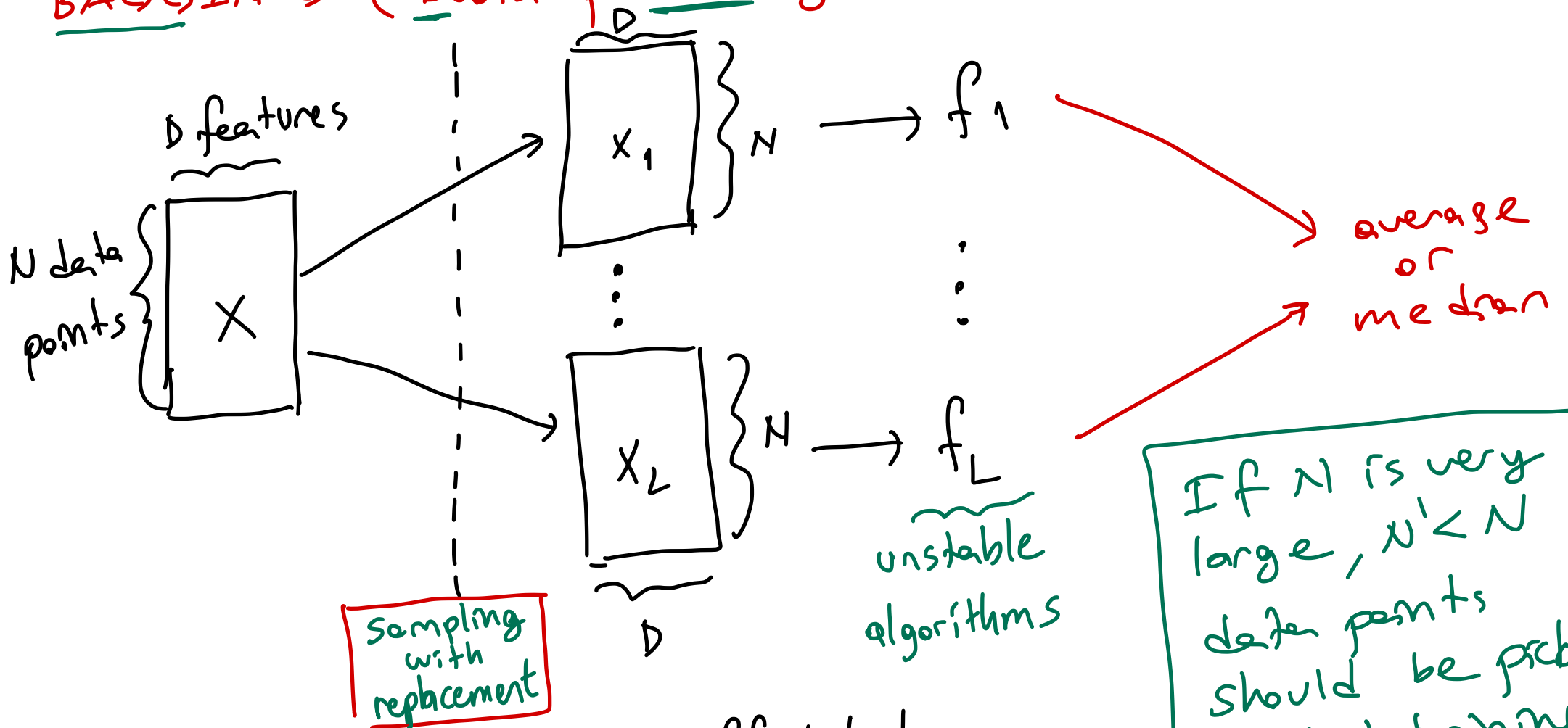# BAGGING (Bootstrap AGGregation)
