

# Chapter 3: Probability

---

STA 101, Summer I 2021, Duke University

Derived from OpenIntro slides, developed by Mine Çetinkaya-Rundel of OpenIntro.  
Edited under the CC BY-SA license.

## **Defining probability**

---

# Random processes

- A *random process* is a situation in which we know what outcomes could happen, but we don't know which particular outcome will happen.
- Examples: coin tosses, die rolls, iTunes shuffle, whether the stock market goes up or down tomorrow, etc.
- It can be helpful to model a process as random even if it is not truly random.

MP3 Players > Stories > iTunes: Just how random is random?

## iTunes: Just how random is random?

By David Braue on 08 March 2007

• Introduction

• Say You, Say What?

• A role for labels?

• The new random

Think that song has appeared in your playlists just a few too many times? David Braue puts the randomness of Apple's song shuffling to the test -- and finds some surprising results.

Quick -- think of a number between one and 20. Now think of another one, and another, and another.

Starting to repeat yourself? No surprise: in practice, many series of random numbers are far less random than you would think.

Computers have the same problem. Although all systems are able to pick random numbers, the method they use is often tied to specific other numbers -- for example, the time -- that means you could get a very similar series of 'random' numbers in different situations.

This tendency manifests itself in many ways. For anyone who uses their iPod heavily, you've probably noticed that your supposedly random 'shuffling' iPod seems to be particularly fond of the Bee Gees, Melissa Etheridge or Pavarotti. Look at a random playlist that iTunes generates for you, and you're likely to notice several songs from one or two artists, while other artists go completely unrepresented.



<http://www.cnet.com.au/>

<itunes-just-how-random-is-random-339274094.htm>

# Probability

- There are several possible interpretations of probability but they (almost) completely agree on the mathematical rules probability must follow.
  - $P(A)$  = Probability of event A
  - $0 \leq P(A) \leq 1$

# Probability

- There are several possible interpretations of probability but they (almost) completely agree on the mathematical rules probability must follow.
  - $P(A)$  = Probability of event A
  - $0 \leq P(A) \leq 1$
- *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.

# Probability

- There are several possible interpretations of probability but they (almost) completely agree on the mathematical rules probability must follow.
  - $P(A)$  = Probability of event A
  - $0 \leq P(A) \leq 1$
- *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: For the same event, two separate people could have different viewpoints and so assign different probabilities.
  - Largely popularized by revolutionary advance in computational technology and methods during the last twenty years.

## Practice

Which of the following events would you be most surprised by?

- (a) exactly 3 heads in 10 coin flips
- (b) exactly 3 heads in 100 coin flips
- (c) exactly 3 heads in 1000 coin flips

## Practice

Which of the following events would you be most surprised by?

- (a) exactly 3 heads in 10 coin flips
- (b) exactly 3 heads in 100 coin flips
- (c) *exactly 3 heads in 1000 coin flips*

## Law of large numbers

*Law of large numbers* states that as more observations are collected, the proportion of occurrences with a particular outcome,  $\hat{p}_n$ , converges to the probability of that outcome,  $p$ .

## Law of large numbers (cont.)

When tossing a *fair* coin, if heads comes up on each of the first 10 tosses, 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 ?

## Law of large numbers (cont.)

When tossing a *fair* coin, if heads comes up on each of the first 10 tosses, 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 ?

- The probability is still 0.5, or there is still a 50% chance that another head will come up on the next toss.

$$P(H \text{ on } 11^{\text{th}} \text{ toss}) = P(T \text{ on } 11^{\text{th}} \text{ toss}) = 0.5$$

## Law of large numbers (cont.)

When tossing a *fair* coin, if heads comes up on each of the first 10 tosses, 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 0.5, or there is still a 50% chance that another head will come up on the next toss.

$$P(H \text{ on } 11^{\text{th}} \text{ toss}) = P(T \text{ on } 11^{\text{th}} \text{ toss}) = 0.5$$

- The coin is not “due” for a tail.

## Law of large numbers (cont.)

When tossing a *fair* coin, if heads comes up on each of the first 10 tosses, 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 0.5, or there is still a 50% chance that another head will come up on the next toss.

$$P(H \text{ on } 11^{\text{th}} \text{ toss}) = P(T \text{ on } 11^{\text{th}} \text{ toss}) = 0.5$$

- The coin is not “due” for a tail.
- The common misunderstanding of the LLN is that random processes are supposed to compensate for whatever happened in the past; this is just not true and is also called *gambler's fallacy* (or *law of averages*).

## Disjoint and non-disjoint outcomes

*Disjoint (mutually exclusive) outcomes:* Cannot happen at the same time.

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

## Disjoint and non-disjoint outcomes

*Disjoint (mutually exclusive) outcomes:* Cannot happen at the same time.

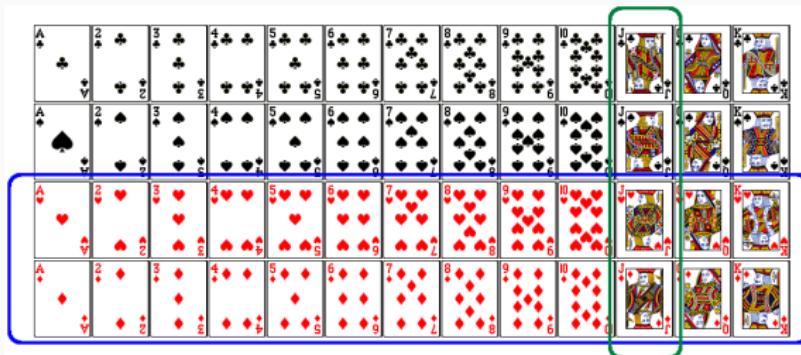
- The outcome of a single coin toss cannot be a head and a tail.
- A student both cannot fail and pass a class.
- A single card drawn from a deck cannot be an ace and a queen.

*Non-disjoint outcomes:* Can happen at the same time.

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

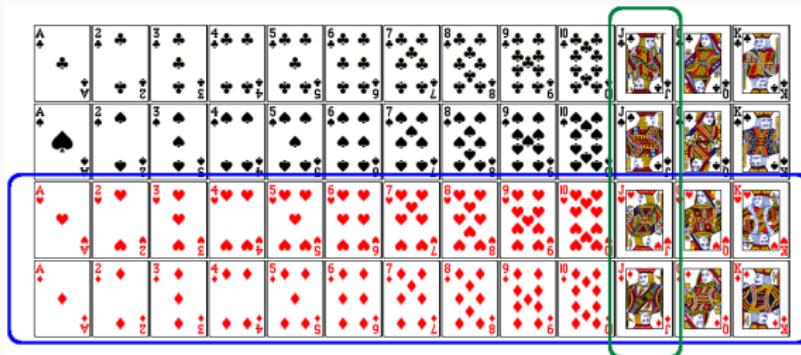
## Union of non-disjoint events

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



## Union of non-disjoint events

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



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

$$= \frac{4}{52} + \frac{26}{52} - \frac{2}{52} = \frac{28}{52}$$

(Can draw a Venn diagram)

## Practice

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

Legalize MJ	Share Parents' Politics		Total
	No	Yes	
No	11	40	51
Yes	36	78	114
Total	47	118	165

- (a)  $\frac{40+36-78}{165}$
- (b)  $\frac{114+118-78}{165}$
- (c)  $\frac{78}{165}$
- (d)  $\frac{78}{188}$
- (e)  $\frac{11}{47}$

## Practice

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

Legalize MJ	Share Parents' Politics		Total
	No	Yes	
No	11	40	51
Yes	36	78	114
Total	47	118	165

- (a)  $\frac{40+36-78}{165}$
- (b)  $\frac{114+118-78}{165}$
- (c)  $\frac{78}{165}$
- (d)  $\frac{78}{188}$
- (e)  $\frac{11}{47}$

## Recap

### General addition rule

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

---

**Note:** For disjoint events  $P(A \text{ and } B) = 0$ , so the above formula simplifies to

$$P(A \text{ or } B) = P(A) + P(B).$$

# Probability distributions

A *probability distribution* lists all possible events and the probabilities with which they occur.

- The probability distribution for the gender of one kid:

Event	Male	Female
Probability	0.5	0.5

# Probability distributions

A *probability distribution* lists all possible events and the probabilities with which they occur.

- The probability distribution for the gender of one kid:

Event	Male	Female
Probability	0.5	0.5

- Rules for probability distributions:

1. The events listed must be disjoint
2. Each probability must be between 0 and 1
3. The probabilities must total 1

# Probability distributions

A *probability distribution* lists all possible events and the probabilities with which they occur.

- The probability distribution for the gender of one kid:

Event	Male	Female
Probability	0.5	0.5

- Rules for probability distributions:

- The events listed must be disjoint
- Each probability must be between 0 and 1
- The probabilities must total 1

- The probability distribution for the genders of two kids:

Event	MM	FF	MF	FM
Probability	0.25	0.25	0.25	0.25

## Sample space and complements

*Sample space* is the collection of all possible outcomes of a trial.

- A couple has one kid, what is the sample space for the gender of this kid?  $S = \{M, F\}$
- A couple has two kids, what is the sample space for the gender of these kids?

## Sample space and complements

*Sample space* is the collection of all possible outcomes of a trial.

- A couple has one kid, what is the sample space for the gender of this kid?  $S = \{M, F\}$
- A couple has two kids, what is the sample space for the gender of these kids?  $S = \{MM, FF, FM, MF\}$

## Sample space and complements

*Sample space* is the collection of all possible outcomes of a trial.

- A couple has one kid, what is the sample space for the gender of this kid?  $S = \{M, F\}$
- A couple has two kids, what is the sample space for the gender of these kids?  $S = \{MM, FF, FM, MF\}$

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

- A couple has one kid. If we know that the kid is not a boy, what is gender of this kid?  $\{ M, F \} \rightarrow$  Boy and girl are *complementary* outcomes.
- A couple has two kids, if we know that they are not both girls, what are the possible gender combinations for these kids?

## Sample space and complements

*Sample space* is the collection of all possible outcomes of a trial.

- A couple has one kid, what is the sample space for the gender of this kid?  $S = \{M, F\}$
- A couple has two kids, what is the sample space for the gender of these kids?  $S = \{MM, FF, FM, MF\}$

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

- A couple has one kid. If we know that the kid is not a boy, what is gender of this kid? { M, **F** } → Boy and girl are *complementary* outcomes.
- A couple has two kids, if we know that they are not both girls, what are the possible gender combinations for these kids? { **MM**, **FF**, **FM**, **MF** }

## Independence

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

# Independence

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

- Knowing that the coin landed on a head on the first toss does not provide any useful information for determining what the coin will land on in the second toss. → Outcomes of two tosses of a coin are independent.

# Independence

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

- Knowing that the coin landed on a head on the first toss does not provide any useful information for determining what the coin will land on in the second toss. → Outcomes of two tosses of a coin are independent.
- Knowing that the first card drawn from a deck is an ace does provide useful information for determining the probability of drawing an ace in the second draw. → Outcomes of two draws from a deck of cards (without replacement) are dependent.

