

Polynomial Regression

Agenda

- ① Polynomial Regression
- ② Generalisation & Occam's Razor
- ③ Underfitting & Overfitting
- ④ Bias Variance Tradeoff

G.D Variants

$$w_j' = w_j - \eta \cdot \frac{\partial L}{\partial w_j}$$

Most Popular

Batch G.D

① iterates all data points -
'gradients / derivative'

② 100% Converge global minima

③ Speed is slow

Mini-Batch G.D

① iterate over 'K' = 32 Batch size data points.

② in Between B.G.D & S.G.D

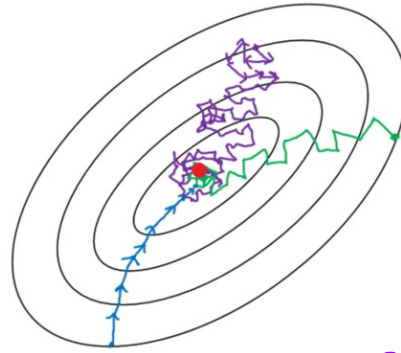
③ Speed is moderate

Stochastic G.D

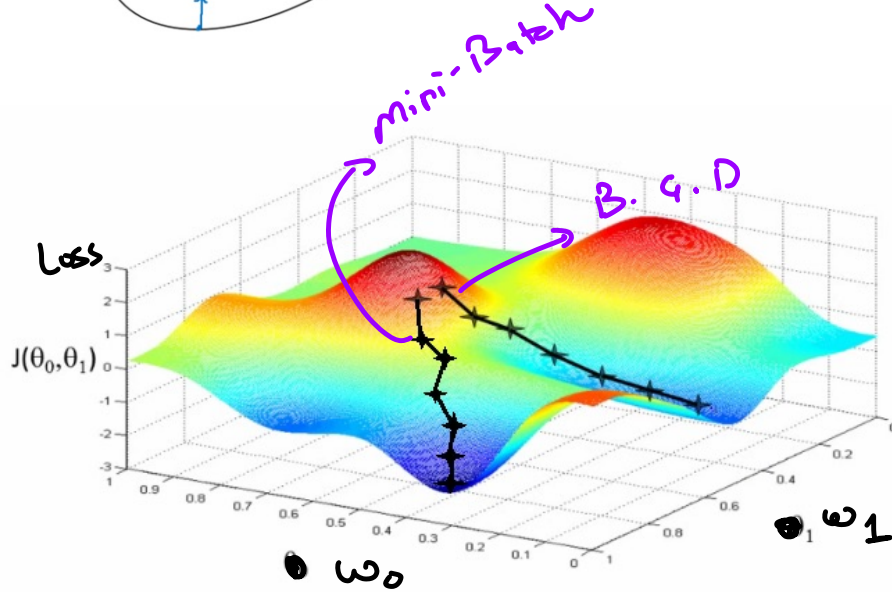
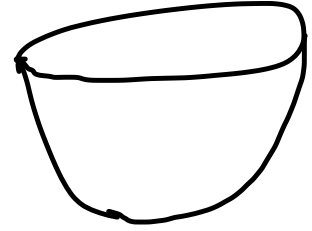
① gradients calculated using 1 data point.

② Sometimes Stucks in local minima.

③ Speed is v. fast



- Batch gradient descent
- Mini-batch gradient Descent
- Stochastic gradient descent

