Sapienza University of Rome

Master in Artificial Intelligence and Robotics

# Machine Learning

A.Y. 2024/2025

Prof. Luca Iocchi

# 13. Multiple learners

Luca Iocchi

with contributions from Valsamis Ntouskos

# Overview

- Combining multiple learners
- Voting
- Bagging
- Boosting
- AdaBoost

Reference

E. Alpaydin. Introduction to Machine Learning. Chapter 17.

C. Bishop. Pattern Recognition and Machine Learning. Chapter 14.

# Multiple learners / Ensemble learning

two options ↗ training in parallel → indipendent each other
↘ training in sequential

General idea: instead of training a complex learner/model, train many different learners/models and then combine their results.

**Committees**: set of models trained on a dataset.

Models can be trained in parallel (*voting* or *bagging*) or in sequence (*boosting*).

# Voting

*we use FULL DATASET*

Given a dataset $D$ ↓

1. use $D$ to train a set of models $y_m(x)$, for $m = 1, \ldots, M$

   *IF WE DON'T WANT USE WEIGHTS WE CAN USE $\frac{1}{M}$*

2. make predictions with

   *must normalize to 1* ↑

$$y_{voting}(x) = \sum_{m=1}^{M} w_m y_m(x) \qquad \text{(regression)}$$

$$y_{voting}(x) = \operatorname*{argmax}_{c} \sum_{m=1}^{M} w_m I(y_m(x) = c) \qquad \begin{array}{l}\text{weighted majority}\\\text{(classification)}\end{array}$$
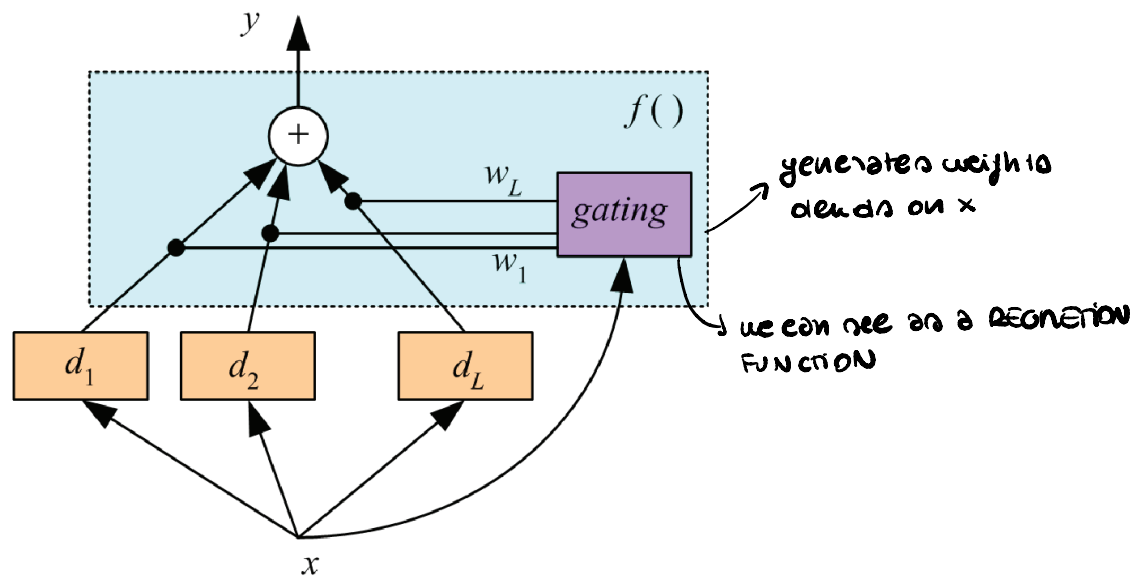
with $w_m \geq 0$, $\sum_m w_m = 1$ (prior probability of each model),
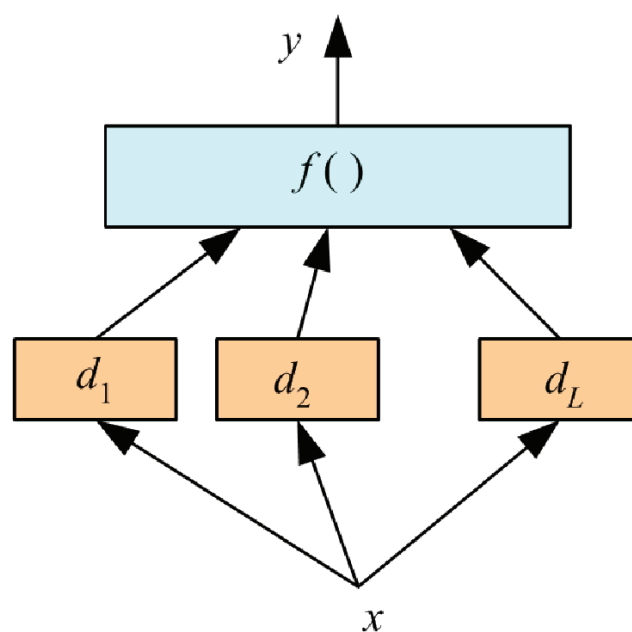$I(e) = 1$ if $e$ is true, 0 otherwise.

