

→ Gradient Descent Variants

→ Polynomial Regression

→ Underfitting & Overfitting

→ Bias - Variance

Repeat until convergence

$$w_j = w_j - n \frac{\partial L}{\partial w_j}$$

$$\hat{y} = w_0 + w_1 u$$
$$= w_0 + \sum_{j=1}^{d=5} w_j u_j$$

$$w_j = w_j - n \sum_m \sum_{i=1}^m (\gamma^{(i)} - \hat{y}^{(i)}) u_j^{(i)}$$
$$= w_0 + w_1 u_1 + w_2 u_2$$
$$- - - w_5 u_5$$

Repeat until convergence

$$\left\{ \begin{array}{l} \check{w}_j = w_j - n \frac{\partial L}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m (\gamma^{(i)} - \hat{\gamma}^{(i)}) x_j^{(i)} \\ \end{array} \right.$$

$\gamma^{(i)}$  = target

$\hat{\gamma}^{(i)}$  = Prediction

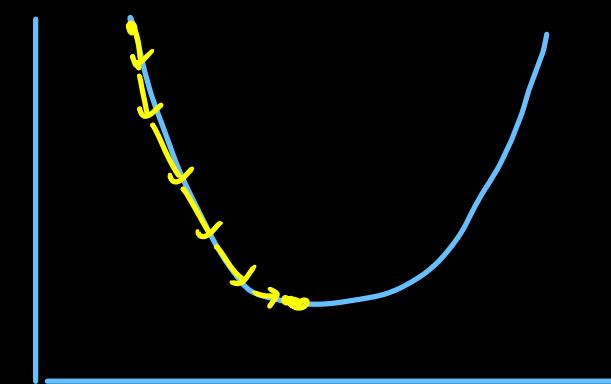
$w_i$  = weight of feature i

$n$  = Learning rate

$$\rightarrow \hat{\gamma}^{(i)} = w_0 + w_1 x_1^{(i)}$$

data point

	$x_1$	$x_2$	$x_3$	$y$
$x_1^{(1)}$	$x_2^{(1)}$	$x_3^{(1)}$	$y^{(1)}$	← 1st sample
$x_1^{(2)}$	$x_2^{(2)}$	$x_3^{(2)}$	$y^{(2)}$	← 2nd sample
$x_1^{(m)}$	$x_2^{(m)}$	$x_3^{(m)}$	$y^{(m)}$	



Batch GD

for iter in range (iteration):

{

    //

$$w_j = w_j - n \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}) x_j^{(i)}$$

}

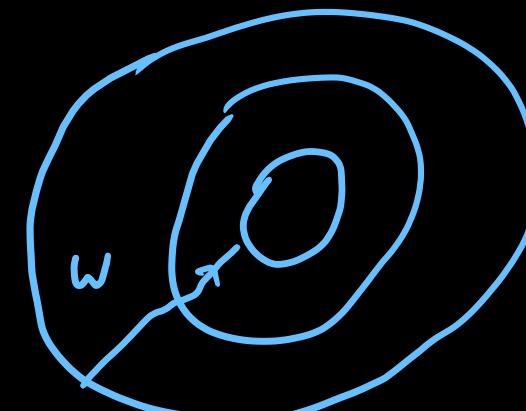
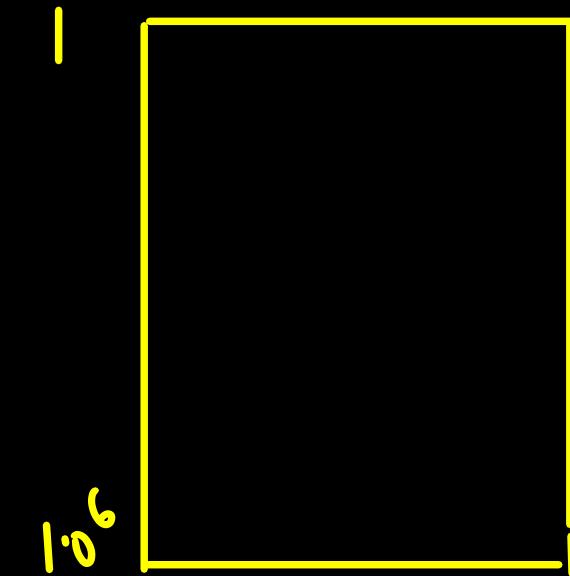
→ Huge memory in RAM

→ Very less updates

1000

11

$m = 10^6$



SGD (Stochastic GD)

for iter in range (iteration):

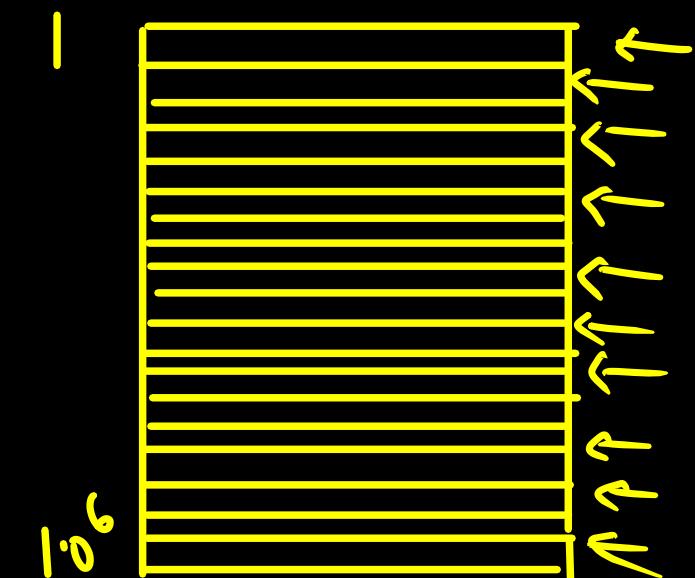
{

for i in range (m):

