

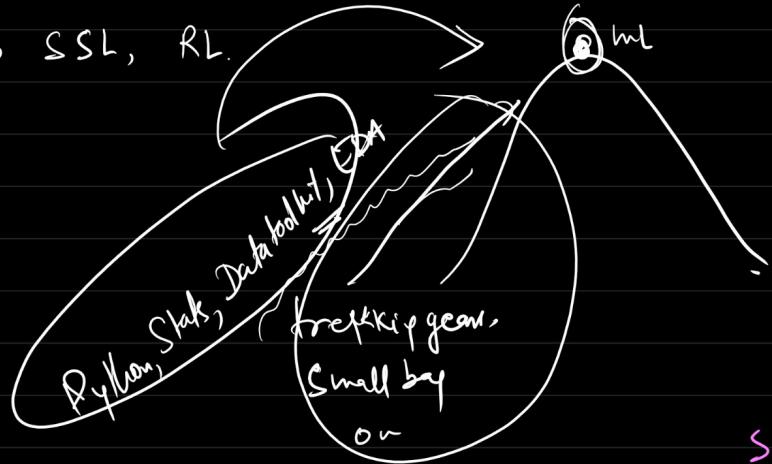
Till now

→ Intro to ML

→ SL, VSL, SSL, RL.

Agenda →

ML-2



SL

ML

	$x_1$ # of hours studied	$x_2$ # marks in internal exam	$y$ Pass / fail.
<del>train</del>	5	80	1
	3	70	0
	1	40	1
mathematical relationship	2	45	0
	3	50	1

How to verify?

→ Test the model on some other data

Unseen data

train-test split

date

take this model to  
business team  
and say

↓  
Verify this model

reason (money, time)

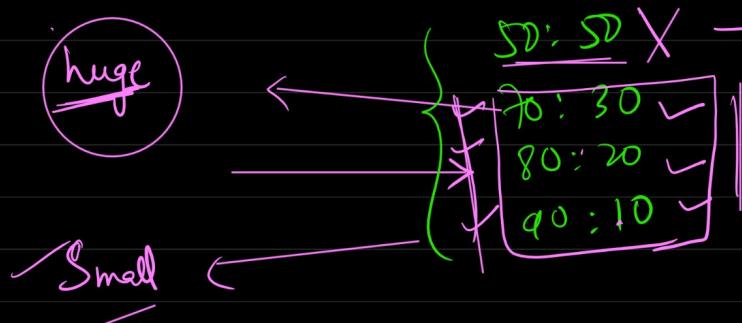
$X_1$ # of hours studied	$X_2$ # marks in internet exam	Pass/Fail.
5	80	1
3	70	0
1	40	1
2	45	0
3	50	1

train [ 5, 3, 1, 2 ]  
 test: [ 3 ]

test → representative of Unseen data.  
 → Used to test model accuracy.

1000 → ratio : train : test split ?

