# All Cheat Sheets
## Machine Learning, Deep Learning, Artificial Intelligence

BY

STANFORD UNIVERSITY

AND

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

COMPILED BY -  NIKHIL YADAV

# Probability–the Science of Uncertainty and Data
by Fabián Kozynski

---

## PROBABILITY

### Probability models and axioms

**Definition (Sample space)** A sample space $\Omega$ is the set of all possible outcomes. The set's elements must be mutually exclusive, collectively exhaustive and at the right granularity.

**Definition (Event)** An event is a subset of the sample space. Probability is assigned to events.

**Definition (Probability axioms)** A probability law $\mathbb{P}$ assigns probabilities to events and satisfies the following axioms:

**Nonnegativity** $\mathbb{P}(A) \geq 0$ for all events $A$.

**Normalization** $\mathbb{P}(\Omega) = 1$.

**(Countable) additivity** For every sequence of events $A_1, A_2, \ldots$ such that $A_i \cap A_j = \varnothing$: $\mathbb{P}\left(\bigcup_i A_i\right) = \sum_i \mathbb{P}(A_i)$.

**Corollaries (Consequences of the axioms)**

- $\mathbb{P}(\varnothing) = 0$.
- For any finite collection of disjoint events $A_1, \ldots, A_n$,
  $\mathbb{P}\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n \mathbb{P}(A_i)$.
- $\mathbb{P}(A) + \mathbb{P}(A^c) = 1$.
- $\mathbb{P}(A) \leq 1$.
- If $A \subset B$, then $\mathbb{P}(A) \leq \mathbb{P}(B)$.
- $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B)$.
- $\mathbb{P}(A \cup B) \leq \mathbb{P}(A) + \mathbb{P}(B)$.

**Example (Discrete uniform law)** Assume $\Omega$ is finite and consists of $n$ equally likely elements. Also, assume that $A \subset \Omega$ with $k$ elements. Then $\mathbb{P}(A) = \frac{k}{n}$.

### Conditioning and Bayes' rule

**Definition (Conditional probability)** Given that event $B$ has occurred and that $P(B) > 0$, the probability that $A$ occurs is

$$\mathbb{P}(A|B) \triangleq \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}.$$

**Remark (Conditional probabilities properties)** They are the same as ordinary probabilities. Assuming $\mathbb{P}(B) > 0$:

- $\mathbb{P}(A|B) \geq 0$.
- $\mathbb{P}(\Omega|B) = 1$
- $\mathbb{P}(B|B) = 1$.
- If $A \cap C = \varnothing$, $\mathbb{P}(A \cup C|B) = \mathbb{P}(A|B) + \mathbb{P}(C|B)$.

**Proposition (Multiplication rule)**

$\mathbb{P}(A_1 \cap A_2 \cap \cdots \cap A_n) = \mathbb{P}(A_1) \cdot \mathbb{P}(A_2|A_1) \cdots \mathbb{P}(A_n|A_1 \cap A_2 \cap \cdots \cap A_{n-1})$.

**Theorem (Total probability theorem)** Given a partition $\{A_1, A_2, \ldots\}$ of the sample space, meaning that $\bigcup_i A_i = \Omega$ and the events are disjoint, and for every event $B$, we have

$$\mathbb{P}(B) = \sum_i \mathbb{P}(A_i)\mathbb{P}(B|A_i).$$

**Theorem (Bayes' rule)** Given a partition $\{A_1, A_2, \ldots\}$ of the sample space, meaning that $\bigcup_i A_i = \Omega$ and the events are disjoint, and if $\mathbb{P}(A_i) > 0$ for all $i$, then for every event $B$, the conditional probabilities $\mathbb{P}(A_i|B)$ can be obtained from the conditional probabilities $\mathbb{P}(B|A_i)$ and the initial probabilities $\mathbb{P}(A_i)$ as follows:

$$\mathbb{P}(A_i|B) = \frac{\mathbb{P}(A_i)\mathbb{P}(B|A_i)}{\sum_j \mathbb{P}(A_j)\mathbb{P}(B|A_j)}.$$

### Independence

**Definition (Independence of events)** Two events are independent if occurrence of one provides no information about the other. We say that $A$ and $B$ are independent if

$$\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B).$$

Equivalently, as long as $\mathbb{P}(A) > 0$ and $\mathbb{P}(B) > 0$,

$$\mathbb{P}(B|A) = \mathbb{P}(B) \qquad \mathbb{P}(A|B) = \mathbb{P}(A).$$

**Remarks**

- The definition of independence is symmetric with respect to $A$ and $B$.
- The product definition applies even if $\mathbb{P}(A) = 0$ or $\mathbb{P}(B) = 0$.

**Corollary** If $A$ and $B$ are independent, then $A$ and $B^c$ are independent. Similarly for $A^c$ and $B$, or for $A^c$ and $B^c$.

**Definition (Conditional independence)** We say that $A$ and $B$ are independent conditioned on $C$, where $\mathbb{P}(C) > 0$, if

$$\mathbb{P}(A \cap B|C) = \mathbb{P}(A|C)\mathbb{P}(B|C).$$

**Definition (Independence of a collection of events)** We say that events $A_1, A_2, \ldots, A_n$ are independent if for every collection of distinct indices $i_1, i_2, \ldots, i_k$, we have

$$\mathbb{P}(A_{i_1} \cap \ldots \cap A_{i_k}) = \mathbb{P}(A_{i_1}) \cdot \mathbb{P}(A_{i_2}) \cdots \mathbb{P}(A_{i_k}).$$

### Counting

This section deals with finite sets with uniform probability law. In this case, to calculate $\mathbb{P}(A)$, we need to count the number of elements in $A$ and in $\Omega$.

**Remark (Basic counting principle)** For a selection that can be done in $r$ stages, with $n_i$ choices at each stage $i$, the number of possible selections is $n_1 \cdot n_2 \cdots n_r$.

**Definition (Permutations)** The number of permutations (orderings) of $n$ different elements is

$$n! = 1 \cdot 2 \cdot 3 \cdots n.$$

**Definition (Combinations)** Given a set of $n$ elements, the number of subsets with exactly $k$ elements is

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}.$$

**Definition (Partitions)** We are given an $n$−element set and nonnegative integers $n_1, n_2, \ldots, n_r$, whose sum is equal to $n$. The number of partitions of the set into $r$ disjoint subsets, with the $i^{\text{th}}$ subset containing exactly $n_i$ elements, is equal to

$$\binom{n}{n_1, \ldots, n_r} = \frac{n!}{n_1! n_2! \cdots n_r!}.$$

**Remark** This is the same as counting how to assign $n$ distinct elements to $r$ people, giving each person $i$ exactly $n_i$ elements.

### Discrete random variables

*Probability mass function and expectation*

**Definition (Random variable)** A random variable $X$ is a function of the sample space $\Omega$ into the real numbers (or $\mathbb{R}^n$). Its range can be discrete or continuous.

**Definition (Probability mass funtion (PMF))** The probability law of a discrete random variable $X$ is called its PMF. It is defined as

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

**Properties**

$p_X(x) \geq 0, \forall x$.

$\sum_x p_X(x) = 1$.

**Example (Bernoulli random variable)** A Bernoulli random variable $X$ with parameter $0 \leq p \leq 1$ ($X \sim \text{Ber}(p)$) takes the following values:

$$X = \begin{cases} 1 & \text{w.p. } p, \\ 0 & \text{w.p. } 1-p. \end{cases}$$

An indicator random variable of an event ($I_A = 1$ if $A$ occurs) is an example of a Bernoulli random variable.

**Example (Discrete uniform random variable)** A Discrete uniform random variable $X$ between $a$ and $b$ with $a \leq b$ ($X \sim \text{Uni}[a, b]$) takes any of the values in $\{a, a+1, \ldots, b\}$ with probability $\frac{1}{b-a+1}$.

**Example (Binomial random variable)** A Binomial random variable $X$ with parameters $n$ (natural number) and $0 \leq p \leq 1$ ($X \sim \text{Bin}(n, p)$) takes values in the set $\{0, 1, \ldots, n\}$ with probabilities $p_X(i) = \binom{n}{i}p^i(1-p)^{n-i}$.

It represents the number of successes in $n$ independent trials where each trial has a probability of success $p$. Therefore, it can also be seen as the sum of $n$ independent Bernoulli random variables, each with parameter $p$.

**Example (Geometric random variable)** A Geometric random variable $X$ with parameter $0 \leq p \leq 1$ ($X \sim \text{Geo}(p)$) takes values in the set $\{1, 2, \ldots\}$ with probabilities $p_X(i) = (1-p)^{i-1}p$.

It represents the number of independent trials until (and including) the first success, when the probability of success in each trial is $p$.

**Definition (Expectation/mean of a random variable)** The expectation of a discrete random variable is defined as

$$\mathbb{E}[X] \triangleq \sum_x x p_X(x).$$

assuming $\sum_x |x| p_X(x) < \infty$.

