

# ML FUNDAMENTAL

DAY 2

MODEL  $\rightarrow$

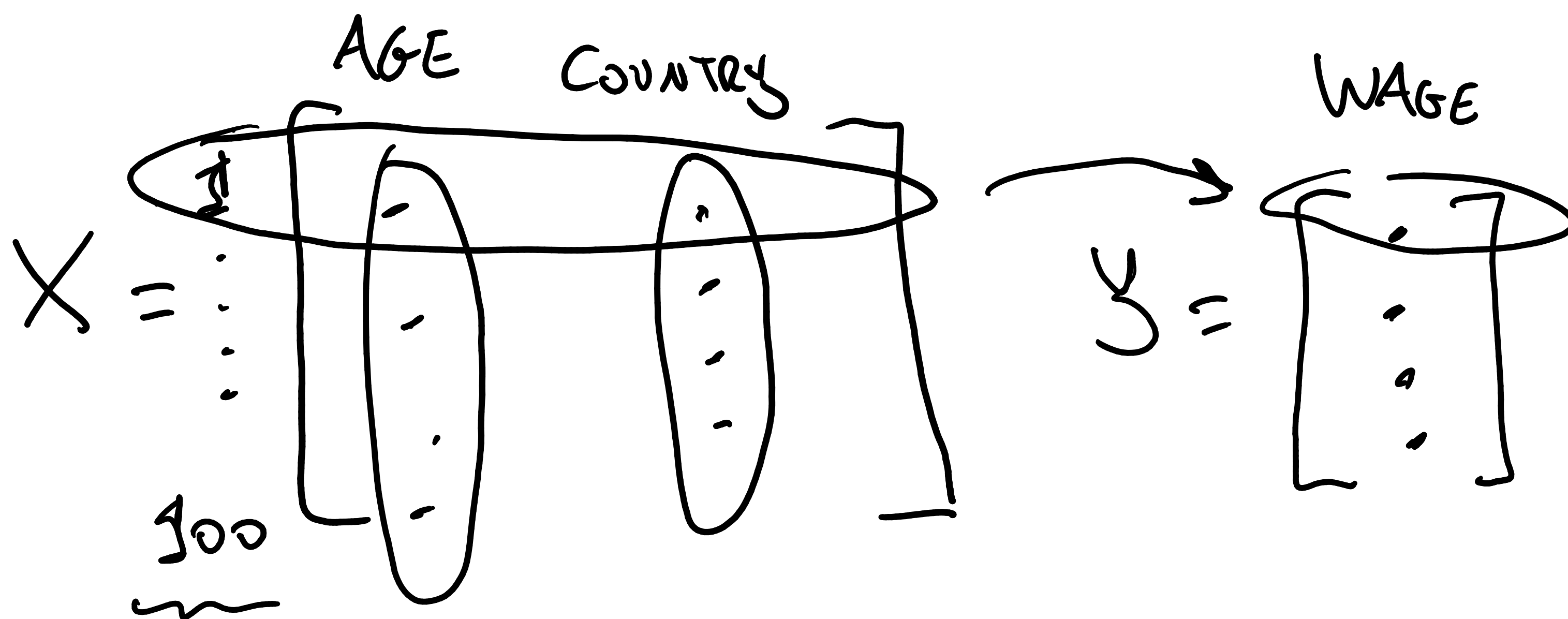
$$f(\cdot) + \varepsilon$$

$y$

TARGET



FEATURES



LINEAR REGN.

$$f(x) = b + w_1 x_1 + \dots + w_p x_p$$

• LOSS FUNCTION

• GRADIENT DESCENT

• OVERFITTING

• REGULAR.

