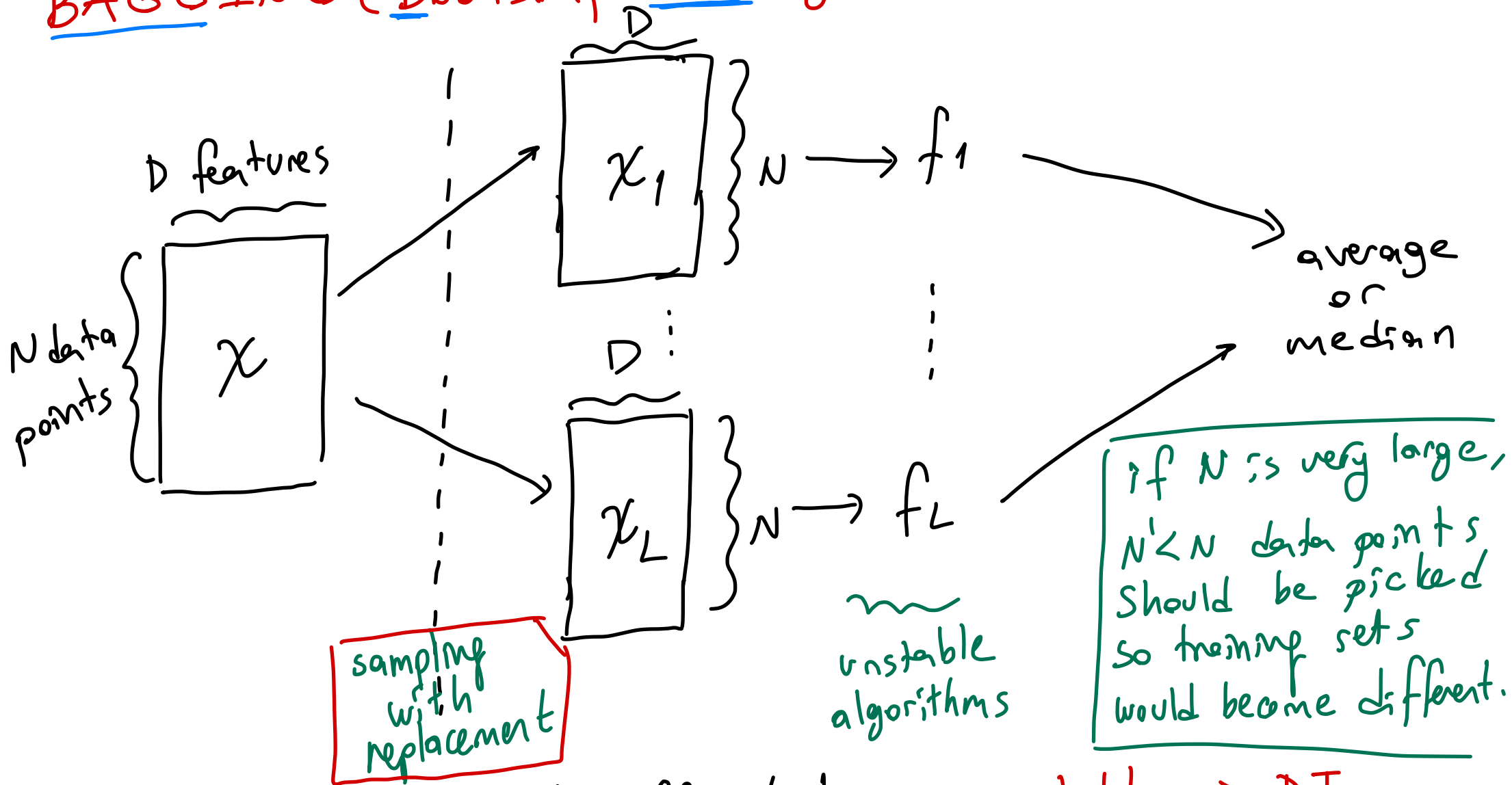


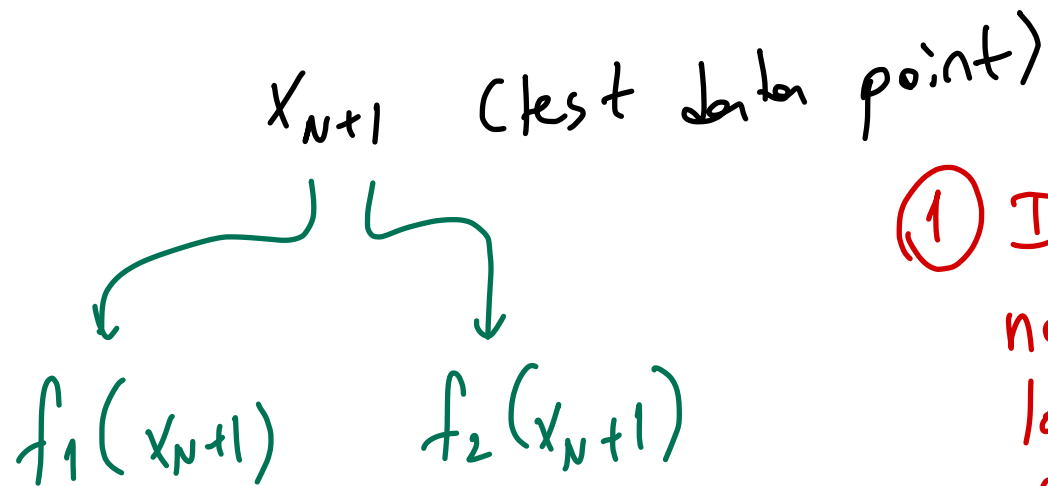
