

Social Science Inquiry II

Week 3: A brief introduction to probability, part II

Molly Offer-Westort

Department of Political Science,
University of Chicago

Winter 2022

Loading packages for this class

```
> library(ggplot2)
> library(gridExtra)
> set.seed(60637)
```

Recall our terms from probability.

Our random process: flipping a fair coin twice.

- ▶ Ω : Sample space. Describes all possible outcomes in our setting.
 - ▶ ω : Generic notation for the realized outcomes in the sample space.
 - ▶ Here, $\Omega = \{HH, HT, TH, TT\}$.
- ▶ Event: a subset of Ω .
 - ▶ We will often use terms like A or B to define events.
 - ▶ Here, the event that we get a head on first flip is $A = \{HT, HH\}$.
- ▶ S : Event space. Describes all subsets of events, including null set.

Full event space

 - ▶ We use this in addition to the sample space, so we can describe all types of events that we can define the probability for.
- ▶ P : Probability measure. An operator that assigns probability to all events in the event space.
 - ▶ Here, the event that we get a head on the first flip, $P[A] = 1/2$.

Random variables

- ▶ A random variable is a mapping X from our sample space Ω , to the Real numbers.

$$X : \Omega \rightarrow \mathbb{R}$$

- ▶ Random variables are ways to quantify random events described by our sample space.
- ▶ We'll mostly work with random variables going forward, but it's important to remember that the random variable is built on the foundations of the sample space – and often, **you'll be the one deciding how that quantification happens.**

For example, with our two coin flips, let $X(\omega)$ be the number of heads in the sequence ω .

Then the random variable, and its probability distribution, can be described as:

ω	$P[\{\omega\}]$	$X(\omega)$
TT	1/4	0
TH	1/4	1
HT	1/4	1
HH	1/4	2

and,

x	$P[X = x]$
0	1/4
1	1/2
2	1/4

We can simulate this in 'R' as well.

```
> X <- c(0, 1, 2)
> probs <- c(0.25, 0.5, 0.25)
> sample(x = X,
+       size = 1,
+       prob = probs)

[1] 0

>

> n <- 1000
> result_n <- sample(x = X,
+                   size = n,
+                   prob = probs,
+                   replace = TRUE)
> table(result_n)

result_n
 0    1    2
247 499 254

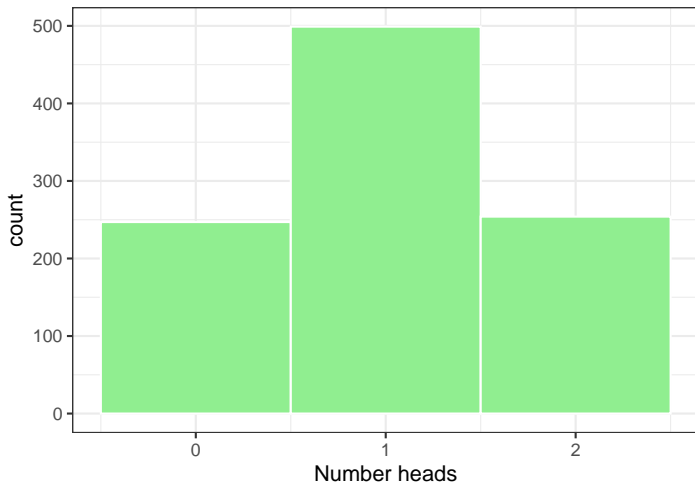
> prop.table(table(result_n))

result_n
 0    1    2
0.247 0.499 0.254

>
```

We can plot a histogram to look at the distribution of results.

```
> ggplot(data.frame(result_n), aes(x = result_n)) +  
+   geom_histogram(bins = 3, color = 'white', fill = 'lightgreen') +  
+   theme_bw() + xlab('Number heads')  
>
```



Probability Mass Function of a discrete random variable

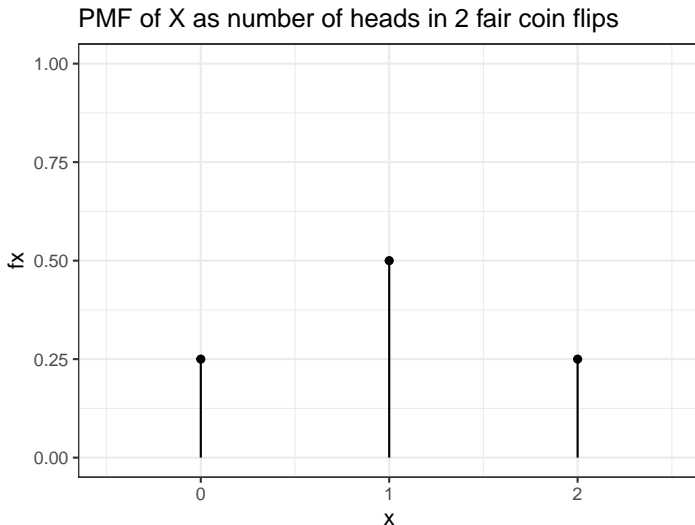
- ▶ A random variable is *discrete* if it takes countably many values.
- ▶ The probability mass function of a discrete RV X tells us the probability we will see an outcome at some value x .

