# Methods for Causal Inference
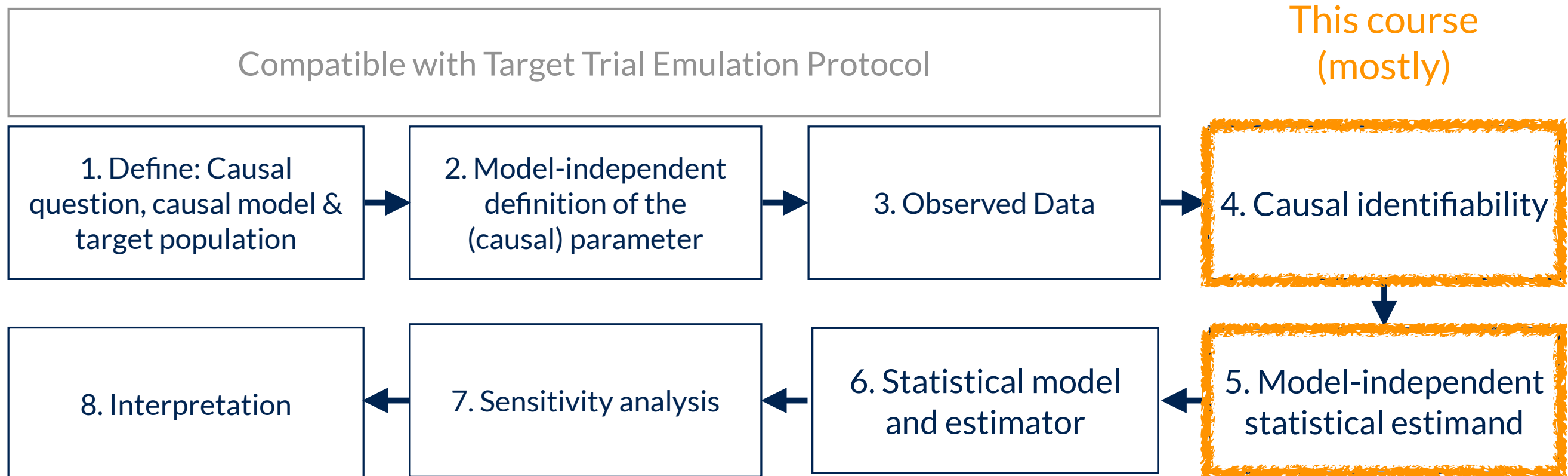# Lecture 2: Basics of probability

Ava Khamseh

School of Informatics
2025-2026

# The Causal Roadmap

Compatible with Target Trial Emulation Protocol

This course (mostly)

1. Define: Causal question, causal model & target population

2. Model-independent definition of the (causal) parameter

3. Observed Data

4. Causal identifiability

5. Model-independent statistical estimand

6. Statistical model and estimator

7. Sensitivity analysis

8. Interpretation

# The Causal Roadmap (summary of steps 1-3)

**Step 1:** Specifying the **causal question(s)** of interest, the study's target population, the exposure(s), outcomes(s), time period and context, as well as the **causal model** (developed with input from subject experts) to describe relationships amongst the variables.

**Step 2:** Define of the causal estimand, which is a mathematical quantity representing the answer to the causal question(s)

**Step 3:** Observed data, assesses questions such as what is the baseline time zero, are the required variables in the causal model of step 1 measured, or measured differently in various data sources. If these essential ingredients are not recorded in the observed data, the causal question in the first step may have to be adapted accordingly.
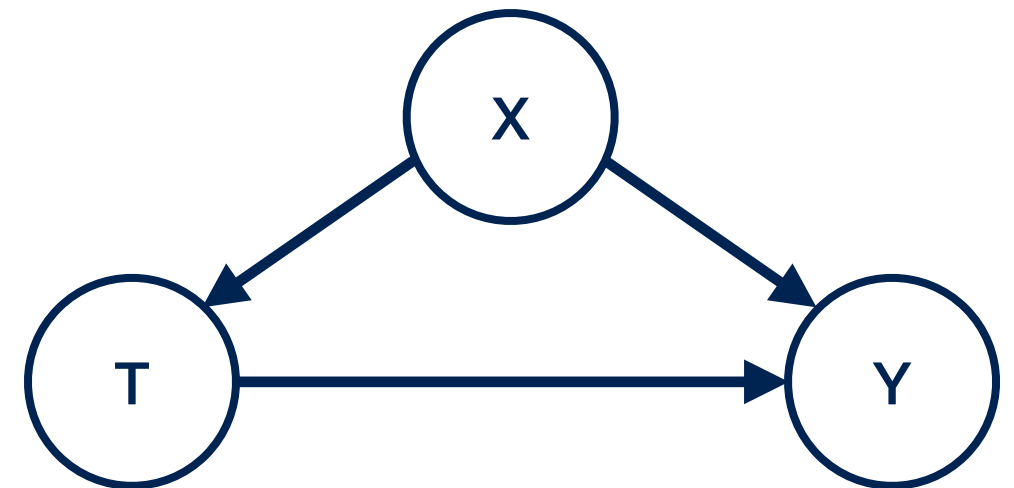
# Causal theory and data

**Disclaimer**: In this course our focus is on steps 4-5. We then use simple models to exemplify 6-7 (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).

# Two main Frameworks for causal identifiablity

- Potential outcomes framework (Neyman-Rubin):

  - Requires a given treatment-outcome pair (known directionality)
  - For causal estimation
  - More familiar to biomedical researchers (this is changing ...)

