# Review of Probability Theory Part 2
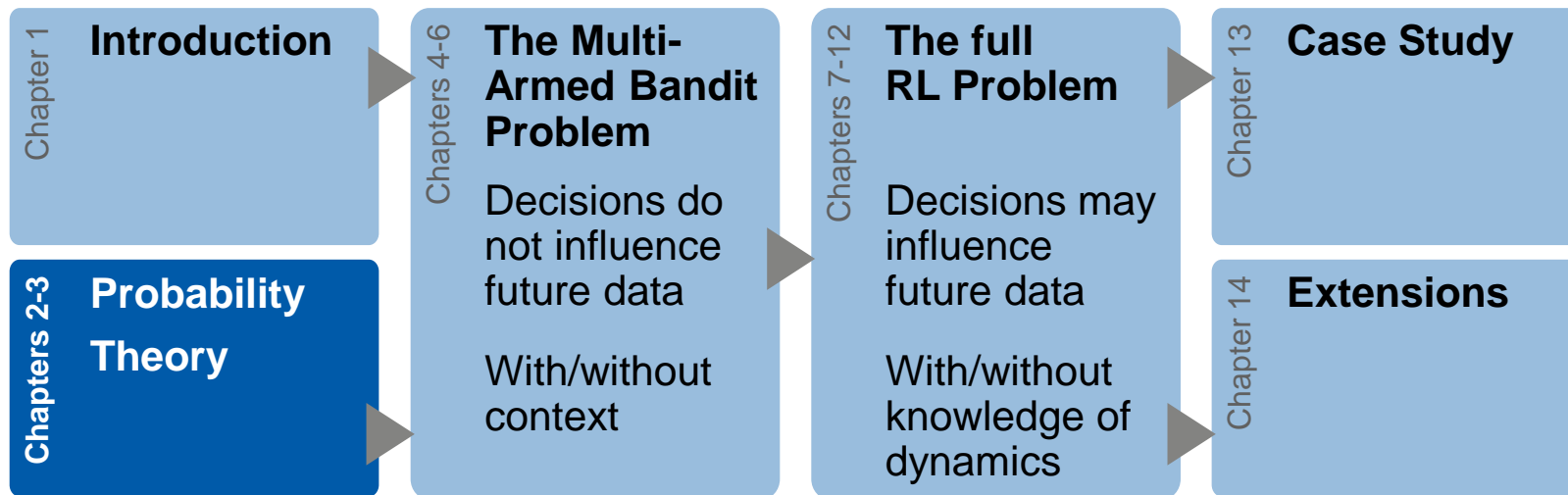
## Fundamentals of Reinforcement Learning

Institut für Nachrichtentechnik

Fachgebiet Kommunikationstechnik

Prof. Dr.-Ing. Anja Klein

Dr. Sabrina Klos & Dr. Andrea Ortiz

# Lecture Overview



| Chapter 1 | **Introduction** |
| Chapters 2-3 | **Probability Theory** |
| Chapters 4-6 | **The Multi-Armed Bandit Problem** <br><br> Decisions do not influence future data <br><br> With/without context |
| Chapters 7-12 | **The full RL Problem** <br><br> Decisions may influence future data <br><br> With/without knowledge of dynamics |
| Chapter 13 | **Case Study** |
| Chapter 14 | **Extensions** |

# Learning Goals

- You can determine characteristics of continuous random variables and relate important examples to their applications.

- You can apply the formulas for multiple random variables and operations on random variables to compute probabilities, distributions, expectation and variance.

- You can distinguish the fundamental concepts of statistics and apply results and formulas for point estimation and confidence intervals.

# Outline

- Continuous Random Variables

- Multiple Random Variables

- Operations on Random Variables

- Statistics

# Recap: Random Variables (RVs)
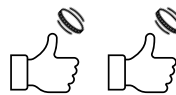
**RVs link sample spaces and events to data**

---

### Definition (Random Variable)

A **random variable (RV)** is a function $X : \Omega \to \mathcal{X}$ that assigns an element of $\mathcal{X}$ to each $\omega \in \Omega$.

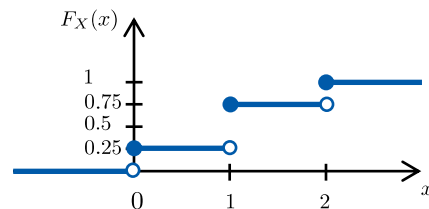- The distribution of an RV $X$ can be completely determined by its **cumulative distribution function (CDF)**

$$F_X(x) := \mathbb{P}(X \leq x).$$

- **Example**: $X(\omega)$: Number of "heads" in 2 coin tosses

$$\mathbb{P}(X = 0) = \mathbb{P}(TT) = \frac{1}{4}$$

$$\mathbb{P}(X = 1) = \mathbb{P}(HT, TH) = \frac{1}{2}$$

$$\mathbb{P}(X = 2) = \mathbb{P}(HH) = \frac{1}{4}$$

$$\Rightarrow \quad F_X(x) = \begin{cases} 0 & x < 0 \\ \frac{1}{4} & 0 \leq x < 1 \\ \frac{3}{4} & 1 \leq x < 2 \\ 1 & x \geq 2 \end{cases}$$



| Continuous RVs | Multiple RVs | Operations on RVs | Statistics |
|---|---|---|---|

