

## Encoding the Categorical features

Categorical features Example: 1) Marital Status : [unmarried, married, divorced, widowed]

2) Colour of a car : [Red, Blue, Sky blue, Grey, Black etc.]

3) Smoking habit : [smoker, non-smoker]

for binary features : 0, 1       $0 \rightarrow \text{non-smoker} ; 0 \rightarrow \text{female}$   
   $1 \rightarrow \text{smoker} ; 1 \rightarrow \text{male}$

for feature having multiple possible values / categories :

Label Encoding :

unmarried	$\rightarrow 0$
married	$\rightarrow 1$
divorced	$\rightarrow 2$
widowed	$\rightarrow 3$

## One-Hot Encoding :-

marital status (MS) }  
 [U, M, D, W]

<u>Observation</u>	<u>MS</u>
1	U
2	D
3	M
4	W

$$U = [1, 0, 0]$$

$$D = [0, 0, 1]$$

$$M = [0, 1, 0]$$

$$W = [0, 0, 0]$$

## Dummy Variables.

MS\_U, MS\_M, MS\_D, MS\_W

	<u>MS_U</u>	<u>MS_M</u>	<u>MS_D</u>	<u>MS_W</u>
1	1	0	0	0
2	0	0	1	0
3	0	1	0	0
4	0	0	0	1

Redundant

n - features

(n-1) : non-redundant

Binary number Encoding:  $(4)_{10} = (100)_2$ ;  $(12)_{10} = (1100)_2$

$$\begin{array}{ccccc} \frac{2^4}{1} & \frac{2^3}{0} & \frac{2^2}{0} & \frac{2^1}{1} & \frac{2^0}{1} \\ (\text{1}) & (\text{0}) & (\text{0}) & (\text{1}) & (\text{1}) \end{array}_2 \rightarrow 1 \times 2^4 + 0 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 = 16 + 0 + 0 + 2 + 1 = (19)_{10}$$

3.7

$$32 \rightarrow 2^5$$

$$(37)_{10} = (100101)_2$$

4 bits

5

$$4 \rightarrow 2^2$$

$$\begin{array}{l} 1111 \\ 2^3 + 2^2 + 2^1 + 2^0 \\ = 2^4 - 1 \end{array}$$

1

$$1 \rightarrow 2^0$$

2-bits

$$\text{max value: } \underline{(2^k - 1)}$$

n - different categories

Label Enc.  $\left\{ \begin{array}{l} 1 \\ 2 \\ 3 \\ \vdots \\ n \end{array} \right.$   $\lceil \log_2(n) \rceil \rightarrow$  bits are required to represent those categories.

### Marital Status

<u>Values</u>	<u>Label Enc</u>	<u>binary Enc</u>	<u>elig<sup>n</sup></u>	$\frac{E_2}{E_1}$	$E_1$
U	0	0 0		0 0	
M	1	0 1		0 1	
W	2	1 0		1 1	
D	3	1 1		1 1	

Multi-class Classification ( $C$ -class) [ $C > 3$ ]

One-vs-Rest (OVR) / One-vs-All (OVA)

train  $C$ - no. of binary classifiers. Each of them will be trained for one class.

## Logistic Regression (binary / Dichotomous)

$\underline{x}$  (features)

$\underline{x_1} \quad \underline{x_2} \quad \underline{x_3} \quad \dots \quad \underline{x_n}$

$\underline{y}$  (target)

$\underline{y}$

logit

$\underline{z}$

$$z = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\text{logit}(x) = \ln\left(\frac{x}{1-x}\right)$$

inverse of sigmoid

Predicted

$$\underline{\hat{y}_p} = \sigma(z)$$

Multinomial Regression / Softmax Regression / Polytomous Regression

Maximum Entropy classifier :

features

$\underline{x_1} \quad \underline{x_2} \quad \underline{x_3} \quad \dots \quad \underline{x_n}$