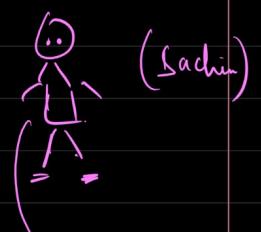
ML  
 ↓ learn patterns  
 ↓ train data



train > test data

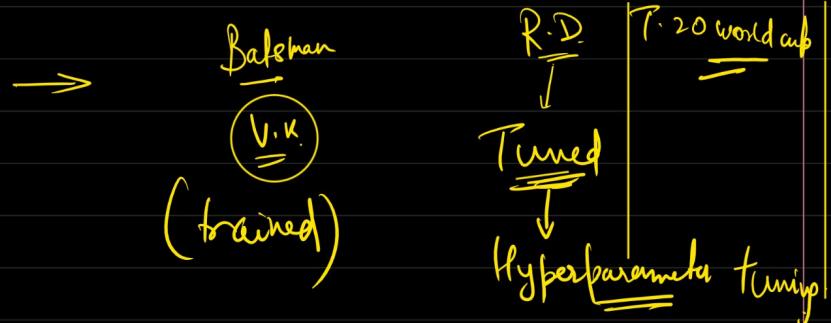
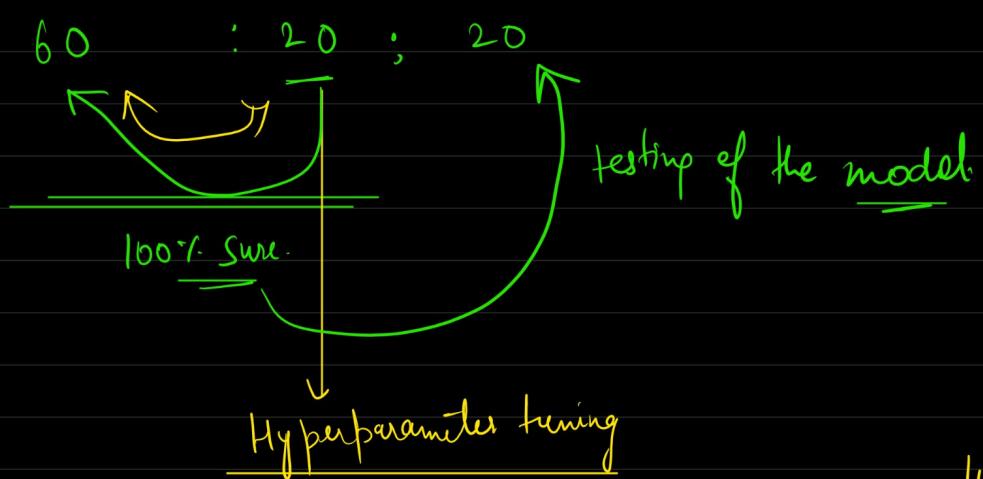


(Permutation  
 2  
 combination)



DS/ML  
 ↘ mock  
 interview  
 practice

appear  
 for interview



Student

parameter - mathematical relationship

eg. Studying from (NCERT)

Solve sample papers

Board exam

training

Hyperparameter tuning

testing of model

eg. guitar tuning

Key takeaway.

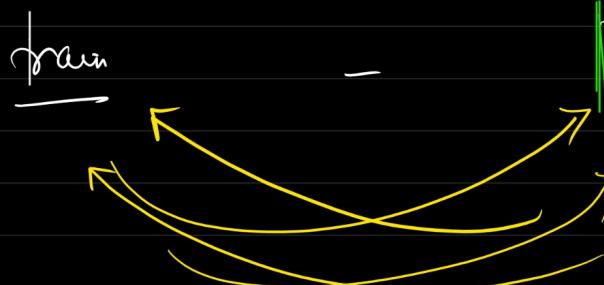
→ train - test

→ train - val - test

60 : 20 : 20

\* Validation date will act as test while training

Hyperparameter tuning ??



Test  
(Unseen data)

→ Ever board exam papers  
are given to students.

→ Jailed  
because  
you are  
cheating.

\* data leakage

You have passed test  
Information of the training  
of the model  $\Rightarrow$  model will  
perform well on both  
train and test.



One day you will deploy  
it.



Model performs  
worst in real world  
Scenario



huge loss  $\Rightarrow$  You  
are fined

\* Machine learning framework

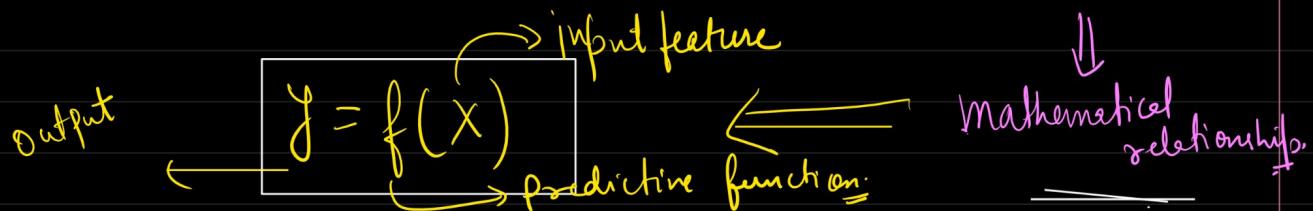


\* framework  $\rightarrow$  set of rules  
a predefined process  
set of tools  
guideline  
standardized process

