

# Lecture 2: Recap of Probability Theory

Isabel Valera

Machine Learning Group  
Department of Mathematics and Computer Science  
Saarland University, Saarbrücken, Germany

21.04.2021

# Outline

- 1 Bibliography
- 2 Introduction
- 3 Discrete Random Variables
- 4 Continuous Random Variables
- 5 Moments
- 6 Bayes' Theorem

# Main references

- Statistics Lab notes by Prof. Wolf
- Bishop - Chapter 1.2

# Outline

- 1 Bibliography
- 2 Introduction**
- 3 Discrete Random Variables
- 4 Continuous Random Variables
- 5 Moments
- 6 Bayes' Theorem

# Why probability theory in ML course

- A key concept in ML is uncertainty.
- Source of uncertainty are diverse and include the noise in the measurements (i.e., in the observed data) and the finite sample size from the underlying data distribution.
- Probability theory gives a theoretical framework to reason under uncertainty, i.e., to quantify and manipulate uncertainty.
- **Frequentist interpretation:** Probability as the frequency or propensity of some event, i.e.,

$$P(A) = \lim_{n \rightarrow \infty} \frac{n_A}{n},$$

where  $n_A$  is the number of times  $A$  happens in  $n$  trials (usually it is assumed that  $n \rightarrow \infty$ ).

- **Bayesian interpretation:**  
Probabilities as quantification of a belief or the uncertainty on unobserved quantities.

# Outline

- 1 Bibliography
- 2 Introduction
- 3 Discrete Random Variables
- 4 Continuous Random Variables
- 5 Moments
- 6 Bayes' Theorem

# Discrete probability

- A random variable is used to represent the outcome of an experiment. When the number of possible outcomes is countable, then we encounter a **discrete random variable**.
- The set of all possible outcomes is called the **sample space**:  
 $\Omega = \{\omega_1, \dots, \omega_n\}$  (e.g., in tossing a coin experiment,  $\Omega = \{H, T\}$ ).
- **Elementary event** is a singleton  $\{\omega_r\}$  of  $\Omega$ , i.e., is an event which cannot be further divided into other events.
- **The set of all possible events** is the power set  $2^\Omega$  (for the coin:  $\{\emptyset, \{H\}, \{T\}, \{H, T\}\}$ ).
- The **probability function**  $P$  maps events  $A \in 2^\Omega$  into the probability of such an event, i.e.,  $P : 2^\Omega \rightarrow [0, 1]$ , such that
  - $P(\emptyset) = 0$  and  $P(\Omega) = 1$ ,
  - $\sum_{\omega_i \in \Omega} P(\{\omega_i\}) = 1$ ,
  - $A \in 2^\Omega \implies P(A) = \sum_{\omega_i \in A} P(\{\omega_i\})$ .
- **Additive rule of probabilities:**  
 Let  $A, B \in 2^\Omega$ , then  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ .

# Example: Binomial distribution I

- An experiment with two possible outcomes  $Y \in \{0, 1\}$  is called **Bernoulli trial** (or binomial trial) and is defined by the “success” probability  $p = P(Y = 1)$ .
- The **binomial distribution** models  $n$  repeated Bernoulli trials where the outcomes are independent (e.g., in a coin toss experiment) and the random variable  $X$  accounts for the number of times we observe “success”  $Y = 1$  (the order does not matter), i.e.,

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k},$$

with the binomial coefficient  $\binom{n}{k} = \frac{n!}{k!(n-k)!}$ .

