

# Machine Learning

(30)

Marks	Placed
24	Placed
	not placed

Classification  
Regression

$x$

$f(x)$

Height

168 cm  
166.5

Weight

82 kg  
81 kg

Size  
of house

1200sqft

price

50L

$f(x)$  Price

Explain whether each of the following situations is a classification or regression problem:

- A company wants to launch a new product and wants to know whether it will turn out to be a success or failure. We have information on the last 100 products this company launched, including if it was a success/failure, price, weight, color, and several other variables. → classification
- We have information on several Bay Area Tech Companies, including size, industry, revenue, average employee salary, and more. We want to know which features influence the average employee salary. → regression
- You are given data of 100 individuals and their sequenced DNA and want to know whether these individuals will exhibit a particular disease based off their genomic mutations. We have information on 10,000 individual genomes and whether or not they exhibit the particular disease.

$$x_1 \quad x_2 \quad \dots \quad x_n \quad f(x) \quad \text{Class } 1$$

You are given reviews of few movies marked as positive, negative or neutral. Classifying reviews of a new movie is an example of

- a. Supervised learning
- b. Unsupervised learning
- c. Semi-Supervised learning
- d. Reinforcement learning

Imagine a newly-born starts to learn walking. It will try to find a suitable policy to learn walking after repeated falling and getting up. Specify what type of machine learning algorithm is best suited to do the same.

- a. Supervised Learning
- b. Unsupervised Learning
- c. Reinforcement Learning
- d. Semi-supervised Learning

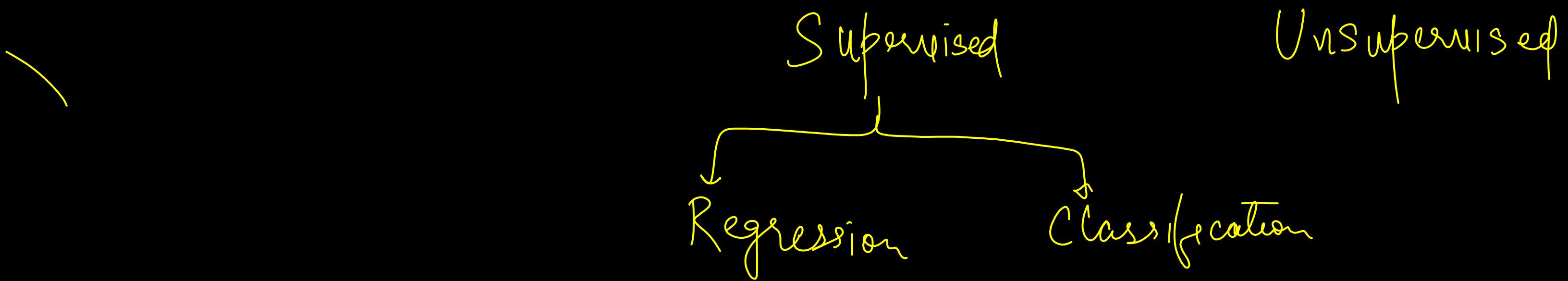
Which of the following are supervised learning problems (Multiple Correct)?

- A. Clustering Spotify users based on their listening history ✓
- B. Weather forecast using data collected by a satellite ✓
- C. Predicting tuberculosis using patient's chest X-Ray ✓
- D. Training a humanoid to walk using a reward system

↓  
Reinforcement

Classify the following as regression or classification tasks

- A. Predicting the outcome of an election → Class
- B. Predicting the weight of a giraffe based on its height → reg
- C. Predicting the emotion conveyed by a sentence → class
- D. Predict the temperature for the next day → reg
- E. Detect pneumonia from chest X-ray image → Class



# Linear Regression

$$\begin{matrix} x \\ y \end{matrix} \left\} \begin{matrix} y = ax + b \end{matrix} \right.$$

$x$	$y$
2	4
3	6
1	2
4	8

$y = 2x$



$x$	$y$	$y'$
2	5	5
3	7	7
1	3	3
4	9	9

$y = 2x + 1$      $y' = 2x + 1$

$x$	$y$	$y'$
2	4	0
3	6	0
1	3	1
4	8	0

Actual value      Predicted value      Error ( $y - y'$ )

$\downarrow$

$y' = 2x \rightarrow$  Predicted value

Best fit line

$x$	$y$	$y'$	Error ( $y - y'$ )	Error <sup>2</sup>	Error
1	3	2	-1	1	1
2	4	4	0	0	0
3	6	6	0	0	0
4	8	8	0	0	0
5	9	10	-1	1	1

$$y = 2x$$

Mean Squared Error =  $\frac{1}{2n} \sum_{i=1}^n (y_i - y'_i)^2$

$$= \frac{1}{2n} \left[ (y_1 - y'_1)^2 + (y_2 - y'_2)^2 + \dots + (y_n - y'_n)^2 \right]$$

Mean Absolute Error

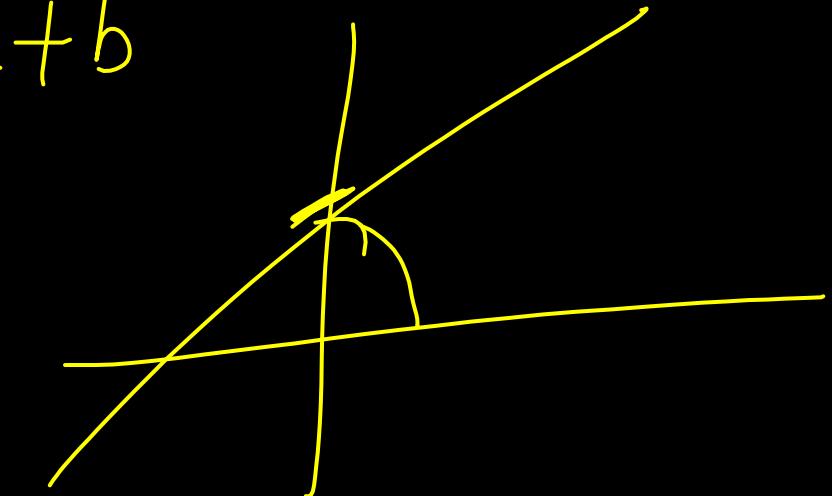
$$= \frac{1}{2n} \sum_{i=1}^n |y_i - y'_i|$$

What is Best Fit Line?

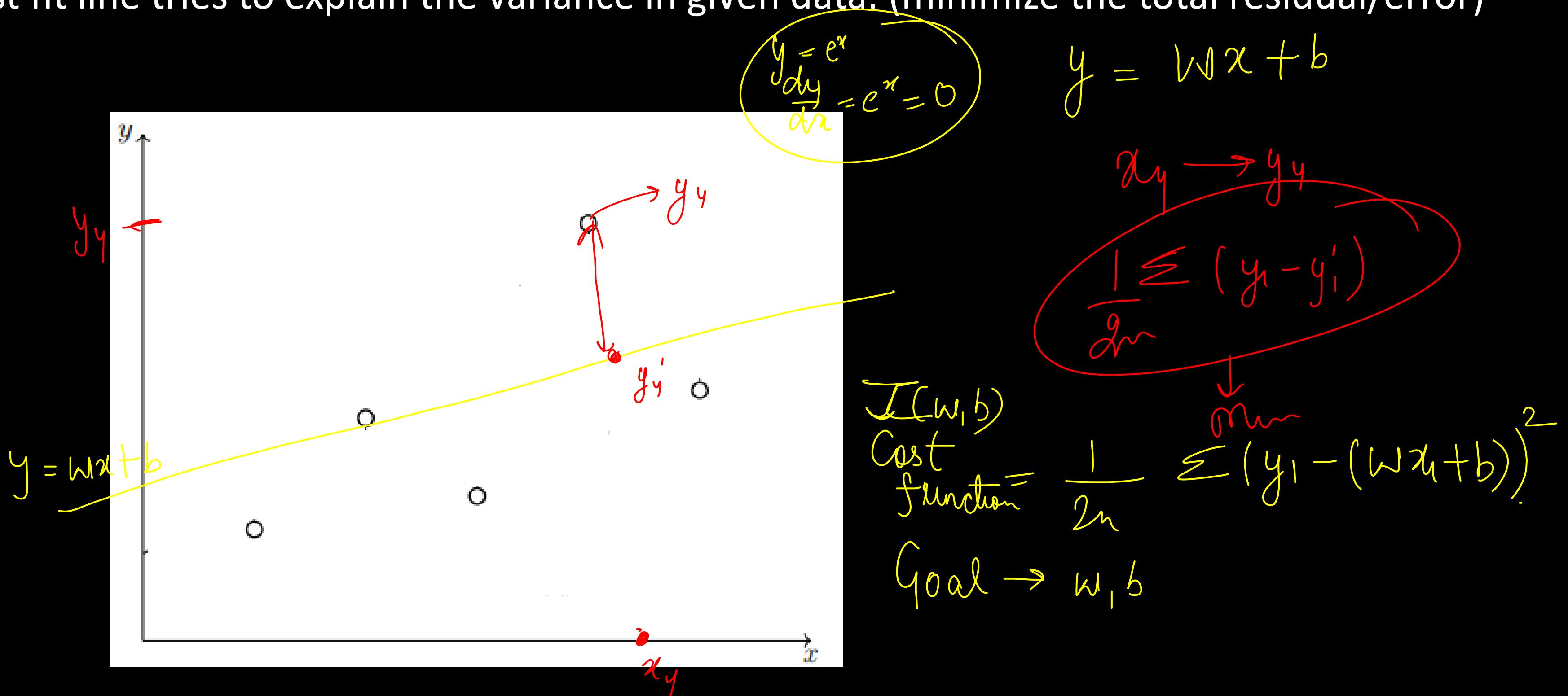
$$\boxed{y_i' = w\bar{x}_i + b}$$

$$\bar{y} = w\bar{x} + b$$

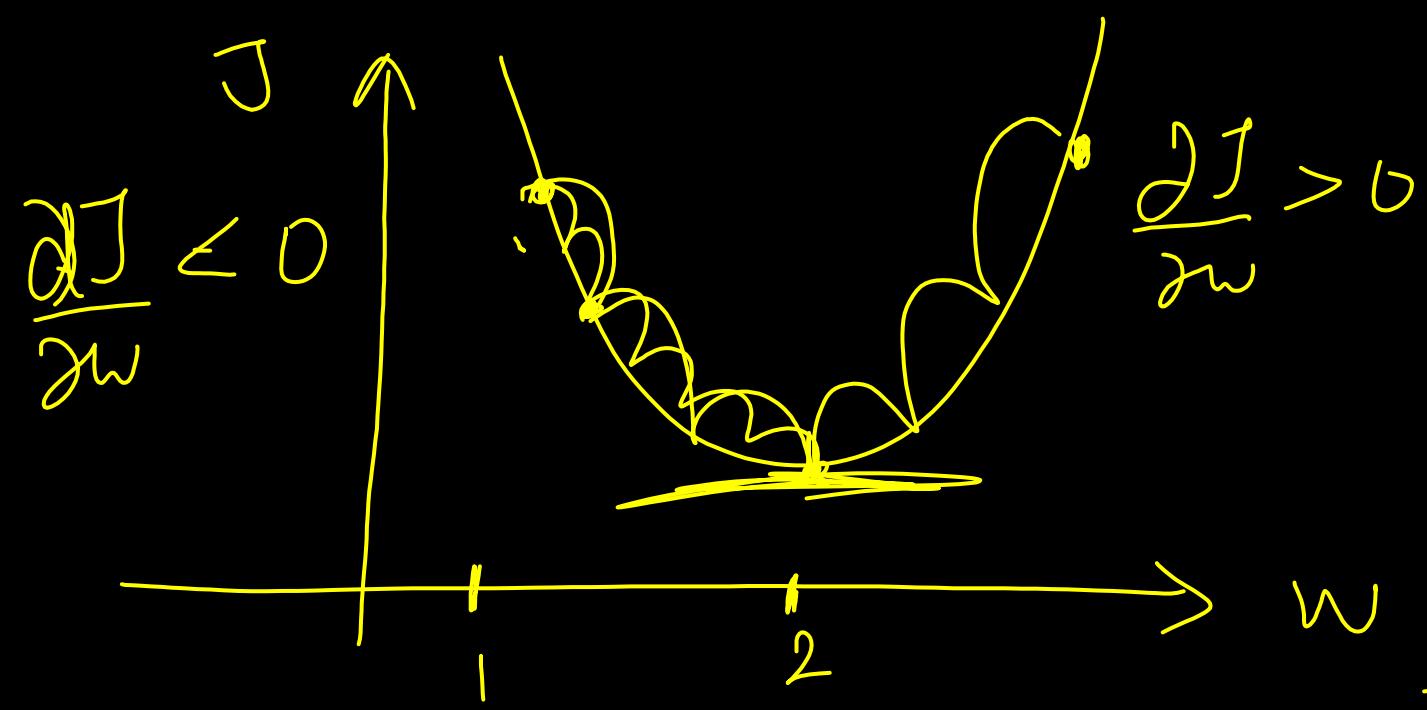
$$\bar{y} = w\bar{x} + b$$



Best fit line tries to explain the variance in given data. (minimize the total residual/error)



## Gradient Descent



$$w = w - \alpha \left( \frac{\partial J}{\partial w} \right)$$

$$\boxed{w = w - \alpha, \frac{\partial J}{\partial w}}$$

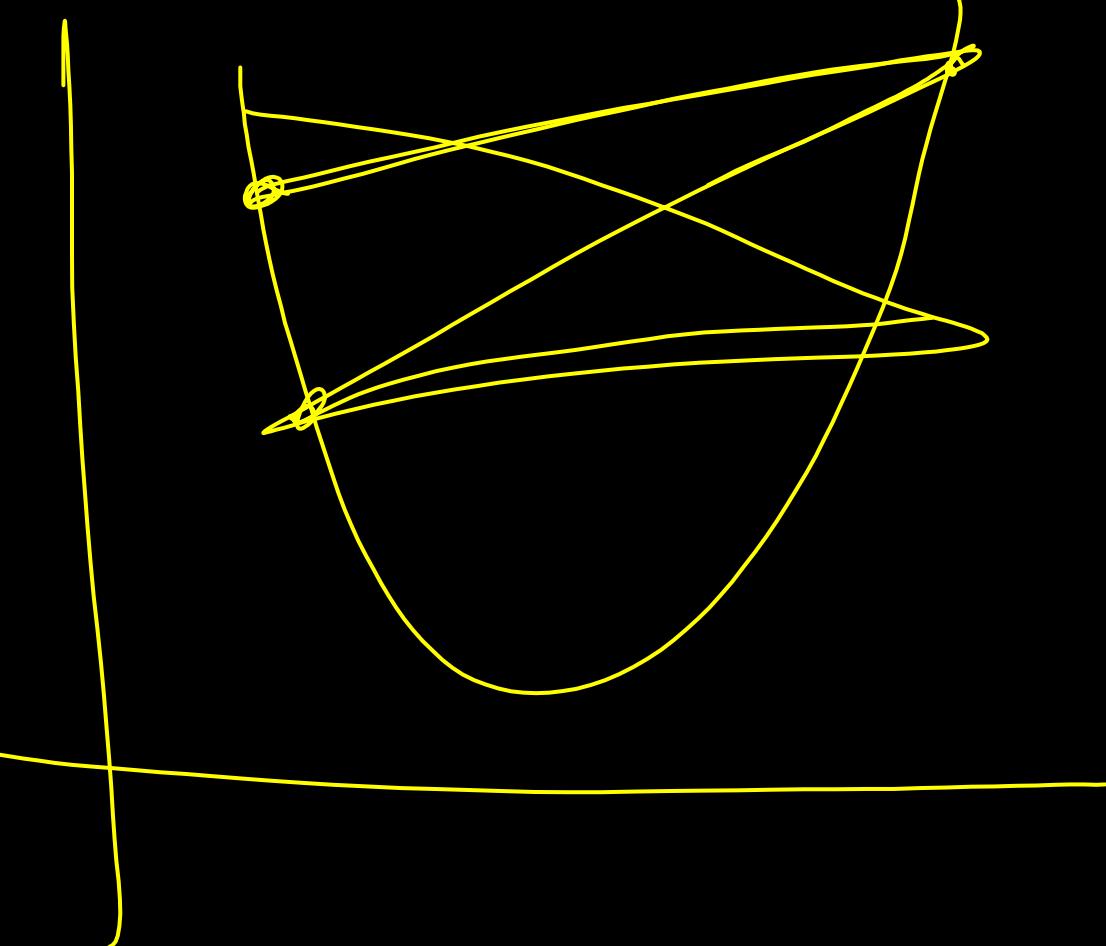
$$= 1 - \alpha (-2)$$

$$= 1 + 0.1 \times 2$$

$$= 1.2$$

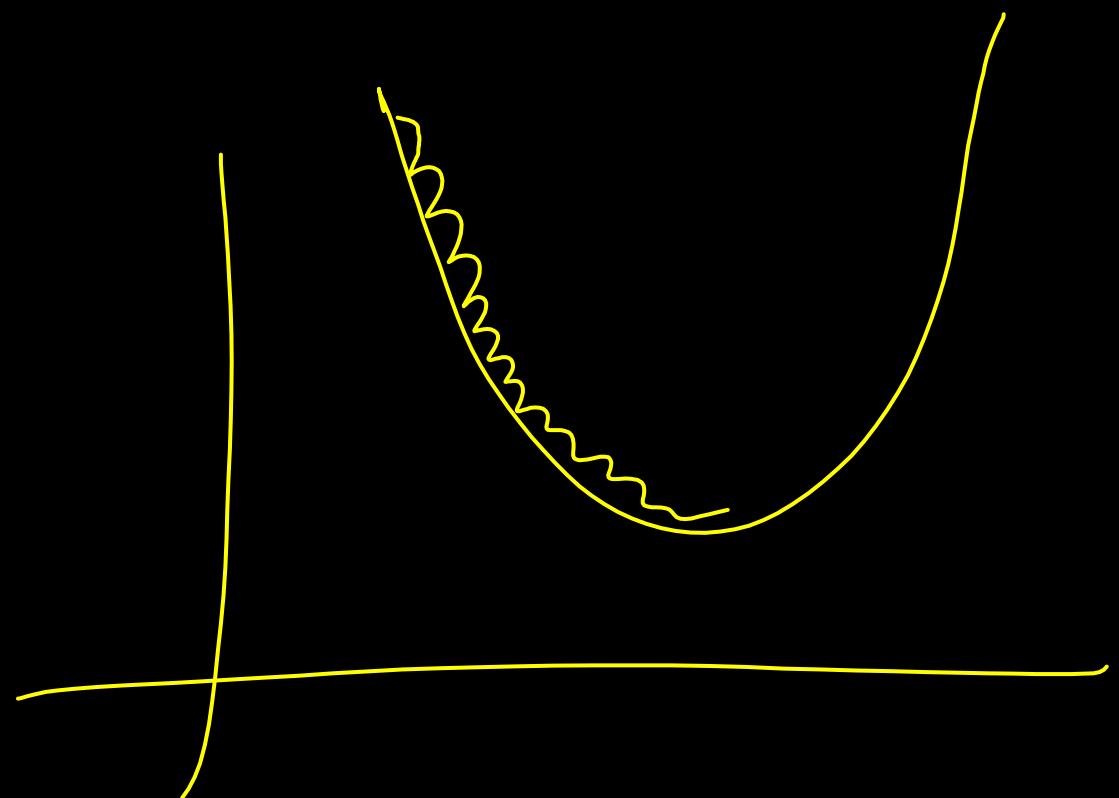
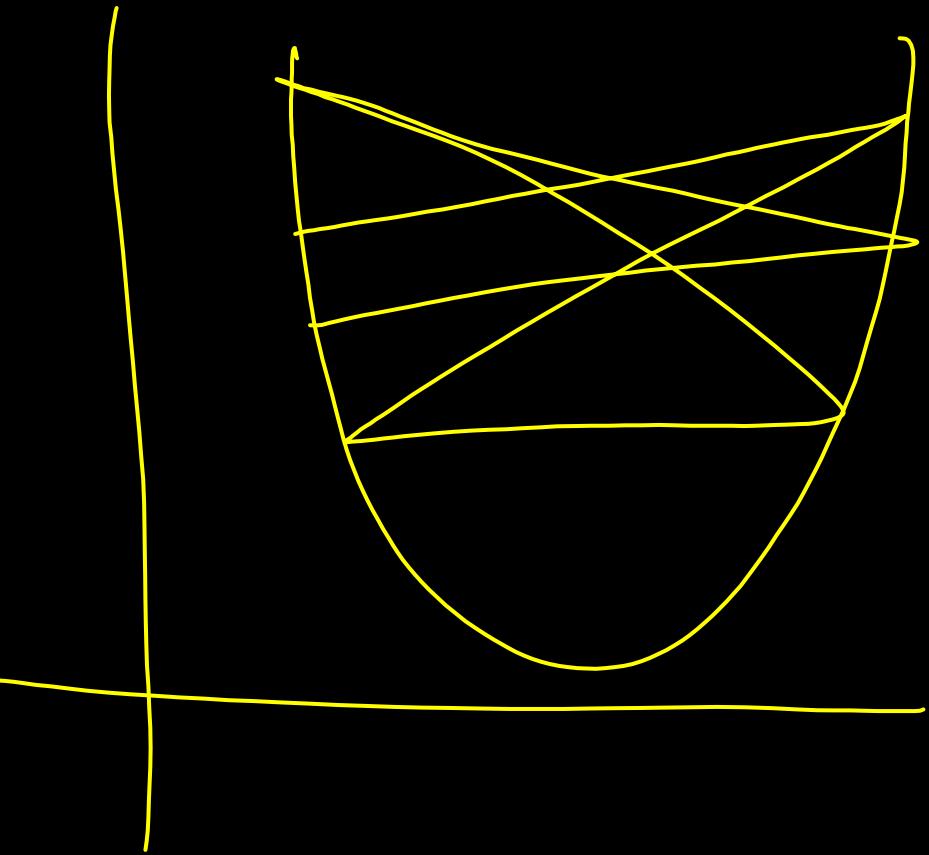
$$\boxed{b = b - \alpha \frac{\partial J}{\partial b}}$$

$$\begin{aligned} w &= 1.2 - 0.1 (1.2) \\ &= 1.08 \end{aligned}$$



## Learning Rate

$$w = w - \alpha \frac{\partial J}{\partial w}, \quad b = b - \alpha \frac{\partial J}{\partial b}$$



$\alpha \rightarrow$  too large.

$\alpha \rightarrow$  too small

Consider the following training set of  $m = 4$  training examples

x	y
1	0.5
2	1
4	2
0	0

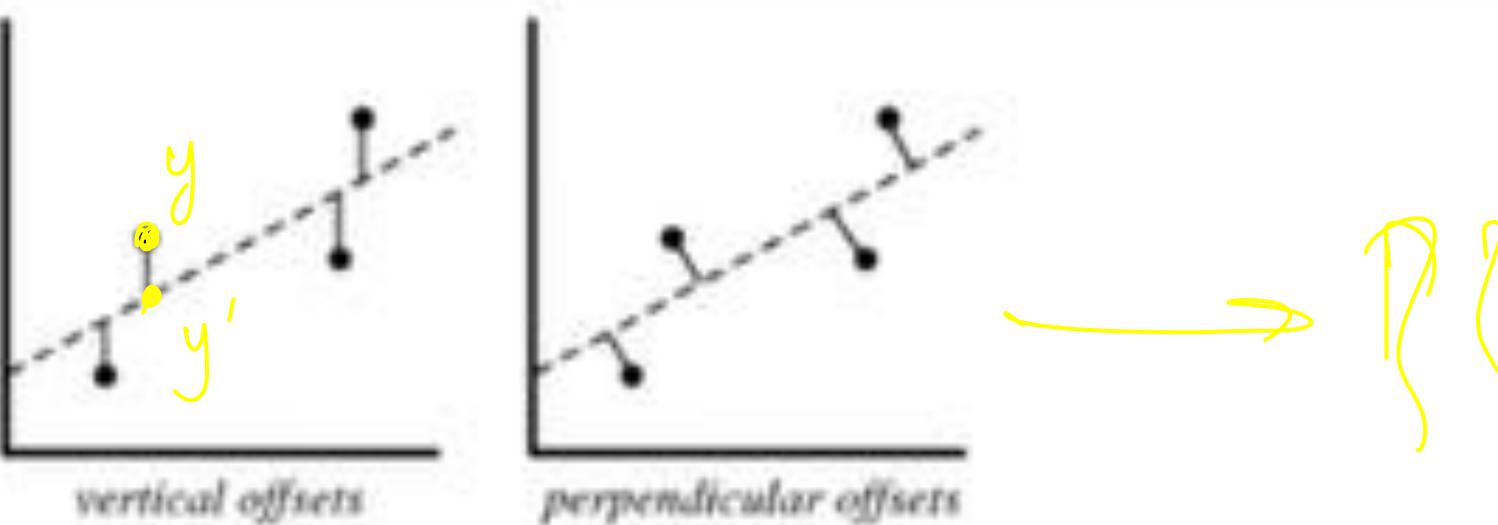
$$y = \frac{1}{2}x + 0$$

Consider the linear regression model  $y = wx + b$ . What are the values of w and b, that you would expect to obtain upon running gradient descent on this model?

- (a)  $b = 0.5$  and  $w = 0$
- (b)  $b = 0.5$  and  $w = 0.5$
- ~~(c)~~  $b = 0$  and  $w = 0.5$
- (d)  $b = 0$  and  $w = 0$

Answer: (c)

Which of the following offsets, do we use in linear regression's least square line fit? Assume the horizontal axis is the independent variable and vertical axis is dependent variable.



- A) Vertical offset
- B) Perpendicular offset
- C) Both, depending on the situation
- D) None of above



PCA

Which of the following if any is a valid cost function in a regression setting and why?

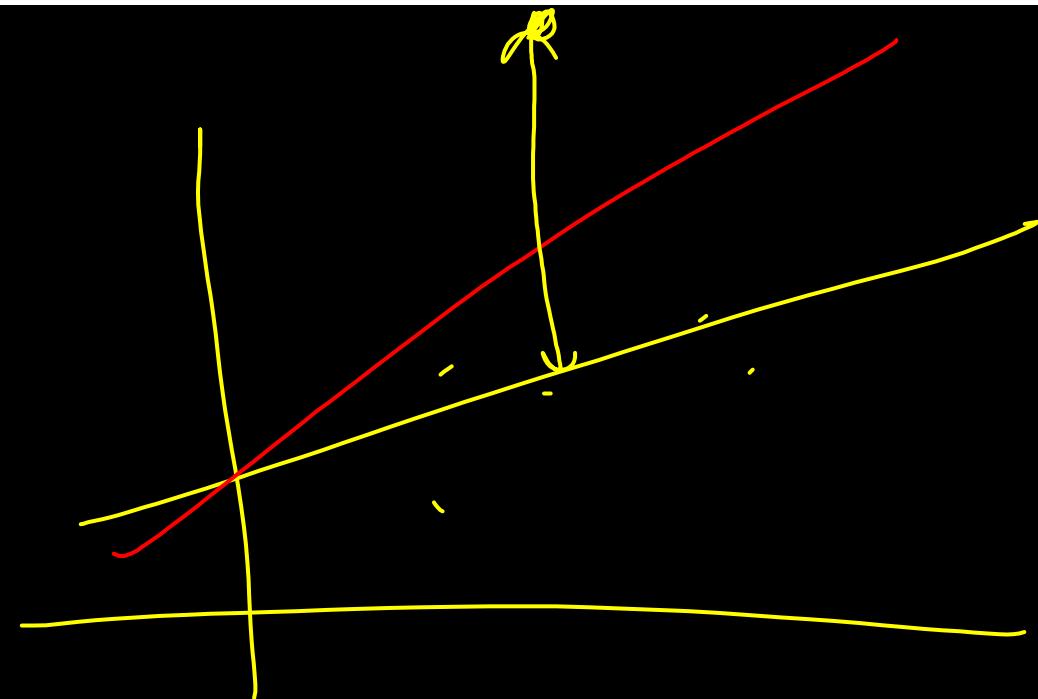
a.  $J(w) = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2$

b.  $J(w) = \frac{1}{2m} \sum_{i=1}^m |f_w(x^{(i)}) - y^{(i)}|$

c.  $J(w) = \frac{1}{2m} \sum_{i=1}^m |f_w(x^{(i)}) - y^{(i)}|$

Which statement is true about outliers in Linear regression?

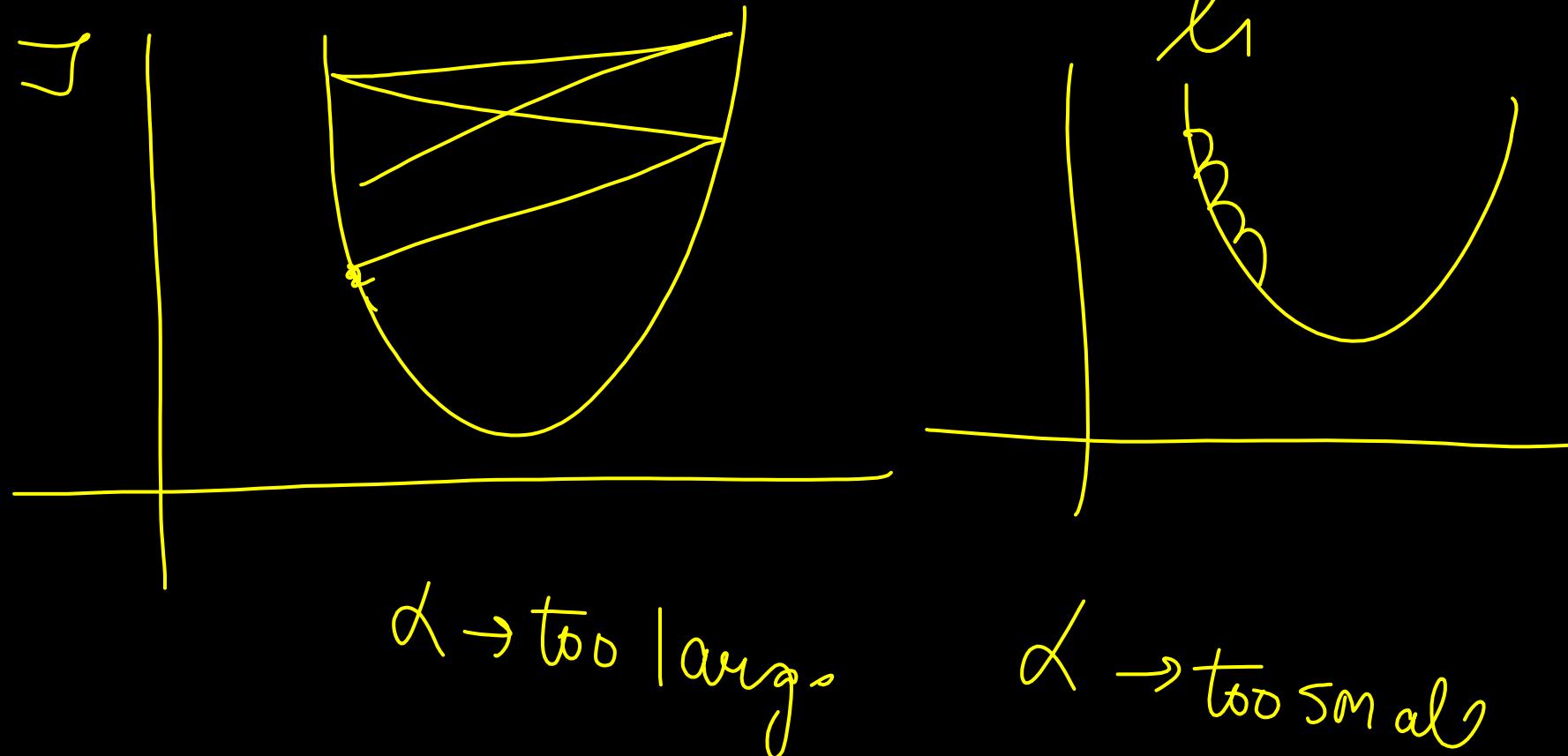
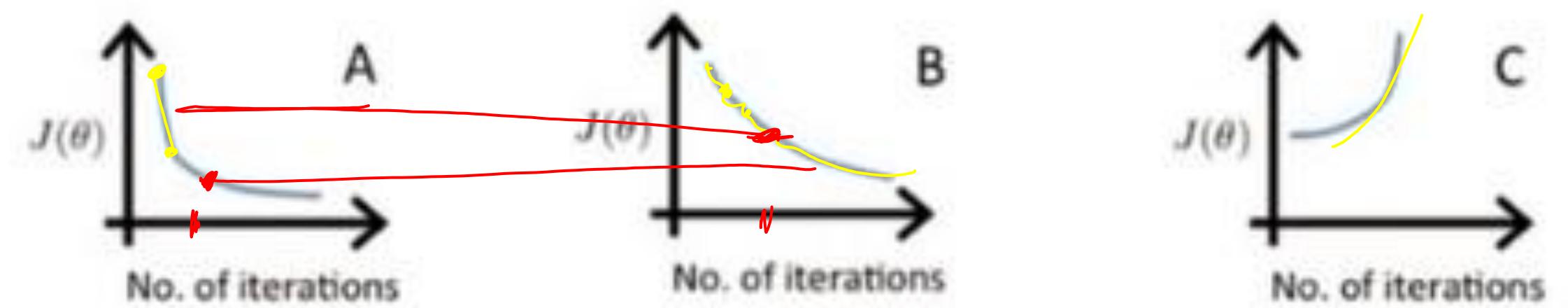
- A) Linear regression is sensitive to outliers
- B) Linear regression is not sensitive to outliers
- C) Can't say
- D) None of these



1, 2, 3, 4, 1000

Suppose  $\alpha_1$ ,  $\alpha_2$  and  $\alpha_3$  are the three learning rates for A, B, C respectively. Which of the following is true about  $\alpha_1$ ,  $\alpha_2$  and  $\alpha_3$ ?

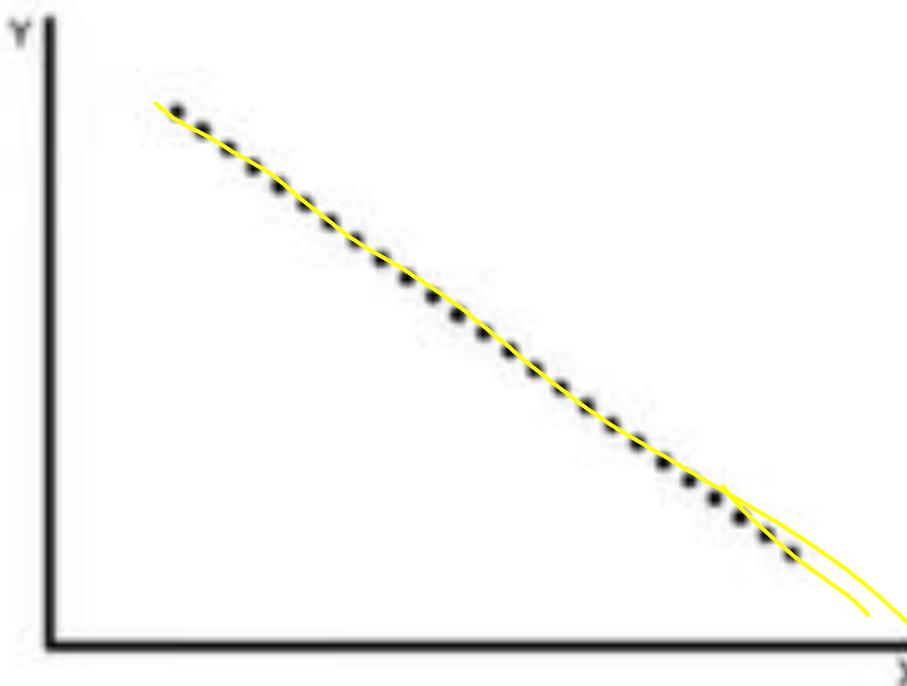
- A)  $\alpha_2 < \alpha_1 < \alpha_3$
- B)  $\alpha_1 > \alpha_2 > \alpha_3$
- C)  $\alpha_1 = \alpha_2 = \alpha_3$
- D) None of these



$$\alpha_2 < \alpha_1 < \alpha_3$$

$\alpha_2$   
Small  
 $\alpha_1$   
 $\alpha_3$   
too large

Consider the following data where one input(X) and one output(Y) is given. What would be the cost for this data if you run a Linear Regression model of the form ( $Y = w_1 * x_1 + b$ )?



- A) Less than 0
- B) Greater than zero
- C) Equal to 0
- D) None of these

The selling price of a house depends on the following factors. For example, it depends on the number of bedrooms, number of kitchen, number of bathrooms, the year the house was built and the square footage of the lot. Given these factors, predicting the selling price of the house is an example of \_\_\_\_\_ task.

- a. Binary Classification
- b. Multilabel Classification
- c. Simple Linear Regression
- d. Multiple Linear Regression

$$x_1 \ x_2 \ \dots \ x_k \ y$$

$$y = w_1 x_1 + w_2 x_2 + \dots + w_k x_k + b$$

## Multiple regression

Allows a response variable Y to be modeled as a linear function of multidimensional feature vector



$$\begin{aligned}
 & y' = \mathbf{w}^T \mathbf{x} \\
 & y' = w_0 x_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 \\
 & = [w_0 \ w_1 \ w_2 \ w_3] \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \\
 & = \mathbf{w}^T \mathbf{x}
 \end{aligned}$$

$$\begin{aligned}
 & y' = w_1 x_1 + w_2 x_2 + w_3 x_3 + b \\
 & w_l = w_l - \alpha \frac{\partial J}{\partial w_l} \\
 & b = b - \alpha \frac{\partial J}{\partial b}
 \end{aligned}
 \longrightarrow \text{Goal} \rightarrow w_1, w_2, w_3, b$$

$$J = \frac{1}{2n} \sum_{i=1}^n (y_i - y'_i)^2$$

In the context of machine, **sparsity** refers to a model where many of the feature weights (coefficients) are exactly zero. In simple terms:

**Sparse models** are models where only a small number of features (input variables) have non-zero coefficients, meaning that the model effectively ignores most of the features.

$$\underbrace{w_1, w_2, \dots, w_k}_{\text{Sparse}} \rightarrow w_1 = 0, w_2 = 0, w_3 \neq 0, w_4 \neq 0 \dots$$

**Dense models**, on the other hand, use most or all of the features with non-zero coefficients.

$$w_1 \neq 0, w_2 \neq 0 \dots$$



$$w_1 \approx 0$$

$$0.006/$$

## L1 Regularization (Lasso)

$$\text{L1 Regularization Term} = \lambda \sum_{i=1}^n |w_i|$$

$w_1 \leftarrow 0$

$\downarrow$   
Sparse

$$\begin{aligned} y' &= w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b \\ \text{Error} &= \frac{1}{2n} (y - y')^2 + \lambda \sum |w_i| \\ &\quad \downarrow \quad \downarrow \\ &\quad \lambda = 1000 \quad \downarrow \\ &\quad 1000 \times \\ &\quad = 1000 \end{aligned}$$

## L2 Regularization (Ridge):

$$\text{L2 Regularization Term} = \lambda \sum_{i=1}^n w_i^2$$

$$\text{Error} = \frac{1}{2n} \sum (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^n w_i^2$$

$w_1 = 0.0061$

$w_1 \neq 0$

$$\hat{y} = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

## Combined L1 and L2 Regularization (Elastic Net):

$$\text{Elastic Net Regularization Term} = \lambda_1 \sum_{i=1}^n |w_i| + \lambda_2 \sum_{i=1}^n w_i^2$$

Here,  $\lambda_1$  and  $\lambda_2$  control the strengths of L1 and L2 regularization, respectively.

Select all the following reasons why we would use LASSO over Ridge?

- A. It can help us identify which features are important
- B. It is faster to learn the weights for LASSO than for Ridge
- C. LASSO usually achieves lower generalization error than Ridge
- D. If there are many features, the model learned using LASSO can make predictions more efficient

Answer: (A), (C)

With Lasso Regression the influence of the hyper parameter lambda, as lambda tends to zero the solution approaches to \_\_\_\_\_.

- a) Zero.
- b) One.
- c) Linear regression.
- d) Infinity.

Answer: (c)

A handwritten note on a black background. At the top left, there is a circled equation  $\lambda = 0$ . To the right, there is a large bracket enclosing a formula:  $\text{Error} = \frac{1}{2n} \sum (y_i - y_i^')^2 + \lambda \sum |w_i|$ . An arrow points from the circled  $\lambda$  towards the zero value in the circled equation above it.

$$\text{Error} = \frac{1}{2n} \sum (y_i - y_i^')^2 + \lambda \sum |w_i|$$

In this Lasso and Ridge regression as alpha value increases, the impact on slope is

a) Slope is fixed for whatever maybe the value.

b) the slope of the regression line reduces and becomes horizontal. ✓

c) the slope of the regression line increases and becomes vertical

d) None of the above.

Answer: (b)

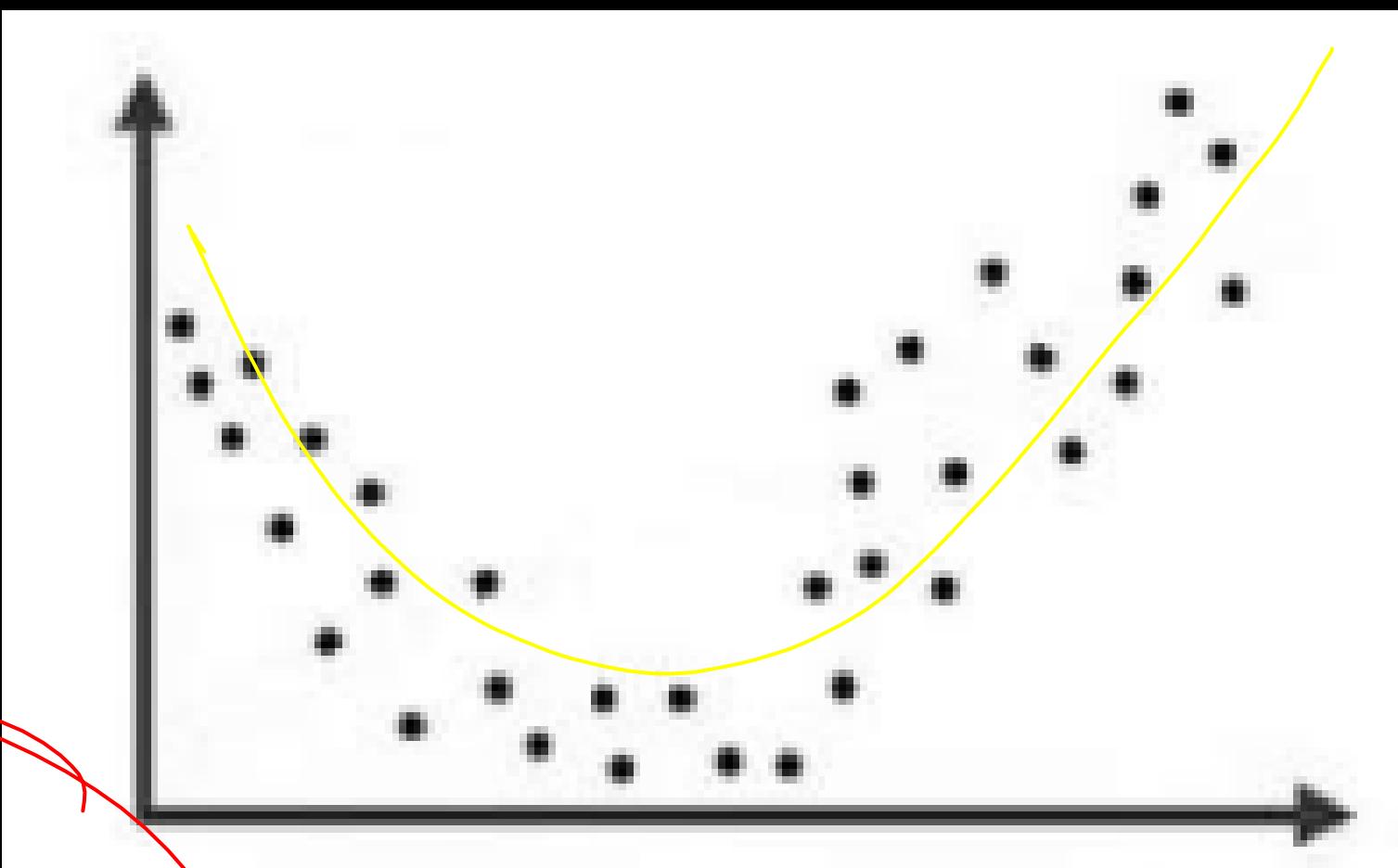
$$\lambda \rightarrow 1000 \\ \omega_1 = 0.001 \quad /$$

$$\lambda = 10,000 \uparrow \\ \omega \downarrow$$

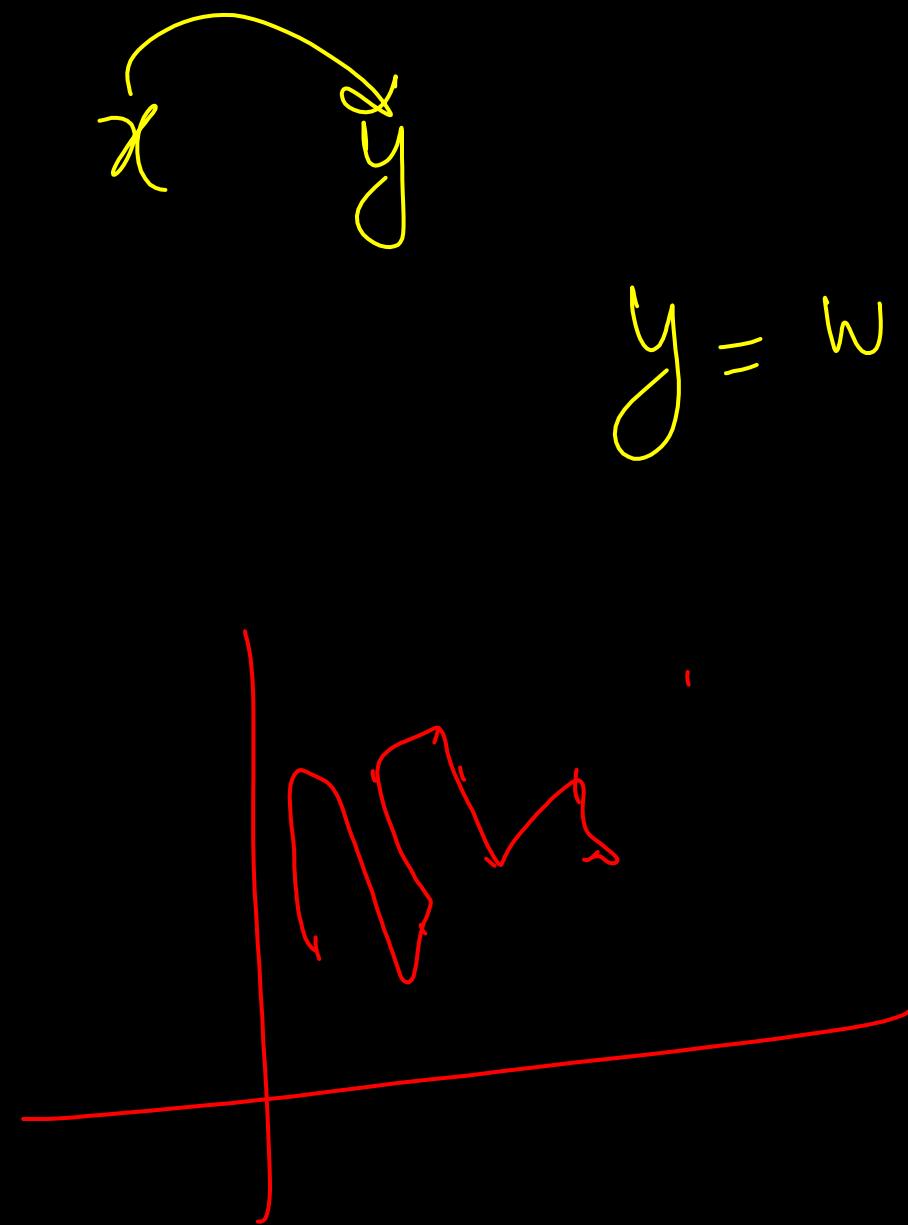
$$\lambda \uparrow \\ \omega_1 \rightarrow 0$$

# Polynomial regression

$$y = w_1x + w_2x^2 + w_3x^3 + w_4x^4 + w_5x^5 + w_6x^6 + w_7x^7 + b$$



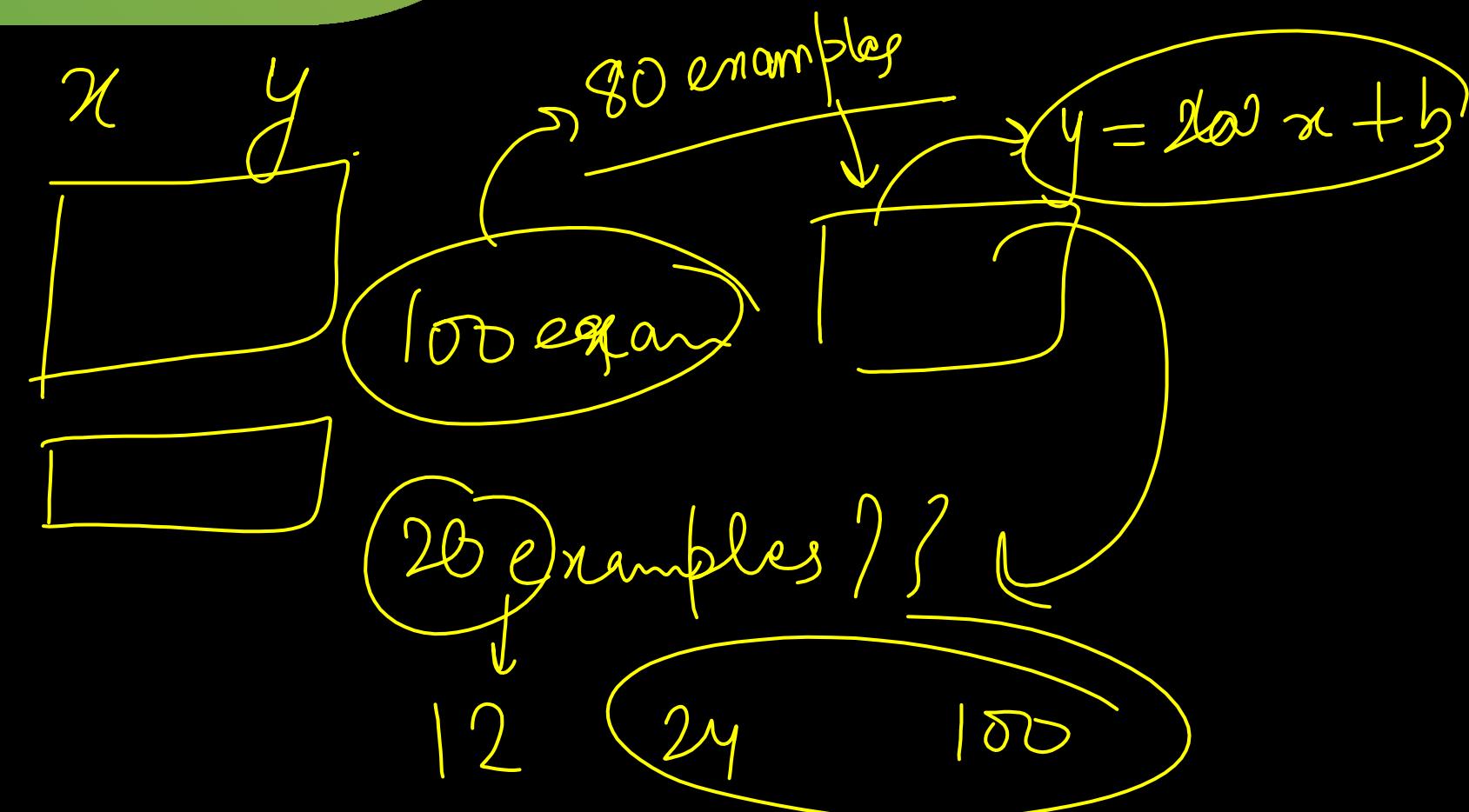
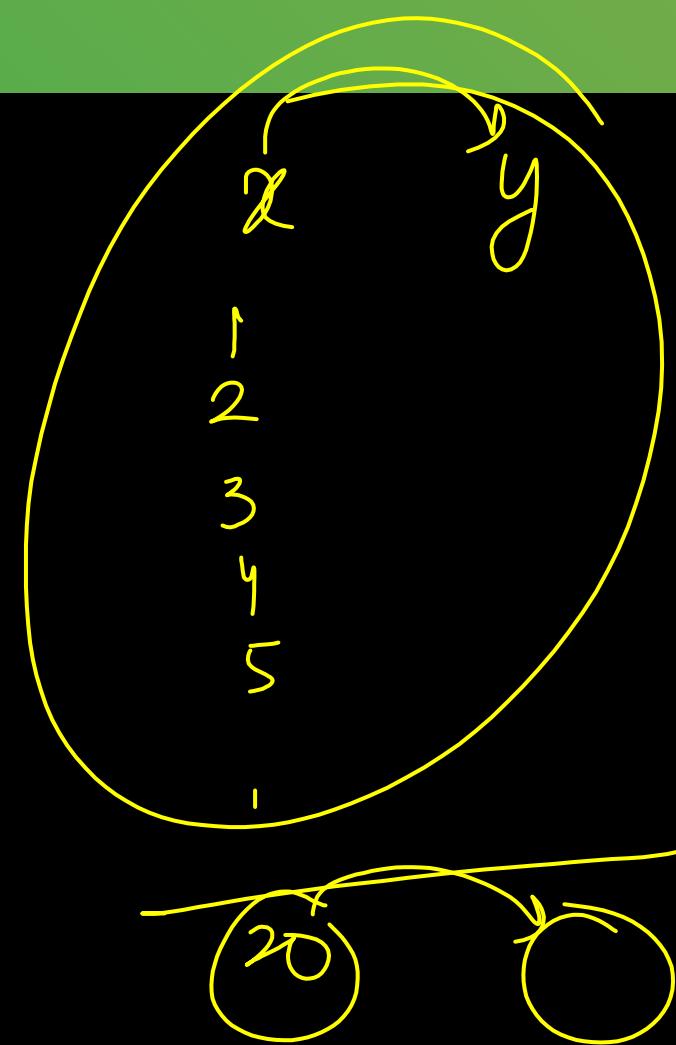
$$y = w_1x + w_2x^2 + b$$



$$y = w_1x + b$$

# Bias Variance Tradeoff

B



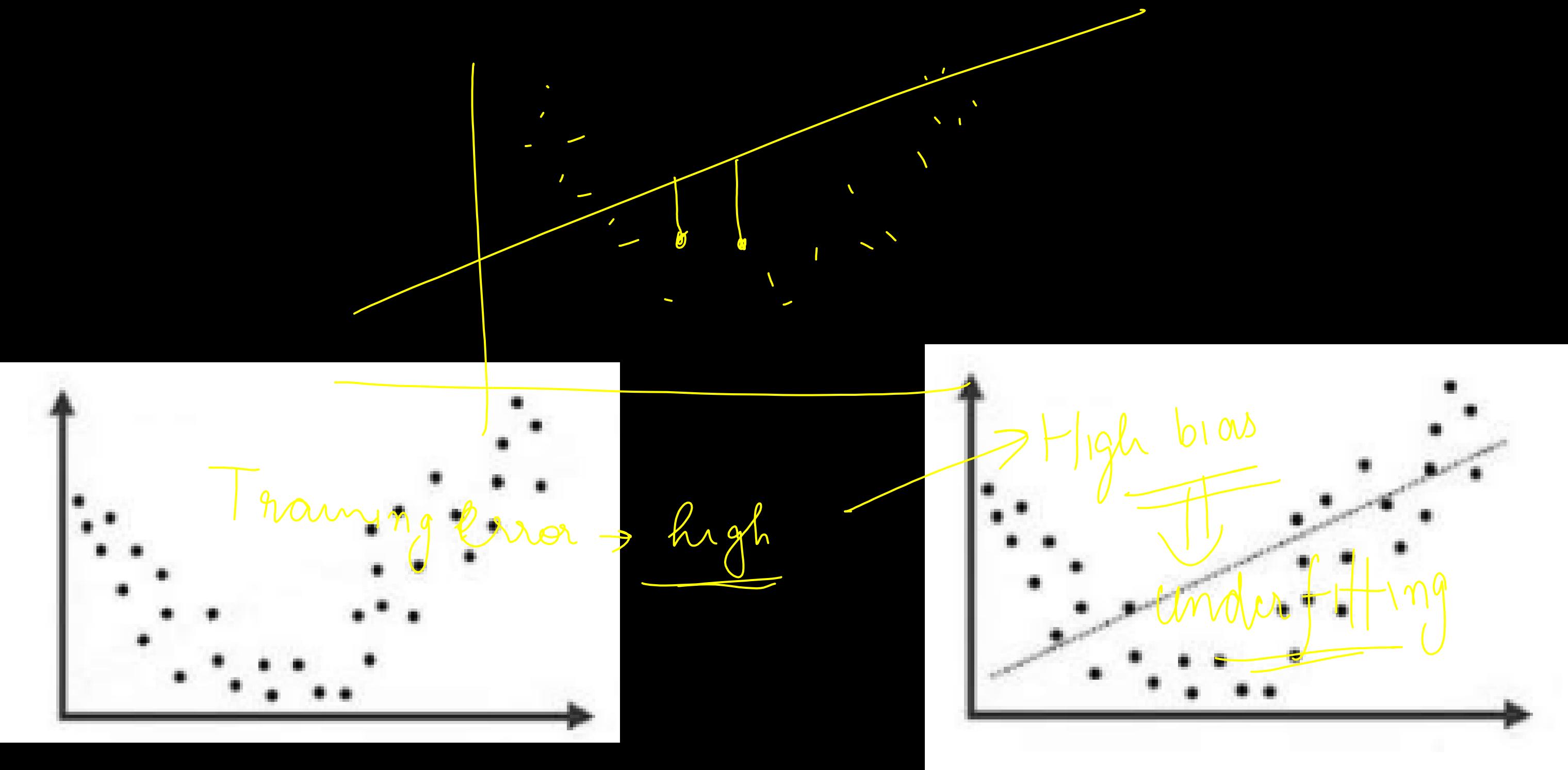
## Training & Test Data

In machine learning, data is split into training data and test data.