# Area of house	# of rooms	Price of house
1100	2	2.1
-	-	-
-	-	-
1200	3	?

y  
train 80 → training of  
model  
test 20

✓ Under finding the pattern



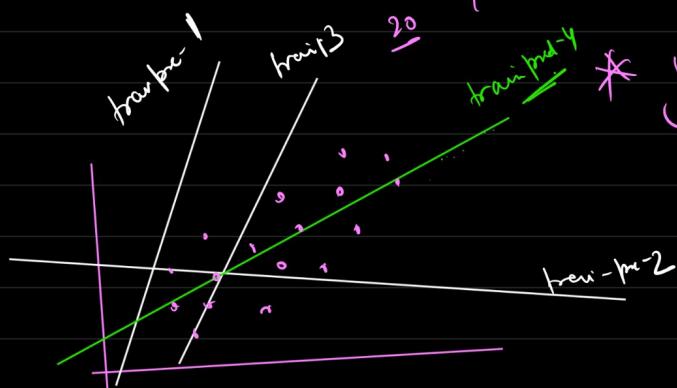
\* Model training → given the data estimate the prediction function by minimizing error =  $\downarrow$   
 ↓  
 evaluation metrics

$y = f(x)$

← training 80 {

# Area of house	# of rooms	Price of house	$y_{\text{train prediction}}$
1100	2	2.1	1.8
-	-	-	2.1
-	-	-	2.3
-	-	-	2.4
1200	3	?	

} Exposed data



\* you want to learn the optimal relationship.

$$y = \underline{f(x)}$$

\* key takeaway

→ train | test.  
 $\rightarrow y = f(x)$

## Overfitting & Underfitting

\* Overfitting

train → Model is trained  $\Rightarrow$  Accuracy is  $\uparrow$  (95%)

test → Model is tested  $\Rightarrow$  Accuracy is  $\uparrow$  (high)

Vishwa Sir. (10 Question)



(unseen)

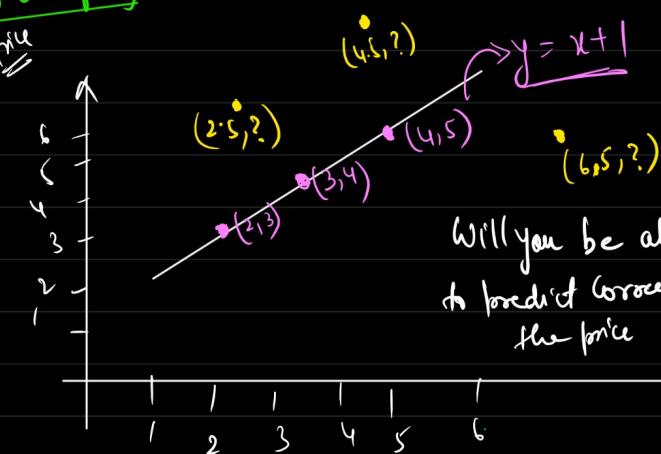
- low  $\downarrow$  Exam  $\rightarrow$  who will perform better in this exam?
- question  
Model performs well on train data but worse on test data.

