

Interview Prep Session 1

[Classical Machine Learning]

Warm Up Questions

Q: If you are training a LR to predict
Fahrenheit using Celcius as independent var.
what will be the value of the intercept?
(Quiz)

→ Automated | Easy

Solⁿ → We know →

$$F = \frac{9}{5}x + 32$$

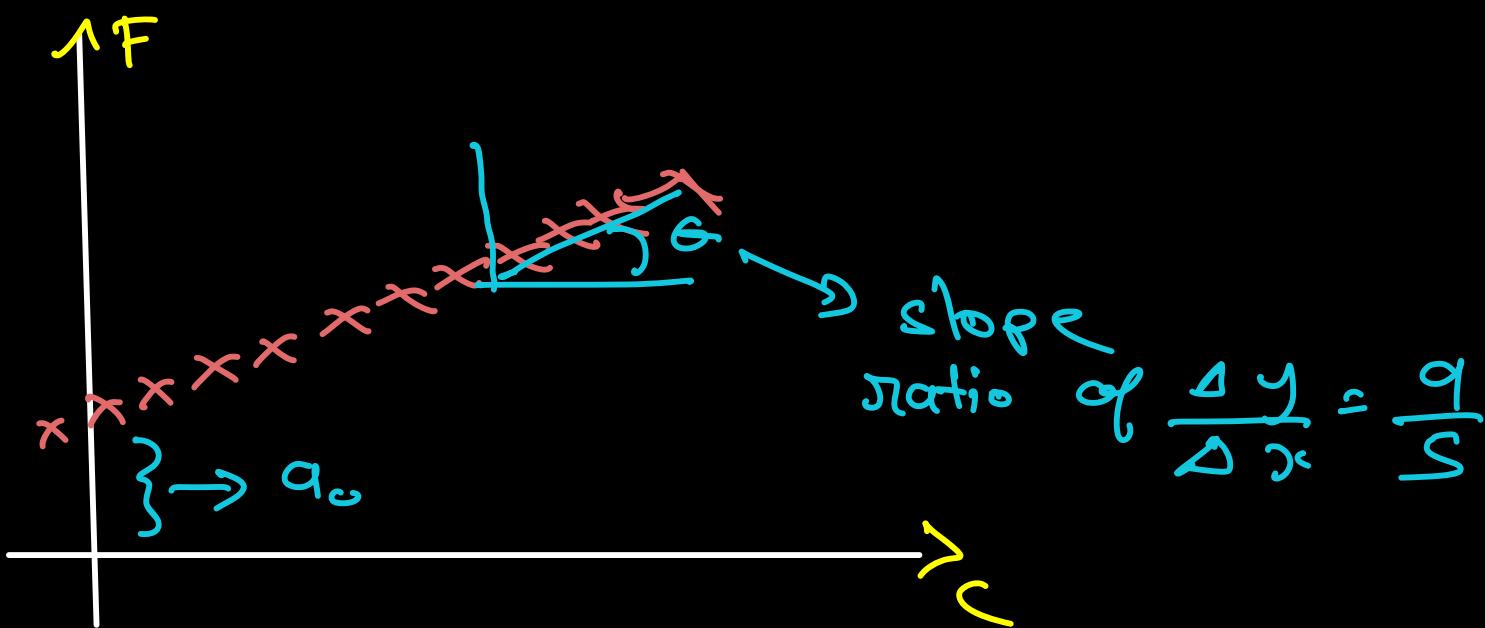
Train a LR, st.

$X = [-, -, -, -, -]$ Readings in C°

$y = [\dots, \cdot, \cdot, \cdot, \cdot, \cdot, \cdot, \cdot]$ Readings in
 F

We need $f(x)$, s.t

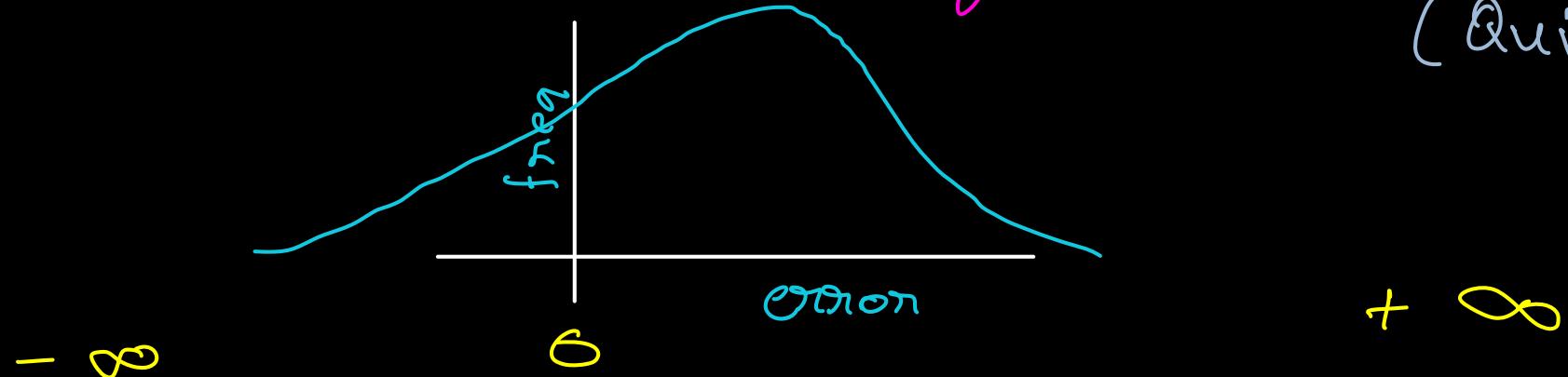
$$X \rightarrow f \rightarrow \hat{y} \sim y = \frac{q}{s} \cdot x + 32$$



$$\therefore \hat{y} = a_1 x + \underbrace{a_0}_{32} \rightarrow \text{intercept}$$

Hence answer = 32.

Q: After training a LR model, the residual distribution has the following production? (Quiz)