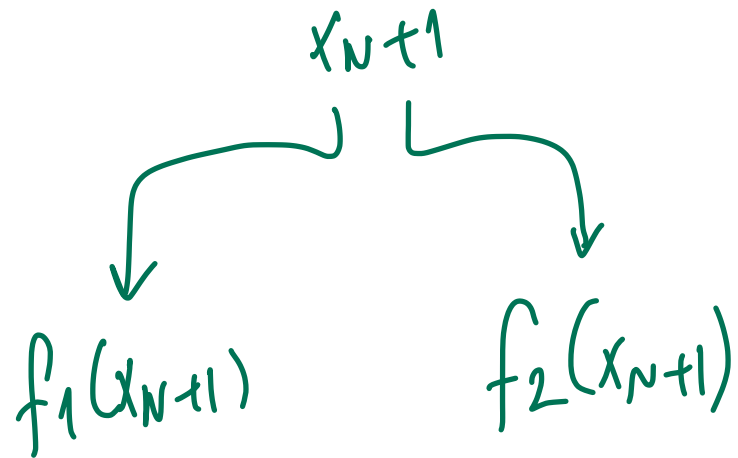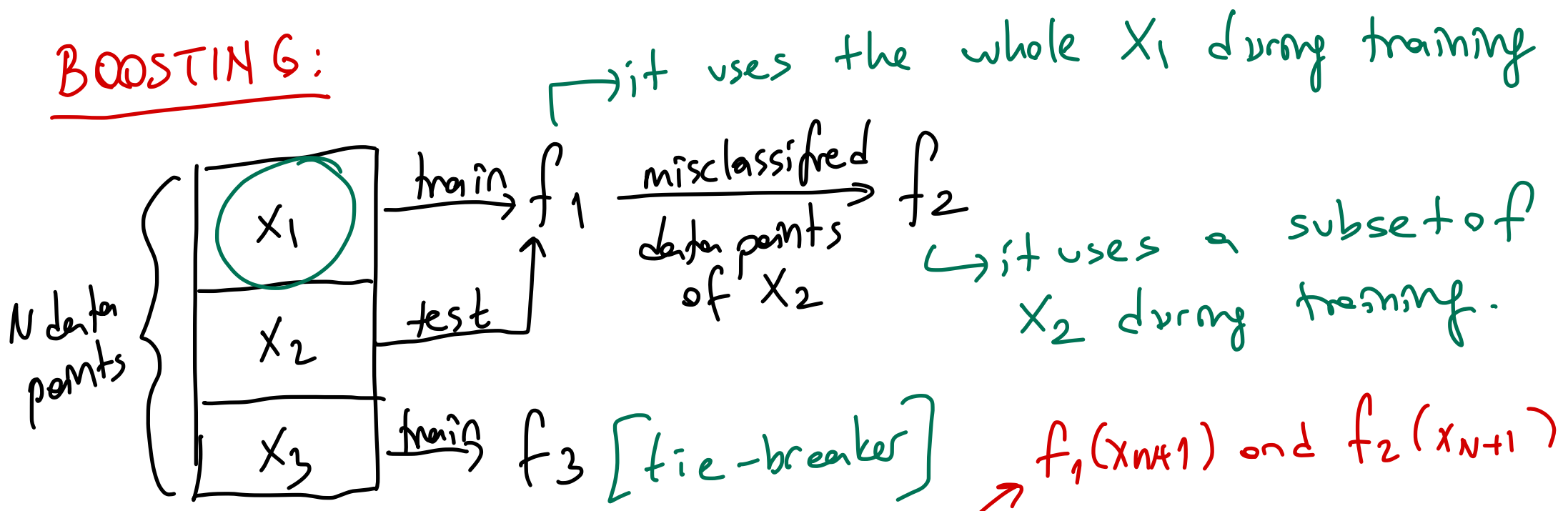