- Structural causal models (Pearl):

  - Causal graphs
  - Structural equations $x = f_x(\epsilon_x), \ t = f_t(x, \epsilon_t), \ y = f_y(x, t, \epsilon_y)$
  - Algorithmic
  - For causal estimation and discovery

  **Assumption: Independent noise terms:** $\epsilon_x \perp\!\!\!\perp \epsilon_t \perp\!\!\!\perp \epsilon_y$

# Overview of the course

- **Lecture 1:** Introduction & Motivation, why do we care about causality? Why deriving causality from observational data is non-trivial.
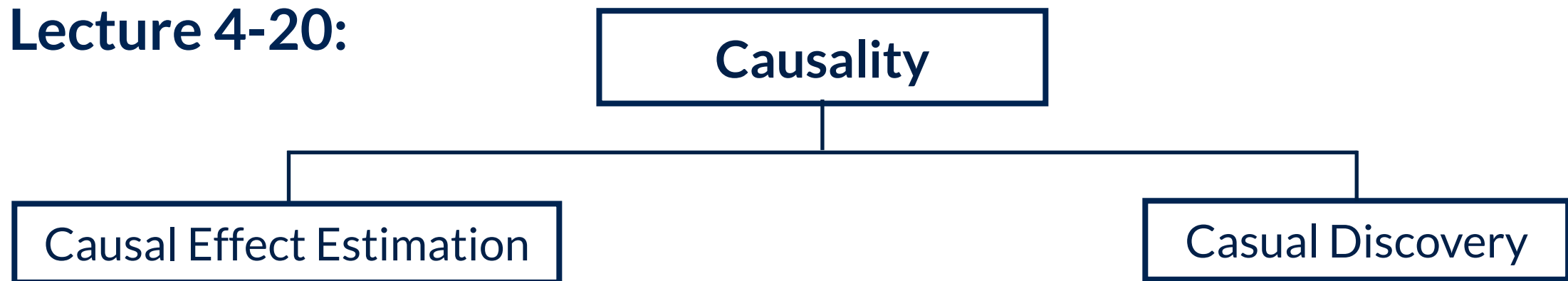
# Overview of the course

- **Lecture 1:** Introduction & Motivation, why do we care about causality? Why deriving causality from observational data is non-trivial.
- **Lecture 2:** Recap of probability theory, variables, events, conditional probabilities, independence, law of total probability, Bayes' rule
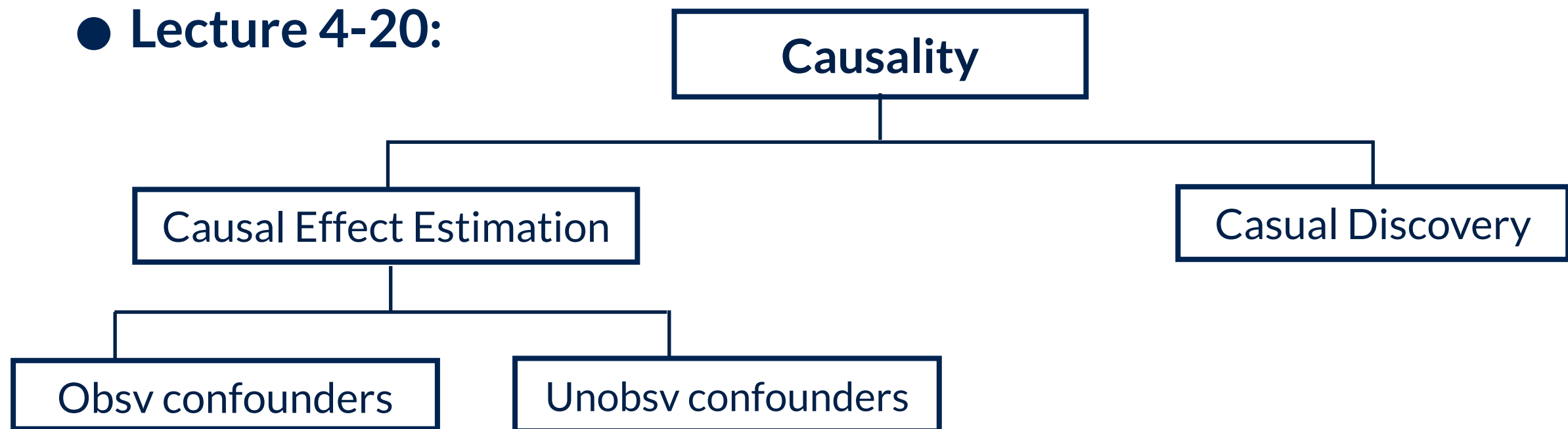
# Overview of the course

- **Lecture 1**: Introduction & Motivation, why do we care about causality? Why deriving causality from observational data is non-trivial.
- **Lecture 2:** Recap of probability theory, variables, events, conditional probabilities, independence, law of total probability, Bayes' rule
- **Lecture 3**: Recap of regression, multiple regression, graphs, SCM

