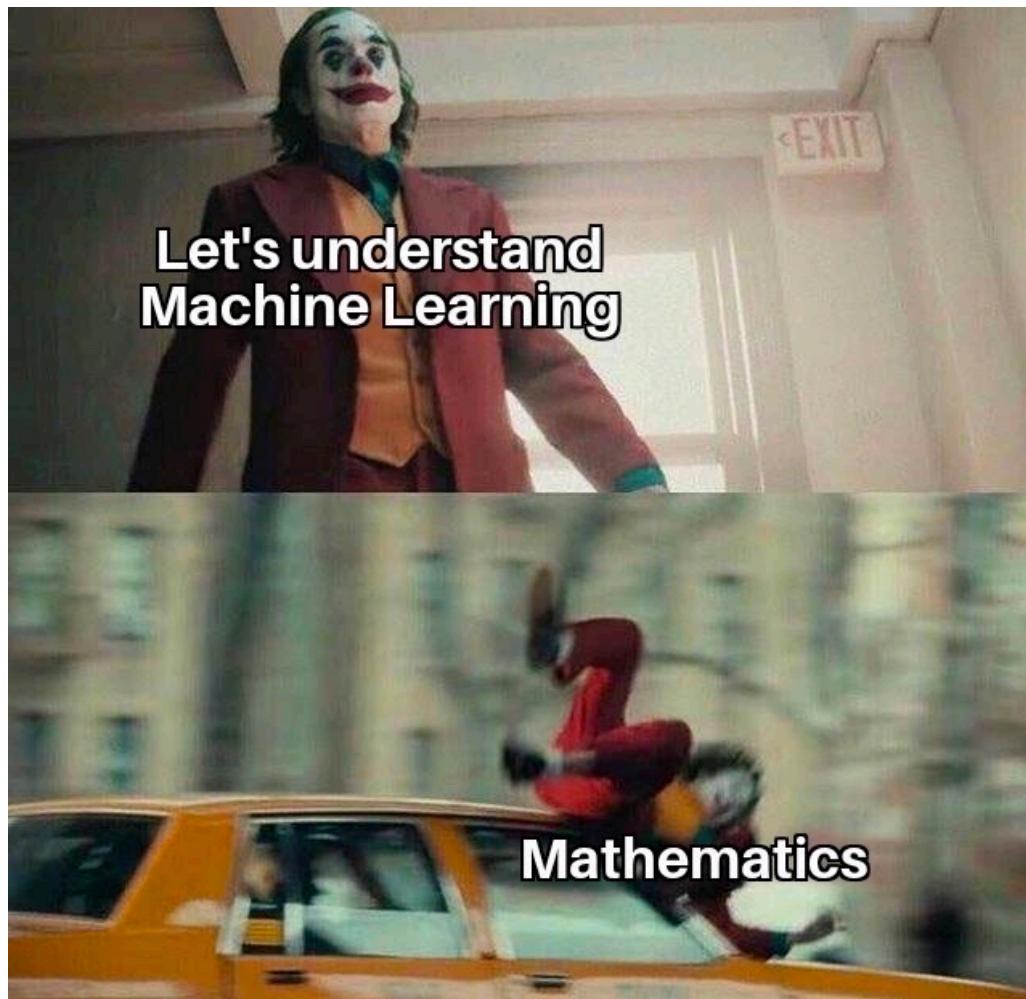


February 21, 2023

DSML : Math for ML.

## Linear Algebra: Halfspaces and Distances.



## Recap:

(a) Vectors:  $\bar{x}, \bar{y} \in \mathbb{R}^d$

$$\bar{x}^T = [x_1 \ x_2 \ \dots \ x_d]$$

$$\bar{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

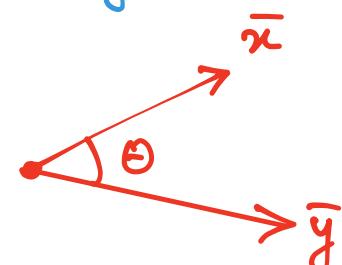
$$\bar{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_d \end{bmatrix}$$

(b) Norm:  $\bar{x} \in \mathbb{R}^d$   
(length, magnitude)

$$\|\bar{x}\| = \sqrt{x_1^2 + x_2^2 + \dots + x_d^2}$$

(c) Dot Product:  $\bar{x} \cdot \bar{y} = \bar{x}^T \bar{y}$

$$= \sum_{i=1}^d x_i \cdot y_i$$



(d) Angle between two vectors:

$$\cos \theta = \frac{\bar{x}^T \bar{y}}{\|\bar{x}\| \cdot \|\bar{y}\|}$$

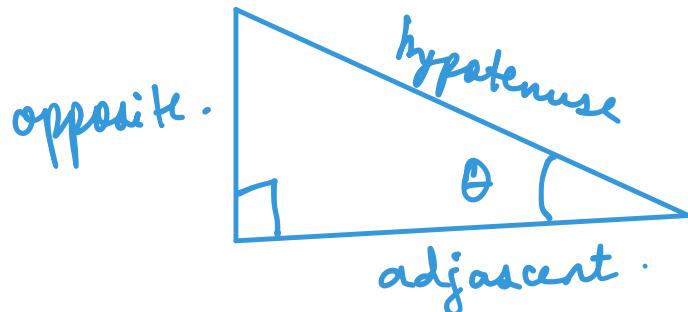
## Unit Vectors:

"Unit vectors represent the direction in which a vector is pointing towards"

\* Vector with magnitude 1.  $\rightarrow$  Unit vector.

$$\|\bar{w}\| = \sqrt{w_1^2 + w_2^2}$$
$$\frac{\bar{w}}{\|\bar{w}\|} = \hat{w} = \begin{bmatrix} \frac{w_1}{\sqrt{w_1^2 + w_2^2}} \\ \frac{w_2}{\sqrt{w_1^2 + w_2^2}} \end{bmatrix} = \sqrt{\frac{w_1^2}{w_1^2 + w_2^2} + \frac{w_2^2}{w_1^2 + w_2^2}} = \sqrt{\frac{w_1^2 + w_2^2}{w_1^2 + w_2^2}} = \sqrt{1} = 1$$
$$\|\hat{w}\| =$$

