

DISCLAIMER

These notes are supposed to explain some major topics covered in the lessons.

While they may not provide an exhaustive overview, they strive to complement the information presented in the slides with some explanations taken from additional readings in various textbooks.

If someone finds some errors or have any suggestion, it's possible to contact me on mmeschini001@gmail.com

MACHINE LEARNING

Why machine learning? Many interesting problems are too complex to admit an algorithmic solution or even a complete description. For these problems only data are available. Machine learning is about using data to solve problems.

Machine learning jargon

MODEL: a mathematical or computational representation of the real world

TRAINING DATA: the portion of the dataset used to train the machine learning model

TEST DATA: the portion of the dataset that is used to evaluate the performance of a trained machine learning model.

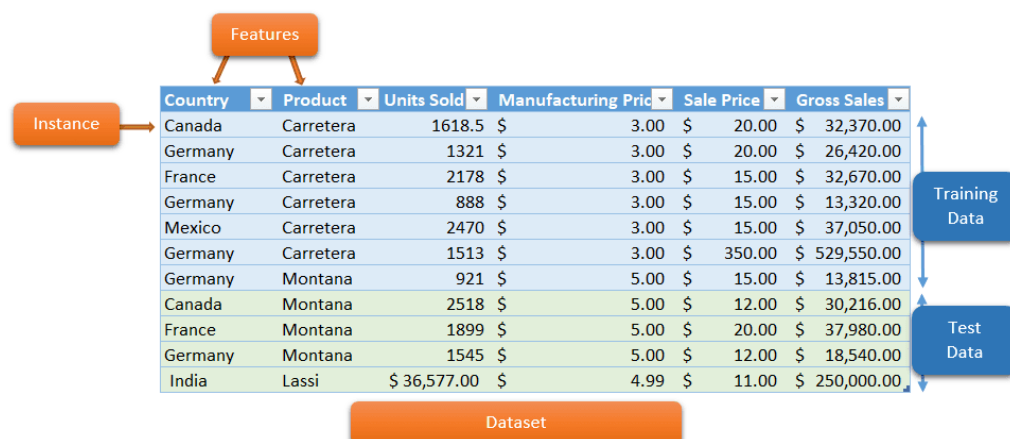
FEATURE/ATTRIBUTE/VARIABLES: an individual measurable property used to train the machine learning model

OBSERVATIONS/SAMPLES/INSTANCES: are individual data points that represent specific records. Typically they correspond to a single entry.

CLASS/LABELS: The output or target variable that a machine learning model is trying to predict.

OVERFITTING: When a machine learning model learns the training data too well, capturing noise and irrelevant details, which can lead to poor generalization to new data. (typically, too complex model)

UNDERFITTING: The opposite of overfitting, where a model is too simple to capture the underlying patterns in the data, resulting in poor performance.



There are three different types of machines learning:

- **SUPERVISED LEARNING:** the algorithm “learns” from the training dataset by iteratively making predictions on the data and adjusting for the correct answer. While supervised learning models tend to be more accurate than unsupervised learning models, they require upfront human intervention to label the data appropriately.
- **UNSUPERVISED LEARNING/CLUSTERING:** Unsupervised learning models, in contrast, work on their own to discover the inherent structure of unlabelled data. Note that they still require some human intervention for validating output variables.

- REINFORCEMENT LEARNING algorithms work by training the agent through interaction with the environment. They learn to optimize their policy over time, improving their decision-making capabilities.

The types of quantity our model want to learn can be of two different types:

- Real values
- Categorical/nominal values: e.g. colour, name. (there isn't a natural ordering)

Depending on the model an the output variables we can distinguish:

		type of output	
		quantitative	nominal
supervised	YES	REGRESSION	CLASSIFICATION
	NO	LOW-DIMENSIONAL MAPPING	CLUSTERING

Possible scenarios

There are three possible scenarios:

- 1) Useful mostly for reasoning: we have a **learning machine** that is **fixed** and also the **data** are **fixed** (we have only the training set). The objective is to find the correct learners.
- 2) Not realizable: we **know the model that generates data** (the probabilities are known), so the objective is to find the optimal learner. Useful for theory.
- 3) Our usual situation: **Data** are **stochastic**, the **learner** is **fixed** and chosen in advance. The objective is to find the best configuration of learning machine which is correct for any realization of the data.

Our model will typically use **inductive learning**. This type of learning is based on examples and data, our objective is to learn a function **for every possible input**, given only a **finite set of examples input** -> **GENERALIZATION**

Induction works only if we have an a-priori knowledge or assumption. EX:

Problem	Bias
Classification	Points in the same class are close to each other
Regression	Nearby points don't have very different outputs
Clustering	There are k compact clusters in the data
Mapping	The data are organized in a low-dimensional shape within the d -dimensional data space

Probability

The probability quantifies the chance of an event or outcome occurring. It provides a way to express uncertainty to which something is likely to happen. Something that may or may not happen is called an **outcome** (ω). All possible outcomes constitute the **sample space** (Ω). An **event** refers to a specific outcome or a subsets of Ω .

We defined as $P(A)$ the probability of event A .

Axioms of probability:

- ① $P(A) \geq 0$
- ② $\sum_{i=1}^N P(\omega_i) = 1$, or $P(\Omega) = 1$
- ③ If A_1 and A_2 are mutually exclusive events (viewed as sets: if they are **disjoint** = have zero intersection), then $P(A_1 \text{ or } A_2) = P(A_1 \cup A_2) = P(A_1) + P(A_2)$

A **random variable** (X) is a numerical variable that doesn't have a fixed value but changes according to a given probability law. The **Cumulative distribution function (CDF)** of a random variable X is the function given by:

$$F_X(x) = P(X \leq x) \quad (\text{Eq.1})$$

Where the right-hand side represents the probability that the random variable X takes on a value less than or equal to x .