$$(w_1, \dots, w_p)$$



TRAINING

$$(2, 1.5, \dots)$$

REGRESSION

$$y \in \mathbb{R}$$

$$y \in \mathbb{N}$$

$$y \in \begin{matrix} A \\ B \end{matrix} \downarrow$$
$$y \in \{0, 1\}$$

$$(y_i - \hat{y}_i)^2$$

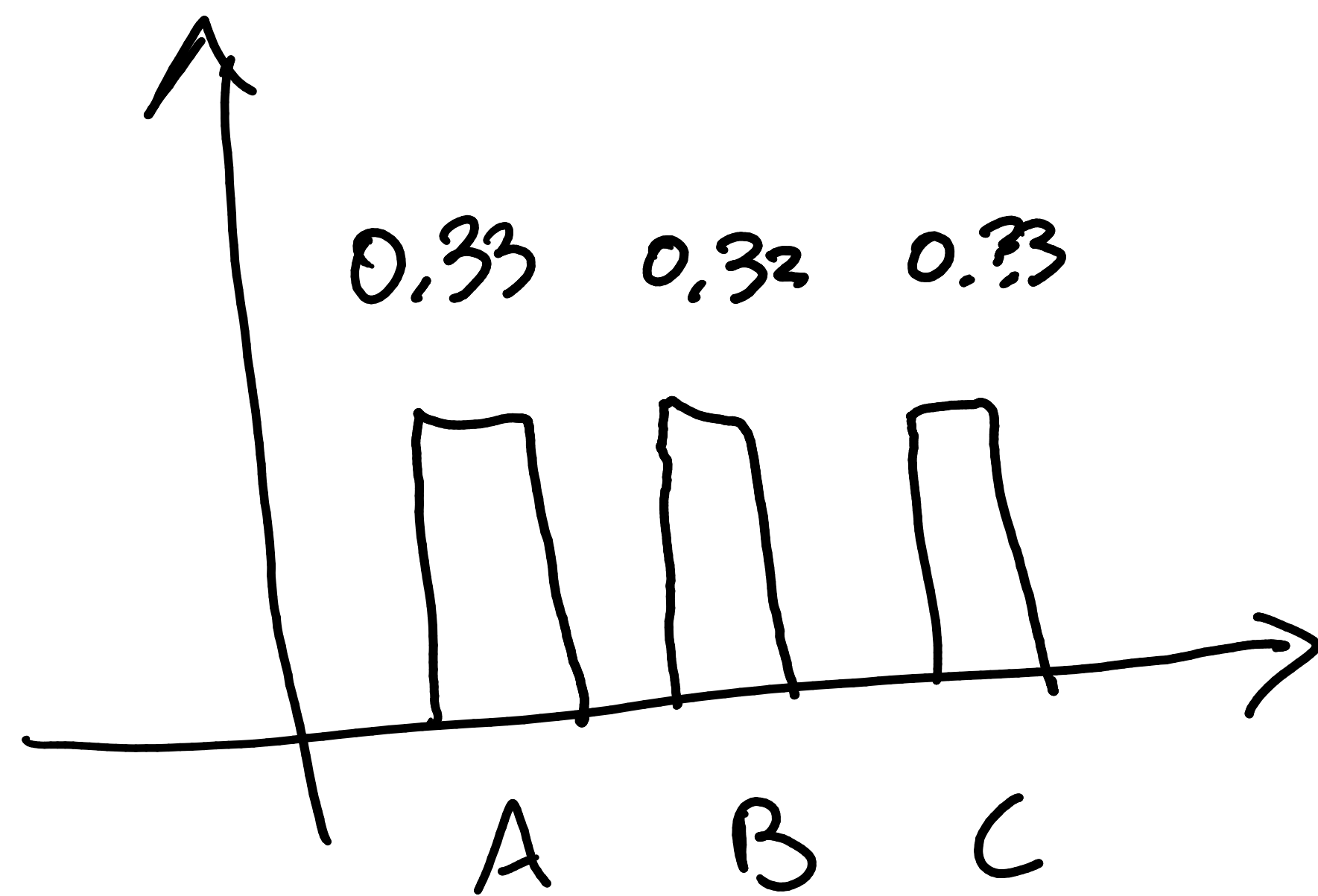
CLASSIFICATION

$$y_1 = 0$$

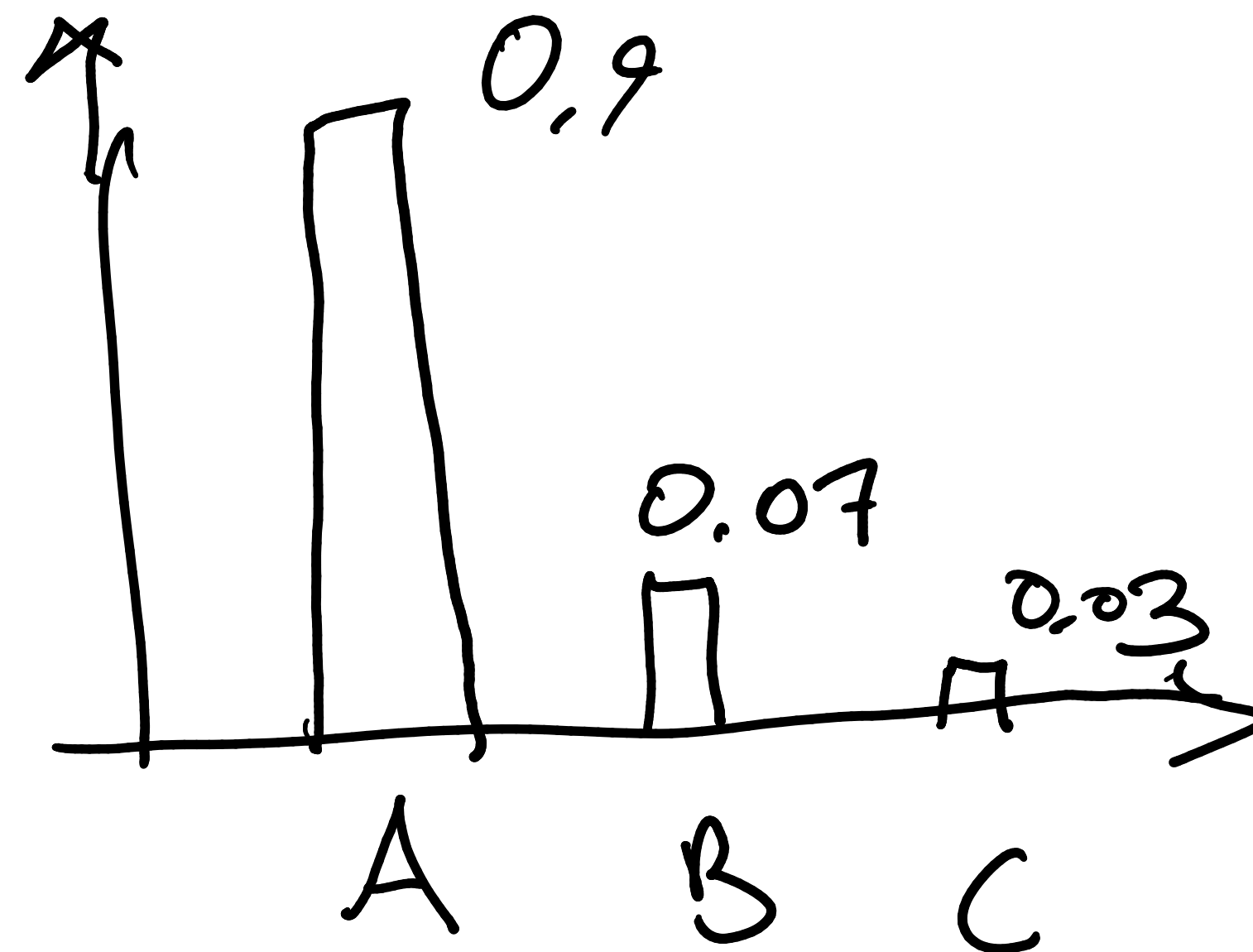
$$\hat{y}_1 = 1$$

$$(0 - 1)^2$$

①



②



$$Y = \begin{cases} A \\ B \\ C \end{cases}$$

$$\hat{f}(x) = 0.93$$

$$Y = \begin{cases} A \rightarrow 0 \\ B \rightarrow 1 \end{cases}$$

$$(y_i - \hat{y}_i)^2 = \underline{1 - 0} = ?$$

+

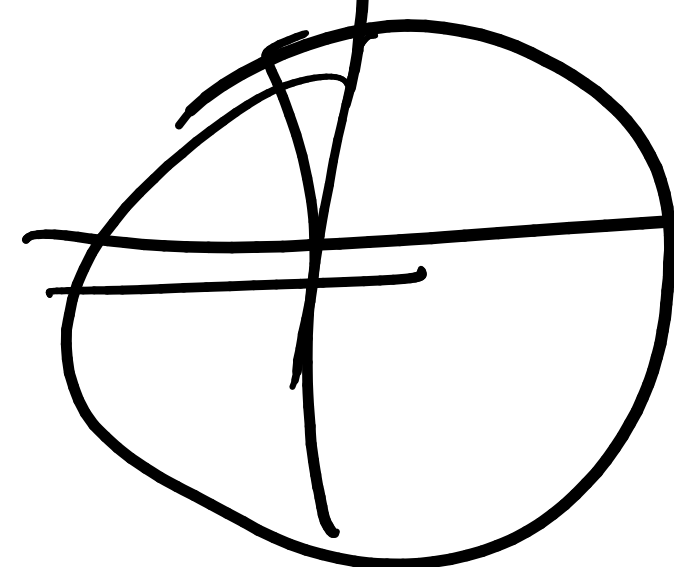
$$\hat{y}_i = \begin{matrix} 1 \\ 0 \end{matrix} \rightarrow p \in [0, 1]$$

VARIANCE

REGR. / CLASS

GB / RF TREES

LOG. REGR.  
LIN. REGRES



BIAS

# Trees

$$X = \left[ \begin{array}{c} \text{ } \end{array} \right]$$

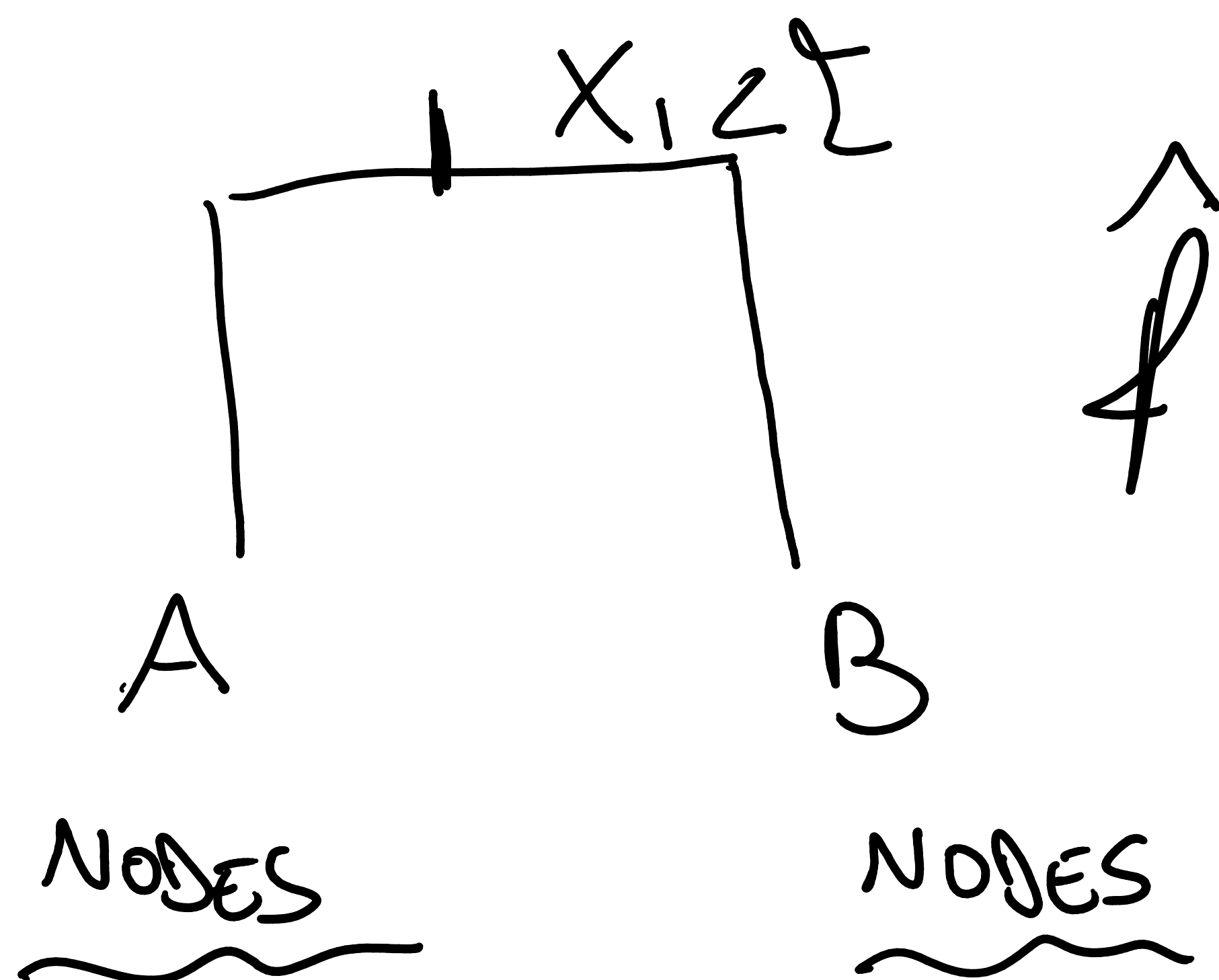
$$X = \left[ \begin{array}{c} x_1 \\ \vdots \\ x_n \end{array} \right]$$