# Recap: Discrete Random Variables (RVs)
**RVs with countably many values**

## Definition (Discrete Random Variable)

A random variable $X$ is **discrete** if it takes only countably many values $\{x_1, x_2, ...\}$.
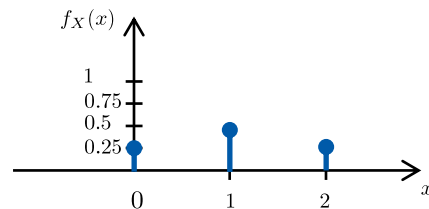
## Definition (Probability Mass Function)

For a discrete random variable $X$, we define the **probability mass function (PMF)** of $X$ by
$$f_X(x) := \mathbb{P}(X = x).$$

- **Example**: $X(\omega)$: Number of "heads" in 2 coin tosses

$$\mathbb{P}(X = 0) = \mathbb{P}(TT) = \frac{1}{4}$$
$$\mathbb{P}(X = 1) = \mathbb{P}(HT, TH) = \frac{1}{2}$$
$$\mathbb{P}(X = 2) = \mathbb{P}(HH) = \frac{1}{4}$$

$$\Rightarrow f_X(x) = \begin{cases} \frac{1}{4} & x = 0 \\ \frac{1}{2} & x = 1 \\ \frac{1}{4} & x = 2 \\ 0 & x \notin \{0, 1, 2\} \end{cases}$$



| Continuous RVs | Multiple RVs | Operations on RVs | Statistics |

# Outline

- **Continuous Random Variables**

- Multiple Random Variables

- Operations on Random Variables

- Statistics

We also consider RVs with an uncountable number of values in $\mathcal{X} = \mathbb{R}$.

**Definition (Continuous Random Variable and Probability Density Function)**

A random variable $X$ is **continuous** if there exists a function $f_X$ such that $f_X(x) \geq 0$ for all $x$, $\int_{-\infty}^{\infty} f_X(x)dx = 1$ and for every $a \leq b$,

$$\mathbb{P}(a < X < b) = \int_a^b f_X(x')dx'.$$

The function $f_X$ is called the **probability density function (PDF).** We have that

$$F_X(x) = \int_{-\infty}^x f_X(x')dx'$$

and $f_X(x) = \dfrac{dF_X(x)}{dx}$ at all points $x$ at which $F_X(x)$ is differentiable.

- **Note:** For a continuous random variable $X$, $\mathbb{P}(X = x) = 0$ for all $x$.
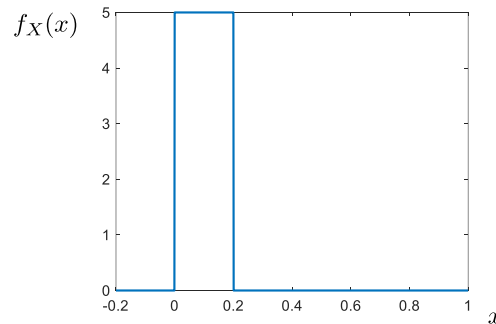
# Important Continuous Random Variables
**Example: Uniform Distribution**

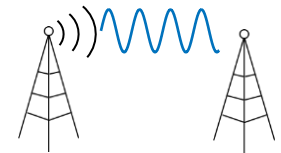## Definition (Uniform Distribution)

A random variable $X$ has a **(continuous) uniform distribution** on the interval $[a, b]$, written $X \sim \mathcal{U}(a, b)$, if it has the PDF

$$f_X(x) = \begin{cases} \frac{1}{b-a} & \text{for } x \in [a, b] \\ 0 & \text{otherwise.} \end{cases}$$



Uniform distribution on $[0, 0.2]$

- **Application:** Can be used to model the belief of a receiver about the unknown phase of a transmitted radio frequency sinusoid in a communications system.
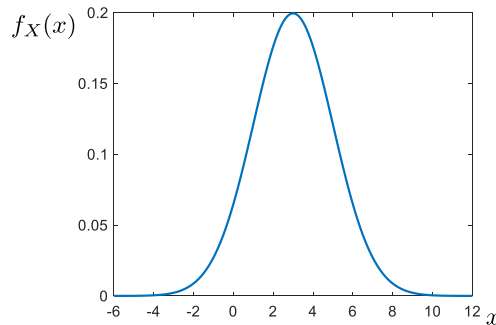
# Important Continuous Random Variables
**Example: Normal/Gaussian Distribution**

---

## Definition (Normal/Gaussian Distribution)

A random variable $X$ has a **Normal/Gaussian distribution** with parameters $\mu \in \mathbb{R}$ and $\sigma > 0$, written $X \sim \mathcal{N}(\mu, \sigma^2)$, if it has the PDF

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



Normal distribution with $\mu = 3, \sigma = 2$

The Normal distribution is very important:

- Many quantities can be approximated by a normal distribution.