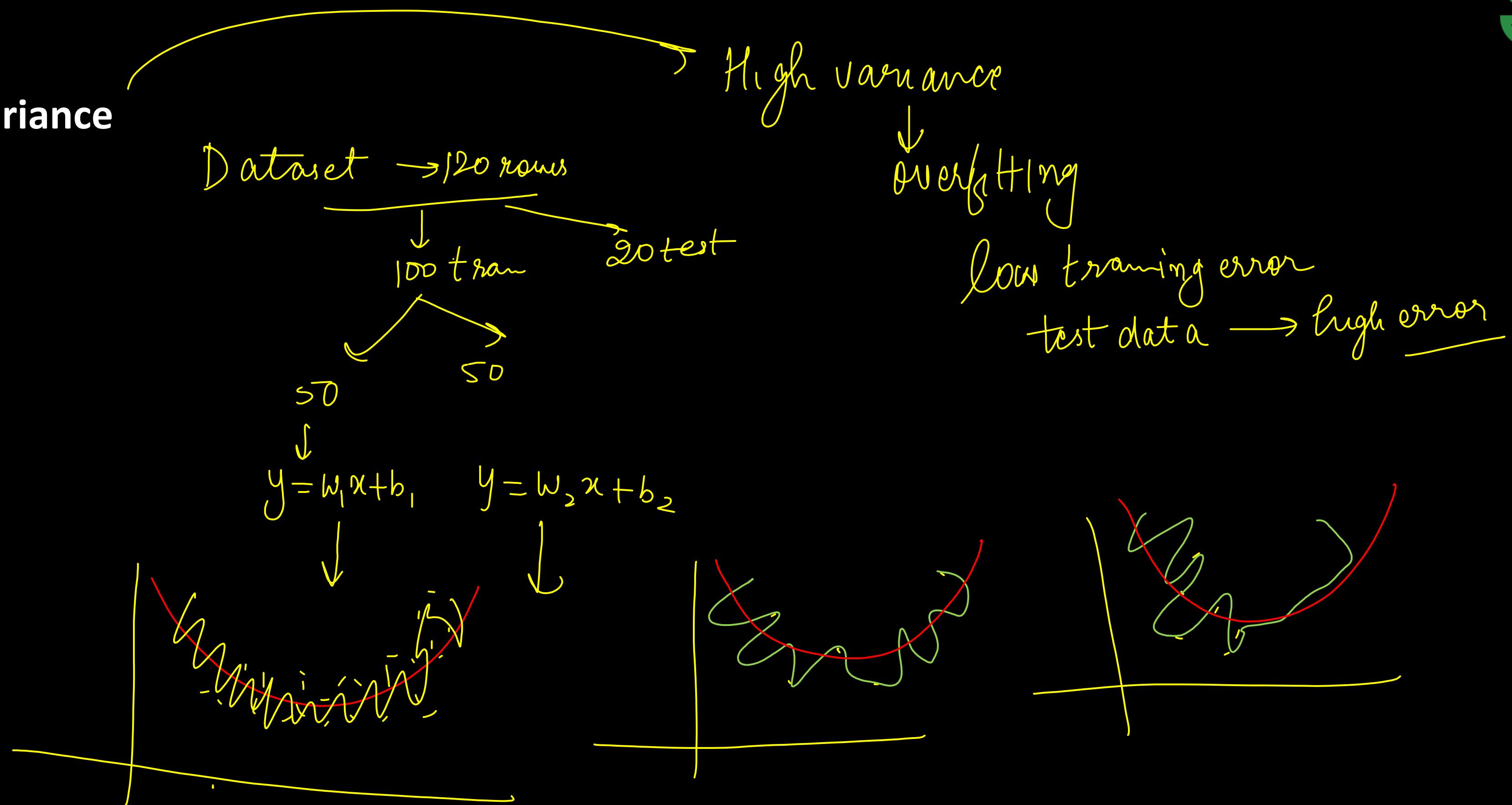
The first split of data, i.e. the initial reserve of data you use to develop your model, provides the training data.

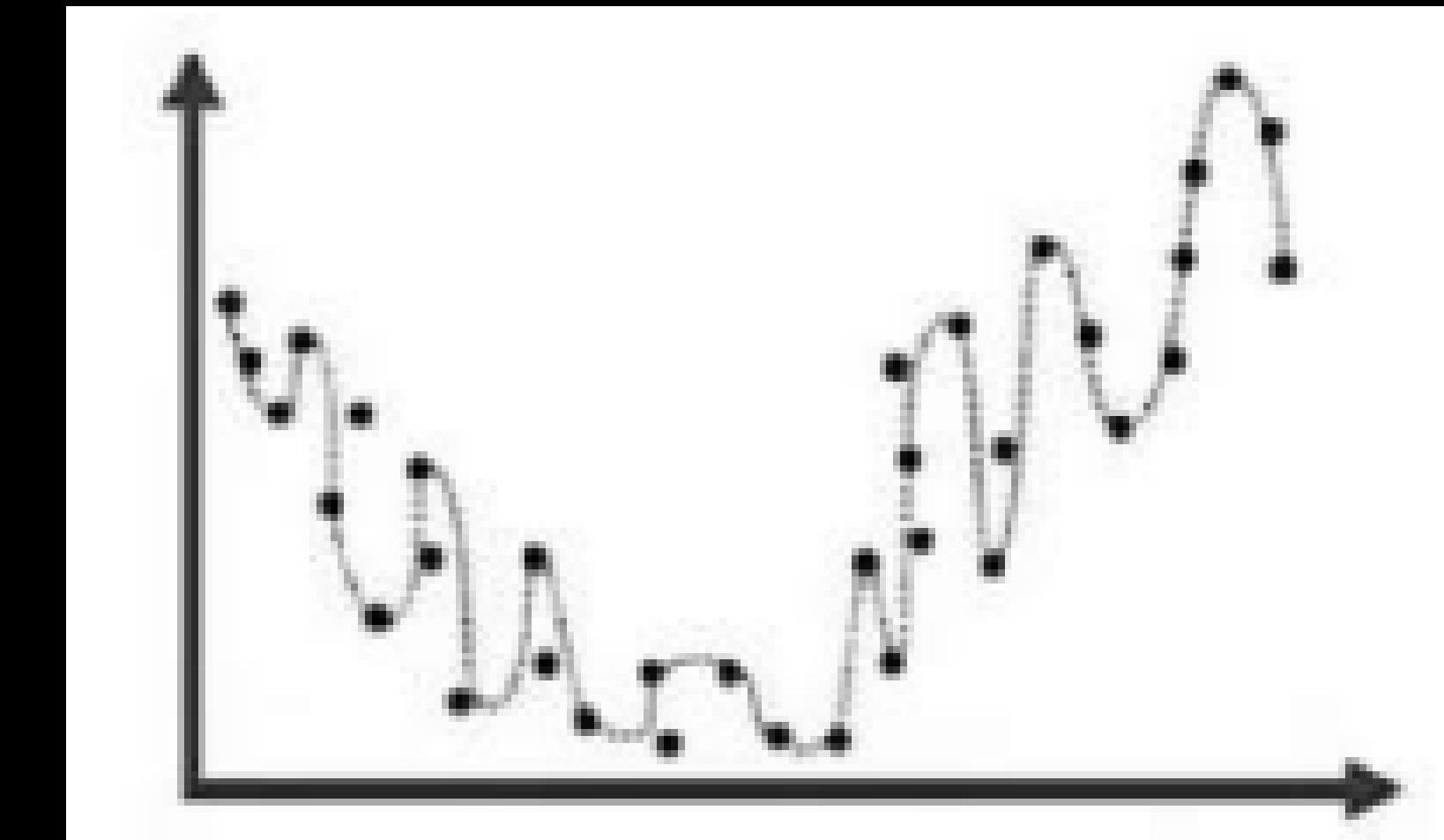
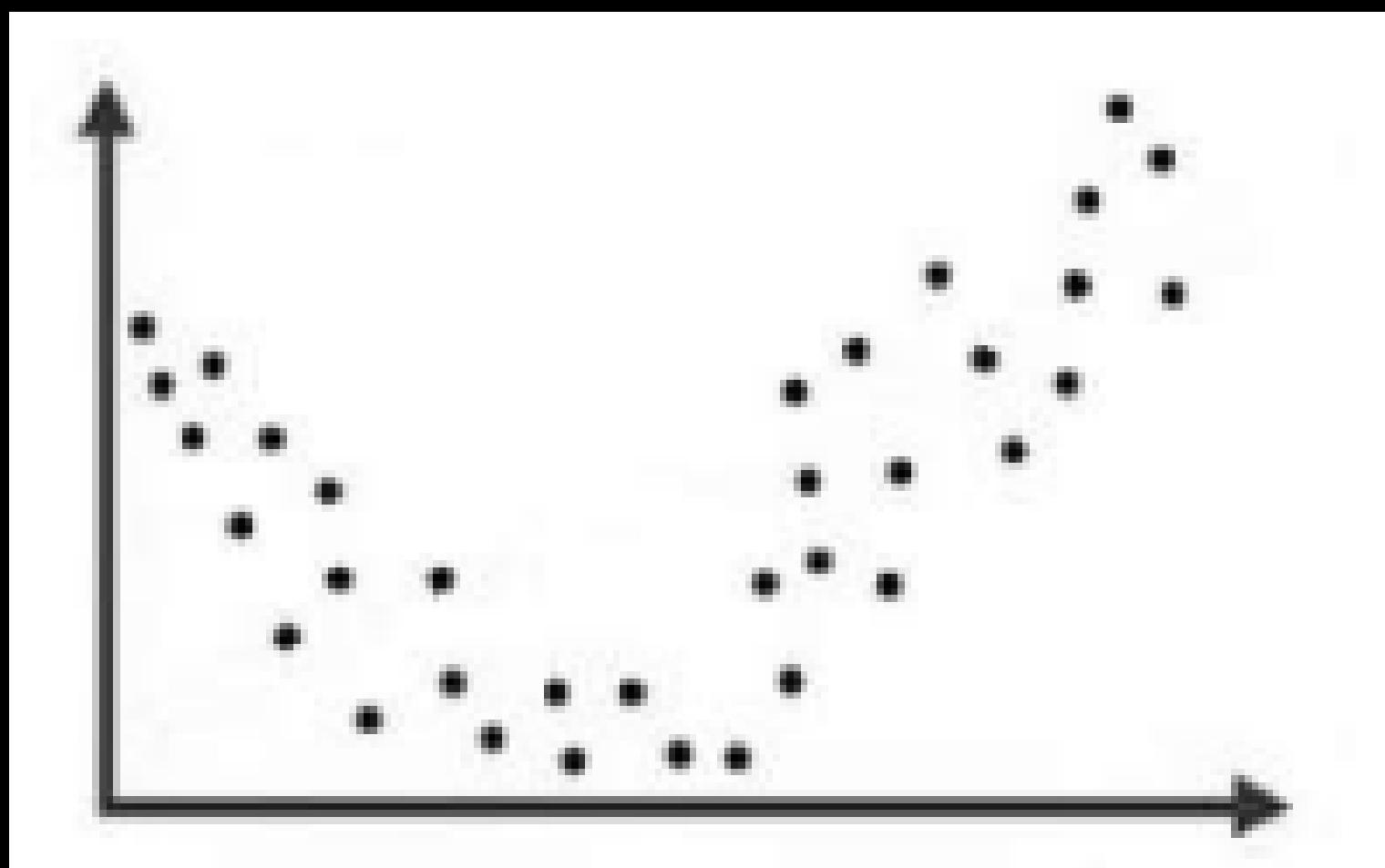
After you have successfully developed a model based on the training data and are satisfied with its accuracy, you can then test the model on the remaining data, known as the test data.

# Bias



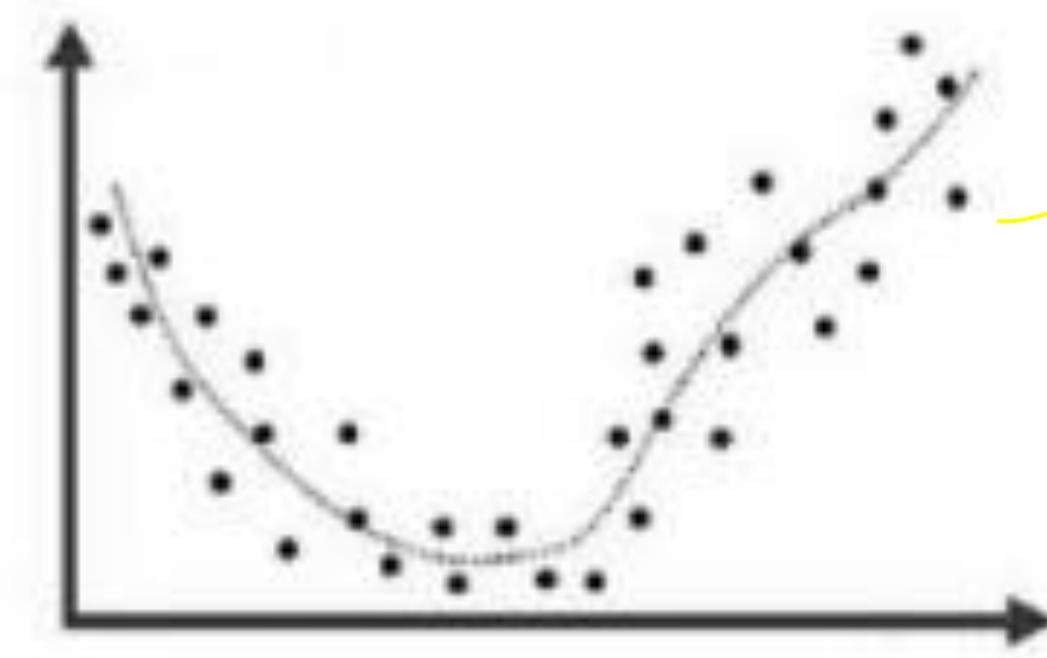
## Variance



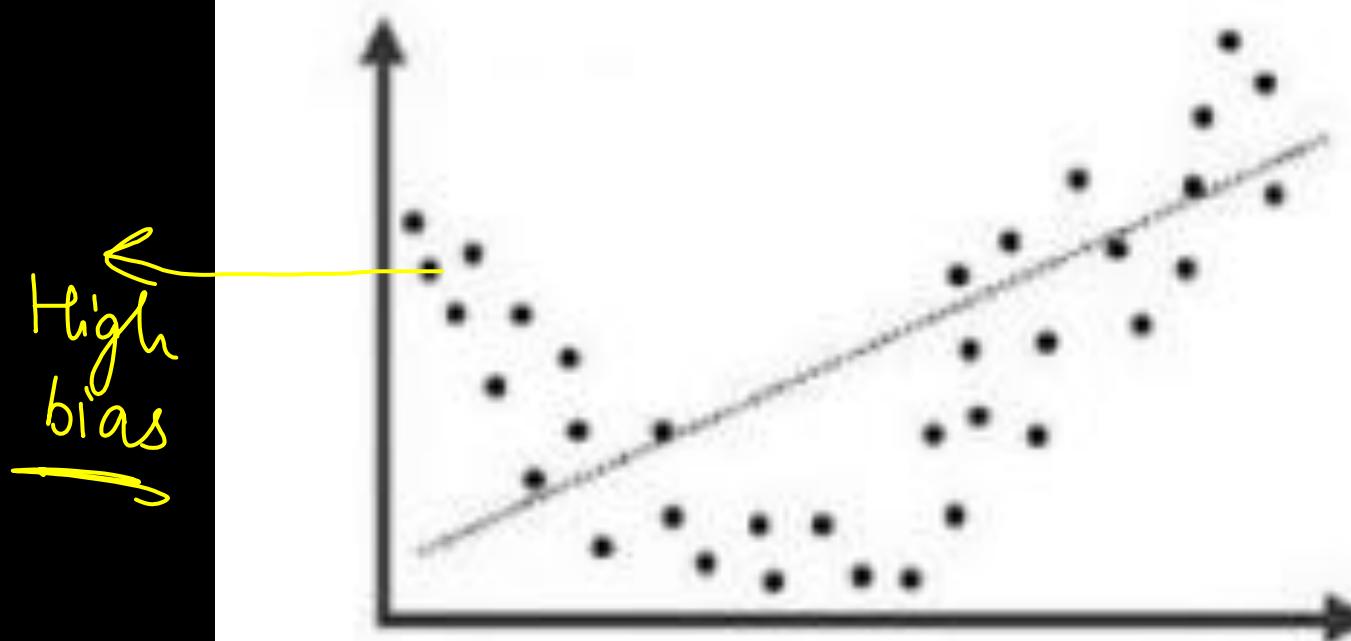




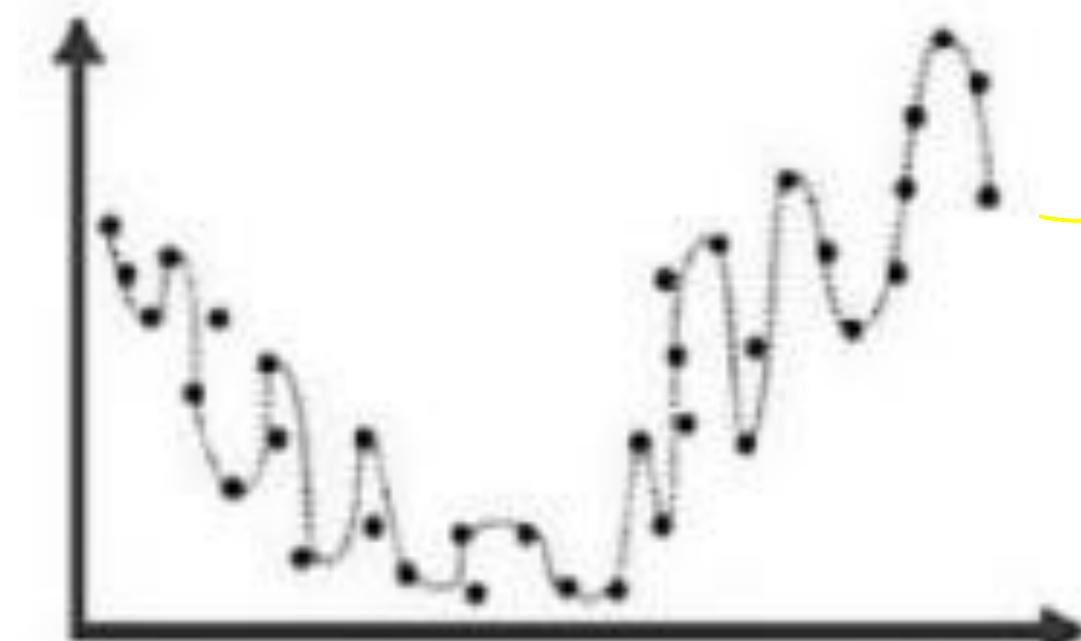
(a) Given dataset



(b) "Just right" model

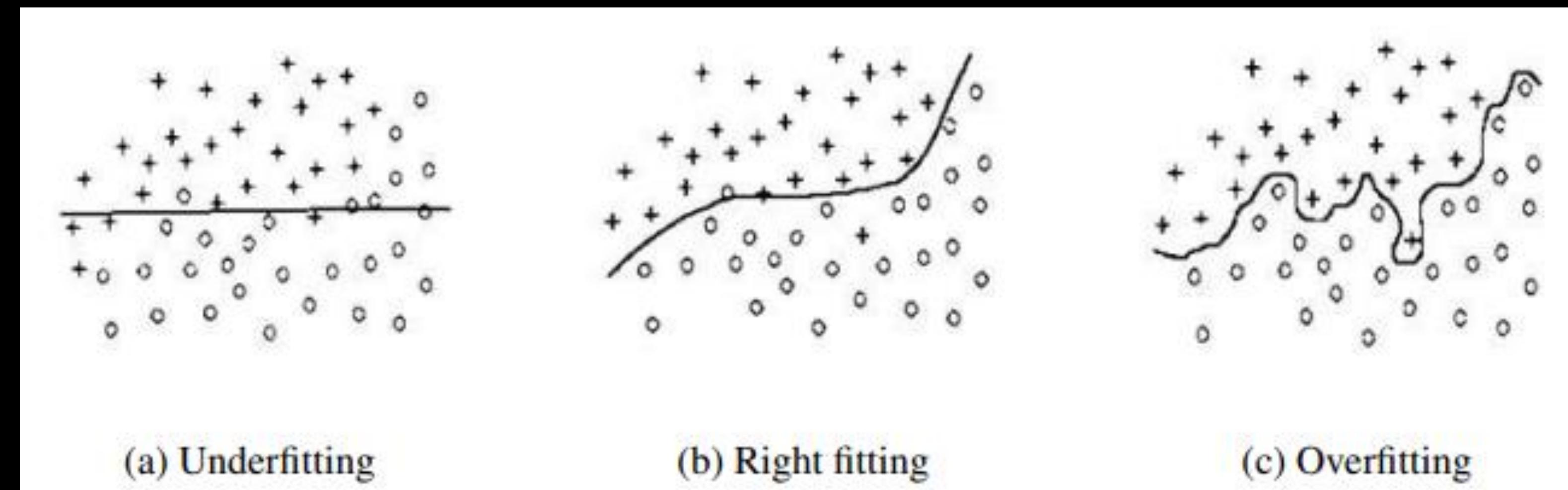


(c) Underfitting model



(d) Overfitting model

High variance,  $\downarrow$   
Low bias



## Tradeoff:

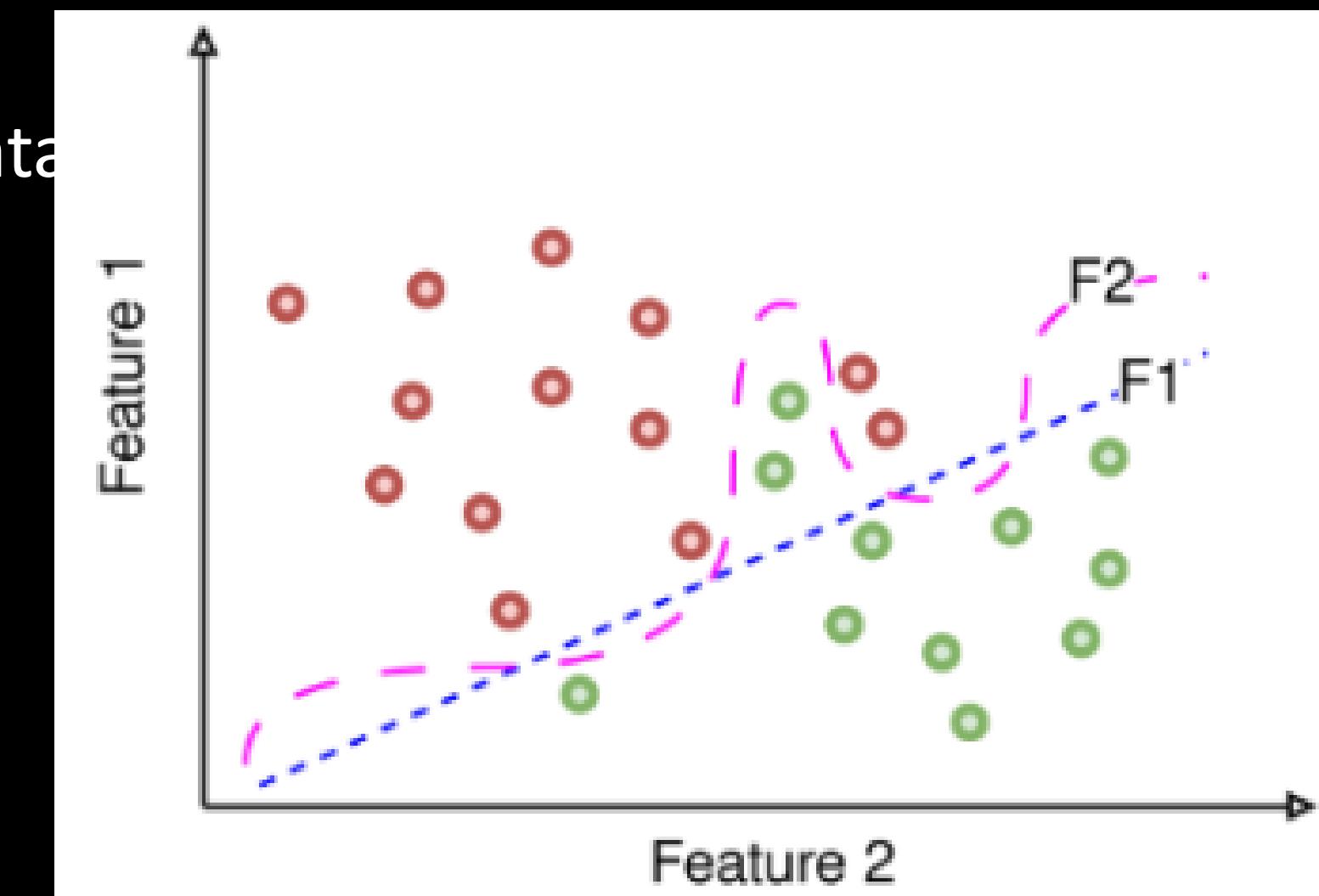
- The bias-variance tradeoff suggests that there is a balance to be struck between bias and variance. As you decrease bias (by increasing model complexity), you typically increase variance, and vice versa.
- The goal is to find the right level of **model complexity** that minimizes the combined error due to bias and variance. The objective is to achieve a model that generalizes well to new, unseen data.

Here is a 2-dimensional plot showing two functions that classify data points into two classes. The red points belong to one class, and the green points belong to another. The dotted blue line ( $F_1$ ) and dashed pink line ( $F_2$ ) represent the two trained functions.

Which of the two functions overfit the training data?

- A. Both functions  $F_1$  &  $F_2$
- B. Function  $F_1$
- C. Function  $F_2$
- D. None of them

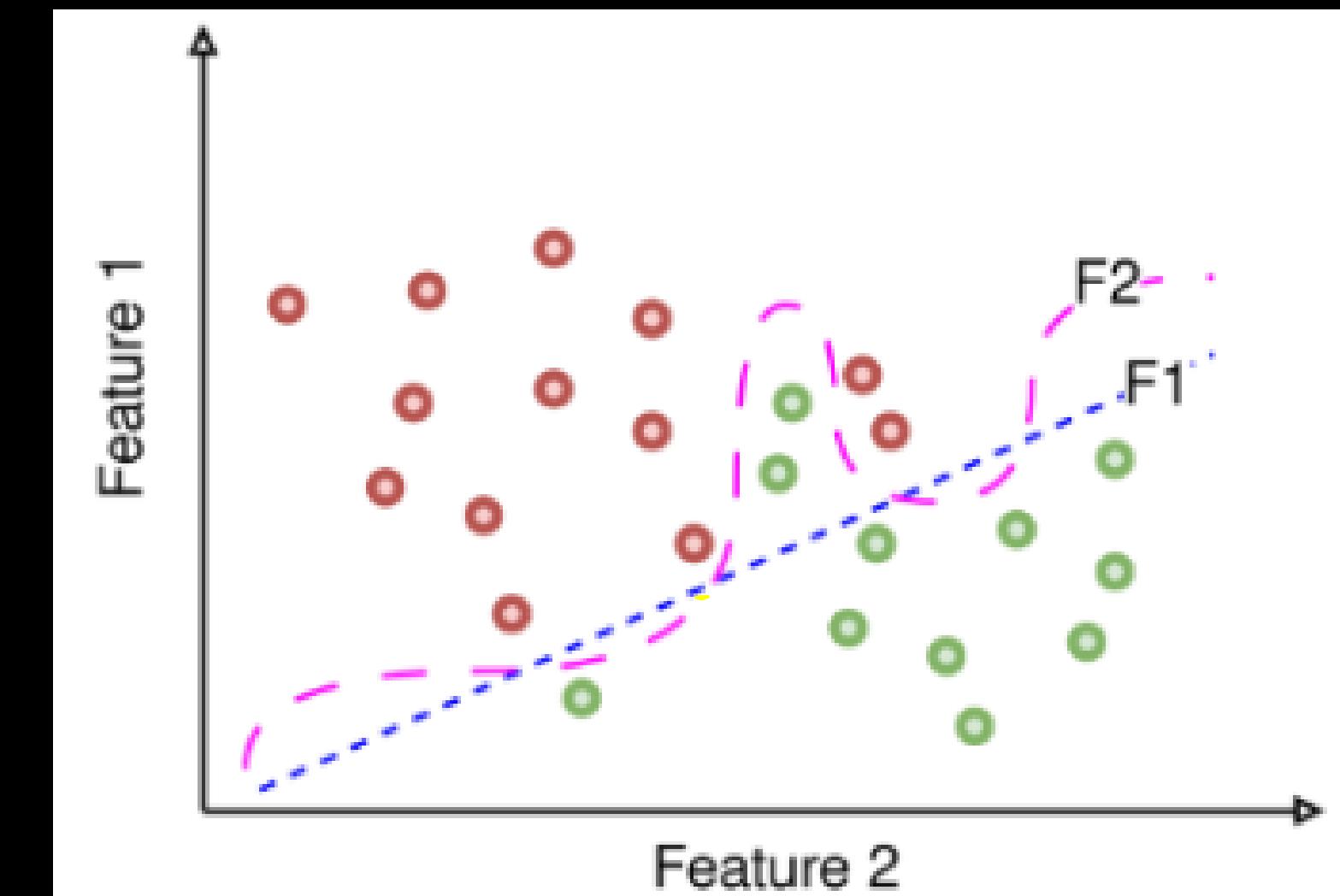
Answer: (C)



Which of the following 2 functions will yield higher training error?

- A. Function F1
- B. Function F2
- C. Both functions F1 & F2 will have the same training error
- D. Cannot be determined

Answer: (A)



Which of the following techniques can help reduce overfitting in a machine learning model?

- a) Increasing the model complexity ✗
- b) Decreasing the amount of training data ✗
- c) Adding more features to the model ✗
- d) Applying regularization techniques ✓

Answer: (d)

$$(w_0 + w_1x_1^2 + w_2x_1 + b + w_3x_2^2 + w_4x_2)$$

After training a model, you observe a significant gap between the training and test performance metrics. Which of the following techniques is most likely to reduce this gap?

- a) Increasing the training dataset size
- b) Decreasing the model complexity
- c) Adding more features to the model
- d) Fine-tuning hyperparameters

Answer: (b),(d)

Suppose, you got a situation where you find that your linear regression model is under fitting the data. In such situation which of the following options would you consider?

- a. You will add more features
- b. You will start introducing higher degree features
- c. You will remove some features
- d. None of the above.

You have generated data from a 3-degree polynomial with some noise. What do you expect of the model that was trained on this data using a 5-degree polynomial as function class?

- a. Low bias, high variance
- b. High bias, low variance.
- c. Low bias, low variance.
- d. High bias, low variance.



# Cross Validation

2.50 pm

## Types of Cross Validation Set

1. Leave-One-Out Cross-Validation (LOOCV)
2. Hold-out cross-validation
3. K-Fold Cross-Validation
4. Stratified K-Fold Cross-Validation
5. Time Series Cross-Validation

## How cross validation helps in bias variance tradeoff?

### **Bias Estimation:**

By performing cross-validation, If the average performance is consistently poor across all folds, it may indicate that the model has high bias.

Cross-validation helps in identifying underfitting by revealing consistent errors in different portions of the dataset.

### **Variance Estimation:**

If there is a significant difference in performance between folds, it suggests that the model is sensitive to the specific data used for training and testing, indicating potential overfitting.

## **Hyperparameter Tuning:**

Cross-validation is commonly used for hyperparameter tuning.

Models with different hyperparameter configurations are trained and evaluated on multiple folds.

This helps in selecting hyperparameter values that balance bias and variance, optimizing the model for better generalization.

Cross validation is a model evaluation method. **Leave-one-out cross validation(LOOCV)** is K-fold cross validation taken to its logical extreme, with K equal to N, the number of data points in the set. That means that N separate times, the function approximator is trained on all the data except for one point and a prediction is made for that point. Thus, it iterates over the other datapoints keeping the rest of the dataset fixed. What can be the major issues in LOOCV?

- a. low variance
- b. high variance
- c. faster run time compared to K-fold cross validation
- d. slower run time compared to K-fold cross validation

**Which of the following cross validation strategies cannot be stratified?**

- a) k-fold cross validation
- b) hold out cross validation
- c) leave one out cross validation
- d) shuffle split cross validation

Which of the following cross validation versions may not be suitable for very large datasets with hundreds of thousands of samples?

- a) k-fold cross-validation
- b) Leave-one-out cross-validation
- c) Holdout method
- d) All of the above

As k increases in k-fold cross-validation method?

- a) The variance of the resulting estimate is reduced as k is increased.
- b) The variance of the resulting estimate is reduced as k is decreased.
- c) None of the above

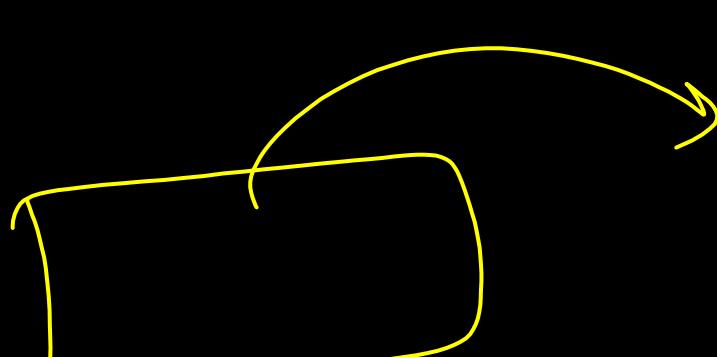
Which of the following is a disadvantage of k-fold cross-validation method?

- a) The variance of the resulting estimate is reduced as k is increased.
- b) This usually does not take longer time to compute
- c) Reduced bias
- d) The training algorithm has to rerun from scratch k times

# KNN

Classification

— PCA  
— t-SNE  
— Place



KNN →  $K=3$

Name	Aptitude	Communication	Class
Karuna	2	5	Speaker
Bhuvna	2	6	Speaker
Gaurav	7	6	Leader
Parul	7	2.5	Intel
Dinesh	8	6	Leader
Jani	4	7	Speaker
Bobby	5	3	Intel
Parimal	3	5.5	Speaker
Govind	8	3	Intel
Susant	6	5.5	Leader
Gouri	6	4	Intel
Bharat	6	7	Leader
Ravi	6	2	Intel
Pradeep	9	7	Leader
Josh	5	4.5	???

$$\sqrt{(5-2)^2 + (4.5-5)^2} =$$

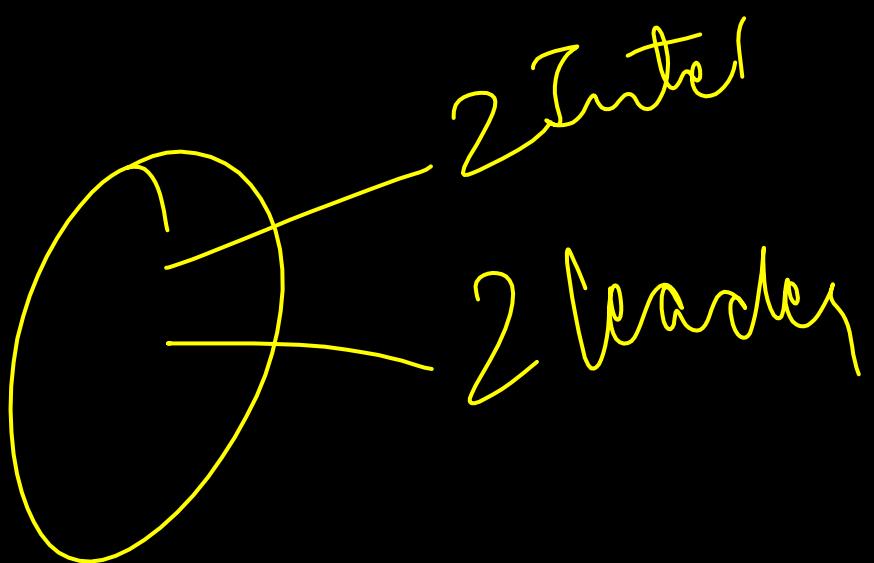
$(x_1, y_1)$  &  $(x_2, y_2)$

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Name	Aptitude	Communication	Class	Distance
Karuna	2	5	Speaker	3.041
Bhuvna	2	6	Speaker	3.354
Parimal	3	5.5	Speaker	2.236
Jani	4	7	Speaker	2.693
Bobby	5	3	Intel	1.500
Ravi	6	2	Intel	2.693
Gouri	6	4	Intel	1.118
Parul	7	2.5	Intel	2.828
Govind	8	3	Intel	3.354
Susant	6	5.5	Leader	1.414
Bharat	6	7	Leader	2.693
Gaurav	7	6	Leader	2.500
Dinesh	8	6	Leader	3.354
Pradeep	9	7	Leader	4.717
Josh	5	4.5	???	→ Intel

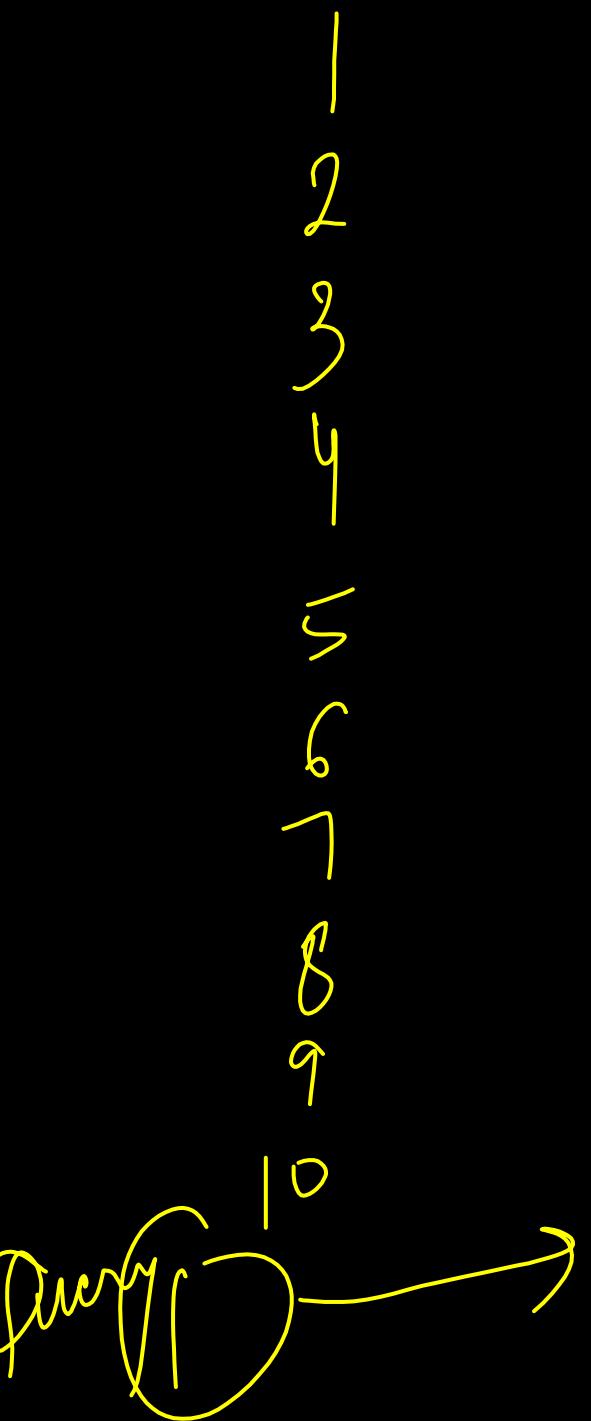
2 Intel & 1 Leader

# Why is the odd value of 'k' preferred over an even value in the k-NN algorithm?



K-Nearest Neighbor is a \_\_\_\_\_, \_\_\_\_\_ algorithm

- a. Non-parametric, eager
- b. Parametric, eager
- c. Non-parametric, lazy
- d. Parametric, lazy



State whether the statement is True/False:

k-NN algorithm does more computation on test time rather than train time.

1. True
2. False

Consider a dataset with the following three data points in a two-dimensional space:

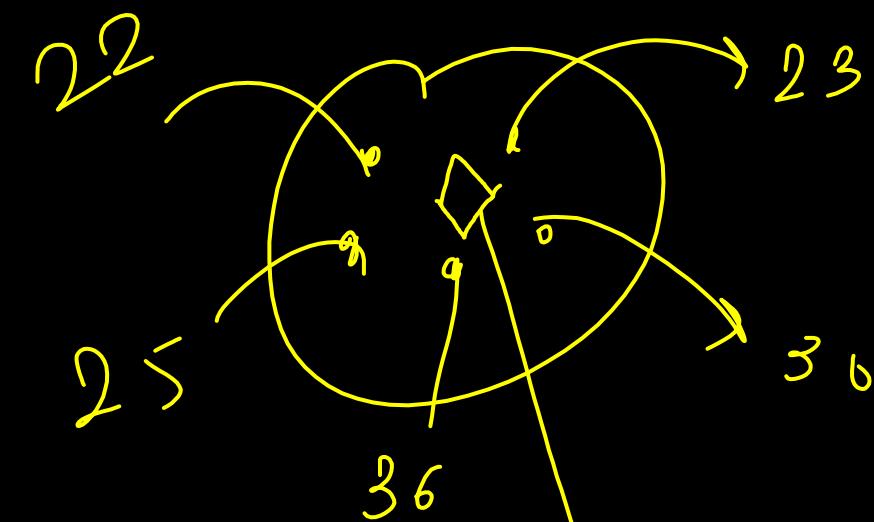
- Data point A: (2, 3), Class: 1
- Data point B: (4, 6), Class: 1
- Data point C: (5, 2), Class: 2

If K=2, what is the predicted class of a new data point at coordinates (3, 4) using the Euclidean distance metric?

$$\left. \begin{aligned} (3, 4) \& A = \sqrt{(3-2)^2 + (4-3)^2} = \sqrt{2} \\ (3, 4) \& B = \sqrt{(3-4)^2 + (4-6)^2} = \sqrt{5} \\ (3, 4) \& C = \sqrt{(3-5)^2 + (4-2)^2} = \sqrt{8} \end{aligned} \right\}$$

# Can k-NN algorithm be used for a regression problem?

Yes



$$\frac{22 + 23 + 25 + 36 + 30}{5}$$

5

What would be the relationship between the training time taken by 1- NN, 2-NN, and 3-NN?

