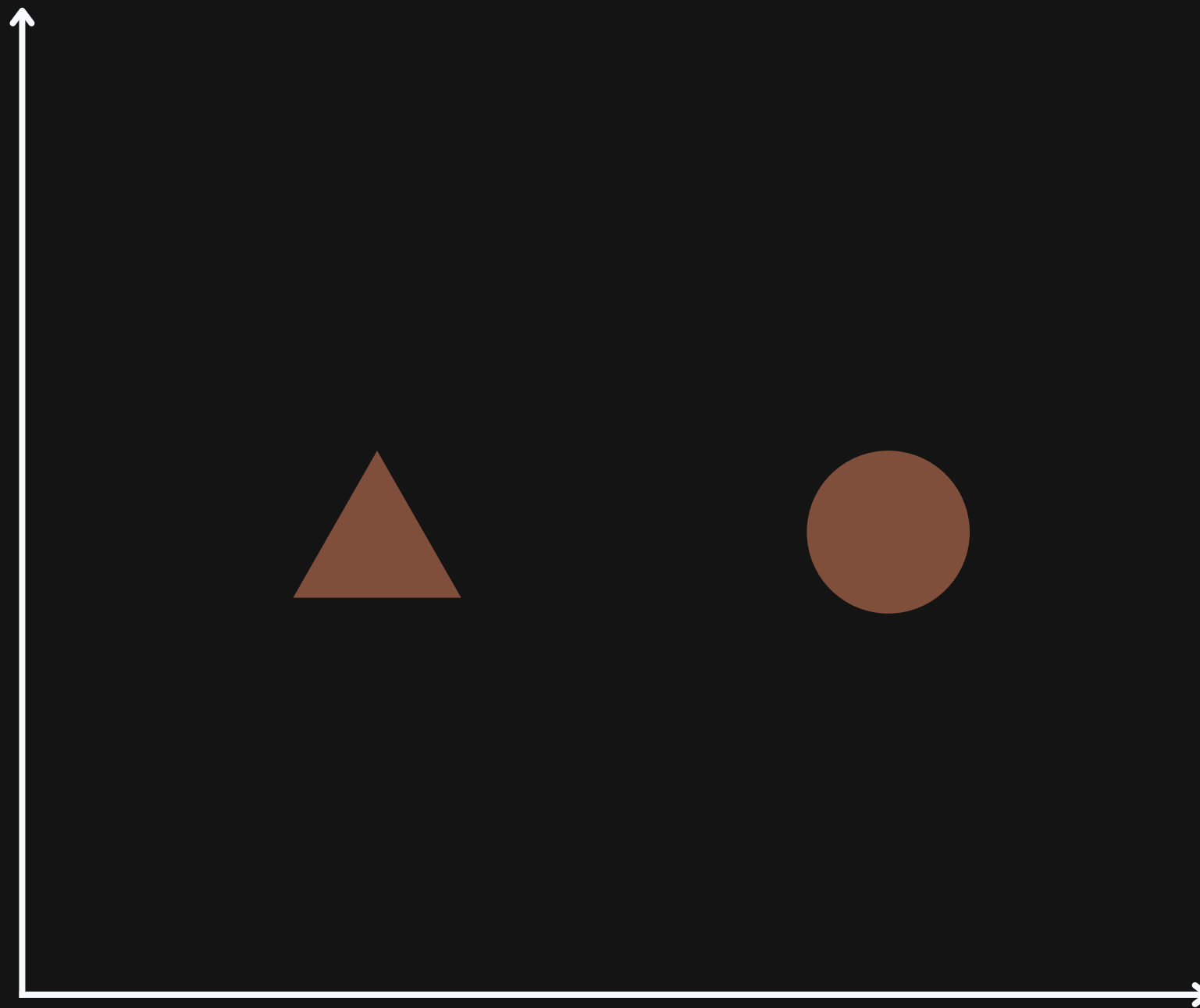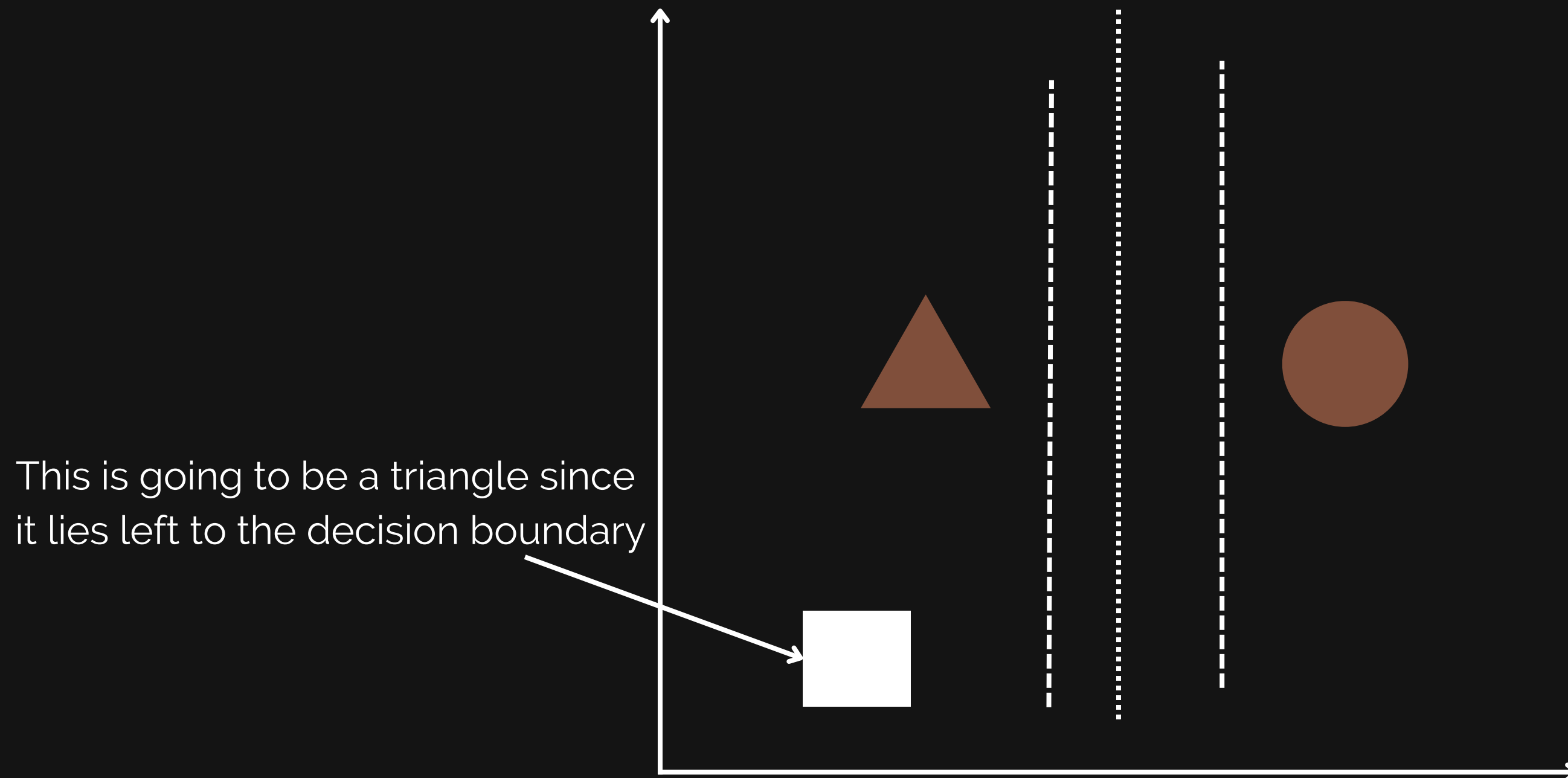# KNN Algo

