

QML-Mod2-Classical Machine Learning

Riccardo Marega

March 2025

Indice

1	Introduction to classical machine learning -22/03/2025	2
1.1	Input data	3
1.2	Machine Learning examples	4
1.2.1	Code overview	5
1.3	Cross validation and hyperparameter tuning	5
1.3.1	Evaluation metrics for classification	6
2	Introduction to deep learning -28/03/2025	7
2.1	Training an artificial neural network	7
2.2	Artificial neural networks: tasks	12
2.2.1	Binary classification	12
3	Advanced topics -29/03/2025	15
3.1	Convolutions from scratch	15
3.1.1	Types of convolutions	17
3.2	Convolutional neural networks for image classification	18
3.3	Convolutional neural networks for classification task	19
3.3.1	Image classification	19
3.3.2	Architectures	19
3.3.3	Transfer learning and fine tuning	21
3.4	Data preparation	22
3.4.1	Data augmentation	22
3.5	Optimizers	23
3.6	Loss function	23
3.7	Metrics	25

1 Introduction to classical machine learning -22/03/2025

Despite maybe we were not aware but we've already trained models, in our life, with machine learning algorithms: an example is given by the tools asking the users to spot cars, traffic signals, etc. in some website before entering them. An other example is given by all the filters that one can apply before taking a picture while using social media.

Artificial Intelligence vs Machine Learning vs Deep Learning

- artificial intelligence: any technique that enables computers to mimic human intelligence. It includes machine learning
- machine learning: a subset of AI that includes techniques that enable machines to improve at tasks with experience. It includes deep learning
- deep learning: a subset of machine learning based on neural networks that permit a machine to train itself to perform tasks

Machine learning teaches computers to do what comes naturally to humans and animals: learning from experience. Machine learning algorithms use computational methods to "learn" information directly from data, without relying on a predetermined equation. These algorithms adaptively improve their performance as the number of samples available for learning increases.

Machine learning includes two kinds of modalities to train the algorithm: supervised and unsupervised learning. In the latter we don't give any label to the algorithm and we let it find itself the correct answer to the given problem.

The choice of using ML algorithms arise when:

- we can't rely on rule-based systems
- the rules depend on too many variables, and many of them overlap or need to be optimized
- scalability becomes an issue

An example of machine learning algorithm is to be found in surgical data science, where AI is at service of surgeons.

The choice of the best algorithm to apply has to be done considering the dimension and the type of data one has to manage.

$$\text{machine learning} \left\{ \begin{array}{l} \text{Unsupervised learning} \left\{ \begin{array}{l} \text{Clustering} \end{array} \right. \\ \text{Supervised learning} \left\{ \begin{array}{l} \text{Classification} \\ \text{Regression} \end{array} \right. \end{array} \right.$$

Supervised learning: the model is trained under the assumption that for each input the label is known.

Unsupervised learning: the model identifies meaningful patterns within the input data, which does not come with any label.

In supervised learning methods, the output is already known, and the goal of training is map inputs to the corresponding outputs. To build a model, the machine learning is fed with a large amount of input data along with their corresponding labels. Supervised learning uses classification and regression techniques to develop predictive models. The difference between classification and regression is that the first one does a prediction of discrete outputs while the second one does a prediction of continuous outputs.

For unsupervised learning case it's not necessary to supervise the model or provide labeled input data. The algorithm begins to learn from the data without guidance. The model uses unlabeled data to identify new information. Since there are no known output values to establish a logical relationship between input and output, specific techniques are used to extract rules, patterns, and grouping of similar types. Machine using unsupervised learning algorithms discover patterns to find meaningful insights. For example, the system can learn to distinguish between dogs and cats by understanding the features and characteristics of each animal.

Every machine learning workflow starts with three key questions:

- What type of data are we working with?
- What do we want to extract or learn from the data?
- What is the context or domain of application?

The main problem with unsupervised learning is that evaluating its performance is not always immediate. We choose unsupervised learning when we aim to understand: the distribution of data, the clustering of data based on similar characteristics and the dimensionality reduction (needed to obtain a more concise and efficient representation).

Note that one can combine supervised and unsupervised learning.

1.1 Input data

Standard algorithms of ML usually take as inputs features and characteristics extracted from the data set. These features are handcrafted or manually designed or selected. Features are relevant pieces of information that help solve the task at hand. **The quality of a machine learning model depends on its ability of extracting features.**

The data is divided in a training set (used exclusively during the training phase) and a validation set (used to monitor the process and optimize model trainings). An algorithm is subject to under-fitting of the data set if it has a poor performance on the training set. This happens because the model fails to learn the relationship between input data and their corresponding data. Otherwise

it is said being subject to overfitting is it works well with the data training but has poor results with data of the validation set. This occurs because the model memorize the data it has seen and is unable to generalize to unseen examples.

For a fully functional training of the model we introduce a third data set called validation set. These data will be used exclusively during the training phase to monitor the learning process and optimize model training. This is called the validation set because it serves to validate the results obtained in the training set. If the performance is poor, we may need to adjust the model's hyperparameters (in the case of ML) and retrain the model until the validation results are satisfactory. When implementing this third dataset the testing set will be exclusively used after training is complete to evaluate model's performance.

Cross validation strategy : one can repeat the whole training of the model swapping every time the data from one set to another (what was first, for example, in the training set now will go in the validation set).

Generalization In Machine Learning, generalization refers to a model's ability to perform well on unseen data, meaning data that was not used during training. A well-generalized model captures the underlying patterns in the data rather than memorizing specific examples, allowing it to make accurate predictions on new inputs. Poor generalization can lead to overfitting (where the model performs well on training data but poorly on new data) or underfitting (where the model fails to learn meaningful patterns from the training data).