## \* Underfitting

train  $\rightarrow$  accuracy  $\downarrow$  50%  
test  $\rightarrow$  accuracy  $\downarrow$  40% } Underfitting

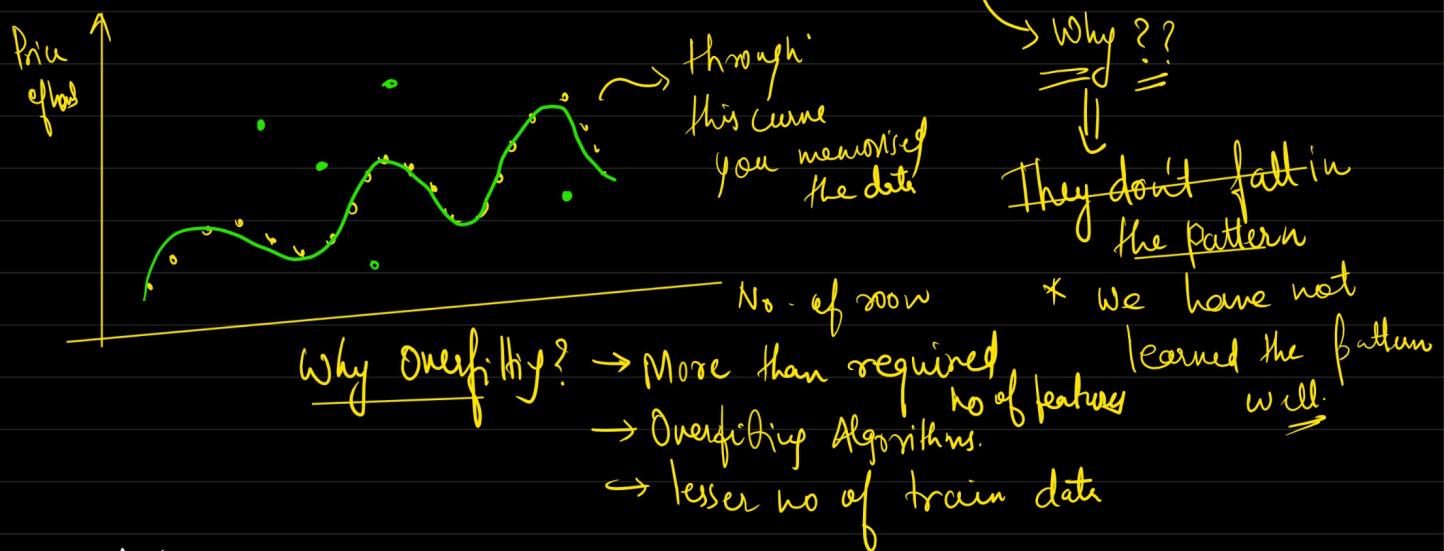
## \* Generalised model

### Expl. of overfitting



Will you be able to predict correctly the price for yellow points?