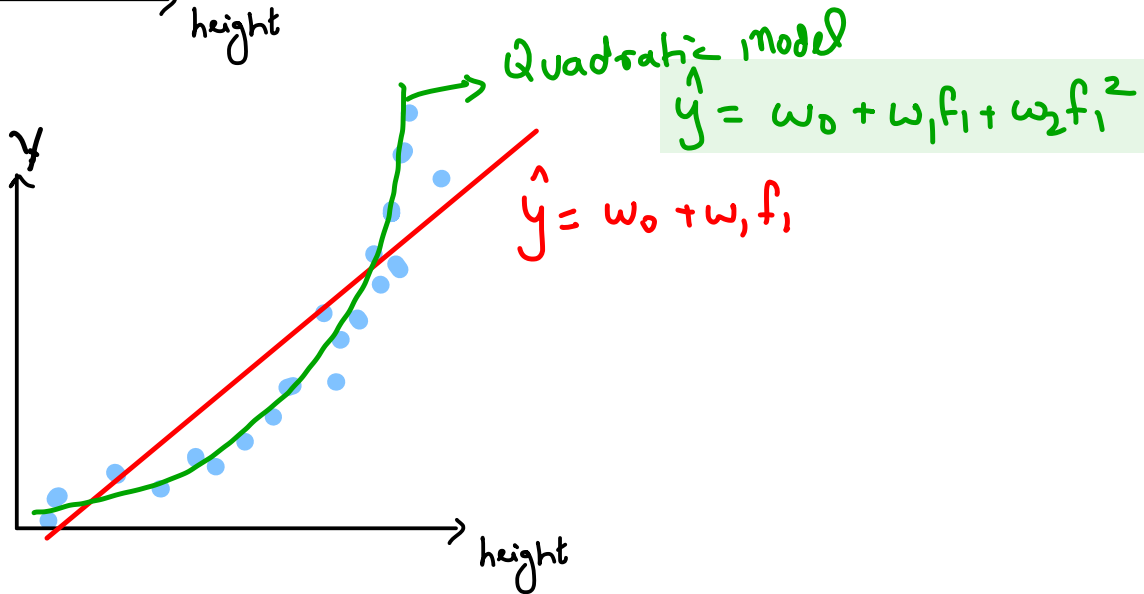
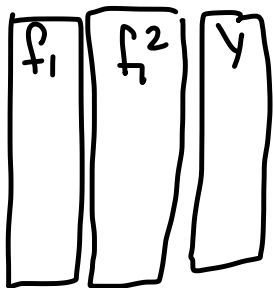
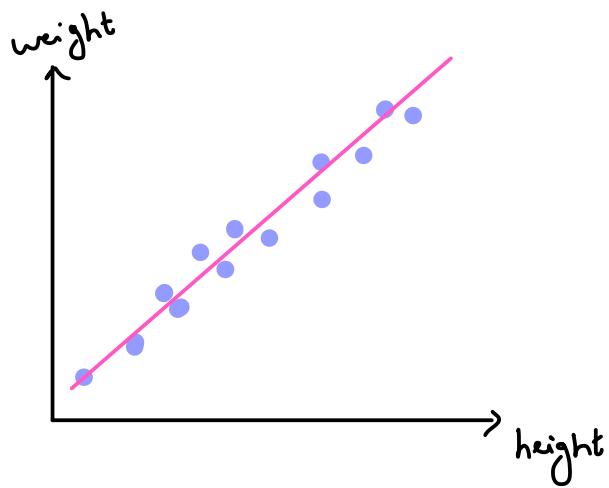


Which statement is true about mini-batch gradient descent?

58 users have participated

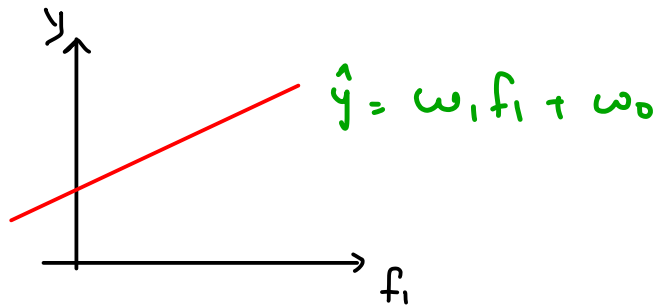
- ☒ **A** It guarantees convergence to the global minima 36%
- ☒ **B** It may converge to a local minima due to the weight fluctuations 57%
- ☒ **C** It requires a very high learning rate. 2%
- ☒ **D** It is not suitable for large datasets (opp.) 5%

