur at the same time. For example, being a freshman and being a sophomore would be considered disjoint events. Jun 20, 2016

Probability: Are disjoint events independent? - Mathematics Stack ...
<https://math.stackexchange.com/questions/.../probability-are-disjoint-events-independent>

union of disjoint events

What is the probability of drawing a Jack or a three from a well shuffled full deck of cards?



$$P(J \text{ or } 3)$$

$$= P(J) + P(3)$$

$$= (4/52) + (4/52)$$

$$\approx 0.154$$

For disjoint events A and B,
 $P(A \text{ or } B) = P(A) + P(B)$

Image source: <http://svg-cards.sourceforge.net/>

union of non-disjoint events

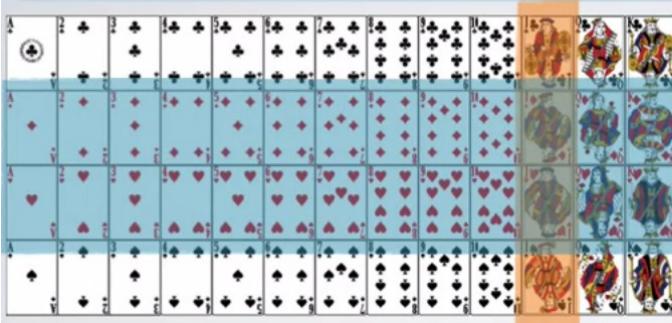
What is the probability of drawing a Jack or a red card from a well shuffled full deck of cards?

$$P(J \text{ or red})$$

$$= P(J) + P(\text{red}) - P(J \text{ and red})$$

$$= (4/52) + (26/52) - (2/52)$$

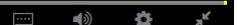
$$\approx 0.538$$



For non-disjoint events A and B,
 $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

Image source: <http://svg-cards.sourceforge.net/>

4:26 / 9:28



What is the probability that a randomly sampled student thinks marijuana should be legalized or they agree with their parents' political views?

Share Parents' Politics

Legalize MJ	No	Yes	Total
No	11	40	51
Yes			
Total			

probability distributions

a **probability distribution** lists all possible outcomes in the sample space, and the probabilities with which they occur.

(40+36-78) / 165
 (114+118-78) / 1

Correct

Use the general

78 / 165

78 / 188

11 / 47

one toss	head	tail
probability	0.5	0.5

two tosses	head - head	tail - tail	head - tail	tail - head
probability	0.25	0.25	0.25	0.25

rules

1. the events listed must be disjoint
2. each probability must be between 0 and 1
3. the probabilities must total 1

complementary events

complementary events are two mutually exclusive events whose probabilities add up to 1.

complementary	
one toss	head
probability	0.5
	tail

complementary	
two tosses	head - head
probability	0.25
	tail - tail
	head - tail
	tail - head

7:53 / 9:28

disjoint vs. complementary

Do the sum of probabilities of two disjoint outcomes always add up to 1?

Not necessarily, there may be more than 2 outcomes in the sample space.

Do the sum of probabilities of two complementary outcomes always add up to 1?

comp

Yes, that's the definition of complementary.

Complementary are always disjoint but disjoint aren't always complementary.

2.1.5 Complement of an event

S
Sample space

A^c
Complement
of outcome A

Rolling a die produces a value in the set $\{1, 2, 3, 4, 5, 6\}$. This set of all possible outcomes is called the **sample space (S)** for rolling a die. We often use the sample space to examine the scenario where an event does not occur.

Let $D = \{2, 3\}$ represent the event that the outcome of a die roll is 2 or 3. Then the **complement** of D represents all outcomes in our sample space that are not in D , which is denoted by $D^c = \{1, 4, 5, 6\}$. That is, D^c is the set of all possible outcomes not already included in D . Figure 2.9 shows the relationship between D , D^c , and the sample space S .

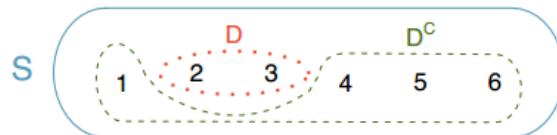


Figure 2.9: Event $D = \{2, 3\}$ and its complement, $D^c = \{1, 4, 5, 6\}$. S represents the sample space, which is the set of all possible outcomes.

independence

two processes are **independent** if knowing the outcome of one provides no useful information about the outcome of the other.

1st toss



2nd toss



$$P(H) = 0.5 \quad P(T) = 0.5$$

outcomes of two tosses of a coin are independent

1st draw



2nd draw



$$P(A) = 3/51 \quad P(J) = 4/51$$

outcomes of two draws from a deck of cards (without replacement) are dependent

deck of cards, without replacement are dependent.

Image sources:

https://commons.wikimedia.org/wiki/File:1913_Eliasberg_Liberty_Head_Nickel.png
<http://openclipart.org/cgi-bin/navigate/recreation/games/cards/white>

If the probability of an event A occurring, given that event B occurred is the same as the probability of event A occurring in the first place, then events A and B are said to be independent. (*as in the case of coin tossing*)

This rule basically says that knowing B tells us nothing about A

Checking for independence:

$P(A | B) = P(A)$, then A and B are independent.

given

In 2013, SurveyUSA interviewed a random sample of 500 NC residents asking them whether they think widespread gun ownership protects law abiding citizens from crime, or makes society more dangerous.

- 58% of all respondents said it protects citizens.
- 67% of White respondents,
- 28% of Black respondents,
- and 64% of Hispanic respondents shared this view.

Opinion on gun ownership and race ethnicity are most likely _____?

- (a) complementary
- (b) mutually exclusive
- (c) independent
- (d) dependent
- (e) disjoint

$$P(\text{protects citizens}) = 0.58$$

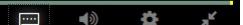
$$P(\text{protects citizens} \mid \text{White}) = 0.67$$

$$P(\text{protects citizens} \mid \text{Black}) = 0.28$$

$$P(\text{protects citizens} \mid \text{Hispanic}) = 0.64$$

Link to poll: <http://www.surveymonkey.com/client/PollReport.aspx?g=a5f160ef-bba9-484b-8579-1101ea26421b>

3:46 / 9:58



$P(\text{citizen who thinks that gun protects given that he is white})$ is 0.67.

Earlier we saw that $P(\text{protects citizens} \mid \text{White}) = 0.67$ and $P(\text{protects citizens} \mid \text{Hispanic}) = 0.64$. Under which sample size would you be more convinced of a real difference between the proportions of Whites and Hispanics who think the widespread gun ownership protects citizens?

n = 500

n = 50,000

Correct

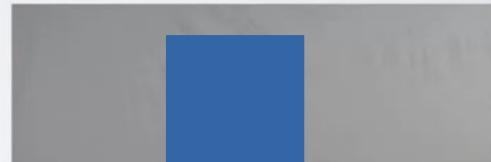
Larger the sample size, the more convincing the evidence from that sample.

determining dependence based on sample data

observed difference
between conditional
probabilities → dependence → hypothesis test

if difference is large, there
is stronger evidence that
the difference is real

if sample size is large, even a small
difference can provide strong
evidence of a real difference



Product rule for independent events:

If A and B are independent, $P(A \text{ and } B) = P(A) \times P(B)$

You toss a coin twice, what is the probability of getting two tails in a row?

$$\begin{aligned}P(\text{two tails in a row}) &= \\&= P(T \text{ on the 1st toss}) \times P(T \text{ on the 2nd toss}) \\&= (1/2) \times (1/2) \\&= 1/4\end{aligned}$$

Note: If A_1, A_2, \dots, A_k are independent, $P(A_1 \text{ and } A_2 \text{ and } \dots \text{ and } A_k) = P(A_1) \times P(A_2) \times \dots \times P(A_k)$