Unstable Algorithm: Highly affected by small changes in the training data set.

Unstable $\Rightarrow$ DT

Stable $\Rightarrow$ k-NN

# BOOSTING:

→ it uses the whole $X_1$ during training



$X_1$ —train→ $f_1$ —$\dfrac{\text{misclassified}}{\text{data points of } X_2}$→ $f_2$

↳ it uses a subset of $X_2$ during training.

$X_2$ —test→

$X_3$ —train→ $f_3$ [tie-breaker]

N data points

$x_{N+1}$

$f_1(x_{N+1})$ and $f_2(x_{N+1})$

$f_1(x_{N+1})$    $f_2(x_{N+1})$

① If they agree on their decisions, no problem and use the predicted class label.

② If they do not agree, use $f_3(x_{N+1})$ as the predicted class label.

**AdaBoost**: modify the probabilities of drawing instances as a function of the error.

$P_{ij}$ = the probability that the data point $x_i$ is selected (used in training) by classifier $f_j$
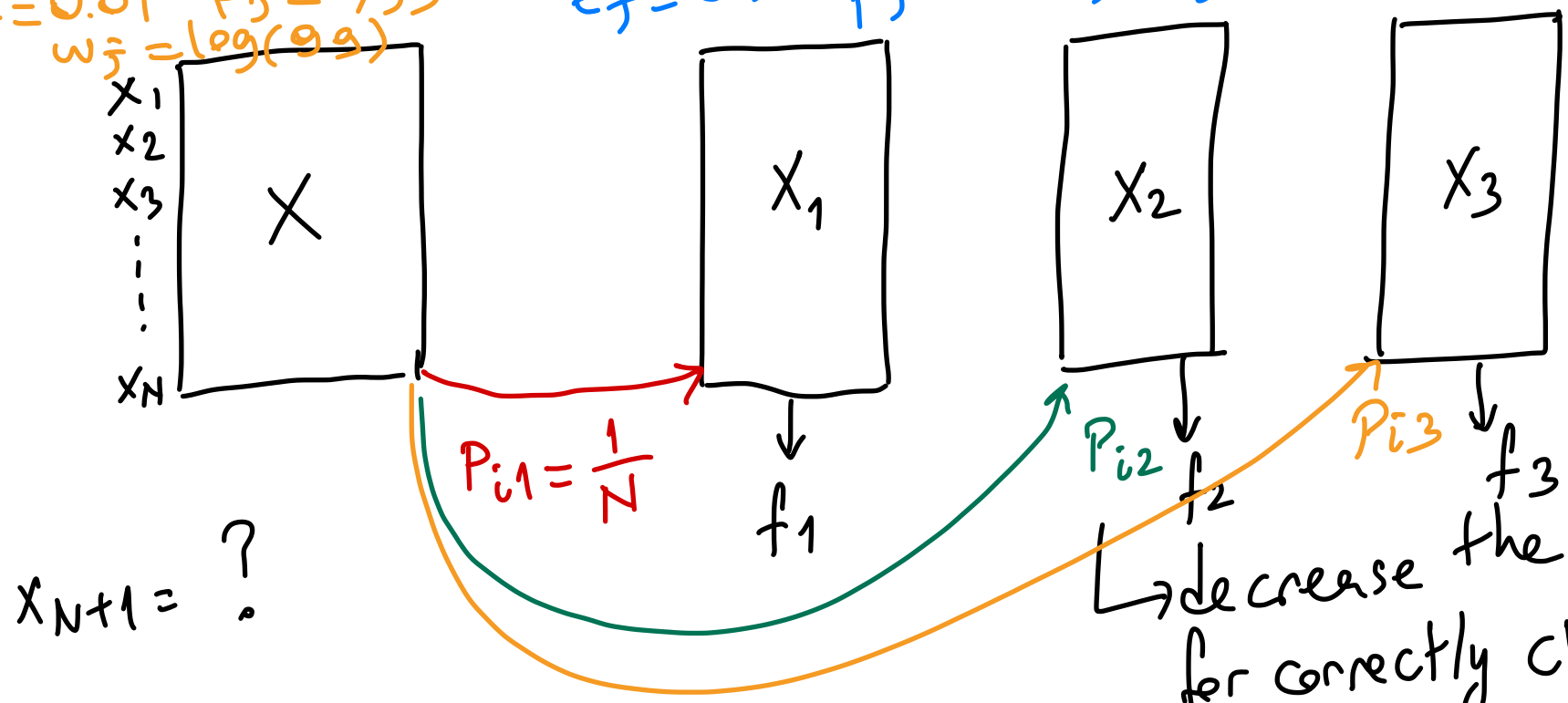
error rate. $\epsilon_j = 0.2$

