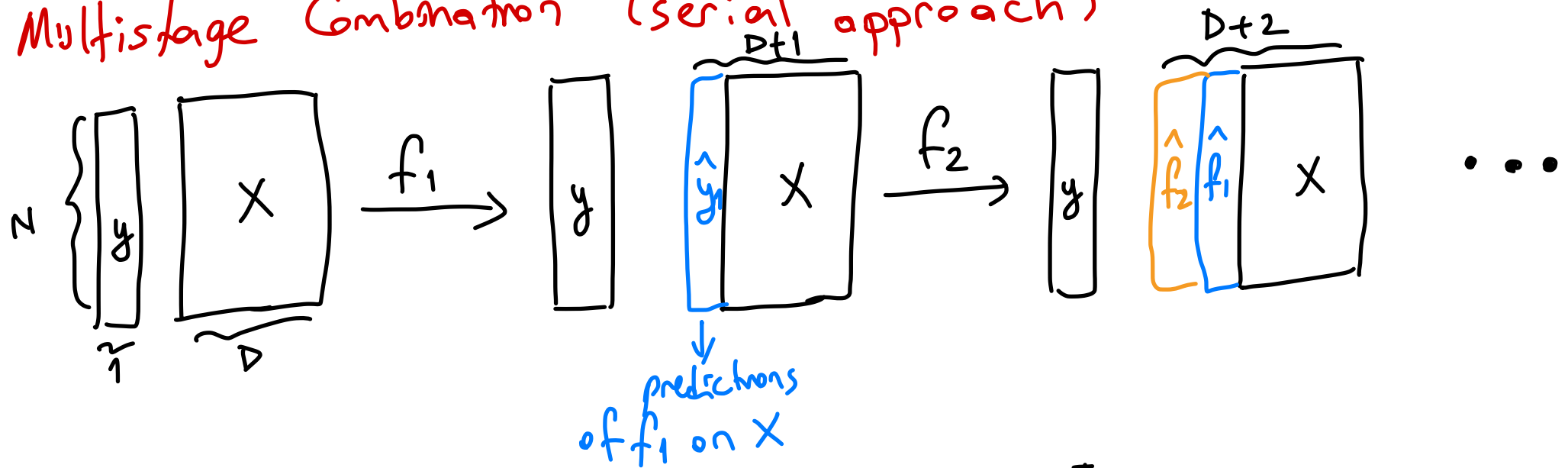


Multistage Combination (serial approach)



$$x_{N+1} \rightarrow f_1(x_{N+1}) \rightarrow f_2([f_1(x_{N+1}), x_{N+1}^T]^T) \rightarrow \dots$$