- The sample space is thus  $\Omega = \{0, 1, \dots, n\}$  and

$$P(\Omega) = \sum_{k=0}^n P(X = k) = \sum_{k=0}^n \binom{n}{k} p^k (1 - p)^{n-k} = (1 - p + p)^n = 1.$$



## Example: Binomial distribution II

- Coin toss:  $\Omega = \{H, T\}$ ,  $P(H) = p$ . Define  $Y : \{H, T\} \rightarrow \{0, 1\}$  by

$$Y = \begin{cases} 1 & \text{if } H, \\ 0 & \text{if } T. \end{cases}$$

$Y$  is a random variable with Bernoulli-distribution:

$$P_Y(Y = 1) = P(H) = p, \text{ and similarly } P_Y(Y = 0) = 1 - p.$$

- Repeat the coin toss independently  $n$  times and denote by  $X$  the number of times we observe head. Let  $\Omega$  be the set of all sequences of  $n$  variables with the alphabet  $\{H, T\}$ , then  $|\Omega| = 2^n$ .  $X$  is a random variable  $X : \Omega \rightarrow \mathbb{Z}$  with distribution

$$P_X(X = k) = P(X^{-1}(k)) = \binom{n}{k} p^k (1 - p)^{n-k}.$$

If  $n = 3$ , then  $X^{-1}(2) = \{HHT, HTH, THH\}$ .

# The Rules of Probability

- There are two fundamental rules of probability theory:

$$\textbf{Sum rule: } P(X) = \sum_Y P(X, Y) \quad (1)$$

$$\textbf{Product rule: } P(X, Y) = P(Y | X)P(X)(= P(X \cap Y)) \quad (2)$$

- Let  $X, Y$  be discrete random variables.  $X$  and  $Y$  are **independent** if,

$$P_{X \times Y}(X = i, Y = j) = P_X(i) P_Y(j), \quad \forall i, j \in \mathbb{Z}.$$

- The **conditional probability**  $P(X = i | Y = j)$  of  $X$  given  $Y = j$  is,

$$P(X = i | Y = j) = \frac{P_{X \times Y}(X = i, Y = j)}{P(Y = j)}, \quad \forall j \text{ with } P(Y = j) > 0.$$

# Example: Oranges v.s Apples from Bishop

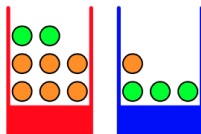


Figure: Figure 1.9 from Bishop

$$P(B = r) = 4/10$$

$$P(B = b) = 6/10$$

$$P(F = a|B = r) = 1/4$$

$$P(F = o|B = r) = 3/4$$

$$P(F = a|B = b) = 3/4$$

$$P(F = o|B = b) = 1/4$$

$$P(F = a) = P(F = a|B = r)P(B = r) + P(F = a|B = b)P(B = b)$$

$$= \frac{1}{4} \times \frac{4}{10} + \frac{3}{4} \times \frac{6}{10} = \frac{11}{20}$$

$$P(B = r|F = o) = \frac{P(F = o|B = r)P(B = r)}{P(F = o)} = \frac{3/4 \times 4/10}{9/20} = \frac{2}{3}$$

# Outline

- 1 Bibliography
- 2 Introduction
- 3 Discrete Random Variables
- 4 Continuous Random Variables**
- 5 Moments
- 6 Bayes' Theorem

# $\sigma$ -algebra

- So far, random variables taking discrete values  $X \in \{1, 2, 3, \dots\}$ , thus  $\Omega$  is a countable set.
- What if we consider continuous variables, e.g.,  $X \in \mathbb{R}$ , and thus  $\Omega = \mathbb{R}$  is uncountable? How do we assign probabilities to all  $2^\Omega$  events?
- If all numbers are equally likely to occur, how do we ensure that  $\sum_{\omega_i \in \Omega} P(\omega_i) = 1$ ?

## Definition ( $\sigma$ -algebra)

A set  $\mathcal{A} \subset 2^\Omega$  is called a  $\sigma$ -**algebra**:

- 1 If  $\emptyset \in \mathcal{A}$  and  $\Omega \in \mathcal{A}$ ,
- 2 If  $A \in \mathcal{A}$ , then also the complement  $A^c$  is contained in  $\mathcal{A}$ ,
- 3 If  $\mathcal{A}$  is closed under **countable** unions, that is if  $A_1, A_2, \dots$  is a sequence of events in  $\mathcal{A}$ , then  $\bigcup_{i=1}^{\infty} A_i \in \mathcal{A}$ .