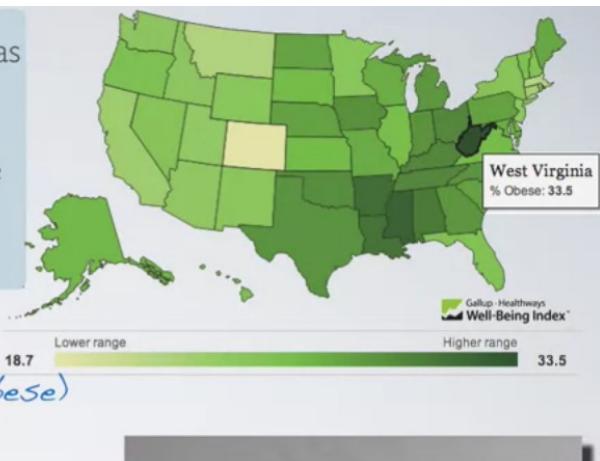
7:07 / 9:58

if A_1, A_2, \dots, A_k all the way through A_k are independent, then probability of all of these events happening at once is simply going to be the product of the individual probabilities of the events

A 2012 Gallup poll suggests that West Virginia has the highest obesity rate among US states, with 33.5% of West Virginians being obese. Assuming that the obesity rate stayed constant, what is the probability that two randomly selected West Virginians are both obese? *independent*

$$P(\text{obese}) = 0.335$$

$$\begin{aligned} P(\text{both obese}) &= P(\text{1st obese}) \times P(\text{2nd obese}) \\ &= 0.335 \times 0.335 \end{aligned}$$

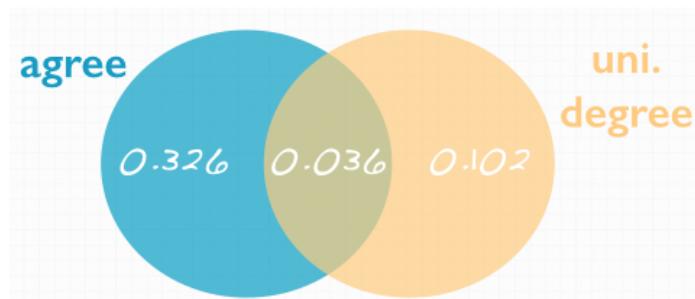


Given the probabilities below and the accompanying Venn diagram, what percent of the world population have a university degree or higher and disagree with the statement about men having more right to a job than women?

$$P(\text{agree}) = 0.362$$

$$P(\text{uni. degree}) = 0.138$$

$$P(\text{agree} \& \text{uni. degree}) = 0.036$$



(This is a long question, please scroll on the right to see the options.)

1 - 0.138 = 0.862 → 86.2%

0.102 → 10.2%

Correct

$$P(\text{uni degree} \& \text{disagree})$$

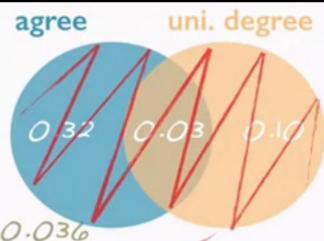
$$= P(\text{uni degree}) - P(\text{uni degree} \& \text{agree})$$

$$= 0.138 - 0.036 = 0.102$$

below are the 2 ways for calculating the probability by using addition rule and venn diagram.

(3) What is the probability that a randomly drawn person has a university degree or higher or agrees with the statement about men having more right to a job than women?

$$P(\text{agree}) = 0.362$$
$$P(\text{uni. degree}) = 0.138$$
$$P(\text{agree} \& \text{uni. degree}) = 0.036$$



$$P(\text{agree or uni. degree})$$

General addition rule:
$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

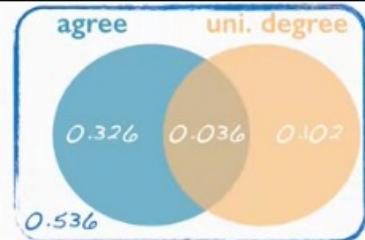
$$\begin{aligned} &= P(\text{agree}) + P(\text{uni. degree}) - P(\text{agree} \& \text{uni. degree}) \\ &= 0.362 + 0.138 - 0.036 \\ &= 0.464 \end{aligned}$$

$$0.326 + 0.036 + 0.102 = 0.464$$

probability in the intersection of the two circles and arrive at the same answer.

(4) What percent of the world population do not have a university degree and disagree with the statement about men having more right to a job than women?

$$\begin{aligned} P(\text{agree}) &= 0.362 \\ P(\text{uni. degree}) &= 0.138 \\ P(\text{agree} \& \text{uni. degree}) &= 0.036 \\ P(\text{agree or uni. degree}) &= 0.464 \end{aligned}$$



$$P(\text{neither agree nor uni. degree})$$

$$\begin{aligned} &= 1 - P(\text{agree or uni. degree}) \\ &= 1 - 0.464 = 0.536 \end{aligned}$$

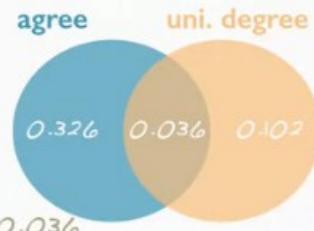
How to check independence ?

(5) Does it appear that the event that someone agrees with the statement is independent of the event that they have a university degree or higher?

$$P(\text{agree}) = 0.362$$

$$P(\text{uni. degree}) = 0.138$$

$$P(\text{agree} \& \text{uni. degree}) = 0.036$$



Product rule for independent events:

If A and B are independent, $P(A \text{ and } B) = P(A) \times P(B)$

$$P(\text{agree} \& \text{uni. degree}) ?= ? P(\text{agree}) \times P(\text{uni. degree})$$

$$0.036 ?= ? 0.362 \times 0.138$$

$$0.036 \neq 0.05$$

They aren't independent as the occurrence of event doesn't give us any hint about the occurrence of second event.

(6) What is the probability that at least 1 in 5 randomly selected people agree with the statement about men having more right to a job than women?

$$P(\text{agree}) = 0.362$$

$$S = \{0, 1, 2, 3, 4, 5\} \longrightarrow S = \{0, \text{at least } 1\}$$

$$P(\text{at least } 1 \text{ agree}) = 1 - P(\text{none agree})$$

$$= 1 - P(D D D D D) \longrightarrow P(\text{disagree})$$

$$= 1 - 0.638^5 \quad = 1 - P(\text{agree})$$

$$= 1 - 0.106 = 0.894 \quad = 1 - 0.362$$

$$= 0.638$$

Which of the following is true? Assume a fair coin.

- Outcomes of one coin toss are disjoint and independent; outcomes of two coin tosses are dependent.
- Outcomes of one coin toss are disjoint and dependent; outcomes of two coin tosses are independent.