## Practice

Between January 9-12, 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. Which of the below is true?

Opinion on gun ownership and race ethnicity are most likely

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

## Practice

Between January 9-12, 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. Which of the below is true?

Opinion on gun ownership and race ethnicity are most likely

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

## Checking for independence

If  $P(A \text{ occurs, given that } B \text{ is true}) = P(A | B) = P(A)$ , then A and B are independent.

## Checking for independence

If  $P(A \text{ occurs, given that } B \text{ is true}) = P(A | B) = P(A)$ , then A and B are independent.

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

## Checking for independence

If  $P(A \text{ occurs, given that } B \text{ is true}) = P(A | B) = P(A)$ , then A and B are independent.

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

$P(\text{randomly selected NC resident says gun ownership protects citizens, given that the resident is white}) =$

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

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

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

## Checking for independence

If  $P(A \text{ occurs, given that } B \text{ is true}) = P(A | B) = P(A)$ , then A and B are independent.

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

$P(\text{randomly selected NC resident says gun ownership protects citizens, given that the resident is white}) =$

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

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

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

$P(\text{protects citizens})$  varies by race/ethnicity, therefore opinion on gun ownership and race ethnicity are most likely dependent.

## Determining dependence based on sample data

- If conditional probabilities calculated based on sample data suggest dependence between two variables, the next step is to conduct a hypothesis test to determine if the observed difference between the probabilities is likely or unlikely to have happened by chance.

## Determining dependence based on sample data

- If conditional probabilities calculated based on sample data suggest dependence between two variables, the next step is to conduct a hypothesis test to determine if the observed difference between the probabilities is likely or unlikely to have happened by chance.
- If the observed difference between the conditional probabilities is large, then there is stronger evidence that the difference is real.

## Determining dependence based on sample data

- If conditional probabilities calculated based on sample data suggest dependence between two variables, the next step is to conduct a hypothesis test to determine if the observed difference between the probabilities is likely or unlikely to have happened by chance.
- If the observed difference between the conditional probabilities is large, then there is stronger evidence that the difference is real.
- If a sample is large, then even a small difference can provide strong evidence of a real difference.

## Determining dependence based on sample data

Practice:

We saw that  $P(\text{protects citizens} \mid \text{White}) = 0.67$  and  $P(\text{protects citizens} \mid \text{Hispanic}) = 0.64$ . Under which condition would you be more convinced of a real difference between the proportions of Whites and Hispanics who think widespread gun ownership protects citizens?  $n = 500$  or  $n = 50,000$

## Determining dependence based on sample data

Practice:

We saw that  $P(\text{protects citizens} \mid \text{White}) = 0.67$  and  $P(\text{protects citizens} \mid \text{Hispanic}) = 0.64$ . Under which condition would you be more convinced of a real difference between the proportions of Whites and Hispanics who think widespread gun ownership protects citizens?  $n = 500$  or  $n = 50,000$

$n = 50,000$

## Product rule for independent events

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

Or more generally,  $P(A_1 \text{ and } \dots \text{ and } A_k) = P(A_1) \times \dots \times P(A_k)$

## Product rule for independent events

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

Or more generally,  $P(A_1 \text{ and } \dots \text{ and } A_k) = P(A_1) \times \dots \times P(A_k)$

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

## Product rule for independent events

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

Or more generally,  $P(A_1 \text{ and } \dots \text{ and } A_k) = P(A_1) \times \dots \times P(A_k)$

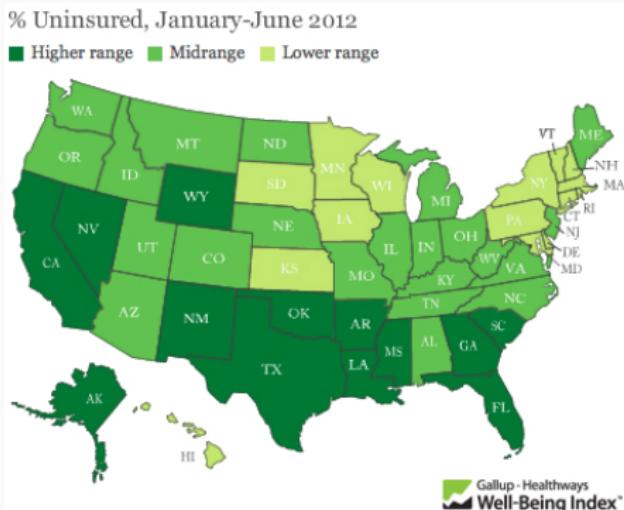
You toss a coin twice, what is the probability of getting two tails in a row?

$$P(\text{T on the first toss}) \times P(\text{T on the second toss}) = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$$

# Practice

A recent Gallup poll suggests that 25.5% of Texans do not have health insurance as of June 2012. Assuming that the uninsured rate stayed constant, what is the probability that two randomly selected Texans are both uninsured?

- (a)  $25.5^2$
- (b)  $0.255^2$
- (c)  $0.255 \times 2$
- (d)  $(1 - 0.255)^2$

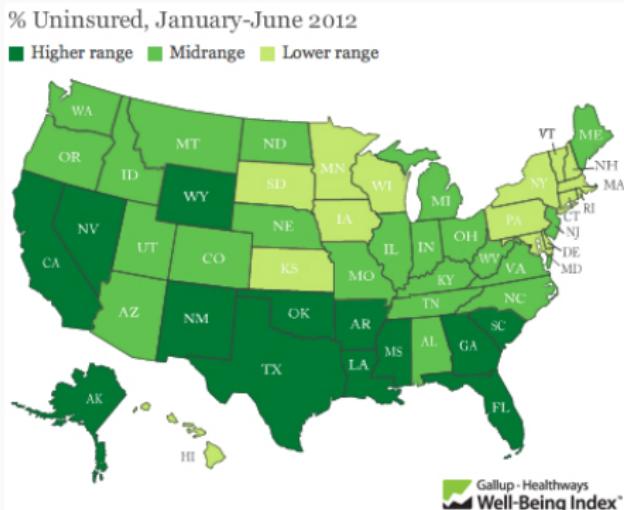


<http://www.gallup.com/poll/156851/uninsured-rate-stable-across-states-far-2012.aspx>

# Practice

A recent Gallup poll suggests that 25.5% of Texans do not have health insurance as of June 2012. Assuming that the uninsured rate stayed constant, what is the probability that two randomly selected Texans are both uninsured?

- (a)  $25.5^2$
- (b)  $0.255^2$
- (c)  $0.255 \times 2$
- (d)  $(1 - 0.255)^2$



<http://www.gallup.com/poll/156851/uninsured-rate-stable-across-states-far-2012.aspx>

## Disjoint vs. complementary

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

## Disjoint vs. complementary

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

*Not necessarily, there may be more than 2 events in the sample space, e.g. party affiliation.*

## Disjoint vs. complementary

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

*Not necessarily, there may be more than 2 events in the sample space, e.g. party affiliation.*

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

## Disjoint vs. complementary

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

*Not necessarily, there may be more than 2 events in the sample space, e.g. party affiliation.*

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

*Yes, that's the definition of complementary, e.g. heads and tails.*

## Putting everything together...

If we were to randomly select 5 Texans, what is the probability that at least one is uninsured?

- If we were to randomly select 5 Texans, the sample space for the number of Texans who are uninsured would be:

$$S = \{0, 1, 2, 3, 4, 5\}$$

- We are interested in instances where at least one person is uninsured:

$$S = \{0, 1, 2, 3, 4, 5\}$$

- So we can divide up the sample space into two categories:

$$S = \{0, \text{at least one}\}$$

## Putting everything together...

Since the probability of the sample space must add up to 1:

$$\text{Prob}(\text{at least 1 uninsured}) = 1 - \text{Prob}(\text{none uninsured})$$

## Putting everything together...

Since the probability of the sample space must add up to 1:

$$\begin{aligned} \text{Prob}(at \ least \ 1 \ uninsured) &= 1 - \text{Prob}(none \ uninsured) \\ &= 1 - [(1 - 0.255)^5] \end{aligned}$$

## Putting everything together...

Since the probability of the sample space must add up to 1:

$$\begin{aligned} \text{Prob(at least 1 uninsured)} &= 1 - \text{Prob}(none uninsured) \\ &= 1 - [(1 - 0.255)^5] \\ &= 1 - 0.745^5 \end{aligned}$$

## Putting everything together...

Since the probability of the sample space must add up to 1:

$$\begin{aligned} \text{Prob(at least 1 uninsured)} &= 1 - \text{Prob}(none uninsured) \\ &= 1 - [(1 - 0.255)^5] \\ &= 1 - 0.745^5 \\ &= 1 - 0.23 \end{aligned}$$

## Putting everything together...

Since the probability of the sample space must add up to 1:

$$\begin{aligned} \text{Prob(at least 1 uninsured)} &= 1 - \text{Prob}(none uninsured) \\ &= 1 - [(1 - 0.255)^5] \\ &= 1 - 0.745^5 \\ &= 1 - 0.23 \\ &= 0.77 \end{aligned}$$

At least 1

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

## Practice

Roughly 20% of undergraduates at a university are vegetarian or vegan. What is the probability that, among a random sample of 3 undergraduates, at least one is vegetarian or vegan?

- (a)  $1 - 0.2 \times 3$
- (b)  $1 - 0.2^3$
- (c)  $0.8^3$
- (d)  $1 - 0.8 \times 3$
- (e)  $1 - 0.8^3$

## Practice

Roughly 20% of undergraduates at a university are vegetarian or vegan. What is the probability that, among a random sample of 3 undergraduates, at least one is vegetarian or vegan?

