# 13 Multiple Learners

Given a problem there is no a best way to solve it. The general idea is to insted of training a complex model/learner, TRAIN MANY MODELS/LEARNERS and combine their results.

You take all the models, you consider all of them. Whenever a new input appears you give it to the models and you combine the various predictions.

There are two different families of approaches:

① PARALLEL (voting or bagging)
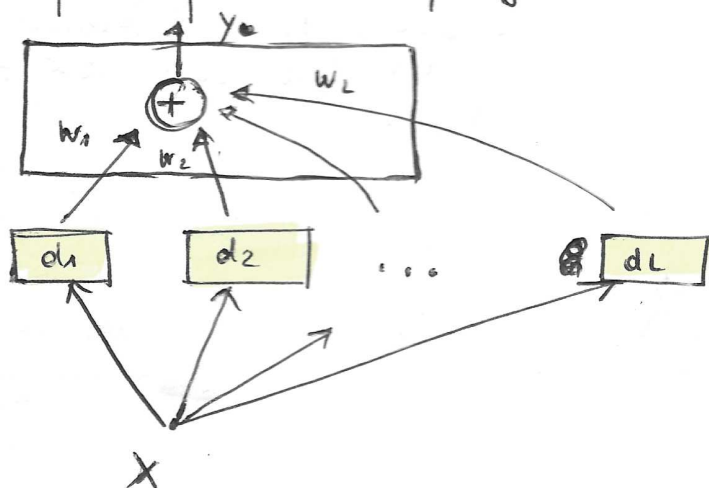
② SEQUENTIAL (boosting)

- ## VOTING

Given a dataset $D$, you use the same dataset to train different models. Then the prediction can be done by considering A WEIGHTED AVG OF THE PREDICTIONS of EACH MODEL.

(REGR.) $$y_{voting}(x) = \sum_{M=1}^{M} w_M y_M(x)$$

$M$ = # classifier/ learners/ models.

(CLASS.) $$y_{voting}(x) = \underset{c}{argmax} \sum_{M=1}^{M} w_m \delta\left(y_M(x) = c\right)$$

$w_M$ = prior probability of models, $\sum w_m = 1$



The scheme. In the voting scheme the weights are fixed, but we have different approach where the weights might depend on the input like in MIXTURE SCHEME.
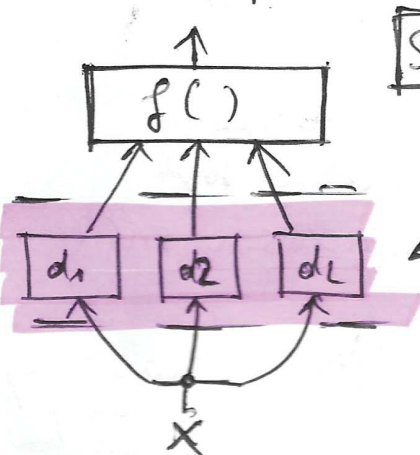
# MIXTURE OF EXPERTS



$f()$

The weights depend on the input, for each input you have different $w_i$ according to the gating function.

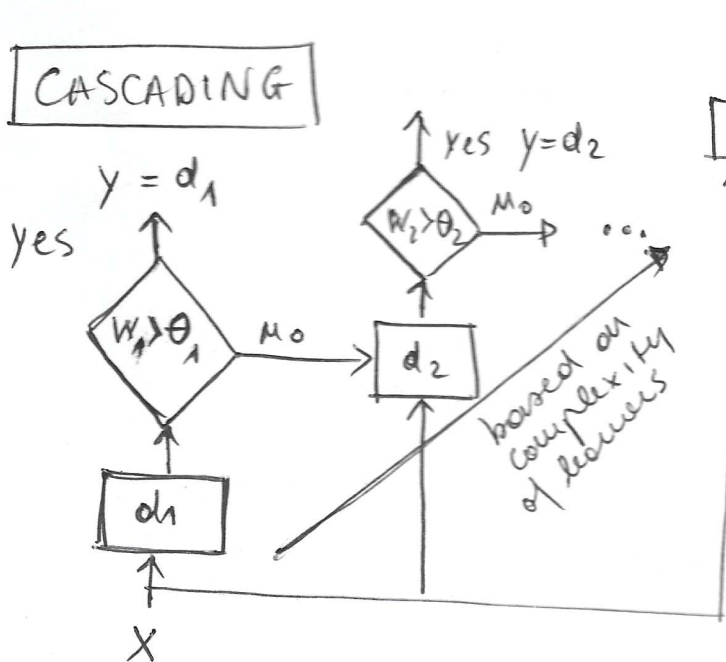This is useful when a subspace can be represented with a particular method.