When an algorithm makes consistent mistakes, we say it has a "bias".

A **bias** is a distortion of the training data which is propagated in the algorithm.

The goals related to the supervised trainings are a low error during training, validation and testing phases. The selection of the algorithm is based on:

- velocity of the training
- storage capacity
- accuracy on new data prediction
- transparency and interpretability

1.2 Machine Learning examples

- **Logistic regression**: the model is trained to predict the probability of a binary choice. Due to its simplicity, logistic regression is commonly used as a starting point for binary classification problems. This algorithm is normally used when the data can be clearly separated by

a single decision boundary and in general as baseline to evaluate more complex classification methods.

- **k nearest neighbor (kNN)**: is a pattern recognition algorithm based on distance of the data. If a subject A has, for example, characteristic close to the one of a subject B, then probably they are associated to the same category. This kind of algorithm is normally used: for establishing baseline learning rules, when memory usage is not a major problem and when prediction speed of the trained model is not a major concern.

A **decision tree** is a supervised learning algorithm used for both classification and regression tasks. It has a hierarchical tree structure composed of a root node, branches, internal nodes and leaf nodes. A decision tree starts with the root node, which has no incoming branches. The outgoing branches from the root lead to internal nodes, also known as decisions nodes. Both root and internal nodes perform evaluations to split the data into homogeneous subset, which are represented by leaf nodes or terminal nodes. Leaf nodes represent all possible outcomes (either continuous or discrete) within the dataset.

ML learning models in training phase learn a set of rules which depends both on the data set and a fixed combination of hyperparameters. The automatic learning of a model is not a single process, indeed it is necessary to experiment different models fixing different values for the hyperparameters.

1.2.1 Code overview

```
1 from sklearn.ensemble import RandomForestClassifier
2
3 rf_clf = RandomForestClassifier(n_estimators=?, max_depth=?)
```

Imports the `RandomForestClassifier` model from the scikit-learn library. This model is a supervised learning algorithm based on decision trees. Creates an instance of the Random Forest classifier with two key parameters:

- `n_estimators=?`: defines the number of trees in the forest (typically between 10 and 1000).
- `max_depth=?`: sets the maximum depth of the trees (can be `None` to allow them to grow until all leaves are pure).

1.3 Cross validation and hyperparameter tuning

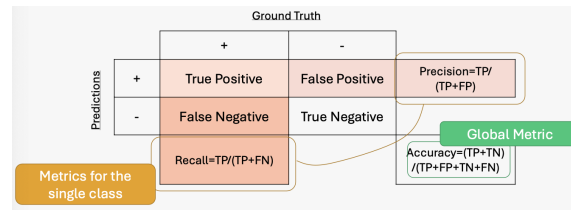
We divide the training and the validation set in multiple sets so that the model can be trained with all the available data. The model's performance for the same model with different hyperparameters combinations strongly depends on the specific data split. Each model is trained and evaluated only once, so its performance is tied to that single evaluation. But this raises an issue: we might get very different results when training and validating on different subsets of the same data. What

if we could split the training and validation sets multiple times, each time using different subsets of data, train and evaluate our models repeatedly, and observe model performance across several rounds of evaluation? → **k-fold Cross Validation**.

Cross validation is statistical method used to estimate the ability of different models of performing automatic learning. This procedure is done by defining a parameter k which represents the number of sets in which every set is subdivided. In this approach the dataset is randomly divided into k groups of roughly equal size. The model is trained on k-1 folds and validated on remaining fold. This process is repeated k times, each time using different fold as the validation set.

We can perform cross validation for hyperparameter tuning in either inner loops or outer loops

1.3.1 Evaluation metrics for classification



Accuracy: $\frac{TP+TN}{TP+FP+FN+TN}$

Recall (Sensitivity/True positive rate): $\frac{TP}{TP+FN}$

Precision: $\frac{TP}{TP+FP}$

Specificity: $\frac{TN}{TN+FP}$

F1 score: $2 \times \frac{Precision \times Recall}{Precision + Recall}$

ROC Curve & AUC (Area under the curve)

- **ROC curve:** Plot true positive rate (recall) vs. False Positive rate ($FPR = \frac{FP}{FP+TN}$) at various threshold.
- **AUC (Area Under Curve)** Probability the classifier ranks a randomly chosen positive higher than a randomly chosen negative (AUC = 0.5: random , AUC = 1: perfect).

2 Introduction to deep learning -28/03/2025

2.1 Training an artificial neural network

Modeling problems the first example is defining the position of a car whose positions is given exactly by $d(t) = d_0 + vt$. This problem is fully solvable. That is not always the case; indeed, a problem could be characterized by a large number of variables. The problem with working with large number of variables is that not always every variable has the same importance as the other. The problem, can be reduced to: $y = \alpha x_1 + \beta x_2 + \gamma x_3$.

Machine learning allows to solve complex problems in which we can easily tell what variables are involved.

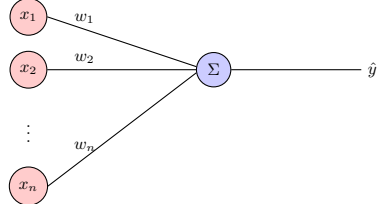
A typical problem that could be modeled is image recognition.

In ML we have a phase that does not appear in deep learning which is feature education: in traditional ML, model require pre-computed features that are manually designed based on domain knowledge. This feature engineering requires human expertise to identify the most relevant attributes for a given task.

Deep learning eliminates the need for manual feature extraction by learning hierarchical representations directly from raw data. Using AI networks, DL can automatically detect patterns at multiple levels of abstraction.

Artificial Neural Networks Biologically, neurons are unique cells that can communicate with one another, thus propagating information.