**Properties (Properties of expectation)**

- If $X \geq 0$ then $\mathbb{E}[X] \geq 0$.
- If $a \leq X \leq b$ then $a \leq \mathbb{E}[X] \leq b$.
- If $X = c$ then $\mathbb{E}[X] = c$.

**Example** Expected value of know r.v.

- If $X \sim \text{Ber}(p)$ then $\mathbb{E}[X] = p$.
- If $X = I_A$ then $\mathbb{E}[X] = \mathbb{P}(A)$.
- If $X \sim \text{Uni}[a, b]$ then $\mathbb{E}[X] = \frac{a+b}{2}$.
- If $X \sim \text{Bin}(n, p)$ then $\mathbb{E}[X] = np$.
- If $X \sim \text{Geo}(p)$ then $\mathbb{E}[X] = \frac{1}{p}$.

**Theorem (Expected value rule)** Given a random variable $X$ and a function $g : \mathbb{R} \to \mathbb{R}$, we construct the random variable $Y = g(X)$. Then
$$\sum_y y p_Y(y) = \mathbb{E}[Y] = \mathbb{E}[g(X)] = \sum_x g(x) p_X(x).$$

**Remark (PMF of $Y = g(X)$)** The PMF of $Y = g(X)$ is
$p_Y(y) = \sum_{x:g(x)=y} p_X(x)$.

**Remark** In general $g(\mathbb{E}[X]) \neq \mathbb{E}[g(X)]$. They are equal if $g(x) = ax + b$.

*Variance, conditioning on an event, multiple r.v.*

**Definition (Variance of a random variable)** Given a random variable $X$ with $\mu = \mathbb{E}[X]$, its variance is a measure of the spread of the random variable and is defined as
$$\mathrm{Var}(X) \overset{\triangle}{=} \mathbb{E}\left[(X - \mu)^2\right] = \sum_x (x - \mu)^2 p_X(x).$$

**Definition (Standard deviation)**
$$\sigma_X = \sqrt{\mathrm{Var}(X)}.$$

**Properties (Properties of the variance)**
- $\mathrm{Var}(aX) = a^2 \mathrm{Var}(X)$, for all $a \in \mathbb{R}$.
- $\mathrm{Var}(X + b) = \mathrm{Var}(X)$, for all $b \in \mathbb{R}$.
- $\mathrm{Var}(aX + b) = a^2 \mathrm{Var}(X)$.
- $\mathrm{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$.

**Example (Variance of known r.v.)**
- If $X \sim \mathrm{Ber}(p)$, then $\mathrm{Var}(X) = p(1 - p)$.
- If $X \sim \mathrm{Uni}[a, b]$, then $\mathrm{Var}(X) = \frac{(b-a)(b-a+2)}{12}$.
- If $X \sim \mathrm{Bin}(n, p)$, then $\mathrm{Var}(X) = np(1 - p)$.
- If $X \sim \mathrm{Geo}(p)$, then $\mathrm{Var}(X) = \frac{1-p}{p^2}$

**Proposition (Conditional PMF and expectation, given an event)** Given the event $A$, with $\mathbb{P}(A) > 0$, we have the following
- $p_{X|A}(x) = \mathbb{P}(X = x | A)$.
- If $A$ is a subset of the range of $X$, then:
$$p_{X|A}(x) \overset{\triangle}{=} p_{X|\{X \in A\}}(x) = \begin{cases} \frac{1}{\mathbb{P}(A)} p_X(x), & \text{if } x \in A, \\ 0, & \text{otherwise.} \end{cases}$$
- $\sum_x p_{X|A}(x) = 1$.
- $\mathbb{E}[X|A] = \sum_x x p_{X|A}(x)$.
- $\mathbb{E}[g(X)|A] = \sum_x g(x) p_{X|A}(x)$.

**Proposition (Total expectation rule)** Given a partition of disjoint events $A_1, \ldots, A_n$ such that $\sum_i \mathbb{P}(A_i) = 1$, and $\mathbb{P}(A_i) > 0$,
$$\mathbb{E}[X] = \mathbb{P}(A_1) \mathbb{E}[X|A_1] + \cdots + \mathbb{P}(A_n) \mathbb{E}[X|A_n].$$

**Definition (Memorylessness of the geometric random variable)** When we condition a geometric random variable $X$ on the event $X > n$ we have memorylessness, meaning that the "remaining time" $X - n$, given that $X > n$, is also geometric with the same parameter. Formally,
$$p_{X-n|X>n}(i) = p_X(i).$$

**Definition (Joint PMF)** The joint PMF of random variables $X_1, X_2, \ldots, X_n$ is
$p_{X_1, X_2, \ldots, X_n}(x_1, \ldots, x_n) = \mathbb{P}(X_1 = x_1, \ldots, X_n = x_n)$.

**Properties (Properties of joint PMF)**
- $\sum_{x_1} \cdots \sum_{x_n} p_{X_1, \ldots, X_n}(x_1, \ldots, x_n) = 1$.
- $p_{X_1}(x_1) = \sum_{x_2} \cdots \sum_{x_n} p_{X_1, \ldots, X_n}(x_1, x_2, \ldots, x_n)$.
- $p_{X_2, \ldots, X_n}(x_2, \ldots, x_n) = \sum_{x_1} p_{X_1, X_2, \ldots, X_n}(x_1, x_2, \ldots, x_n)$.

**Definition (Functions of multiple r.v.)** If $Z = g(X_1, \ldots, X_n)$, where $g : \mathbb{R}^n \to \mathbb{R}$, then $p_Z(z) = \mathbb{P}(g(X_1, \ldots, X_n) = z)$.

**Proposition (Expected value rule for multiple r.v.)** Given $g : \mathbb{R}^n \to \mathbb{R}$,
$$\mathbb{E}[g(X_1, \ldots, X_n)] = \sum_{x_1, \ldots, x_n} g(x_1, \ldots, x_n) p_{X_1, \ldots, X_n}(x_1, \ldots, x_n).$$

**Properties (Linearity of expectations)**
- $\mathbb{E}[aX + b] = a\mathbb{E}[X] + b$.
- $\mathbb{E}[X_1 + \cdots + X_n] = \mathbb{E}[X_1] + \cdots + \mathbb{E}[X_n]$.

*Conditioning on a random variable, independence*

**Definition (Conditional PMF given another random variable)** Given discrete random variables $X, Y$ and $y$ such that $p_Y(y) > 0$ we define
$$p_{X|Y}(x|y) \overset{\triangle}{=} \frac{p_{X,Y}(x, y)}{p_Y(y)}.$$

**Proposition (Multiplication rule)** Given jointly discrete random variables $X, Y$, and whenever the conditional probabilities are defined,
$$p_{X,Y}(x, y) = p_X(x) p_{Y|X}(y|x) = p_Y(y) p_{X|Y}(x|y).$$

**Definition (Conditional expectation)** Given discrete random variables $X, Y$ and $y$ such that $p_Y(y) > 0$ we define
$$\mathbb{E}[X|Y = y] = \sum_x x p_{X|Y}(x|y).$$

Additionally we have
$$\mathbb{E}[g(X)|Y = y] = \sum_x g(x) p_{X|Y}(x|y).$$

**Theorem (Total probability and expectation theorems)** If $p_Y(y) > 0$, then
$$p_X(x) = \sum_y p_Y(y) p_{X|Y}(x|y),$$
$$\mathbb{E}[X] = \sum_y p_Y(y) \mathbb{E}[X|Y = y].$$

**Definition (Independence of a random variable and an event)** A discrete random variable $X$ and an event $A$ are independent if $\mathbb{P}(X = x \text{ and } A) = p_X(x) \mathbb{P}(A)$, for all $x$.

**Definition (Independence of two random variables)** Two discrete random variables $X$ and $Y$ are independent if $p_{X,Y}(x, y) = p_X(x) p_Y(y)$ for all $x, y$.

**Remark (Independence of a collection of random variables)** A collection $X_1, X_2, \ldots, X_n$ of random variables are independent if
$$p_{X_1, \ldots, X_n}(x_1, \ldots, x_n) = p_{X_1}(x_1) \cdots p_{X_n}(x_n), \ \forall \, x_1, \ldots, x_n.$$

**Remark (Independence and expectation)** In general, $\mathbb{E}[g(X, Y)] \neq g(\mathbb{E}[X], \mathbb{E}[Y])$. An exception is for linear functions: $\mathbb{E}[aX + bY] = a\mathbb{E}[X] + b\mathbb{E}[Y]$.

**Proposition (Expectation of product of independent r.v.)** If $X$ and $Y$ are discrete independent random variables,
$$\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y].$$

**Remark** If $X$ and $Y$ are independent, $\mathbb{E}[g(X)h(Y)] = \mathbb{E}[g(X)]\mathbb{E}[h(Y)]$.

**Proposition (Variance of sum of independent random variables)** IF $X$ and $Y$ are discrete independent random variables,
$$\mathrm{Var}(X + Y) = \mathrm{Var}(X) + \mathrm{Var}(Y).$$

## Continuous random variables

*PDF, Expectation, Variance, CDF*

**Definition (Probability density function (PDF))** A probability density function of a r.v. $X$ is a non-negative real valued function $f_X$ that satisfies the following
- $\int_{-\infty}^{\infty} f_X(x)\mathrm{d}x = 1$.
- $\mathbb{P}(a \leq X \leq b) = \int_a^b f_X(x)\mathrm{d}x$ for some random variable $X$.