$x_1^{(1)} \quad x_2^{(1)} \quad x_3^{(1)} \quad \dots \quad x_n^{(1)}$

$x_1^{(i)} \quad x_2^{(i)} \quad x_3^{(i)} \quad \dots \quad x_n^{(i)}$

One-hot Encoded target

$\underline{y_1} \quad \underline{y_2} \quad \dots \quad \underline{y_C}$

0 1 0 0 0

0 0 0 1 0

Logit

$$\ln\left(\frac{x}{1-x}\right)$$

$$\text{sigmoid}\left(\ln\left(\frac{x}{1-x}\right)\right) \\ = z$$

features       $x_1$      $x_2$      $x_3$     ...     $x_n$

targets       $y_1$      $y_2$     ...     $y_c$

Logits       $z_1$      $z_2$     ...     $z_c$

$$z_k = \theta_{k,0} + \theta_{k,1} x_1 + \theta_{k,2} x_2 + \cdots + \theta_{k,n} x_n$$

$$z_k = \vec{\theta}_k^T \vec{x} + \theta_{k,0} \quad \left[ \vec{\theta}_k = [\theta_{k,1}, \theta_{k,2}, \dots, \theta_{k,n}]^T \right]$$

$i^{th}$  observation

$x_1$      $x_2$      $x_3$     ...     $x_n$

$x_1^{(i)}$      $x_2^{(i)}$      $x_3^{(i)}$     ...     $x_n^{(i)}$

$z_1$      $z_2$     ...     $z_k$     ...     $z_c$

Logits

$z_1$      $z_2$     ...     $z_k$     ...     $z_c$

(.)    (.) ...  $\vec{\theta}_k^T \vec{x} + \theta_{k,0}$  .. (.)

2	-0.5	1	0.7
---	------	---	-----

$$\begin{matrix} z_1 & z_2 & z_3 & \dots & z_c \end{matrix} \quad \sigma([z_1, z_2, z_3]) = (\hat{y}_1, \hat{y}_2, \hat{y}_3) \\ \begin{matrix} 0 & 2 & -1 \end{matrix} \quad \text{s.t. } \boxed{\hat{y}_1 + \hat{y}_2 + \hat{y}_3 = 1}$$

$$0.5 \quad 1.2 \quad 0.1 \quad \text{and } 0 \leq \hat{y}_i \leq 1$$

$$\vdots \quad \vdots \quad \vdots \quad \sigma(z) = \hat{y} \quad \text{s.t. } \sum_i \hat{y}_i = 1 \quad \& \quad \boxed{0 \leq \hat{y}_i \leq 1}$$

$$[0, 2, -1] \xrightarrow{e^x} [e^0, e^2, e^{-1}] \xrightarrow{\text{normalize}} \left[ \frac{e^0}{e^0 + e^2 + e^{-1}}, \frac{e^2}{e^0 + e^2 + e^{-1}}, \frac{e^{-1}}{e^0 + e^2 + e^{-1}} \right]$$

$$\hat{y}_1 \quad \hat{y}_2 \quad \hat{y}_3$$

$$z_1, z_2, \dots, z_c$$

$$\hat{y}_k = \sigma(z_k) = \frac{e^{z_k}}{\sum_{j=1}^c e^{z_j}}$$

} Softmax function.

$$\sigma(z_k) = \frac{e^{z_k}}{\sum_{j=1}^c e^{z_j}}$$

$$z = z_1 - z_2$$

$\frac{z_1}{1}$	$\frac{z_2}{0}$	$\frac{z_1 - z_2}{1}$
0	1	(-1)

if  $c = 2 \}$  binary classification

$$\sigma(z_1) = \frac{e^{z_1}}{e^{z_1} + e^{z_2}}$$

$$= \frac{1}{1 + e^{-(z_1 - z_2)}}$$

$$\sigma(z_2) = \frac{e^{z_2}}{e^{z_1} + e^{z_2}}$$

$$= \frac{1}{1 + e^{(z_1 - z_2)}}$$

predicted

$$\underline{\hat{y}_1} \quad \underline{\hat{y}_2} \dots \underline{\hat{y}_k} \dots \underline{\hat{y}_c} \quad | \quad \underline{y_1} \quad \underline{y_2} \dots \underline{y_k} \dots \underline{y_c}$$