Correct

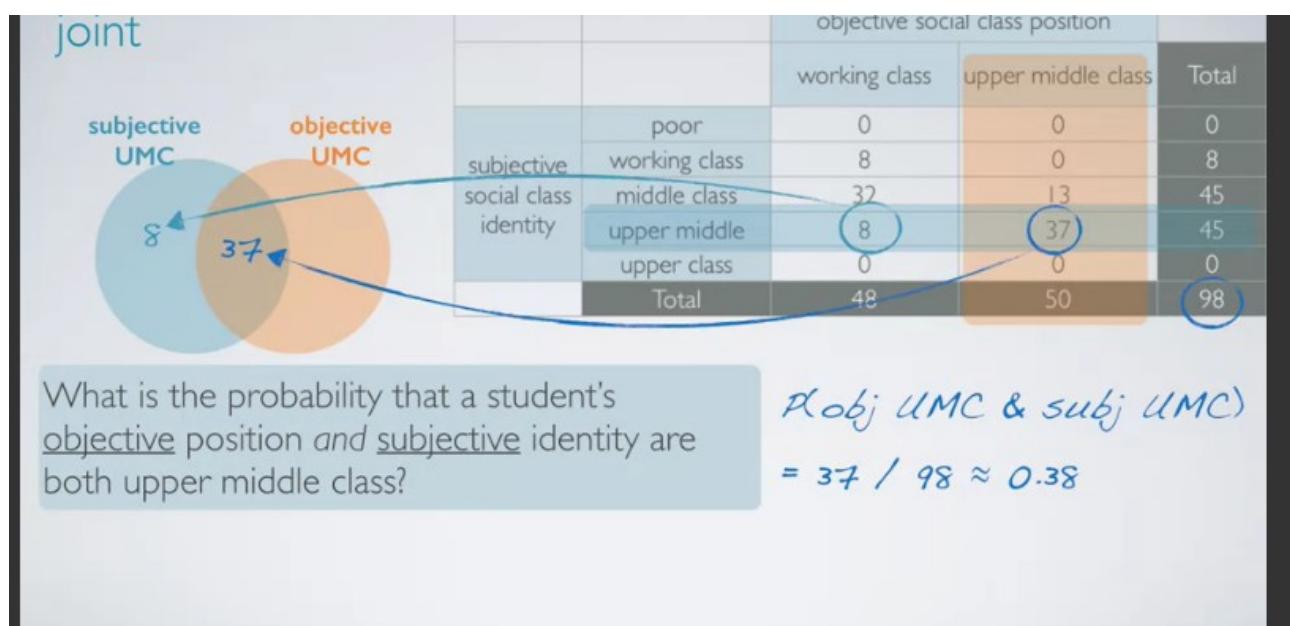
For a single coin toss, if we know that it's heads, then we know that it's not tails, hence dependent. For two consecutive coin tosses, knowing the result of the first toss does not tell us anything about the result of the second toss, hence independent.

- Outcomes of one coin toss are disjoint and dependent; outcomes of two coin tosses are also dependent.

We can generalize this to say that disjoint events with non-zero probability are always dependent on each other. Because if we know that one happened, we know that the other one cannot happen.

Conditional probability

joint probability:



The term joint probability comes from the fact that we're considering the students who are at the **intersection** of the two events of interest.

CONDITIONAL :

read it from the book too.

conditional

		objective social class position		Total
		working class	upper middle class	
subjective social class identity	poor	0	0	0
	working class	8	0	8
	middle class	32	13	45
	upper middle	8	37	45
	upper class	0	0	0
	Total	48	50	98

What is the probability that a student who is objectively in the working class associates with upper middle class?

$P(\text{subj UMC} | \text{obj WC}) = 8 / 48 \approx 0.17$

the conditional probabilities.

When the counts are not given we use bayes theorem to calculate the conditional probabilities.

Bayes' theorem:

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

		objective social class position		Total
		working class	upper middle class	
subjective social class identity	poor	0	0	0
	working class	8	0	8
	middle class	32	13	45
	upper middle	8	37	45
	upper class	0	0	0
	Total	48	50	98

$P(\text{subj UMC} | \text{obj WC}) = \frac{P(\text{subj UMC} \& \text{ obj WC})}{P(\text{obj WC})} = \frac{8 / 98}{48 / 98}$

48 out of 98 students are working class based on their objective categorization.

The numerator is the joint probability while the denominator is the marginal .

The American Community Survey is an ongoing survey that provides data every year to give communities the current information they need to plan investments and services.

The 2010 American Community Survey estimates that 14.6% of Americans live below the poverty line, 20.7% speak a language other than English at home, and 4.2% fall into both categories.

Based on this information, what percent of Americans live below the poverty line given that they speak a language other than English at home?

$$P(\text{below PL} \mid \text{speak non-Eng}) = ?$$

$$= \frac{P(\text{below PL} \& \text{speak non-Eng})}{P(\text{speak non-Eng})} = \frac{0.042}{0.207} \approx 0.2$$

Bayes' theorem:

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

at, for the part of the public that speaks
a language other than English at home.

Community Survey 1-Year Estimates,
People by Language Spoken at Home.

8:46 / 12:40



roughly 20% of Americans who speak a language other than English at home also live below the poverty line.

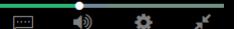
independence and conditional probabilities

Generically, if $P(A|B) = P(A)$ then the events A and B are said to be independent.

- ▶ **Conceptually:** Giving B doesn't tell us anything about A.
- ▶ **Mathematically:** If events A and B are independent, $P(A \text{ and } B) = P(A) \times P(B)$. Then,

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

11:13 / 12:40



IMPORTANT EXAMPLE:

example

		major		Total
		social science	non-social science	
gender	female	30	20	50
	male	30	20	50
	Total	60	40	100

$$P(SS) = 60 / 100 = 0.6$$

$$P(SS | F) = 30 / 50 = 0.6$$

$$P(SS | M) = 30 / 50 = 0.6$$

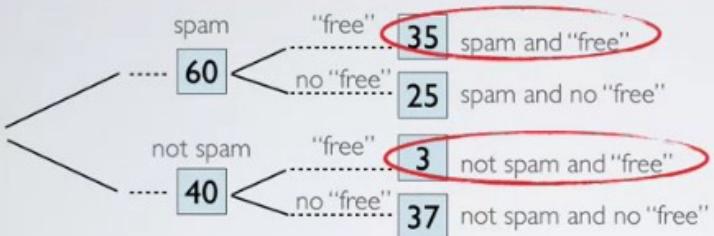
So what we're seeing here is that all of these probabilities are exactly the same. So this goes back to probability of a given b. If that equals probability of a, then we know that the events are independent. In this case, probability of social science equals probability of social science given female or social science given male. So we would determine that the two variables, gender and major are independent of each other, given this hypothetical distribution.

PROBABILITY TREES

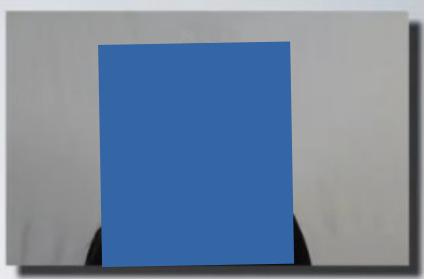
they are useful when we have to calculate the reverse of what we are given.

$$P(A | B) \rightarrow P(B | A)$$

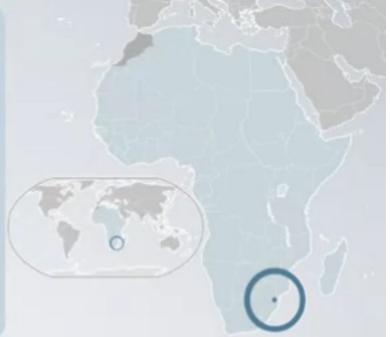
You have 100 emails in your inbox: 60 are spam, 40 are not. Of the 60 spam emails, 35 contain the word "free". Of the rest, 3 contain the word "free". If an email contains the word "free", what is the probability that it is spam?



$$P(\text{spam} | \text{"free"}) = \frac{35}{35 + 3} = 0.92$$



As of 2009, Swaziland had the highest HIV prevalence in the world. 25.9% of this country's population is infected with HIV. The ELISA test is one of the first and most accurate tests for HIV. For those who carry HIV, the ELISA test is 99.7% accurate. For those who do not carry HIV, the test is 92.6% accurate. If an individual from Swaziland has tested positive, what is the probability that he carries HIV?

