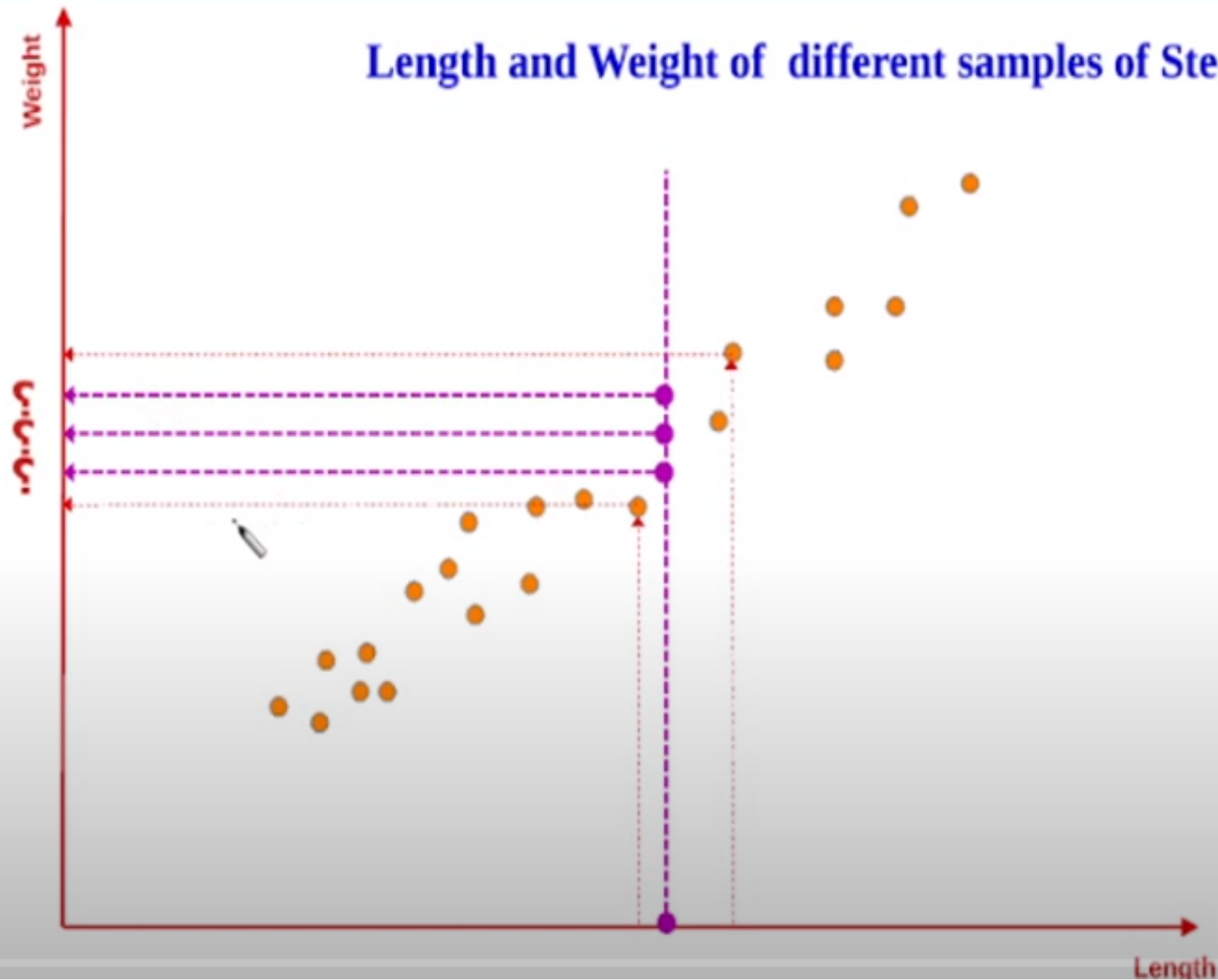


## Length and Weight of different samples of Steel Bars

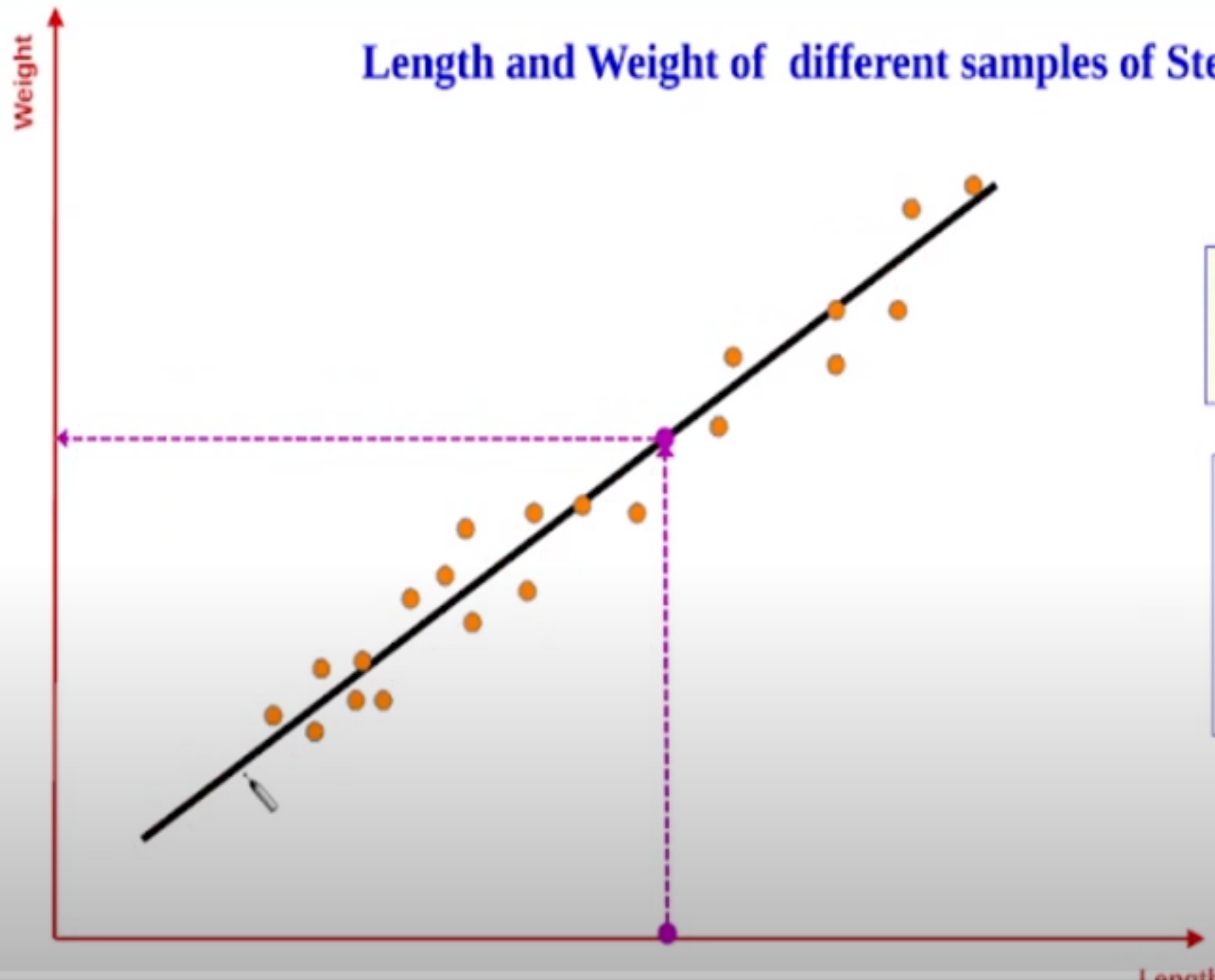


**Model**

**I/P: Length**

**O/P: Weight**

## Length and Weight of different samples of Steel Bars



Model

I/P: Length

O/P: Weight

Model (Line)