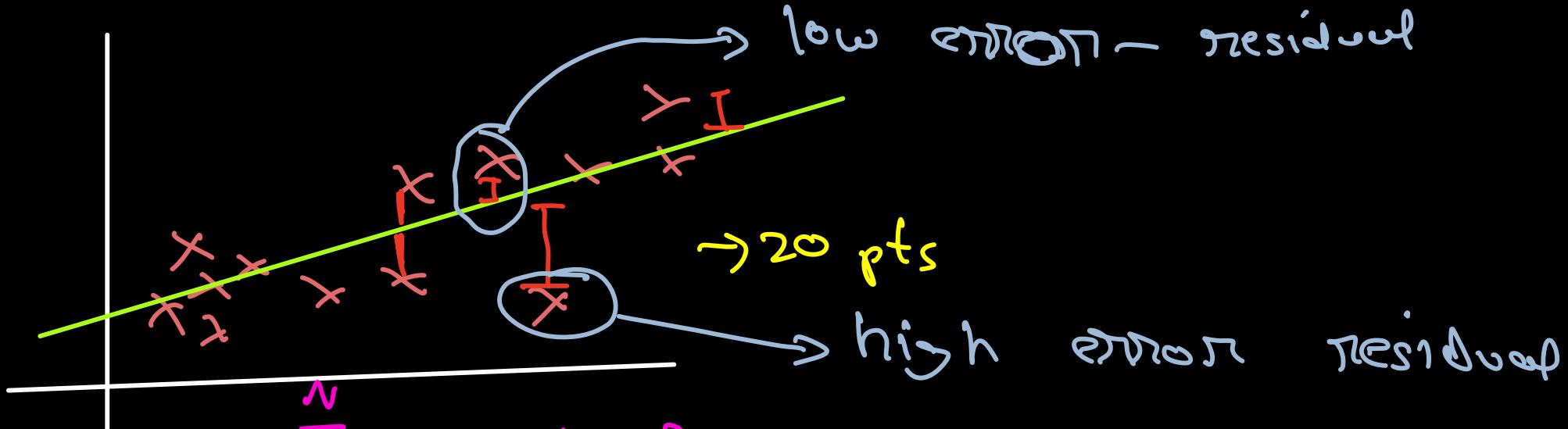
→ Automated | Easy

Solⁿ → No

→ The error dist is not normal, assumptions violated

→ The mean of residue is non-zero, +ve, indicating high tendency to Over predict

Concept: Regression Residuals



$$MSE = \sum_{i=0}^n (y_i - \hat{y}_i)^2$$

residual

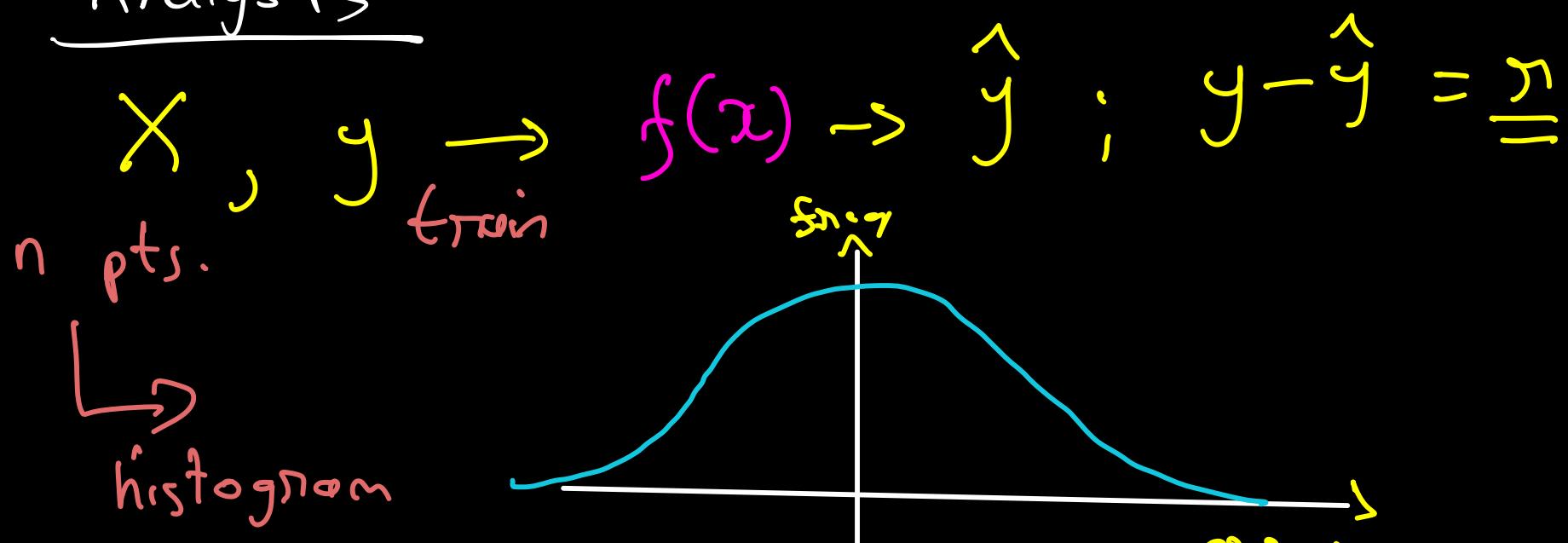
$$MSE = \sum_{i=0}^n r_i^2$$

$R = \underbrace{y_i - \hat{y}_i}$
 $r_i = y_i - \hat{y}_i$

Importance of Residuals

→ We just saw → all error metrics are computed on residuals

Analysis

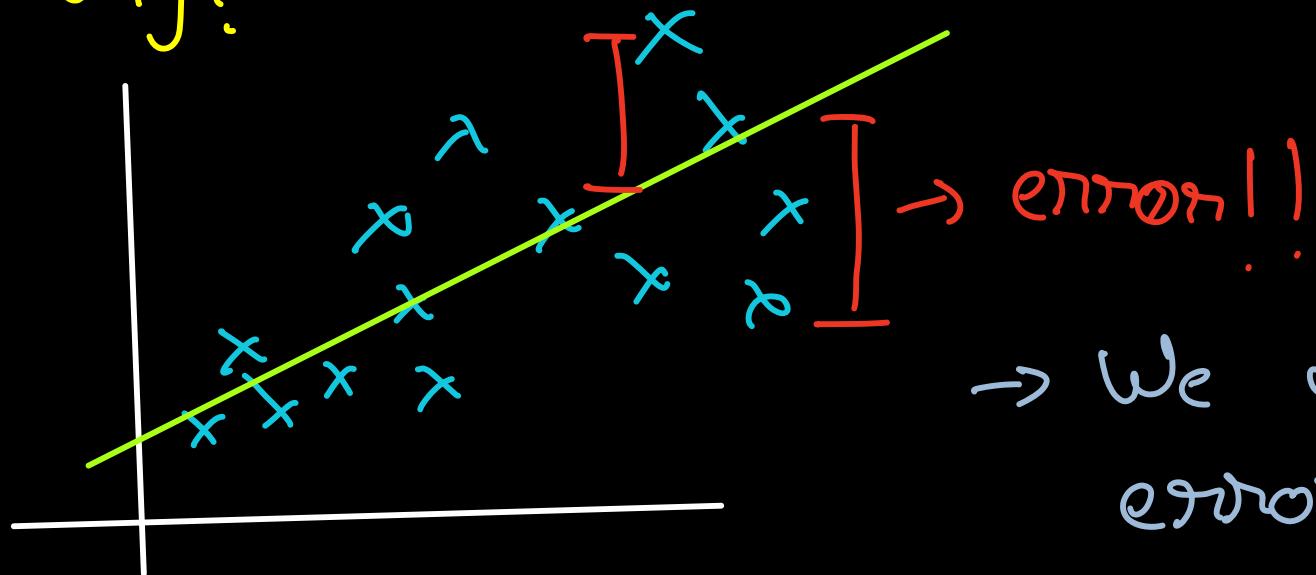


→ Ideally for LR → dist should be normal

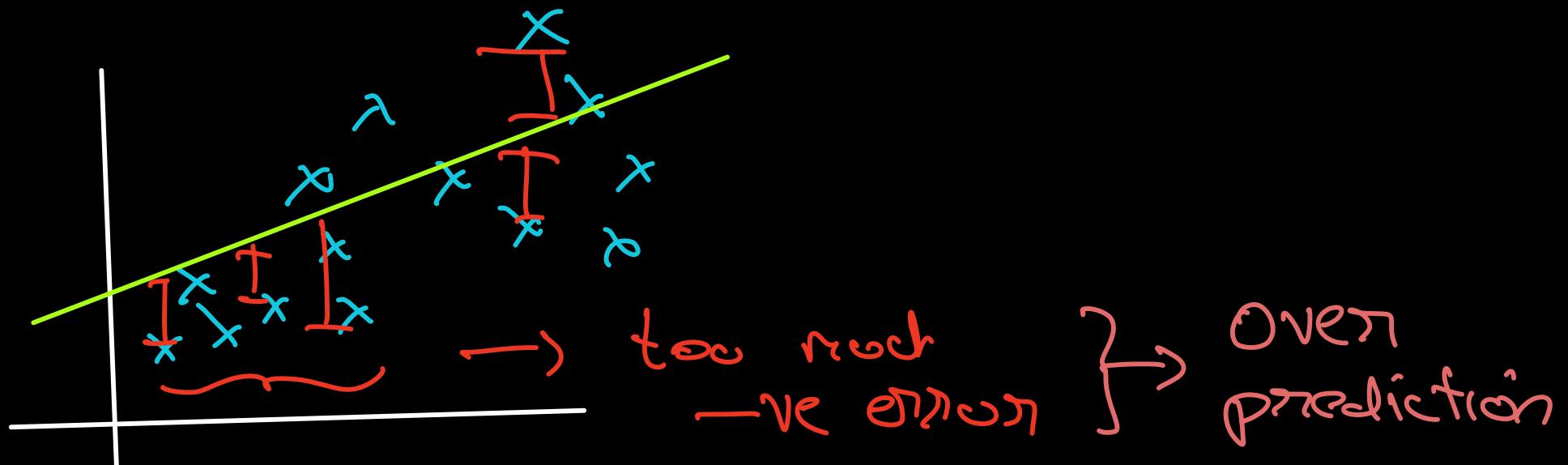
→ Mean should be 0

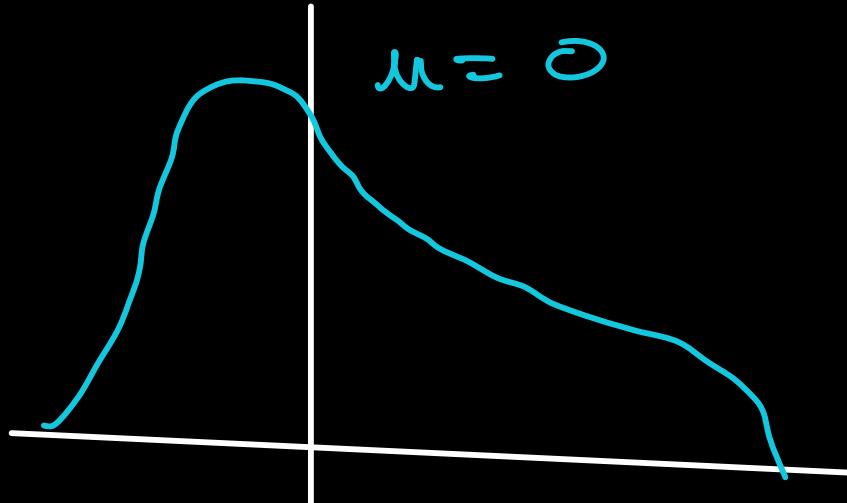
→ Variance : lower the better

why?



But overall
-ve error
~ +ve error





→ We want to predict
"error rate" in test
set.

→ In this model

→ if $\text{error} > 0$, error is low
if $\text{error} < 0$, error is high

⇒ asymmetry !!.

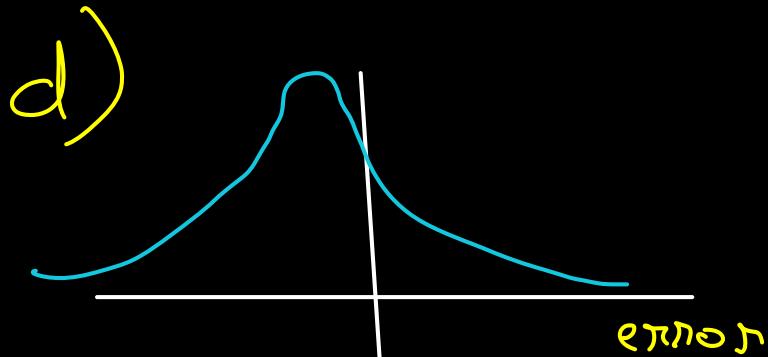
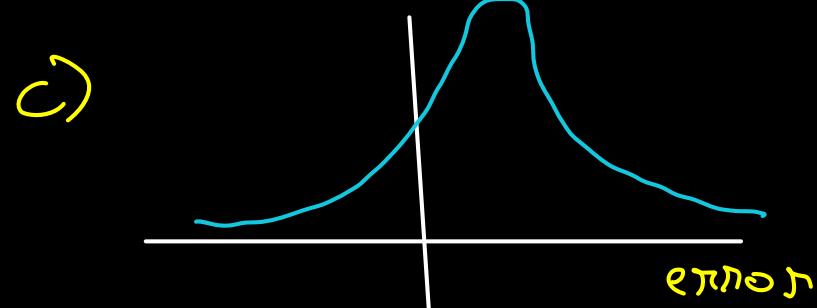
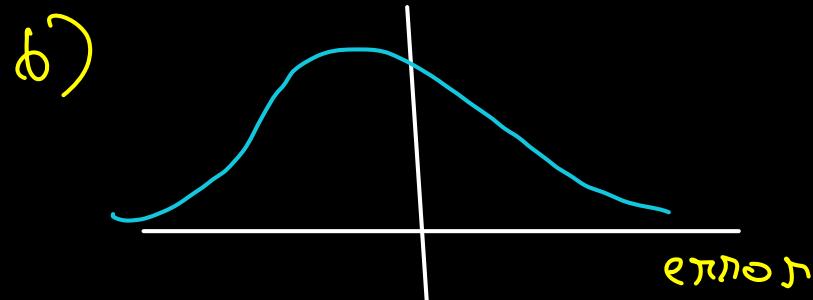
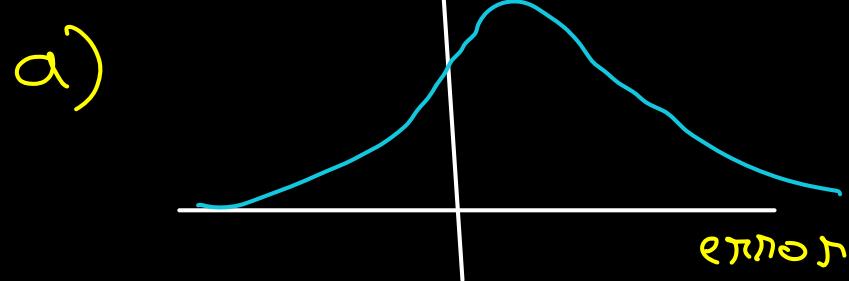
To make decisions we need confidence

→ my prediction is

$$53.5 \pm \underbrace{3}_{\text{at}} \quad 1\sigma$$

Q If a retail company has 50% profit margin for each sale, and 80% salvage price for unsold goods, while predicting Order Quantity, which of the following models seem best? (Quiz)

→ Automated / Live / Medium



Salvage price = Discounted price @ unsold goods
are sold (Company loses money)

The company wants max profits

if over-produce \rightarrow need to sell at

salvage price

\hookrightarrow Loss = 20%

if under-produce \rightarrow Loosing potential profit

Hence, better to $y - \hat{y} > 0$ \hookrightarrow Loss = 50%

overproduced \rightarrow Next, choose one with low variance

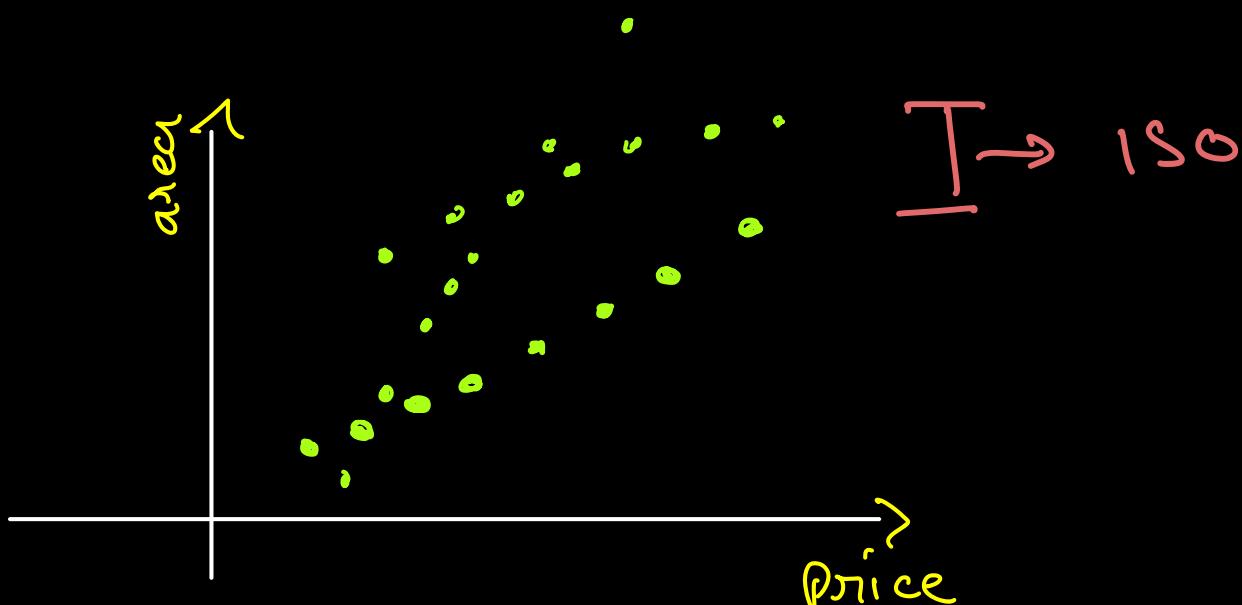
Hence Ans \rightarrow d

(more accurate)

Q You have been given house prices for a certain city. The only variables available are latitude, longitude and square footage. One scatter plot is shown below. How would you approach to build a price prediction model?