# Voting



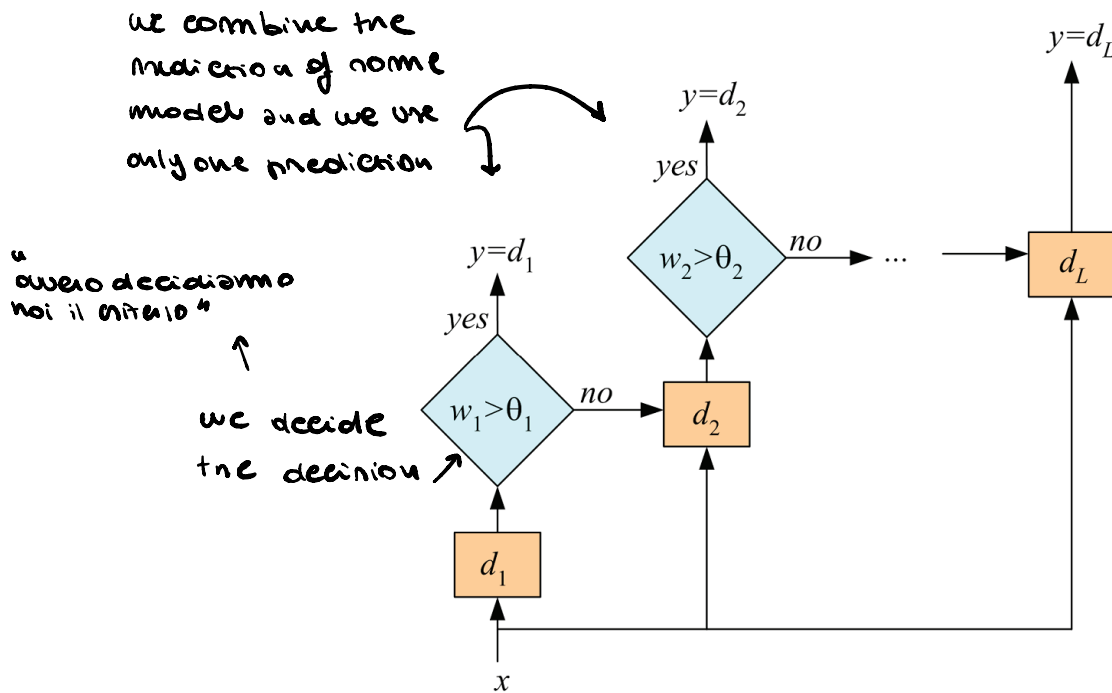*IN THIS APPROACH WEIGHTS ARE COSTANT*

# Mixture of experts



Non linear gating function $f$ depending on input

# Stacking



Combination function $f$ is also learned

# Cascading

*we combine the prediction of some models and we use only one prediction*

*"ovvero decidiamo noi il criterio"*

*we decide the decision*



$$y = d_L$$

$$y = d_2$$

yes

$$y = d_1$$

yes

no

$$w_2 > \theta_2$$

no

...