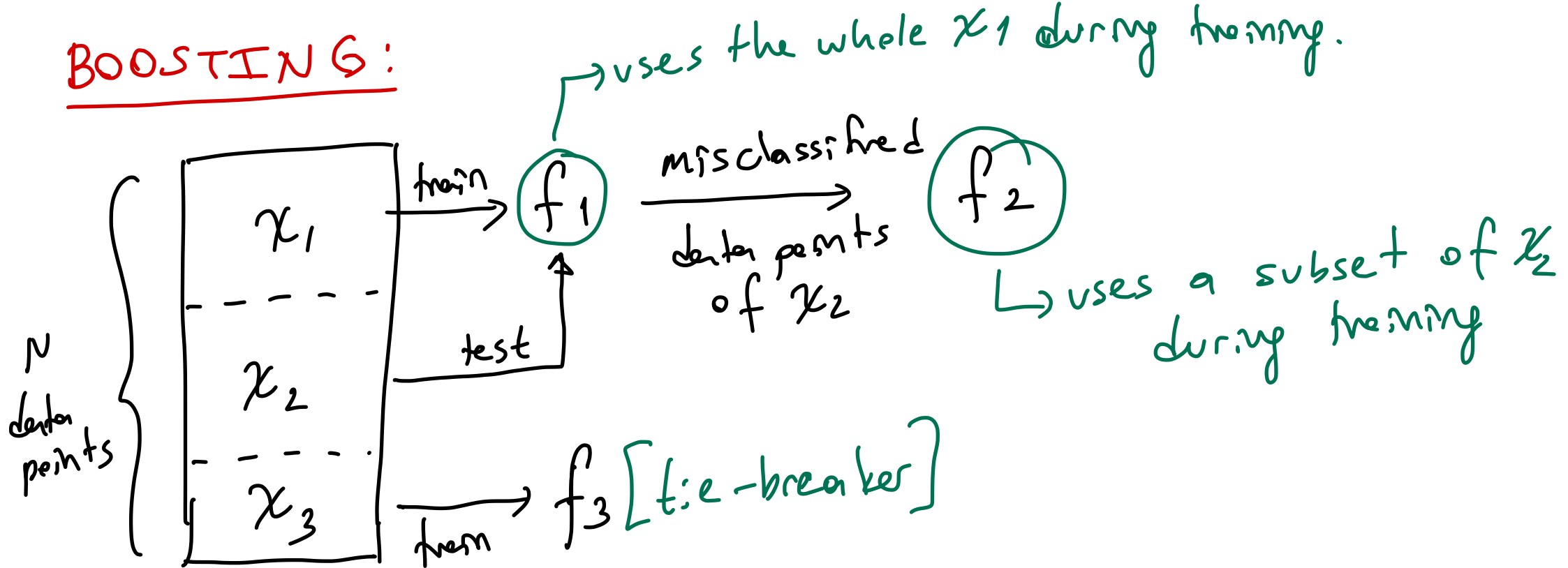
# BAGGING (Bootstrap AGGregation)



