







Introduction to Probability

Statistics is everywhere

Berkeley, CA 10 Day Weather

7:54 am PST [Print](#)

| DAY | | DESCRIPTION | HIGH / LOW | PRECIP | WIND | HUMIDITY |
|-----------------|---|---------------|------------|--------|-----------|----------|
| TODAY FEB 13 |  | Partly Cloudy | 56°/45° | 10% | SW 8 mph | 78% |
| FRI FEB 14 |  | Mostly Sunny | 62°/45° | 10% | WSW 5 mph | 75% |
| SAT FEB 15 |  | Partly Cloudy | 64°/48° | 10% | SW 7 mph | 71% |
| SUN FEB 16 |  | Partly Cloudy | 62°/47° | 10% | W 11 mph | 71% |
| MON FEB 17 |  | Sunny | 66°/45° | 0% | N 16 mph | 36% |
| TUE FEB 18 |  | Sunny | 64°/44° | 0% | N 11 mph | 34% |

Different definitions

According to the U.S. National Weather Service (NWS), PoP is the probability of exceedance that more than 0.01 inches (0.25 mm) of precipitation will fall in a single spot, averaged over the forecast area. This can be expressed mathematically as

$$PoP = C \times A$$

where C is the confidence that any form of precipitation (e.g., snow or rain) will occur somewhere in the forecast area and A is the percent of the area that will receive measurable precipitation, if it occurs at all. For instance, if there is a 100% probability of rain covering one half of a city, and a 0% probability of rain on the other half of the city, the POP for the city would be 50%. A 50% chance of a rainstorm covering the entire city would also lead to a PoP of 50%.

The PoP measure is meaningless unless it is associated with a period of time. NWS forecasts commonly use PoP defined over 12-hour periods (PoP12), though 6-hour periods (PoP6) and other measures are also published. A “daytime” PoP12 means from 6 am to 6 pm.

Different definitions

Environment Canada reports a chance of precipitation (COP) that is defined as “The chance that measurable precipitation (0.2 mm of rain or 0.2 cm of snow) will fall on any random point of the forecast region during the forecast period.” The values are rounded to 10% increments, but are never rounded to 50%.

Learning objectives for today

- Why does probability matter?
- Terminology for probability: sample space (discrete vs. continuous), event, discrete probability model, continuous probability model
- Some basic rules for probability
- Definition of a random variable
- Probabilities in Public Health: risk, odds, case fatality rate

Why is probability important:

Although we don't always explicitly discuss probability, we think about probability all the time:

- Probability of rain today
- Probability of getting a job with a college degree
- Probability of a question showing up on the next exam

Statistics can be misleading, knowing the basic rules of probability and understanding how they are generated can help you to interpret statistics clearly, think critically about information and draw relevant conclusions.

Why is probability important:

- Misleading statistics can be used to make and defend policy decisions You should know how to interpret those statistics for yourself
- Determining probability can be difficult and not intuitive. Our gut instinct about the probability of an event may be way off and lead to poor decision making in policy, medical or personal settings.
- Using predictive models to calculate probabilities sounds like an objective process, however these models can encode bias and discrimination.

Example. 1: Misleading statistics

- Kirstjen Nielsen, former Homeland Security Secretary stated, about folks at the US-Mexico border:
“Again, let’s just pause to think about this statistic: **314 percent** increase in adults showing up with kids that are not a family unit. . . Those are traffickers, those are smugglers, that is MS-13, those are criminals, those are abusers.”
- Nielsen was speaking about a relative increase in the probability of the event of “adults with kids who are not their own at the US-Mexico border”. The relative increase is very large (314%). However, how often did the event happen in the first place?
- In a Washington Post analysis¹, the increase was from 0.19% in 2017 to 0.61% in 2018. Thus the actual chance of the event happening is very small and increased by $0.61\% - 0.19\% = 0.42$ percentage points.
- Takeaway: looking at the increase in absolute percentage points provides a different interpretation than the increase on the relative scale.