In general you have that the output of the multiple leans are combined in a function whose parameters are learned during the process ( seems/reminds a neural network ).

# STACKING



in a given layer we don't have omogeneous units like in NNs, but eterogeneous ~~xxxxx~~ units.

# CASCADING



$y = d_1$

yes

$W_1 > \theta_1$ No $\rightarrow$ $d_2$

yes $y = d_2$

$W_2 > \theta_2$ No $\rightarrow$ ...

based on complexity of learners

$y = d_L$

$d_L$

The output in this approach is used in a different way. In all before approaches, you query all learners, IN CASCADING YOU QUERY ONE LEARNER AT TIME.

The result can come just from one model: THIS A SELECTION SCHEME. It is based on threshold, confidence.

## BAGGING

You can generate subsets of the dataset and train each model with a different subset. I take my dataset and I generate a set of sub-set of the dataset ( You can even use bagging with same learners). In general, this is better than training any individual model.
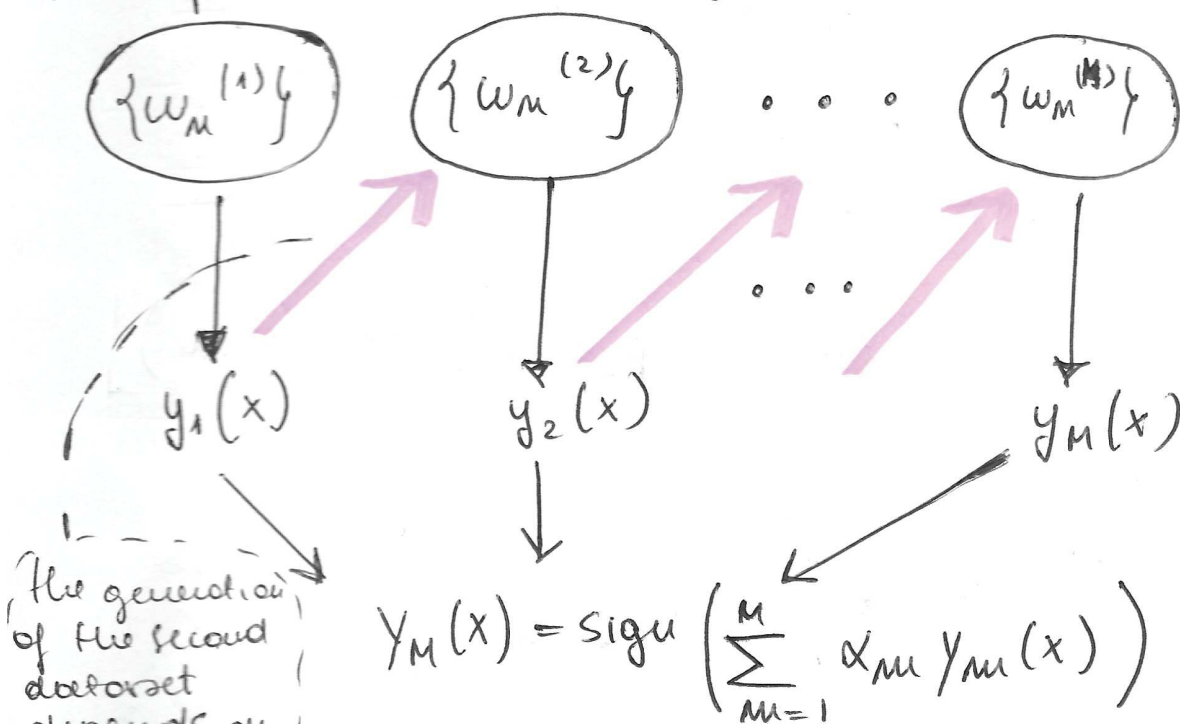
$$Y_{bagging}(x) = \frac{1}{M} \sum_{m=1}^{M} Y_m(x)$$ ← do the average

To generate Bootstrap data sets, you can USE RANDOM SAMPLING WITH REPLACEMENT.