$$\rightarrow w_j = w_j - \eta (y^{(i)} - \hat{y}^{(i)}) u_j^{(i)}$$

→ Very large update for  $w_j$



Chunks = Batch size

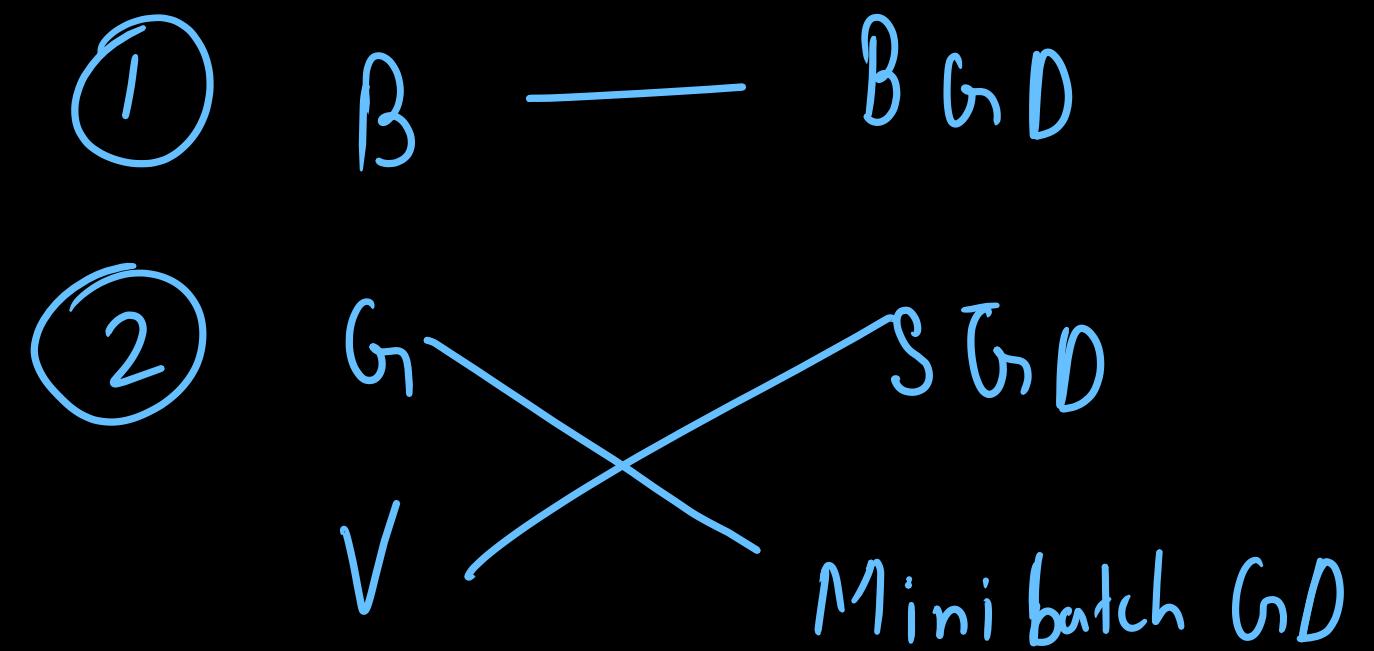
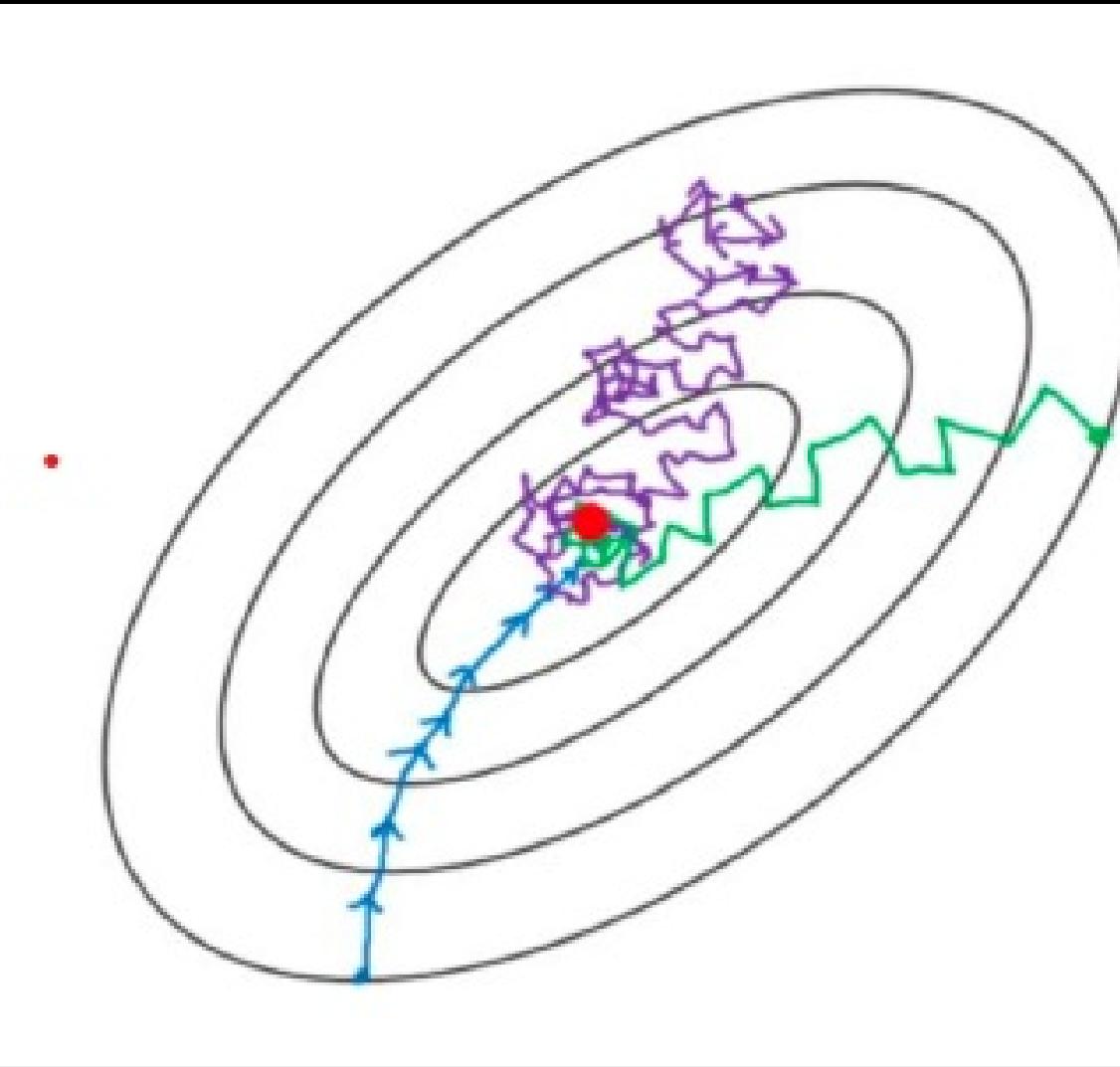
Mini batch

```
for iter in range (iteration):  
{  
    for i in range (No of batches):  
    {
```