- It has convenient mathematical properties.

- **Application:** Can be used to model e.g., noise in wireless communication channels and thermal noise in electronic circuits.

# Important Continuous Random Variables
**Example: Exponential Distribution**

## Definition (Exponential Distribution)

A random variable $X$ has an **Exponential** distribution with parameter $\lambda > 0$, written $X \sim \mathrm{Exp}(\lambda)$, if it has the PDF

$$f_X(x) = \lambda e^{-\lambda x} \quad \text{for} \quad x > 0 \,.$$



Exponential distribution with $\lambda = 2$

- **Application:** Can be used to describe the waiting time for a memoryless process, e.g., the interarrival times between independent accesses to a server.

$T = t_1 - t_2$

# Outline

- Continuous Random Variables

- **Multiple Random Variables**

- Operations on Random Variables

- Statistics

# Joint Distributions

**Joint distribution functions characterize the joint distribution of multiple RVs**

## Definition (Joint Distribution and Joint Density Function)

- For $n$ random variables $X_1, X_2, \ldots, X_n$, the function

$$F_{X_1, X_2, \ldots, X_n}(x_1, x_2, \ldots, x_n) = \mathbb{P}(X_1 \leq x_1, X_2 \leq x_2, \ldots, X_n \leq x_n)$$

  is called the **joint (cumulative) distribution function (joint CDF)**.

- In the discrete case, we define the **joint (probability) mass function (joint PMF)** by

$$f_{X_1, X_2, \ldots X_n}(x_1, \ldots, x_n) := \mathbb{P}(X_1 = x_1, \ldots, X_n = x_n).$$

- In the continuous case, we call a function $f_{X_1, X_2, \ldots X_n}(x_1, x_2, \ldots, x_n)$ a **joint (probability) density function (joint PDF)** if

  i.    $f_{X_1, X_2, \ldots X_n}(x_1, x_2, \ldots, x_n) \geq 0$ for all $(x_1, x_2, \ldots, x_n)$

  ii.    $\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f_{X_1, \ldots, X_n}(x_1, \ldots, x_n) dx_1 \ldots dx_n = 1$

  iii.    For any $A \subset \mathbb{R}^n$: $\mathbb{P}((X_1, X_2, \ldots, X_n) \in A) = \int_A f_{X_1, X_2, \ldots X_n}(x_1, x_2, \ldots, x_n) \, dx_1 dx_2 \ldots dx_n$

# Marginal Distributions

**Marginal distribution functions characterize the distribution of one of multiple RVs**

## Definition (Marginal Distribution Function)

- Let $F_{X_1,X_2,\ldots,X_n}(x_1, x_2, \ldots, x_n)$ denote the joint distribution of $X_1, X_2, \ldots, X_n$. The **marginal distribution function** of $X_i$ is given by

$$F_{X_i}(x_i) = \lim_{\substack{x_j \to \infty \\ j=1,\ldots,i-1,i+1,\ldots,n}} F_{X_1,X_2\ldots,X_n}(x_1, x_2 \ldots, x_n).$$

- In the discrete case, the **marginal mass function** for $X_i$ is defined by

$$f_{X_i}(x_i) := \mathbb{P}(X_i = x_i) = \sum_{x_1} \cdots \sum_{x_{i-1}} \sum_{x_{i+1}} \cdots \sum_{x_n} f_{X_1,\ldots,X_n}(x_1, \ldots, x_n).$$

- In the continuous case, we obtain the **marginal density function** by

$$f_{X_i}(x_i) = \int_{\mathbb{R}^{n-1}} f_{X_1,X_2,\ldots,X_n}(x_1, x_2, \ldots, x_n)\, dx_1 \ldots dx_{i-1} dx_{i+1} \ldots dx_n.$$

- Marginals can be defined for any subset of the RVs $X_1, X_2, \ldots, X_n$.

# Joint and Marginal Distributions
**Example: A bivariate distribution for two discrete RVs**

Here is a bivariate distribution for two discrete RVs X,Y each taking values 0 or 1:

|       | Y=0 | Y=1 |     |
|-------|-----|-----|-----|
| X=0   | 1/9 | 2/9 | 1/3 |
| X=1   | 2/9 | 4/9 | 2/3 |
|       | 1/3 | 2/3 | 1   |

- The inner part of the table shows the joint mass function $f_{X,Y}(x,y) = \mathbb{P}(X = x, Y = y).$

- The row totals show the marginal mass function of X $\quad f_X(x) = \mathbb{P}(X = x) = \sum_y f_{X,Y}(x,y).$

