



THE UNIVERSITY  
*of* EDINBURGH

# Methods for Causal Inference

## Lecture 2: Basics of probability

---

Ava Khamseh

School of Informatics  
2023-2024

# Causal theory and data

Requires 5 steps:

1. Definition of Causation
2. Clearly formulating causal **assumptions** and creating the **causal model**
3. Linking the structure of causal model to features of data
4. **Estimating**, given the causal model and data
5. **Uncertainty quantification**, e.g., confidence/credible interval

# Causal theory and data

Requires 5 steps:

1. Definition of Causation
2. Clearly formulating causal **assumptions** and creating the **causal model**
3. Linking the structure of causal model to features of data
4. **Estimating**, given the causal model and data
5. **Uncertainty quantification**, e.g., confidence/credible interval

**Disclaimer:** In this course our focus is on 1-3. We then use simple models to exemplify 4-5 (taking model assumptions as ‘true’),

i.e., we do not discuss valid **statistical inference**.

For causal/statistical inference please refer to the course:

**Targeted Causal Learning** (code: MATH11238).



# Causal theory and data

Requires 5 steps:

1. Definition of Causation
2. Clearly formulating causal **assumptions** and creating the **causal model**
3. Linking the structure of causal model to features of data
4. **Estimating**, given the causal model and data
5. **Uncertainty quantification**, e.g., confidence/credible interval

## Defining causation:

A variable  $X$  is a **cause** of a variable  $Y$  if  $Y$  in any way relies on  $X$  for its value. (Intuitively:  $X$  is a cause of  $Y$  if  $Y$  listens to  $X$  and decides its value in response to what it hears)

**Pre-requisites:** Elementary concepts from probability theory, statistics, graph theory

# Basics of probability

Most causal statements are uncertain: “drinking causes liver disease”, does not mean every person who consumes alcohol is certain to have liver disease



Need language and laws of probability.

# Basics of probability

Most causal statements are uncertain: “drinking causes liver disease”, does not mean every person who consumes alcohol is certain to have liver disease



Need language and laws of probability.

**(Random) variables:** Any property or descriptor that can take multiple values, e.g., age ( $x=40$ ), sex ( $x'=F$ ), family history of disease ( $x''=0$ ), ... .

# Basics of probability

Most causal statements are uncertain: “drinking causes liver disease”, does not mean every person who consumes alcohol is certain to have liver disease



Need language and laws of probability.

**(Random) variables:** Any property or descriptor that can take multiple values, e.g., age ( $x=40$ ), sex ( $x'=F$ ), family history of disease ( $x''=0$ ), ... .

**Events:** An event is any assignment of a **value or set of values** to a variable or set of variables.

# Basics of probability

Most causal statements are uncertain: “drinking causes liver disease”, does not mean every person who consumes alcohol is certain to have liver disease



Need language and laws of probability.

**(Random) variables:** Any property or descriptor that can take multiple values, e.g., age ( $x=40$ ), sex ( $x'=F$ ), family history of disease ( $x''=0$ ), ... .

**Events:** An event is any assignment of a **value or set of values** to a variable or set of variables.

**Example:** Individual  $> 40$  and recovered from covid  $y=0$ , event is  $(x > 40, y=0)$ . So variables are ‘age’ and ‘recovery status’ with values  $> 40$  and  $0$ .

Can ask what is the probability of an event, e.g., what is  $P(x > 40, y=0)$ ?



# Basics of probability

Most causal statements are uncertain: “drinking causes liver disease”, does not mean every person who consumes alcohol is certain to have liver disease



Need language and laws of probability.

**(Random) variables:** Any property or descriptor that can take multiple values, e.g., age ( $x=40$ ), sex ( $x'=F$ ), family history of disease ( $x''=0$ ), ... .

**Events:** An event is any assignment of a **value or set of values** to a variable or set of variables.

**Discrete** (binary/categorical): Are being treated or not, have a disease or not, ...

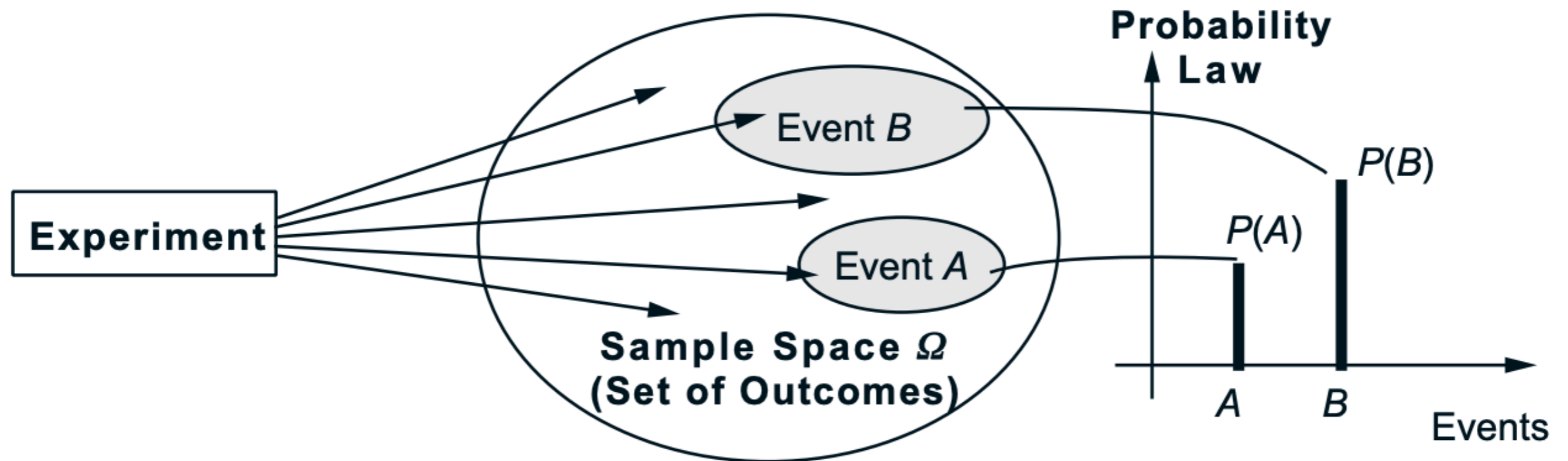
**Continuous** (can take infinite set of values): age, weight, ...

Drug (yes/no) vs dose of drug (categorical). Sun intake (time is continuous)

# Basics of probability

For probabilistic modelling (of a random experiment) we need to:

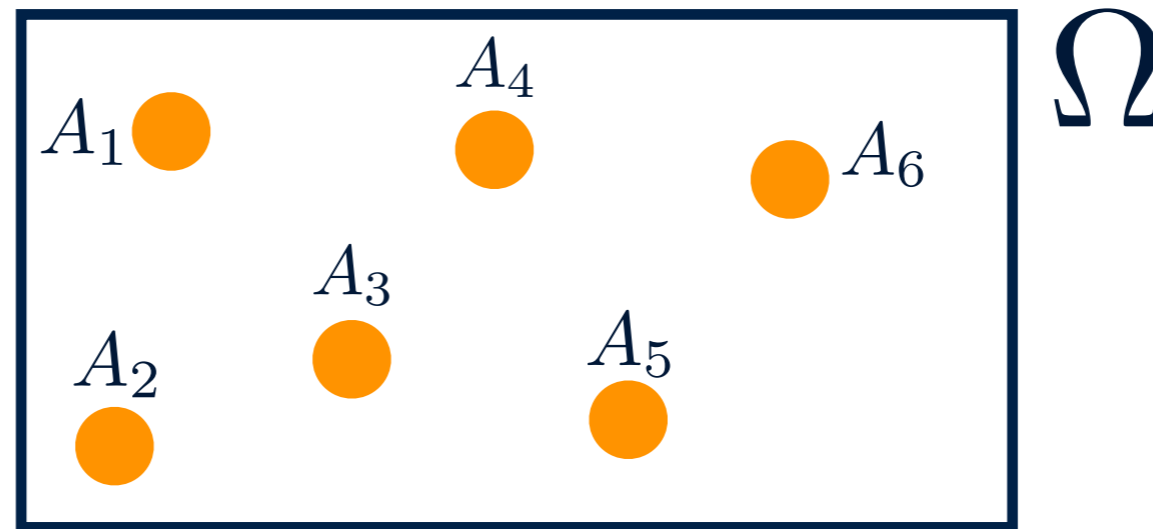
- Describe possible outcomes: **sample space**
- **Event**: A subset of sample space
- Describe beliefs about likelihood of these events: **probability law**



# Sample space

The sample space is the set of all possible outcomes of the experiment:

e.g. Rolling a dice



Outcomes must be:

- **Mutually Exclusive:** If I tell you, after the experiment, that  $A_1$  happened, then it should not be possible that  $A_6$  also happened.
- **Collectively Exhaustive:** Collectively, all the outcomes in  $\Omega$  exhaust all possibilities.

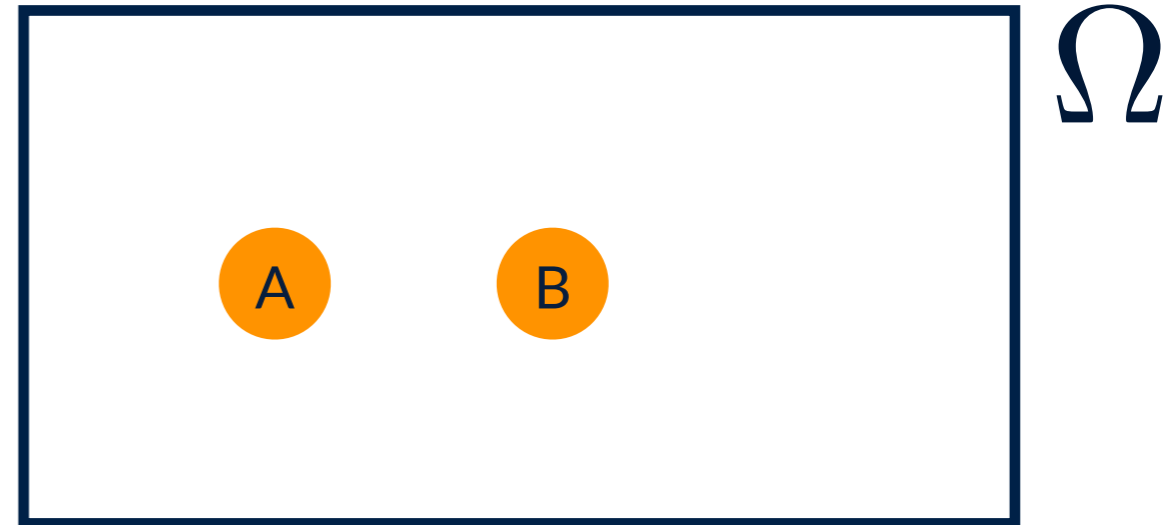
# Probability Axioms

Non-negativity:  $P(A) \geq 0$

Normalisation:  $P(\Omega) = 1$

For any two **mutually exclusive events** (i.e. A and B cannot co-occur) we have:

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



# Probability Axioms

Non-negativity:  $P(A) \geq 0$

Normalisation:  $P(\Omega) = 1$

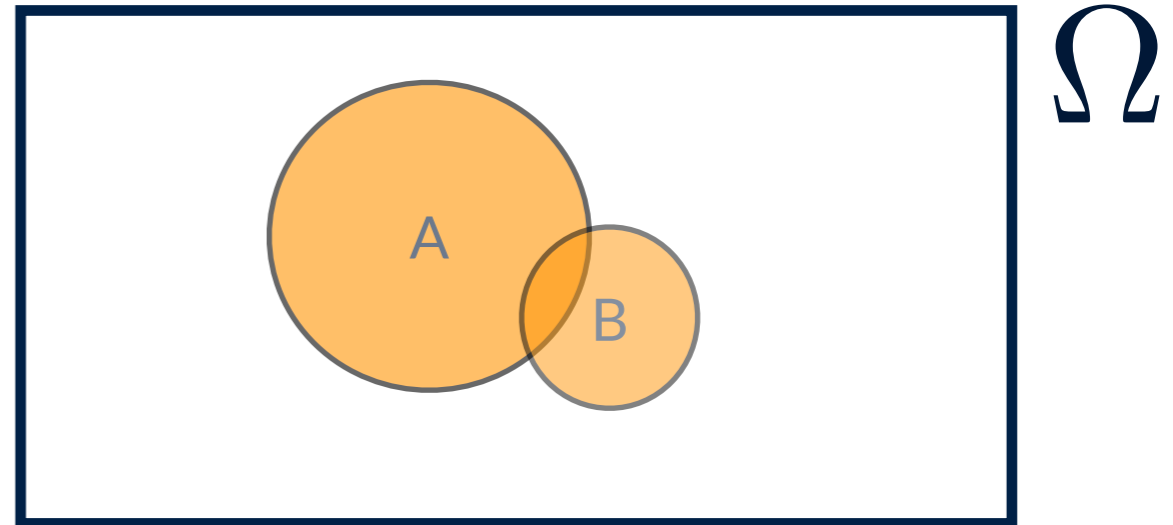
For any two **mutually exclusive events** (i.e. A and B cannot co-occur) we have:

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

As a consequence, take any two events A and B (they may overlap!), then:

$$P(A) = P(A \text{ and } B) + P(A \text{ and 'not } B')$$

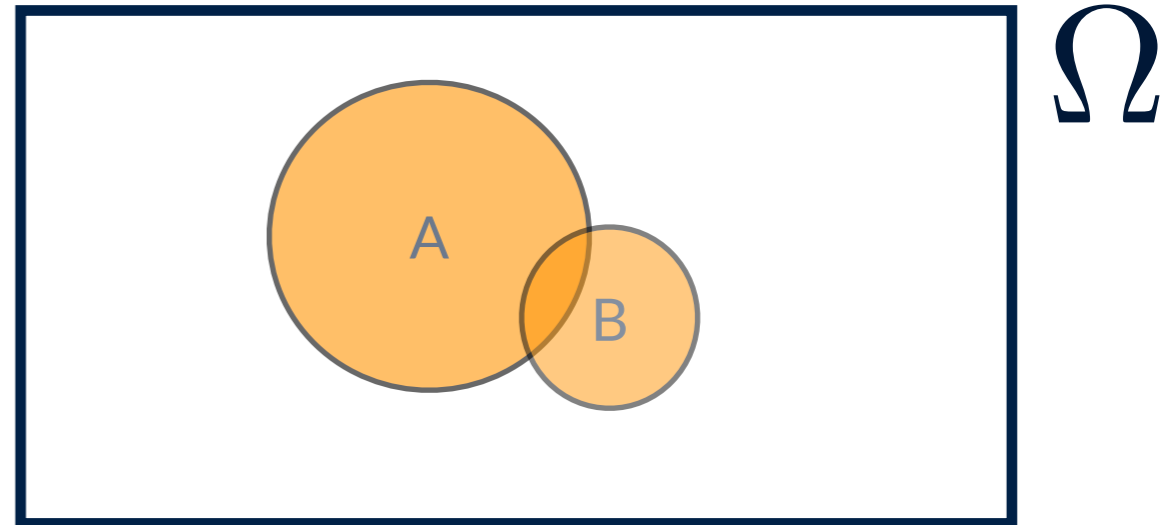
Why is the last statement true?



# Probability Axioms

- Non-negativity:  $P(A) \geq 0$
- Normalisation:  $P(\Omega) = 1$
- For any two **mutually exclusive events** (i.e. A and B cannot co-occur) we have:

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



As a consequence, take any two events A and B (they may overlap!), then:

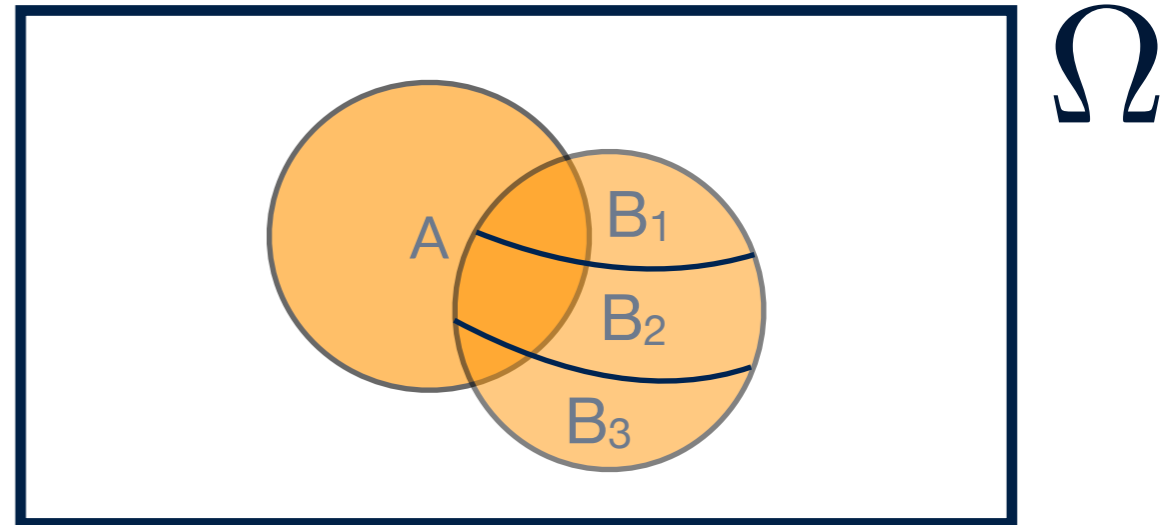
$$P(A) = P(A \text{ and } B) + P(A \text{ and 'not } B')$$

Mutually exclusive: If A is true, either “A and B” or “A and not B” must be true.

# Probability Axioms

- Non-negativity:  $P(A) \geq 0$
- Normalisation:  $P(\Omega) = 1$
- For any two **mutually exclusive events** (i.e. A and B cannot co-occur) we have:

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



Corollary:  $B_1, B_2, B_3$ , are exclusive, and together form all of  $B$ . Then,

$$P(A \text{ and } B) = P(A \text{ and } B_1) + P(A \text{ and } B_2) + P(A \text{ and } B_3)$$

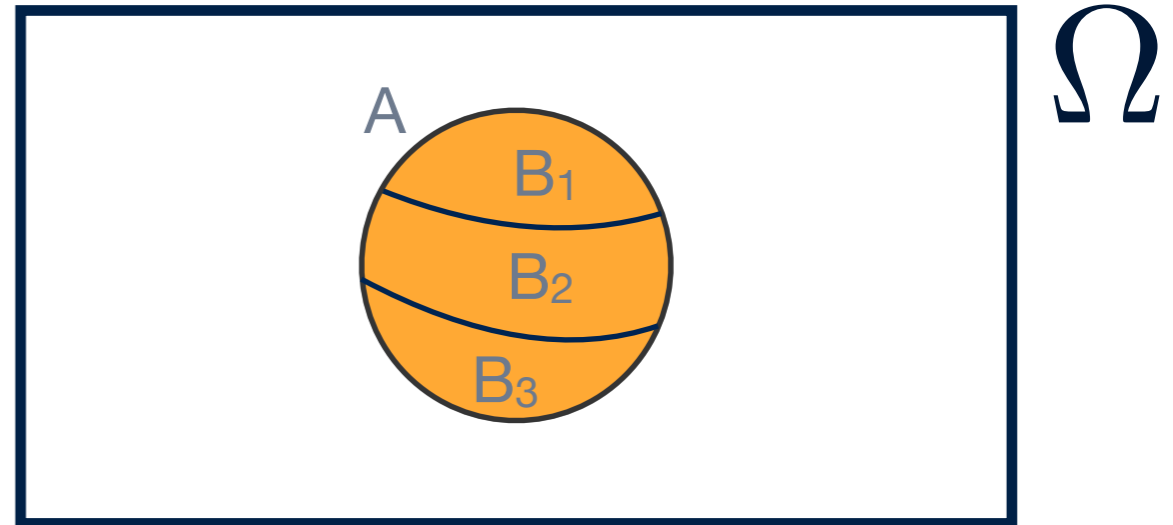
Generalise for (exhaustive, mutually exclusive) **partitions** of  $B$ :

$$P(A \text{ and } B) = \sum_{i=1}^n P(A \text{ and } B_i) \quad \text{where } B_i \cap B_j = \emptyset, \bigcup_{i=1}^n B_i = B$$

# Probability Axioms

- Non-negativity:  $P(A) \geq 0$
- Normalisation:  $P(\Omega) = 1$
- For any two **mutually exclusive events** (i.e. A and B cannot co-occur) we have:

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



Corollary: Let  $B_i, i=1, \dots, n$  be mutually exclusive and exhaustive partitions of  $B$ , and let  $A=B$  (complete overlap). Then,

$$P(A) = P(A \text{ and } A) = P(A \text{ and } B) = \sum_{i=1}^n P(A \text{ and } B_i) \quad \text{where } B_i \cap B_j = \emptyset, \bigcup_{i=1}^n B_i = B$$

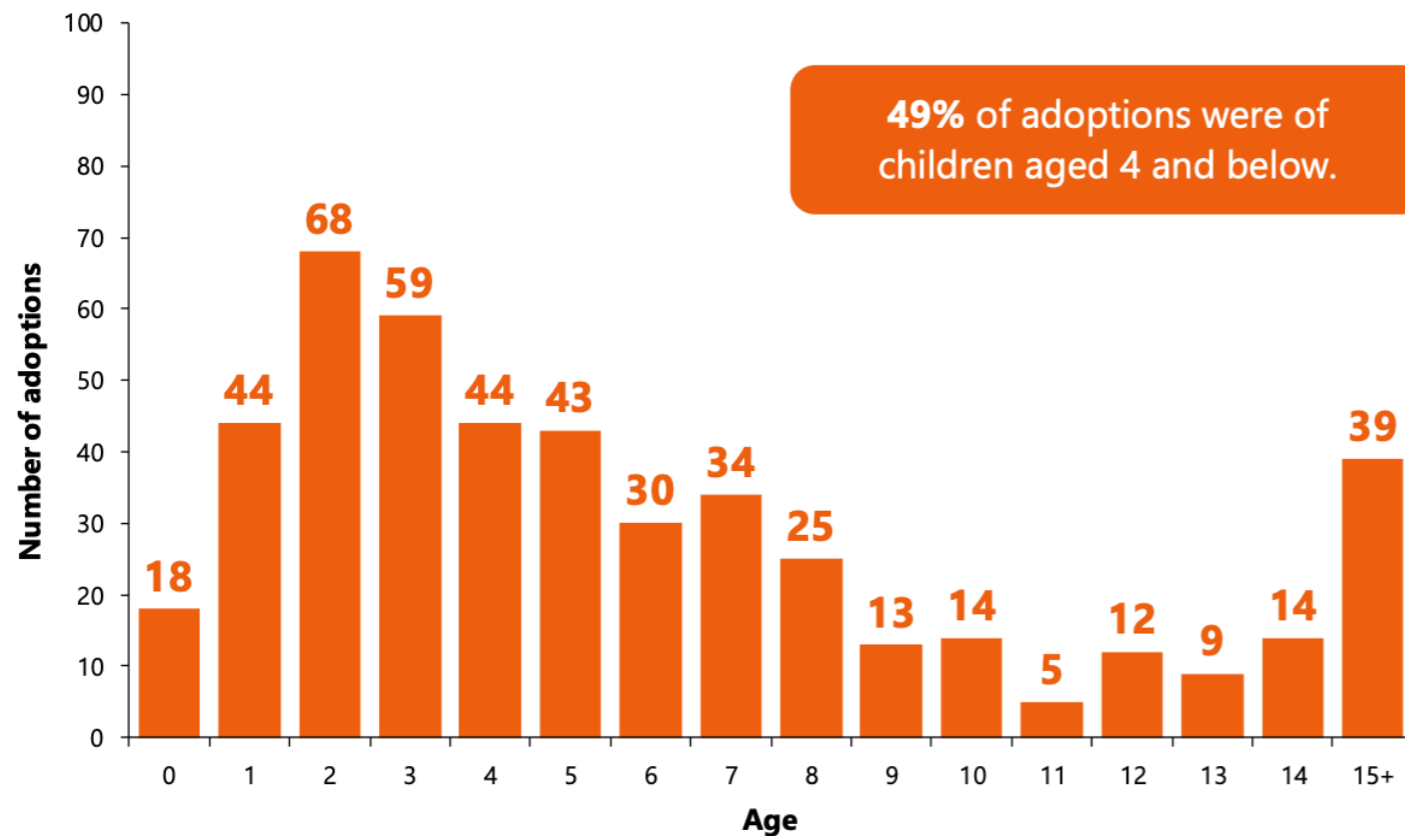
See later: “marginalisation”