1. 1-NN > 2-NN > 3-NN
2. 1-NN < 2-NN < 3-NN
3. 1-NN ~ 2-NN ~ 3-NN ✓
4. None of these

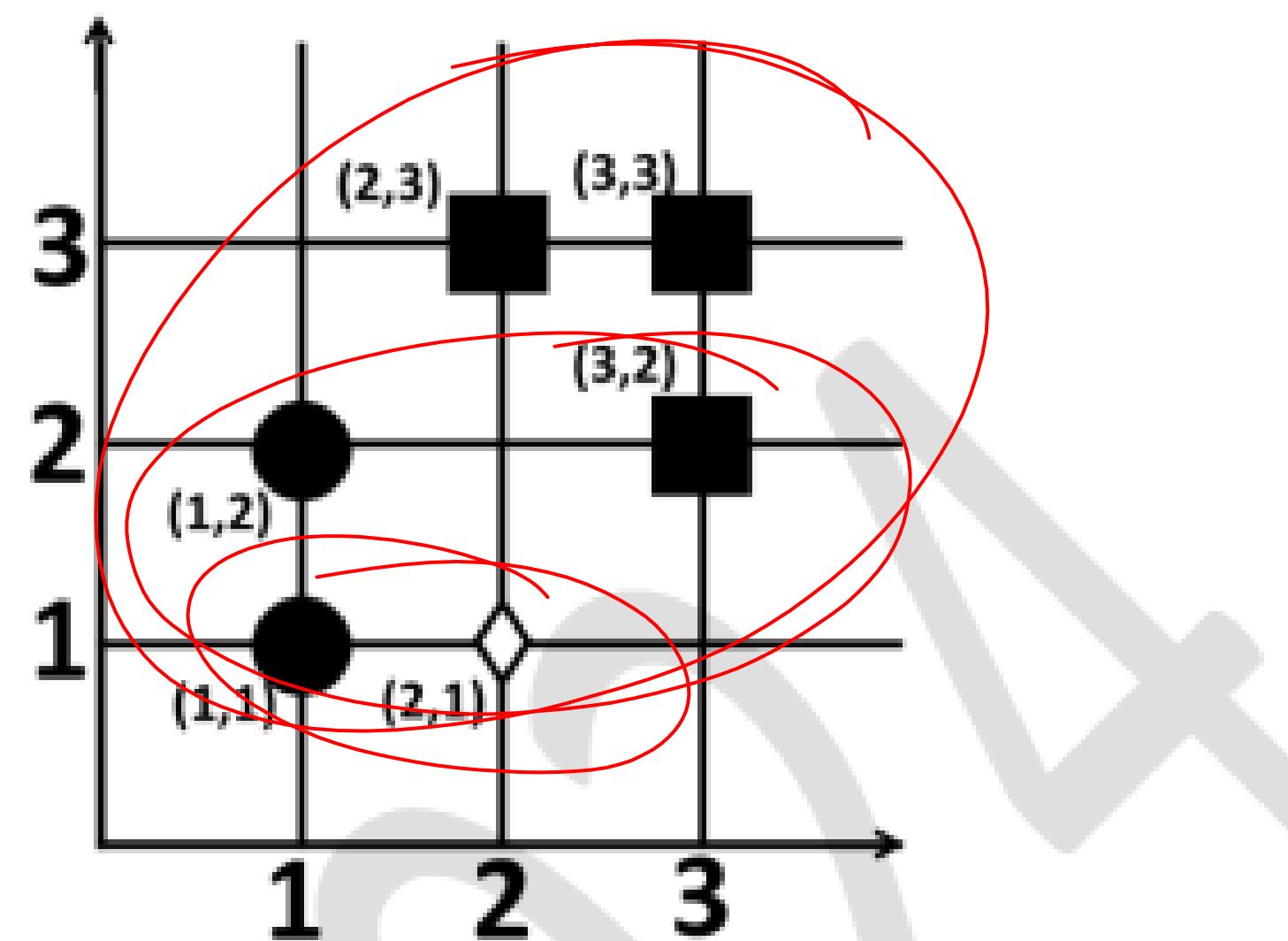
Given the two-dimensional dataset consisting of 5 data points from two classes (circles and squares) and assume that the Euclidean distance is used to measure the distance between two points. The minimum odd value of  $k$  in  $k$ -nearest neighbor algorithm for which the diamond ( $\diamond$ ) shaped data point is assigned the label square is \_\_\_\_\_.

$$K=1$$

$$K=3$$

$$K=5$$

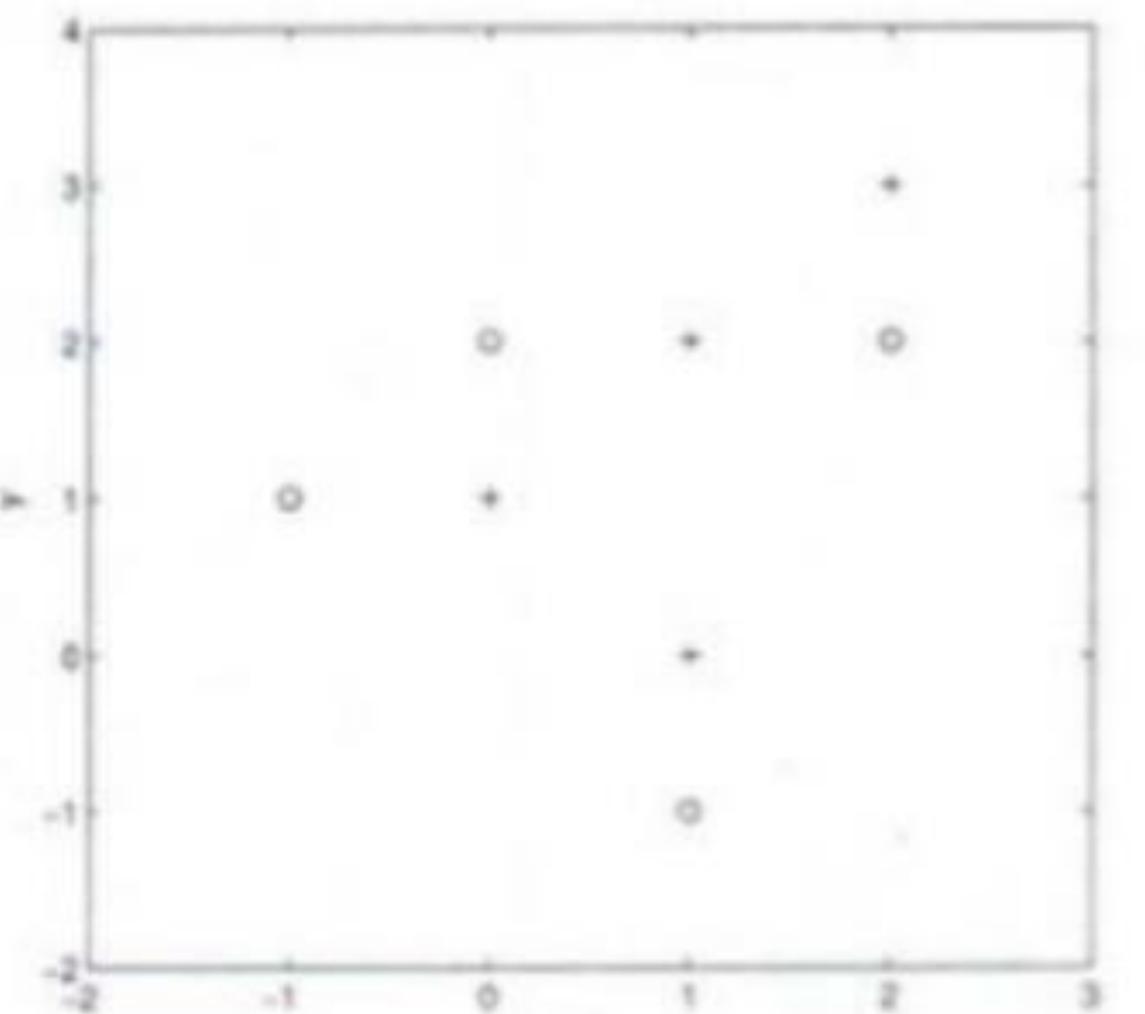
$\diamond$   
 $\diamond$   
□



Suppose, you have given the following data where x and y are the 2 input variables and Class is the dependent variable.

x	y	Class
-1	1	-
0	1	+
0	2	-
1	-1	-
1	0	+
1	2	+
2	2	-
2	3	+

Below is a scatter plot which shows the above data in 2D space.



Suppose, you want to predict the class of new data point  $x=1$  and  $y=1$  using Euclidian distance in 3-NN. In which class this data point belong to?

- a. + Class
- b. - Class
- c. Can't Say
- d. None of these

Consider a set of five training examples given as  $((x_i, y_i), c_i)$  values, where  $x_i$  and  $y_i$  are the two attribute values (positive integers) and  $c_i$  is the binary class label:

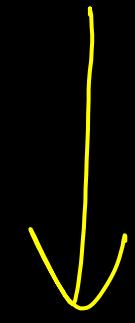
$$\left\{ \begin{array}{l} ((1,1),-1) \rightarrow 2+5 = 7 \\ ((1,7),+1) \rightarrow 2+1 = 3 \\ ((3,3),+1) \rightarrow 0+3 = 3 \\ ((5,4),-1) \rightarrow 2+2 = 4 \\ ((2,5),-1) \rightarrow 1+1 = 2 \end{array} \right.$$

$$\begin{aligned} & (x_1, y_1) \neq (x_2, y_2) \\ & |y_2 - y_1| + |x_2 - x_1| \end{aligned}$$

Classify a test example at coordinates  $(3,6)$  using a k-NN classifier with  $k=3$  and Manhattan distance defined by  $d((u,v),(p,q)) = |u-p| + |v-q|$

+  
→

Large value of  $K$



underfitting



high bias

Small value of  $K$



overfitting



high variance & low bias

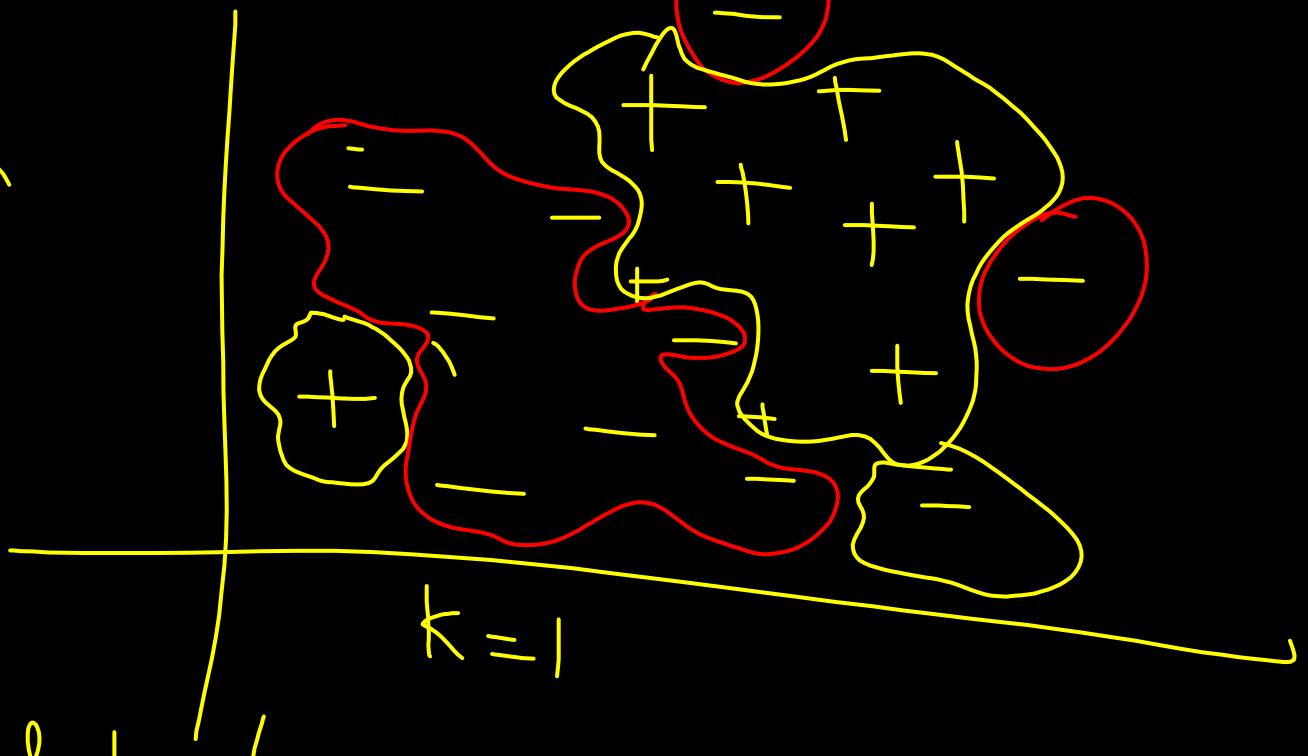
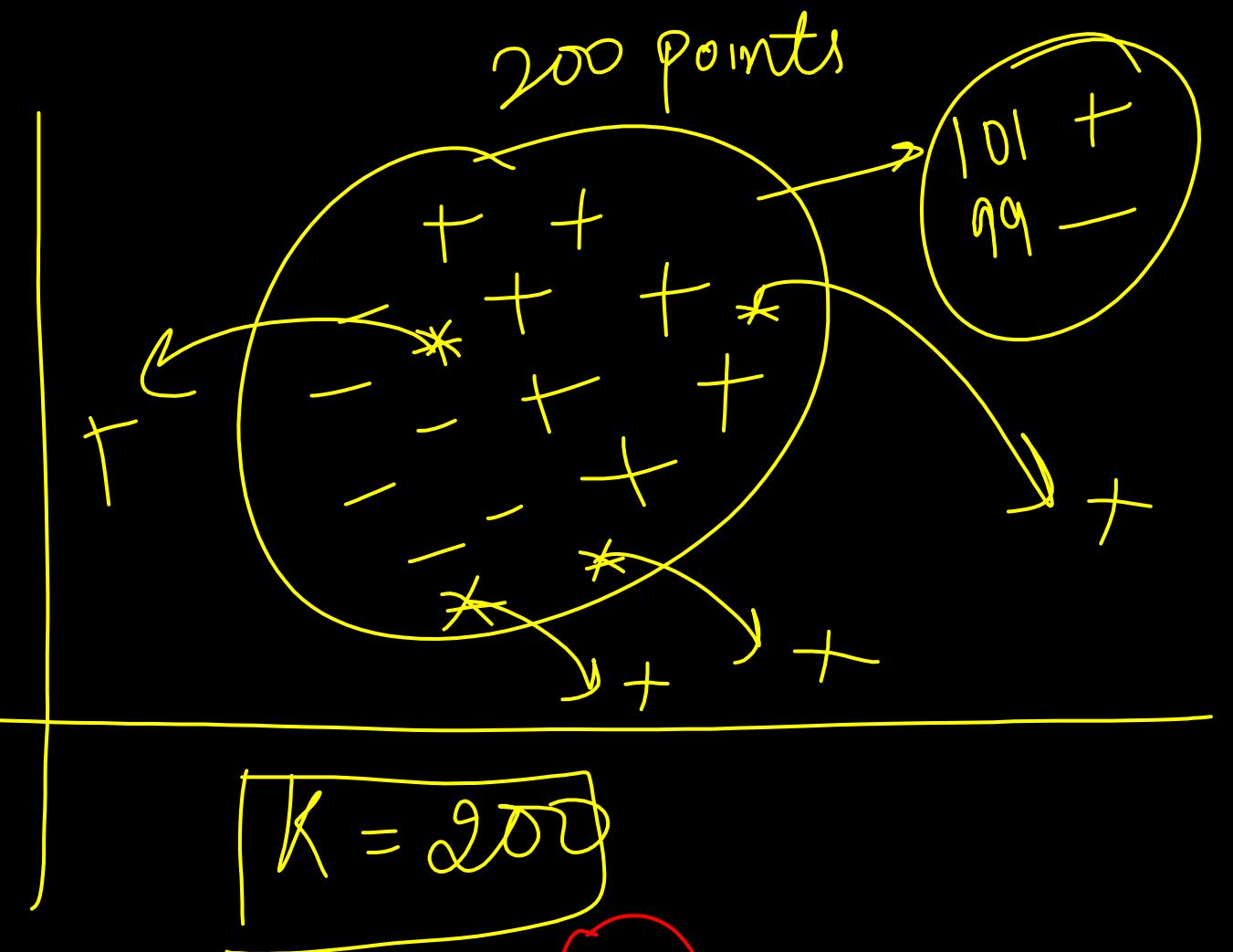
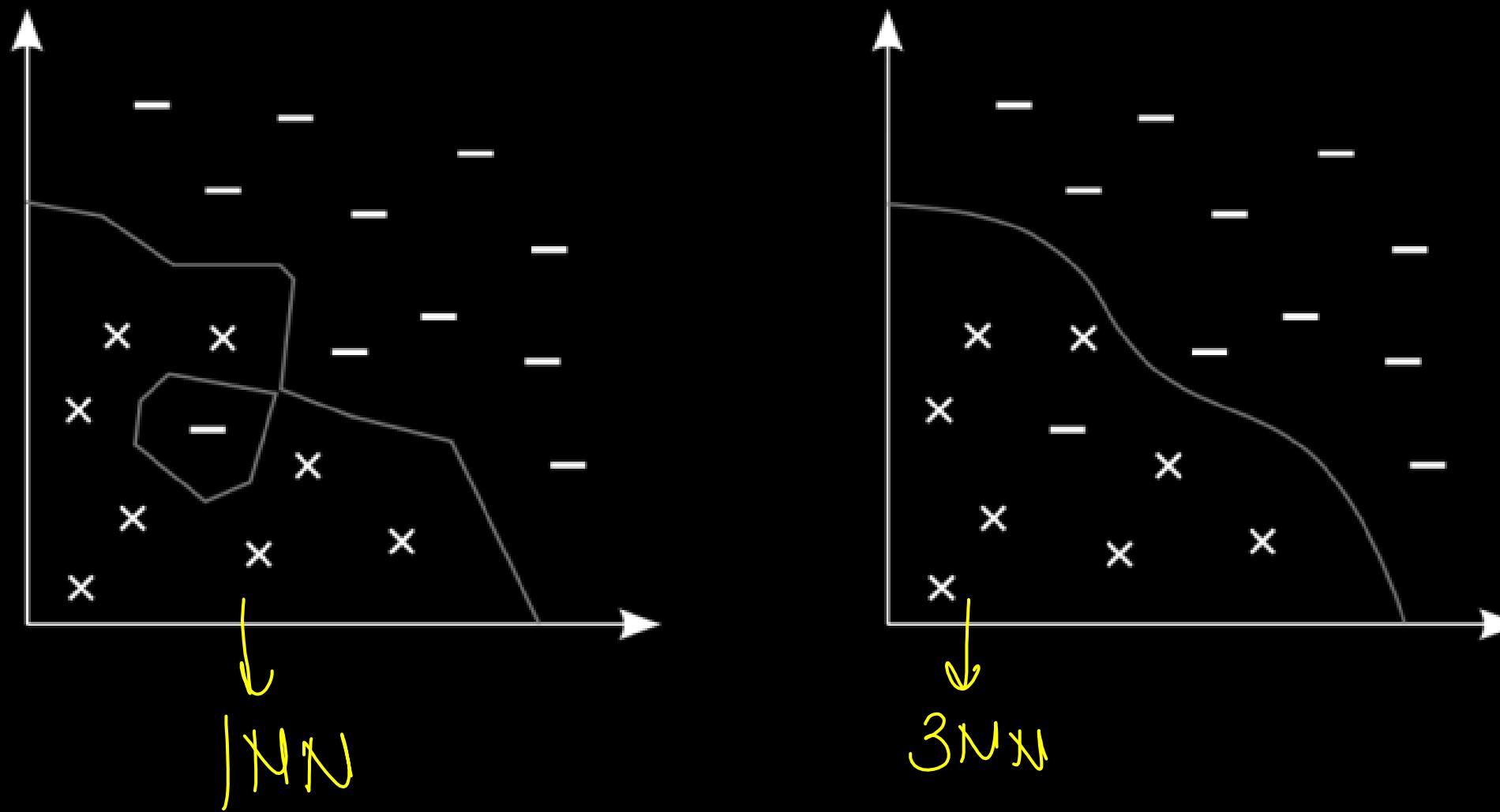


Figure illustrates decision boundaries for two nearest – neighbor classifiers. Determine which one of the boundaries belongs to the 1-NN and which one belong to 3-NN?



You have been given the following 2 statements. Find out which of these options is/are true in case of k-NN?

1. In case of very large value of k, we may include points from other classes into the neighborhood.
  2. In case of too small value of k, the algorithm is very sensitive to noise.
- 
- a. 1 is True and 2 is False
  - b. 1 is False and 2 is True
  - c. Both are True
  - d. Both are False

What is the effect of increasing the value of K in KNN on the bias and variance of the model?

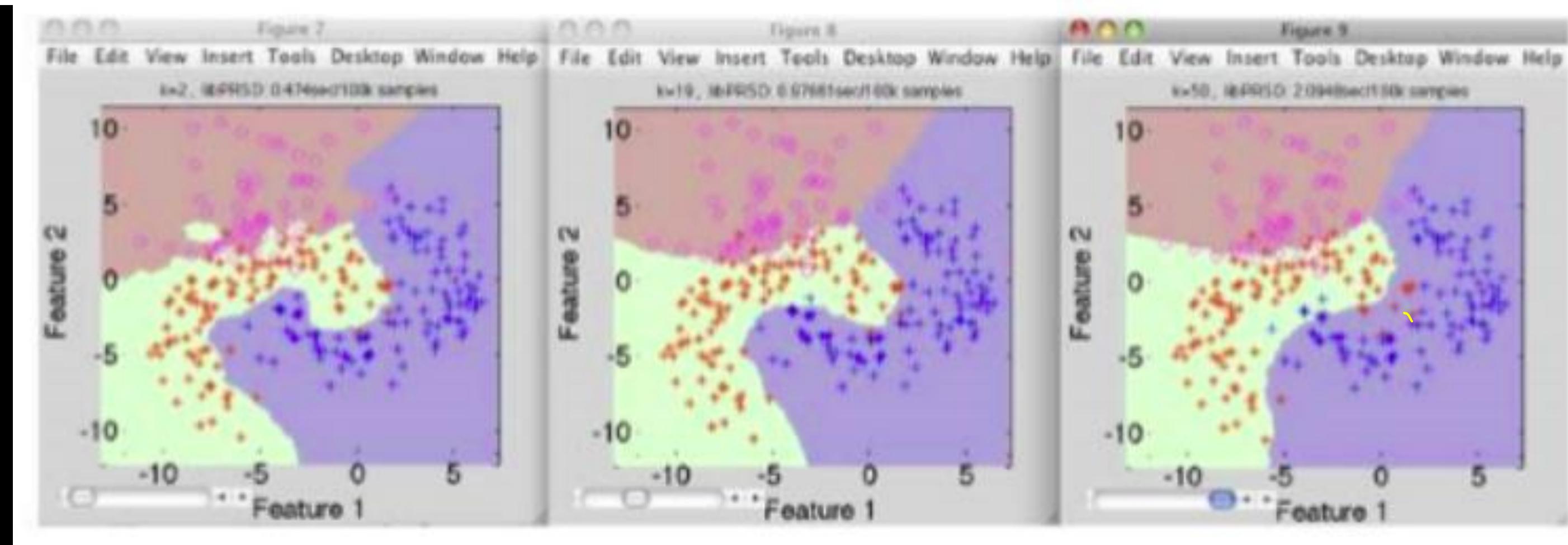
- a) Increasing K increases bias and decreases variance.
- b) Increasing K decreases bias and increases variance.
- c) Increasing K has no effect on bias or variance.
- d) Increasing K decreases both bias and variance.

Overfitting → underfitting

When you find many noises in data, which of the following options would you consider in kNN?

- 1. Increase the value of k
- 2. Decrease the value of k
- 3. Noise does not depend on k
- 4. K = 0

Suppose you are given the following images (1 represents the left image, 2 represents the middle and 3 represents the right). Now your task is to find out the value of k in k-NN in each of the images shown below. Here  $k_1$  is for 1<sup>st</sup>,  $k_2$  is for 2<sup>nd</sup> and  $k_3$  is for 3rd figure.



- a.  $k_1 > k_2 > k_3$
- b.  $k_1 < k_2 > k_3$
- c.  $k_1 < k_2 < k_3$
- d. None of these

Given the following dataset, for  $k = 3$ , use KNN regression to find the prediction for a new data-point  $(2,3)$  (*Use Euclidean distance measure for finding closest points*)

X1	X2	Y
2	5	3.4
5	5	5
3	3	3
6	3	4.5
2	2	2
4	1	2.8

2      3

$$\sqrt{0+2^2} = \sqrt{4} = 2$$

$$\sqrt{3+2^2} = \sqrt{13}$$

$$\sqrt{1}$$

$$\sqrt{4}$$

$$\sqrt{1}$$

$$\sqrt{8}$$

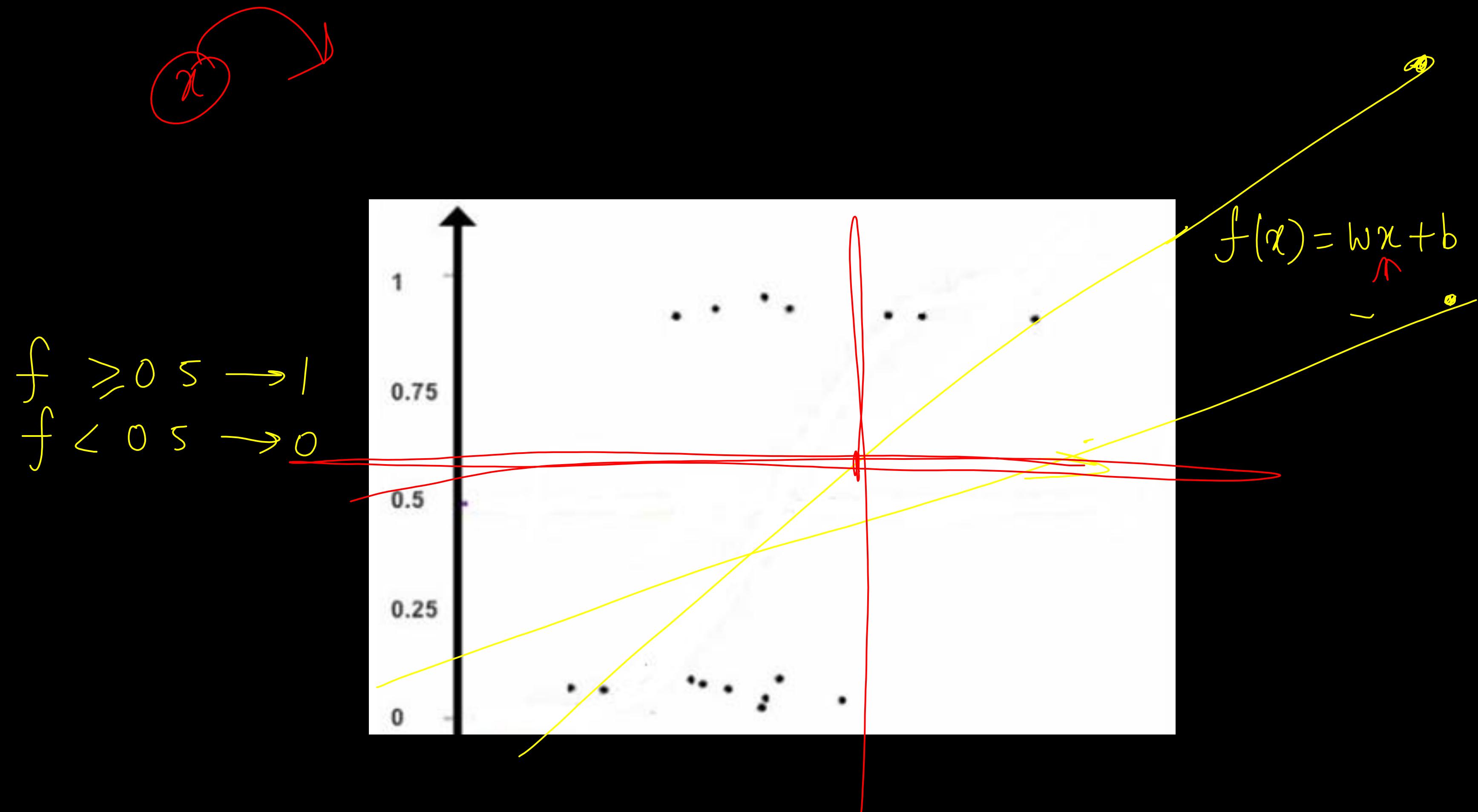
- A. 2.0
- B. 2.6
- C. 2.8
- D. 3.2

$$\frac{3.4 + 3 + 2}{3} = \frac{8.4}{3} = 2.8$$

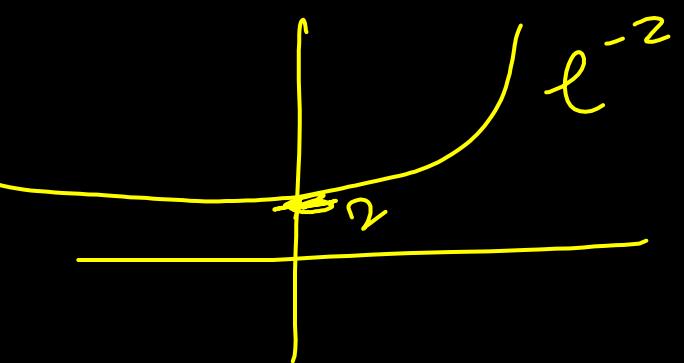
Answer: (C)

# Classification

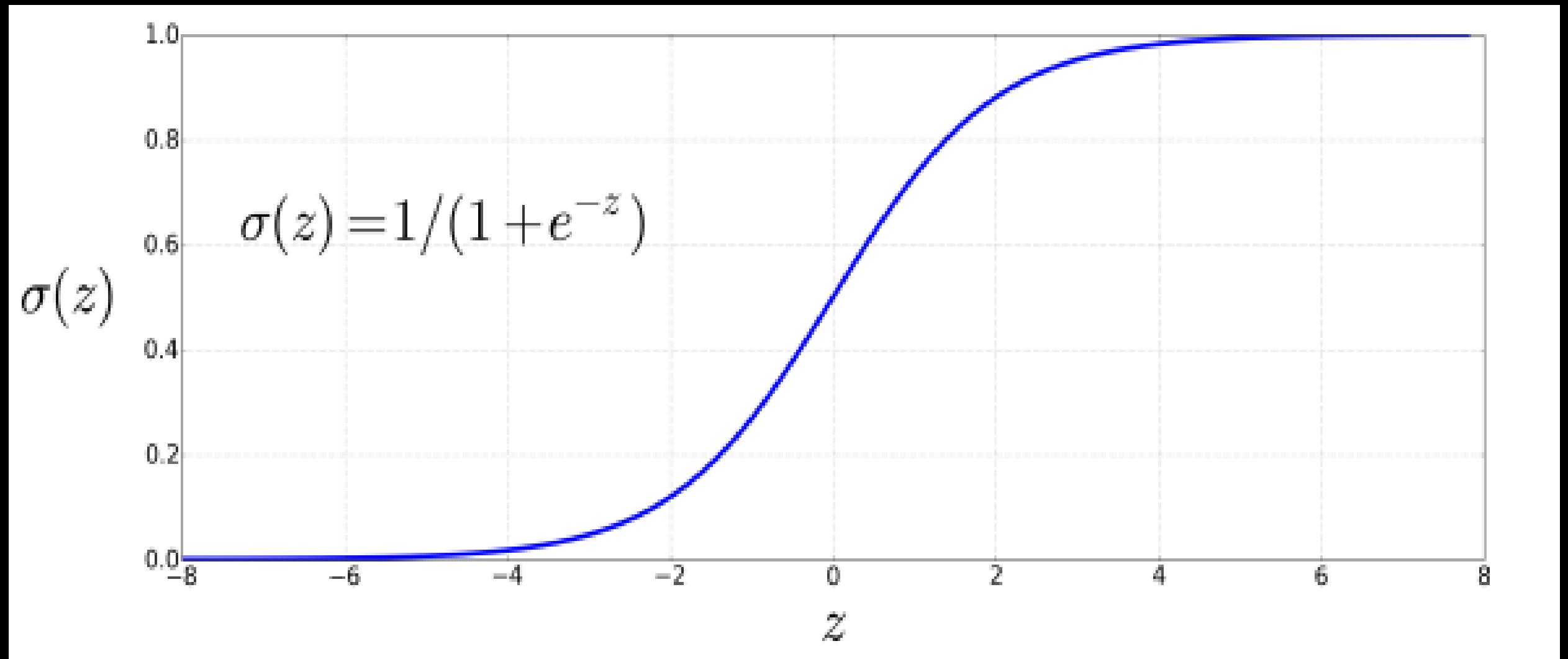
# Logistic Regression



## Sigmoid Function



$$z = w\chi + b$$



$$\begin{aligned}\sigma(z) &= \frac{1}{1+e^{-z}} \\ &= \frac{1}{1+e^{-(wz+b)}}\end{aligned}$$

$$\begin{aligned}\sigma(z) &= P(y=1) \\ 1 - \sigma(z) &= P(y=0)\end{aligned}$$

$$w=2, b=1$$

$$\begin{aligned}\sigma(z) &= \frac{1}{1+e^{-2z+1}} \\ &= \frac{1}{1+e^{-1}}\end{aligned}$$

$0 < \sigma(z) < 1$  (Output between 0 and 1)

$$\begin{aligned}
 P(y=1) &= \sigma(\mathbf{w} \cdot \mathbf{x} + b) \\
 &= \frac{1}{1 + \exp(-(\mathbf{w} \cdot \mathbf{x} + b))}
 \end{aligned}$$

$$\begin{aligned}
 P(y=0) &= 1 - \sigma(\mathbf{w} \cdot \mathbf{x} + b) \\
 &= 1 - \frac{1}{1 + \exp(-(\mathbf{w} \cdot \mathbf{x} + b))} \\
 &= \frac{\exp(-(\mathbf{w} \cdot \mathbf{x} + b))}{1 + \exp(-(\mathbf{w} \cdot \mathbf{x} + b))}
 \end{aligned}$$

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

$$\sigma'(z) = \frac{e^{-z}}{(1+e^{-z})^2} = \sigma(z)[1-\sigma(z)]$$

Goal  $\mathbf{w} \cdot \mathbf{x} + b$   
 $\mathbf{x}$        $y - y'$   
 $|$        $|$   
 $0$        $0$

## Decision boundary

$$\text{decision}(x) = \begin{cases} 1 & \text{if } P(y=1|x) > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

## Loss Function

$$\hat{y} = \frac{1}{1 + e^{-z}}$$

$$L(\hat{y}, y) = -[y \log \hat{y} + (1-y) \log(1-\hat{y})]$$

$$\hat{y} = \frac{1}{1 + e^{-z}}.$$

~~$-y \log \hat{y}$~~

$$L(\hat{y}, y) = -[y \log \sigma(w \cdot x + b) + (1-y) \log(1 - \sigma(w \cdot x + b))]$$

$$y=1 \quad , \quad y=0$$

$$L(\hat{y}, y) = -\log \hat{y} \quad J(\hat{y}, y) = -\log(1-\hat{y})$$

$$J = \frac{1}{m} \sum L(y, \hat{y})$$

l

$$z = w\alpha + b$$

## Gradient Descent Implementation

$$L(\hat{y}, y) = -[y \log \sigma(w \cdot x + b) + (1 - y) \log (1 - \sigma(w \cdot x + b))]$$

$$w = w - \alpha \frac{\partial L}{\partial w}$$

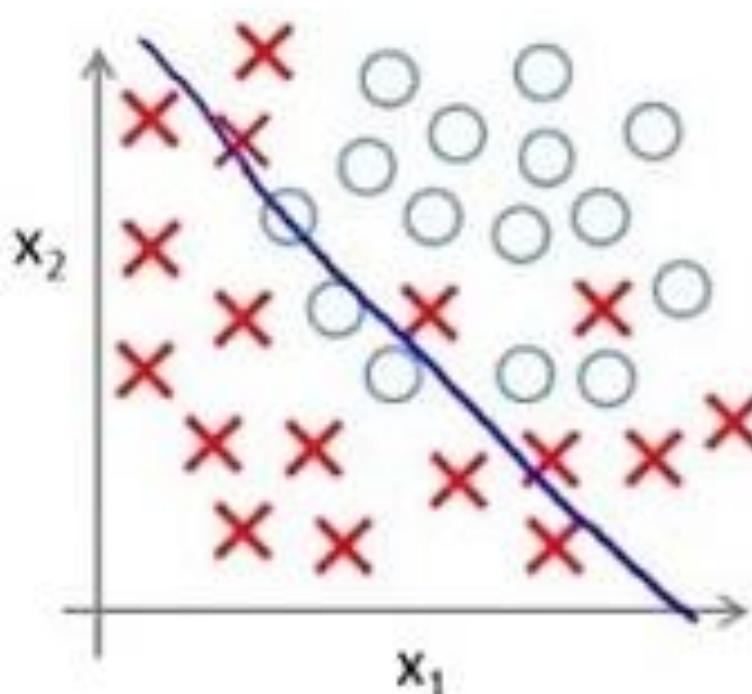
$$b = b - \alpha \frac{\partial L}{\partial b}$$

$$\begin{aligned} \frac{\partial L}{\partial w} &= - \left[ y \frac{1}{\sigma(z)} \times \cancel{\sigma(z)} (1 - \sigma(z)) \alpha + (1 - y) \frac{1}{1 - \cancel{\sigma(z)}} (-\sigma(z)) (1 - \cancel{\sigma(z)}) \alpha \right] \\ &= - \left[ \cancel{y} - \cancel{y} \cancel{\sigma(z)} - \alpha \sigma(z) + \cancel{y} \cancel{\sigma(z)} \right] \\ &\equiv -\alpha (y - \sigma(z)) \\ &= (y - \sigma(z)) \alpha - (y - \hat{y}) \alpha \\ \frac{\partial L}{\partial b} &= (y - \hat{y}) \end{aligned}$$

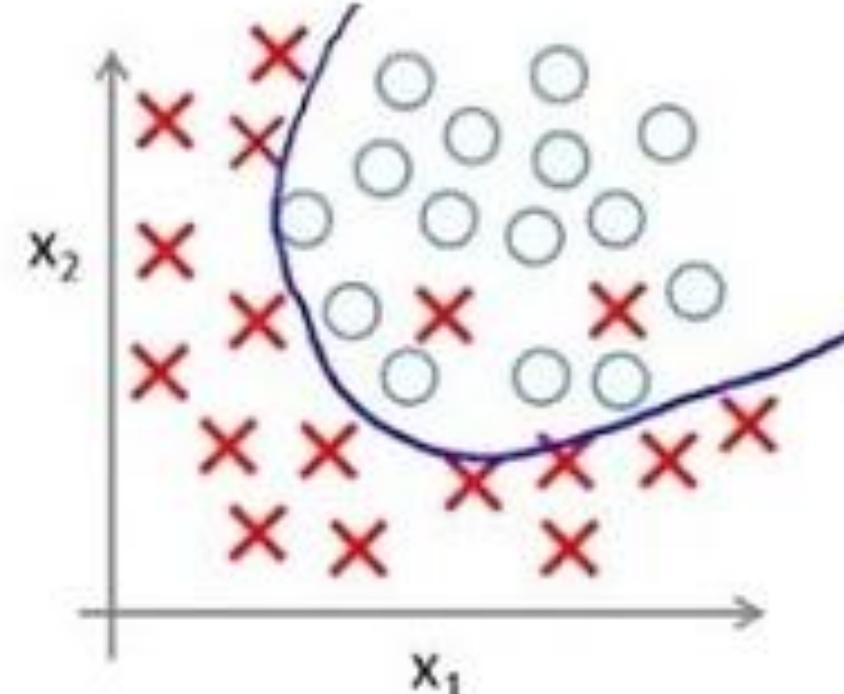
If  $g(z)$  is the sigmoid function, then its derivative with respect to  $z$  may be written in term of  $g(z)$  as

- a.  $g(z)(g(z)-1)$
- b.  $g(z)(1+g(z))$
- c.  $-g(z)(1+g(z))$
- d.  $g(z)(1-g(z))$

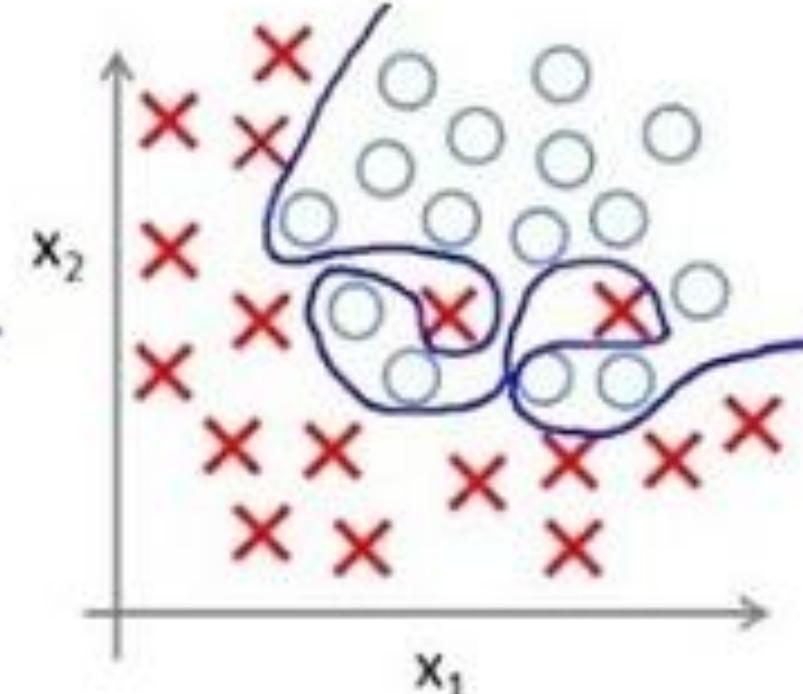
Below are the three scatter plot(A,B,C left to right) and hand drawn decision boundaries for logistic regression.



A



B

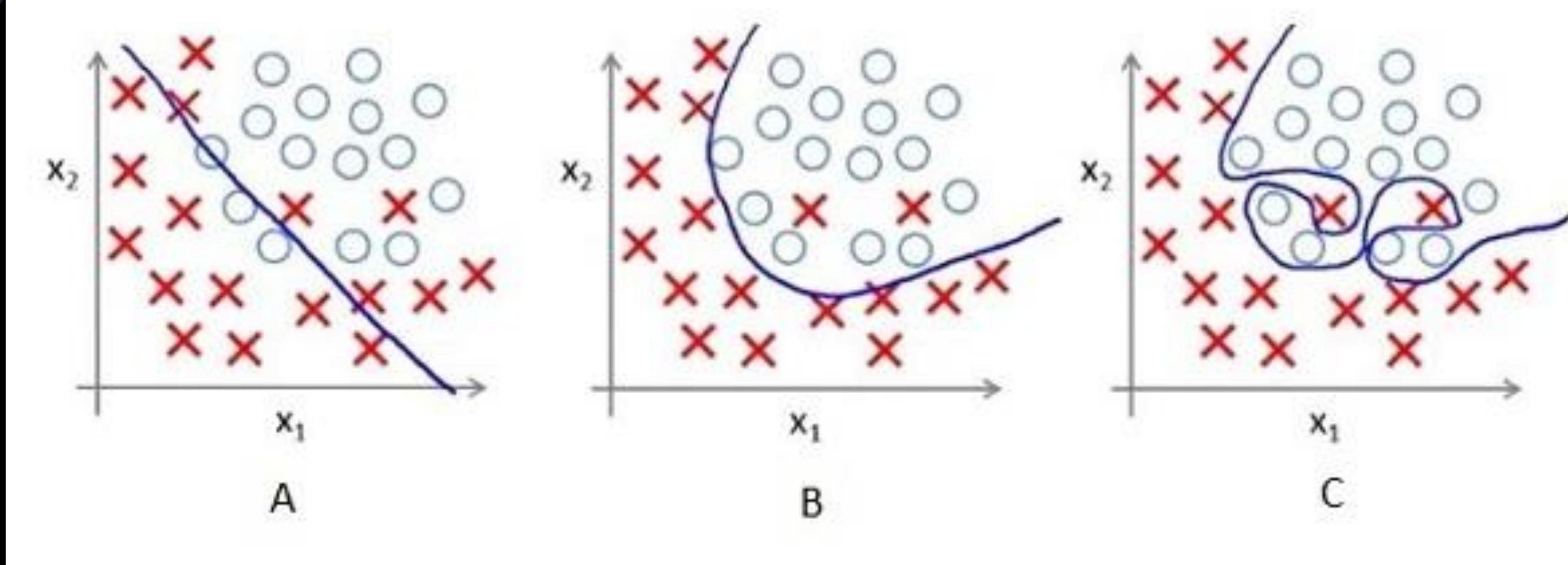


C

**Which of the following above figure shows that the decision boundary is overfitting the training data?**

- A) A
- B) B
- C) C
- D) None of these

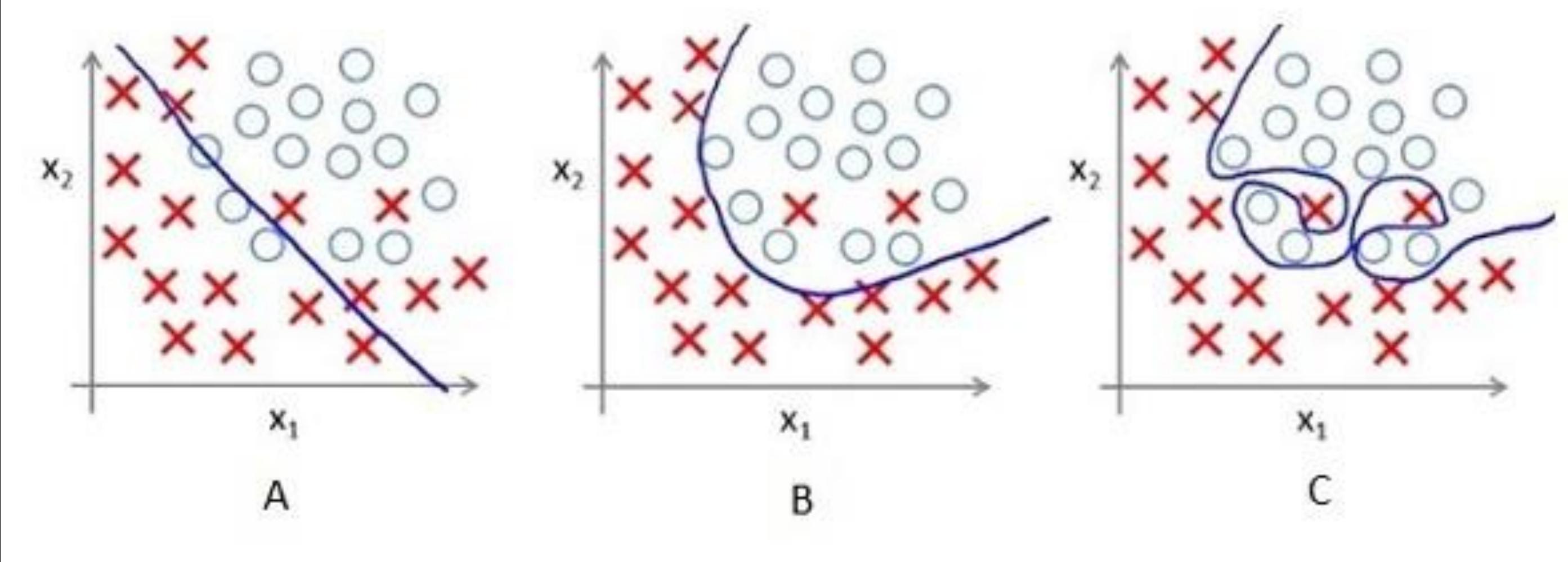
Answer: (C)



**What do you conclude after seeing this visualization?**

1. The training error in first plot is maximum as compare to second and third plot.
2. The best model for this regression problem is the last (third) plot because it has minimum training error (zero).
3. The second model is more robust than first and third because it will perform best on unseen data.
4. The third model is overfitting more as compare to first and second.
5. All will perform same because we have not seen the testing data.

Answer: (1,3 &4)



**Suppose, above decision boundaries were generated for the different value of regularization. Which of the above decision boundary shows the maximum regularization?**

- A) A
- B) B
- C) C
- D) All have equal regularization

Answer: (A)

Which of the following are symptoms of a logistic regression model being overfit? Select all that apply

- (a) Large estimated coefficients
- (b) Good generalization to unseen data
- (c) Simple decision boundary
- (d) Complex decision boundary

Answer: (a), (d)

## How will the bias change on using high(infinite) regularization?

- A) Bias will be high
- B) Bias will be low
- C) Can't say
- D) None of these

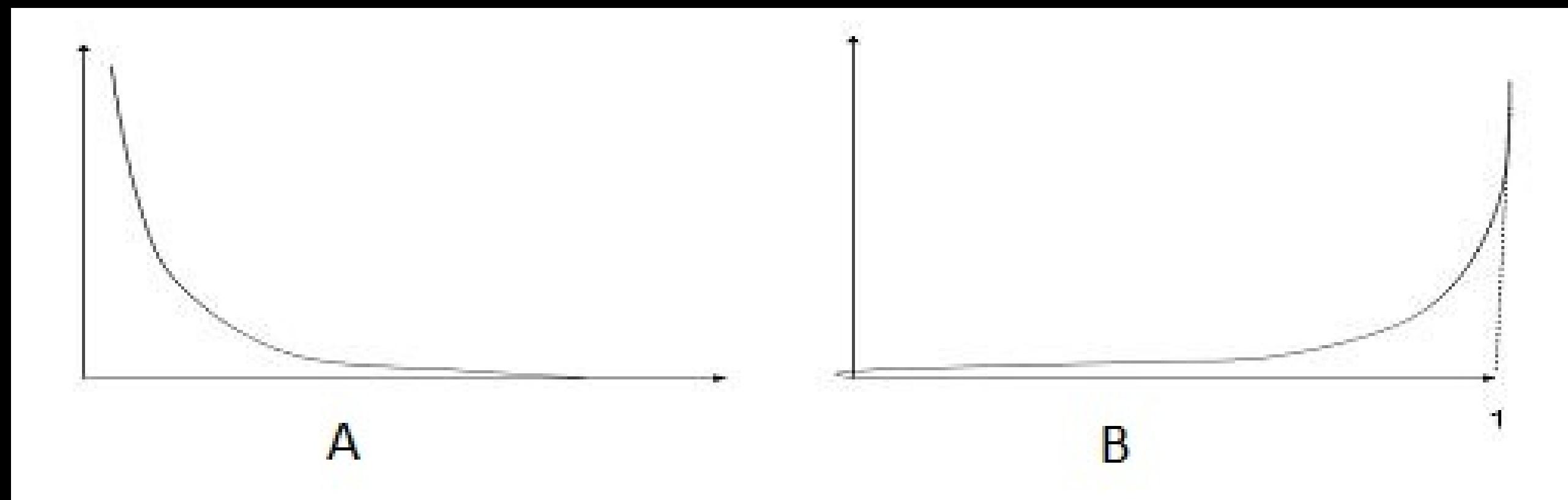
Answer: (A)

Suppose you are using a Logistic Regression model on a huge dataset. One of the problems you may face on such huge data is that Logistic regression will take very long time to train.

- A) Decrease the learning rate and decrease the number of iteration
- B) Decrease the learning rate and increase the number of iteration
- C) Increase the learning rate and increase the number of iteration
- D) Increase the learning rate and decrease the number of iteration

**Solution:** D

Which of the following image is showing the cost function for  $y = 1$ . ~



- A) A
- B) B
- C) Both
- D) None of these

logistic regression

$J(\hat{y}) = -\log \hat{y}$

A hand-drawn yellow graph of the logistic (sigmoid) function  $y = 1 / (1 + e^{-x})$ . The curve is S-shaped, passing through the point (0, 0.5). It approaches y=1 as x goes to positive infinity and y=0 as x goes to negative infinity. A vertical dashed line is drawn at x=0, and a horizontal dashed line is drawn at y=0.5.