$$f(x) = P[X = x]$$

For our coin flip example,

$$f(x) = \begin{cases} 1/4 & x = 0 \\ 1/2 & x = 1 \\ 1/4 & x = 2 \\ 0 & \text{otherwise} \end{cases}$$

Illustrating the PMF of a discrete RV



Note that the probabilities sum to 1. This is one of the foundational axioms of probability.
Social Science Inquiry II, Winter 2022

Molly Offer-Westort

Cumulative Distribution Functions

The cumulative distribution function of X tells us the probability we will see an outcome less than or equal to some value x .

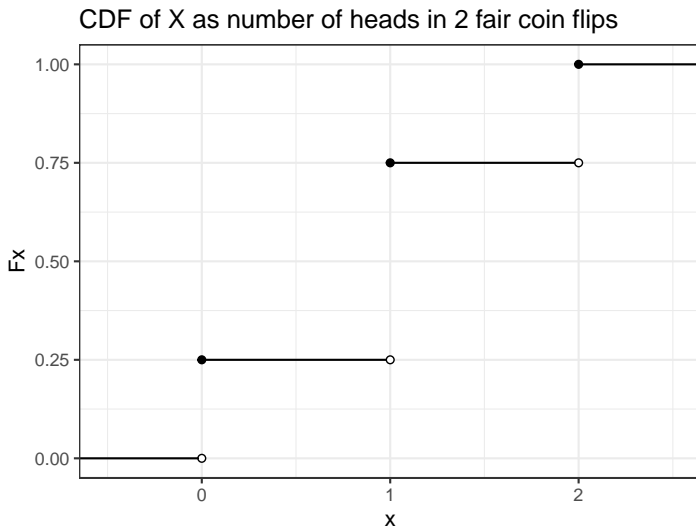
$$F(x) = P[X \leq x]$$

For our coin flip example,

$$F(x) = \begin{cases} 0 & x < 0 \\ 1/4 & 0 \leq x < 1 \\ 3/4 & 1 \leq x < 2 \\ 1 & x \geq 2 \end{cases}$$

CDFs are really useful, because if we know the CDF, we can fully describe the distribution of *any* random variable.

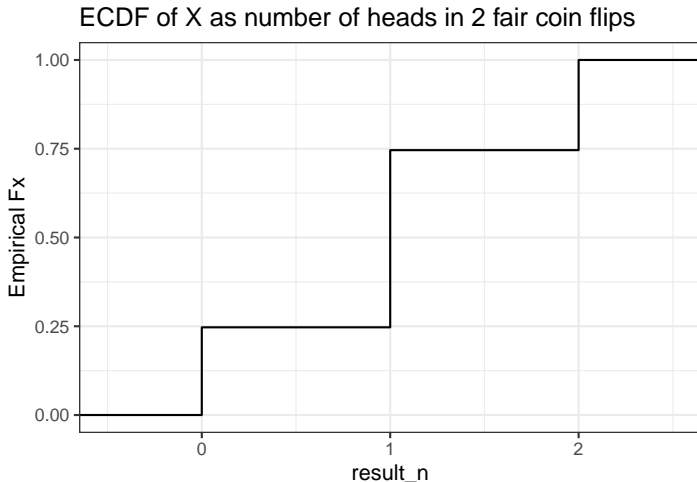
Illustrating the CDF of a RV



Illustrating the CDF of a RV

And we can use ggplot2 to see what the *Empirical* CDF looks like

```
> ggplot(data.frame(result_n), aes(x = result_n)) +  
+   stat_ecdf() +  
+   coord_cartesian(xlim = c(-0.5, 2.5)) +  
+   ylab('Empirical Fx') +  
+   ggtitle('ECDF of X as number of heads in 2 fair coin flips') + theme_bw()  
>
```



Joint and conditional relationships

Bivariate relationships

We often care about how random variables vary with each other

- ▶ age and voter turnout
- ▶ sex and income
- ▶ education and earnings

Just like with univariate random variables, we can describe these bivariate relationships by their distributions

Joint PMF of discrete random variables

$$f(x, y) = P[X = x, Y = y]$$

Returning to our example of flipping two fair coins

- ▶ Let X be 1 if we get *at least one heads*, and 0 otherwise
- ▶ Let Y be 1 if we get *two heads* in our two coin flips, and 0 otherwise