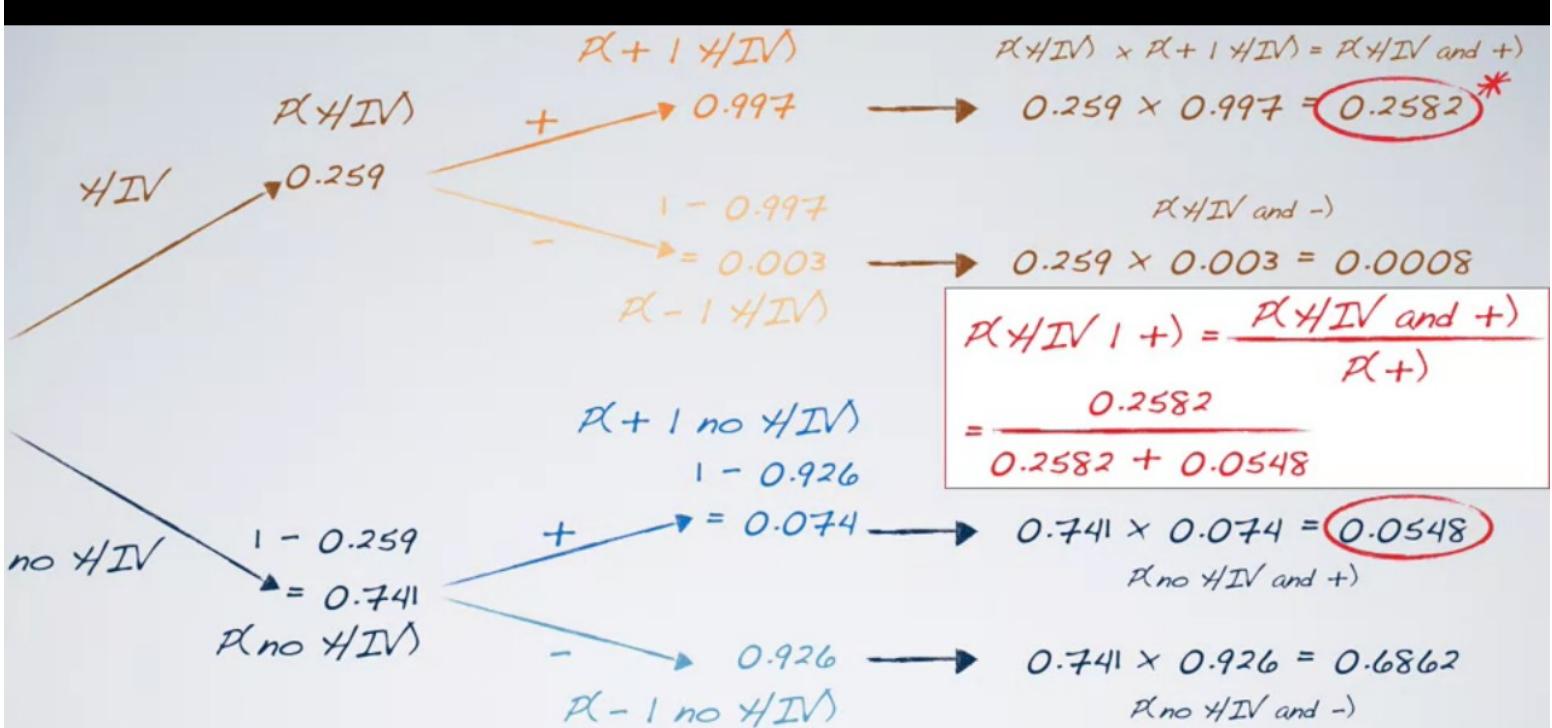


$$P(HIV) = 0.259$$

$$P(+ | HIV) = 0.997 \quad P(- | \text{no HIV}) = 0.926$$

~~tree diagram!~~

$$P(HIV | +) = ?$$



Since we're saying or and these are disjoined probabilities to get

9:56 / 10:31

Note that the marginal probability [$P(+)$] is the sum of [$P(+ \text{ and hiv}) + P(+ \text{ and no hiv})$].

BAYESIAN INFERENCE

posterior

- The probability we just calculated is also called the **posterior probability**.
 $P(H_1: \text{good die on the Right} | \text{you rolled } \geq 4 \text{ with the die on the Right})$
- Posterior probability is generally defined as $P(\text{hypothesis} | \text{data})$.
- It tells us the probability of a hypothesis we set forth, given the data we just observed.
- It depends on both the prior probability we set and the observed data.
- This is different than what we calculated at the end of the randomization test on gender discrimination – the probability of observed or more extreme data given the null hypothesis being true, i.e. $P(\text{data} | \text{hypothesis})$, also called a **p-value**.

the hypothesis, which we had called a p-value.

12:13 / 14:27

updating the prior

- ▶ In the Bayesian approach, we evaluate claims iteratively as we collect more data.
- ▶ In the next iteration (roll) we get to take advantage of what we learned from the data.
- ▶ In other words, we **update** our prior with our posterior probability from the previous iteration.

updated:

P(H ₁ : good die on the Right)	P(H ₂ : good die on the Left)
0.6	0.4

to be the probability of the competing hypothesis.

recap

- ▶ Take advantage of prior information, like a previously published study or a physical model.
- ▶ Naturally integrate data as you collect it, and update your priors.
- ▶ Avoid the counter-intuitive definition of a p-value:
 $P(\text{observed or more extreme outcome} \mid H_0 \text{ is true})$
- ▶ Instead base decisions on the posterior probability:
 $P(\text{hypothesis is true} \mid \text{observed data})$
- ▶ A good prior helps, a bad prior hurts, but the prior matters less the more data you have.
- ▶ More advanced Bayesian techniques offer flexibility not present in Frequentist models.

continue your studies with statistics and work with more advanced Bayesian models.

if you didn't have a great prior to begin with, as you collect my, more data, you're going to be able to converge to the right probabilities.

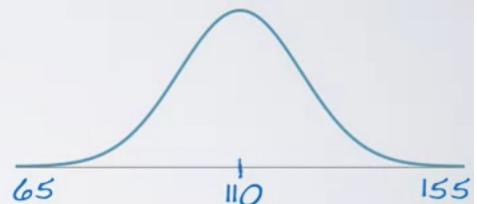
NORMAL DISTRIBUTION

we calculate the STD with min,max and mean.

Practice

A doctor collects a large set of heart rate measurements that approximately follow a normal distribution. He only reports 3 statistics, the mean = 110 beats per minute, the minimum = 65 beats per minute, and the maximum = 155 beats per minute. Which of the following is most likely to be the standard deviation of the distribution?

- (a) 5 $\rightarrow 110 \pm (3 \times 5) = (95, 125)$
- (b) 15** $\rightarrow 110 \pm (3 \times 15) = (65, 155)$
- (c) 35 $\rightarrow 110 \pm (3 \times 35) = (5, 215)$
- (d) 90 $\rightarrow 110 \pm (3 \times 90) = (-160, 380)$



which would give us that within three standard deviations would

▶ 5:25 / 17:19

[play] [volume] [settings] [exit]

We can calculate Z scores for any distribution.

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

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

After all, every distribution will have a mean and a standard deviation,

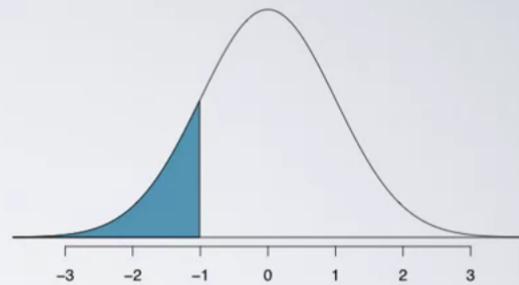
▶ 8:47 / 17:19

[play] [volume] [settings] [exit]

Percentiles can be calculated for dist of any shape but if it is not symmetric then we have to use calculus.