$$Y = \left[ \begin{array}{c} \text{ } \end{array} \right]$$

If  $x_1 \in A$  then  $\hat{Y} = A$

ELSE

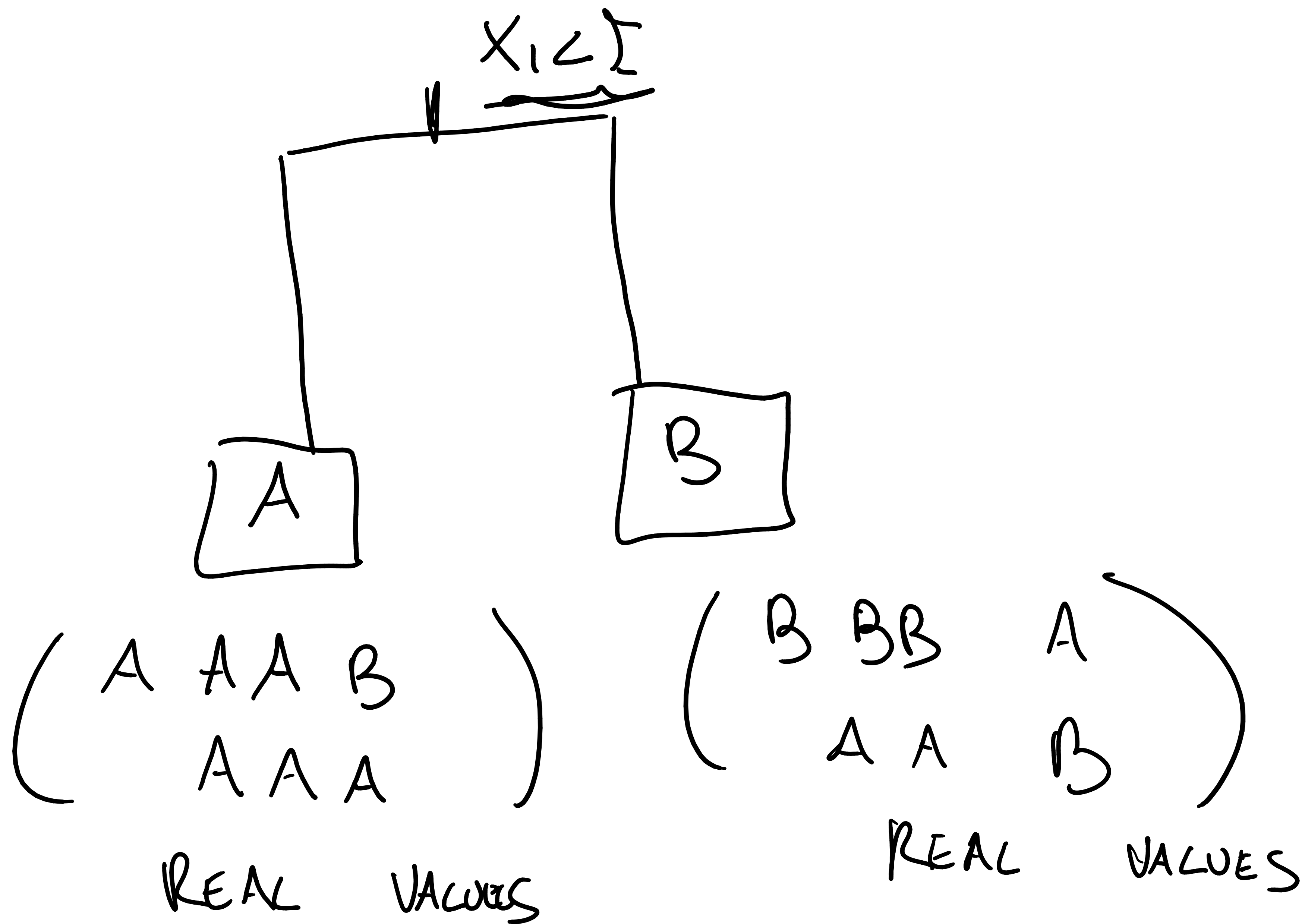
$$\hat{Y} = B$$



$$\hat{x} = 20$$

$$X = \begin{matrix} & X_1 \\ \begin{matrix} 1 \\ \vdots \\ 100 \end{matrix} & \begin{bmatrix} 30 \\ 21 \\ \vdots \\ 1 \end{bmatrix} \end{matrix} \rightarrow \hat{f} \rightarrow \hat{y} = \begin{bmatrix} B \\ B \\ A \end{bmatrix} \quad \hat{y} = \begin{bmatrix} \end{bmatrix}$$