Then the joint probability distribution can be described as:

ω	$P[\{\omega\}]$	$X(\omega)$	$Y(\omega)$
TT	1/4	0	0
TH	1/4	1	0
HT	1/4	1	0
HH	1/4	1	1

or, considering the joint PMF,

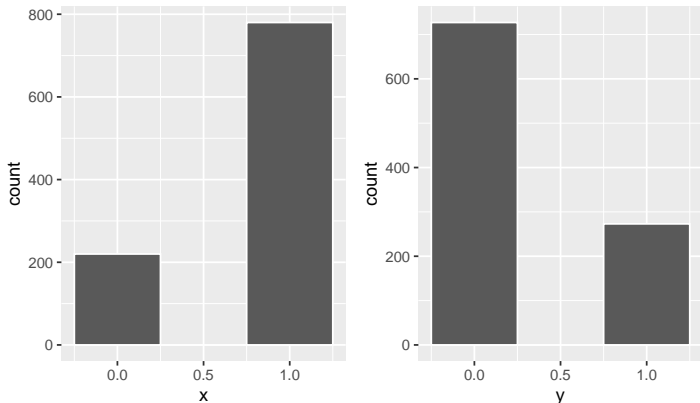
$$f(x, y) = \begin{cases} 1/4 & x = 0, y = 0 \\ 1/2 & x = 1, y = 0 \\ 1/4 & x = 1, y = 1 \\ 0 & \text{otherwise} \end{cases}$$


```

> Omega <- c('HH', 'HT', 'TH', 'TT')
> probs <- c(0.25, 0.25, 0.25, 0.25)
> result_n <- sample(x = Omega,
+                   size = n,
+                   prob = probs,
+                   replace = TRUE)
> result_mat <- data.frame(omega = result_n,
+                          x = ifelse(result_n == 'TT', 0, 1),
+                          y = ifelse(result_n == 'HH', 1, 0))
> options <- list(theme(panel.grid.minor = element_blank()), scale_x_continuous(breaks = 0, 1),
+                  theme_bw())
> p1 <- ggplot(result_mat) + geom_histogram(aes(x = x), bins = 3, position = 'identity')
> p2 <- ggplot(result_mat) + geom_histogram(aes(x = y), bins = 3, position = 'identity')
>

```

```
> grid.arrange(p1, p2, ncol = 2)
```



Seeing X and Y plotted side by side doesn't really give us a full picture of their relationship.

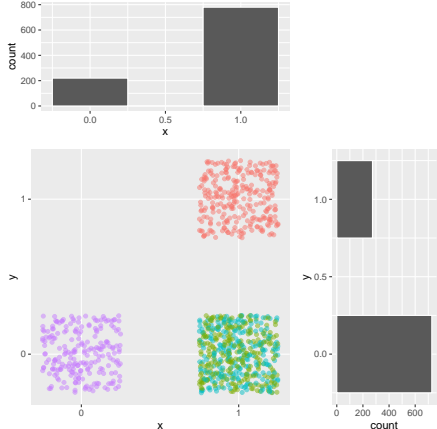
These are the *marginal* distributions of X and Y , i.e., their distributions where we *marginalize* or sum over the distribution of the other random variable.

Marginal distributions

$$f_X(x) = P[X = x] = \sum_y P[X = x, Y = y] = \sum_y f_{X,Y}(x, y)$$

	$Y = 0$	$Y = 1$	
$X = 0$	1/4	0	1/4
$X = 1$	1/2	1/4	3/4
	3/4	1/4	

Notational aside: we can subscript X in f_X to denote that it is the mass function of X specifically, as X and Y have different probability mass functions. But often we will just omit the subscript for convenience.



Plotting X and Y jointly gives us a better understanding of their joint relationship.

Conditional distributions

We are also often interested in conditional relationships.

$$f_{Y|X}(y|x) = P[Y = y|X = x] = \frac{P[X = x, Y = y]}{P[X = x]} = \frac{f_{X,Y}(x, y)}{f_X(x)}$$

For example,

$$f_{Y|X}(y|x) = \begin{cases} 1 & x = 0, y = 0 \\ 2/3 & x = 1, y = 0 \\ 1/3 & x = 1, y = 1 \\ 0 & \text{otherwise} \end{cases}$$

Here, what is the probability of observing two heads, conditional on having observed at least one heads?

Summarizing single variable distributions

Expectation

$$E[X] = \sum_x xf(x)$$

- ▶ Expectation is an *operator* on a random variable; it maps the distribution of X to a specific number.
- ▶ Specifically, the expectation operator tells us about the mean, or average value of X across its distribution.