Reference: https://www.washingtonpost.com/news/politics/wp/2018/06/18/how-to-mislead-with-statistics-dhs-secretary-nielsen-edition/?noredirect=on&utm_term=.9193534ee80c

Example 2: Calculating probabilities in medical settings can be difficult and not intuitive

- Suppose that there is test for a specific type of cancer that has a 90% chance of a positive screening test result for cancer if the individual truly has cancer and a 90% chance of testing negative for cancer when the individual does not have it.
 - 1% of patients in the population have the cancer being tested for.
 - What is the chance that a patient has cancer given that they test positive?
- a) Between 0% - 24.9%
 - b) Between 25.0% - 49.9%
 - c) Between 50.0% - 74.9%
 - d) Between 75.0% - 100%

Example 2: Calculating probabilities in medical settings can be difficult and not intuitive

Many people choose 4), but the true answer is 1)! Why do we get this so wrong?

Video link (2 mins): [click here](#)

What is probability?

Fundamental components

Probability can be thought of in three fundamental components:

- A random experiment
- All possible outcomes of that experiment
- An **event** or events of interest

Sample space

We refer to the entire set of possible outcomes as the **Sample Space**

All possible outcomes of a sample space (S) together have a probability of 1. $P(S) = 1$

Sample space

- **Discrete sample space**
 - e.g., Marital status: $S = \{\text{Single, married, divorced, widowed}\}$
 - careful, discrete spaces remind us of nominal, ordinal, or discrete variables
 - anything that is countable with “gaps” between the events
 - the notation is important: $S = \{\text{elements in the space}\}$
- **Continuous sample space**
 - e.g., The interval $[0, 1]$: $S = \{\text{all numbers between 0 and 1}\}$
 - continuous sample spaces remind us of continuous variables only
 - the events are not countable (i.e., we cannot list the numbers between 0 and 1, there are infinite)

Defining probability

The probability or occurrence of an event A often called the probability of A and denoted as $P(A)$ is the ratio of the number of outcomes where event A occurs to the total number of possible outcomes.

For a coin what is the probability of heads?

- **Probability model:** Description of random phenomena. Consists of sample space S and a way of assigning probabilities to events

Frequentist definition

- Probability corresponds to frequency over many repetitions
- For example, the probability of a coin landing heads on a single toss is 0.5: if we toss a coin numerous times, we would expect one-half of the tosses to land heads, the more times we toss the coin the closer we expect the fraction of tosses to come to exactly half.

From B&M: How common is the common cold?

- Suppose that there are 100,000 people in your community. You work for your community’s public health office and want to estimate the number of people who had a common cold. You are not able to sample everyone but suppose you could randomly call people in your community and ask them “Did you have a cold yesterday?” and then calculate the proportion of the sample who had a cold.

From B&M: How common is the common cold?

- Here are the dimensions, a data frame for the **whole population**, and the mean of the variable `had_cold_yesterday`:

```
dim(cold_data)
```

```
## [1] 100000      2
```

```
cold_data %>% summarize(population_mean = mean(had_cold_yesterday))
```

```
## population_mean
## 1 0.11214
```

- Note that the mean of a 0/1 variable is called a **proportion**. This is because the mean is the number of individuals with a cold (coded as `had_cold_yesterday = 1`) divided by the total number of individuals.

From B&M page 216: How common is the common cold?

How big should your sample size be?

- We want to sample enough people such that the proportion of those with colds in the sample is close to the proportion of those with colds in the population
- Let's take samples of size 5, 100, and so on using `dplyr`'s `sample_n` function:

```
sample_5 <- dplyr::sample_n(tbl = cold_data, size = 5)
sample_100 <- dplyr::sample_n(tbl = cold_data, size = 100)
sample_1000 <- dplyr::sample_n(tbl = cold_data, size = 1000)
sample_10000 <- dplyr::sample_n(tbl = cold_data, size = 10000)
sample_100000 <- dplyr::sample_n(tbl = cold_data, size = 100000)
```

Now estimate the proportion of those with a cold in the random samples

```
## sample_mean_n5
## 1 0
## sample_mean_n100
## 1 0.08
## sample_mean_n1k
## 1 0.112
## sample_mean_n10k
## 1 0.1091
## sample_mean_n100k
## 1 0.11214
```

Estimate the proportion of those with a cold in the random samples

- What do you notice about the proportion estimates?
- Do they approach the true estimate as the sample size increases?

