\* Performance prediction of models :- Accuracy → metric
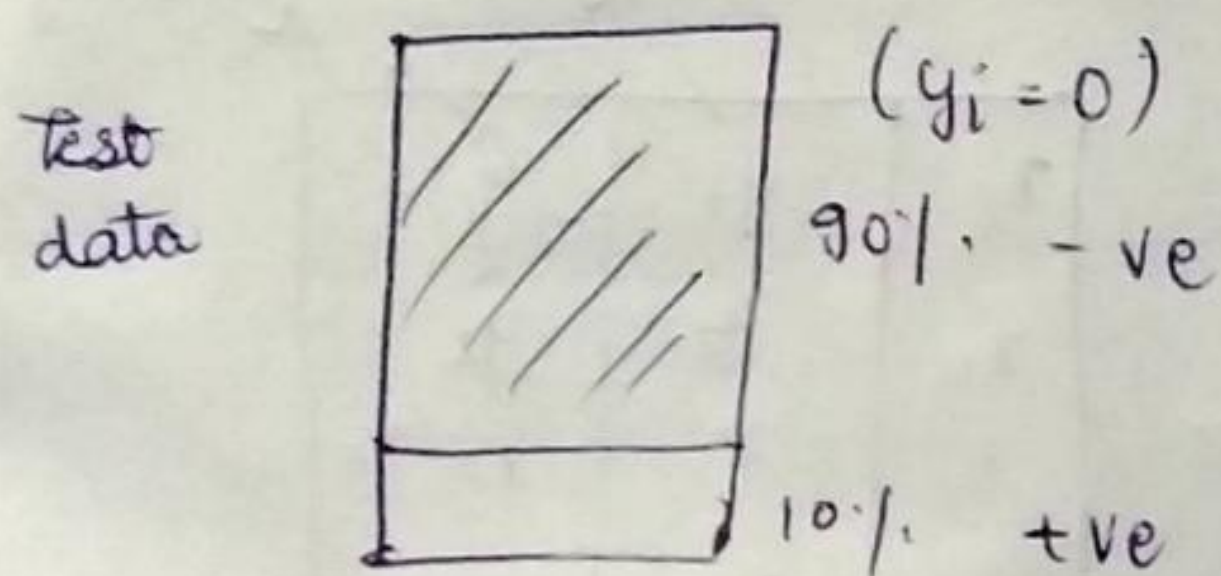
$$\downarrow$$

classification, regression

(KNN)

Accuracy = $\dfrac{\text{no. of correctly classified pts}}{\text{total no. of points in } D_{test}}$

easy to
understand
measure

Performance is always measured on test data.

① Imbalanced data :→



Test
data

$(y_i = 0)$

90% -ve

10% +ve

dumb
model    $x_q$ → -ve

{return -ve}

$\rightsquigarrow$ model → accuracy = 90% = 0.9

high accuracy

So, accuracy is not a good measure to predict if
a model is performing good or not when
dataset is imbalanced.

②

| | x | y | $M_1$ | $M_2$ | $\hat{y}_1$ | $\hat{y}_2$ |
|---|---|---|---|---|---|---|
| +ve | $x_1$ | 1 | 0.9 | 0.6 | 1 | 1 |
| | $x_2$ | 1 | 0.8 | 0.65 | 1 | 1 |
| -ve | $x_3$ | 0 | 0.1 | 0.45 | 0 | 0 |
| | $x_4$ | 0 | 0.15 | 0.48 | 0 | 0 |
| | | | | | | |