**Definition (Continuous random variable)** A random variable $X$ is continuous if its probability law can be described by a PDF $f_X$.

**Remark** Continuous random variables satisfy:
- For small $\delta > 0$, $\mathbb{P}(a \leq X \leq a + \delta) \approx f_X(a)\delta$.
- $\mathbb{P}(X = a) = 0$, $\forall a \in \mathbb{R}$.

**Definition (Expectation of a continuous random variable)** The expectation of a continuous random variable is
$$\mathbb{E}[X] \overset{\triangle}{=} \int_{-\infty}^{\infty} x f_X(x)\mathrm{d}x.$$
assuming $\int_{-\infty}^{\infty} |x| f_X(x)\mathrm{d}x < \infty$.

**Properties (Properties of expectation)**
- If $X \geq 0$ then $\mathbb{E}[X] \geq 0$.
- If $a \leq X \leq b$ then $a \leq \mathbb{E}[X] \leq b$.
- $\mathbb{E}[g(X)] = \int_{-\infty}^{\infty} g(x) f_X(x)\mathrm{d}x$.
- $\mathbb{E}[aX + b] = a\mathbb{E}[X] + b$.

**Definition (Variance of a continuous random variable)** Given a continuous random variable $X$ with $\mu = \mathbb{E}[X]$, its variance is
$$\mathrm{Var}(X) = \mathbb{E}\left[(X - \mu)^2\right] = \int_{-\infty}^{\infty} (x - \mu)^2 f_X(x)\mathrm{d}x.$$

It has the same properties as the variance of a discrete random variable.

**Example (Uniform continuous random variable)** A Uniform continuous random variable $X$ between $a$ and $b$, with $a < b$, ($X \sim \mathrm{Uni}(a, b)$) has PDF
$$f_X(x) = \begin{cases} \frac{1}{b-a}, & \text{if } a < x < b, \\ 0, & \text{otherwise.} \end{cases}$$

We have $\mathbb{E}[X] = \frac{a+b}{2}$ and $\mathrm{Var}(X) = \frac{(b-a)^2}{12}$.

**Example (Exponential random variable)** An Exponential random variable $X$ with parameter $\lambda > 0$ ($X \sim Exp(\lambda)$) has PDF

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x}, & \text{if } x \geq 0, \\ 0, & \text{otherwise.} \end{cases}$$

We have $E[X] = \frac{1}{\lambda}$ and $\text{Var}(X) = \frac{1}{\lambda^2}$.

**Definition (Cumulative Distribution Function (CDF))** The CDF of a random variable $X$ is $F_X(x) = \mathbb{P}(X \leq x)$.
In particular, for a continuous random variable, we have

$$F_X(x) = \int_{-\infty}^{x} f_X(x)\mathrm{d}x,$$

$$f_X(x) = \frac{\mathrm{d}F_X(x)}{\mathrm{d}x}.$$

**Properties (Properties of CDF)**

- If $y \geq x$, then $F_X(y) \geq F_X(x)$.
- $\lim_{x \to -\infty} F_X(x) = 0$.
- $\lim_{x \to \infty} F_X(x) = 1$.

**Definition (Normal/Gaussian random variable)** A Normal random variable $X$ with mean $\mu$ and variance $\sigma^2 > 0$ ($X \sim \mathcal{N}(\mu, \sigma^2)$) has PDF

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

We have $E[X] = \mu$ and $\text{Var}(X) = \sigma^2$.

**Remark (Standard Normal)** The standard Normal is $\mathcal{N}(0, 1)$.

**Proposition (Linearity of Gaussians)** Given $X \sim \mathcal{N}(\mu, \sigma^2)$, and if $a \neq 0$, then $aX + b \sim \mathcal{N}(a\mu + b, a^2\sigma^2)$.

Using this $Y = \frac{X-\mu}{\sigma}$ is a standard gaussian.

*Conditioning on an event, and multiple continuous r.v.*

**Definition (Conditional PDF given an event)** Given a continuous random variable $X$ and event $A$ with $P(A) > 0$, we define the conditional PDF as the function that satisfies

$$\mathbb{P}(X \in B|A) = \int_B f_{X|A}(x)\mathrm{d}x.$$

**Definition (Conditional PDF given $X \in A$)** Given a continuous random variable $X$ and an $A \subset \mathbb{R}$, with $P(A) > 0$:

$$f_{X|X \in A}(x) = \begin{cases} \frac{1}{\mathbb{P}(A)} f_X(x), & x \in A, \\ 0, & x \notin A. \end{cases}$$

**Definition (Conditional expectation)** Given a continuous random variable $X$ and an event $A$, with $P(A) > 0$:

$$\mathbb{E}[X|A] = \int_{-\infty}^{\infty} f_{X|A}(x)\mathrm{d}x.$$

**Definition (Memorylessness of the exponential random variable)** When we condition an exponential random variable $X$ on the event $X > t$ we have memorylessness, meaning that the "remaining time" $X - t$ given that $X > t$ is also geometric with the same parameter i.e.,

$$\mathbb{P}(X - t > x|X > t) = \mathbb{P}(X > x).$$

**Theorem (Total probability and expectation theorems)** Given a partition of the space into disjoint events $A_1, A_2, \ldots, A_n$ such that $\sum_i \mathbb{P}(A_i) = 1$ we have the following:

$$F_X(x) = \mathbb{P}(A_1)F_{X|A_1}(x) + \cdots + \mathbb{P}(A_n)F_{X|A_n}(x),$$
$$f_X(x) = \mathbb{P}(A_1)f_{X|A_1}(x) + \cdots + \mathbb{P}(A_n)f_{X|A_n}(x),$$
$$\mathbb{E}[X] = \mathbb{P}(A_1)\mathbb{E}[X|A_1] + \cdots + \mathbb{P}(A_n)\mathbb{E}[X|A_n].$$

**Definition (Jointly continuous random variables)** A pair (collection) of random variables is jointly continuous if there exists a joint PDF $f_{X,Y}$ that describes them, that is, for every set $B \subset \mathbb{R}^n$

$$\mathbb{P}((X, Y) \in B) = \iint_B f_{X,Y}(x, y)\mathrm{d}x\mathrm{d}y.$$

**Properties (Properties of joint PDFs)**

- $f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x, y)\mathrm{d}y.$

- $F_{X,Y}(x, y) = \mathbb{P}(X \leq x, Y \leq y) = \int_{-\infty}^{x}\left[\int_{-\infty}^{y} f_{X,Y}(u, v)\mathrm{d}v\right]\mathrm{d}u.$

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

**Example (Uniform joint PDF on a set $S$)** Let $S \subset \mathbb{R}^2$ with area $s > 0$, then the random variable $(X, Y)$ is uniform over $S$ if it has PDF

$$f_{X,Y}(x, y) = \begin{cases} \frac{1}{s}, & (x, y) \in S, \\ 0, & (x, y) \notin S. \end{cases}$$

*Conditioning on a random variable, independence, Bayes' rule*

**Definition (Conditional PDF given another random variable)** Given jointly continuous random variables $X, Y$ and a value $y$ such that $f_Y(y) > 0$, we define the conditional PDF as

$$f_{X|Y}(x|y) \triangleq \frac{f_{X,Y}(x, y)}{f_Y(y)}.$$

Additionally we define $\mathbb{P}(X \in A|Y = y) \int_A f_{X|Y}(x|y)\mathrm{d}x.$

**Proposition (Multiplication rule)** Given jointly continuous random variables $X, Y$, whenever possible we have

$$f_{X,Y}(x, y) = f_X(x)f_{Y|X}(y|x) = f_Y(y)f_{X|Y}(x|y).$$

**Definition (Conditional expectation)** Given jointly continuous random variables $X, Y$, and $y$ such that $f_Y(y) > 0$, we define the conditional expected value as

$$\mathbb{E}[X|Y = y] = \int_{-\infty}^{\infty} x f_{X|Y}(x|y)\mathrm{d}x.$$

Additionally we have

$$\mathbb{E}[g(X)|Y = y] = \int_{-\infty}^{\infty} g(x)f_{X|Y}(x|y)\mathrm{d}x.$$

**Theorem (Total probability and total expectation theorems)**

$$f_X(x) = \int_{-\infty}^{\infty} f_Y(y)f_{X|Y}(x|y)\mathrm{d}y,$$

$$\mathbb{E}[X] = \int_{-\infty}^{\infty} f_Y(y)\mathbb{E}[X|Y = y]\mathrm{d}y.$$

**Definition (Independence)** Jointly continuous random variables $X, Y$ are independent if $f_{X,Y}(x, y) = f_X(x)f_Y(y)$ for all $x, y$.