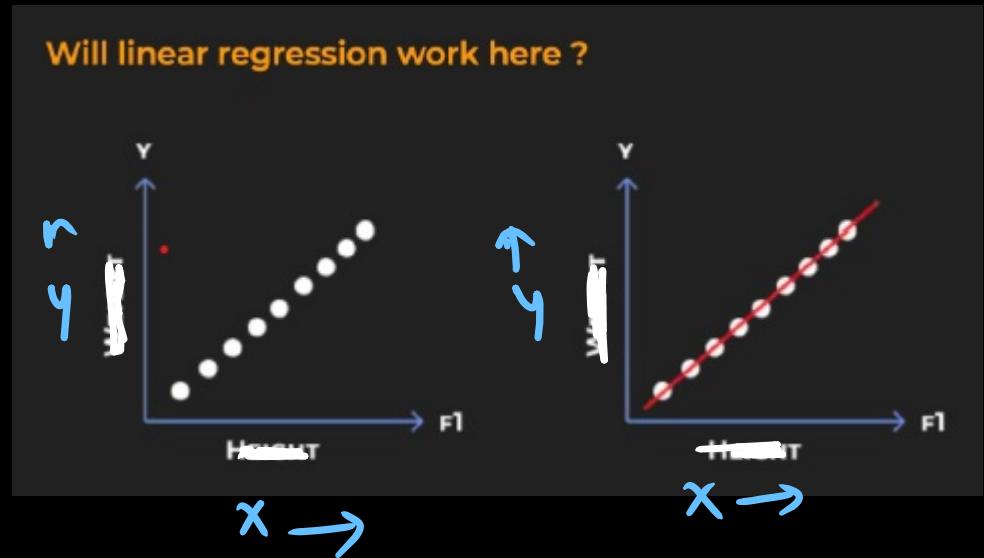
$$w_j = w_j - \eta \frac{1}{\text{Batch size}} \sum_{i=1}^{\text{Batch size}} (y^{(i)} - \hat{y}^{(i)}) x_j^{(i)}$$