# Intervals

$$P(\text{age} > 4) = 1 - P(\text{age} \leq 4) = 1 - 0.49 = 0.51$$

Figure 7.2: Age at adoption, Scotland, 2018



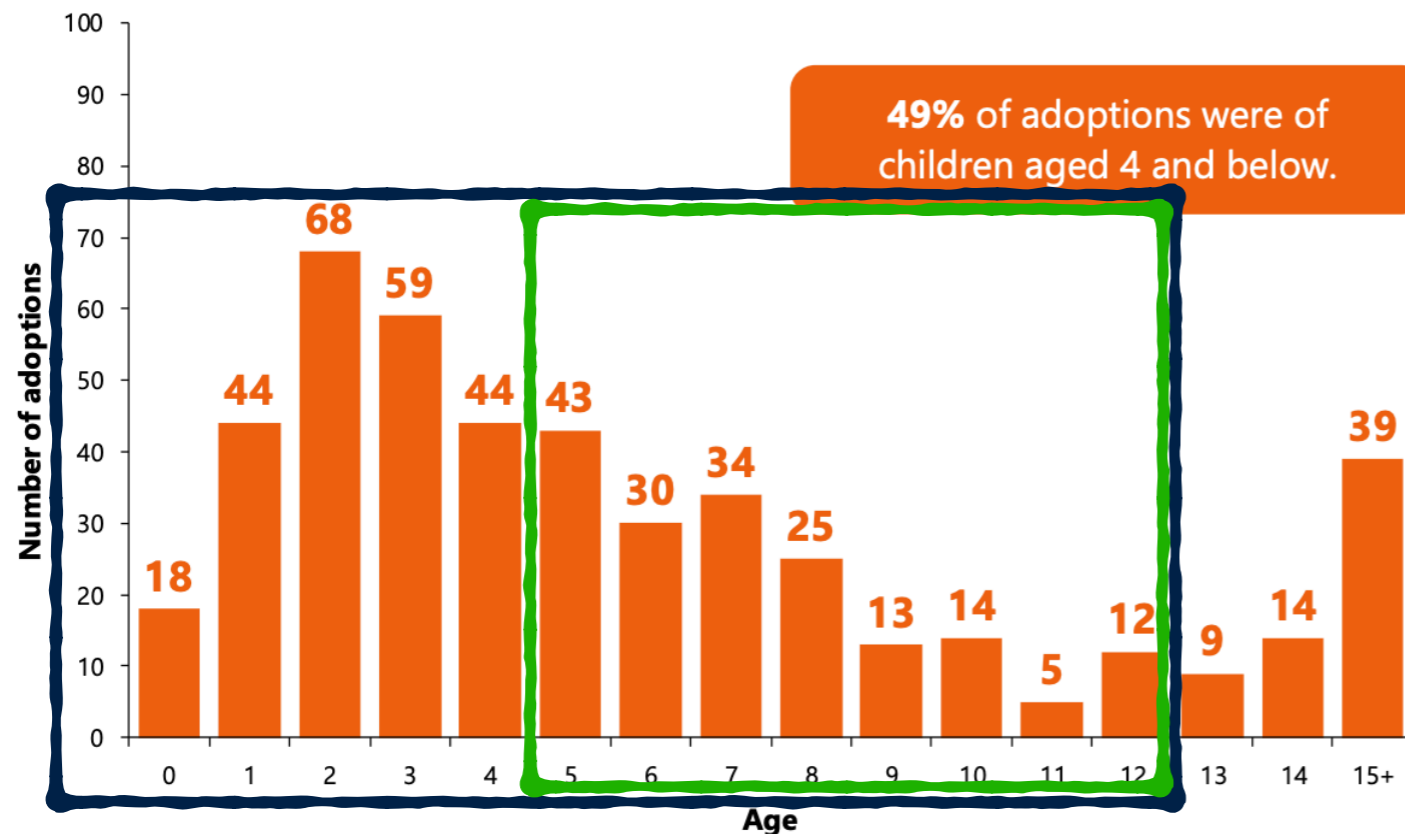
Total = 471

# Intervals

$$P(\text{age} > 4) = 1 - P(\text{age} \leq 4) = 1 - 0.49 = 0.51$$

$$P(4 < \text{age} \leq 12) = (43+30+34+25+13+14+5+12) / 471 = 0.37$$

Figure 7.2: Age at adoption, Scotland, 2018



Total = 471

# Law of Total probability: Example

Assuming 'no multi-tasking', the event:

“Passing the causality exam AND not being on your phone during the lectures”

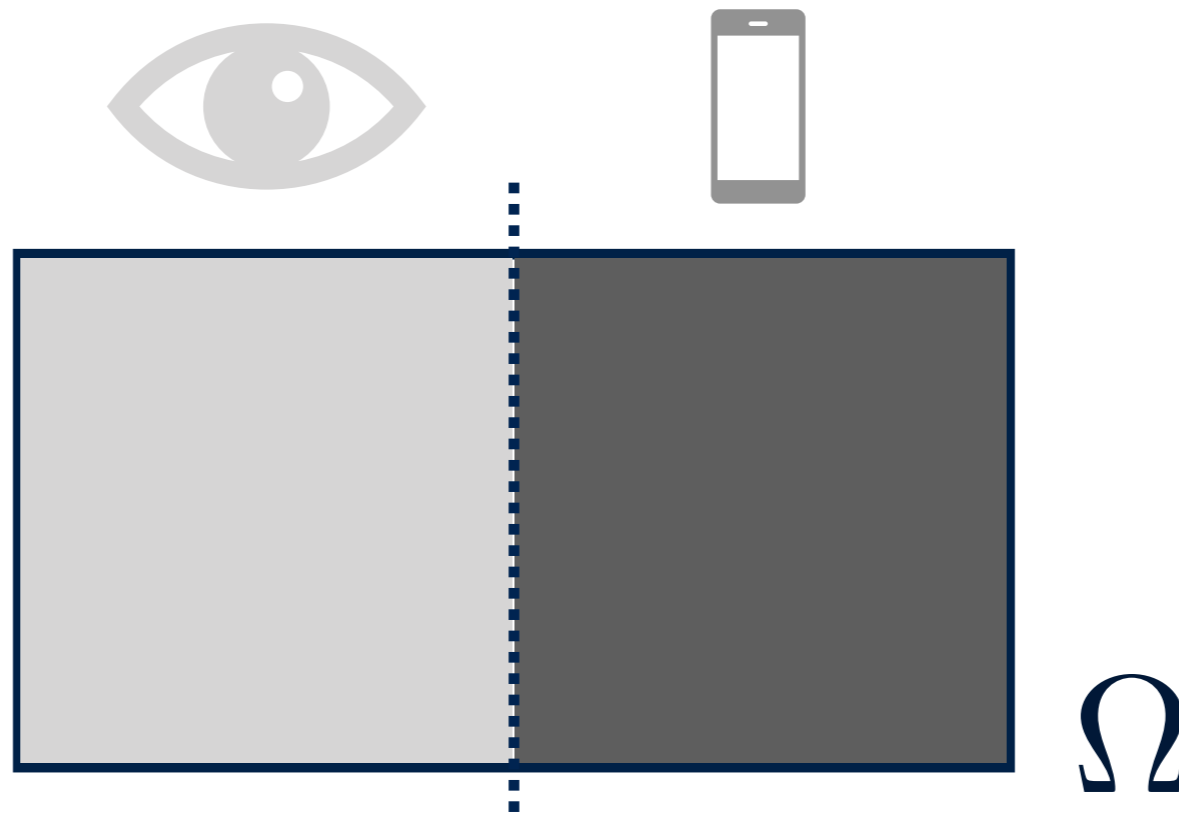
is **mutually exclusive** from

“Passing the causality exam AND being entirely on your phone during the lectures”

**P(passing the causality exam) =**

P(passing the exam, being entirely on your phone during the lecture) +

P(passing the exam, fully paying attention during the lecture)



# Law of Total probability: Example

Assuming 'no multi-tasking', the event:

“Passing the causality exam AND not being on your phone during the lectures”

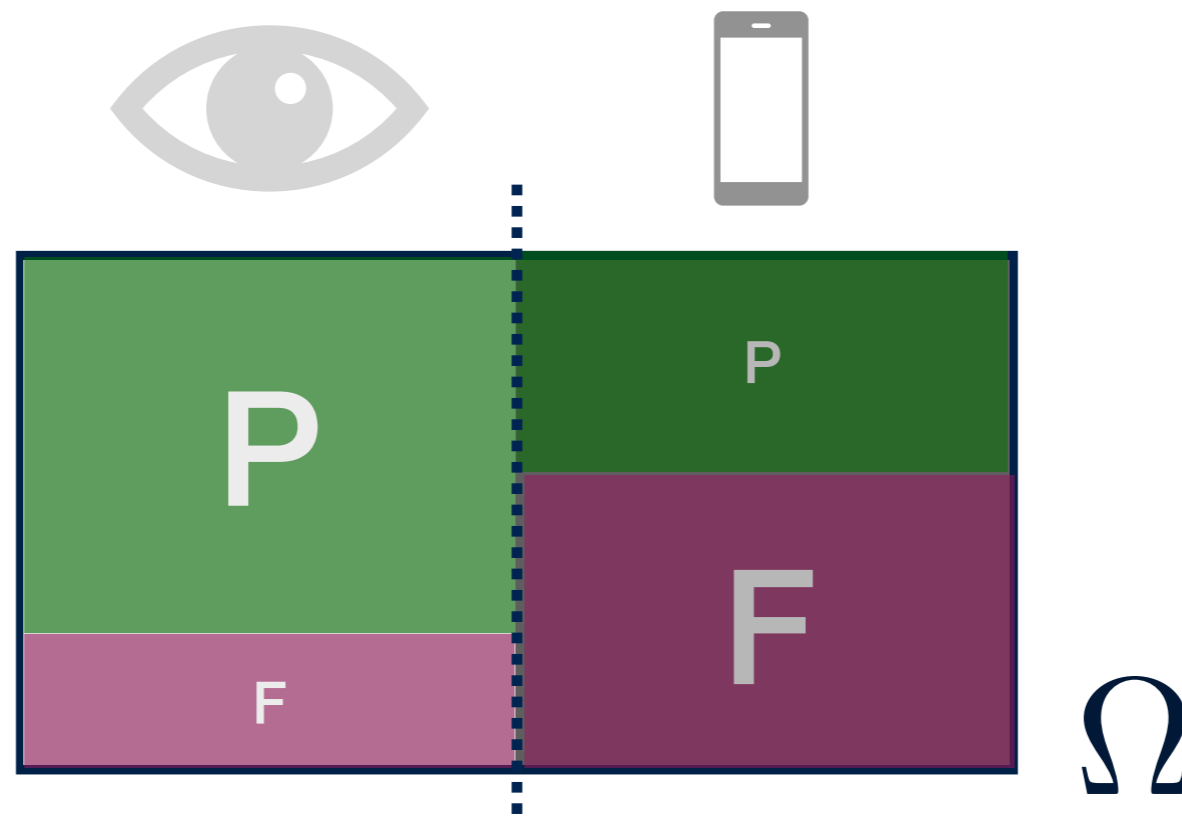
is **mutually exclusive** from

“Passing the causality exam AND being entirely on your phone during the lectures”

**P(passing the causality exam) =**

P(passing the exam, being entirely on your phone during the lecture) +

P(passing the exam, fully paying attention during the lecture)



# Law of Total probability: Example

Assuming 'no multi-tasking', the event:

“Passing the causality exam AND not being on your phone during the lectures”