An artificial neuron receives n inputs (x_1, \dots, x_n), each scaled by a factor (weights) and sums them all:



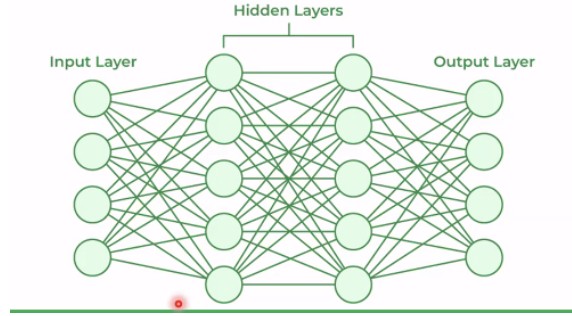
To map all the space is necessary to add a factor b called bias, such that:

$$y = \sum_{i=1}^n w_i \times x_i + b = w_1 \times x_1 + \dots + w_n \times x_n + b.$$

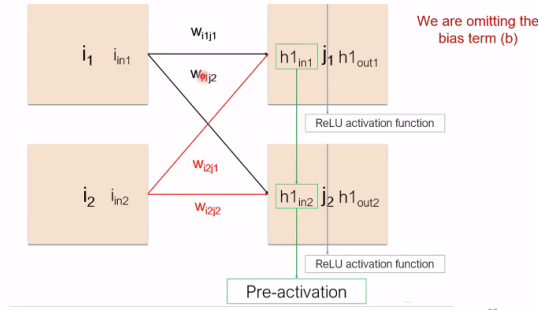
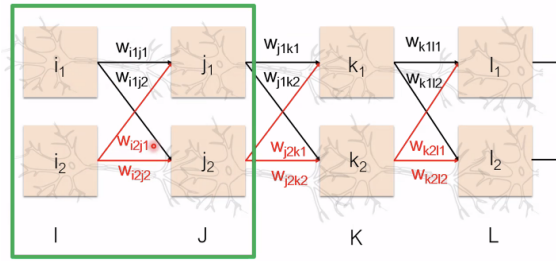
Weights and bias are the parameters of the neural network. Note that very few models are correctly modeled using linear combinations of the variables; indeed, we introduce an activate factor such that:

$$y = \sum_{i=1}^n g(w_i \times x_i + b) = g(w_1 \times x_1 + \dots + w_n \times x_n + b).$$

Typical activation functions are: sigmoid, tanh, ReLU ($\max(0, x)$), Leaky ReLU ($\max(0.1x, x)$). The sigmoid and tanh functions limit the output, indeed, the output is always confined between 0 and 1.



This is what a simple neural network looks like: an input layer, a bunch of hidden layers, and an output layer. Each layer extracts information from the input and transmits it to the next one for further data processing.



$$\begin{bmatrix} i_{in1} & i_{in2} \end{bmatrix} \times \begin{bmatrix} w_{i1j1} & w_{i1j2} \\ w_{i2j1} & w_{i2j2} \end{bmatrix} = \begin{bmatrix} h1_{in1} & h1_{in2} \end{bmatrix}$$

$$\begin{bmatrix} h1_{out1} & h1_{out2} \end{bmatrix} = ReLU(\begin{bmatrix} h1_{in1} & h1_{in2} \end{bmatrix})$$

What we are building is a neural architecture. All these parameters will have to be adjusted. This process is exactly what is called "train of the neural network". The training process continues by making mistakes.

Let's consider the problem of predicting the cost of a house: y is the real cost and \hat{y} is the predicted one. y is also called ground-truth value and is used to supervise the training. The idea under the training process lays in computing the error, and what we want indeed is to try to minimize this error. Error minimization is an optimization problem: I need a parameter configuration that explains the problem in the best possible way (i.e. i get as close as possible to the ground-truth value)

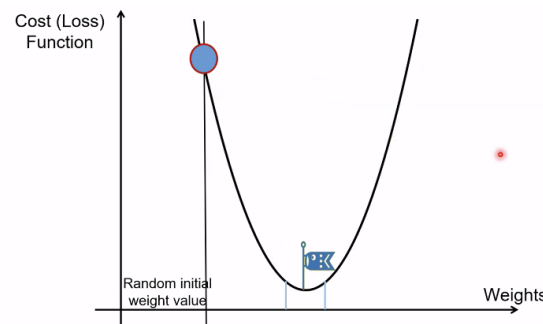
$$error = |y - \hat{y}|.$$

To optimize the parameters of a neural network we use the **gradient method**. What we always know when computing an algorithm is whether the error is increasing or decreasing.

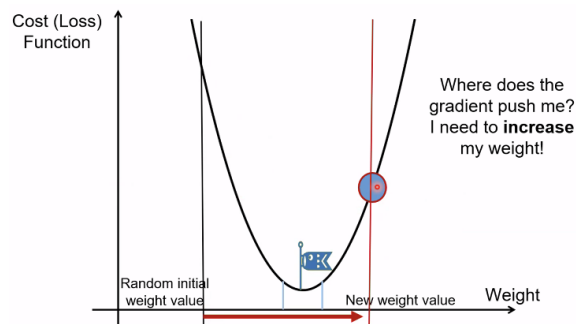
Gradient descent is the most used optimization technique to minimize the error. It means that, for each step, I compute the error and change the parameters in order to minimize it according to its dependence on each parameter.

$$error = f(w, b)$$

this function is called **cost function**. First order derivative (gradient) is a measure of dependence between each independent variable (w and b) and the dependent one (error).