**Proposition (Expectation of product of independent r.v.)** If $X$ and $Y$ are independent continuous random variables,

$$\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y].$$

**Remark** If $X$ and $Y$ are independent,
$\mathbb{E}[g(X)h(Y)] = \mathbb{E}[g(X)]\mathbb{E}[h(Y)].$

**Proposition (Variance of sum of independent random variables)** If $X$ and $Y$ are independent continuous random variables,

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y).$$

**Proposition (Bayes' rule summary)**

- For $X, Y$ discrete: $p_{X|Y}(x|y) = \frac{p_X(x)p_{Y|X}(y|x)}{p_Y(y)}$.

- For $X, Y$ continuous: $f_{X|Y}(x|y) = \frac{f_X(x)f_{Y|X}(y|x)}{f_Y(y)}$.

- For $X$ discrete, $Y$ continuous: $p_{X|Y}(x|y) = \frac{p_X(x)f_{Y|X}(y|x)}{f_Y(y)}$.

- For $X$ continuous, $Y$ discrete: $f_{X|Y}(x|y) = \frac{f_X(x)p_{Y|X}(y|x)}{p_Y(y)}$.

**Derived distributions**

**Proposition (Discrete case)** Given a discrete random variable $X$ and a function $g$, the r.v. $Y = g(X)$ has PMF

$$p_Y(y) = \sum_{x:g(x)=y} p_X(x).$$

**Remark (Linear function of discrete random variable)** If $g(x) = ax + b$, then $p_Y(y) = p_X\left(\frac{y-b}{a}\right)$.

**Proposition (Linear function of continuous r.v.)** Given a continuous random variable $X$ and $Y = aX + b$, with $a \neq 0$, we have

$$f_Y(y) = \frac{1}{|a|} f_X\left(\frac{y-b}{a}\right).$$

**Corollary (Linear function of normal r.v.)** If $X \sim \mathcal{N}(\mu, \sigma^2)$ and $Y = aX + b$, with $a \neq 0$, then $Y \sim \mathcal{N}(a\mu + b, a^2\sigma^2)$.

**Example (General function of a continuous r.v.)** If $X$ is a continuous random variable and $g$ is any function, to obtain the pdf of $Y = g(X)$ we follow the two-step procedure:

1. Find the CDF of $Y$: $F_Y(y) = \mathbb{P}(Y \leq y) = \mathbb{P}(g(X) \leq y)$.

2. Differentiate the CDF of $Y$ to obtain the PDF: $f_Y(y) = \frac{\mathrm{d}F_Y(y)}{\mathrm{d}y}$.

**Proposition (General formula for monotonic $g$)** Let $X$ be a continuous random variable and $g$ a function that is monotonic wherever $f_X(x) > 0$. The PDF of $Y = g(X)$ is given by

$$f_Y(y) = f_X(h(y))\left|\frac{\mathrm{d}h}{\mathrm{d}y}(y)\right|.$$

where $h = g^{-1}$ in the interval where g is monotonic.

## Sums of independent r.v., covariance and correlation

**Proposition (Discrete case)** Let $X, Y$ be discrete independent random variables and $Z = X + Y$, then the PMF of $Z$ is

$$p_Z(z) = \sum_x p_X(x) p_Y(z - x).$$

**Proposition (Continuous case)** Let $X, Y$ be continuous independent random variables and $Z = X + Y$, then the PDF of $Z$ is

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(x) f_Y(z - x) \mathrm{d}x.$$

**Proposition (Sum of independent normal r.v.)** Let $X \sim \mathcal{N}(\mu_x, \sigma_x^2)$ and $Y \sim \mathcal{N}(\mu_y, \sigma_y^2)$ independent. Then $Z = X + Y \sim \mathcal{N}(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$.

Definition (Covariance) We define the covariance of random variables $X, Y$ as

$$\mathrm{Cov}(X, Y) \triangleq \mathbb{E}\left[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])\right].$$

Properties (Properties of covariance)

- If $X, Y$ are independent, then $\mathrm{Cov}(X, Y) = 0$.
- $\mathrm{Cov}(X, X) = \mathrm{Var}(X)$.
- $\mathrm{Cov}(aX + b, Y) = a\,\mathrm{Cov}(X, Y)$.
- $\mathrm{Cov}(X, Y + Z) = \mathrm{Cov}(X, Y) + \mathrm{Cov}(X, Z)$.
- $\mathrm{Cov}(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$.

**Proposition (Variance of a sum of r.v.)**

$$\mathrm{Var}(X_1 + \cdots + X_n) = \sum_i \mathrm{Var}(X_i) + \sum_{i \neq j} \mathrm{Cov}(X_i, X_j).$$

Definition (Correlation coefficient) We define the correlation coefficient of random variables $X, Y$, with $\sigma_X, \sigma_Y > 0$, as

$$\rho(X, Y) \triangleq \frac{\mathrm{Cov}(X, Y)}{\sigma_X \sigma_Y}.$$

Properties (Properties of the correlation coefficient)

- $-1 \leq \rho \leq 1$.
- If $X, Y$ are independent, then $\rho = 0$.
- $|\rho| = 1$ if and only if $X - \mathbb{E}[X] = c(Y - \mathbb{E}[Y])$.
- $\rho(aX + b, Y) = \mathrm{sign}(a)\rho(X, Y)$.

## Conditional expectation and variance, sum of random number of r.v.

Definition (Conditional expectation as a random variable) Given random variables $X, Y$ the conditional expectation $\mathbb{E}[X|Y]$ is the random variable that takes the value $\mathbb{E}[X|Y = y]$ whenever $Y = y$.

Theorem (Law of iterated expectations)

$$\mathbb{E}\left[\mathbb{E}[X|Y]\right] = \mathbb{E}[X].$$

Definition (Conditional variance as a random variable) Given random variables $X, Y$ the conditional variance $\mathrm{Var}(X|Y)$ is the random variable that takes the value $\mathrm{Var}(X|Y = y)$ whenever $Y = y$.

Theorem (Law of total variance)

$$\mathrm{Var}(X) = \mathbb{E}\left[\mathrm{Var}(X|Y)\right] + \mathrm{Var}\left(\mathbb{E}[X|Y]\right).$$

Proposition (Sum of a random number of independent r.v.) Let $N$ be a nonnegative integer random variable. Let $X, X_1, X_2, \ldots, X_N$ be i.i.d. random variables. Let $Y = \sum_i X_i$. Then

$$\mathbb{E}[Y] = \mathbb{E}[N]\mathbb{E}[X],$$

$$\mathrm{Var}(Y) = \mathbb{E}[N]\,\mathrm{Var}(X) + (\mathbb{E}[X])^2\,\mathrm{Var}(N).$$

---

## CONVERGENCE OF RANDOM VARIABLES

### Inequalities, convergence, and the Weak Law of Large Numbers

Theorem (Markov inequality) Given a random variable $X \geq 0$ and, for every $a > 0$ we have

$$\mathbb{P}(X \geq a) \leq \frac{\mathbb{E}[X]}{a}.$$

Theorem (Chebyshev inequality) Given a random variable $X$ with $\mathbb{E}[X] = \mu$ and $\mathrm{Var}(X) = \sigma^2$, for every $\epsilon > 0$ we have

$$\mathbb{P}\left(|X - \mu| \geq \epsilon\right) \leq \frac{\sigma^2}{\epsilon^2}.$$

Theorem (Weak Law of Large Number (WLLN)) Given a sequence of i.i.d. random variables $\{X_1, X_2, \ldots\}$ with $\mathbb{E}[X_i] = \mu$ and $\mathrm{Var}(X_i) = \sigma^2$, we define

$$M_n = \frac{1}{n} \sum_{i=1}^{n} X_i,$$

for every $\epsilon > 0$ we have

$$\lim_{n \to \infty} \mathbb{P}\left(|M_n - \mu| \geq \epsilon\right) = 0.$$

Definition (Convergence in probability) A sequence of random variables $\{Y_i\}$ converges in probability to the random variable $Y$ if

$$\lim_{n \to \infty} \mathbb{P}\left(|Y_i - Y| \geq \epsilon\right) = 0,$$

for every $\epsilon > 0$.

Properties (Properties of convergence in probability) If $X_n \to a$ and $Y_n \to b$ in probability, then

- $X_n + Y_n \to a + b$.
- If $g$ is a continuous function, then $g(X_n) \to g(a)$.
- $\mathbb{E}[X_n]$ does not always converge to $a$.

## The Central Limit Theorem

Theorem (Central Limit Theorem (CLT)) Given a sequence of independent random variables $\{X_1, X_2, \ldots\}$ with $\mathbb{E}[X_i] = \mu$ and $\mathrm{Var}(X_i) = \sigma^2$, we define

$$Z_n = \frac{1}{\sigma \sqrt{n}} \sum_{i=1}^{n} (X_i - \mu).$$

Then, for every $z$, we have

$$\lim_{n \to \infty} \mathbb{P}(Z_n \leq z) = \mathbb{P}(Z \leq z),$$

where $Z \sim \mathcal{N}(0, 1)$.

