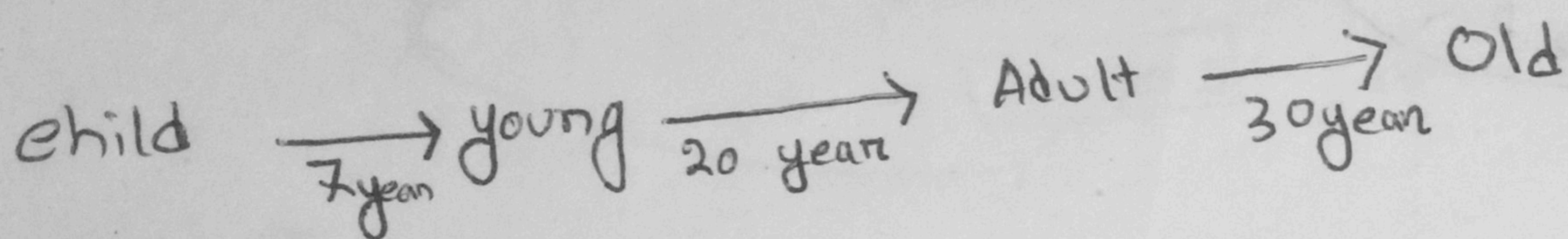


1

Machine learning

It is a branch of Artificial Intelligence
 Use of data and algorithms to imitate the way that human
 learn, gradually improving its accuracy

IBM



model₁
 (Algo) → m₂ → m₃ → m₄ (data driven)

Q. if input X rules define

if ($x_1 \cdot x_2 == 0$)
 #even

else
 #odd

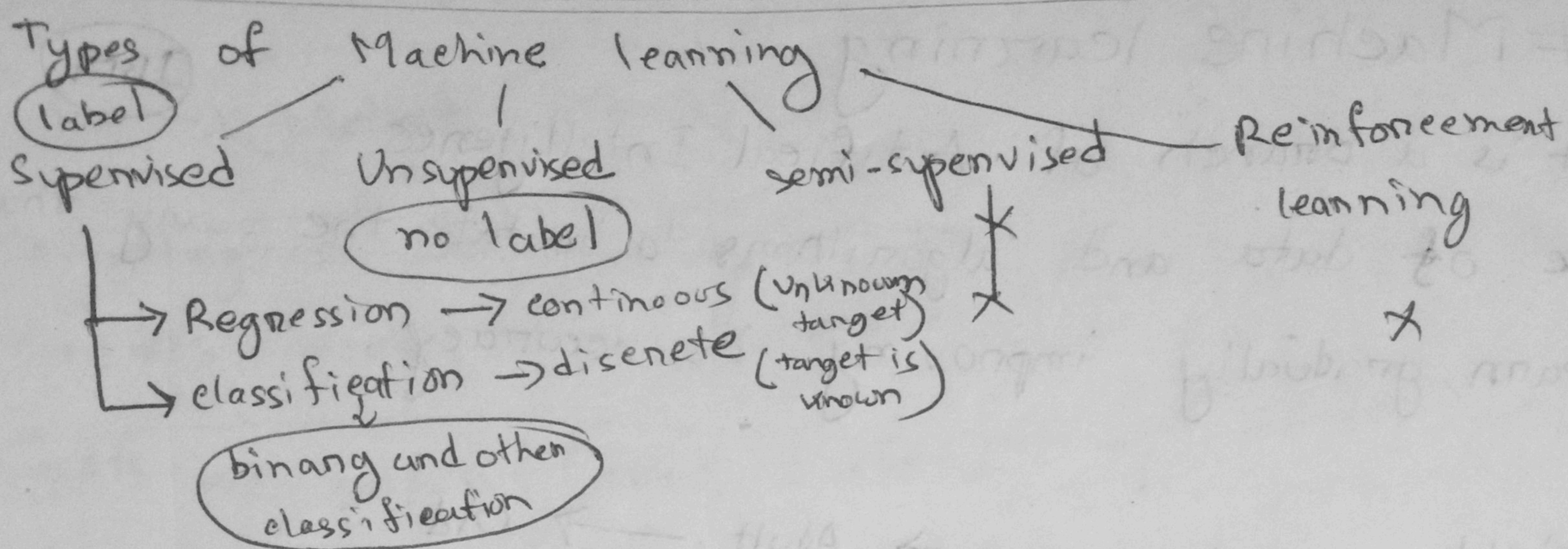
Floor	Location	Size	target	
			Rent	Cost
7	Banani	1500	50,000	
4	Tongi	2200	55,000	
2	Banani	1200	30,000	
5	Dhamondi	1500	40,000	
3	Ottawa	1800	??	

size > 1500 X
 location X
 floor X
 → Not maintaining
 rules

When we can not set
 the rules → we
 can use ML
 (Machine learning)

11

2



Example

Cat or Dog

	Features	Target
image 1		Cat
image 2		Dog
image 3		Cat
image 4		Dog

label → labeled image

student exam pass/fail

study hour	decision
0h	fail
1h	fail
2h	pass
3h	pass
4h	pass

0/1

Classification

target

CGPA [2-4] & Continuous

3.14 3.42 2.92

2.03 1.23 4.00

Rent east (0-30,000) any value

12K

18K

target rent cost

possible
no idea

continuous → Regression

discrete → classification

↳ cat or dog

3

Regression

Size	Rent cost (e^{-10})
1200	25,000 (5×10^3)
1500	30,000 (6×10^3)
1800	35,000 (7×10^3)
2000	? ($10,000$)

linear regression

$$f(x) = 2x + 3$$

↓
dependent (target)

$\boxed{\text{Rent cost} = 20 \times \text{Size}}$ hypothesis

$$\boxed{R = 20}$$

Randomly

Prediction

$$R = 20 \times 1200 = 24,000$$

$$R = 20 \times 1500 = 30,000$$

$$R = 20 \times 1800 = 36,000$$

$$25,000 \quad L = 1000$$

$$30,000 \quad L = 0$$

$$35,000 \quad L = 1000$$

$$\frac{2000}{3}$$

$$= 666.66$$

$$\underline{\text{Loss}}$$

$$\text{total loss} = 2000$$

$$\text{Loss} = \frac{1}{N} \sum (\text{Actual} - \text{Prediction})$$

Loss $\downarrow \rightarrow$ Accuracy \uparrow

$$\text{Loss} = 2000 \quad (3 \text{ sample})$$

We want to minimize loss

sample	each	$\sum L$
$m_1 \quad N \quad 3$	1000	3000
$m_2 \quad N \quad 1000$	5	5000