# Overview of the course

- **Lecture 1**: Introduction & Motivation, why do we care about causality? Why deriving causality from observational data is non-trivial.
- **Lecture 2:** Recap of probability theory, variables, events, conditional probabilities, independence, law of total probability, Bayes' rule
- **Lecture 3**: Recap of regression, multiple regression, graphs, SCM
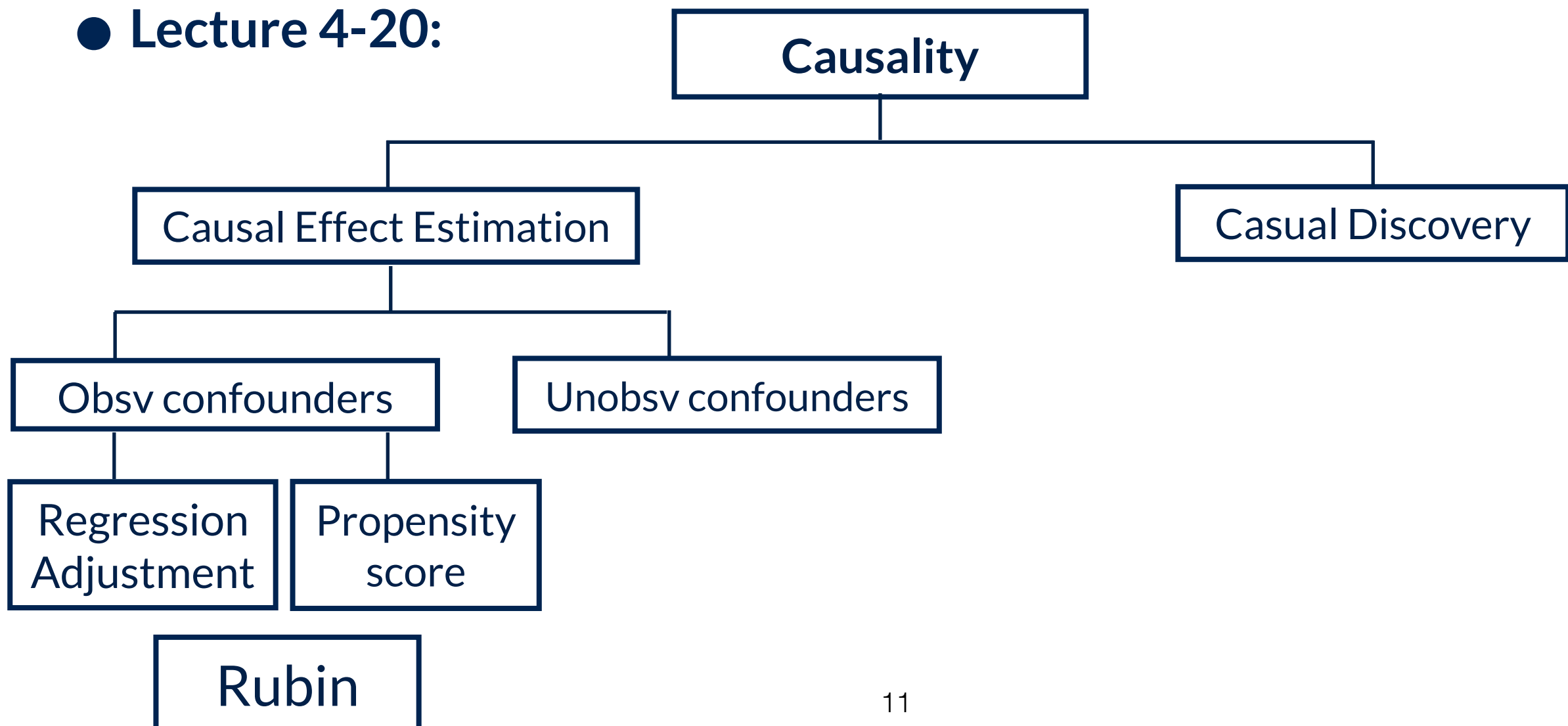- **Lecture 4-20:**

# Overview of the course

- **Lecture 1**: Introduction & Motivation, why do we care about causality? Why deriving causality from observational data is non-trivial.
- **Lecture 2:** Recap of probability theory, variables, events, conditional probabilities, independence, law of total probability, Bayes' rule
- **Lecture 3**: Recap of regression, multiple regression, graphs, SCM
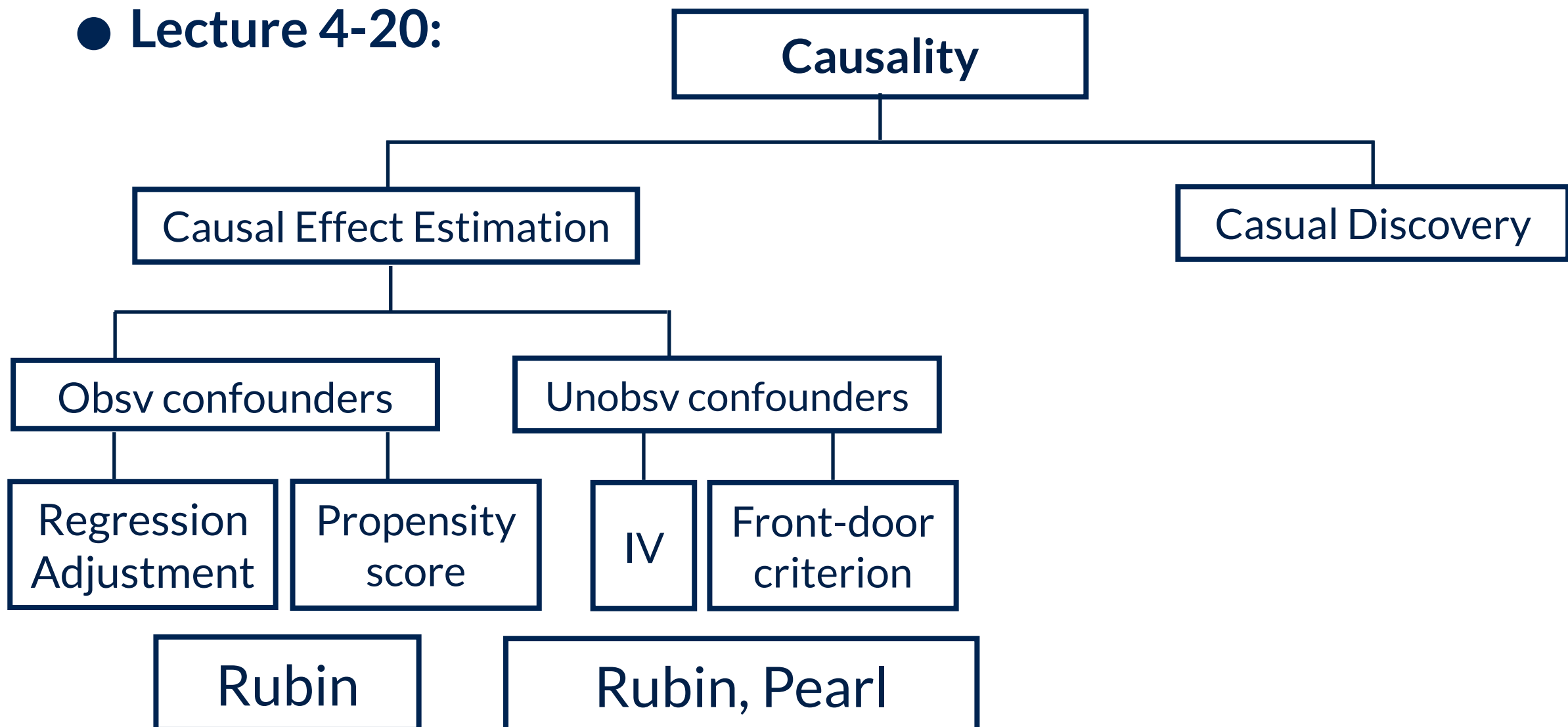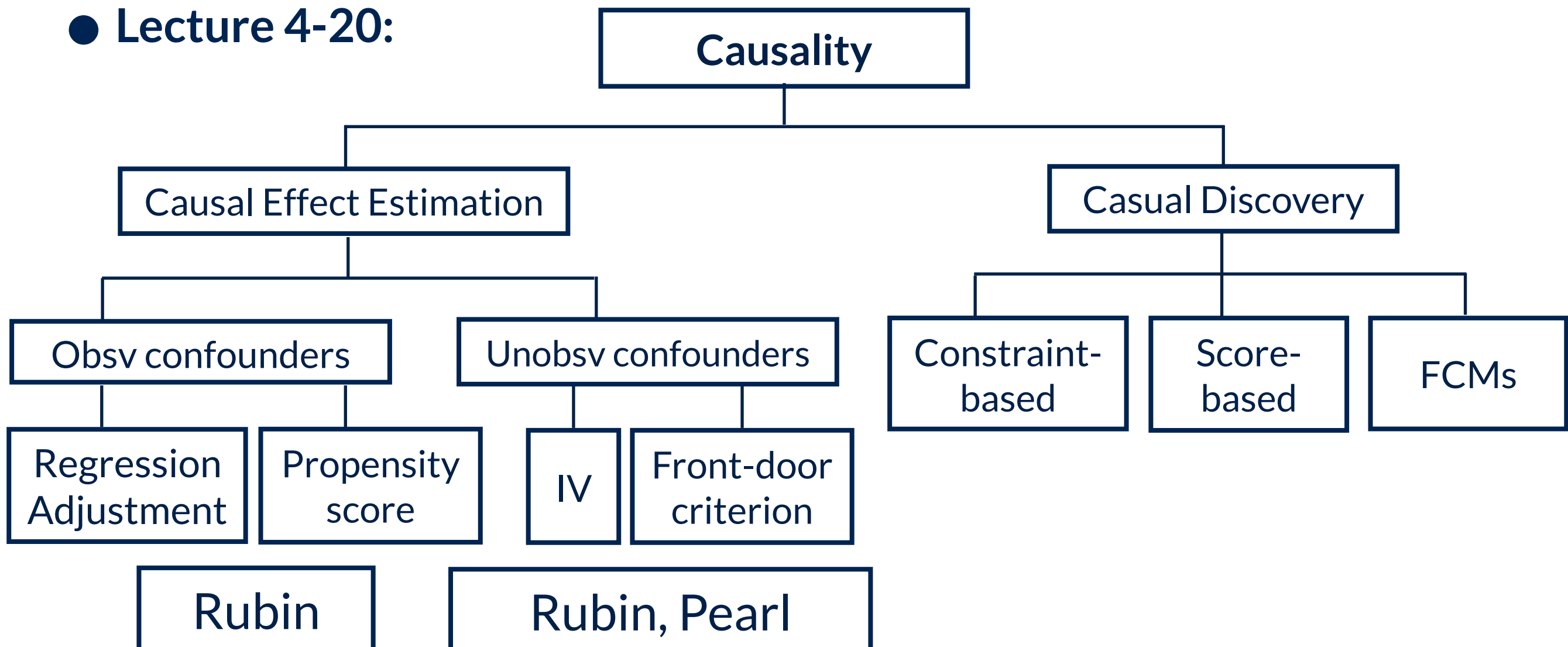- **Lecture 4-20:**