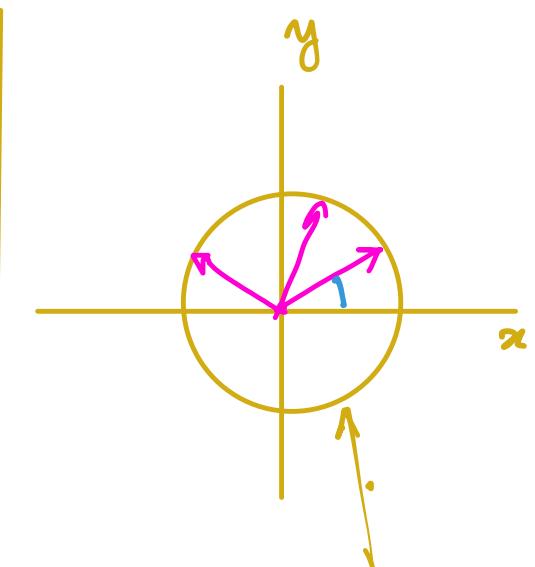
## Basic Trigonometry



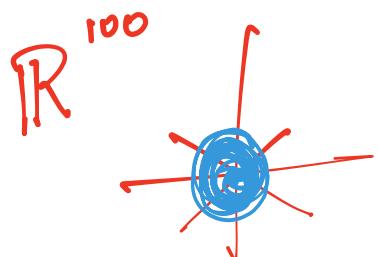
$$\sin(\theta) = \frac{\text{opposite}}{\text{hypotenuse}}$$

$$\cos(\theta) = \frac{\text{adjacent}}{\text{hypotenuse}}$$

$$\tan(\theta) = \frac{\text{opposite}}{\text{adjacent}}$$

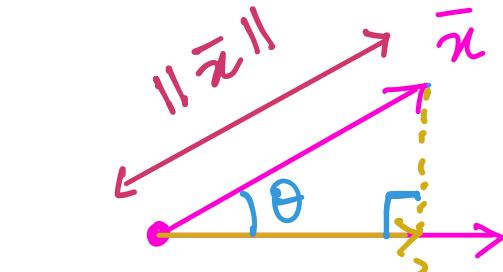


circle  
with  
radius  
1.



## Projection :

$$\bar{x}^T \bar{y} = \bar{x} \cdot \bar{y}$$



$\bar{P}$  → projection of  $\bar{x}$   
on  $\bar{y}$ .  
(shadow)

$$\cos(\theta) = \frac{\|\bar{P}\|}{\|\bar{x}\|}$$

$$\frac{\bar{x}^T \bar{y}}{\|\bar{x}\| \cdot \|\bar{y}\|} = \frac{\|\bar{P}\|}{\|\bar{x}\|}$$

$$\|\bar{P}\| = \|\bar{x}\| \cdot \frac{\bar{x}^T \bar{y}}{\|\bar{x}\| \cdot \|\bar{y}\|}$$

$$\|\bar{P}\| = \frac{\bar{x}^T \bar{y}}{\|\bar{y}\|} = \frac{\bar{x} \cdot \bar{y}}{\|\bar{y}\|}$$

$$= \bar{x}^T \left( \frac{\bar{y}}{\|\bar{y}\|} \right) = \bar{x}^T \hat{\bar{y}}$$

not useful large dimensional

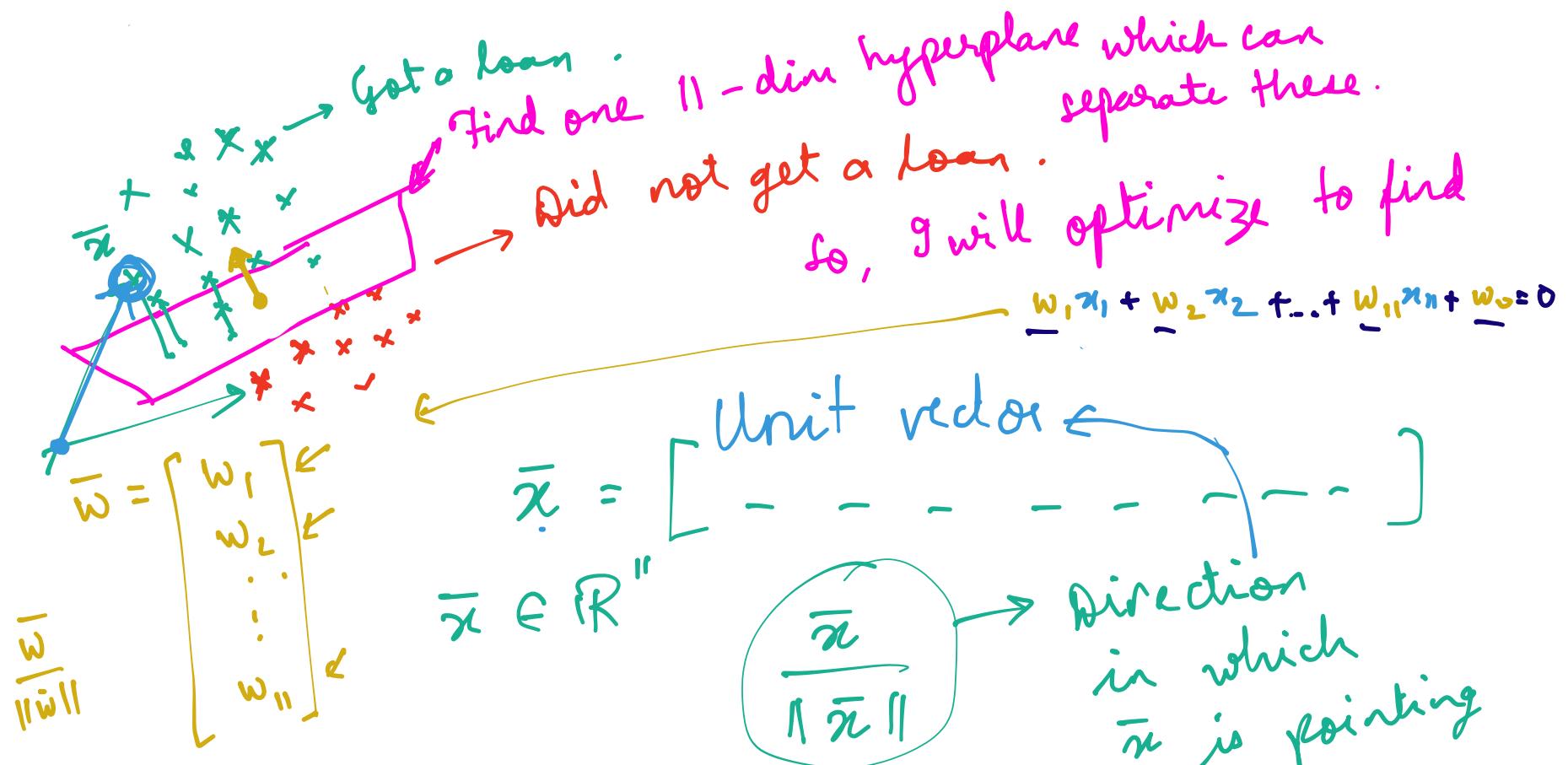
Classification problems.