Let us say we have L base learners

$f_j(x)$ f_1 f_2 \dots f_L

$$\hat{y} = f(f_1, f_2, \dots, f_L | \Phi)$$

\rightarrow Combination function

\rightarrow Combination parameters

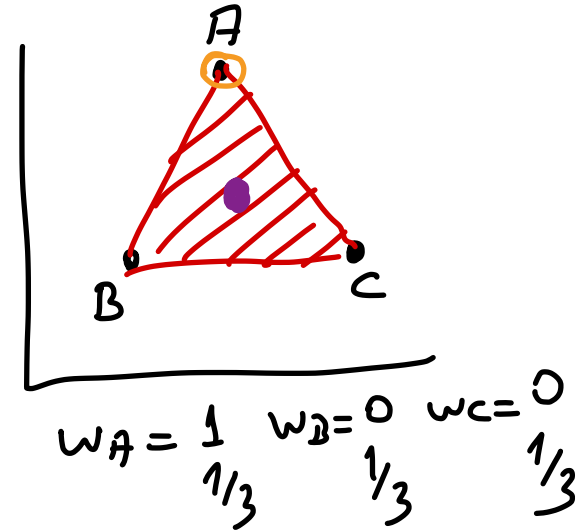
VOTING

$$\hat{y}_i = \sum_{j=1}^L w_j f_j(x_i)$$

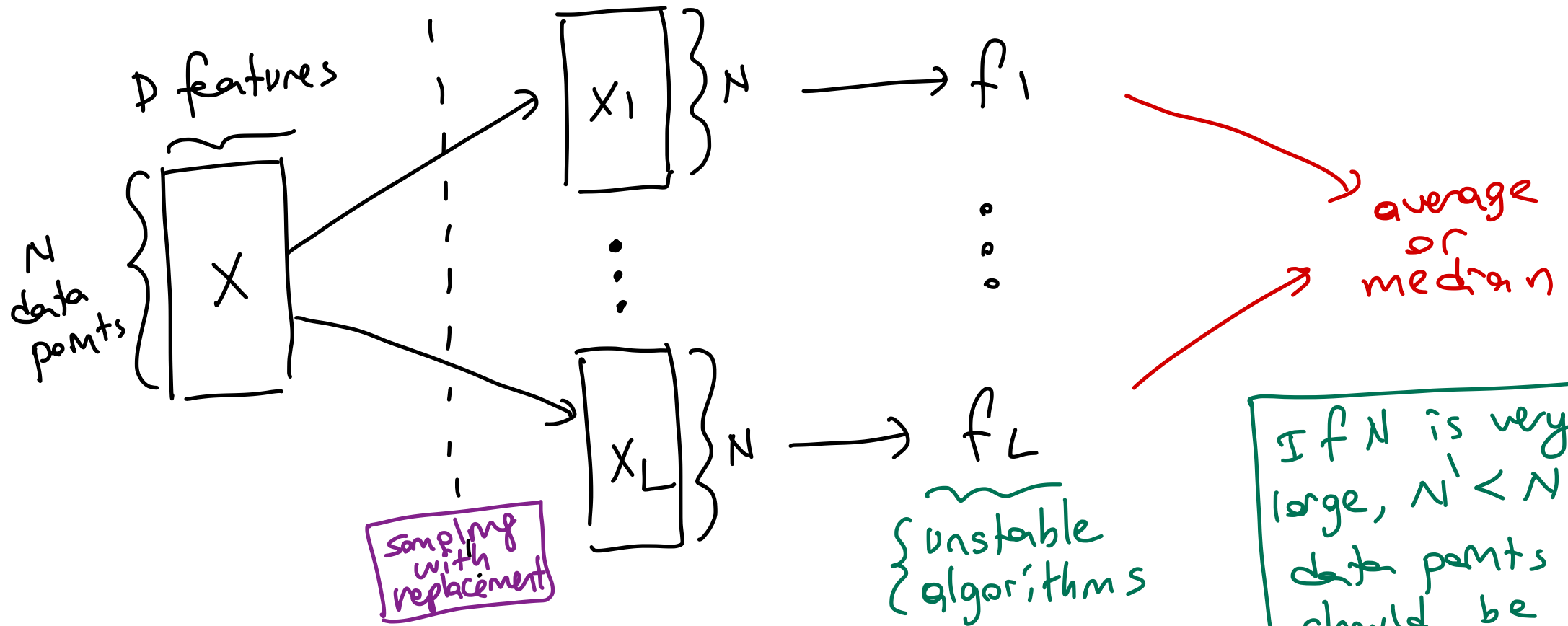
linear opinion models,
ensembles

convex combination $\Rightarrow w_j \geq 0 \quad \forall j$
 $\sum_{j=1}^L w_j = 1$

linear combination $\Rightarrow w_j \in \mathbb{R} \quad \forall j$



BAGGING (Bootstrap AGGregation)