↑ Dataset \rightarrow Loss \downarrow

$$\text{Loss} = \frac{1}{N} \sum | \text{Actual} - \text{Predicted} |$$

$$\Rightarrow = \frac{3000}{3}$$

$$\Rightarrow = \frac{5000}{5000}$$

state of the art

$$2018 \rightarrow \text{Loss} = 1000 \times$$

$$2020 \rightarrow \text{Loss} = 100 \times$$

$$2023 \rightarrow \text{Loss} = 10 \times$$

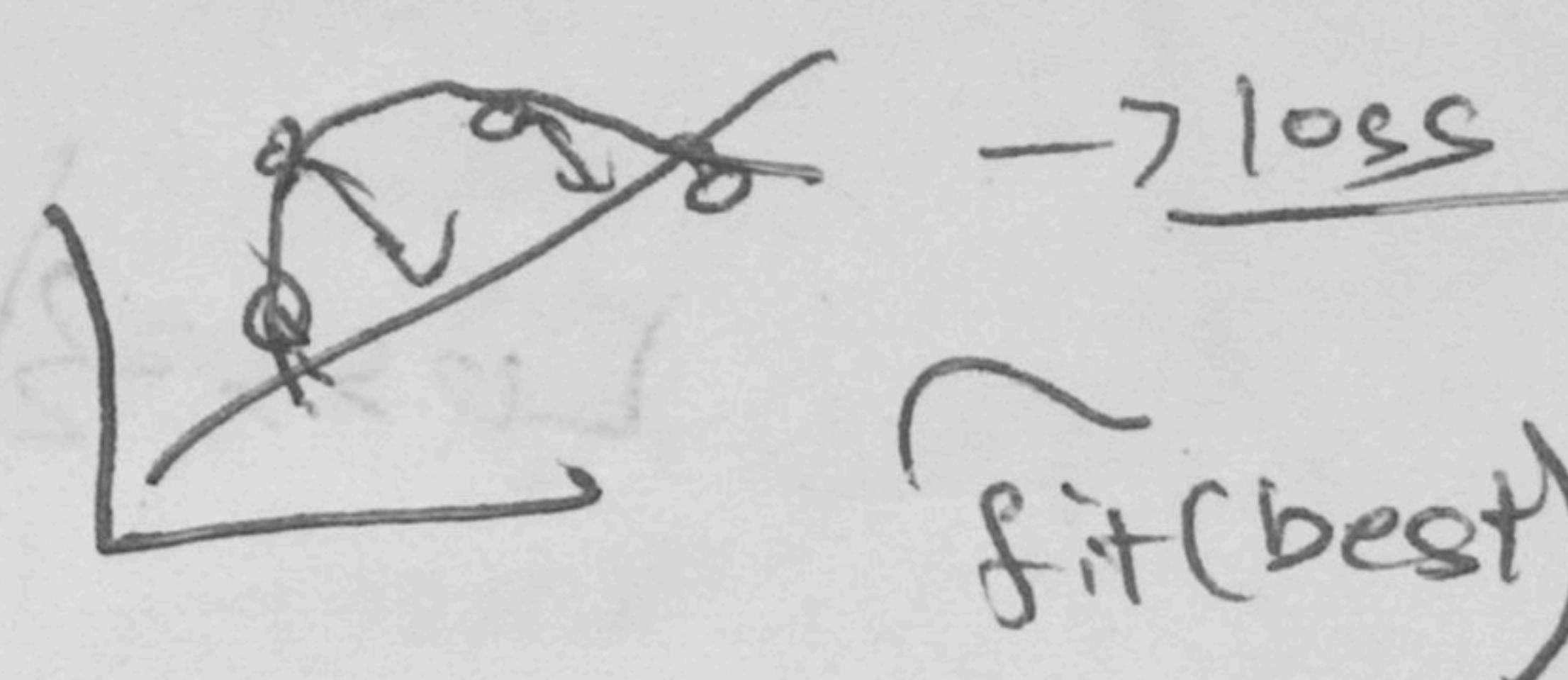
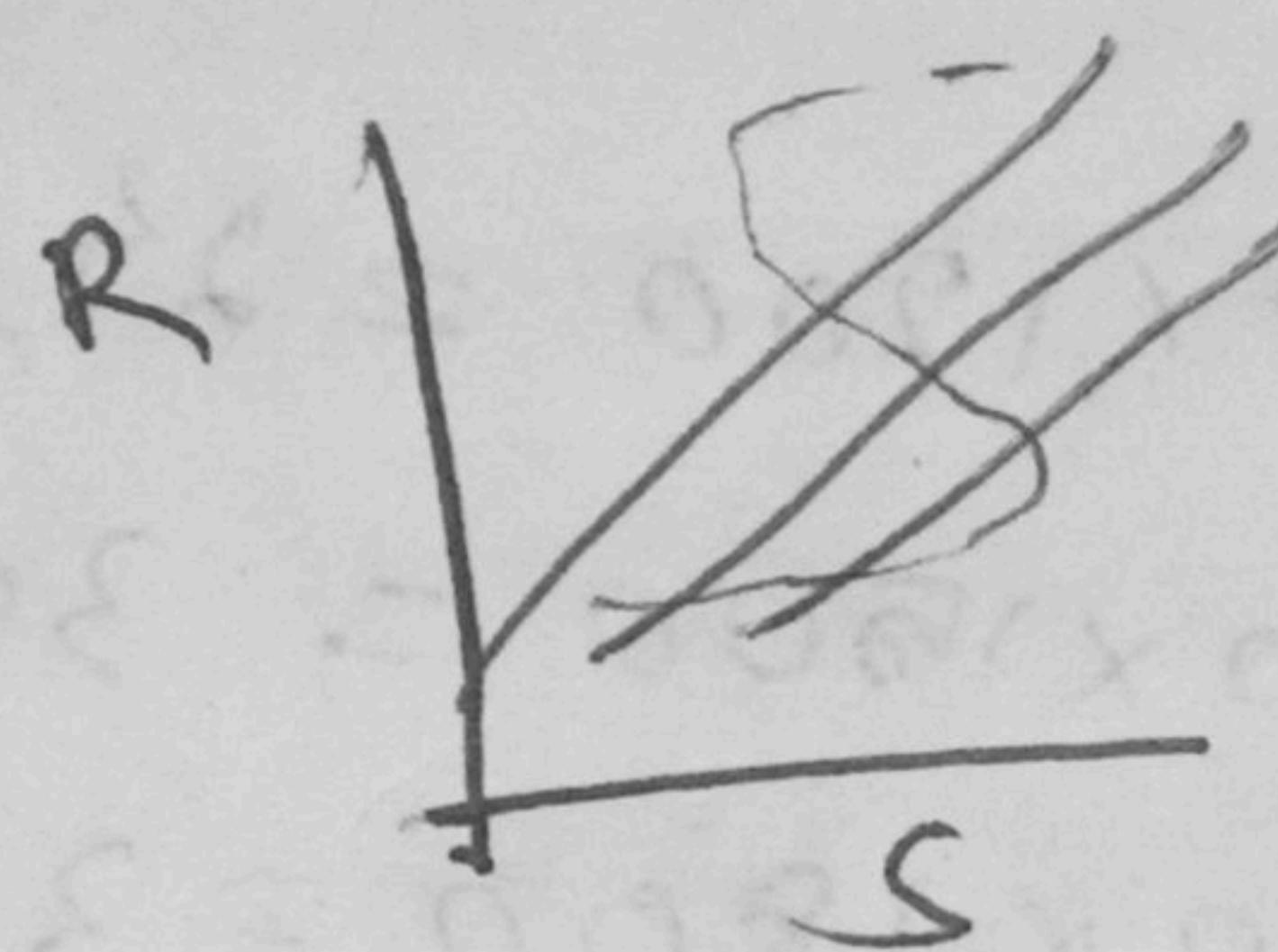
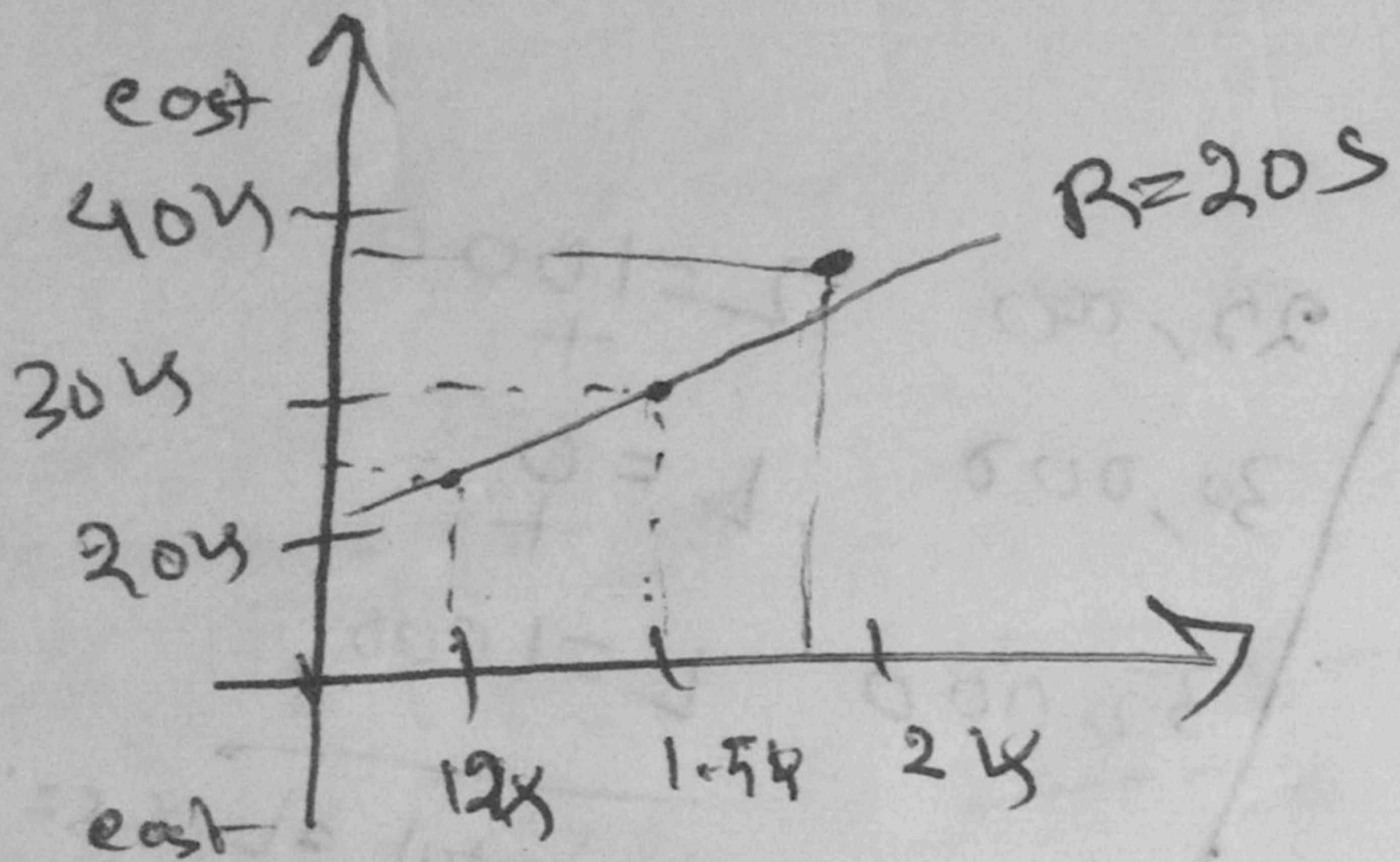
$$R=20S$$