$n$

	Loan_ID	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History	Property_Area
0	LP001002	Male	No	0	Graduate	No	5849	0.0	NaN	360.0	1.0	Urban
1	LP001003	Male	Yes	1	Graduate	No	4583	1508.0	128.0	360.0	1.0	Rural
2	LP001005	Male	Yes	0	Graduate	Yes	3000	0.0	66.0	360.0	1.0	Urban
3	LP001006	Male	Yes	0	Not Graduate	No	2583	2358.0	120.0	360.0	1.0	Urban
4	LP001008	Male	No	0	Graduate	No	6000	0.0	141.0	360.0	1.0	Urban

$$X = \{\bar{x}_i\}_{i=1}^n$$

$\bar{x}_i \rightarrow$  feature vectors.



## Shifting 2-D line.

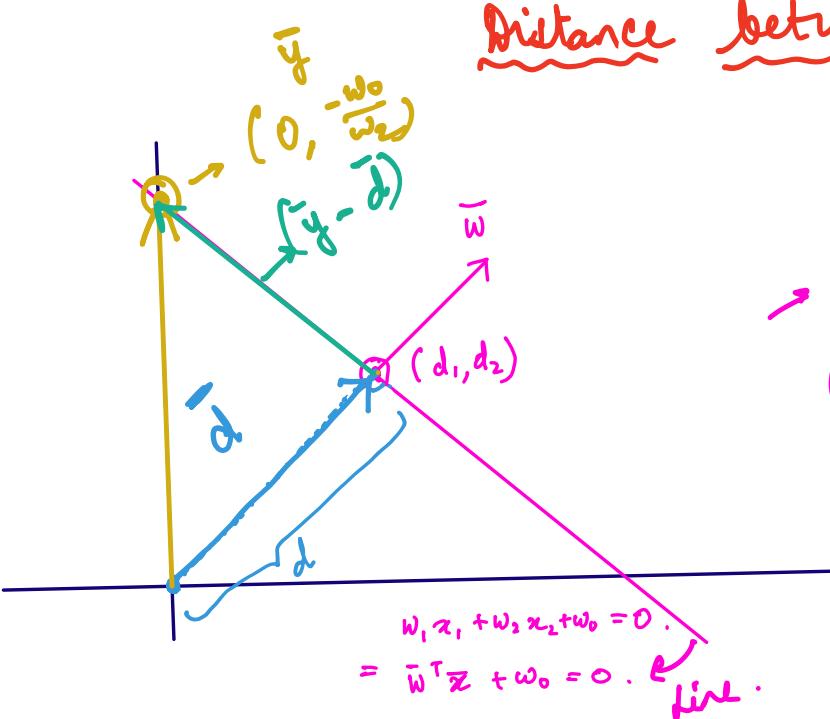
$$l : w_1 x_1 + w_2 x_2 + w_0 = 0$$

Q1] How to shift this line by 'a' units to the right?

$$w_1(x_1 - a) + w_2 x_2 + w_0 = 0$$

Q2] How to shift this line by 'b' units to the top?

$$w_1 x_1 + w_2(x_2 - b) + w_0 = 0$$



Distance between origin and a line

$$\textcircled{1} \quad \|\bar{d}\| = \sqrt{d_1^2 + d_2^2}$$

$$\textcircled{2} \quad w_1 d_1 + w_2 d_2 + w_0 = 0.$$

$$\textcircled{3} \quad \hat{w} = \frac{\bar{w}}{\|\bar{w}\|}, \quad \bar{d} = \begin{bmatrix} d_1 \\ d_2 \end{bmatrix} = \begin{bmatrix} k w_1 \\ k w_2 \end{bmatrix}$$

Substitute \textcircled{3} in \textcircled{2}  $\rightarrow$

$$d_1 = \frac{k \cdot w_1}{\|\bar{w}\|}, \quad d_2 = \frac{k \cdot w_2}{\|\bar{w}\|}$$

