

Marco Casu

MACHINE LEARNING



SAPIENZA
UNIVERSITÀ DI ROMA

Faculty of Information Engineering, Computer Science and Statistics
Department of Computer, Control and Management Engineering
Master's degree in Artificial Intelligence and Robotics

This document summarizes and presents the topics for the Machine Learning course for the Master's degree in Artificial Intelligence and Robotics at Sapienza University of Rome. The document is free for any use. If the reader notices any typos, they are kindly requested to report them to the author.



CONTENTS

1	Introduction	3
1.1	Basic Definition of a ML Problem	3
1.2	Types of Machine Learning Problems	5
1.3	Performance Evaluation	7
1.3.1	Concept Learning	8
1.3.2	Error Estimation and Performance Metric	10
1.3.3	Unbiased Estimation	11
1.3.4	The Cross Validation Algorithm and others Performance Metrics	11

CHAPTER

1

INTRODUCTION

1.1 Basic Definition of a ML Problem

In this chapter we will introduce the basics of what is a machine learning problem, giving a mathematical definition. informally, with machine learning we define the use of knowledge (data) to improve the performance of a given program, using past experiences.

Generally, we use machine learning to solve problems with no deterministic solutions, trying to find an approximate one (such as recognizing what animal is represented in a given photo).

Usually, a machine learning problem consists in three main component:

- T : the given task
- P : a performance metric
- E : the past experiences (the data)

Let's see an example, we want to model a program capable of playing *Checkers*.



The task T is to play the game, the performance metric P measure the ratio of win in a tournament, and the experiences is given by the past match. We can create this matches by letting the computer play against himself, or by making it play against a human. We can consider two types of target functions that models the behavior of the model:

- $ChooseMove : Board \rightarrow Move$
- $V : Board \rightarrow \mathbb{R}$

Board is the set of the possible board configurations, move is the set of feasible moves. The image of the function V should represents the validity of a board in the following sense:

- if $b \in Boards$ is a configuration that represents the win of the player, then $V(b) = 1$
- else if $b \in Boards$ is a configuration that represents the win of the opponent, then $V(b) = -1$
- else if $b \in Boards$ is a configuration that represents a draw, then $V(b) = 0$
- else, $V(b) = V(b')$ where b' is the best final board state that can be achieved starting from b and playing optimally until the end of the game.

The function V can be used to predict the next move by considering all possible boards that can be obtained from the feasible moves. We focus on the function V , this should be an optimal model, but is not computable, because we cannot tell if a play is "optimal" (this is the goal of the model), so we want to consider a function that approximate the behavior of V .

We consider a function $\hat{V} : Board \rightarrow \mathbb{R}$ defined as follows:

$$\hat{V}(b) = \sum_{i=0}^6 w_i f_i(b) \quad (1.1)$$

where $\mathbf{w} = (w_i)$ is real coefficients, and the functions f_i represents some features of the given board:

- $f_1(b)$: number of black pieces on b
- $f_2(b)$: number of red pieces on b
- ecc...

We don't know if some features are useful to predict which move is optimal, is important that this features model the knowledge of the *domain* (in this case, the game of Checkers).

with *learning the function* \hat{V} , we intend to finds the coefficients \mathbf{w} which make the function \hat{V} more similar to V as possible. There are various method to find these coefficients.

Let's introduce some notation:

- with V we define the **target function** (always unknown and uncomputable).
- with \hat{V} we define the **learned function**, the approximation of V that we want to find.
- with $V_{train}(b)$ we define the value of V obtained at b , where b is a part of a data set that is given. We will use the values of V_{train} on the data set to synthesize \hat{V} .
- with D we define the given data set:

$$D = \bigcup_{i=1}^n \{b_i, V_{train}(b_i)\} \quad (1.2)$$

n is the number of the available data.

The iterative method given in the Algorithm 1 is an informal example of how we can find the values for \mathbf{w} . In this case c is a small constant, usually in $(0, 1]$, to moderate the rate of learning.