→ Live | Medium

(Open ended)



$\rightarrow \text{Sol}^n \rightarrow$ train a simple regression using
any kind of model SVM, DT, LR.
incorrect \rightarrow bad performance.

2) Use lat-long info to get neighbours.
 \rightarrow Google \rightarrow Assign 3 neigh.

X neigh	area	Y price	"Target encoding"
Say A	1500	300k	
B	1900	150k	
C	2250	900k	
A	1300	200k	\rightarrow rows of $n = A$
		150k	$\xrightarrow{\text{avg}} 500k$

\downarrow

N	A
500K	170°
350K	22°
1000K	130°
.	1
.	1

$$P_i \quad \hat{y} = a_0 + a_1 N + a_2 A$$

lets assume $a_1 = 1$

For $N = \underline{A}$ (1st res)

$$\hat{y} = a_0 + 500K + a_L A$$

For $N = \underline{B}$ (2nd res)

$$\hat{y} = a_0 + 350K + a_2 A$$

By design $\Rightarrow A > B \Rightarrow \underline{\underline{150K}}$

Use clustering (ML based on Manual)

→ Different LR for each NeighbourLoc.

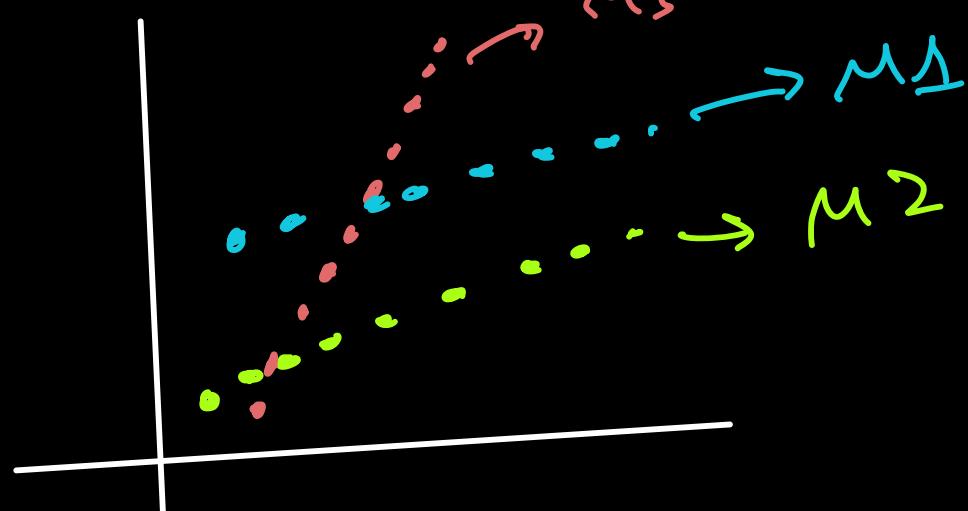
$$x = (A, 1000) \rightarrow N=A$$

Input now

hence → select M1

if $M_1 = \hat{a}_0 + \hat{a}_1 A = \hat{y}$

$N=B, M_2 = \hat{a}_0^2 + \hat{a}_1^2 A = \hat{y}$



Each model is
highly accurate

Q: For a KNN classifier, as we increase the value of K, does the model move towards . . .

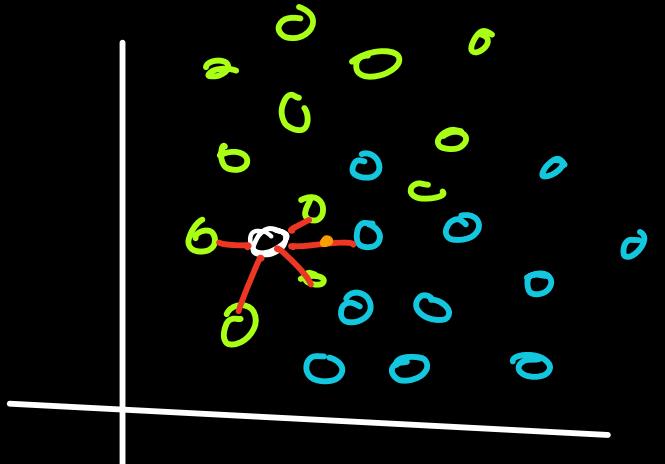
(Quiz)

→ overfitting

→ underfitting

→ Automated + Live | Easy

Sol^x.



→ check dist from
K nearest neighbours

K = defined by DS

↳ hyper parameter

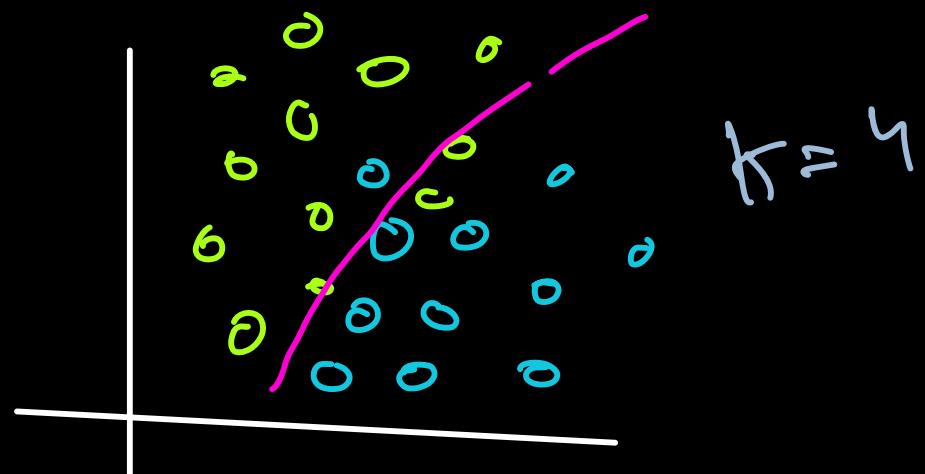
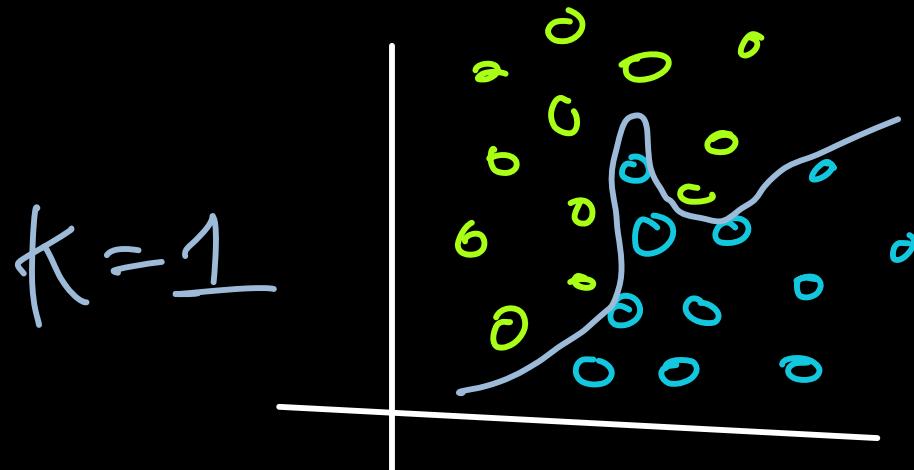
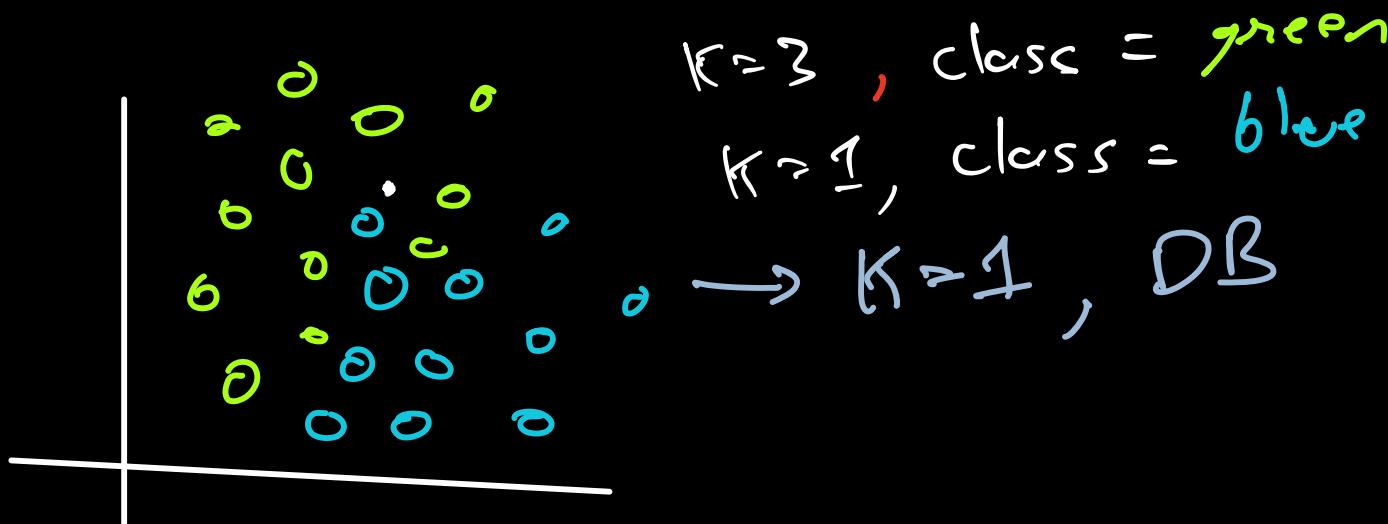
let K = 5

$\underbrace{5+1}_{\text{voting}} \rightarrow \text{green}$

As we increase K , what happens?

Let $K = 100M$

L_s will always take majority class
→ underfit



Q: Your colleague has trained a Tree based regression model, to predict the marks of students for a very hard test.

The max marks for all exams are 100.