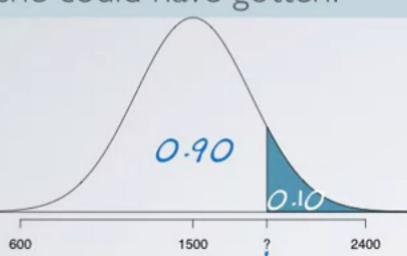
percentiles

- when the distribution is normal, Z scores can be used to calculate percentiles
- percentile** is the percentage of observations that fall below a given data point
- graphically, percentile is the area below the probability distribution curve to the left of that observation.



And for the purposes of this course,
we're not going to be using calculus, so

A friend of yours tells you that she scored in the top 10% on the SAT. What is the lowest possible score she could have gotten?



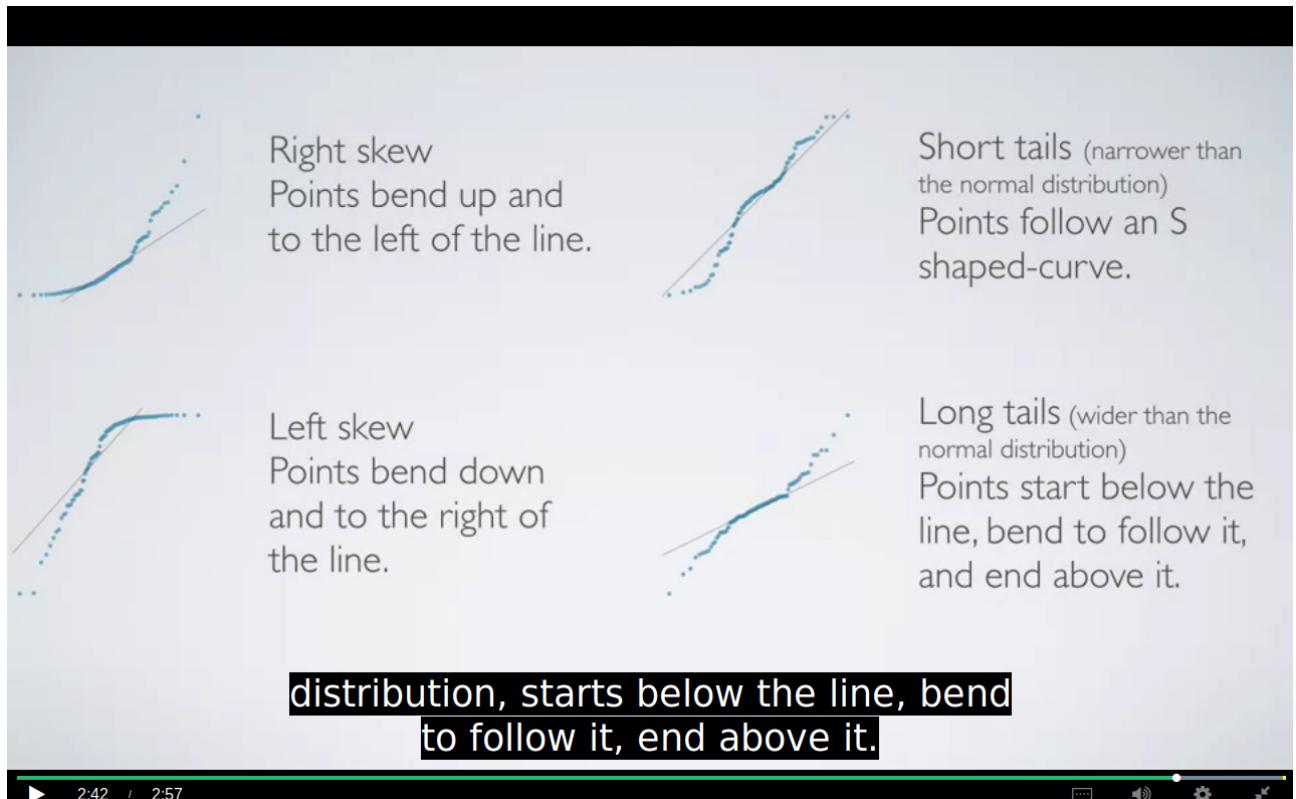
$$Z = 1.28 = \frac{X - 1500}{300}$$
$$X = (1.28 \times 300) + 1500 = 1884$$

Z	Second decimal place of Z									
	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.5000	0.5040	0.5080	0.5120	0.5160	0.5199	0.5239	0.5279	0.5319	0.5359
0.1	0.5398	0.5438	0.5478	0.5517	0.5557	0.5596	0.5636	0.5675	0.5714	0.5753
0.2	0.5793	0.5832	0.5871	0.5910	0.5948	0.5987	0.6026	0.6064	0.6103	0.6141
0.3	0.6179	0.6217	0.6255	0.6293	0.6331	0.6368	0.6406	0.6443	0.6480	0.6517
0.4	0.6554	0.6591	0.6628	0.6664	0.6700	0.6736	0.6772	0.6808	0.6844	0.6879
0.5	0.6915	0.6950	0.6985	0.7019	0.7054	0.7088	0.7123	0.7157	0.7190	0.7224
0.6	0.7257	0.7291	0.7324	0.7357	0.7389	0.7422	0.7454	0.7486	0.7517	0.7549
0.7	0.7580	0.7611	0.7642	0.7673	0.7704	0.7734	0.7764	0.7794	0.7823	0.7852
0.8	0.7881	0.7910	0.7939	0.7967	0.7995	0.8023	0.8051	0.8078	0.8106	0.8133
0.9	0.8159	0.8186	0.8212	0.8238	0.8264	0.8289	0.8315	0.8340	0.8365	0.8389
1.0	0.8413	0.8438	0.8461	0.8485	0.8508	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8643	0.8665	0.8686	0.8708	0.8729	0.8749	0.8770	0.8790	0.8810	0.8830
1.2	0.8849	0.8869	0.8888	0.8907	0.8925	0.8944	0.8962	0.8980	0.8997	0.9015
1.3	0.9032	0.9049	0.9066	0.9082	0.9099	0.9115	0.9131	0.9147	0.9162	0.9177
1.4	0.9192	0.9207	0.9222	0.9236	0.9251	0.9265	0.9279	0.9292	0.9306	0.9319

R
> qnorm(0.90, 1500, 300)
[1] 1884.465

Practice

Shape of dist using Normal probability plot.



2:42 / 2:57



Distribution of number of hours of sleep college students get has a mean of 7 hours. 62% of students sleep between 6 and 8 hours, 92% of students sleep between 5 and 9 hours, and 95% sleep between 4 and 10 hours. Which of the following is true?

- The distribution is more variable than a normal distribution with mean 7 and standard deviation 1

Correct

In a normal distribution with mean 7 and standard deviation 1, we would expect 68% of the data to be between 6-8, 95% of the data between 5-9, and 99.7% of the data between 4-10 hours. In this data set a lower percentage than expected fall within these ranges, and hence the data are more spread out, i.e. more variable, than $N(\text{mean} = 7, \text{SD} = 1)$.

- The distribution is less variable than a normal distribution with mean 7 and standard deviation 1
 The distribution is nearly normal

Suppose weights of checked baggage of airline passengers follow a nearly normal distribution with mean 45 pounds and standard deviation 3.2 pounds. Most airlines charge a fee for baggage that weigh in excess of 50 pounds.

$$\text{baggage} \sim N(\text{mean} = 45, \text{SD} = 3.2)$$

What percent of airline passengers are expected to incur this fee?