$$y = 20R + 3$$

$$y = 2u^2 + 3$$

$$y = u^{1/2}$$

which is linear $\rightarrow y = mu + c$



$$\boxed{R=20S} \quad \text{loss} = -650 \times \downarrow \text{Reduce}$$

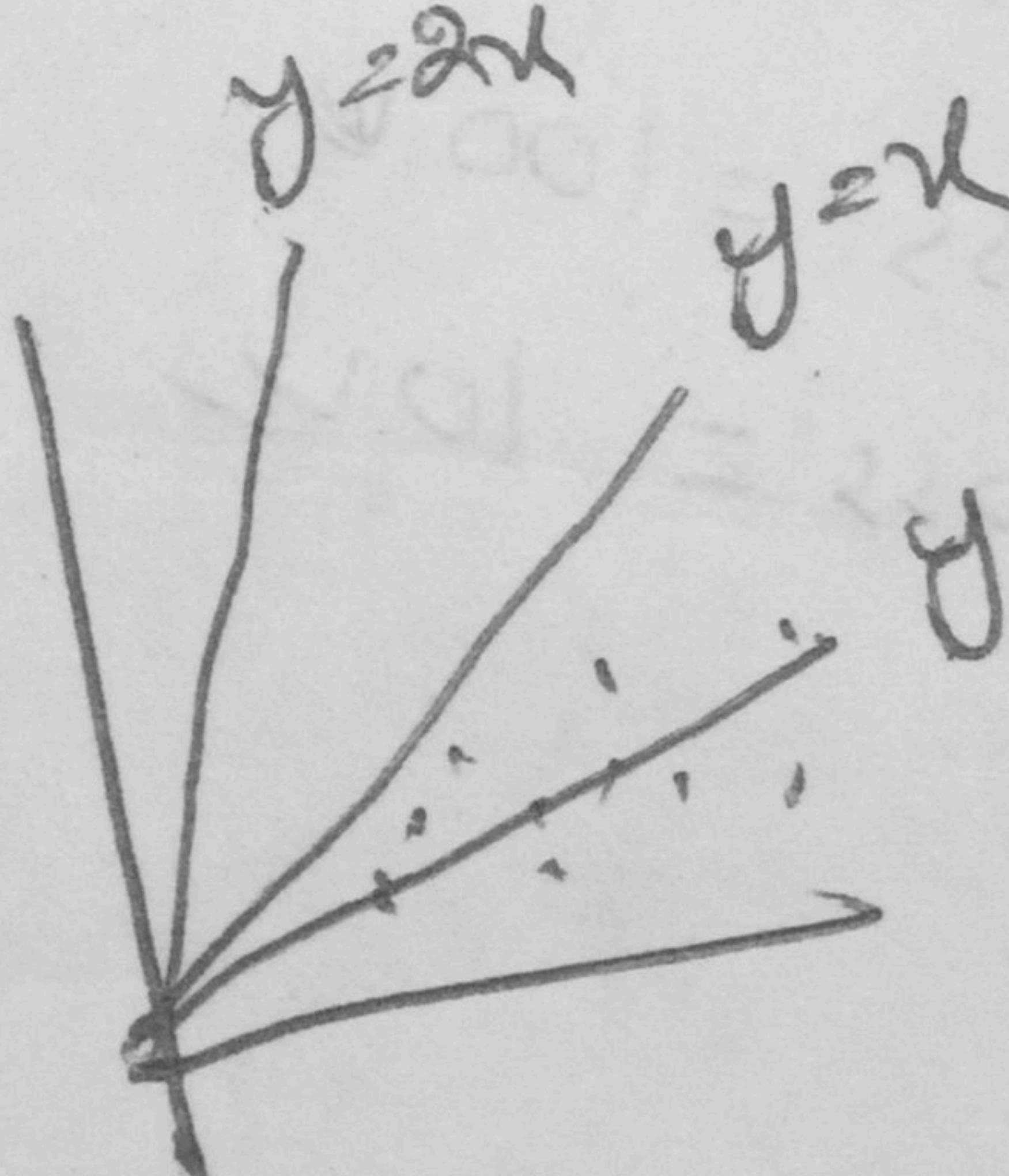
$$\min(\text{loss})$$

$$R = \underline{\underline{20S}} \quad \begin{array}{l} \xrightarrow{\text{tangent}} \text{weight} \\ \xrightarrow{\text{y} = w^T x \rightarrow \text{feature}} \end{array}$$

gradient descent δ

$$w_{\text{new}} = w_{\text{old}} - \alpha \frac{\partial}{\partial w} (\text{loss})$$

partial derivative



(Best fits with line)

Size	y price	\hat{y} prediction
1000	20000	21000
1200	25000	27000
1800	35000	31000

A
MAE (mean absolute error)

$$MAE = \frac{1}{N} (y - \hat{y})$$

$$= \sum \frac{1}{N} |y - \hat{y}|$$

$$= \frac{1000 + 2000 + 4000}{3}$$

$$= \frac{7000}{3}$$

$$\sqrt{\frac{RMSE^2}{BDF}} = \sqrt{\frac{MSE}{BDF}}$$

MSE (Mean square error)