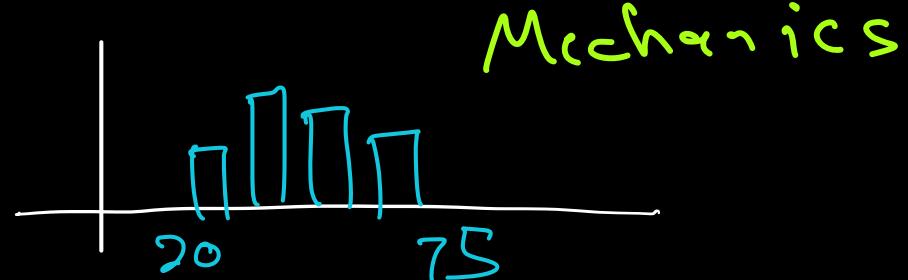
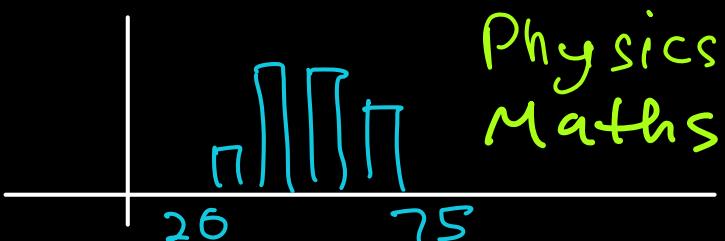
→ Live | Medium (Open ended)

Physics	Maths
65	68
72	67
35	40

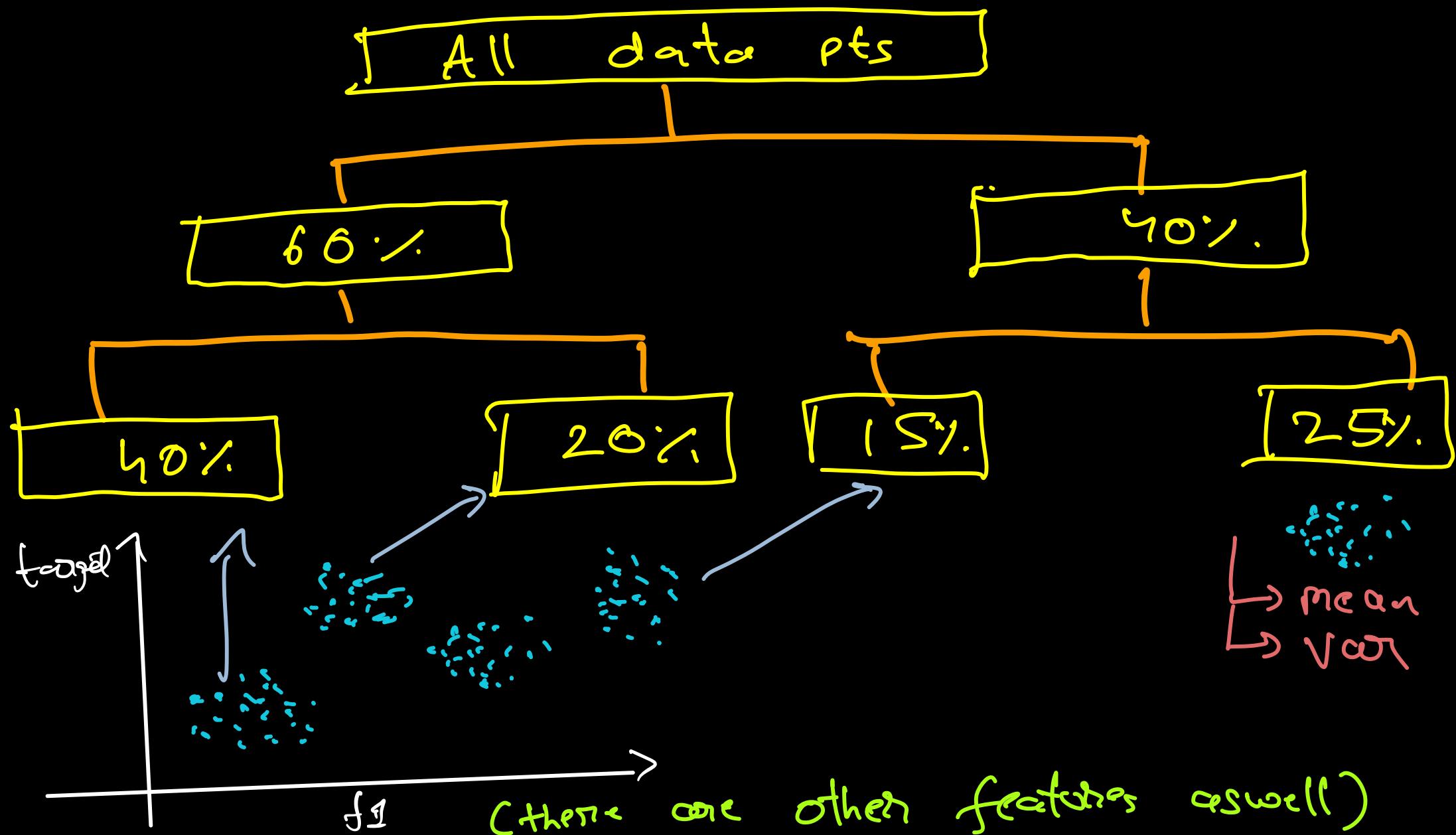
Mechanics → target

74 The model MAPE
70 is 1%. Can we
43 put this into
 production?

Distributions:



Let's understand how a DT works



25% Output \rightarrow mean: \rightarrow more predictable
if Σ

mean
var \rightarrow Variance \downarrow low

40% 20% 15% 25%
 $\downarrow u_1$ $\downarrow u_2$ $\downarrow u_3$ $\downarrow u_4$

Q: What will be the O/P range??

$\rightarrow [u_1, u_2, u_3, u_4]$

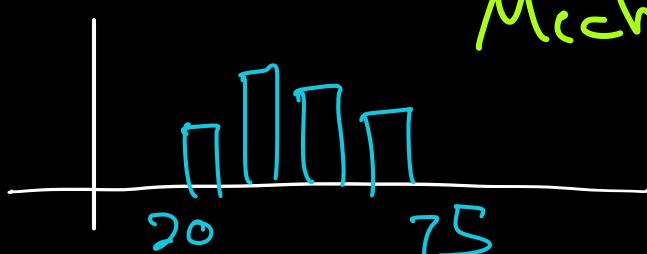
\Rightarrow Yes, there will be only 4 unique values in regression O/P in this case

Had you guys realized this?

For regressions, we use RFR \rightarrow so we get many (but finite) unique values because it averages many trees.

Tree based Regression can never o/p anything outside of training data 'target' range.

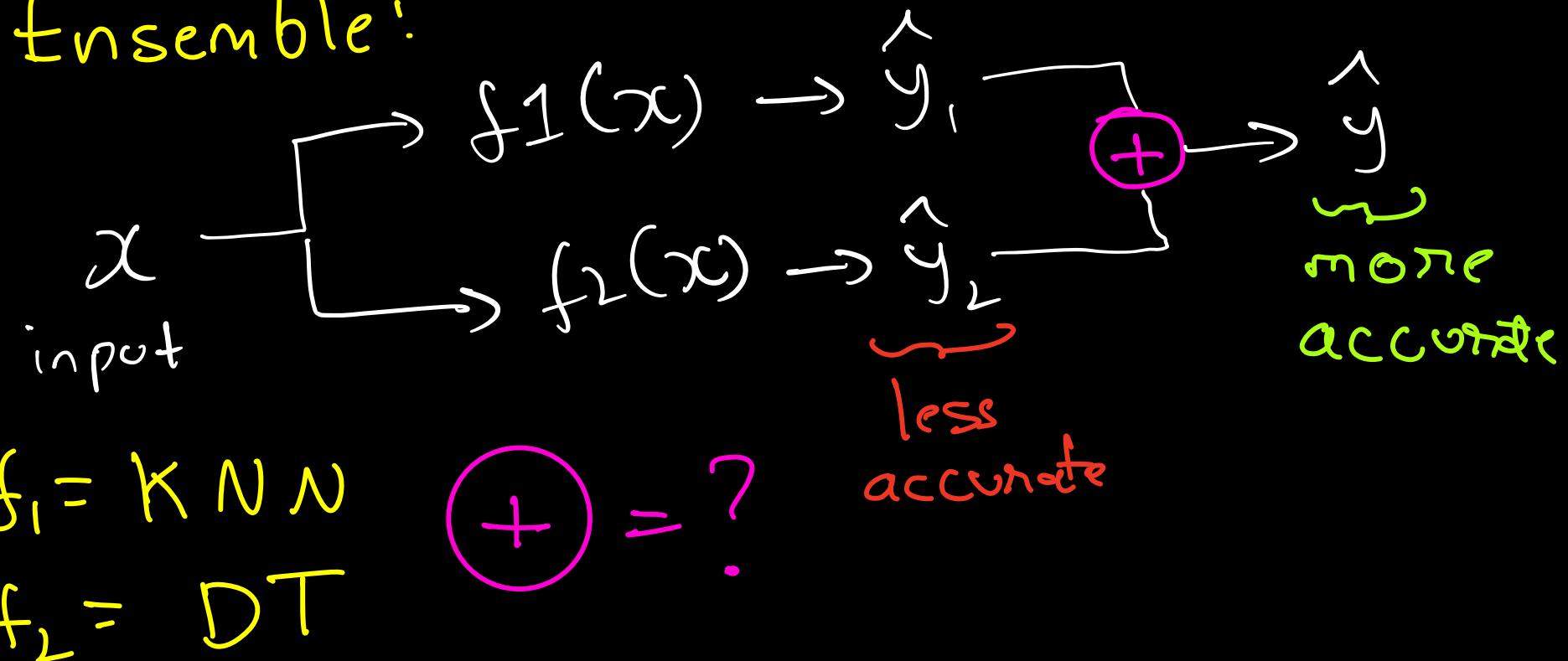
Target has a range of only 20 - 75 in training data



Mechanics
Linear regression would be better

Q: How do you create ensemble model from KNN and DT ? (open ended)

Ensemble:



Soln

Since KNN is purely a classification model, we need to create ensemble for



2 blue 3 green

Output \rightarrow green

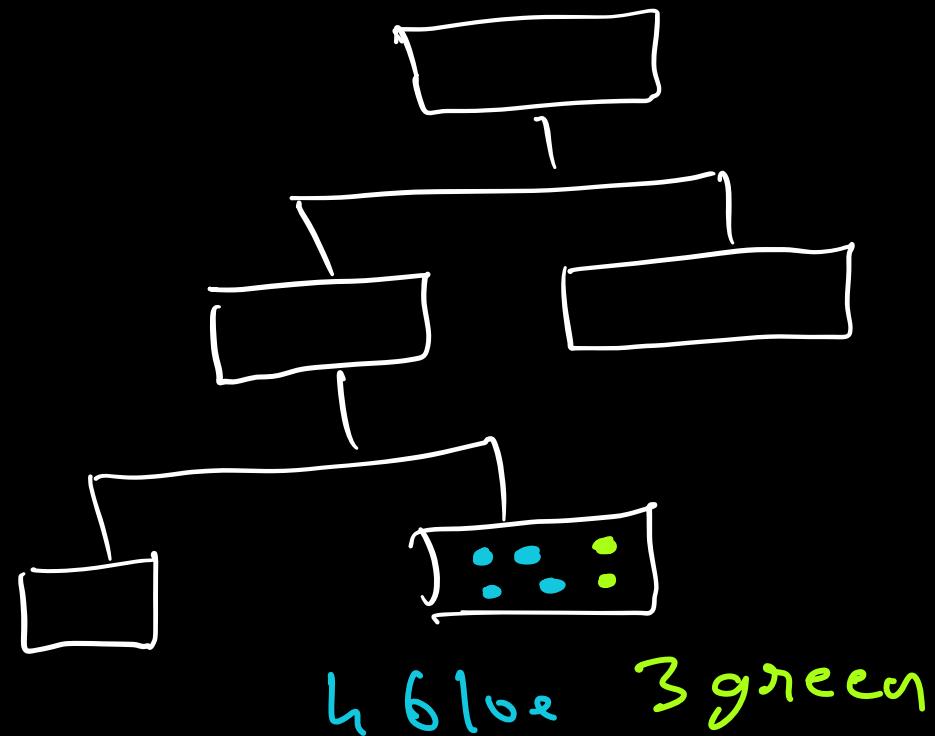
Voting? = Tie

\downarrow
keep odd number.

Better Idea \rightarrow Use Probabilities

KNN Prob:

$$\frac{2}{5} \quad \frac{3}{5}$$



4 blue 3 green

Output \rightarrow Blue

DT	Prob:
$\frac{4}{6}$	$\frac{2}{6}$

Adding up !

