

Probability Theory

Maximilian Soelch

Technische Universität München

Probability Theory

Probability Theory is the study of uncertainty. Uncertainty is all around us.

Mathematical probability theory is based on *measure theory*—we do not work at this level.

Slides are mostly based on *Review of Probability Theory* by Arian Meleki and Tom Do.

Probability Theory

The basic problem that we study in probability theory:

Given a data generating process,
what are the properties of the outcomes?

The basic problem of statistics (or better statistical *inference*) is the *inverse* of probability theory:

Given the outcomes,
what can we say about the process that generated the data?

Statistics uses the formal language of probability theory.

Basic Elements of Probability

Sample space Ω :

The set of all outcomes of a random experiment.

e.g. rolling a die: $\Omega = \{1, 2, 3, 4, 5, 6\}$

e.g. rolling a die twice: $\Omega' = \Omega \times \Omega = \{(1, 1), (1, 2), \dots, (6, 6)\}$

Set of events \mathcal{F} (event space):

A set whose elements $A \in \mathcal{F}$ (*events*) are subsets of Ω .

\mathcal{F} (σ -field) must satisfy

- ▶ $\emptyset \in \mathcal{F}$,
- ▶ $A \in \mathcal{F} \Rightarrow \Omega \setminus A \in \mathcal{F}$,
- ▶ $A_1, A_2, \dots \in \mathcal{F} \Rightarrow \cup_i A_i \in \mathcal{F}$.

e.g. "die outcome is even" \rightsquigarrow event $A = \{\omega \in \Omega : \omega \text{ even}\} = \{2, 4, 6\}$

e.g. (smallest) σ -field that contains A : $\mathcal{F} = \{\emptyset, \{1, 3, 5\}, \{2, 4, 6\}, \{1, 2, 3, 4, 5, 6\}\}$

Basic Elements of Probability ctd.

Probability measure $P : \mathcal{F} \rightarrow [0, 1]$

with *Axioms of Probability*:

- ▶ $P(A) \geq 0$ for all $A \in \mathcal{F}$,
- ▶ $P(\Omega) = 1$,
- ▶ If A_1, A_2, \dots are disjoint events ($A_i \cap A_j = \emptyset, i \neq j$) then

$$P(\cup_i A_i) = \sum_i P(A_i).$$

e.g. for rolling a die $P(A) = \frac{|A|}{|\Omega|}$

The triple (Ω, \mathcal{F}, P) is called a probability space.

Important Properties

The three axioms from the previous slide suffice to show:

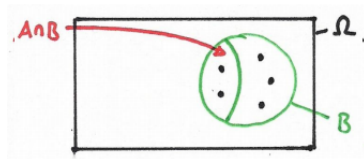
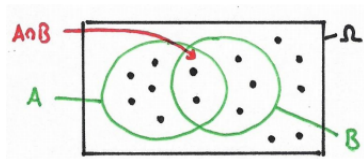
- ▶ If $A \subseteq B \Rightarrow P(A) \leq P(B)$
- ▶ $P(A \cap B) (\equiv P(A, B)) \leq \min(P(A), P(B))$
- ▶ $P(A \cup B) \leq P(A) + P(B)$
- ▶ $P(\Omega \setminus A) = 1 - P(A)$
- ▶ If A_1, A_2, \dots, A_k are sets of *disjoint* events such that $\cup_i A_i = \Omega$ then $\sum_k P(A_k) = 1$ (*Law of total probability*).

Conditional Probability

Let $A, B \subseteq \Omega$ be two events with $P(B) \neq 0$, then:

$$P(A | B) := \frac{P(A \cap B)}{P(B)}.$$

$P(A | B)$ is the probability of A conditioned on B and represents the probability of A , if it is known that B was observed.