AI Club
IIT MADRAS
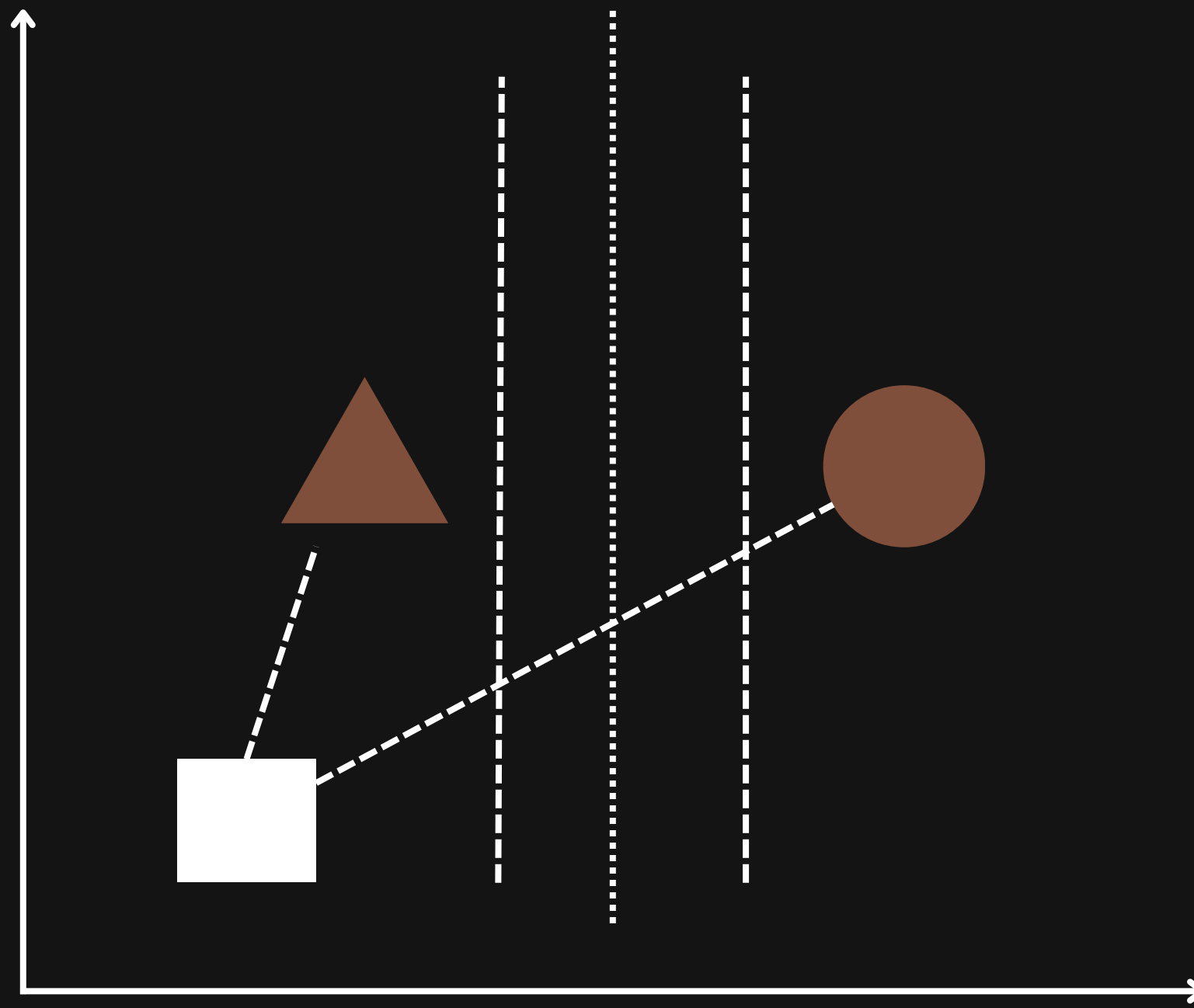
Let us begin with a simplest case of a classification problem.Let us try to classify the instances into  two classes.