# No. of rooms (x)	Price of house (y) (in rupees) (y)
2	3
3	4
4	5
5	5



## Underfitting

train & test accuracy  $\downarrow$

model has not learned well

Why underfitting

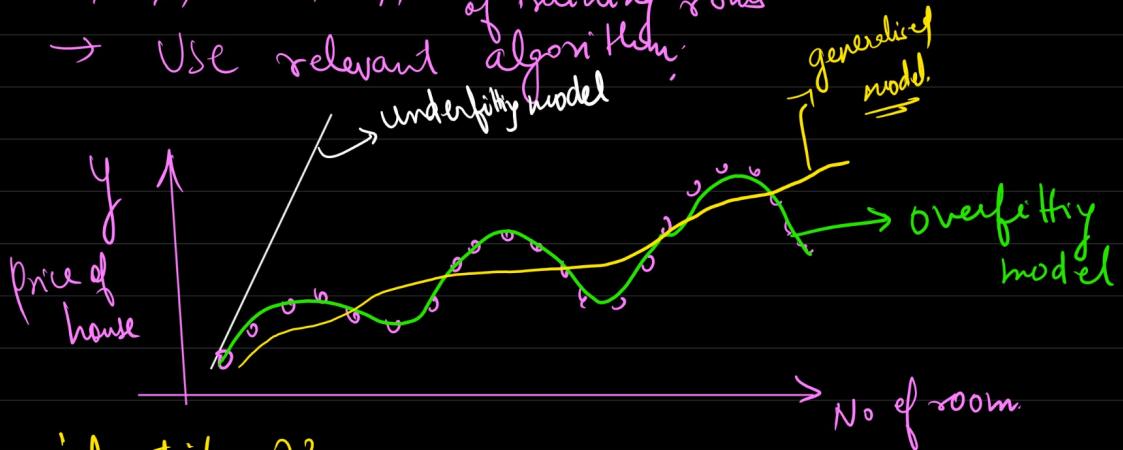
## Generalised model

→ give required no of features

→ " " " of training rows

→ Use relevant algorithm.

- less no of required features
- less no of train data
- Underfitting algo thrills



## \* How to identify ??

\* Underfitting ≈ Model is not performing well while training.

\* Overfitting ≈ train acc↑, test ↓ (difference more than 5%)

\* Generalised model ≈ train↑ test↑

$$\boxed{65, 64} \text{ diff is } \approx \text{ range } 5\%$$

## Bias | Variance

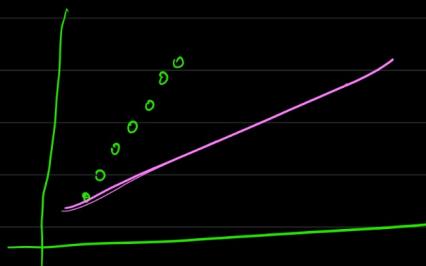
\* training error is also known as bias.

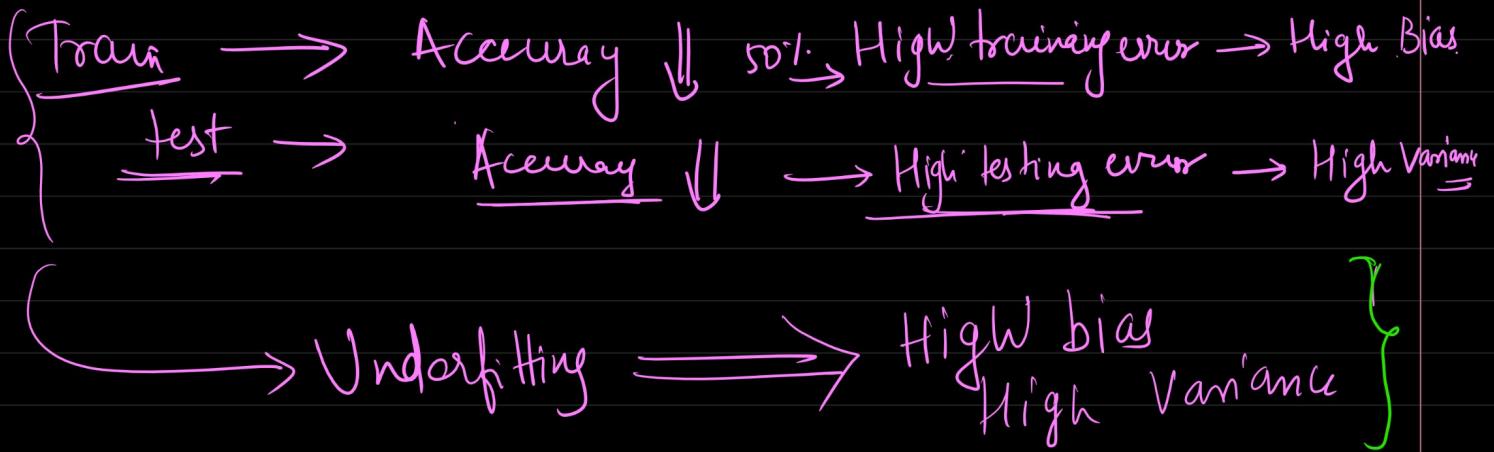
High training error means high bias

\* testing error is also known as Variance.  
High testing error means high variance