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

unstable  $\Rightarrow$  DT  
stable  $\Rightarrow$  k-NN

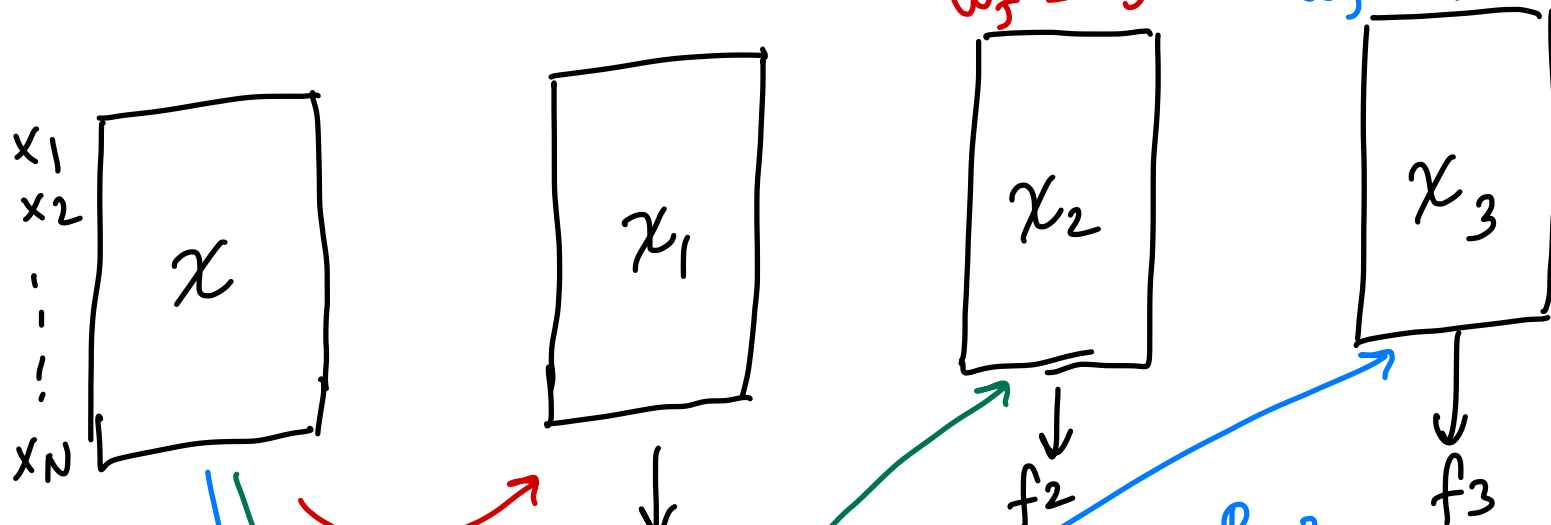
# BOOSTING:



- ① If they agree on their decisions, no problem and use the predicted label.
- ② If they do not agree, use  $f_3(x_{N+1})$  as the predicted label.