- The column totals show the marginal mass function of Y $\quad f_Y(y) = \mathbb{P}(Y = y) = \sum_x f_{X,Y}(x,y).$

Here is a bivariate distribution for two discrete RVs X,Y each taking values 0 or 1:

|  | Y=0 | Y=1 |  |
|---|---|---|---|
| X=0 | 1/9 | 2/9 | 1/3 |
| X=1 | 2/9 | 4/9 | 2/3 |
|  | 1/3 | 2/3 | 1 |

- The inner part of the table shows the joint mass function $f_{X,Y}(x,y) = \mathbb{P}(X=x, Y=y).$

- The row totals show the marginal mass function of X $\quad f_X(x) = \mathbb{P}(X=x) = \sum_y f_{X,Y}(x,y).$

- The column totals show the marginal mass function of Y $\quad f_Y(y) = \mathbb{P}(Y=y) = \sum_x f_{X,Y}(x,y).$

We read the probabilities from the joint and marginal mass functions in the table:

|  | Y=0 | Y=1 |  |
|---|---|---|---|
| X=0 | 1/9 | 2/9 | 1/3 |
| X=1 | 2/9 | 4/9 | 2/3 |
|  | 1/3 | 2/3 | 1 |

$$\mathbb{P}(X = 1, Y = 1) = f_{X,Y}(1,1) = 4/9$$

$$\mathbb{P}(X = 1) = f_X(1) = 2/3$$

$$\mathbb{P}(Y = 1) = f_Y(1) = 2/3$$

# Independent Random Variables

**RVs are independent iff joint density is product of marginal densities**

TECHNISCHE UNIVERSITÄT DARMSTADT

---

### Definition (Independent Random Variables)

The random variables $X_1, X_2, \ldots, X_n$ are said to be **independent** if for every $A_1, A_2, \ldots, A_n$

$$\mathbb{P}(X_1 \in A_1, X_2 \in A_2, \ldots, X_n \in A_n) = \prod_{i=1}^{n} \mathbb{P}(X_i \in A_i) \quad \Leftrightarrow \quad f(x_1, x_2, \ldots x_n) = \prod_{i=1}^{n} f_{X_i}(x_i)$$

where $f$ is the PMF in the discrete case and the PDF in the continuous case.

---

### Definition (i.i.d. Random Variables)

If $X_1, X_2, \ldots, X_n$ are independent and all $X_i$ have the same marginal distribution $F$, we say that $X_1, X_2, \ldots, X_n$ are **independent and identically distributed (i.i.d.)**.

- We can also see the i.i.d. RVs $X_1, X_2, \ldots, X_n$ as a random sample of size $n$ from distribution $F$.

    → This idea is important for statistical inference.

# Independent Random Variables

**Example: Check independence in a bivariate distribution for two discrete RVs**

Consider again the bivariate distribution for two discrete RVs X,Y each taking values 0 or 1:

|       | Y=0 | Y=1 |     |
|-------|-----|-----|-----|
| X=0   | 1/9 | 2/9 | 1/3 |
| X=1   | 2/9 | 4/9 | 2/3 |
|       | 1/3 | 2/3 |  1  |

**Joint mass function:**

$$f_{X,Y}(0,0) = 1/9$$
$$f_{X,Y}(0,1) = 2/9$$
$$f_{X,Y}(1,0) = 2/9$$
$$f_{X,Y}(1,1) = 4/9$$

**Marginal mass function of X:**

$$f_X(0) = 1/3$$
$$f_X(1) = 2/3$$

**Marginal mass function of Y:**

$$f_Y(0) = 1/3$$
$$f_Y(1) = 2/3$$

$$\Rightarrow$$

$$f_{X,Y}(0,0) = f_X(0)f_Y(0)$$
$$f_{X,Y}(0,1) = f_X(0)f_Y(1)$$
$$f_{X,Y}(1,0) = f_X(1)f_Y(0)$$
$$f_{X,Y}(1,1) = f_X(1)f_Y(1)$$

$\Rightarrow$ X and Y are independent.

# Conditional Distributions
**We can condition an RV on the value of another RV**

For simplicity, consider two RVs $X_1$ and $X_2$ with a joint distribution $F_{X_1, X_2}$.

We are interested in the distribution of $X_1$ for a given value of $X_2$.