{ Test set }          $\hat{y}$ = Predicted value

$M_1$ & $M_2$ return a
prob score

KNN

{ $x_q$ → prob ($y_q = 1$) }

$(0 \le \phi \le 1)$
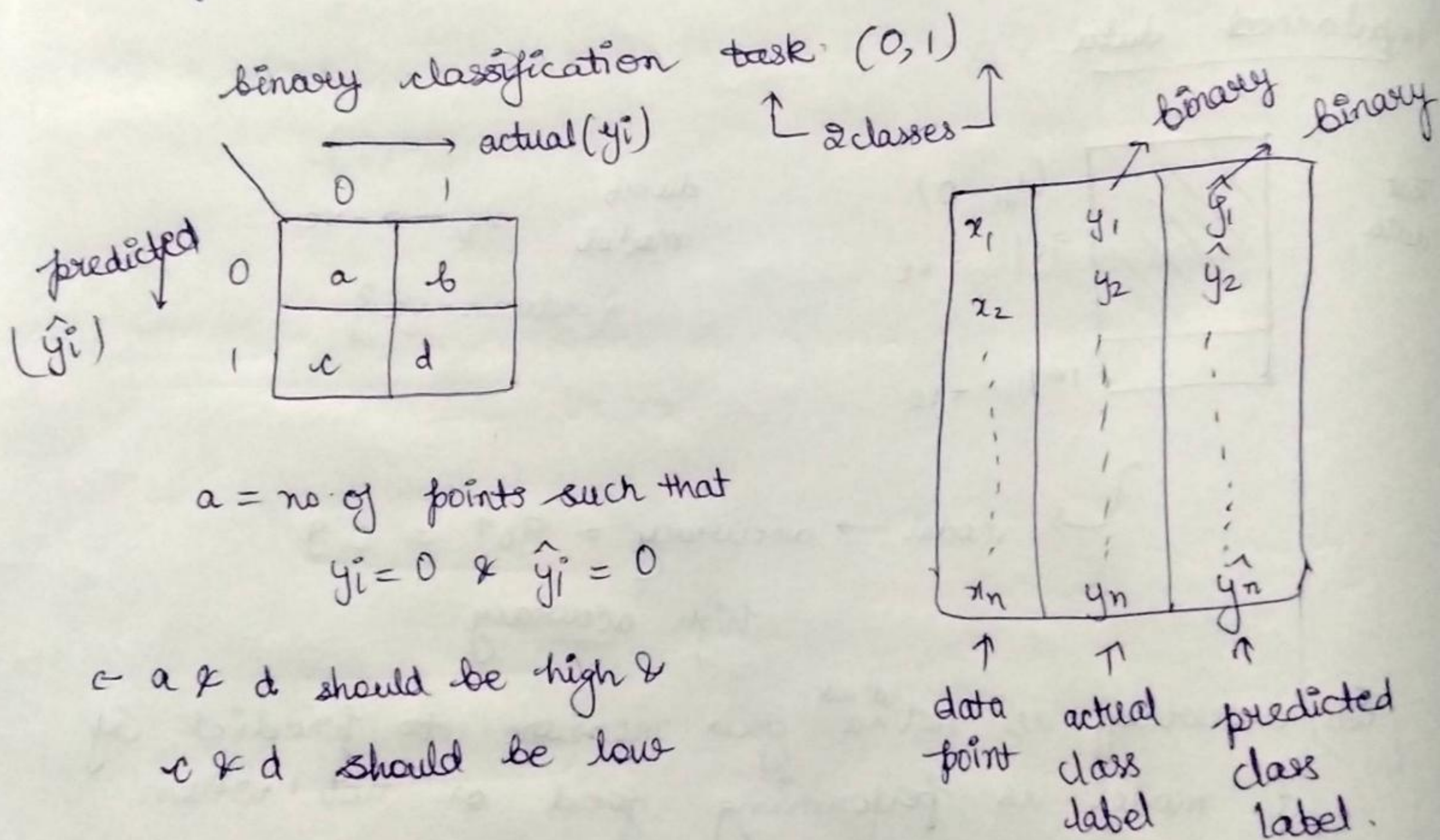
⟨ ⇒ predicted class labels are exactly the same in $M_1$ & $M_2$.

⇒ $M_1$ is better than $M_2$ by looking at probability scores.

accuracy ⟶ cannot use prob-score
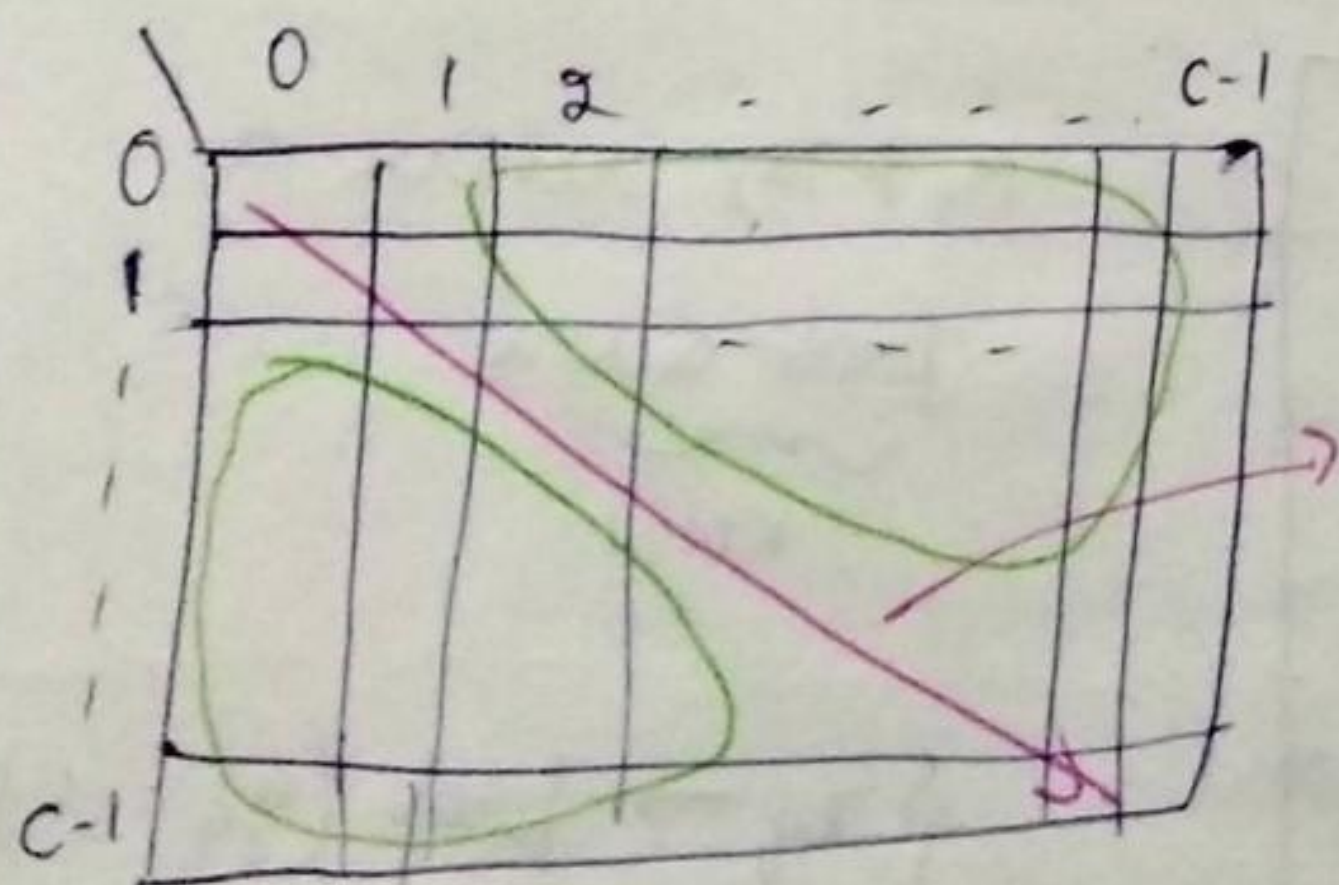
$\hat{y_1}$ & $\hat{y_2}$ → by looking at $\hat{y_1}$ & $\hat{y_2}$, $M_1$ & $M_2$ are having same accuracies.