[End Quiz Now](#)



f_1	y

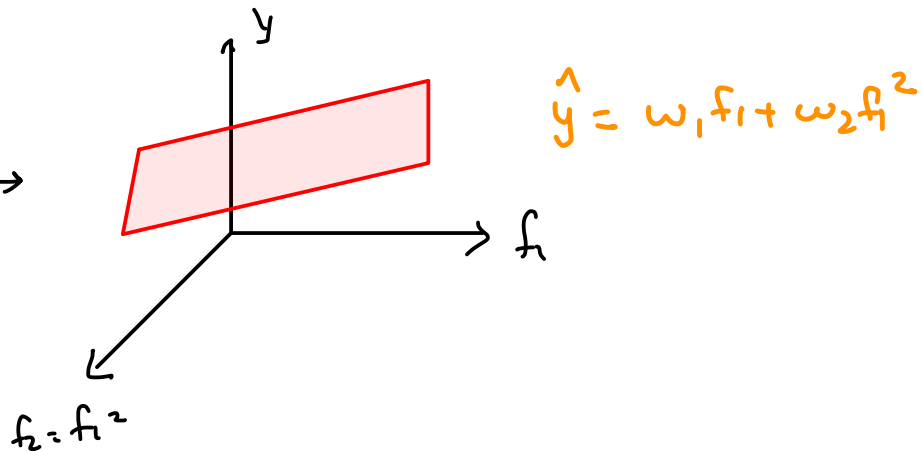
$L.R \rightarrow$



↓ transform Data

f_1	f_1^2	y

$L.R \rightarrow$



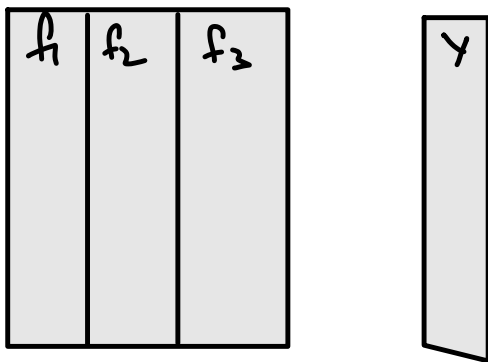
Does $f_2 = f_1^2$ cause m.c?
↓
linear comb. of other features.

$$f_2 = \alpha \cdot f_1 + \beta \rightarrow \text{m.c}$$

$$f_2 = \underbrace{f_1 \times f_1}_{\rightarrow \text{Non-linear}} \rightarrow \text{No. m.c}$$

f_1	f_1^2	y
-------	---------	-----

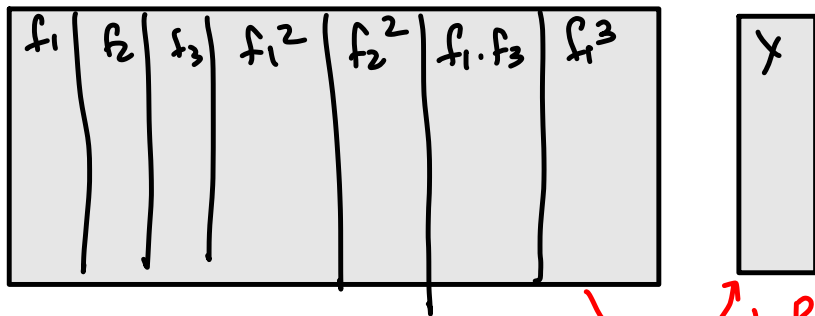
"Create Non-linear features"



$$\Rightarrow R_2 = 0.65$$

$L.R \rightarrow$

F.E (Non-linear feature) / Polyⁿ features



$$\Rightarrow R_2 = 0.79$$

$L.R \rightarrow$

POLYNOMIAL Regression

→ Linear Regression + (non-linear features)
↓
Polyⁿ features.

$$\hat{y} = w_0 + w_1 f_1 + w_2 f_2 + w_3 f_1^2 + w_4 f_2^2 + w_5 f_1^3 + w_6 f_1^4$$