\textcircled{4} To get  $(d_1, d_2)$ ,  
substitute  $k$  in \textcircled{3}

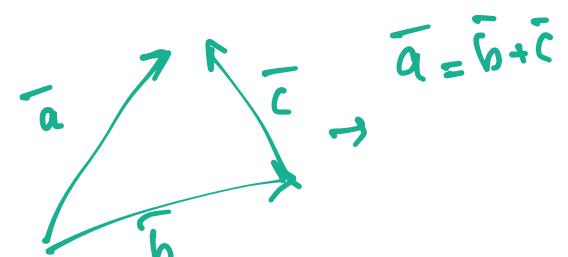
$$\bar{d} = \begin{bmatrix} -w_0 w_1 \\ \frac{(w_1^2 + w_2^2)}{\|\bar{w}\|^2} \end{bmatrix}$$

$$\begin{bmatrix} -w_0 w_2 \\ \frac{(w_1^2 + w_2^2)}{\|\bar{w}\|^2} \end{bmatrix}$$

$$k \frac{w_1^2}{\|\bar{w}\|} + k \frac{w_2^2}{\|\bar{w}\|} + w_0 = 0.$$

$$k \left( \frac{w_1^2 + w_2^2}{\|\bar{w}\|} \right) = -w_0$$

$$k = \frac{-w_0 \times \|\bar{w}\|}{(w_1^2 + w_2^2)}$$



$$\bar{y} = \left( 0, -\frac{\omega_0}{\omega_2} \right)$$

$$\omega_1 x_1 + \omega_2 x_2 + \omega_0 = 0.$$

$$\bar{d} = \begin{bmatrix} -\frac{\omega_0 \omega_1}{(\omega_1^2 + \omega_2^2)} \\ -\frac{\omega_0 \omega_2}{(\omega_1^2 + \omega_2^2)} \end{bmatrix}$$

$$x_2 = -\frac{\omega_1}{\omega_2} x_1 - \underbrace{\frac{\omega_0}{\omega_2}}_{-}$$

$$d^T (\bar{y} - \bar{d}) = \frac{\omega_0 \omega_1}{(\omega_1^2 + \omega_2^2)} \times \frac{-\omega_0 \omega_1}{\omega_1^2 + \omega_2^2}$$

$$+ \left( -\frac{\omega_0}{\omega_2} + \frac{\omega_0 \omega_2}{(\omega_1^2 + \omega_2^2)} \right) \times \frac{-\omega_0 \omega_2}{\omega_1^2 + \omega_2^2}.$$

$$\left( \frac{-\omega_0 (\omega_1^2 + \omega_2^2) + \omega_0 \omega_2^2}{\omega_2 (\omega_1^2 + \omega_2^2)} \right) - \frac{\omega_0 \omega_2}{\omega_1^2 + \omega_2^2}$$

$$d^T (\bar{y} - \bar{d}) = \frac{-w_0^2 w_1^2}{(w_1^2 + w_2^2)^2} + \frac{w_0^2}{(w_1^2 + w_2^2)^2} \underbrace{(w_1^2 + w_2^2) - w_0^2 w_2^2}_{(w_1^2 + w_2^2)^2}$$

$$= \frac{-w_0^2 w_1^2 + w_0^2 w_1^2 + w_0^2 w_2^2 - w_0^2 w_2^2}{(w_1^2 + w_2^2)^2}$$

$$= \frac{0}{(w_1^2 + w_2^2)^2} = \frac{0}{\cdot}$$

Distance between a line and the origin.

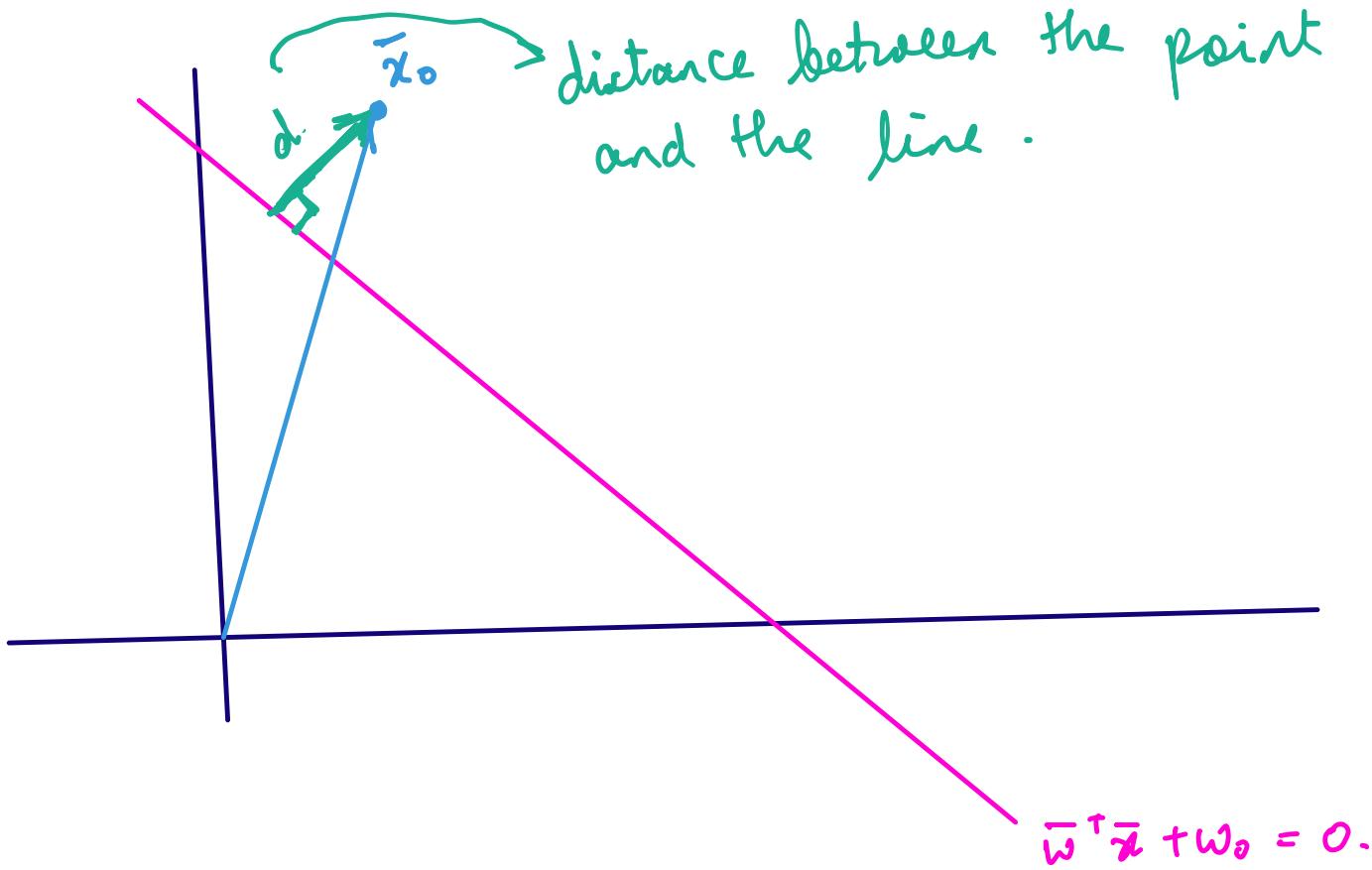
$$\bar{d} = \begin{bmatrix} -w_0 w_1 \\ \frac{(w_1^2 + w_2^2)}{w_0} \\ -w_0 w_2 \\ \frac{(w_1^2 + w_2^2)}{w_0} \end{bmatrix}$$