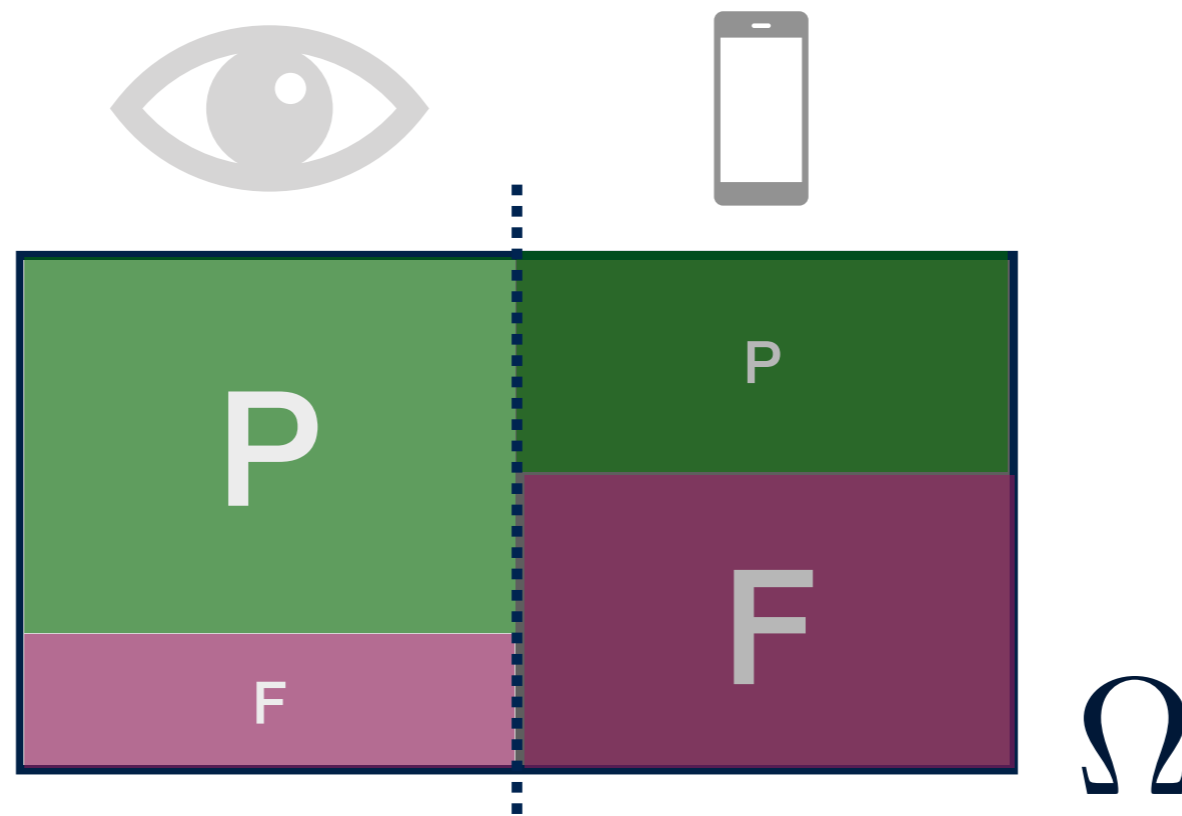
is **mutually exclusive** from

“Passing the causality exam AND being entirely on your phone during the lectures”

**P(passing the causality exam) =**

P(passing the exam, being entirely on your phone during the lecture) +

P(passing the exam, fully paying attention during the lecture)



# Conditional Probability

The probability that event A occurs, given that we know some other event B has occurred. (Think of filtering the data based on the value of some variable)

$P(X = x)$  vs  $P(X = x|Y = y)$ : The probability of  $X=x$  can drastically change depending on the knowledge  $Y=y$

# Conditional Probability

The probability that event A occurs, given that we know some other event B has occurred. (Think of filtering the data based on the value of some variable)

$P(X = x)$  vs  $P(X = x|Y = y)$ : The probability of  $X=x$  can drastically change depending on the knowledge  $Y=y$

**Example:**  $P(\text{lung cancer} | \text{smoker})$  vs

$P(\text{lung cancer} | \text{smoker, socio-economic status})$

Given that the patient is a smoker, does knowing their socio-economic status add further information to the probability of lung cancer?

# Conditional Probability

The probability that event A occurs, given that we know some other event B has occurred. (Think of filtering the data based on the value of some variable)

$P(X = x)$  vs  $P(X = x|Y = y)$ : The probability of  $X=x$  can drastically change depending on the knowledge  $Y=y$

**Example:**  $P(\text{lung cancer} | \text{smoker})$  vs

$P(\text{lung cancer} | \text{smoker, socio-economic status})$

Given that the patient is a smoker, does knowing their socio-economic status add further information to the probability of lung cancer?

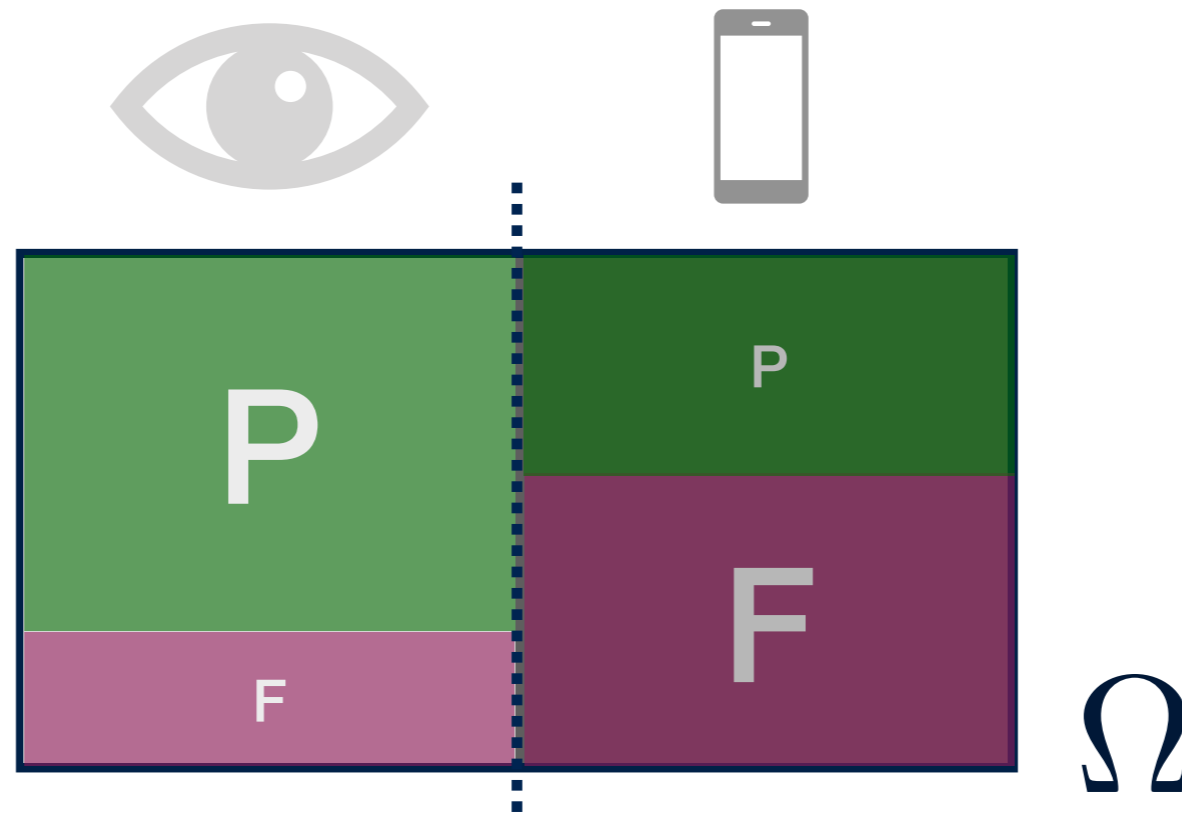
Relation between “**joint**”, “**conditional**”, and “**marginal**” probabilities:

$$P(X, Y) = P(X|Y)P(Y)$$



# Conditional Law of Total probability: Example

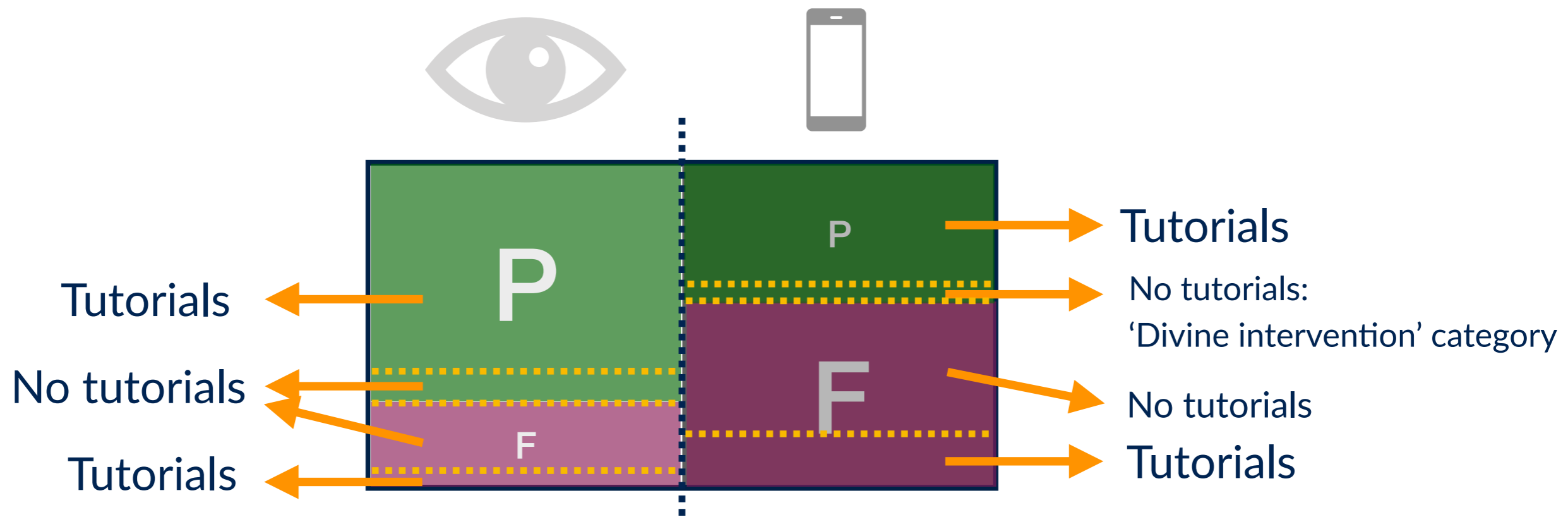
$P(\text{passing the causality exam} \mid \text{paying attention}) >$   
 $P(\text{passing the causality exam} \mid \text{being on your phone})$



# Conditional Law of Total probability: Example

$P(\text{passing the causality exam} \mid \text{fully paying attention during the lecture}) =$   
 $P(\text{passing the exam, attending tutorials} \mid \text{attention in lecture}) +$   
 $P(\text{passing the exam, not attending tutorials} \mid \text{attention in lecture})$

$P(\text{passing the causality exam} \mid \text{being on one's phone during the lectures}) =$   
 $P(\text{passing the exam, attending tutorials} \mid \text{being on phone during lecture}) +$   
 $P(\text{passing the exam, not attending tutorials} \mid \text{being on phone lecture})$



# Bayes' Rule

$A_1, A_2, \dots, A_n$  are disjoint events forming a **partition** of the sample space and  $P(A_i) > 0, \forall A_i$ . Then, for any event  $B, P(B) > 0$ , Bayes' rule states:

$$P(A_i|B) = \frac{P(A_i \cap B)}{P(B)} = \frac{P(A_i)P(B|A_i)}{P(B)}$$

# Bayes' Rule

$A_1, A_2, \dots, A_n$  are disjoint events forming a **partition** of the sample space and  $P(A_i) > 0, \forall A_i$ . Then, for any event  $B, P(B) > 0$ , Bayes' rule states:

$$\begin{aligned} P(A_i|B) &= \frac{P(A_i \cap B)}{P(B)} = \frac{P(A_i)P(B|A_i)}{P(B)} \\ &= \frac{P(A_i)P(B|A_i)}{P(A_1 \cap B) + \dots + P(A_n \cap B)} \\ &= \frac{P(A_i)P(B|A_i)}{P(A_1)P(B|A_1) + \dots + P(A_n)P(B|A_n)} \end{aligned}$$

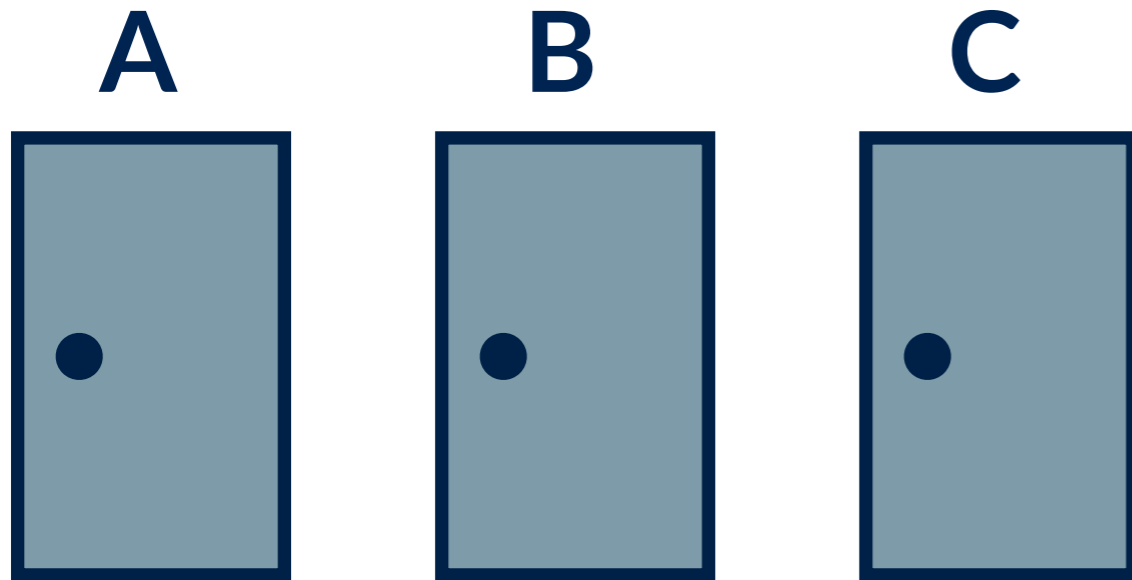
# Bayes' Rule

$A_1, A_2, \dots, A_n$  are disjoint events forming a **partition** of the sample space and  $P(A_i) > 0, \forall A_i$ . Then, for any event  $B, P(B) > 0$ , Bayes' rule states:

$$\begin{aligned} P(A_i|B) &= \frac{P(A_i \cap B)}{P(B)} = \frac{P(A_i)P(B|A_i)}{P(B)} \\ &= \frac{P(A_i)P(B|A_i)}{P(A_1 \cap B) + \dots + P(A_n \cap B)} \\ &= \frac{P(A_i)P(B|A_i)}{P(A_1)P(B|A_1) + \dots + P(A_n)P(B|A_n)} \end{aligned}$$