<u>21.2</u>
⊛ <u>Confusion matrix</u> → doesnot / cannot process prob-scores

binary classification task (0,1) $\quad \lfloor$ 2classes $\rfloor$     binary   binary

⟶ actual ($y_i$)

| | 0 | 1 |
|---|---|---|
| 0 | a | b |
| 1 | c | d |

predicted ($\hat{y_i}$)

| $x_1$ | $y_1$ | $\hat{y_1}$ |
|---|---|---|
| $x_2$ | $y_2$ | $\hat{y_2}$ |
| ⋮ | ⋮ | ⋮ |
| $x_n$ | $y_n$ | $\hat{y_n}$ |

↑ data point   ↑ actual class label   ↑ predicted class label.

$a$ = no. of points such that
$$y_i = 0 \ \& \ \hat{y_i} = 0$$

⟵ a & d should be high &
c & d should be low

multiclass - classification → (c-Classes)



principal diagonal values should be high and off-diagonal values should be small for a sensible model.

actual →

| | 0 | 1 |
|---|---|---|
| predicted 0 | true negative | false negative |
| 1 | false positive | true positive |

N    P

True positive (TP)

↓ └→ what is the predicted label?

are your predicted label matching with actual? (T/F)

N = sum of false positive and true negative

P = sum of false negative and true positive

n = total no. of points (N + P)

True positive Rate (TPR) = $\dfrac{\text{True positives}}{P}$

True negative Rate (TNR) = $\dfrac{\text{True negatives}}{N}$

false ~~negative~~ positive Rate (FPR) = $\dfrac{\text{false positives}}{N}$

False negative Rate (FNR) = $\dfrac{\text{false negatives}}{P}$

for ex :→

→ actual

| | 0 | 1 |
|---|---|---|
| predicted 0 | 850 (TN) | 6 (FN) |
| 1 | 50 (FP) | 94 (TP) |

N = 900    P = 100

Test :→ 900 -ve ⎫ imbalanced
       100 +ve ⎭

↑ TPR = 94%.

↑ TNR = $\dfrac{850}{900}$

FPR = $\dfrac{50}{900}$ ↓

FNR = 6% ↓

for a sensible model, TPR & TNR should be high &
FPR & FNR should be low.

dumb model :→ all -ve

→ actual



|  | | 0 | 1 |
|---|---|---|---|
| predicted | 0 | 900 TN | 100 FN |
| ↓ | 1 | 0 EP | 0 TP |

900 = N     P = 100

900 - ve
100 +ve

$x_q$ :→ -ve : '0'

totally low  ←✗  TPR = 0%.          FPR = 0%. ✓

✓ TNR = 100%.          FNR = 100%. ✗

for er:→   To diagnose a cancer patient if he/she has
cancer or not.

Domain specific

→ act



| pred. | 0 | 1 |
|---|---|---|
| ↓ 0 | TN | FN |
| 1 | FP | TP |

our objective is not to
miss any cancerous
patient.

V. low FNR should be there as, FNR
means to say that patient does
not have cancer but it actually has
cancer.

If FPR is high → that's okay also, as if a
patient doesn't have cancer and it predicted
to be cancer patient then he/she will go
through more powerful tests to prove him

non-cancerous.

, interpretable in english
sentences.

21.3
★ Precision, Recall & $F_1$-Score :→



$$Precision = \frac{TP}{TP + FP}$$

of all the points the model
declared / predicted to be +ve,
what %age of them are actually +ve

Recall :→ same as TPR $= \frac{TP}{P}$

of all the actual positive points, what %age of them,
are predicted positive by the model.

we want precision to be high & Recall should be high.