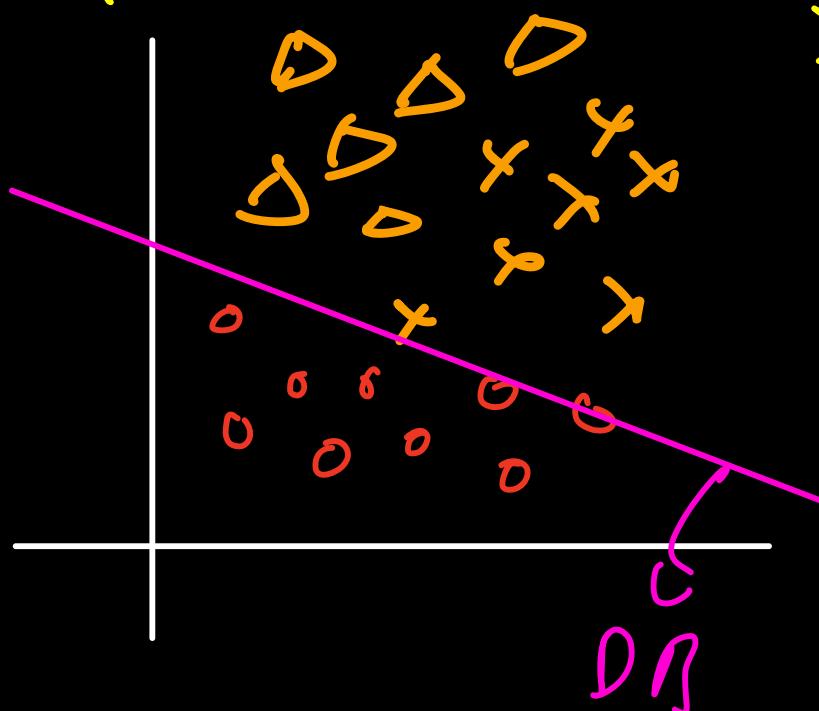
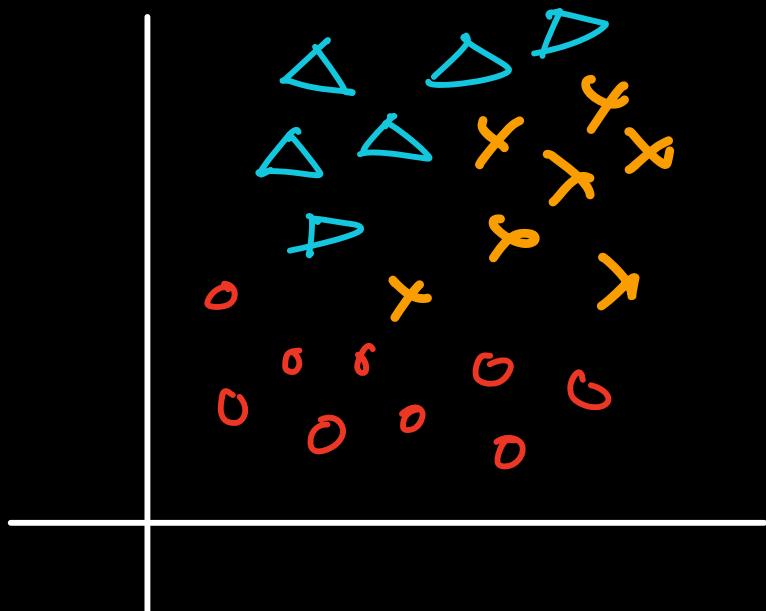
blue $(\frac{2}{5} + \frac{4}{6})_2 = \frac{16}{30} = \frac{8}{15}$

green $(\frac{3}{5} + \frac{2}{6})_2 = \frac{14}{30} = \frac{7}{15}$

↳ Output \rightarrow Blue

Q: How can you use Log Reg for a multiclass problem statement?
→ Live Easy (Quiz)

Soln One-vs-Rest framework.



if $M1 = 1$
↳ Red
elif $M2 = 1$
↳ Blue
else
↳ Orange

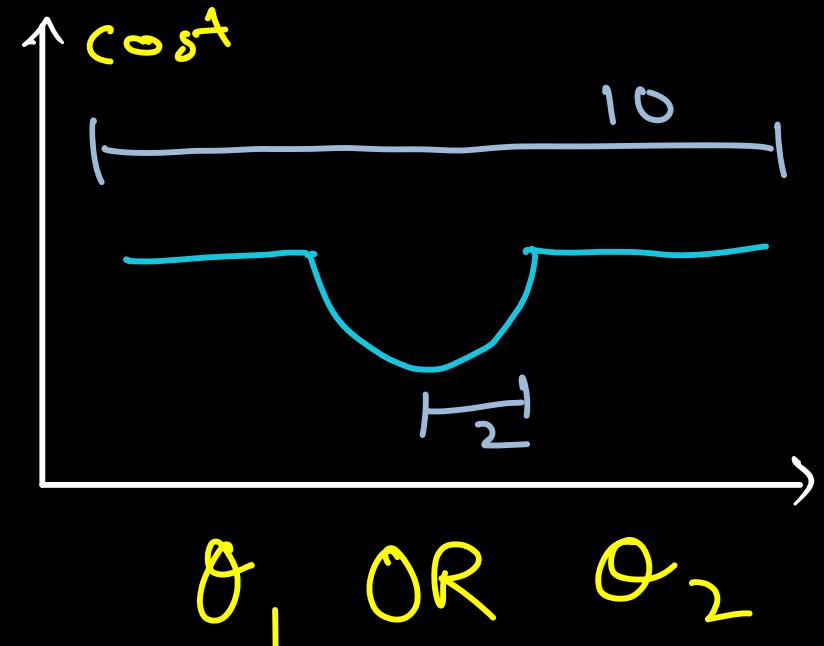
Make 2 LogReg → Red vs Rest
blue vs Rest

Q: What is the probability of finding the optimum value using simple G.D (θ_1, θ_2) if the cost function looks like this?

↳ 2D view w.r.t θ_1

↳ same view w.r.t θ_2

→ Automated | Medium
(Quiz)



So (\rightarrow) If you init on flat surface, G.D will
 $10 \times 10 = 100$ be

$$\pi(2)^2 = 4\pi \rightarrow \text{prob} = \frac{4\pi}{100} \quad \underline{\text{stuck}}$$

Q1. Your teammate tried adding a new feature to the training data and has reported that R^2 has dropped.

What is your suggestion? (Quiz)

→ Live | Medium

→ New feature is not important, should remove it.

→ Colleague has made a mistake, R^2 cannot reduce in this case

→ Feature could be correlated with another feature. → Deep dive

→ Do hyper-param tuning first, then take a decision

Soln

$$R^2 = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2} ; \bar{y} = \text{mean}$$

→ How much better is the prediction compared to just the mean

Now, let's recall LR.

$$\min_{\beta} SSE = \min_{\beta} \sum_{i=1}^n (y - \hat{y})^2$$

$$= \min_{\beta} \sum_{i=1}^n (y - \beta_0 - \beta_1 x_{i1} - \beta_2 x_{i2} - \dots - \beta_n x_{in})$$

↳ G.D will make $\vec{\beta}$ such that

SSE is min $\rightarrow S_1$

Now when we add one more Var.

$$= \min_{\beta} \sum_{i=1}^n (y - \beta_0 - \beta_1 x_{i1} - \dots - \beta_n x_{in} - \underbrace{\beta_{n+1} x_{int}}_{\downarrow})$$

G.D returns $\rightarrow S_2$

This is extra

$$S_2 \leq S_1$$

Because, if $\beta_{n+1} = 0$,
 $S_1 = S_2$

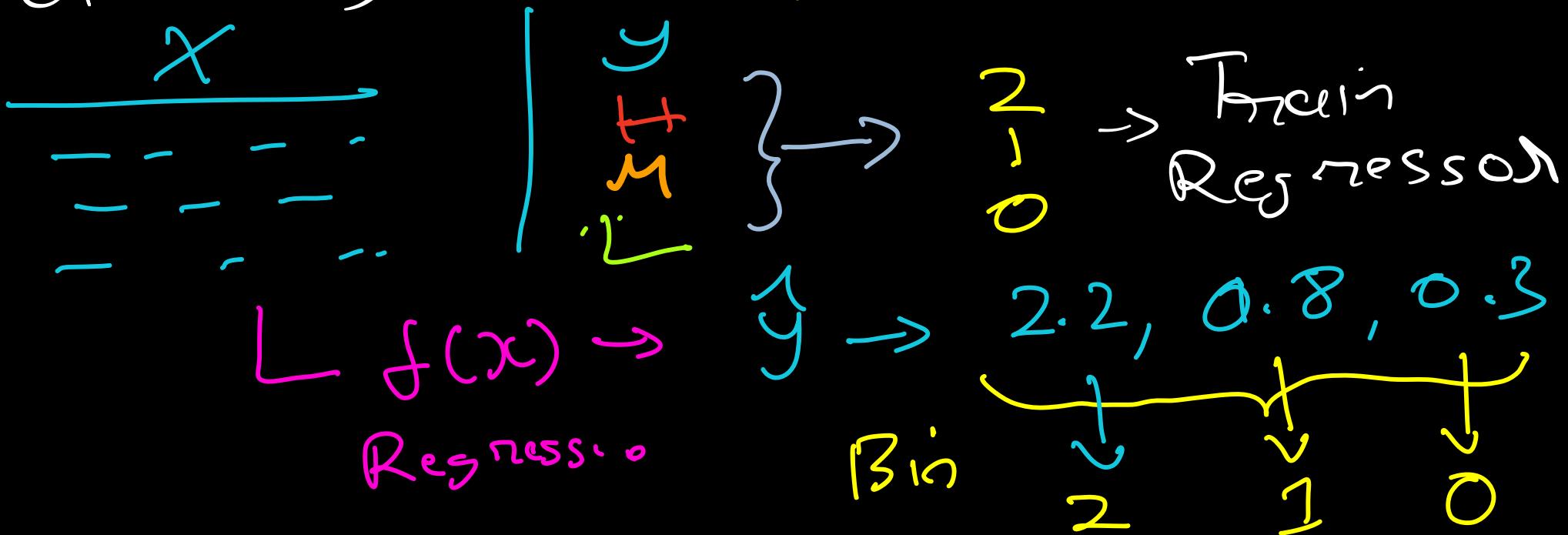
OR, $\beta_{n+1} \neq 0$; in a way that $S_2 < S_1$

Q: How would you approach a cancer severity classification problem? HML?

(open ended)

→ Live | Medium

Sol → 1) Use Regression and bin



Binning may have issues at the boundaries

Sol-2) Classification for Ordinal Classes

$$X \rightarrow f_1(x) \xrightarrow{\text{No}} \hat{y} = 0$$

$$\boxed{\hat{y} > 1}$$

$$\xrightarrow{\text{Yes}} f_2(x) \xrightarrow{\text{No}} \hat{y} = 1$$

$$\boxed{\hat{y} > 2}$$

It might be

easier to predict

$\hat{y} > 1$ as a classification

compared to predicting exact value of \hat{y} .

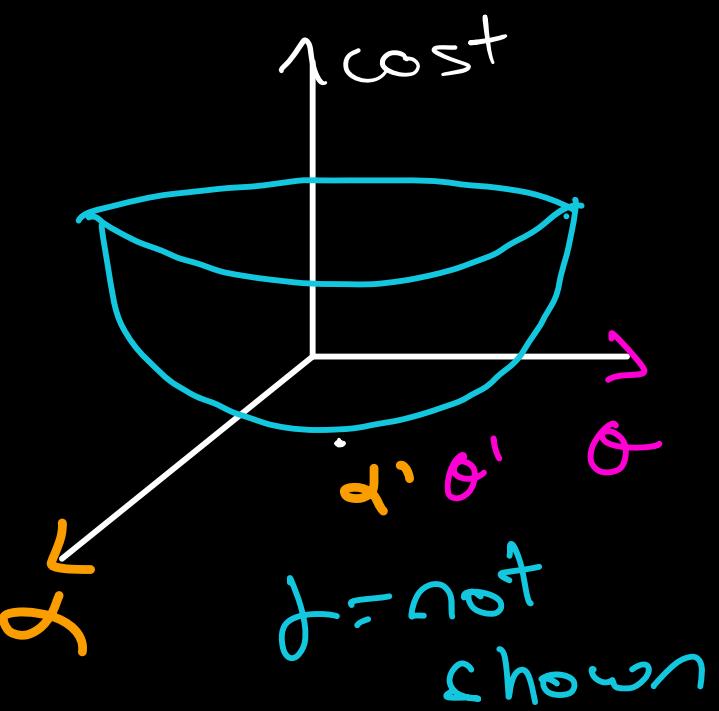
Q: Can you design a Gradient Descent algorithm, to optimize a cost function with params $\alpha, \beta, \gamma, \theta$, when β is a boolean variable and can only be 0 or 1?