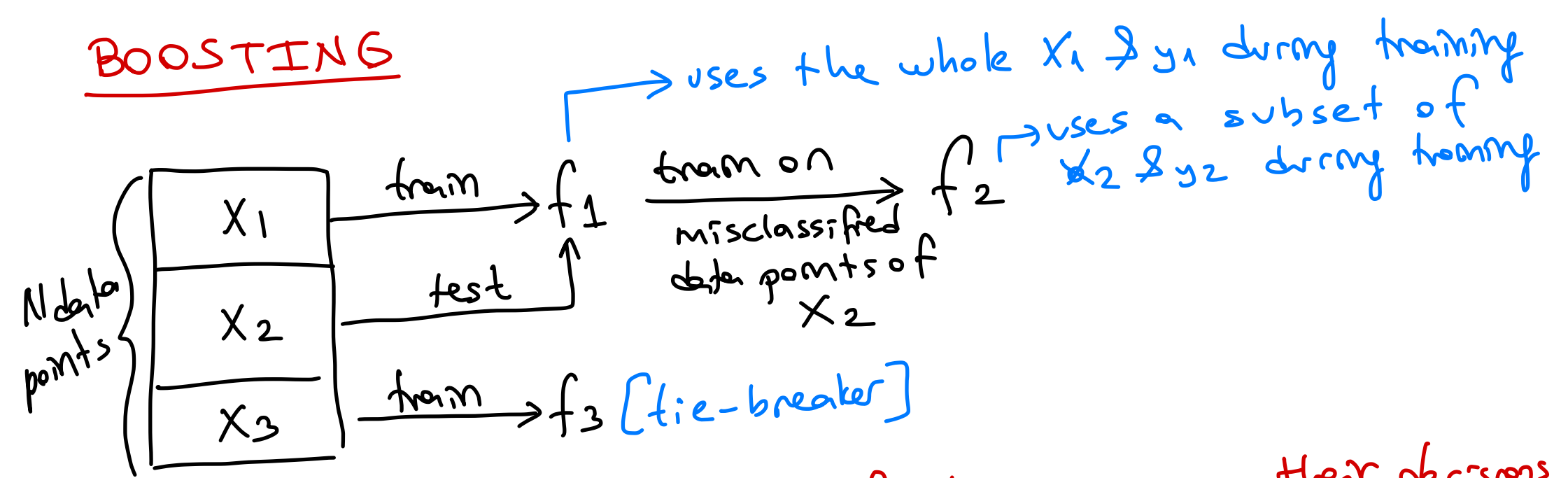


Unstable Algorithm: Highly affected by data set.

Small changes in the training data set.
unstable \Rightarrow DT
stable \Rightarrow KNN

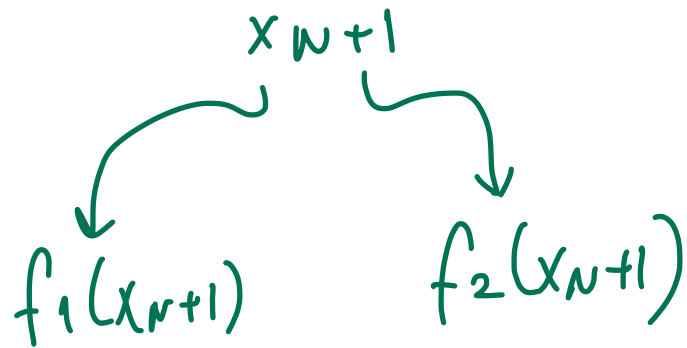
If N is very large, $N' < N$ data points should be picked so that training sets would become different enough.

BOOSTING



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

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



Ada Boost: 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 .

$$w_j = \log \left[\frac{1}{\beta_j} \right]$$

$$\beta_j = \frac{\epsilon_j}{1 - \epsilon_j}$$

...