(a)  $1 - 0.2 \times 3$

(b)  $1 - 0.2^3$

(c)  $0.8^3$

(d)  $1 - 0.8 \times 3$

(e)  $1 - 0.8^3$

$$\begin{aligned}P(\text{at least 1 from veg}) &= 1 - P(\text{none veg}) \\&= 1 - (1 - 0.2)^3 \\&= 1 - 0.8^3 \\&= 1 - 0.512 = 0.488\end{aligned}$$

## Conditional probability

---

## Relapse

Researchers randomly assigned 72 chronic users of cocaine into three groups: desipramine (antidepressant), lithium (standard treatment for cocaine) and placebo. Results of the study are summarized below.

	no		total
	relapse	relapse	
desipramine	10	14	24
lithium	18	6	24
placebo	20	4	24
total	48	24	72

What is the probability that a patient did not relapse?

	relapse	no relapse	total
desipramine	10	14	24
lithium	18	6	24
placebo	20	4	24
total	48	24	72

## Marginal probability

What is the probability that a patient relapsed?

	relapse	no relapse	total
desipramine	10	14	24
lithium	18	6	24
placebo	20	4	24
total	48	24	72

## Marginal probability

What is the probability that a patient relapsed?

	relapse	no relapse	total
desipramine	10	14	24
lithium	18	6	24
placebo	20	4	24
total	48	24	72

$$P(\text{relapsed}) = \frac{48}{72} \approx 0.67$$

## Joint probability

What is the probability that a patient received the antidepressant (desipramine) and relapsed?

	no relapse		total
	relapse	no relapse	
desipramine	10	14	24
lithium	18	6	24
placebo	20	4	24
total	48	24	72

## Joint probability

What is the probability that a patient received the antidepressant (desipramine) and relapsed?

	relapse	no relapse	total
desipramine	10	14	24
lithium	18	6	24
placebo	20	4	24
total	48	24	72

$$P(\text{relapsed and desipramine}) = \frac{10}{72} \approx 0.14$$

## Conditional probability

### Conditional probability

The conditional probability of the outcome of interest  $A$  given condition  $B$  is calculated as

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

# Conditional probability

## Conditional probability

The conditional probability of the outcome of interest  $A$  given condition  $B$  is calculated as

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

	relapse	no relapse	total	$P(\text{relapse} \text{desipramine})$
desipramine	10	14	24	$= \frac{P(\text{relapse and desipramine})}{P(\text{desipramine})}$
lithium	18	6	24	
placebo	20	4	24	
total	48	24	72	

# Conditional probability

## Conditional probability

The conditional probability of the outcome of interest  $A$  given condition  $B$  is calculated as

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

	relapse	no relapse	total	$P(\text{relapse} \text{desipramine})$
desipramine	10	14	24	$= \frac{P(\text{relapse and desipramine})}{P(\text{desipramine})}$
lithium	18	6	24	$= \frac{10/72}{24/72}$
placebo	20	4	24	
total	48	24	72	

# Conditional probability

## Conditional probability

The conditional probability of the outcome of interest  $A$  given condition  $B$  is calculated as

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

	relapse	no relapse	total	$P(\text{relapse} \text{desipramine})$
desipramine	10	14	24	$= \frac{P(\text{relapse and desipramine})}{P(\text{desipramine})}$
lithium	18	6	24	$= \frac{10/72}{24/72}$
placebo	20	4	24	$= \frac{10}{24}$
total	48	24	72	

# Conditional probability

## Conditional probability

The conditional probability of the outcome of interest  $A$  given condition  $B$  is calculated as

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

	relapse	no relapse	total	$P(\text{relapse} \text{desipramine})$
desipramine	10	14	24	$= \frac{P(\text{relapse and desipramine})}{P(\text{desipramine})}$
lithium	18	6	24	$= \frac{10/72}{24/72}$
placebo	20	4	24	$= \frac{10}{24}$
total	48	24	72	$= 0.42$

## Conditional probability (cont.)

If we know that a patient received the antidepressant (desipramine), what is the probability that they relapsed?

	no relapse		
	relapse	no relapse	total
desipramine	10	14	24
lithium	18	6	24
placebo	20	4	24
total	48	24	72

## Conditional probability (cont.)

If we know that a patient received the antidepressant (desipramine), what is the probability that they relapsed?

		no relapse		total
		relapse	no relapse	
desipramine	relapse	10	14	24
	no relapse	18	6	24
lithium	20	4	24	
placebo	48	24	72	
total				

$$P(\text{relapse} \mid \text{desipramine}) = \frac{10}{24} \approx 0.42$$

## Conditional probability (cont.)

If we know that a patient received the antidepressant (desipramine), what is the probability that they relapsed?

		no relapse		total
		relapse	no relapse	
desipramine	relapse	10	14	24
	no relapse	18	6	24
lithium	20	4	24	
placebo	48	24	72	
total				

$$P(\text{relapse} \mid \text{desipramine}) = \frac{10}{24} \approx 0.42$$

$$P(\text{relapse} \mid \text{lithium}) = \frac{18}{24} \approx 0.75$$

$$P(\text{relapse} \mid \text{placebo}) = \frac{20}{24} \approx 0.83$$

## Conditional probability (cont.)

If we know that a patient relapsed, what is the probability that they received the antidepressant (desipramine)?

	relapse	no relapse	total
desipramine	10	14	24
lithium	18	6	24
placebo	20	4	24
total	48	24	72

## Conditional probability (cont.)

If we know that a patient relapsed, what is the probability that they received the antidepressant (desipramine)?

	relapse	no relapse	total
desipramine	10	14	24
lithium	18	6	24
placebo	20	4	24
total	48	24	72

$$P(\text{desipramine} \mid \text{relapse}) = \frac{10}{48} \approx 0.21$$

## Conditional probability (cont.)

If we know that a patient relapsed, what is the probability that they received the antidepressant (desipramine)?

	relapse	no relapse	total
desipramine	10	14	24
lithium	18	6	24
placebo	20	4	24
total	48	24	72

$$P(\text{desipramine} \mid \text{relapse}) = \frac{10}{48} \approx 0.21$$

$$P(\text{lithium} \mid \text{relapse}) = \frac{18}{48} \approx 0.375$$

$$P(\text{placebo} \mid \text{relapse}) = \frac{20}{48} \approx 0.42$$

## General multiplication rule

- Earlier: if  $A$  and  $B$  are *independent*, then

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

## General multiplication rule

- Earlier: if  $A$  and  $B$  are *independent*, then

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

- If  $A$  and  $B$  are any two outcomes or events, then

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

Note that this formula is simply the conditional probability formula, rearranged.

## General multiplication rule

- Earlier: if  $A$  and  $B$  are *independent*, then

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

- If  $A$  and  $B$  are any two outcomes or events, then

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

Note that this formula is simply the conditional probability formula, rearranged.

- It is useful to think of  $A$  as the outcome of interest and  $B$  as the condition.

## Independence and conditional probabilities

Consider the following (hypothetical) distribution of gender and major of students in a class:

	social science	non-social science	total
female	30	20	50
male	30	20	50
total	60	40	100

## Independence and conditional probabilities

Consider the following (hypothetical) distribution of gender and major of students in a class:

	social science	non-social science	total
female	30	20	50
male	30	20	50
total	60	40	100

- The probability that a randomly selected student is a social science major is

## Independence and conditional probabilities

Consider the following (hypothetical) distribution of gender and major of students in a class:

	social science	non-social science	total
female	30	20	50
male	30	20	50
total	60	40	100

- The probability that a randomly selected student is a social science major is  $\frac{60}{100} = 0.6$ .

## Independence and conditional probabilities

Consider the following (hypothetical) distribution of gender and major of students in a class:

	social science	non-social science	total
female	30	20	50
male	30	20	50
total	60	40	100

- The probability that a randomly selected student is a social science major is  $\frac{60}{100} = 0.6$ .
- The probability that a randomly selected student is a social science major given that they are female is

## Independence and conditional probabilities

Consider the following (hypothetical) distribution of gender and major of students in a class:

	social science	non-social science	total
female	30	20	50
male	30	20	50
total	60	40	100

- The probability that a randomly selected student is a social science major is  $\frac{60}{100} = 0.6$ .
- The probability that a randomly selected student is a social science major given that they are female is  $\frac{30}{50} = 0.6$ .

## Independence and conditional probabilities

Consider the following (hypothetical) distribution of gender and major of students in a class:

	social science	non-social science	total
female	30	20	50
male	30	20	50
total	60	40	100

- The probability that a randomly selected student is a social science major is  $\frac{60}{100} = 0.6$ .
- The probability that a randomly selected student is a social science major given that they are female is  $\frac{30}{50} = 0.6$ .
- Since  $P(SS|M)$  also equals 0.6, major of students in this class does not depend on their gender:

$$P(SS | F) = P(SS | M) = P(SS).$$

## Independence and conditional probabilities (cont.)

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

## Independence and conditional probabilities (cont.)

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  $A$ .

## Independence and conditional probabilities (cont.)

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: We know that if events  $A$  and  $B$  are independent,  $P(A \text{ and } B) = P(A) \times P(B)$ . Then,

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

## Breast cancer screening

- American Cancer Society estimates that about **1.7%** of women have breast cancer.  
*<http://www.cancer.org/cancer/cancerbasics/cancer-prevalence>*
- Susan G. Komen For The Cure Foundation states that mammography correctly identifies about **78%** of women who truly have breast cancer.  
*http://ww5.komen.org/BreastCancer/AccuracyofMammograms.html*
- An article published in 2003 suggests that **up to 10%** of all mammograms result in false positives for patients who do not have cancer.  
*<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1360940>*

---

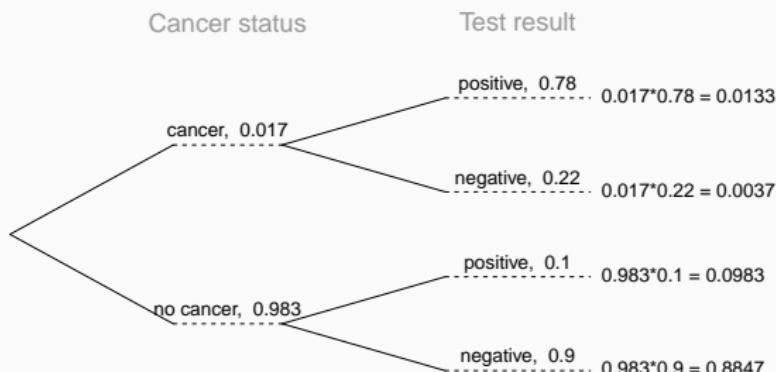
**Note:** These percentages are approximate, and very difficult to estimate.

## Inverting probabilities

When a patient goes through breast cancer screening there are two competing claims: patient has cancer and patient doesn't have cancer. If a mammogram yields a positive result, what is the probability that patient actually has cancer?

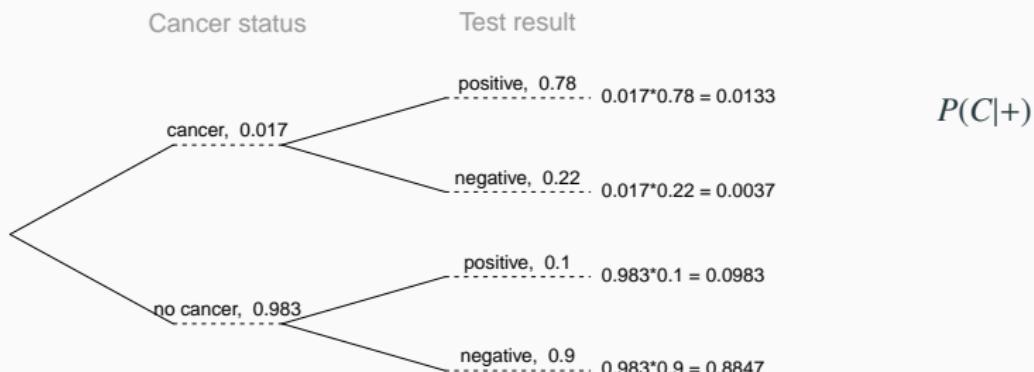
# Inverting probabilities

When a patient goes through breast cancer screening there are two competing claims: patient has cancer and patient doesn't have cancer. If a mammogram yields a positive result, what is the probability that patient actually has cancer?