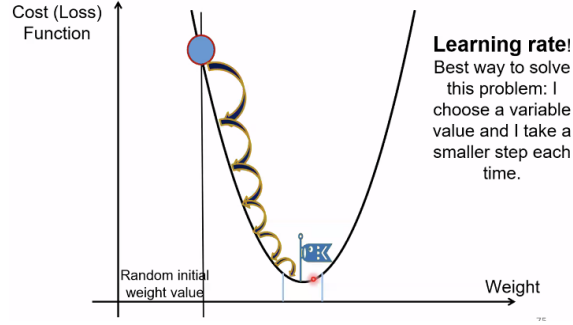


Initializing the parameters randomly gives better results.



and so on till we reach the minimum.

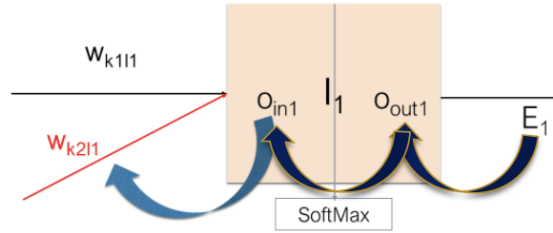
In general the first step to reach the minimum is the bigger among all the steps that will be done. In this case we talk about learning rate. While "learning" we have to decrease the learning rate (it should take smaller steps) in order to reach the target.



Mathematically, training a neural network means to change its parameters N times in order to minimize the cost function (error). Each update can be expressed as

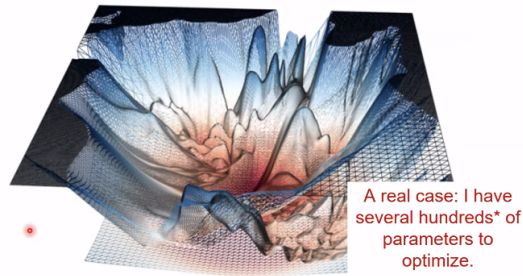
$$w_{t+1} = w_t - \eta \times \nabla error$$

where η is the learning rate and $\nabla error$ is the gradient of the cost function with respect to w . A typical loss function is the MSE. The only constraint that we impose to the loss functions is to be differentiable everywhere. The error can be computed only from the output of the layers:



$$\frac{dE_1}{dw_{k1l1}} = \frac{dE_1}{dO_{out1}} \times \frac{dO_{out1}}{dO_{in1}} \times \frac{dO_{in1}}{dw_{ikl1}}$$

For a forward pass of information we have an error back-propagation. This kind of neural network is called fully connected.

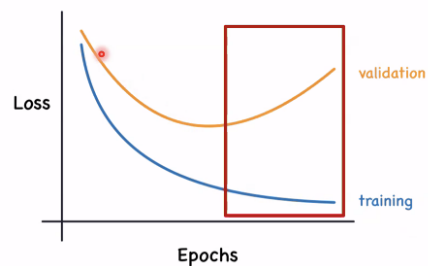


Usually, what we need is a dataset, i.e. a set of (x,y) samples where x is the input and y is the output. The dataset is divided into three sub-sets: a training set, a validation set and a test set. The training process involves parameters optimization on the same training samples several times, or epochs.

AI does not really "learns": it just finds patterns. It's like a student who only studies math by doing every exercise on the book but never reads a page of theory. How can we make sure that the student is learning and not just memorizing? During the test, the student has no way to learn new knowledge because it has no supervision (correct answers). It will be up to the teacher to evaluate the results from the test and assess the student's knowledge. Therefore, our neural network needs to perform well both on the training set and on the two test sets. What happens if it doesn't? The loss function on the training set has a very nice trend: as the network is iteratively optimized on the training sample, its error on these samples decreases. However, the loss function on the validation test doesn't seem as good: while the student improves with those "training" exercises, they make a lots of mistakes during the exam.

Two things are never to be expected:

- no error or 100% accuracy
- better results in validation



However, we should minimize the distance between training and validation performance: **over-fitting**.

The easiest task for which ML and DL are used are regression and classification: regression is used to predict a number from a continuous set of numbers (e.g. prediction of the price of a house) while classification is used to predict a number from a discrete set of numbers.

The number of neurons in the final layer will depend on the task: for the regression case we need to predict one number \rightarrow one neuron is needed, while in case of classification we need to label the inputs as one of the possible N classes $\rightarrow N$ neurons are needed. The optimal number of neurons in the hidden layers (as well as the number of the hidden layers) cannot be assessed a priori. The more the neurons and the layers, the more abstract information we can extract from the input data.

2.2 Artificial neural networks: tasks

Check out the following site: [playground tensorflow](#)

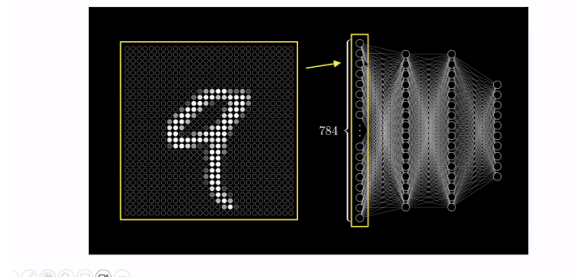
2.2.1 Binary classification

Classification is the most studied problem in DL. To classify means to assign a label to the input data among a closed set of labels.

The **MNIST** dataset is one of the most famous one: our network needs to classify the input image as one of the 10 possible labels (the digits 0-9).

How do we feed these images to the neural network?

Pixel by pixel: MNIST images are $28 \times 28 = 784$ pixels.



Pixels are integer values between 0 (black) and 255 (white) that are processed by neural network as numbers.

Check out: The stilwell brain

In the lecture of today we are going to classify whether a digit (lower than 10) is odd or even. Therefore, ours is a binary classification task. For classification, a perfect loss function is (binary) cross-entropy.