The **probability mass function** is a function that gives the probability that a **discrete** random variable is exactly equal to some value.

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

The **probability density function** is a function whose value at any given sample in the sample space can be interpreted as providing a relative likelihood that the value of the random variable would be equal to that sample.

A function f_X such that

$$P_X(\hat{x}) = \int_{-\infty}^{\hat{x}} f_X(x) dx$$

is the **probability density function** of X .

We define the **conditional probability** $P(E | F)$ as the probability of an event E given the knowledge that another event F has occurred.

Two events are **independent** if the outcome of one event does not influence the outcome of the second event.

Bayesian decision theory

Bayesian decision theory is a fundamental statistical approach to the problem of pattern classification. This approach is based on quantifying the trade-offs between various classification decisions using probability and the costs that accompany such decisions. It makes the assumption that the decision problem is posed in probabilistic terms, and that **all the relevant probability values are known**.

We define as **state of nature** (ω) (t in professor slide) the different possible states or scenarios that a system could be in. (e.g. salmon and sea bass are two different state of nature). We define a **priori probability** refers to the probability of an event occurring based on prior knowledge. It is represented by the likelihood in the occurrence of an event before any new evidence is taken into account. Then, is important to define the **class-conditional probability density function** expressed as $P(\mathbf{x}|\omega)$ that represents the distribution of the random variable \mathbf{x} depending on the state of nature.

The probability density of finding a pattern that is in category ω_j and has a feature value \mathbf{x} can be written in two ways: $P(\omega_j, \mathbf{x}) = P(\omega_j|\mathbf{x}) * P(\mathbf{x}) = P(\mathbf{x}|\omega_j) * P(\omega_j)$ rearranging these leads us to the **Bayes' formula**:

$$P(\omega_j|\mathbf{x}) = \frac{p(\mathbf{x}|\omega_j)P(\omega_j)}{p(\mathbf{x})}, \quad \text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}.$$

Bayes' formula shows that by observing the value of \mathbf{x} we can convert the prior probability $P(\omega_j)$ to the **a posteriori probability** $P(\omega_j|\mathbf{x})$, that is the probability of the state of nature being ω_j given that the feature value \mathbf{x} has been measured.

Oss: the evidence can be viewed as a scale factor that guarantees that the posterior probability sum to one.

In this context, is useful to introduce the **decision rule** that represents a criterion (a set of boundaries) used to minimize the cost of our decision. To evaluate each decision, we use the **expected loss**, also called **risk**. The **loss function** (λ) states exactly how costly each action is and is used to convert a probability determination into a decision. For a generic decision y_i the risk is:

$$R(y_i) = \sum_{j=1}^c \lambda(y_i, t_j)P(t_j)$$

Finally, the **conditional risk** represent the expected loss associated a decision under a specific state of nature.

$$R(y_i|\mathbf{x}) = \sum_{j=1}^c \lambda(y_i, t_j)P(t_j|\mathbf{x})$$

$R(y_i|\mathbf{x})$ is the **conditional risk** of decision y_i
when we have the **experimental observation** \mathbf{x}

Classification is a decision problem where there is no decision to take and we want only to recognize the state of nature. In short, the model tries to predict the correct class of a given input data. A classifier is a rule $y()$ that **receives an observation \mathbf{x} and output a class $y(\mathbf{x})$** .

Naive Bayes classifier

A Naive Bayes classifier is a simple and probabilistic machine learning algorithm that is often used for classification tasks. It is based on Bayes' theorem and makes a "naïve" assumption that features used to describe data are conditionally independent.

Naïve Bayes is a conditional probability model: it assigns probabilities $P(C_k | x_1, \dots, x_n)$ for each of the K possible classes C_k (= state of natures) given a problem instance to be classified represented by a vector $\mathbf{x} = (x_1, \dots, x_n)$ encoding some n features. From the Bayes' theorem

$$p(C_k | \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} | C_k)}{p(\mathbf{x})}$$

In practice, there is interest only in the numerator of that fraction, because the denominator does not depend on C and the values of the features x_i are given, so that the denominator is effectively constant. The numerator is equivalent to the [joint probability model](#)

$$p(C_k, x_1, \dots, x_n)$$

which can be rewritten as follows, using the [chain rule](#) for repeated applications of the definition of [conditional probability](#):

$$\begin{aligned} p(C_k, x_1, \dots, x_n) &= p(x_1, \dots, x_n, C_k) \\ &= p(x_1 | x_2, \dots, x_n, C_k) p(x_2, \dots, x_n, C_k) \\ &= p(x_1 | x_2, \dots, x_n, C_k) p(x_2 | x_3, \dots, x_n, C_k) p(x_3, \dots, x_n, C_k) \\ &= \dots \\ &= p(x_1 | x_2, \dots, x_n, C_k) p(x_2 | x_3, \dots, x_n, C_k) \dots p(x_{n-1} | x_n, C_k) p(x_n | C_k) p(C_k) \end{aligned}$$

Now the "naive" [conditional independence](#) assumptions come into play: assume that all features in \mathbf{x} are [mutually independent](#), conditional on the category C_k . Under this assumption,

$$p(x_i | x_{i+1}, \dots, x_n, C_k) = p(x_i | C_k).$$

Thus, the joint model can be expressed as

$$\begin{aligned} p(C_k | x_1, \dots, x_n) &\propto p(C_k, x_1, \dots, x_n) \\ &= p(C_k) p(x_1 | C_k) p(x_2 | C_k) p(x_3 | C_k) \dots \\ &= p(C_k) \prod_{i=1}^n p(x_i | C_k), \end{aligned}$$

The decision rule for minimizing the probability of error is:

$$\text{Decide } \omega_1 \text{ if } P(\omega_1 | x) > P(\omega_2 | x); \text{ otherwise decide } \omega_2$$

The loss function for this classifier is **zero-one**: $\lambda(y | t) = \begin{cases} 0 & \text{if } y = t \\ 1 & \text{if } y \neq t \end{cases}$

So, all types of errors have the same cost (=1) and correct classifications don't have a cost.