$$RMSE = \sqrt{\frac{1}{N} (y - \hat{y})^2}$$

$$= \sqrt{\frac{1000^2 + 2000^2 + 4000^2}{3}}$$

$$= \sqrt{\frac{21 \times 10^6}{3}}$$

$$= \dots$$

R²

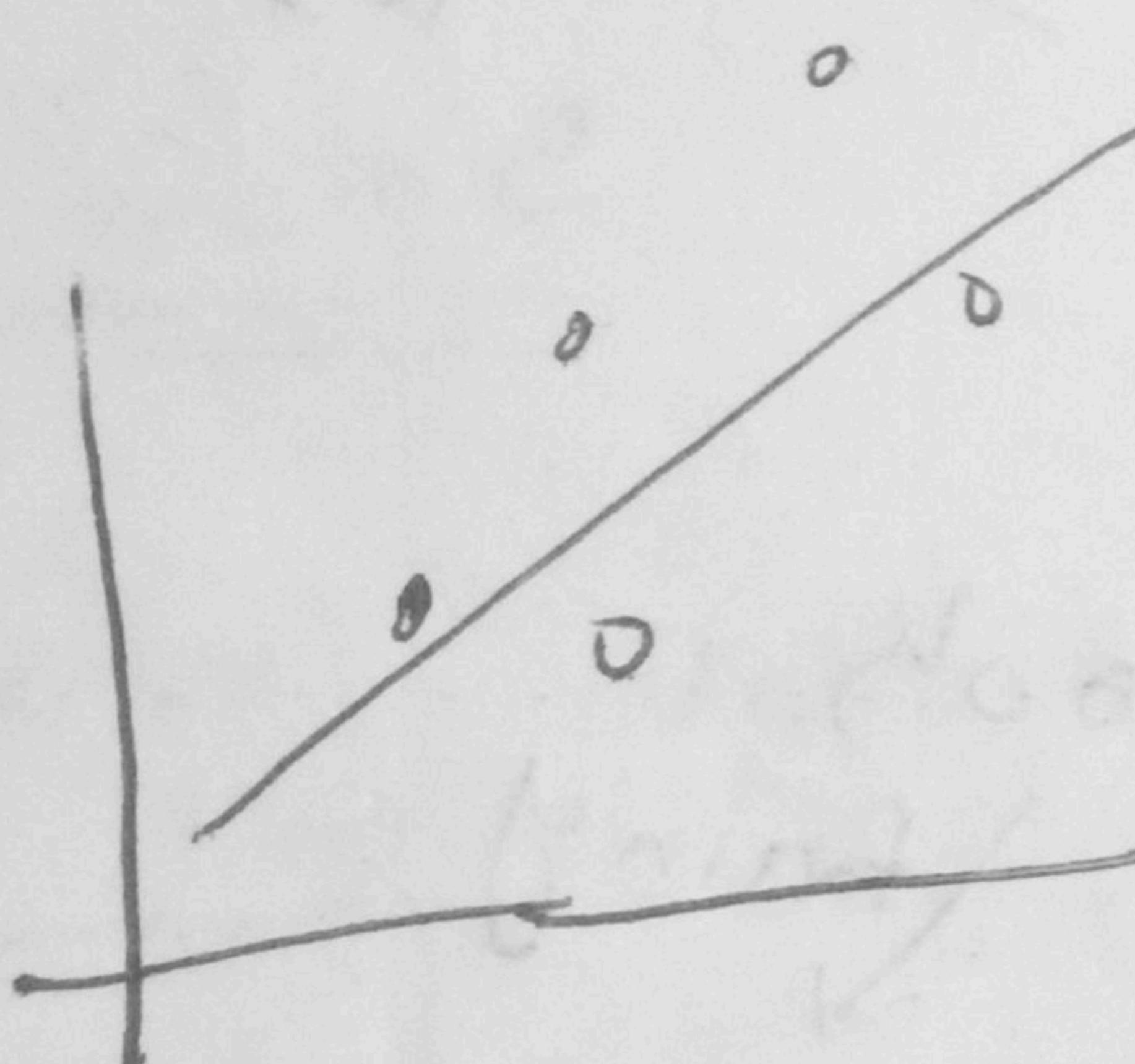


Fig 1

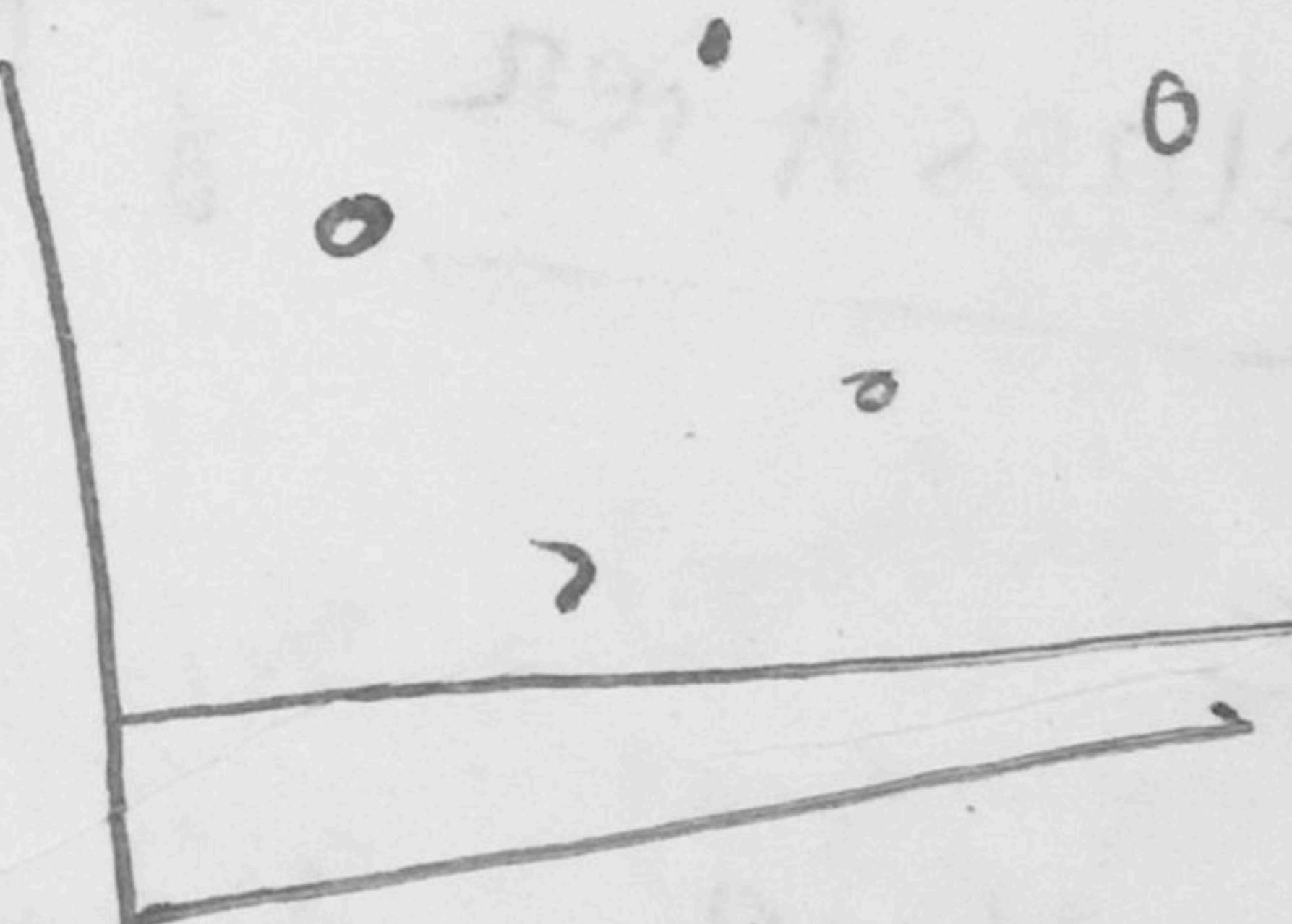


Fig 2

$$0 \leq R^2 \leq 1$$

$R^2 = 1 \Rightarrow$ perfectly fitted

$R^2 = 0 \Rightarrow$ poorly fitted

$$R^2 = \frac{\sum (y - \bar{y})^2}{\sum (y - \hat{y})^2}$$

$$R^2 = 1 - \frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2}$$

\bar{y} - Average

\hat{y} - Predicted

y - Actual

$R^2 = 0.8$
$R^2 = 0.2$

Classification

Index	Weather	Outlook	Play
1	hot	sunny	Y
2	cold	sunny	Y
3	cold	rainy	N
4	hot	sunny	Y
5	cold	sunny	Y