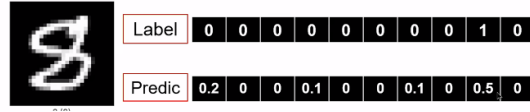


Figura 1: Note: 0.6 instead of 0.5

I want my prediction to be a probability distribution. The sum of all my elements needs to be 1, so that I can "compare" it with the label (also a probability distribution). In the case of a odd vs even prediction the binary classification will give more precise results.

Loss function: Binary cross entropy

$$-\sum_{j=1}^M y_j \log(p(y_j))$$

$$-\sum_{i=1}^N y_i \log(p(y_i)) + (1 - y_i) \log(1 - p(y_i))$$

Final activation: softmax

In order to obtain a probability distribution, I need to use the softmax activation function in the outer layer.

$$\begin{bmatrix} 1.3 \\ 5.1 \\ 2.2 \\ 0.7 \\ 1.1 \end{bmatrix} = \left[\frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \right] = \begin{bmatrix} 0.02 \\ 0.9 \\ 0.05 \\ 0.01 \\ 0.02 \end{bmatrix}$$

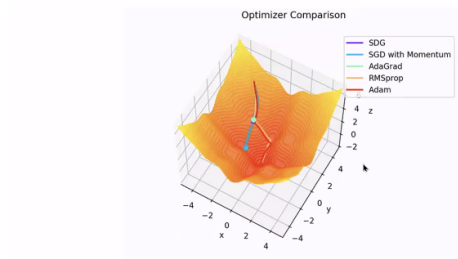
Metric: accuracy

$$\begin{bmatrix} \text{True positive (TP)} & \text{False Positive (FP)} \\ \text{False Negative (FN)} & \text{True Negative (TN)} \end{bmatrix}$$

where

- Recall = $\frac{\sum TP}{\sum TP + FN}$
- Precision = $\frac{\sum TP}{\sum TP + FP}$
- Accuracy = $\frac{\sum TP + TN}{\sum TP + FP + FN + TN}$

Optimizer: Adam (Adaptive moment estimator)



3 Advanced topics -29/03/2025

3.1 Convolutions from scratch

Edge detection is an old but gold problem in computer vision that involves detecting edges in an image to determinate object boundaries and thus separate the object of interest.

One of the most popular techniques for edge detection is the Conny Edge Detection algorithm.

An edge is a point of rapid change of intensity of the image function. The gradient points in the direction of the most rapid increase




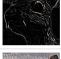
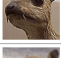
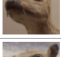
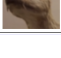
$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix}.$$

Using filters (aka matrices) is possible defining the edges in an image. An example is given by:

$$\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$$

A **convolution** (of images) is simply an elementary multiplication of two matrices followed by a sum:

- take two matrices (which both have the same dimension)
- multiply them, element by element
- add up the elements

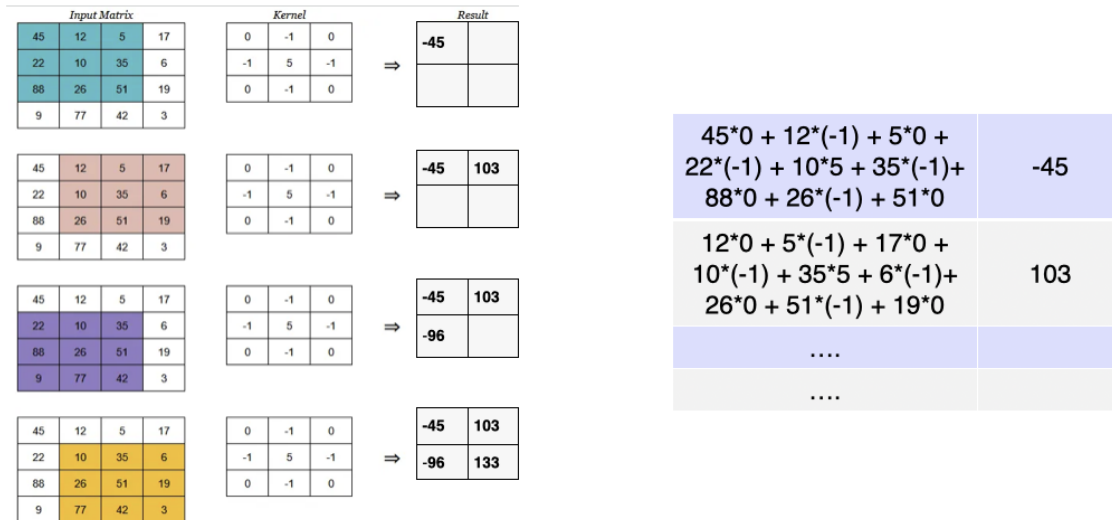
Originale	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
Edge-Detect	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	
Sharpen	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	
Blur	$\begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$	
	$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	

An image is just a multidimensional matrix, but unlike traditional matrices, images (RGB) can also have a depth. The kernel should be thought as a small matrix that is used for blurring, sharpening,

edge detection, and other image processing and functions. It is common to define the kernel by hand to achieve various image processing functions, edge detection: all these operations are hand-defined forms of kernels designed specifically to perform a particular function. The question the arises: is there a way to automatically learn these type of filters? And even use these filters for image classification and object detection? Of course there is: **CNN**.

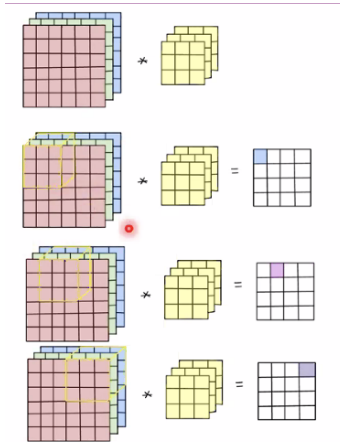
Kernel Most of the kernels we usually see are $N \times N$ square matrices. We use an odd kernel dimension to ensure that there is a valid integer coordinate in the center of the image. In image processing, a convolution requires three components: an input image, a kernel matrix to apply to the input image and an output image to store the results of the input image convolved with the kernel.