Linear Regression

When only data are available is possible to find the correct output using a regression model.

Regression means approximating a functional dependency based on measured data.

It is a supervised problem given that we have observations and target

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{pmatrix} \quad \begin{pmatrix} t_1 \\ t_2 \\ t_3 \\ \vdots \\ t_N \end{pmatrix}$$

Observations Target

We are looking for a linear model $y(x)$ that predicts t given x .

Since the model is linear the form will be: $y(x) = w_0 + wx$

From the training set is possible to estimate the parameters w . Difficulty is possible to find values that are good for all points, so the solution is to find two values that minimize the average error.

The most popular estimation method is the **square error loss**, in which we pick the coefficients w to minimize the following loss function:

$$\lambda_{SE}(y, t) = (y - t)^2 \quad \xrightarrow{\text{So our objective function will be}} \quad J_{MSE} = \frac{1}{N} \sum_{l=1}^N (y_l - t_l)^2$$

- This function when building a model (=training) is a function of the **parameters** of the model and the data are fixed.
- When using a model (=inference) is function of the **data**, while the parameters are now fixed.

The least squares solution of the one-dimensional linear regression problem (for the model $y(x) = wx$) is:

$$w = \frac{\sum_{l=1}^N x_l t_l}{\sum_{l=1}^N x_l^2}$$

While for the model $y(x) = w_0 + w_1 x$ where is also present the offset is:

$$\bar{x} = \frac{1}{N} \sum_{l=1}^N x_l \quad \bar{t} = \frac{1}{N} \sum_{l=1}^N t_l$$

Note:

w_1 is called: slope, gain

w_0 is called intercept, offset, bias

$$w_1 = \frac{\sum_{l=1}^N (x_l - \bar{x})(t_l - \bar{t})}{\sum_{l=1}^N (x_l - \bar{x})^2}$$

$$w_0 = \bar{t} - w_1 \bar{x}$$

Multidimensional linear regression problem

The data is now composed of d-dimensional vectors $\mathbf{x}_1 = [x_{1,1}, x_{1,2}, \dots, x_{1,d}]$ and we can organize them into a N x d matrix. Now we have also d parameters:

$$X = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \vdots \\ \mathbf{x}_N \end{pmatrix} = \begin{pmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,d} \\ x_{2,1} & x_{2,2} & \dots & x_{2,d} \\ x_{3,1} & x_{3,2} & \dots & x_{3,d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N,1} & x_{N,2} & \dots & x_{N,d} \end{pmatrix} \quad \mathbf{w} = \begin{pmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_d \end{pmatrix}$$

The linear model takes the d inputs of each observation and combines them by using the d parameters to produce one output:

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,d} \\ x_{2,1} & x_{2,2} & \dots & x_{2,d} \\ x_{3,1} & x_{3,2} & \dots & x_{3,d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N,1} & x_{N,2} & \dots & x_{N,d} \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_d \end{pmatrix} = \begin{pmatrix} x_{1,1}w_1 + x_{1,2}w_2 + \dots + x_{1,d}w_d \\ x_{2,1}w_1 + x_{2,2}w_2 + \dots + x_{2,d}w_d \\ x_{3,1}w_1 + x_{3,2}w_2 + \dots + x_{3,d}w_d \\ \vdots \\ x_{N,1}w_1 + x_{N,2}w_2 + \dots + x_{N,d}w_d \end{pmatrix} = X\mathbf{w}$$

The affine case let us incorporate the additive parameter w_0 by adding one constant column to the data matrix

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} 1 & x_{1,1} & x_{1,2} & \dots & x_{1,d} \\ 1 & x_{2,1} & x_{2,2} & \dots & x_{2,d} \\ 1 & x_{3,1} & x_{3,2} & \dots & x_{3,d} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N,1} & x_{N,2} & \dots & x_{N,d} \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_d \end{pmatrix} = \begin{pmatrix} w_0 + x_{1,1}w_1 + x_{1,2}w_2 + \dots + x_{1,d}w_d \\ w_0 + x_{2,1}w_1 + x_{2,2}w_2 + \dots + x_{2,d}w_d \\ w_0 + x_{3,1}w_1 + x_{3,2}w_2 + \dots + x_{3,d}w_d \\ \vdots \\ w_0 + x_{N,1}w_1 + x_{N,2}w_2 + \dots + x_{N,d}w_d \end{pmatrix} = X\mathbf{w}$$

Our final goal is to make this model's prediction \mathbf{y} as similar as possible to the measured outputs for each observation

$$\mathbf{t} = \begin{pmatrix} t_1 \\ t_2 \\ t_3 \\ \vdots \\ t_N \end{pmatrix}$$

The solution in this case is: $\mathbf{w} = (X^T X)^{-1} X^T \mathbf{t}$

Optimisation

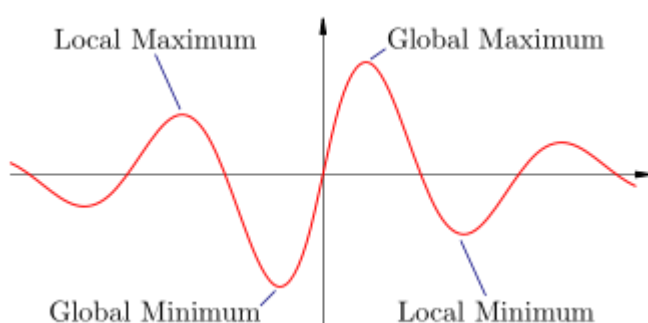
The target of optimisation is finding the extrema of an objective function $J(\mathbf{w})$, where $J: S \subset \mathbb{R}^m \rightarrow \mathbb{R}$. An extremum is a point $\mathbf{w}^* \in S$ that may be a maximum or minimum.