The function $error(b)$ is computable only on the sample b given in the dataset B , the goal of the method is to converge to a local minimum for the function

$$\sum_{b_i \in D} error(b_i). \quad (1.3)$$

Once we synthesize \hat{V} with this method, we can make the model play against human and track the result to use it as an additional dataset. The diagram in image 1.1 represents the process of designing an artificial intelligence agent.

Algorithm 1 LMS weight update rule

Require: D, V_{train}, k
 initialize \mathbf{w} with small random values
for k times **do**
 select a sample b from D
 $error(b) = V_{train}(b) - \hat{V}(b)$
 for each feature i **do**
 $w_i \leftarrow w_i + c \cdot f_i(b) \cdot error(b)$
 end for
end for

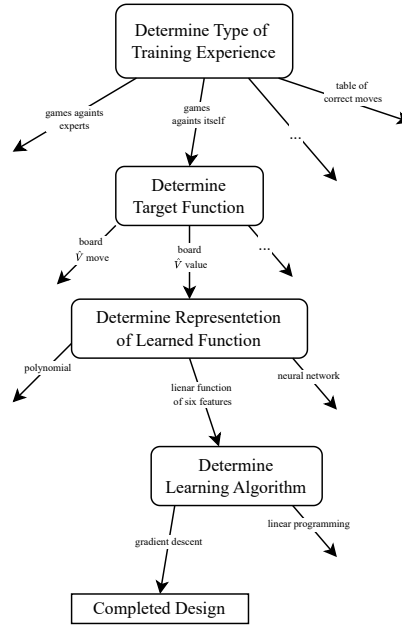


Figure 1.1: Design Choices

1.2 Types of Machine Learning Problems

There are different types of machine learning problems:

- Supervised Learning
 - Classification
 - Regression
- Unsupervised Learning
- Reinforcement Learning

Definition 1 Given a function $f : X \rightarrow Y$, and a training set $X_D \subset X$ containing information about f , **learning** the function f means computing an approximated function \hat{f} such that is much close as possible to f on X

$$\hat{f}(x) \simeq f(x), \quad x \in X \quad (1.4)$$

more formally, we want to find \hat{f} that minimize the following integral:

$$\int_X |\hat{f}(x) - f(x)| dx \quad (1.5)$$

if $X \subseteq \mathbb{R}^n$ and is uncountable, and

$$\sum_{x \in X} |\hat{f}(x) - f(x)| \quad (1.6)$$

if X is countable.

This is not a simple problem, f is not computable, so the difference $f - \hat{f}$, and X is usually uncountable or a big set, way bigger than the training set X_D .

Machine learning problems can be classified in terms of the input data set D , given a target function $f : X \rightarrow Y$ a problem is

- a **supervised learning** problem if $D = \bigcup_{i=1}^n \{(x_i, y_i)\} \subset X \times Y$.
- an **unsupervised learning** problem if $D = \bigcup_{i=1}^n \{(x_i)\} \subset X$.
- a **reinforcement learning** problem, the condition on the input dataset will be discussed later.

The problems can also be classified in terms of the target function $f : X \rightarrow Y$

$$X = \begin{cases} A_1 \times \cdots \times A_m, & A_i \text{ finite set} & \textbf{(Finite Discrete Problem)} \\ \mathbb{R}^n & & \textbf{(Continuous)} \end{cases}$$