Hence, the distance between the origin and a line  $\bar{w}^T \bar{x} + w_0 = 0$  is given by

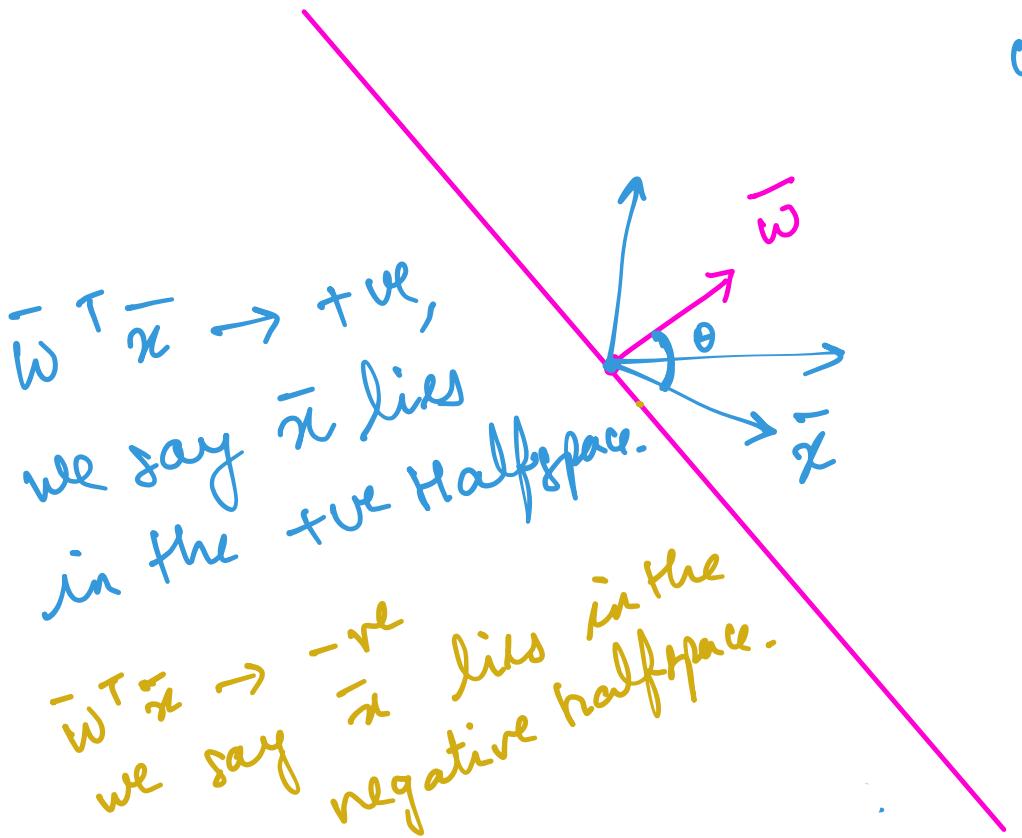
$$d = \frac{w_0}{\|\bar{w}\|}$$

$$\begin{aligned} \|\bar{d}\| &= \sqrt{\frac{w_0^2 w_1^2}{(w_1^2 + w_2^2)^2} + \frac{w_0^2 w_2^2}{(w_1^2 + w_2^2)^2}} \\ &= \sqrt{\frac{w_0^2 (w_1^2 + w_2^2)}{(w_1^2 + w_2^2)^2}} = \frac{w_0}{\sqrt{w_1^2 + w_2^2}} = \frac{w_0}{\|\bar{w}\|} \end{aligned}$$

Distance between a point and a line.