Multiplication law

Let A_1, \dots, A_n be events with $P(A_1 \cap \dots \cap A_n) \neq 0$. Then:

$$\begin{aligned} P(A_1 \cap \dots \cap A_n) &= \prod_{i=1}^n P\left(A_i \mid \bigcap_{j<i} A_j\right) \\ &= P(A_1) \cdot P(A_2 \mid A_1) \cdot P(A_3 \mid A_1 \cap A_2) \cdot \dots \cdot P(A_n \mid A_1 \cap \dots \cap A_{n-1}). \end{aligned}$$

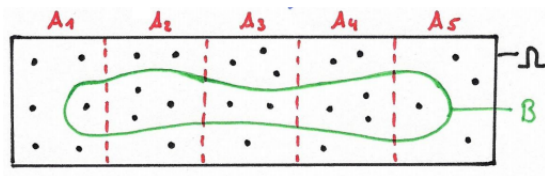
Law of total probability (revisited)

Let B be an event and Φ a partition of Ω with $P(A) > 0$ for all $A \in \Phi$.

Then:

$$P(B) = \sum_{A \in \Phi} P(B \cap A) = \sum_{A \in \Phi} P(A) \cdot P(B | A).$$

Graphical representation for a 5-partition $\Phi = \{A_1, \dots, A_5\}$ of Ω :



Bayes' rule

Let A and B be two events with $P(A), P(B) \neq 0$. Then:

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

Bayes' rule applies the multiplication rule twice to set $P(A | B)$ and $P(B | A)$ in relation:

$$P(B | A) \cdot P(A) = \underbrace{P(A \cap B)}_{=P(A,B)} = P(A | B) \cdot P(B).$$

Bayes' rule is always used if the conditional probability $P(A | B)$ is easy to calculate or given, but the conditional probability $P(B | A)$ is searched for.

Independence

Two events A , B are called *independent* if and only if

$$P(A, B) = P(A)P(B),$$

or equivalently $P(A | B) = P(A)$. What does that mean in words?

Two events A , B are called *conditionally independent* given a third event C if and only if

$$P(A, B | C) = P(A | C)P(B | C).$$

Random variables

We are usually only interested in some aspects of a random experiment.

Random variable $X : \Omega \rightarrow \mathbb{R}$ (actually not every function is allowed ...).

A random variable is usually just denoted by an upper case letter X (instead of $X(\omega)$). The value a random variable may take is denoted by the corresponding lower-case letter.

For a discrete random variable

$$P(X = x) := P(\{\omega \in \Omega : X(\omega) = x\}).$$

For a continuous random variable

$$P(a \leq X \leq b) := P(\{\omega \in \Omega : a \leq X(\omega) \leq b\}).$$

Note the usage of P here.

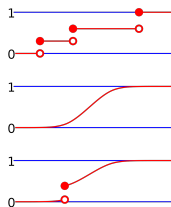
Cumulative distribution function – CDF

A probability measure P is specified by a *cumulative distribution function* (CDF), a function $F_X : \mathbb{R} \rightarrow [0, 1]$:

$$F_X(x) \equiv P(X \leq x).$$

Properties:

- ▶ $0 \leq F_X(x) \leq 1$
- ▶ $\lim_{x \rightarrow -\infty} F_X(x) = 0$
- ▶ $\lim_{x \rightarrow \infty} F_X(x) = 1$
- ▶ $x \leq y \rightarrow F_X(x) \leq F_X(y)$



Let X have CDF F_X and Y have CDF F_Y . If $F_X(x) = F_Y(x)$ for all x , then $P(X \in A) = P(Y \in A)$ for all (measurable) A .

We call X and Y *identically distributed* (or *equal in distribution*).

Probability density function—PDF

For some continuous random variables, the CDF $F_X(x)$ is continuous on \mathbb{R} . The *probability density function* is then defined as the piecewise derivative

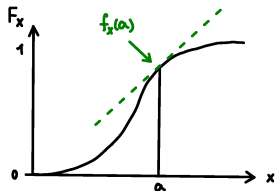
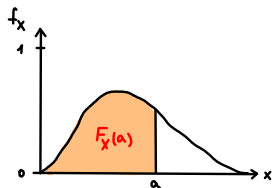
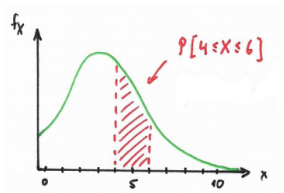
$$f_X(x) = \frac{dF_X(x)}{dx},$$

and X is called continuous.

$$P(x \leq X \leq x + \Delta x) \approx f_X(x) \cdot \Delta x.$$

Properties:

- ▶ $f_X(x) \geq 0$
- ▶ $\int_{x \in A} f_X(x) dx = P(X \in A)$
- ▶ $\int_{-\infty}^{\infty} f_X(x) dx = 1$



Probability mass function—PMF

X takes on only a *countable* set of possible values (*discrete* random variable).