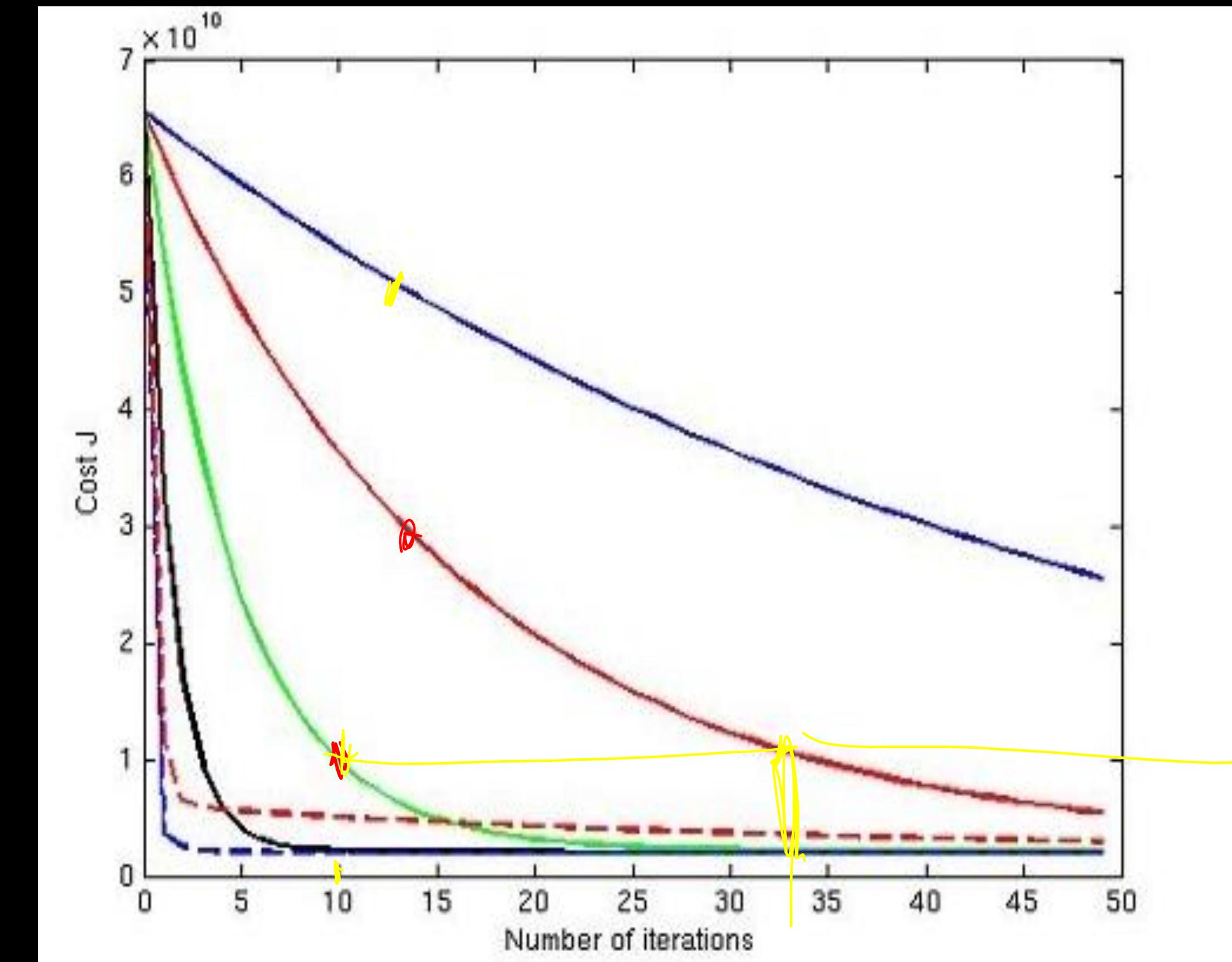
Solution: A

Imagine, you have given the below graph of logistic regression which shows the relationships between cost function and number of iteration for 3 different learning rate values (different colors are showing different curves at different learning rates).

1. The learning rate for blue is  $\alpha_1$
2. The learning rate for red is  $\alpha_2$
3. The learning rate for green is  $\alpha_3$

- A)  $\alpha_1 > \alpha_2 > \alpha_3$
- B)  $\alpha_1 = \alpha_2 = \alpha_3$
- C)  $\alpha_1 < \alpha_2 < \alpha_3$
- D) None of these

**Solution: C**



# Naïve Bayes

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

$$P(C_1|X) = \frac{P(X|C_1) P(C_1)}{P(X)}$$

$$\asymp P(X|C_1) P(C_1) = P(x_1, x_2, x_3|C_1) P(C_1)$$

$$P(C_2|X) = \frac{P(X|C_2) P(C_2)}{P(X)}$$

$$\asymp P(X|C_2) P(C_2) = P(x_1|C_2) P(x_2|C_2) P(x_3|C_2) P(C_2)$$

$$P(C_3|X) = \frac{P(X|C_3) P(C_3)}{P(X)}$$

$$\asymp P(X|C_3) P(C_3) = P(x_1|C_3) P(x_2|C_3) P(x_3|C_3) P(C_3)$$

$C_1 = \text{FISH}$ ,  $C_2 = \text{Animal}$ ,  $C_3 = \text{Bird}$

Sl. No.	Swim	Fly	Crawl	Class
1	Fast	No	No	Fish
2	Fast	No	Yes	Animal
3	Slow	No	No	Animal
4	Fast	No	No	Animal
5	No	Short	No	Bird
6	No	Short	No	Bird
7	No	Rarely	No	Animal
8	Slow	No	Yes	Animal
9	Slow	No	No	Fish
10	Slow	No	Yes	Fish
11	No	Long	No	Bird
12	Fast	No	No	Bird

$$P(C_1|X) = P(x_1|C_1) P(x_2|C_1) P(x_3|C_1) P(C_1)$$

$$= \frac{2}{3} \times \frac{0}{3} \times \frac{2}{3} \times \frac{3}{12} = 0$$

$$P(C_2|X) = P(x_1|C_2) P(x_2|C_2) P(x_3|C_2) P(C_2)$$

$$= \frac{1}{5} \times \frac{1}{5} \times \frac{3}{5} \times \frac{1}{12} = \frac{1}{50}$$

$$P(C_3|X) = P(x_1|C_3) P(x_2|C_3) P(x_3|C_3) P(C_3)$$

$$= \frac{0}{4} \times \frac{0}{4} \times \frac{4}{4} \times \frac{4}{12} = 0$$

$$P(C_1) = \frac{3}{12}, P(C_2) = \frac{5}{12}, P(C_3) = \frac{4}{12}$$

X

The test instance is  $(\text{Slow}, \text{Rarely}, \text{No}) \rightarrow C_2 \rightarrow \text{Animal}$

$$C_1 = \text{Yes}, C_2 = \text{No} \quad P(C_1) = \frac{5}{10}, P(C_2) = \frac{5}{10}$$



Use naive Bayes algorithm to determine whether a red domestic SUV car is a stolen car or not using the following data:

Example no.	Colour	Type	Origin	Whether stolen
1	red	sports	domestic	yes
2	red	sports	domestic	no
3	red	sports	domestic	yes
4	yellow	sports	domestic	no
5	yellow	sports	imported	yes
6	yellow	SUV	imported	no
7	yellow	SUV	imported	yes
8	yellow	SUV	domestic	no
9	red	SUV	imported	no
10	red	sports	imported	yes

$$P(C_1|X) = P(x_1|C_1) P(x_2|C_1) P(x_3|C_1) P(C_1)$$

$$= \frac{3}{5} \times \frac{1}{5} \times \frac{2}{5} \times \frac{1}{2}$$

$$= \frac{3}{125}$$

$$P(C_2|X) = P(x_1|C_2) P(x_2|C_2) P(x_3|C_2) P(C_2)$$

$$= \frac{2}{5} \times \frac{3}{5} \times \frac{3}{5} \times \frac{1}{2}$$

$$= \frac{9}{125}$$

$\times (Red \quad SUV \quad domestic)$  ??  $\rightarrow C_2 \rightarrow \text{Not}$

Consider the following dataset showing the result whether a person has passed or failed the exam based on various factors. Suppose the factors are independent to each other. We want to classify a new instance with Confident=Yes, Studied=Yes, and Sick=No.

Confident	Studied	Sick	Result
Yes	No	No	Fail
Yes	No	Yes	Pass
No	Yes	Yes	Fail
No	Yes	No	Pass
Yes	Yes	Yes	Pass

Yes Yes No →

(True/False) Naïve Bayes algorithm can be used when some input features are in continuous range?

Salary

75,000

50,000

1 lac

Salary

0 - 50

50 - 600

100 - 150

What is the naive assumption in a Naive Bayes Classifier?

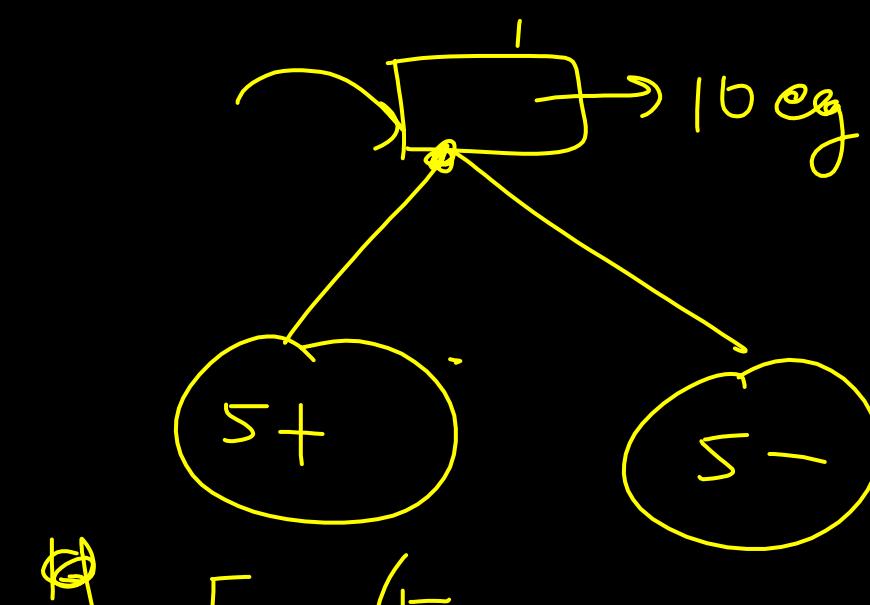
- a. All the classes are independent of each other
- b. All the features of a class are independent of each other 
- c. The most probable feature for a class is the most important feature to be considered for classification 
- d. All the features of a class are conditionally dependent on each other. 

# Decision Tree

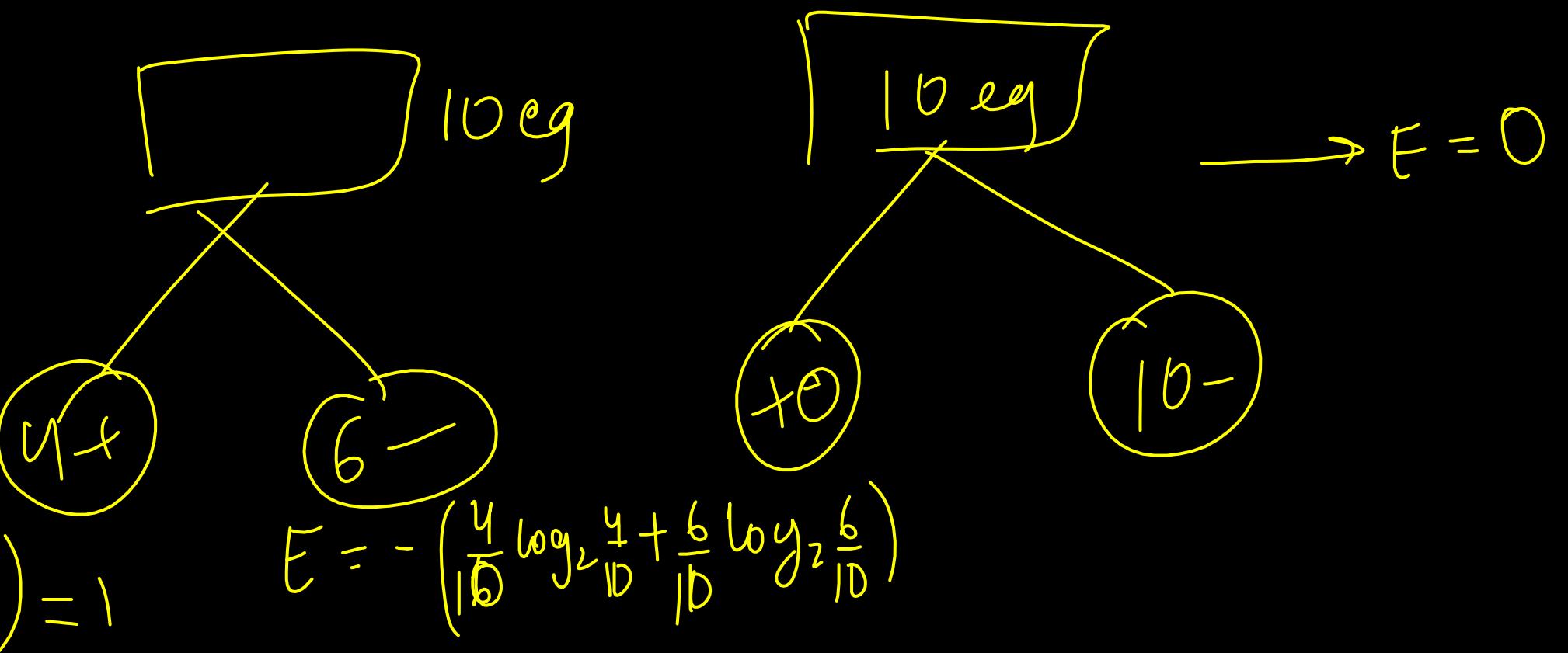
Salary	Credit score	Loan Approved/NA
H	1	A
M	2	NA
L	3	.
	4	.
	5	.
	6	.
	7	.
	8	.
	9	.
	10	10 eq

Entropy

$$= - \sum p_i \log_2 p_i$$



$$E = - \left( \frac{5}{10} \log_2 \frac{5}{10} + \frac{5}{10} \log_2 \frac{5}{10} \right) = 1$$



$$E = - \left( \frac{4}{10} \log_2 \frac{4}{10} + \frac{6}{10} \log_2 \frac{6}{10} \right)$$

$$\rightarrow E = 0$$

## Entropy

Measure of “impurity” in a dataset

## Information gain

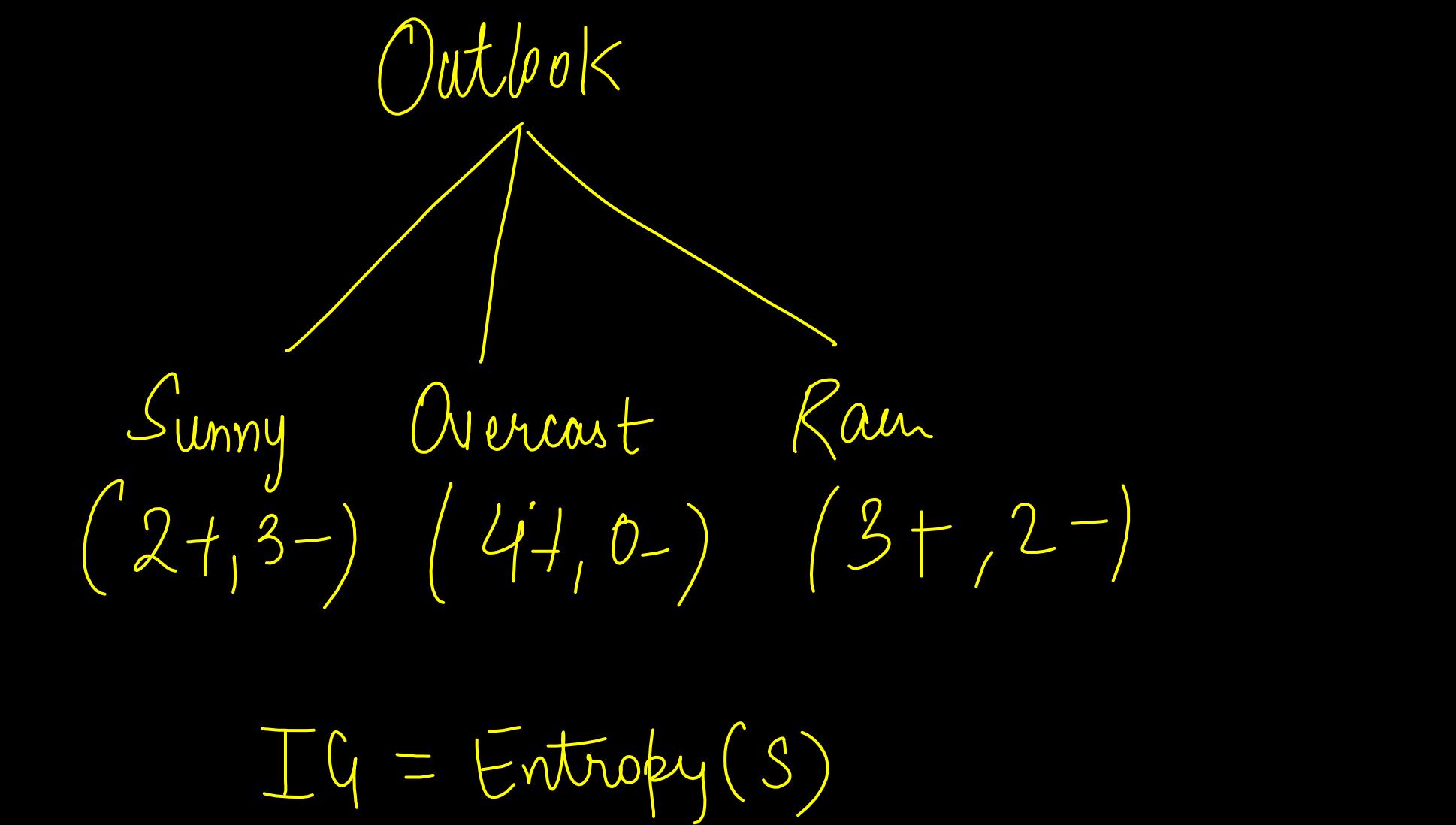
How much information gain I am getting from particular node.

$x_1$      $x_2$      $x_3$     output

Day	Outlook	Temp	Humidity	Wind	Tennis?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Day	Outlook	Temp	Humidity	Wind	Tennis?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Day	Outlook	Temp	Humidity	Wind	Tennis?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No



$$= 0.940 - \left[ \frac{5}{14} \times 0.971 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.971 \right]$$

$$= 0.247$$

$$IG(\text{Wind}) = 0.048$$

$$\underline{IG(\text{Temp})} = 0.029$$

$$IG(\text{Hum}) = 0.151$$

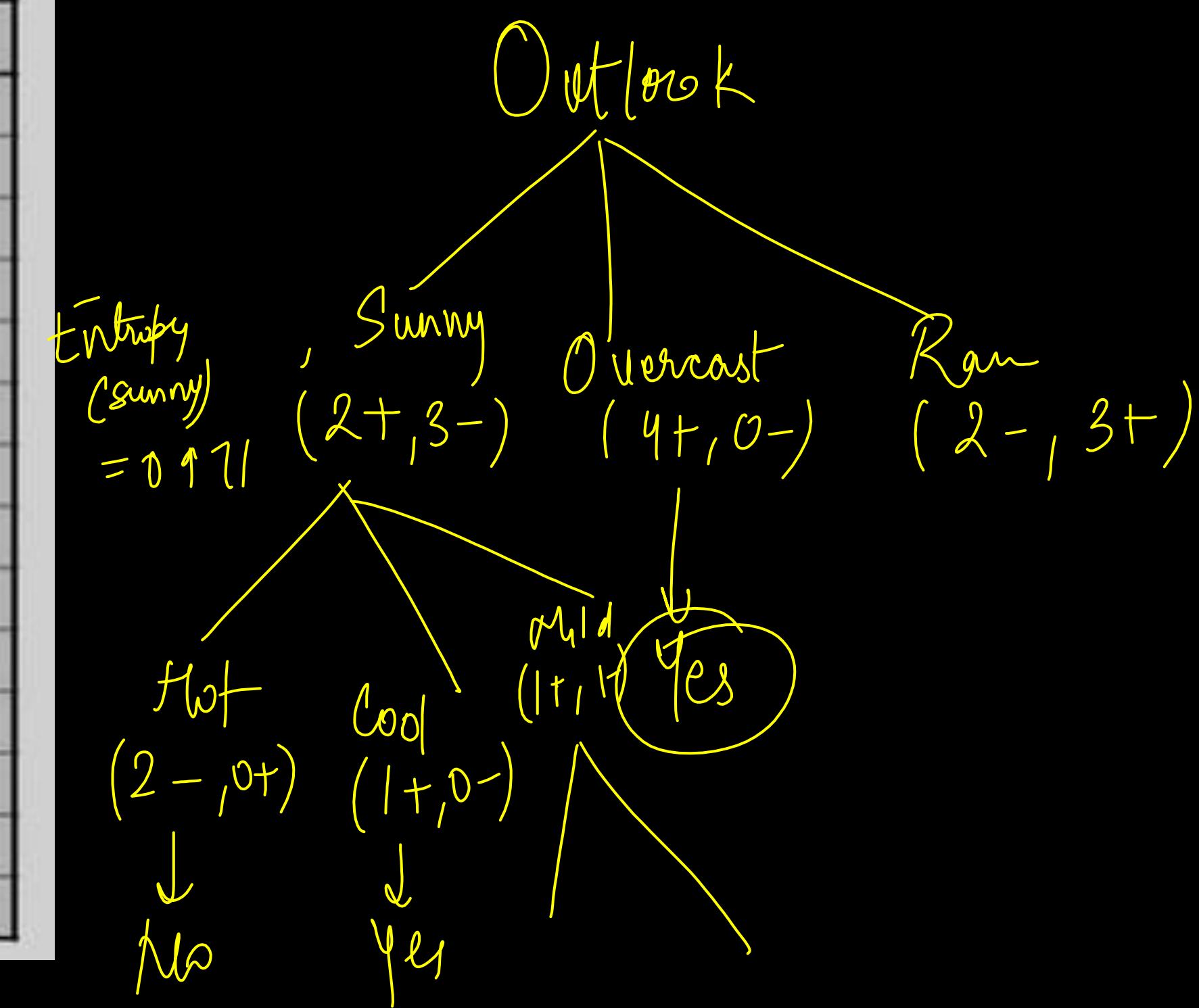
$$-\text{Entropy}(S) = -\left(\frac{5}{14} \log_2 \frac{5}{14} + \frac{9}{14} \log_2 \frac{9}{14}\right) = 0.940$$

$$\text{Entropy}(S, \text{Sunny}) = -\left(\frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5}\right) = 0.971$$

$$\text{Entropy}(S, \text{Overcast}) = 0$$

$$\text{Entropy}(S, \text{Rain}) = 0.971$$

Day	Outlook	Temp	Humidity	Wind	Tennis?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No



Day	Outlook	Temp	Humidity	Wind	Tennis?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Day	Outlook	Temp	Humidity	Wind	Tennis?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Day	Outlook	Temp	Humidity	Wind	Tennis?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Consider the dataset, S given below:

Elevation	Road Type	Speed Limit	Speed
steep	Uneven	Yes	Slow
steep	Smooth	Yes	Slow
flat	Uneven	No	Fast
steep	Smooth	No	Fast

Elevation, Road Type and speed Limit are the features and Speed is the target label that we want to predict.

$$\text{Entropy} = -\left(\frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{4} \log_2 \frac{2}{4}\right)$$

$$I(S) = 1 - \left[ \frac{3}{4} \times 0.918 + \frac{1}{4} \times 0 \right] = 0.3$$

Elevation

```

graph TD
    E[Elevation] --> S1[Steep]
    E --> S2[flat]
    S1 --> S1L["(2S, 1F)"]
    S2 --> S2L["(1F, 0S)"]

```

$$0.918 = \text{Entropy}(S, \text{steep}) = -\left[\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3}\right]$$

$$\text{Entropy}(S, \text{flat}) = -\left[\frac{1}{4} \log_2 1 + 0\right] = 0$$

Find the entropy of the dataset, S as given above:

- a. 0.5
- b. 0
- c. 1
- d. 0.7

Find the information Gain if the dataset is split at the feature "Elevation":

- a. 1
- b. 0
- c. 0.675
- d. 0.325

Find the feature on which the parent node must be chosen to split the dataset,  $S$  based on information gain:

- a. Speed Limit
- b. Road Type
- c. Elevation

## GATE DA 2024

Details of ten international cricket games between two teams “Green” and “Blue” are given in Table C. This table consists of matches played on different pitches, across formats along with their winners. The attribute Pitch can take one of two values: spin-friendly (represented as  $S$ ) or pace-friendly (represented as  $F$ ). The attribute Format can take one of two values: one-day match (represented as  $O$ ) or test match (represented as  $T$ ).

A cricket organization would like to use the information given in Table C to develop a decision-tree model to predict outcomes of future games between these two teams.

To develop such a model, the computed  $\text{InformationGain}(C, \text{Pitch})$  with respect to the Target is \_\_\_\_\_ (rounded off to two decimal places). 

C

4G, 6B

Match Number	Pitch	Format	Winner (Target)
1	S	T	Green →
2	S	T	Blue
3	F	O	Blue
4	S	O	Blue
5	F	T	Green →
6	F	O	Blue
7	S	O	Green →
8	F	T	Blue
9	F	O	Blue
10	S	O	Green →

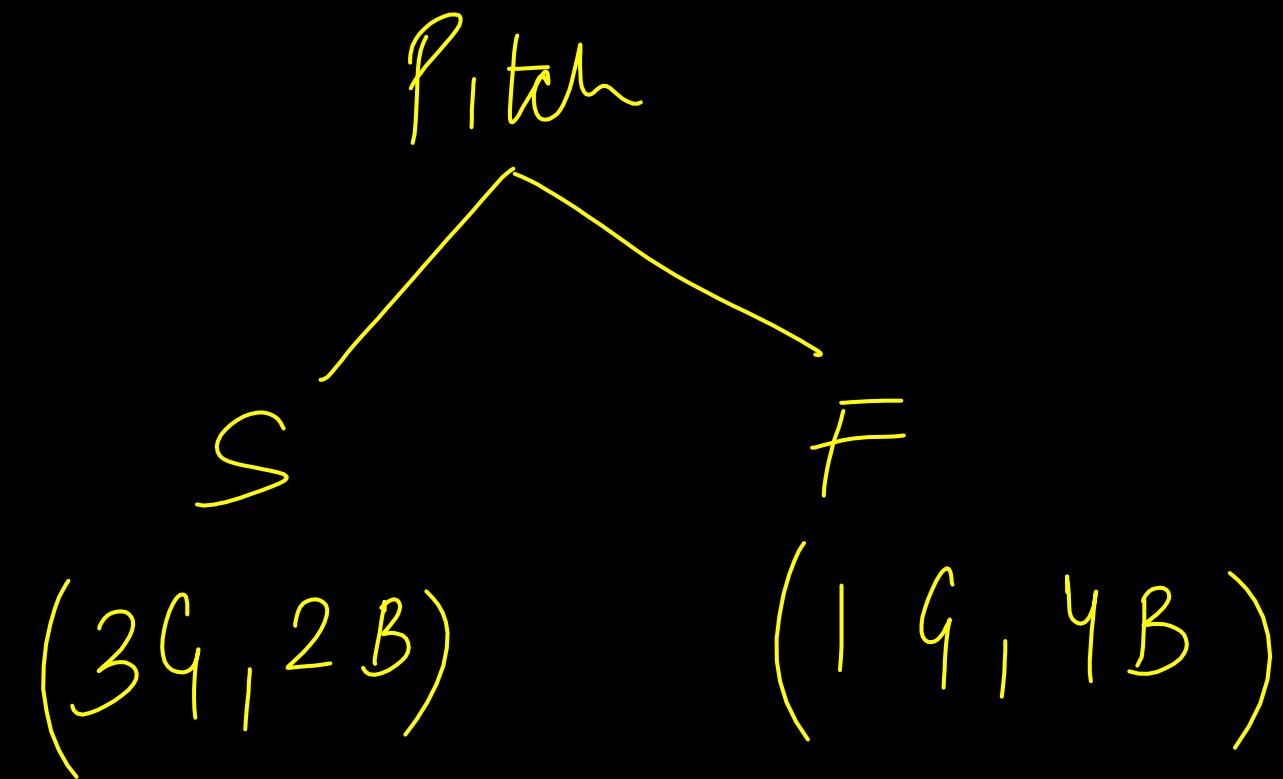
Information Gain (G, Pitch)

$$= 0.970 - \left[ \frac{5}{10} \times 0.970 + \frac{5}{10} \times 0.721 \right]$$

$$= 0.125$$

$$\text{Entropy}(C) = - \left[ \frac{4}{10} \log_2 \frac{4}{10} + \frac{6}{10} \log_2 \frac{6}{10} \right]$$

$$= 0.970$$



$$\text{Entropy}(\text{Pitch}, S) = - \left[ \frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5} \right]$$

$$\begin{aligned} \text{Entropy}(\text{Pitch}, F) &= - \left[ \frac{1}{5} \log_2 \frac{1}{5} + \frac{4}{5} \log_2 \frac{4}{5} \right] \\ &= 0.721 \end{aligned}$$

**True/False. Can we use a decision tree for both Regressions as well as Classification?**

a) True



b) False

Find the decision tree.

$x_1$	1	3	4	6	10	15	2	7	16	0
$x_2$	12	23	21	10	27	23	35	12	27	17
$y$	10.1	15.3	11.5	13.9	17.8	23.1	12.7	43.0	17.6	14.9

→

T

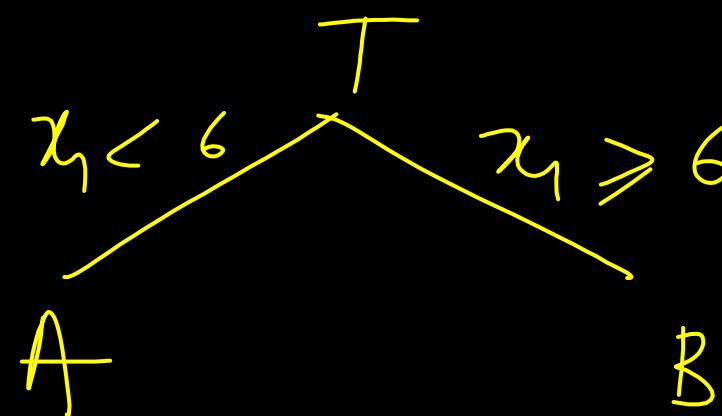
$x_1$	1	3	4	6	10	15	2	7	16	0
$x_2$	12	23	21	10	27	23	35	12	27	17
$y$	10.1	15.3	11.5	13.9	17.8	23.1	12.7	43.0	17.6	14.9

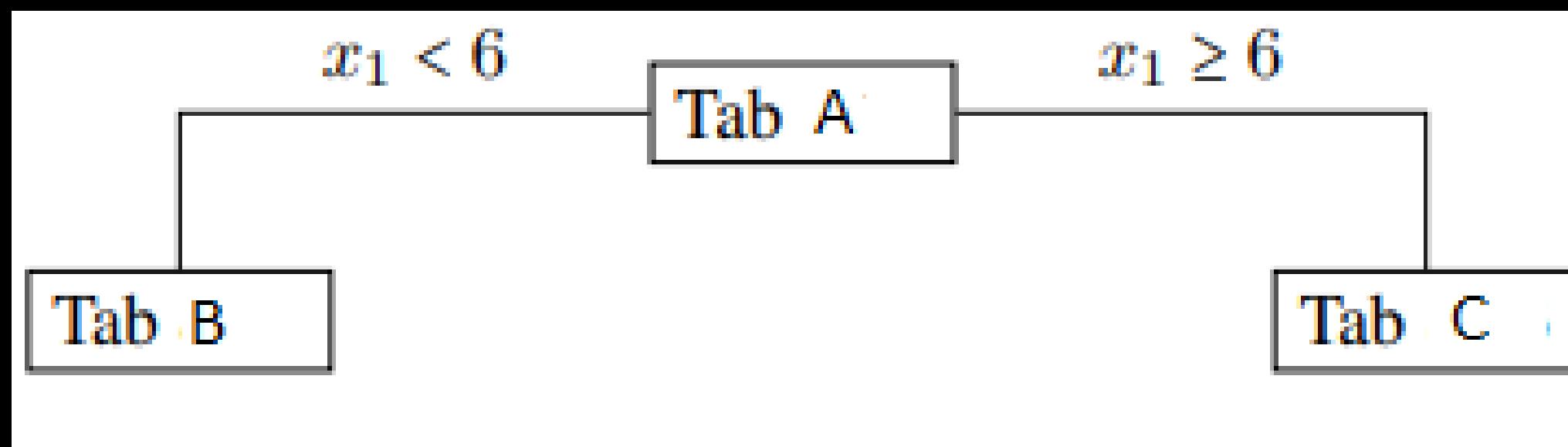
A

$x_1$	1	3	4	2	0
$x_2$	12	23	21	35	17
$y$	10.1	15.3	11.5	12.7	14.9

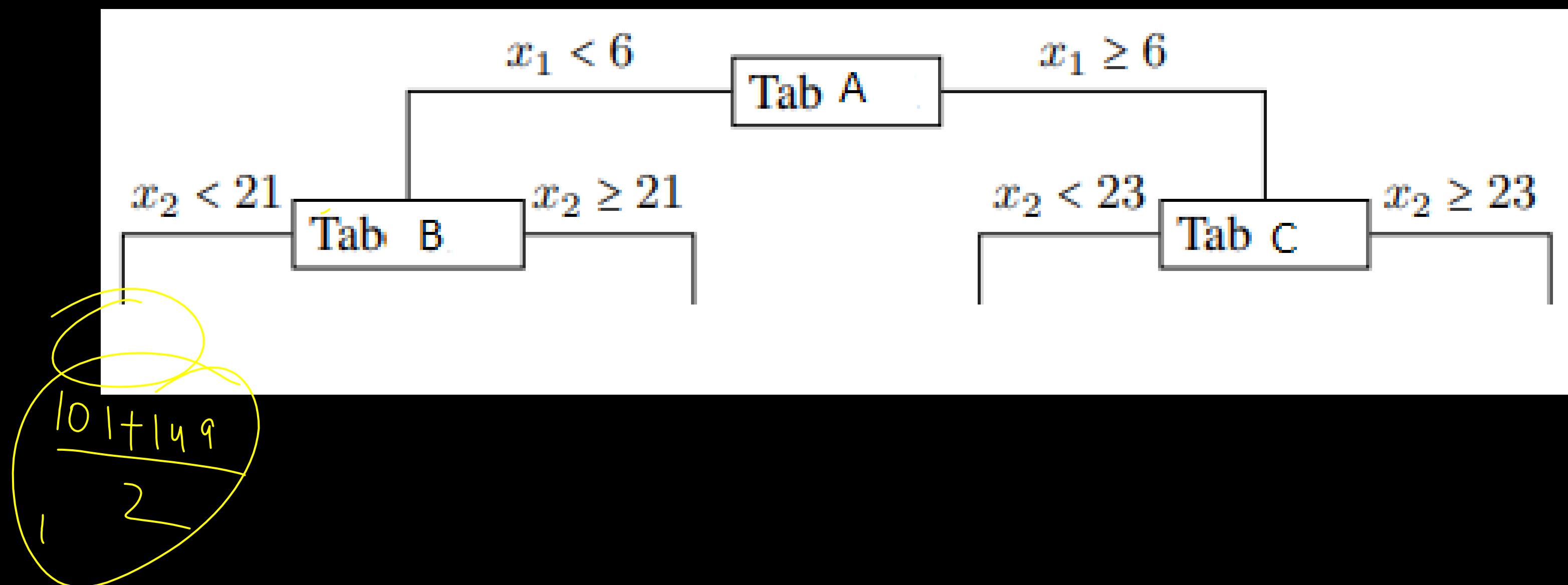
B

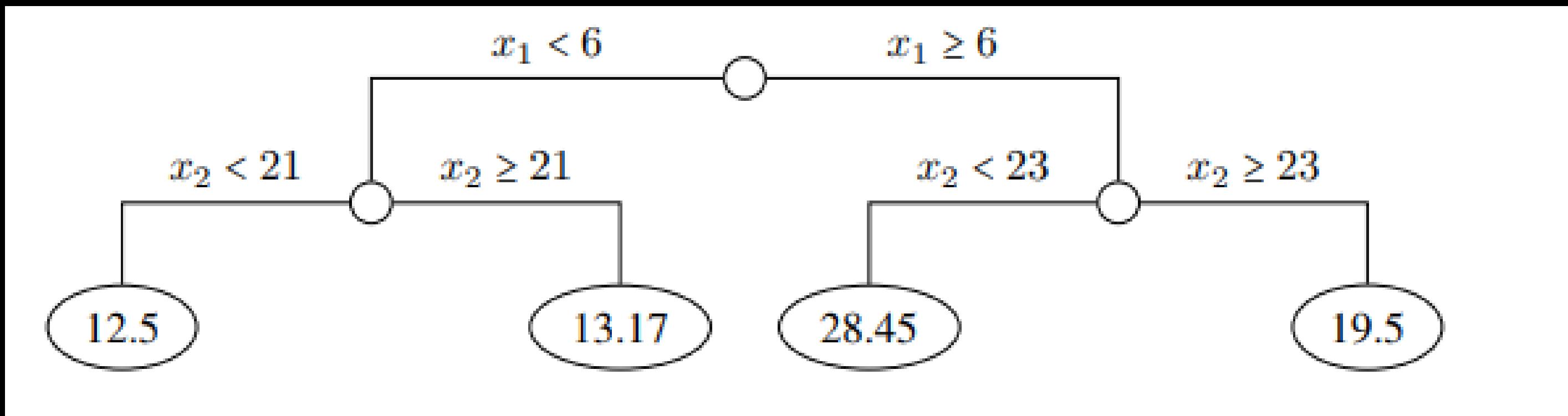
$x_1$	6	10	15	7	16
$x_2$	10	27	23	12	27
$y$	13.9	17.8	23.1	43.0	17.6





$$x_1 = 5, m = 10$$



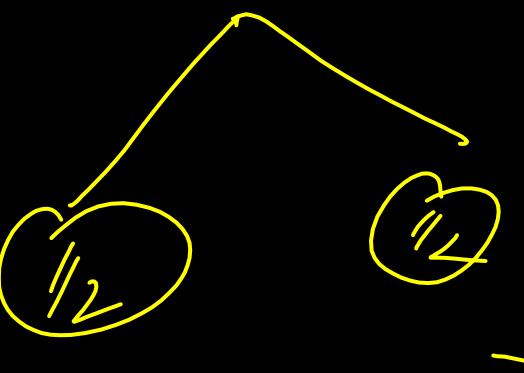


**Decision tree is a \_\_\_\_\_.**

- a) Non-linear ML technique.
- b) Non-Parametric technique.
- c) Supervised Learning technique.
- d) All of the above.

**Entropy value for the data sample that has 50-50 split belonging to two categories is**

- a. 1
- b. 0
- c. None



Binary class

$$\underline{0 < E < 1}$$

Which of the following is true for a decision tree?

- a. A decision tree is an example of a linear classifier.
- b. The entropy of a node typically decreases as we go down a decision tree.
- c. Entropy is a measure of purity.
- d. An attribute with lower mutual information should be preferred to other attributes.

Day	Outlook	Temp	Humidity	Wind	Tennis?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Find the Gini Index of tennis.

$$\frac{5}{14}, \frac{9}{14}$$

$$G I = 1 - \sum p_i^2$$
$$= 1 - \left[ \left( \frac{5}{14} \right)^2 + \left( \frac{9}{14} \right)^2 \right]$$

$$\rightarrow \text{Gini split index} = \frac{\sum |S_v|}{|S|} GI(S_{1,v})$$

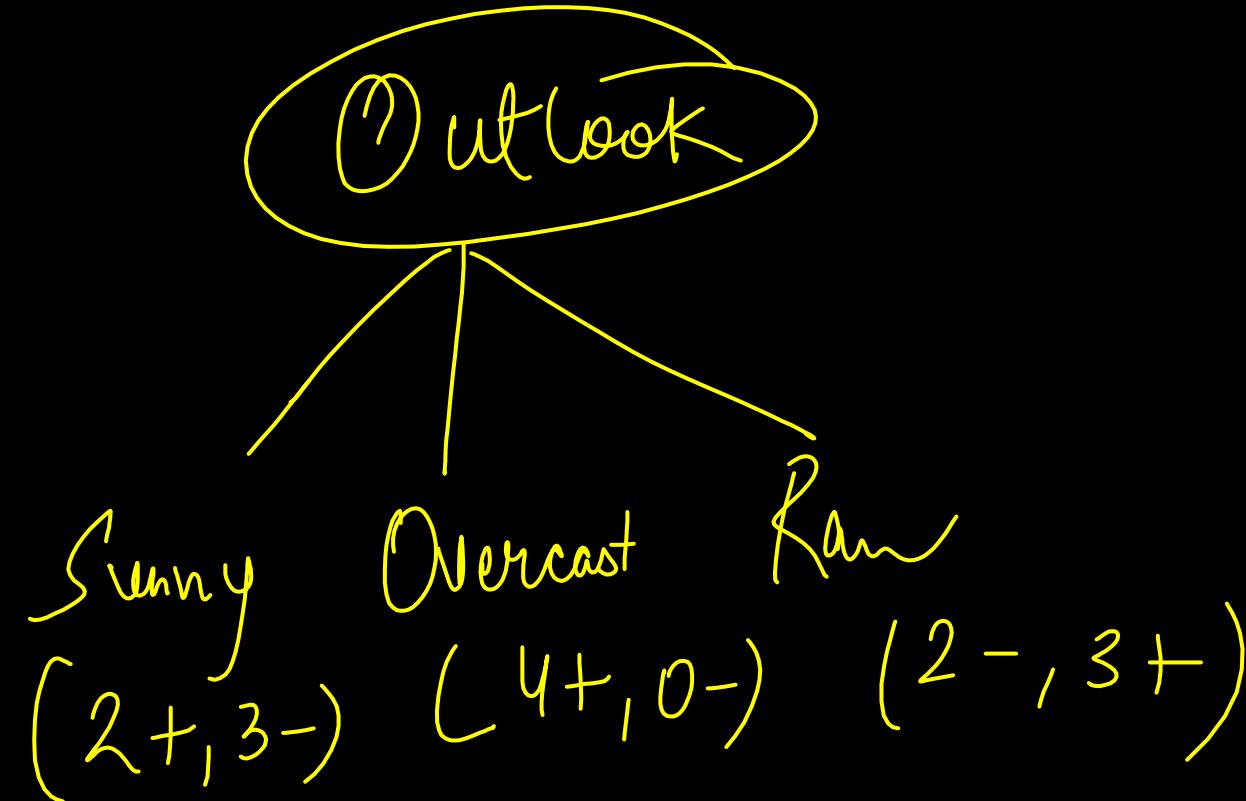
choose feature  
which has minimum  
Gini split index.

Find the Gini Split Index of Outlook.

$$GI(S_1, \text{sunny}) = 1 - \left[ \left( \frac{2}{5} \right)^2 + \left( \frac{3}{5} \right)^2 \right] = 0.48$$

$$\begin{aligned} GI(S_1, \text{Overcast}) &= 1 - \left[ \left( \frac{4}{4} \right)^2 \right] \\ &= 0 \end{aligned}$$

$$GI(S_1, \text{Rain}) = 0.48$$



$$\begin{aligned} \text{Gini Split Index}(\text{Outlook}) &= \sum \frac{|S_v|}{|S|} \times GI(S_1, v) \\ &= \frac{5}{14} \times 0.48 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.48 \\ &= 0.342 \end{aligned}$$

# Overfitting in decision trees

Causes:

1. The tree has too many levels (i.e., it's too deep).
2. The tree splits down to very specific details of the training set

To deal with an overfitted Decision Tree, you can use one or more of the following methods:

**Early stopping**

**Feature selection**



**Pruning:** Pruning is a technique used to simplify the decision tree by removing branches that do not improve the tree's performance on new data.

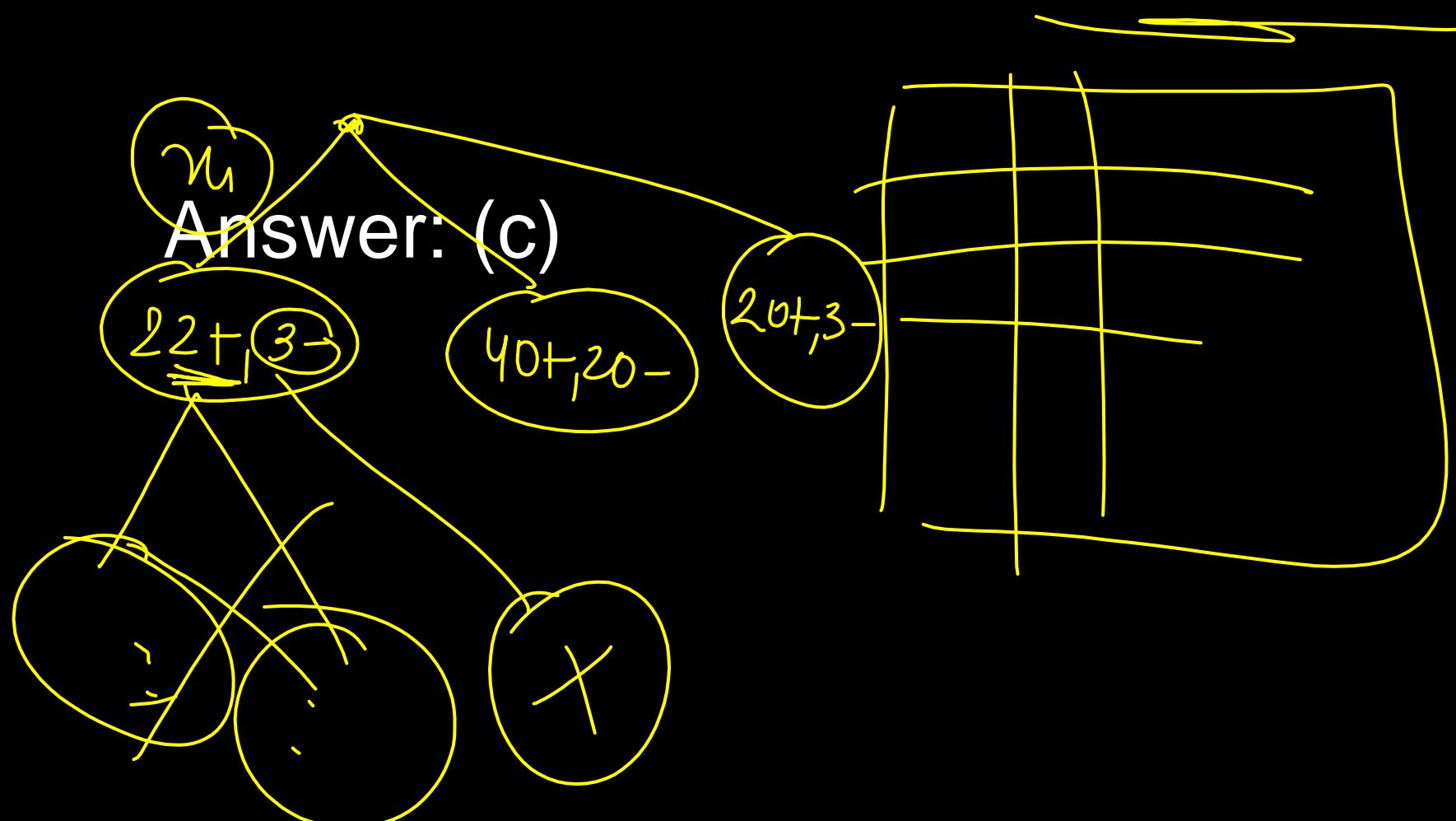
Pruning involves building the full decision tree and then removing branches that do not improve the tree's performance on a validation set. The idea is to remove branches that are unlikely to be useful for predicting new data, thereby simplifying the tree and reducing the risk of overfitting.

Pruning is a technique that reduces the size of decision trees by removing sections of the tree that provide little power to classify instances. This is done in order to avoid:

- a. overfitting
- b. underfitting

Which of the following statements is not true about Pruning?