The process of "sliding" a convolutional kernel over an image and storing the output decreases the spatial dimensions of the output data. This decrease in spatial dimension is simply a side effect of applying convolutions to images.



However, in most cases, we want the output image to be the same size as the input image. To ensure this, we apply padding. The fact that the output is smaller than the input does not seem to be a big problem: we did not lose much data because most of the important features are located in the central area of the input. The only case when losing this information is a real problem is when much information is concentrated on border of the image. Padding could be done using zero elements or just copies of the border.

For an image we have:



An RGB image is represented as a $6 \times 6 \times 3$ volume, where the 3 correspond to the 3 color channels (RGB). To detect edges or other features in this image, one could convolve the $6 \times 6 \times 3$ with a 3-D filter. Then also the filter itself will have 3 levels corresponding to the red, green and blue channels. So, the filter also has a height, a width and a number of channels. The number of channels in the image must match the number of channels in the filter.

In convolutional neural networks, convolutional layers are not only applied to input data, such as pixel values, but can also be applied to the output layers. The sequence of convolutional layers allows a hierarchical breakdown of the input. Consider that filters operating directly on the raw pixel values will learn to extract low level features from the starting image, such as lines. Filters operating on the output of the first convolutional layer can extract features that are combinations of lower level features, such as features that comprise multiple lines to extract shapes. This process continues until very deep layers extract faces, animals, houses and so on. The abstraction of characteristics to ever higher orders increases with network depth.

Kernel size The kernel size defines the convolution field of view.

Padding The padding defines how the edge of a sample is handled. A convolution with padding will keep the spatial dimensions of the output equal to those of the input, while convolutions without padding will crop some of the edges if kernel is larger than 1.

Stride The stride defines the size of the kernel step when passing through the image. Although the default setting is usually 1, you can use a stride of 2 to downsample an image (this is used to improve the features of the network).

3.1.1 Types of convolutions

Dilated/Atrous convolution The atrous convolutions introduce another parameter called the rate of expansion. This parameter defines the distance between values of a kernel. This way you get

a wider field of view at the same computational cost. That type of convolution is used to understand large partial context.

Spatial separable convolutions A separable spatial convolution simply splits one kernel into smaller kernels. With fewer multiplications, the computational complexity decreases (the network has fewer parameters to learn).

Depth-wise separable convolutions We have an RGB input image (with 3 channels). After convolutions, a feature map can have more channels (as usually happens). Each channel can be as a particular interpretation of the image. Similarly to spatial separable convolution, a deep separable convolution divides a kernel into two separate kernels that perform two convolutions: the **deep convolution** and the **point convolution**.

Width multiplier allows scaling the network's width to control the number of parameters and computation.

Resolution multiplier enables adjusting the input image resolution to further reduce computational requirements.

3.2 Convolutional neural networks for image classification

The visual cortex has hierarchical structure: LGB (lateral geniculate body) -> simple cells -> complex cells -> lower order hypercomplex cells -> higher order hypercomplex cells

We can think that these complex cells are performing an aggregation of activations using functions such as the maximum, the sum. In this way, these cells can recognize edges and orientations.

Deep neural networks are normally organized in alternate repetitions of linear and non-linear operators. The reason for having multiple layers of this type is to build a hierarchical representation of the data.

The local pixels assemble to perform simple patterns like oriented edges. These borders are in turn combined to form patterns that are even more abstract. We can continue to build above these hierarchical representations until we arrive at the objects that we observe in the real world.

This compositional and hierarchical nature that we observe in the natural world is therefore not only the result of our visual perception, but it is something true in the physical world. At the lowest level of description, we have elementary particles, which are composed to form atoms, more atoms form molecules, and we continue to increase this process until we form materials, parts of objects and finally complete objects in the physical world.

3.3 Convolutional neural networks for classification task

3.3.1 Image classification

is the task of assigning to an input image a label belonging to pre-set set of categories.

Pro convolutional neural networks The 3D element that flows along the 3 image channels (RGB) is called convolutional kernel. The convolutional kernel slides over the image (in the case of the input layer) and extract features or features map. Unlike fully connected networks where the input was a carrier, CNNs operate on volumes (multi-channel image).

The pixels of the feature map with the same color are from the same kernel. Both types of networks learn by updating the weights which, in the case of fully connected networks, are the values of the connections whereas, in the case of convolutional neural networks, they are the values of the kernels and connections. In both types of networks, the neuron receives an input which is a combination of weighted inputs. This combination of weighted inputs represents the overall level of neuron excitation and is given as input to an activation function which produces some limited output.