Corollary (Normal approximation of a binomial) Let $X \sim Bin(n, p)$ with $n$ large. Then $S_n$ can be approximated by $Z \sim \mathcal{N}(np, np(1 - p))$.

Remark (De Moivre-Laplace 1/2 approximation) Let $X \sim Bin$, then $\mathbb{P}(X = i) = \mathbb{P}\left(i - \frac{1}{2} \leq X \leq i + \frac{1}{2}\right)$ and we can use the CLT to approximate the PMF of $X$.

# Super VIP Cheatsheet: Machine Learning

Afshine Amidi and Shervine Amidi
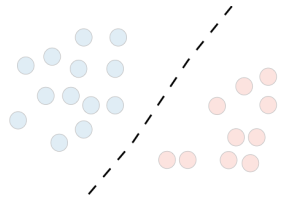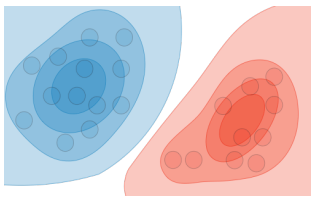
October 6, 2018

## Contents

# 1 Supervised Learning

## 1.1 Introduction to Supervised Learning

Given a set of data points $\{x^{(1)}, ..., x^{(m)}\}$ associated to a set of outcomes $\{y^{(1)}, ..., y^{(m)}\}$, we want to build a classifier that learns how to predict $y$ from $x$.

❏ **Type of prediction** – The different types of predictive models are summed up in the table below:

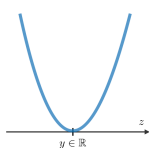|  | Regression | Classifier |
|---|---|---|
| **Outcome** | Continuous | Class |
| **Examples** | Linear regression | Logistic regression, SVM, Naive Bayes |

❏ **Type of model** – The different models are summed up in the table below:

|  | Discriminative model | Generative model |
|---|---|---|
| **Goal** | Directly estimate $P(y\|x)$ | Estimate $P(x\|y)$ to deduce $P(y\|x)$ |
| **What's learned** | Decision boundary | Probability distributions of the data |
| **Illustration** | | |
| **Examples** | Regressions, SVMs | GDA, Naive Bayes |

## 1.2 Notations and general concepts

❏ **Hypothesis** – The hypothesis is noted $h_\theta$ and is the model that we choose. For a given input data $x^{(i)}$, the model prediction output is $h_\theta(x^{(i)})$.

❏ **Loss function** – A loss function is a function $L : (z,y) \in \mathbb{R} \times Y \longmapsto L(z,y) \in \mathbb{R}$ that takes as inputs the predicted value $z$ corresponding to the real data value $y$ and outputs how different they are. The common loss functions are summed up in the table below:
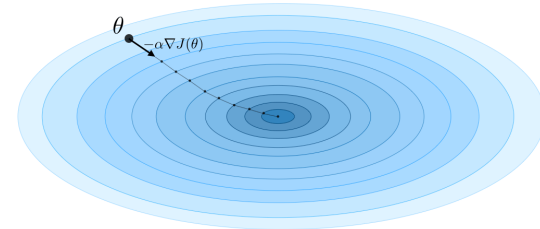
| Least squared | Logistic | Hinge | Cross-entropy |
|---|---|---|---|
| $\frac{1}{2}(y - z)^2$ | $\log(1 + \exp(-yz))$ | $\max(0, 1 - yz)$ | $-\Big[y \log(z) + (1 - y) \log(1 - z)\Big]$ |
| | | | |
| Linear regression | Logistic regression | SVM | Neural Network |

❏ **Cost function** – The cost function $J$ is commonly used to assess the performance of a model, and is defined with the loss function $L$ as follows:

$$J(\theta) = \sum_{i=1}^{m} L(h_\theta(x^{(i)}), y^{(i)})$$

❏ **Gradient descent** – By noting $\alpha \in \mathbb{R}$ the learning rate, the update rule for gradient descent is expressed with the learning rate and the cost function $J$ as follows:

$$\theta \longleftarrow \theta - \alpha \nabla J(\theta)$$



*Remark: Stochastic gradient descent (SGD) is updating the parameter based on each training example, and batch gradient descent is on a batch of training examples.*

❏ **Likelihood** – The likelihood of a model $L(\theta)$ given parameters $\theta$ is used to find the optimal parameters $\theta$ through maximizing the likelihood. In practice, we use the log-likelihood $\ell(\theta) = \log(L(\theta))$ which is easier to optimize. We have:

$$\theta^{\text{opt}} = \arg \max_\theta L(\theta)$$

❏ **Newton's algorithm** – The Newton's algorithm is a numerical method that finds $\theta$ such that $\ell'(\theta) = 0$. Its update rule is as follows:

$$\theta \leftarrow \theta - \frac{\ell'(\theta)}{\ell''(\theta)}$$

*Remark: the multidimensional generalization, also known as the Newton-Raphson method, has the following update rule:*

$$\theta \leftarrow \theta - \left(\nabla_\theta^2 \ell(\theta)\right)^{-1} \nabla_\theta \ell(\theta)$$

## 1.3 Linear models

### 1.3.1 Linear regression

We assume here that $y|x;\theta \sim \mathcal{N}(\mu, \sigma^2)$

❏ **Normal equations** – By noting $X$ the matrix design, the value of $\theta$ that minimizes the cost function is a closed-form solution such that:

$$\theta = (X^T X)^{-1} X^T y$$

❐ **LMS algorithm** – By noting $\alpha$ the learning rate, the update rule of the Least Mean Squares (LMS) algorithm for a training set of $m$ data points, which is also known as the Widrow-Hoff learning rule, is as follows:

$$\forall j, \quad \theta_j \leftarrow \theta_j + \alpha \sum_{i=1}^{m} \left[ y^{(i)} - h_\theta(x^{(i)}) \right] x_j^{(i)}$$

*Remark: the update rule is a particular case of the gradient ascent.*

❐ **LWR** – Locally Weighted Regression, also known as LWR, is a variant of linear regression that weights each training example in its cost function by $w^{(i)}(x)$, which is defined with parameter $\tau \in \mathbb{R}$ as:

$$w^{(i)}(x) = \exp\left( -\frac{(x^{(i)} - x)^2}{2\tau^2} \right)$$

### 1.3.2 Classification and logistic regression

❐ **Sigmoid function** – The sigmoid function $g$, also known as the logistic function, is defined as follows:

$$\forall z \in \mathbb{R}, \quad g(z) = \frac{1}{1 + e^{-z}} \in ]0,1[$$

❐ **Logistic regression** – We assume here that $y|x;\theta \sim \text{Bernoulli}(\phi)$. We have the following form:

$$\phi = p(y = 1|x;\theta) = \frac{1}{1 + \exp(-\theta^T x)} = g(\theta^T x)$$

*Remark: there is no closed form solution for the case of logistic regressions.*

❐ **Softmax regression** – A softmax regression, also called a multiclass logistic regression, is used to generalize logistic regression when there are more than 2 outcome classes. By convention, we set $\theta_K = 0$, which makes the Bernoulli parameter $\phi_i$ of each class $i$ equal to:

$$\phi_i = \frac{\exp(\theta_i^T x)}{\displaystyle\sum_{j=1}^{K} \exp(\theta_j^T x)}$$

### 1.3.3 Generalized Linear Models

❐ **Exponential family** – A class of distributions is said to be in the exponential family if it can be written in terms of a natural parameter, also called the canonical parameter or link function, $\eta$, a sufficient statistic $T(y)$ and a log-partition function $a(\eta)$ as follows:

$$p(y;\eta) = b(y)\exp(\eta T(y) - a(\eta))$$

*Remark: we will often have $T(y) = y$. Also, $\exp(-a(\eta))$ can be seen as a normalization parameter that will make sure that the probabilities sum to one.*

Here are the most common exponential distributions summed up in the following table:

| Distribution | $\eta$ | $T(y)$ | $a(\eta)$ | $b(y)$ |
|---|---|---|---|---|
| Bernoulli | $\log\left(\frac{\phi}{1-\phi}\right)$ | $y$ | $\log(1 + \exp(\eta))$ | $1$ |
| Gaussian | $\mu$ | $y$ | $\frac{\eta^2}{2}$ | $\frac{1}{\sqrt{2\pi}}\exp\left(-\frac{y^2}{2}\right)$ |
| Poisson | $\log(\lambda)$ | $y$ | $e^\eta$ | $\frac{1}{y!}$ |
| Geometric | $\log(1-\phi)$ | $y$ | $\log\left(\frac{e^\eta}{1-e^\eta}\right)$ | $1$ |