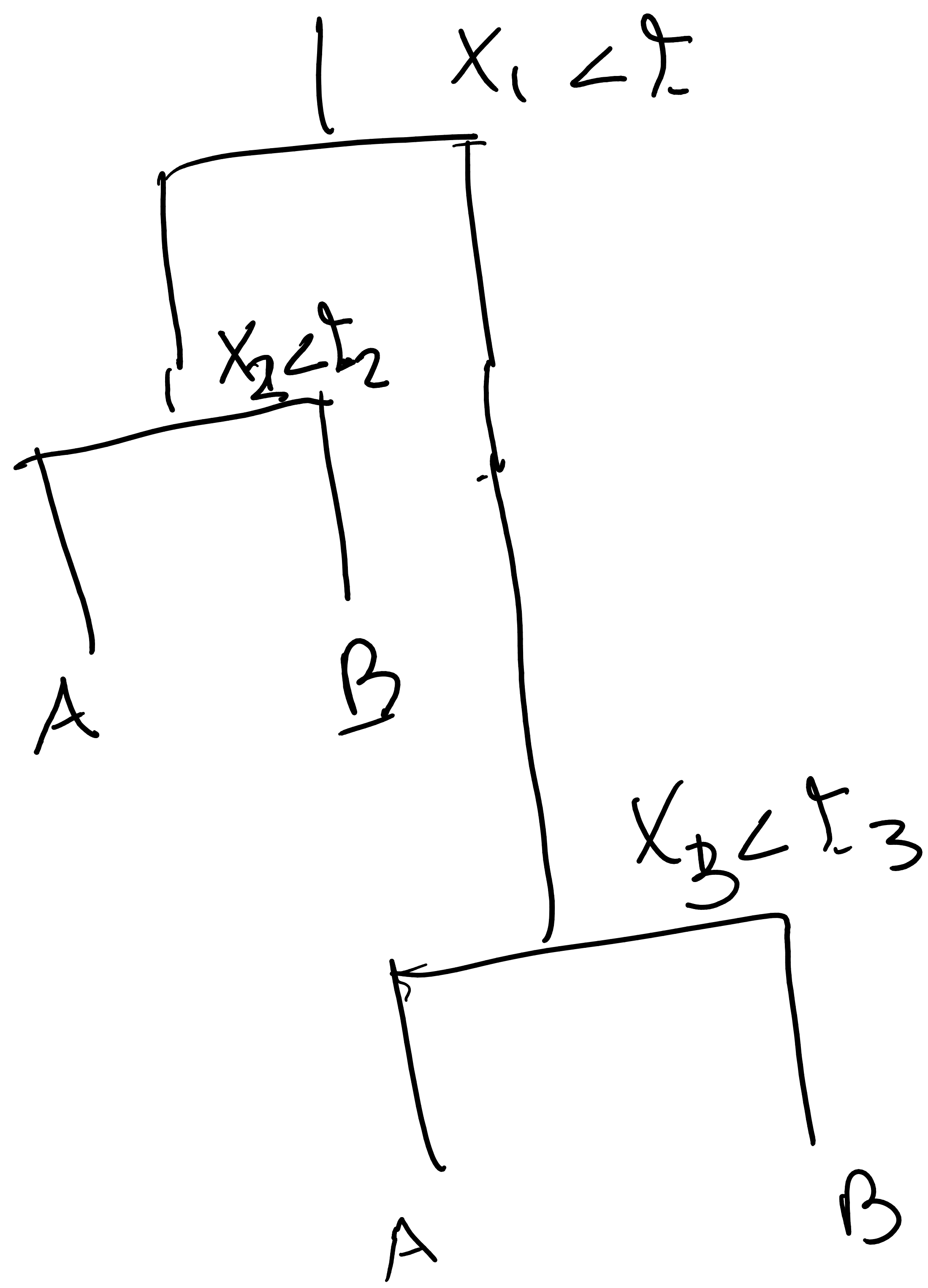


men! (Entropy\_left + Entropy\_right)

GREEDY

ALG.

~~Handwritten signature~~



1  
SPECIFICATION

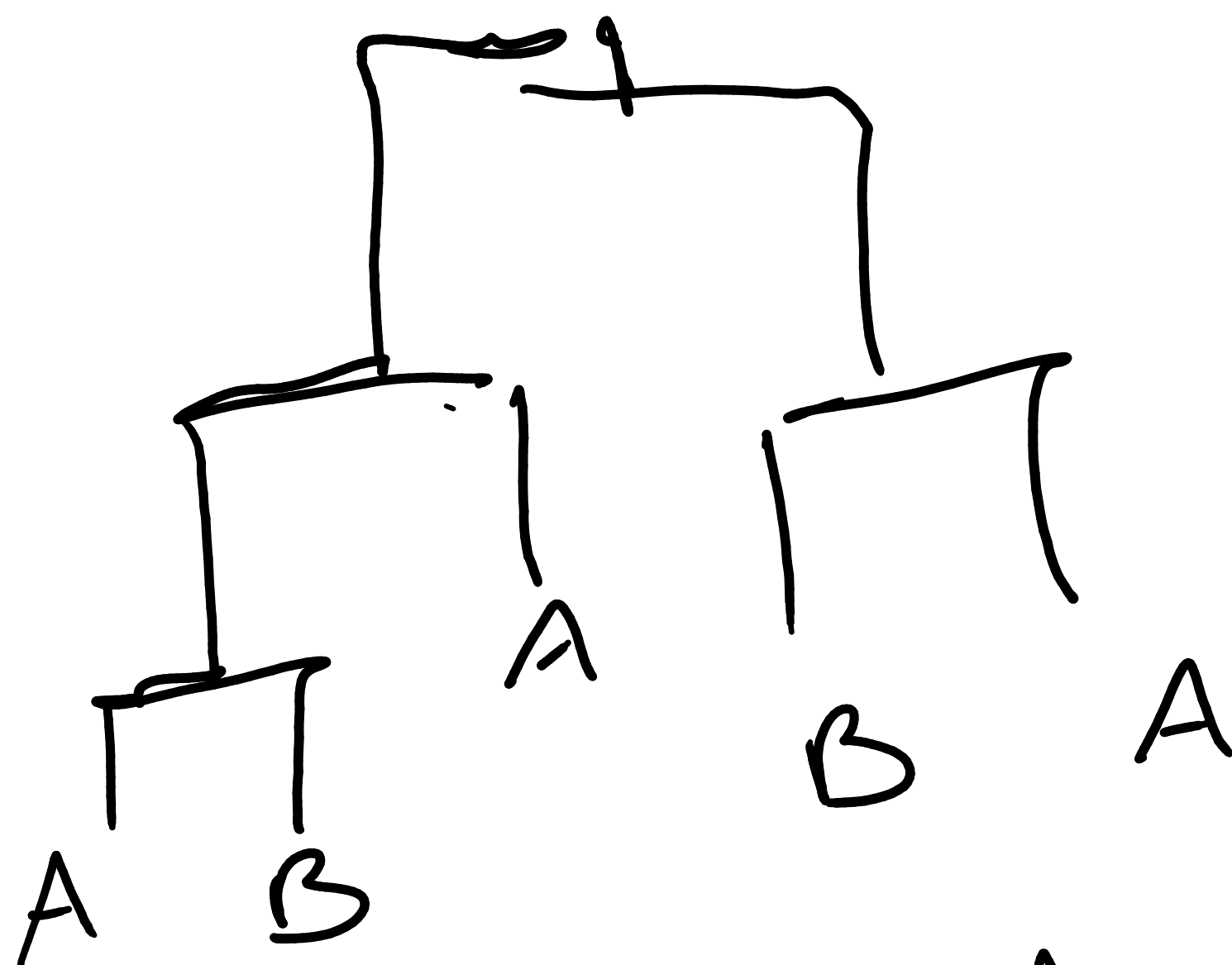
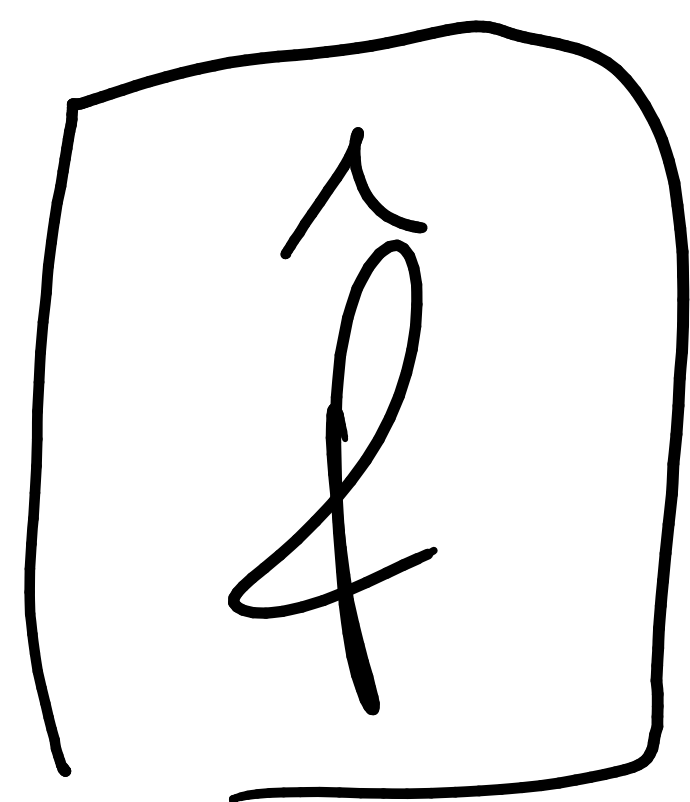
2  
TRAINING