A point \mathbf{w}^* is a **minimum** if there is a neighbourhood $R \subseteq S$, where the following holds:

$$J(\mathbf{w}) \geq J(\mathbf{w}^*) \quad \forall \mathbf{w} \in R$$

A point \mathbf{w}^* is a **maximum** if \leq is used.

An extremum is **relative/local** if R is strictly contained in S , or in other words if it is a minimum only in a local neighbourhood. An extremum is **absolute/global** if $R = S$. Finally, an extremum is **isolated** if we can draw a sphere around \mathbf{w}^* and there are no other points there with function values greater (for a maximum) or smaller (for a minimum) than the value at that point.



Typical objective function in machine learning

- **Cost function:** which are inversely related to how much we like a solution (higher cost, less quality of the solution).
- **Risk (expected loss):** $J(\mathbf{w}) = \int_{\mathcal{X}} \lambda(t, \mathbf{y}(\mathbf{w})) p(\mathbf{x}) d\mathbf{x}$
- **Likelihood:** which is a common function that we have to maximize.

We define as **optimal solution** an extremum \mathbf{w}^* of J , and the optimal value will be $J(\mathbf{w}^*)$. Then we can have a **feasible solution** which is any point \mathbf{w} which satisfies all hypotheses of the optimisation problem. Then another important aspect is the convexity of a function.

Convex sets

A set $S \subset \mathbb{R}^m$ is convex if and only if, for any $\theta \in [0, 1]$,

$$\forall \mathbf{v}, \mathbf{w} \in S \Rightarrow \theta \mathbf{v} + (1 - \theta) \mathbf{w} \in S$$

more generally if for any $\theta_1 > 0, \dots, \theta_n > 0$ such that $\sum_k \theta_k = 1$

$$\forall \mathbf{v}_1, \dots, \mathbf{v}_n \in S \Rightarrow \sum_k \theta_k \mathbf{v}_k \in S$$

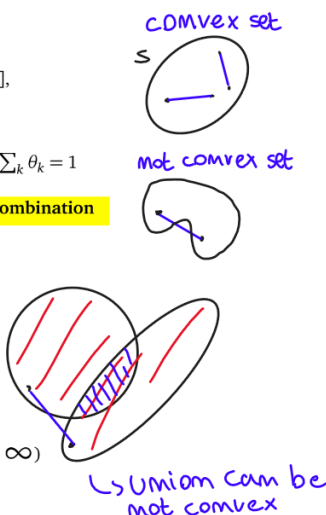
Convex combination

Properties:

- $\mathbf{v} \in \mathbb{R}^m$ (a single point) is convex
- $\emptyset = \{ \}$ (the empty set) is convex
- \mathbb{R}^m is convex

and if S_1 and S_2 are convex, then

- $S_1 \cap S_2$ is convex
- $S_1 \cup S_2$ IS NOT NECESSARILY convex ($0 + 0 = \infty$)



Convex functions

A function $J: S \subseteq \mathbb{R}^m \rightarrow \mathbb{R}$ is convex if S is a convex set and if $\forall \mathbf{v}, \mathbf{w} \in S$, and with $0 \leq \theta \leq 1$:

$$J(\theta \mathbf{v} + (1 - \theta) \mathbf{w}) \leq \theta J(\mathbf{v}) + (1 - \theta) J(\mathbf{w})$$

i.e., if its **epigraph** is a convex set

more generally if for any $\theta_1 > 0, \dots, \theta_n > 0$ such that $\sum_k \theta_k = 1$

$$\forall \mathbf{v}_1, \dots, \mathbf{v}_n \in S \Rightarrow J\left(\sum_k \theta_k \mathbf{v}_k\right) \leq \sum_k \theta_k J(\mathbf{v}_k)$$

J is concave if $-J$ is convex

The convexity is a good thing because it ensures the **uniqueness of extrema** and the **convergence of iterative algorithms**. While for non-convex problems we cannot be sure whether an extremum is absolute or relative.

Types of optimisation problems

The optimisation problems can be of different types:

- **Discrete:** if $S \in \mathbb{N}^m$, or more generally there are countable variables
- **Continuous:** $S \in \mathbb{R}^m$
- **Constrained:** if we want to optimize an objective function with respect to some variables in the presence of constraints on those variables.

Constraints are expressed by **constraint functions**:

- **Equality constraints:** $f_1(\mathbf{w}) = 0, f_2(\mathbf{w}) = 0, \dots, f_k(\mathbf{w}) = 0$
- **Inequality constraints:** $f_1(\mathbf{w}) \geq 0, f_2(\mathbf{w}) \geq 0, \dots, f_k(\mathbf{w}) \geq 0$

- **Unconstrained**
- **Smooth:** if the objective function is differentiable
- **Nonsmooth:** if the objective function is not differentiable
- **Linear** (typically needs a constraint, otherwise the objective can increase/decrease indefinitely)
- **Convex/Nonconvex**

Useful for this type of problems is the gradient, that indicates the direction of maximum increase.

$$\nabla J(\mathbf{w}) = \begin{bmatrix} \frac{\partial J(\mathbf{w})}{\partial w_1} \\ \frac{\partial J(\mathbf{w})}{\partial w_2} \\ \vdots \\ \frac{\partial J(\mathbf{w})}{\partial w_m} \end{bmatrix}$$

- Derivative \rightarrow rate of growth of a function of a scalar variable
- Negative sign \rightarrow decreasing
- Gradient length (norm) \rightarrow rate of maximum growth
- Direction \rightarrow direction of maximum growth