Ada Boost: modify the probabilities of drawing instances as a function of the error.

$P_{i,j}$  = the probability that the instance  $x_i$  is selected by classifier  $f_j$ .



$$\begin{aligned} \epsilon_j &= 0.20 \\ \beta_j &= 0.20/0.80 \\ w_j &= \log(4) \end{aligned}$$

$$\begin{aligned} \epsilon_j &= 0.50 \\ \beta_j &= 0.50/0.50 \\ w_j &= \log(1) = 0 \end{aligned}$$

$$\begin{aligned} \epsilon_j &= 0.01 \\ \beta_j &= 0.01/0.99 \\ w_j &= \log[99] \end{aligned}$$

$$\begin{aligned} w_j &= \log[1/\beta_j] \\ \beta_j &= \frac{\epsilon_j}{1-\epsilon_j} \\ \epsilon_j &= \text{error rate.} \end{aligned}$$

$x_{N+1} \Rightarrow ?$

decrease the probabilities for correctly classified data points  
increase the probabilities for incorrectly classified data points

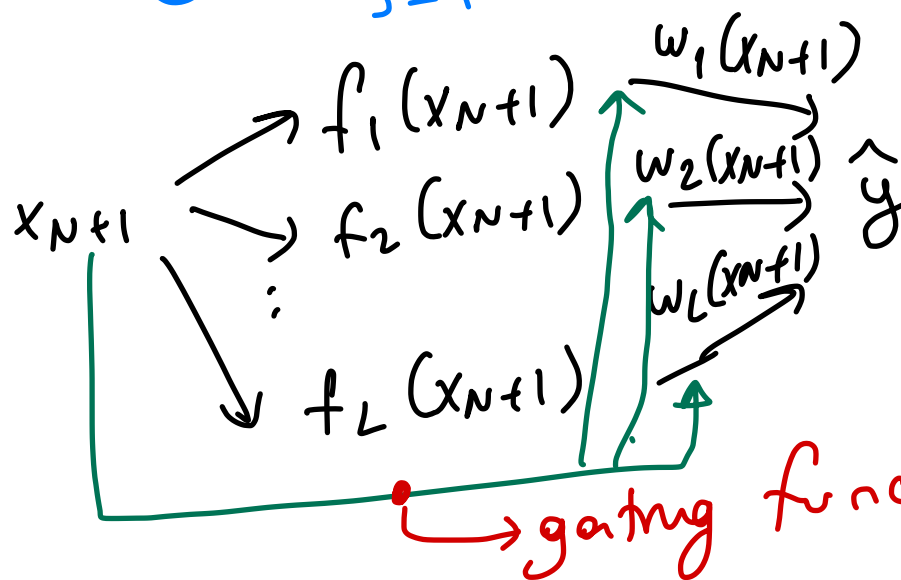
$$f(x_{N+1}) = w_1 f_1(x_{N+1}) + w_2 f_2(x_{N+1}) + \dots + w_L f_L(x_{N+1})$$

# Mixture of Experts (MoE):

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

Constant over the input space.

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



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

## Competitive:

- $\rightarrow w_1, w_2, \dots, w_L$  are usually sparse (i.e., mostly zero)
- $\rightarrow$  one or some of them are non zero.

Softmax

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

## Cooperative:

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

Sigmoid

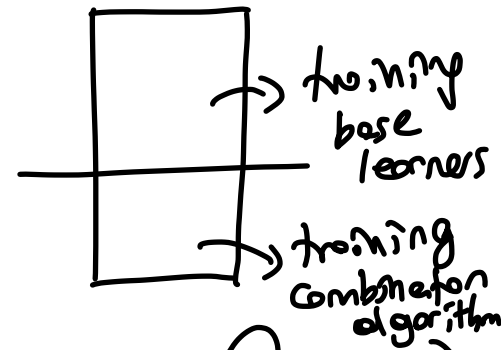
$$w_j = \frac{1}{1 + \exp[-(v_j^T \cdot x + v_{j0})]}$$