❐ **Assumptions of GLMs** – Generalized Linear Models (GLM) aim at predicting a random variable $y$ as a function fo $x \in \mathbb{R}^{n+1}$ and rely on the following 3 assumptions:

$$(1) \quad \boxed{y|x;\theta \sim \text{ExpFamily}(\eta)} \qquad (2) \quad \boxed{h_\theta(x) = E[y|x;\theta]} \qquad (3) \quad \boxed{\eta = \theta^T x}$$

*Remark: ordinary least squares and logistic regression are special cases of generalized linear models.*
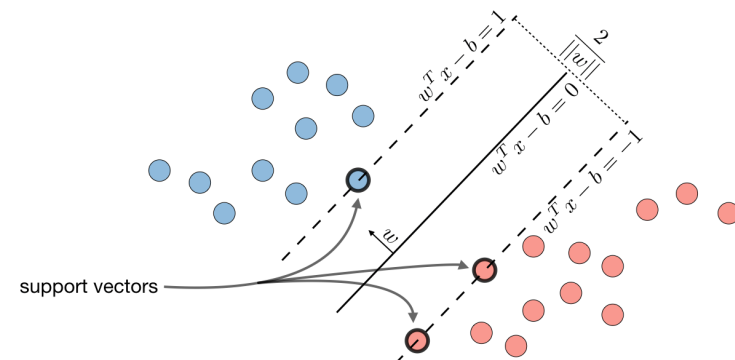
### 1.4 Support Vector Machines

The goal of support vector machines is to find the line that maximizes the minimum distance to the line.

❐ **Optimal margin classifier** – The optimal margin classifier $h$ is such that:

$$h(x) = \text{sign}(w^T x - b)$$

where $(w, b) \in \mathbb{R}^n \times \mathbb{R}$ is the solution of the following optimization problem:

$$\min \frac{1}{2}||w||^2 \qquad \text{such that} \qquad y^{(i)}(w^T x^{(i)} - b) \geqslant 1$$



support vectors

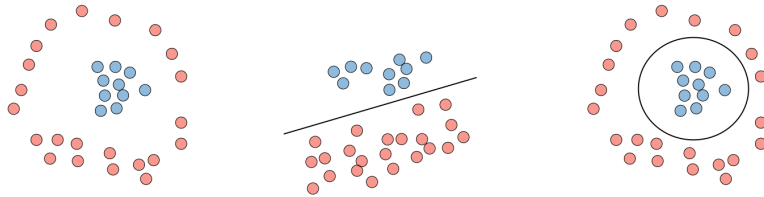*Remark: the line is defined as $\boxed{w^T x - b = 0}$.*

❐ **Hinge loss** – The hinge loss is used in the setting of SVMs and is defined as follows:

$$L(z,y) = [1 - yz]_+ = \max(0, 1 - yz)$$

❐ **Kernel** – Given a feature mapping $\phi$, we define the kernel $K$ to be defined as:

$$K(x,z) = \phi(x)^T \phi(z)$$

In practice, the kernel $K$ defined by $K(x,z) = \exp\left(-\frac{||x-z||^2}{2\sigma^2}\right)$ is called the Gaussian kernel and is commonly used.



Non-linear separability ⟹ Use of a kernel mapping $\phi$ ⟹ Decision boundary in the original space

*Remark: we say that we use the "kernel trick" to compute the cost function using the kernel because we actually don't need to know the explicit mapping $\phi$, which is often very complicated. Instead, only the values $K(x,z)$ are needed.*

❐ **Lagrangian** – We define the Lagrangian $\mathcal{L}(w,b)$ as follows:

$$\mathcal{L}(w,b) = f(w) + \sum_{i=1}^{l} \beta_i h_i(w)$$

*Remark: the coefficients $\beta_i$ are called the Lagrange multipliers.*

## 1.5   Generative Learning

A generative model first tries to learn how the data is generated by estimating $P(x|y)$, which we can then use to estimate $P(y|x)$ by using Bayes' rule.

### 1.5.1   Gaussian Discriminant Analysis

❐ **Setting** – The Gaussian Discriminant Analysis assumes that $y$ and $x|y = 0$ and $x|y = 1$ are such that:

$$y \sim \text{Bernoulli}(\phi)$$

$$x|y = 0 \sim \mathcal{N}(\mu_0, \Sigma) \quad \text{and} \quad x|y = 1 \sim \mathcal{N}(\mu_1, \Sigma)$$

❐ **Estimation** – The following table sums up the estimates that we find when maximizing the likelihood:

| $\widehat{\phi}$ | $\widehat{\mu}_j \quad (j = 0,1)$ | $\widehat{\Sigma}$ |
|---|---|---|
| $\dfrac{1}{m}\sum_{i=1}^{m} 1_{\{y^{(i)}=1\}}$ | $\dfrac{\sum_{i=1}^{m} 1_{\{y^{(i)}=j\}} x^{(i)}}{\sum_{i=1}^{m} 1_{\{y^{(i)}=j\}}}$ | $\dfrac{1}{m}\sum_{i=1}^{m} (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T$ |

### 1.5.2   Naive Bayes

❐ **Assumption** – The Naive Bayes model supposes that the features of each data point are all independent:

$$P(x|y) = P(x_1, x_2, ...|y) = P(x_1|y)P(x_2|y)... = \prod_{i=1}^{n} P(x_i|y)$$

❐ **Solutions** – Maximizing the log-likelihood gives the following solutions, with $k \in \{0,1\}$, $l \in [\![1,L]\!]$

$$P(y = k) = \frac{1}{m} \times \#\{j|y^{(j)} = k\} \quad \text{and} \quad P(x_i = l|y = k) = \frac{\#\{j|y^{(j)} = k \text{ and } x_i^{(j)} = l\}}{\#\{j|y^{(j)} = k\}}$$

*Remark: Naive Bayes is widely used for text classification and spam detection.*

## 1.6   Tree-based and ensemble methods

These methods can be used for both regression and classification problems.

❐ **CART** – Classification and Regression Trees (CART), commonly known as decision trees, can be represented as binary trees. They have the advantage to be very interpretable.

❐ **Random forest** – It is a tree-based technique that uses a high number of decision trees built out of randomly selected sets of features. Contrary to the simple decision tree, it is highly uninterpretable but its generally good performance makes it a popular algorithm.
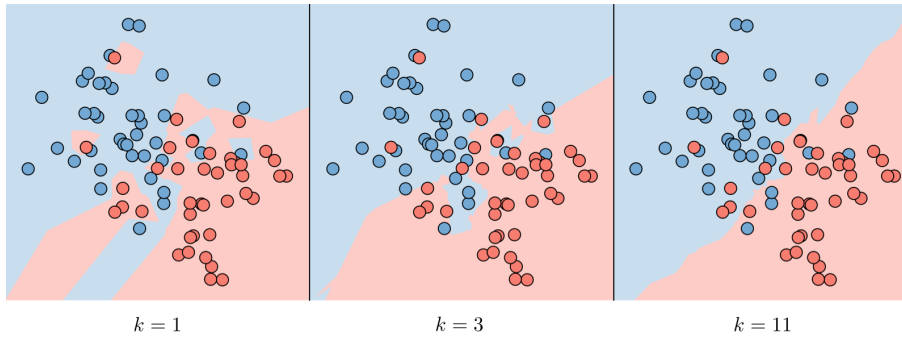
*Remark: random forests are a type of ensemble methods.*

❐ **Boosting** – The idea of boosting methods is to combine several weak learners to form a stronger one. The main ones are summed up in the table below:

| Adaptive boosting | Gradient boosting |
|---|---|
| - High weights are put on errors to improve at the next boosting step<br>- Known as Adaboost | - Weak learners trained on remaining errors |

## 1.7   Other non-parametric approaches

❐ **$k$-nearest neighbors** – The $k$-nearest neighbors algorithm, commonly known as $k$-NN, is a non-parametric approach where the response of a data point is determined by the nature of its $k$ neighbors from the training set. It can be used in both classification and regression settings.
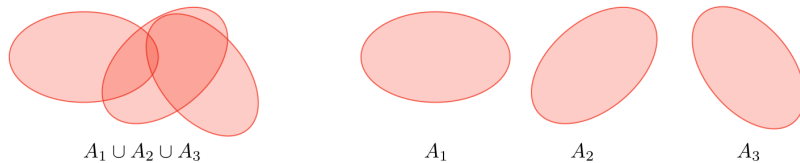
*Remark: The higher the parameter $k$, the higher the bias, and the lower the parameter $k$, the higher the variance.*

$$k = 1 \qquad\qquad k = 3 \qquad\qquad k = 11$$

## 1.8 Learning Theory

❒ **Union bound** – Let $A_1, ..., A_k$ be $k$ events. We have:

$$P(A_1 \cup ... \cup A_k) \leqslant P(A_1) + ... + P(A_k)$$



$$A_1 \cup A_2 \cup A_3 \qquad\qquad A_1 \qquad\quad A_2 \qquad\quad A_3$$

❒ **Hoeffding inequality** – Let $Z_1, .., Z_m$ be $m$ iid variables drawn from a Bernoulli distribution of parameter $\phi$. Let $\widehat{\phi}$ be their sample mean and $\gamma > 0$ fixed. We have:

$$P(|\phi - \widehat{\phi}| > \gamma) \leqslant 2\exp(-2\gamma^2 m)$$

*Remark: this inequality is also known as the Chernoff bound.*

❒ **Training error** – For a given classifier $h$, we define the training error $\widehat{\epsilon}(h)$, also known as the empirical risk or empirical error, to be as follows:

$$\widehat{\epsilon}(h) = \frac{1}{m} \sum_{i=1}^{m} 1_{\{h(x^{(i)}) \neq y^{(i)}\}}$$