1	32
33	32
65	32
97	32
.	
.	
106	

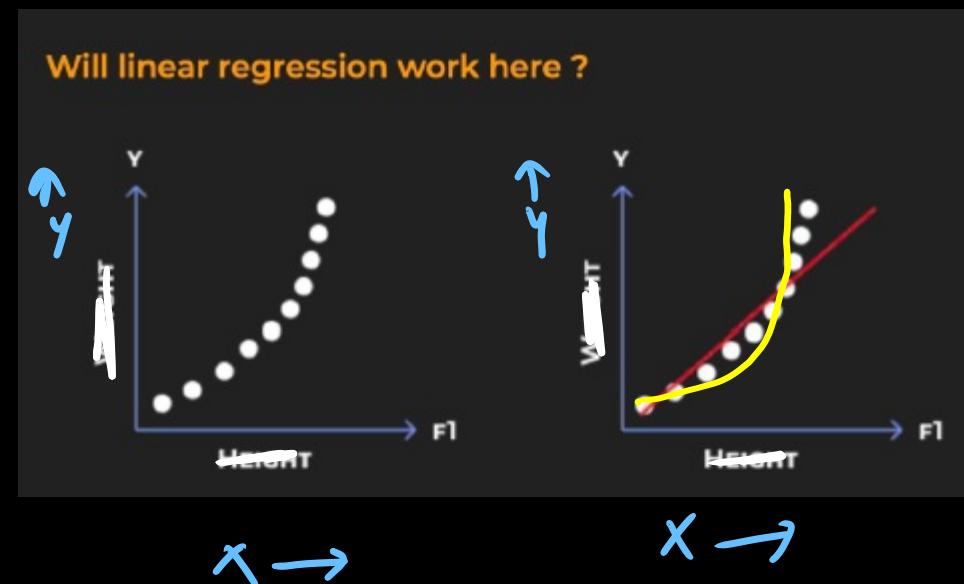
- ① BGD ✓
- ② SGD ✓
- ③ Mini batch GD →



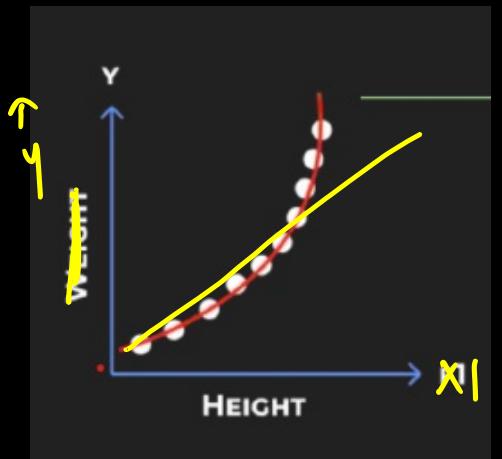
# Linear Reg



$$\hat{y} = \omega_0 + \sum \omega_j x_j$$



## Polynomial Reg



## Linear Reg

$$\hat{y} = w_0 + w_1 x_1 \rightarrow \text{linear}$$

$$\hat{y} = w_0 + \sum_{j=1}^d w_j x_j$$

Poly Reg

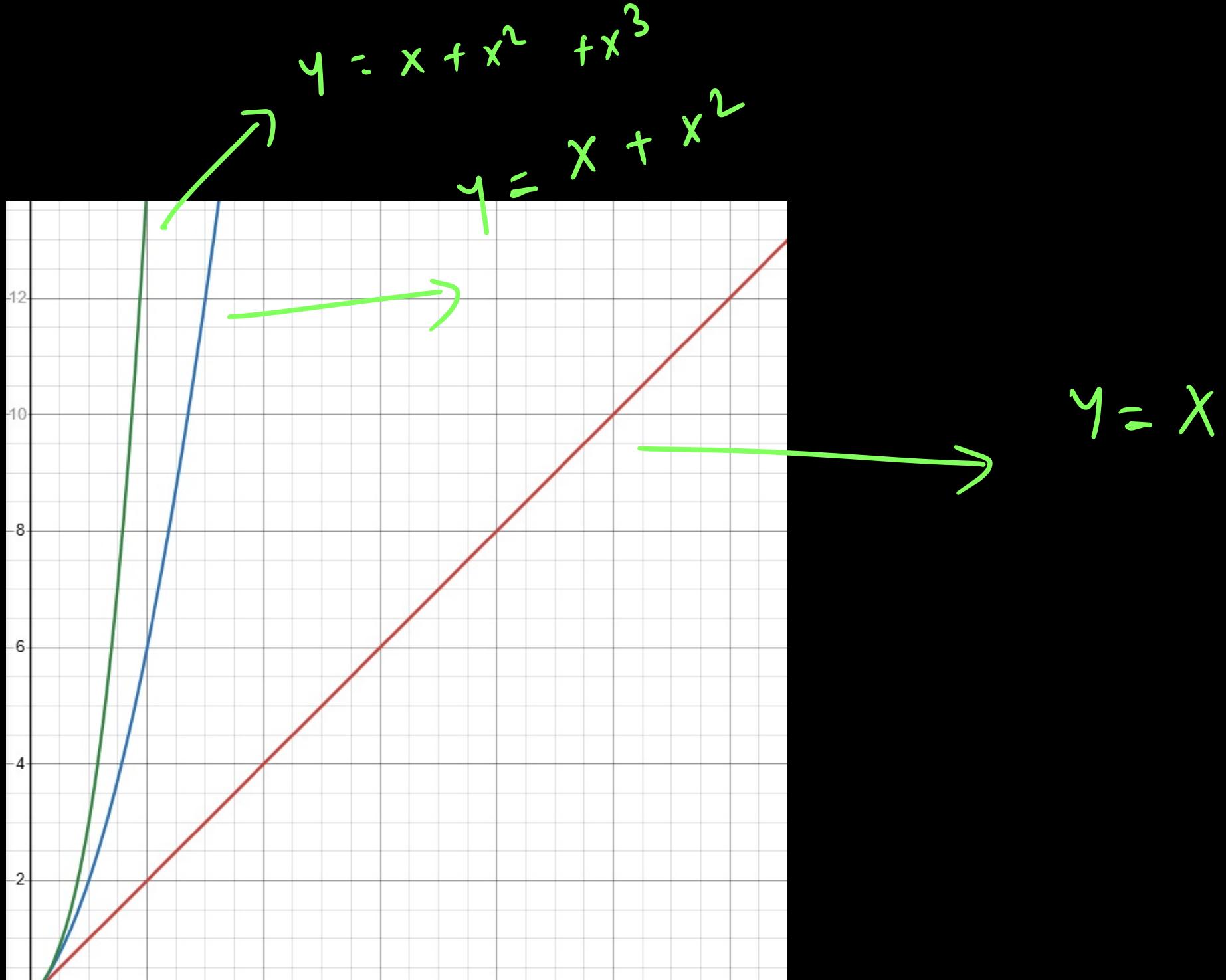
$$\checkmark \quad \hat{y} = w_0 + w_1 x_1 + w_{12} \underbrace{x_1^2}_{\text{quadratic}}$$

{  
2 = Quadratic  
3 = Cubic  
 $n > 3 \rightarrow \text{Poly}$

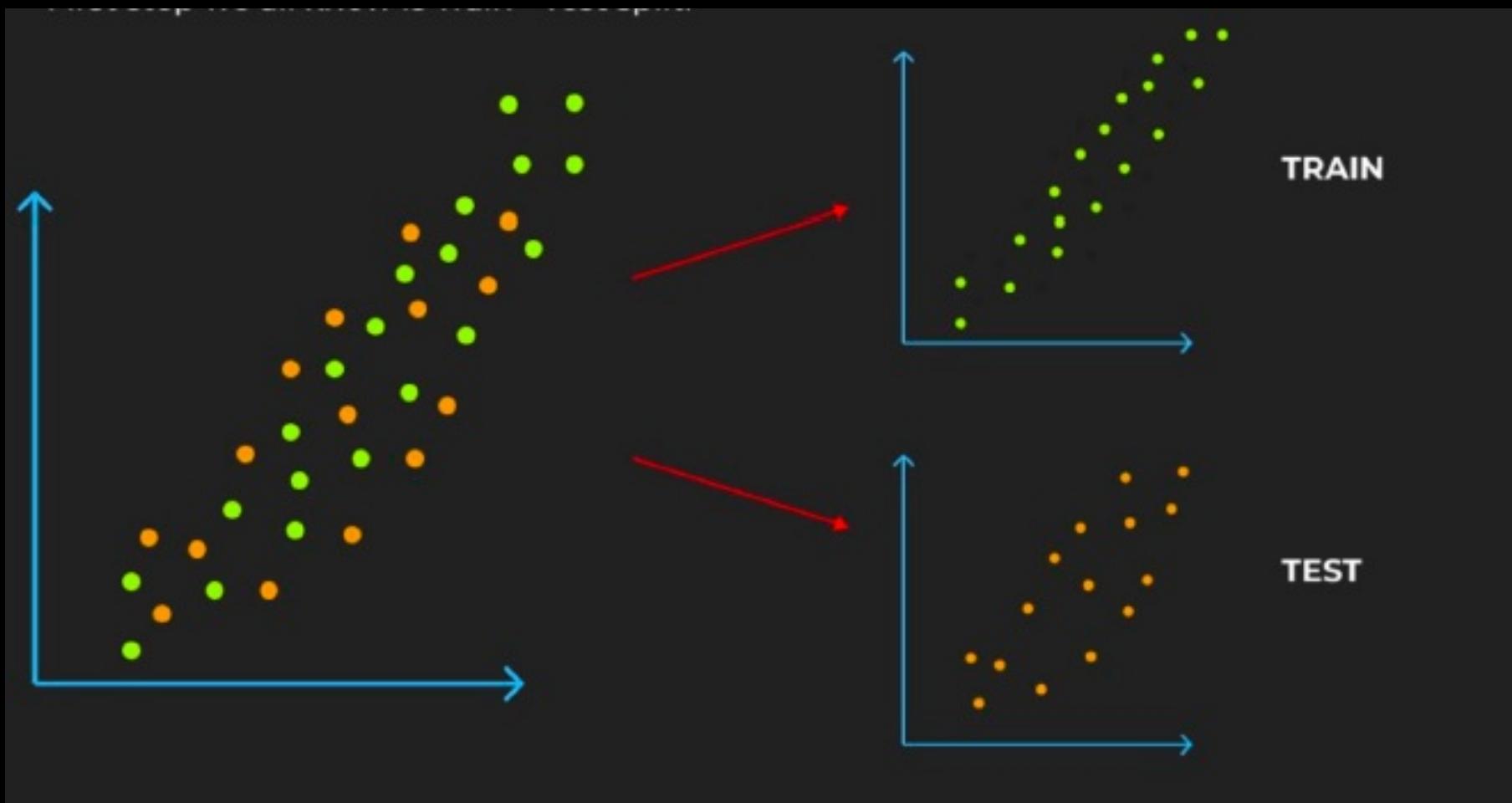
$$\hat{y} = w_0 + w_1 x_1 + w_2 x_2 - w_d x_d$$

$$+ \underbrace{w_{12} x_1^2 + w_{22} x_2^2 - w_d x_d^2}_{\text{quadratic terms}}$$