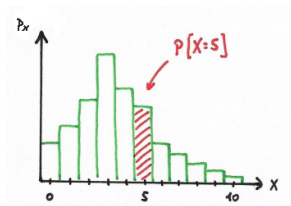
A probability mass function $p_X : \Omega \rightarrow [0, 1]$ is a simple way to *represent* the probability measure associated with X :

$$p_X(x) = P(X = x)$$

(Note: We use the probability measure P on the random variable X)

Properties:

- ▶ $0 \leq p_X(x) \leq 1$
- ▶ $\sum_x p_X(x) = 1$
- ▶ $\sum_{x \in A} p_X(x) = P(X \in A)$



Transformation of Random Variables

Given a (continuous) random variable X and a strictly monotonic (increasing *or* decreasing) function s , what can we say about $Y = s(X)$?

$$f_Y(y) = f_X(t(y)) |t'(y)|,$$

where t is the inverse of s .

Expectation

For any measurable function $g : \mathbb{R} \rightarrow \mathbb{R}$, we define the *expected value*:

$$\mathbb{E}[g(X)] = \sum_x g(x)p_X(x) \quad \text{discrete}$$

$$\mathbb{E}[g(X)] = \int_{-\infty}^{\infty} g(x)f_X(x) dx \quad \text{continuous}$$

Special case: $\mathbb{E}[X]$, i.e., $g(x) = x$, is called the *mean* of X .

Properties:

- ▶ $\mathbb{E}[a] = a$ for any *constant* $a \in \mathbb{R}$
- ▶ $\mathbb{E}[af(X)] = a\mathbb{E}[f(X)]$ for any constant $a \in \mathbb{R}$
- ▶ $\mathbb{E}[f(X) + g(X)] = \mathbb{E}[f(X)] + \mathbb{E}[g(X)]$

For any $A \in \mathbb{R}$: $\mathbb{E}[\mathbb{I}_A(X)] = P(X \in A)$

Variance and Standard Deviation

Variance measures the concentration of a random variable's distribution around its mean.

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

Properties:

- ▶ $\text{Var}(a) = 0$ for any constant $a \in \mathbb{R}$.
- ▶ $\text{Var}(af(X)) = a^2 \text{Var}(f(X))$ for any constant $a \in \mathbb{R}$.

$\sigma(X) = \sqrt{\text{Var}(X)}$ is called the *standard deviation* of X .

Entropy

The *Shannon entropy* or just *entropy* of a discrete random variable X is

$$H[X] \equiv - \sum_x \mathbb{P}(X = x) \log \mathbb{P}(X = x) = -\mathbb{E}[\log \mathbb{P}(X)].$$

Given two probability mass functions p_1 and p_2 , the *Kullback-Leibler divergence* (or *relative entropy*) between p_1 and p_2 is

$$\text{KL}(p_1 \parallel p_2) \equiv - \sum_x p_1(x) \log \frac{p_2(x)}{p_1(x)}$$

Note that the KL divergence is **not** symmetric.

Bernoulli distribution

A Bernoulli-distributed random variable $X \sim \text{Ber}(\mu)$, $\mu \in [0, 1]$ models the outcome of an experiment. It is positive with a probability of μ and negative with a probability of $1 - \mu$.

$$p_X(x) = \begin{cases} \mu, & \text{if } x = 1 \\ 1 - \mu, & \text{if } x = 0 \\ 0 & \text{else} \end{cases}$$