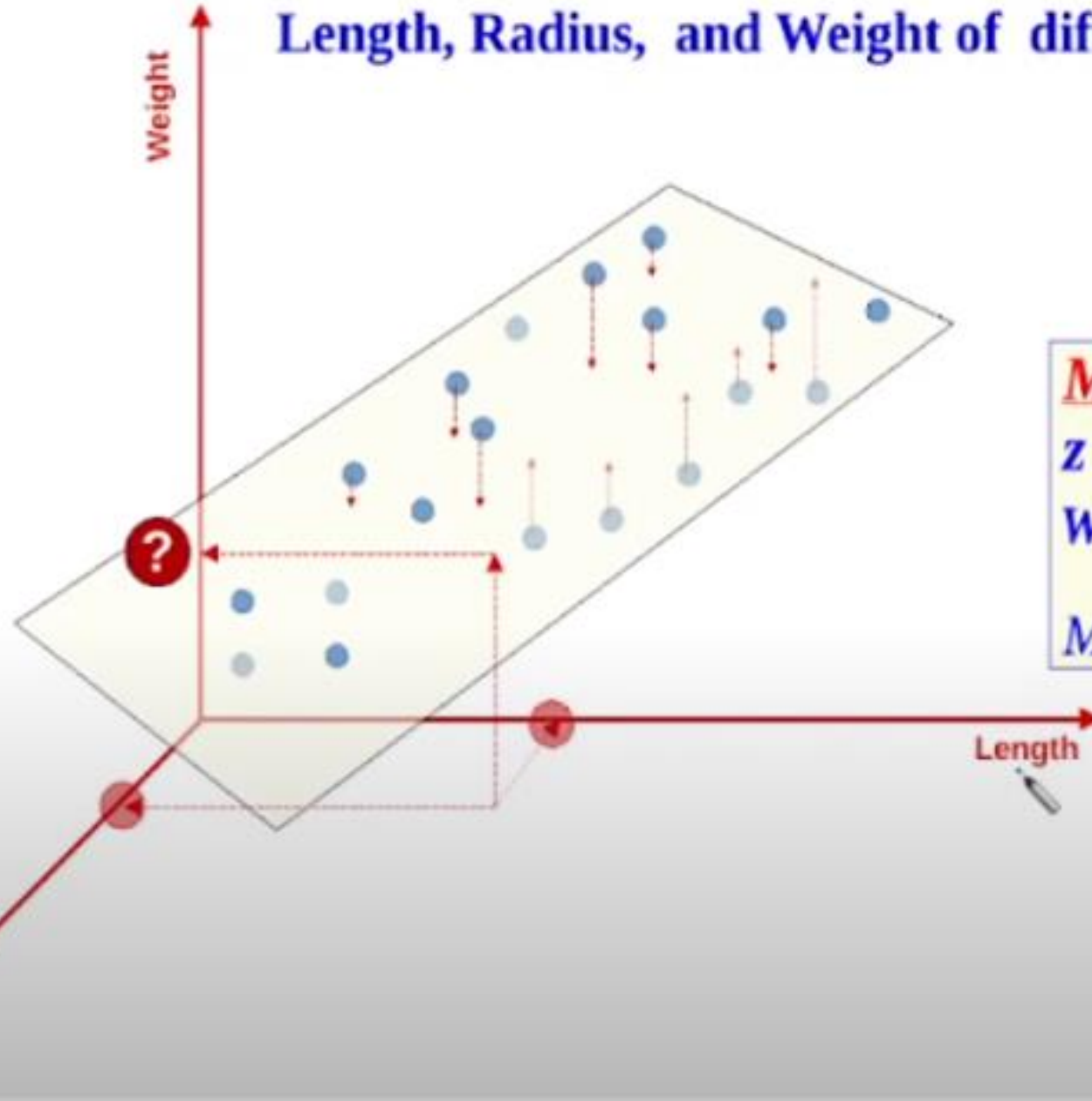
$$y = ax + b$$

$$\text{Weight} = a * \text{Length} + b$$

Model Parameters :  $a, b$

## Regression – Two Dimensional Case

Length, Radius, and Weight of different samples of Steel Bars



Model

I/P: Length, Radius

O/P: Weight

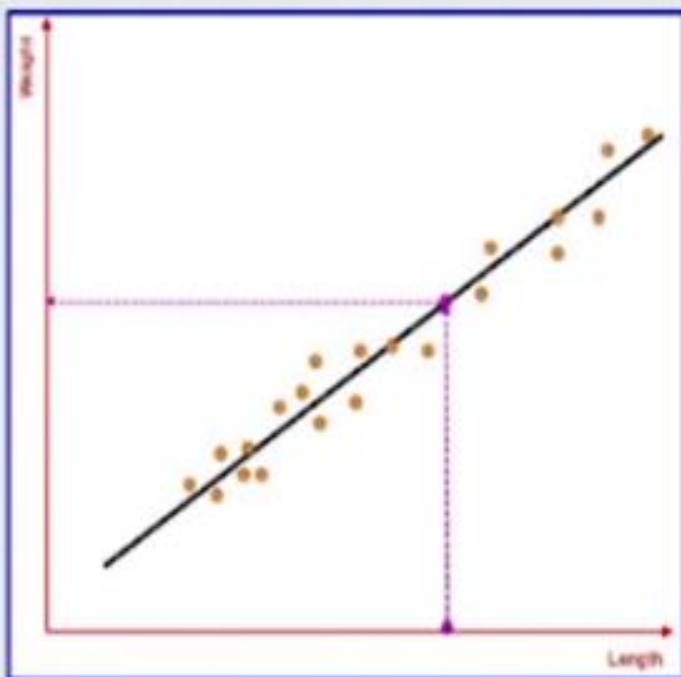
Model (Plane)

$$z = ax + by + c$$

$$\text{Weight} = a * \text{Length} + b * \text{Radius} + c$$

Model Parameters:  $a, b, c$

### Input Vector: 1-D



#### Model

I/P: Length [x]

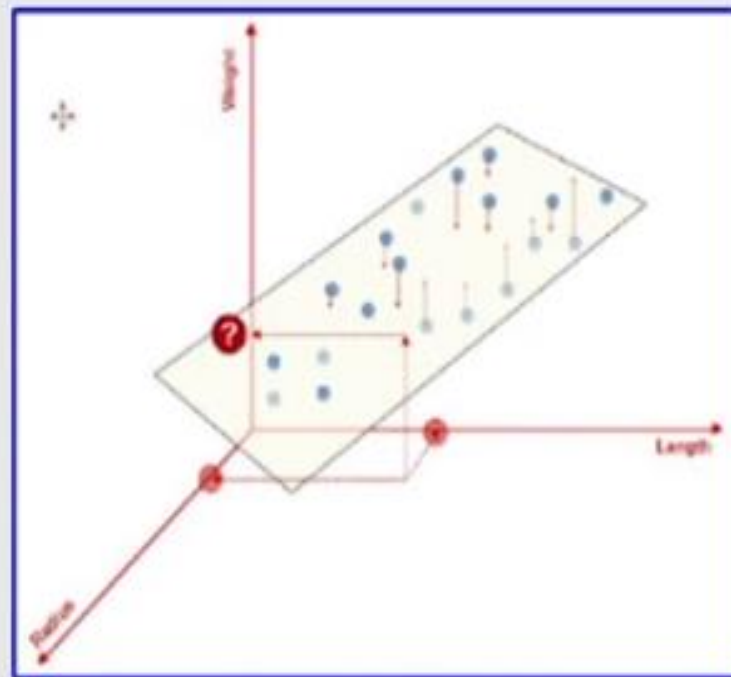
O/P: Weight [y]