$\beta_j = \dfrac{\epsilon_j}{1-\epsilon_j} = \dfrac{0.2}{0.8}$

$\epsilon_j = 0.01 \quad \beta_j = 1/99$
$w_j = \log(99)$

$\epsilon_j = 0.5 \quad \beta_j = 1 \quad w_j = \log(1) = 0$

$w_j = \log\left[\dfrac{1}{\beta_j}\right] = \log(4)$

$x_1$
$x_2$
$x_3$
$\vdots$
$x_N$

X

$X_1$

$X_2$

$X_3$

$P_{i1} = \dfrac{1}{N}$

$P_{i2}$

$P_{i3}$

$f_1$

$f_2$

$f_3$

↳ decrease the probabilities for correctly classified data points

↳ increase the probabilities for incorrectly classified data points

$x_{N+1} = ?$

$f(x_{N+1}) = \widehat{w_1} f_1(x_{N+1}) + \dots + \widehat{w_L} f_L(x_{N+1})$

based on their error rate.

# Mixture of Experts (MoE):

→ Constant over the input space.

$$Voting \Rightarrow \hat{y} = \sum_{J=1}^{L} \boxed{w_j} f_j(x_{N+1})$$

$$MoE \Rightarrow \hat{y} = \sum_{J=1}^{L} w_j(x_{N+1}) f_j(x_{N+1})$$

→ $w_j$'s will be assigned by the gating function

## Cooperative

- $w_1, w_2, \ldots, w_L$ are assumed to be independent.



$f_1(x_{N+1})$   $w_1(x_{N+1})$
$f_2(x_{N+1})$   $w_2(x_{N+1})$   $\hat{y}$
$\vdots$
$f_L(x_{N+1})$   $w_L(x_{N+1})$

$x_{N+1}$

gating function

## Competitive

- $w_1, w_2, \ldots, w_L$ are producing sperse weights (mostly zero)

- one or some of them are nonzero.