Then, the Hessian matrix $H_J(\mathbf{w}) : \mathbb{R}^m \rightarrow \mathbb{R}^m \times \mathbb{R}^m$ s.t. $H_J(\mathbf{w})_{ij} = \frac{\partial^2 J(\mathbf{w})}{\partial w_i \partial w_j}$ (is the gradient of the

$$H = \begin{pmatrix} \frac{\partial^2 J}{\partial w_1^2} & \frac{\partial^2 J}{\partial w_1 \partial w_2} & \frac{\partial^2 J}{\partial w_1 \partial w_3} & \cdots & \frac{\partial^2 J}{\partial w_1 \partial w_m} \\ \frac{\partial^2 J}{\partial w_2 \partial w_1} & \frac{\partial^2 J}{\partial w_2^2} & \frac{\partial^2 J}{\partial w_2 \partial w_3} & \cdots & \frac{\partial^2 J}{\partial w_2 \partial w_m} \\ \frac{\partial^2 J}{\partial w_3 \partial w_1} & \frac{\partial^2 J}{\partial w_3 \partial w_2} & \frac{\partial^2 J}{\partial w_3^2} & \cdots & \frac{\partial^2 J}{\partial w_3 \partial w_m} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 J}{\partial w_m \partial w_1} & \frac{\partial^2 J}{\partial w_m \partial w_2} & \frac{\partial^2 J}{\partial w_m \partial w_3} & \cdots & \frac{\partial^2 J}{\partial w_m^2} \end{pmatrix}$$

The matrix H is **symmetric**

Finally, the **Taylor polynomials** centred around w_0 :

$$J(w) \approx J(w_0) + J'(w)|_{w=w_0} (w - w_0) + \frac{1}{2} J''(w)|_{w=w_0} (w - w_0)^2$$

Or for the multidimensional case:

$$J(\mathbf{w}) \approx J(\mathbf{w}_0) + \nabla J(\mathbf{w})|_{\mathbf{w}=\mathbf{w}_0} (\mathbf{w} - \mathbf{w}_0) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_0)^T H|_{\mathbf{w}=\mathbf{w}_0} (\mathbf{w} - \mathbf{w}_0)$$

To check whether a given point \mathbf{w} is a minimum...

- 1 Necessary conditions of extremum: **gradient = 0 in the point of interest**
- 2 Necessary and sufficient condition of local convexity: **Hessian ≥ 0 at the point of interest**
- 3 Sufficient conditions of local minimum: **gradient = 0 and Hessian ≥ 0 at the point of interest**
- 4 Necessary and sufficient condition of convexity: **Hessian ≥ 0 everywhere**
- 5 Sufficient conditions of global minimum: **gradient = 0 and Hessian ≥ 0 everywhere**

Basic optimisation methods

An important observation is that in general we don't know everything about the objective function (realistic problems are typically too complex). On the following part, it will be proposed some situations and some possible solutions (algorithms) to solve the problems.

- We only know the function $J(\mathbf{w})$, so it's not possible to use Taylor expansion and the only feasible solution is to **direct search** for minimum. Is possible to use different techniques: simulated annealing, genetic algorithms, particle swarm optimisation, ant colony optimisation. PRO: simple, CONS: we need infinite time available.
- We know the function $J(\mathbf{w})$ and the gradient $\nabla J(\mathbf{w})$, in this case we can use first-order Taylor expansions and then find minimum in local region.

Gradient descent algorithm

- 1 Initialize: set $i = 0$; select some $\mathbf{w}_0 = \mathbf{w}_{\text{start}}$
 - 2 Compute $\nabla J(\mathbf{w}_i)$
 - 3 Select the appropriate step size η_i
 - 4 Compute the step $\Delta \mathbf{w}_i = -\eta_i \nabla J(\mathbf{w}_i)$
 - 5 Perform step $\mathbf{w}_{i+1} \leftarrow \mathbf{w}_i + \Delta \mathbf{w}_i$
 - 6 Compute convergence test. If necessary, iterate from step 2.
- We know the function $J(\mathbf{w})$, the gradient $\nabla J(\mathbf{w})$ and also the Hessian H . It's possible to use the second-order Taylor polynomial and then the Newton-Raphson method. This method attempts to solve the problem of finding the minimum of a function by constructing a sequence $\{x_k\}$ from an initial guess x_0 that converges to a minimum by using a sequence of second-order Taylor approximations of J around the iterates.
 - Start at x_0
 - At each step τ compute the update as

$$x_{\tau+1} = x_{\tau} - \frac{J'(x_{\tau})}{J''(x_{\tau})}$$

PRO: simple, much faster than gradient descent. CONS: need to compute the Hessian

- The Quasi-Newton methods are the most popular, they use the gradient to approximate the second-order information.

Statistics

Statistics consists in collecting, organizing, analyzing and interpreting data. Statistical models are mathematical and are based on probability theory. It is also defined as the science of using empirical data to create models based on probability.

Terminology

- **Model**, an approximation to reality that retains only those characteristics that are useful to a certain purpose.
- **Population**: a collection of people, items, events about which we want to make inferences.
- **Sample**: a subset of the population.
- **Observation**: An individual element of the sample.

An important assumption is that our samples will be independent and identically distributed.

Statistical estimates

They typically measures central tendency, assuming that there is some “middle value” around which random fluctuations occur.