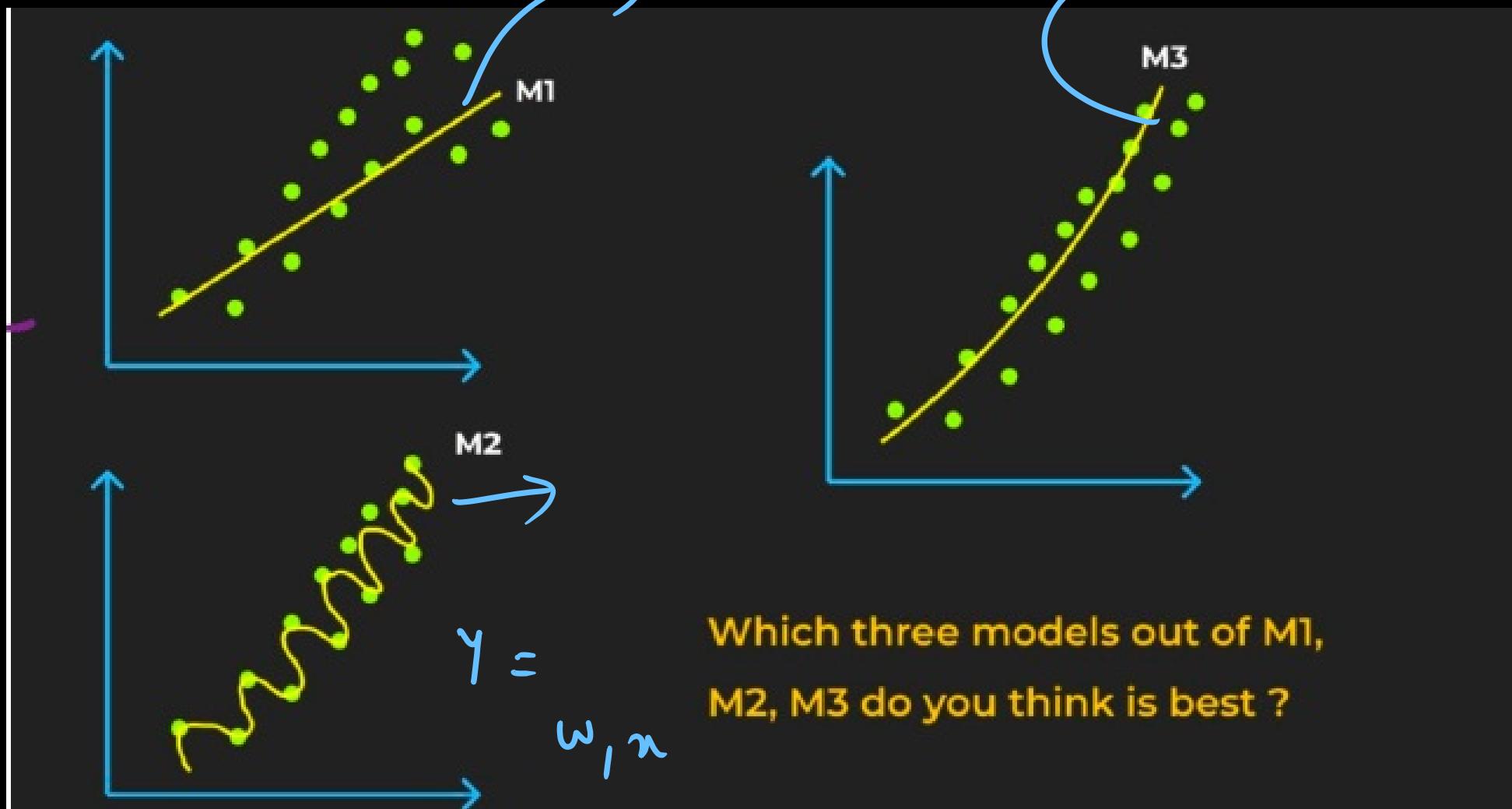
$$+ \underbrace{w_{13} x_1^3 + w_{23} x_2^3 - w_d x_d^3}_{\text{cubic terms}}$$



→ Break until 22:07 PM



$$\gamma = \omega_1 n_1 + \omega_0$$

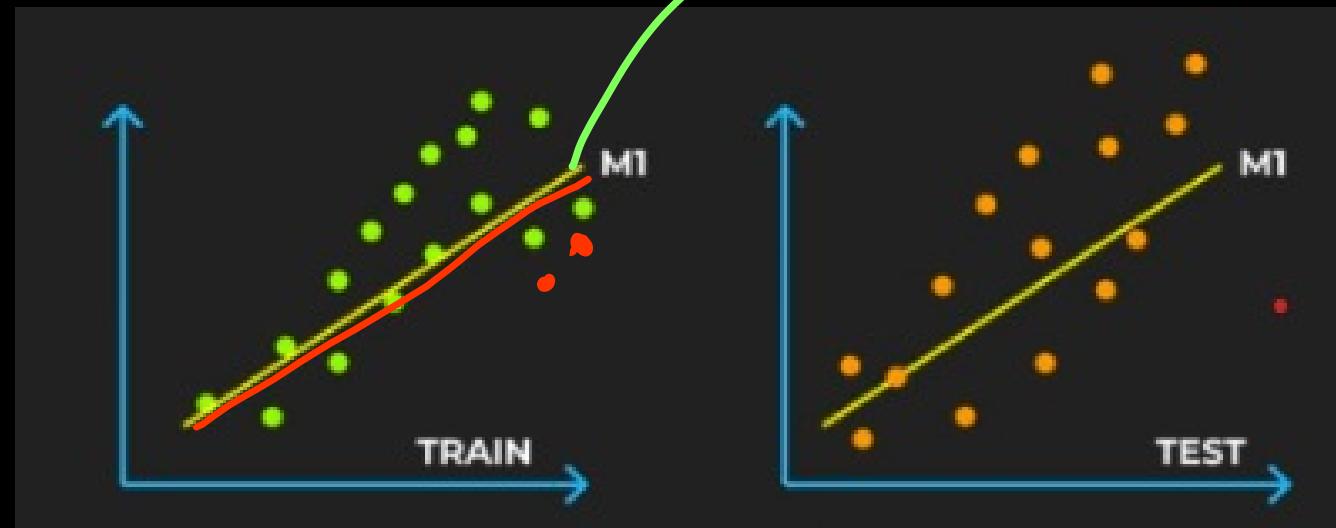


$$\gamma = \omega_1 n_1 + \omega_{11} n_1^2$$

$$+ \omega_{13} n_1^3 + \omega_{14} n_1^4$$

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

Underfitting



$$y = w_0 + w_1 x_1$$

$\rightarrow \text{Error} (y - \hat{y})$

$\hookrightarrow$  High bias & low  
Variance

> High

Sensitivity

w.r.t. change  
in data

2 > Low

$\rightarrow$  Student did not study text book (training data)