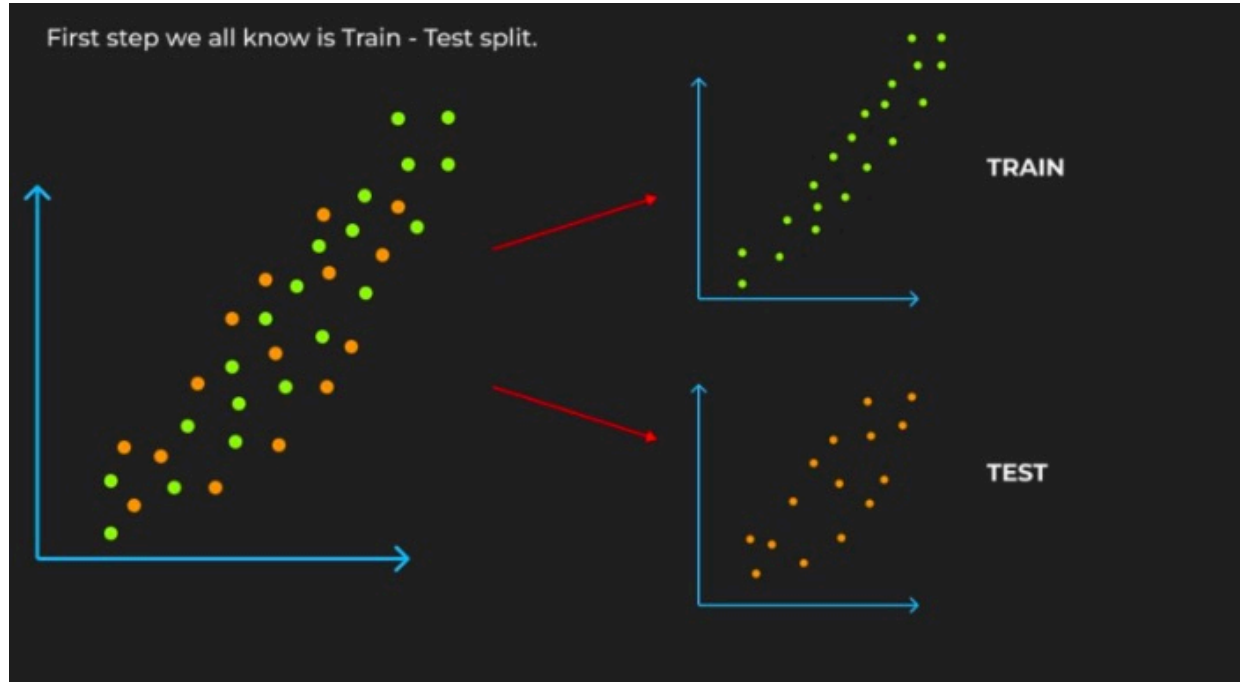
What metric should be used during Polynomial Regression ?

77 users have participated

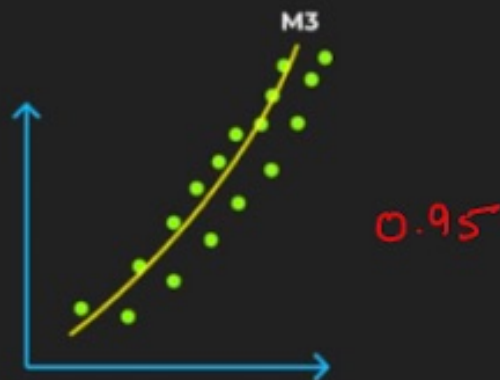
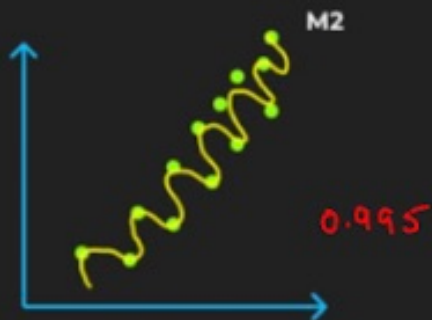
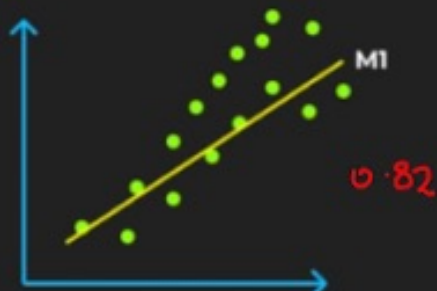
A	R-sq	13%
✓ B	Adj R-sq	84%
C	Doesnt matter	1%
D	Use a different metric	1%

[End Quiz Now](#)

Generalisation / Occam Razor



Train



Which three models out of M1, M2, M3 do you think is best ?