$$d_L$$

$$w_1 > \theta_1$$

no

$$d_2$$

$$d_1$$

$$x$$

Cascade learners based on confidence thresholds

*one problem could be __overfitting__*

# Bagging

*PARALLEL METHOD*

*L'INTERSEZIONE È VUOTA*

*PARTITION → $D_1 \cap D_2 = \emptyset$*

*!! IS NOT PARTITION !!*

*randomly sample generated by uniform distribution*

Given a dataset $D$,

1. generate $M$ bootstrap data sets $D_1, \ldots, D_M$, with $D_i \subset D$

2. use each bootstrap data set $D_m$ to train a model $y_m(x)$, for $m = 1, \ldots, M$

*THIS REDUCE OVERFITTING BECAUSE doesn't know all dataset to the model*

3. make predictions with a voting scheme

$$y_{bagging}(x) = \frac{1}{M} \sum_{m=1}^{M} y_m(x)$$

In general, this is better than training any individual model.

Bootstrap data sets chosen with *random sampling with replacement*

WE TRAIN MODEL SEQUEN-
TIALY

## Boosting: general approach
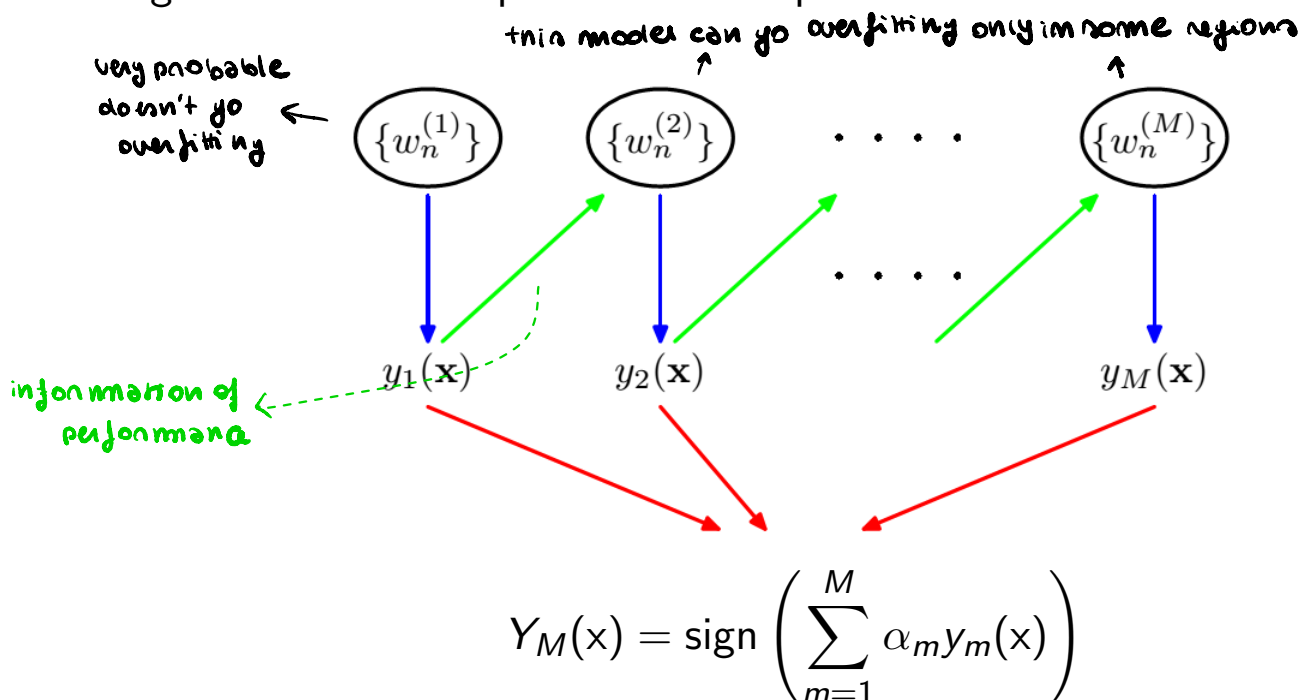
in general in better
respect the other

↳ the second model try to
improve performance in region of    ⟹ same samples are more
input of previous model                           important than other