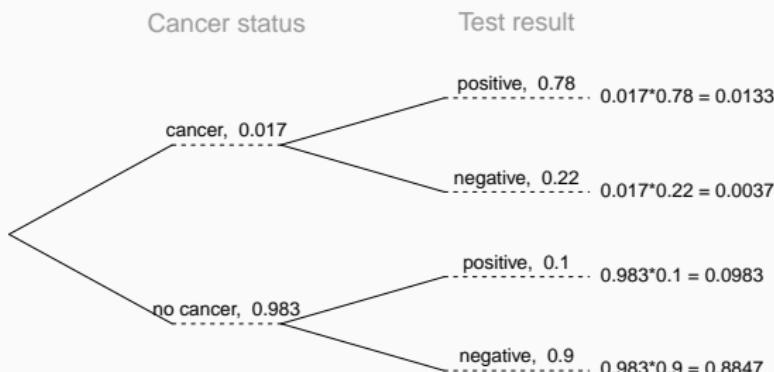
# Inverting probabilities

When a patient goes through breast cancer screening there are two competing claims: patient has cancer and patient doesn't have cancer. If a mammogram yields a positive result, what is the probability that patient actually has cancer?



# Inverting probabilities

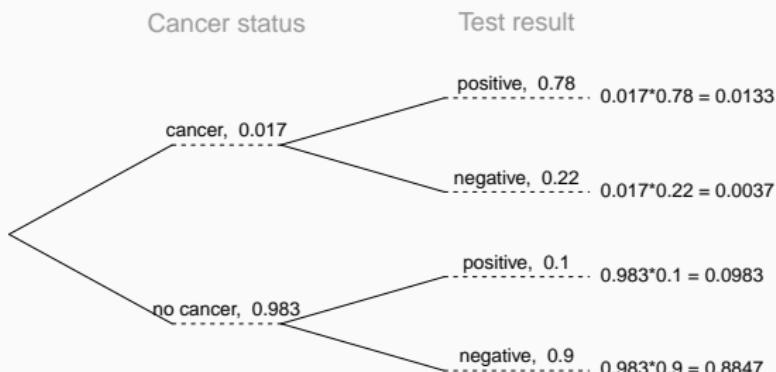
When a patient goes through breast cancer screening there are two competing claims: patient has cancer and patient doesn't have cancer. If a mammogram yields a positive result, what is the probability that patient actually has cancer?



$$P(C|+)$$
$$= \frac{P(C \text{ and } +)}{P(+)}$$

# Inverting probabilities

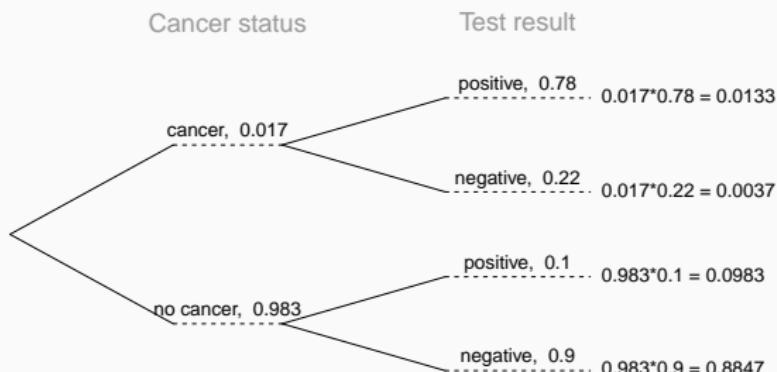
When a patient goes through breast cancer screening there are two competing claims: patient has cancer and patient doesn't have cancer. If a mammogram yields a positive result, what is the probability that patient actually has cancer?



$$\begin{aligned}P(C|+) &= \frac{P(C \text{ and } +)}{P(+)} \\&= \frac{0.0133}{0.0133 + 0.0983}\end{aligned}$$

# Inverting probabilities

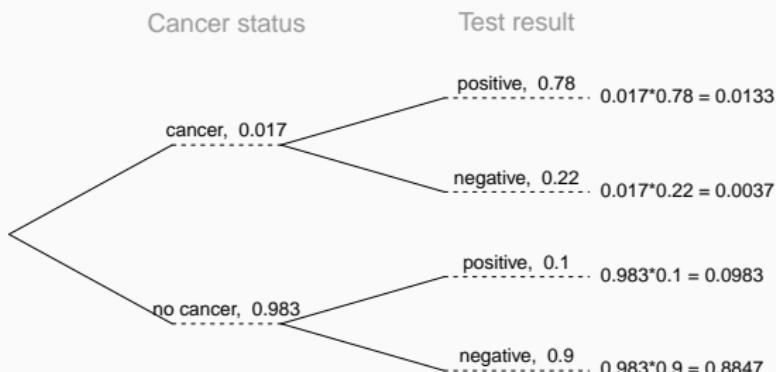
When a patient goes through breast cancer screening there are two competing claims: patient has cancer and patient doesn't have cancer. If a mammogram yields a positive result, what is the probability that patient actually has cancer?



$$\begin{aligned}P(C|+) &= \frac{P(C \text{ and } +)}{P(+)} \\&= \frac{0.0133}{0.0133 + 0.0983} \\&= 0.12\end{aligned}$$

# Inverting probabilities

When a patient goes through breast cancer screening there are two competing claims: patient has cancer and patient doesn't have cancer. If a mammogram yields a positive result, what is the probability that patient actually has cancer?



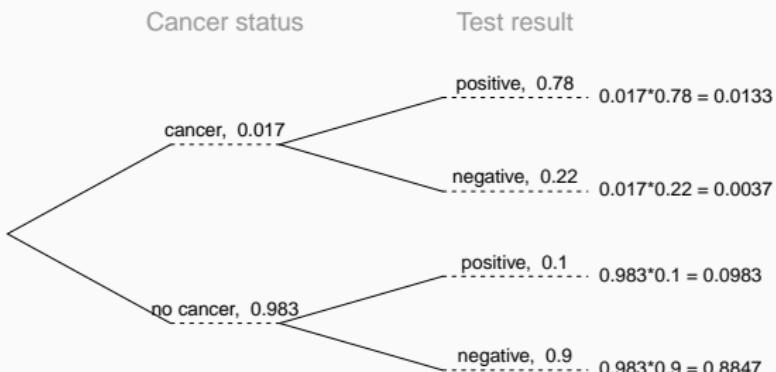
$$\begin{aligned}P(C|+) &= \frac{P(C \text{ and } +)}{P(+)} \\&= \frac{0.0133}{0.0133 + 0.0983} \\&= 0.12\end{aligned}$$

**Note:** Tree diagrams are useful for inverting probabilities: we are given  $P(+|C)$  and asked for  $P(C|+)$ .

## Practice

Suppose a woman who gets tested once and obtains a positive result wants to get tested again. In the second test, what should we assume to be the probability of this specific woman having cancer?

- (a) 0.017
- (b) 0.12
- (c) 0.0133
- (d) 0.88

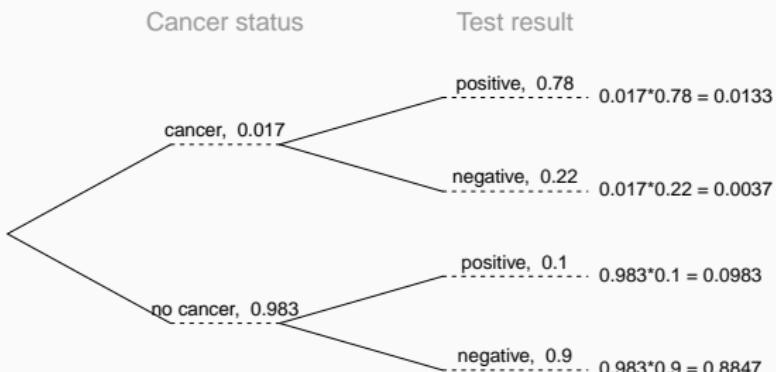


$$P(C | +) = 0.12.$$

## Practice

Suppose a woman who gets tested once and obtains a positive result wants to get tested again. In the second test, what should we assume to be the probability of this specific woman having cancer?

- (a) 0.017
- (b) 0.12
- (c) 0.0133
- (d) 0.88



$$P(C | +) = 0.12.$$

## Practice

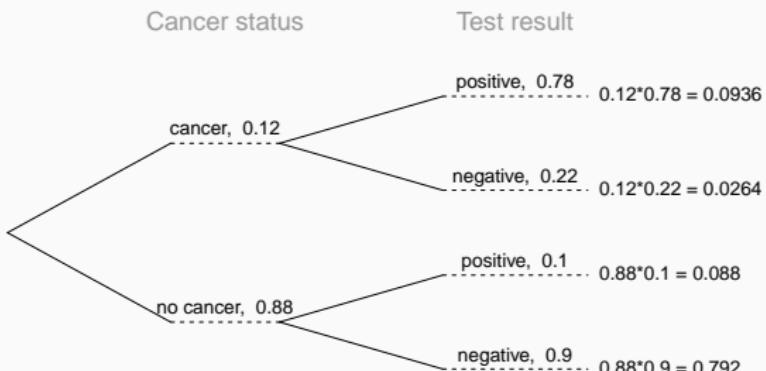
What is the probability that this woman has cancer if this second mammogram also yielded a positive result?

- (a) 0.0936
- (b) 0.088
- (c) 0.48
- (d) 0.52

## Practice

What is the probability that this woman has cancer if this second mammogram also yielded a positive result?

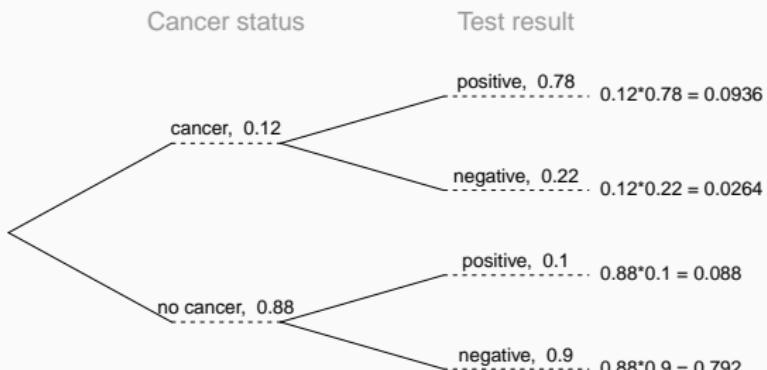
- (a) 0.0936
- (b) 0.088
- (c) 0.48
- (d) 0.52



## Practice

What is the probability that this woman has cancer if this second mammogram also yielded a positive result?

- (a) 0.0936
- (b) 0.088
- (c) 0.48
- (d) 0.52



$$P(C|+) = \frac{P(C \text{ and } +)}{P(+)} = \frac{0.0936}{0.0936 + 0.088} = 0.52$$

## Bayes' Theorem

---

- The conditional probability formula we have seen so far is a special case of the Bayes' Theorem, which is applicable even when events have more than just two outcomes.