Note: For random variables, we often write  $P(X, Y)$ , instead of  $P(X \cap Y)$

# Monte Hall Problem & Application of Bayes' Rule

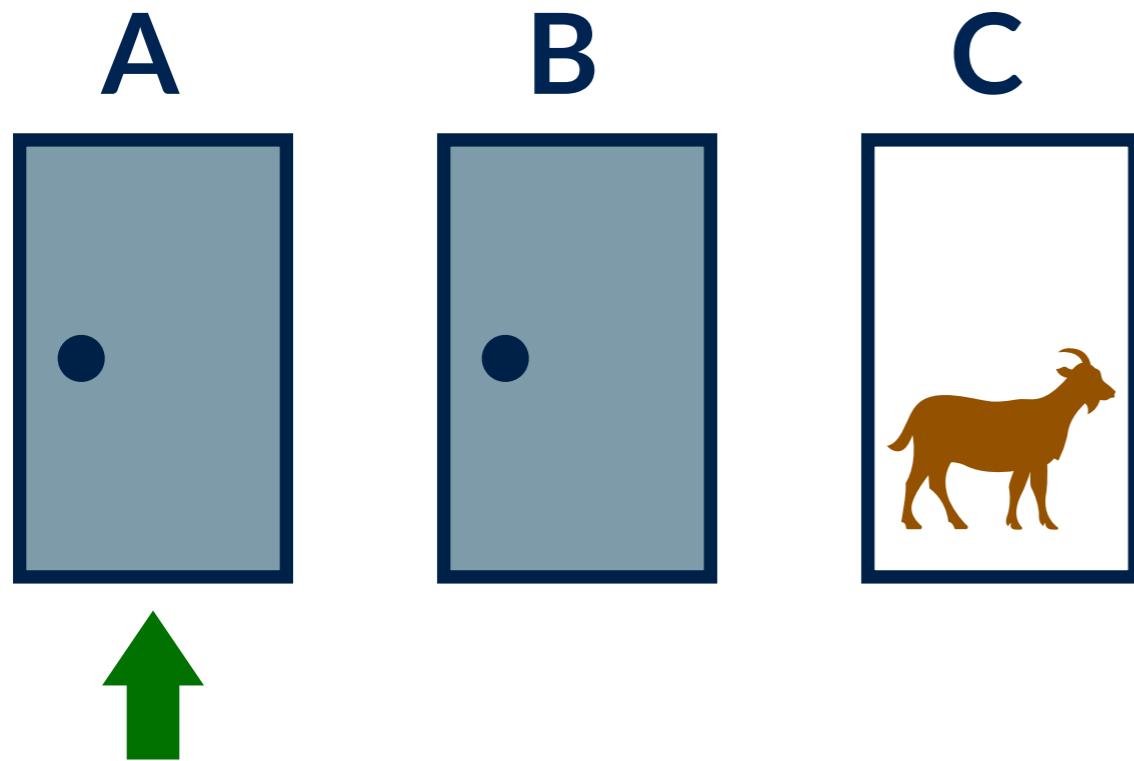


$X$  = Door chosen by player

$Y$  = Door hiding the car

$Z$  = Door opened by host

# Monte Hall Problem & Application of Bayes' Rule



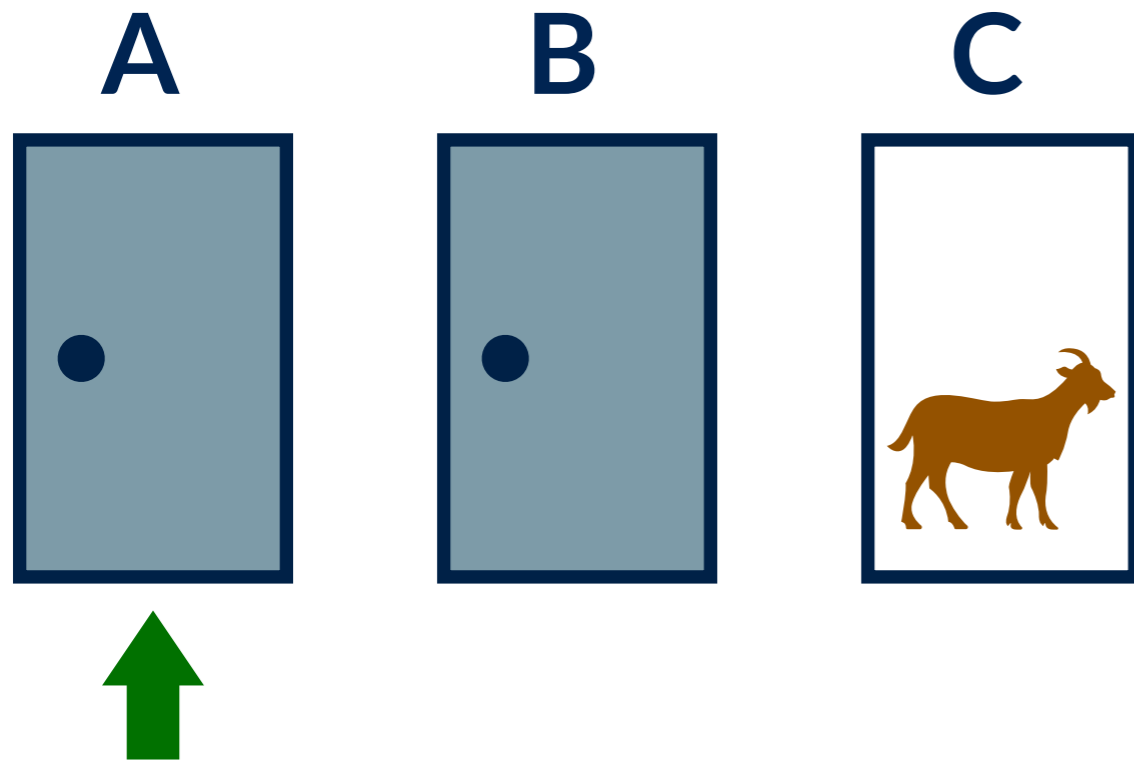
$X$  = Door chosen by player

$Y$  = Door hiding the car

$Z$  = Door opened by host

Prove that switching doors improves our chance of winning the car.

# Monte Hall Problem & Application of Bayes' Rule



$X$  = Door chosen by player

$Y$  = Door hiding the car

$Z$  = Door opened by host

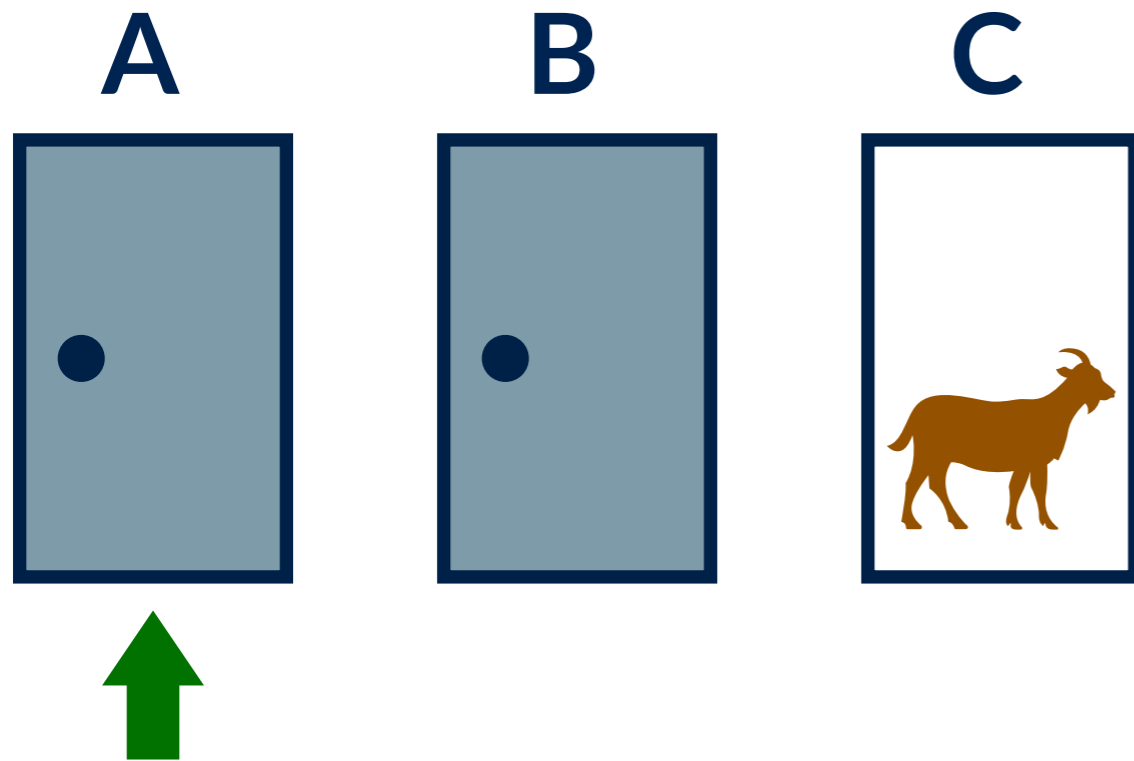
Prove that switching doors improves our chance of winning the car.

Note the assumptions:

1. The host will not open the door we have chosen
2. **The host will never open a door with a car behind**
3. Given a choice of doors, the host will choose at **random** (whilst 2)
4. Given no info, the car is equally likely to be behind any door



# Monte Hall Problem & Application of Bayes' Rule



$X$  = Door chosen by player

$Y$  = Door hiding the car

$Z$  = Door opened by host

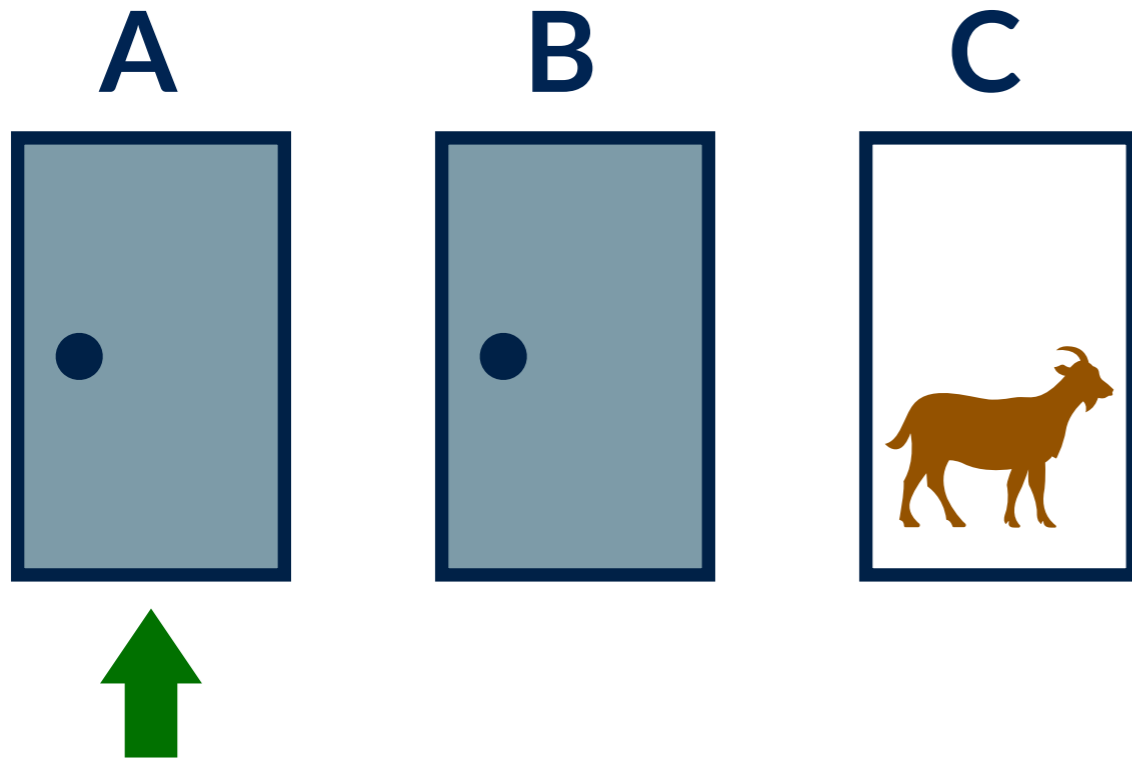
Prove that switching doors improves our chance of winning the car.