Predict
Yes → No
Outcome

Zero-R classifier

↳ count majority

yes → 4

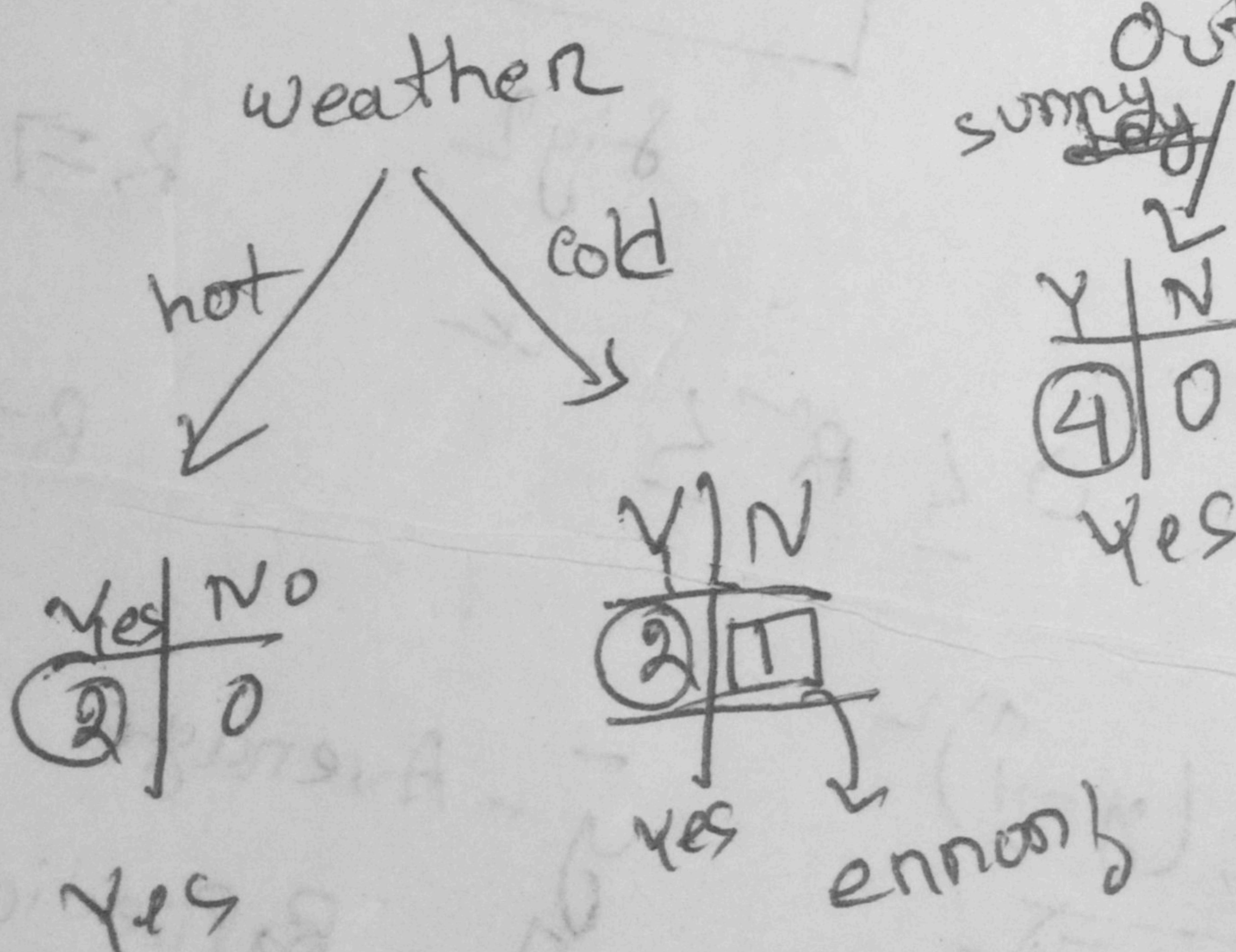
No → 1

prediction All → Yes.

One-R classifier of Outlook

outlook
sunny
Yes
Rainy
No

Feature



outlook

sunny	Rainy
Y N	Y N
④ 0	0 ①
Yes	No

If $Y = N$
Random choice

Evaluation metric

weather	outlook		play
	S	R	
S	S	R	Y
S	R	R	N
R	S	S	Y
R	R	R	N

ZenoR

prediction \rightarrow "Yes"

$$\text{Accuracy} = \frac{1+1+0+1+1}{5} = \frac{4}{5} = 0.8 = 80\%$$

1000 apple

950 apple good }
50 apple bad }

$$\frac{950}{1000} = 95\%$$

→ model → "good" factory

Confusion metric

Actual value		predicted value	
		P	N
predicted value	P	TP 2	FP 2
	N	FN 3	TN 6

$$\text{Total} = (2+2+3+6)$$

$$\text{precision} = \frac{2}{2+2} = \frac{1}{2}$$

positive - (good) P
negative - (bad) N

$$\textcircled{1} \quad \text{precision} = \frac{TP}{TP+FP}$$

$$\textcircled{2} \quad \text{recall} = \frac{TP}{TP+FN}$$

PP $\downarrow \rightarrow \uparrow$ precision

FP $\downarrow \rightarrow \uparrow$ precision

FP = 0 $\rightarrow \uparrow$ precision
(100%)

FN $\downarrow \rightarrow \uparrow$ recall

FN $\uparrow \rightarrow \downarrow$ recall

FN = 0 $\rightarrow \uparrow$ recall
(100%)

8

precision \rightarrow FP
 recall \rightarrow FN

(medical test) concern
 (Negative)

which metric should I use precision or recall?

Ans: Based on

@ Target variable (P, N)

For PP on FN

should be lesser than base line

value, for example (spam email detection)

$$F_1 \text{ score} = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

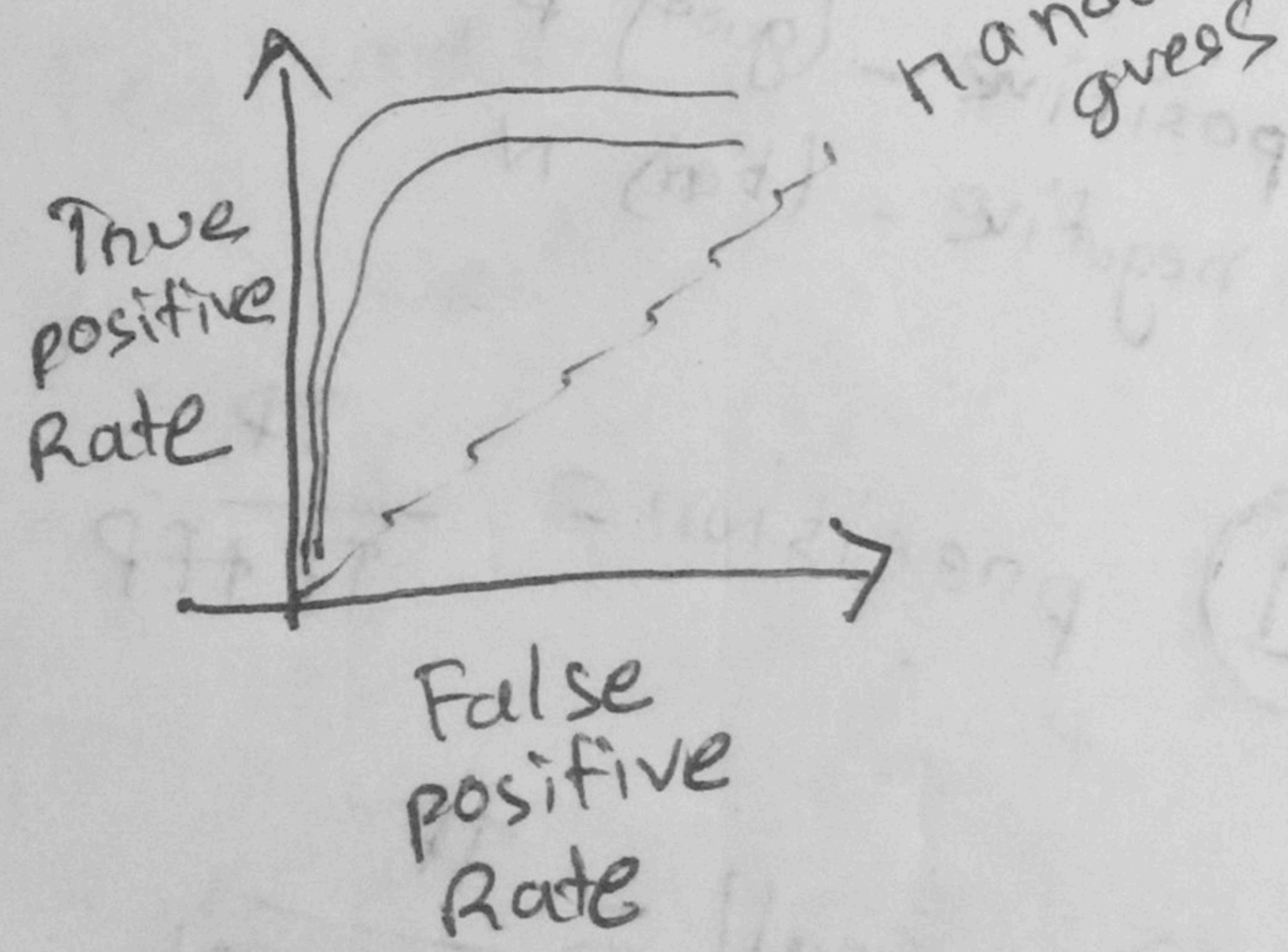
F_1 measure \rightarrow harmonic mean of precision and recall
 \hookrightarrow (close to smallest value)

\hookrightarrow close to $\frac{\text{mean}}{\text{positive non negative values}}$

Mean

ROC curve: Receiver characteristics curve

ROC



\rightarrow We want

\rightarrow Average

Under Random guess line, is considered 1000 values

AUC

curve:

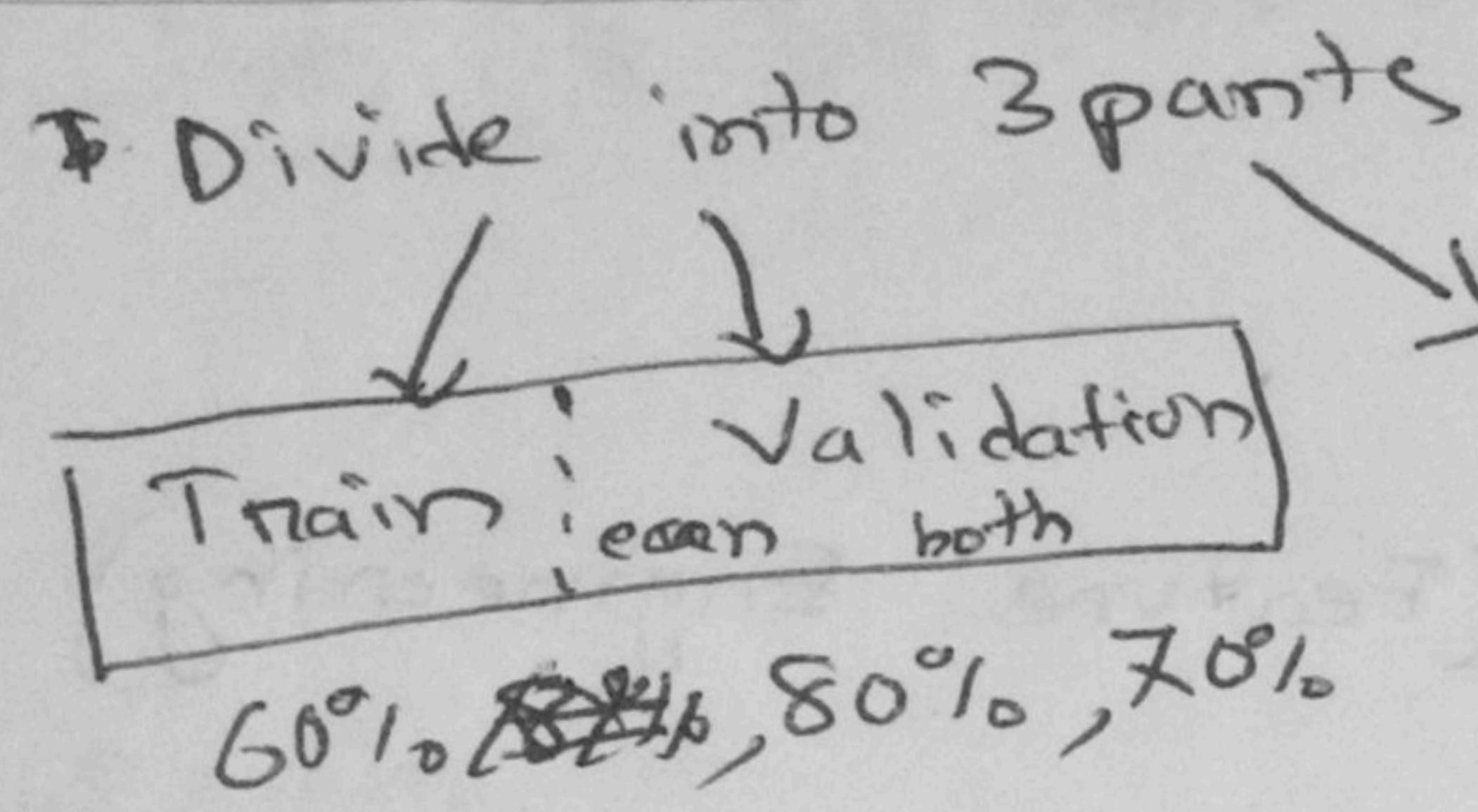
\hookrightarrow Area under curve

\Rightarrow must be high



\Rightarrow must be low w

F1	F2	Target
1	1	1
1	2	1
2	1	2
2	2	2
3	1	1
3	2	2
4	1	1
4	2	2
5	1	1
5	2	2



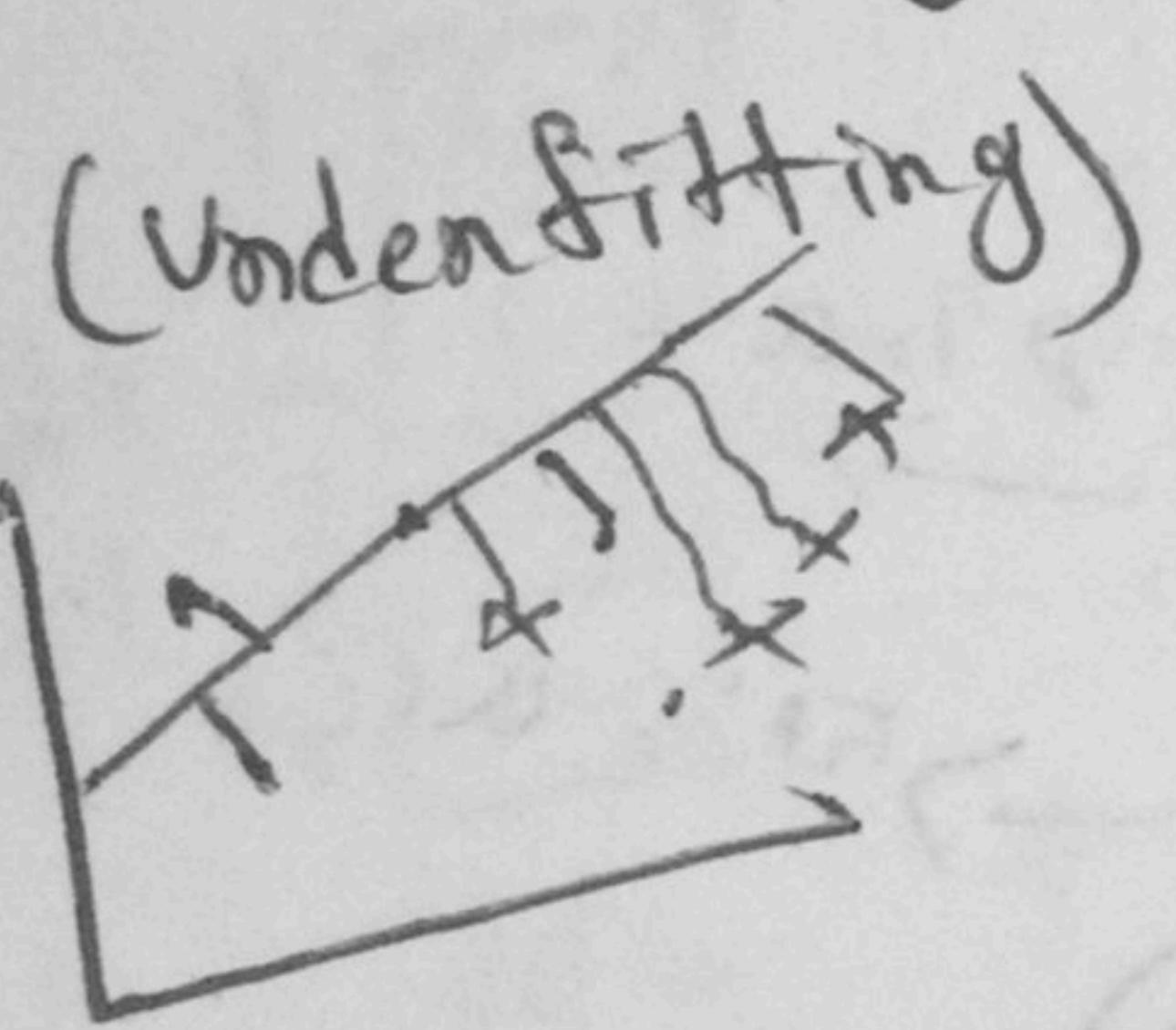
9

Train → Model → Train
Fitting Accuracy

Val → Model → Validation
Validation Accuracy

Test → Model → Test
Test Accuracy
must be unseen

Overshifting and Underfitting

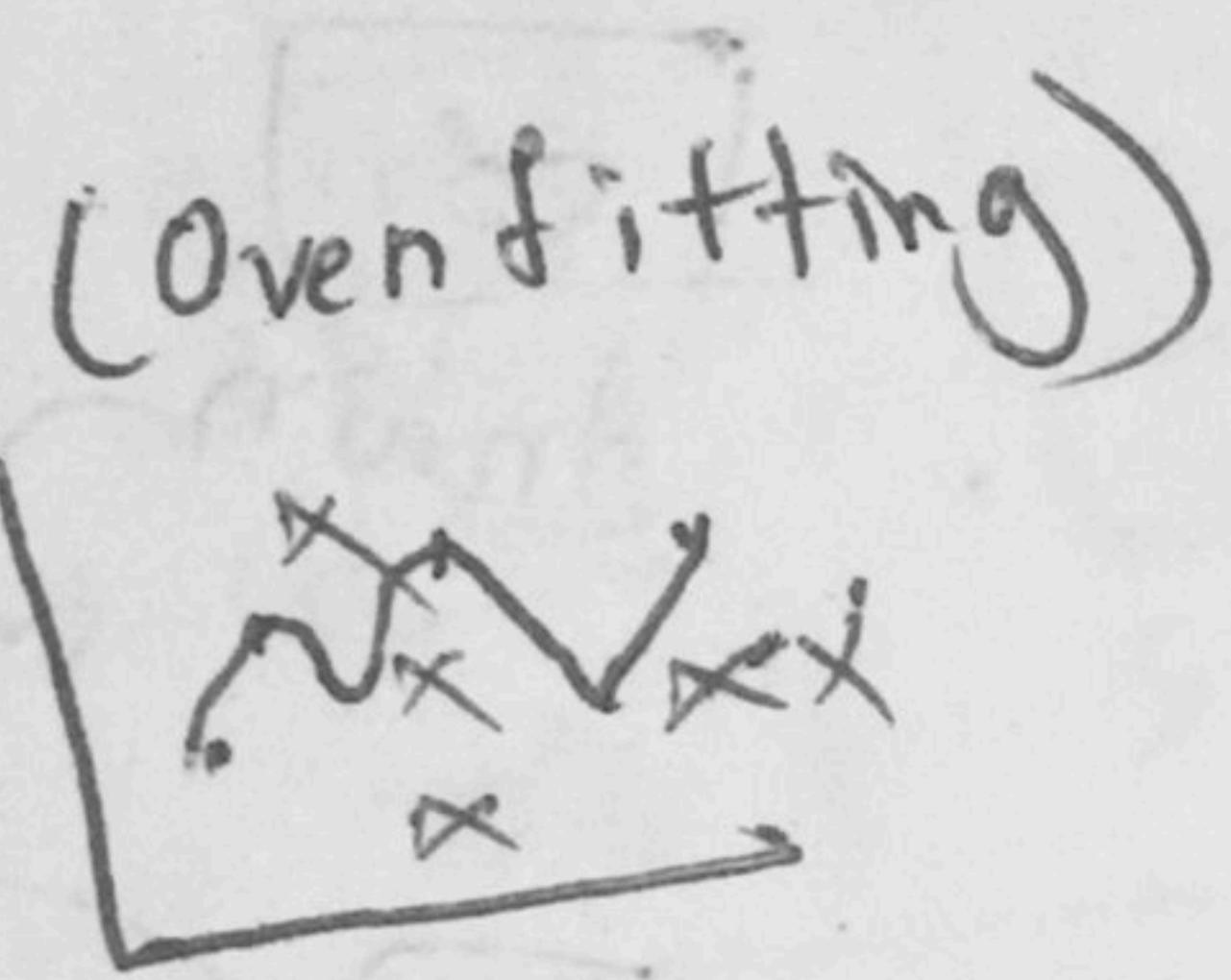


train → Low accuracy
test → low accuracy

(High bias, low variance)

Bias → Inaccuracy

Variance → Inconsistent (Overshifting)



train → High Accuracy

test → low Accuracy

(Mixed bias, High variance)

(Underfitting)



train → good Accuracy

test → good Accuracy

Perfect fit

(low Bias and low Variance)

Bias or Variance (Both cases)

Train acc → 90% | 80%
Test acc → 50% | 70%

underfit

good

30%	90%	85%	95%
28%	82%	83%	72%
underfit	overfit	overfit	perfect
good	good	overshifting	overshifting

if test accuracy is greater than train accuracy,
then different dataset, wrong dataset split, ~~or~~ High complexity of model