MODEL M1 Clearly M1 model is not fitting well on both training and test data.



Failed to
capture
patterns.

This situation is called as Underfitting.

- Imagine M1 as a student who did not study at all for exam, so bound to fail.



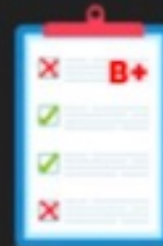
MODEL M2 M2 is performing great on training but only decent on testing.



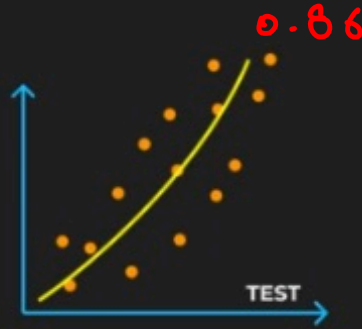
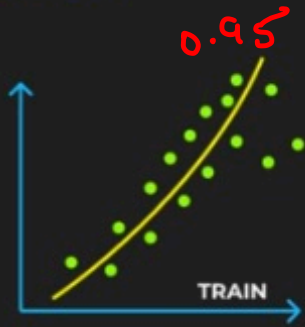
Overcaptured
the pattern
due to noise.

This is an example of Overfitting.

- Imagine M2 as a student who studied hard bt instead of understanding better, just cramped the Q/A



MODEL M3 M2 is performing great on training but only decent on testing.

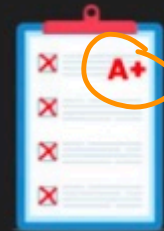


Balanced between overfitting & underfitting & underfitting

- Imagine M3 as a student who studied smartly and understood the concept.

Able to see
the generic

(Pattern)



To summarize,

		TRAINING	TESTING
<div>SIMPLEST ↑ ↓ COMPLEX</div>	UNDERFIT	POOR	POOR
	PERFECTLY FIT	GREAT	GREAT
	OVERFIT	BEST	DECENT

$$y = w_0 + w_1 f_1 + w_2 f_2 + w_3 f_1^2 + w_4 f_2^2 + w_5 f_1 f_2 + w_6 f_1^3 + \dots$$

If there are several ML model, Generalization states that:

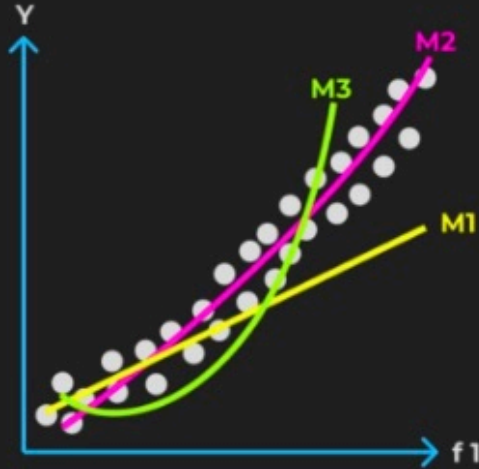
- Always choose the models that learns and understand the data inside out so that it can make good predictions on unseen data.

OCCAM's RAZOR

There is another rule that we follow while choosing best model.

**"There are many solutions to the problem,
always choose the simpler one"**





Now suppose given these 3 models :

M1 Linear: f_1

M2 Quad: f_1, f_1^2

M3 Cubic: f_1, f_1^2, f_1^3

Which one should we choose ?