# Overview of the course

- **Lecture 1**: Introduction & Motivation, why do we care about causality? Why deriving causality from observational data is non-trivial.
- **Lecture 2:** Recap of probability theory, variables, events, conditional probabilities, independence, law of total probability, Bayes' rule
- **Lecture 3**: Recap of regression, multiple regression, graphs, SCM
- **Lecture 4-20:**

```
                        ┌──────────────┐
                        │  Causality   │
                        └──────┬───────┘
              ┌────────────────┴──────────────────┐
    ┌─────────────────────────┐        ┌────────────────────┐
    │ Causal Effect Estimation │        │ Casual Discovery   │
    └────────────┬─────────────┘        └────────────────────┘
       ┌─────────┴──────────┐
┌──────────────┐    ┌──────────────────┐
│ Obsv confounders │ │ Unobsv confounders │
└──────┬───────────┘ └──────────────────┘
   ┌───┴────────┐
┌────────────┐ ┌────────────┐
│ Regression │ │ Propensity │
│ Adjustment │ │   score    │
└────────────┘ └────────────┘
   ┌────────────┐
   │   Rubin    │
   └────────────┘
```

# Overview of the course

- **Lecture 1**: Introduction & Motivation, why do we care about causality? Why deriving causality from observational data is non-trivial.
- **Lecture 2:** Recap of probability theory, variables, events, conditional probabilities, independence, law of total probability, Bayes' rule
- **Lecture 3**: Recap of regression, multiple regression, graphs, SCM
- **Lecture 4-20:**

```
                        ┌──────────────┐
                        │  Causality   │
                        └──────┬───────┘
          ┌────────────────────┴────────────────────┐
┌──────────────────────────┐          ┌──────────────────────┐
│ Causal Effect Estimation │          │  Casual Discovery    │
└────────────┬─────────────┘          └──────────────────────┘
      ┌───────┴────────────────┐
┌──────────────────┐   ┌──────────────────────┐
│ Obsv confounders │   │  Unobsv confounders  │
└───────┬──────────┘   └──────────┬───────────┘
    ┌───┴────────┐          ┌──────┴─────────┐
┌───────────┐ ┌──────────┐ ┌────┐ ┌────────────┐
│Regression │ │Propensity│ │ IV │ │ Front-door │
│Adjustment │ │  score   │ │    │ │ criterion  │
└───────────┘ └──────────┘ └────┘ └────────────┘
      ┌────────┐              ┌──────────────┐
      │ Rubin  │              │ Rubin, Pearl │
      └────────┘              └──────────────┘
```

# Overview of the course

- **Lecture 1**: Introduction & Motivation, why do we care about causality? Why deriving causality from observational data is non-trivial.
- **Lecture 2:** Recap of probability theory, variables, events, conditional probabilities, independence, law of total probability, Bayes' rule
- **Lecture 3**: Recap of regression, multiple regression, graphs, SCM
- **Lecture 4-20:**

# Lecture 2: Recap of probability theory

# Causal theory and data

**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

$$\Omega$$

A box containing outcomes $A_1$, $A_2$, $A_3$, $A_4$, $A_5$, $A_6$.