#### Model ( Line )

$$y = ax + b$$

$$\text{Weight} = a * \text{Length} + b$$

### Input Vector: 2-D



#### Model

I/P: Length, Radius [x,y]

O/P: Weight [z]

#### Model ( Plane )

$$z = ax + by + c$$

$$\text{Weight} = a * \text{Length} + b * \text{Radius} + c$$

### Input Vector: d-D

#### General Case

Input:  $\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ \dots \\ x_d \end{bmatrix}$  Output: (y)

#### Model

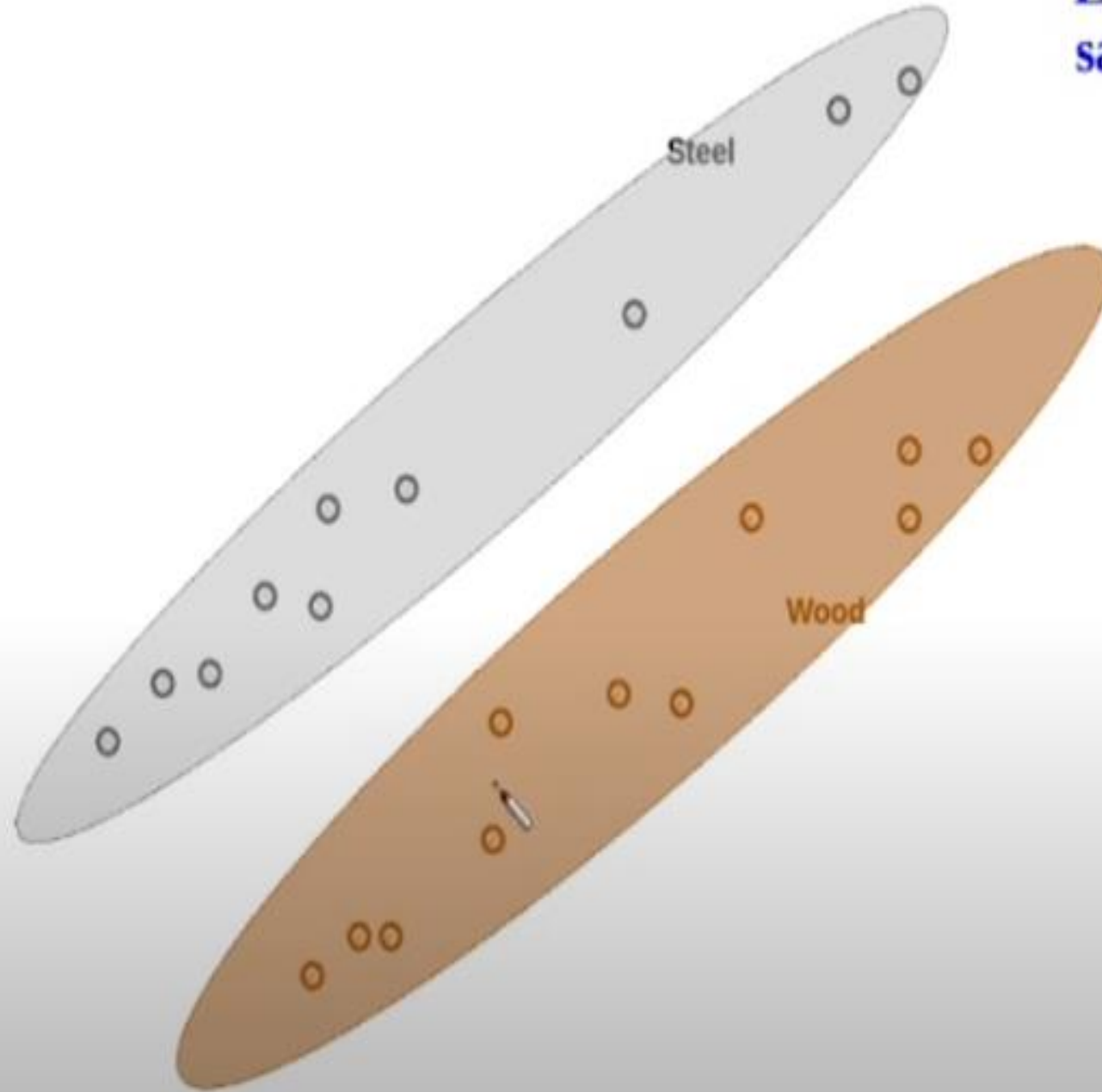
I/P:  $[x_1, x_2, x_3, \dots, x_d]$

O/P: [Y]