→ Live Hand (open ended)

Solⁿ

Problem → Gr. D needs cost funcⁿ to be differentiable w.r.t $\alpha, \beta, \gamma, \theta$

But β is not differentiable $\frac{\partial}{\partial \beta} = \text{not defined}$



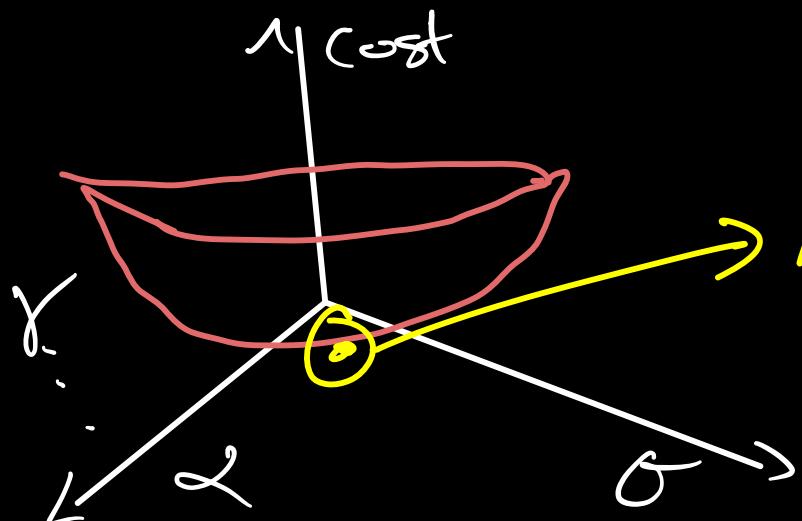
Q: Any Ideas?

1) Assume β to be continuous from $0 \rightarrow 1$. Use GD. where β is constrained.

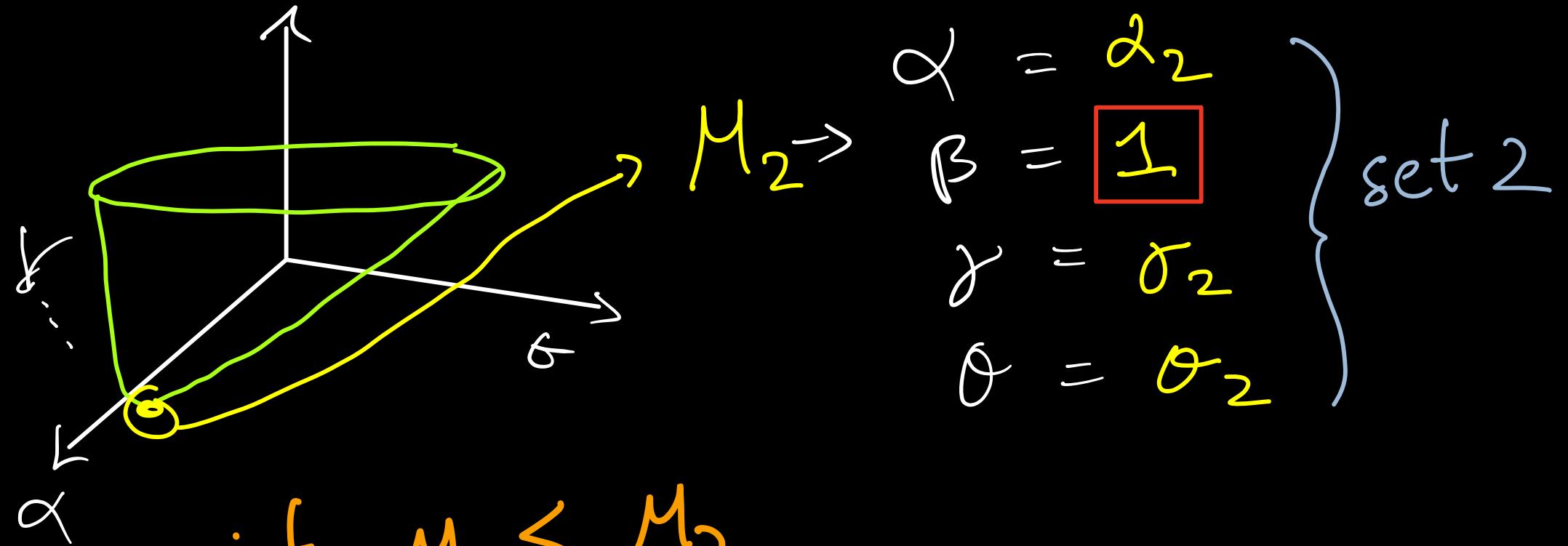
\hookrightarrow IF, GD outputs $\beta = 0.09$
 $\hookrightarrow \beta = 0$

if, $\beta = 0.65 \rightarrow \beta = 1$ (sub optimal)

2) Fix $\beta = 0$, Solve GD with α, γ, θ .



$$\left. \begin{array}{l} \alpha = \alpha_1 \\ \beta = 0 \\ \gamma = \gamma_1 \\ \theta = \theta_1 \end{array} \right\} \text{set 1}$$



if $M_1 < M_2$

$$\text{out} = \{\alpha_1, 0, \gamma_1, \theta_1\}$$

else

$$\text{out} = \{\alpha_2, 1, \gamma_2, \theta_2\}$$

Q : Based on product attributes , how would you model , a product assortment problem at Amazon? Please ask any further data points that you require .

(open ended)

(long disc)

→ Live | Hard

Soln : First , the interviewer wants to know how much you can imagine about product data. [long discussion]

→ Applications :

→ Search query → Filters

→ UI segmentation

→ Clustering: But, these clusters need to have names, so they can be filtered / UI.

→ Classification :

→ Using business logic

→ Using Unsupervised learning + human.

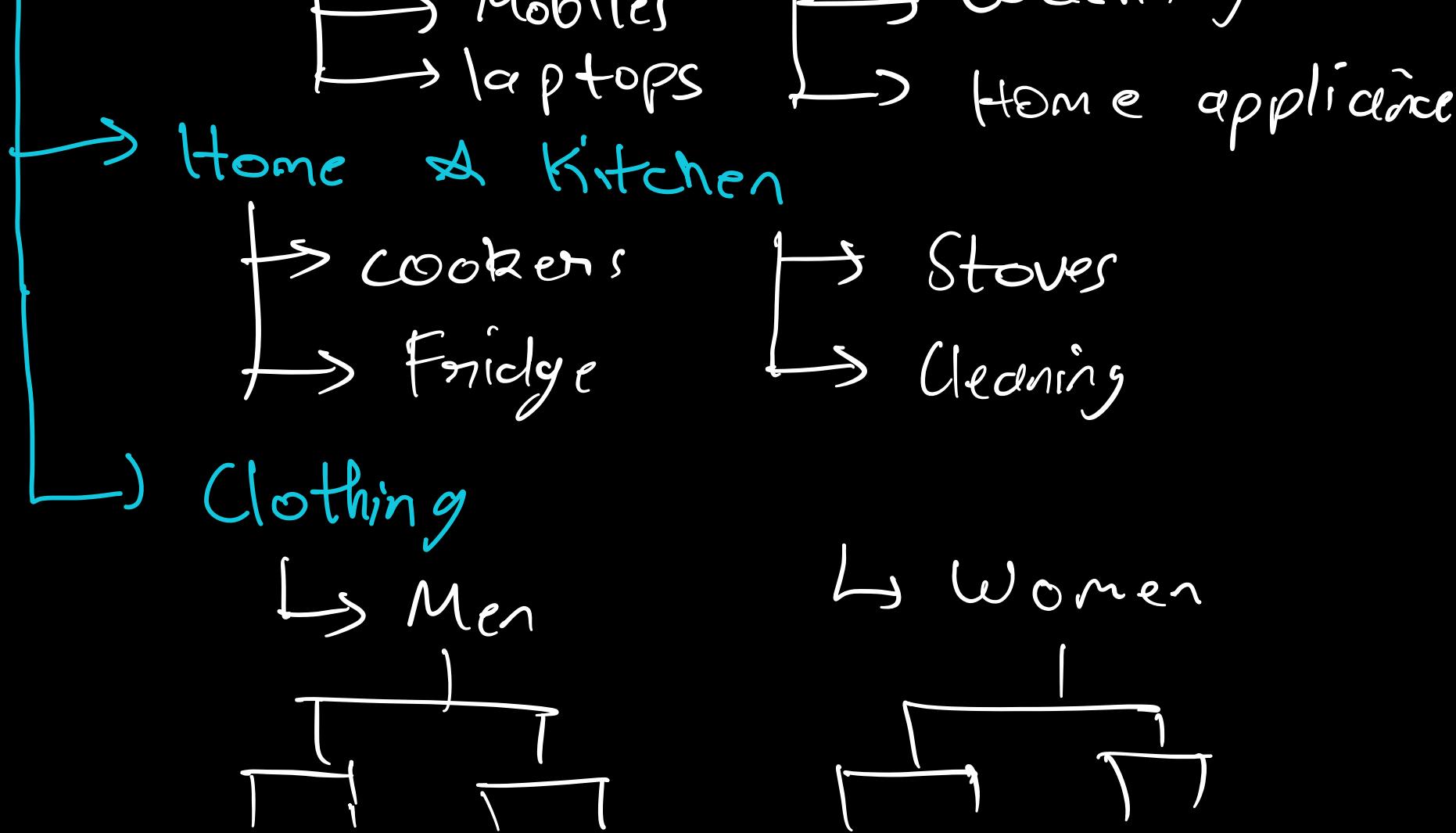
Product Hierarchy

Do not miss the point about P-H.

Eg. →

→ Electronics

| → Mobiles → Laptops → Laptops Mech



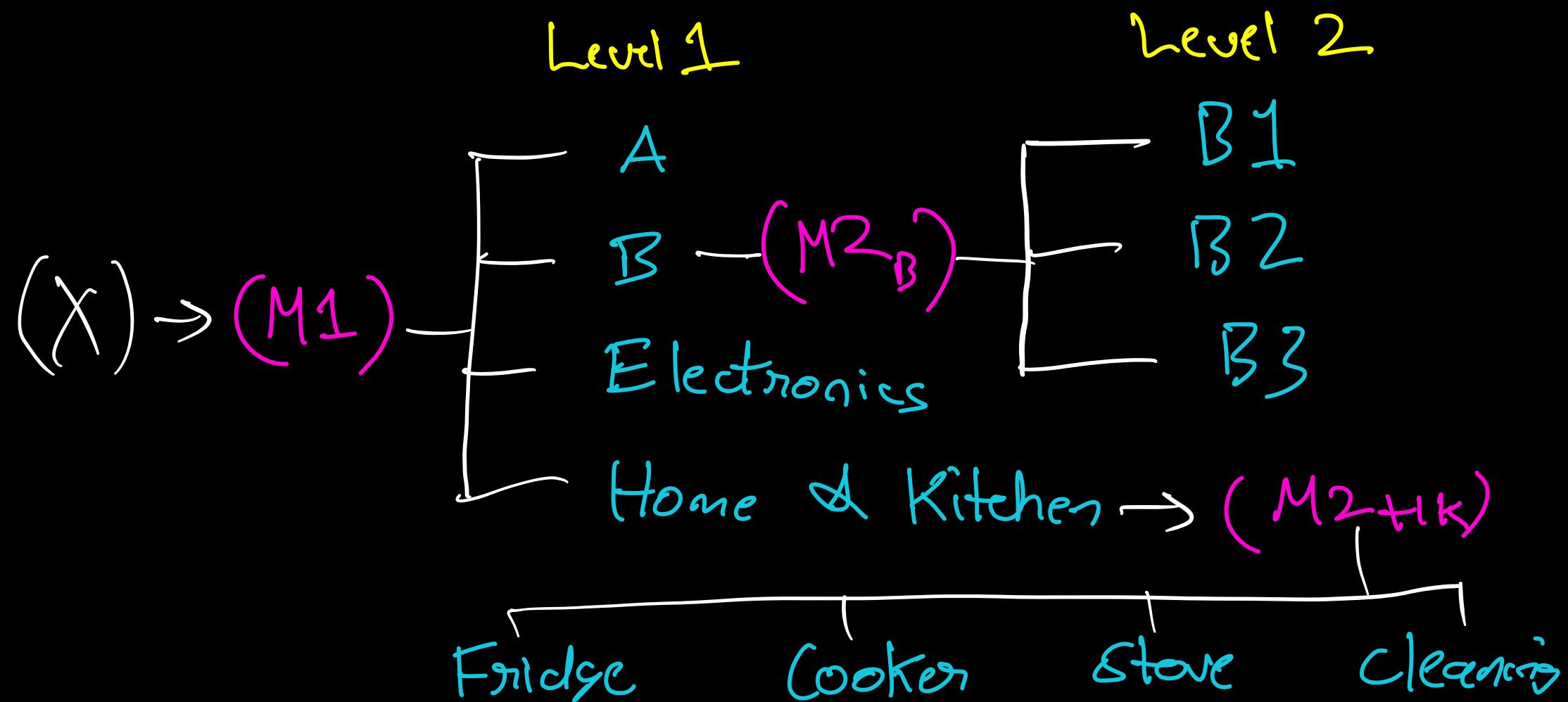
Q: So how do you solve a
hierarchical, multiclass, multilabel problem?