Other layers in CNNs: Pooling Pooling provides a form of translation invariance: small shift in the input can lead to the same output. Statistics over neighboring features to reduce the size of the feature maps:

- separate the image into non-overlapping subimages
- select the maximum/average/... in each layer

3.3.2 Architectures

- LeNet-5 1998

Layer		Feature Map	Size	Kernel Size	Stride	Activation
Input	Image	1	32x32	-	-	-
1	Convolution	6	28x28	5x5	1	tanh
2	Average Pooling	6	14x14	2x2	2	tanh
3	Convolution	16	10x10	5x5	1	tanh
4	Average Pooling	16	5x5	2x2	2	tanh
5	Convolution	120	1x1	5x5	1	tanh
6	FC	-	84	-	-	tanh
Output	FC	-	10	-	-	softmax

- Alexnet 2012

Layer		Feature Map	Size	Kernel Size	Stride	Activation
Input	Image	1	227x227x3	-	-	-
1	Convolution	96	55x55x96	11x11	4	relu
	Max Pooling	96	27 x 27 x 96	3x3	2	relu
2	Convolution	256	27 x 27 x 256	5x5	1	relu
	Max Pooling	256	13 x 13 x 256	3x3	2	relu
3	Convolution	384	13 x 13 x 384	3x3	1	relu
4	Convolution	384	13 x 13 x 384	3x3	1	relu
5	Convolution	256	13 x 13 x 256	3x3	1	relu
	Max Pooling	256	6 x 6 x 256	3x3	2	relu
6	FC	-	9216	-	-	elu
7	FC	-	4096	-	-	relu
8	FC	-	4096	-	-	relu
Output	FC	-	1000	-	-	Softmax

In deep learning application, the ReLU activation feature is among the most popular. ImageNet is a large dataset of multimedia data annotated manually and divided into 1000 categories.

- **Googlenet inception 2014**

The inception module relies on several convolutions with reduced kernel size to drastically lower the number of parameters.

- **VGG16 2014**

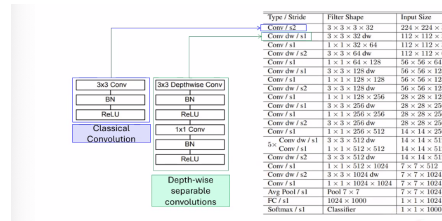
Characterized by a total of 16 layers with weights, that is, of parameters which are implemented.

- **ResNet50 2014**

Particularity: skip connections.

- **MobileNet V1**

MobileNet is designed to be a lightweight architecture optimized for mobile and embedded devices. It uses a new type of convolutional layer, known as Depthwise Separate convolution which comprises a depthwise convolution and a pointwise convolution.



Batch normalization:

$$\hat{x} = \frac{x - \mu}{\sigma},$$

where μ is the mean of x in mini-batch and σ is the std of x in mini-batch.

- **DensNet**

Densnet emphasizes strong feature reuse within the network through dense connections. Each

layer receives inputs from all preceding layers in a dense block. Dense connections include direct connections between layers within dense blocks, reducing the number of parameters and enabling better gradient flow. Growth rate controls the number of output feature maps produced by each layer within a dense block. Compound scaling: scaled the network's depth, width and resolution uniformly to find an optimal balance model size and accuracy.

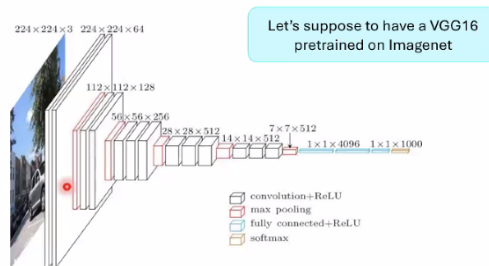
All previous CNNs have been trained on ImageNet to classify 1000 classes. This means that all CNN weights are available online. It's possible to use the pre-trained CNNs.

3.3.3 Transfer learning and fine tuning

It consists in taking the characteristics learned on a problem and exploiting them on a new similar problem. Transfer learning is usually used for tasks where the dataset has too little data to train a complete model from scratch. But how can we implement transfer learning?

All levels of a previously trained model are considered trainable. These levels are frozen so that the information contained in them is not destroyed during future training cycles. New trainable layers are added on top of the frozen layers. Pass the new data set into the "new" network and record the output of one (or more) levels from the base model (this operation is called feature extraction). This result is used as input data for a new "smaller" model to be trained. The training of this model is called **fine-tuning**.

How many layers of the original model do I freeze? It depends on the task.



I can try to freeze all layers up to the last convolutional layer and train a mini-network composed of: convolutional block and 3 fully connected blocks

Transfer learning is about "transferring" the representation learned during training of a CNN to another problem. For example, one can use pretrained CNN features to initialize the weights of a new CNN, developed for a different task. Fine tuning is about making fine adjustments. For example, during transfer learning, you can unfreeze some of (or all) the pre-trained CNN layers and let it adapt more to the task at hand.

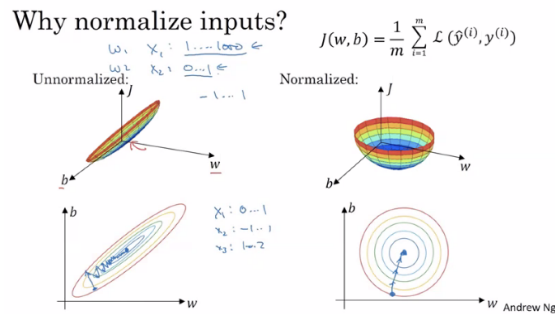
3.4 Data preparation

With the hold-out cross validation, i have:

- 70% training set (development of the model)
- 20% validation set (tuning and selection of the model)
- 10% testing set (reporting of results)

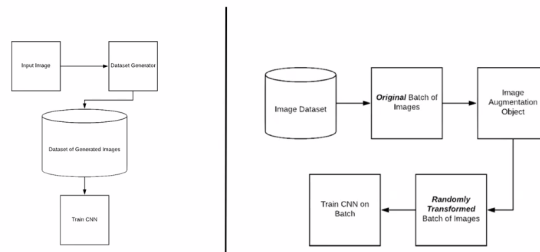
This is often done, unless my data are small size. I can't build statistics on my results, if I have few testing samples.

The data needs to be normalized in order to obtain better results.



3.4.1 Data augmentation

The aim is to increase artificially the dimension of the training set applying a series of transformations.