Notational aside: it is common to write the expectation of a distribution as μ .

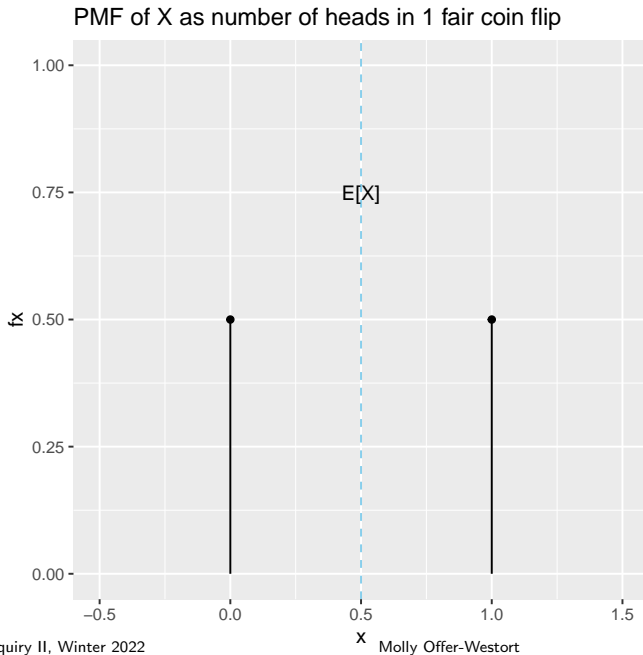
Let's flip a single coin, and let X be 1 if we get a head, and 0 otherwise.

$$f(x) = \begin{cases} 1/2 & x = 0 \\ 1/2 & x = 1 \\ 0 & \text{otherwise} \end{cases}$$

Mathematically,

$$\begin{aligned} E[X] &= \sum_x xf(x) \\ &= 0 \times \frac{1}{2} + 1 \times \frac{1}{2} \\ &= \frac{1}{2} \end{aligned}$$

Visually,



Spread of a distribution

We often describe the spread of a distribution by its variance

$$\text{Var}[X] = E[(X - E[X])^2]$$

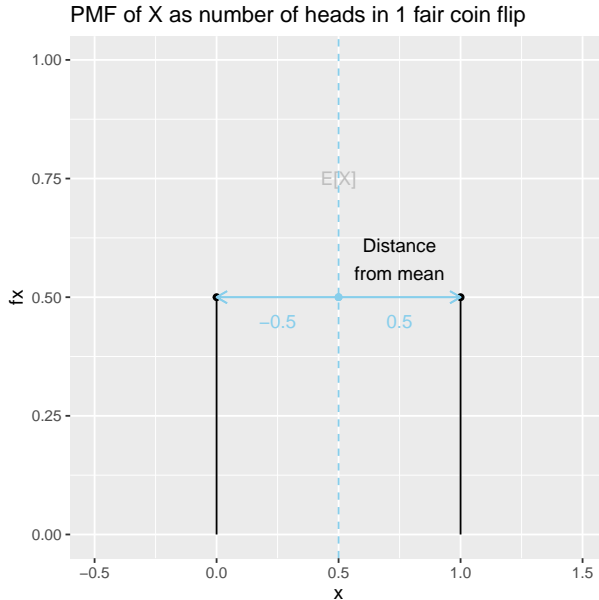
Or equivalently,

$$= E[X^2] - E[X]^2$$

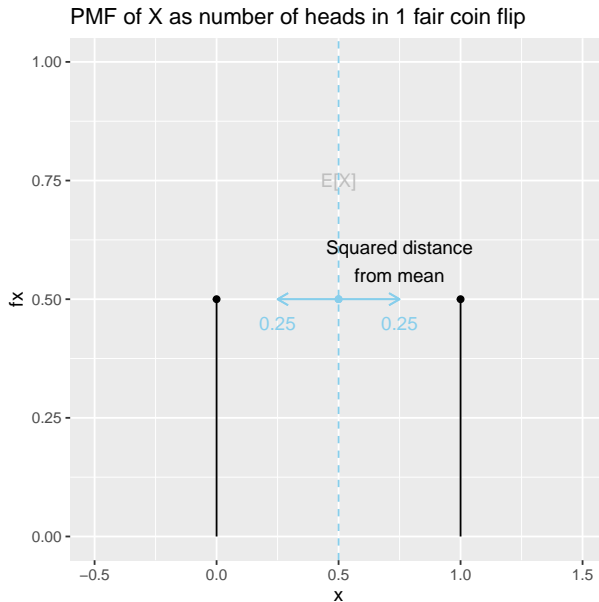
The standard deviation is the square root of the variance.