#### Model ( Hyperplane )

$$y = w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_d x_d + w_{d+1}$$

Weight ↑



## Length and Weight of different samples of Steel and Wood Bars

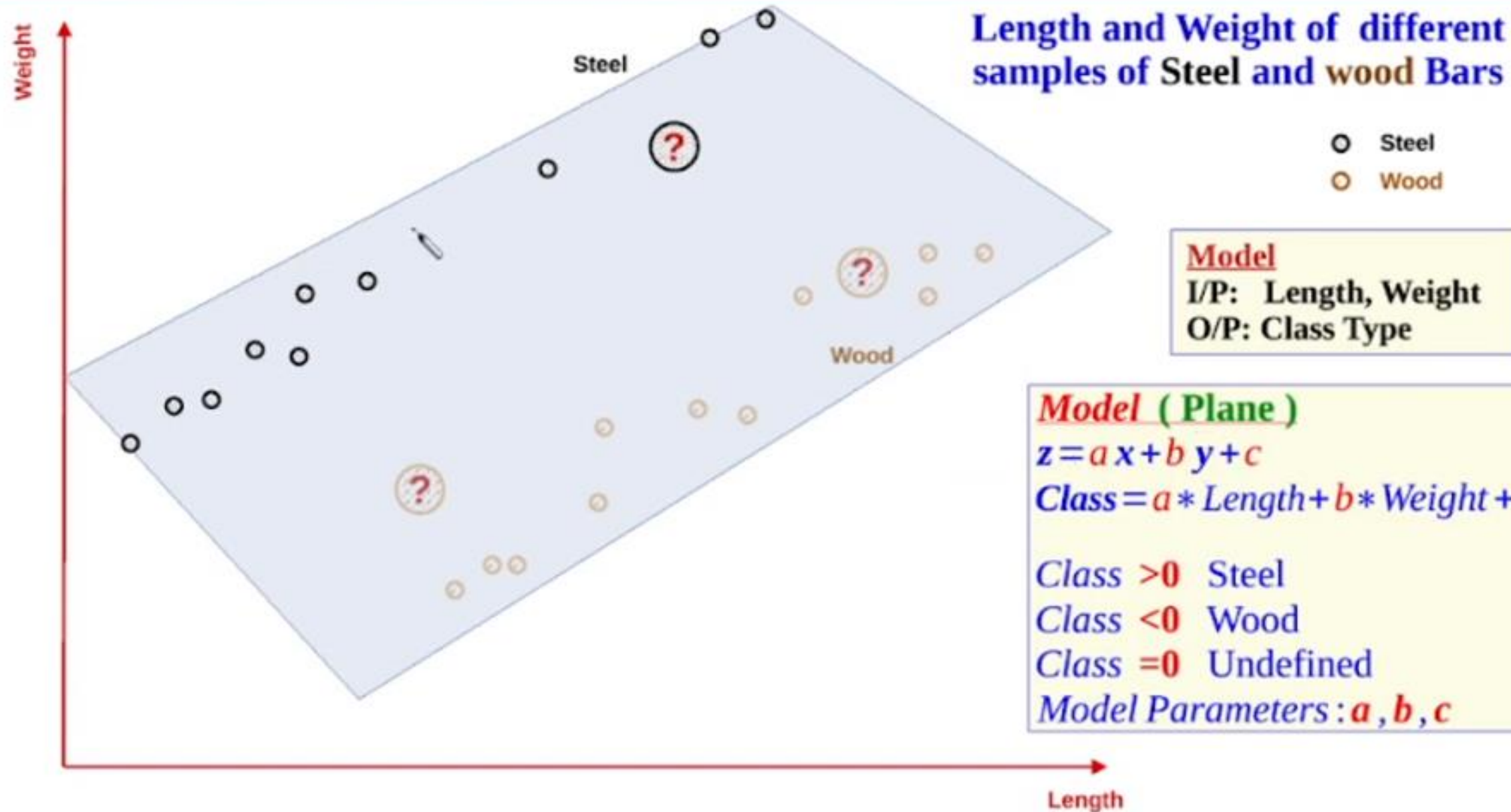
- Steel
- Wood

### Model

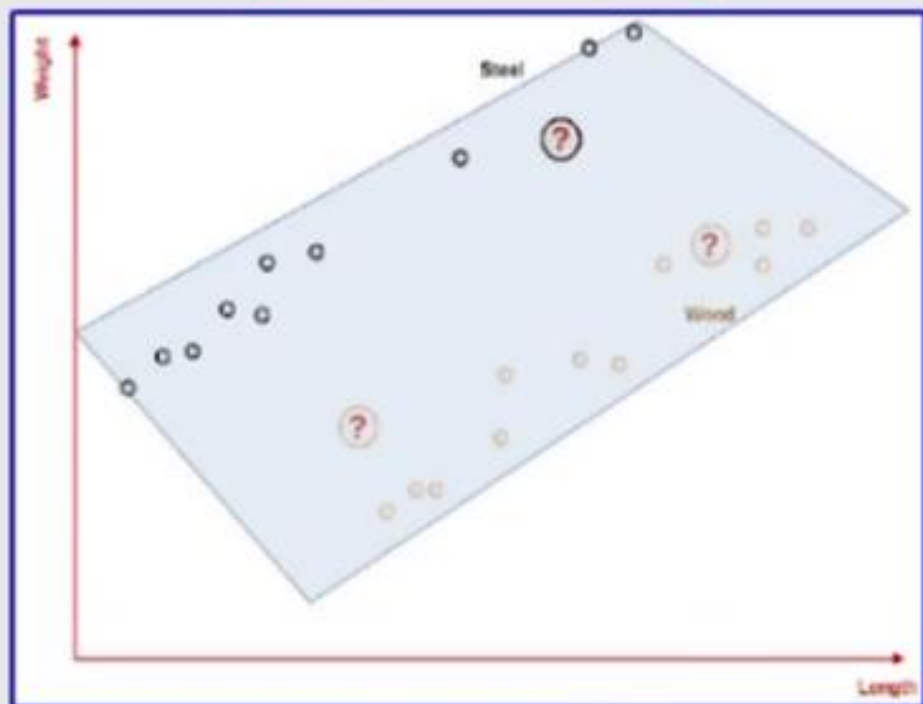
I/P: Length, Weight  
O/P: Class Type



## Classification – Two Dimensional Case



## Input Vector: 2-D



### Model

I/P: Length, Weight

O/P: Class Type

### Model ( Plane )

$$z = ax + by + c$$

$$\text{Class} = a * \text{Length} + b * \text{Weight} + c$$

Class  $> 0$  Steel    Class  $< 0$  Wood

Class  $= 0$  Undefined

## Input Vector: d-D

13

### General Case

Input:

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ \dots \\ x_d \end{bmatrix}$$

Output: (y)

### Model

I/P:  $[x_1, x_2, x_3, \dots, x_n]$

O/P:  $[y]$

### Model ( Hyperplane )

$$y = w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_d + w_{d+1}$$

$y > 0 \Rightarrow$  Class 1     $y < 0 \Rightarrow$  Class 2     $y = 0 \Rightarrow$  Undefined

Model Parameters:  $w_1, w_2, w_3, \dots, w_{d+1}$

## Same Equation for Linear Regression and Classification

### Model

I/P:  $(x_1 \ x_2 \ x_3 \ \dots \ x_n)$

O/P:  $(y)$

### General Case - Regression

### Model

$$y = w_1 x_1 + w_2 x_2 + w_2 x_2 + \dots + w_n x_d + b$$

Model Parameters:  $w_1, w_2, w_3, \dots, w_d, b$

$b$  is called "bias" term as it is independent of  $x$

### General Case - Classification

### Model

$$y = w_1 x_1 + w_2 x_2 + w_2 x_2 + \dots + w_d x_d + b$$

$y > 0 \Rightarrow$  Class 1     $y < 0 \Rightarrow$  Class 2     $y = 0 \Rightarrow$  Undefined

Model Parameters:  $w_1, w_2, w_3, \dots, w_d, b$

$$\text{Input: } \underline{X} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ \dots \\ x_d \end{pmatrix} \quad \underline{W} = \begin{pmatrix} W_1 \\ W_2 \\ W_3 \\ \dots \\ \dots \\ W_d \end{pmatrix} \quad \text{Output: } (y)$$



## Vector Notation for Linear Regression and Classification (One Sample)

15

### Model

I/P:  $(x_1 \ x_2 \ x_3 \ \dots \ x_n)$       O/P:  $(y)$

### Model

$$y = w_1 x_1 + w_2 x_2 + w_2 x_2 + \dots + w_d x_d + b$$

$$y = \underline{W}^T \underline{X} + b$$

$$y = \begin{pmatrix} w_1 & w_2 & \dots & w_d \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_d \end{pmatrix} + b$$

$$\text{Input: } \underline{X} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ \dots \\ x_d \end{pmatrix} \quad \underline{W} = \begin{pmatrix} w_1 \\ w_2 \\ w_3 \\ \dots \\ \dots \\ w_d \end{pmatrix} \quad \text{Output: } (y)$$

# Vector Notation for Linear Regression and Classification (One Sample)

16

## Model

I/P:  $(x_1 \ x_2 \ x_3 \ \dots \ x_n)$     O/P:  $(y)$

## Model

$$y = w_1 x_1 + w_2 x_2 + w_2 x_2 + \dots + w_n x_n + b$$

$$y = \underline{W}^T \underline{X} + b$$

$$y = \begin{pmatrix} w_1 & w_2 & \dots & w_n \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix} + b$$