# Probability measure

## Definition (Probability measure)

A **probability measure** defined on a  $\sigma$ -algebra  $\mathcal{A}$  of  $\Omega$  is a function  $P : \mathcal{A} \rightarrow [0, 1]$  that satisfies:

- 1  $P(\Omega) = 1$ ,
- 2 For every countable sequence  $(A_n)_{n \geq 1}$  of elements of  $\mathcal{A}$ , pairwise disjoint (that is  $A_m \cap A_n = \emptyset$  whenever  $m \neq n$ ), one has

$$P\left(\bigcup_{n=1}^{\infty} A_n\right) = \sum_{n=1}^{\infty} P(A_n).$$

- Any discrete **probability space**  $(\Omega, 2^{\Omega}, P)$  is a probability measure, since  $2^{\Omega}$  is a  $\sigma$ -algebra and  $P$  is a probability measure.

# Borel $\sigma$ -Algebra

Let  $C \subset 2^\Omega$ . The  $\sigma$ -algebra generated by  $C$  is the smallest  $\sigma$ -algebra containing  $C$ .

## Definition (Borel $\sigma$ -algebra)

The **Borel  $\sigma$ -algebra**  $\mathcal{B}$  in  $\mathbb{R}^d$  is the  $\sigma$ -algebra generated by the open sets in  $\mathbb{R}^d$ .

## Lebesgue Measure on $\mathbb{R}^d$

- The Lebesgue measure  $\mu : \mathcal{B} \rightarrow \mathbb{R}_+$  is now just the usual measure of volume. For the one-dimensional case, we have

$$\mu([a, b[) = b - a,$$

- A set  $A \in \mathcal{B}$  has **measure zero** if  $\mu(A) = 0$ . Any countable set of points has Lebesgue measure zero.

**Warning:** The Lebesgue measure works on its own (larger)  $\sigma$ -algebra but the difference is for our purposes negligible.

# Probability on continuous spaces

In the case  $\Omega = \mathbb{R}^d$  we will work with measures which have a density with respect to the **Lebesgue measure**.

Let  $\mathcal{B}$  be the Borel  $\sigma$ -algebra in  $\mathbb{R}^d$ . A probability measure  $P$  on  $(\mathbb{R}^d, \mathcal{B})$  has a **density**  $p$  if  $p$  is a non-negative (Borel measurable) function on  $\mathbb{R}^d$  satisfying for all  $A \in \mathcal{B}$  that:

$$P(A) = \int_A p(x) dx = \int_A p(x_1, \dots, x_d) dx_1 \dots dx_d,$$

where  $dx = dx_1 \dots dx_d$ .

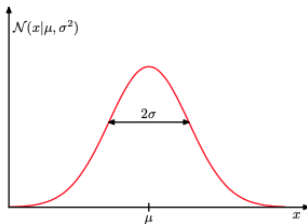
- This implies:  $P(\mathbb{R}^d) = \int_{\mathbb{R}^d} p(x) dx = 1$ .
- Observation: Not all probability measures on  $\mathbb{R}^d$  have a density.



# Example of a probability measure with density

The **Gaussian distribution** or normal distribution on  $\mathbb{R}$  has two parameters  $\mu$  (mean) and  $\sigma^2$  (variance). The associated density function is denoted by  $\mathcal{N}(\mu, \sigma^2)$  and defined as:

$$p(X = x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}.$$



$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x | \mu, \sigma^2) x \, dx = \mu$$

$$\mathbb{E}[x^2] = \int_{-\infty}^{\infty} \mathcal{N}(x | \mu, \sigma^2) x^2 \, dx = \mu^2 + \sigma^2$$

$$\text{var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma^2$$

Figure: Figure 1.13 from Bishop

# Other densities

- **Multivariate Gaussian**  $\mathcal{N}(\mu, \Sigma)$  is uniquely determined by the mean  $\mu \in \mathbb{R}^d$  and the covariance matrix  $\Sigma \in \mathbb{R}^{d \times d}$  (positive-definite) as

$$p(x) = \frac{1}{(2\pi)^{\frac{d}{2}} |\det \Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