Notational aside: it is common to write the variance of a distribution as σ^2 , or the standard deviation as σ .

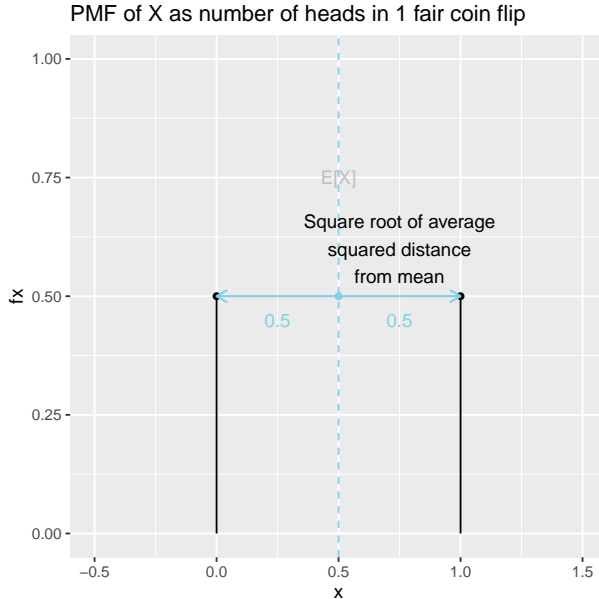
The variance is the average squared distance from the mean. The standard deviation is the square root of this.



The variance is the average squared distance from the mean. The standard deviation is the square root of this.



The variance is the average squared distance from the mean. The standard deviation is the square root of this.

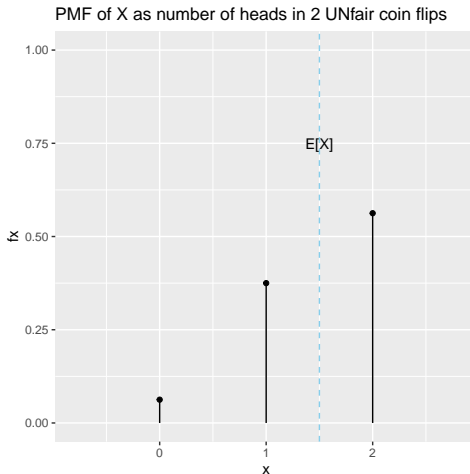


Let's take another example, where we flip a coin twice, and let X be the number of heads. However, let's say our coin is *not* fair, and the probability of getting a heads is 0.75.

The random variable's probability distribution is then:

$$f(x) = \begin{cases} 1/16 & x = 0 \\ 3/8 & x = 1 \\ 9/16 & x = 2 \\ 0 & \text{otherwise} \end{cases}$$

Let's take a look at the mean.

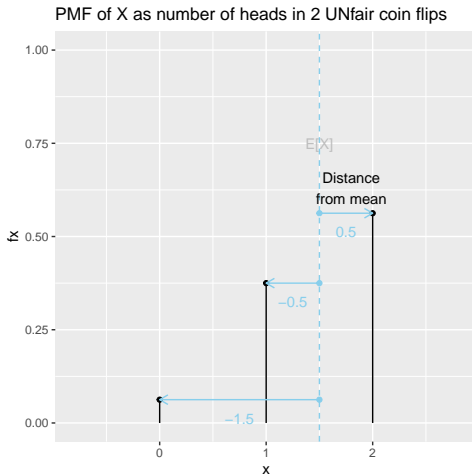


$$E[X] = \sum_x xf(x)$$

$$= 0 \times \frac{1}{16} + 1 \times \frac{3}{8} + 2 \times \frac{9}{16}$$

$$= \frac{24}{16} = 1.5$$

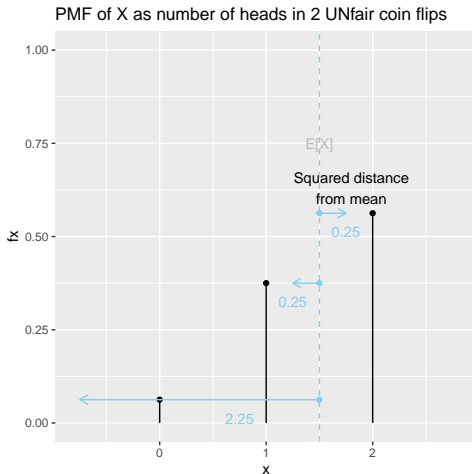
And the spread.