This transformations could be rotations, flips or a change of the brightness of the image itself. Be careful to consistency in the application of data augmentation techniques. I have to develop an algorithm for monitoring, for example, the movement of preterm infants. Suppose that the position

and orientation of the room are fixed in relation to the child. A data augmentation that would be non sense to apply would be a rotation of the images.

Note that some bias of the original adding set persists.

3.5 Optimizers

Tensor A tensor is an N-dimensional array of data (non è vero : _ ().

SDG is an optimization algorithm used to train machine learning algorithms. The algorithm's task is to find a set of parameters within the model that minimizes the loss or error function. Optimization is a type of research process and you can think of this research as learning. The optimization algorithm is called "gradient descent", where "gradient" refers to the calculation of an error gradient or slope of the error and "descent" refers to the movement along that slope towards a minimum level of error. The algorithm is iterative: **use the error to update the internal parameters of the model (back propagation)**. A sample is individual data to which, in the case of supervised learning, a label is attached. This is used to compare the forecast and calculate an error. The batch size is a hyperparameter that defines the number of samples to be analyzed before updating the internal parameters of the model. A training dataset can be split into one or more batches. When the batch has sample size, the learning algorithm is called stochastic gradient descent. When the batch size is greater than a sample and less than the size of the training data set, the learning algorithm is called mini-batch gradient descent. The latter is the most common implementation used in the field of deep learning. Mini-batch gradient descent, the most common batch size are 32, 64 and 128.

In the batch gradient descent we use all the training data in a single iteration.

A training epoch means that the learning algorithm has made a pass through the training data set, in which the examples have been separated into randomly selected batch size = 32 groups.

Momentum is a technique used to help the optimizer go faster towards the optimal solution.

3.6 Loss function

For continuous outputs

$$MSE = \frac{1}{n} \sum_{i=1}^N (t_i - p_i)^2$$

while for categorical output we'll implement accuracy (used only for evaluation) and cross entropy (used during training)

$$L_{CE} = - \sum_{i=1}^{c=2} t_i \log p_i = -t_1 \log(p_1) - (1 - t_1) \log(p_2).$$

The prediction is a probability vector: it represents the predicted probabilities of all classes with sum equal to 1. In a neural network, this prediction is usually made with the last layer activated

by softmax function. We calculate the cross-entropy loss for the image. Loss is a measure of the performance of the model. The lower it is, the better. During learning, the model aims to achieve the lowest possible loss. The target represents the probability for all classes and is one-hot encoded vector, that is has 1 in a single position and 0 in all others. We'll start by calculating the loss for each class separately and then add it up. The loss for each class is calculated as follows:

$$Loss = -p(X) \ln(q(X))$$

where the first term is the probability of class X in target while the second is the probability of class X in predictions. If the probability for the target is 0, the loss is 0.

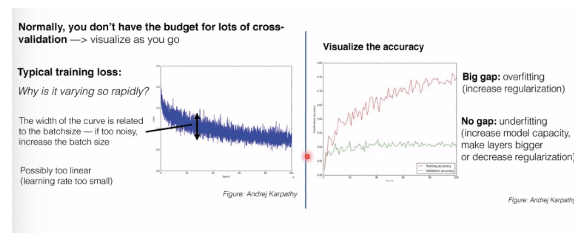
The loss is 0 if the prediction is 1. The loss tends to infinity if the prediction is 0.

Cross entropy = $\sum_X p(X) \ln q(X)$. This is the cross entropy formula which can be used as a loss function for any two probability vectors. If we want to get the loss for our branch or for the whole set of data? You add up the losses of individual images. If our goal is one-hot encoded vector, we can actually forget the targets and predictions for all other classes and only calculate the loss for the hot class. Cross-entropy loss also works for distributions that are not one-hot vectors. In summary: **General formula used to calculate the loss between two probability vectors.**

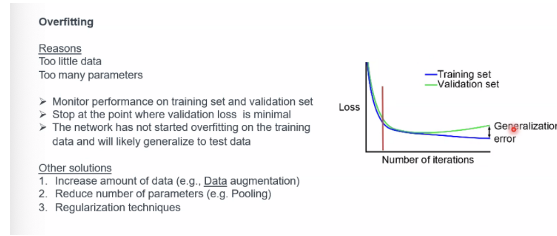
The further we get from our target the more the error grows.

Binary classification = $-[p(X) \ln q(X) + (1 - p(X)) \ln(1 - q(X))]$ we use binary cross entropy, a specific case of cross entropy where our target is 0 or 1. It can be calculated with the cross entropy formula if you convert the target into a one-hot encoded vector, such as [0,1] or [1,0] and forecasts respectively. We can calculate it with the previous formula.

Multi-label classification = $\sum_x \text{Binary cross entropy}_x$ Our target can represent multiple classes (or even zero) at the same time. We calculate the binary cross entropy for each class separately and then add it up to get the complete loss.



Loss function- learning curves



3.7 Metrics

The **F1 score** is a popular metric for evaluating the performance of a classification model. In the case of multi-class classification, for the calculation, for example, of the F1 score, averaging methods are used, which translate into a series of different averages (macro, weighted, micro) in the rating report. In order to assess the performance of a model in a comprehensive way, we should examine both recall and accuracy. The F1 is a useful metric that considers both.

Instead of having multiple F1 scores per class, it would be better to average for a single number that describes overall performance.

The weighted average F1 score takes into account the support of each class. The weight refers to the proportion of the support for each class in relation to the sum of all the support values.

Micro average F1 score calculates the global mean F1 score by counting the sums of the true positives, false negatives and false positives:

$$\text{Micro-average F1 score} = \frac{TP}{TP + 0.5(FP + FN)}.$$

The micro-average F1 score calculates the percentage of correctly ranked observations on all observations.