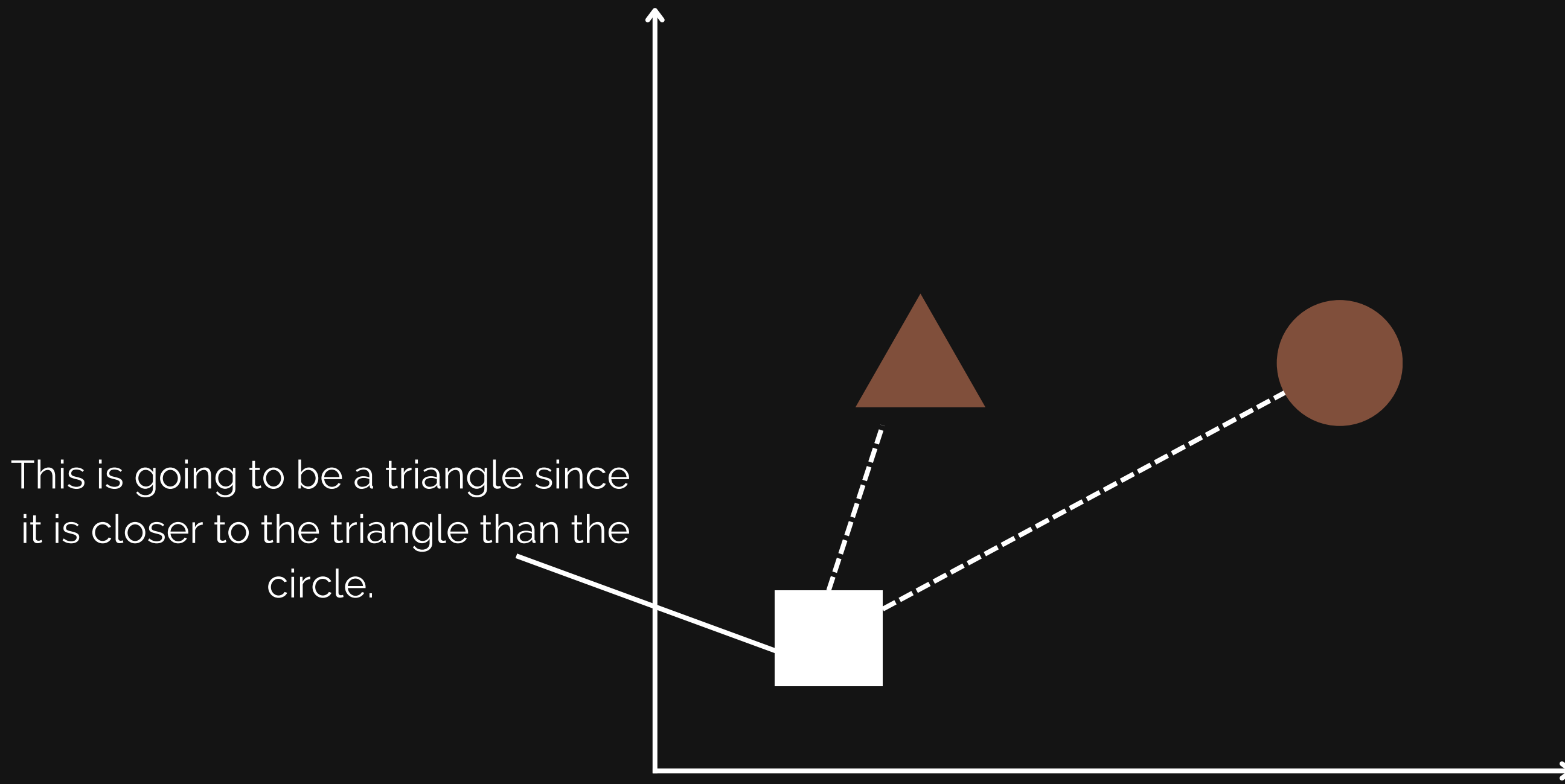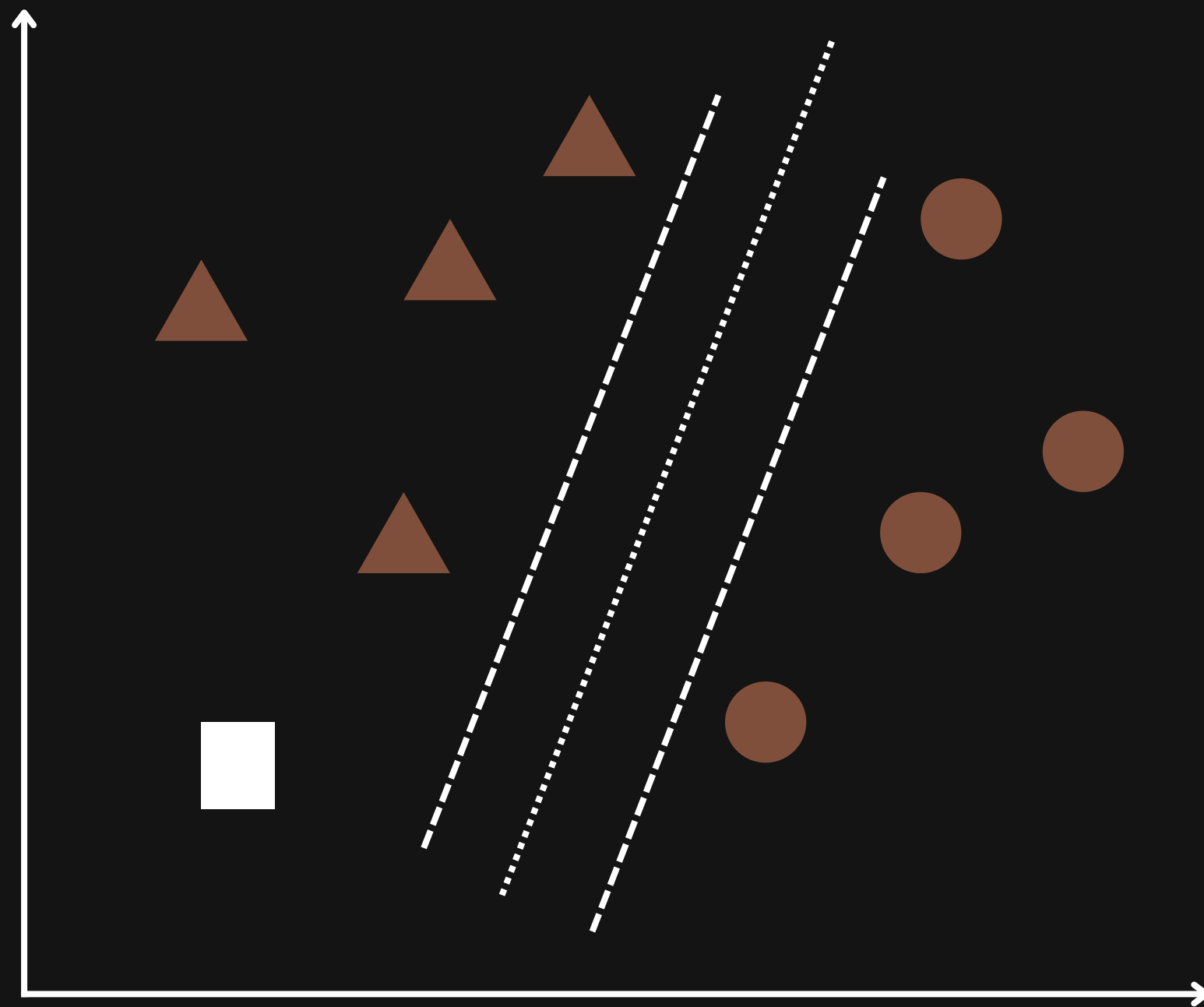This is going to be a triangle since it lies left to the decision boundary

Let us try to classify it using a SVM