- **Expectation**, used under scenario 2: know the model that generates data

$$\text{continuous } \mathbf{x}: E\{X\} = \int_X \mathbf{x} f_X(\mathbf{x}) d\mathbf{x} \quad \text{discrete } \mathbf{x}: E\{X\} = \sum_i \mathbf{x}_i F_X(\mathbf{x}_i)$$

- **Mean (empirical expectation)**, used under scenarios 1 and 3

$$\text{mean}\{X\} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$$

- **Median**: the middle value after sorting \mathbf{x} in increasing order

$$\begin{array}{ll} 1, 3, 3, \mathbf{6}, 7, 8, 9 & 2, 2, 3, \mathbf{5}, \mathbf{6}, 8, 8, 9 \\ \text{Median} = \mathbf{6} & \text{Median} = (5 + 6) \div 2 \\ & = \mathbf{5,5} \end{array}$$

- **Mode**: the value that occurs more frequently
- **Variance**: how much the values vary around the mean

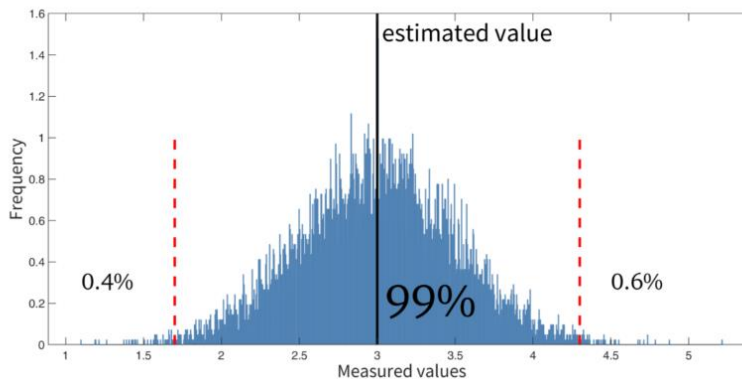
$$\text{theoretical } \sigma^2 = E\{(\mathbf{x} - E\{X\})^2\} \quad \text{empirical } \sigma^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \text{mean}\{X\})^2$$

- **Standard deviation**: $\sigma = \sqrt{\sigma^2}$
- **Percentile**: measure the position of a specific value within a data set. First, the data are sorted in ascending order, then the percentile of a given data point is computed by finding what percentage of data are lower than that point. For example, the 50th percentile is the value below which 50% of data are lower

- **Confidence intervals:** represent the probability that a sample will be in a given interval. Mathematically: given ϵ , the probability that the actual values will differ from the estimated location by more than ϵ is: $\delta = \Pr(|\hat{\mu} - \mu| > \epsilon)$

$[\hat{\mu} - \epsilon, \hat{\mu} + \epsilon]$ is a **confidence interval** at **level** $1 - \delta$

Example



In this case: $\epsilon = 1.3$ and $\delta = 0.99$

So, with probability 99% the values will be between 1.7 and 4.3

Non-parametric models

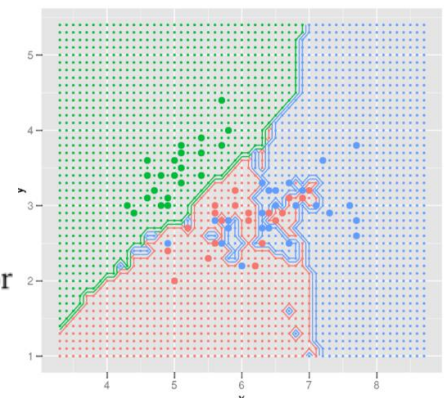
Non-parametric models are a class of models that do not make strong assumptions about the underlying data distribution or the functional form of the relationships between variables. Unlike parametric models, which have a fixed number of parameters and assume a specific data distribution (e.g., linear regression assumes a linear relationship between variables), non-parametric models are more flexible and adaptive. They can capture complex patterns and are particularly useful when data relationships are not well-understood or when assumptions about linearity and distribution are not met. The negative side is that the complexity (model size, number of parameters) is not pre-defined, but depends on the data

Nearest-neighbour classifier

- We are given a training set $X = \{x_1, x_2, \dots, x_n\}$.
- Then we have a point x^* that we want to classify.
- We search for the nearest point in the training set: $q = \arg \min_i ||x_i - x^*||$; the output will be the class of the nearest point: $y = t_q$

Properties

- + $O(1)$ “learning” complexity (just store the training set!)
- $O(nd)$ classification complexity
- + (sort of) works even for regression, $t \in \mathbb{R}$
- + Theoretical guarantee: For $n \rightarrow \infty$, error rate $\leq 2 \cdot \text{Bayes error}$
- May “overfit”, i.e., decision regions may have jagged borders



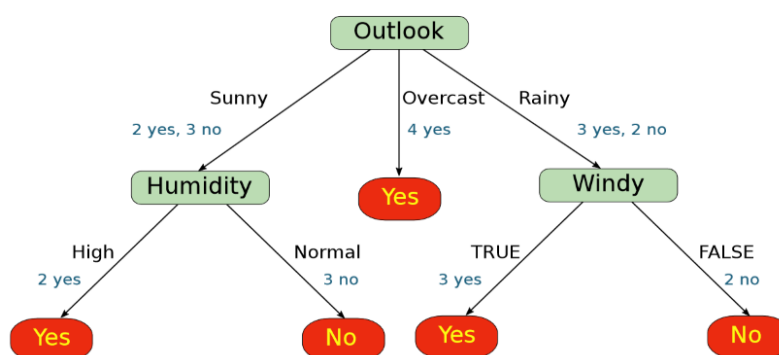
The **K-Nearest-neighbour classifier** is another version of the algorithm in which instead of only one point we select k nearest point. So, $\{n_1, \dots, n_k\} = \text{top } k \ ||x_i - x^*||$ and the output is $y = \text{mode}\{t_{n_1}, \dots, t_{n_k}\}$ -> it is a majority voting algorithm. Other variants use weighted nearest neighbours, condensed NN and approximated NN-

Decision trees

Assumption: categorical features

- Select one feature.
- Partition of the training set according to its possible values (**levels**).
- **For each** of the parts of the training set so obtained (one per level):
 1. If it contains data from one class only, then decide for that class.
 2. Else if it contains data from more than one class and all features have been used, then decide for the majority class.
 3. Else apply again the same procedure using a different feature.