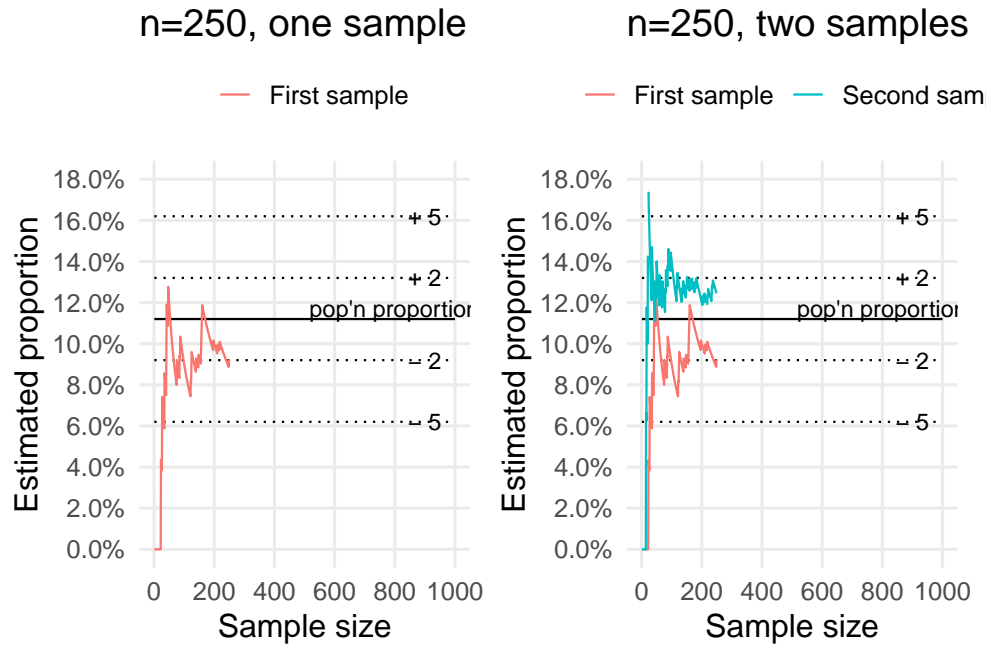
How many people should we sample?

- So we know we should sample more than five people, but feasibly can't sample everyone.
- We need to sample *enough* people to reasonably estimate the true chance of having a cold. But how many is enough?

Look at the sample's estimated proportion as a function of sample size

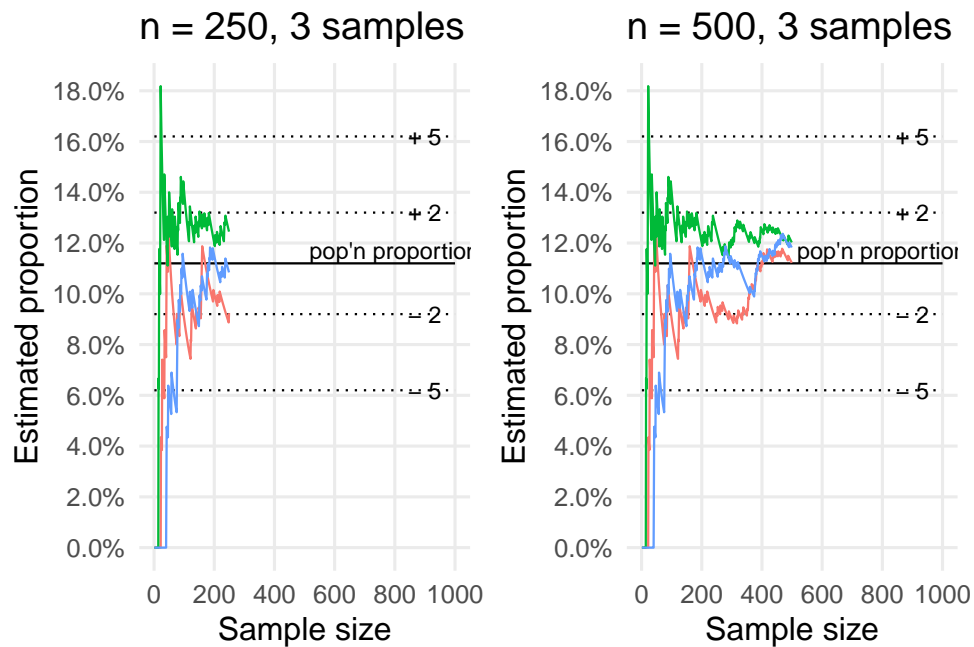
- To help us decide I first sampled one person and took the mean of that sample.
- Then I added another person and took the mean of that sample of size 2 ($n=2$). ... and so on, until I had 5000 people.
- The plot on the next slide shows the estimated proportion vs. the sample size.