classifier.Clearly, a line that goes midway between the two figures would be an ideal classifier
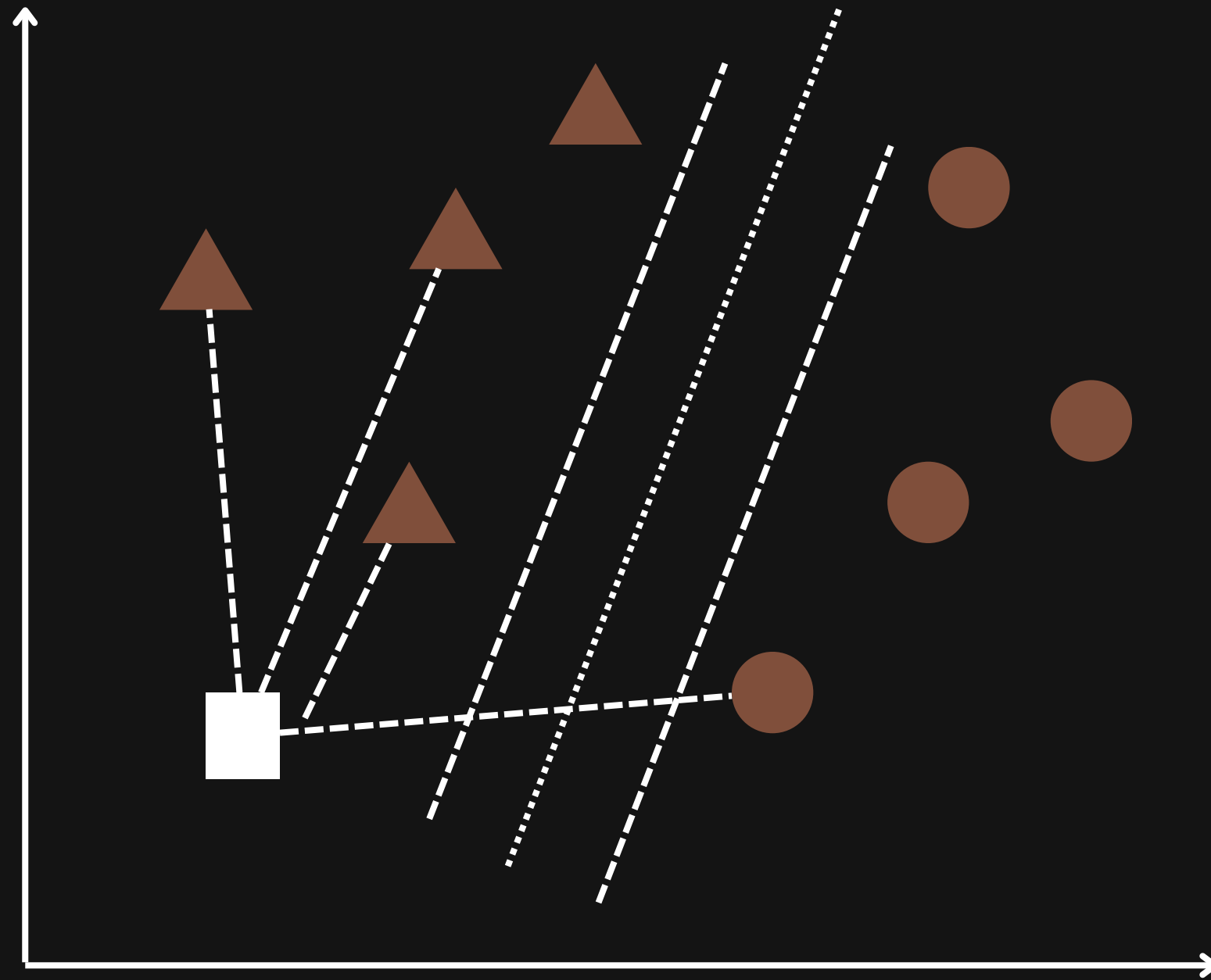
If the new point introduced falls into the left half plane,it is going to be closer to the triangle lying in the left plane.Infact it is intuitively easy to see that all points that lie to the left of the boundary lie closer to the triangle than the circle.