SOLUTION:

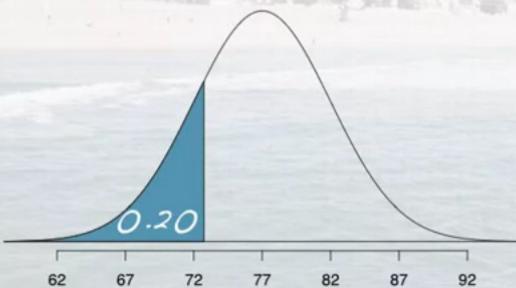
`> 1-pnorm(50,45,3.2)`

`[1] 0.05908512`

The average daily high temperature in June in LA is 77 °F with a standard deviation of 5 °F. Suppose that the temperatures in June closely follow a normal distribution.

How cold are the coldest 20% of the days during June in LA?

temperatures $\sim N(\text{mean} = 77, \text{SD} = 5)$

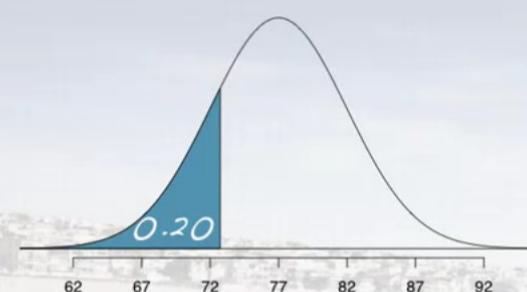


```
R
> qnorm(0.20, mean = 77, sd = 5)
[1] 72.79
```

than 72, 72.79 degrees Fahrenheit.

https://gallery.shinyapps.io/dist_calc/

another approach for solving using Z table



Second decimal place of Z							Z
0.06	0.05	0.04	0.03	0.02	0.01	0.00	
0.0003	0.0003	0.0003	0.0003	0.0003	0.0003	0.0003	-3.4
0.0004	0.0004	0.0004	0.0004	0.0005	0.0005	0.0005	-3.3
0.0006	0.0006	0.0006	0.0006	0.0006	0.0007	0.0007	-3.2
0.0008	0.0008	0.0008	0.0009	0.0009	0.0009	0.0010	-3.1
0.0011	0.0011	0.0012	0.0012	0.0013	0.0013	0.0013	-3.0
0.0015	0.0016	0.0016	0.0017	0.0018	0.0018	0.0019	-2.9
0.0021	0.0022	0.0023	0.0023	0.0024	0.0025	0.0026	-2.8
0.1685	0.1711	0.1736	0.1762	0.1788	0.1814	0.1841	-0.9
0.1949	0.1977	0.2005	0.2033	0.2061	0.2090	0.2119	-0.8
0.2236	0.2266	0.2296	0.2327	0.2358	0.2389	0.2420	-0.7

$$Z = -0.84 = \frac{x - 77}{5}$$

$$x = (-0.84 \times 5) + 77 = 72.8$$

once again get to roughly the same answer,
72.8 degrees Fahrenheit.

https://gallery.shinyapps.io/dist_calc/

Test yourself: True/False: In a right skewed distribution the Z score of the median is positive.

the answer is false bcz in Right skewed distribution the mean is > median.

BINOMIAL DISTRIBUTION

Bernoulli random variables

- ▶ each person in Milgram's experiment can be thought of as a **trial**
- ▶ a person is labeled a **success** if she refuses to administer a severe shock, and **failure** if she administers such shock
- ▶ since only 35% of people refused to administer a shock, **probability of success** is $p = 0.35$.

of the analysis we're going to focus

▶ 1:55 / 17:13



Suppose we randomly select four individuals to participate in this experiment. What is the probability that exactly 1 of them will refuse to administer the shock?

- ▶ Four individuals:
(A) Anthony
(B) Brittany
(C) Clara
(D) Dorian
- ▶ Multiple scenarios where "exactly 1 refuses"

Scenario 1: $\frac{0.35}{(\text{A) refuse}} \times \frac{0.65}{(\text{B) shock}} \times \frac{0.65}{(\text{C) shock}} \times \frac{0.65}{(\text{D) shock}} = 0.0961$

OR

Scenario 2: $\frac{0.65}{(\text{A) shock}} \times \frac{0.35}{(\text{B) refuse}} \times \frac{0.65}{(\text{C) shock}} \times \frac{0.65}{(\text{D) shock}} = 0.0961$

OR

Scenario 3: $\frac{0.65}{(\text{A) shock}} \times \frac{0.65}{(\text{B) shock}} \times \frac{0.35}{(\text{C) refuse}} \times \frac{0.65}{(\text{D) shock}} = 0.0961$

OR

Scenario 4: $\frac{0.65}{(\text{A) shock}} \times \frac{0.65}{(\text{B) shock}} \times \frac{0.65}{(\text{C) shock}} \times \frac{0.35}{(\text{D) refuse}} = 0.0961$

4 x 0.0961 = 0.3844

of one scenario with the number of scenarios to arrive at the same answer.

▶ 5:45 / 17:13



binomial distribution

the **binomial distribution** describes the probability of having exactly k successes in n independent Bernoulli trials with probability of success p

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

of scenarios x $P(\text{single scenario})$

\downarrow

"n choose k"

\downarrow

probability of success
to the power of
number of successes

\downarrow

$p^k(1-p)^{(n-k)}$

\downarrow

probability of failure
to the power of
number of failures

n minus k factorial.
Let's give a couple examples here.

How many scenarios yield
1 success in 4 trials?

$$n = 4 \quad k = 1$$

$$\binom{4}{1} = \frac{4!}{1! \times (4-1)!}$$
$$= \frac{4 \times \cancel{3} \times \cancel{2} \times 1}{1 \times \cancel{3} \times \cancel{2} \times 1} = 4$$

$$SSFFFFFF$$

$$SFSFFFFF$$

$$SFFSFFFF$$

How many scenarios yield
2 successes in 9 trials?

$$n = 9 \quad k = 2$$

$$\binom{9}{2} = \frac{9!}{2! \times 7!}$$

$$= \frac{9 \times 8 \times \cancel{7}!}{2 \times 1 \times \cancel{7}!} = 36$$

R

> choose(9, 2)

[1] 36

Binomial distribution:

If p represents probability of success, $(1-p)$ represents probability of failure, n represents number of independent trials, and k represents number of successes

$$P(k \text{ successes in } n \text{ trials}) = \binom{n}{k} p^k (1-p)^{(n-k)}$$

$$\text{where } \binom{n}{k} = \frac{n!}{k!(n-k)!}$$

binomial conditions:

1. the trials must be independent
2. the number of trials, n , must be fixed
3. each trial outcome must be classified as a success or a failure
4. the probability of success, p , must be the same for each trial

Same data: 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). We are interested in finding the probability that among a random sample of 10 employees, what is the probability that 8 of them are engaged at work.

Without doing any calculations, would you expect this probability to be pretty low or pretty high?

- pretty low

Correct

Among a sample of 10, we would only expect $10 \times 0.13 = 1.3$ people to be engaged at work, so getting a random sample where 8 out of 10 are engaged would be pretty unlikely.

- pretty high

practice

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?

$$n = 10$$

$$p = 0.13$$

$$1 - p = 0.87$$

$$k = 8$$

$$\begin{aligned}P(K = 8) &= \binom{10}{8} 0.13^8 \times 0.87^2 \\&= \frac{10!}{8! \times 2!} \times 0.13^8 \times 0.87^2 \\&= \frac{10 \times 9 \times 8!}{8! \times 2 \times 1} \times 0.13^8 \times 0.87^2 \\&= 45 \times 0.13^8 \times 0.87^2 \\&= 0.00000278\end{aligned}$$

employees to be engaged than eight if the probability of success is only 13%.