Estimated proportion vs. sample size for $n = 250$



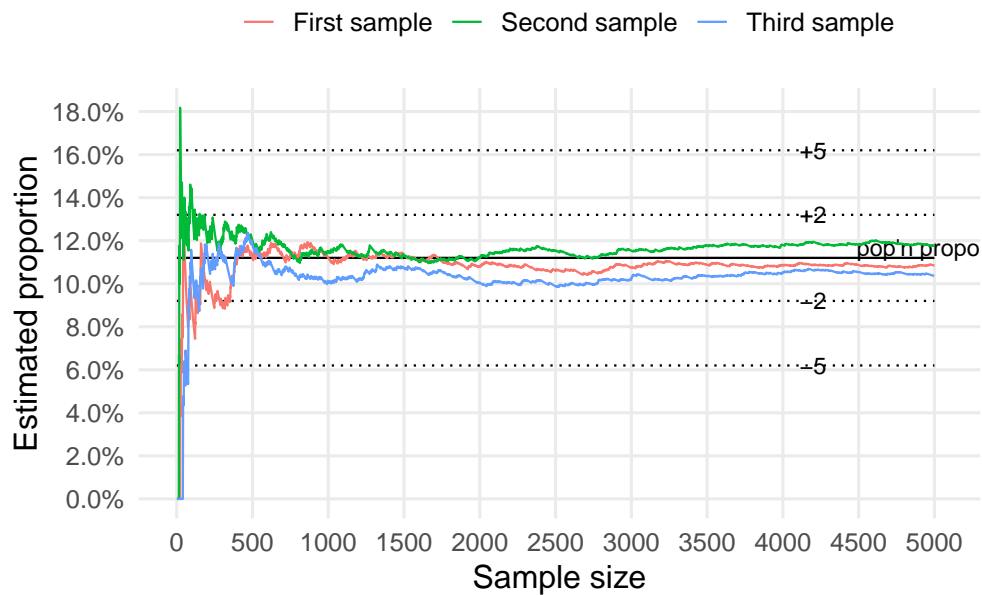
Estimated proportion vs. sample size for $n = 500$

- Increase the sample size how the estimate becomes closer to the true value
- Add in a third sample to compare how different samples perform in the short vs. the long run



Estimated proportion vs. sample size for $n = 5000$

$n = 5000$, 3 samples



Summary of the example on estimating the proportion with a cold

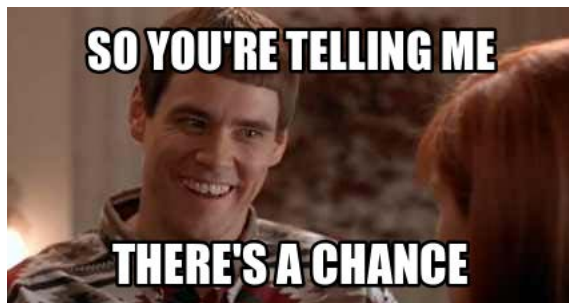
- As the sample size increases, the estimated proportion becomes closer to the true proportion
- Random samples of the same size will provide different estimates of the true proportion, but will be closer to each other (and the true value) if they are “large enough”

Some vocabulary, definitions, and rules

Rule of Range

Probabilities are numbers between 0 and 1.