Ex 806
Test Performance

M1	M2	M3
HIGH	LOW	LOW

Params


w_0, w_1, w_2

Params

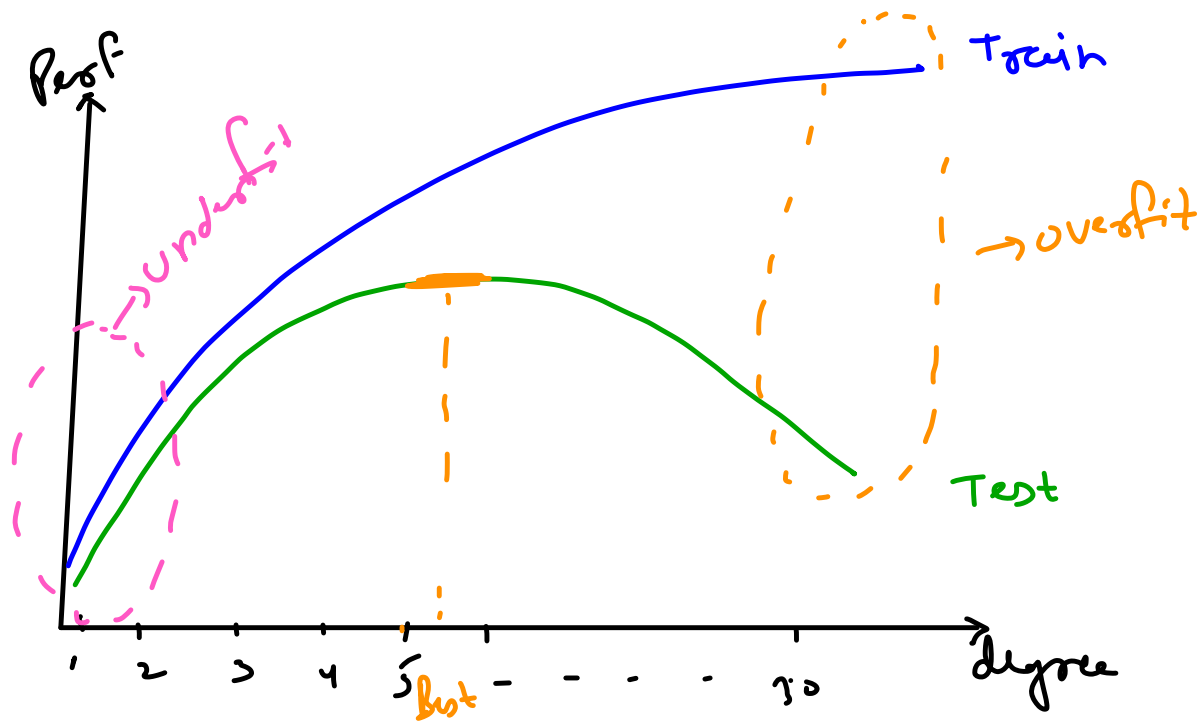
w_0, w_1, w_2, w_3

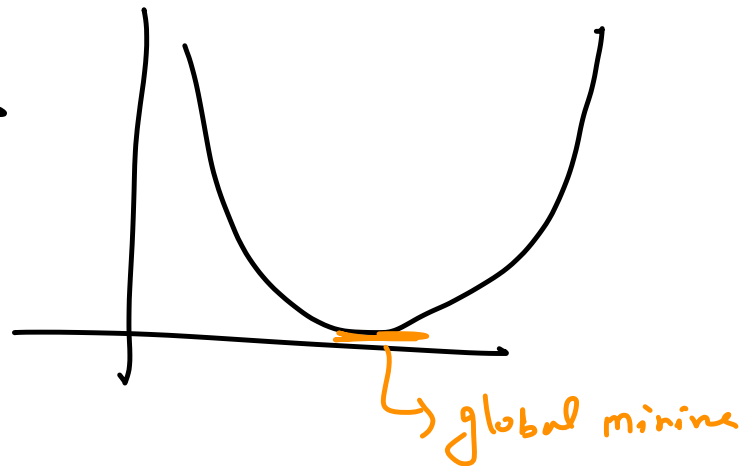
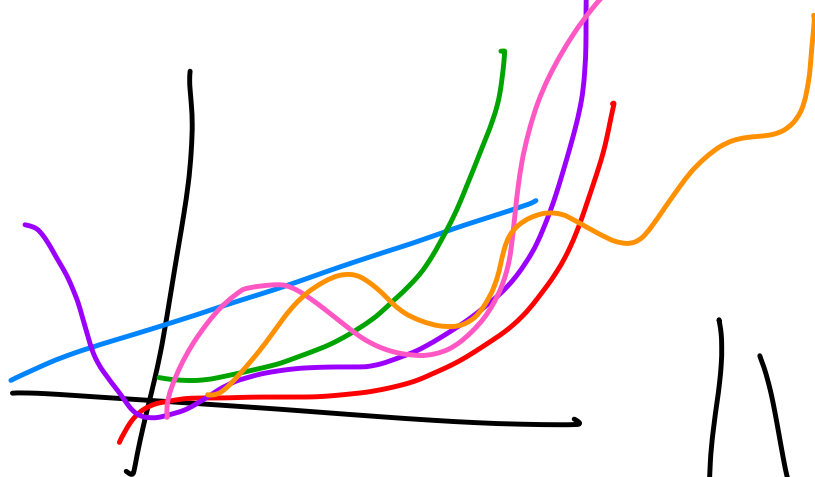
Why is Occam's Razor important in machine learning?