[worldwide-employees-engaged-work.aspx](#)

▶ 13:38 / 17:13



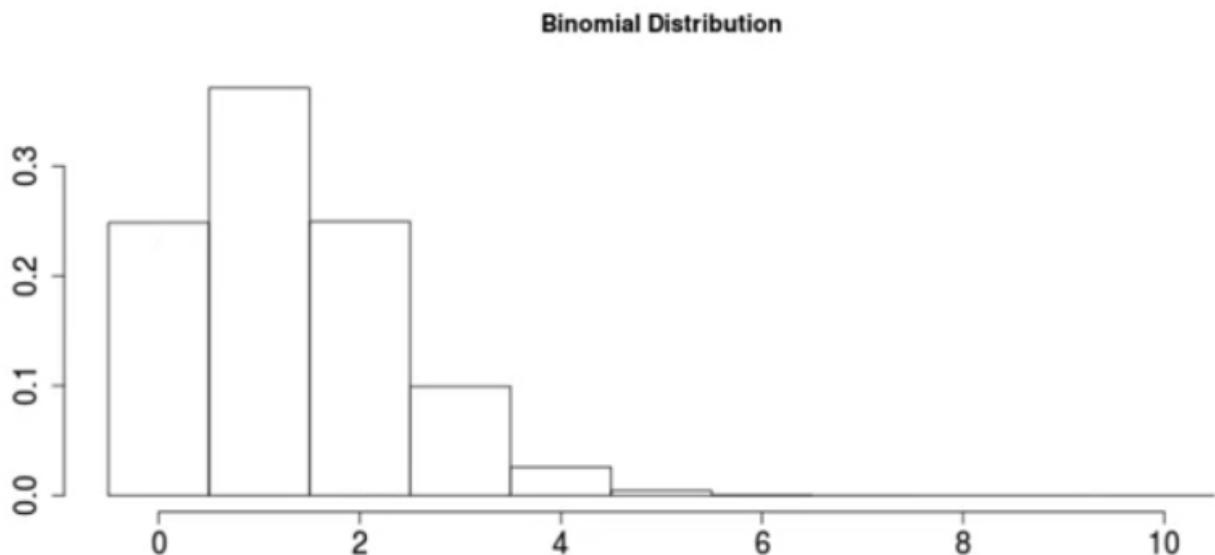
R command:

```
> choose(10,8) * (0.13^8) * (0.87^2)
```

[1] 2.77842e-06

> dbinom(8,size=10,p=0.13)

plot for above question using applet.



$$P(X = 8) = 2.78 \times 10^{-6}$$

The probability for 8 is very low therefore there is no bar above it. The height of the bar shows probability.

Among a random sample of 100 employees, how many would you expect to be engaged at work? Remember: $p = 0.13$.

$$\mu = 100 \times 0.13 = 13$$

Expected value (mean) of binomial distribution: $\mu = np$

Standard deviation of binomial distribution: $\sigma = \sqrt{np(1 - p)}$

$$\sigma = \sqrt{100 \times 0.13 \times 0.87} = 3.36$$

This means that 13 out of 100 employees are expected to be at engaged at work, give or take approximately 3.36. Note that the mean on the standard deviation of a binomial, might not always be whole numbers, and that's alright. These values represent what we would expect to see on average.

A 2012 Gallup survey suggests that 26.2% of Americans are obese. Which of the following is false?

- Among a random sample of 1,000 Americans, we would expect 262 to be obese.
- Random samples of 1,000 Americans where there are at most 230 are obese people would be considered unusual.
- The standard deviation of number of obese Americans in random samples of 1,000 is roughly 14.
- Random samples of 1,000 Americans where at least 300 are obese would not be considered unusual.

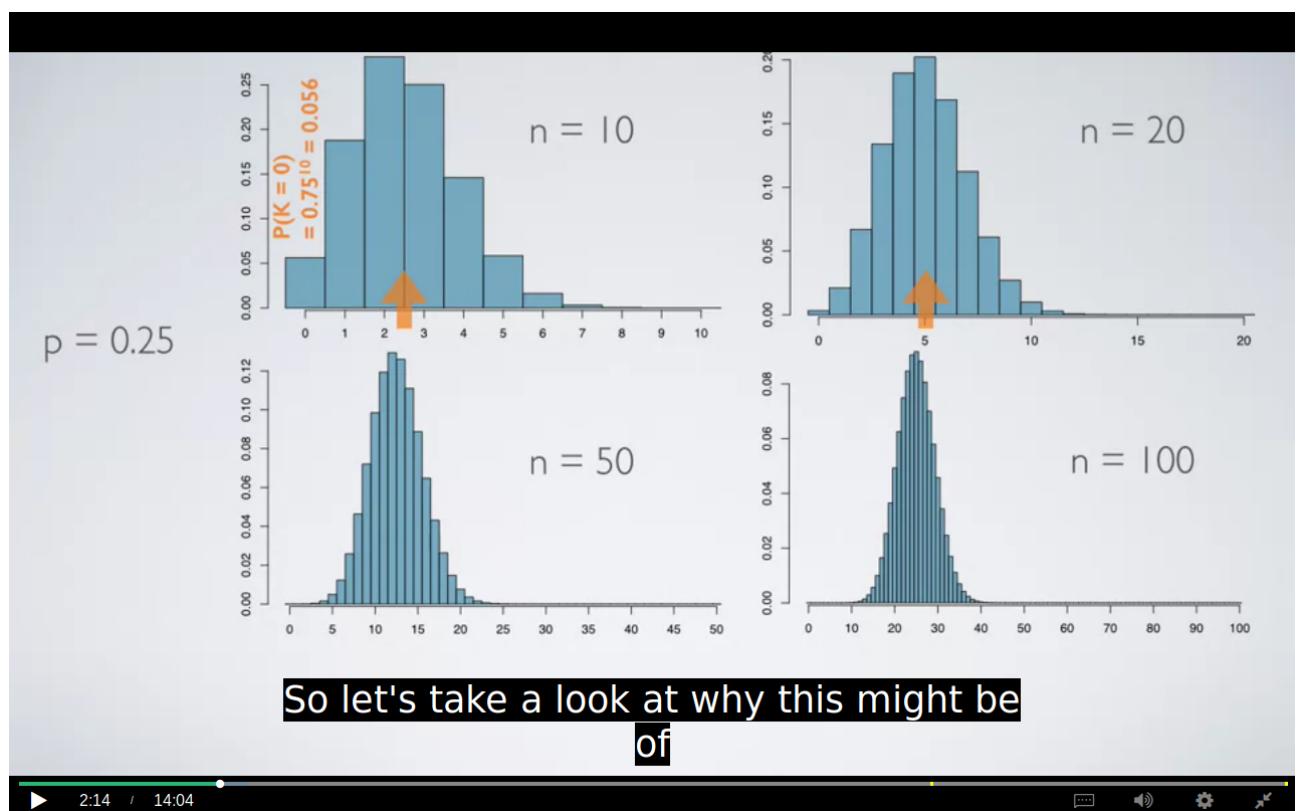
Correct

$$\mu = 1000 \times 0.262 = 262$$

$$\sigma = \sqrt{1000 \times 0.262 \times 0.738} \approx 14$$

Range of "usual" observations: $262 \pm 2 \times 14 = (234, 290)$, anything beyond these would be considered unusual.

NORMAL APPROXIMATION TO BINOMIAL:



THE CENTER OF THE DISTRIBUTION CHANGES AS WE ARE INCREASING THE NO OF TRIALS.