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



What is the probability of a certain event? Of an impossible event?

Complement

The total number of outcomes in a random experiment (the sample space) can always divide into two mutually exclusive groups:

outcomes where A occurs $P(A)$

outcomes where A does not occur $P(\bar{A})$

These are called **complementary events** These must cover the entire set of possibilities which sums to 1

Thus $P(A) + P(\bar{A}) = 1$

Complement

The probability of the complement is 1 minus the probability of the event occurring.

- $P(A \text{ does not occur}) = 1 - P(A)$
- Shorthand: $P(\bar{A}) = 1 - P(A)$ or $P(A^c) = 1 - P(A)$ or $P(A') = 1 - P(A)$

Composite Events

We started by defining probability in terms of one event, but we can expand this to think about the probability of more than one event. For example, let A and B be two separate events. A **composite event** would then be the event which describes the outcomes of both A and B.

The composite event where both A and B occur is also referred to as the **intersect** of A and B.

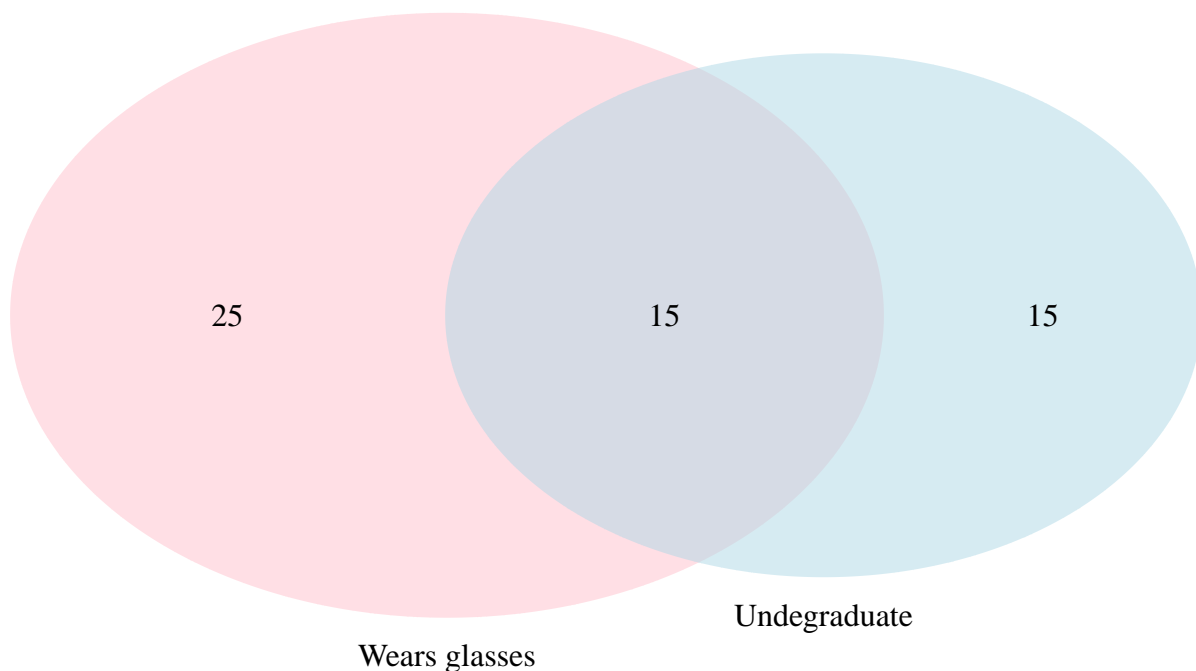
This is sometimes written as $P(AB)$ or $P(A \cap B)$

The composite event where either A or B occurs is also referred to as the **union** of A and B.

This is sometimes written as $P(A \text{ or } B)$ or $P(A \cup B)$

Composite Events

Imagine we have 100 students in a classroom.