Final answer:  $d = \frac{\bar{w}^T \bar{x}_0 + w_0}{\|\bar{w}\|}$



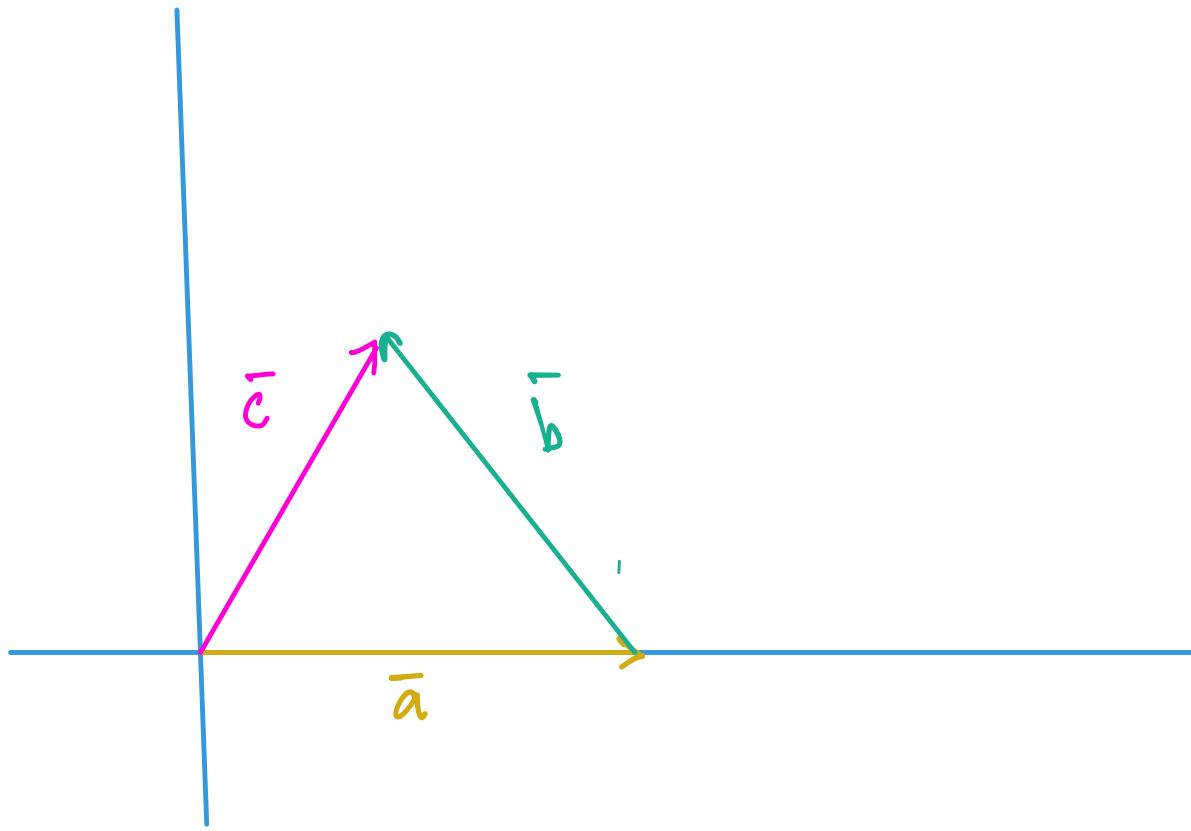
sign  $(\bar{w}^T \bar{x})$

$\rightarrow +ve$ when $0^\circ \leq \theta \leq 90^\circ;$ $270^\circ \leq \theta \leq 360^\circ$	$\rightarrow -ve$ when $90^\circ \leq \theta \leq 270^\circ$
--	---

$$\cos \theta = \frac{\bar{w}^T \bar{x}}{\|\bar{w}\| \cdot \|\bar{x}\|}$$

$\cos \theta \rightarrow +ve$  when  
 $0^\circ \leq \theta \leq 90^\circ,$   
 $270^\circ \leq \theta \leq 360^\circ$   
  
 $\rightarrow -ve$  when  
 $90^\circ \leq \theta \leq 270^\circ$

$$\bar{w}^T \bar{x} + w_0$$



$$\bar{c} = \bar{a} + \bar{b}$$

$$\bar{b} = \bar{c} - \bar{a}.$$

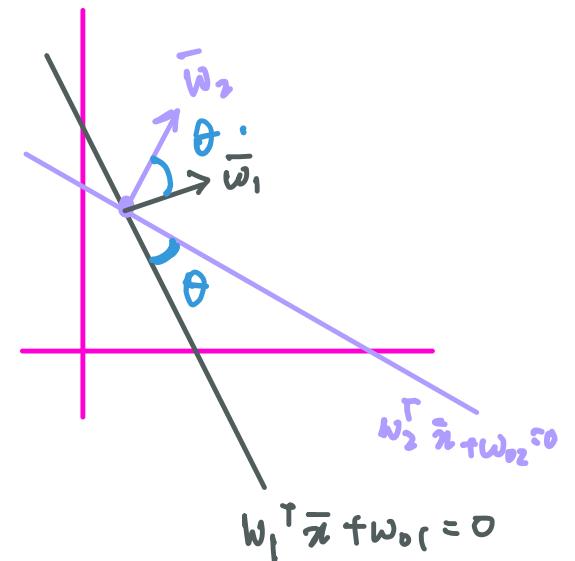
$$w_1 x_1 + w_2 x_2 + \dots + w_d x_d + w_0 = 0.$$

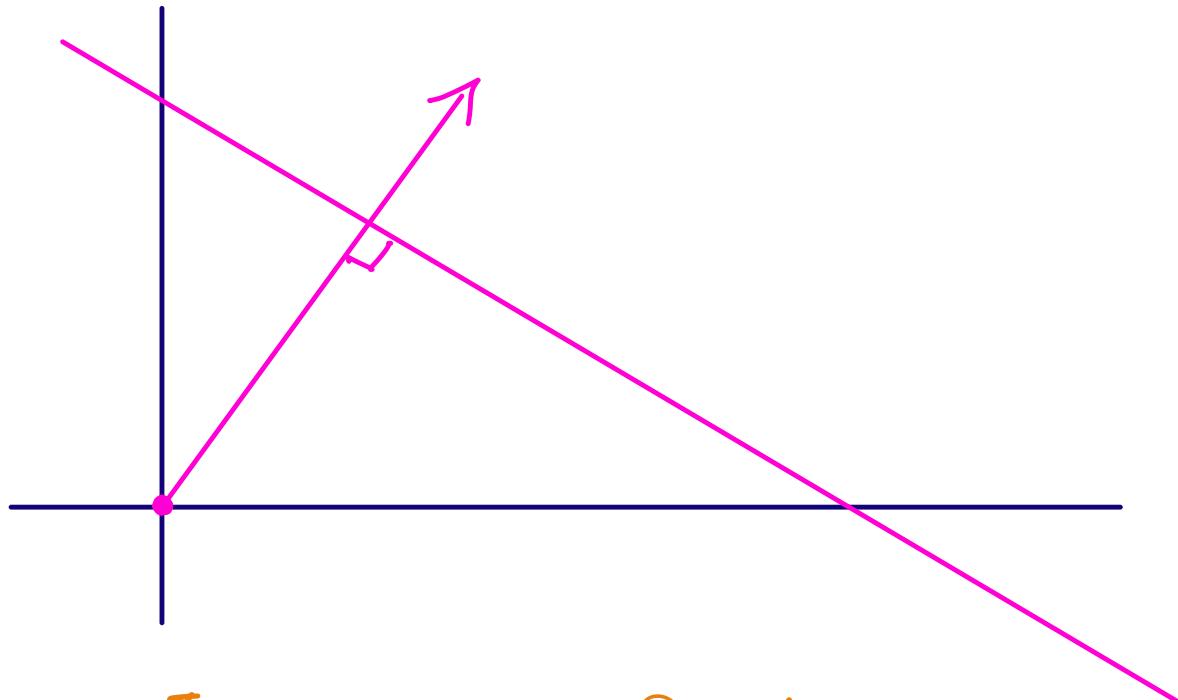
$$\Rightarrow \bar{w}^T \bar{x} + w_0 = 0$$