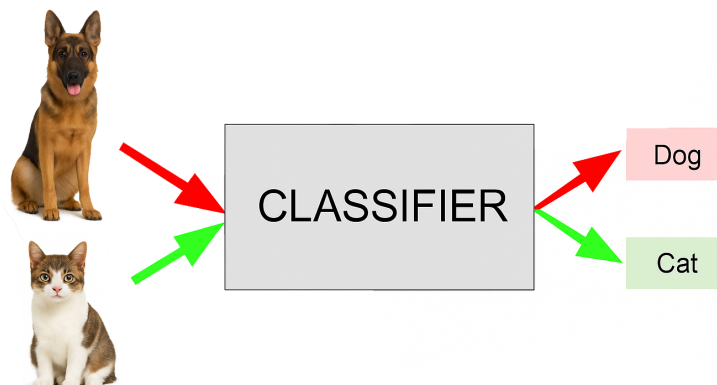
$$Y = \begin{cases} \mathbb{R}^k & \textbf{(Regression)} \\ \{C_1, C_2 \times C_k\} & \textbf{(Classification)} \end{cases}$$

special case (**Concept Learning**):

$$X = A_1 \times \cdots \times A_m, \quad A_i \text{ finite set}$$

$$Y = \{0, 1\}$$

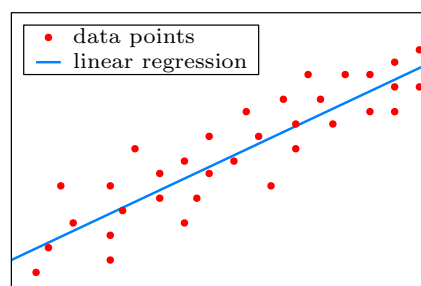
Classification problems is also known as *Pattern Recognition Problems*, the goal is to return the class to which a specific instance belong.



Some examples are:

- face/object/character recognition
- speech/sound recognition
- medical diagnosis
- document classification.

Regression problems consists in approximating real valued functions.

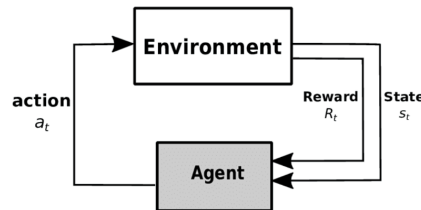


Unsupervised learning discovers data patterns without a specific output. The main goal is to understand what's normal in a dataset. *Clustering* is a key technique that groups similar data points. Applications include customer segmentation, image compression, and bioinformatics motif learning.

Reinforcement learning consists in learning a policy, which is a strategy that tells the agent what action to take in any given situation (or state) to maximize its long-term reward. This is often represented as a function that maps a state to an action. Unlike supervised learning, where the model is trained with labeled examples, in reinforcement learning agents learn through a process of trial and error.

They don't have a correct answer to guide them. Instead, they receive sparse and time-delayed rewards. This means the feedback (the reward) for a good action might not come immediately, and it might be a simple, numerical value (like +1 for winning a game or -1 for losing). Some examples are:

- Game playing: Think of programs that have learned to play chess, Go, or video games better than humans.
- Robotic tasks: A robot learning to navigate a room, pick up an object, or perform a specific manufacturing task.
- Any dynamical system with an unknown or partially known model: This is a broad category that includes things like optimizing traffic flow, managing a power grid, or controlling financial trading strategies.



In concept learning, the output consists in only two classes, the target function $c : X \rightarrow \{0, 1\}$ maps any kind of input in one of two distinct values. The following example model a program that predict if is a good day to play Tennis, the input set is

$$X = Day \times Outlook \times Temperature \times Humidity \times Wind$$

the output set is $PlayTennis = \{Yes, No\}$. An example of data samples:

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

1.3 Performance Evaluation

Usually, we call **Hypothesis** a possible learned function h , and we define H as the **Hypothesis space**, such space contains all possible function that can be learnt (all possible approximation of the target function). In this terms, a learning problem is described as a search in the hypothesis space using the given dataset D , that aims to find the best possible approximation h^* :

$$h^* = \arg \max_{h \in H} Performance(h, D) \quad (1.7)$$



1.3.1 Concept Learning

We focus now on the concept learning (binary classification), let's consider a target function

$$c : X \rightarrow \{0, 1\} \quad (1.8)$$

let $D \subset \{X, Y\}$ to be the sample set

$$D = \bigcup_{i=1}^n \{(x_i, c(x_i))\} \quad (1.9)$$

we denote X_D the point of X in the sample set. For now, we assume that the sample set does not have noise:

- noise dataset $D = \bigcup_{i=1}^n \{(x_i, c(x_i) + \varepsilon_i)\}$, $\varepsilon_i \in \mathbb{R}$
- perfect dataset $D = \bigcup_{i=1}^n \{(x_i, c(x_i))\}$

in general, the dataset is noisy. We want to find a function \hat{f} that approximate c , as we just said, the hypothesis space H is the set of all hypothesis h . Each hypothesis is a possible solution to the problem

(1.7), it is possible to compare an hypothesis h with c on a sample in X_D .

Definition 2 Given a target function c and a set of sample point X_D , an hypothesis h is **consistent** if

$$h(x) = c(x), \quad \forall x \in X_D. \quad (1.10)$$

The **Version Space** $VS_{H,D}$ is the set of all consistent hypothesis:

$$VS_{H,D} = \{h : h(x) = c(x), \quad \forall x \in X_D\} \subset H. \quad (1.11)$$

A solution that does not lie in the version space is probably not a good solution.

Let's consider an example, let c to be the following target function

$$c : \mathbb{N} \rightarrow \{-, +\} \quad (1.12)$$

the dataset is

$$D = \{(1, +), (3, +), (5, +), (6, -), (8, -), (10, -)\} \quad (1.13)$$

given the hypothesis space H , we consider four different hypothesis

$$\begin{aligned} h_1(n) = + &\iff n \text{ is odd} \\ h_2(n) = + &\iff n \leq 5 \\ h_3(n) = + &\iff n \text{ is prime} \\ h_4(n) = + &\iff n \in \{1, 3, 5\} \end{aligned}$$

if $n = 11$ we have

$$\begin{aligned} h_1(11) &= + \\ h_2(11) &= - \\ h_3(11) &= + \\ h_4(11) &= - \end{aligned}$$

we can't tell which of the four hypothesis is better than the others. Let's now consider a new hypothesis space H' , defined as the power set of H

$$H' = \mathcal{P}(H) = \{I : I \subseteq H\} \quad (1.14)$$

The set H' contains more information than H , let $\theta \in H'$, we define the value of θ on x in a different way

$$\theta = \bigcup_i \{h_i\} \quad (1.15)$$

$$\theta(x) = \text{maj} \bigcup_i \{h_i(x)\} \quad (1.16)$$

where the function maj returns the element that is more frequent in a set. It's important to know that, if you have a set where two different items appear the exact same number of times, and no other item appears more often than either of them, then the majority element can't be identified, and the value of the function maj is undefined.

With the hypothesis space H' , the point $n = 11$ can't be classified

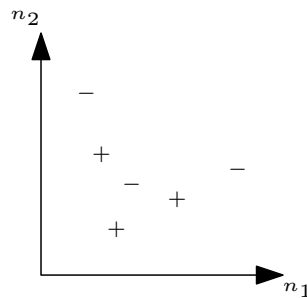
$$\theta(11) = \text{maj}\{h_1(11), h_2(11), h_3(11), h_4(11)\} = \text{maj}\{+, -, +, -\}. \quad (1.17)$$

Even if H' is more powerful than H , this hypothesis space can't classify an element that can be classified in H , this problem is called overfitting.