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:

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

$\Omega$

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:

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

$\Omega$

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, …,n be mutually exclusive and exhaustive partitions of B, and let A=B (complete overlap). Then,

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

Total Law of probability. See later: "marginalisation"

# Intervals

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

Figure 7.2: Age at adoption, Scotland, 2018



**49%** of adoptions were of children aged 4 and below.

Total = 471

# Intervals

**P(age > 4)** = 1 - P(age <= 4) = 1- 0.49 = 0.51

**P( 4 < age <= 12)** = (43+30+34+25+13+14+5+12) / 471 = 0.37

**Figure 7.2: Age at adoption, Scotland, 2018**



Total = 471

# 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(lung cancer | smoker) vs

P(lung cancer | 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(lung cancer | smoker) vs
         P(lung cancer | 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:

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

# Bayes' Rule

$A_1, A_2, ..., A_n$ 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

$A_1, A_2, ..., A_n$ 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:

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

$$= \frac{P(A_i)P(B|A_i)}{P(A_1 \cap B) + \cdots + P(A_n \cap B)}$$

$$= \frac{P(A_i)P(B|A_i)}{P(A_1)P(B|A_1) + \cdots + P(A_n)P(B|A_n)}$$

# Bayes' Rule

$A_1, A_2, ..., A_n$ 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:

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

$$= \frac{P(A_i)P(B|A_i)}{P(A_1 \cap B) + \cdots + P(A_n \cap B)}$$

$$= \frac{P(A_i)P(B|A_i)}{P(A_1)P(B|A_1) + \cdots + P(A_n)P(B|A_n)}$$

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

# Monte Hall Problem & Application of Bayes' Rule

A  B  C

X = Door chosen by player

Y = Door hiding the car

Z = Door opened by host

# Monte Hall Problem & Application of Bayes' Rule

A      B      C

X = Door chosen by player
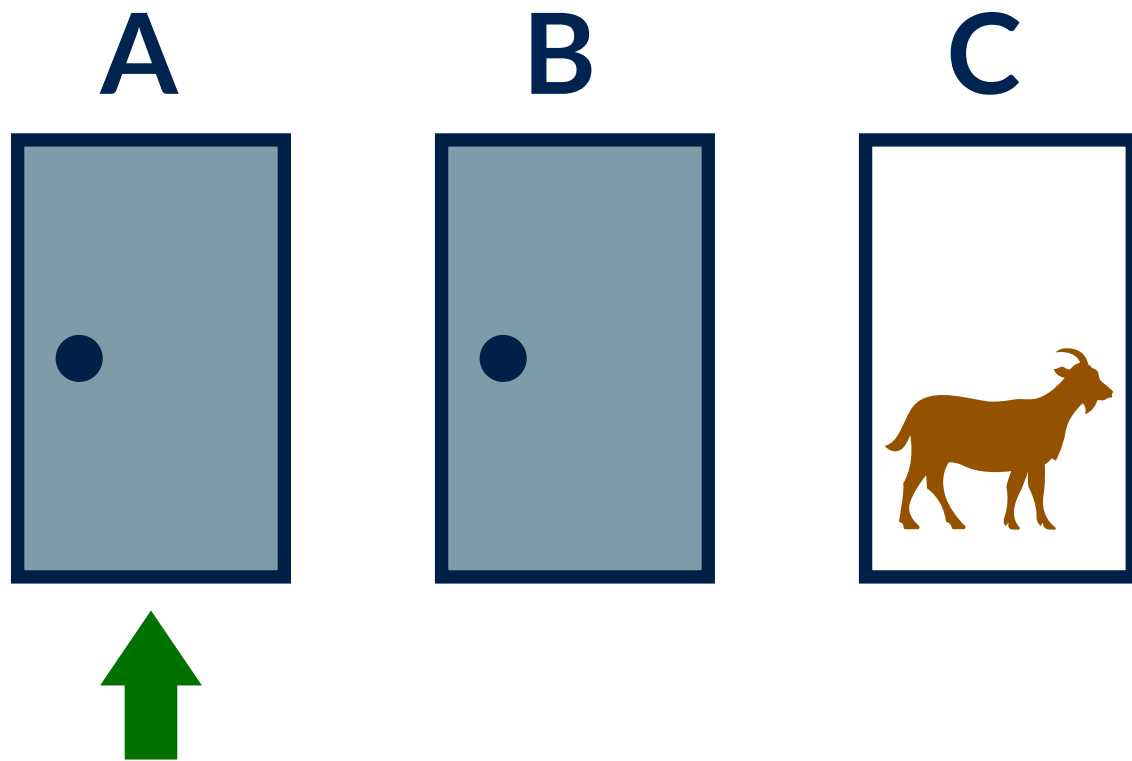
Y = Door hiding the car