- a) It removes the sections of the tree that provide little power to classify instances
- b) It is a technique in machine learning and search algorithms to reduce the size of the decision trees
- c) It increases the complexity of the final classifier
- d) It improves the predictive accuracy by the reduction of overfitting

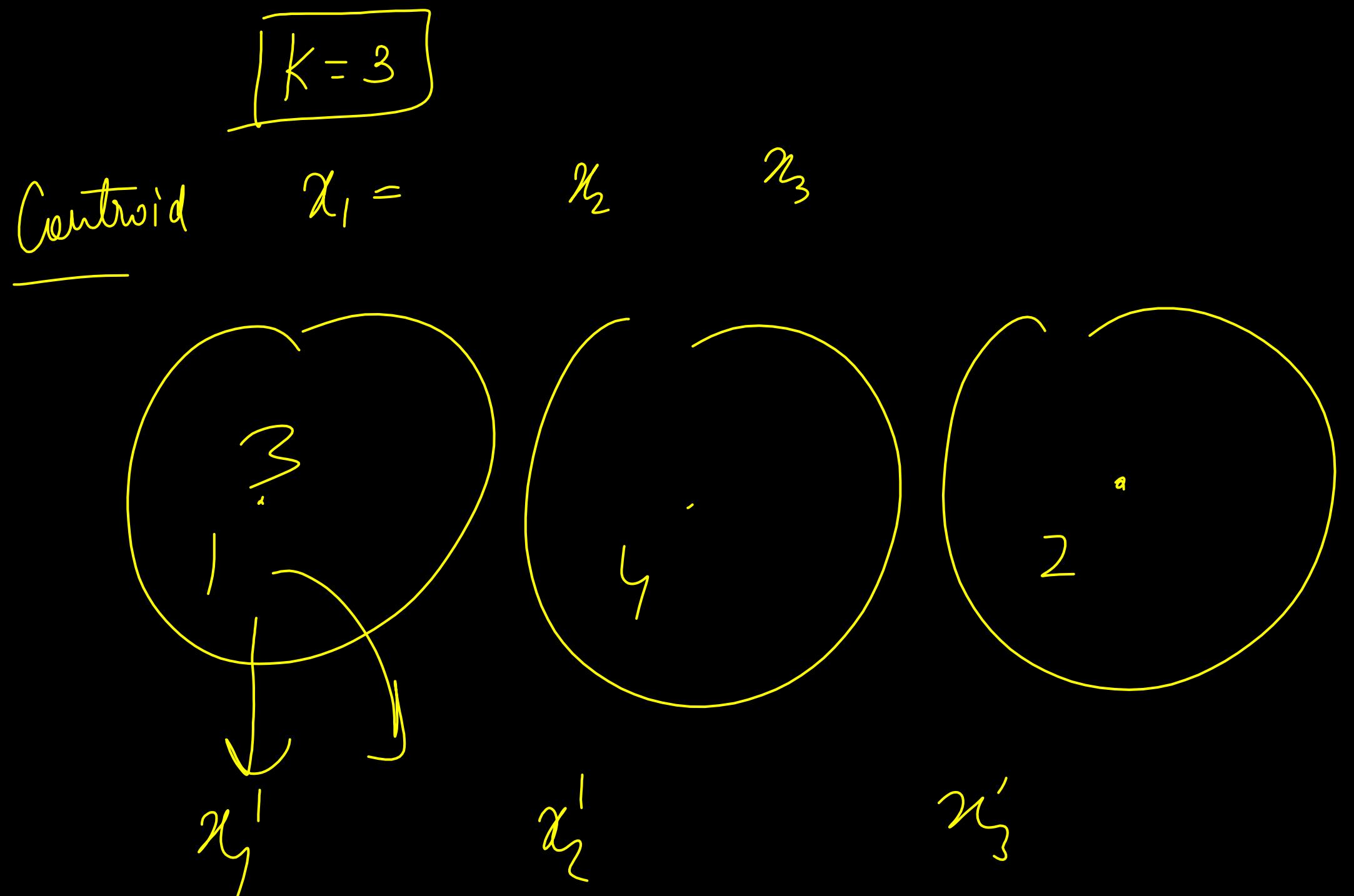


# Clustering Algorithm – K Means

K-Means

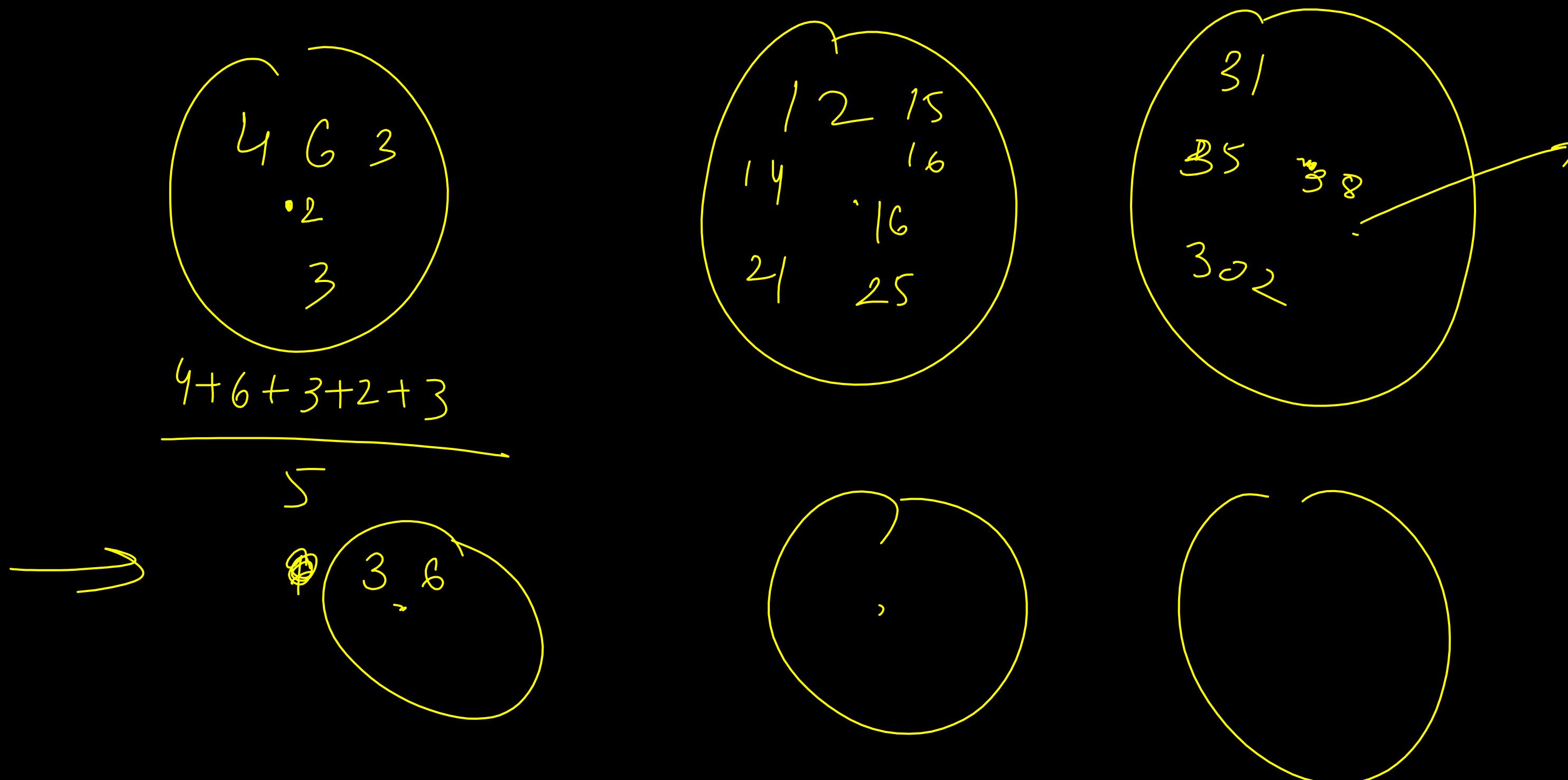
Health, En , Temp , Weather Medicine

## K-Means Clustering

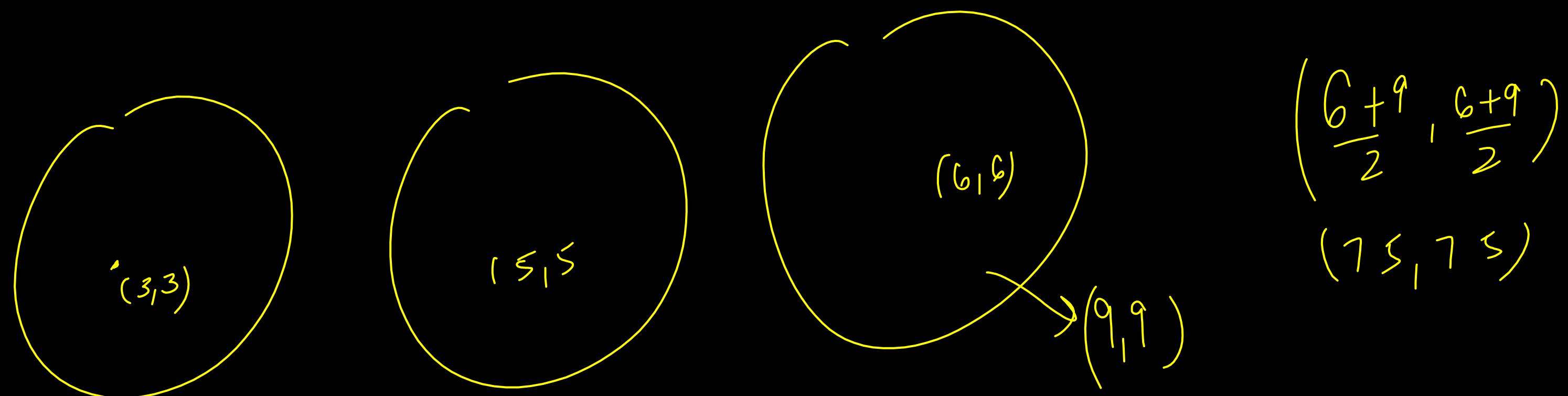


**Problem** Apply k-means algorithm for given data with  $k = 3$ . Use  $C_1(2)$ ,  $C_2(16)$  and  $C_3(38)$  as initial centers.

Data: 2, 4, 6, 3, 31, 12, 15, 16, 38, 35, 14, 21, 3, 25, 30 2.



Given the dataset:  $(1, 1), (3, 3), (4, 4), (5, 5), (6, 6), (9, 9), (0, 3), (3, 0)$  and assuming the initial centroids for ( $K = 3$  – means clustering) to be  $C_1 = (3, 3)$ ,  $C_2 = (5, 5)$  and  $C_3 = (6, 6)$ . One iteration for K-means clustering, will update  $C_3$  to  $(\underline{\quad}, \underline{\quad})$



During a research work, you found seven observations as described with the data points below. You want to create three clusters from these observations using K-means algorithm.

After first iteration, the clusters C1, C2, C3 has following observations:

C1:  $\{(2,2), (4,4), (6,6)\}$

C2:  $\{(0,4), (4,0)\}$

C3:  $\{(5,5), (9,9)\}$

If you want to run a second iteration then what will be the cluster centroids?

K-mean clustering algorithm has clustered the given 8 observations into 3 clusters after 1st iteration as follows:

$$C_1 : \{(3,3), (5,5), (7,7)\}$$

$$C_2 : \{(0,6), (6,0), (3,0)\}$$

$$C_3 : \{(8,8), (4,4)\}$$

What will be the Manhattan distance for observation (4,4) from cluster centroid  $C_1$  in second iteration?

Answer: (2)

## GATE DA 2024

Euclidean distance based  $k$ -means clustering algorithm was run on a dataset of 100 points with  $k = 3$ . If the points  $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$  and  $\begin{bmatrix} -1 \\ 1 \end{bmatrix}$  are both part of cluster 3, then which **ONE** of the following points is necessarily also part of cluster 3?

(A)	$\begin{bmatrix} 0 \\ 0 \end{bmatrix}$
(B)	$\begin{bmatrix} 0 \\ 2 \end{bmatrix}$
(C)	$\begin{bmatrix} 2 \\ 0 \end{bmatrix}$
(D)	$\begin{bmatrix} 0 \\ 1 \end{bmatrix}$

Which of the following can be possible termination conditions in K-Means?

1. For a fixed number of iterations.
2. Assignment of observations to clusters does not change between iterations. Except for cases with a bad local minimum.
3. Centroids do not change between successive iterations.
4. All of the above

## What is true about K-Mean Clustering?

- 1. K-means is extremely sensitive to cluster center initializations
  - 2. Bad initialization can lead to Poor convergence speed
  - 3. Bad initialization can lead to bad overall clustering
- 
- a. 1 and 2
  - b. 1 and 3
  - c. All of the above
  - d. 2 and 3

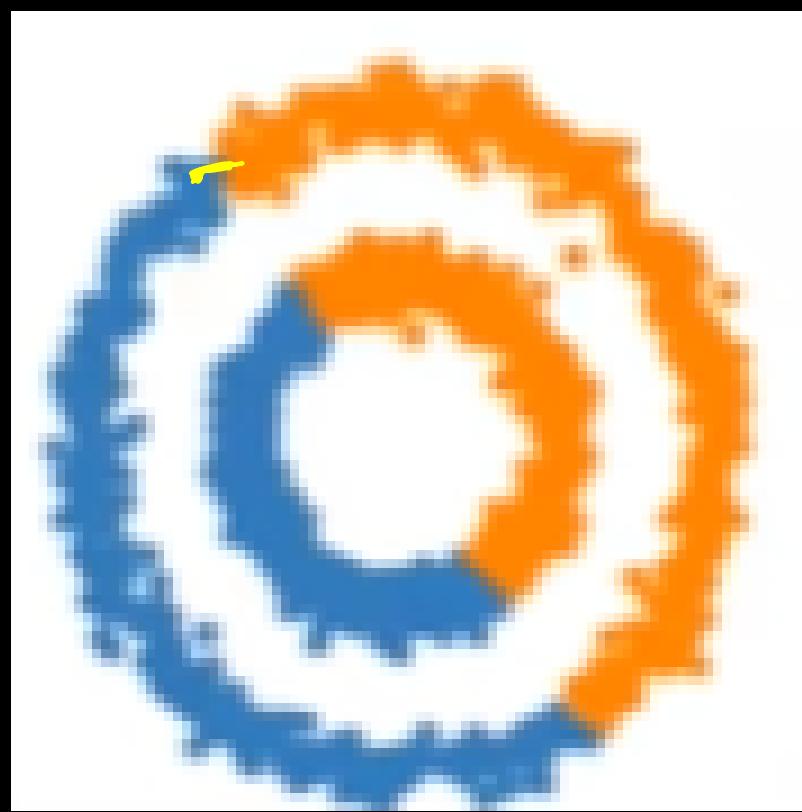
**Which of the following clustering algorithms is the most sensitive to outliers?**

- a. K-means clustering algorithm
- b. K-medians clustering algorithm
- c. K-modes clustering algorithm
- d. None of the above

## Limitations

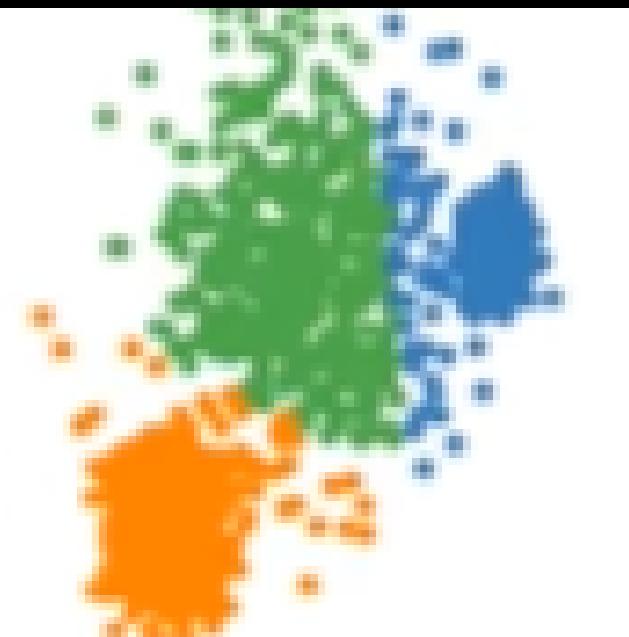
### 1. Assumes Spherical Clusters:

K-Means performs well when clusters are spherical, equally sized, and have similar variances.



## 2. Sensitive to Outliers:

Outliers can significantly impact cluster assignments and centroid positions.

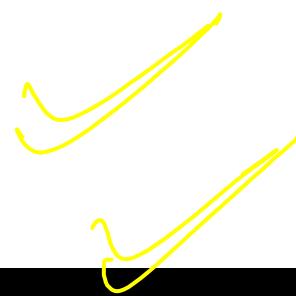


## 3. Requires Specifying $K$ :

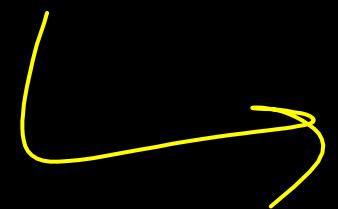
Determining the appropriate number of clusters ( $K$ ) can be challenging.

In which of the following cases will K-Means clustering fail to give good results? (Choose all the correct answers)

- a. Data points with outliers
- b. Data points with round shapes
- c. Data points with non-convex shapes
- d. Data points with different densities

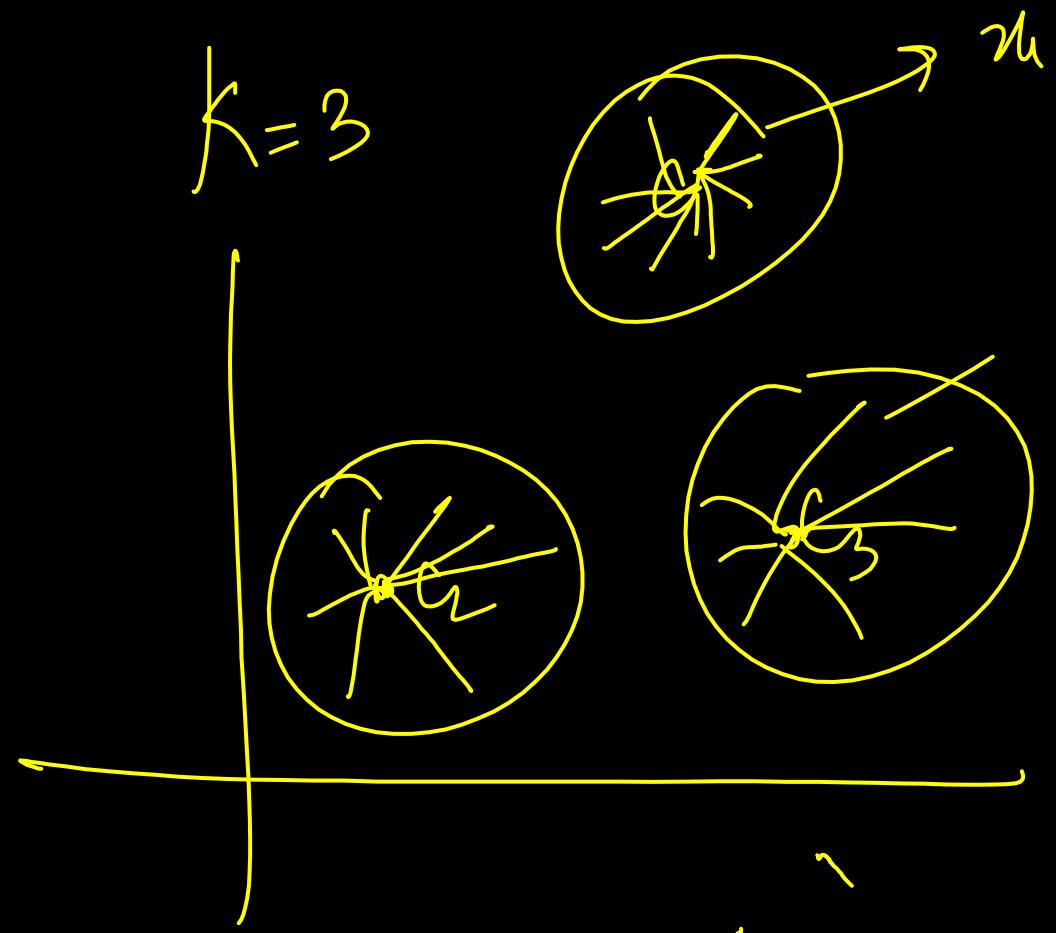


## SSE (or WCSS) in K Means



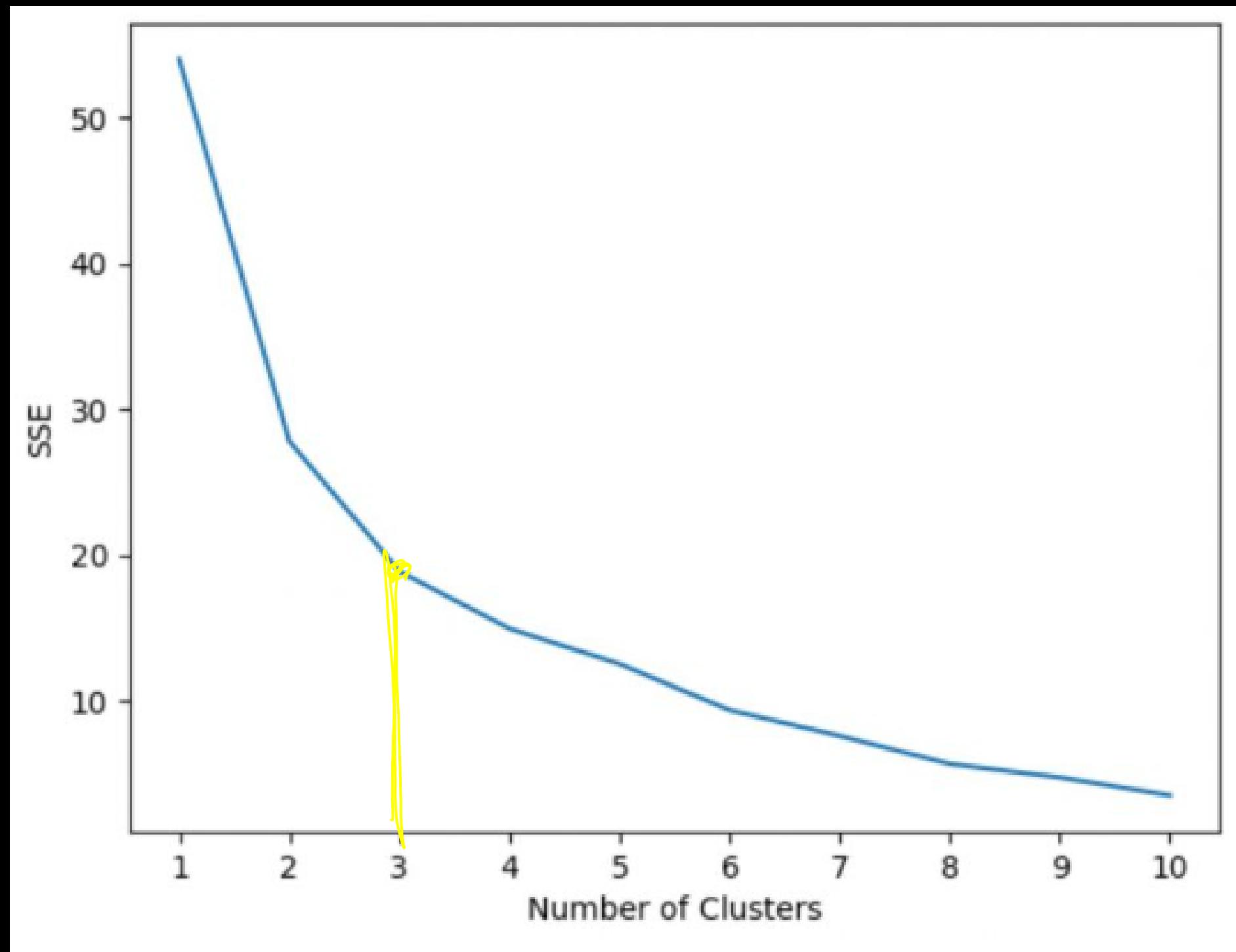
$$\text{WCSS} = \sum_{k=1}^K \sum_{x \in C_k} \|x - \mu_k\|^2$$

$$SSE = \sum_{k \in \text{part}} \sum_{l \in C} |x_k - c_l|^2$$



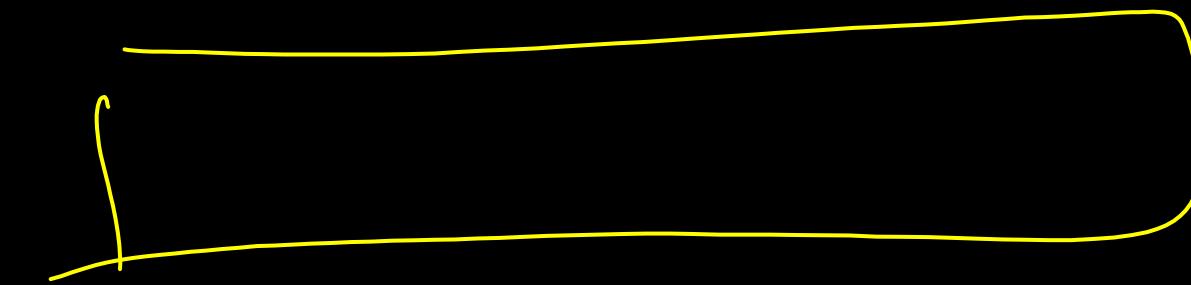
## Elbow method

To determine the optimal number of clusters, we'll create a plot that displays the number of clusters along with the SSE (sum of squared errors) of the model.

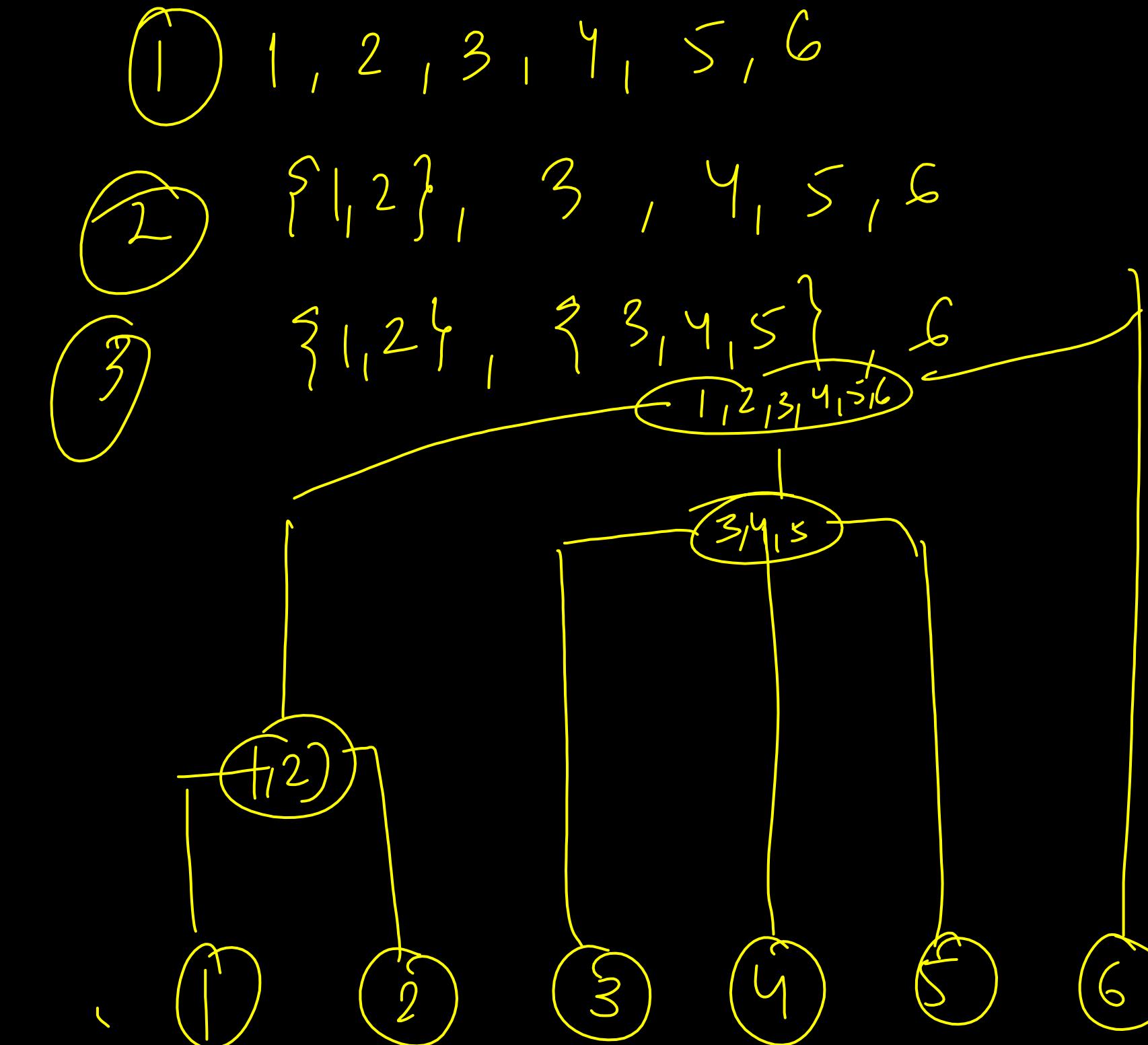
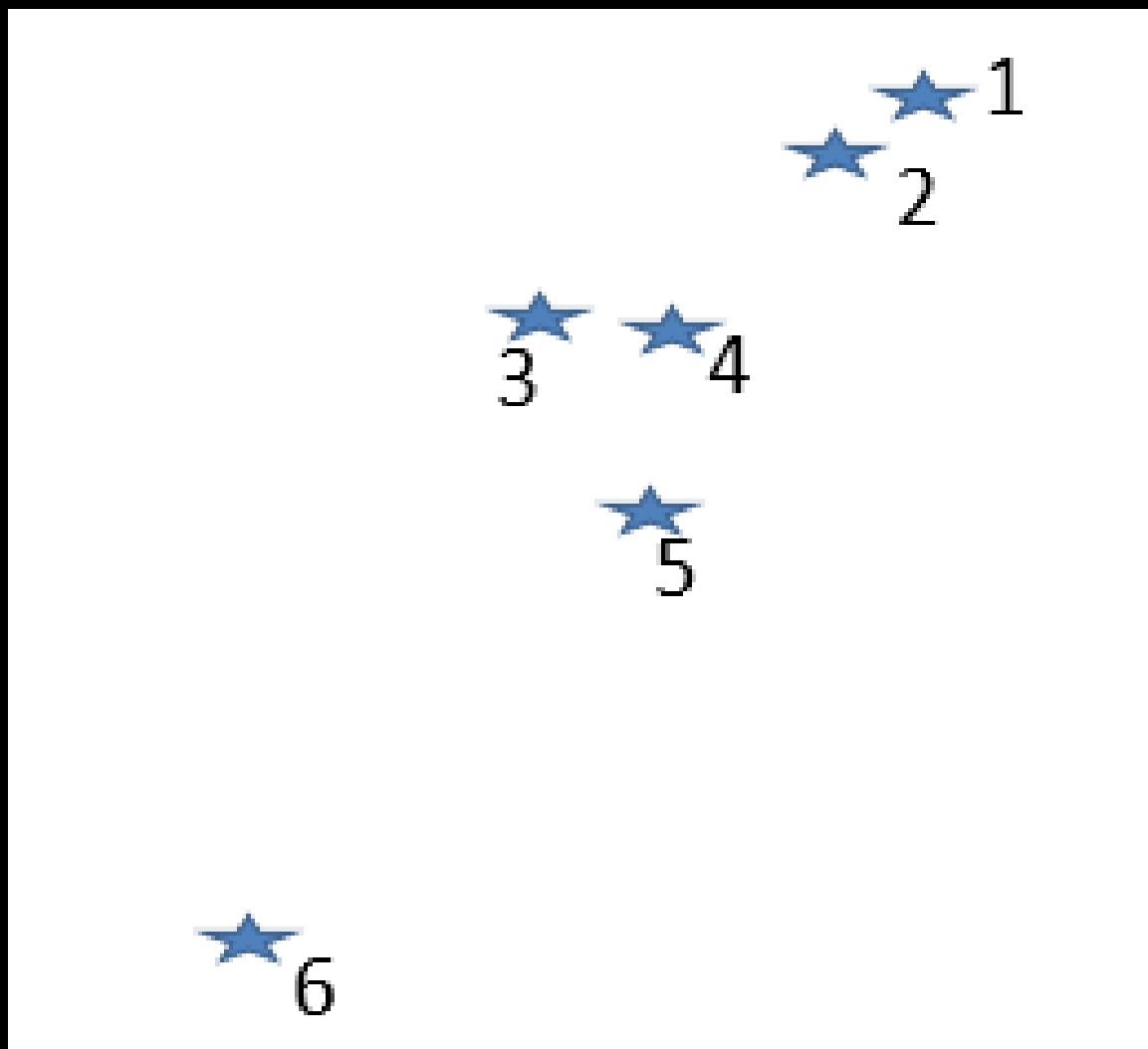


$k = 3$  clusters

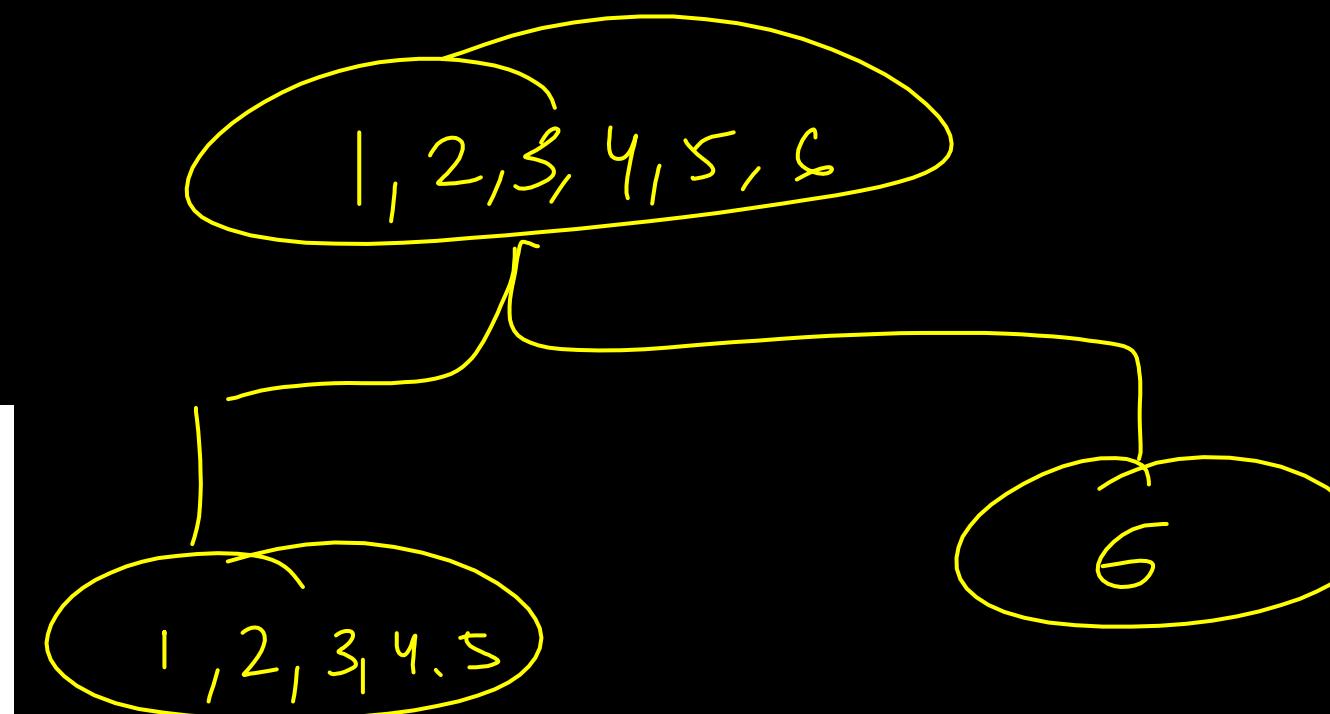
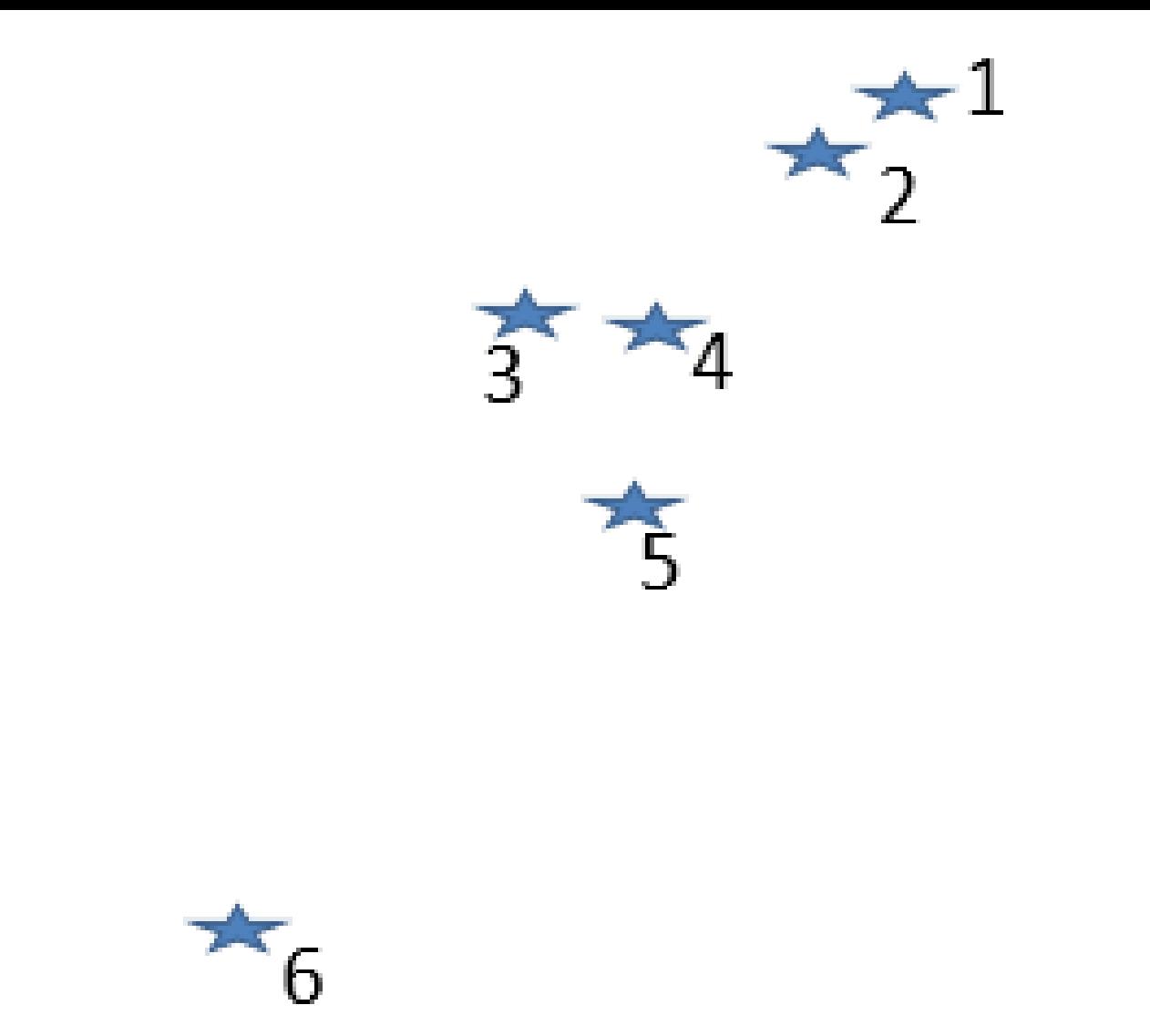
# Hierarchical Clustering



## Agglomerative clustering (bottom up)



## Divisive clustering (top down)



## Measures of dissimilarity

The decision regarding whether two clusters are to be merged or not is taken based on the measure of dissimilarity between the clusters.

## Measures of distance between data points

### Numeric data

Name	Formula
Euclidean distance	$\ \bar{x} - \bar{y}\ _2 = \sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2}$
Squared Euclidean distance	$\ \bar{x} - \bar{y}\ _2^2 = (x_1 - y_1)^2 + \dots + (x_n - y_n)^2$
Manhattan distance	$\ \bar{x} - \bar{y}\ _1 =  x_1 - y_1  + \dots +  x_n - y_n $
Maximum distance	$\ \bar{x} - \bar{y}\ _\infty = \max\{ x_1 - y_1 , \dots,  x_n - y_n \}$

## Non-numeric data

the Levenshtein distance between “kitten” and “sitting” is 3

)

## Measures of dissimilarity

### Measures of distance between groups of data points

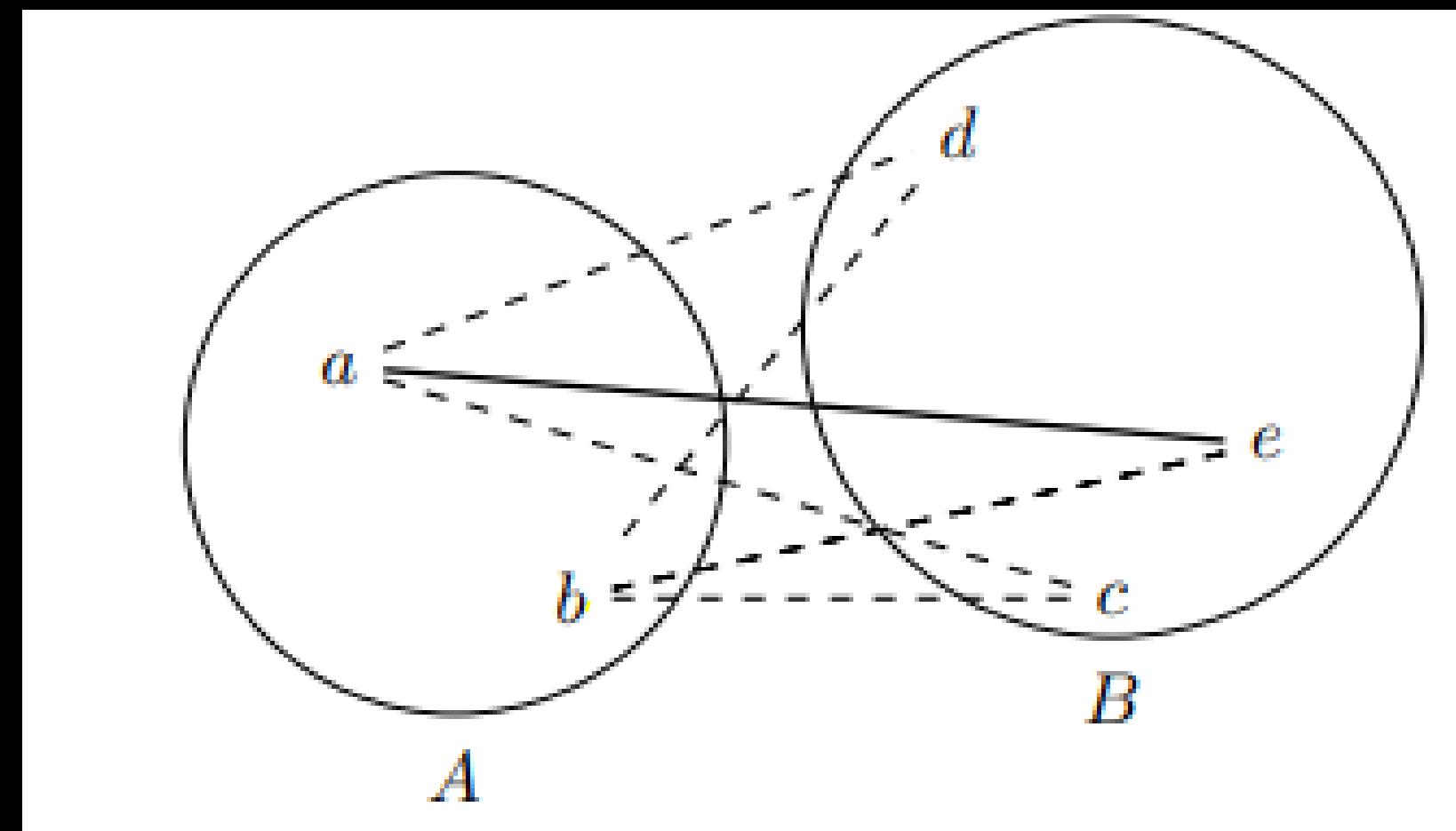
Let  $d(x, y)$  = distance between  $x$  and  $y$  (say Euclidean distance formula)

$d(A, B)$  = distance between the groups A and B.

The following are some of the different methods in which  $d(A, B)$  is defined.

1.  $d(A, B) = \max\{d(x, y) : x \in A, y \in B\}.$

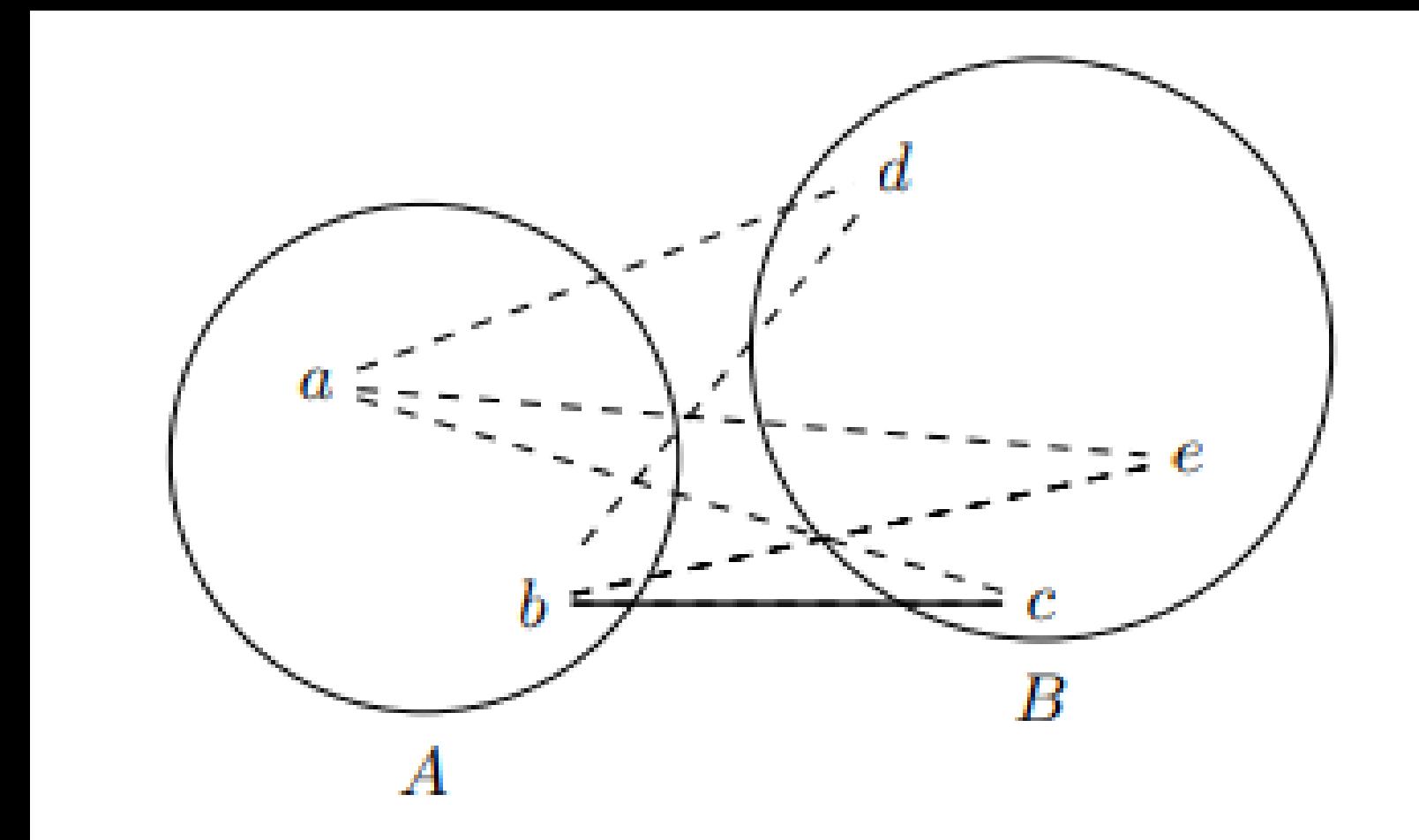
Agglomerative hierarchical clustering using this measure of dissimilarity is known as complete linkage clustering. The method is also known as farthest neighbour clustering.



$$2. d(A, B) = \min\{d(x, y) : x \in A, y \in B\}.$$

Agglomerative hierarchical clustering using this measure of dissimilarity is known as single linkage clustering.

The method is also known as nearest neighbour clustering.

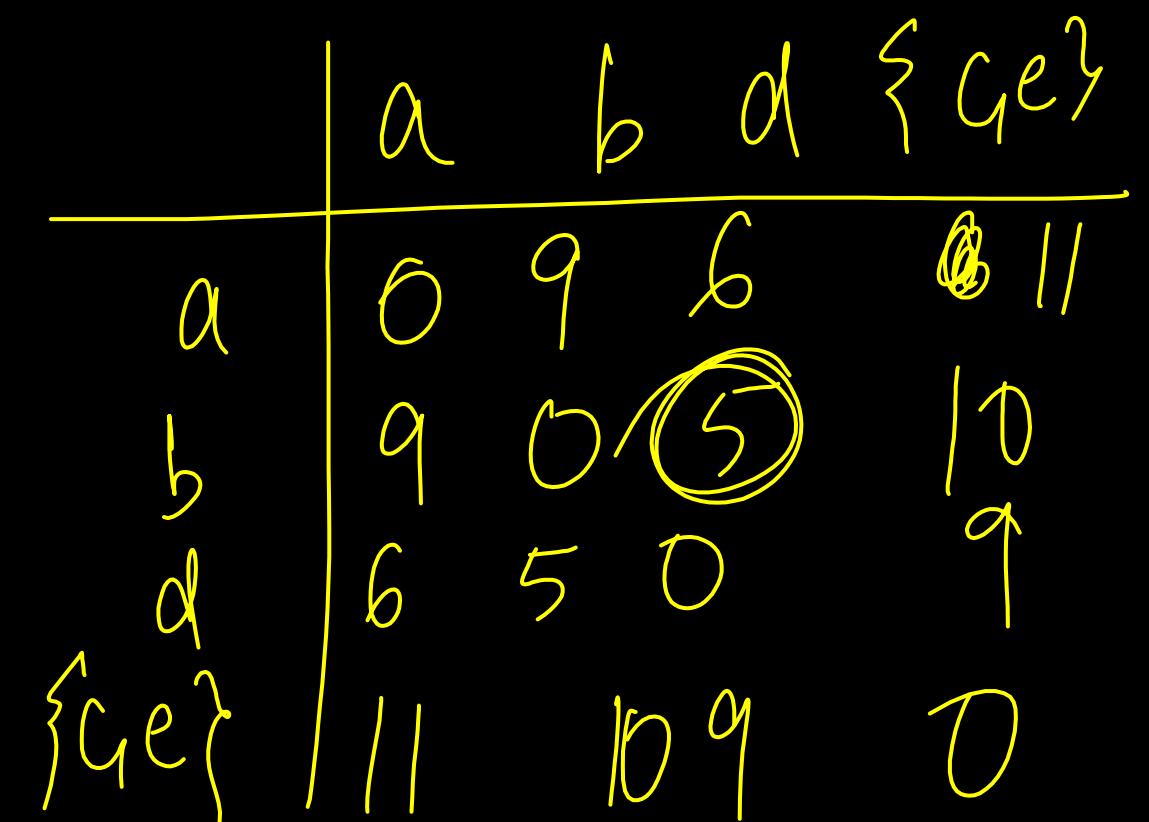
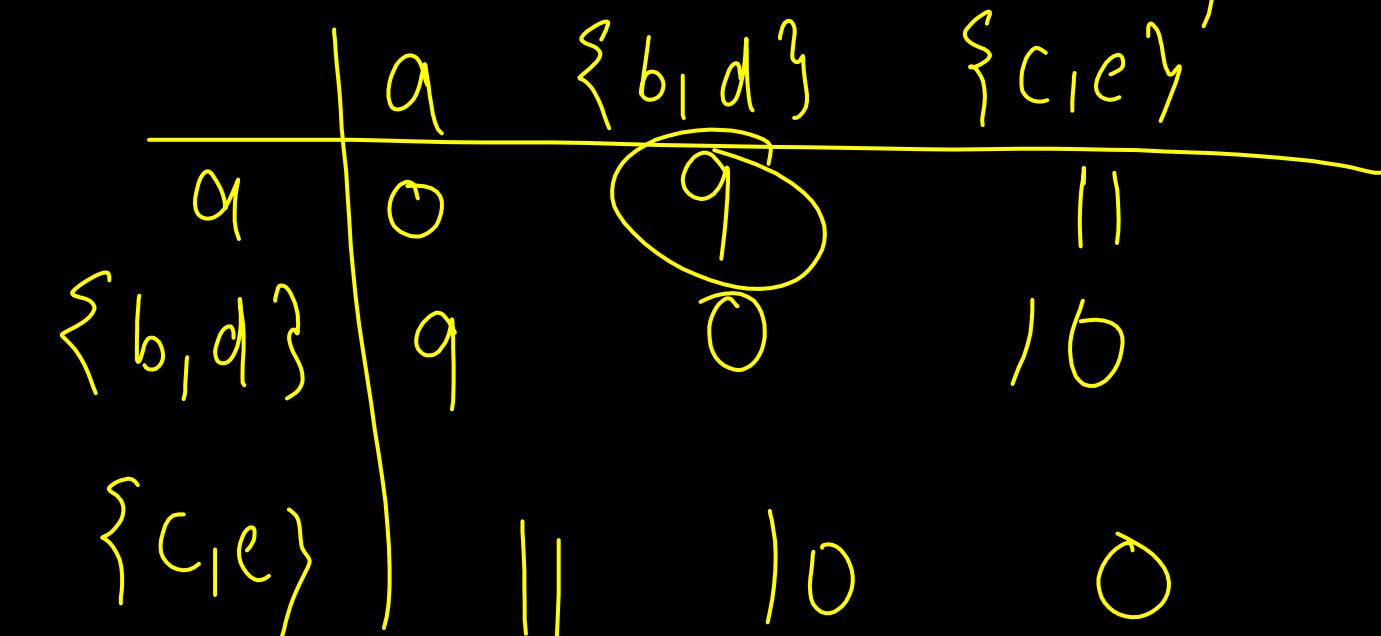


3.  $d(A, B) = \frac{1}{|A||B|} \sum_{x \in A, y \in B} d(x, y)$  where  $|A|, |B|$  are respectively the number of elements in  $A$  and  $B$ .