## Bayes' Theorem

- The conditional probability formula we have seen so far is a special case of the Bayes' Theorem, which is applicable even when events have more than just two outcomes.
- Bayes' Theorem:*

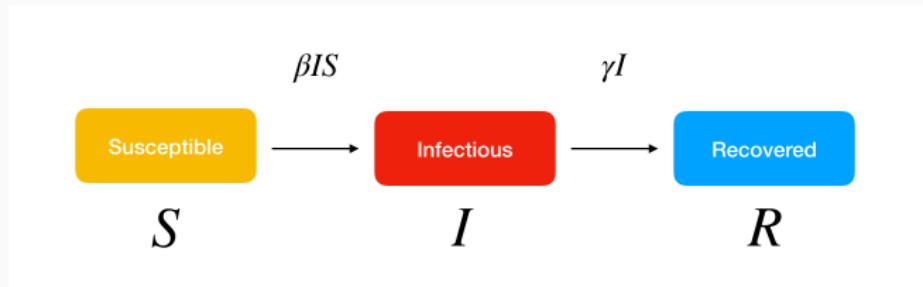
$P(\text{outcome } A_1 \text{ of variable 1} \mid \text{outcome } B \text{ of variable 2})$

$$= \frac{P(B|A_1)P(A_1)}{P(B|A_1)P(A_1) + P(B|A_2)P(A_2) + \cdots + P(B|A_k)P(A_k)}$$

where  $A_2, \dots, A_k$  represent all other possible outcomes of variable 1.

## Application activity: Inverting probabilities

A common epidemiological model for the spread of diseases is the SIR model, where the population is partitioned into three groups: Susceptible (S), Infectious (I), and Recovered (R).



Imagine a population amidst an epidemic that

- Susceptible: 60%
- Infectious: 10%
- Recovered: 30%

## Application activity: Inverting probabilities

Imagine a population amidst an epidemic that

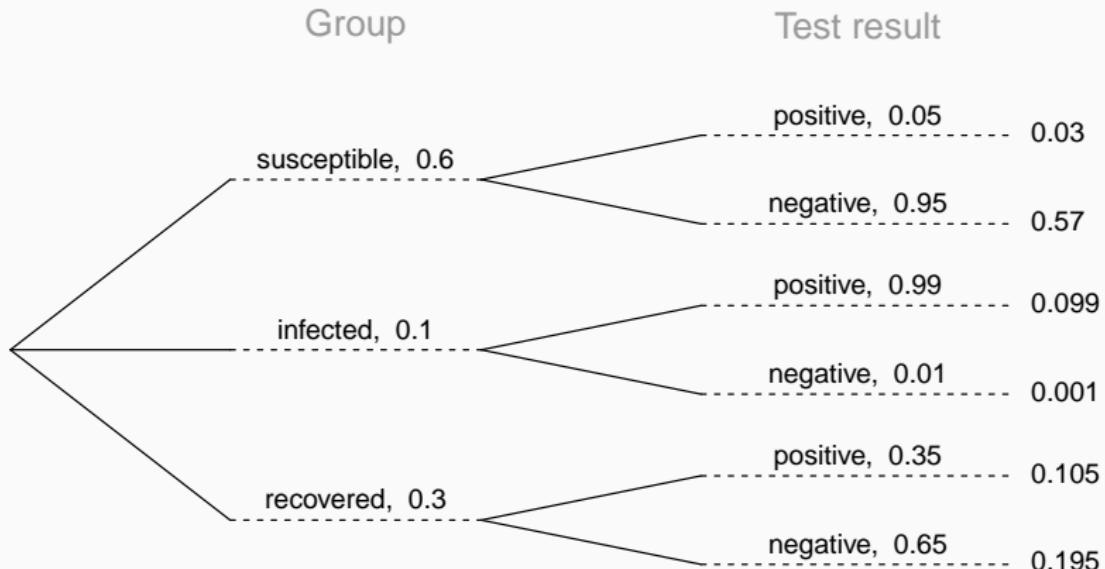
- Susceptible: 60%
- Infectious: 10%
- Recovered: 30%

Suppose there is a test (e.g., the COVID nucleic acid test) that is not 100% accurate, and the accuracy on different groups is

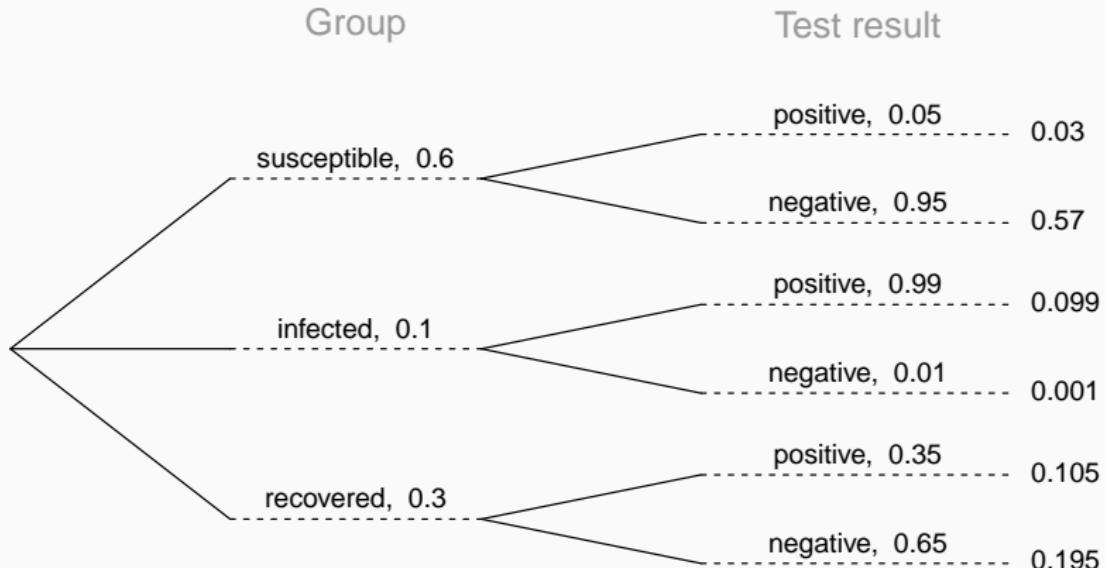
- Susceptible: 95% accurate
- Infectious: 99% accurate
- Recovered: 65% accurate

Draw a probability tree to reflect the information given above. If the individual has tested positive, what is the probability that they are actually infected?

## Application activity: Inverting probabilities (cont.)



## Application activity: Inverting probabilities (cont.)



$$P(\text{inf}|+) = \frac{P(\text{inf and } +)}{P(+)} = \frac{0.099}{0.03 + 0.099 + 0.105} \approx 0.423$$

## **Sampling from a small population**

---

## Sampling with replacement

When sampling *with replacement*, you put back what you just drew.

## Sampling with replacement

When sampling *with replacement*, you put back what you just drew.

- Imagine you have a bag with 5 red, 3 blue and 2 orange chips in it. What is the probability that the first chip you draw is blue?

5 ● , 3 ● , 2 ●

## Sampling with replacement

When sampling *with replacement*, you put back what you just drew.

- Imagine you have a bag with 5 red, 3 blue and 2 orange chips in it. What is the probability that the first chip you draw is blue?

5 ● , 3 ● , 2 ●

$$Prob(1^{st} \text{ chip } B) = \frac{3}{5 + 3 + 2} = \frac{3}{10} = 0.3$$

## Sampling with replacement

When sampling *with replacement*, you put back what you just drew.

- Imagine you have a bag with 5 red, 3 blue and 2 orange chips in it. What is the probability that the first chip you draw is blue?

5 ● , 3 ● , 2 ●

$$Prob(1^{\text{st}} \text{ chip } B) = \frac{3}{5 + 3 + 2} = \frac{3}{10} = 0.3$$

- Suppose you did indeed pull a blue chip in the first draw. If drawing with replacement, what is the probability of drawing a blue chip in the second draw?

## Sampling with replacement

When sampling *with replacement*, you put back what you just drew.

- Imagine you have a bag with 5 red, 3 blue and 2 orange chips in it. What is the probability that the first chip you draw is blue?

5 ● , 3 ● , 2 ●

$$Prob(1^{st} \text{ chip } B) = \frac{3}{5 + 3 + 2} = \frac{3}{10} = 0.3$$

- Suppose you did indeed pull a blue chip in the first draw. If drawing with replacement, what is the probability of drawing a blue chip in the second draw?

1<sup>st</sup> draw: 5 ● , 3 ● , 2 ●

## Sampling with replacement

When sampling *with replacement*, you put back what you just drew.

- Imagine you have a bag with 5 red, 3 blue and 2 orange chips in it. What is the probability that the first chip you draw is blue?

5 ● , 3 ● , 2 ●

$$\text{Prob}(1^{\text{st}} \text{ chip } B) = \frac{3}{5 + 3 + 2} = \frac{3}{10} = 0.3$$

- Suppose you did indeed pull a blue chip in the first draw. If drawing with replacement, what is the probability of drawing a blue chip in the second draw?

1<sup>st</sup> draw: 5 ● , 3 ● , 2 ●

2<sup>nd</sup> draw: 5 ● , 3 ● , 2 ●

## Sampling with replacement

When sampling *with replacement*, you put back what you just drew.

- Imagine you have a bag with 5 red, 3 blue and 2 orange chips in it. What is the probability that the first chip you draw is blue?

5 ● , 3 ● , 2 ●

$$\text{Prob}(1^{\text{st}} \text{ chip } B) = \frac{3}{5 + 3 + 2} = \frac{3}{10} = 0.3$$

- Suppose you did indeed pull a blue chip in the first draw. If drawing with replacement, what is the probability of drawing a blue chip in the second draw?

1<sup>st</sup> draw: 5 ● , 3 ● , 2 ●

2<sup>nd</sup> draw: 5 ● , 3 ● , 2 ●

$$\text{Prob}(2^{\text{nd}} \text{ chip } B | 1^{\text{st}} \text{ chip } B) = \frac{3}{10} = 0.3$$

## Sampling with replacement (cont.)

- Suppose you actually pulled an orange chip in the first draw. If drawing with replacement, what is the probability of drawing a blue chip in the second draw?

## Sampling with replacement (cont.)

- Suppose you actually pulled an orange chip in the first draw. If drawing with replacement, what is the probability of drawing a blue chip in the second draw?

1<sup>st</sup> draw: 5 ● , 3 ● , 2 ●

## Sampling with replacement (cont.)

- Suppose you actually pulled an orange chip in the first draw. If drawing with replacement, what is the probability of drawing a blue chip in the second draw?

1<sup>st</sup> draw: 5 ● , 3 ● , 2 ●

2<sup>nd</sup> draw: 5 ● , 3 ● , 2 ●

## Sampling with replacement (cont.)

- Suppose you actually pulled an orange chip in the first draw. If drawing with replacement, what is the probability of drawing a blue chip in the second draw?