> **Definition (Conditional Probability Mass and Density Functions)**
>
> - For $X_1, X_2$ discrete and $f_{X_2}(x_2) > 0$, the **conditional probability mass function** is
> $$f_{X_1 | X_2}(x_1 | x_2) := \mathbb{P}(X_1 = x_1 | X_2 = x_2) = \frac{\mathbb{P}(X_1 = x_1, X_2 = x_2)}{\mathbb{P}(X_2 = x_2)} = \frac{f_{X_1, X_2}(x_1, x_2)}{f_{X_2}(x_2)}.$$
>
> - For $X_1, X_2$ continuous and $f_{X_2}(x_2) > 0$, the **conditional probability density function** is
> $$f_{X_1 | X_2}(x_1 | x_2) := \frac{f_{X_1, X_2}(x_1, x_2)}{f_{X_2}(x_2)} \quad \text{and} \quad \mathbb{P}(a_1 \leq X_1 \leq b_1 | X_2 = x_2) = \int_{a_1}^{b_1} f_{X_1 | X_2}(x_1 | x_2) \, dx_1 \,.$$

# Outline

- Continuous Random Variables

- Multiple Random Variables

- **Operations on Random Variables**

- Statistics

# Expectation of an RV
**The expectation of an RV is its "average" value**

In applications, the full distribution of an RV is usually inaccessible.

→ We therefore consider certain summary functions.

---

**Definition (Expectation)**

- The **expected value**, or **mean**, or **first moment**, of a discrete RV $X$ is defined as

$$\mathbb{E}[X] = \sum_{x \in \mathcal{X}} x\, f_X(x),$$

where $f_X(x)$ is the PMF of $X$.

- The **expected value**, or **mean**, or **first moment**, of a continuous RV $X$ is defined as

$$\mathbb{E}[X] = \int_{-\infty}^{\infty} x\, f_X(x)\, dx,$$

where $f_X(x)$ is the PDF of $X$.

---

- Generalization to multiple random variables is straightforward.
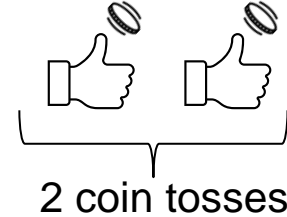
- **Random Variable:** Let $X(\omega)$ be the number of "heads"

  in the sequence $\omega \in \Omega$ of two coin tosses,

  where $\Omega = \{H, T\} \times \{H, T\}$.

2 coin tosses

**What is the expected value of this RV?**

# Answer

**Sum up PMF weighted by values of RV**

- **Random Variable:** Let $X(\omega)$ be the number of "heads"

  in the sequence $\omega \in \Omega$ of two coin tosses,
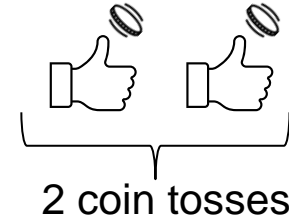
  where $\Omega = \{H, T\} \times \{H, T\}$.

2 coin tosses

We already found the PMF:

$$f_X(x) = \begin{cases} \frac{1}{4} & x = 0 \\ \frac{1}{2} & x = 1 \\ \frac{1}{4} & x = 2 \\ 0 & x \notin \{0, 1, 2\} \end{cases}$$



$$\Rightarrow \quad \mathbb{E}(X) = \sum_x x \, f_X(x)$$
$$= 0 \cdot f_X(0) + 1 \cdot f_X(1) + 2 f_X(2)$$
$$= 0 \cdot \frac{1}{4} + 1 \cdot \frac{1}{2} + 2 \frac{1}{4}$$
$$= 1$$

# Properties of the Expectation
**The expectation is a linear, monotone operator**

---

## Theorem (The rule of the lazy statistician)

For an RV $X$ with density $f_X(x)$ and a function $g$, define the new RV $Y = g(X)$. Then

$$\mathbb{E}[Y] := \mathbb{E}[g(X)] = \int_{-\infty}^{\infty} g(x)\, f_X(x)\, dx \,.$$

## Theorem (Properties of Expectation)

Let $X, Y$ be general RVs with $\mathbb{E}[X], \mathbb{E}[Y] < \infty$.

- **Linearity:** $\mathbb{E}[\alpha X + \beta Y] = \alpha \mathbb{E}[X] + \beta \mathbb{E}[Y]$ for constants $\alpha, \beta$.

- **Monotonicity:** If $X \leq Y$ $(F_X(x) \geq F_Y(x), \forall x)$, then also $\mathbb{E}[X] \leq \mathbb{E}[Y]$.

- For $X, Y$ independent: $\mathbb{E}[X \cdot Y] = \mathbb{E}[X] \cdot \mathbb{E}[Y]$

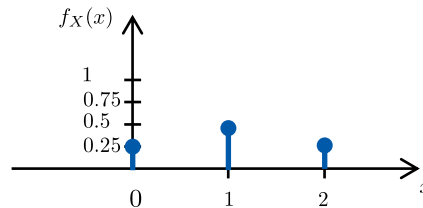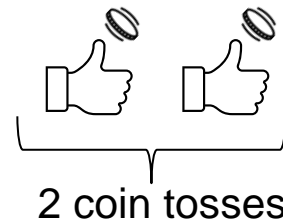# Properties of the Expectation
**Example: Expected profit in game with two coin tosses**

- **Random Variable:** Let $X(\omega)$ be the number of "heads"

  in the sequence $\omega \in \Omega$ of two coin tosses,

  where $\Omega = \{H, T\} \times \{H, T\}$.