↳ one-hot enc.  
target

predicted (softmax)      target      |       $\hat{y}$        $y$

<u><math>\hat{y}_1</math></u>	<u><math>\hat{y}_2</math></u>	<u><math>\hat{y}_3</math></u>	<u><math>y_1</math></u>	<u><math>y_2</math></u>	<u><math>y_3</math></u>		<u><math>\hat{y}</math></u>	<u><math>y</math></u>
0.2	0.7	0.1	0	1	0		- $[y \log \hat{y} + (1-y) \log(1-\hat{y})]$	
0.8	0.05	0.15	0	0	1			(BCE)

$$\hat{y}_1 \hat{y}_2 \dots \hat{y}_k \dots \hat{y}_c \quad | \quad y_1 y_2 \dots y_k \dots y_c$$

$$\hat{y}_1^{(i)} \hat{y}_2^{(i)} \dots \hat{y}_k^{(i)} \dots \hat{y}_c^{(i)} \quad | \quad y_1^{(i)} y_2^{(i)} \dots y_k^{(i)} \dots y_c^{(i)}$$

$$\text{Loss}^{(i)} = - \sum_{k=1}^c y_k^{(i)} \log(\hat{y}_k^{(i)})$$

$\frac{\partial}{\partial \hat{y}_1} \text{Loss}^{(i)} = \frac{y_1^{(i)}}{\hat{y}_1^{(i)}} - 1$        $\frac{\partial}{\partial \hat{y}_2} \text{Loss}^{(i)} = \frac{y_2^{(i)}}{\hat{y}_2^{(i)}} - 1$   
 $\frac{\partial}{\partial \hat{y}_3} \text{Loss}^{(i)} = \frac{y_3^{(i)}}{\hat{y}_3^{(i)}} - 1$        $\frac{\partial}{\partial \hat{y}_4} \text{Loss}^{(i)} = \frac{y_4^{(i)}}{\hat{y}_4^{(i)}} - 1$   
 $\vdots$        $\vdots$

$$\text{Loss}^{(i)} = - \sum_{k=1}^c y_k^{(i)} \log(\hat{y}_k^{(i)}) \quad \rightarrow \text{i}^{\text{th}} \text{ instance.}$$

total loss = average of all the losses

$$J = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^c y_k^{(i)} \log(\hat{y}_k^{(i)})$$