1<sup>st</sup> draw: 5 ●, 3 ●, 2 ●

2<sup>nd</sup> draw: 5 ●, 3 ●, 2 ●

$$Prob(2^{\text{nd}} \text{ chip } B | 1^{\text{st}} \text{ chip } O) = \frac{3}{10} = 0.3$$

## Sampling with replacement (cont.)

- Suppose you actually pulled an orange chip in the first draw. If drawing with replacement, what is the probability of drawing a blue chip in the second draw?

1<sup>st</sup> draw: 5 ●, 3 ●, 2 ●

2<sup>nd</sup> draw: 5 ●, 3 ●, 2 ●

$$Prob(2^{\text{nd}} \text{ chip } B | 1^{\text{st}} \text{ chip } O) = \frac{3}{10} = 0.3$$

- If drawing with replacement, what is the probability of drawing two blue chips in a row?

## Sampling with replacement (cont.)

- Suppose you actually pulled an orange chip in the first draw. If drawing with replacement, what is the probability of drawing a blue chip in the second draw?

1<sup>st</sup> draw: 5 ●, 3 ●, 2 ●

2<sup>nd</sup> draw: 5 ●, 3 ●, 2 ●

$$Prob(2^{\text{nd}} \text{ chip } B | 1^{\text{st}} \text{ chip } O) = \frac{3}{10} = 0.3$$

- If drawing with replacement, what is the probability of drawing two blue chips in a row?

1<sup>st</sup> draw: 5 ●, 3 ●, 2 ●

2<sup>nd</sup> draw: 5 ●, 3 ●, 2 ●

## Sampling with replacement (cont.)

- Suppose you actually pulled an orange chip in the first draw. If drawing with replacement, what is the probability of drawing a blue chip in the second draw?

1<sup>st</sup> draw: 5 ●, 3 ●, 2 ●

2<sup>nd</sup> draw: 5 ●, 3 ●, 2 ●

$$Prob(2^{\text{nd}} \text{ chip } B | 1^{\text{st}} \text{ chip } O) = \frac{3}{10} = 0.3$$

- If drawing with replacement, what is the probability of drawing two blue chips in a row?

1<sup>st</sup> draw: 5 ●, 3 ●, 2 ●

2<sup>nd</sup> draw: 5 ●, 3 ●, 2 ●

$$\begin{aligned} Prob(1^{\text{st}} \text{ chip } B) \cdot Prob(2^{\text{nd}} \text{ chip } B | 1^{\text{st}} \text{ chip } B) &= 0.3 \times 0.3 \\ &= 0.3^2 = 0.09 \end{aligned}$$

## Sampling with replacement (cont.)

- When drawing with replacement, probability of the second chip being blue does not depend on the color of the first chip since whatever we draw in the first draw gets put back in the bag.

$$\text{Prob}(B|B) = \text{Prob}(B|O)$$

- In addition, this probability is equal to the probability of drawing a blue chip in the first draw, since the composition of the bag never changes when sampling with replacement.

$$\text{Prob}(B|B) = \text{Prob}(B)$$

- When drawing with replacement, draws are independent.*

## Sampling without replacement

When drawing *without replacement* you do not put back what you just drew.

## Sampling without replacement

When drawing *without replacement* you do not put back what you just drew.

- Suppose you pulled a blue chip in the first draw. If drawing without replacement, what is the probability of drawing a blue chip in the second draw?

## Sampling without replacement

When drawing *without replacement* you do not put back what you just drew.

- Suppose you pulled a blue chip in the first draw. If drawing without replacement, what is the probability of drawing a blue chip in the second draw?

1<sup>st</sup> draw: 5 ● , 3 ● , 2 ●

## Sampling without replacement

When drawing *without replacement* you do not put back what you just drew.

- Suppose you pulled a blue chip in the first draw. If drawing without replacement, what is the probability of drawing a blue chip in the second draw?

1<sup>st</sup> draw: 5 ● , 3 ● , 2 ●

2<sup>nd</sup> draw: 5 ● , 2 ● , 2 ●

## Sampling without replacement

When drawing *without replacement* you do not put back what you just drew.

- Suppose you pulled a blue chip in the first draw. If drawing without replacement, what is the probability of drawing a blue chip in the second draw?

1<sup>st</sup> draw: 5 ● , 3 ● , 2 ●

2<sup>nd</sup> draw: 5 ● , 2 ● , 2 ●

$$Prob(2^{\text{nd}} \text{ chip } B | 1^{\text{st}} \text{ chip } B) = \frac{2}{9} = 0.22$$

## Sampling without replacement

When drawing *without replacement* you do not put back what you just drew.

- Suppose you pulled a blue chip in the first draw. If drawing without replacement, what is the probability of drawing a blue chip in the second draw?

1<sup>st</sup> draw: 5 ● , 3 ● , 2 ●

2<sup>nd</sup> draw: 5 ● , 2 ● , 2 ●

$$Prob(2^{\text{nd}} \text{ chip } B | 1^{\text{st}} \text{ chip } B) = \frac{2}{9} = 0.22$$

- If drawing without replacement, what is the probability of drawing two blue chips in a row?

## Sampling without replacement

When drawing *without replacement* you do not put back what you just drew.

- Suppose you pulled a blue chip in the first draw. If drawing without replacement, what is the probability of drawing a blue chip in the second draw?

1<sup>st</sup> draw: 5 ● , 3 ● , 2 ●

2<sup>nd</sup> draw: 5 ● , 2 ● , 2 ●

$$\text{Prob}(2^{\text{nd}} \text{ chip } B | 1^{\text{st}} \text{ chip } B) = \frac{2}{9} = 0.22$$

- If drawing without replacement, what is the probability of drawing two blue chips in a row?

1<sup>st</sup> draw: 5 ● , 3 ● , 2 ●

2<sup>nd</sup> draw: 5 ● , 2 ● , 2 ●

## Sampling without replacement

When drawing *without replacement* you do not put back what you just drew.

- Suppose you pulled a blue chip in the first draw. If drawing without replacement, what is the probability of drawing a blue chip in the second draw?

1<sup>st</sup> draw: 5 ●, 3 ●, 2 ●

2<sup>nd</sup> draw: 5 ●, 2 ●, 2 ●

$$Prob(2^{\text{nd}} \text{ chip } B | 1^{\text{st}} \text{ chip } B) = \frac{2}{9} = 0.22$$

- If drawing without replacement, what is the probability of drawing two blue chips in a row?

1<sup>st</sup> draw: 5 ●, 3 ●, 2 ●

2<sup>nd</sup> draw: 5 ●, 2 ●, 2 ●

$$Prob(1^{\text{st}} \text{ chip } B) \cdot Prob(2^{\text{nd}} \text{ chip } B | 1^{\text{st}} \text{ chip } B) = 0.3 \times 0.22 \\ = 0.066$$

## Sampling without replacement (cont.)

- When drawing without replacement, the probability of the second chip being blue given the first was blue is not equal to the probability of drawing a blue chip in the first draw since the composition of the bag changes with the outcome of the first draw.

$$\text{Prob}(B|B) \neq \text{Prob}(B)$$

## Sampling without replacement (cont.)

- When drawing without replacement, the probability of the second chip being blue given the first was blue is not equal to the probability of drawing a blue chip in the first draw since the composition of the bag changes with the outcome of the first draw.

$$\text{Prob}(B|B) \neq \text{Prob}(B)$$

- When drawing without replacement, draws are not independent.*

## Sampling without replacement (cont.)

- When drawing without replacement, the probability of the second chip being blue given the first was blue is not equal to the probability of drawing a blue chip in the first draw since the composition of the bag changes with the outcome of the first draw.

$$\text{Prob}(B|B) \neq \text{Prob}(B)$$

- When drawing without replacement, draws are not independent.*
- This is especially important to take note of when the sample sizes are small. If we were dealing with, say, 10,000 chips in a (giant) bag, taking out one chip of any color would not have as big an impact on the probabilities in the second draw.

## Practice

In most card games cards are dealt without replacement. What is the probability of being dealt an ace and then a 3? Choose the closest answer.

- (a) 0.0045
- (b) 0.0059
- (c) 0.0060
- (d) 0.1553

## Practice

In most card games cards are dealt without replacement. What is the probability of being dealt an ace and then a 3? Choose the closest answer.

- (a) 0.0045
- (b) 0.0059
- (c) **0.0060**
- (d) 0.1553

$$P(\text{ace then 3}) = \frac{4}{52} \times \frac{4}{51} \approx 0.0060$$

## Random variables

---

# Random variables

- A *random variable* is a numeric quantity whose value depends on the outcome of a random event
  - We use a capital letter, like  $X$ , to denote a random variable
  - The values of a random variable are denoted with a lowercase letter, in this case  $x$
  - For example,  $P(X = x)$
- There are two types of random variables:
  - *Discrete random variables* often take only integer values
    - Example: Number of credit hours, Difference in number of credit hours this term vs last
  - *Continuous random variables* take real (decimal) values
    - Example: Cost of books this term, Difference in cost of books this term vs last

# Expectation

- We are often interested in the average outcome of a random variable.
- We call this the *expected value* (mean), and it is a weighted average of the possible outcomes

$$\mu = E(X) = \sum_{i=1}^k x_i P(X = x_i)$$

## Expected value of a discrete random variable

In a game of cards you win \$1 if you draw a heart, \$5 if you draw an ace (including the ace of hearts), \$10 if you draw the king of spades and nothing for any other card you draw. Write the probability model for your winnings, and calculate your expected winning.

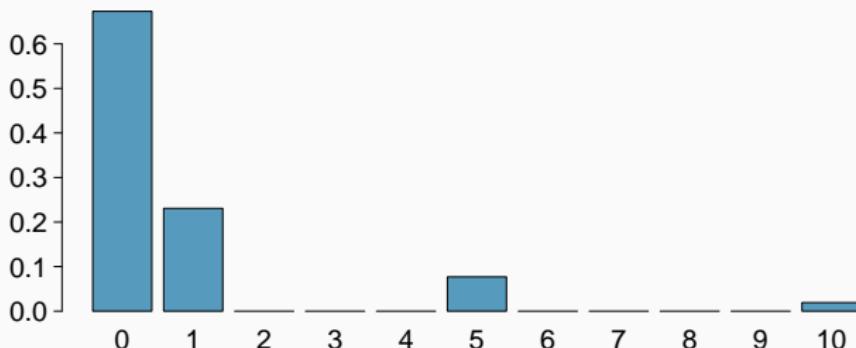
## Expected value of a discrete random variable

In a game of cards you win \$1 if you draw a heart, \$5 if you draw an ace (including the ace of hearts), \$10 if you draw the king of spades and nothing for any other card you draw. Write the probability model for your winnings, and calculate your expected winning.