## Note:

“b” may be added at begin of the vector  
or at end of the vector  
Same for “1” added to input vector

$$\text{Input: } \underline{X} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ \dots \\ x_d \end{pmatrix} \quad \underline{W} = \begin{pmatrix} w_1 \\ w_2 \\ w_3 \\ \dots \\ \dots \\ w_d \end{pmatrix} \quad \text{Output: } (y)$$

$$y = \begin{pmatrix} b & w_1 & w_2 & \dots & w_d \end{pmatrix} \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ \dots \\ x_d \end{pmatrix}$$

$$y = \underline{W}^T \underline{X}$$

$$y = \underline{W}^T \underline{X}$$

Where  $\underline{X}$  is Augmented vector  
 $\underline{W}$  is  $d+1$  dimension

Model Parameter is  $\underline{W}$

# Matrix Notation for Linear Regression and Classification (N- Samples)

17

N Training Samples

$$\begin{array}{cccc}
 \mathbf{X}_1 & \mathbf{X}_2 & \dots & \mathbf{X}_n \\
 \left( \begin{array}{c} x_{11} \\ x_{12} \\ \dots \\ x_{1d} \end{array} \right) & \left( \begin{array}{c} x_{21} \\ x_{22} \\ \dots \\ x_{2d} \end{array} \right) & \left( \begin{array}{c} \dots \\ \dots \\ \dots \\ \dots \end{array} \right) & \left( \begin{array}{c} x_{n1} \\ x_{n2} \\ \dots \\ x_{nd} \end{array} \right) \\
 \mathbf{y}_1 & \mathbf{y}_2 & \dots & \mathbf{y}_n
 \end{array}$$

Remember

$$\mathbf{y} = \mathbf{W}^T \mathbf{X} = \mathbf{X}^T \mathbf{W}$$

In General:  $N \gg d$

$$\begin{array}{c}
 \left( \begin{array}{c} y_1 \\ y_2 \\ \dots \\ y_n \end{array} \right)_{N \times 1} \\
 +
 \end{array}
 =
 \begin{array}{c}
 \left( \begin{array}{cccccc}
 1 & x_{11} & x_{12} & \dots & x_{1d} \\
 1 & x_{21} & x_{22} & \dots & x_{2d} \\
 \dots & \dots & \dots & \dots & \dots \\
 1 & x_{n1} & x_{n2} & \dots & x_{nd}
 \end{array} \right)_{N \times d+1}
 \left( \begin{array}{c} 1 \\ w_2 \\ w_2 \\ \dots \\ w_d \end{array} \right)_{d+1 \times 1}
 \end{array}$$

$$\mathbf{Y} = \mathbf{X} \mathbf{W}$$

Where:  $\mathbf{X}_{N \times d+1}$  is a **matrix**  
 $\mathbf{Y}_{N \times 1}$  is labels **vector**

## Calculation of Model Parameter “W”

$$\underline{Y} = \underline{X} \underline{W}$$

Where:  $\underline{X}_{N \times d+1}$  is a **matrix**  $\Rightarrow$  **Given**

$\underline{Y}_{N \times 1}$  is labels **vector**  $\Rightarrow$  **Given**

$\underline{W}$  is Model Parameter  $\Rightarrow$  **Required**



$$\underline{W} = \underline{X}^{-1} \underline{Y} \quad \text{😊}$$

$$\underline{W} = \left( \underline{X}_{N \times d} \right)^{-1} \underline{Y}$$

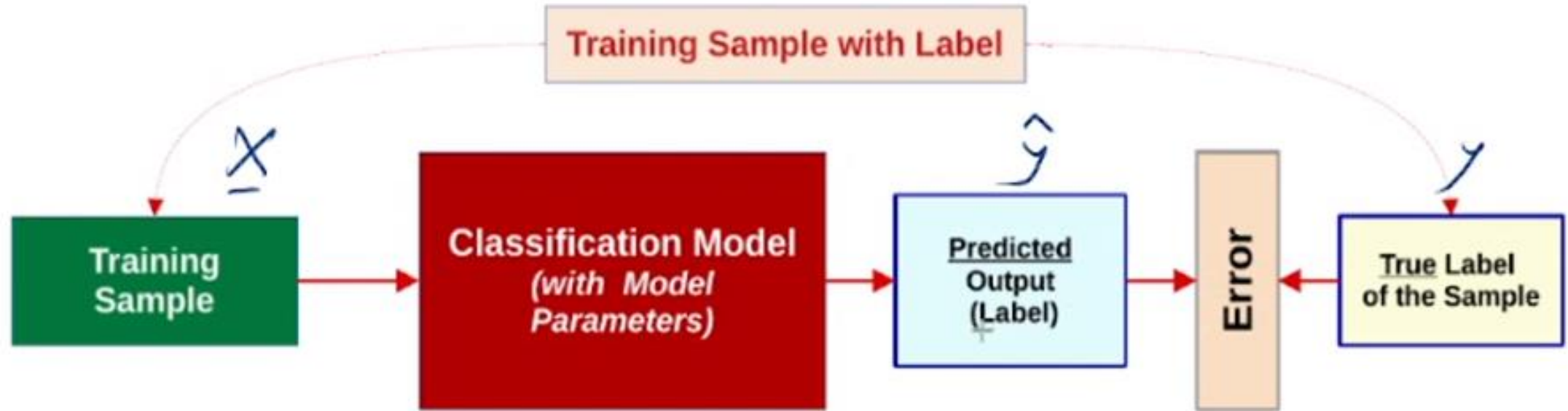
$\underline{X}_{N \times d}$  is non-square matrix

*There is NO Inverse*





## How to obtain the BEST model parameters (W)



Try to obtain model parameters by minimizing the **Error** (Loss Function)

### Note:

There are three terms that are commonly used interchangeably (*though there is a difference*)

**Loss** function , **Error** function, and **Objective** function (to be discussed later)

# Concept of LOSS Function



Let  $\underline{Y} = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix}_{N \times 1}$  **Original Labels**  
from Training Samples

Let  $\hat{\underline{Y}} = \begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \dots \\ \hat{y}_n \end{pmatrix}_{N \times 1}$  **Predicted Labels**  
from Model  $\hat{Y} = \underline{X} \underline{W}$