Z = Door opened by host

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

# Monte Hall Problem & Application of Bayes' Rule

**A**  **B**  **C**

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

**A**   **B**   **C**

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
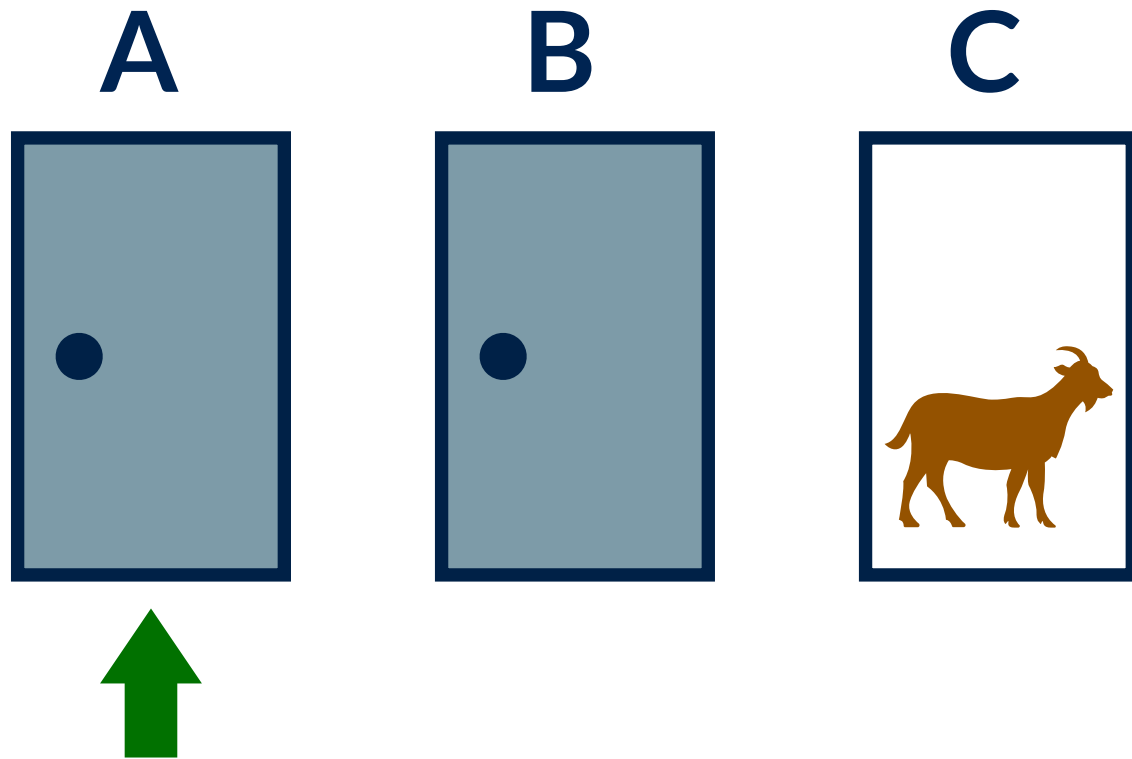
**A**  **B**  **C**

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

**A**  **B**  **C**

$X$ = Door chosen by player

$Y$ = Door hiding the car

$Z$ = Door opened by host

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

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 choses randomly (assumption 3), we get 1/2.

# Monte Hall Problem & Application of Bayes' Rule

**A**    **B**    **C**

$X$ = Door chosen by player

$Y$ = Door hiding the car

$Z$ = Door opened by host

1/3

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

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

**A**  **B**  **C**

$X$ = Door chosen by player

$Y$ = Door hiding the car

$Z$ = Door opened by host

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

**Total law of prob**

**Product rule**

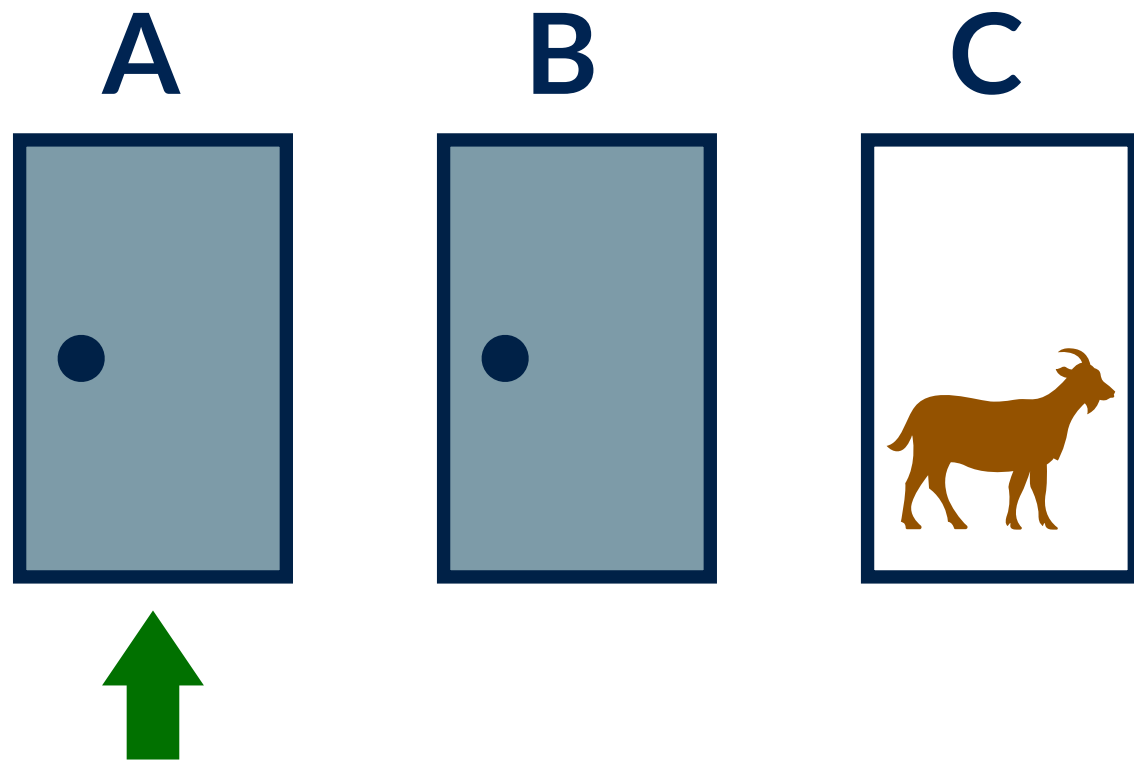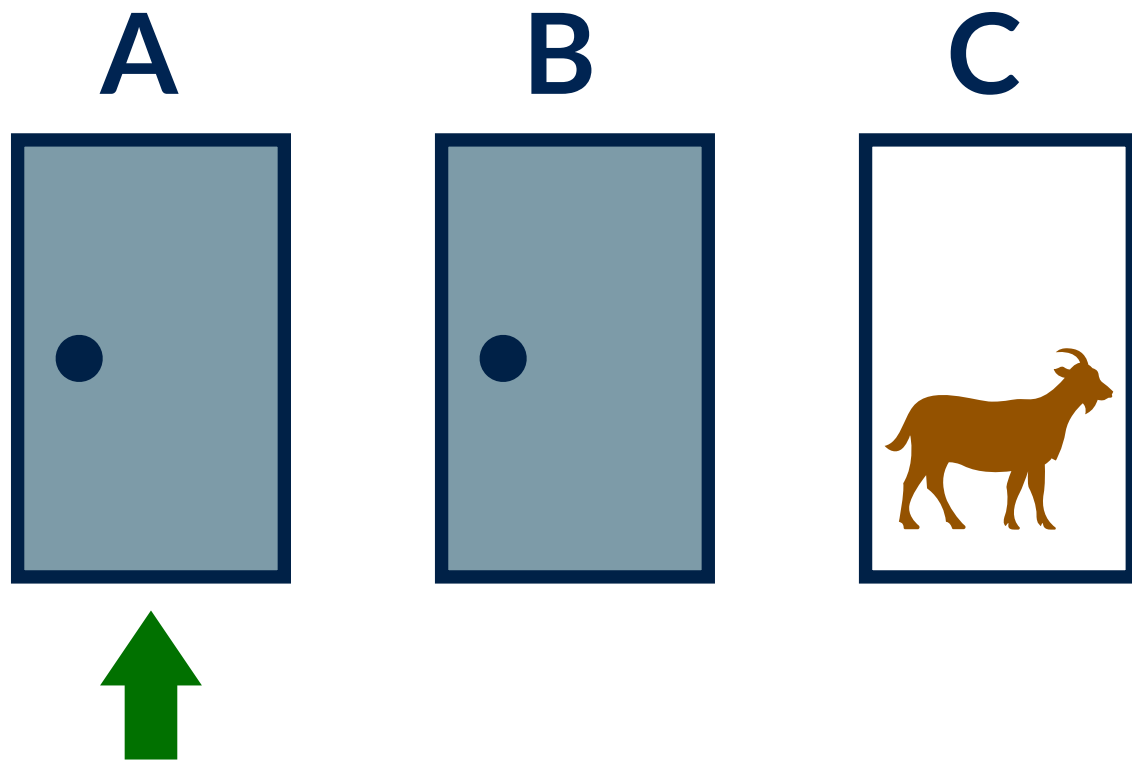$$P(Z = C | X = A) = \sum_{d=A.B.C} P(Z = C, Y = d | X = A) = \sum_{d=A.B.C} P(Z = C | X = A, Y = d)P(Y = d)$$