  2 coin tosses

- **Game:** After the two coin tosses, you are paid a profit of $2^{X(\omega)}$.

- **Expected Profit:** Set $Y = g(X) := 2^X$ and apply rule of the lazy statistician.

$$\mathbb{E}(Y) = \sum_x g(x)\, f_X(x)$$

$$= 2^0 \cdot f_X(0) + 2^1 \cdot f_X(1) + 2^2 \cdot f_X(2)$$

$$= 1 \cdot \tfrac{1}{4} + 2 \cdot \tfrac{1}{2} + 4 \cdot \tfrac{1}{4}$$

$$= \tfrac{9}{4}$$

# Variance
**The variance measures the "spread" of a distribution**

> **Definition (Variance)**
>
> For an RV $X$ with $\mathbb{E}[X], \mathbb{E}[X^2] < \infty$, the variance is defined as
> $$\mathbb{V}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2].$$

- The variance can also be written as $\mathbb{V}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$.

- It is a measure of the spread of a distribution around its mean.

- The standard deviation is related to the variance via $\mathrm{std}[X] = \sqrt{\mathbb{V}[X]}$.

- The variance is often denoted by $\sigma^2$ and the standard deviation by $\sigma$.

# Outline

- Continuous Random Variables

- Multiple Random Variables

- Operations on Random Variables

- **Statistics**

# Statistical Inference
**The process of using data to infer the distribution that generated the data**
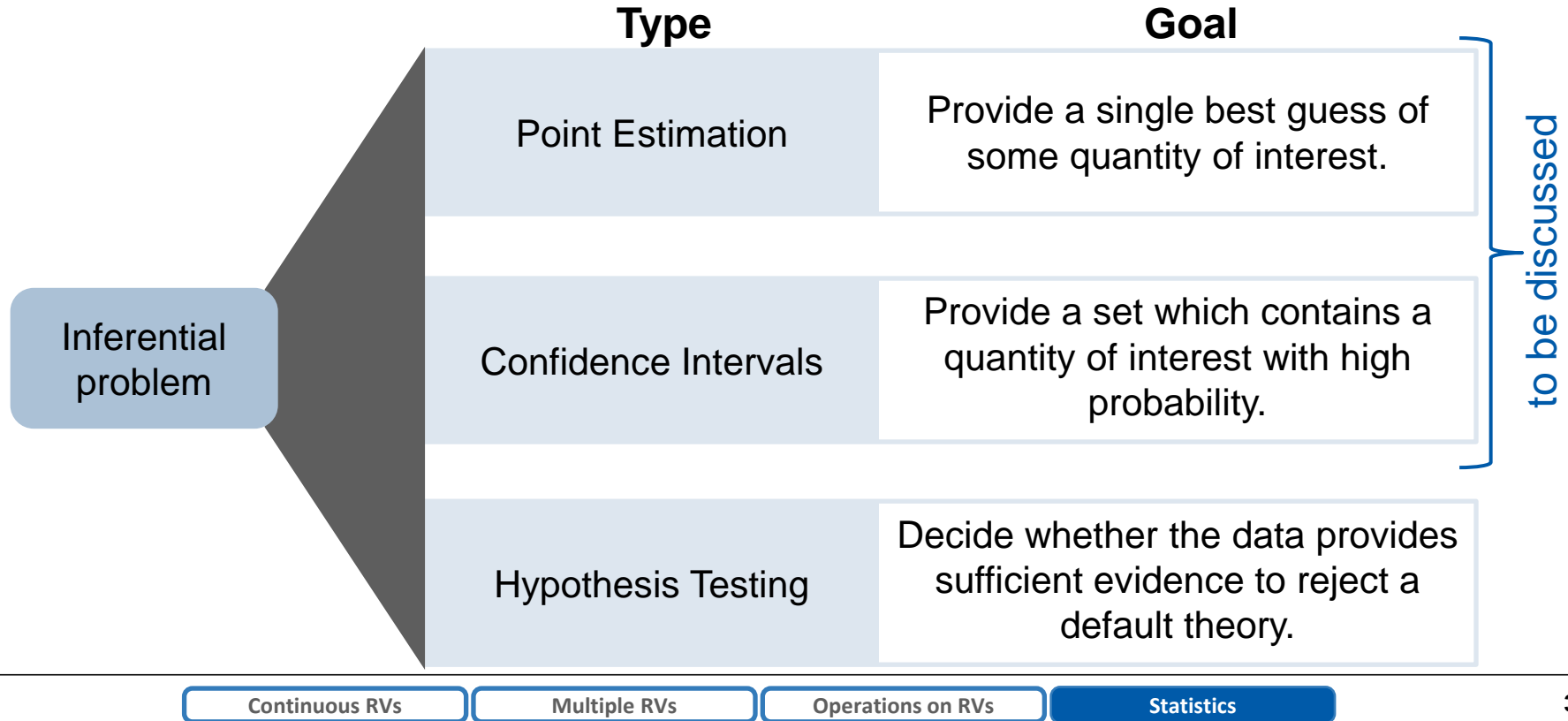
Basic statistical inference problem:

> We observe $X_1, X_2, \ldots, X_n \sim F$ i.i.d.
> How can we **infer** (or **estimate** or **learn**) the distribution $F$ or some features of $F$ ?