(polygon[GRID.polygon.286], polygon[GRID.polygon.287], polygon[GRID.polygon.288], polygon[GRID.polygon.289])

Composite Events

There are 40 students who wear glasses, and 30 who are undergraduates.

Based on the Venn diagram, what is $P(\text{Glasses})$?

What is the complement $P(\overline{\text{Glasses}})$

What is the union of these two events $P(\text{Glasses} \cup \text{Undergraduate})$

What is the intersect of these two events $P(\text{Glasses} \cap \text{Undergraduate})$

Composite Events

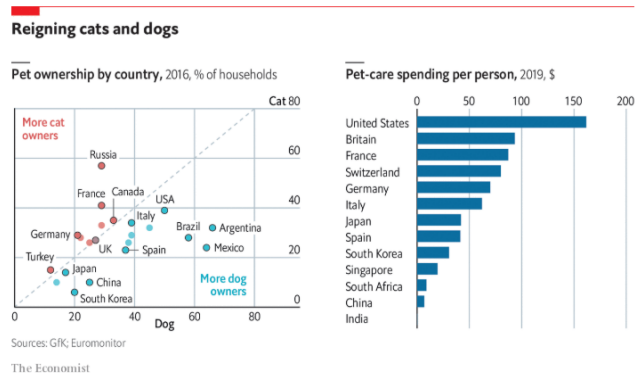
If we toss a coin twice, what are the possible composite events in our sample space?

What is the probability of tossing a combination of one Heads and one Tails?

What is the complement of that?

Probability of pet ownership

What is the probability space here?



from the economist article here

Disjoint events

If two events have a joint probability of 0 (i.e., no overlap in their event spaces $P(A \cap B) = 0$) then they are **disjoint** and the probability of either event occurring is the summation of their individual probabilities.

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

Disjoint events are also described as **mutually exclusive** meaning that it is not possible to have both events

Discrete probability models

Discrete probability model

- A probability model with a sample space made up of a list of individual outcomes is called discrete
- To assign probabilities in a discrete model, list the probabilities of all the individual outcomes. These probabilities must be numbers between 0 and 1 and must sum to 1. The probability of any event is the sum of the probabilities of the outcomes making up the event.

Discrete probability model example

For example, we could survey a sample of people and ask them their marital status. Based on this survey we can calculate the portion of each event in the sample space:

| Single | Married | Divorced | Widowed |
|--------|---------|----------|---------|
| 47% | 30% | 18% | 5% |

This is a discrete probability model shown in a table. How else could you display these data?

Continuous probability model

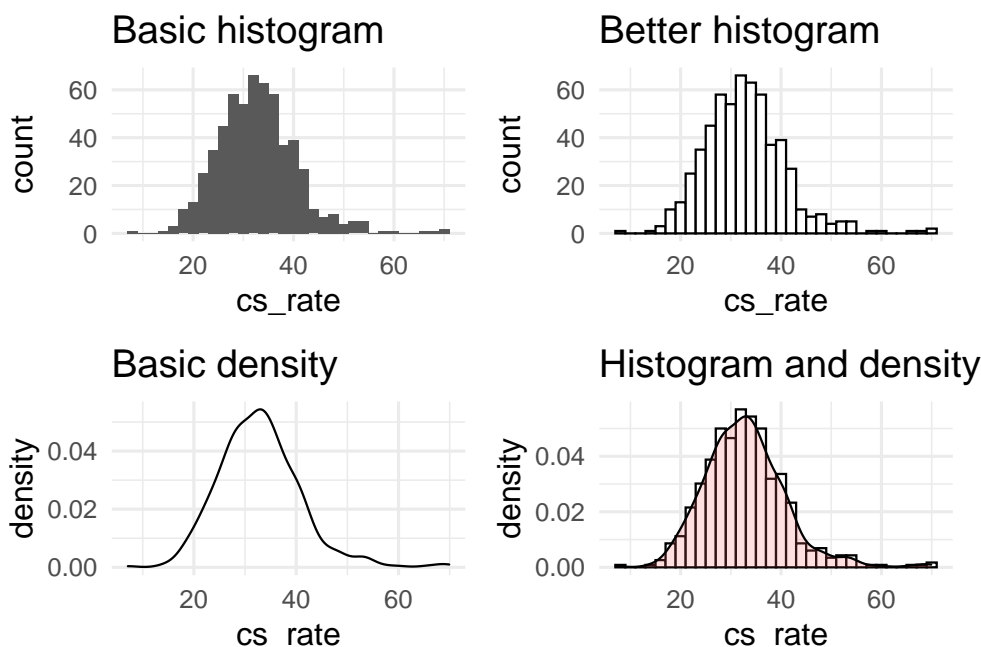
- A continuous probability model assigns probabilities as areas under a **density** curve. The area under the curve and between a range of specified values on the horizontal axis is the probability of an outcome in that range.
- What is a density curve?