For calculations the following equation is more useful:

$$\text{Ber}(x \mid \mu) = \mu^x \cdot (1 - \mu)^{1-x}$$

Binomial distribution

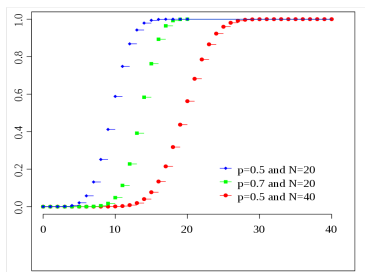
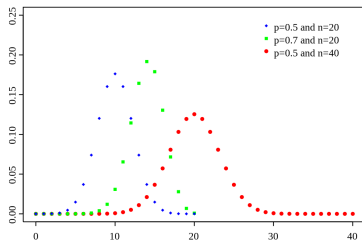
A Binomial random variable

$X \sim \text{Bin}(N, \mu)$, $N \geq 1$, $\mu \in [0, 1]$

shows the number of successes by performing N trials, where each trial is independent from the others. The success probability is μ .

For $x \in \{0, 1, \dots, N\}$:

$$\text{Bin}(x | N, \mu) = \binom{N}{x} \cdot \mu^x \cdot (1-\mu)^{N-x}$$



Poisson distribution

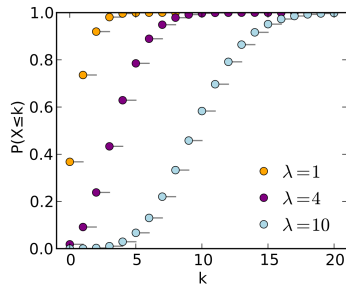
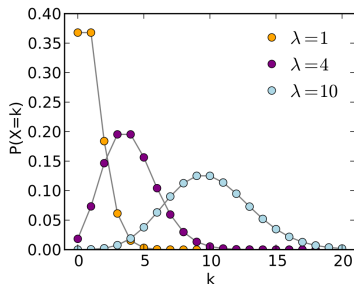
A Binomial random variable with large N and small μ can be approximated by a Poisson random variable $X \sim \text{Poi}(\lambda)$.

For $\lambda = N\mu$ and as $N \rightarrow \infty$:

$$X \sim \text{Bin}(N, \mu) \rightarrow X \sim \text{Poi}(\lambda)$$

For $x \in \mathbb{N}_0$:

$$\text{Poi}(x | \lambda) = \frac{e^{-\lambda} \cdot \lambda^x}{x!}$$

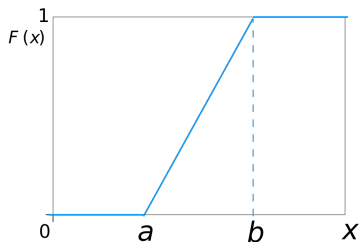
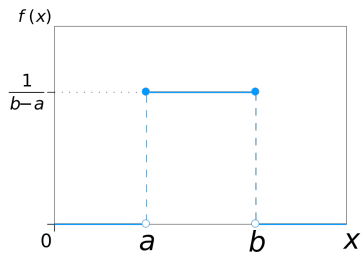


Uniform distribution

A uniformly distributed random variable $X \sim U(a, b)$, $a, b \in \mathbb{R}$, $a < b$, takes any value on the interval $[a, b]$ with equal probability.

For $x \in [a, b]$:

$$U(x | a, b) = \frac{1}{b - a}$$

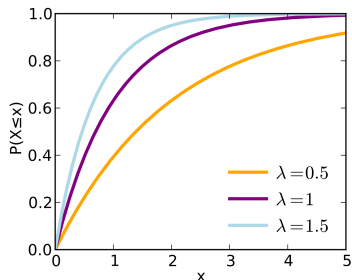
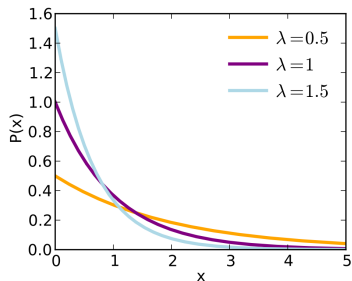


Exponential distribution

An exponentially distributed random variable $X \sim \text{Exp}(\lambda)$, $\lambda > 0$ can be referred to as the latency until an event (“success”) occurs the first time. λ corresponds to the expected number of successes in one unit of time.