Need to show (given the we have selected A and host has shown us C):

$$P(Y = A | X = A, Z = C) < P(Y = B | X = A, Z = C)$$

Is the car more likely to be behind B than A, i.e. switching improves our chance.

# Monte Hall Problem & Application of Bayes' Rule



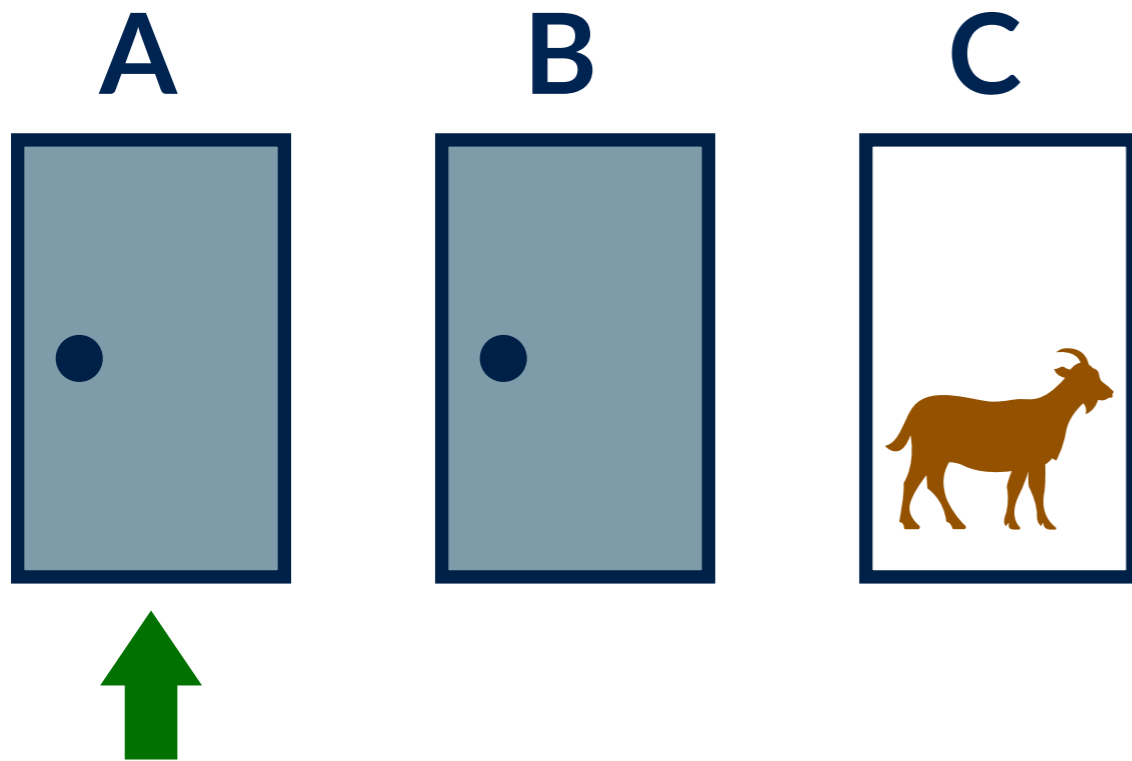
$X$  = Door chosen by player

$Y$  = Door hiding the car

$Z$  = Door opened by host

$$P(Y = A | X = A, Z = C) = \frac{P(Z = C | X = A, Y = A)P(Y = A | X = A)}{P(Z = C | X = A)}$$

# Monte Hall Problem & Application of Bayes' Rule



$X$  = Door chosen by player

$Y$  = Door hiding the car

$Z$  = Door opened by host

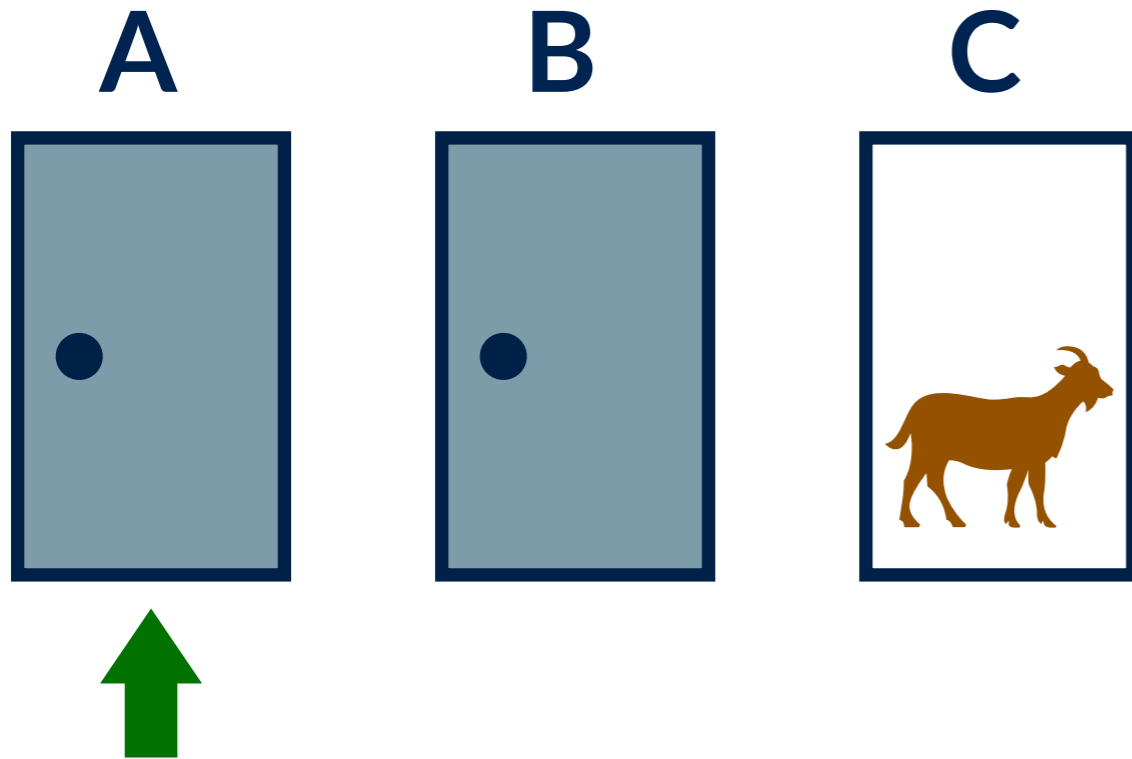
$$P(Y = A | X = A, Z = C) = \frac{P(Z = C | X = A, Y = A) P(Y = A | X = A)}{P(Z = C | X = A)}$$

1/2

Given we choose A ( $X=A$ ), and the car is in A ( $Y=A$ ), then the host is allowed to choose either B or C, as neither has the car behind it.

Since the host chooses randomly (assumption 3), we get 1/2.

# Monte Hall Problem & Application of Bayes' Rule



$X$  = Door chosen by player

$Y$  = Door hiding the car

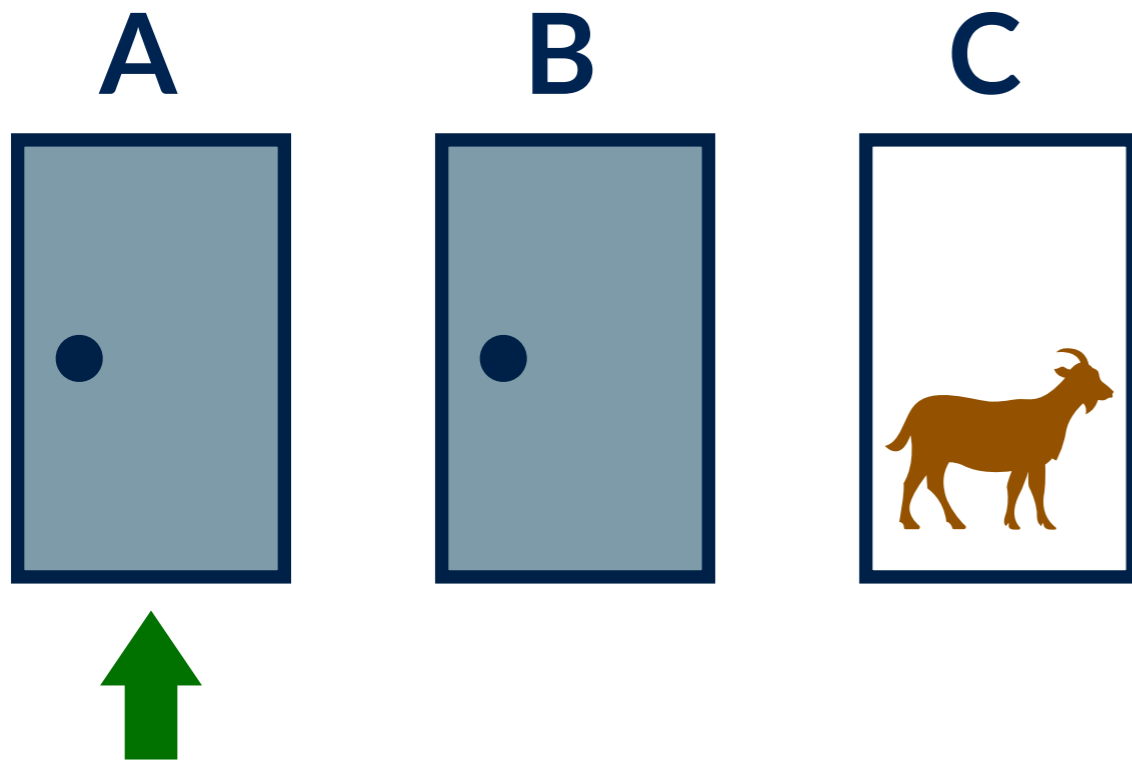
$Z$  = Door opened by host

$$P(Y = A | X = A, Z = C) = \frac{P(Z = C | X = A, Y = A) P(Y = A | X = A)}{P(Z = C | X = A)}$$

$1/3$

Given we choose A ( $X=A$ ), what is the probability that the car is behind A?  
With no further information, this is equal to  $1/3$ .