sigmoid

$$w_j = \frac{1}{1 + \exp\left[-(v_j^T x + v_{j0})\right]}$$

softmax

$$w_j = \frac{\exp(v_j^T \cdot x + v_{j0})}{\sum_{k=1}^{L} \exp(v_k^T \cdot x + v_{k0})}$$
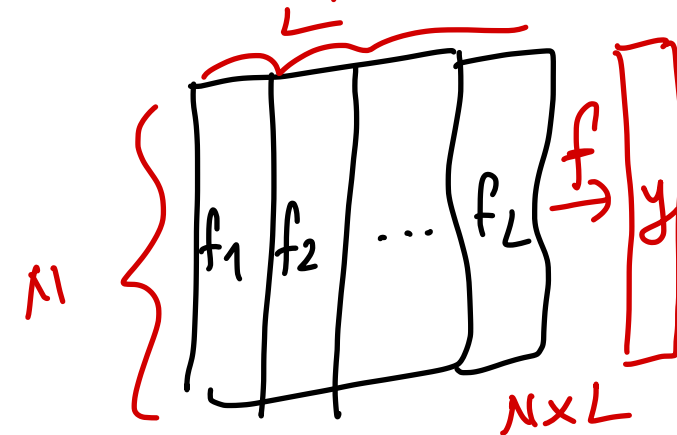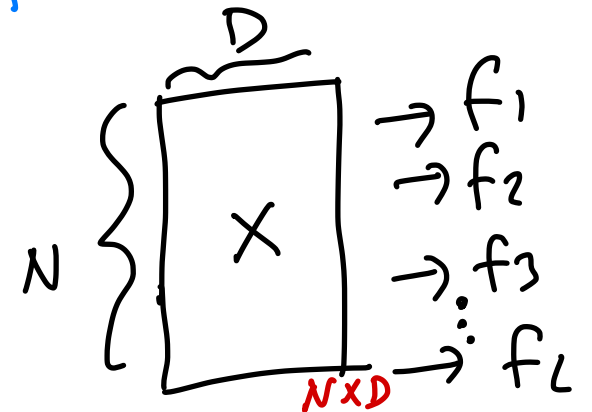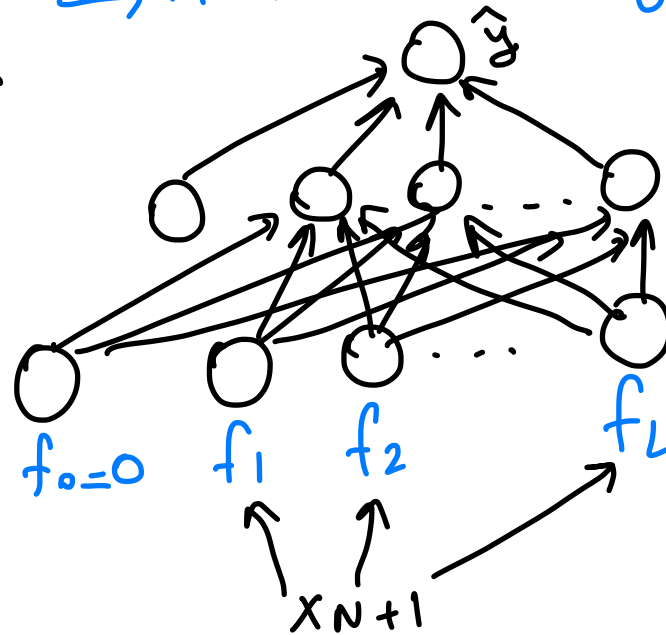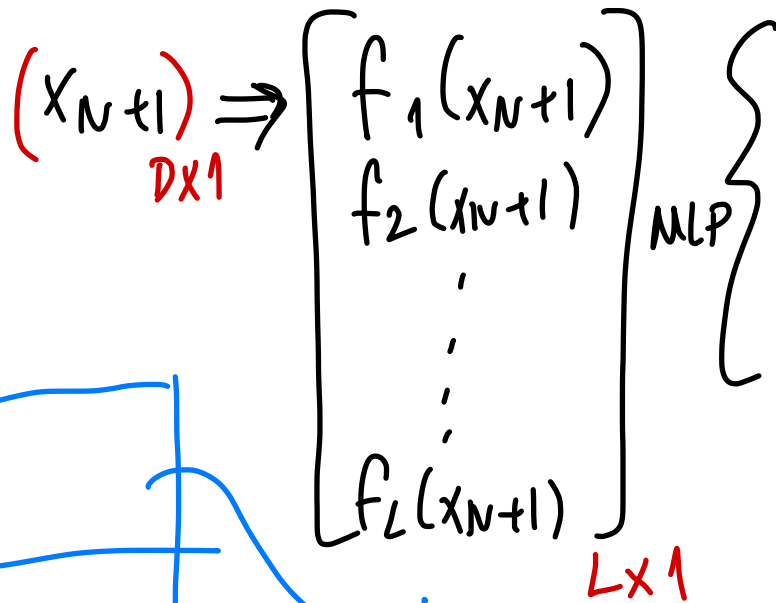
# Stacked Generalization

$$\text{Voting} \Rightarrow \hat{y} = \sum_{j=1}^{L} w_j f_j(x_{N+1})$$

$$\text{MoE} \Rightarrow \hat{y} = \sum_{j=1}^{L} w_j(x_{N+1}) f_j(x_{N+1})$$

$$\begin{array}{l}\text{Stacked} \\ \text{Generalization}\end{array} \Rightarrow \hat{y} = f\left( f_1(x_{N+1}), f_2(x_{N+1}), \dots, f_L(x_{N+1}) \right)$$