- **Laplace distribution**  $\text{Laplace}(\mu, b)$  is given by

$$p(x) = \frac{1}{2b} e^{-\frac{1}{b}|x-\mu|}$$

- **Gamma distribution**  $\Gamma(\alpha, \beta)$  given by:

$$p(x) = \frac{x^{\alpha-1} \beta^\alpha e^{-\beta x}}{\Gamma(\alpha)}, \text{ where } \Gamma(\cdot) \text{ is the Gamma function.}$$

# Cumulative distribution function

- The **(cumulative) distribution function** of a probability measure  $P$  on  $(\mathbb{R}, \mathcal{B})$  is the function

$$F(x) = P\left(X \in (-\infty, x]\right) = P\left(X \leq x\right) = \int_{-\infty}^x p(t)dt.$$

If the distribution function  $F$  is sufficiently differentiable, then

$$p(x) = \left. \frac{\partial F}{\partial x} \right|_x.$$

- The distribution function of  $P$  on  $(\mathbb{R}^d, \mathcal{B})$  is the function

$$F(x_1, \dots, x_d) = P\left(X_1 \leq x_1, \dots, X_d \leq x_d\right).$$

If the distribution function  $F$  is sufficiently differentiable, then

$$p(x_1, \dots, x_d) = \left. \frac{\partial^d F}{\partial x_1 \dots \partial x_d} \right|_{x_1, \dots, x_d}.$$

# Quantile

**Quantiles:** Quantiles are only defined for distributions on  $\mathbb{Z}$  and  $\mathbb{R}$ .

## Definition

The  $\alpha$ -**quantile** of a probability measure on  $\mathbb{Z}$  or  $\mathbb{R}$  is the real number  $q_\alpha$  such that

$$F(q_\alpha) = P([-\infty, q_\alpha]) = \alpha.$$

The **median** is the  $\frac{1}{2}$ -quantile.

- Median and mean agree if the distributions are symmetric (and unimodal).
- The median is more robust to changes of the probability measure.



# Joint density and marginals

Let  $X = (X_1, X_2)$  be a  $\mathbb{R}^2$ -valued random variable with density  $p_X$  on  $\mathbb{R}^2$ . Then the densities  $p_{X_1}$  of  $X_1$  and  $p_{X_2}$  of  $X_2$  are given as

$$p_{X_1}(x_1) = \int_{\mathbb{R}} p_X(x_1, x_2) dx_2, \quad p_{X_2}(x_2) = \int_{\mathbb{R}} p_X(x_1, x_2) dx_1.$$

- $p_X(x_1, x_2)$  denotes the **joint density**.
- $p_{X_1}$  and  $p_{X_2}$  are called **marginal densities** of  $X$  and are associated to the probability measures of  $X_1$  respectively  $X_2$ .

Observation: The joint measure can in general not be reconstructed from the knowledge of the marginal densities (only if  $X_1$  and  $X_2$  are independent).

# Independence and conditional density

Let  $X, Y$  be  $\mathbb{R}$ -valued random variables with joint-density  $p_{X \times Y}$  and marginal densities  $p_X$  and  $p_Y$ , then  $X$  and  $Y$  are **independent** if

$$p_{X \times Y}(x, y) = p_X(x) p_Y(y), \quad \forall x, y \in \mathbb{R}.$$

The **conditional density**  $p(x|Y = y)$  of  $X$  given  $Y = y$  is defined as,