## BOOSTING a sequential approach | More efficient, better!

In a sequential approach, each model is ~~trained~~ trained on a dataset that depends on the training of previous iteration.

$\{w_m^{(1)}\}$    $\{w_m^{(2)}\}$   . . . .   $\{w_m^{(M)}\}$

$y_1(x)$     $y_2(x)$    . . .    $y_M(x)$

_The generation of the second dataset depends on the other training_

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

Base classifiers are trained in sequence using a weighted data set when weights are based on the performance of previous classifiers. WE TRY TO SPECIALIZE LEARNERS, you give more importance (in the next phase training) to misclassified data by previous learner.

# Ada Boost Algorithm

Given $D = \{ (x_i, t_i)_{i=1}^{N} \}$, where $x_M \in X$, $t_M \in \{+1, -1\}$

1. Initialize $w_M^{(1)} = 1/N$, $M = 1, \ldots, N$.

2. For $M = 1, \ldots, M$

   Train a weak learner $Y_{M}(x)$ by minimizing the weighted error function:
   
   → (any learner better than random (accuracy ≥ 0.51))

   $$J_M = \sum_{M=1}^{N} W_M^{(M)} I \left( Y_{M}(x_M) \neq t_M \right)$$

   WEIGHTED ERROR FUNCTION

   where $I(e) = \begin{cases} 1 & \text{if } e \text{ true} \\ 0 & \text{otherwise} \end{cases}$

   Evaluate! 
   
   RESIDUAL ERROR
   
   how good has been the classifier
   
   $$\varepsilon_M = \frac{\sum_{M=1}^{N} W_M^{(M)} I \left( Y_{M}(x_M) \neq t_M \right)}{\sum_{M=1}^{N} W_M^{(M)}}$$

   and $\alpha_M = \ln \left[ \frac{1 - \varepsilon_M}{\varepsilon_M} \right]$

   update $w_M^{(M+1)} = w_M^{(M)} \exp \left[ \alpha_M I \left( Y_{M}(x_M) \neq t_M \right) \right]$

3. Output! the final classifier

   $$Y_M(\bar{x}') = \text{sign} \left( \sum_{M=1}^{M} \alpha_M Y_M(x) \right)$$

Ada Boost IS OPTIMIZING AN EXPONENTIAL ERROR FUNCTION. Consider the error function!

$$E = \sum_{M=1}^{N} \exp \left[ -t_M f_M(x_M) \right]$$

Where

$$f_M(x) = \frac{1}{2} \sum_{i=1}^{M} \alpha_M Y_M(x), \quad t_M \in \{-1, +1\}$$

GOAL: minimize $E$ wrt $\alpha_m, y_m(x)$

$$\forall \ m = 1, \ldots, M.$$

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 wrt $\alpha_M$ and $y_M(x)$

> The idea is to consider at each step, some parameters fixed

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

$$E = \sum_{i=1}^{N} \exp\left[-t_m f_{M-1}(x_m) - \tfrac{1}{2} t_m \alpha_M y_M(x_m)\right]$$

$$= \sum_{m=1}^{N} w_m^{(M)} \exp\left[-\tfrac{1}{2} t_m \alpha_M y_M(x_m)\right]$$

with $w_m^{(M)} = \exp\left[-t_m f_{(M-1)}(x_m)\right]$ constant since we are optimizing wrt $\alpha_M$ and $y_M(x)$.

From sequential minimization of $E$ we obtain:

$$w_m^{(m+1)} = w_m^{(m)} \exp\left[\alpha_m \, \mathbb{I}\left(y_m(x_m) \neq t_m\right)\right]$$

predictions are made with:

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

which is equivalent to:

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

WITH THE ASSUMPTION OF OPTIMIZING WRT $\alpha_M$ and $y_M$ WE SIMPLIFY A LOT THE PROBLEM.

(+) FAST, simple and easy to program.

(+) No requirements to learners (best than random)

(+) No parameter to tune, except for M. A theoretical result says that performance improve with an higher M

(+) practical evidence to be good (better performance wrt to individual learners).

(−) Performance depends on data and the base learners (can fail with insufficient data or when base learners are too weak)

(−) Sensitive to noise.