Event	$X$	$P(X)$	$X P(X)$
Heart (not ace)	1	$\frac{12}{52}$	$\frac{12}{52}$
Ace	5	$\frac{4}{52}$	$\frac{20}{52}$
King of spades	10	$\frac{1}{52}$	$\frac{10}{52}$
All else	0	$\frac{35}{52}$	0
Total			$E(X) = \frac{42}{52} \approx 0.81$

## Expected value of a discrete random variable (cont.)

Below is a visual representation of the probability distribution of winnings from this game:



## Variability

We are also often interested in the variability in the values of a random variable.

$$\sigma^2 = \text{Var}(X) = \sum_{i=1}^k (x_i - E(X))^2 P(X = x_i)$$

$$\sigma = SD(X) = \sqrt{\text{Var}(X)}$$

## Variability of a discrete random variable

For the previous card game example, how much would you expect the winnings to vary from game to game?

## Variability of a discrete random variable

For the previous card game example, how much would you expect the winnings to vary from game to game?

$X$	$P(X)$	$X P(X)$	$(X - E(X))^2$	$P(X) (X - E(X))^2$
1	$\frac{12}{52}$	$1 \times \frac{12}{52} = \frac{12}{52}$	$(1 - 0.81)^2 = 0.0361$	$\frac{12}{52} \times 0.0361 = 0.0083$
5	$\frac{4}{52}$	$5 \times \frac{4}{52} = \frac{20}{52}$	$(5 - 0.81)^2 = 17.5561$	$\frac{4}{52} \times 17.5561 = 1.3505$
10	$\frac{1}{52}$	$10 \times \frac{1}{52} = \frac{10}{52}$	$(10 - 0.81)^2 = 84.4561$	$\frac{1}{52} \times 84.4561 = 1.6242$
0	$\frac{35}{52}$	$0 \times \frac{35}{52} = 0$	$(0 - 0.81)^2 = 0.6561$	$\frac{35}{52} \times 0.6561 = 0.4416$
		$E(X) = 0.81$		

## Variability of a discrete random variable

For the previous card game example, how much would you expect the winnings to vary from game to game?

$X$	$P(X)$	$X P(X)$	$(X - E(X))^2$	$P(X) (X - E(X))^2$
1	$\frac{12}{52}$	$1 \times \frac{12}{52} = \frac{12}{52}$	$(1 - 0.81)^2 = 0.0361$	$\frac{12}{52} \times 0.0361 = 0.0083$
5	$\frac{4}{52}$	$5 \times \frac{4}{52} = \frac{20}{52}$	$(5 - 0.81)^2 = 17.5561$	$\frac{4}{52} \times 17.5561 = 1.3505$
10	$\frac{1}{52}$	$10 \times \frac{1}{52} = \frac{10}{52}$	$(10 - 0.81)^2 = 84.4561$	$\frac{1}{52} \times 84.4561 = 1.6242$
0	$\frac{35}{52}$	$0 \times \frac{35}{52} = 0$	$(0 - 0.81)^2 = 0.6561$	$\frac{35}{52} \times 0.6561 = 0.4416$
		$E(X) = 0.81$		$V(X) = 3.4246$

## Variability of a discrete random variable

For the previous card game example, how much would you expect the winnings to vary from game to game?

$X$	$P(X)$	$X P(X)$	$(X - E(X))^2$	$P(X) (X - E(X))^2$
1	$\frac{12}{52}$	$1 \times \frac{12}{52} = \frac{12}{52}$	$(1 - 0.81)^2 = 0.0361$	$\frac{12}{52} \times 0.0361 = 0.0083$
5	$\frac{4}{52}$	$5 \times \frac{4}{52} = \frac{20}{52}$	$(5 - 0.81)^2 = 17.5561$	$\frac{4}{52} \times 17.5561 = 1.3505$
10	$\frac{1}{52}$	$10 \times \frac{1}{52} = \frac{10}{52}$	$(10 - 0.81)^2 = 84.4561$	$\frac{1}{52} \times 84.4561 = 1.6242$
0	$\frac{35}{52}$	$0 \times \frac{35}{52} = 0$	$(0 - 0.81)^2 = 0.6561$	$\frac{35}{52} \times 0.6561 = 0.4416$
		$E(X) = 0.81$		$V(X) = 3.4246$
				$SD(X) = \sqrt{3.4246} = 1.85$

## Linear combinations

- A *linear combination* of random variables  $X$  and  $Y$  is given by

$$aX + bY$$

where  $a$  and  $b$  are some fixed numbers.

## Linear combinations

- A *linear combination* of random variables  $X$  and  $Y$  is given by

$$aX + bY$$

where  $a$  and  $b$  are some fixed numbers.

- The average value of a linear combination of random variables is given by

$$E(aX + bY) = a \times E(X) + b \times E(Y)$$

## Calculating the expectation of a linear combination

On average you take 10 minutes for each statistics homework problem and 15 minutes for each chemistry homework problem. This week you have 5 statistics and 4 chemistry homework problems assigned. What is the total time you expect to spend on statistics and chemistry homework for the week?

## Calculating the expectation of a linear combination

On average you take 10 minutes for each statistics homework problem and 15 minutes for each chemistry homework problem. This week you have 5 statistics and 4 chemistry homework problems assigned. What is the total time you expect to spend on statistics and chemistry homework for the week?

$$\begin{aligned}E(S + S + S + S + S + C + C + C + C) &= 5 \times E(S) + 4 \times E(C) \\&= 5 \times 10 + 4 \times 15 \\&= 50 + 60 \\&= 110 \text{ min}\end{aligned}$$

## Linear combinations

- The variability of a linear combination of two independent random variables is calculated as

$$V(aX + bY) = a^2 \times V(X) + b^2 \times V(Y)$$

## Linear combinations

- The variability of a linear combination of two independent random variables is calculated as

$$V(aX + bY) = a^2 \times V(X) + b^2 \times V(Y)$$

- The standard deviation of the linear combination is the square root of the variance.

## Linear combinations

- The variability of a linear combination of two independent random variables is calculated as

$$V(aX + bY) = a^2 \times V(X) + b^2 \times V(Y)$$

- The standard deviation of the linear combination is the square root of the variance.

---

**Note:** If the random variables are not independent, the variance calculation gets a little more complicated and is beyond the scope of this course.

## Calculating the variance of a linear combination

The standard deviation of the time you take for each statistics homework problem is 1.5 minutes, and it is 2 minutes for each chemistry problem. What is the standard deviation of the time you expect to spend on statistics and physics homework for the week if you have 5 statistics and 4 chemistry homework problems assigned? Suppose that the time it takes to complete each problem is independent of another.

## Calculating the variance of a linear combination

The standard deviation of the time you take for each statistics homework problem is 1.5 minutes, and it is 2 minutes for each chemistry problem. What is the standard deviation of the time you expect to spend on statistics and physics homework for the week if you have 5 statistics and 4 chemistry homework problems assigned? Suppose that the time it takes to complete each problem is independent of another.

$$\begin{aligned}V(5S + 4C) &= 5 \times V(S) + 4 \times V(C) \\&= 5 \times 1.5^2 + 4 \times 2^2 \\&= 27.25\end{aligned}$$

## Practice

A casino game costs \$5 to play. If the first card you draw is red, then you get to draw a second card (without replacement). If the second card is the ace of clubs, you win \$500. If not, you don't win anything, i.e. lose your \$5. What is your expected profits/losses from playing this game? Remember: profit/loss = winnings - cost.

- (a) A profit of 5¢
- (c) A loss of 25¢
- (b) A loss of 10¢
- (d) A loss of 30¢

## Practice

A casino game costs \$5 to play. If the first card you draw is red, then you get to draw a second card (without replacement). If the second card is the ace of clubs, you win \$500. If not, you don't win anything, i.e. lose your \$5. What is your expected profits/losses from playing this game? Remember: profit/loss = winnings - cost.

- (a) A profit of 5¢
- (c) A loss of 25¢
- (b) **A loss of 10¢**
- (d) A loss of 30¢

Event	Win	Profit: $X$	$P(X)$	$X \times P(X)$
Red, A♣	500	$500 - 5 = 495$	$\frac{26}{52} \times \frac{1}{51} = 0.0098$	$495 \times 0.0098 = 4.851$
Other	0	$0 - 5 = -5$	$1 - 0.0098 = 0.9902$	$-5 \times 0.9902 = -4.951$
				$E(X) = -0.1$

## Fair game

A *fair* game is defined as a game that costs as much as its expected payout, i.e. expected profit is 0.

## Fair game

A *fair* game is defined as a game that costs as much as its expected payout, i.e. expected profit is 0.

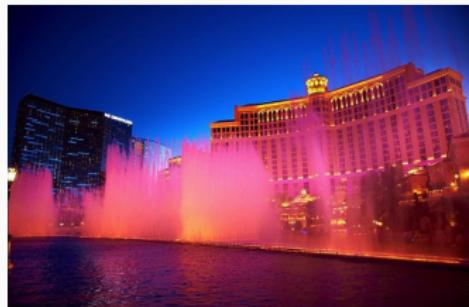
Do you think casino games in Vegas cost more or less than their expected payouts?

## Fair game

A *fair* game is defined as a game that costs as much as its expected payout, i.e. expected profit is 0.

Do you think casino games in Vegas cost more or less than their expected payouts?

*If those games cost less than their expected payouts, it would mean that the casinos would be losing money on average, and hence they wouldn't be able to pay for all this:*



*Image by Moyan\_Brenn on Flickr [http://www.flickr.com/photos/aigle\\_dore/5951714693](http://www.flickr.com/photos/aigle_dore/5951714693).*

## Simplifying random variables

Random variables do not work like normal algebraic variables:

$$X + X \neq 2X$$

## Simplifying random variables

Random variables do not work like normal algebraic variables:

$$X + X \neq 2X$$

$$\begin{aligned} E(X + X) &= E(X) + E(X) & Var(X + X) &= Var(X) + Var(X) \text{ (assuming independence)} \\ &= 2E(X) & &= 2 Var(X) \end{aligned}$$

$$\begin{aligned} E(2X) &= 2E(X) & Var(2X) &= 2^2 Var(X) \\ & & &= 4 Var(X) \end{aligned}$$

## Simplifying random variables

Random variables do not work like normal algebraic variables:

$$X + X \neq 2X$$

$$\begin{aligned} E(X + X) &= E(X) + E(X) & Var(X + X) &= Var(X) + Var(X) \text{ (assuming independence)} \\ &= 2E(X) & &= 2 Var(X) \end{aligned}$$

$$\begin{aligned} E(2X) &= 2E(X) & Var(2X) &= 2^2 Var(X) \\ & & &= 4 Var(X) \end{aligned}$$

$E(X + X) = E(2X)$ , but  $Var(X + X) \neq Var(2X)$ .

## Adding or multiplying?

A company has 5 Lincoln Town Cars in its fleet. Historical data show that annual maintenance cost for each car is on average \$2,154 with a standard deviation of \$132. What is the mean and the standard deviation of the total annual maintenance cost for this fleet?

## Adding or multiplying?

A company has 5 Lincoln Town Cars in its fleet. Historical data show that annual maintenance cost for each car is on average \$2,154 with a standard deviation of \$132. What is the mean and the standard deviation of the total annual maintenance cost for this fleet?