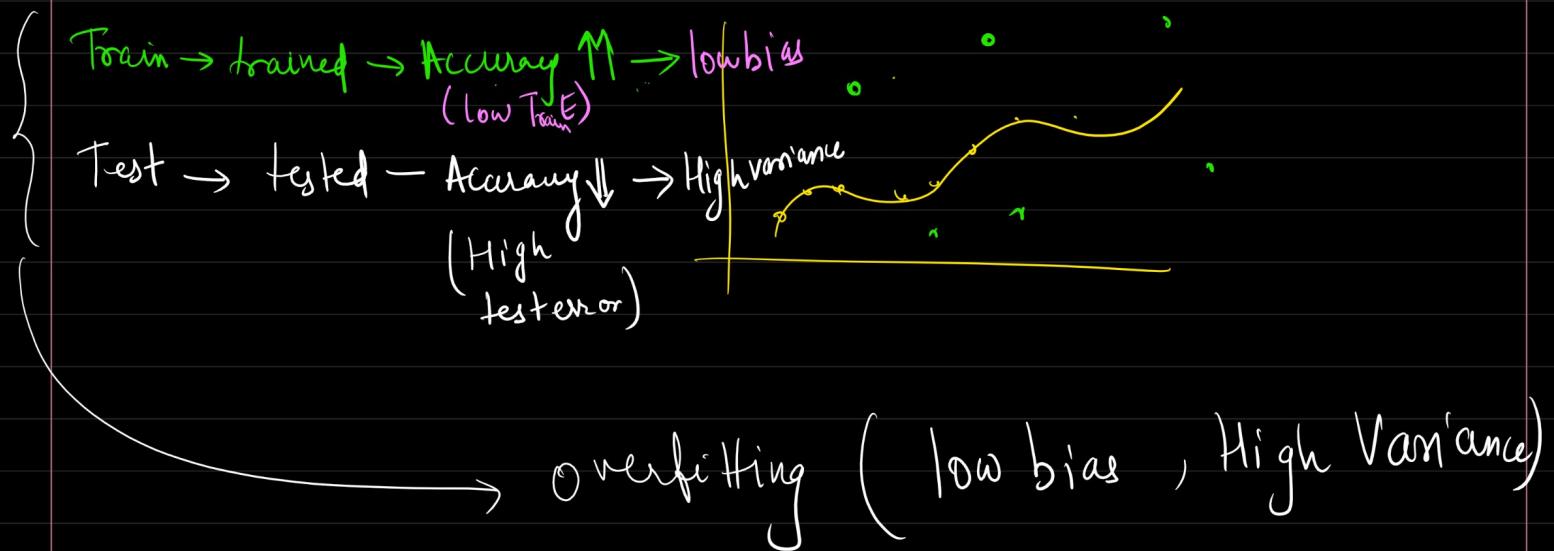
bias → assumption | prejudice      Why ??

Because you assumed the dist of data wrong





Variance → Why Variance is analogous to High testing error?  
 → test date varies from the pattern  
 that's why high testing error