# Monte Hall Problem & Application of Bayes' Rule



$X$  = Door chosen by player

$Y$  = Door hiding the car

$Z$  = Door opened by host

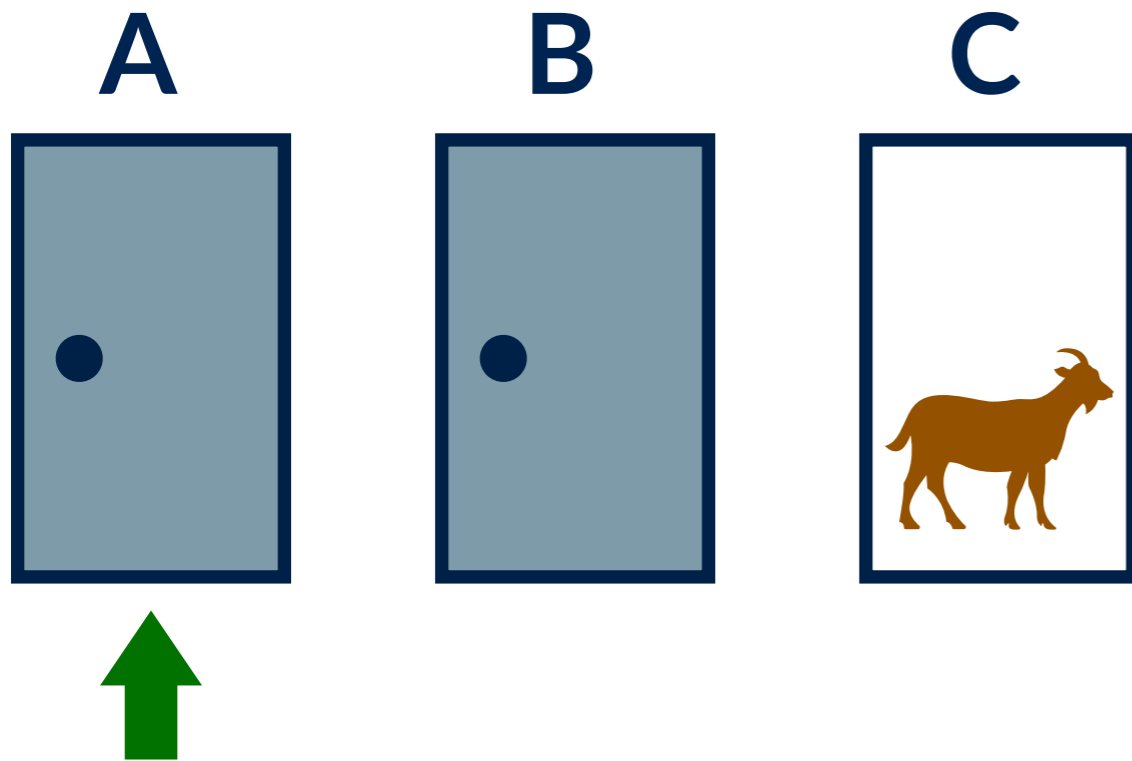
$$P(Y = A | X = A, Z = C) = \frac{P(Z = C | X = A, Y = A) P(Y = A | X = A)}{P(Z = C | X = A) \cdot 1/2}$$

Total law of prob

Product rule

$$P(Z = C | X = A) = \sum_{d=A,B,C} P(Z = C, Y = d | X = A) = \sum_{d=A,B,C} P(Z = C | X = A, Y = d) P(Y = d)$$

# Monte Hall Problem & Application of Bayes' Rule



$X$  = Door chosen by player

$Y$  = Door hiding the car

$Z$  = Door opened by host

$$P(Y = A | X = A, Z = C) = \frac{P(Z = C | X = A, Y = A) P(Y = A | X = A)}{P(Z = C | X = A) \cdot \frac{1}{2}}$$

Total law of prob

Product rule

$$P(Z = C | X = A) = \sum_{d=A,B,C} P(Z = C, Y = d | X = A) = \sum_{d=A,B,C} P(Z = C | X = A, Y = d) P(Y = d)$$

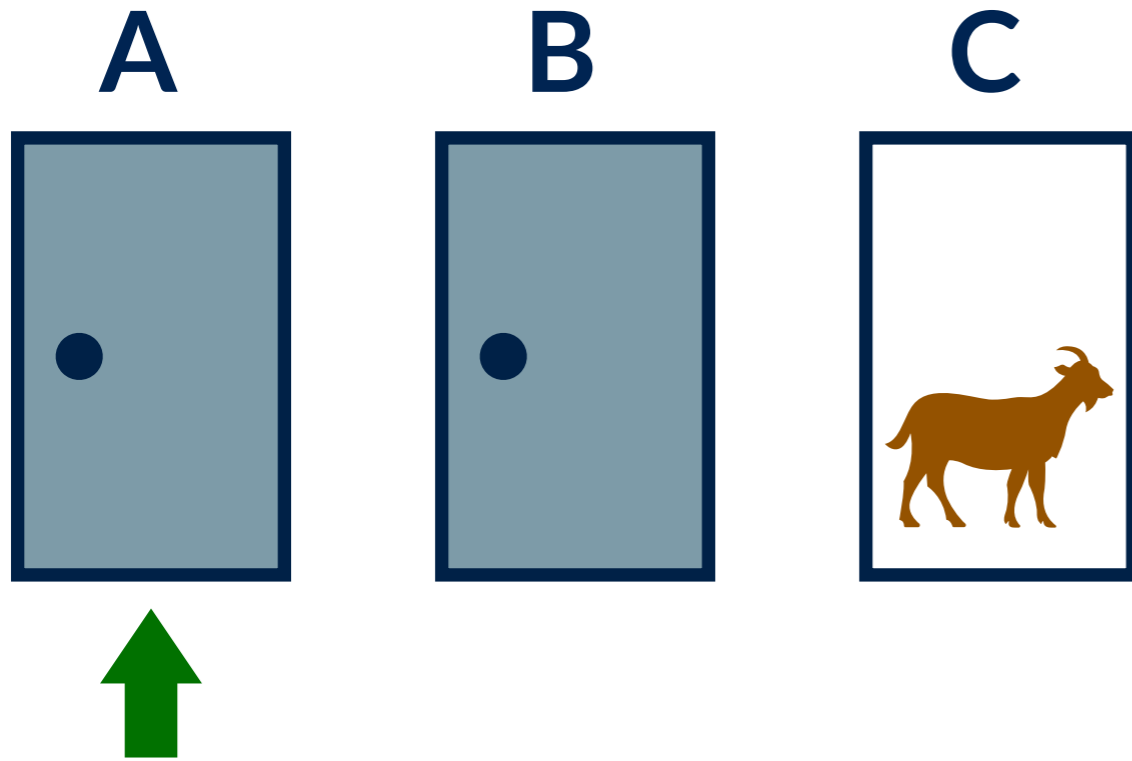
$$= \frac{1}{3} \left( P(Z = C | X = A, Y = A) + P(Z = C | X = A, Y = B) + P(Z = C | X = A, Y = C) \right)$$

1/2 as above

1: Given we chose A and car is behind B, host is **forced** to choose C (Assumption 2)

0: Given we chose A and car is behind C, the host cannot choose C (Assumption 2)

# Monte Hall Problem & Application of Bayes' Rule



$X$  = Door chosen by player

$Y$  = Door hiding the car

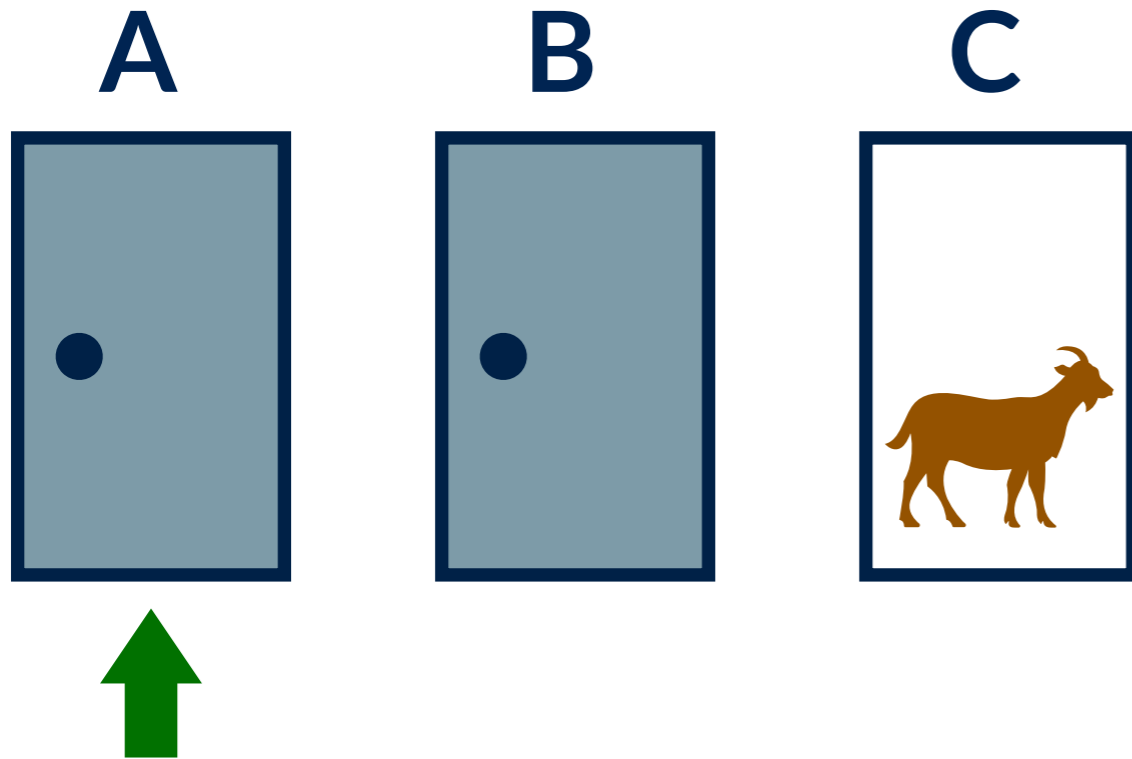
$Z$  = Door opened by host

$1/2$

$1/3$

$$P(Y = A | X = A, Z = C) = \frac{P(Z = C | X = A, Y = A) P(Y = A | X = A)}{P(Z = C | X = A) \quad 1/2}$$

# Monte Hall Problem & Application of Bayes' Rule



$X$  = Door chosen by player

$Y$  = Door hiding the car

$Z$  = Door opened by host

**Importance:** Incorporating knowledge about the process that generated the data.  
The first step towards **causal inference**.

'Host could have opened', 'he was forced to open', 'randomly opened', 'about to open', ...



# Independence

X and Y are independent events:  $P(X,Y) = P(X)P(Y)$

Equivalently:  $P(X|Y) = P(X)$  (where  $P(Y)$  is non-zero, otherwise  $P(X|Y)$  not defined)

Conditional independence:  $P(X,Y|Z) = P(X|Z)P(Y|Z)$

Equivalently:  $P(X|Y,Z) = P(X|Z)$  (again, for  $P(Y,Z)$  non-zero)

Independence of several events:

**Remark:** Pairwise independence does not imply independence

Example: 2 independent fair coin tosses ( $p_1, p_2 = 0.5$ )

Consider 3 events:

H1 = first coin is a head

H2 = second coin is a head

J = the two tosses have the same results

	H	T
H	X	
T		X

# Independence

X and Y are independent events:  $P(X,Y) = P(X)P(Y)$

Equivalently:  $P(X|Y) = P(X)$  (where  $P(Y)$  is non-zero, otherwise  $P(X|Y)$  not defined)

Conditional independence:  $P(X,Y|Z) = P(X|Z)P(Y|Z)$

Equivalently:  $P(X|Y,Z) = P(X|Z)$  (again, for  $P(Y,Z)$  non-zero)

Independence of several events:

**Remark:** Pairwise independence does not imply independence

Example: 2 independent fair coin tosses ( $p_1, p_2 = 0.5$ )

H1 & H2: independent coin tosses

$P(H1,H2) = P(H1|H2)P(H2) = 0.5 \times 0.5 = P(H1)P(H2)$

	H	T
H	X	
T		X

# Independence

X and Y are independent events:  $P(X,Y) = P(X)P(Y)$

Equivalently:  $P(X|Y) = P(X)$  (where  $P(Y)$  is non-zero, otherwise  $P(X|Y)$  not defined)

Conditional independence:  $P(X,Y|Z) = P(X|Z)P(Y|Z)$

Equivalently:  $P(X|Y,Z) = P(X|Z)$  (again, for  $P(Y,Z)$  non-zero)

Independence of several events:

**Remark:** Pairwise independence does not imply independence

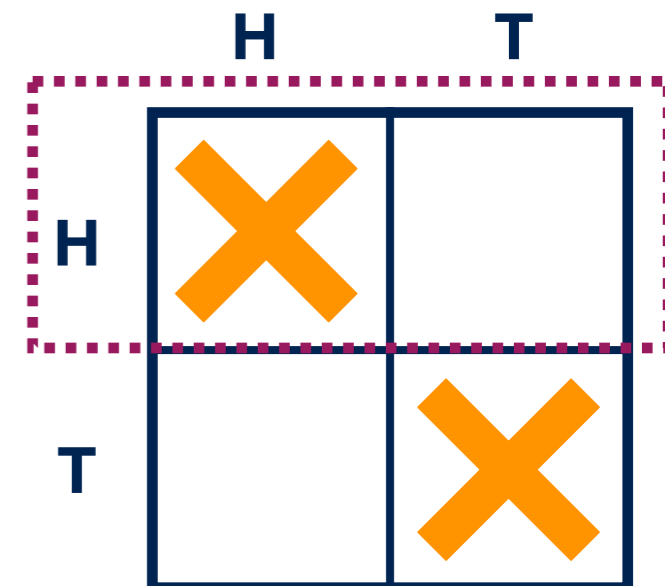
Example: 2 independent fair coin tosses ( $p_1, p_2 = 0.5$ )

H1 & H2: independent coin tosses

$P(H1,J) = P(J | H1)P(H1) =$

Given H1, what is the probability of J  
(i.e second toss also being a head)

So:  $P(J | H1) = 0.5$



# Independence

X and Y are independent events:  $P(X,Y) = P(X)P(Y)$

Equivalently:  $P(X|Y) = P(X)$  (where  $P(Y)$  is non-zero, otherwise  $P(X|Y)$  not defined)

Conditional independence:  $P(X,Y|Z) = P(X|Z)P(Y|Z)$

Equivalently:  $P(X|Y,Z) = P(X|Z)$  (again, for  $P(Y,Z)$  non-zero)

Independence of several events:

**Remark:** Pairwise independence does not imply independence

Example: 2 independent fair coin tosses ( $p_1, p_2 = 0.5$ )

H1 & H2: independent coin tosses

$P(H1,J) = P(J | H1)P(H1) = 0.5 \times 0.5 = P(J)P(H1)$

Given H1, what is the probability of J

(i.e second toss also being a head)

So:  $P(J | H1) = 0.5$

	H	T
H	X	
T		X

# Independence

X and Y are independent events:  $P(X,Y) = P(X)P(Y)$

Equivalently:  $P(X|Y) = P(X)$  (where  $P(Y)$  is non-zero, otherwise  $P(X|Y)$  not defined)

Conditional independence:  $P(X,Y|Z) = P(X|Z)P(Y|Z)$

Equivalently:  $P(X|Y,Z) = P(X|Z)$  (again, for  $P(Y,Z)$  non-zero)

Independence of several events:

**Remark:** Pairwise independence does not imply independence

Example: 2 independent fair coin tosses ( $p_1, p_2 = 0.5$ )

H1 & H2: independent coin tosses

$P(H_2, J) = P(J | H_2)P(H_2) = 0.5 \times 0.5 = P(J)P(H_2)$

So pair-wise independent. BUT ...