$\nearrow D_1 : w_1^T \bar{x} + w_{01} = 0$

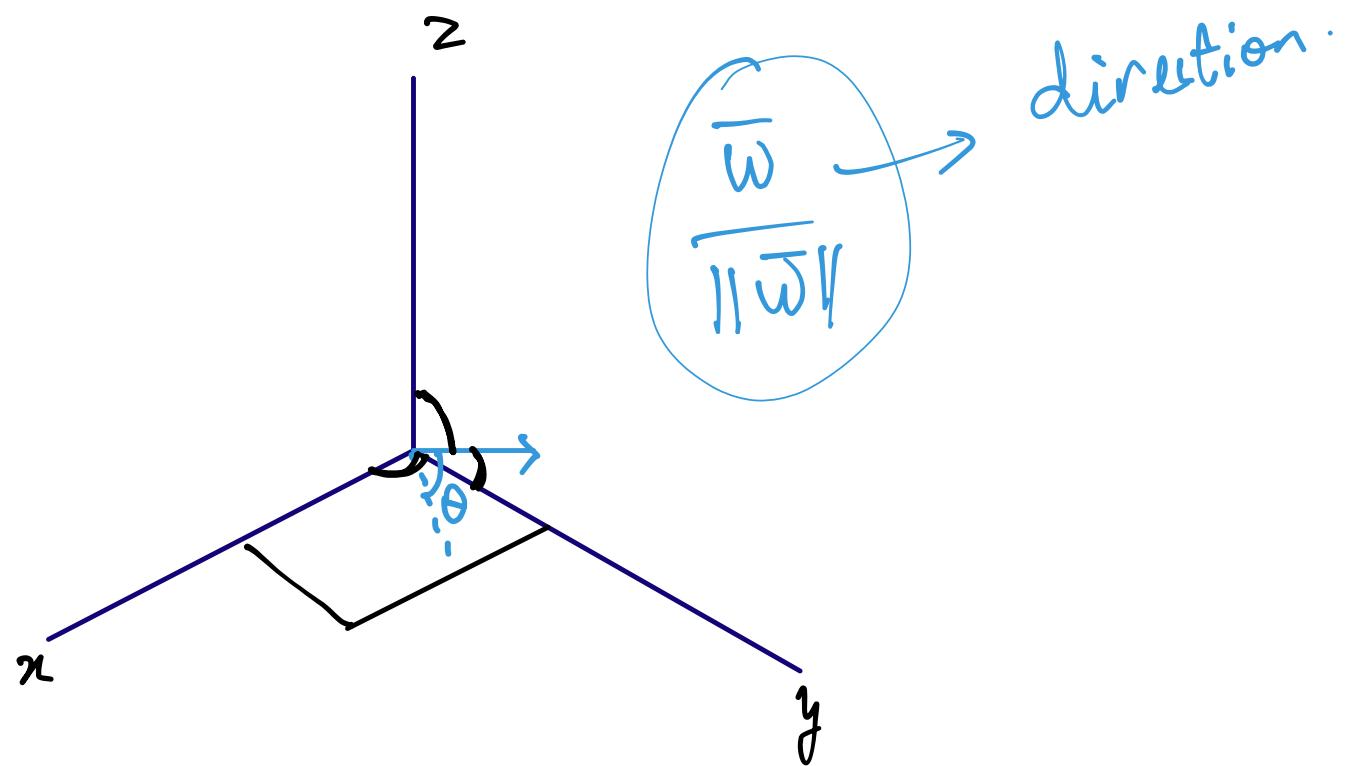
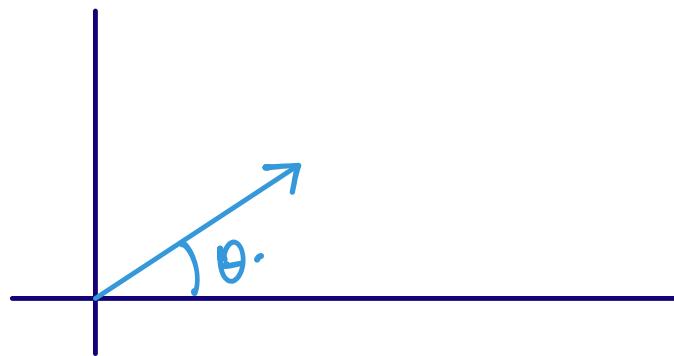
$\nearrow D_2 : w_2^T \bar{x} + w_{02} = 0.$

$$\cos(\theta) = \frac{\bar{w}_1^T \bar{w}_2}{\|\bar{w}_1\| \cdot \|\bar{w}_2\|}$$





If  $\bar{w}^\top \bar{x} + w_0 = 0$  is a  
D-dim hyperplane, then,  $\bar{w}$  is  $\perp^r$   
to all points lying on the d-Dim  
hyperplane.



direction.

$$\bar{x}^T \underbrace{w}_{\substack{\text{vector} \\ \text{matrix}}} \bar{y}$$

← another vector.

$$(\bar{x}^T w \bar{y})^T = \bar{y}^T w^T \bar{x}$$

$$\bar{x}^T \bar{y}$$

✓

### Problem Description:

Given two arrays, representing two continuous features from a certain dataset, your task is to calculate the **Spearman rank correlation coefficient** between these two arrays.