Density curves

- Density curves are also known as **probability density functions**
- You can think of density curves as **smoothed histograms**.

Density curves using `geom_density()`

- Recall the data on cesarean delivery rates across hospitals in the US. We can use these data to also make a density plot (also called density curve):



Density curves

- From this plot, we can see that the density curve approximates the shape of the histogram very well.
- Remember because there are infinitely many cesarean delivery rates that could be observed between 0 and 1 it is impossible to assign a finite probability to any specific number.
- If we did, we could do this infinitely and their summed probability would surpass 100%.
- Instead, the density curve is used to determine the probability of an observed event within a specific range.

Density curves

You could use the density curve to calculate:

- $P(CS < 0.20)$
- $P(0.20 < CS < 0.40)$
- $P(CS < 0.2 \text{ or } CS > 0.4) = P(CS < 0.2) + P(CS > 0.4)$ because these events are independent
- $P(CS > 0.4) = 1 - P(CS \leq 0.4)$
- The calculations can be interpreted as either:
 - the **proportion** of hospitals with cesarean delivery rates in the specified range
 - the **probability** that a randomly chosen hospital will have cesarean delivery rate in the specified range.

Random variables

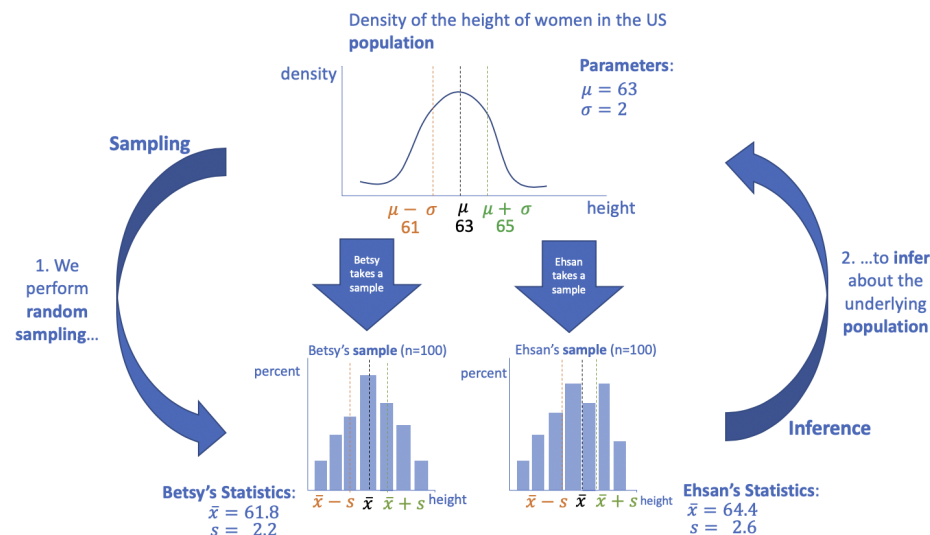
- A random variable is a variable whose value is a numerical outcome of a random phenomenon
- Random variables are represented by capital letters, most popularly X .
- A lower case letter represents a particular value for the random variable has been taken. For example $P(X = x)$ asks, what is the probability that random variable X takes the value x ?
- For continuous random variables, we ask $P(X < x)$ for example, because $P(X = x) = 0$ for continuous random variables.
 - Why is this the case?