# Stacked Generalization

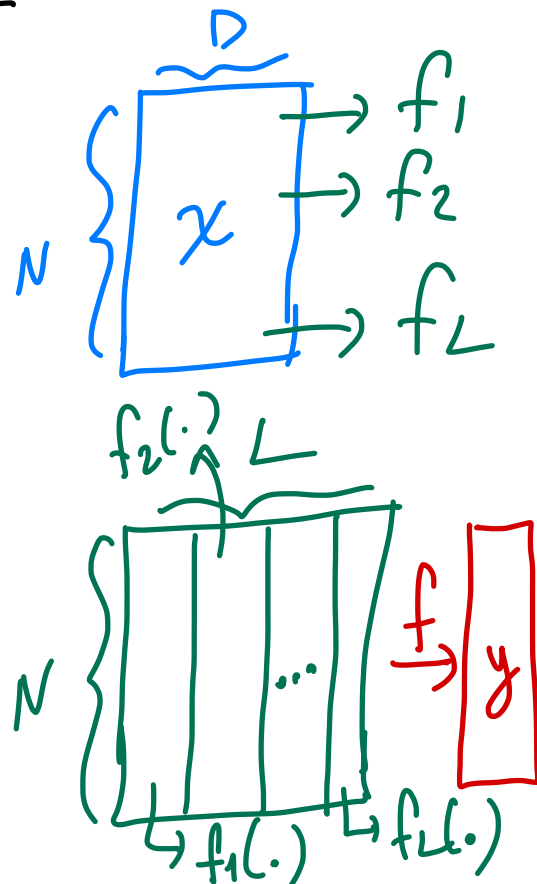
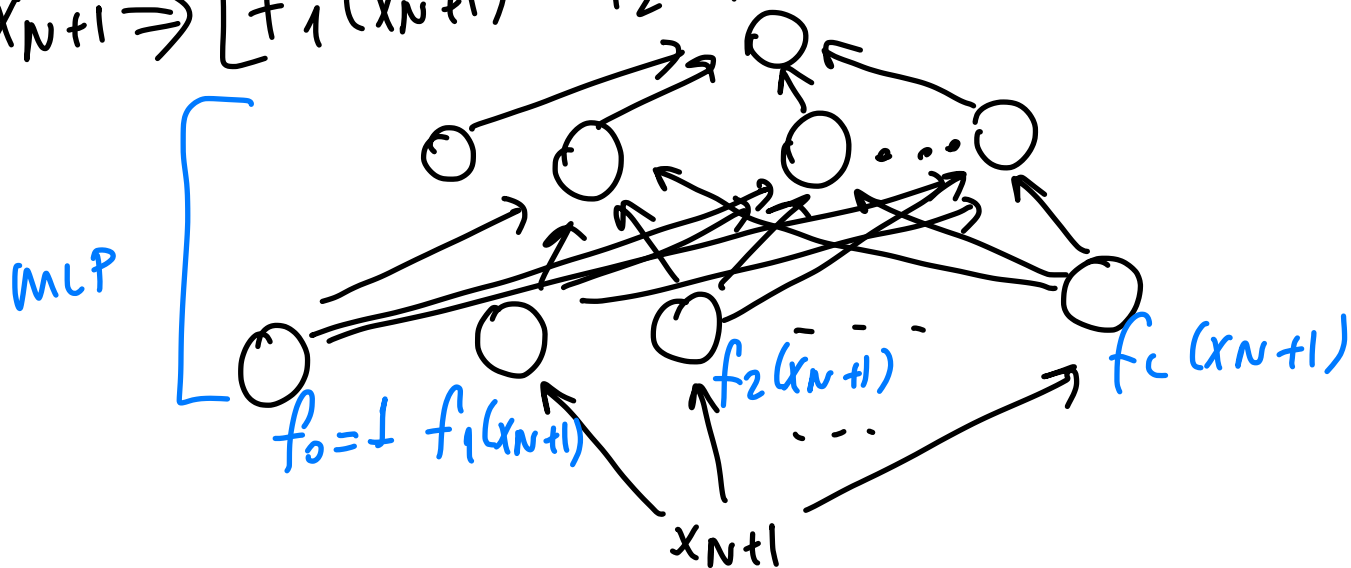
Voting  $\Rightarrow \hat{y} = \sum_{j=1}^L \underline{w_j} f_j(x_{N+1})$