$$p(x|y) = \frac{p(x, y)}{p(y)}, \quad \forall y \text{ with } p(y) > 0.$$

# Example: Joint, marginals and conditionals

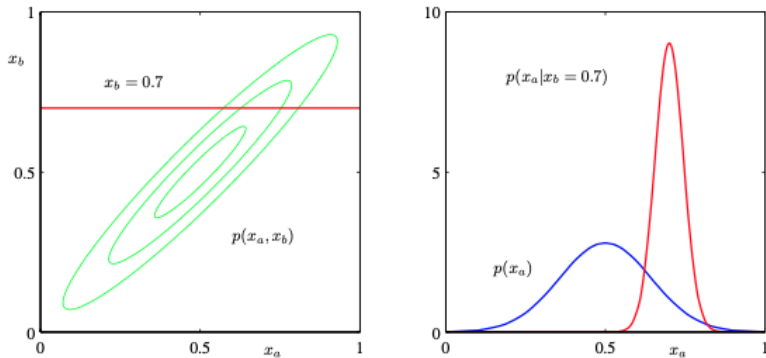


Figure 2.9 from Bishop



# Transformation of Random Variables

## Theorem

Let  $X = (X_1, \dots, X_d)$  have joint density  $p_X$ . Let  $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$  be continuously differentiable and injective, with non-vanishing Jacobian. Then  $Y = g(X)$  has density

$$p_Y(y) = p_X(g^{-1}(y)) |\det J_{g^{-1}}(y)|$$

- The Jacobian  $J_g(x)$  of a function  $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$  at value  $x$  is the  $d \times d$ - matrix

$$J_g(x)_{ij} = \left. \frac{\partial g_i}{\partial x_j} \right|_x, \quad i, j = 1, \dots, d$$

- This result allows us to generate samples from complicated densities from simple ones.

# Example: Sampling from an exponential distribution

$$p_\lambda(y) = \lambda \exp(-\lambda y), \text{ for } y \geq 0.$$

1. We can first sample from a uniform distribution on  $[0, 1]$ .
2. Apply a function  $g : [0, 1] \rightarrow \mathbb{R}_+$  (resp.  $g^{-1}$ ) such that

$$p_\lambda(y) = \lambda \exp(-\lambda y) = p_X(g^{-1}(y)) \left| \frac{\partial g^{-1}}{\partial y} \right| = \left| \frac{\partial g^{-1}}{\partial y} \right|.$$

**General case:** complicated differential equation.

**This case:**  $g^{-1}(y) = \exp(-\lambda y) \implies g(x) = -\frac{\log(x)}{\lambda}$

- $X_i$  samples from the uniform distribution on  $[0, 1]$ ,
- $Y_i = g(X_i) = -\frac{\log(X_i)}{\lambda}$  are samples from the exponential distribution.

# Outline

- 1 Bibliography
- 2 Introduction
- 3 Discrete Random Variables
- 4 Continuous Random Variables
- 5 Moments**
- 6 Bayes' Theorem

# Expectation

The **expected value** or **expectation** of a  $\mathbb{R}^d$ -valued random variable  $X$  is defined as

$$(\mathbb{E}[X])_i = \int_{\mathbb{R}^d} x_i p(x) dx = \int_{\mathbb{R}^d} x_i p(x_1, \dots, x_d) dx_1 \dots dx_d,$$

and for a discrete random variable  $X$  taking values in  $\mathbb{Z}$  it is defined as,

$$\mathbb{E}[X] = \sum_{n=-\infty}^{\infty} n P(X = n).$$

## Expectation of functions of random variables

We can also define the expectation of functions of random variables.

$$\mathbb{E}[f(X)] = \int_{\mathbb{R}^d} f(x) p(x) dx = \int_{\mathbb{R}^d} f(x_1, \dots, x_d) p(x_1, \dots, x_d) dx_1 \dots, dx_d.$$

# Variance, Covariance and Correlation

The **variance**  $\text{Var}[X]$  (also  $\sigma^2(X)$ ) of an  $\mathbb{Z}$ - or  $\mathbb{R}$ -valued random variable  $X$  is defined as

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

The standard deviation of  $X$  is  $\sigma(X) = \sqrt{\text{Var}[X]}$ .