Now let's consider the following target function

$$c : \mathbb{N}^2 \rightarrow \{+, -\} \quad (1.18)$$

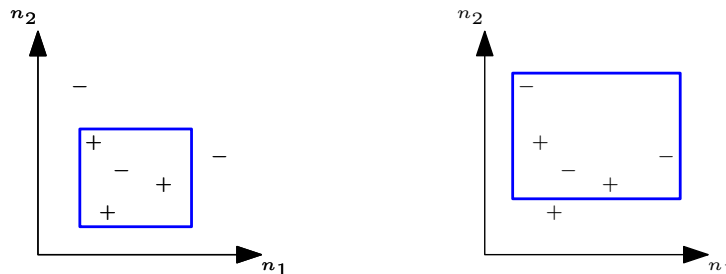
with the following dataset D that can be plotted on a 2D plane:



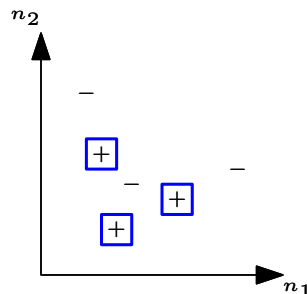
We consider H as the set of the functions that, assigns at all the points in a specific rectangle the $+$ value, and the $-$ value for all the points outside.

$$h \in H \iff (\exists R := \{(x, y) : a \leq x \leq b \wedge c \leq y \leq d\} : h(x, y) = + \iff (x, y) \in R)$$

For that data samples, does not exists a consistent hypothesis, because does not exists a rectangle that contains all the $+$ value point, letting outside the $-$ value points.



If we consider the hypothesis space $H' = \mathcal{P}(H)$, geometrically, we can represent a function by a finite number of rectangles, in this case consistent hypothesis are allowed, but no $-$ value points can be predicted.



Even in this example, a powerful representation of the hypothesis space lead to less generalization. This is known (as previously said) as **overfitting**, when the hypothesis space is not powerful enough, the problem could not be represented well, this is known as **underfitting**.

If n is a number that determines how big an hypothesis space, we can observe the following trend about the performance of a model.

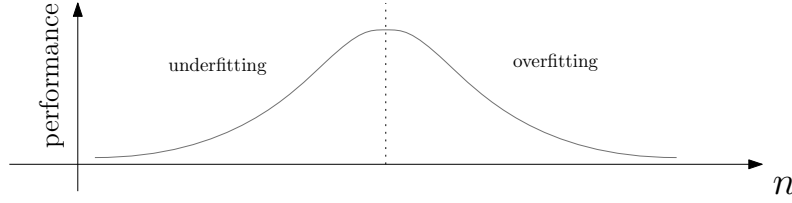


Figure 1.2: overfitting and underfitting trend

1.3.2 Error Estimation and Performance Metric

We want to define a metric for the performance of a model, to estimate how a given hypothesis h is good respect to the others. We have to introduce some statistical methods, let's consider a target function

$$f : X \rightarrow Y \quad (1.19)$$

where Y is countable (classification problem). Let \mathcal{D} to be a probability distribution over X .

Definition 3 If X is a set, \mathcal{D} is a probability distribution on X if

$$\forall x \in X, \quad \mathcal{D}(x) \in [0, 1] \quad (1.20)$$

and

$$\begin{aligned} \sum_{x \in X} \mathcal{D}(x) &= 1 \text{ if } X \text{ is countable} \\ \int_X \mathcal{D}(x) dx &= 1 \text{ if } X \text{ is uncountable} \end{aligned}$$

Definition 4 Given a target function $f : X \rightarrow Y$, an hypothesis h and a probability distribution \mathcal{D} on X , the **true error** is the probability that h will misclassify an instance $x \in X$ drawn according to the distribution \mathcal{D} :

$$\text{error}_{\mathcal{D}}(h) = \mathbb{P}_{x \sim \mathcal{D}}(f(x) \neq h(x)) \quad (1.21)$$

Note: With $x \sim \mathcal{D}$, we mean that x was extracted from X according to the probability distribution \mathcal{D} . We remind that, if $\mathcal{D}(x) = p$, then, the probability of extracting x from X by taking a random element is p .

Since we can't compute f for all $x \in X$, the true error can't be computed. We need to define an estimation of the error given by the sample set extracted from X .