# Monte Hall Problem & Application of Bayes' Rule

**A**

**B**

**C**

$X$ = 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

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

**A**  **B**  **C**



X = Door chosen by player

Y = Door hiding the car

Z = Door opened by host

$$P(Y = A | X = A, Z = C) = \frac{\overset{1/2}{P(Z = C | X = A, Y = A)}\overset{1/3}{P(Y = A | X = A)}}{P(Z = C | X = A)\ \underset{1/2}{}}$$

# Monte Hall Problem & Application of Bayes' Rule

A      B      C

X = Door chosen by player
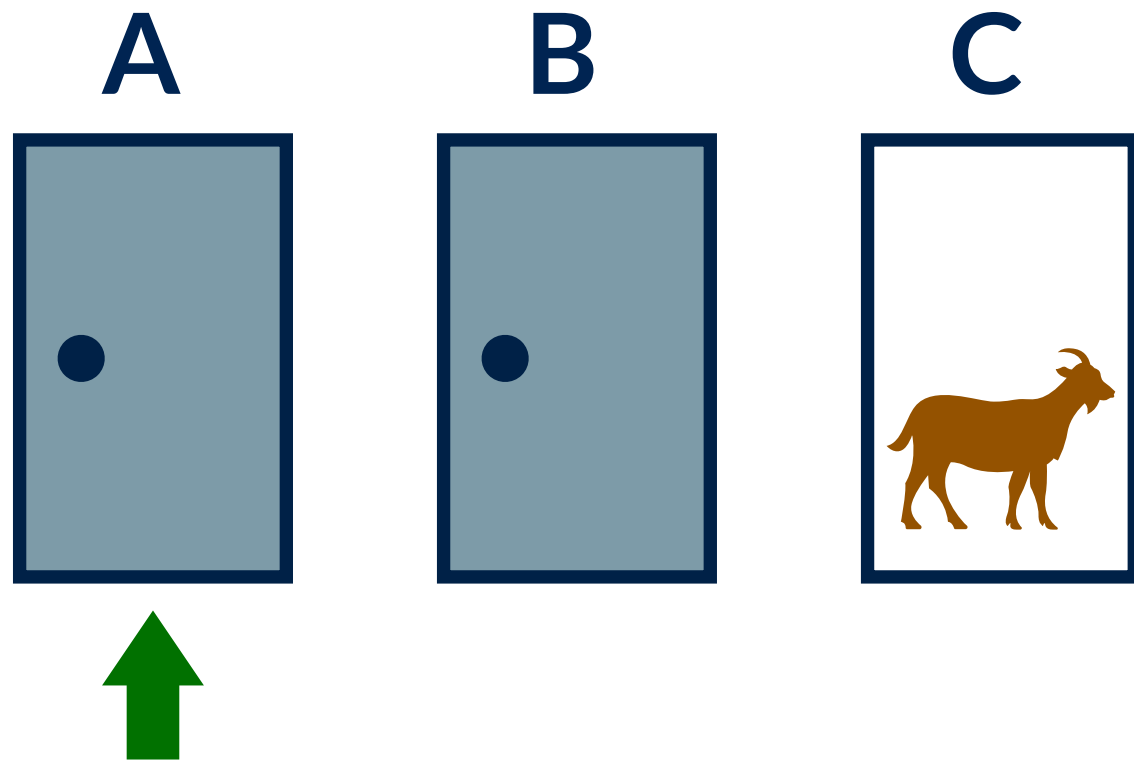
Y = Door hiding the car

Z = Door opened by host

$$P(Y = A | X = A, Z = C) = \frac{\overset{1/2}{P(Z = C | X = A, Y = A)} \overset{1/3}{P(Y = A | X = A)}}{P(Z = C | X = A) \underset{1/2}{}}$$

$$P(Y = B | X = A, Z = C) = 1 - P(Y = A | X = A, Z = C) - P(Y = C | X = A, Z = C)$$

$$= 1 - \frac{1}{3} - 0 = 2/3$$

# Monte Hall Problem & Application of Bayes' Rule

**A**  **B**  **C**

$X$ = 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 (p1, p2 = 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

# 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 (p1, p2 = 0.5)

H1 & H2: independent coin tosses

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

# 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 (p1, p2 = 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



49

# 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 (p1, p2 = 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: $P(J | H1) = 0.5$

# 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 (p1, p2 = 0.5)

H1 & H2: independent coin tosses

P(H2,J) = P(J | H2)P(H2) = 0.5 x 0.5 = P(J)P(H2)

So pair-wise independent. BUT ...