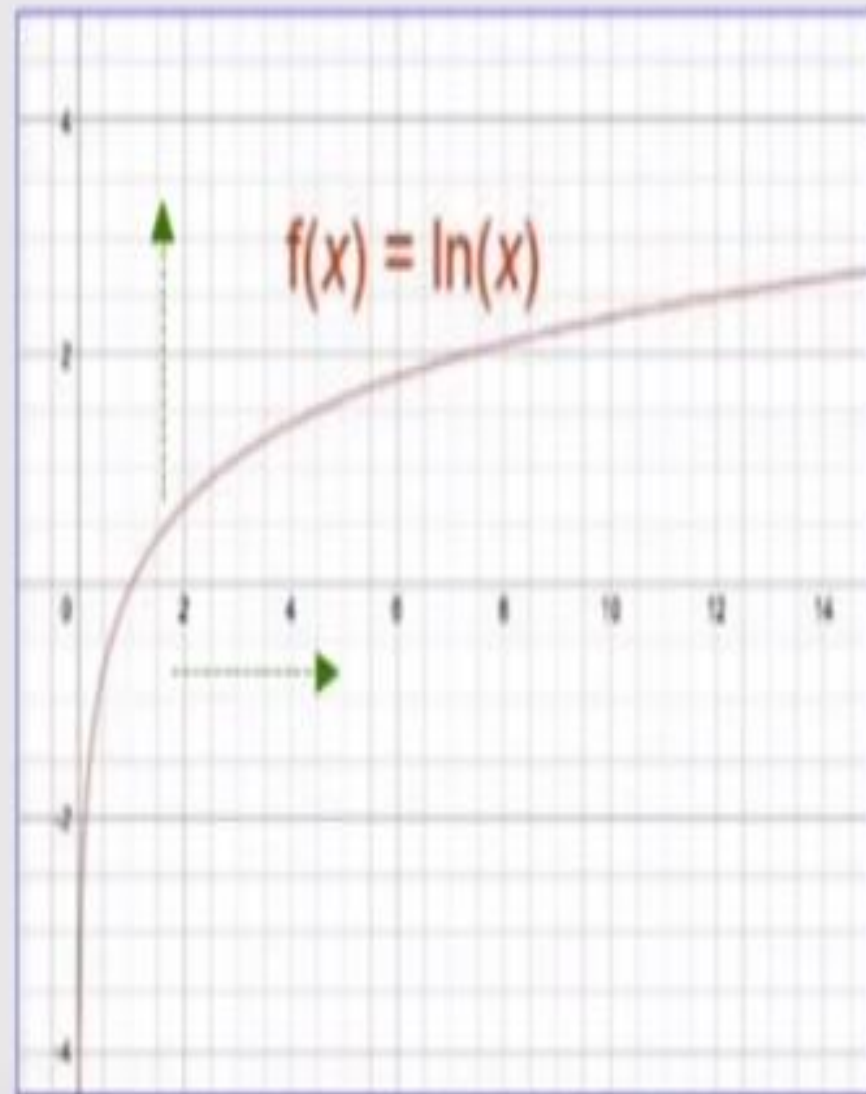
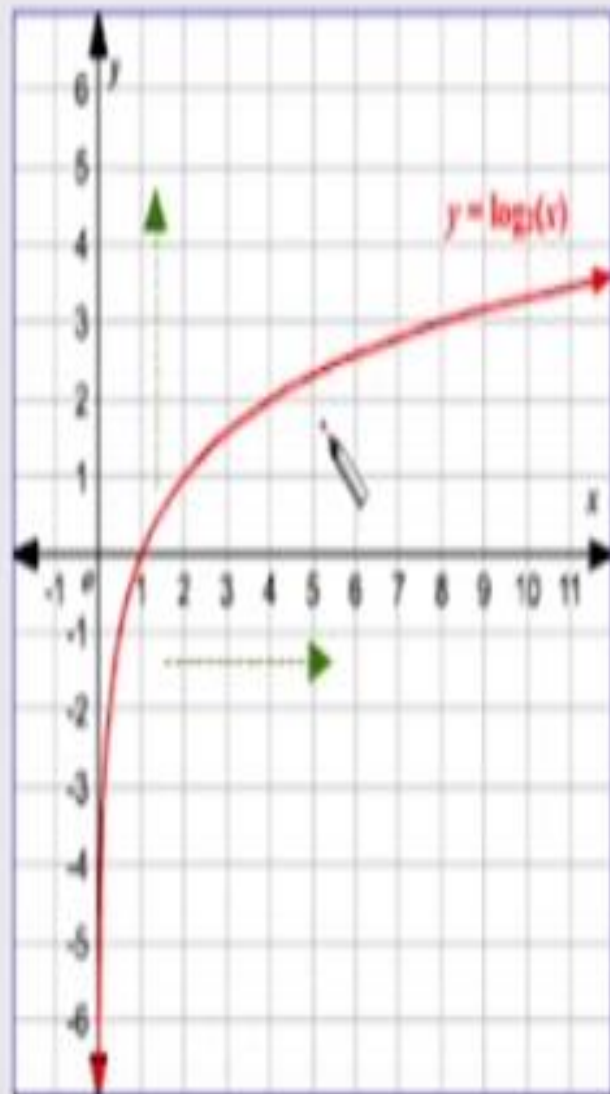
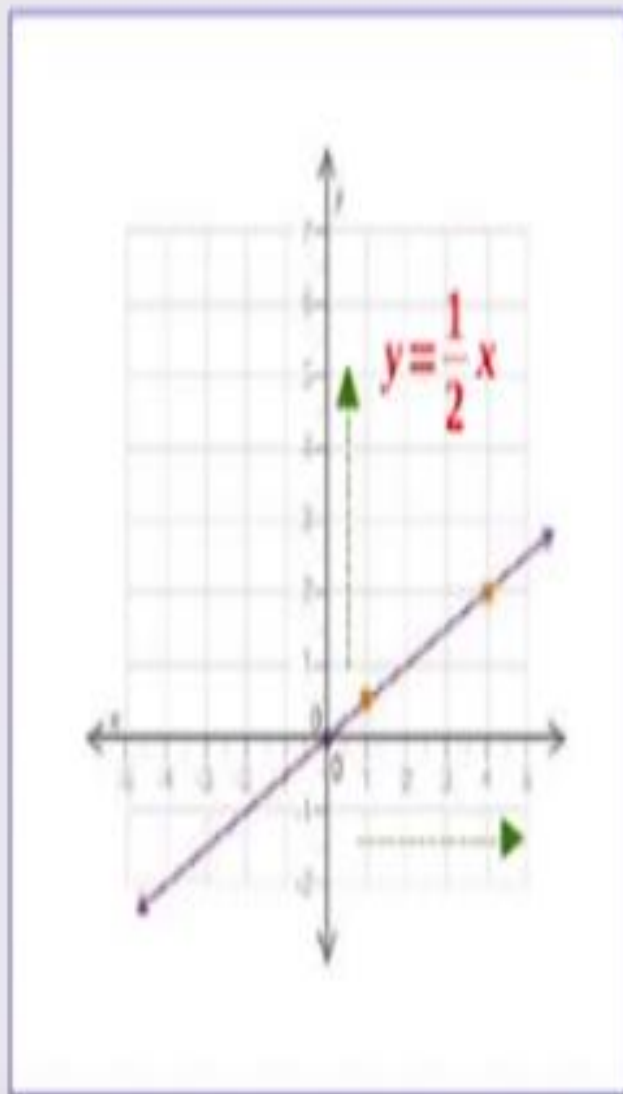
Objective is  $\hat{\underline{Y}} = \underline{Y}$

If  $\hat{\underline{Y}} \neq \underline{Y}$   $\Rightarrow$  There is an **Error (Loss)**

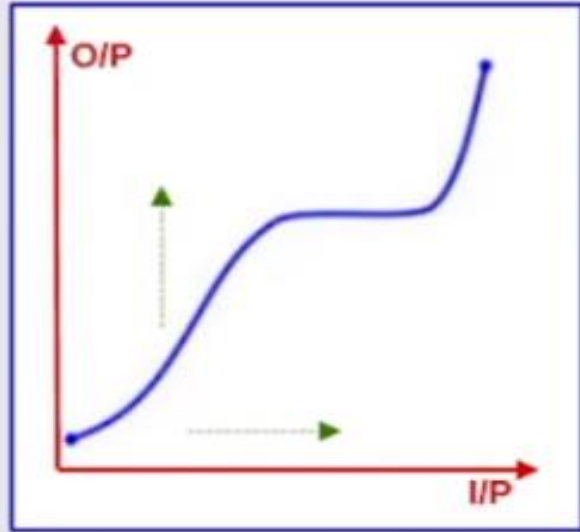
Define **Loss** function as squared distance between  $\hat{\underline{Y}}$  and  $\underline{Y}$