3  
EVALUA.  
(Pred.)

TREE

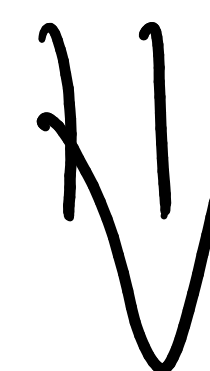
ENTR,  
ENTR

ACCURACY



ACCURACY  $\hat{y}$   $\hat{y}$

$\hat{f}(X_{TEST})$



$\hat{y} = \begin{pmatrix} A \\ A \\ \vdots \\ B \end{pmatrix}$



# ENSEMBLING OF MODELS

$\hat{f}_1, \dots, \hat{f}_B$

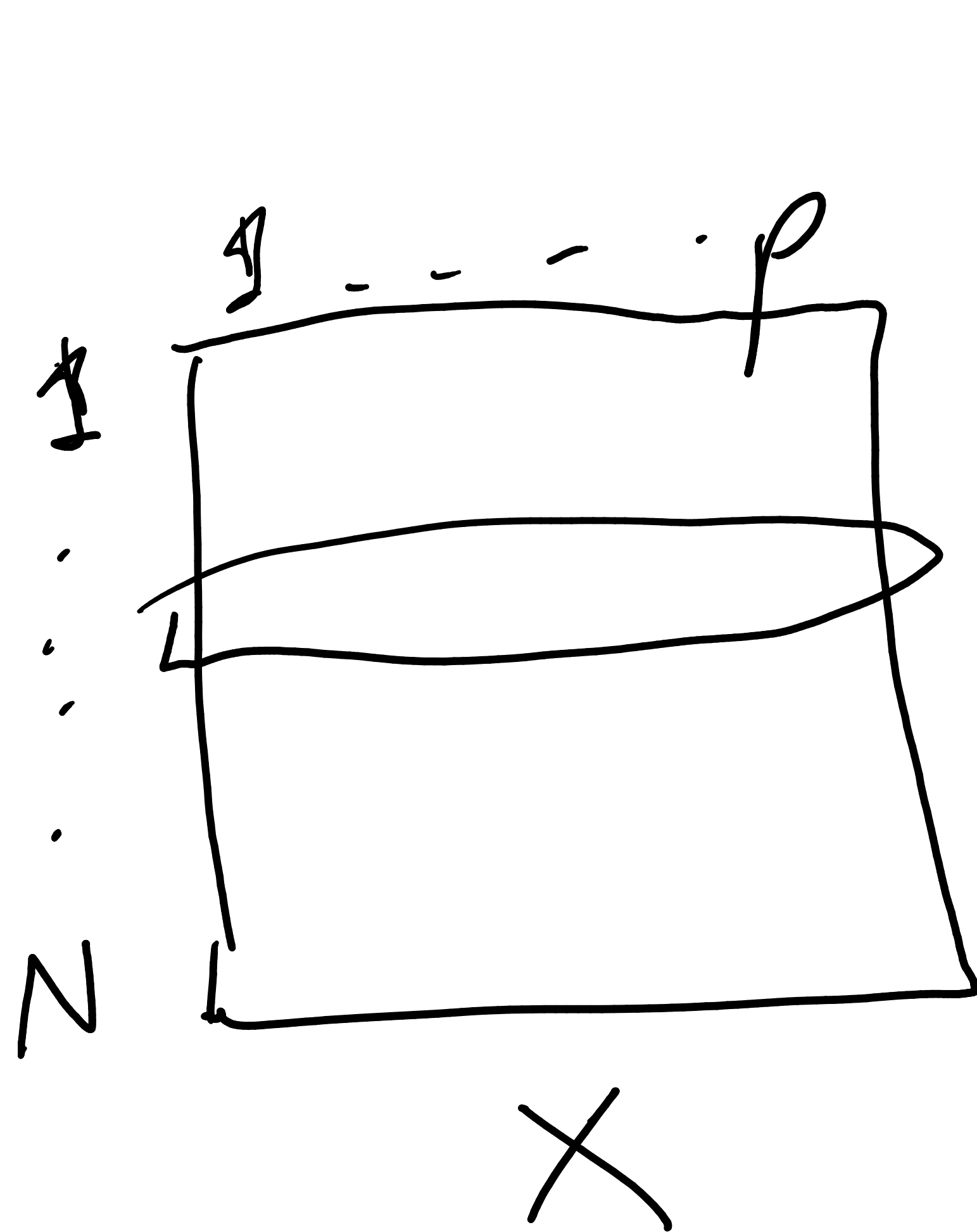
PARALLEL

$\hat{f}_1, \dots, \hat{f}_B$

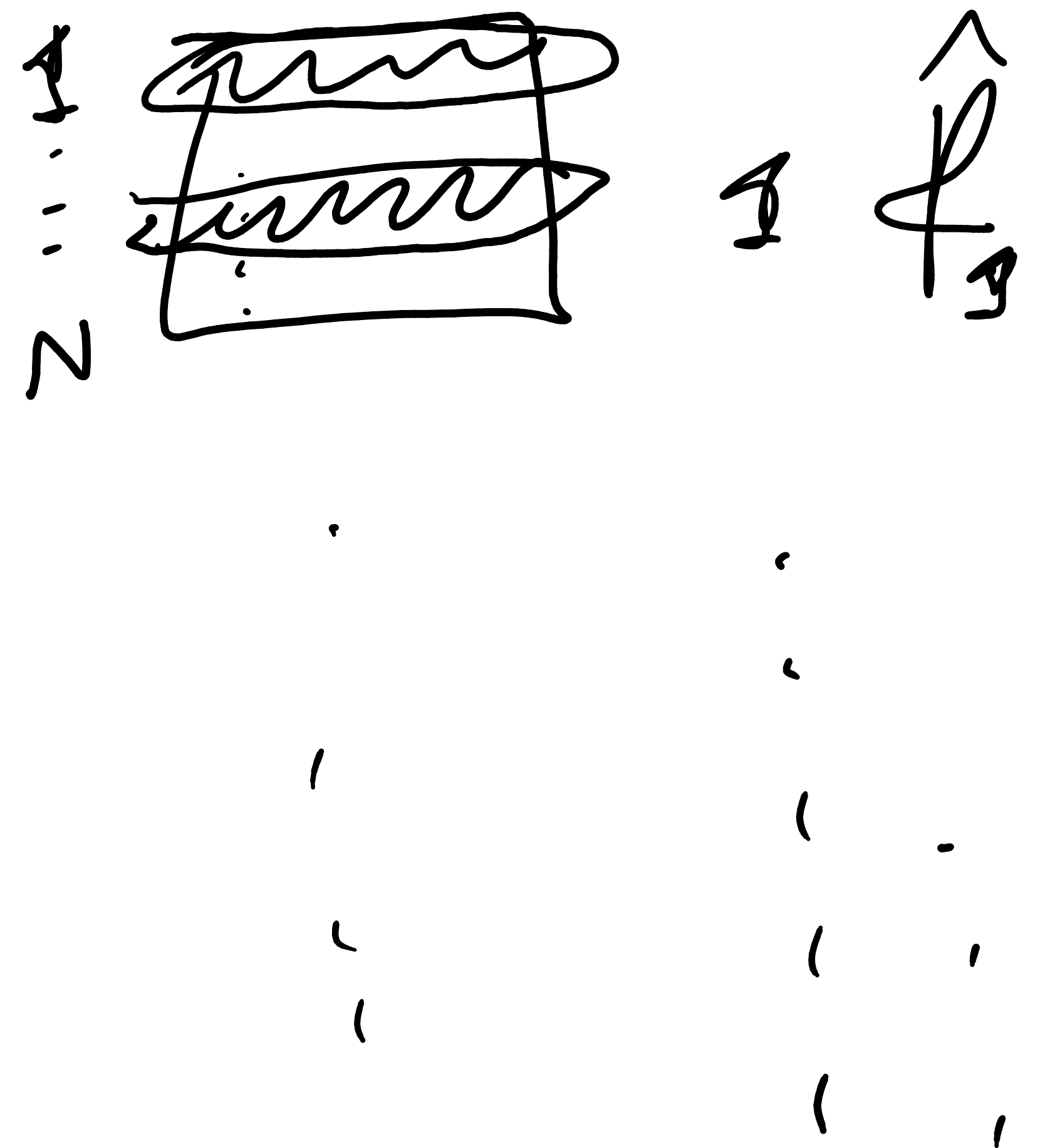
SEQUENTIAL

$\hat{f}_B(\hat{f}_{B-1}(\dots(\hat{f}_1)\dots))$

$\hat{f}_1 \rightarrow \hat{f}_2 \rightarrow \dots$

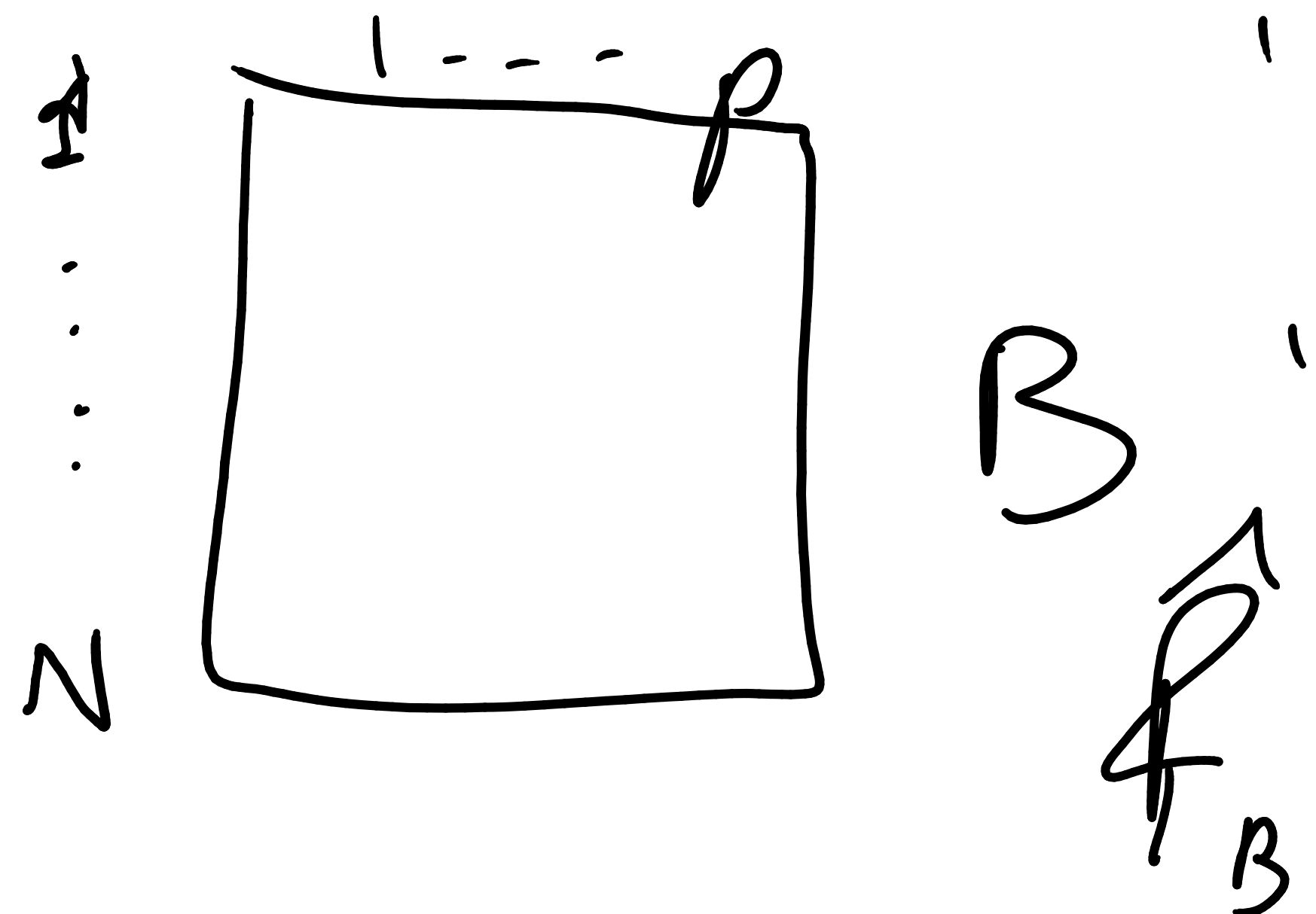


SAMPLING WITH REPLACEMENT  
 BOOTSTRAP



BAGGING

(BOOTSTRAPPING AND AGGREGATING)



$$(X_1, \dots, X_n)$$

$$X_i \sim N(\mu, \sigma^2)$$

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

# RANDOM FOREST

- BAGGING ( $1, \dots, B$  datasets)
- FOR EACH DATASET:
  - FIT A TREE
    - FOR EACH split, use subset of FEATURES