↳ "avero il meccanico cerca di migliorare
il modello precedente con i pesi
del modello precedente"

Main points:

- Base classifiers (*weak learners*) trained sequentially
- Each classifier trained on weighted data
- Weights depend on performance of previous classifiers
- Points misclassified by previous classifiers are given greater weight
- Predictions based on weighted majority of votes

## Boosting: general approach

Base classifiers are trained in sequence using a weighted data set where
weights are based on performance of previous classifiers.

this model can go overfitting only in some regions

very probable
doesn't go
overfitting

information of
performance



$$Y_M(x) = \text{sign}\left(\sum_{m=1}^{M} \alpha_m y_m(x)\right)$$

# AdaBoost

Given $D = \{(x_1, t_1), \ldots, (x_N, t_N)\}$, where $x_n \in X$, $t_n \in \{-1, +1\}$

1. Initialize $w_n^{(1)} = 1/N$, $n = 1, ..., N$.

2. For $m = 1, ..., M$:

   - Train a weak learner $y_m(x)$ by minimizing the weighted error function:

   STANDARD ERROR-FUNCTION → FOR BINARY CLASSIFICATION

   $$J_m = \sum_{n=1}^{N} w_n^{(m)} I(y_m(x_n) \neq t_n), \text{ with } I(e) = \begin{cases} 1 & \text{if } e \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

   - Evaluate: $\epsilon_m = \dfrac{\sum_{n=1}^{N} w_n^{(m)} I(y_m(x_n) \neq t_n)}{\sum_{n=1}^{N} w_n^{(m)}}$ and $\alpha_m = \ln\left[\dfrac{1 - \epsilon_m}{\epsilon_m}\right]$

   USE FOR PREDICTION OF MODEL
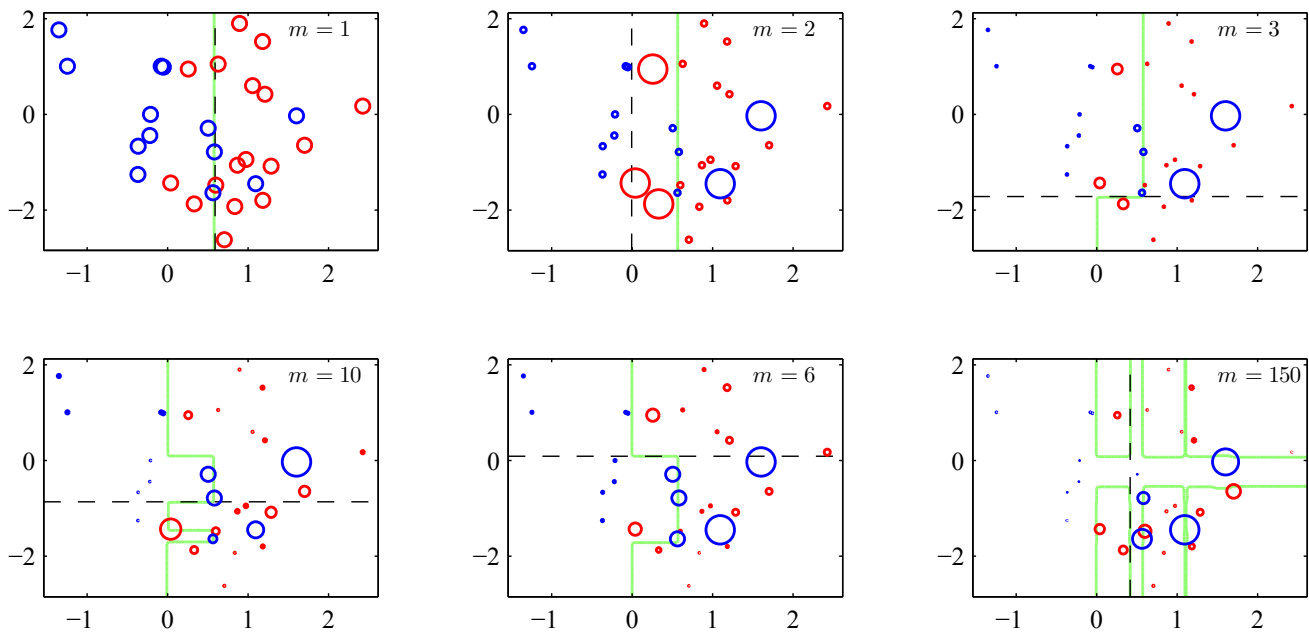
   - Update the data weighting coefficients:

   $$w_n^{(m+1)} = w_n^{(m)} \exp[\alpha_m I(y_m(x_n) \neq t_n)]$$

# AdaBoost

3. Output the final classifier

$$Y_M(x) = \text{sign}\left(\sum_{m=1}^{M} \alpha_m y_m(x)\right)$$

# AdaBoost

# Exponential error minimization

AdaBoost can be explained as the sequential minimization of an exponential error function.

Consider the error function

ADABOOST MINIMIZE
THIS ERROR FUNCTION

$$E = \sum_{n=1}^{N} \exp[-t_n f_M(x_n)],$$

where

$$f_M(x) = \frac{1}{2} \sum_{m=1}^{M} \alpha_m y_m(x), \quad t_n \in \{-1, +1\}$$

**Goal**:

$$\text{minimize } E \text{ w.r.t. } \alpha_m, y_m(x), \ m = 1, \dots, M$$

# Exponential error minimization

**Sequential minimization.** Instead of minimizing $E$ globally

- assume $y_1(x), \ldots, y_{M-1}(x)$ and $\alpha_1, \ldots, \alpha_{M-1}$ fixed;
- minimize w.r.t. $y_M(x)$ and $\alpha_M$.

Making $y_M(x)$ and $\alpha_M$ explicit we have:

$$E = \sum_{n=1}^{N} \exp\left[-t_n f_{M-1}(x_n) - \frac{1}{2} t_n \alpha_M y_M(x_n)\right]$$

$$= \sum_{n=1}^{N} w_n^{(M)} \exp\left[-\frac{1}{2} t_n \alpha_M y_M(x_n)\right],$$

with $w_n^{(M)} = \exp[-t_n f_{M-1}(x_n)]$ constant as we are optimizing w.r.t. $\alpha_M$ and $y_M(x)$.

# Exponential error minimization

From sequential minimization of $E$, we obtain

$$w_n^{(m+1)} = w_n^{(m)} \exp[\alpha_m I(y_m(x_n) \neq t_n)] \text{ and } \alpha_m = \ln\left[\frac{1 - \epsilon_m}{\epsilon_m}\right]$$

predictions are made with

$$\text{sign}(f_M(x)) = \text{sign}\left(\frac{1}{2} \sum_{m=1}^{M} \alpha_m y_m(x)\right)$$

which is equivalent to

$$Y_M(x) = \text{sign}\left(\sum_{m=1}^{M} \alpha_m y_m(x)\right)$$

thus proving that AdaBoost minimizes such error function.

# AdaBoost Remarks

Advantages:

- fast, simple and easy to program
- no prior knowledge about base learner is required
- no parameters to tune (except for $M$)
- can be combined with any method for finding base learners
- theoretical guarantees given sufficient data and base learners with moderate accuracy

Issues:

- Performance depends on data and the base learners
  (can fail with insufficient data or when base learners are too weak)
- Sensitive to noise

# Summary

RANDOM FOREST → are simple a set of DECISION TREES

- Instead of designing a learning algorithm that is accurate over the entire space one can focus on finding base learning algorithms that only need to be better than random
- Combined learners theoretically outperforms any individual learner
- AdaBoost practically outperforms many other base learners in many problems
- Ensembles of small DNN outperform very deep NN in some cases