70 users have participated

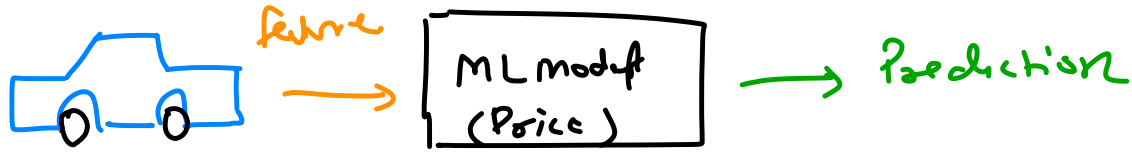
- | | | |
|---|--|-----|
| A | It helps in selecting the model that fits the training data perfectly. | 21% |
| B | It encourages the use of complex models. | 1% |
|  C | It helps in avoiding overfitting by favoring simpler models. | 76% |
| D | It promotes the use of large datasets for training models. | 1% |

[End Quiz Now](#)



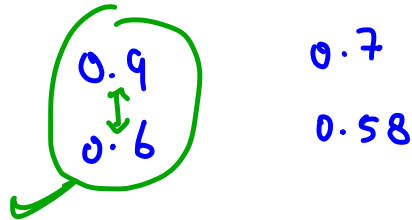


Inference.

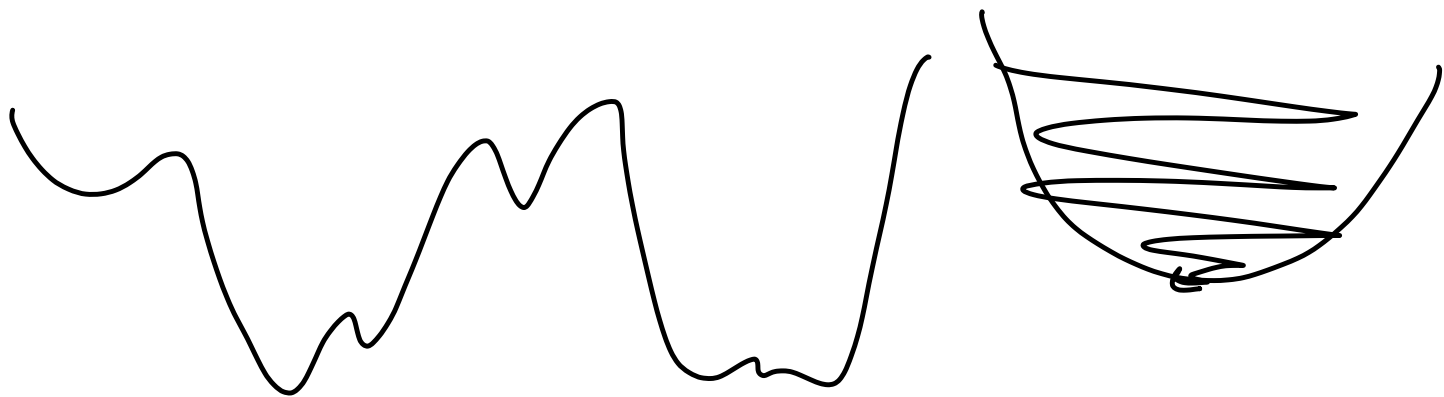


Pipeline.

1. Preprocess. → Encoding, Cat → Numerical
2. Scaling



improve overfitting



GD f momentum