- This task is known as **statistical inference** or **learning**.

- Statistics is deeply connected with **machine learning**.

# Fundamental Concepts in Statistical Inference
**Many inferential problems can be identified as being one of 3 types**

| Type | Goal | |
|------|------|---|
| Point Estimation | Provide a single best guess of some quantity of interest. | |
| Inferential problem → Confidence Intervals | Provide a set which contains a quantity of interest with high probability. | to be discussed |
| Hypothesis Testing | Decide whether the data provides sufficient evidence to reject a default theory. | |

Continuous RVs   Multiple RVs   Operations on RVs   **Statistics**

# Point Estimation

**Point estimators provide a single best guess of some quantity of interest**

---

### Definition (Point estimator)

Assume we have $X_1, X_2, \ldots, X_n$ i.i.d. samples from a distribution $F_\theta$ from a class of candidate distributions defined by some parameter vector $\theta \in \Theta$.

- A **point estimator** $\hat{\theta}_n$ of $\theta$ is a function
$$\hat{\theta}_n = g(X_1, \ldots, X_n).$$

- We define the **bias** of $\hat{\theta}_n$ to be
$$\mathrm{bias}[\hat{\theta}_n] = \mathbb{E}[\hat{\theta}_n] - \theta.$$

- We call $\hat{\theta}_n$ **unbiased** if
$$\mathbb{E}[\hat{\theta}_n] = \theta.$$

- We call $\hat{\theta}_n$ **consistent** if
$$\hat{\theta}_n \longrightarrow \theta \quad \text{for } n \to \infty.$$

- Point estimators often have a limiting Normal distribution.

# Point Estimation

**Two important point estimators are sample mean and sample variance**

---

**Definition (Sample Mean and Sample Variance)**

If $X_1, X_2, \ldots, X_n$ are random variables, then we define the **sample mean** to be

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$$

and the **sample variance** to be

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^{n} (X_i - \bar{X}_n)^2 \ .$$

- The sample mean is an **unbiased** and **consistent** estimator of the true expected value.

- The sample variance is an **unbiased** and **consistent** estimator of the true variance.

1+1= → Exercise 2

# Limit Theorems

**Theorems describe the limiting behaviour of sequences of random variables**

The behavior of the sample mean for a large number of samples is described by two important theorems.

Let $\mu = \mathbb{E}[X] < \infty$ and $\sigma^2 = \text{Var}[X] < \infty$ denote expected value and variance of $X$ with $X \sim F$.

### Theorem (Weak Law of Large Numbers (WLLN))

Let $X_1, X_2, \ldots, X_n$ be i.i.d. random variables from $F$. Then

$$\bar{X}_n \longrightarrow \mu \quad \text{for } n \to \infty.$$

### Theorem (Central limit theorem (CLT))

Let $X_1, X_2, \ldots, X_n$ be i.i.d. random variables from $F$. Then

$$Z_n = \frac{\bar{X}_n - \mu}{\sqrt{\text{Var}[\bar{X}_n]}} \longrightarrow Z \sim \mathcal{N}(0, 1) \quad \text{for } n \to \infty.$$

# Point Estimation

**Example: How to estimate the probability of heads in coin tossing**

- **Experiment:** Consider tossing a coin for which the probability of heads is $p$.

- **Random Variable:** Let $X_i$ be the outcome of a single coin toss, where

$$X_i(\omega) = \begin{cases} 1, \omega = H \\ 0, \omega = T. \end{cases}$$

**What is the distribution of this RV and what is its expected value?**

- **Experiment:** Consider tossing a coin for which the probability of heads is $p$.

- **Random Variable:** Let $X_i$ be the outcome of a single coin toss, where

$$X_i(\omega) = \begin{cases} 1, \omega = H \\ 0, \omega = T. \end{cases}$$

- **Distribution:** ?

- **Expected Value:** ?

# Answer

**The RV is Bernoulli distributed with expected value p**

- **Experiment:**  Consider tossing a coin for which the probability of heads is $p$.

- **Random Variable:** Let $X_i$ be the outcome of a single coin toss, where

$$X_i(\omega) = \begin{cases} 1, \omega = H \\ 0, \omega = T. \end{cases}$$

- **Distribution:**  $X_i \sim \text{Bernoulli}(p)$

- **Expected Value:** $\mathbb{E}(X_i) = \sum_x x\, f_{X_i}(x)$

$$= 0 \cdot f_{X_i}(0) + 1 \cdot f_{X_i}(1)$$

$$= 0 \cdot \mathbb{P}[X_i = 0] + 1 \cdot \mathbb{P}[X_i = 1]$$

$$= p$$

- **Experiment:** Consider tossing a coin for which the probability of heads is $p$.