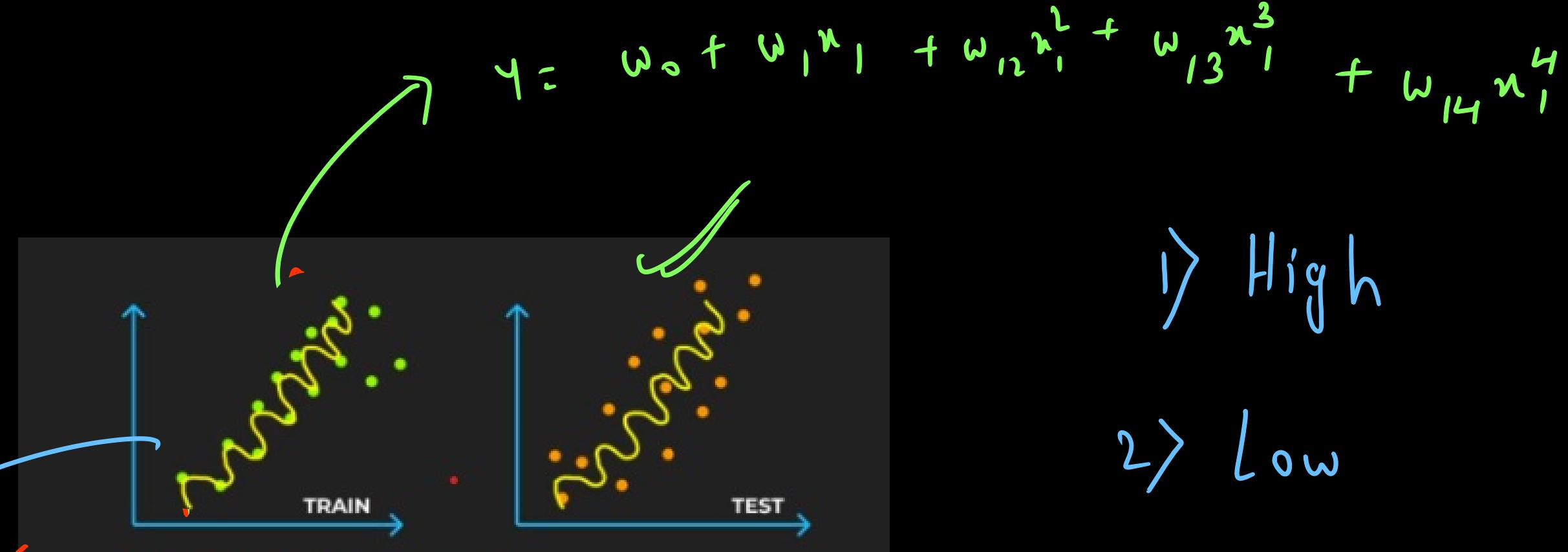
✓ High train error

✓ High test error

Overfitting

Error ( $y - \hat{y}$ )  
||  
low  
bias

Variance  
High

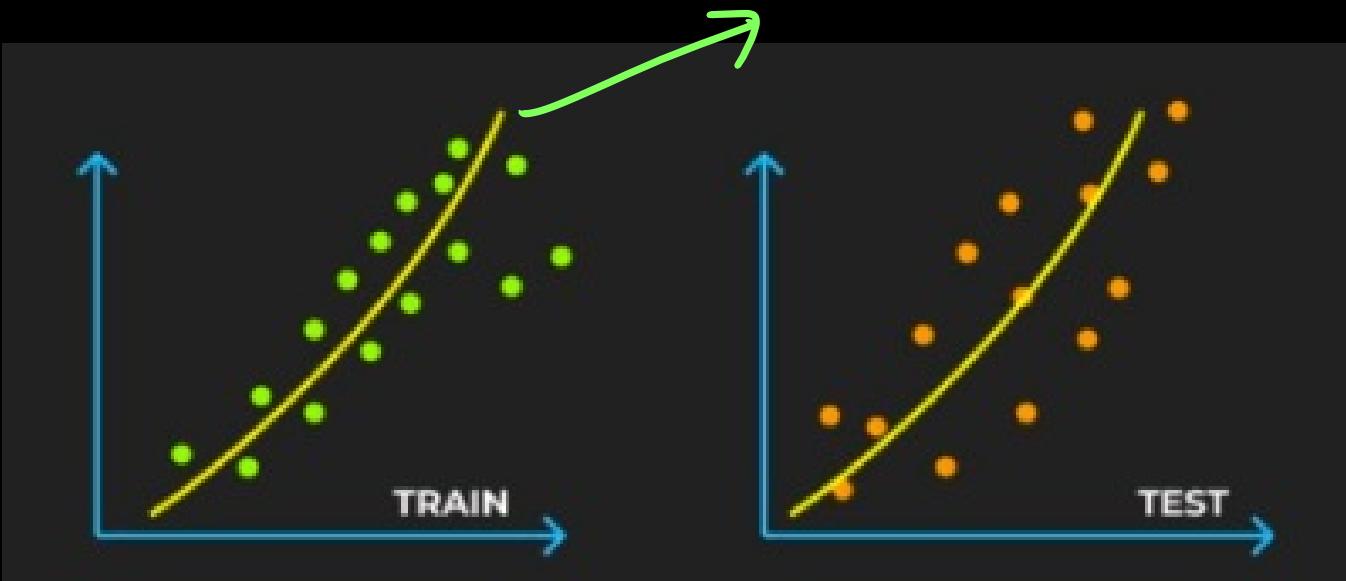


→ Student has memorized the text book (-train data)