**Problem** Given the dataset  $\{a, b, c, d, e\}$  and the following distance matrix, construct a dendrogram by complete linkage hierarchical clustering using the agglomerative method.

- ①  $a, b, c, d, e$
- ②  $a, b, d, \{c, e\}$
- ③  $a, \{b, d\}, \{c, e\}$
- ④  $\{a, b, d\}, \{c, e\}$
- ⑤  $\{a, b, c, d, e\}$

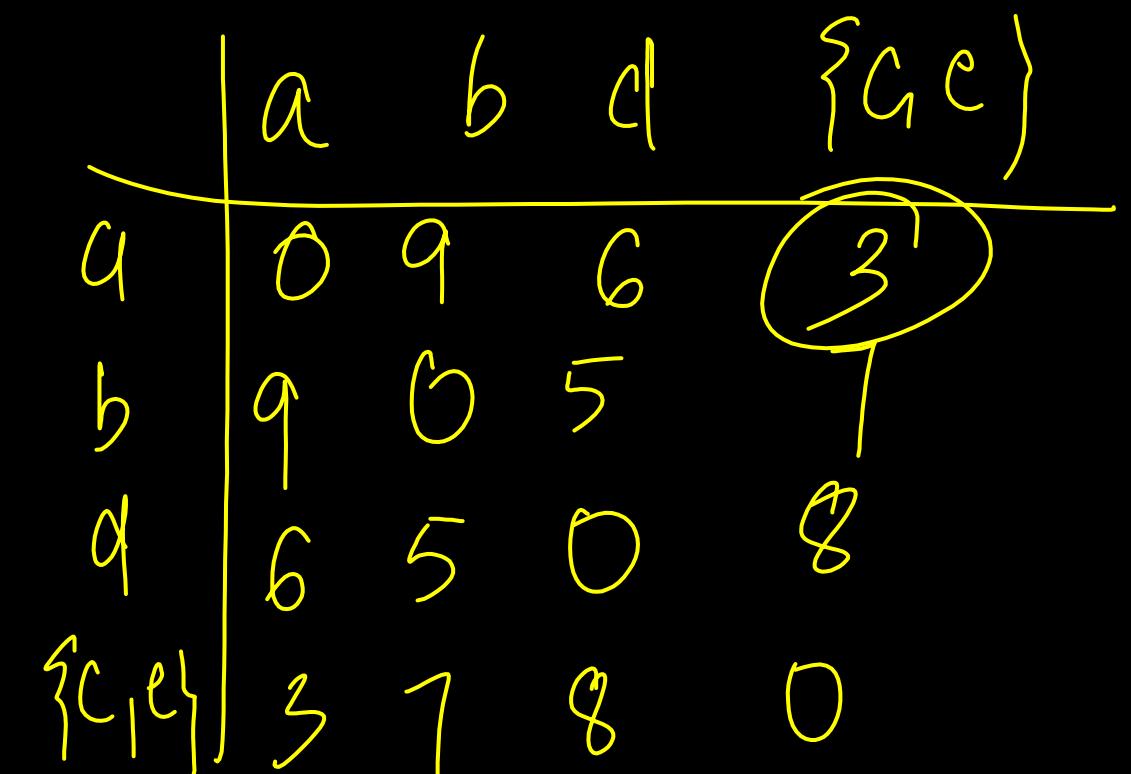
	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>
<i>a</i>	0	9	3	6	11
<i>b</i>	9	0	7	5	10
<i>c</i>	3	7	0	9	2
<i>d</i>	6	5	9	0	8
<i>e</i>	11	10	2	8	0



**Problem** Given the dataset  $\{a, b, c, d, e\}$  and the following distance matrix, construct a dendrogram by single linkage hierarchical clustering using the agglomerative method.

- ①  $a, b, c, d, e$
- ②  $a, b, d, \{c, e\}$
- ③  $a, b, \{d, \{c, e\}\}$
- ④

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>
<i>a</i>	0	9	3	6	11
<i>b</i>	9	0	7	5	10
<i>c</i>	3	7	0	9	2
<i>d</i>	6	5	9	0	8
<i>e</i>	11	10	2	8	0



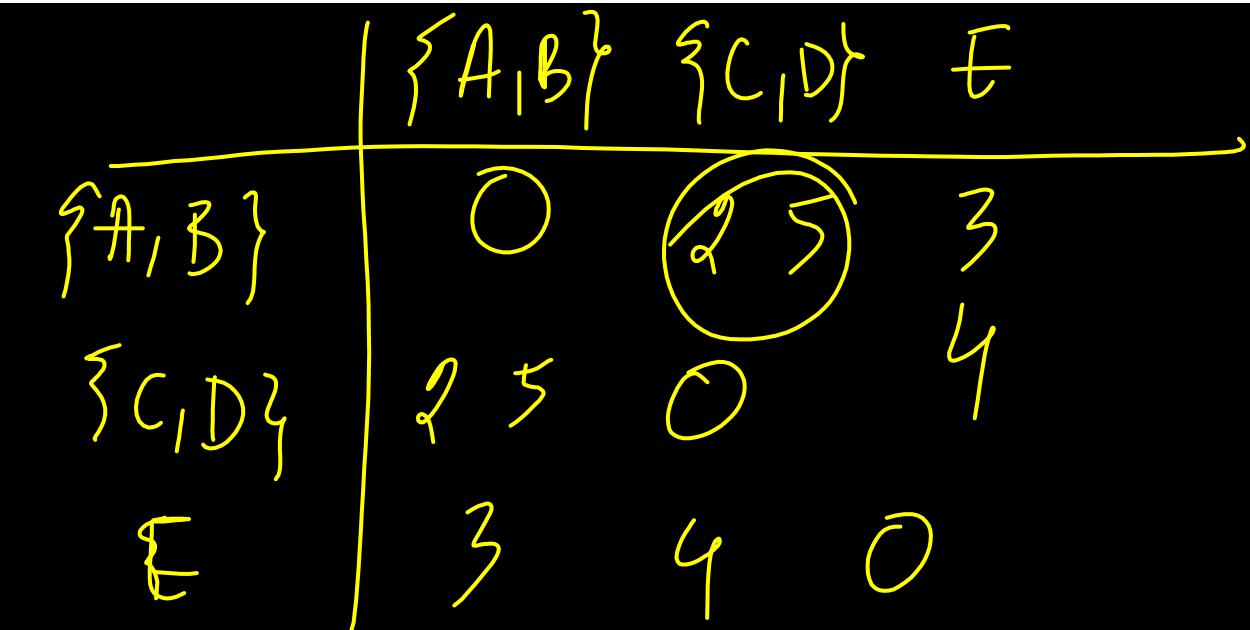
Given the following distance matrix, construct the dendrogram using agglomerative clustering with average linkage.

$$\frac{2+2+2+1}{9 \times 2} =$$

$$\frac{3+3}{2 \times 1}$$

	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	1	5
D	2	4	1	0	3
E	3	3	5	3	0

- ① A, B, C, D, E
- ② {A, B}, {C, D}, {E}
- ③ {A, B, C, D}, {E}
- ④ {A, B, C, D, E}

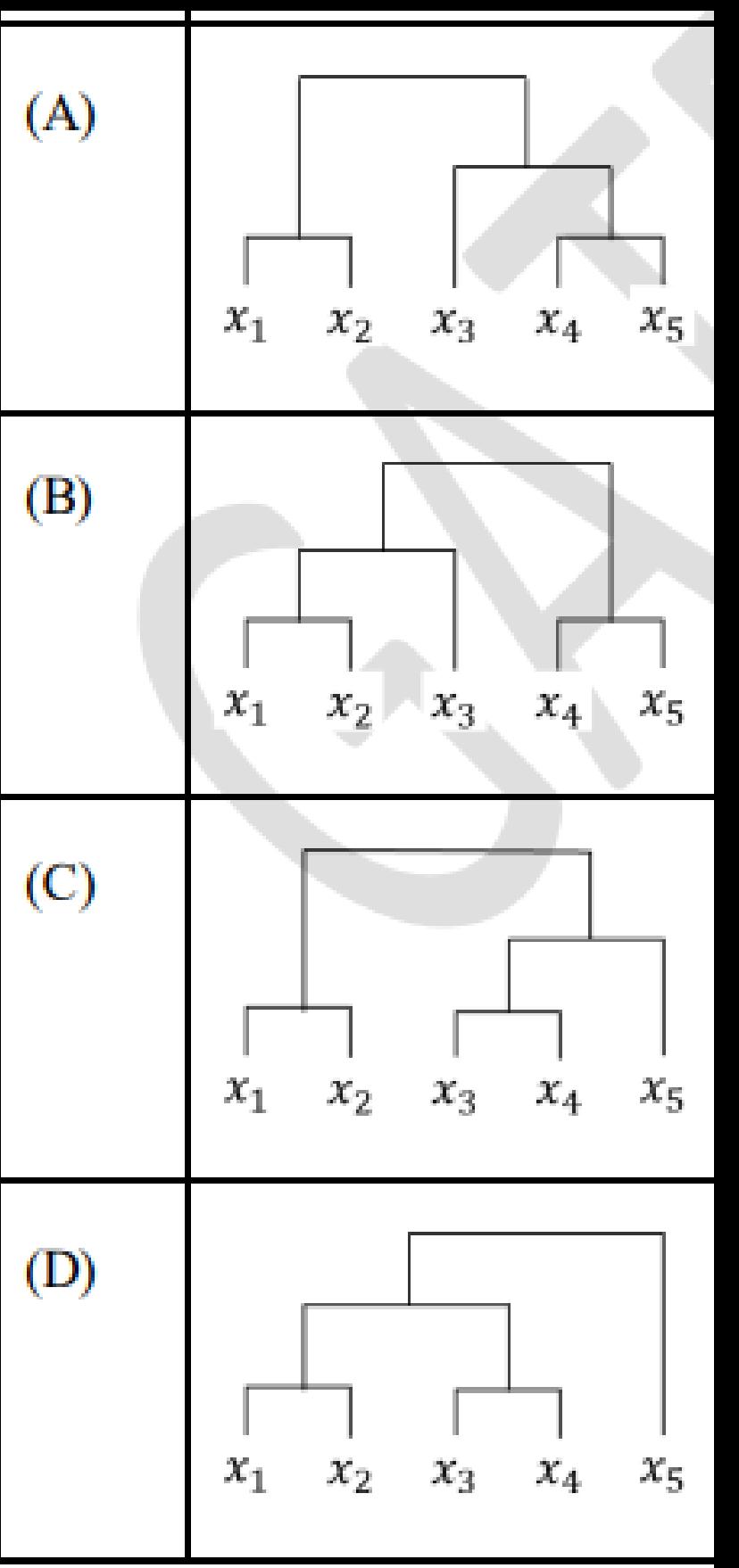


$$\frac{5+3}{2 \times 1} =$$

Consider the table below, where the  $(i, j)^{\text{th}}$  element of the table is the distance between points  $x_i$  and  $x_j$ . Single linkage clustering is performed on data points,  $x_1, x_2, x_3, x_4, x_5$ .

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$
$x_1$	0	1	4	3	6
$x_2$	1	0	3	5	3
$x_3$	4	3	0	2	5
$x_4$	3	5	2	0	1
$x_5$	6	3	5	1	0

Which **ONE** of the following is the correct representation of the clusters produced?



# Dimensionality reduction

## Two ways of Dimensionality Reduction

1. Feature Selection 
2. Feature Extraction

## Feature Selection – 3 Methods

### 1. Filter Method

- Correlation
- Chi-Square Test, etc.

### 2. Wrapper Method

- Forward Selection
- Backward Selection
- Bi-directional Elimination

### 3. Embedded Method

- LASSO
- Elastic Net
- Ridge Regression, etc.

## Feature Extraction

Feature extraction is the process of transforming the space containing many dimensions into space with fewer dimensions.

This approach is useful when we want to keep the whole information but use fewer resources while processing the information

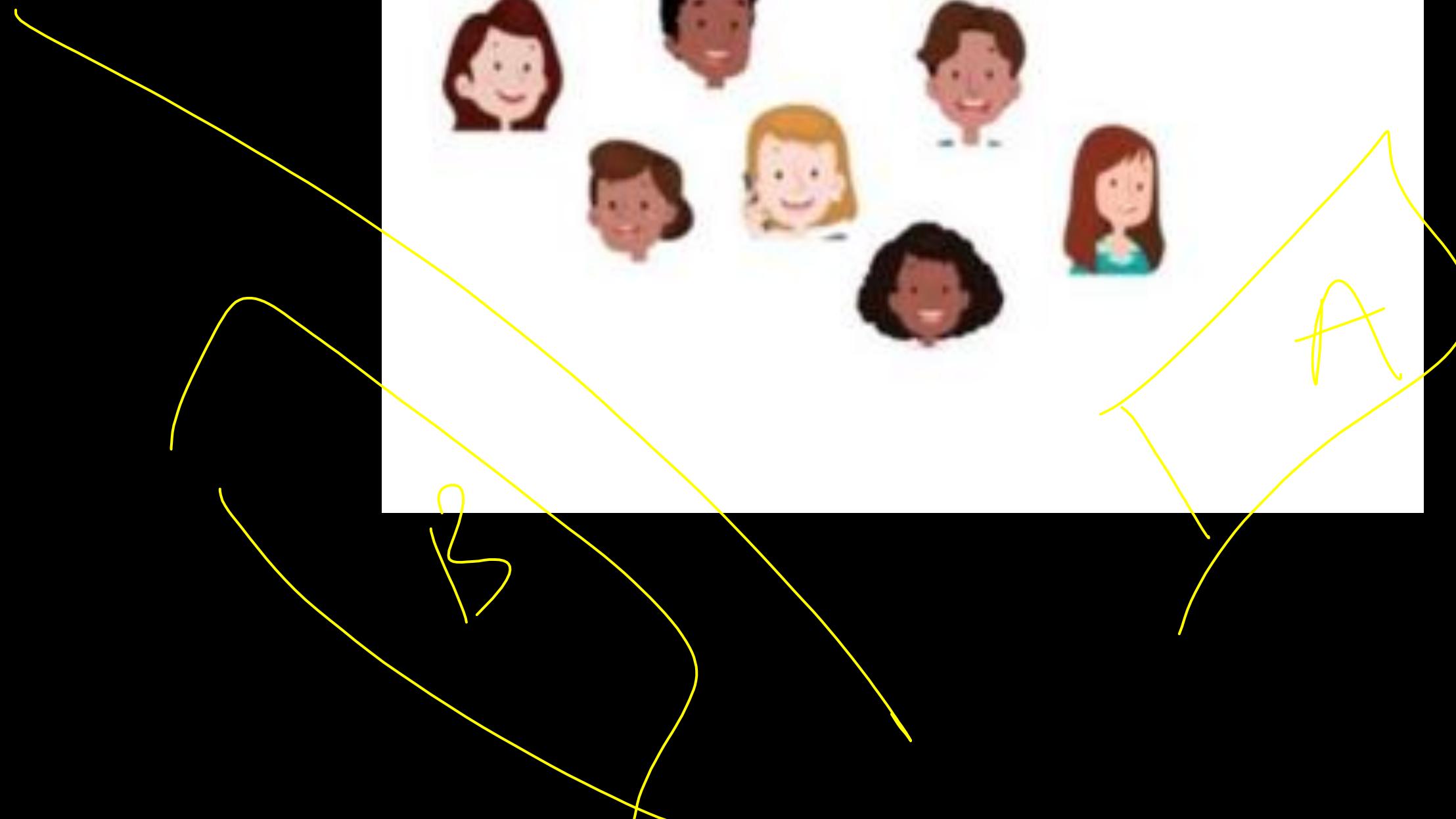
$$\begin{matrix} y \\ \downarrow \\ x \end{matrix}$$

Some common feature extraction techniques are:

1. Principal Component Analysis (PCA) 
2. Linear Discriminant Analysis (LDA) 
3. Kernel PCA
4. Quadratic Discriminant Analysis (QDA)etc

## PCA

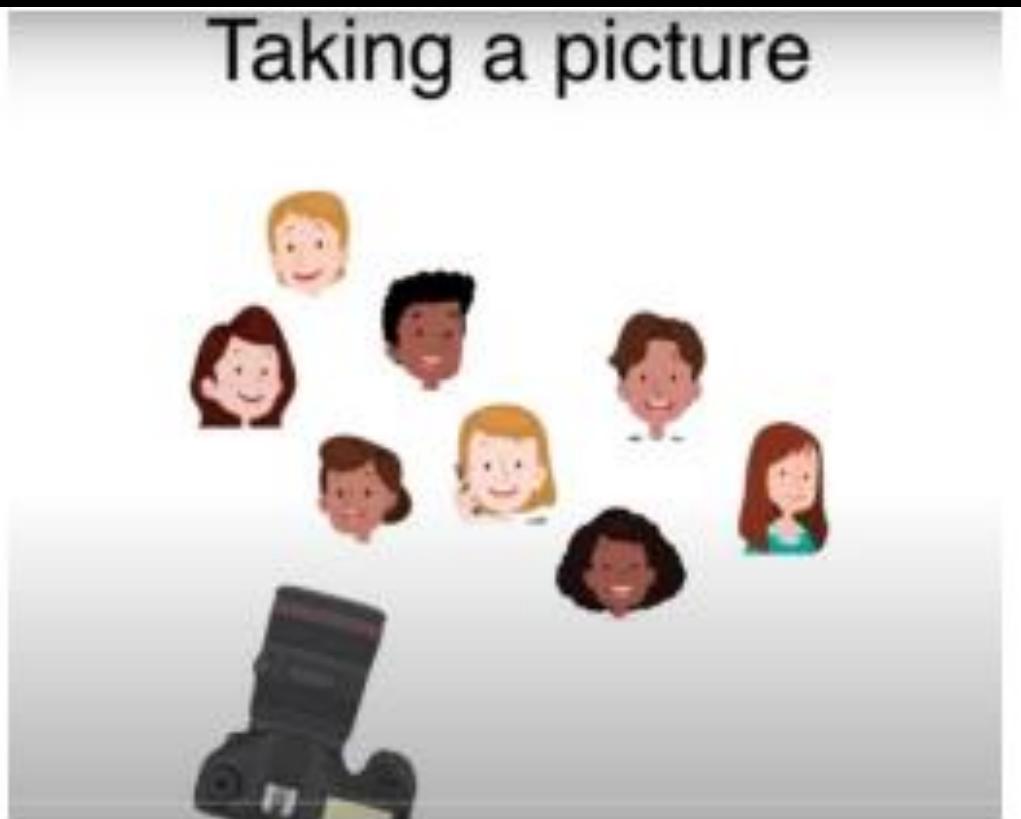
# Taking a picture



Taking a picture

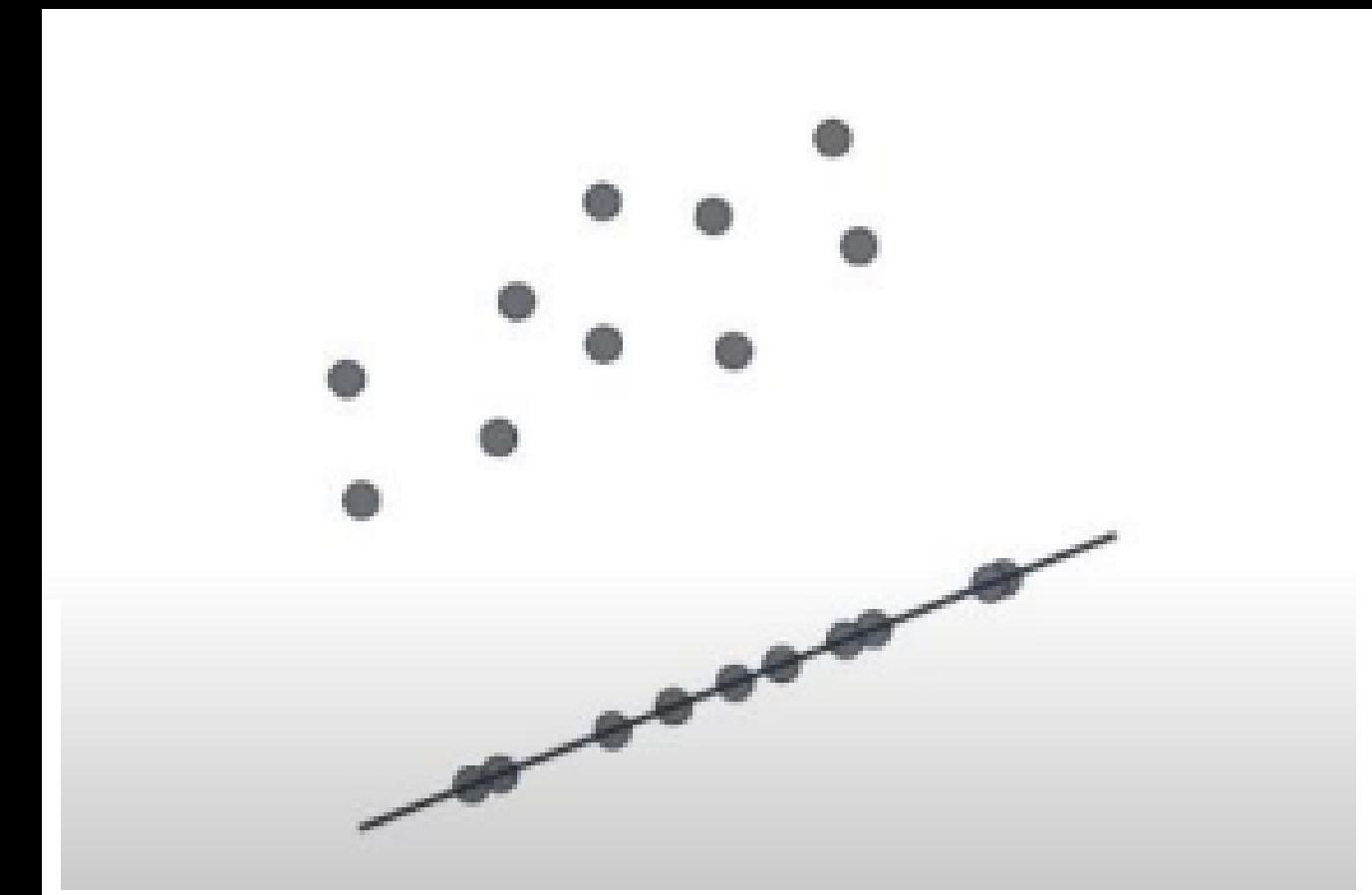
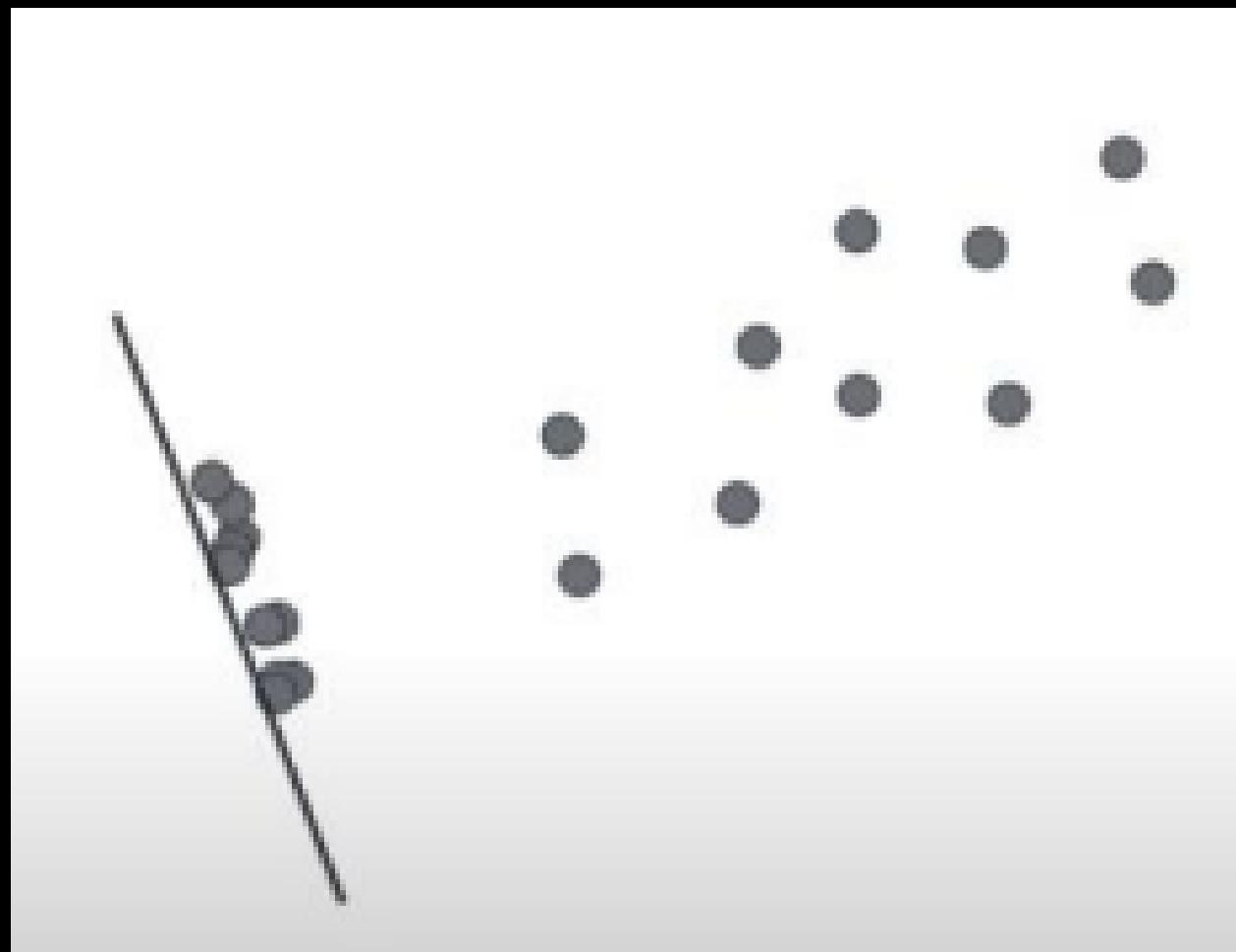


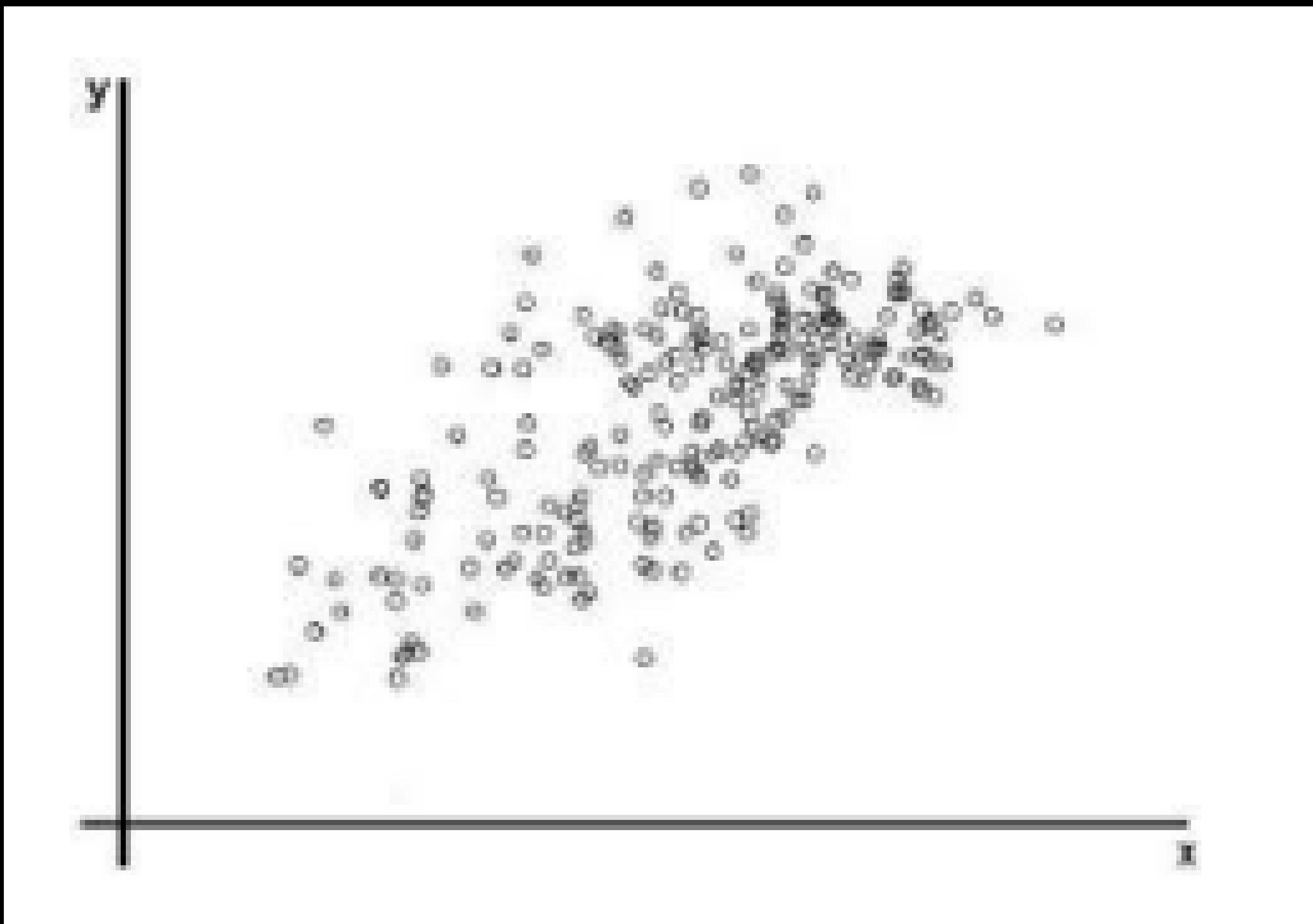
Taking a picture

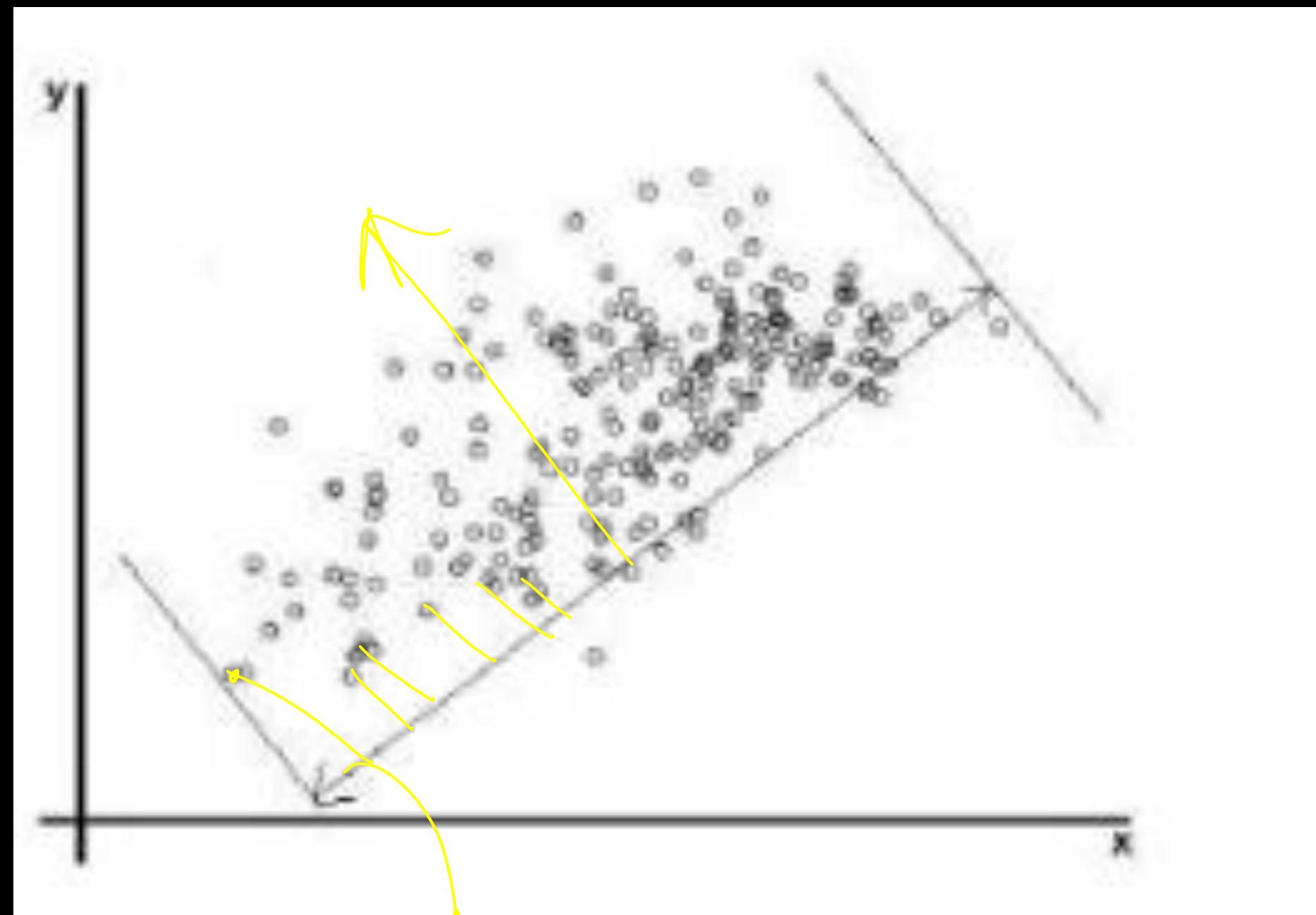


Taking a picture

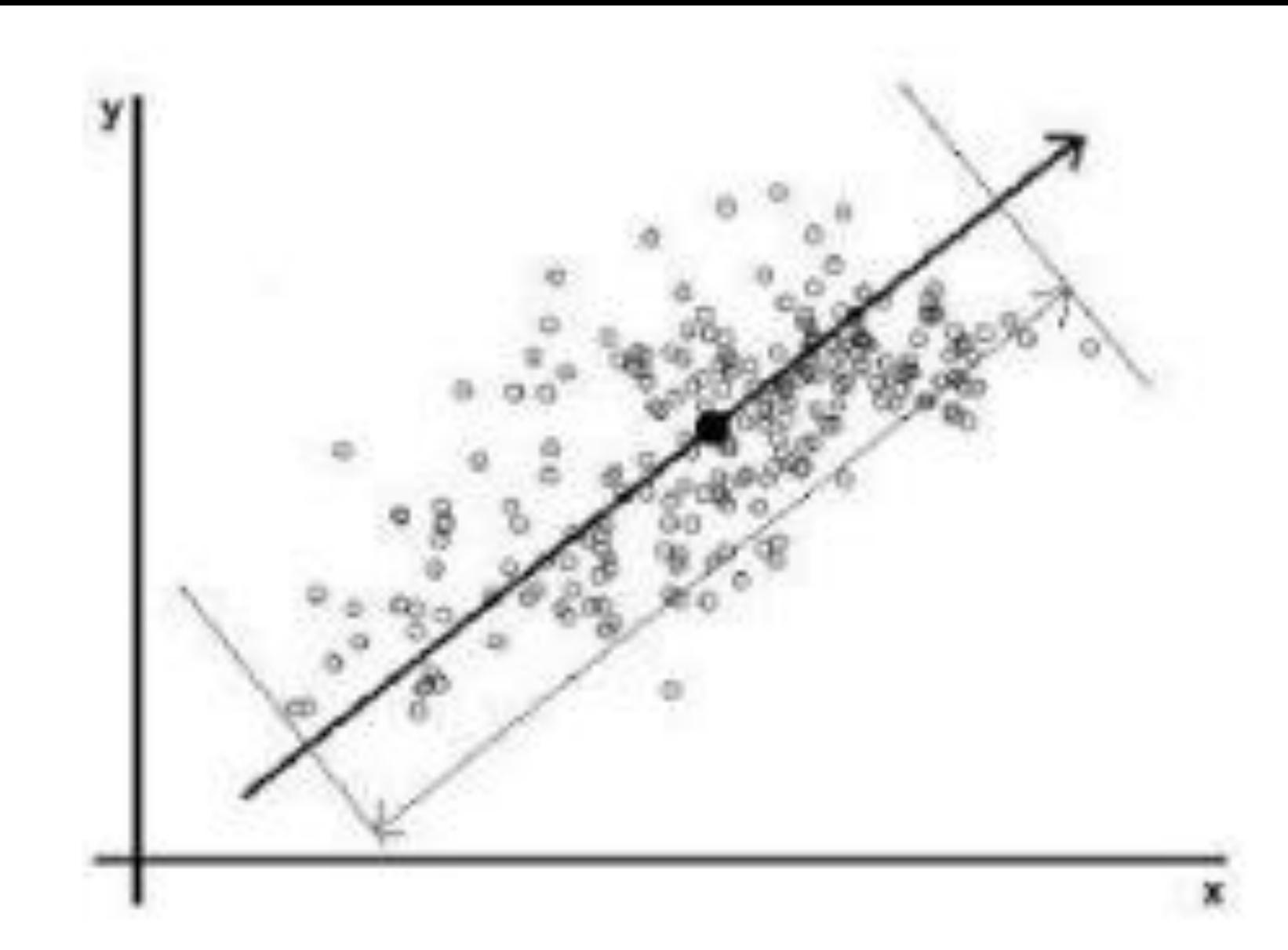


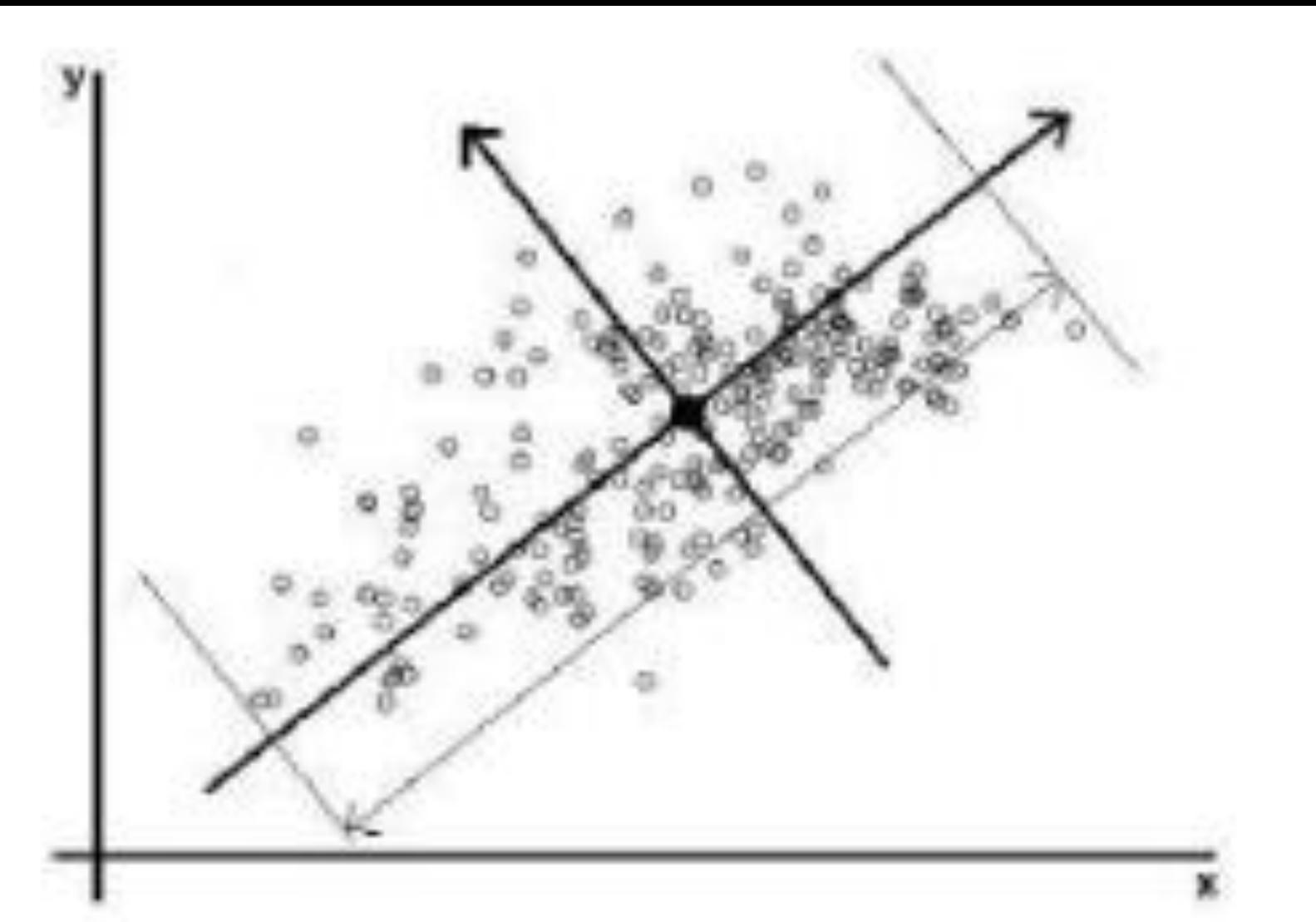


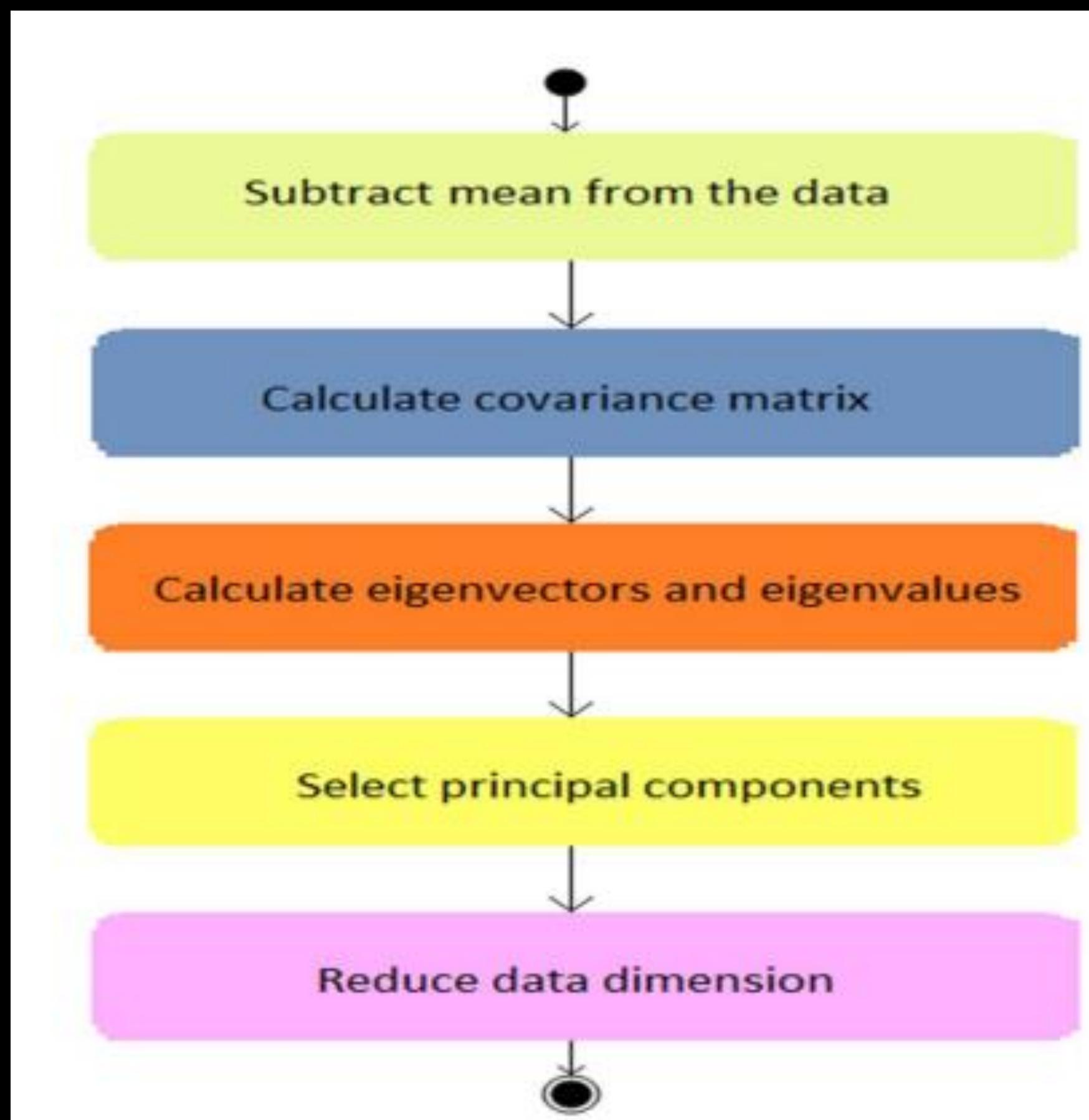




IPC8







Feature	Example 1	Example 2	Example 3	Example 4	
$X_1$	$x_{11}$	$x_{12}$	$x_{13}$	$x_{14}$	$\bar{x}_1 = 8$
$X_2$	$x_{21}$	$x_{22}$	$x_{23}$	$x_{24}$	$\bar{x}_2 = 8.5$

$$\text{Cov}(x_1, x_1) = \text{Var} x_1 = \frac{(x_i - \bar{x})^2}{n} = \frac{(4-8)^2 + (8-8)^2 + (13-8)^2 + (7-8)^2}{4} = 14$$

$$\text{Cov}(x_2, x_2) = 23$$

$$\text{Cov}(x_1, x_2) = \frac{1}{n-1} \sum (x_{1k} - \bar{x}_1)(x_{2k} - \bar{x}_2) = -11$$

$$\text{Cov} = \begin{bmatrix} \text{Cov}(x_1, x_1) & \text{Cov}(x_1, x_2) \\ \text{Cov}(x_2, x_1) & \text{Cov}(x_2, x_2) \end{bmatrix}$$

$$e_1^T \left( \begin{bmatrix} 4-8 & 8-8 & 13-8 & 7-8 \\ 11-8.5 & 4-8.5 & 5-8.5 & 14-8.5 \end{bmatrix} \right)$$

$$\text{Cov} = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}$$

Eigen vector  
↓  
unit vector  $\rightarrow e_1$

The most popularly used dimensionality reduction algorithm is Principal Component Analysis (PCA). Which of the following is/are true about PCA?

1. PCA is an unsupervised method of optimization
2. It searches for the directions that data have the largest variance
3. Maximum number of principal components  $\leq$  number of features
4. All principal components are orthogonal to each other

Answer: 1,2,3&4

What will happen when eigenvalues are roughly equal?

- A. PCA will perform outstandingly
- B. PCA will perform badly
- C. Can't say
- D. None of above

**Solution: (B)**

What happens when you get features in lower dimensions using PCA?

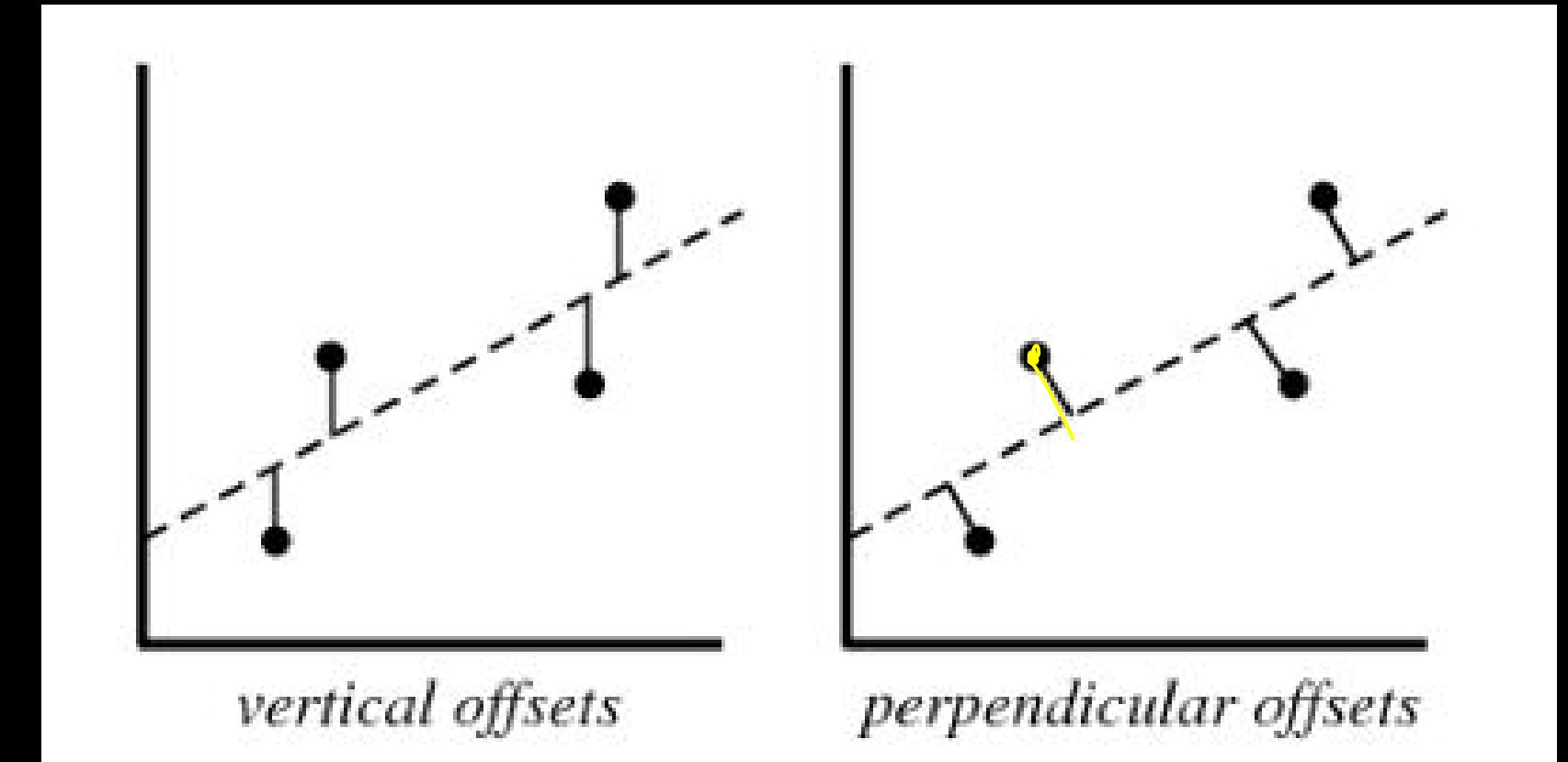
1. The features will still have interpretability
2. The features will lose interpretability
3. The features must carry all information present in the data
4. The features may not carry all information present in the data

Answer: 2 & 4

Which of the following offset do we consider in PCA?