Calculating the area under the curve is much more simpler then calculating the individual probabilities therefore we use normal approximation

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,

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

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

What is the minimum required n for a binomial distribution with $p = 0.25$ to closely follow a normal distribution?

$$n \times 0.25 \geq 10$$

$$n \geq 10 / 0.25$$

$$\textcircled{n \geq 40}$$

$$n \times 0.75 \geq 10$$

$$n \geq 10 / 0.75$$

$$n \geq 13.33$$

with p equals 0.25, to closely follow a normal distribution.

What is the probability that in a random sample of 1000 people exactly 600 plan to get health insurance through a government health insurance exchange?

- (a) 0.243, same as $P(K = 6)$
- (b) less than 0.243**
- (c) more than 0.243

$$P = 0.56$$

$$n_1 = 10 \quad \mu_1 = 10 \times 0.56 = 5.6$$

$$\Delta = 6 - 5.6 = 0.4$$

$$n_2 = 1000 \quad \mu_2 = 1000 \times 0.56 = 560 > \text{dbinom}(600, 1000, 0.56)$$

$$\Delta = 600 - 560 = 40$$

R

[1] 0.00098

4:58 / 9:59

the diff between the expected and desired is much lower in the case of $n=10$ while it is much larger in $n=1000$ therefore acc to law of large number the desired outcome is pretty unusual.

What is the probability that at least 60 out of a random sample of 100 uninsured Americans plan to get health insurance through a government health insurance exchange?

R

```
> sum(dbinom(60:100, size = 100, p = 0.56))
[1] 0.241
```

Alternate Solution for same question using 0.5 correction by normal approximation.

$$Z = \frac{60 - 56}{4.96} \approx 0.81$$

$$P(Z > 0.81) = 0.209$$

The probability was .20 which differs from the one calculated with the binomial formula.

$$Z = \frac{59.5 - 56}{4.96} \approx 0.71$$

$$P(Z > 0.71) = 0.239$$

The probability is almost equal to the binomial after using 0.5 correction.

Important question

3. You are about to take a multi-day tour through a national park which is famous for its wildlife. The tour guide tells you that on any given day there's a 61% chance that a visitor will see at least one "big game" animal, and a 39% chance they'll see no big game animals; when the tour guide says "big game", he refers to either a moose or a bear. The guide assures you that big game sightings on a single day are independent of any other day's sightings. Given the information from the tour guide, which of the following calculations cannot be performed using a binomial distribution?

-  Calculate the probability that you see at least 4 big game animals on the first day of a 5-day trip.

Correct

This question refers to the following learning objective: Determine if a random variable is binomial using the four conditions.

- The trials are independent.
- The number of trials, n , is fixed.
- Each trial outcome can be classified as a success or failure.
- The probability of a success, p , is the same for each trial.

You cannot calculate the probability of seeing at least 4 big game animals on a particular day because not enough information was given in the problem. In other words, the problem gave information on the probability of a success and failure where "success" means seeing any (i.e. at least 1) big game animals on a particular day - we have no information on the exact number of big game animals seen on a particular day.

6. About 30% of human twins are identical, and the rest are fraternal. Identical twins are necessarily the same sex, half are males and the other half are females. One-quarter of fraternal twins are both male, one-quarter both female, and one-half are mixes: one male, one female. You have just become a parent of twins and are told they are both girls. Given this information, what is the probability that they are identical?

33%

50%

46%

Correct

This question refers to the following learning objective: Distinguish between marginal and conditional probabilities. Construct tree diagrams to calculate conditional probabilities and probabilities of intersection of non-independent events using Bayes' theorem: $P(A|B) = P(A \text{ and } B) / P(B)$

A good strategy is to first decide what quantity we're being asked to calculate. Then we'll use Bayes Theorem to write the quantity in terms of other quantities which we hope to be given directly in the problem. So, we're asked to calculate $P(I|G)$ where I is the event of known twins being identical and G the event that both twins are girls. Then use Bayes' Theorem to write $\frac{P(G|I)P(I)}{P(G|I)P(I) + P(G|F)P(F)}$. Let F be the event that known twins are fraternal. It turns out that we're given all of the quantities on the right-hand side of that equation, so the result is $P(I|G) = \frac{.5 \times .3}{.5 \times .3 + .25 \times .7}$. Your answer may vary slightly due to rounding.

CLT:

INCLUDE THE BELOW WITH CLT TOPIC:

4.5 Hen eggs. The distribution of the number of eggs laid by a certain species of hen during their breeding period has a mean of 35 eggs with a standard deviation of 18.2. Suppose a group of researchers randomly samples 45 hens of this species, counts the number of eggs laid during their breeding period, and records the sample mean. They repeat this 1,000 times, and build a

distribution of sample means.

- (a) What is this distribution called?
- (b) Would you expect the shape of this distribution to be symmetric, right skewed, or left skewed? Explain your reasoning.
- (c) Calculate the variability of this distribution and state the appropriate term used to refer to this value.
- (d) Suppose the researchers' budget is reduced and they are only able to collect random samples of 10 hens. The sample mean of the number of eggs is recorded, and we repeat this 1,000 times, and build a new distribution of sample means. How will the variability of this new distribution compare to the variability of the original distribution?

ANS 4.5:

- (a) We are building a distribution of sample statistics, in this case the sample mean. Such a distribution is called a sampling distribution.
- (b) Because we are dealing with the distribution of sample means, we need to check to see if the Central Limit Theorem applies. Our sample size is greater than 30, and we are told that random sampling is employed. With these conditions met, we expect that the distribution of the sample mean will be nearly normal and therefore symmetric.
- (c) Because we are dealing with a sampling distribution, we measure its variability with the standard error. $SE = 18.2 / \sqrt{45} = 2.713$.
- (d) The sample means will be more variable with the smaller sample size.

4.13 Waiting at an ER, Part I. A hospital administrator hoping to improve wait times decides to estimate the average emergency room waiting time at her hospital. She collects a simple random sample of 64 patients and determines the time (in minutes) between when they checked in to the ER until they were first seen by a doctor. A 95% confidence interval based on this sample is (128 minutes, 147 minutes), which is based on the normal model for the mean. Determine whether the following statements are true or false, and explain your reasoning.

- (a) This confidence interval is not valid since we do not know if the population distribution of the ER wait times is nearly Normal.
- (b) We are 95% confident that the average waiting time of these 64 emergency room patients is between 128 and 147 minutes.
- (c) We are 95% confident that the average waiting time of all patients at this hospital's emergency room is between 128 and 147 minutes.
- (d) 95% of random samples have a sample mean between 128 and 147 minutes.
- (e) A 99% confidence interval would be narrower than the 95% confidence interval since we need to be more sure of our estimate.
- (f) The margin of error is 9.5 and the sample mean is 137.5.
- (g) In order to decrease the margin of error of a 95% confidence interval to half of what it is now, we would need to double the sample size.

4.13 (a) False. Provided the data distribution is not very strongly skewed ($n = 64$ in this sample, so we can be slightly lenient with the skew), the sample mean will be nearly normal, allowing for the method normal approximation described.

(b) False. Inference is made on the population parameter, not the point estimate. The point estimate is always in the confidence interval.

(c) True.

(d) False. The confidence interval is not about a sample mean.

(e) False. To be more confident that we capture the parameter, we need a wider interval. Think about needing a bigger net to be more sure of catching a fish in a murky lake.

(f) True. Optional explanation: This is true since the normal model was used to model the sample mean. The margin of error is half the width of the interval, and the sample mean is the midpoint of the interval.

(g) False. In the calculation of the standard error, we divide the standard deviation by the square root of the sample size. To cut the SE (or margin of error) in half, we would need to sample $2^2 = 4$ times the number of people in the initial sample.