- **Random Variable:** Let $X_i$ be the outcome of a single coin toss, where

$$X_i(\omega) = \begin{cases} 1, \omega = H \\ 0, \omega = T. \end{cases} \Rightarrow X_i \sim \text{Bernoulli}(p), \mathbb{E}(X_i) = p.$$

- **Point Estimation:** ?

# Answer

**Use fraction of heads after n tosses as point estimator**

- **Experiment:** Consider tossing a coin for which the probability of heads is $p$.

- **Random Variable:** Let $X_i$ be the outcome of a single coin toss, where

$$X_i(\omega) = \begin{cases} 1, \omega = H \\ 0, \omega = T. \end{cases} \Rightarrow X_i \sim \text{Bernoulli}(p), \mathbb{E}(X_i) = p.$$

- **Point Estimation:** A possible point estimator $\hat{p}_n$ for parameter $p$ is the fraction of heads after $n$ coin tosses, given by the (unbiased and consistent) sample average, i.e.,

$$\hat{p}_n := \bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i.$$

# Likelihood-based inference

**Maximum likelihood is most common method for parameter estimation**

Idea: If $X \sim F_\theta$, we can understand this as a distribution conditional on parameter $\theta$.

---

### Definition (Likelihood Function)

Let $X_1, X_2, \ldots, X_n$ i.i.d. be continuous random variables with PDF $f(X_i|\theta)$.
The **likelihood function** $L_n(\theta)$ is a defined as the joint conditional

$$L_n(\theta) \equiv f(x_1, \ldots, x_n \mid \theta) = \prod_{i=1}^{n} f(x_i \mid \theta).$$

---

### Definition (Maximum Likelihood Estimator (MLE))

The **maximum likelihood estimator (MLE)** is defined as the value of $\theta$ that maximizes the likelihood, i.e.,

$$\hat{\theta}_n = \arg\max_\theta L_n(\theta) = \arg\max_\theta \log L_n(\theta).$$

1+1= ✏ → Exercise 2

# Confidence Intervals
**Confidence intervals contain a quantity of interest with high probability**

> **Definition (Confidence Interval)**
>
> Assume we have $X_1, X_2, \ldots, X_n$ i.i.d. samples from a distribution $F_\theta$ from a class of candidate distributions defined by some (one-dimensional) parameter $\theta$. Let $\alpha \in [0, 1]$.
> A $1 - \alpha$ **confidence interval** for $\theta$ is an interval
> $$C_n = (a, b)$$
> where $a = a(X_1, ..., X_n)$ and $b = b(X_1, ..., X_n)$ are functions of $X_1, X_2, \ldots, X_n$ such that
> $$\mathbb{P}[\theta \in C_n] \geq 1 - \alpha \quad \forall \theta \in \Theta.$$

- I.e., $C_n = (a, b)$ traps $\theta$ with probability $1 - \alpha$. Note that $C_n = (a, b)$ is random and $\theta$ is fixed!

- (Approximate) confidence intervals can often be constructed based on point estimators with limiting Normal distribution.

- If $\theta$ is a vector, we use a **confidence set** (e.g., a sphere or an ellipse) instead of an interval.

# Hoeffding's Inequality

**This inequality is useful for constructing confidence intervals**

---

**Theorem (Hoeffding's Inequality)**

Let $X_1, X_2, \ldots, X_n$ be i.i.d. RVs with values in $[0, 1]$ and expected value $\mathbb{E}[X]$ and let

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$$

be the sample mean. Then, for any $u > 0$,

$$\mathbb{P}[\mathbb{E}[X] \geq \bar{X}_n + u] \leq e^{-2nu^2}$$

and

$$\mathbb{P}[|\bar{X}_n - \mathbb{E}[X]| \geq u] \leq 2e^{-2nu^2}.$$

- There exist also variants of Hoeffding's inequality for independent random variables with bounded supports $a_i \leq X_i \leq b_i, i = 1, ..., n.$

- We can use Hoeffding's inequality to construct confidence intervals of samples of i.i.d. RVs.

1+1= → Exercise 2    → Chapter 5

# Learning Goals

- You can determine the characteristics of continuous random variables and relate important examples to their applications.

  → Probability Density Function; Uniform / Gaussian / Exponential Distribution.

- You can apply the formulas for multiple random variables and operations on random variables to compute probabilities, distributions, expectation and variance.

  → Joints; Marginals; Independence; Conditionals; Expectation and its properties; Variance.

- You can distinguish the fundamental concepts of statistics and apply results and formulas for point estimation and confidence intervals.

  → Point Estimators and their properties; Sample Mean and Variance; Limit Theorems; Maximum Likelihood; Confidence Intervals; Hoeffding's Inequality.

# Lecture Overview
**Next week, we'll study a simplified version of RL**