Low train error      High test error

→ Perfect fit (Best fit)

$$y = \omega_0 + \omega_1 u_1 + \omega_{12} u_1^2$$



1) High

2) Low

→ Student has studied text book smartly  
& can apply concept

Low train error

Low test error

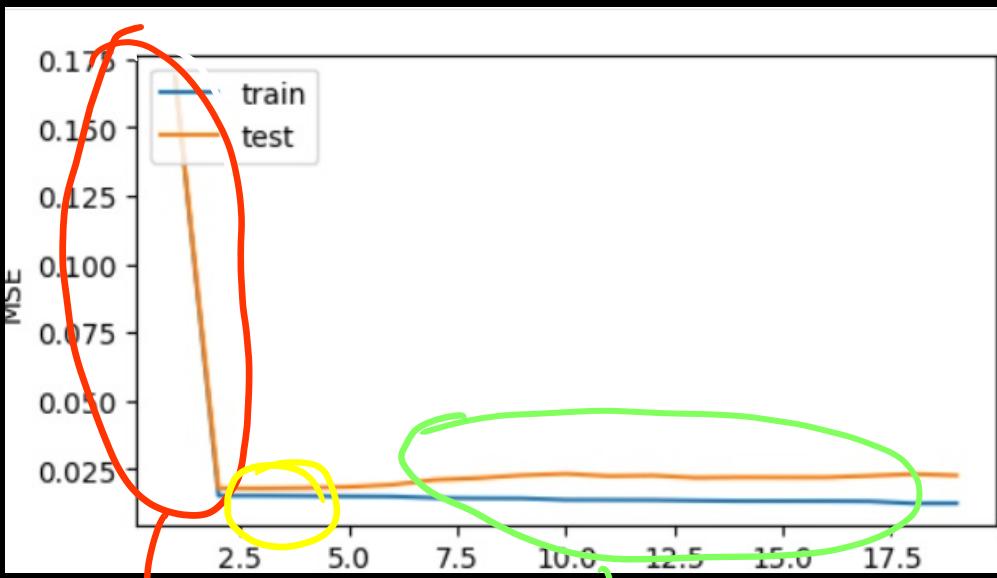
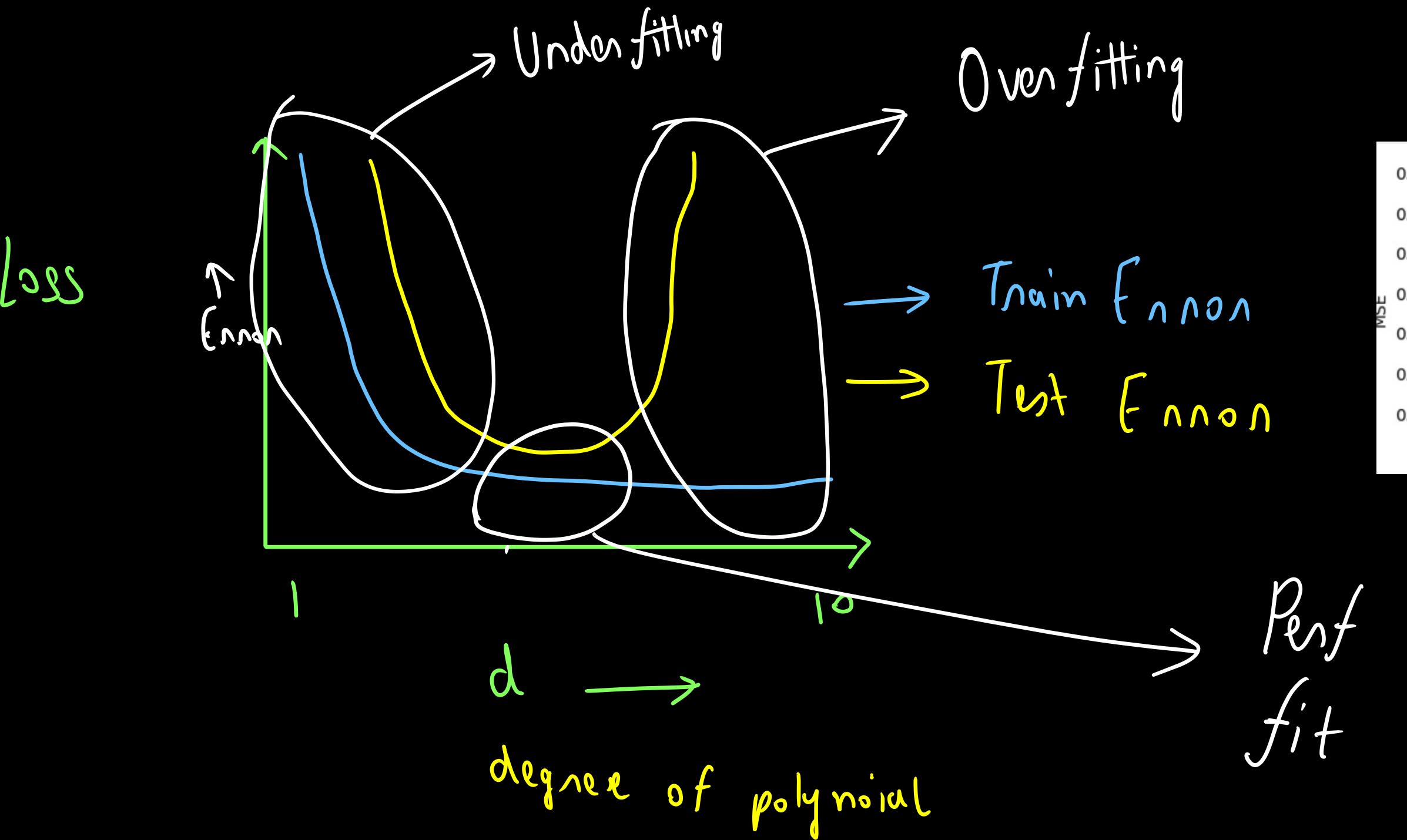
$$y = w_0 + w_1 x$$