Variance = average squared distance from the mean

$$\text{Var}[X] = E[(X - E[X])^2]$$

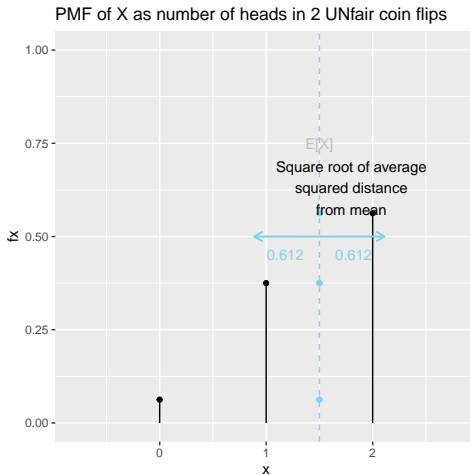
And the spread.



Variance = average squared distance from the mean

$$\begin{aligned}\text{Var}[X] &= E[(X - E[X])^2] \\ &= 2.25 \times \frac{1}{16} + 0.25 \times \frac{3}{8} + 0.25 \times \frac{9}{16} = 0.375\end{aligned}$$

And the spread.



SD = square root of variance

$$= \sqrt{0.375} = 0.612$$

Applications

- ▶ Coin flips are a pretty trivial example of a random event \rightarrow random variable.
- ▶ But often, as researchers, our job is to map events that happen in the world to variables in our data sets.

Presidential Daily Briefing

The President's Daily Brief

The PDB is an intelligence report created by the Central Intelligence Agency (CIA), which synthesizes multiple information streams from around the world—including local media, reports from U.S. outposts abroad, and clandestine activities—into a succinct and accessible overview of new developments in global affairs. The document is delivered to the president and a narrow circle of senior officials each morning, with follow-up oral briefings taking place at a reader's request. The PDB is widely considered to be the premier document created by the intelligence community. In his memoirs, former CIA director George Tenet characterized the PDB as “our most important product” (Tenet, 2007, 30).

Presidential Daily Briefing

What goes into the data set?

DAILY BRIEF
7 AUGUST 1965

1. Vietnam

A Soviet cargo ship, the Polotsk,

50X1



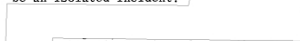
is en route to Haiphong. The ship unloaded military cargo in Indonesia; there is no evidence it is carrying such cargo now. 50X1

2. South Vietnam

There has been no significant change in the situation at Duc Co in Pleiku Province, where South Vietnamese airborne troops are trying to eliminate Viet Cong harassment of a government paramilitary camp.

3. Communist China

The loss of two Chinese Nationalist patrol craft on 5 August in an encounter with Chinese Communist naval vessels off the mainland coast at the southern end of the Taiwan Strait appears so far to be an isolated incident. 50X1



There is no sign of any other significant Communist military reaction to what seems to have been a Nationalist incursion into Communist-controlled waters. (See map)

- ▶ Unit of observation: "entries" in daily briefings
- ▶ Record date
- ▶ President
- ▶ Pages in briefing
- ▶ Number of maps
- ▶ How to code redactions?

Presidential Daily Briefing

What goes into the data set?

The *event* that happens is a certain amount of the briefing is redacted before it's made public.
How is this encoded in a variable?

DAILY BRIEF
7 AUGUST 1965

1. Vietnam

A Soviet cargo ship, the Polotsk,

50X1



is en route to Haiphong. The ship unloaded military cargo in Indonesia; there is no evidence it is carrying such cargo now.

50X1

2. South Vietnam

There has been no significant change in the situation at Duc Co in Pleiku Province, where South Vietnamese airborne troops are trying to eliminate Viet Cong harassment of a government paramilitary camp.

3. Communist China

The loss of two Chinese Nationalist patrol craft on 5 August in an encounter with Chinese Communist naval vessels off the mainland coast at the southern end of the Taiwan Strait appears so far to be an isolated incident.

50X1



There is no sign of any other significant Communist military reaction to what seems to have been a Nationalist incursion into Communist-controlled waters. (See map)

Presidential Daily Briefing

```
> file <- "https://raw.githubusercontent.com/UChicago-pol-methods/IntroQSS-F21/main/  
> df_pdb <- read.csv(file, as.is = TRUE)  
> head(df_pdb)
```

	PDBid	date	President	Total_Pgs	maps	Redaction_total
1	17061961	6/17/61	1	8	2	26
2	19061961	6/19/61	1	7	1	23
3	20061961	6/20/61	1	5	0	8
4	21061961	6/21/61	1	5	0	17
5	22061961	6/22/61	1	6	1	13
6	23061961	6/23/61	1	9	1	22

```
> sapply(df_pdb[,c('Total_Pgs', 'maps', 'Redaction_total')],  
+       function(x) c('mean' = mean(x, na.rm = TRUE),  
+                     'var' = var(x, na.rm = TRUE)))
```

	Total_Pgs	maps	Redaction_total
mean	9.44898	0.6460584	8.693077
var	11.98529	0.7410201	22.452908

Summarizing joint distributions

Covariance

$$\text{Cov}[X, Y] = E[(X - E[X])(Y - E[Y])]$$

Covariance is how much X and Y vary together.