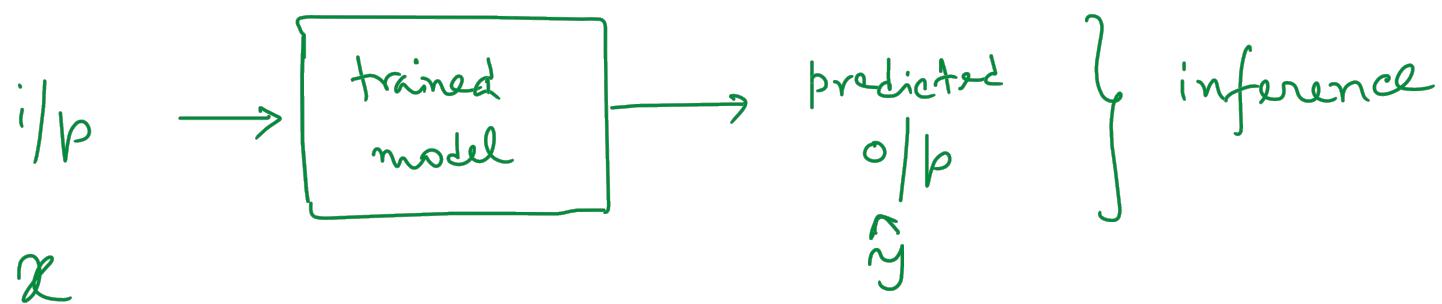
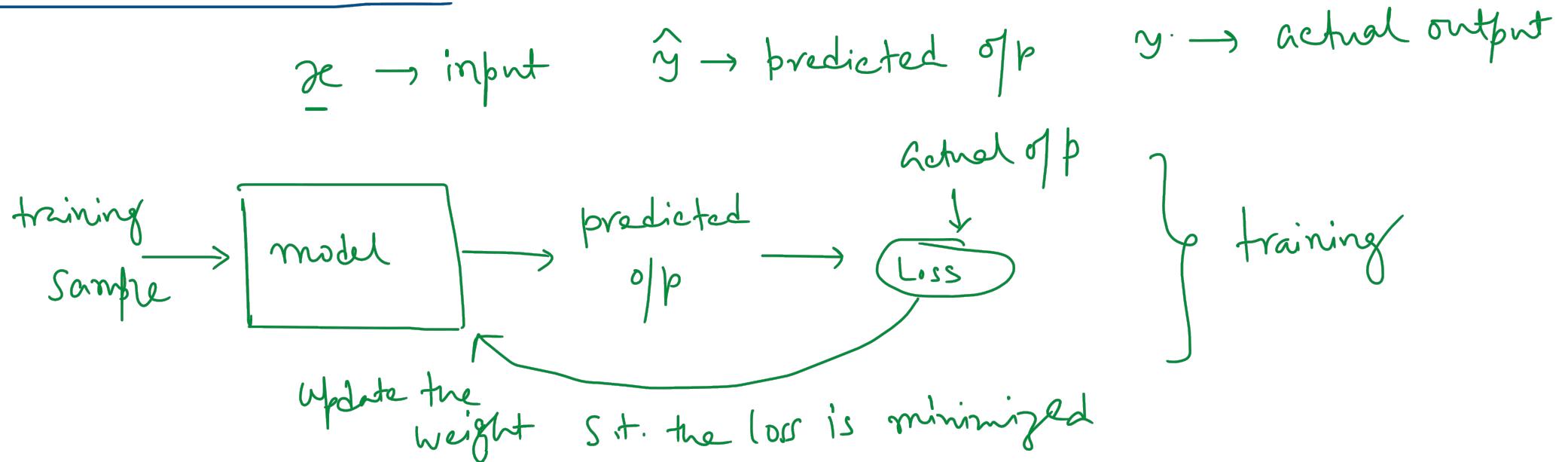
Loss function for  
multi-class classification.

Categorical Cross Entropy Loss (CCE)

$$\sum_{k=1}^c \} \text{ all the classes}$$

$$\sum_{i=1}^m \} \text{ all the observations.}$$

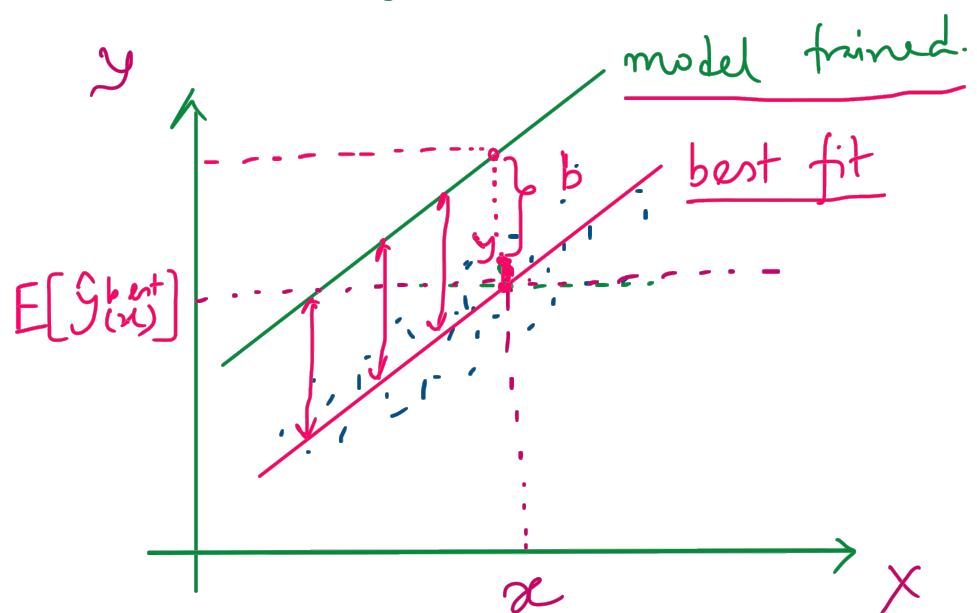
## Bias & Variance



Bias:  $\text{Bias}(\text{model}) = \mathbb{E}[\hat{y}] - \underline{y}$

$\mathbb{E}[\hat{y}] \rightarrow$  the expected (average) prediction of the model for input  $x$  over different training dataset.