Precision ↑        Recall ↑
  (0-1)            (0-1)

$F_1$-Score is combination of Precision & Recall.

$$F_1\text{-Score} = \left( 2 * \frac{precision * recall}{precision + recall} \right)$$

imp for
kaggle competitions

harmonic mean
of precision &
recall

┌─────────────────────────────────────────┐
│ $F_1$-score ↑ if Precision ↑ & recall ↑  │
└─────────────────────────────────────────┘

# Reciever Operating characteristic curve & AUC

(ROC) curve. ↳ designed by electronics & radio engineers during second world war.

| $\hat{y}_{z_2}=0.9$ | $\hat{y}_{z_1}=0.95$ $x$ | | $y$ | $\hat{y}$ | |
|---|---|---|---|---|---|
| 1 | 1 | $x_1$ | 1 | 0.95 | $\rightarrow \tau_1$ |
| 1 | 0 | $x_2$ | 1 | 0.92 | $\rightarrow \tau_2$ |
| 0 | 0 | $x_3$ | 0 | 0.80 | |
| 0 | 0 | $x_4$ | 1 | 0.76 | |
| 0 | 0 | $x_5$ | 1 | 0.71 | |

model → 1, 0

↳ gives score (like prob-score)

① Sort your data in decreasing order of $\hat{y}$

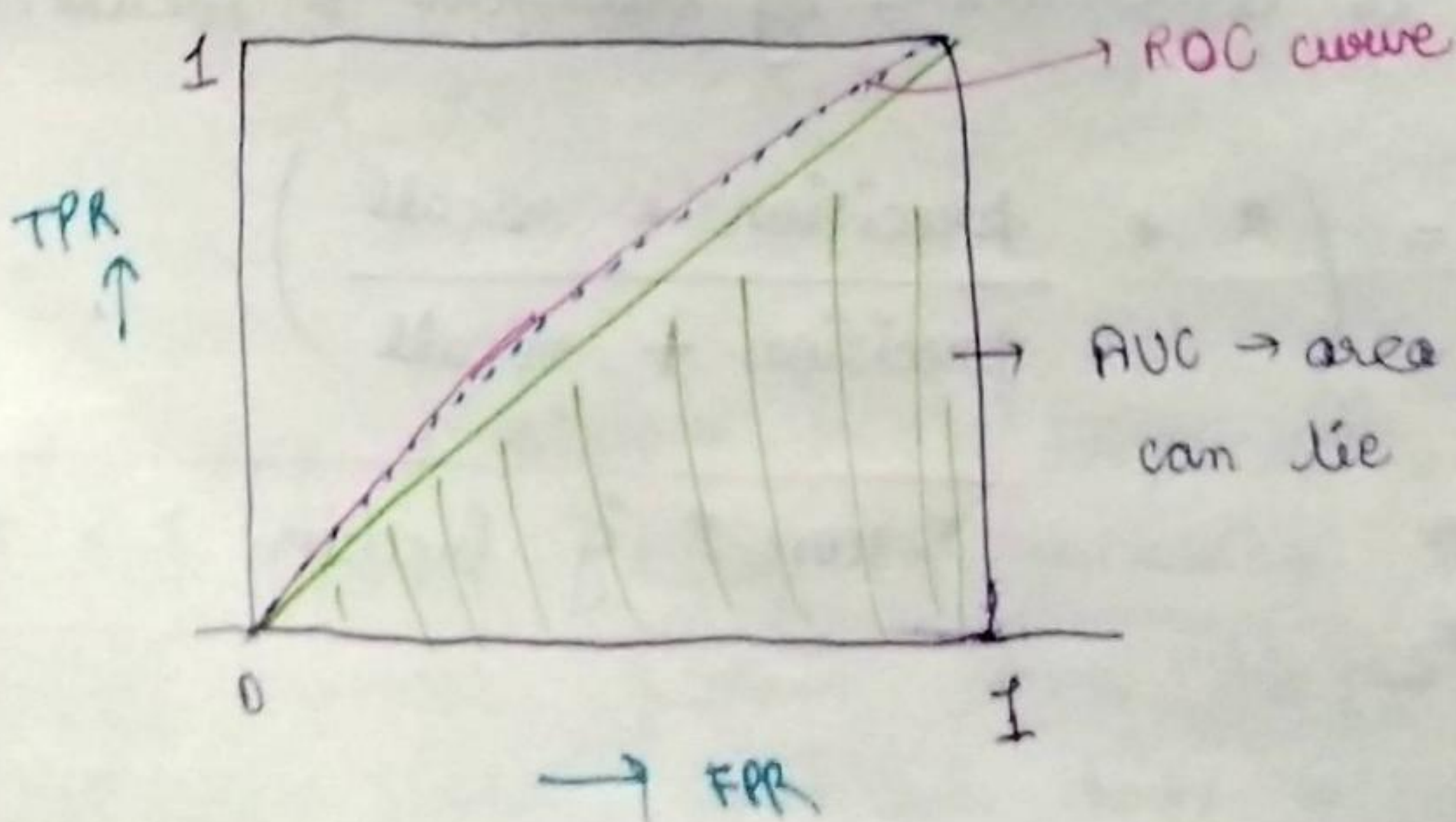② Thresholding → Take any value
($\tau$)
↓
tow

@ $\tau_1 = 0.95$

if $\hat{y} \geq \tau_1$ . ①

else ⓪

For each $\tau_1, \tau_2, \tau_3, \ldots \ldots \tau_n$
↓       ↓ $FPR_2, TPR_2$ and so on.
we can easily get

$FPR_1, TPR_1$



→ ROC curve

AUC → area under curve
can lie b/w ⓪ to ①
↓           ↓
trouble      v-good.