For $x \in \mathbb{R}_0^+$:

$$\text{Exp}(x \mid \lambda) = \lambda \cdot e^{-\lambda x}$$

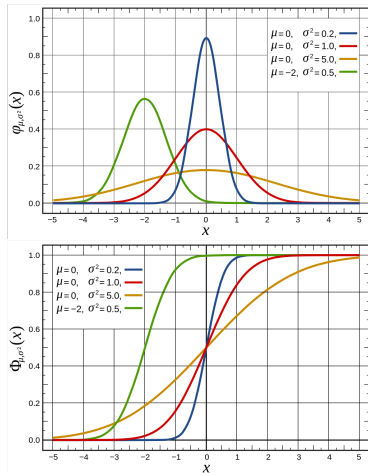


Normal/Gaussian distribution

A Normal or Gaussian random variable $X \sim \mathcal{N}(\mu, \sigma^2)$, $\mu, \sigma \in \mathbb{R}$, $\sigma > 0$ has approximately the same distribution as the sum of many independently, arbitrarily but identically distributed random variables.

For $x \in \mathbb{R}$:

$$\mathcal{N}(x \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

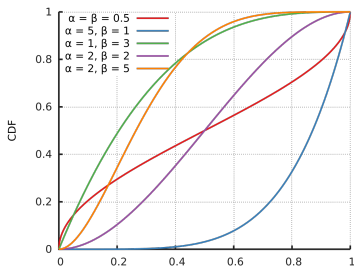
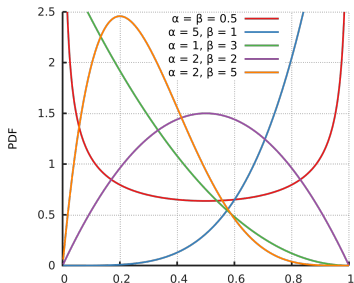


Beta distribution

Random variables $X \sim \text{Beta}(a, b)$, $a, b > 0$, following a Beta distribution can often be seen as the success probability for a binary event.

For $x \in [0, 1]$:

$$\text{Beta}(x \mid a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$$



Gamma distribution

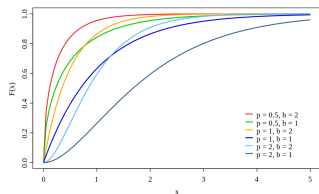
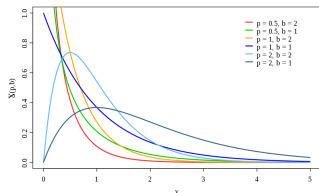
Random variables $X \sim \text{Gamma}(a, b)$ following a Gamma distribution are governed by the parameters $a, b > 0$. For $x > 0$:

$$\text{Gamma}(x | a, b) = \frac{1}{\Gamma(a)} b^a x^{a-1} e^{-bx}$$

$$\Gamma(a) = \int_0^{\infty} t^{a-1} e^{-t} dt$$

$$\Gamma(n+1) = n! \text{ for } n \in \mathbb{N}_0^+$$

The Gamma distribution is the conjugate prior for the precision (inverse variance) of a univariate Gauss distribution.