- AGGREGATE pred.



Predictions WITH

ENSEMBLING

(CLASSIFICATION)

$$X^* = (x_1^*, \dots, x_p^*)$$

$$\begin{array}{lcl} \hat{f}_1(x^*) = 0.7 \Rightarrow B \\ \vdots \\ \vdots = 0.3 \Rightarrow A \\ \vdots \\ \hat{f}_B(x^*) = 0.9 \Rightarrow B \end{array}$$

~~# A~~

# B

A  
AVERAGE IN  
CASE OF REGRESSION

↓  
[ FINAL  
PREDICTION ]

<u>Depth</u>	<u># OF TREES</u>
$\uparrow \Rightarrow \text{BIAS} \downarrow$	$\uparrow \Rightarrow \text{VARIANCE} \downarrow$

# BOOSTING

$X_{\text{TRAIN}}$

$y_{\text{TRAIN}}$

$$\hat{f}_1(X_{\text{TRAIN}}) = \hat{y}_1, \quad \underline{e}_1 = |\hat{y}_1 - y_{\text{TRAIN}}| \quad \Delta$$

$$\hat{f}_2(X_{\text{TRAIN}}) = \hat{e}_1, \quad \underline{e}_2 = |\hat{e}_1 - e_1| \quad \Delta$$

$\vdots$

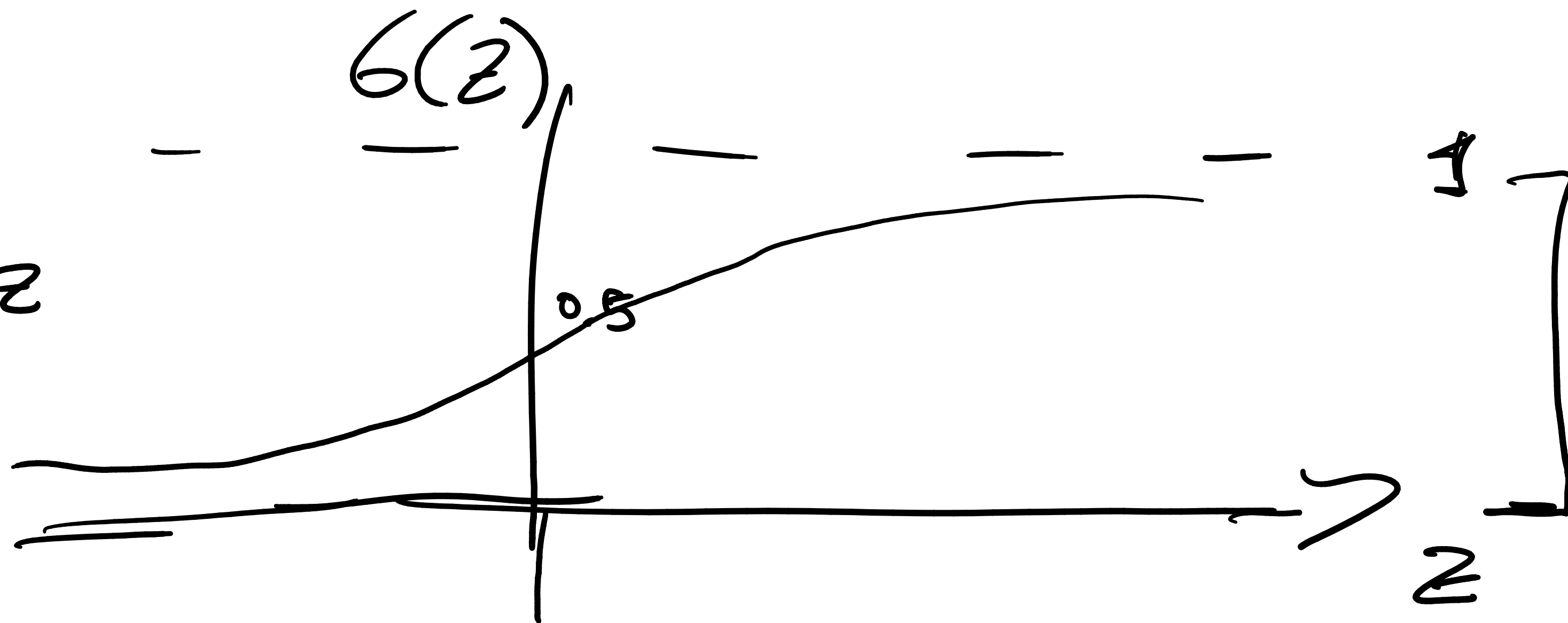
$$\hat{f}_B(X_{\text{TRAIN}}) = \hat{e}_{B-1} \quad \Delta$$