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

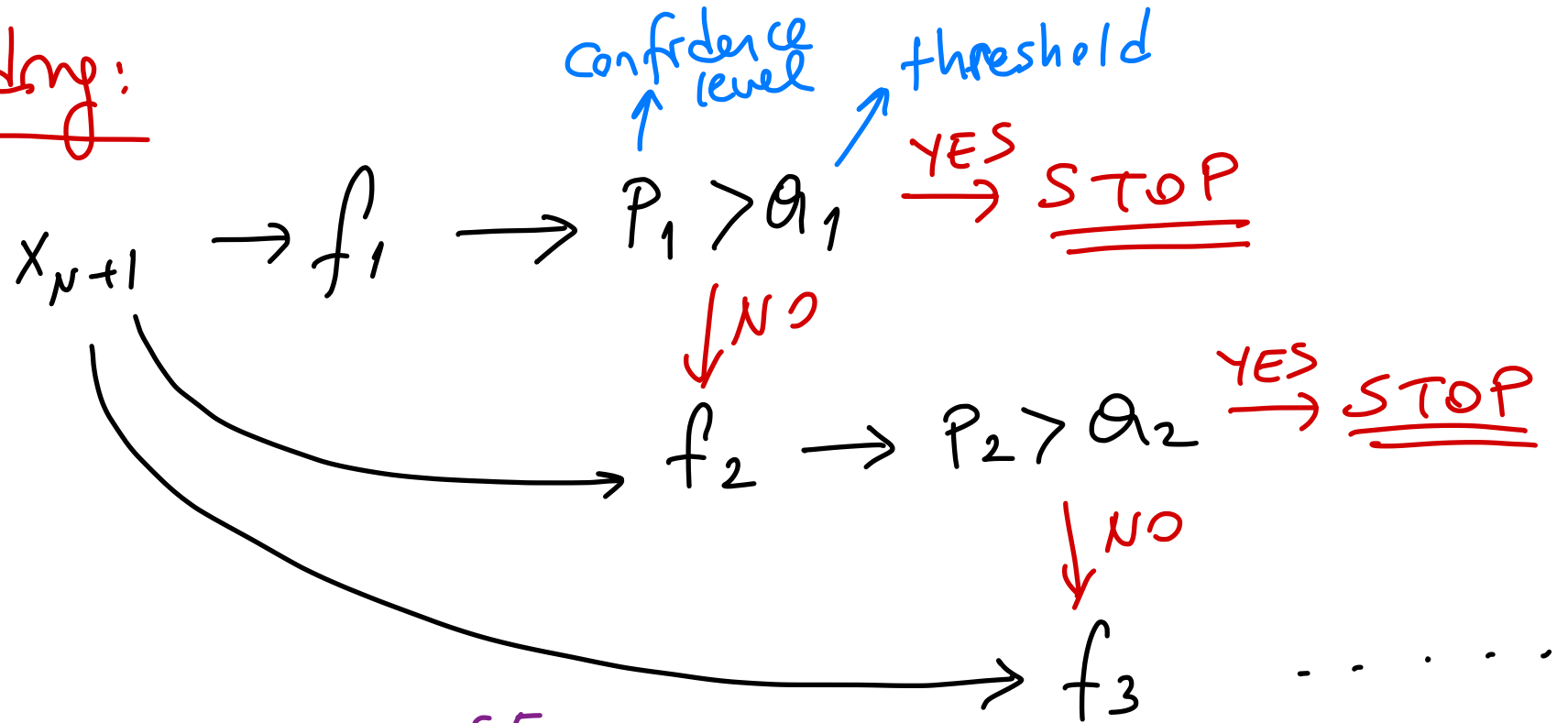
Stacked Generalization  $\Rightarrow \hat{y} = \underbrace{f}_{\text{nonlinear algorithm}}(f_1(x_{N+1}), f_2(x_{N+1}), \dots, f_L(x_{N+1}))$