TPR ↑

→ FPR

$\boxed{AUC}$ :→

① Imbalanced data → AOC can be high with a dumb model too / Simple model.

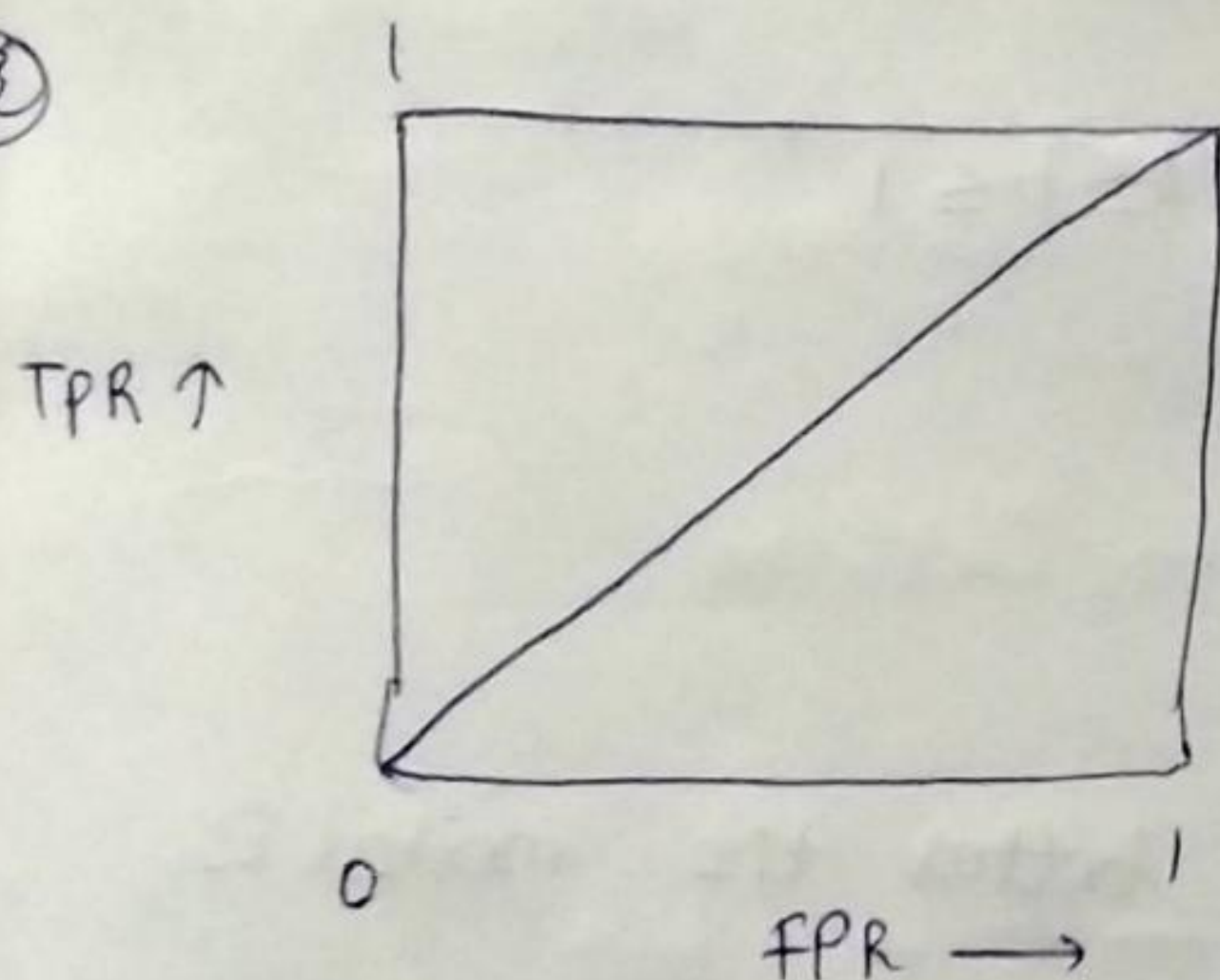② AUC is not dependent on the $\hat{y}$ scores but it depends on the decreasing order.

|       |   | $M_1$ | $M_2$ |
|-------|---|-------|-------|
| $x_1$ | 1 | 0.95  | 0.2   |
| $x_2$ | 1 | 0.92  | 0.1   |
| $x_3$ | 0 | 0.80  | 0.08  |
| $x_4$ | 1 | 0.76  | 0.07  |
| $x_5$ | 1 | 0.71  | 0.06  |