# LOGISTIC REGRESSION

$$\underbrace{y} = b + w_1 x_1 + \dots + w_p x_p$$

$$y \in \mathbb{R} \Rightarrow y \in [0, 1]$$

$$G(z) = \frac{1}{1 + e^{-z}}$$



# LOGISTIC REGR.

$$p = \frac{1}{1 + e^{-(b + w_1 x_1 + \dots + w_p x_p)}}$$

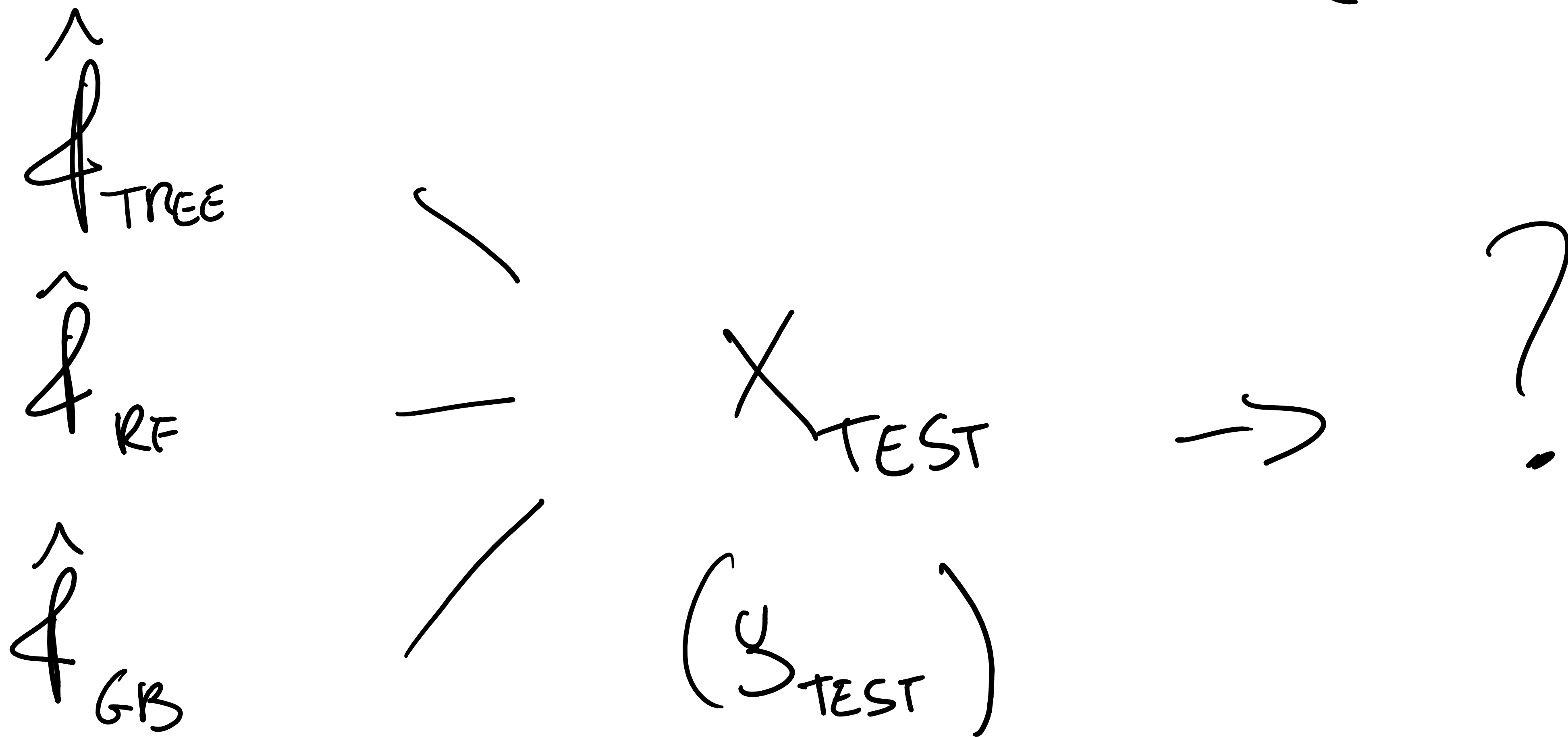
Loss FUNCTION : BINARY CROSS ENTROPY

$\hat{y}$	0	1
0	$-\log(1-0) = 0$	$-\log(1-1) = +\infty$
1		$-\log(1) = 0$

$$\hat{p} = 1$$

$$\hat{p} = 0$$

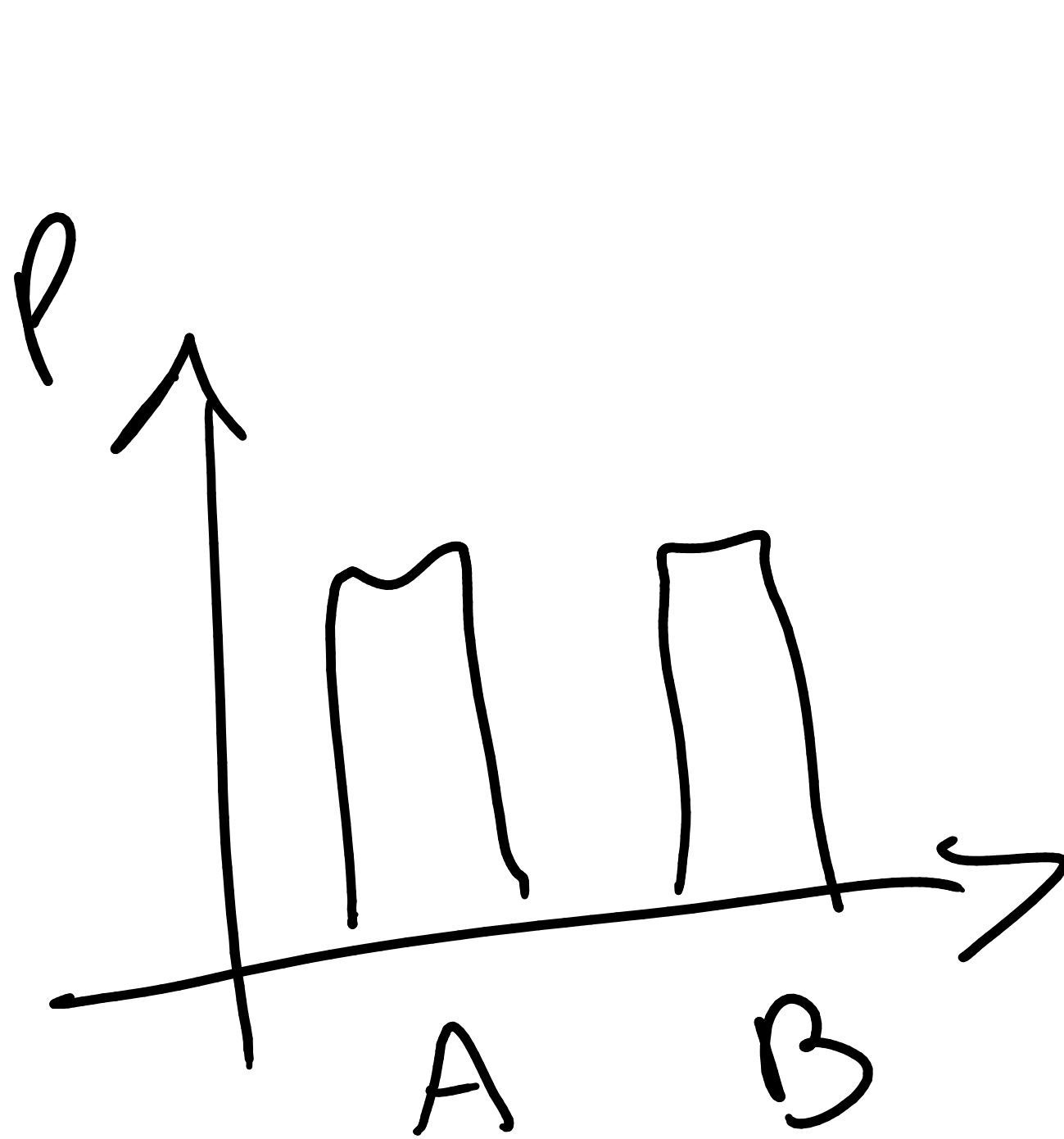
# EVALUATION METRICS (CLASS.)