$x_{N+1} \Rightarrow [f_1(x_{N+1}) \quad f_2(x_{N+1}) \quad \dots \quad f_L(x_{N+1})]^T$



Cascading:



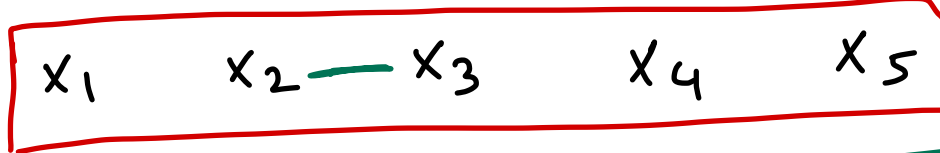
$$0.95 \quad 0.90 \quad 0.85 \\ \alpha_1 > \alpha_2 > \alpha_3 > \dots$$

decreasing thresholds

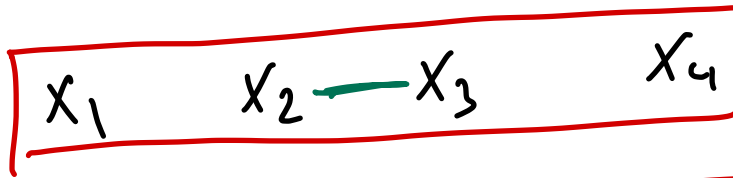
Example 1:  $x_{N+1} \xrightarrow{f_1} P_1 = 0.98 \Rightarrow$  we are confident enough, let's stop.  
since  $0.98 > 0.95$

Example 2:  $x_{N+1} \xrightarrow{f_1} P_1 = 0.92 \xrightarrow{f_2} P_2 = 0.91 \Rightarrow$  we are confident enough, let's stop.  
since  $0.92 < 0.95$  since  $0.91 > 0.90$

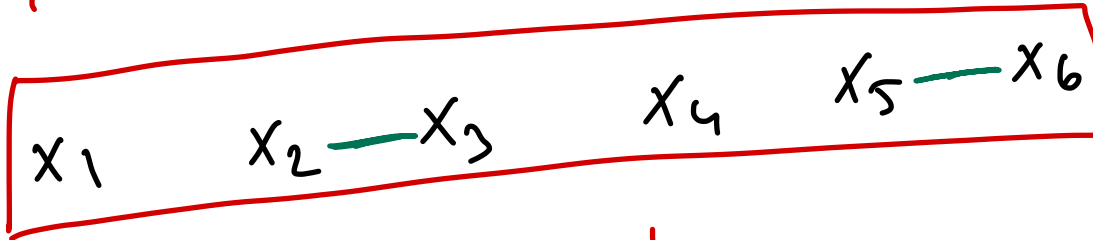
A 1



A 2



A 3



$A^*$

$C_{ij}$  = # of clustering algorithms that put  $x_i$  &  $x_j$  into the same cluster

$C =$

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$	$x_9$	$x_{10}$	$x_{11}$
$x_1$	3										
$x_2$		3	3								
$x_3$			3								
$x_4$				3							
$x_5$					3	2					
$x_6$											
$x_7$											
$x_8$											
$x_9$											
$x_{10}$											
$x_{11}$											

11 x 11

Clustering on the  $C$  matrix would give us  $A^*$ .