$\underline{y} \rightarrow$  the true value for input  $x$ .



High bias: Our predicted o/p will be widely different from actual o/p.

Low bias: Our predicted o/p is close to the actual o/p.

Variance:

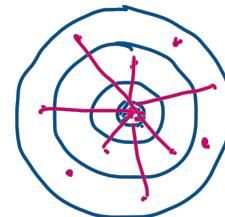
$$\mathbb{E}[(\hat{y} - \mathbb{E}(\hat{y}))^2]$$

Bias

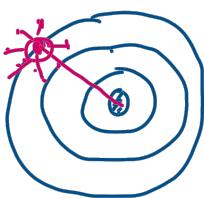
- Error between avg. model prediction & ground truth (actual value)
- $\mathbb{E}[\hat{y}] - y$
- We want a low bias because more bias means more deviant the model is from actual o/p

Variance

- Average variability of the model prediction
- $\mathbb{E}[(\hat{y} - \mathbb{E}[\hat{y}])^2]$
- We want to reduce variance because we want our model to be robust, i.e. not affected by slight variation in i/p data.

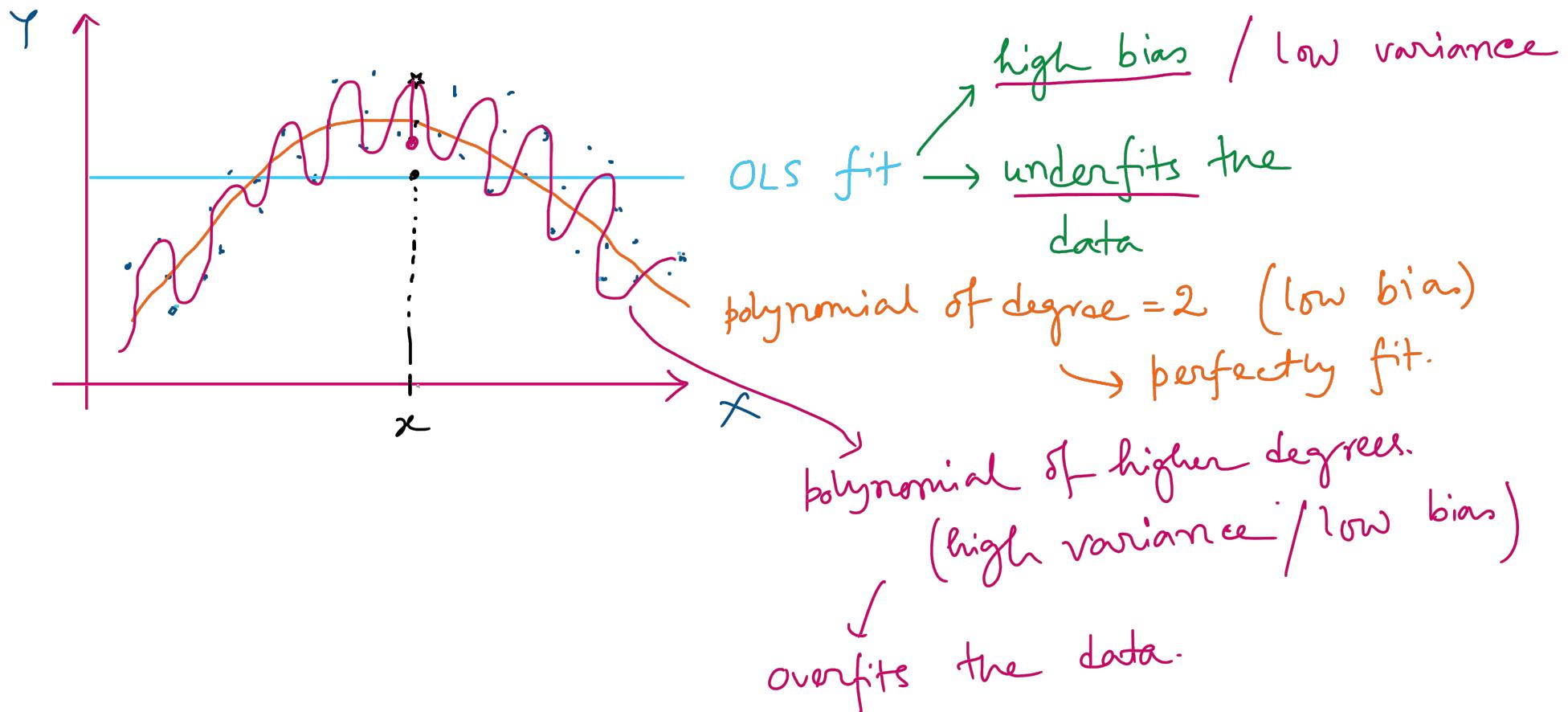


low bias  
high variance



high bias  
low variance

## Problems with high bias & high variance

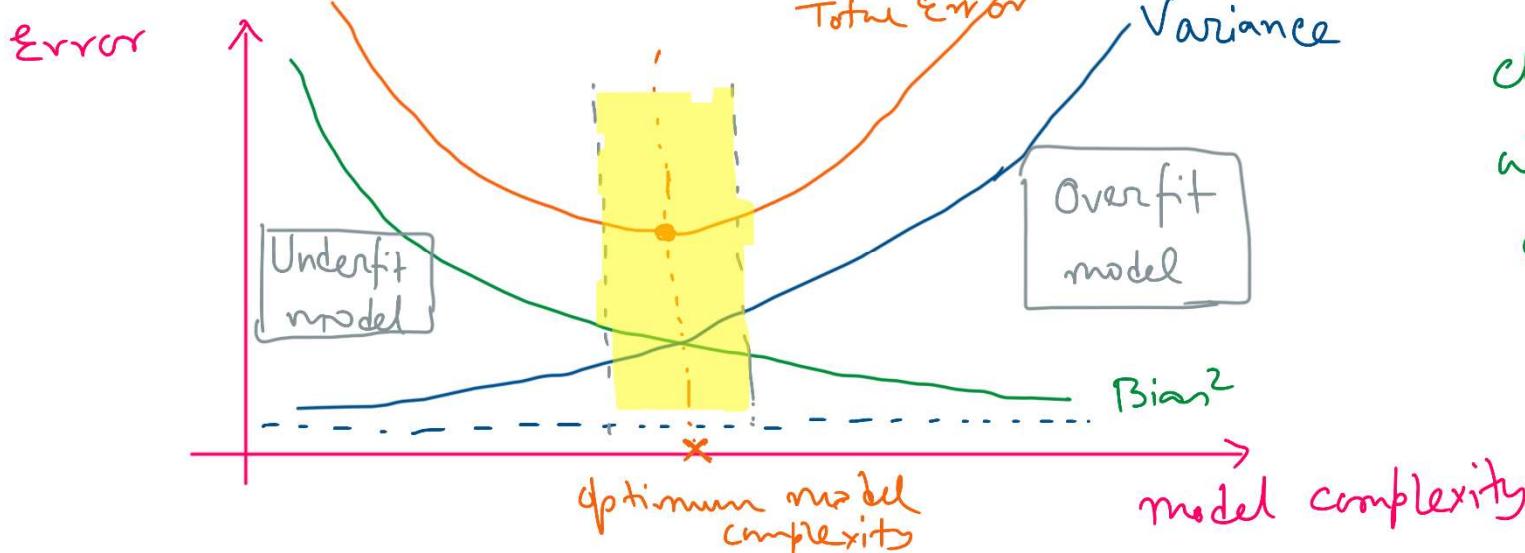


## High Bias

- 1) Simpler models
- 2) Underfits the data
- 3) High error in both training & test dataset

(MSE)

$$\text{Total Error} = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$



## High Variance

- 1) Complex model.
- 2) Overfits the data
- 3) Low error in training / High error in testing dataset

We should always choose the model which has lowest error in test dataset