- ▶ If covariance is positive, when the value of X is large (relative to its mean), the value of Y will also tend to be large (relative to its mean)
- ▶ If covariance is negative, when the value of X is large (relative to its mean), the value of Y will tend to be small (relative to its mean)

Correlation

$$\rho[X, Y] = \frac{\text{Cov}[X, Y]}{\sigma[X]\sigma[Y]}$$

Rescaled version of covariance

- ▶ positive when covariance is positive
- ▶ negative when covariance is negative

$$-1 \leq \rho[X, Y] \leq 1$$

What relationship would you expect to see between number of pages in a Presidential Daily Brief and number of redactions?

Number of pages and number of redactions have a positive linear relationship.
Covariance (and correlation) is positive.

Following content if time, or in handout.

Continuous random variables

- ▶ So far, our coin flip example was for a *discrete* random variable.
- ▶ A random variable is *continuous* if it has a continuous density function
- ▶ Practically, we will treat RVs as discrete if they have countably many outcomes, and RVs as continuous if the number of values they can take on is only constrained by our measurement tool.

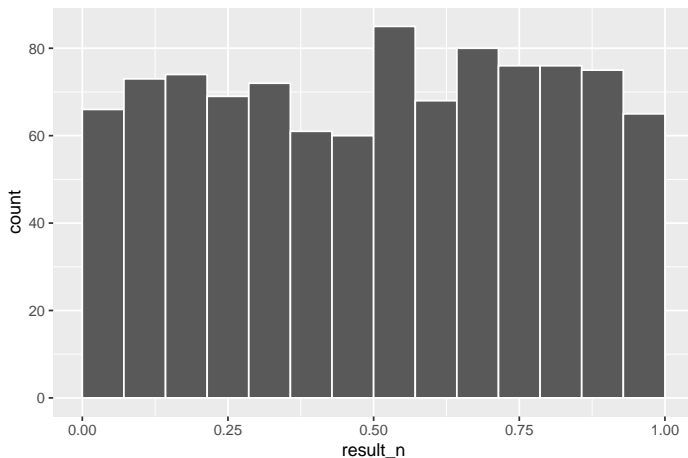
Uniform distribution

- ▶ If you take a draw from the standard uniform distribution, you are equally likely to draw any number between zero and one.
- ▶ We can simulate this in R. R allows you to sample from a number of canonical distributions; to see which distributions are available, search '?Distributions'.

```
> runif(n = 1, min = 0, max = 1)
```

```
[1] 0.4894512
```

We can again sample from the distribution many times, and plot a histogram to look at the distribution of results.



Probability Density Function of continuous random variables

- ▶ Discrete random variables have non-zero mass on specific points, but for continuous random variables, $P[X = x] = 0$. Instead of mass, we refer to *density* for continuous variables.
- ▶ The *probability density function* $f(x)$ for a continuous random variable gives the slope of the CDF at any given point. This means that we can integrate the area under the PDF to get the relative probability of being between two points.

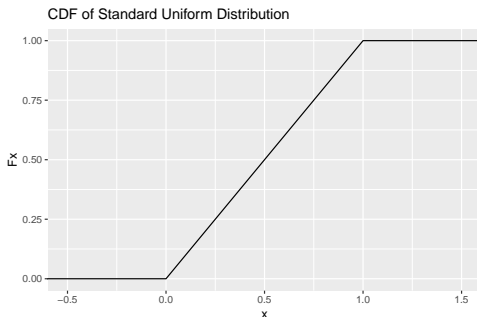
$$P[a < X < b] = \int_a^b f(x)dx$$

Illustrating the CDF of a continuous RV

We start by showing the CDF of the standard uniform distribution, to illustrate how the PDF relates to the CDF. The CDF for the standard uniform distribution is:

$$F(x) = \begin{cases} 0 & x < 0 \\ x & 0 \leq x \leq 1 \\ 1 & x > 1 \end{cases}$$

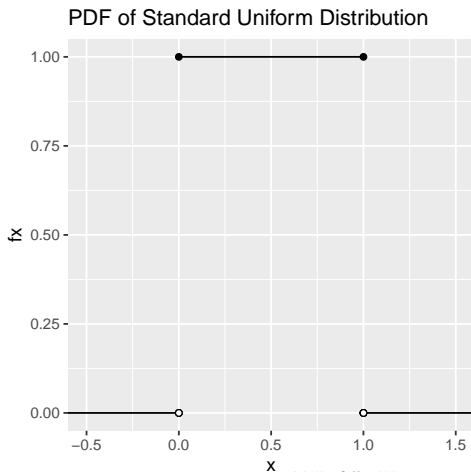
Notice that the slope is 1 between 0 and 1.