The covariance matrix  $\Sigma$  of an  $\mathbb{R}^d$ -valued random variable  $X$  is given as  $\Sigma_{ij} = \text{Cov}(X_i, X_j)$  or in matrix form

$$\Sigma = \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^T].$$

The **covariance**  $\text{Cov}(X, Y)$  of two  $\mathbb{R}$ -valued random variables  $X$  and  $Y$  is defined as,

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

The **correlation**  $\text{Corr}(X, Y)$  of two  $\mathbb{R}$ -valued random variables  $X$  and  $Y$  is then defined as,

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Cov}(X, X) \text{Cov}(Y, Y)}} = \frac{\text{Cov}(X, Y)}{\sigma(X)\sigma(Y)}.$$

# Properties

- The expectation and variance have the following properties

$\forall a, b \in \mathbb{R}$ ,

$$\mathbb{E}[aX + b] = a \mathbb{E}[X] + b, \quad \mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y],$$

$$\text{Var}[aX + b] = a^2 \text{Var}[X],$$

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] + 2 \text{Cov}(X, Y).$$

- Correlation is a measure of **linear dependence**, and satisfies  $-1 \leq \text{Corr}(X, Y) \leq 1$ . If  $X$  and  $Y$  are linearly dependent, that is  $Y = aX + b$  with  $a, b \in \mathbb{R}$ , then

$$\text{Corr}(X, Y) = \text{Corr}(X, aX + b) = \frac{a}{|a|} = \begin{cases} 1, & \text{if } a > 0, \\ 0, & \text{if } a = 0, \\ -1, & \text{if } a < 0. \end{cases}$$

In words, linearly dependent random variables achieve maximal correlation.

# Conditional expectation

Let  $X, Y$  be two  $\mathbb{R}$ -valued random variables. The **conditional expectation**  $\mathbb{E}[X|Y = y]$  of  $X$  given  $Y = y$  is defined for  $y$  with  $p(y) > 0$  as the quantity

$$\mathbb{E}[X|Y = y] = \int_{\mathbb{R}} x p(x|y) dx.$$

The **conditional expectation**  $\mathbb{E}[X|Y]$  of  $X$  given  $Y$  is a random variable  $h(Y)$  with values

$$h(y) = \mathbb{E}[X|Y = y].$$

Important properties of the conditional expectation are:

- $\mathbb{E}[X|Y] = \mathbb{E}[X]$ , if  $X$  and  $Y$  are **independent**,
- $\mathbb{E}[\mathbb{E}[X|Y]] = \mathbb{E}[X]$ ,
- $\mathbb{E}[f(Y)|Y] = f(Y)$  and  $\mathbb{E}[f(Y)g(X)|Y] = f(Y)\mathbb{E}[g(X)|Y]$ .

# Outline

- 1 Bibliography
- 2 Introduction
- 3 Discrete Random Variables
- 4 Continuous Random Variables
- 5 Moments
- 6 Bayes' Theorem**



# Law of total probability

Assume that we have a finite or countably infinite number of events

$$\mathcal{A} = \{A_1, A_2, A_3, \dots\} \text{ and } \Omega = A_1 \cup A_2 \cup A_3 \cup \dots$$

## Definition

A collection of events  $(A_n)_{n \geq 1}$  is called a **partition** of  $\Omega$  if  $A_n \in \mathcal{A}$  for each  $n$ , they are pairwise disjoint,  $A_n \cap A_m = \emptyset$  for  $m \neq n$ ,  $P(A_n) > 0$  for each  $n$ , and  $\cup_n A_n = \Omega$ .

## Theorem (Law of total probability)

Let  $(A_n)_{n \geq 1}$  be a finite or countable partition of  $\Omega$ . Then if  $B \in \mathcal{A}$ ,

$$P(B) = \sum_n P(B|A_n)P(A_n).$$

# Bayes' theorem

## Theorem (Bayes' theorem)

Let  $A, B$  be two events and  $P(B) > 0$ , then

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)},$$

- The above definition follows from the definition of conditional probability.
- Implication: Let  $(A_n)_{n \geq 1}$  be a finite or countable partition of  $\Omega$ , and suppose  $P(B) > 0$ . Then

$$P(A_i|B) = \frac{P(B|A_i)P(A_i)}{\sum_n P(B|A_n)P(A_n)}.$$