**Remember**  
 $\|a\|_2 = a^T a$

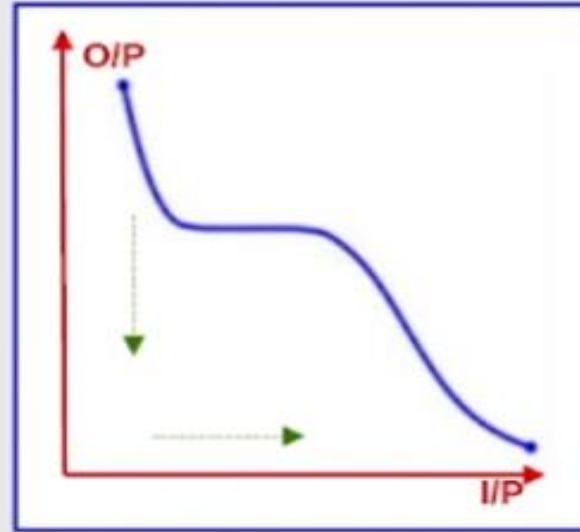
$$\text{Loss} = \|\underline{Y} - \hat{\underline{Y}}\|_2 = (\underline{Y} - \hat{\underline{Y}})^T (\underline{Y} - \hat{\underline{Y}})$$



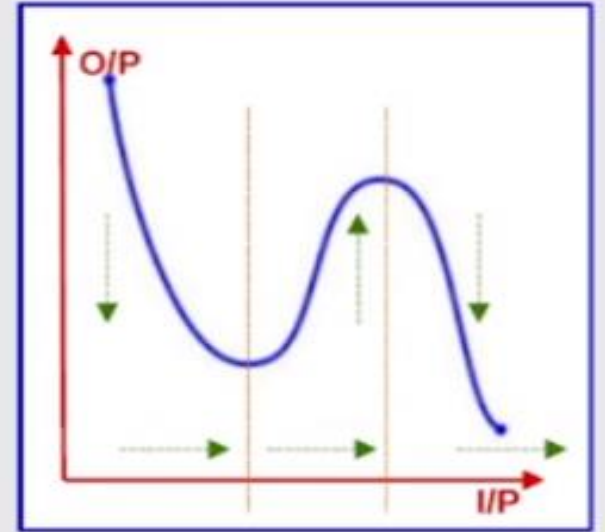
# Monotonic Functions



**Monotonically  
Increase**



**Monotonically  
Decrease**



**Not  
Monotonic**

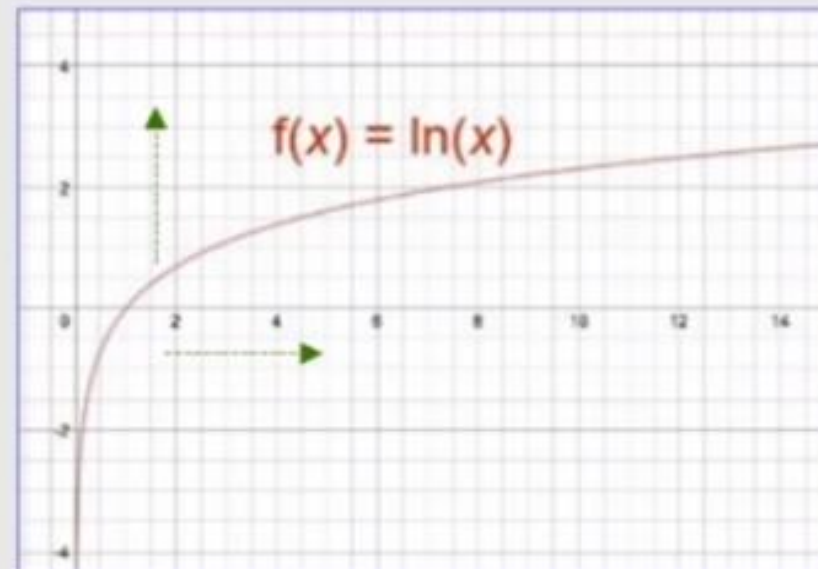
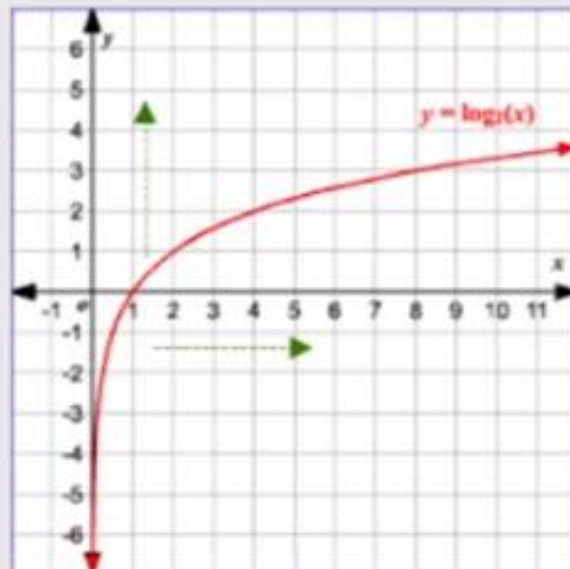
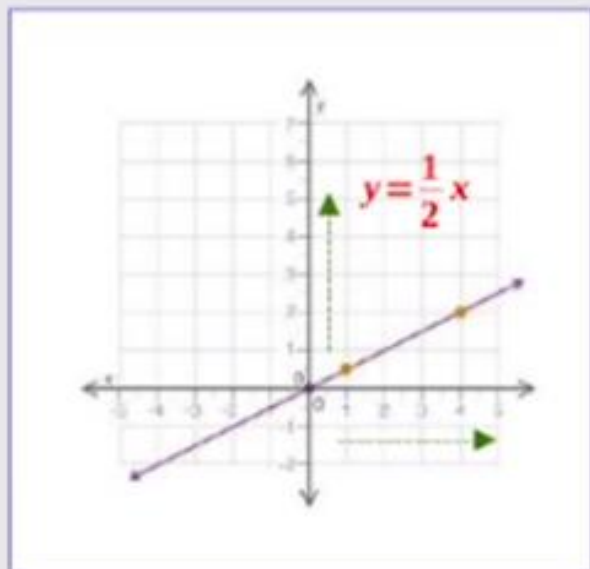
For Monotonically **increase** function: When Input **increases** → Output **increases**

For Monotonically **decrease** function: When Input **increases** → Output **decreases**



# Some Monotonically increase Functions

24



Apply monotonically increase function  $\Phi$  on Loss function  $L(x) \Rightarrow \Phi(L(x))$   
when  $L(x)$  is minimized  $\Rightarrow \Phi(L(x))$  is minimized

**Conclusion:** instead of using **Loss** function  $L(x)$  during training, one can use  $\Phi(L(x))$   
where  $\Phi$  is any monotonically increase function

We will use  $\frac{1}{2} L(x)$  instead of  $L(x)$  during minimization (to Simplify Math operations)