❒ **Probably Approximately Correct (PAC)** – PAC is a framework under which numerous results on learning theory were proved, and has the following set of assumptions:

- the training and testing sets follow the same distribution

- the training examples are drawn independently

❒ **Shattering** – Given a set $S = \{x^{(1)}, ..., x^{(d)}\}$, and a set of classifiers $\mathcal{H}$, we say that $\mathcal{H}$ shatters $S$ if for any set of labels $\{y^{(1)}, ..., y^{(d)}\}$, we have:
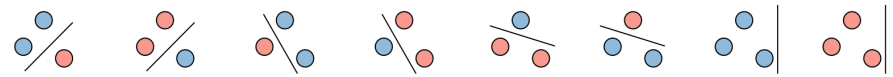
$$\exists h \in \mathcal{H}, \quad \forall i \in [\![1, d]\!], \quad h(x^{(i)}) = y^{(i)}$$

❒ **Upper bound theorem** – Let $\mathcal{H}$ be a finite hypothesis class such that $|\mathcal{H}| = k$ and let $\delta$ and the sample size $m$ be fixed. Then, with probability of at least $1 - \delta$, we have:

$$\epsilon(\widehat{h}) \leqslant \left( \min_{h \in \mathcal{H}} \epsilon(h) \right) + 2\sqrt{\frac{1}{2m} \log\left(\frac{2k}{\delta}\right)}$$

❒ **VC dimension** – The Vapnik-Chervonenkis (VC) dimension of a given infinite hypothesis class $\mathcal{H}$, noted VC($\mathcal{H}$) is the size of the largest set that is shattered by $\mathcal{H}$.

*Remark: the VC dimension of $\mathcal{H} = \{set\ of\ linear\ classifiers\ in\ 2\ dimensions\}$ is 3.*



❒ **Theorem (Vapnik)** – Let $\mathcal{H}$ be given, with VC($\mathcal{H}$) = $d$ and $m$ the number of training examples. With probability at least $1 - \delta$, we have:

$$\epsilon(\widehat{h}) \leqslant \left( \min_{h \in \mathcal{H}} \epsilon(h) \right) + O\left( \sqrt{\frac{d}{m} \log\left(\frac{m}{d}\right) + \frac{1}{m} \log\left(\frac{1}{\delta}\right)} \right)$$

## 2    Unsupervised Learning

### 2.1    Introduction to Unsupervised Learning

❏ **Motivation** – The goal of unsupervised learning is to find hidden patterns in unlabeled data $\{x^{(1)},...,x^{(m)}\}$.

❏ **Jensen's inequality** – Let $f$ be a convex function and $X$ a random variable. We have the following inequality:

$$E[f(X)] \geqslant f(E[X])$$

### 2.2    Clustering

#### 2.2.1    Expectation-Maximization

❏ **Latent variables** – Latent variables are hidden/unobserved variables that make estimation problems difficult, and are often denoted $z$. Here are the most common settings where there are latent variables:

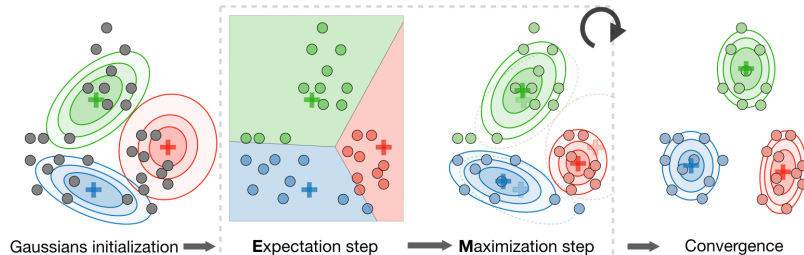| Setting | Latent variable $z$ | $x\|z$ | Comments |
|---|---|---|---|
| Mixture of $k$ Gaussians | Multinomial($\phi$) | $\mathcal{N}(\mu_j,\Sigma_j)$ | $\mu_j \in \mathbb{R}^n, \phi \in \mathbb{R}^k$ |
| Factor analysis | $\mathcal{N}(0,I)$ | $\mathcal{N}(\mu + \Lambda z, \psi)$ | $\mu_j \in \mathbb{R}^n$ |

❏ **Algorithm** – The Expectation-Maximization (EM) algorithm gives an efficient method at estimating the parameter $\theta$ through maximum likelihood estimation by repeatedly constructing a lower-bound on the likelihood (E-step) and optimizing that lower bound (M-step) as follows:

- E-step: Evaluate the posterior probability $Q_i(z^{(i)})$ that each data point $x^{(i)}$ came from a particular cluster $z^{(i)}$ as follows:

$$Q_i(z^{(i)}) = P(z^{(i)}|x^{(i)};\theta)$$

- M-step: Use the posterior probabilities $Q_i(z^{(i)})$ as cluster specific weights on data points $x^{(i)}$ to separately re-estimate each cluster model as follows:

$$\theta_i = \underset{\theta}{\arg\max} \sum_i \int_{z^{(i)}} Q_i(z^{(i)}) \log\left(\frac{P(x^{(i)},z^{(i)};\theta)}{Q_i(z^{(i)})}\right) dz^{(i)}$$



Gaussians initialization ➡ **E**xpectation step ➡ **M**aximization step ➡ Convergence
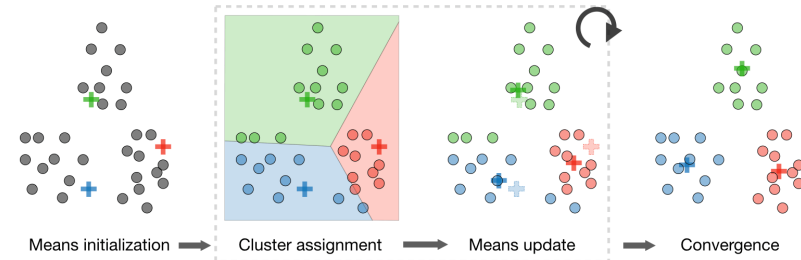
#### 2.2.2    $k$-means clustering

We note $c^{(i)}$ the cluster of data point $i$ and $\mu_j$ the center of cluster $j$.

❏ **Algorithm** – After randomly initializing the cluster centroids $\mu_1,\mu_2,...,\mu_k \in \mathbb{R}^n$, the $k$-means algorithm repeats the following step until convergence:

$$c^{(i)} = \underset{j}{\arg\min}||x^{(i)} - \mu_j||^2 \quad \text{and} \quad \mu_j = \frac{\sum_{i=1}^{m} 1_{\{c^{(i)}=j\}} x^{(i)}}{\sum_{i=1}^{m} 1_{\{c^{(i)}=j\}}}$$



Means initialization ➡ Cluster assignment ➡ Means update ➡ Convergence

❏ **Distortion function** – In order to see if the algorithm converges, we look at the distortion function defined as follows:

$$J(c,\mu) = \sum_{i=1}^{m} ||x^{(i)} - \mu_{c^{(i)}}||^2$$

#### 2.2.3    Hierarchical clustering

❏ **Algorithm** – It is a clustering algorithm with an agglomerative hierarchical approach that build nested clusters in a successive manner.

❏ **Types** – There are different sorts of hierarchical clustering algorithms that aims at optimizing different objective functions, which is summed up in the table below:

| Ward linkage | Average linkage | Complete linkage |
|---|---|---|
| Minimize within cluster distance | Minimize average distance between cluster pairs | Minimize maximum distance of between cluster pairs |

#### 2.2.4    Clustering assessment metrics

In an unsupervised learning setting, it is often hard to assess the performance of a model since we don't have the ground truth labels as was the case in the supervised learning setting.

❏ **Silhouette coefficient** – By noting $a$ and $b$ the mean distance between a sample and all other points in the same class, and between a sample and all other points in the next nearest cluster, the silhouette coefficient $s$ for a single sample is defined as follows:

$$s = \frac{b - a}{\max(a,b)}$$

❒ **Calinski-Harabaz index** – By noting $k$ the number of clusters, $B_k$ and $W_k$ the between and within-clustering dispersion matrices respectively defined as

$$B_k = \sum_{j=1}^{k} n_{c^{(i)}}(\mu_{c^{(i)}} - \mu)(\mu_{c^{(i)}} - \mu)^T, \qquad W_k = \sum_{i=1}^{m}(x^{(i)} - \mu_{c^{(i)}})(x^{(i)} - \mu_{c^{(i)}})^T$$

the Calinski-Harabaz index $s(k)$ indicates how well a clustering model defines its clusters, such that the higher the score, the more dense and well separated the clusters are. It is defined as follows:

$$s(k) = \frac{\text{Tr}(B_k)}{\text{Tr}(W_k)} \times \frac{N - k}{k - 1}$$

## 2.3 Dimension reduction

### 2.3.1 Principal component analysis

It is a dimension reduction technique that finds the variance maximizing directions onto which to project the data.