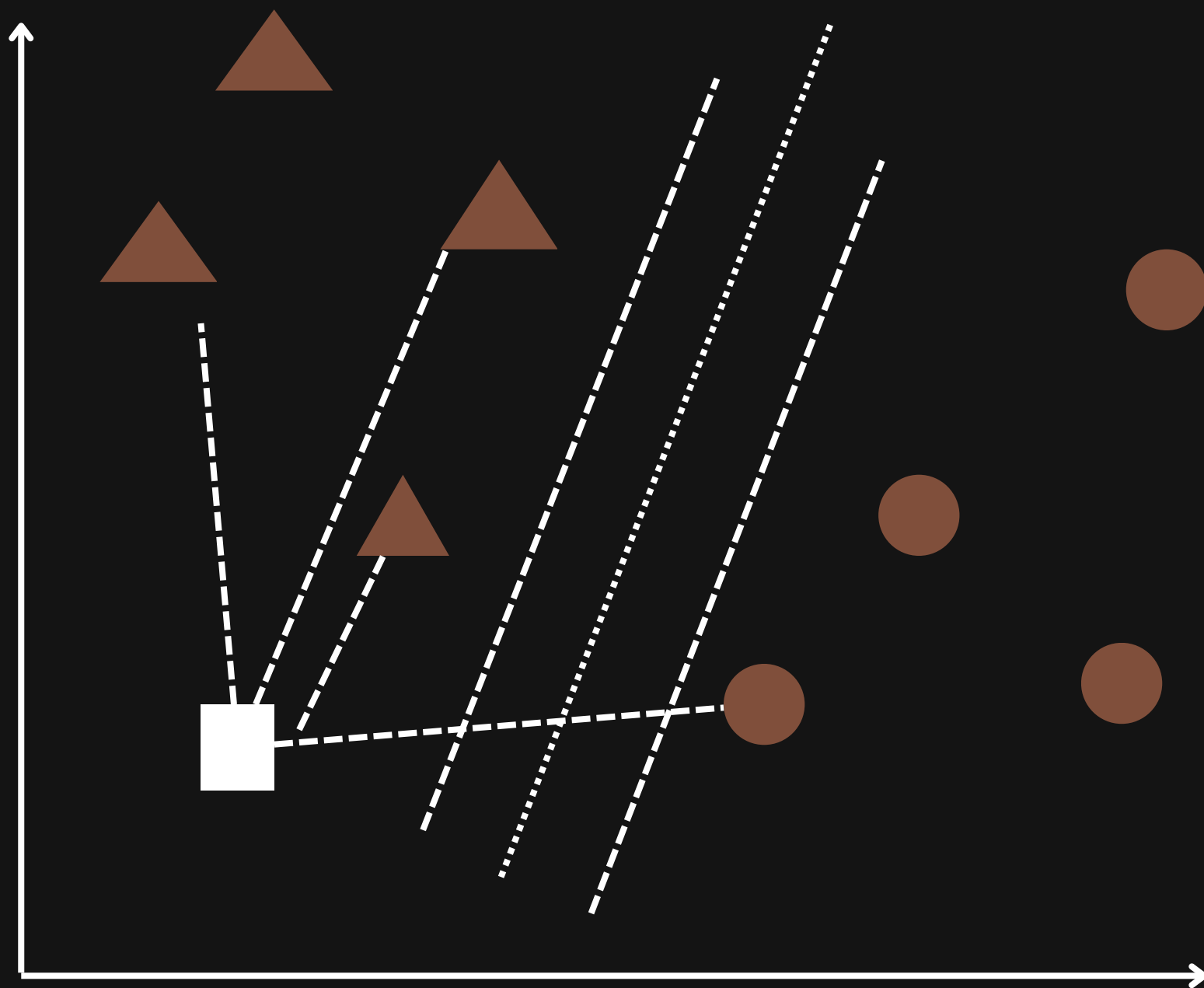This is going to be a triangle since it is closer to the triangle than the circle.

So another way of thinking to classify the unknown object as follows
Find out the nearest known object closer to it;And then classify it under the same label.;

Here we have a slightly more complex example.We clearly can see that the newly introduced point is a triangle,using SVM classifier,but can we use the alternative approach to figure out its label?

Ii should seem to work;Afterall the points on the left of the decision boudary seems to  have more  no of triangles as its "neighbours" than circles and hence can be classified as a triangle.