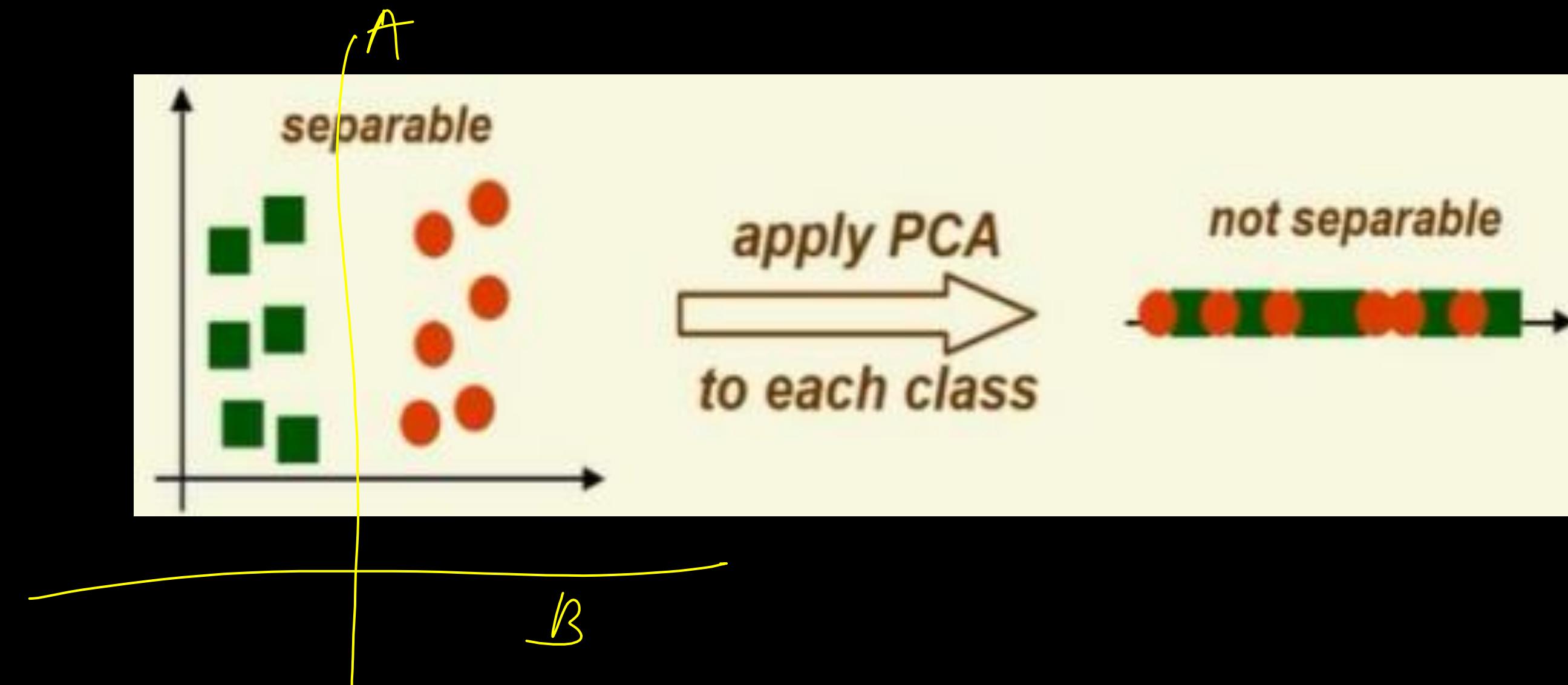
- A. Vertical offset
- B. Perpendicular offset
- C. Both
- D. None of these

**Solution: (B)**



PCA projects the data in the directions of maximum variance.

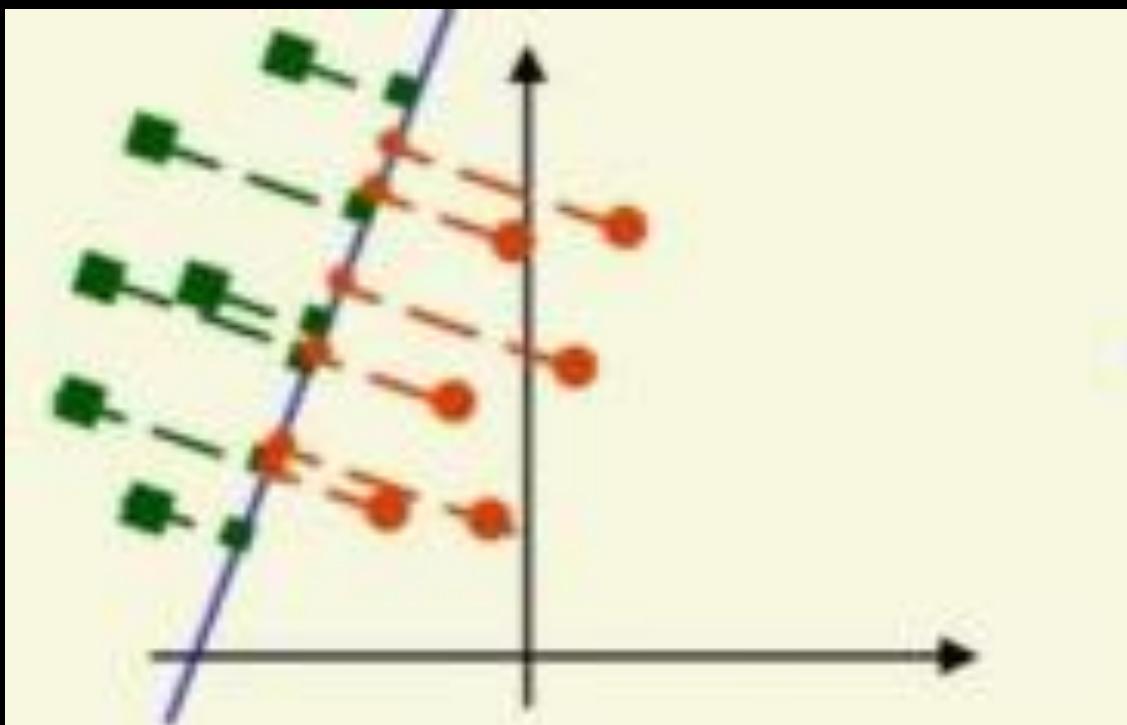
However the directions of maximum variance may be useless for classification



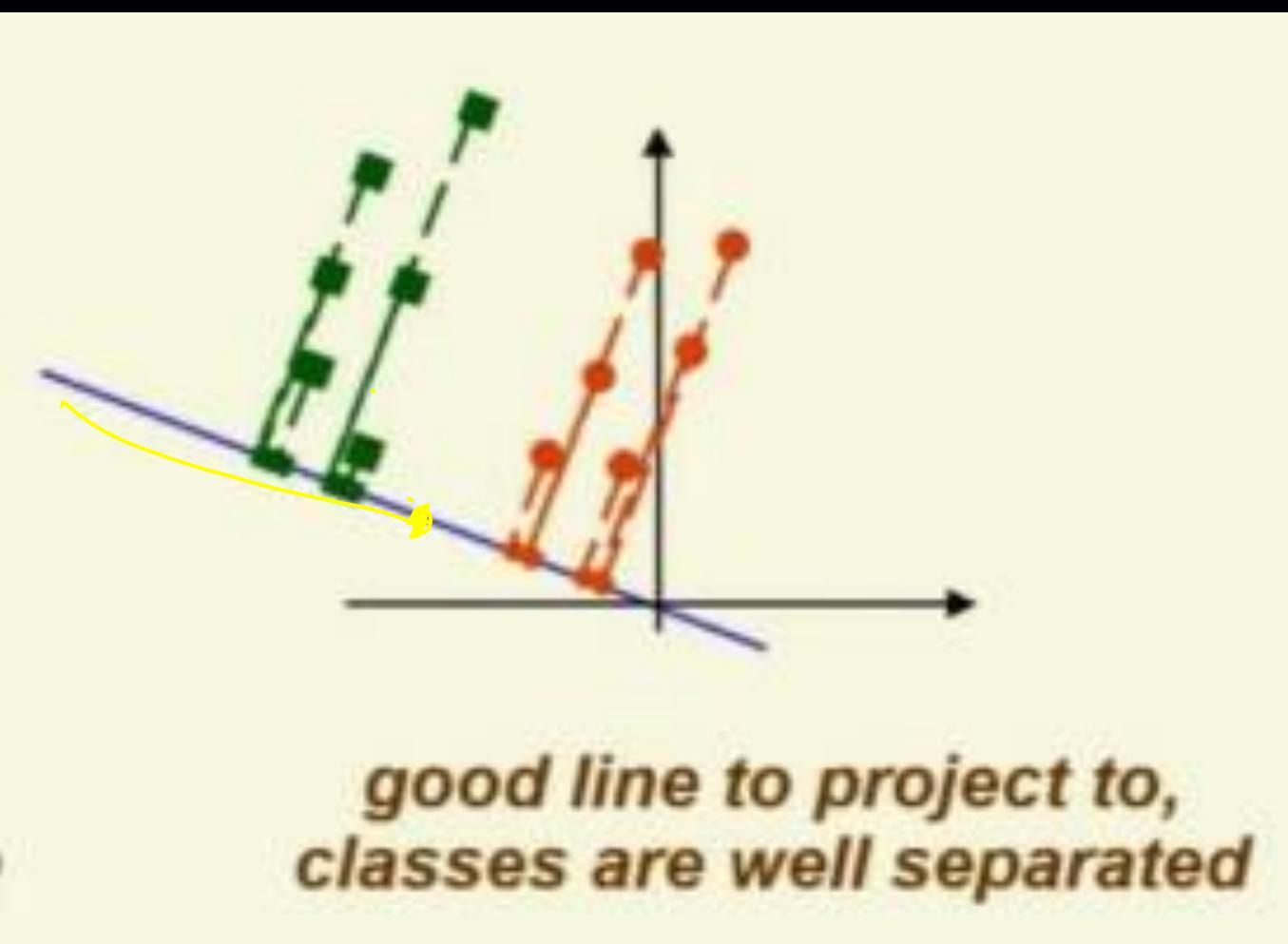
In such condition LDA which is also called as Fisher LDA works well.

## Fisher Linear Discriminant

- Main idea: find projection to a line s.t. samples from different classes are well separated

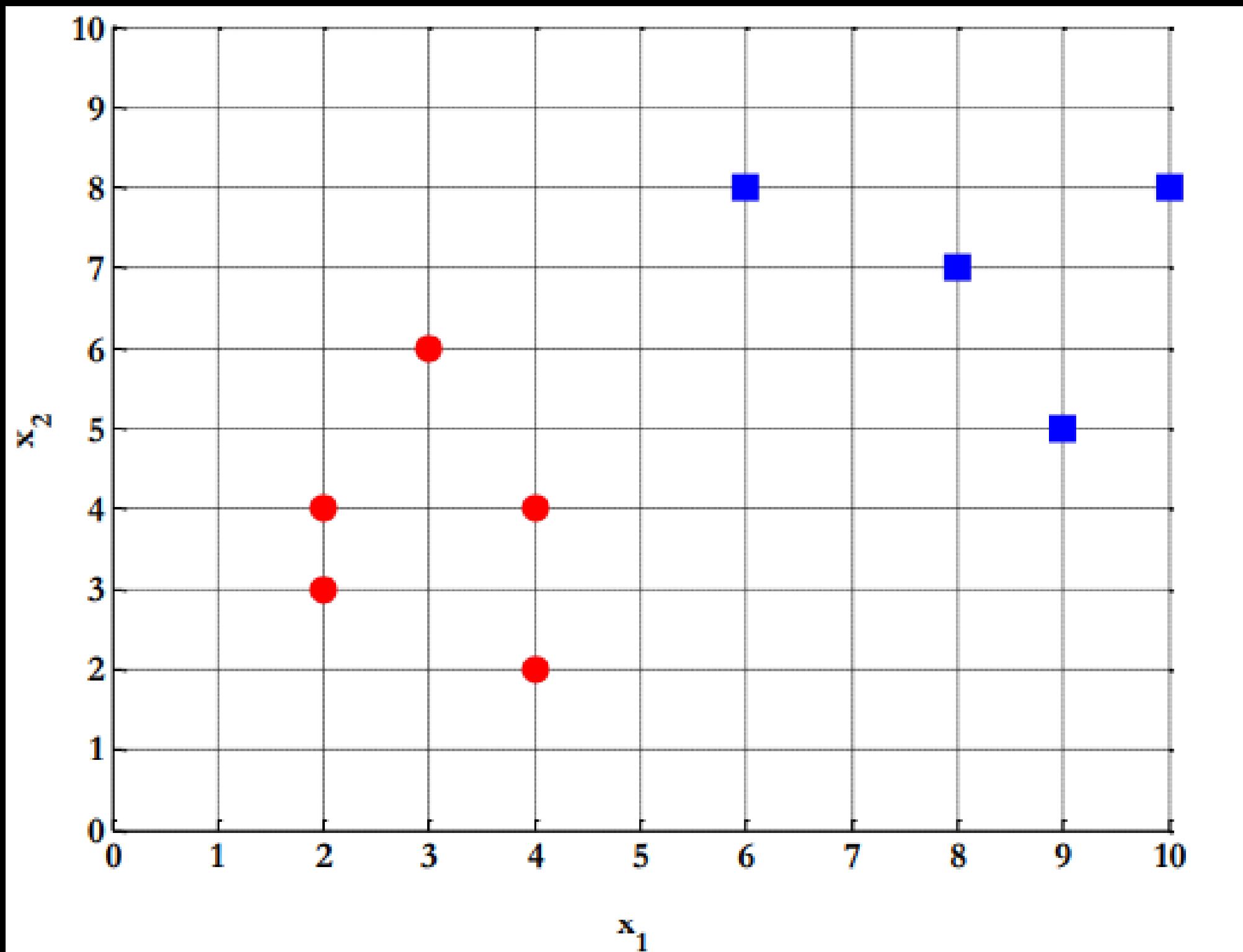


*bad line to project to,  
classes are mixed up*



Samples for class  $\omega_1$  :  $\mathbf{X}_1 = (\mathbf{x}_1, \mathbf{x}_2) = \{(4,2), (2,4), (2,3), (3,6), \underline{(4,4)}\}$

Sample for class  $\omega_2$  :  $\mathbf{X}_2 = (\mathbf{x}_1, \mathbf{x}_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$



The classes mean are :

$$\mu_1 = \frac{1}{N_1} \sum_{x \in \omega_1} x = \frac{1}{5} \left[ \binom{4}{2} + \binom{2}{4} + \binom{2}{3} + \binom{3}{6} + \binom{4}{4} \right] = \begin{pmatrix} 3 \\ 3.8 \end{pmatrix}$$

$$\mu_2 = \frac{1}{N_2} \sum_{x \in \omega_2} x = \frac{1}{5} \left[ \binom{9}{10} + \binom{6}{8} + \binom{9}{5} + \binom{8}{7} + \binom{10}{8} \right] = \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix}$$

Covariance matrix of the first class:

$$\begin{aligned} S_1 &= \sum_{x \in \omega_1} (x - \mu_1)(x - \mu_1)^T = \left[ \begin{pmatrix} 4 \\ 2 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 + \left[ \begin{pmatrix} 2 \\ 4 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 \\ &\quad + \left[ \begin{pmatrix} 2 \\ 3 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 + \left[ \begin{pmatrix} 3 \\ 6 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 + \left[ \begin{pmatrix} 4 \\ 4 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 \\ &= \begin{pmatrix} 1 & -0.25 \\ -0.25 & 2.2 \end{pmatrix} \end{aligned}$$

Covariance matrix of the second class:

$$\begin{aligned} S_2 &= \sum_{x \in \omega_2} (x - \mu_2)(x - \mu_2)^T = \left[ \begin{pmatrix} 9 \\ 10 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^2 + \left[ \begin{pmatrix} 6 \\ 8 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^2 \\ &\quad + \left[ \begin{pmatrix} 9 \\ 5 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^2 + \left[ \begin{pmatrix} 8 \\ 7 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^2 + \left[ \begin{pmatrix} 10 \\ 8 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^2 \\ &= \begin{pmatrix} 2.3 & -0.05 \\ -0.05 & 3.3 \end{pmatrix} \end{aligned}$$

Within-class scatter matrix:

$$\begin{aligned} S_w &= S_1 + S_2 = \begin{pmatrix} 1 & -0.25 \\ -0.25 & 2.2 \end{pmatrix} + \begin{pmatrix} 2.3 & -0.05 \\ -0.05 & 3.3 \end{pmatrix} \\ &= \begin{pmatrix} 3.3 & -0.3 \\ -0.3 & 5.5 \end{pmatrix} \end{aligned}$$

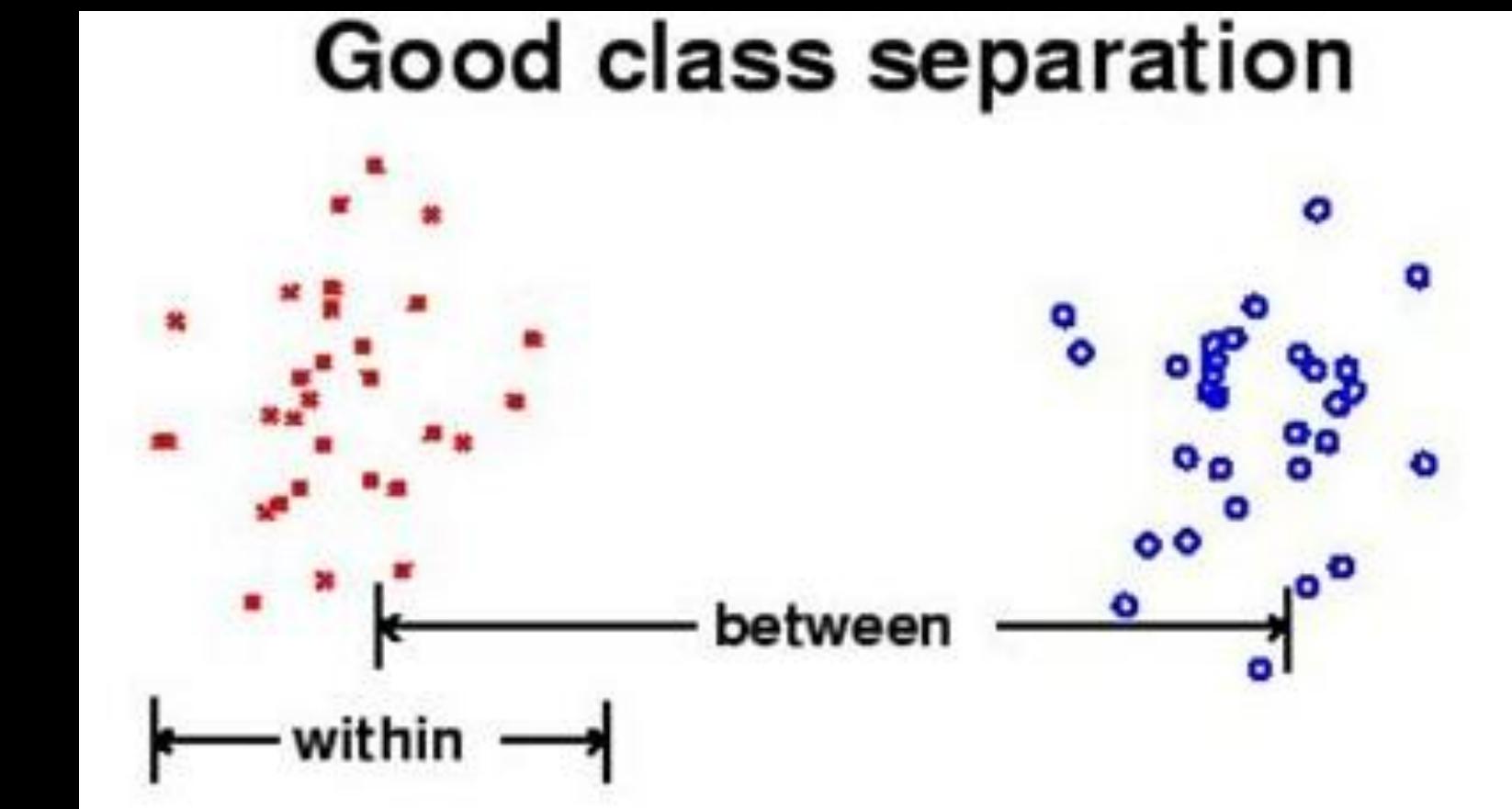
$$\begin{aligned}w^* &= S_W^{-1}(\mu_1 - \mu_2) = \begin{pmatrix} 3.3 & -0.3 \\ -0.3 & 5.5 \end{pmatrix}^{-1} \left[ \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right] \\&= \begin{pmatrix} 0.3045 & 0.0166 \\ 0.0166 & 0.1827 \end{pmatrix} \begin{pmatrix} -5.4 \\ -3.8 \end{pmatrix} \\&= \begin{pmatrix} 0.9088 \\ 0.4173 \end{pmatrix}\end{aligned}$$

Samples for class  $\omega_1$  :  $\mathbf{X}_1 = (x_1, x_2) = \{(4,2), (2,4), (2,3), (3,6), (4,4)\}$

Sample for class  $\omega_2$  :  $\mathbf{X}_2 = (x_1, x_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$

$x_1$	4	2	2	3	4	9	6	9	8	10
$x_2$	2	4	3	6	4	10	8	5	7	8
1 <sup>st</sup> LD	4.46	3.48	3.06	5.2	5.3	12.35	8.8	10.2	10.19	12.42

Which of the following is true about LDA?



- A. LDA aims to maximize the distance between classes and minimize the within-class distance.
- B. LDA aims to minimize both distances between classes and the distance within the class.
- C. LDA aims to minimize the distance between classes and maximize the distance within the class.
- D. LDA aims to maximize both distances between classes and the distance within the class.

**Solution: (A)**

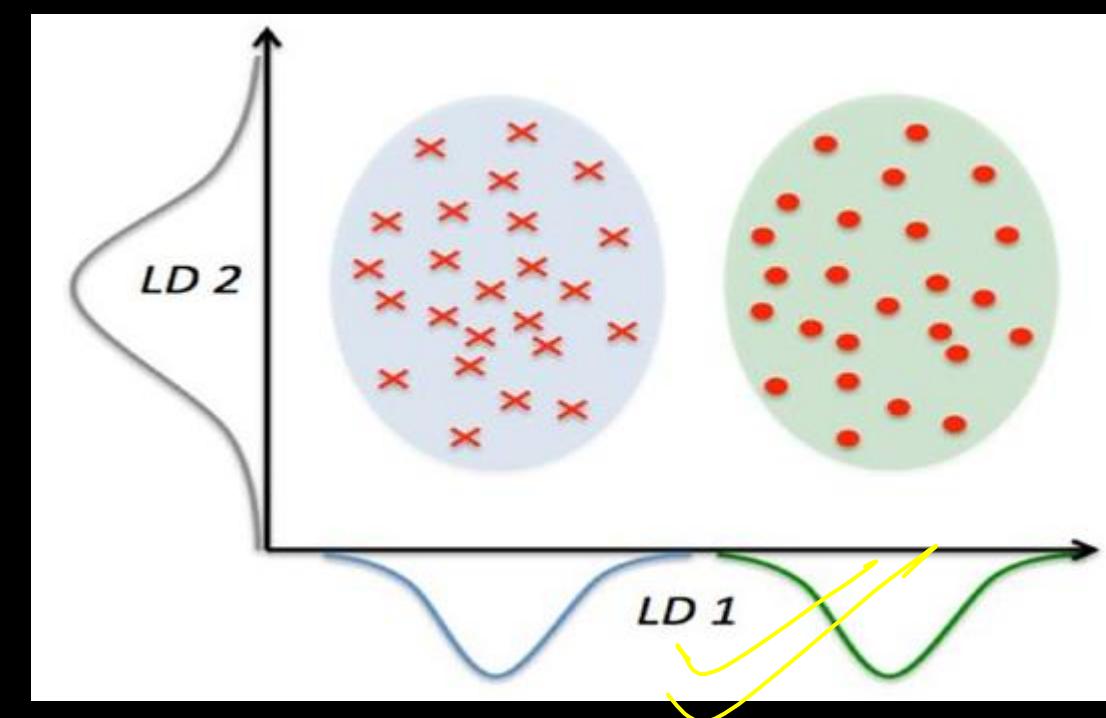
Which of the following comparison(s) are true about PCA and LDA?

1. Both LDA and PCA are linear transformation techniques
2. LDA uses supervised learning, whereas PCA uses unsupervised learning
3. PCA maximizes the variance of the data, whereas LDA maximizes the separation between different classes.
4. None of these

Answer 1,2,&3

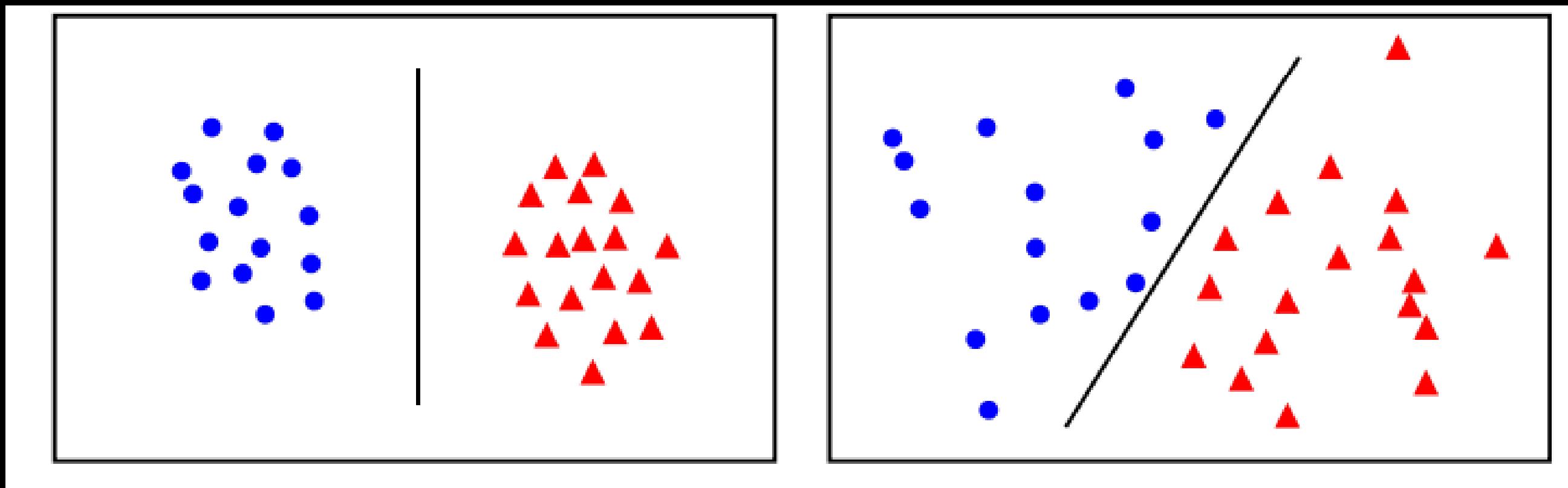
In LDA, the idea is to find the line that best separates the two classes. In the given image, which of the following is a good projection?

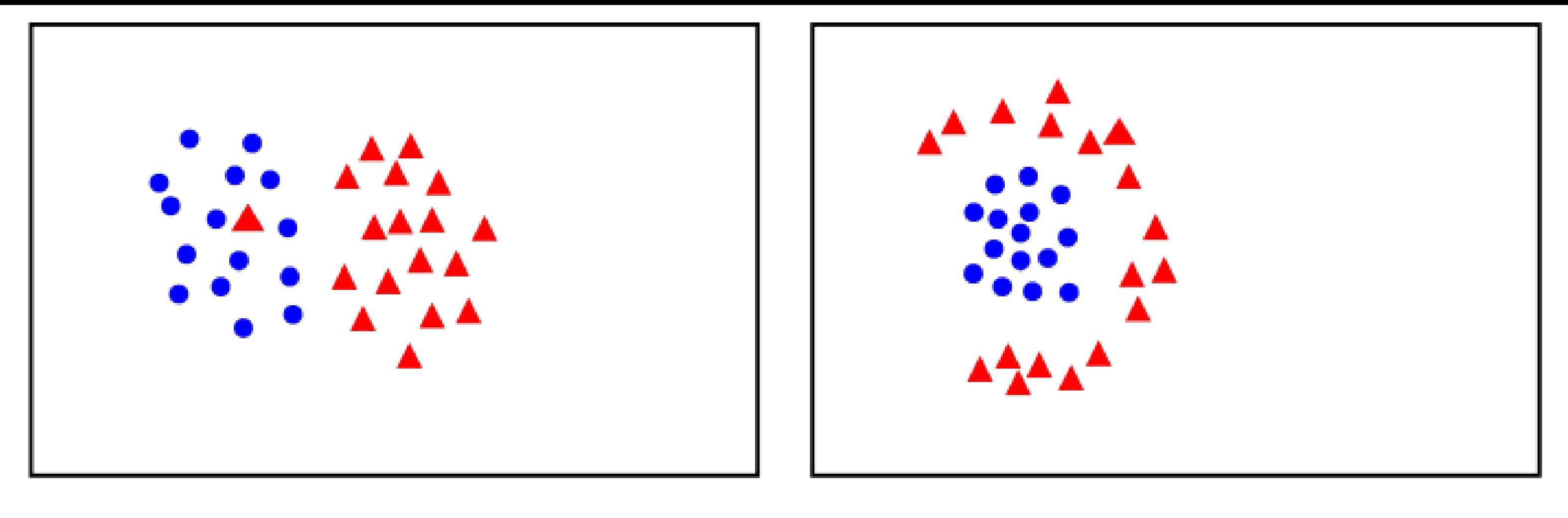
- A. LD1
- B. LD2
- C. Both
- D. None of these

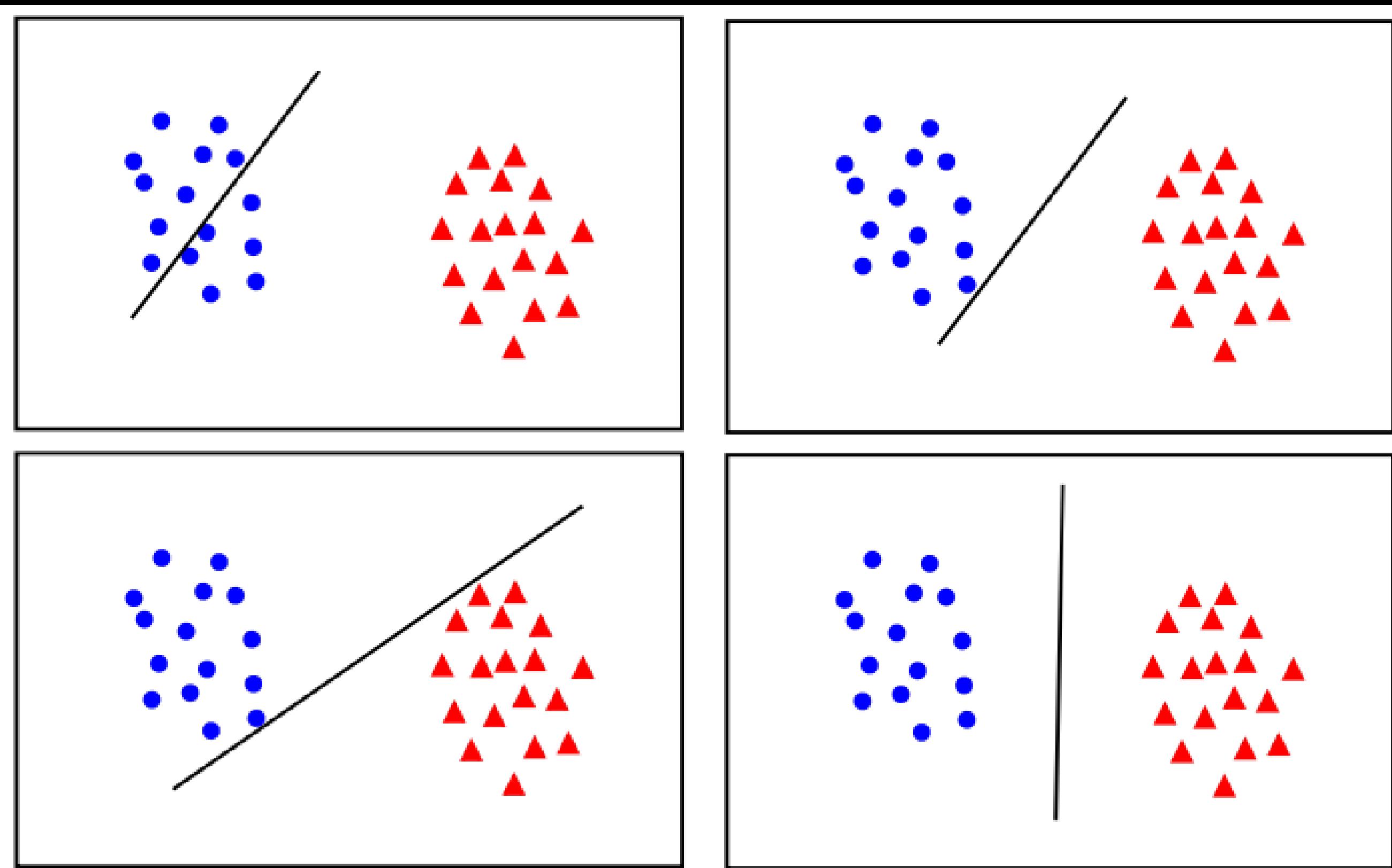


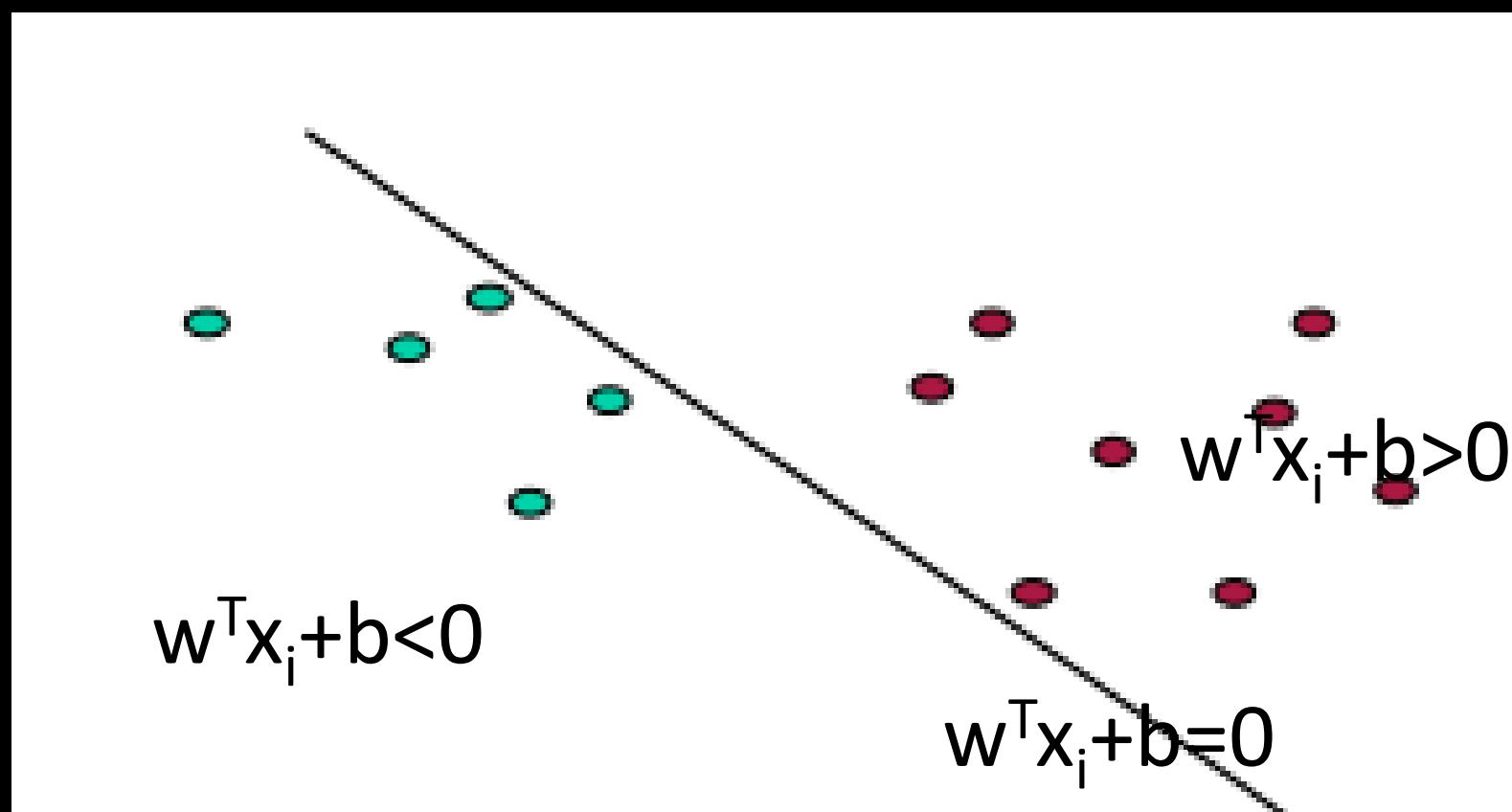
**Solution:** (A)

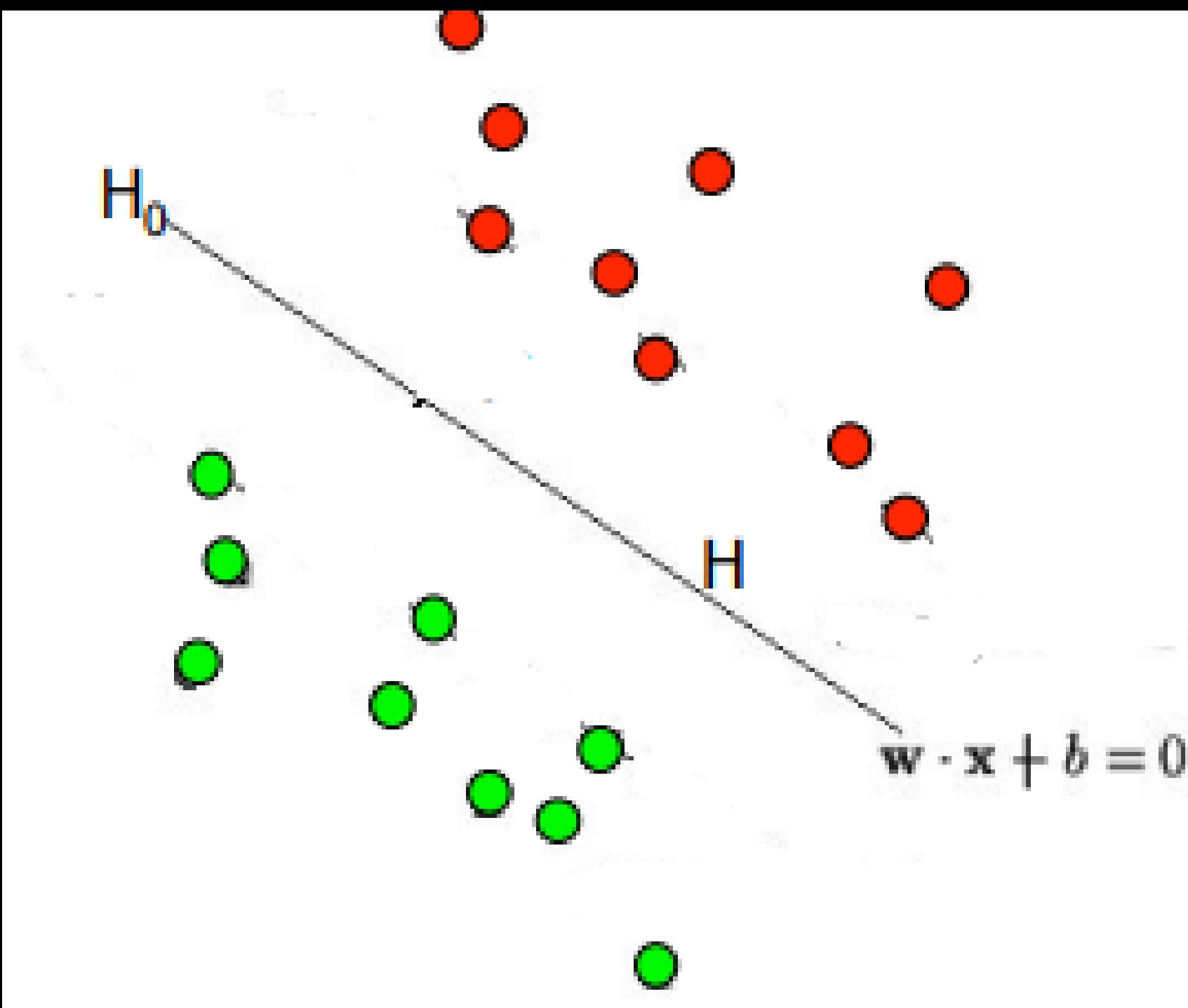
# SVM

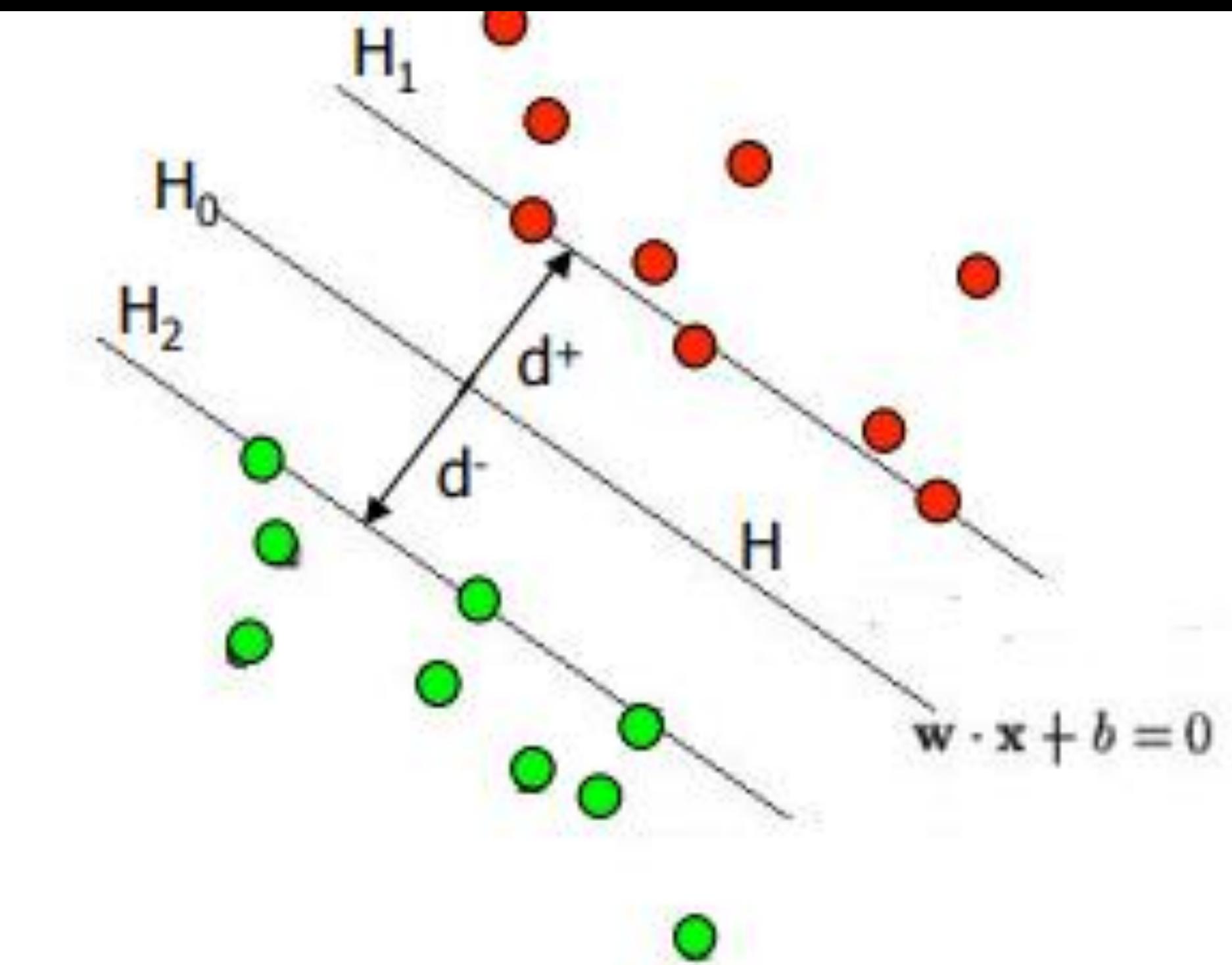


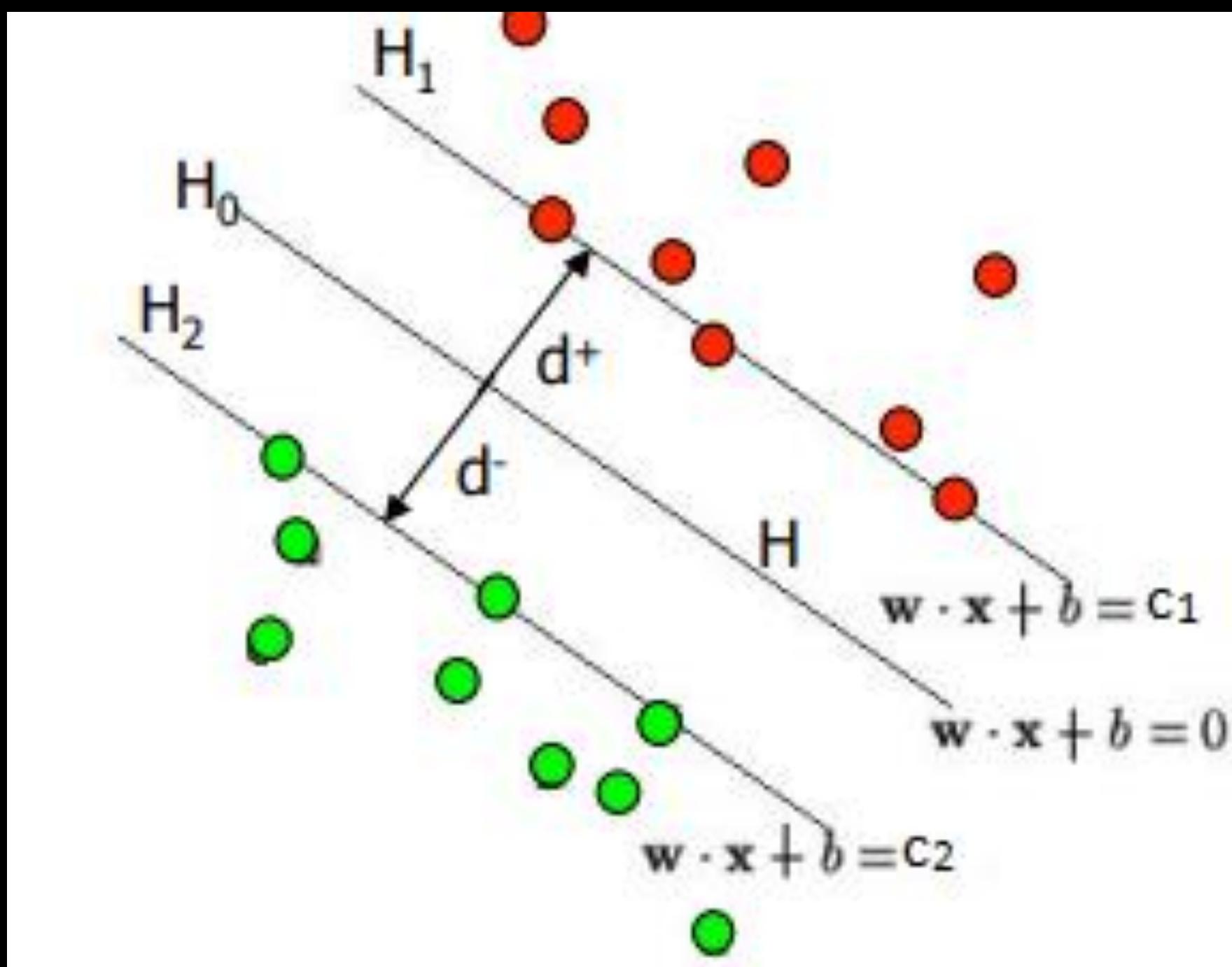


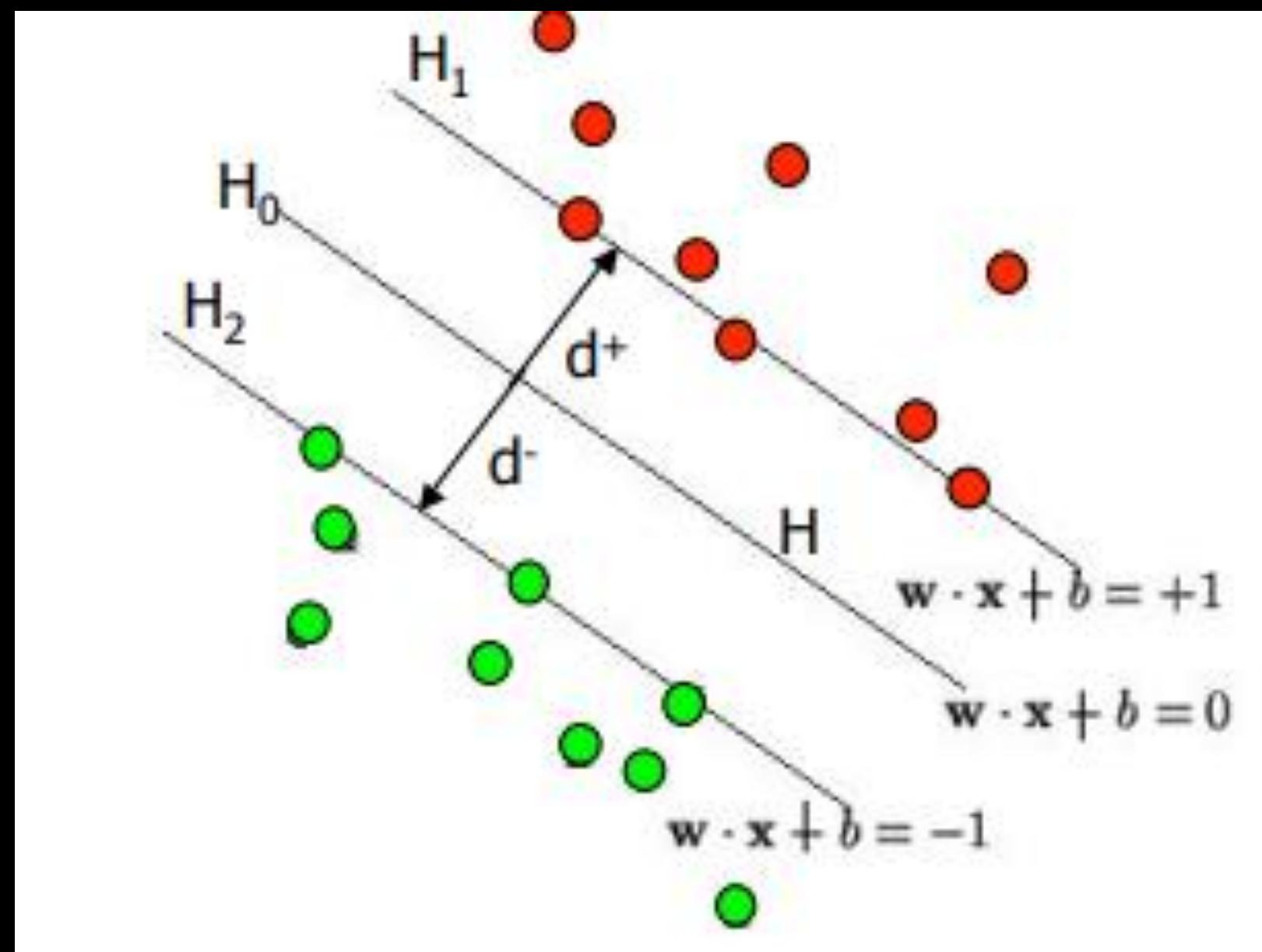






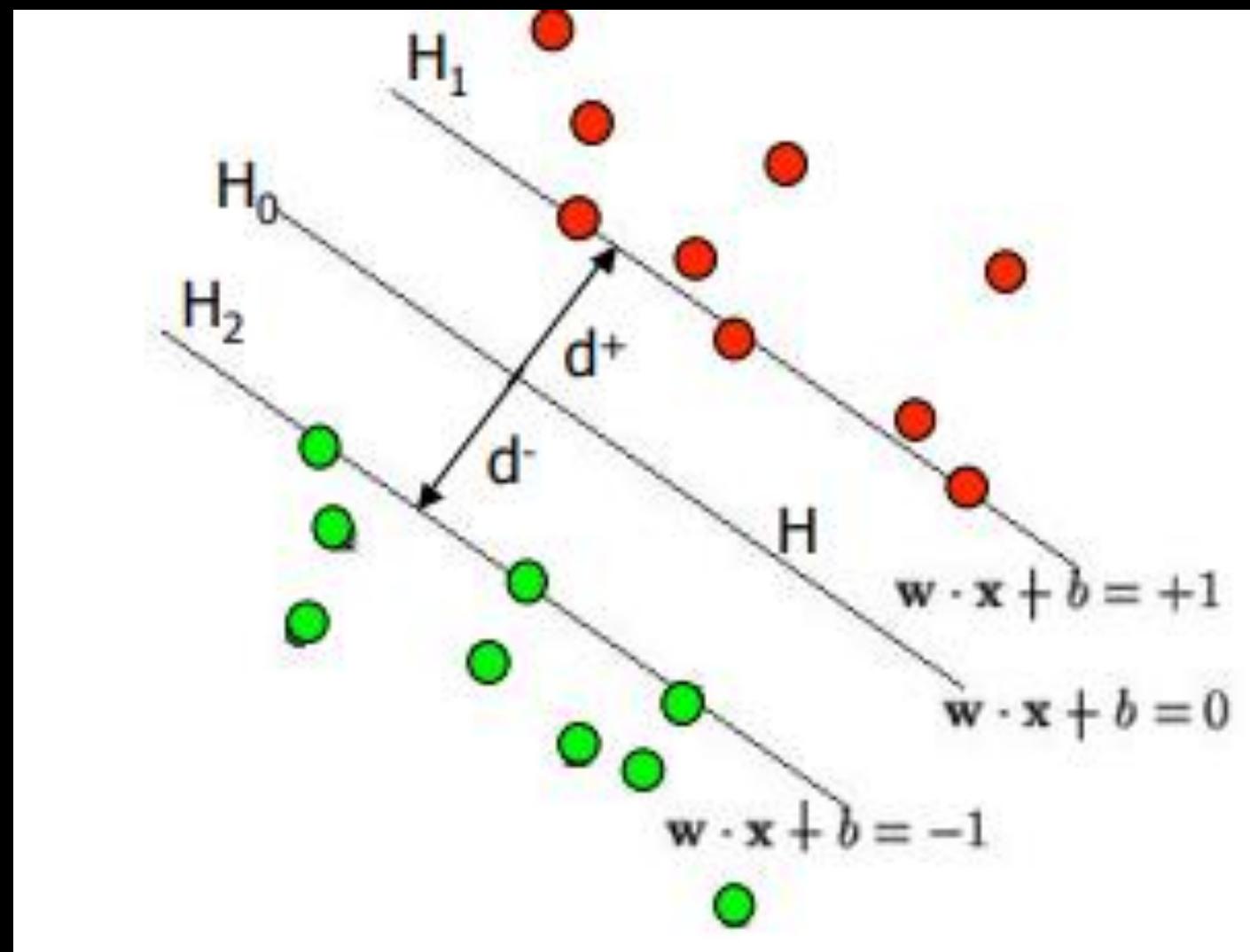






Recall the distance from a point  $(x_0, y_0)$  to a line:  $Ax + By + c = 0$  is:

$$\frac{|Ax_0 + By_0 + c|}{\sqrt{A^2 + B^2}}$$



The distance between  $H_0$  and  $H_1$  is then:

$$\frac{|wx + b|}{||w||} = \frac{1}{||w||}$$

The distance between H2 and H1 is  $\frac{2}{||w||}$

## Optimization Problem of SVM

To summarize, the SVM optimization problem can be written as:

Objective:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

Subject to the constraints:

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \quad \forall i$$

$$\begin{aligned} & \underset{\mathbf{w}, b}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{w}\|^2 \\ & \text{subject to} \quad y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1 \geq 0, \quad 1 \leq i \leq n. \end{aligned}$$

The vector  $\vec{w}$  and the scalar  $b$  are given by

$$\vec{w} = \sum_{i=1}^N \alpha_i y_i \vec{x}_i$$
$$b = \frac{1}{2} \left( \min_{i:y_i=+1} (\vec{w} \cdot \vec{x}_i) + \max_{i:y_i=-1} (\vec{w} \cdot \vec{x}_i) \right)$$

where  $\vec{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_N)$  is a vector which maximizes

$$\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^N \alpha_i \alpha_j y_i y_j (\vec{x}_i \cdot \vec{x}_j)$$

subject to

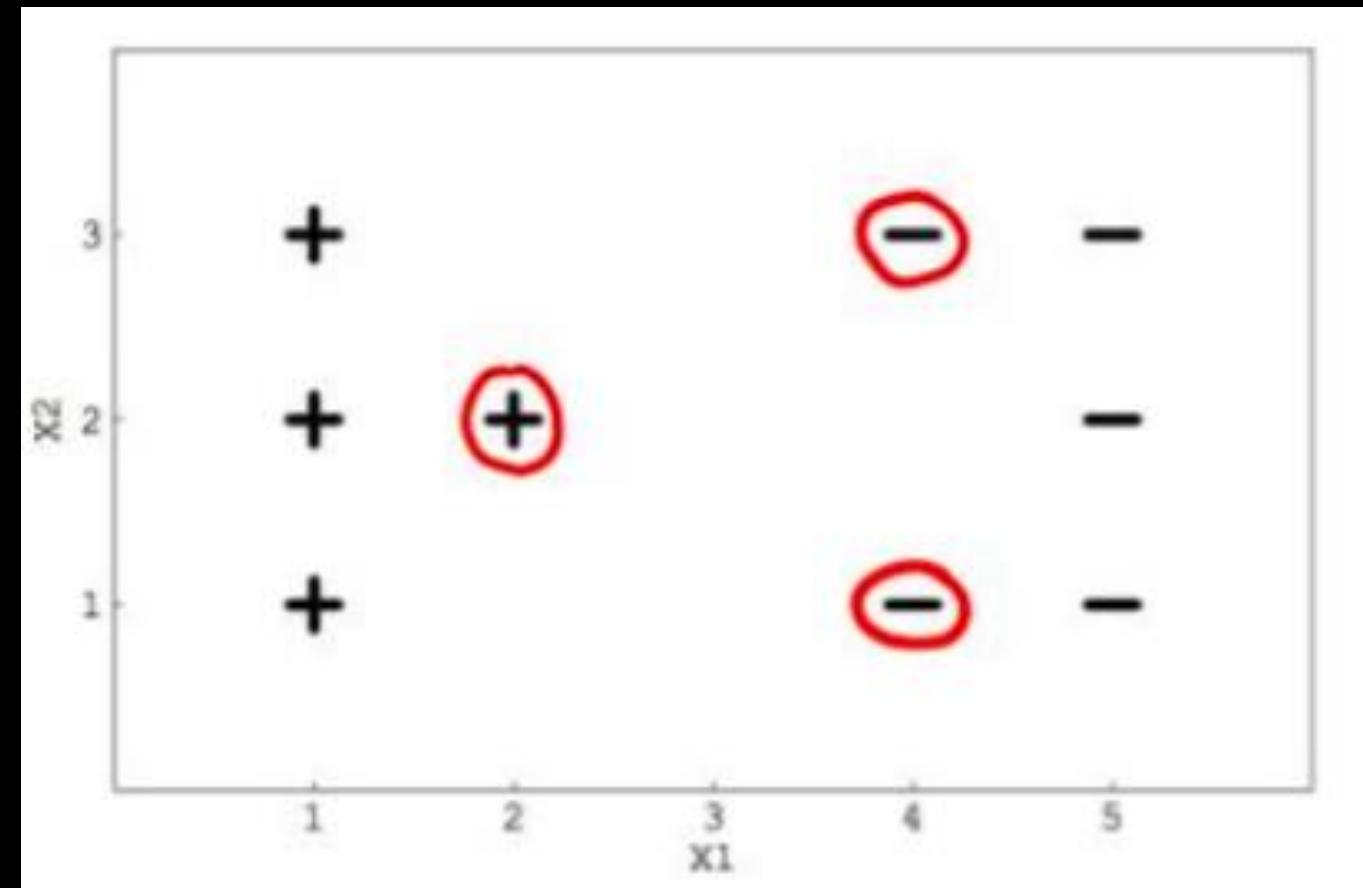
$$\sum_{i=1}^N \alpha_i y_i = 0$$

$$\alpha_i > 0 \text{ for } i = 1, 2, \dots, N.$$

Using the SVM algorithm, find the SVM classifier for the following data.

Example no.	$x_1$	$x_2$	Class
1	2	1	+1
2	4	3	-1

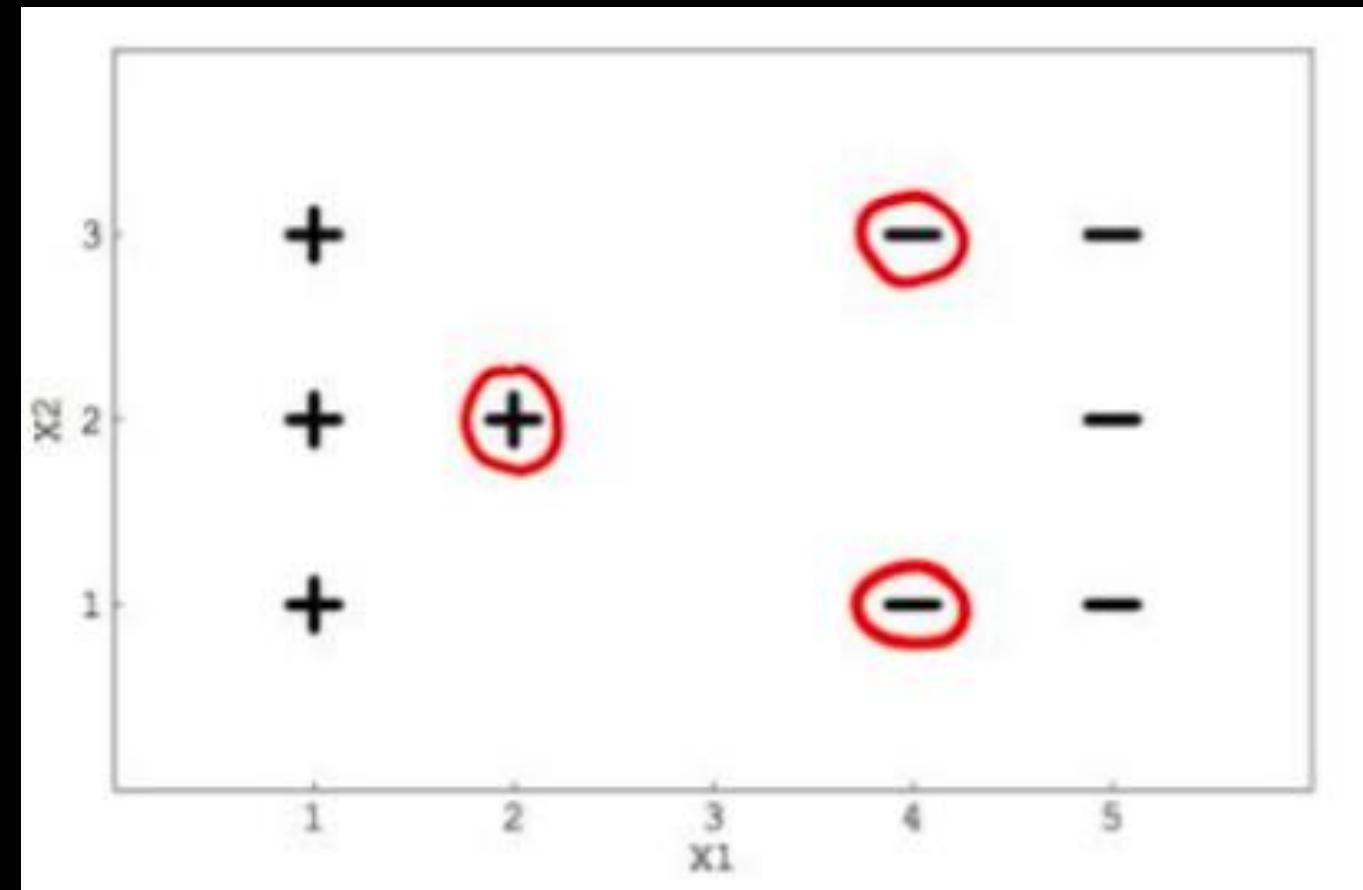
Suppose you are using a Linear SVM classifier with 2 class classification problem. Consider the following data in which the points circled red represent support vectors.



Will the decision boundary change if any of the red points are removed?

- a. Yes
- b. No

Suppose you are using a Linear SVM classifier with 2 class classification problem. Consider the following data in which the points circled red represent support vectors.

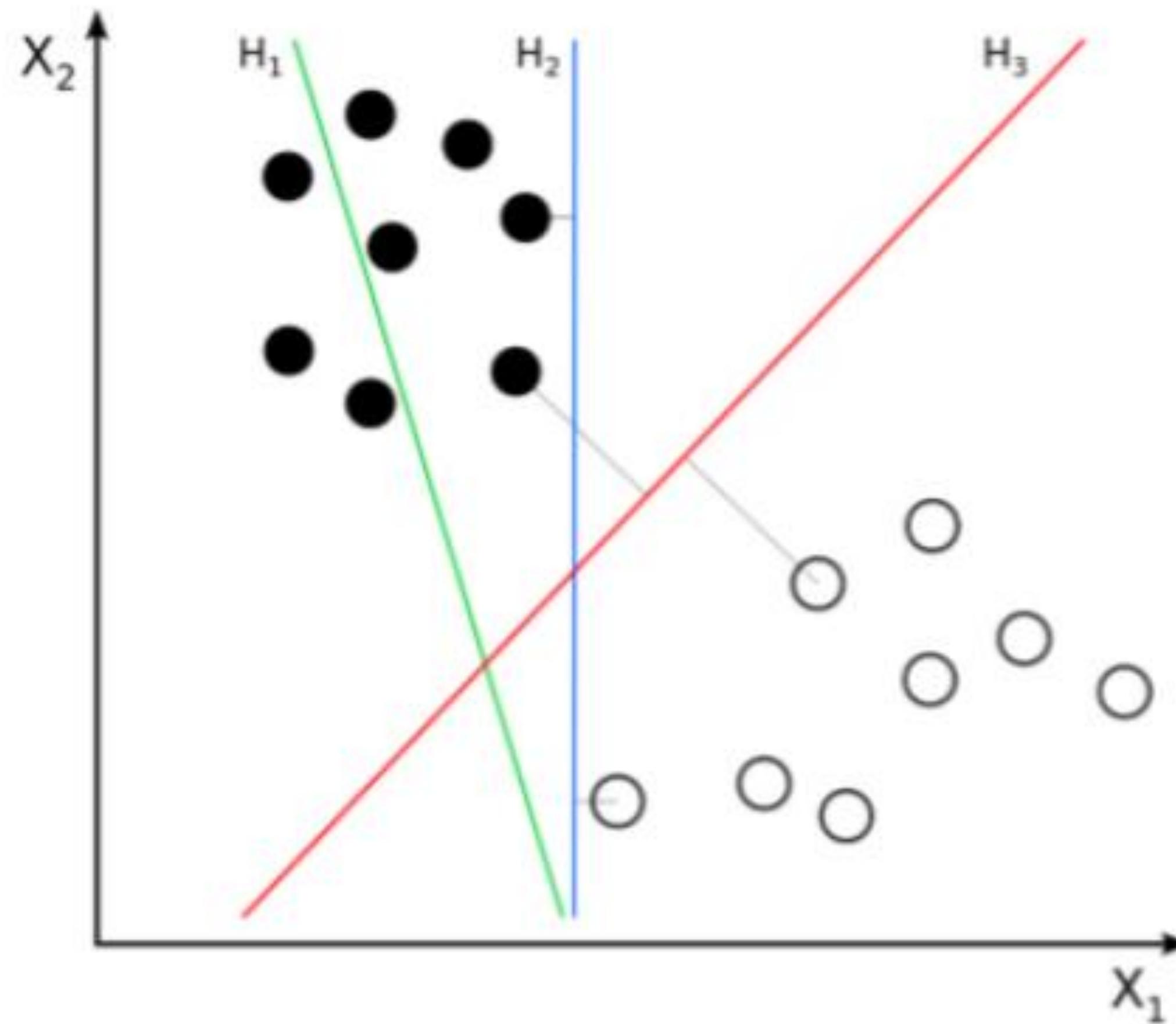


If you remove the non-red circled points from the data, the decision boundary will change?

- A) True
- B) False

**Solution:** B

Consider the data-points in the figure below.



Let us assume that the black-colored circles represent positive class whereas the white-colored circles represent negative class. Which of the following among H1, H2 and H3 is the maximum-margin hyperplane?

- (a) H1
- (b) H2
- (c) H3

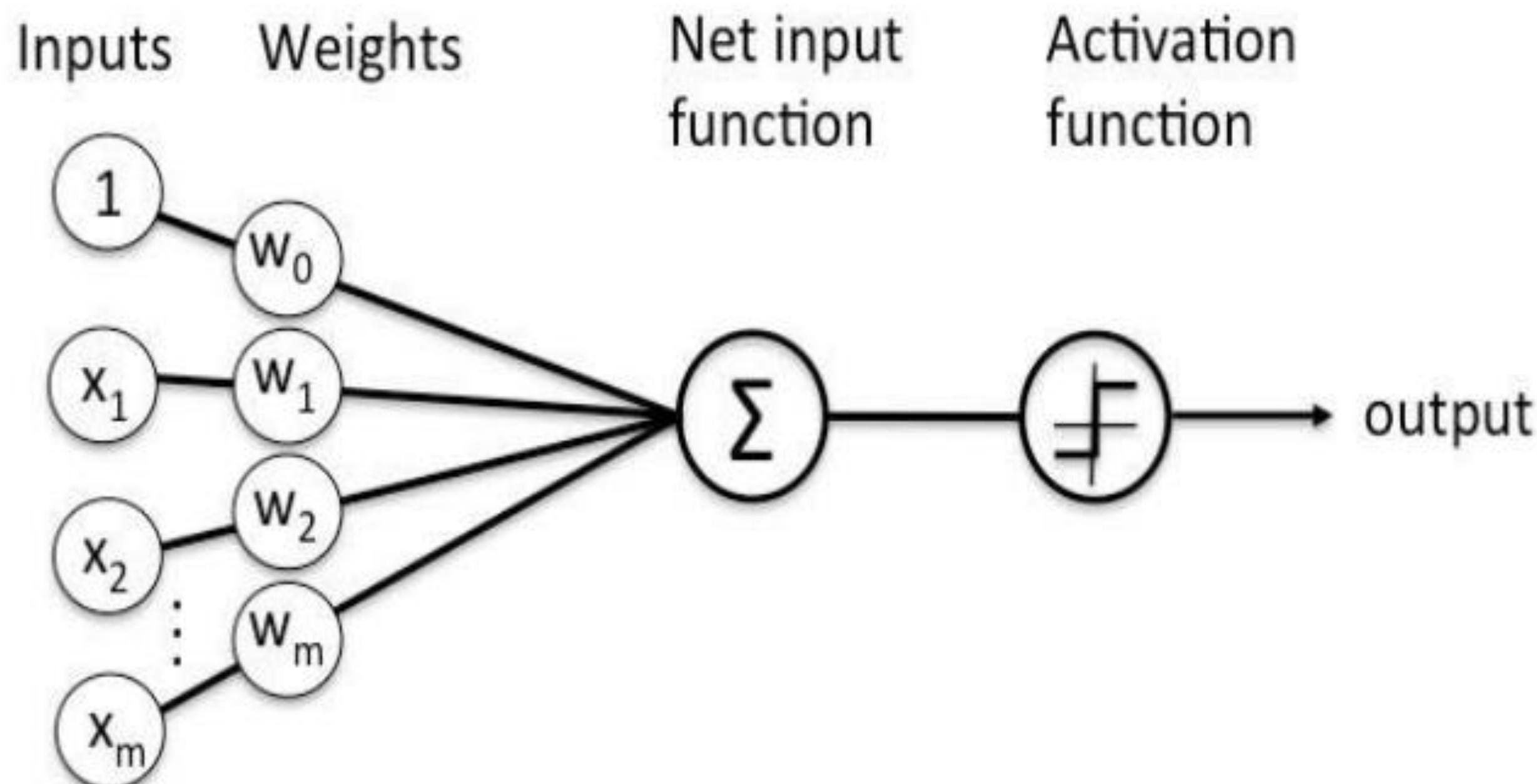
The soft margin SVM is more preferred than the hard-margin svm when:

1. The data is linearly separable
2. The data is noisy and contains overlapping point



# Perceptron

# Perceptron



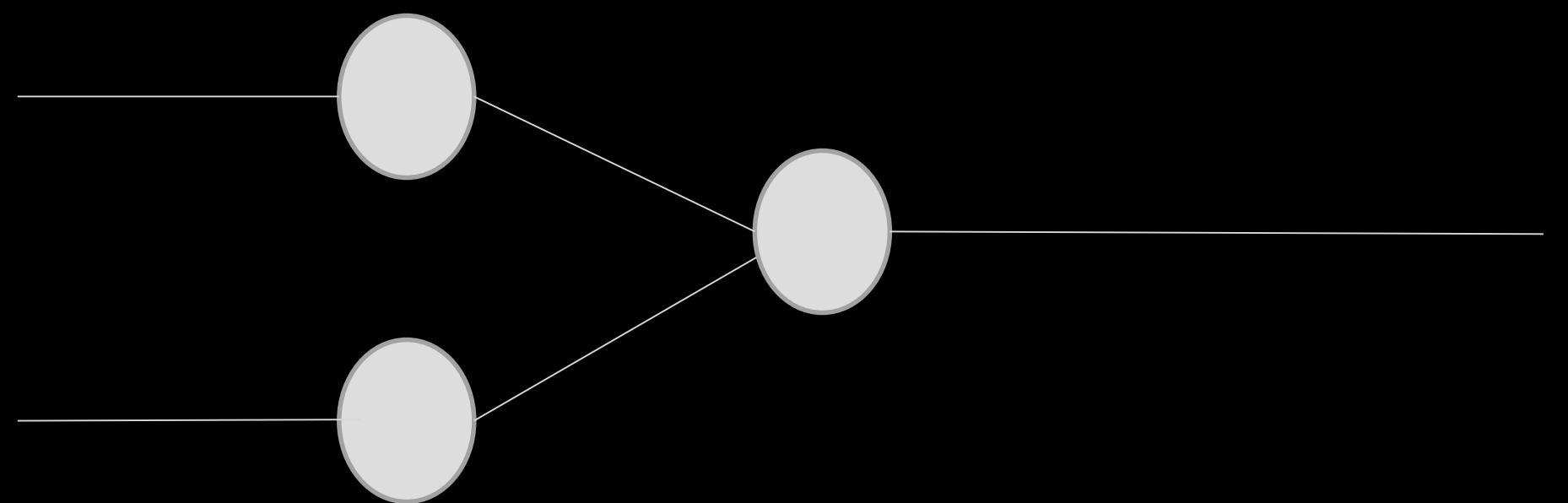
# Forward Propagation

## Activation Function

1. Sigmoid Function
2. Threshold Activation Function
3. Unit Step Function
4. Relu
5. Linear Activation Function
6. Hyperbolic tangent function

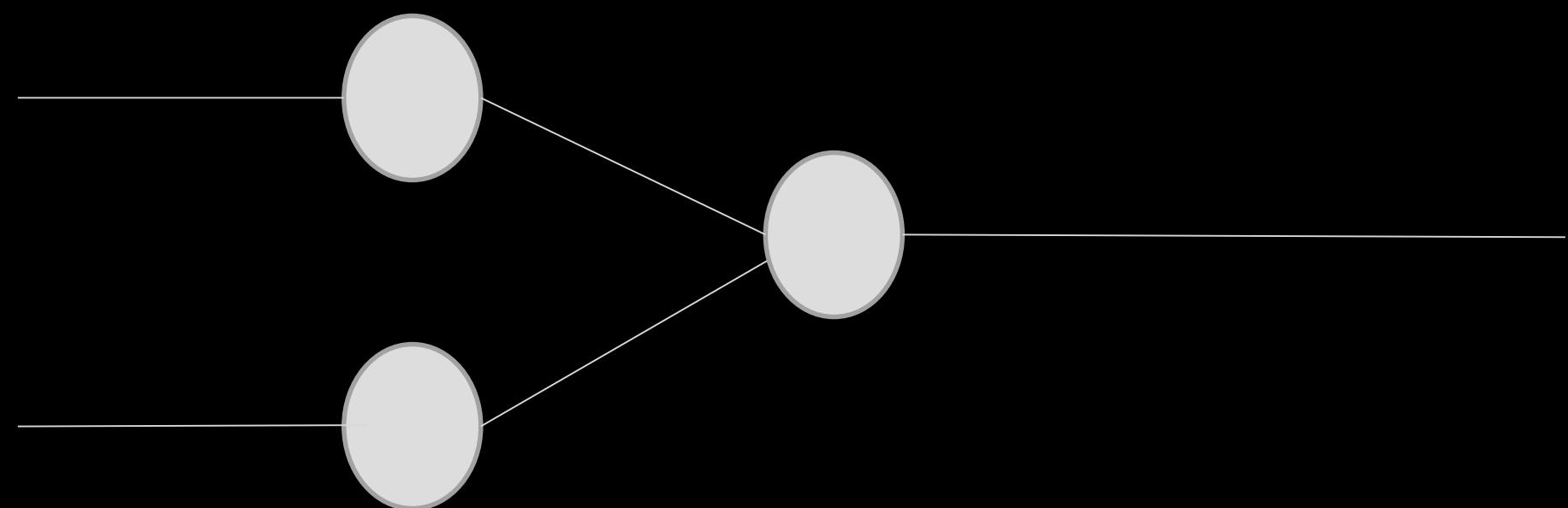
**AND**

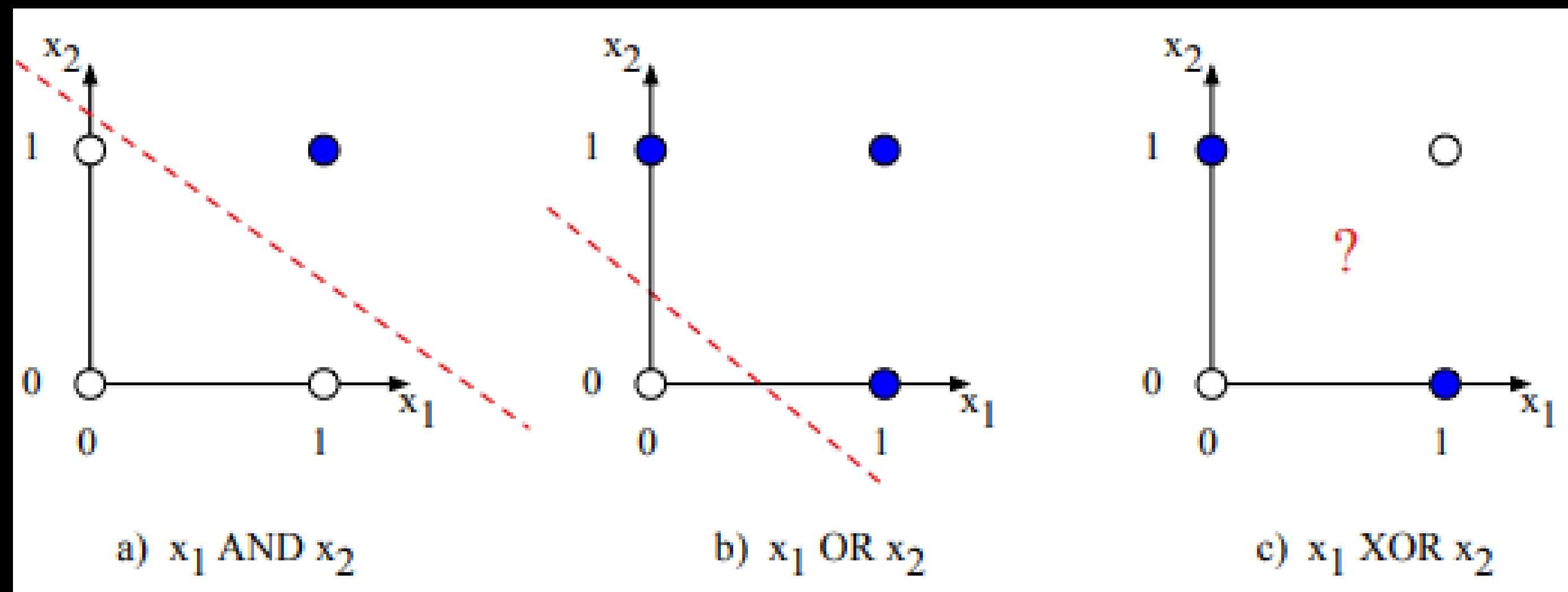
x1	x2	y
0	0	0
0	1	0
1	0	0
1	1	1



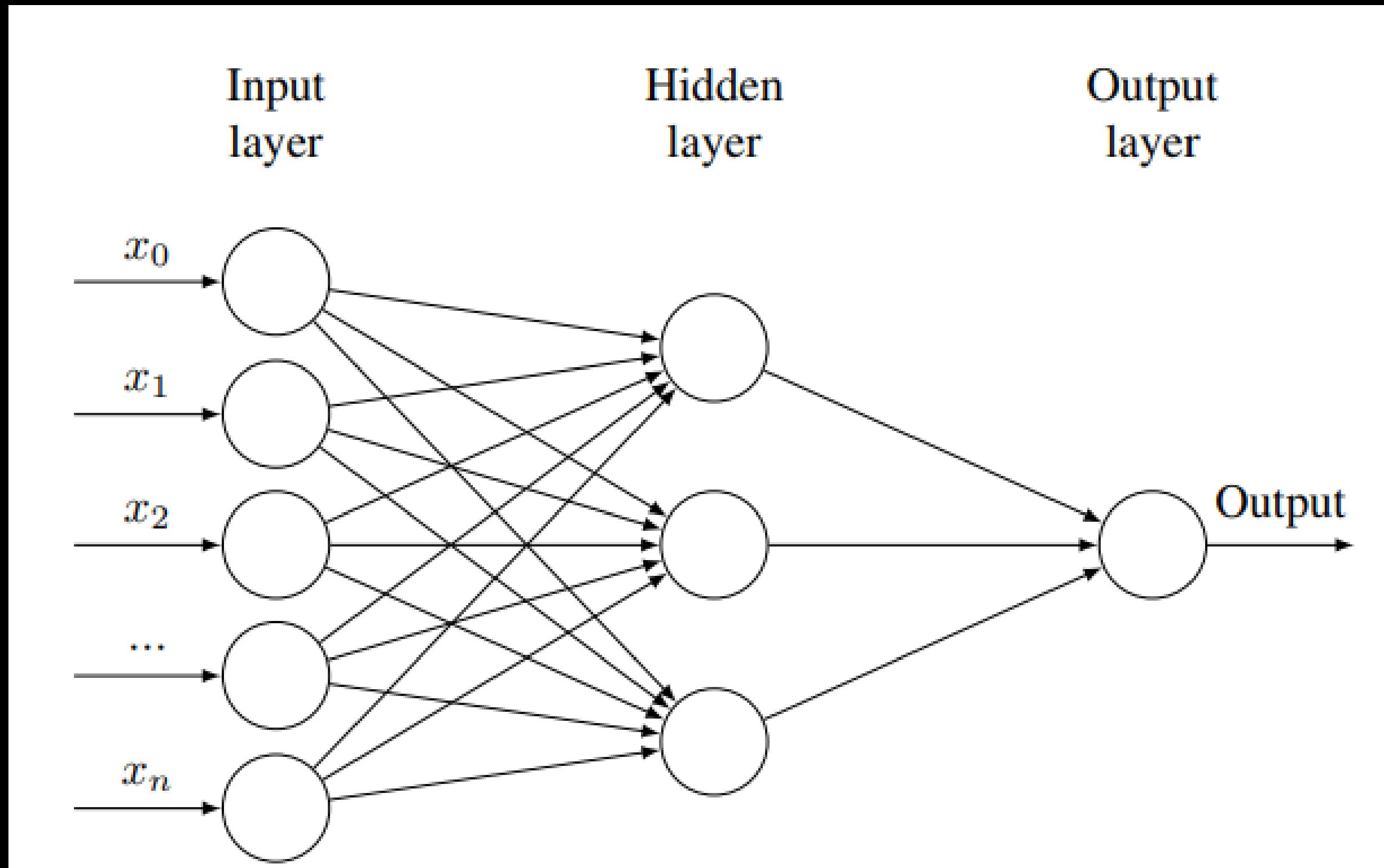
**XOR**

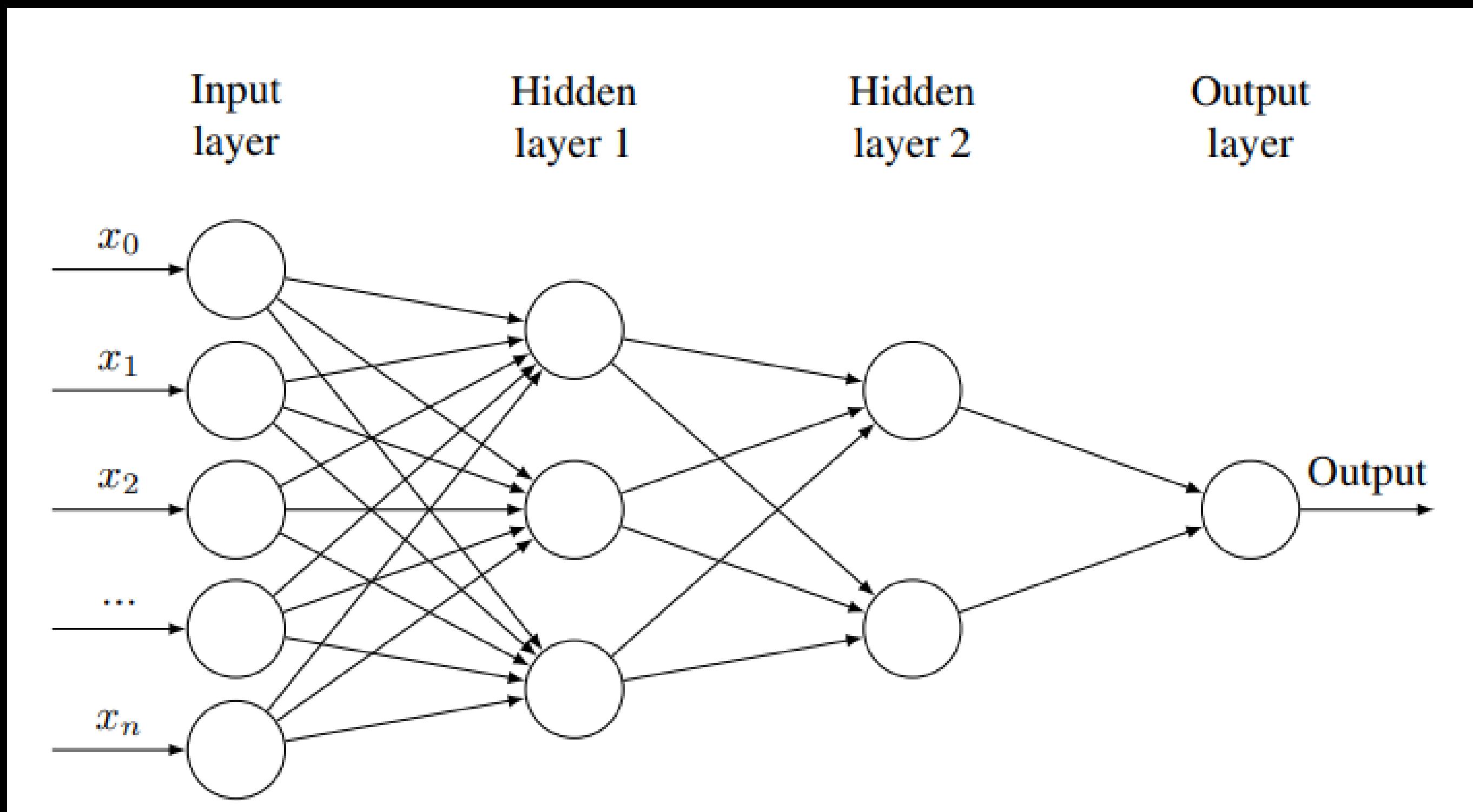
x1	x2	y
0	0	0
0	1	1
1	0	1
1	1	0



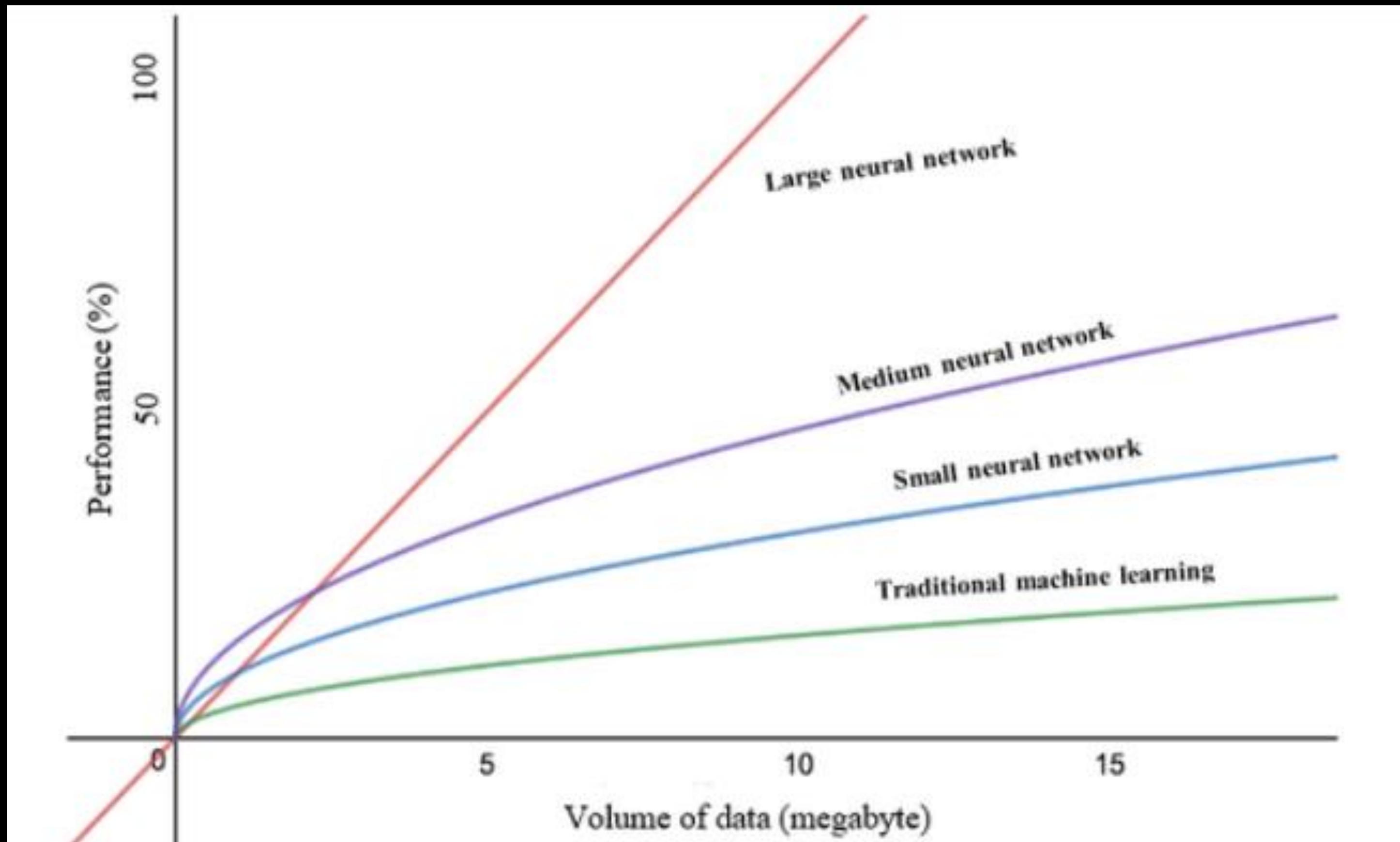


# Multi Layer Neural Network





# Backpropagation



Q.43

Consider a Multi-Layer Perceptron (MLP) model with one hidden layer and one output layer. The hidden layer has 10 neurons, and the output layer has 3 neurons. The input to the MLP is a 5-dimensional vector. Each neuron is connected to every neuron in the previous layer, and a bias term is included for each neuron. The activation function used is the sigmoid function. The total number of trainable parameters in this MLP model is \_\_\_\_\_.

Consider the two class classification task that consists of the following points:

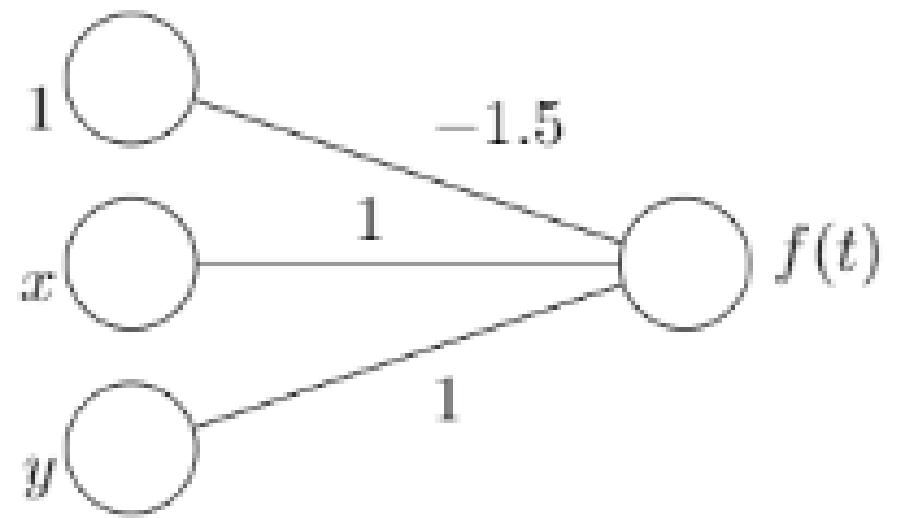
$$\text{Class } C_1 : [1 \ 1.5] \ [1 \ -1.5]$$

$$\text{Class } C_2 : [-2 \ 2.5] \ [-2 \ -2.5]$$

The decision boundary between the two classes using single perceptron is given by:

- A.  $x_1 + x_2 + 1.5 = 0$
- B.  $x_1 + x_2 - 1.5 = 0$
- C.  $x_1 + 1.5 = 0$
- D.  $x_1 - 1.5 = 0$

Consider a single perception with weights as given in the following figure:



and  $f(t)$  is defined as

$$f(t) \begin{cases} 1, & t > 0 \\ 0, & t \leq 0 \end{cases}$$

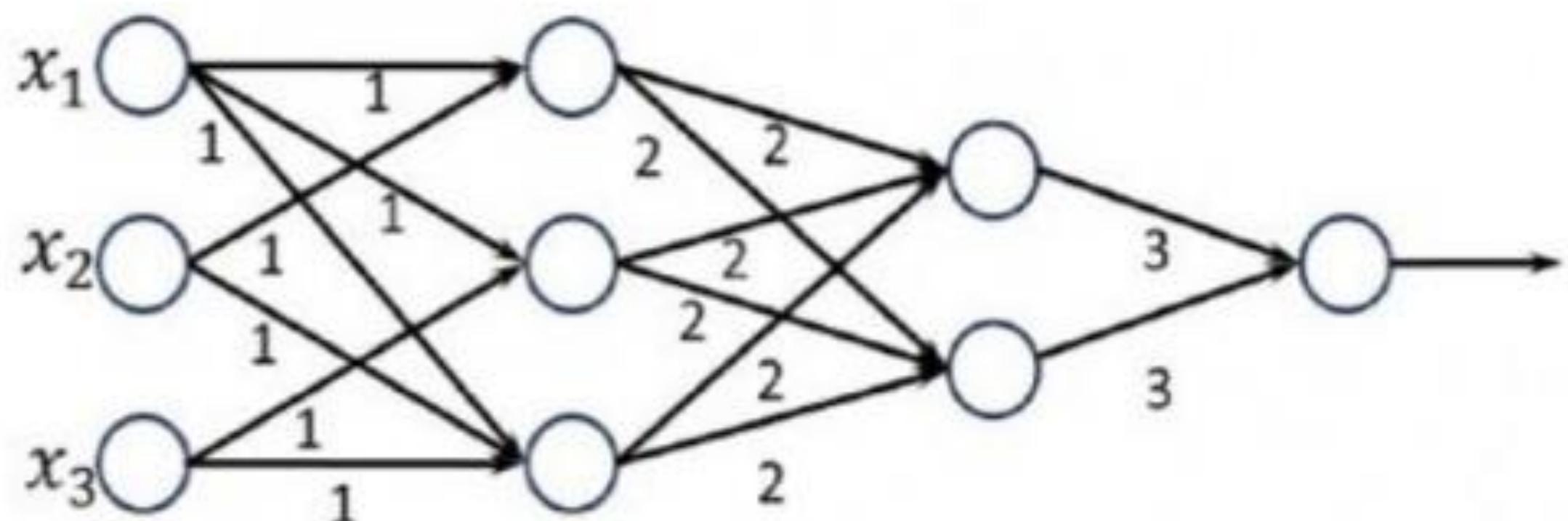
The above perception can solve

Answer: (B)

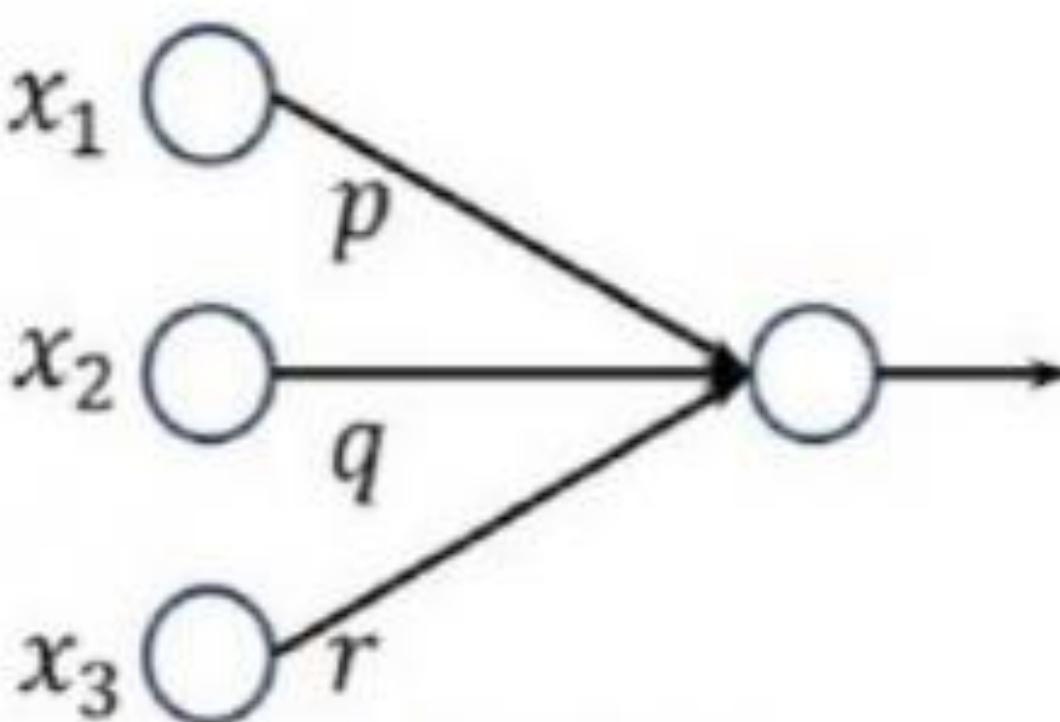
- A. OR problem
- B. AND problem
- C. XOR problem
- D. All of the above

Consider the two neural networks (NNs) shown in Figures 1 and 2, with *ReLU* activation ( $\text{ReLU}(z) = \max\{0, z\}, \forall z \in \mathbb{R}$ ). The connections and their corresponding weights are shown in the Figures. The biases at every neuron are set to 0.

For what values of  $p, q, r$  in Figure 2 are the two NNs equivalent, when  $x_1, x_2, x_3$  are positive?



**Figure 1**



**Figure 2**

Note:  $\mathbb{R}$  denotes the set of real numbers.

- A.  $p = 36, q = 24, r = 24$
- B.  $p = 24, q = 24, r = 36$
- C.  $p = 18, q = 36, r = 24$
- D.  $p = 36, q = 36, r = 36$

**Given a feedforward neural network with a large number of hidden layers, which of the following issues is it most likely to face?**

- A. Overfitting and vanishing gradient
- B. Underfitting and exploding gradient
- C. Bias-variance trade-off and linear separability
- D. None of the above

**Answer:** A

Thank you