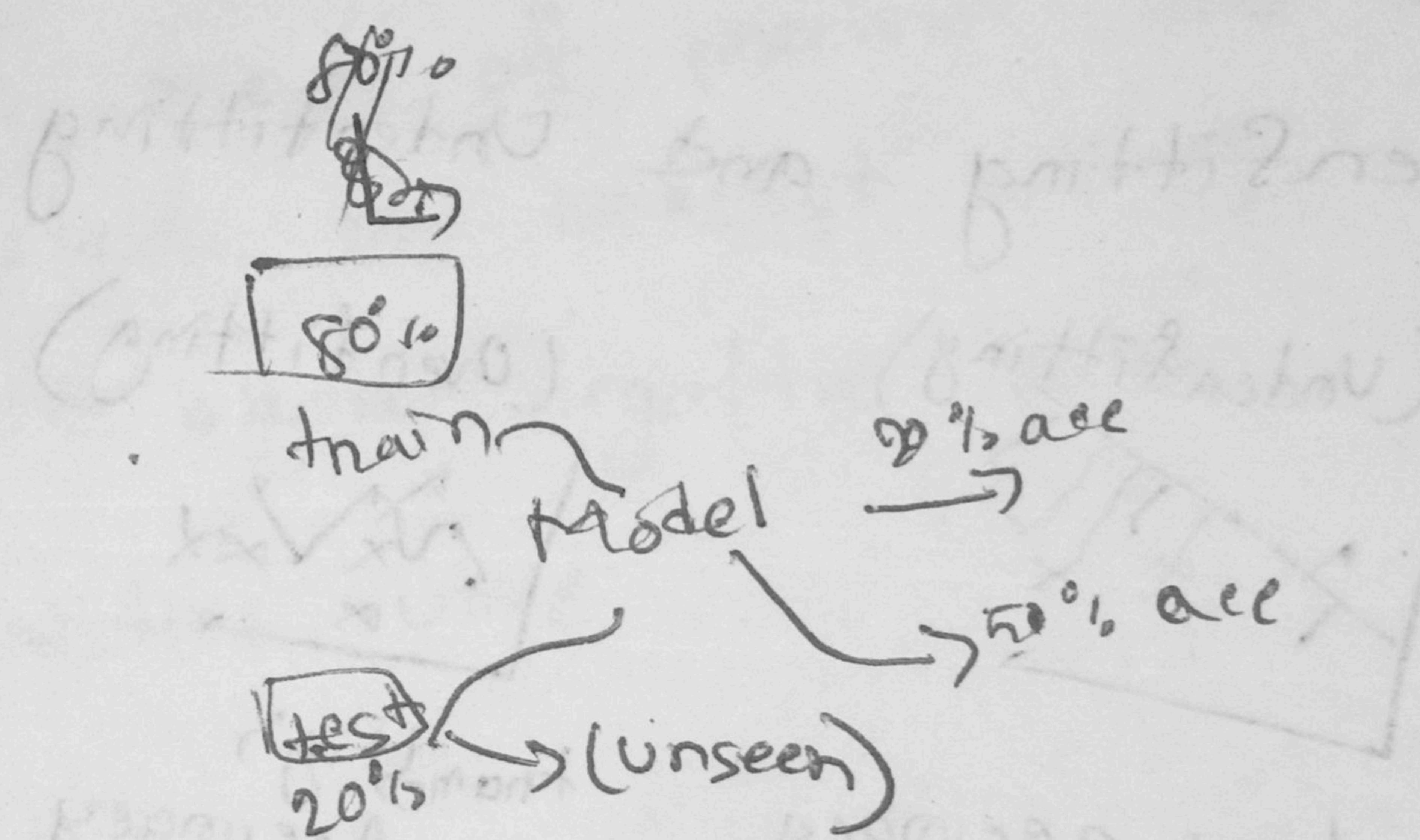
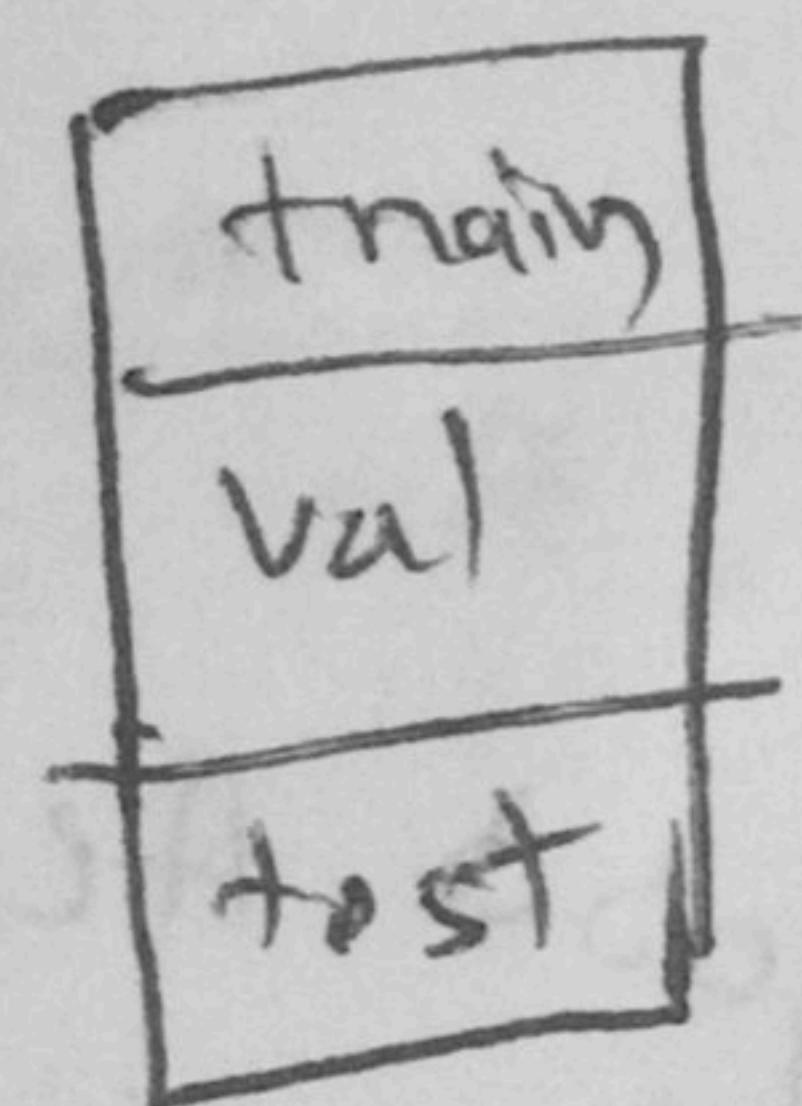
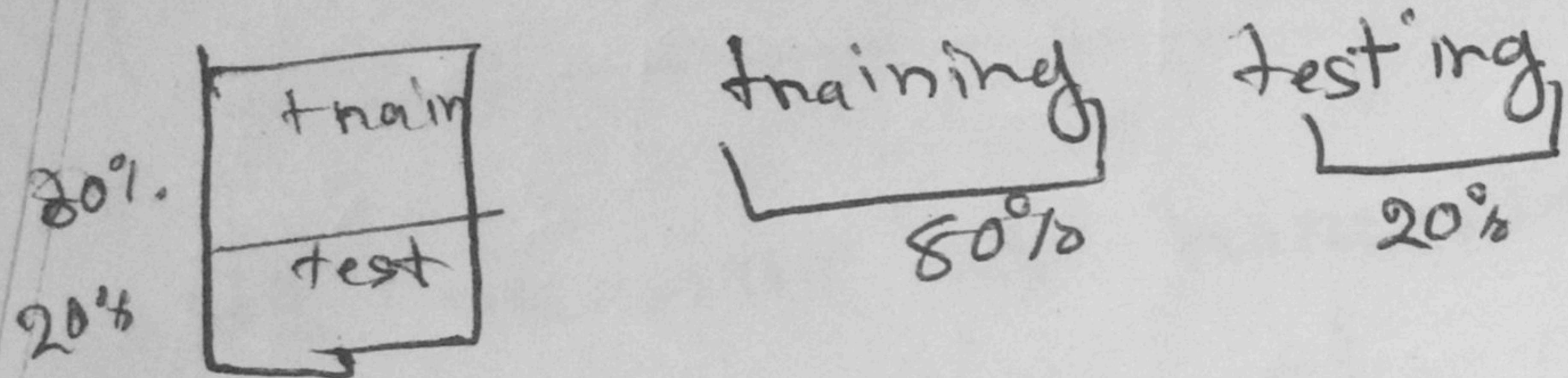
How to tackle

Underfitting

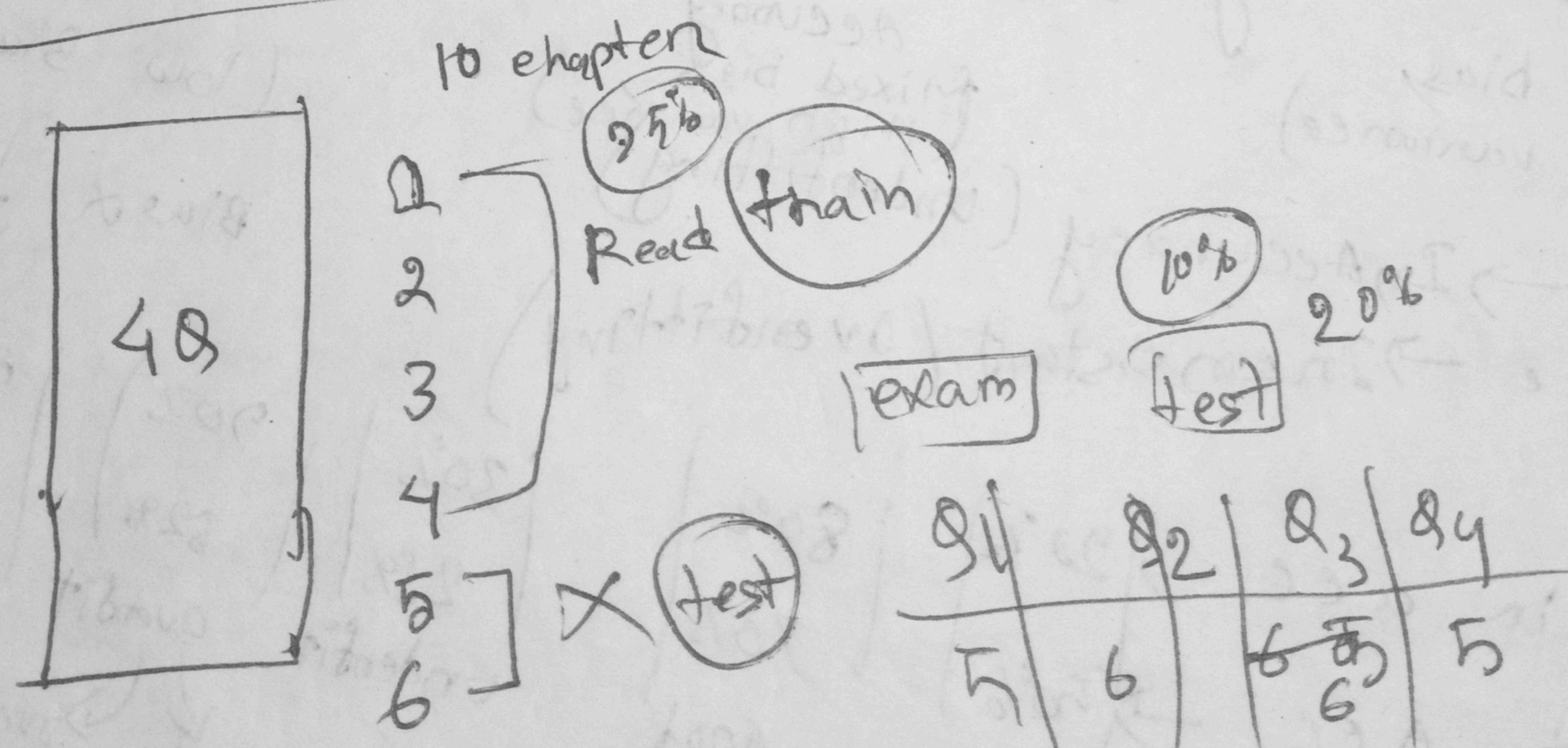
- ① Training time ↑
- ② Increase the features ↑ (Feature Engineering)

Oversetting

- ① Cross validation (Mandatory for cross validation)
- ② Regularization
- ③ Dropout
(Lasso, Ridge, elasticNet)



math program



LCI

II