$AUC(M_1) = AUC(M_2)$

as AUC only depends on order.
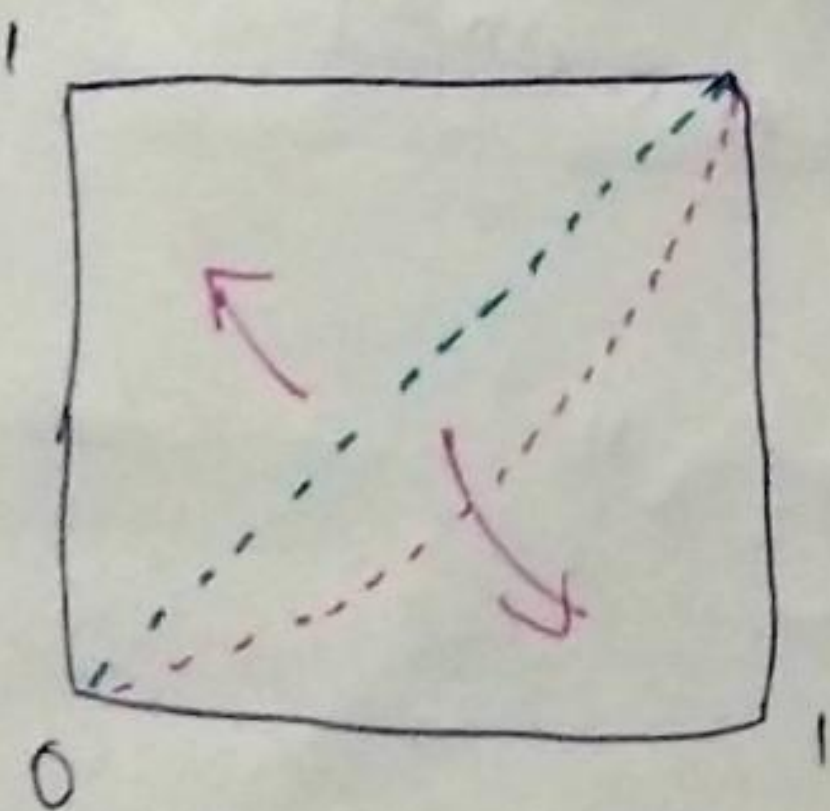
and not at all depend on actual value

③



TPR ↑

0          FPR ⟶

random model

↳ $x_q$ → 1 or 0

$\left\{ \begin{array}{c} AUC\ (random - model) \\ = 0.5 \end{array} \right\}$

④ Model M :→



1

0

$\hat{y}_i = 0 \longrightarrow 1$  } swa
$\hat{y}_i = 1 \longrightarrow 0$  } -pping

let $AUC(M) = 0.2$ → means you did something wrong in modelling

↳ worse than random model.

AUC →

0.5 to 1 ⟶ ✓

0.5 → random model

0.0 to 0.5 → worst
↳ swap class labels, $AUC = 1 - 0.2$

21.5 $\rightarrow$ (0 to $\alpha$)

Log-loss :$\rightarrow$ uses prob-scores.

as small
as
possible

Binary classification :$\rightarrow$

| $x$ | $y$ | $\hat{y} = P$ | |
|---|---|---|---|
| $x_1$ | 1 | 0.9 | $\rightarrow -\log(0.9)$ = 0.0457 |
| $x_2$ | 1 | 0.6 | |
| $x_3$ | 0 | 0.1 | $\rightarrow -\log(0.9)$ |
| $x_4$ | 0 | 0.4 | $\hookrightarrow 0.0457$ |

$-\log(0.6)$ = 0.22 $\leftarrow x_2$

Test set of n-
points

$-\log(0.6)$ = 0.22 $\leftarrow x_4$

$$\text{log-loss} = -\frac{1}{n}\sum_{i=1}^{n}\left\{\left(\log(P_i)\times y_i\right) + \left((1-y_i)\times \log(1-P_i)\right)\right\}$$

log-loss :$\rightarrow$ average of negative log prob of correct class
label.



$0 \leq p \leq 1$

smaller the log loss, better the model is.

Multiclass log loss :$\rightarrow$

$$-\frac{1}{n}\sum_{i=1}^{n}\sum_{j=1}^{c} y_{ij} \log(P_{ij}) \rightarrow \text{probability that } x_i \in \text{class } j$$

$y_{ij} = 1$ if $x_i \in$ class $j$

otherwise 0

Best case of log loss is 0, but other values are not easy
to make sense.

## $R^2$ (or) coefficient of determination :→

For regression

Test :→ $\quad\quad y_i \in R$

$i = 1$ to $n$

$\quad\quad\quad x_i \,,\, y_i \,,\, \hat{y}_i \quad\quad\quad , \quad \boxed{e_i \Rightarrow y_i - \hat{y}_i}$

$\quad\quad\quad\quad\quad \uparrow \quad\quad \uparrow \quad\quad\quad\quad\quad\quad \downarrow \text{error}$

$\quad\quad\quad\quad\quad$ actual

$\quad\quad\quad\quad\quad$ output  model predicted

$\quad\quad\quad\quad\quad\quad\quad\quad$ output

$$\boxed{SS_{\text{Total}} = \sum_{i=1}^{n} \left(y_i - \bar{y}\right)^2}$$

$\underbrace{SS}_{\text{Sum of squares}}$ Total

$\quad\quad\quad\quad\quad\quad\quad\quad \underset{\text{actual}}{\uparrow} \quad \underset{\substack{\text{mean} \\ \text{of actual} \\ \text{class labels}}}{}$

, here $\bar{y} = \dfrac{1}{n} \sum_{i=1}^{n} y_i$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \bar{y} \rightarrow$ avg. value of $y_i$'s

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ in test data.

For regression,

$\quad\quad\quad$ the simplest model you can build is

$\quad\quad\quad\quad\quad$ called the average model $\rightarrow$ simple

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ mean model

$\quad\quad\quad\quad\quad x_q \rightarrow$ mean$(y_i) = \bar{y} = y_q$

$SS_{\text{total}} \rightarrow$ sum of squared errors using simple-mean

$\quad\quad\quad\quad\quad\quad\quad$ model.