	H	T
H	X	
T		X

# Independence

X and Y are independent events:  $P(X,Y) = P(X)P(Y)$

Equivalently:  $P(X|Y) = P(X)$  (where  $P(Y)$  is non-zero, otherwise  $P(X|Y)$  not defined)

Conditional independence:  $P(X,Y|Z) = P(X|Z)P(Y|Z)$

Equivalently:  $P(X|Y,Z) = P(X|Z)$  (again, for  $P(Y,Z)$  non-zero)

Independence of several events:

**Remark:** Pairwise independence does not imply independence

Example: 2 independent fair coin tosses ( $p_1, p_2 = 0.5$ )

H1 & H2: independent coin tosses

$P(H1,H2,J) = P(H1 | H2,J) P(H2,J) = 1 \times 0.25 = 0.25$

	H	T
H	X	
T		X

# Independence

X and Y are independent events:  $P(X,Y) = P(X)P(Y)$

Equivalently:  $P(X|Y) = P(X)$  (where  $P(Y)$  is non-zero, otherwise  $P(X|Y)$  not defined)

Conditional independence:  $P(X,Y|Z) = P(X|Z)P(Y|Z)$

Equivalently:  $P(X|Y,Z) = P(X|Z)$  (again, for  $P(Y,Z)$  non-zero)

Independence of several events:

**Remark:** Pairwise independence does not imply independence

Example: 2 independent fair coin tosses ( $p_1, p_2 = 0.5$ )

H1 & H2: independent coin tosses

$P(H1,H2,J) = P(H1 | H2,J) P(H2,J) = 1 \times 0.25 = 0.25$

However,  $P(H1)P(H2)P(J)=0.5 \times 0.5 \times 0.5=0.125 \neq$

i.e. not jointly independent

	H	T
H	X	
T		X

# Expected values

The probability distribution of a random variable  $X$  provides us with probabilities of all possible values of  $X$ .

Summarise information, with some loss of information, represented by:

The **expected value** or **mean**:

$$\mathbb{E}[X] = \sum_x x P(X = x)$$

For a dice:  $(1 \times 1/6) + (2 \times 1/6) + (3 \times 1/6) + (4 \times 1/6) + (5 \times 1/6) + (6 \times 1/6) = 3.5$



# Expected values

The probability distribution of a random variable  $X$  provides us with probabilities of all possible values of  $X$ .

Summarise information, with some loss of information, represented by:

The **expected value** or **mean**:

$$\mathbb{E}[X] = \sum_x x P(X = x)$$

For a dice:  $(1 \times 1/6) + (2 \times 1/6) + (3 \times 1/6) + (4 \times 1/6) + (5 \times 1/6) + (6 \times 1/6) = 3.5$

The expected value of any function of  $X$ , e.g.  $g(x)$ :

$$\mathbb{E}[g(X)] = \sum_x g(x) P(X = x)$$

Dice:  $(1 \times 1/6) + (4 \times 1/6) + (9 \times 1/6) + (16 \times 1/6) + (25 \times 1/6) + (36 \times 1/6) = 15.17$

# Expected values

The probability distribution of a random variable  $X$  provides us with probabilities of all possible values of  $X$ .

Summarise information, with some loss of information, represented by:

The **expected value** or **mean**:

$$\mathbb{E}[X] = \int x P(x) dx$$

for a continuous variable  $X$ .

# Variance

The **variance** of a random variable  $X$ , denoted  $\text{Var}(X)$  or  $\sigma_X^2$  :

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

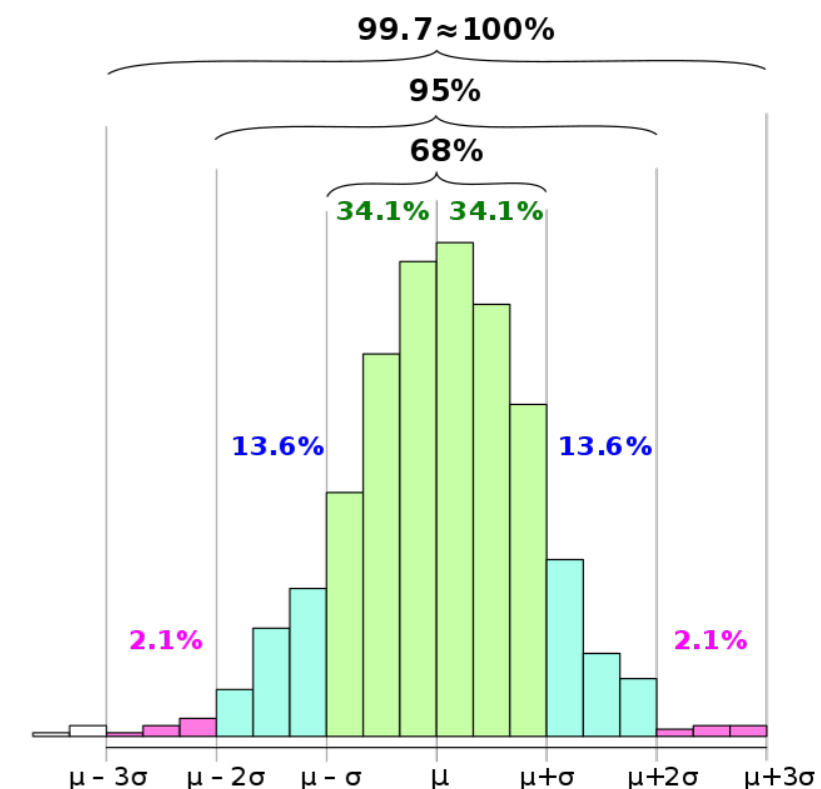
and can be calculated as

$$\text{var}(X) = \sum_x (X - \mathbb{E}[X])^2 p_X(x)$$

(Integral of continuous variables ), and measure how “spread out” the values of  $X$  in a data set are relative to their mean.

The **standard deviation**  $\sigma_X$ , (has the same units as  $X$ ).

For a normal distribution,  $\sim 2/3$  of the population values of  $X$  fall within one  $\sigma_X$ , 95% fall between  $2\sigma_X$ , etc.



# Covariance

The degree to which two random variables  $X$  and  $Y$  co-vary (degree associated):

$$\sigma_{XY} = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

and measures a specific way  $X$  and  $Y$  co-vary, i.e., linearly. When normalised, it yields the correlation coefficient (**Pearson correlation**):

$$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

a dimensionless quantity between -1 and 1.

# Covariance

The degree to which two random variables  $X$  and  $Y$  co-vary (degree associated):

$$\sigma_{XY} = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

and measures a specific way  $X$  and  $Y$  co-vary, i.e., linearly. When normalised, it yields the correlation coefficient (**Pearson correlation**):

$$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

a dimensionless quantity between -1 and 1.

When  $X$  and  $Y$  are independent, then  $\rho_{XY} = 0$ .

The reverse is not true!

(e.g.  $\rho_{XY}$  may be zero, but not linear-correlation, hence dependence exists.

This requires more complex methods of demonstrating if  $P(Y|X) = P(Y)$  )

# Anscombe's Quartet

Group of 4 datasets with nearly identical simple descriptive statistical properties:

- Mean and sample variance of X
- Mean and sample variance of Y
- Correlation between X and Y
- Linear regression line (coefficient the same up to 2 or 3 decimal places)
- $R^2$  coefficient

A note on  $R^2$ : A measure for goodness-of-fit

$$R^2 = 1 - \frac{\sum_i (y_i - f_i)^2}{\sum_i (y_i - \bar{y})^2}, \quad y_i = f(x_i), \quad \bar{y} = \frac{1}{n} \sum_i y_i$$

If the fit  $y=f(x)$  is a perfect fit, the numerator is zero,  $R^2 = 1$ , and  $R^2 = 0$  implies the fit  $f(x)$  is no better than baseline average  $\bar{y}$ .

Negative values corresponds to models worse than the baseline average.

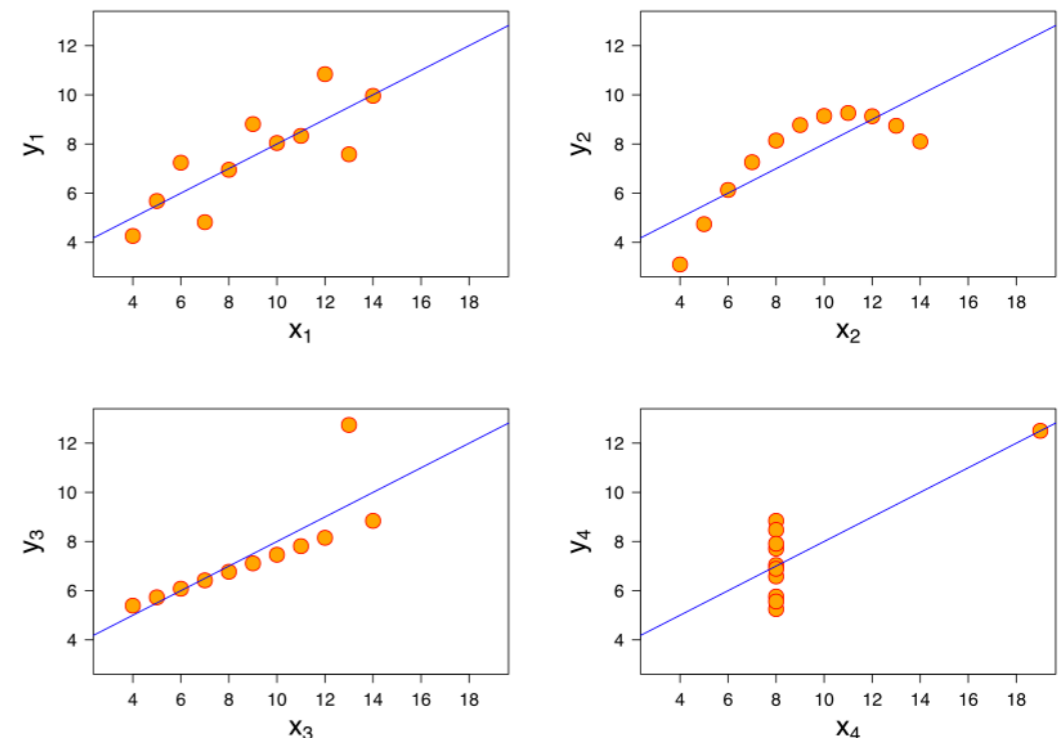
# Anscombe's Quartet

Group of 4 datasets with nearly identical simple descriptive statistical properties:

- Mean and sample variance of  $X$
- Mean and sample variance of  $Y$
- Correlation between  $X$  and  $Y$
- Linear regression line (coefficient the same up to 2 or 3 decimal places)
- $R^2$  coefficient

Yet, very different distributions, which can be observed by plotting the graphs

Same Pearson correlation, but,  
different dependence structure  
( $X$  causes  $Y$ , but in different ways)



# Next time

First of the two causal frameworks:

- Potential Outcomes (due to Neyman-Rubin)
- Study our first causal question

Next, we estimate:

- Answer to causal question
- Uncertainty on this answer (under model assumptions)





THE UNIVERSITY  
*of* EDINBURGH

# Methods for Causal Inference

## Lecture 2: Basics of probability

---

Ava Khamseh

School of Informatics  
2023-2024