Variation Q: How would you approach

a classification problem with classes as follows? A1 A2 A3, B1 - B5, C1 - C3, etc?



Multiple Levels of Models



M1 → Classification | M2_B → ORDINAL
M2_{Hf} → U | Classifier

Follow up Question: Won't this make
too many models if there are 3-4 layers
→ Yes, decide based on, How many products
each model is classifying. Eg: 100K per
model in last layer.
→ Layer depth can also vary: Eg:
3 layers for electronics, 5 for clothing
FUQ: What are the accuracy trade-offs?

Concerns / Priorities ?

Sd \rightarrow L1 needs more accuracy, since

\rightarrow It aggregates more data pts.

\rightarrow If User goes in wrong category
he is misjudged from the start.

FQ: What changes would you make if
one product can belong to multiple cat?

\rightarrow Each model is a binary model to
identify own class (one-vs-rest). Input
needs to be passed through the entire tree.

Q: You are working on a binary classif^n model. When you changed threshold from 0.5 to 0.55, the confusion matrix flipped. What can you say abt this model?
(discussion)

→ Live | Hard

100	200
50	900



$T = 0.5$

102	56
197	895

$T = 0.55$

Soln → The model is not confident

Use of threshold:

$X \rightarrow f(x) \rightarrow \text{probability} \rightarrow \text{thresh} \rightarrow \text{label}$

(trained model) $(0 - 1) > T$

Ideally, when the model returns class '0'
it should be 100% confident. Realistically,
a good model will be $>90\%$ confident
for most data points.

Hence, changing threshold from 50% \rightarrow 55%
should not alter the predictions.

This confusion matrix implies, most predictions are around 50 - 60% confidence.

Unstable Model

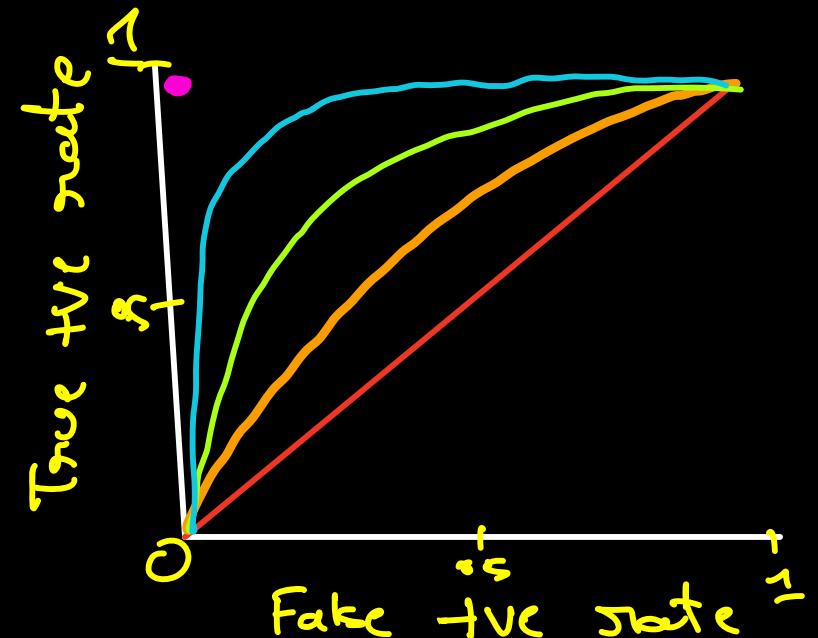
→ Slight changes in I/P will cause different predictions

FUQ: How can you best represent this?

→ ROC Curve.

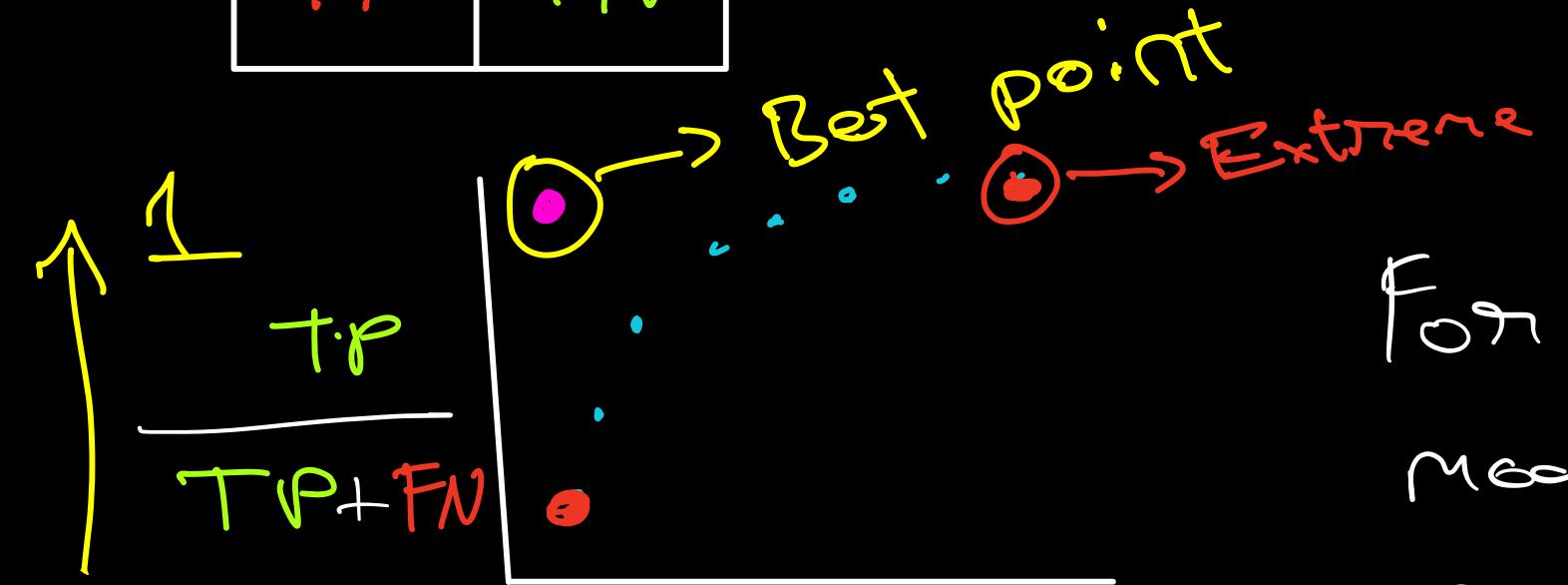
AUC ≈ 0.6
(guess)

- = perfect
- = Great
- = Better
- = Bad
- = Random



TP	FP
FN	TN

Green should be more
Red should be less



For a good model, $TPR = 1$
and $FPR = 0$.

$$\text{At highest } T, \quad \frac{\text{TP}}{\text{TP} + \text{FN}} = 1 \quad \frac{\text{FP}}{\text{FP} + \text{TN}} = 0$$

$T_1 \rightarrow (TPR_1, FPR_1)$
 $T_2 \rightarrow (TP_2, FPR_2)$
 \vdots

Best

Q: Predicting House prices : The dist of prices are shown below: (Assume LR)

What metrics would you use to evaluate your predictions?

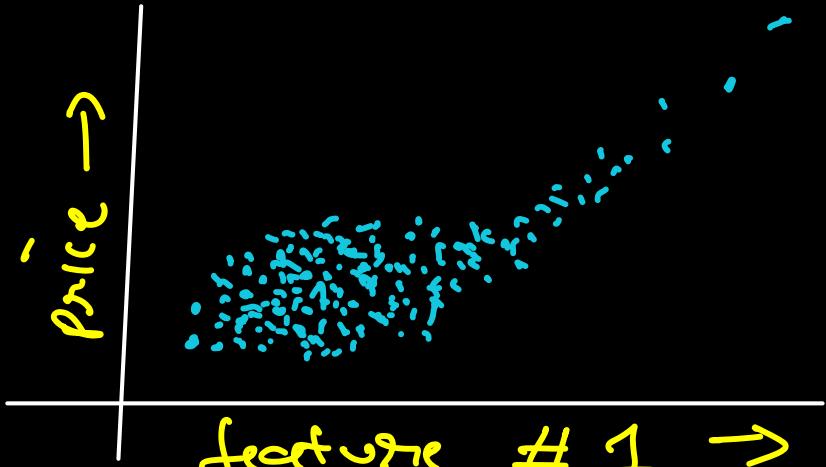
- a) MAE b) MAPE
- c) RMSE d) MSLE

FUQ: Any other suggestions while modelling?