Definition 5 Given a target function $f : X \rightarrow Y$, an hypothesis h and a sample (finite) set $\mathcal{S} \subset X$, the **sample error** is defined as follows:

$$\text{error}_{\mathcal{S}}(h) = \frac{1}{|\mathcal{S}|} \sum_{x \in \mathcal{S}} \delta(f(x) \neq h(x)) \quad (1.22)$$

where

- $\delta(\phi) = 1$ if ϕ is true
- $\delta(\phi) = 0$ if ϕ is false.

The sample error can be computed since, for each $x \in \mathcal{S}$, we know the value of $f(x)$. Is an approximation of the true error that depends from the sample set \mathcal{S} .

Since $\text{error}_{\mathcal{S}}(h)$ depends from the choice of \mathcal{S} , is not fixed, but we can model it as a random process, by extracting n random values from X according to the distribution \mathcal{D} to construct \mathcal{S} , and we can evaluate the expected value of $\text{error}_{\mathcal{S}}(h)$ denoted

$$\mathbb{E}(\text{error}_{\mathcal{S}}(h)). \quad (1.23)$$



We want to formalize this expected value, in this case the sample error is a random variable that assign to each subset \mathcal{S} of X of size n a number between 0 and 1:

$$\text{error}_{\mathcal{S}}(h) : \Omega \rightarrow [0, 1] \quad (1.24)$$

$$\Omega = \{\mathcal{S} : \mathcal{S} \subset X \wedge |\mathcal{S}| = n\} \quad (1.25)$$

in this context, n is fixed. The probability of getting a certain value γ from this random variable, is the probability to extract from X a subset \mathcal{S} of n items such that, the sample error is γ , and this depends from the probability distribution \mathcal{D} on X . The expected value $\mathbb{E}(\text{error}_{\mathcal{S}}(h))$ is now well defined.

1.3.3 Unbiased Estimation

Definition 6 The **bias** is defined as the expected value of the difference between the sample error and the true error:

$$\mathbb{E}(\text{error}_{\mathcal{S}}(h) - \text{error}_{\mathcal{D}}(h)). \quad (1.26)$$

We want to find an unbiased sample error, in such case:

$$\text{bias} = 0 \implies \mathbb{E}(\text{error}_{\mathcal{S}}(h)) = \text{error}_{\mathcal{D}}(h) \quad (1.27)$$

To compute an unbiased estimation, we need to train an evaluate an hypothesis h on different set, let D to be the dataset, we have to split D in two disjoint set

$$D = T \cup S \quad (1.28)$$

$$T \cap S = \emptyset \quad (1.29)$$

Usually, $\frac{|T|}{|D|} \simeq \frac{2}{3}$. We use T to train our learning function to get h , and then, we calculate the sample error on S

$$\text{error}_S(h) = \frac{1}{|S|} \sum_{x \in S} \delta(f(x) \neq h(x)) \quad (1.30)$$

is ideal to choose T and S such that they have similar probability distribution over the features, in this case the random variable error_S is an unbiased estimator for the true error $\text{error}_{\mathcal{D}}$.

1.3.4 The Cross Validation Algorithm and others Performance Metrics

There exists an algorithm to estimate the expected value of the sample error, more the dataset D is large, more the estimation will be precise. The algorithm 2 estimate the expected value of the error, with L is denoted a fixed learning algorithm:

- $h = L(T)$, is the result of the learning algorithm L applied on the training set T

Algorithm 2 K-Fold Cross Validation

Require: D, k, h, L

partition D in k disjoint sets S_1, S_2, \dots, S_k

for $i = 1, 2, \dots, k$ **do**

$T_i \leftarrow D \setminus S_i$

$h_i \leftarrow L(T_i)$

$\delta_i = \text{error}_{S_i}(h_i)$

end for

return $\text{error}_{L,D} = \frac{1}{k} \sum_{i=1}^k \delta_i$

We define the accuracy of a learning algorithm L as

$$\text{accuracy} = 1 - \text{error}_{L,D}. \quad (1.31)$$

The cross-validation algorithm can be used to compare the accuracy of two different learning methods L_a and L_b , as shown in algorithm 3.