The mean and standard deviation revisited

- In Chapter 2, we learned about how to calculate the mean (\bar{x}) and standard deviation (s) of a sample.
- We can also calculate the mean (μ) and standard deviation (σ) of a population.
- The mean and standard deviation of a population are represented using different notation to remind us that we are describing a population *parameter* vs. the *sample* mean and *sample* standard deviation that are *statistics* used to describe samples.

Putting it all together

The figure illustrates the difference between a (hypothetical) underlying distribution of heights among women in the US and its mean and s.d. vs. that of the sampled distribution.



Linking probabilities to public health

Risk

- Generally speaking, **risk** is another word for the probability or chance of an event occurring. In epidemiology and public health, we often use the word risk to represent the risk of some adverse health outcome among a group of individuals.
- For example (according to the American Cancer Society) the risk of developing cancer at some point in your life is roughly 1 in 3

Other risk definitions

- **Attack rate** of a virus: the risk (probability) of becoming afflicted during an infectious period, like the flu season. If the attack rate of influenza during a specific flu season was 10% or 10 per 100, this would imply that 10 out of every 100 individuals develop influenza during the epidemic period.
- **Case fatality rate**: the risk (probability) of dying among individuals with a specified condition. For example, the case fatality rate of measles is 1.5 per 1000 cases or 0.15%. (Examples from Epidemiology an Introduction by KJ Rothman, 2002:p.28)
- Note that these definitions use the word “rate” but not in the same way that we usually define rates in epidemiology. The attack rate and case fatality rate are risks, not rates!


Case fatality rate

News regarding case fatality due to COVID in February 2020,

Article link

Case fatality rate

More recently from data published in JAMA in late March

 JAMA Network™

From: **Case-Fatality Rate and Characteristics of Patients Dying in Relation to COVID-19 in Italy**

JAMA. 2020;323(18):1775-1776. doi:10.1001/jama.2020.4683

| | Italy as of March 17, 2020 | | China as of February 11, 2020 | |
|---------------|----------------------------|------------------------------------|-------------------------------|------------------------------------|
| | No. of deaths (% of total) | Case-fatality rate, % ^b | No. of deaths (% of total) | Case-fatality rate, % ^b |
| All | 1625 (100) | 7.2 | 1023 (100) | 2.3 |
| Age groups, y | | | | |
| 0-9 | 0 | 0 | 0 | 0 |
| 10-19 | 0 | 0 | 1 (0.1) | 0.2 |
| 20-29 | 0 | 0 | 7 (0.7) | 0.2 |
| 30-39 | 4 (0.3) | 0.3 | 18 (1.8) | 0.2 |
| 40-49 | 10 (0.6) | 0.4 | 38 (3.7) | 0.4 |
| 50-59 | 43 (2.7) | 1.0 | 130 (12.7) | 1.3 |
| 60-69 | 139 (8.6) | 3.5 | 309 (30.2) | 3.6 |
| 70-79 | 578 (35.6) | 12.8 | 312 (30.5) | 8.0 |
| ≥80 | 850 (52.3) | 20.2 | 208 (20.3) | 14.8 |

^a Data from China are from Chinese Center for Disease Control and Prevention.⁴ Age was not available for 1 patient.

^b Case-fatality rate calculated as number of deaths/number of cases.

Table Title:
Case-Fatality Rate by Age Group in Italy and China^a

Date of download: 7/16/2020

Copyright 2020 American Medical Association.
All Rights Reserved.

Odds

The **odds** is another commonly used measure in epidemiology that is a ratio of the probability of the adverse event over the probability of adverse event not occurring.

Sometimes, the popular press uses the word “odds” when we would use the word risk or probability.

Comparing risks and odds

In public health and epidemiology we are often presenting risks and odds not as overall (marginal) probabilities, but also as probabilities among groups (conditional probabilities) compared to each other. Next lecture we will talk more about conditional probabilities.

Comic Relief

From xkcd.com