Overview: probability distributions

Distribution	Notation	Param.	Co-dom.	PMF / PDF	Mean	Variance
Bernoulli*	$\text{Ber}(\mu)$	$\mu \in [0, 1]$	$x \in \{0, 1\}$	$\mu^x(1 - \mu)^{1-x}$	μ	$\mu(1 - \mu)$
Binomial*	$\text{Bin}(N, \mu)$	$N \geq 1, \mu \in [0, 1]$	$x \in \{0, 1, \dots, N\}$	$\binom{N}{x} \mu^x (1 - \mu)^{N-x}$	$N\mu$	$N\mu(1 - \mu)$
Poisson*	$\text{Poi}(\lambda)$	$\lambda > 0$	$x \in \mathbb{N}_0^+$	$\frac{e^{-\lambda} \lambda^x}{x!}$	λ	λ
Uniform	$U(a, b)$	$a, b \in \mathbb{R}, a < b$	$x \in [a, b]$	$\frac{1}{b-a}$	$\frac{a+b}{2}$	$\frac{1}{12}(b-a)^2$
Exponential	$\text{Exp}(\lambda)$	$\lambda > 0$	$x \in \mathbb{R}_0^+$	$\lambda e^{-\lambda x}$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$
Normal/Gauss	$\mathcal{N}(\mu, \sigma^2)$	$\mu \in \mathbb{R}, \sigma > 0$	$x \in \mathbb{R}$	$\frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$	μ	σ^2
Beta	$\text{Beta}(a, b)$	$a, b > 0$	$x \in [0, 1]$	$\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$	$\frac{a}{a+b}$	$\frac{ab}{(a+b)^2(a+b+1)}$
Gamma	$\text{Gamma}(a, b)$	$a, b > 0$	$x \in \mathbb{R}_0^+$	$\frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}$	$\frac{a}{b}$	$\frac{a}{b^2}$

*Discrete distributions

With the gamma function $\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt$, with the property that $\Gamma(n+1) = n!$ for $n \in \mathbb{N}_0^+$.

Two random variables—*Bivariate* case

Two random variables X and Y can interact. We need to consider them simultaneously for statistical analysis. To this end, we introduce the *joint cumulative distribution function* of X and Y :

$$F_{XY}(x, y) = P(X \leq x, Y \leq y).$$

$F_X(x)$ and $F_Y(y)$ are the *marginal cumulative distribution function* of $F_{XY}(x, y)$.

Properties:

- ▶ $0 \leq F_{XY}(x, y) \leq 1$
- ▶ $\lim_{x, y \rightarrow -\infty} F_{XY}(x, y) = 0$
- ▶ $\lim_{x, y \rightarrow \infty} F_{XY}(x, y) = 1$
- ▶ $F_X(x) = \lim_{y \rightarrow \infty} F_{XY}(x, y)$

Two *continuous* random variables

Most properties can be defined analogously to the univariate case.

Joint probability density function:

$$f_{XY}(x, y) = \frac{\partial^2 F_{XY}(x, y)}{\partial x \partial y}.$$

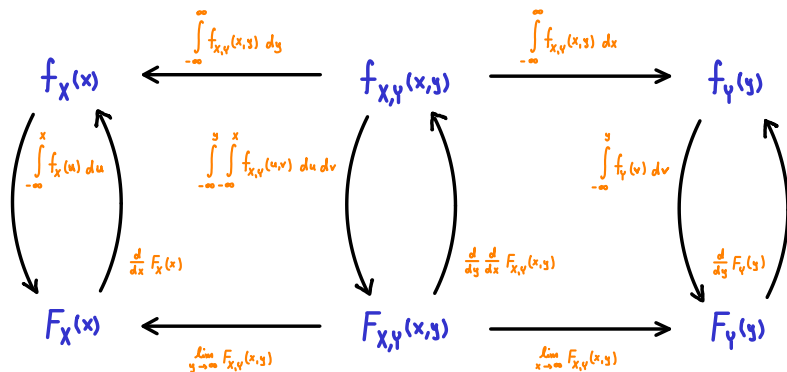
Properties:

- ▶ $f_{XY}(x, y) \geq 0$
- ▶ $\int \int_A f_{XY}(x, y) \, dx dy = P((X, Y) \in A)$
- ▶ $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{XY}(x, y) \, dx dy = 1$

If we remove the effect of one of the random variables, we yield the *marginal probability density function* or *marginal density* for short:

$$f_X(x) = \int_{-\infty}^{\infty} f_{XY}(x, y) \, dy.$$

Relations between $f_{X,Y}$, f_X , f_Y , $F_{X,Y}$, F_X and F_Y



Two *discrete* random variables

Joint probability mass function:

$$p_{XY}(x, y) = P(X = x, Y = y).$$

Properties:

- ▶ $0 \leq p_{XY}(x, y) \leq 1$
- ▶ $\sum_x \sum_y p_{XY}(x, y) = 1$

In order to get the *marginal probability mass function* $p_X(x)$, we need to *sum out* all possible y (*marginalization*):

$$p_X(x) = \sum_y p_{XY}(x, y).$$

Conditional distributions/Bayes' rule

	discrete	continuous
Definition	$p_{Y X}(y x) = \frac{p_{XY}(x,y)}{p_X(x)}$	$f_{Y X}(y x) = \frac{f_{XY}(x,y)}{f_X(x)}$
Bayes' rule	$p_{Y X}(y x) = \frac{p_{X Y}(x y)p_Y(y)}{p_X(x)}$	$f_{Y X}(y x) = \frac{f_{X Y}(x y)f_Y(y)}{f_X(x)}$
Probabilities	$p_{Y X}(y x) = P(Y = y X = x)$	$P(Y \in A X = x) = \int_A f_{Y X}(y x) dy$

Independence

Two random variables X , Y are *independent* if $F_{XY}(x, y) = F_X(x)F_Y(y)$ for all values x and y .

Equivalently:

- ▶ $p_{XY}(x, y) = p_X(x) p_Y(y)$
- ▶ $p_{Y|X}(y | x) = p_Y(y)$
- ▶ $f_{XY}(x, y) = f_X(x)f_Y(y)$
- ▶ $f_{Y|X}(y | x) = f_Y(y)$

Independent and identically distributed—i.i.d.

If two random variables X and Y are called *identically distributed* it means that the following holds:

$$\begin{aligned}f_X(x) &= f_Y(x), \\F_X(x) &= F_Y(x).\end{aligned}$$

As a consequence (among many others):

$$\begin{aligned}E[X] &= E[Y], \\ \text{Var}(X) &= \text{Var}(Y).\end{aligned}$$

It does **not** mean that $X = Y$ is true! X and Y following the same distribution does not imply that they always provide the same values! If X and Y are also independent, we call them *independent and identically distributed* (i.i.d.).

Expectation and covariance

Given two random variables X , Y and $g : \mathbb{R}^2 \rightarrow \mathbb{R}$.

- ▶ $E[g(X, Y)] := \sum_x \sum_y g(x, y) p_{XY}(x, y)$.
- ▶ $E[g(X, Y)] := \int_{-\infty}^{\infty} g(x, y) f_{XY}(x, y) dx dy$.

Covariance

- ▶ $\text{Cov}(X, Y) := E[(X - E[X])(Y - E[Y])] = E[XY] - E[X]E[Y]$.
- ▶ When $\text{Cov}(X, Y) = 0$, X and Y are *uncorrelated*.
- ▶ *Pearson* correlation coefficient $\rho(X, Y)$:

$$\rho(X, Y) := \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} \in [-1, 1].$$

- ▶ $E[f(X, Y) + g(X, Y)] = E[f(X, Y)] + E[g(X, Y)]$.
- ▶ $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2 \text{Cov}(X, Y)$.
- ▶ If X and Y are independent, then $\text{Cov}(X, Y) = 0$.
- ▶ If X and Y are independent, then $E[f(X)g(Y)] = E[f(X)]E[g(Y)]$.

Multiple random variables—Random vectors

Generalize previous ideas to more than two random variables. Putting all these random variables together in one vector \mathbf{X} , a *random vector* ($\mathbf{X} : \Omega \rightarrow \mathbb{R}^n$). The notions of joint CDF and PDF apply equivalently, e.g.

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

Expectation of a continuous random vector for $g : \mathbb{R}^n \rightarrow \mathbb{R}$:

$$E[g(\mathbf{X})] = \int_{\mathbb{R}^n} g(x_1, x_2, \dots, x_n) f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n) dx_1 dx_2 \dots dx_n.$$

If $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ then the expected value of g is the *element-wise* values of the output vector:

$$E[g(\mathbf{X})] = \begin{bmatrix} E[g_1(\mathbf{X})] \\ E[g_2(\mathbf{X})] \\ \vdots \\ E[g_m(\mathbf{X})] \end{bmatrix}.$$