Outlook	Temperature	Humidity	Windy	Play
overcast	hot	high	FALSE	yes
overcast	cool	normal	TRUE	yes
overcast	mild	high	TRUE	yes
overcast	hot	normal	FALSE	yes
rainy	mild	high	FALSE	yes
rainy	cool	normal	FALSE	yes
rainy	cool	normal	TRUE	no
rainy	mild	normal	FALSE	yes
rainy	mild	high	TRUE	no
sunny	hot	high	FALSE	no
sunny	hot	high	TRUE	no
sunny	mild	high	FALSE	no
sunny	cool	normal	FALSE	yes
sunny	mild	normal	TRUE	yes



Root Node: The top node in the tree, representing the initial input data.

Internal Nodes: These nodes represent questions or conditions on which the data is split. Internal nodes have branches to other internal nodes or leaf nodes.

Leaf Nodes: The terminal nodes that provide the final output or decision. In classification, leaf nodes represent class labels, and in regression, they represent predicted numerical values.

Properties

- + Interpretable (only one feature at a time)
- + Very quick learning
- Suboptimal (a **greedy** algorithm)
- Splits on features with many values are more likely
- With few data splits may be on irrelevant features (noise)
- Not suitable for quantitative (e.g., real-valued) features

Some improvements are the pruning of the tree according to some criterion (for irrelevant splits) and the use for real-valued features by divide value range into sub-ranges (quantization) and use that sub-ranges as categorical features.

DEF: an **ensemble classifier** is composed of **many classifiers**, a **learning procedure** that encourages diversity and an **aggregation rule**.

Random forest

It is a set of random forest is a set of random trees, trained on randomly resampled training sets. It is a type of **ensemble classifier**. The model is based on the following algorithm:

- Receive one training set X of cardinality n patterns (n observations).
- Create b new training sets of size n by randomly sampling from X with replacement.
- Train one random tree on each of the b new training sets.
- Inference: given one query pattern, classify the pattern with each tree and make the final decision by majority voting among the b individual decisions.

Random resampling with replacement (Bootstrap)

The second point of the algorithm is explained by this:

Given one training set X of cardinality n , a bootstrap sample is another training set $X^{(i)}$ obtained by randomly sampling with replacement n times from X .

$$\begin{array}{lcl}
 X^{(1)} & = & [A A C D E] \\
 X^{(2)} & = & [B C C D E] \\
 X & = & [A B C D E] \\
 X^{(3)} & = & [A A B D D] \\
 X^{(4)} & = & [A B C E E] \\
 X^{(5)} & = & [A A A C C]
 \end{array}$$

OSS: The decision tree induction rule will tend to select always the most discriminating features, so trees will resemble each other and will not be very independent. To solve this problem is possible to random sampling also for features.

Evaluation of classifiers

It is possible to decompose the expected square error objective J_{mse} as the sum of three terms:

$$J_{mse} = B + V + N$$

- **Bias (B):** the expected difference between the error of the best model and the error of our model (*it suggests problem with model*)
- **Variance (V):** the variance of our model's error (*it suggests problem with data*)
- **Noise (N):** a residual, irreducible quantity.

The bias can be lowered by using a **more adaptable** model -> can cause **overfitting**

The variance can be lowered by using a **less adaptable** model -> cause **underfitting**

So to sum up, low bias = a good performance on the training set. While a low variance = reliable performance in inference (good generalisation). The statistical problem of model selection consists in choosing the model that generalizes best. There are two possible strategies:

1. Evaluate quality empirically (ex post) and use this for designing a new and better learner.
2. Estimate quality theoretically (ex ante) and use this for designing the learner once for all.

Empirical estimates of generalization

Typically, the designer trains the machine using a **training set**, then he evaluates the quality using an independent **validation set**, if the quality is not sufficient, the structure of the model is changed. Finally, the costumer evaluates the model on the **test set**. An important parameter to

considers is the size of the datasets, because small training set leads to overfitting, while with small validation or test set it's not possible to be sure about the performance of the model.

A way to obtain more data to use is given by the resampling methods. Resampling means generating new samples from a given available sample and estimates the confidence intervals, from this is possible to generalize.

Another way is the **cross-validation**: the samples are splitted into a training set and a test set, then the model is evaluated. This can be done during iterative training, when the test value of the model starts to grow, we are overfitting (stop the procedure).

An alternative to this method is the **leave-one-out cross-validation**, in this case the samples are splitted in n training sets of size $n-1$ and corresponding n test sets of size 1. The model is then evaluated on each test obtaining the results for each test. The final estimate of our model is the average of all these one-point estimate. (possible to use the leave-k-out)

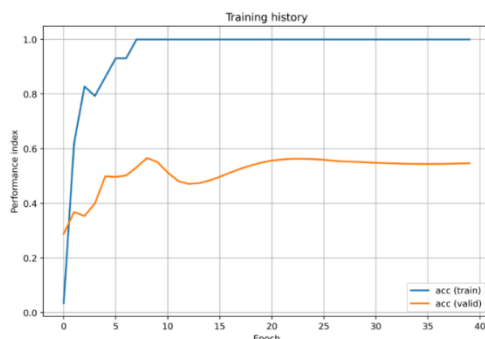


Figure 1: example of cross-validation

Leave-one-out validation algorithm

Input: Data set X of n rows

- 1: **for** $l: 1 \dots n$ **do**
- 2: Remove x_l from X , obtaining training set X_l with $n-1$ rows
- 3: Train learning machine M_l
- 4: Compute validation cost $R(x_l)$ using machine M_l
- 5: Add $R(x_l)$ to set S
- 6: **end for**
- 7: Train final learning machine M on X
- 8: Compute statistics on S (average cost, standard deviation of cost, min/max/percentiles of cost, ...)