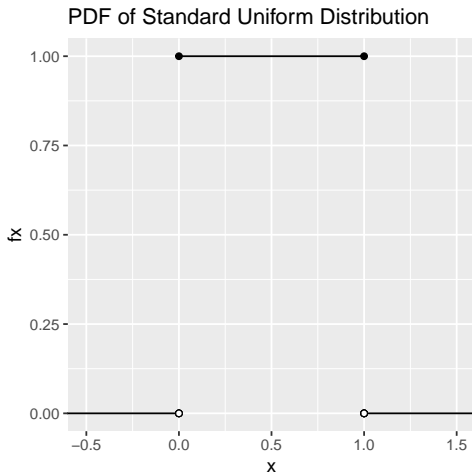
Illustrating the PDF of a continuous random variable

The PDF for the standard uniform distribution is:

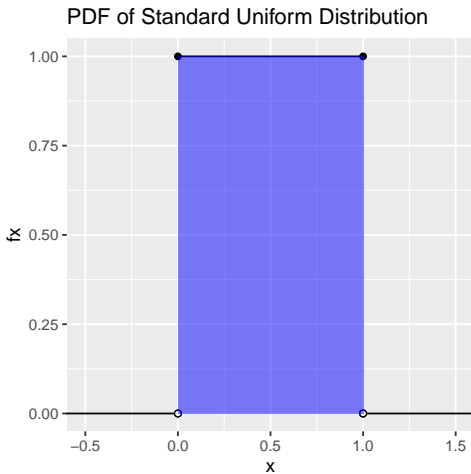
$$f(x) = \begin{cases} 1 & 0 < x < 1 \\ 0 & \text{otherwise.} \end{cases}$$



Notice that if we take the area under the density curve, the total area will sum to 1.



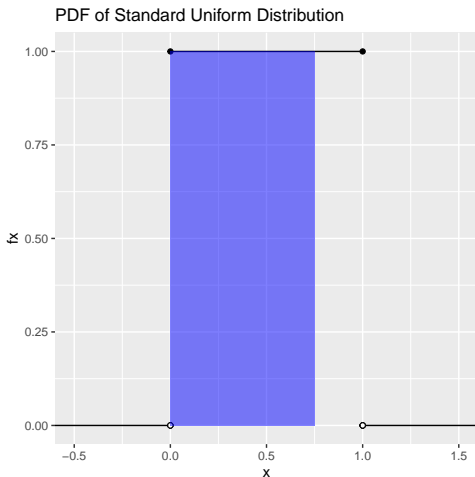
Notice that if we take the area under the density curve, the total area will sum to 1. Relative density gives us relative probability.



If we want to get the probability X is between 0 and 0.75,

$$P[0 \leq x \leq .75] = \int_0^{.75} f(x) dx$$

we take the area under the density curve between 0 and 0.75 – which is also 0.75. (Notice that we don't need to use calculus here.)

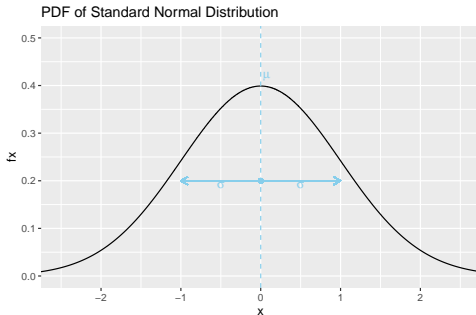


Normal distribution

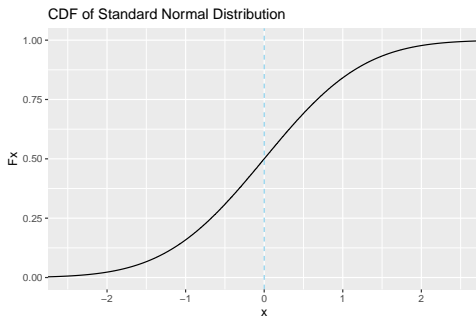
The Normal distribution is frequently used in probability and statistics, because it is a useful approximation to many natural phenomena.

$$f(x) = \frac{1}{\sqrt{\sigma^2 2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

It is defined by two parameters, μ , the center of the distribution, and σ , which defines the distribution's standard deviation, or spread. The distribution is often notated $\mathcal{N}(\mu, \sigma^2)$. It has a bell curve shape, with more density around the middle, and less density at more extreme values.



Normal distribution



References I