# 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 ($p1, p2 = 0.5$)

H1 & H2: independent coin tosses

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

# 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 (p1, p2 = 0.5)

H1 & H2: independent coin tosses

P(H1,H2,J) = P(H1 | H2,J) P(H2,J) = 1 x 0.25 = 0.25

However, P(H1)P(H2)P(J)=0.5x0.5x0.5=0.125  $\neq$

i.e. not jointly independent

# 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: (1x1/6) + (2x1/6) + (3x1/6) + (4x1/6) + (5x1/6) + (6x1/6) = 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: (1x1/6) + (2x1/6) + (3x1/6) + (4x1/6) + (5x1/6) + (6x1/6) = 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: (1x1/6) + (4x1/6) + (9x1/6) + (16x1/6) + (25x1/6) + (36x1/6) = 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 Var(X) or $\sigma_X^2$:
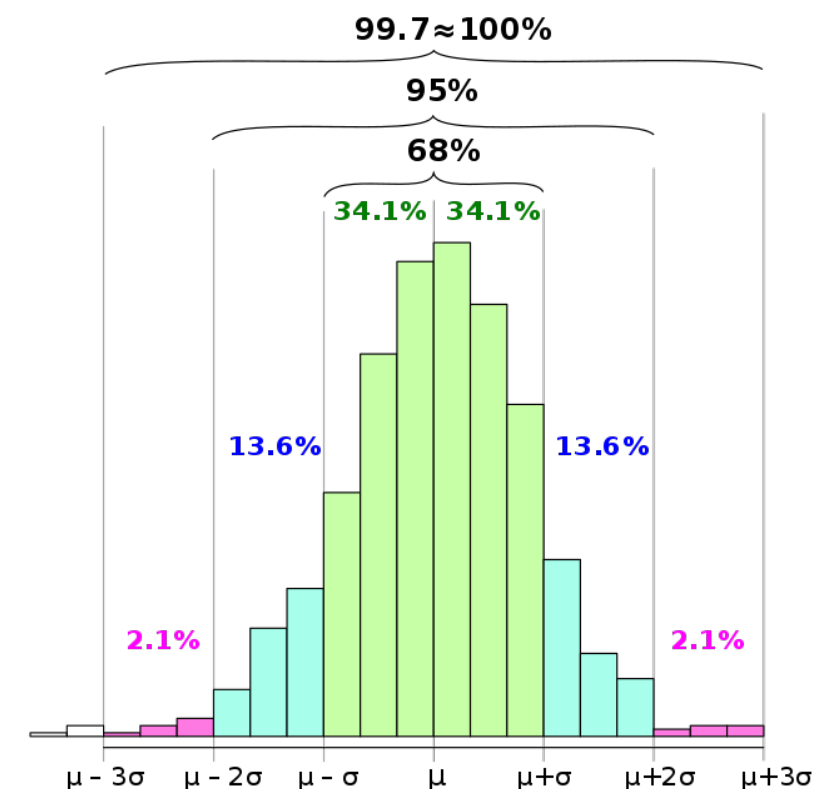
$$var(X) = \mathbb{E}[(X - \mathbb{E}[X])^2]$$

and can be calculated as

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

(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, ~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 (<span style="color:red">Pearson correlation</span>):

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

a dimensionless quantity between -1 and 1.

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

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

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 $\bar{y}$.
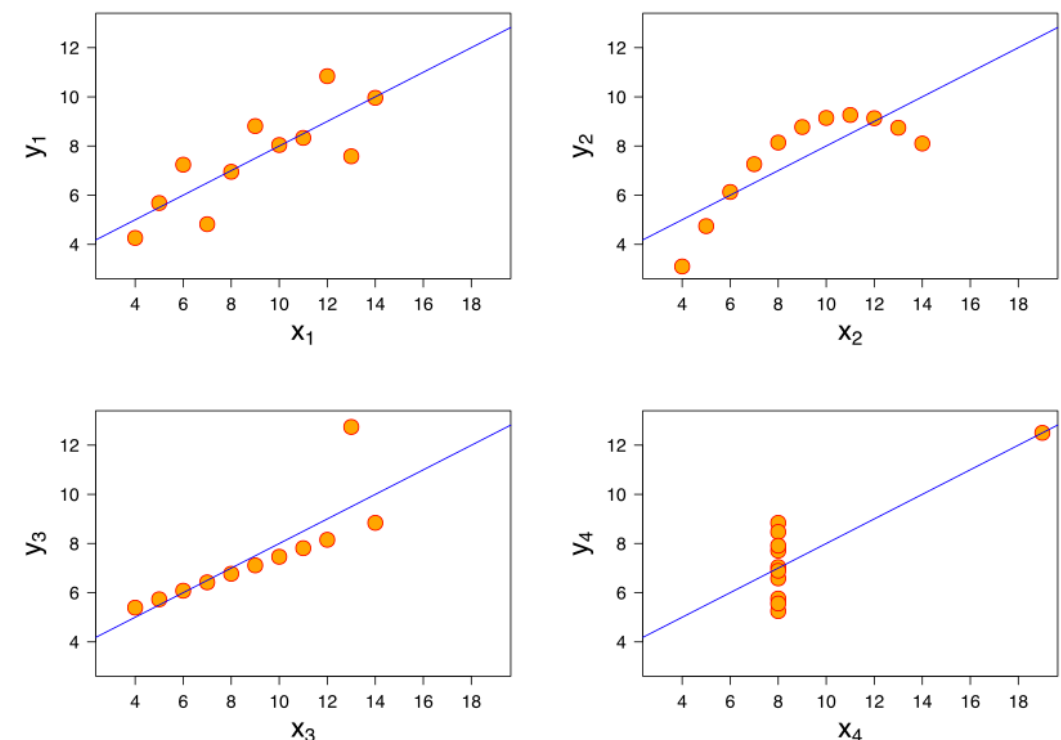Negative values corresponds to models worse than the baseline average.

# 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)