## Minimization of Loss function to obtain **BEST** model parameters



$$Loss = \frac{1}{2}(\underline{Y} - \hat{\underline{Y}})^T(\underline{Y} - \hat{\underline{Y}}) = \frac{1}{2}(\underline{Y} - \underline{X}\underline{W})^T(\underline{Y} - \underline{X}\underline{W})$$

$$Loss = \frac{1}{2}[\underline{Y}^T \underline{Y} + (\underline{X}\underline{W})^T(\underline{X}\underline{W}) - \underline{Y}^T(\underline{X}\underline{W}) - (\underline{X}\underline{W})^T \underline{Y}]$$

$$Loss = \frac{1}{2}[\underline{Y}^T \underline{Y} + (\underline{X}\underline{W})^T(\underline{X}\underline{W}) - 2(\underline{X}\underline{W})^T \underline{Y}]$$

### Remember

**X** is a matrix, **W** is a vector,  
**XW** is a vector

$$\underline{Y}^T(\underline{X}\underline{W}) = (\underline{X}\underline{W})^T \underline{Y}$$

$$\frac{\partial Loss}{\partial \underline{W}} = 0 \quad \text{Get Best } \underline{W} \text{ that minimizes the loss}$$

## Minimization of Loss function to obtain BEST model parameters



$$Loss = \frac{1}{2} (\underline{Y} - \hat{\underline{Y}})^T (\underline{Y} - \hat{\underline{Y}}) = \frac{1}{2} (\underline{Y} - \underline{X} \underline{W})^T (\underline{Y} - \underline{X} \underline{W})$$

$$Loss = \frac{1}{2} [\underline{Y}^T \underline{Y} + (\underline{X} \underline{W})^T (\underline{X} \underline{W}) - \underline{Y}^T (\underline{X} \underline{W}) - (\underline{X} \underline{W})^T \underline{Y}]$$

$$Loss = \frac{1}{2} [\underline{Y}^T \underline{Y} + (\underline{X} \underline{W})^T (\underline{X} \underline{W}) - 2 (\underline{X} \underline{W})^T \underline{Y}]$$

**Remember**

$\underline{X}$  is a matrix,  $\underline{W}$  is a vector,  
 $\underline{X} \underline{W}$  is a vector

$$\underline{Y}^T (\underline{X} \underline{W}) = (\underline{X} \underline{W})^T \underline{Y}$$

$$\frac{\partial Loss}{\partial \underline{W}} = 0 \quad \text{Get Best } \underline{W} \text{ that minimizes the loss}$$

$$\frac{\partial Loss}{\partial \underline{W}} = 0 = \frac{1}{2} [(\underline{X} \underline{W})^T (\underline{X}) + (\underline{X})^T (\underline{X} \underline{W}) - 2 (\underline{X})^T \underline{Y}]$$

**Remember**

$$\underline{X}^T (\underline{X} \underline{W}) = (\underline{X} \underline{W})^T \underline{X}$$

$$\frac{\partial Loss}{\partial \underline{W}} = 0 = \frac{1}{2} [2 \underline{X}^T \underline{X} \underline{W} - 2 \underline{X}^T \underline{Y}] \Rightarrow (\underline{X}^T \underline{X}) \underline{W} = \underline{X}^T \underline{Y} \Rightarrow \underline{W} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y}$$





# Minimum Sum Squared Error (MSSE) Technique (Least Squares)

$$\underline{W} = \left( X_{N \times d} \right)^{-1} \underline{Y}$$



$X_{N \times d}$  is non-square matrix  
*There is NO Inverse*

$$\underline{W} = \left( X^T X \right)^{-1} X^T \underline{Y}$$

Pseudo Inverse of Non-Square Matrix

$$\left( X_{N \times d} \right)^{-1} = \left( X^T X \right)^{-1} X^T$$

*Remember*

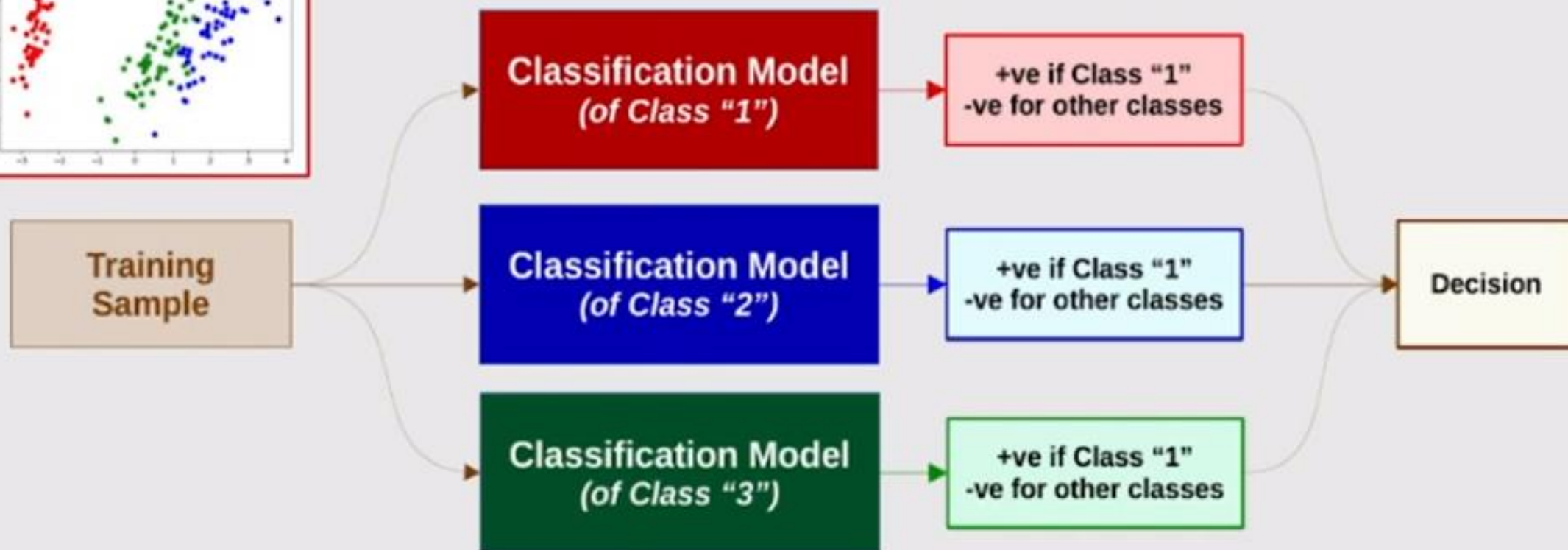
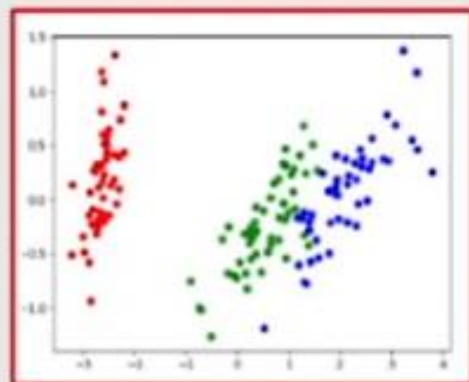
$$\left( X_{d \times N}^T X_{N \times d} \right)^{-1} = \left( \dots \right)_{d \times d}^{-1}$$

In General:  $d \lll N$





## One vs All (for Multi Class Classification)



### Decision Rules

**Number of Positive output (Only One)**

→ Input is properly Classified

**Number of Positive output (More than One)**

→ can NOT define the class

**NO Positive output**

→ Not one of the classes (new class)