Output: Machine M characterized by statistics over S

Finally, the last method is the **bootstrapping**.

Now is possible to consider indexes to measure the classification quality. Some of them are:

- **Contingency tables:** each entry of the matrix indicates the number of co-occurrences of pairs of events from events set E_1, E_2, \dots, E_k and F_1, F_2, \dots, F_h . They are useful to summarize and analyse the relationship between two or more categorical variables.

$$\begin{matrix} & E_1 & E_2 & \dots & E_k \\ \begin{matrix} F_1 \\ F_2 \\ \vdots \\ F_h \end{matrix} & \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1k} \\ c_{21} & c_{22} & \dots & c_{2k} \\ \vdots & \vdots & & \vdots \\ c_{h1} & c_{h2} & \dots & c_{hk} \end{bmatrix}
 \end{matrix}$$

Each row and column indicates one category

The entry c_{ij} represents the number of experiments where events E_i and F_j have been observed together. It's possible to find the frequencies using: $f_{ij} = c_{ij}/N$

- **Confusion matrix (C)** : it is a contingency table for **classification outputs vs actual classes**.

		Output class		
		Non-bot	Bot	
Target class	Non-bot	56.3%	3.7%	93.8%
	Bot	0.8%	39.1%	98.0%
		98.5%	91.3%	95.5%

The entry c_{ij} represents the number of times that the classifier decided w_j when in fact the true class was w_i

Figure 2: Example with classification model for bot vs human visitors in a web site

With the confusion matrix is possible to find the number of correct classification $n_{correct} = trC$ and then the accuracy $f_{correct} = \frac{n_{correct}}{n}$.

The dichotomic case (only two classes)

		Output	
		0	1
Target	0	NR CORRECT NEG	NR FALSE POS
	1	NR FALSE NEG	NR CORRECT POS

- Cases 00 and 11 are correct classifications
- Case 01 is called **false positive**
- Case 10 is called **false negative**

Note: in statistics event 01 is also called error of the first type and event 10 is called error of the second type

In this situation we can't trust the accuracy alone, in fact suppose to have the following example
Recognize publications written by T. Bayes in the whole Web

There is only one document by Bayes and (suppose) 10^{12} other documents on the Web.

Case 1: All correctly recognized.

$$\text{Accuracy} = \frac{10^{12} + 1}{10^{12} + 1} = 1$$

Case 2: We always classify as "not-by-Bayes".

$$\text{Accuracy} = \frac{10^{12}}{10^{12} + 1} = \frac{1000000000000}{1000000000001} \approx 1$$

! We got one of the two classes completely wrong (100% error), but the accuracy was almost perfect. So, it's useful to introduce new indexes:

- **Sensitivity/true positive rate:** probability to correct recognize class 1

$$\frac{\text{prob. true pos}}{\text{prob. true pos} + \text{prob. false neg}} = \frac{P(11)}{P(11) + P(10)} \approx \frac{c_{11}}{c_{11} + c_{10}}$$

- **Specificity/true negative rate:** probability to correctly recognize class 0

$$\frac{\text{prob. true neg}}{\text{prob. true neg} + \text{prob. false pos}} = \frac{P(00)}{P(00) + P(01)} \approx \frac{c_{00}}{c_{00} + c_{01}}$$

A good classifier has both high sensitivity and high specificity.

- **Precision/positive predictive values:** fraction of positive outputs that was actually positive

$$\frac{P(11)}{P(11) + P(01)} \approx \frac{c_{11}}{c_{11} + c_{01}}$$

- **Recall = sensitivity**
- **F measure:** it is a combined index that express with a single number the evaluation

$$F = 2 \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

$$= \frac{2P(11)^2}{2P(11) + P(10) + P(01)}$$

$$\approx \frac{2c_{11}^2}{2c_{11} + c_{10} + c_{01}}$$

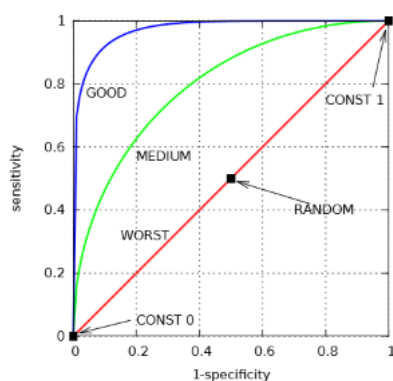
A generalized version exists:

$$F_\beta = (1 + \beta^2) \frac{\text{precision} \times \text{recall}}{\beta^2 \text{precision} + \text{recall}}$$

OSS: Neither precision/recall nor F consider the probability of being correct on the negative class. This is an advantage when it is easy to be correct on the negative class.

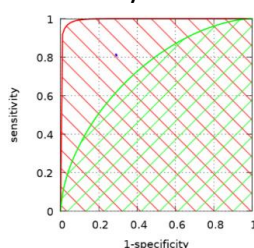
Receiver operating characteristic curve (ROC)

It is a graphic indicator for a binary classifier whose operation depends on some parameter θ (ex: a threshold) and changing this parameter changes the performance of the classifier itself. We plot sensitivity as the abscissa and false positive rate (1-specificity) as the ordinate for varying values of θ .



A good classifier is near the top-left corner, this is also the one enclosing the largest area.

The **Area Under the Curve (AUC)** is another quality index to synthesize the information in one value.



NOTE: these indexes cannot directly be used as objective function because they cannot be expressed as a loss function. For instance:

$$F1 = \frac{2c_{11}^2}{2c_{11} + c_{10} + c_{01}} = \sum_{l=1}^n \lambda(y_l, t_l) \Rightarrow \lambda(y_l, t_l) = ???$$

NOTE: if I want to evaluate my solution according to two incompatible criteria the problem becomes multi-objective and some external criteria will be needed.