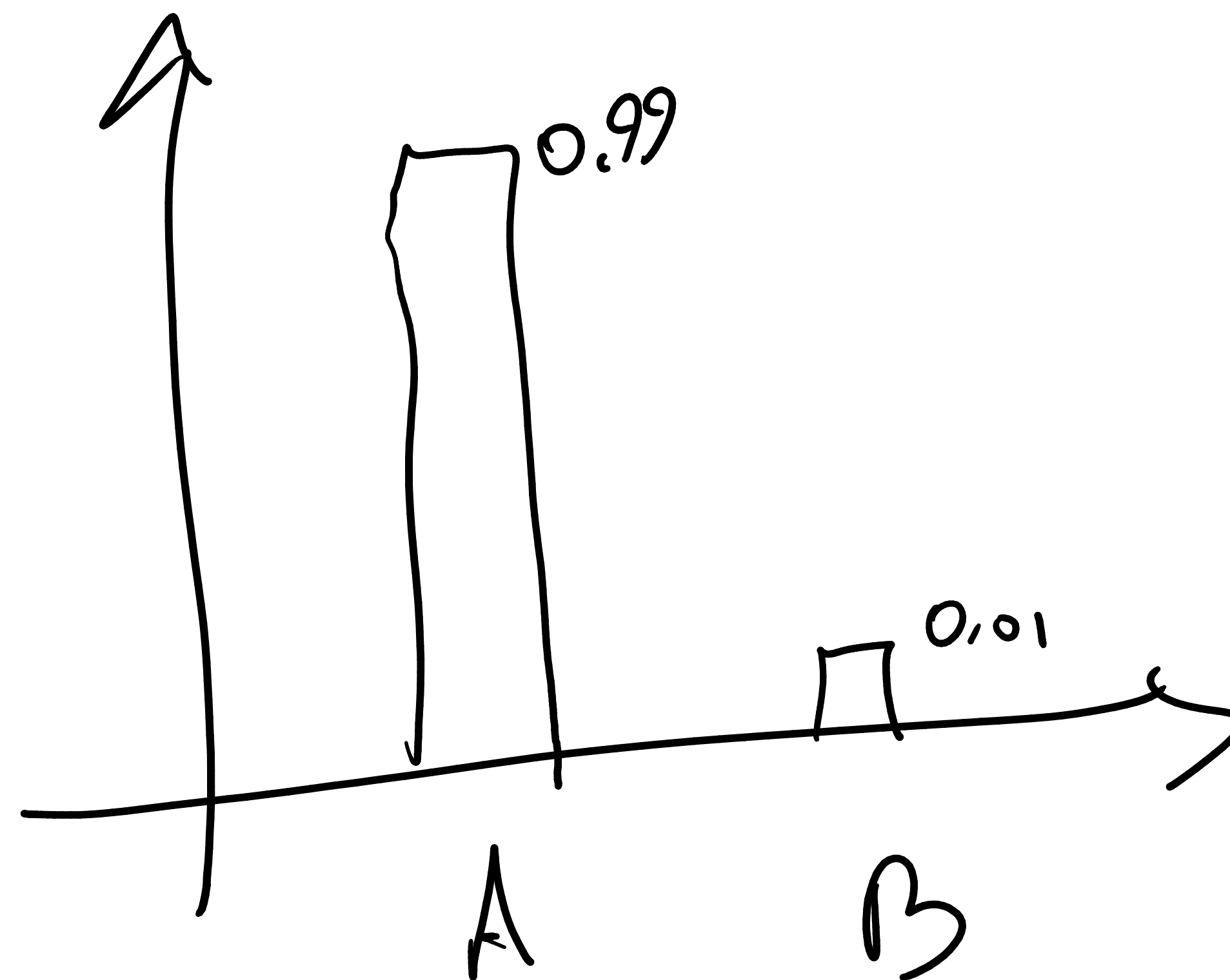
---

$Y = \begin{cases} A \\ B \end{cases}$  BINARY CLASS.

• UNBALANCED



SAMPLE



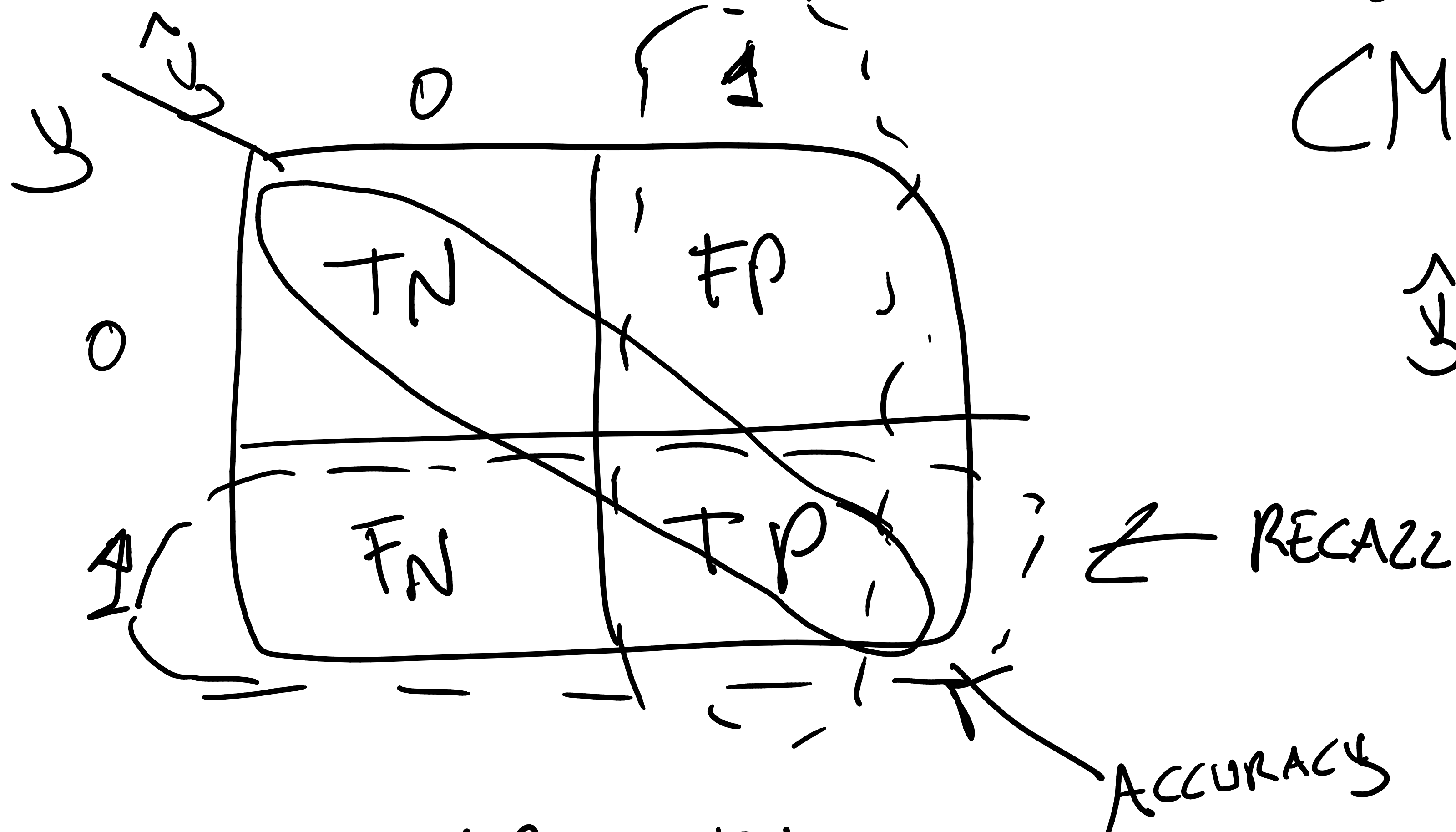
— CONFUSION

MATRIX

PRECISION

$$CM(I)$$

$$\hat{y} = \begin{cases} 0 & \hat{p} < \tau \\ 1 & \hat{p} \geq \tau \end{cases}$$



• ACCURACY :  $\frac{TP + TN}{N}$

• PRECISION :  $\frac{TP}{TP + FP}$

• RECALL :  $\frac{TP}{FN + TP}$

$$F_1\text{-score} = \frac{2}{\frac{1}{\text{PREC}} + \frac{1}{\text{REC}}}$$