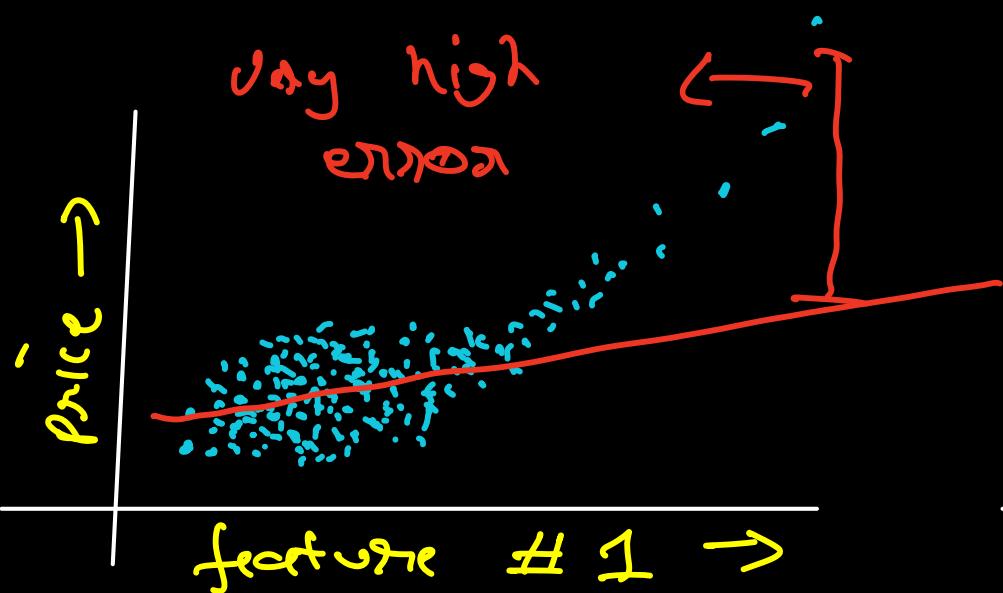
$$\text{MAPE} = \frac{1}{N} \sum_{i=1}^N \left| \frac{(y_i - \hat{y}_i)}{y_i} \right| \quad \mid \quad \text{MSLE} = \frac{1}{N} \sum_{i=1}^N \left[\ln(y_i) - \ln(\hat{y}_i) \right]^2$$

[Note: Many times interviewer will give you a new metric / formula and see if you can make sense]

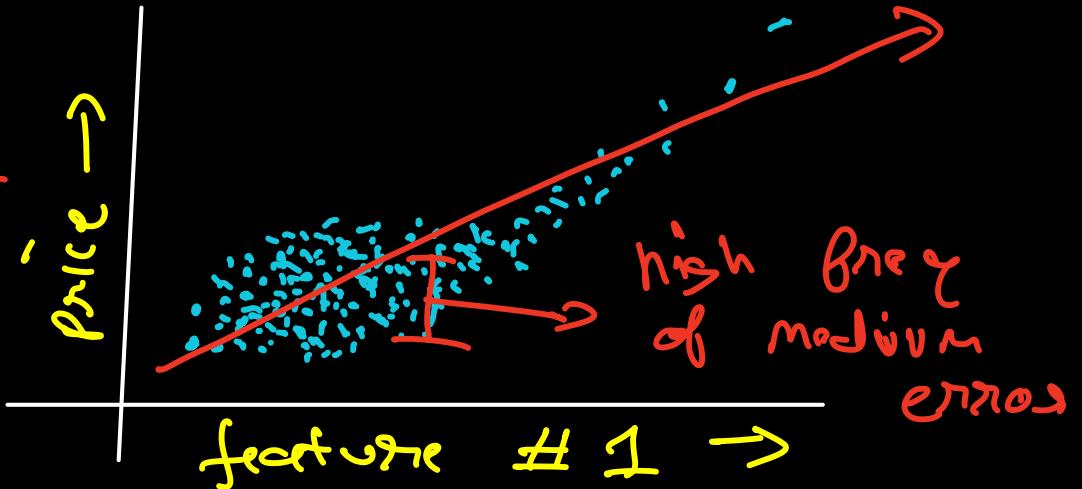
→ Notice that
density of pts
is higher closer
to the y-axis.



feature #1 →



feature #1 →



feature #1 →

Trade off.

MAE \rightarrow X' : because it's not at the same scale.

RMSE \rightarrow X : As we see we will need to trade-off somewhere:

Eg:

	y	\hat{y}	e
case 1	{ 100K 5M	150K 300K	50K <u>4.7M</u>
case 2	{ 100K 5M	200K 7M	100K <u>1.1M</u>
			case-2 ✓
	which one of them is more wrong?		will total error $\frac{1}{N} \sum (\dots)$ make sense?

Percentage Error:

$$\% \text{ error} = \frac{(y - \hat{y})}{y} = \frac{100}{100} = 100\% \text{ error}$$

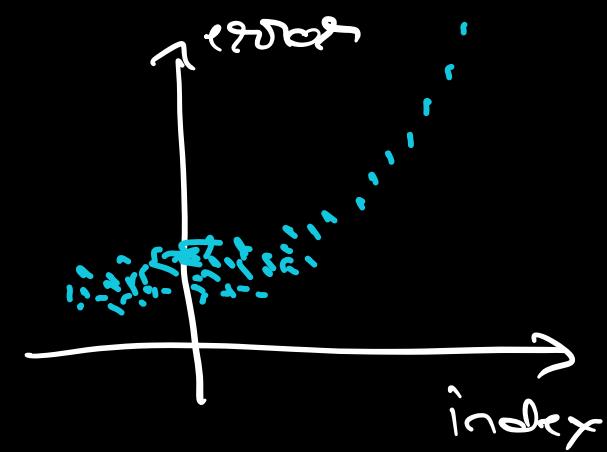
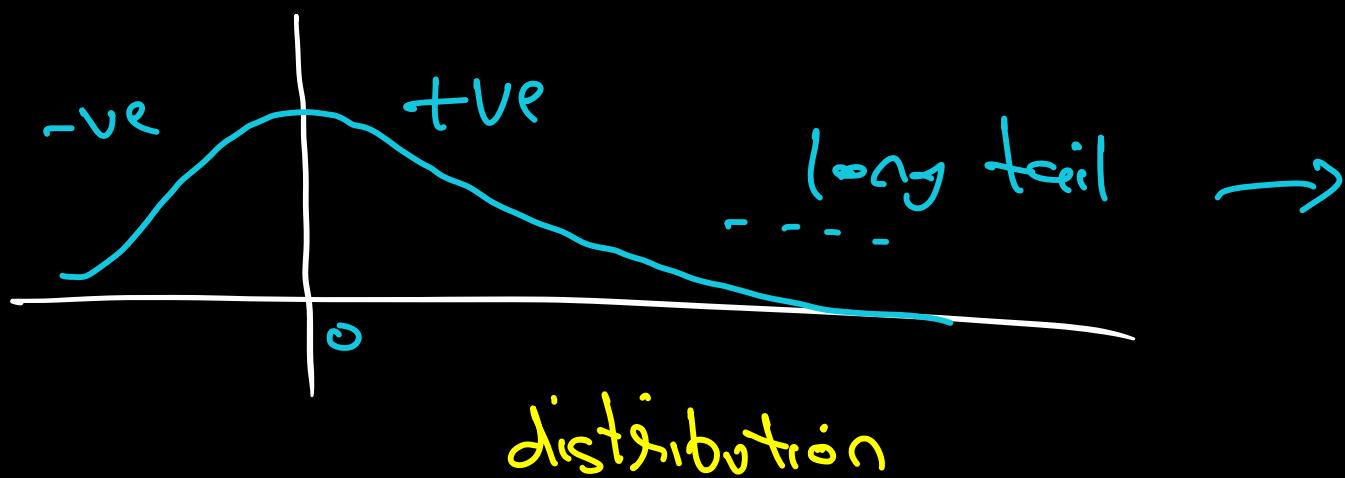
(2x value)

$$= \frac{1n}{Sn} = 20\% \text{ error.}$$

So for such a varying scale look - 10M percentage error, i.e. MAPE makes more sense.

Now, moving Further:

Let's look at the dist of error for case - 1



Scatter

Notice the non-linear scatter plot. And let's look at MSLE:

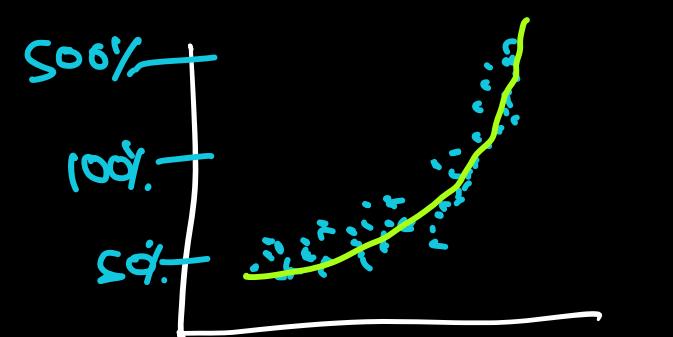
$$\frac{1}{N} \sum_{i=1}^N [\ln(y) - \ln(\hat{y})]^2 = \frac{1}{N} \sum_{i=1}^N [\ln(\frac{y}{\hat{y}})]^2$$

$\therefore \boxed{\ln(a/b) = \ln(a) - \ln(b)}$ percentage

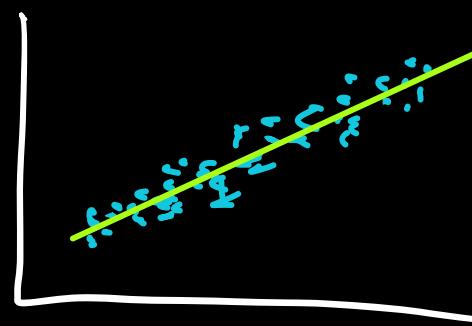
$$\ln(1) = 0 \quad \xleftarrow{\text{No error.}} \quad \leftarrow \frac{y}{\hat{y}} = 1$$

best pred

Let's say % error is rising fast



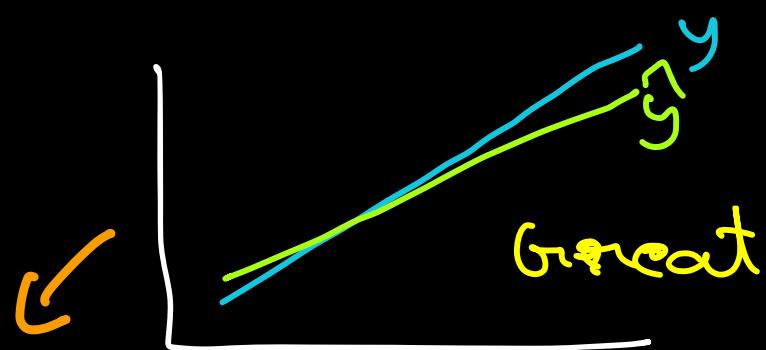
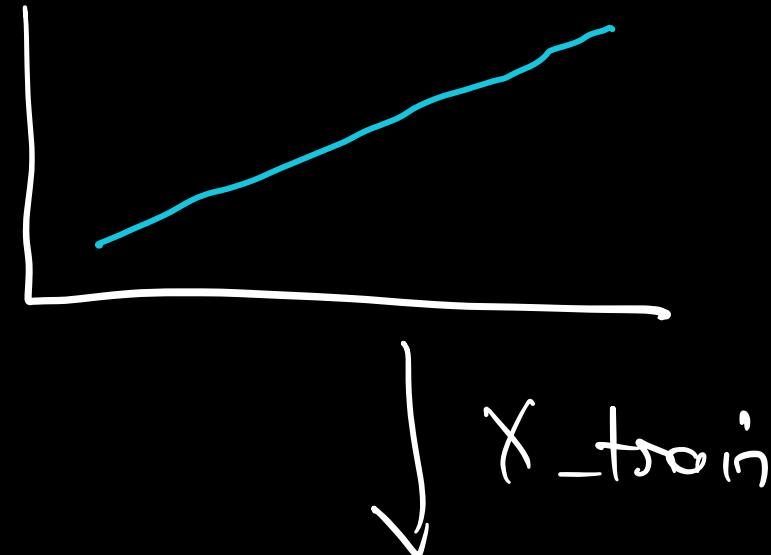
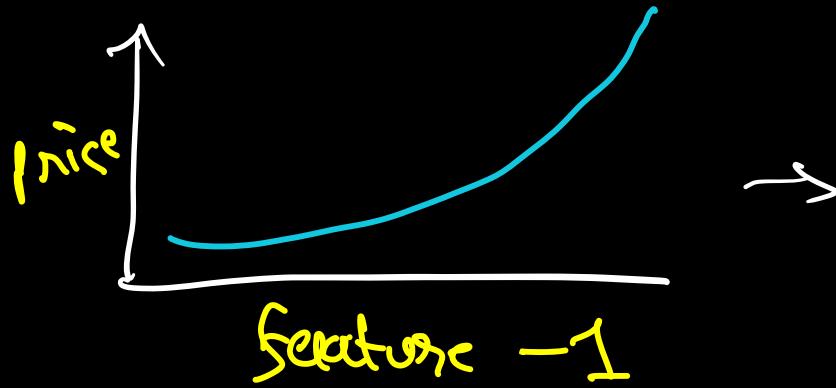
\log



Reduces error for high error case.

Best way:

Step -1: Take log of price



$\leftarrow LR()$

Measure

RMSE



$\rightarrow X \rightarrow LR(x) \rightarrow \hat{y} \rightarrow \underbrace{c^{\hat{y}}}_{\text{Anti log}} \rightarrow \# \underline{\text{pred}}$

Anti log