$\underset{\text{residues}}{}$

$$\boxed{SS_{\text{res}} = \sum_{i=1}^{n} \left(y_i - \hat{y}_i\right)^2 = \sum_{i=1}^{n} e_i^2}$$

$\quad\quad\quad\quad$ residue $= e_i = y_i - \hat{y}_i$

$\quad\quad\quad\quad\quad\quad\quad\quad \underset{\text{actual}}{\downarrow} \quad \underset{\text{predicted}}{\searrow}$

$$R^2 = \left(1 - \frac{SS_{res}}{SS_{tot}}\right)$$

Case 1 :→ $SS_{res} = 0 \leftarrow \boxed{e_i = 0} \rightarrow R^2 = 1$

phenomenal
model

$\downarrow$

Best value

Case 2 :→ if $SS_{res} < SS_{tot}$ ; $R^2 = 0$ to $1$

Case 3 :→ if $SS_{res} = SS_{tot}$ ; $R^2 = 1-1 = 0 \rightarrow$ model is
same as simple
mean model.

Case 4 :→ if $SS_{res} > SS_{tot}$ ; $R^2 = 1 - (greater > 1) = -ve$
model is worse than simple
mean model.

21.7

## Median absolute deviation of errors :→

$$SS_{res}^2 = \sum_{i=1}^{n} e_i^2 \qquad \text{if one } e_3 \text{ is very large}$$

$R^2$ is not very robust to outliers.

$x_i \longrightarrow y_i, \hat{y}_i, e_i \qquad , \ |e_i|_s \rightarrow 0 \rightarrow great$

$|e_i|_s \longrightarrow large \rightarrow not\ so\ good$.

if $e_i \rightarrow$ random variable,

mean $\leftarrow$ median $(e_i) \leftarrow$ central value of errors $\rightarrow$ small