50% accurate

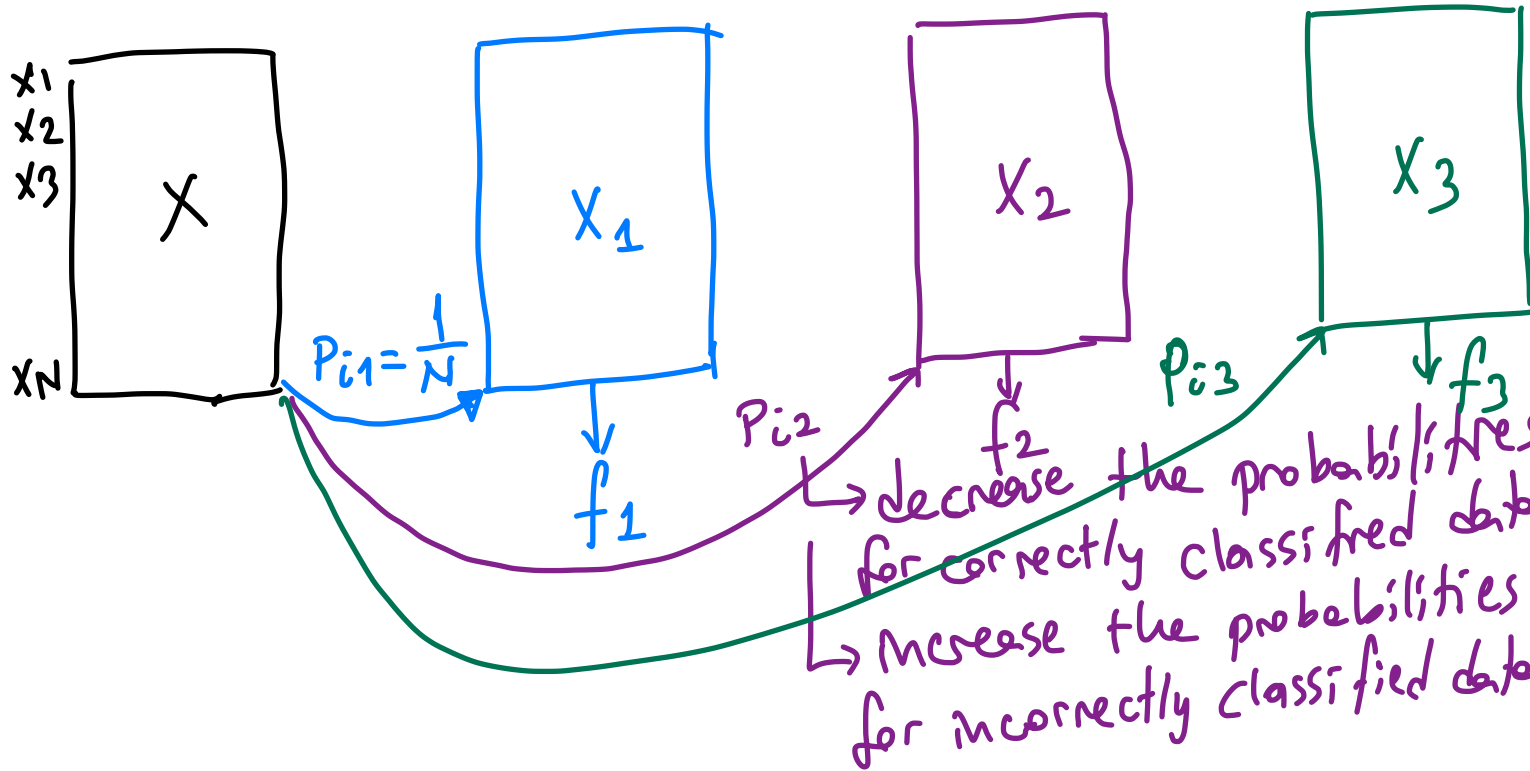
$$\beta_j = \frac{0.5}{1 - 0.5}$$

$$\log(1/\beta_j) = 0$$

99% accurate

$$\beta_j = \frac{0.01}{0.99}$$

$$\log(1/\beta_j) = \log(99)$$



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

based on their error rate

80% accurate $\rightarrow \log(4)$

Mixture of Experts (MoE)

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

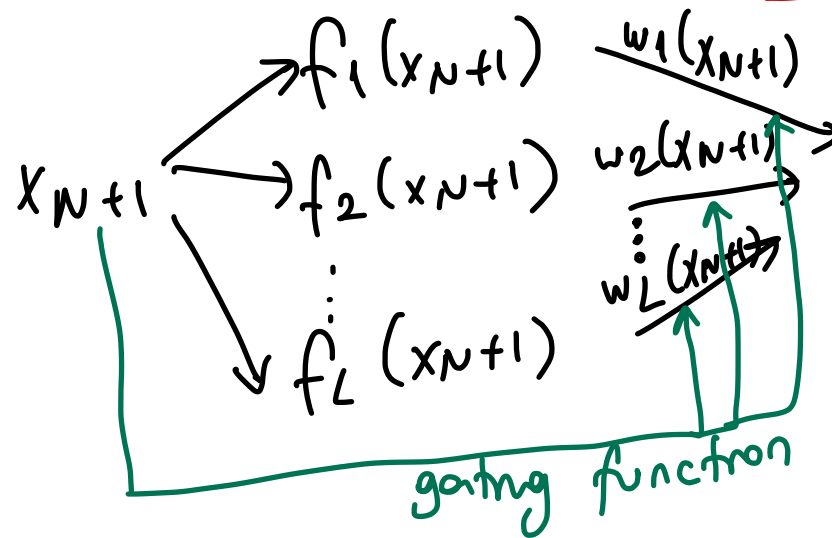
MoE $\Rightarrow \hat{y}_{N+1} = \sum_{j=1}^L \underbrace{w_j(x_{N+1})}_{\substack{\text{w}_j\text{'s will assigned by} \\ \text{the gating function}}} f_j(x_{N+1})$

Competitive

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

\Downarrow
sigmoid
 \Downarrow

$$w_j(x) = \frac{1}{1 + \exp[-V_j^T x - V_{j0}]}$$



- w_1, w_2, \dots, w_L are producing sparse weights mostly zero

- one or some of them are nonzero

\Downarrow
softmax

$$w_j(x) = \frac{\exp(V_j^T x + V_{j0})}{\sum_{k=1}^L \exp(V_k^T x + V_{k0})}$$

A1

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}	x_{11}
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------

A2

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}	x_{11}
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------

A3

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}	x_{11}
-------	-------	-------	-------	-------	-------	-------	-------	-------	----------	----------

C_{ij} = # of clustering algorithms that put x_i and 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	3	3	3	2	1	0	0	0	0	0
x_2	3	3	3	3	2	1	0	0	0	0	0
x_3	3	3	3	3	2						
x_4	3	3	3	3	2						
x_5	2	2	2	2	3						
x_6						1					
x_7							1				
x_8								1			
x_9									1		
x_{10}										1	
x_{11}											1

3
11x11

⇒ clustering on the C matrix* would give you A.

Consensus clustering