❒ **Eigenvalue, eigenvector** – Given a matrix $A \in \mathbb{R}^{n \times n}$, $\lambda$ is said to be an eigenvalue of $A$ if there exists a vector $z \in \mathbb{R}^n \backslash \{0\}$, called eigenvector, such that we have:

$$Az = \lambda z$$

❒ **Spectral theorem** – Let $A \in \mathbb{R}^{n \times n}$. If $A$ is symmetric, then $A$ is diagonalizable by a real orthogonal matrix $U \in \mathbb{R}^{n \times n}$. By noting $\Lambda = \text{diag}(\lambda_1,...,\lambda_n)$, we have:

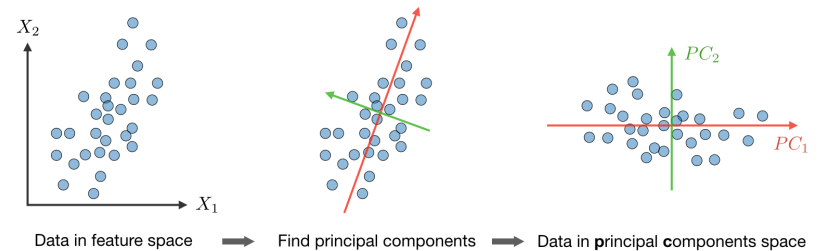$$\exists \Lambda \text{ diagonal}, \quad A = U\Lambda U^T$$

*Remark: the eigenvector associated with the largest eigenvalue is called principal eigenvector of matrix $A$.*

❒ **Algorithm** – The Principal Component Analysis (PCA) procedure is a dimension reduction technique that projects the data on $k$ dimensions by maximizing the variance of the data as follows:

- Step 1: Normalize the data to have a mean of 0 and standard deviation of 1.

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{\sigma_j} \quad \text{where} \quad \mu_j = \frac{1}{m}\sum_{i=1}^{m} x_j^{(i)} \quad \text{and} \quad \sigma_j^2 = \frac{1}{m}\sum_{i=1}^{m}(x_j^{(i)} - \mu_j)^2$$

- Step 2: Compute $\Sigma = \frac{1}{m}\sum_{i=1}^{m} x^{(i)}x^{(i)T} \in \mathbb{R}^{n \times n}$, which is symmetric with real eigenvalues.

- Step 3: Compute $u_1, ..., u_k \in \mathbb{R}^n$ the $k$ orthogonal principal eigenvectors of $\Sigma$, i.e. the orthogonal eigenvectors of the $k$ largest eigenvalues.

- Step 4: Project the data on $\text{span}_{\mathbb{R}}(u_1,...,u_k)$. This procedure maximizes the variance among all $k$-dimensional spaces.



Data in feature space ➡ Find principal components ➡ Data in **p**rincipal **c**omponents space

### 2.3.2 Independent component analysis

It is a technique meant to find the underlying generating sources.

❒ **Assumptions** – We assume that our data $x$ has been generated by the $n$-dimensional source vector $s = (s_1,...,s_n)$, where $s_i$ are independent random variables, via a mixing and non-singular matrix $A$ as follows:

$$x = As$$

The goal is to find the unmixing matrix $W = A^{-1}$ by an update rule.

❒ **Bell and Sejnowski ICA algorithm** – This algorithm finds the unmixing matrix $W$ by following the steps below:

- Write the probability of $x = As = W^{-1}s$ as:

$$p(x) = \prod_{i=1}^{n} p_s(w_i^T x) \cdot |W|$$

- Write the log likelihood given our training data $\{x^{(i)}, i \in [\![1,m]\!]\}$ and by noting $g$ the sigmoid function as:

$$l(W) = \sum_{i=1}^{m}\left(\sum_{j=1}^{n} \log\left(g'(w_j^T x^{(i)})\right) + \log|W|\right)$$

Therefore, the stochastic gradient ascent learning rule is such that for each training example $x^{(i)}$, we update $W$ as follows:
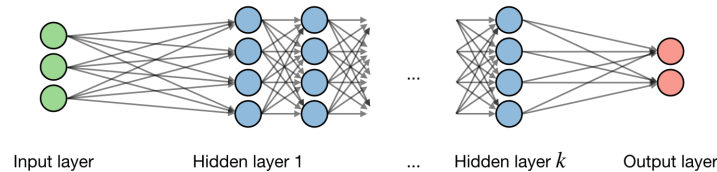
$$W \longleftarrow W + \alpha\left(\begin{pmatrix} 1 - 2g(w_1^T x^{(i)}) \\ 1 - 2g(w_2^T x^{(i)}) \\ \vdots \\ 1 - 2g(w_n^T x^{(i)}) \end{pmatrix} x^{(i)T} + (W^T)^{-1}\right)$$

## 3    Deep Learning

### 3.1    Neural Networks

Neural networks are a class of models that are built with layers. Commonly used types of neural networks include convolutional and recurrent neural networks.

❒ **Architecture** – The vocabulary around neural networks architectures is described in the figure below:



| Input layer | Hidden layer 1 | ... | Hidden layer $k$ | Output layer |

By noting $i$ the $i^{th}$ layer of the network and $j$ the $j^{th}$ hidden unit of the layer, we have:

$$z_j^{[i]} = {w_j^{[i]}}^T x + b_j^{[i]}$$

where we note $w$, $b$, $z$ the weight, bias and output respectively.

❒ **Activation function** – Activation functions are used at the end of a hidden unit to introduce non-linear complexities to the model. Here are the most common ones:

| Sigmoid | Tanh | ReLU | Leaky ReLU |
|---------|------|------|------------|
| $g(z) = \dfrac{1}{1 + e^{-z}}$ | $g(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$ | $g(z) = \max(0,z)$ | $g(z) = \max(\epsilon z, z)$ <br> with $\epsilon \ll 1$ |
|  |  |  |  |

❒ **Cross-entropy loss** – In the context of neural networks, the cross-entropy loss $L(z,y)$ is commonly used and is defined as follows:

$$L(z,y) = -\Big[ y \log(z) + (1 - y) \log(1 - z) \Big]$$

❒ **Learning rate** – The learning rate, often noted $\eta$, indicates at which pace the weights get updated. This can be fixed or adaptively changed. The current most popular method is called Adam, which is a method that adapts the learning rate.

❒ **Backpropagation** – Backpropagation is a method to update the weights in the neural network by taking into account the actual output and the desired output. The derivative with respect to weight $w$ is computed using chain rule and is of the following form:

$$\frac{\partial L(z,y)}{\partial w} = \frac{\partial L(z,y)}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial w}$$

As a result, the weight is updated as follows:

$$w \longleftarrow w - \eta \frac{\partial L(z,y)}{\partial w}$$

❒ **Updating weights** – In a neural network, weights are updated as follows:

- Step 1: Take a batch of training data.

- Step 2: Perform forward propagation to obtain the corresponding loss.

- Step 3: Backpropagate the loss to get the gradients.

- Step 4: Use the gradients to update the weights of the network.

❒ **Dropout** – Dropout is a technique meant at preventing overfitting the training data by dropping out units in a neural network. In practice, neurons are either dropped with probability $p$ or kept with probability $1 - p$.

### 3.2    Convolutional Neural Networks

❒ **Convolutional layer requirement** – By noting $W$ the input volume size, $F$ the size of the convolutional layer neurons, $P$ the amount of zero padding, then the number of neurons $N$ that fit in a given volume is such that:

$$N = \frac{W - F + 2P}{S} + 1$$

❒ **Batch normalization** – It is a step of hyperparameter $\gamma, \beta$ that normalizes the batch $\{x_i\}$. By noting $\mu_B, \sigma_B^2$ the mean and variance of that we want to correct to the batch, it is done as follows:

$$x_i \longleftarrow \gamma \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

It is usually done after a fully connected/convolutional layer and before a non-linearity layer and aims at allowing higher learning rates and reducing the strong dependence on initialization.

### 3.3    Recurrent Neural Networks

❒ **Types of gates** – Here are the different types of gates that we encounter in a typical recurrent neural network:

| Input gate | Forget gate | Output gate | Gate |
|------------|-------------|-------------|------|
| Write to cell or not? | Erase a cell or not? | Reveal a cell or not? | How much writing? |

❒ **LSTM** – A long short-term memory (LSTM) network is a type of RNN model that avoids the vanishing gradient problem by adding 'forget' gates.