Note that we have 5 cars each with the given annual maintenance cost ( $X_1 + X_2 + X_3 + X_4 + X_5$ ), not one car that had 5 times the given annual maintenance cost ( $5X$ ).

## Adding or multiplying?

A company has 5 Lincoln Town Cars in its fleet. Historical data show that annual maintenance cost for each car is on average \$2,154 with a standard deviation of \$132. What is the mean and the standard deviation of the total annual maintenance cost for this fleet?

Note that we have 5 cars each with the given annual maintenance cost ( $X_1 + X_2 + X_3 + X_4 + X_5$ ), not one car that had 5 times the given annual maintenance cost ( $5X$ ).

$$E(X_1 + X_2 + X_3 + X_4 + X_5) = E(X_1) + E(X_2) + E(X_3) + E(X_4) + E(X_5)$$

## Adding or multiplying?

A company has 5 Lincoln Town Cars in its fleet. Historical data show that annual maintenance cost for each car is on average \$2,154 with a standard deviation of \$132. What is the mean and the standard deviation of the total annual maintenance cost for this fleet?

Note that we have 5 cars each with the given annual maintenance cost ( $X_1 + X_2 + X_3 + X_4 + X_5$ ), not one car that had 5 times the given annual maintenance cost ( $5X$ ).

$$\begin{aligned}E(X_1 + X_2 + X_3 + X_4 + X_5) &= E(X_1) + E(X_2) + E(X_3) + E(X_4) + E(X_5) \\&= 5 \times E(X) = 5 \times 2,154 = \$10,770\end{aligned}$$

## Adding or multiplying?

A company has 5 Lincoln Town Cars in its fleet. Historical data show that annual maintenance cost for each car is on average \$2,154 with a standard deviation of \$132. What is the mean and the standard deviation of the total annual maintenance cost for this fleet?

Note that we have 5 cars each with the given annual maintenance cost ( $X_1 + X_2 + X_3 + X_4 + X_5$ ), not one car that had 5 times the given annual maintenance cost ( $5X$ ).

$$\begin{aligned}E(X_1 + X_2 + X_3 + X_4 + X_5) &= E(X_1) + E(X_2) + E(X_3) + E(X_4) + E(X_5) \\&= 5 \times E(X) = 5 \times 2,154 = \$10,770\end{aligned}$$

$$Var(X_1 + X_2 + X_3 + X_4 + X_5) = Var(X_1) + Var(X_2) + Var(X_3) + Var(X_4) + Var(X_5)$$

## Adding or multiplying?

A company has 5 Lincoln Town Cars in its fleet. Historical data show that annual maintenance cost for each car is on average \$2,154 with a standard deviation of \$132. What is the mean and the standard deviation of the total annual maintenance cost for this fleet?

Note that we have 5 cars each with the given annual maintenance cost ( $X_1 + X_2 + X_3 + X_4 + X_5$ ), not one car that had 5 times the given annual maintenance cost ( $5X$ ).

$$\begin{aligned} E(X_1 + X_2 + X_3 + X_4 + X_5) &= E(X_1) + E(X_2) + E(X_3) + E(X_4) + E(X_5) \\ &= 5 \times E(X) = 5 \times 2,154 = \$10,770 \end{aligned}$$

$$\begin{aligned} Var(X_1 + X_2 + X_3 + X_4 + X_5) &= Var(X_1) + Var(X_2) + Var(X_3) + Var(X_4) + Var(X_5) \\ &= 5 \times V(X) = 5 \times 132^2 = 87,120 \end{aligned}$$

## Adding or multiplying?

A company has 5 Lincoln Town Cars in its fleet. Historical data show that annual maintenance cost for each car is on average \$2,154 with a standard deviation of \$132. What is the mean and the standard deviation of the total annual maintenance cost for this fleet?

Note that we have 5 cars each with the given annual maintenance cost ( $X_1 + X_2 + X_3 + X_4 + X_5$ ), not one car that had 5 times the given annual maintenance cost ( $5X$ ).

$$\begin{aligned} E(X_1 + X_2 + X_3 + X_4 + X_5) &= E(X_1) + E(X_2) + E(X_3) + E(X_4) + E(X_5) \\ &= 5 \times E(X) = 5 \times 2,154 = \$10,770 \end{aligned}$$

$$\begin{aligned} Var(X_1 + X_2 + X_3 + X_4 + X_5) &= Var(X_1) + Var(X_2) + Var(X_3) + Var(X_4) + Var(X_5) \\ &= 5 \times V(X) = 5 \times 132^2 = 87,120 \end{aligned}$$

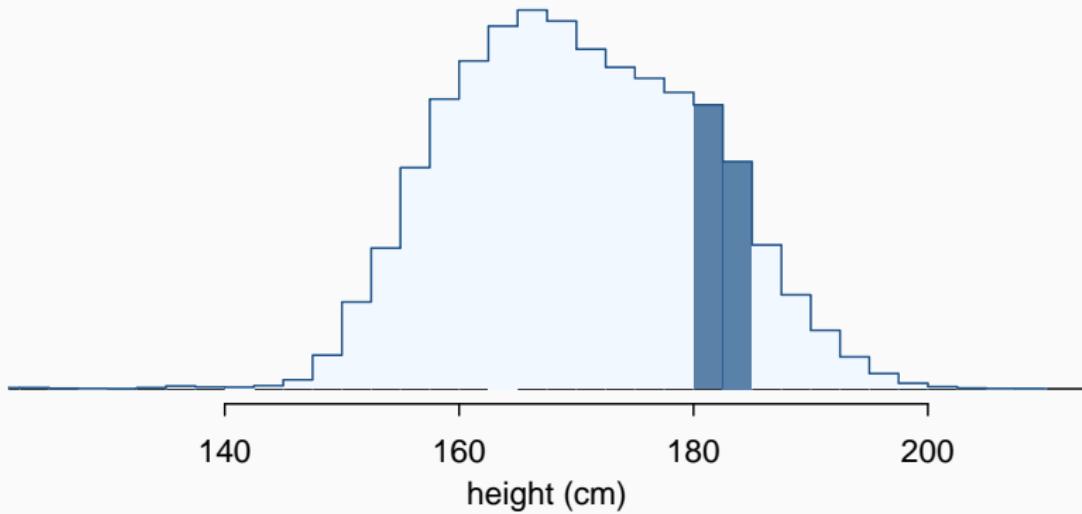
$$SD(X_1 + X_2 + X_3 + X_4 + X_5) = \sqrt{87,120} = \$295.16$$

## **Continuous distributions**

---

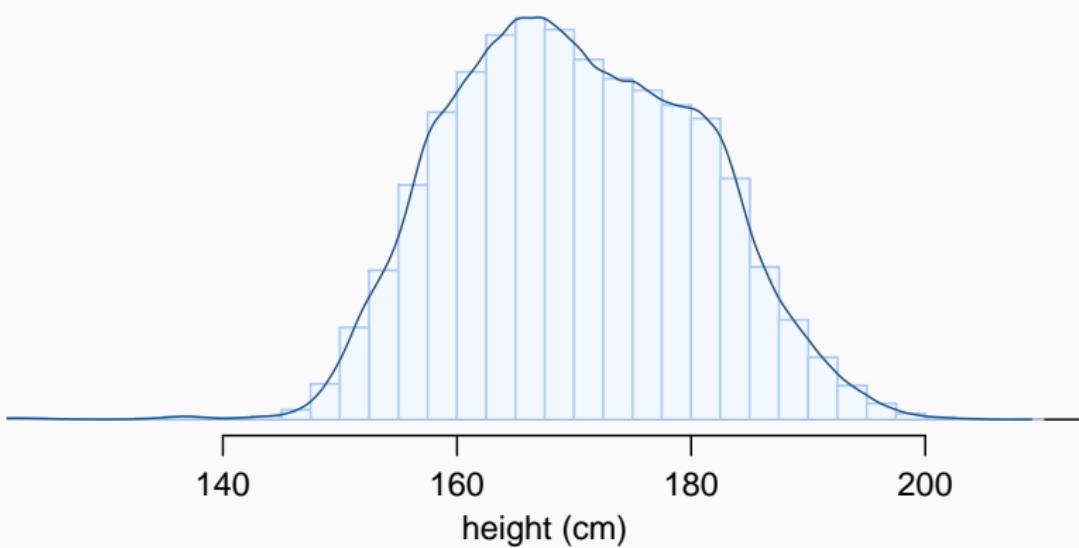
## Continuous distributions

- Below is a histogram of the distribution of heights of US adults.
- The proportion of data that falls in the shaded bins gives the probability that a randomly sampled US adult is between 180 cm and 185 cm (about 5'11" to 6'1").



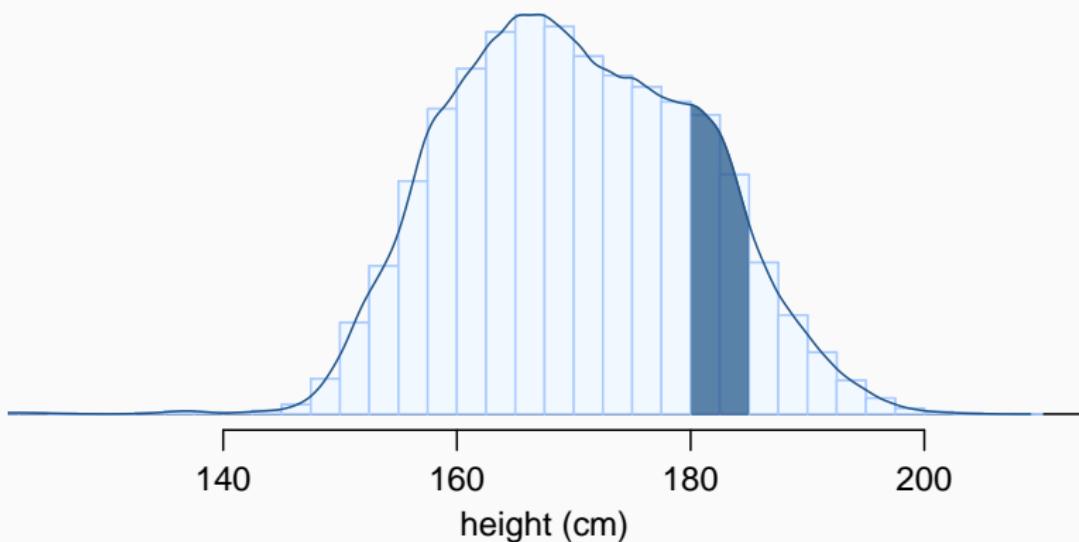
## From histograms to continuous distributions

Since height is a continuous numerical variable, its *probability density function* is a smooth curve.



## Probabilities from continuous distributions

Therefore, the probability that a randomly sampled US adult is between 180 cm and 185 cm can also be estimated as the shaded area under the curve.



## By definition...

Since continuous probabilities are estimated as “the area under the curve”, the probability of a person being exactly 180 cm (or any exact value) is defined as 0.