acting
like these

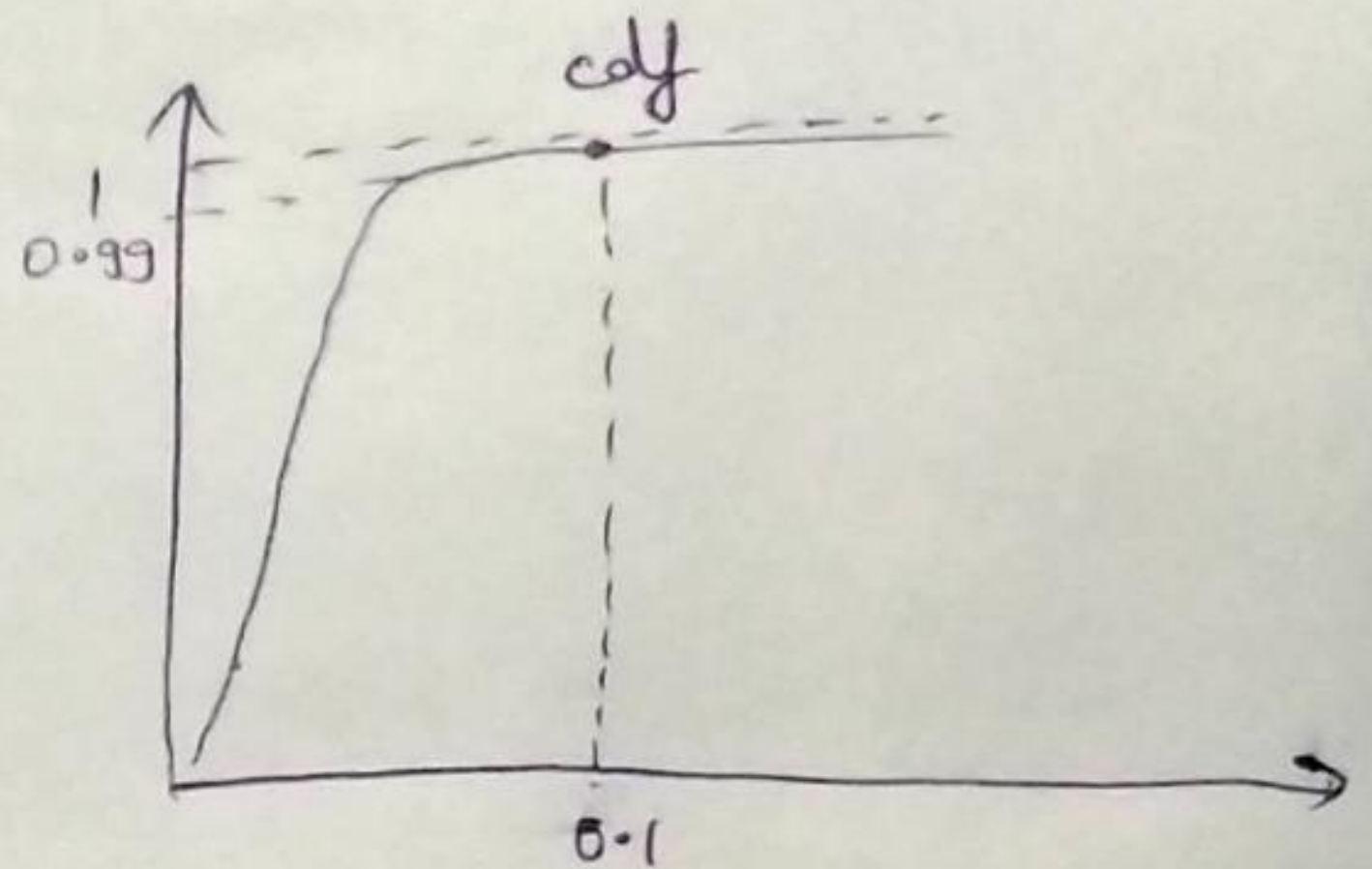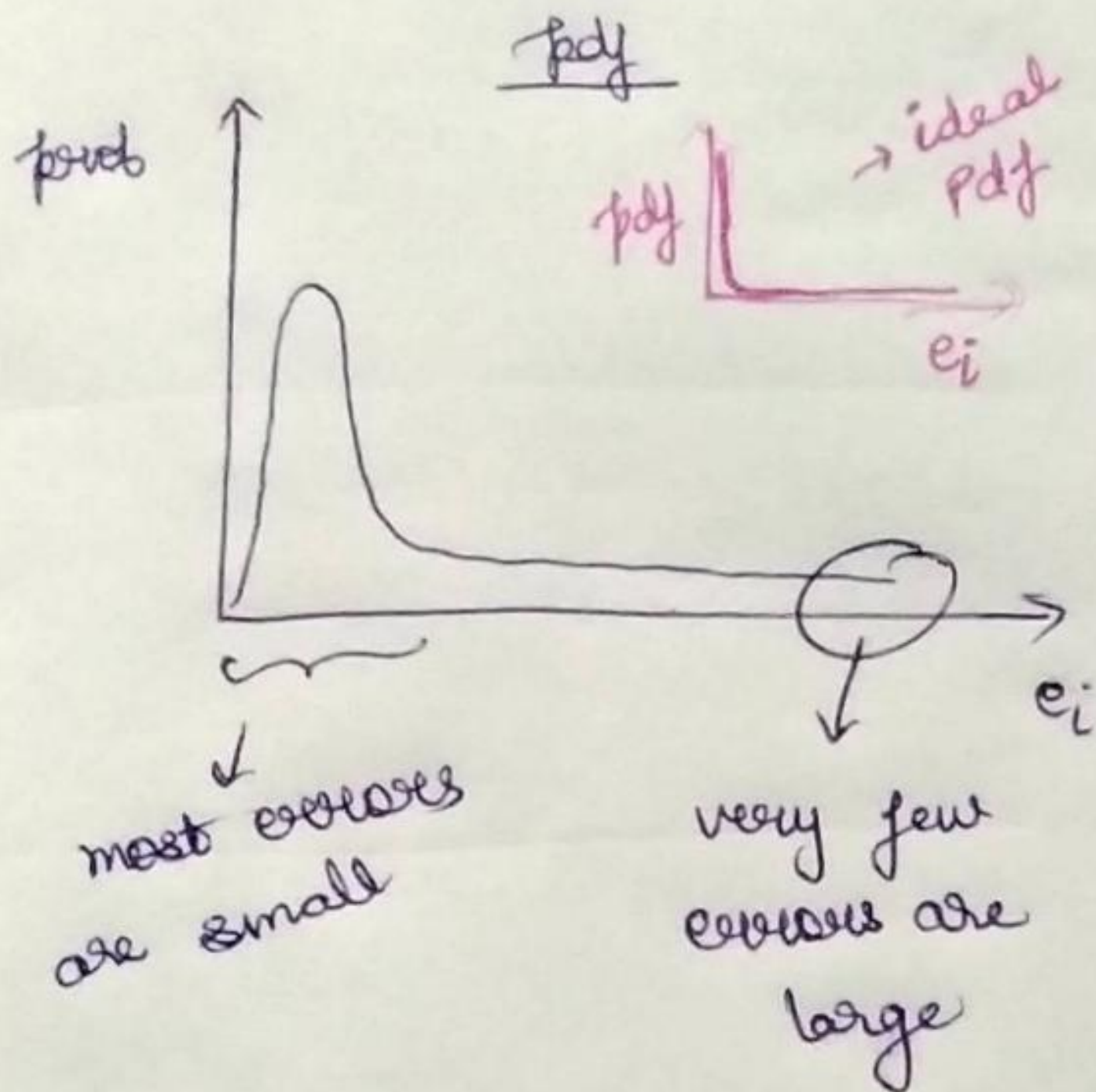std.
dev $\leftarrow$ MAD $(e_i) =$ median $\left( \underbrace{|e_i - \underbrace{median(e_i)}_{deviation}|}_{abs} \right) \rightarrow$ small.

mean or (median) of eis → used to understand if
std-dev or MAD the errors are small
or large.
↓
robust to outliers

21.8 Distribution of errors :-

pdf

prob

pdf → ideal pdf

$e_i$

↓ (most errors are small)

very few errors are large

$e_i$

cdf

0.99

0.1

99% of errors are < 0.1
1% errors are ⩾ 0.1

Models $M_1$ & $M_2$

$M_1$  $M_2$

0.95
0.8

0.1  $e_i$

$M_2$ cdf is below $M_1$

$M_1$ :→ 95% errors are below 0.1

$M_2$ :→ 80% errors are below 0.1.

So, $M_1$ is a better model for regression than $M_2$