Algorithm 3 Accuracy Comparator

Require: D, k, h, L
 partition D in k disjoint sets S_1, S_2, \dots, S_k
for $i = 1, 2, \dots, k$ **do**
 $T_i \leftarrow D \setminus S_i$
 $h_a \leftarrow L_a(T_i)$
 $h_b \leftarrow L_b(T_i)$
 $\delta_i = \text{error}_{S_i}(h_a) - \text{error}_{S_i}(h_b)$
end for
return $\bar{\delta} = \frac{1}{k} \sum_{i=1}^k \delta_i$

Now that we defined the sample error, we can give a formal definition of overfitting. Let h to be an hypothesis for a model, h is overfitting if exists an hypothesis h' such that

$$\text{error}_S(h) < \text{error}_S(h') \quad (1.32)$$

$$\wedge \quad (1.33)$$

$$\text{error}_D(h) > \text{error}_D(h') \quad (1.34)$$

Let's consider other performance metrics in binary classification. Let $f : X \rightarrow \{-, +\}$ to be the target function and let D to be a sample set such that, 90% of elements in D is of class $+$. An hypothesis that always returns $+$ will have an accuracy of 90%, in this scenario the dataset is unbalanced, so the accuracy is not a good performance metric for the model. We can consider a table that counts the number of points in the sample set that are well classified or misclassified by an hypothesis:

true class	predicted class	
	+	-
+	true positive	false negative
-	false positive	true negative

we can define two additional metrics that is useful when we deal with binary classification:

- the **recall** is the ability of the hypothesis to avoid false negatives and is defined as follows

$$\frac{\text{true positive}}{\text{true positive} + \text{false negative}} \quad (1.35)$$

- the **precision** is the ability of the hypothesis to avoid false positives and is defined as follows

$$\frac{\text{true positive}}{\text{true positive} + \text{false positive}} \quad (1.36)$$

the importance of these metrics depend on the application.

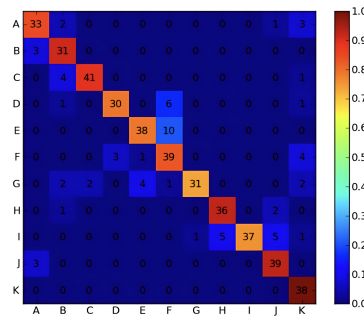


Figure 1.3: an example of a confusion matrix

For the classification problems we can define an extension of the previous table, called **confusion matrix**, and report how many instances of class C_i are classified in class C_j , the main diagonal contains the accuracy for each class. An example is shown in figure 1.3.

We consider now some performance metrics for the regression problems. Let $f : X \rightarrow \mathbb{R}^d$ to be the target function, and let \hat{f} to be the learned function, for each sample $(x_i, f(x_i))$ in the dataset, we can compute the euclidian distance

$$|\hat{f}(x_i) - f(x_i)|. \quad (1.37)$$

Let n to be the number of samples, the three main metrics are the following:

- Mean Absolute Error

$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{f}(x_i) - f(x_i)| \quad (1.38)$$

- Mean Squared Error

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{f}(x_i) - f(x_i))^2 \quad (1.39)$$

- Root Mean Squared Error

$$RMSE = \sqrt{MSE} \quad (1.40)$$

The cross validation algorithm can be extended for the regression problems as shown in algorithm 4.

Algorithm 4 K-Fold Cross Validation for Regression

Require: D, k, h, L

partition D in k disjoint sets S_1, S_2, \dots, S_k

for $i = 1, 2, \dots, k$ **do**

$T_i \leftarrow D \setminus S_i$

$h_i \leftarrow L(T_i)$

$\delta_i = MAE_{S_i}(h_i)$

end for

return $\text{error}_{L,D} = \frac{1}{k} \sum_{i=1}^k \delta_i$