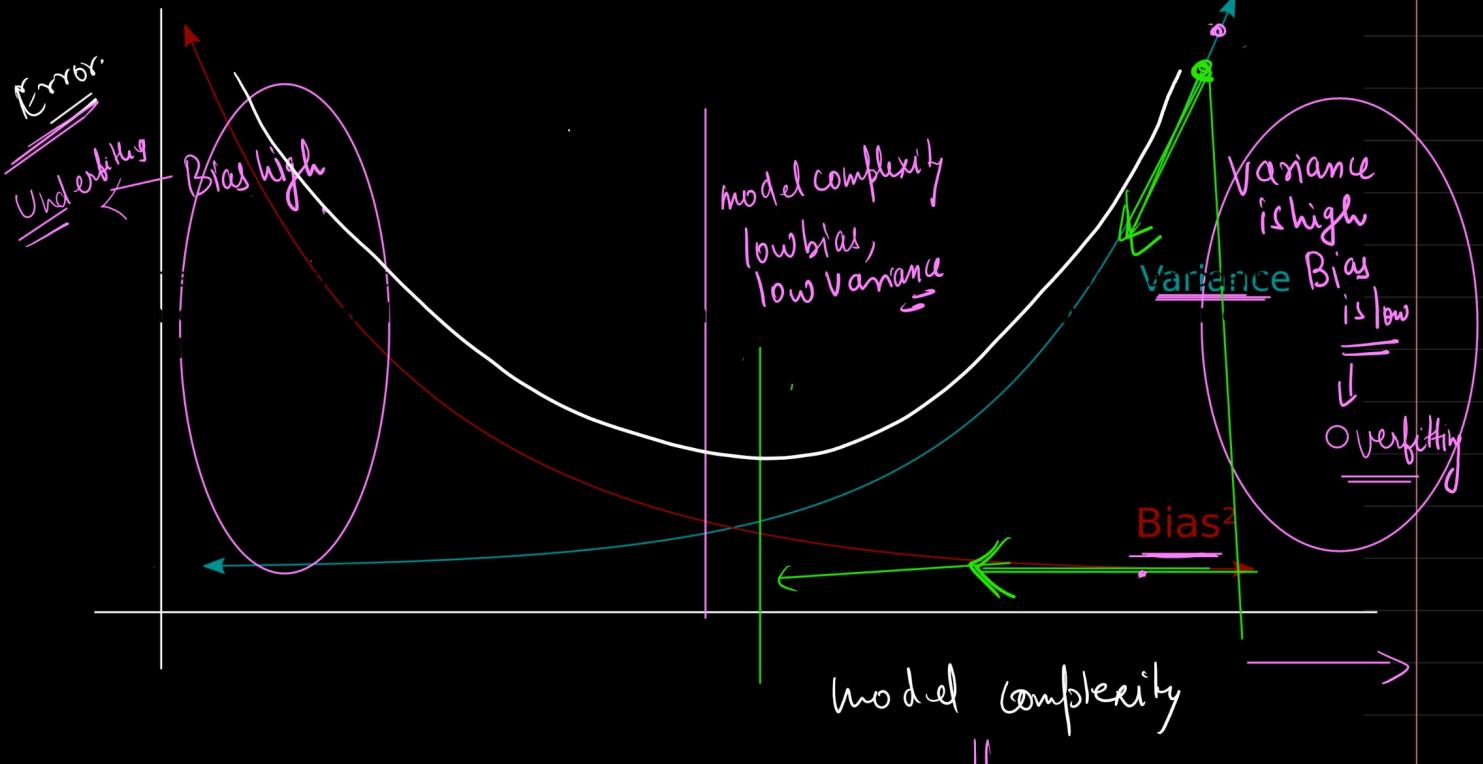


\* Generalised model  
 $(\text{low bias, low variance})$       train acc  $\uparrow$       (low bias)  
 $\qquad\qquad\qquad$       test acc  $\uparrow$       (low Variance)

## Bias - Variance tradeoff

In statistics and machine learning, the bias-variance tradeoff describes the relationship between a model's complexity, the accuracy of its predictions, and how well it can make predictions on previously unseen data that were not used to train the model.

$$\text{total error} = \text{Bias} + \text{Variance}$$



- \* In order to reduce Variance (reduce the testing error), introduce some bias in training of model - by selecting relevant features or selecting the right algorithms.
  - \* In order to reduce high bias  $\rightarrow$  Use some more d.p's, more feature engineering, Use some other MI algorithms.
- Model complexity  $\rightarrow$  You have used more than required features.