$\hookrightarrow$ nonlinear algorithm

$(x_{N+1}) \Rightarrow$ _DX1_
$$\begin{bmatrix} f_1(x_{N+1}) \\ f_2(x_{N+1}) \\ \vdots \\ f_L(x_{N+1}) \end{bmatrix} \text{MLP}$$
_LX1_



$f_0 = 0 \quad f_1 \quad f_2 \quad \dots \quad f_L$

$x_{N+1}$

$N \left\{ \begin{array}{|c|} \hline X \\ \hline \end{array} \right.$ $D$ → $f_1$ → $f_2$ → $f_3$ → $f_L$ _NXD_

_L NXD_

$N \left\{ \begin{array}{|c|c|c|c|} \hline f_1 & f_2 & \dots & f_L \\ \hline \end{array} \right. \xrightarrow{f} y$ _NXL_

→ training base learners

→ training the combination algorithm

# Cascading:

Confidence level.

$$x_{N+1} \rightarrow f_1 \rightarrow P_1 > \theta_1 \xrightarrow{\text{YES}} \text{STOP}$$

threshold

$\downarrow$ NO

$$f_2 \rightarrow P_2 > \theta_2 \xrightarrow{\text{YES}} \text{STOP}$$

$\downarrow$ NO

$$f_3 \; \bullet \; \bullet \; \bullet$$

$$\underset{0.95}{\theta_1} > \underset{0.90}{\theta_2} > \underset{0.85}{\theta_3} > \bullet \; \bullet \; \bullet \bullet$$

decreasing thresholds

$$x_{N+1} \xrightarrow{f_1} P_1 = 0.98 \implies \text{we are confident enough, let's STOP.}$$

since $0.98 > 0.95$

$$x_{N+1} \xrightarrow{f_1} P_1 = 0.92 \xrightarrow{f_2} P_2 = 0.91 \implies \text{we are confident enough,}$$

since $0.92 < 0.95$    since $0.91 > 0.90$    let's STOP.

A1: $x_1 \quad x_2 \overset{\checkmark}{-} x_3 \quad x_4 \quad x_5 \overset{\times}{-} x_6 \quad x_7 \quad x_8 \quad | x_9 \quad x_{10} \quad x_{11}$

A2: $x_1 \quad x_2 \overset{\checkmark}{-} x_3 \quad x_4 \quad | x_5 \overset{\checkmark}{-} x_6 \quad x_7 \quad x_8 \quad | x_9 \quad x_{10} \quad x_{11}$

A3: $x_1 \quad x_2 \overset{\checkmark}{-} x_3 \quad x_4 \quad x_5 \overset{\checkmark}{-} x_6 \quad | x_7 \quad x_8 \quad x_9 \quad | x_{10} \quad x_{11}$

A※

$c_{ij}$ = # of clustering algorithms that put $x_i$ and $x_j$ into the same cluster.

$C_{5,6} = 2$

$C_{2,3} = 3$

$$C = \begin{array}{c c} & \begin{array}{c c c c c c c c c c c} x_1 & x_2 & x_3 & x_4 & x_5 & x_6 & x_7 & x_8 & x_9 & x_{10} & x_{11} \end{array} \\ \begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \\ x_8 \\ x_9 \\ x_{10} \\ x_{11} \end{array} & \left[ \begin{array}{c c c c c c} 3 & 3 & 3 & 3 & 2 & 1 \\ 3 & 3 & 3 & 3 & 2 & 1 \\ 3 & 3 & 3 & 3 & 2 & 1 \\ 3 & 3 & 3 & 3 & 2 \\ 2 & 2 & 2 & 2 & 3 & 2 \\ & & & & & \\ & & & & & 3 \\ & & & & & \quad 3 \end{array} \right] \end{array}$$

11 × 11

Clustering on the C matrix would give you A※.