	TRAINING	TESTING
SIMPLEST	UNDERFIT	POOR
	PERFECTLY FIT	GREAT
COMPLEX	OVERFIT	BEST

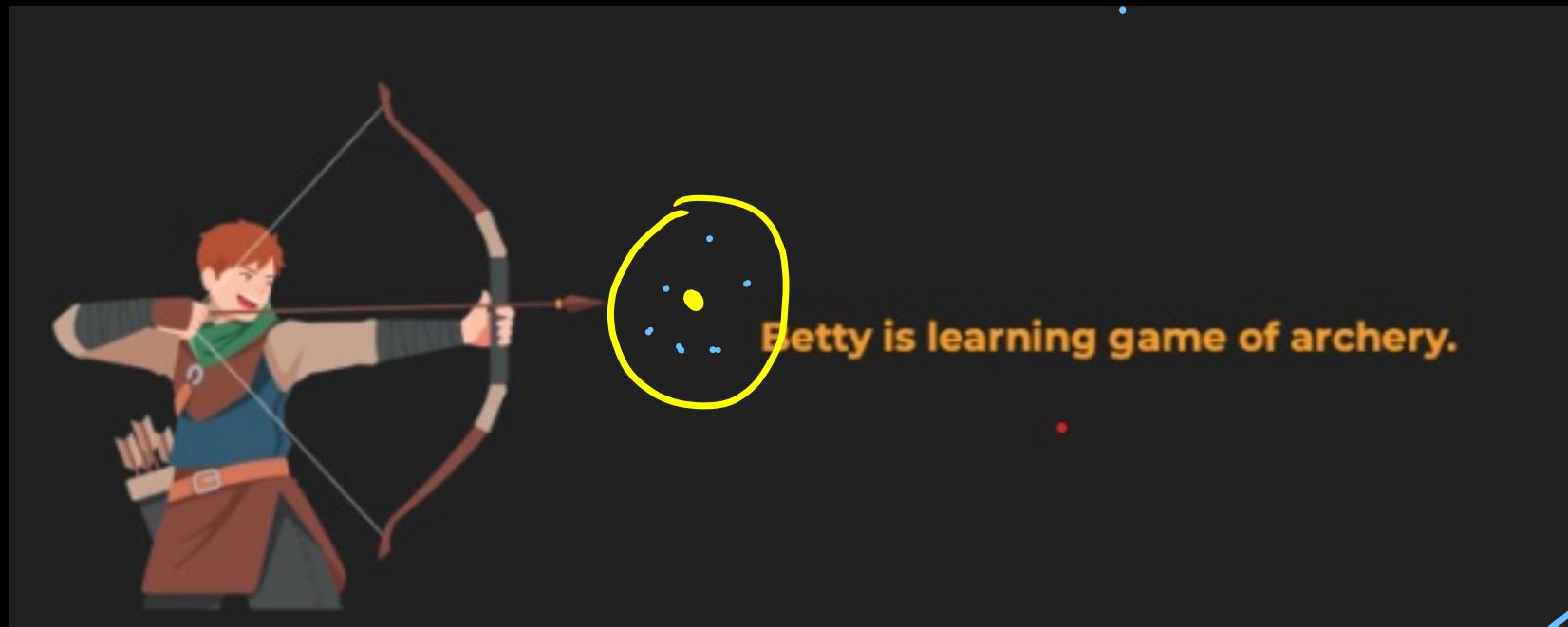
high train  
& test

low train  
& high  
test

$$y = w_0 + w_1 x \quad \dots \quad w_{15} x^5$$



## Bias & Variance

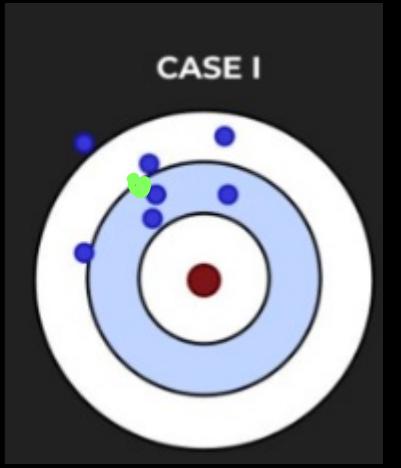


$f$  = Actual target

$\hat{f}$  = Where you hit

Bias =  $\text{Mean}(\hat{f}) - f$

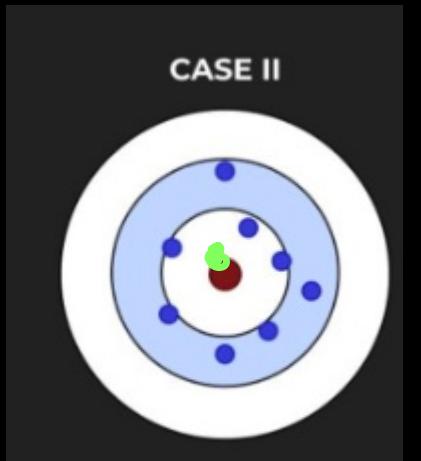
$$\begin{aligned}\text{Variance} &= \frac{\text{Spread of } \hat{f}}{m} \\ &= \frac{1}{m} \sum_{i=1}^m (\hat{f}_i - \bar{\hat{f}})^2\end{aligned}$$



→ Bias - H

↑ High  
↓ Low

→ Var - H



→ Bias - L

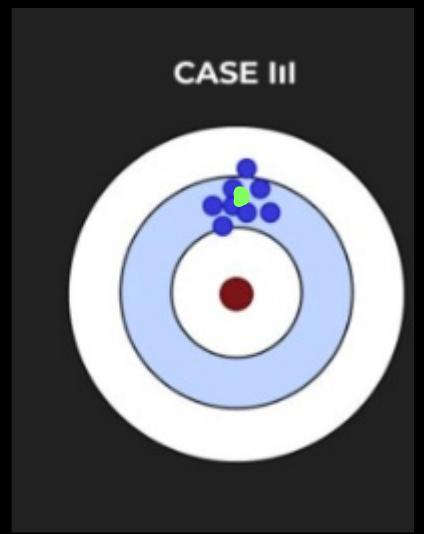
→ Var - H

$f$  = Actual target

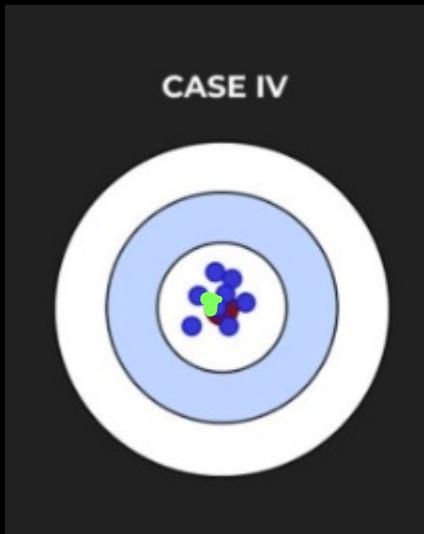
$\hat{f}$  = Where you hit

$$\text{Bias} = \text{Mean}(\hat{f}) - f$$

$$\begin{aligned}\text{Variance} &= \frac{\text{Spread of } \hat{f}}{m} \\ &= \frac{1}{m} \sum_{i=1}^m (\hat{f}_i - \bar{\hat{f}})^2\end{aligned}$$



→ Bias - H      ↗ high  
 → Var - L      ↘ low



→ Bias - L  
 → Var - L

$\checkmark$   
 $f$  = Actual target

$\hat{f}$  = Where you hit

$$\text{Bias} \quad \checkmark = \text{Mean}(\hat{f}) - f$$

$$\begin{aligned} \text{Variance} &= \frac{\text{Spread of } \hat{f}}{m} \\ &= \frac{1}{m} \sum_{i=1}^m (\hat{f}_i - \bar{\hat{f}})^2 \end{aligned}$$