Independence of more than two random variables

The random variables X_1, \dots, X_n are independent if **for all** subsets $I = \{i_1, \dots, i_k\} \subset \{1, \dots, N\}$ and all $(x_{i_1}, \dots, x_{i_k})$

$$f_{X_{i_1}, \dots, X_{i_k}}(x_{i_1}, \dots, x_{i_k}) = f_{X_{i_1}}(x_{i_1}) \cdot \dots \cdot f_{X_{i_k}}(x_{i_k}),$$

or equivalently

$$F_{X_{i_1}, \dots, X_{i_k}}(x_{i_1}, \dots, x_{i_k}) = F_{X_{i_1}}(x_{i_1}) \cdot \dots \cdot F_{X_{i_k}}(x_{i_k}),$$

hold.

If there exists a combination of values so that the equations above do not hold then the random variables are not independent.

For better distinction, this notion of independence is sometimes called *mutual independence*.

Covariance matrix

For a random vector $\mathbf{X} : \Omega \rightarrow \mathbb{R}^n$, the *covariance matrix* Σ is the $n \times n$ square *symmetric, positive definite* matrix whose entries are

$$\Sigma_{ij} = \text{Cov}(X_i, X_j).$$

$$\Sigma = \text{E}[(\mathbf{X} - \text{E}[\mathbf{X}])(\mathbf{X} - \text{E}[\mathbf{X}])^T] = \text{E}[\mathbf{X}\mathbf{X}^T] - \text{E}[\mathbf{X}]\text{E}[\mathbf{X}]^T$$

Multinomial distribution

The multivariate version of the Binomial is called a *Multinomial*, $\mathbf{X} \sim \text{Multinomial}(N, \boldsymbol{\mu})$. We have $k \geq 1$ mutually exclusive events with a success probability of μ_k (such that $\sum_{i=1}^k \mu_k = 1$). We draw N times independently.

$$p_{\mathbf{X}}(x_1, x_2, \dots, x_k) = \binom{N}{x_1 \ x_2 \ \dots \ x_k} \mu_1^{x_1} \mu_2^{x_2} \dots \mu_k^{x_k}.$$

where

$$\begin{aligned} \sum_k x_k &= N, \\ \binom{N}{x_1 \ x_2 \ \dots \ x_k} &= \frac{N!}{x_1! \ x_2! \ \dots \ x_k!}, \\ \mathbb{E}[\mathbf{X}] &= (N\mu_1, N\mu_2, \dots, N\mu_k), \\ \text{Var}(X_i) &= N\mu_i(1 - \mu_i), \\ \text{Cov}(X_i, X_j) &= -N\mu_i\mu_j. \end{aligned}$$

Example: An urn with n balls of $k \geq 1$ different labels, drawn $N \geq 1$ times with replacement and probabilities $\mu_k = \frac{\#k}{n}$. The marginal distribution of X_i is $\text{Bin}(n, \mu_i)$.

Multivariate Gaussian

The multivariate version of the Gaussian $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is very similar to the univariate, except that it allows for dependencies between the individual components. $\boldsymbol{\mu} \in \mathbb{R}^k$ is the mean vector, the positive definite, symmetric $\boldsymbol{\Sigma} \in \mathbb{R}^{k \times k}$ the covariance matrix

$$f_{\mathbf{X}}(x_1, x_2, \dots, x_k) = \frac{1}{\sqrt{(2\pi)^k \det \boldsymbol{\Sigma}}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right),$$

where

$$\mathbb{E}[\mathbf{X}] = \boldsymbol{\mu},$$

$$\text{Var}(X_i) = \boldsymbol{\Sigma}_{ii},$$

$$\text{Cov}(X_i, X_j) = \boldsymbol{\Sigma}_{ij}.$$

The marginal distribution of X_i is $\mathcal{N}(\mu_i, \boldsymbol{\Sigma}_{ii})$.

Notation in the lecture

Consider

$$p_X(x), \quad x \in \mathbb{R} \quad \text{vs} \quad p_X(y), \quad y \in \mathbb{R}$$

$p_X(x), f_X(x), p_{XY}(x, y), f_{XY}(x, y)$ are written as $p(x)$ or $p(x, y)$.
Likewise $p_Y(y)$ is written as $p(y)$.