### Input Format:

The input contains two arrays, passed to the function `calculateSRCC()`

### Output Format:

Return the calculated value, rounded off up to 3 decimal places

### Input Constraints:

`1 <= len(arr) <= 10000`, for both the input arrays  
`-2 <= arr[i] <= 2`

### Sample Input:

```
-0.5722404243297754 1.8421092051183359 0.06532025415596275 -0.021812955390535325 0.2386780529341842
-0.04628825915644745 0.2073075701391909 1.0467470550072373 -0.31174784507404363 0.13821676157233972
```

### Sample Output:

0.6

### Output Explanation:

Using the Spearman correlation coefficient formula on both the arrays, we can easily calculate the value of the coefficient.

### Code Constraints:

You aren't allowed to use `scipy`, `stats`, or any other libraries for inbuild correlation coefficient calculation

**Note:** Feel free to use `NumPy` to compute the helper functions.

Handwritten notes:

- Three arrays are shown:
  - [1, 11, 7, -3]
  - [0, 1, 2, 3]
  - [-3, 0, 2, 1]
- A handwritten formula below the arrays is:  $\text{rank} = \text{arr.argsort().argsort()}$ .

Sun Pharmaceutical Industries claims that a person's IQ improves after they use the Donepezil drug.

To test this claim a trial was conducted considering **20** patients. An IQ test was conducted for these patients before giving the drug and an IQ test was conducted for the same set of patients after the drug the recorded results are shown below.

```
IQ_before=[101,124,89,57,135,98,69,105,114,106,97,121,93,116,102,71,88,108,144,99]
```

```
IQ_after=[113,127,89,70,127,104,69,127,115,99,104,120,95,129,106,71,94,112,154,96]
```

Perform an appropriate test to test the claim at **90%** confidence.

1 Spearman rank Correlation ▾

2 Percentage ▾

3 Customer tip ▾

Q 1	Grocery delivery ▾	
Q 2	Drill sergeant ▾	
Q 3	Correlation ▾	
Q 4	Suitable for T-test ▾	
Q 5	Train & Test data ▾	
Q 6	IQ test ▾	
Q 7	Cloud computing company ▾	
Q 8	Dating website ▾	
Q 9	Coin Flip ▾	

Q 10 Prakruthi cafe ▾

Needs to  
be fixed.

A drill sergeant claims that his battalion is able to complete the toughest obstacle course in **5 minutes or less.**

**35 soldiers** in the battalion are picked at random and asked to undergo the course, which they complete with an **average time of 4.78 minutes** and a **standard deviation of 1.8 minutes**.

To test the drill sergeant's claim, an alpha of **0.05** was selected, and the critical z-score was **-1.64**. The calculated z-score for the sample was **-0.12**. What conclusion can be drawn from this?

should be -0.72.

$$H_0 : \mu = 5$$

$$H_a : \mu \leq 5.$$

$$n = 35.$$

$$m_{obs} = 4.78$$

$$\sigma = 1.8.$$

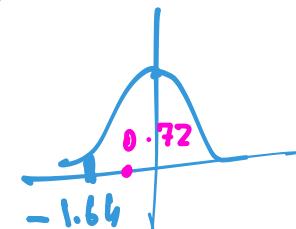
$$-1.64.$$

less than

$$-0.72.$$

$$z = \frac{4.78 - 5}{1.8 / \sqrt{35}}.$$

=



X correct (marked) answer : Test conclusively shows alt. hypothesis is true. ↑ wrong ans.

A cloud computing company offers storage services and owns data centers in four geographical regions to reduce the latency of its services.

The company is carrying out a study to compare the load on each data center, measured by the number of million 'fetch' requests per day.

- $H_0$ : The load is the same among all the regions.
- $H_1$ : The load in at least one of the regions is different from the others. ↗

It was found that the population variances across the four regions were approximately the same, and the data in each group followed an approximately Normal distribution.

The mean square between the groups (i.e, **MSB**) is **180.23**, and the mean square within the groups (i.e, **MSW**) is **8.159**, Calculate the **F-test statistic** while the critical F-score was calculated as **3.49** at the alpha=0.05 significance level. What can you conclude about the test?

$$F_{\text{stat}} = \frac{\text{MSB}'}{\text{MSW}} = \frac{180.23}{8.159} = \underline{22.03} \gg \underline{3.49}$$

\* Mistake,  
correct answer : Reject  $H_0$ .

#### **Problem Description:**

A Chinese restaurant collects data on their customer's tips given to waiters. The dataset includes information about the bill, the tip paid, and customer details such as the sex of the customer, whether the customer smokes, and the day and time that the customer visited the restaurant.

We wish to study the relationship between the tips paid and customer details. Specifically, the restaurant claims that the **female customers tip more than the male customers**. Test the claim with the help of the given dataset at a **90%** confidence interval.

**Note:** Calculate the **p-value** with the help of obtained z-score and return the status as '**Reject**' or '**Fail to reject**'.

#### **Input Description:**

A data frame.

\* No issues here .

#### **Output Description:**

Status ('Reject' or 'Fail to reject')

#### **Sample Input:**

	total_bill	tip	sex	smoker	day	time
0	16.99	1.01	Female	No	Sun	Dinner
1	10.34	1.66	Male	No	Sun	Dinner
2	21.01	3.50	Male	No	Sun	Dinner
3	23.68	3.31	Male	No	Sun	Dinner
4	24.59	3.61	Female	No	Sun	Dinner

#### **Sample Output:**

Reject

#### **Sample Explanation:**

From the given data frame, first, calculate the z-score using two sample z-test.

Based on the z-score calculate p\_value and return the status as 'Reject' or 'Fail to reject'.