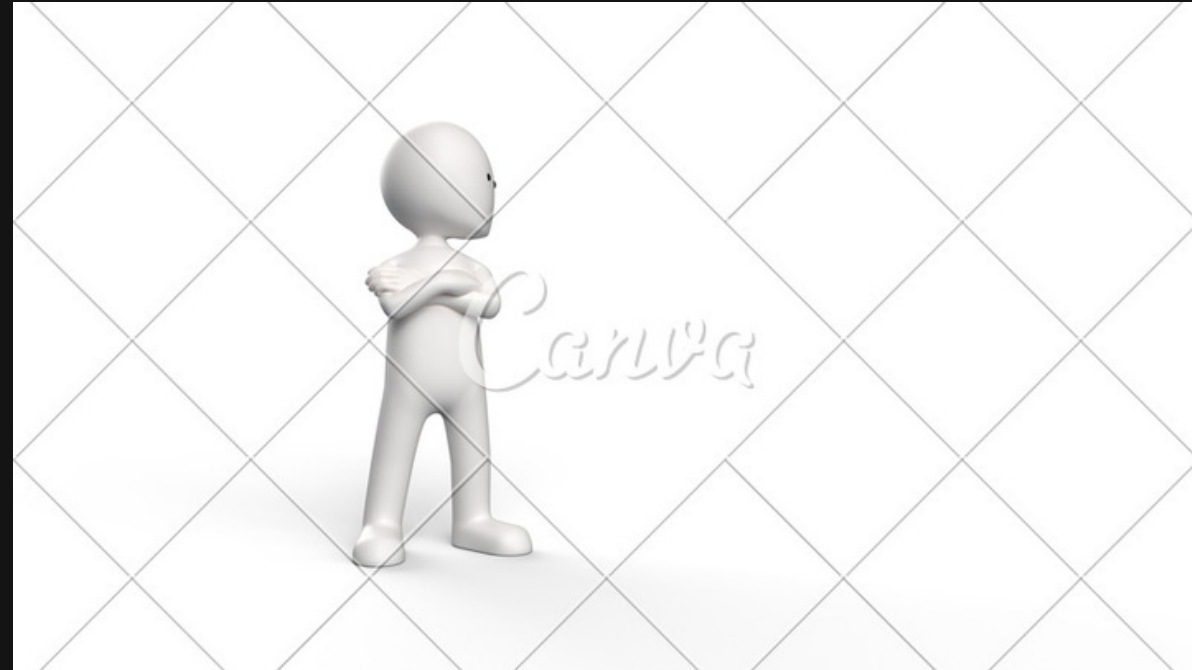
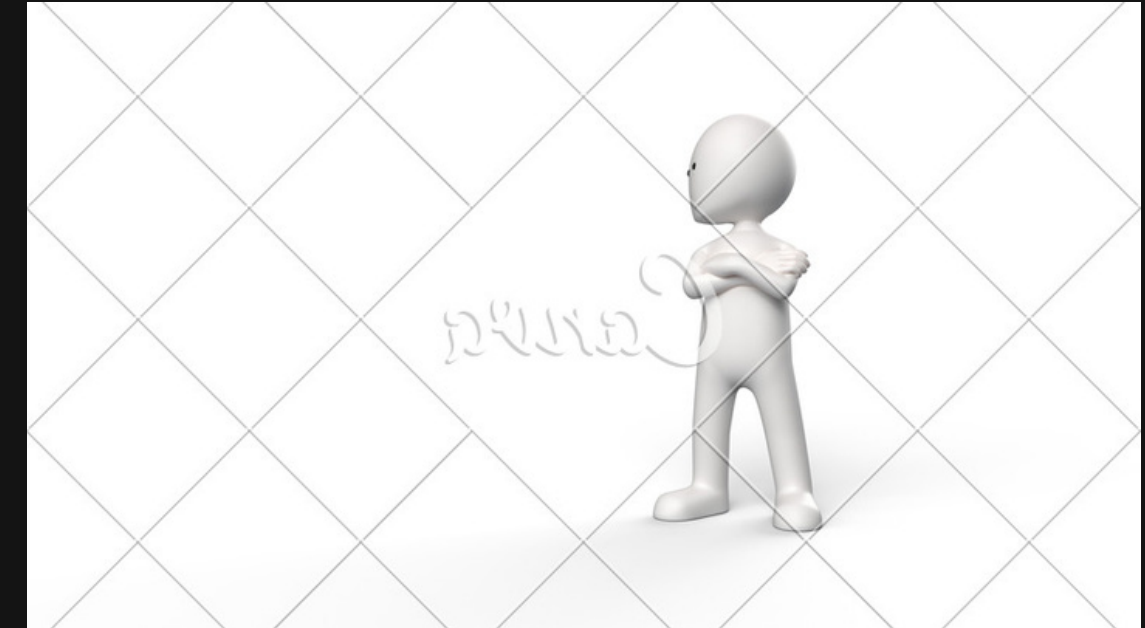But in this case,there is an ambiguity;How many neighbours are we to consider?

# S-MAN

# K-MAN





"The man has to be french since he is within the territorial region of france."

"The man has to be french since the most of the guys he is interacting with is french."

Note the two distinct ways in which the two men derive their conclusion;Both could be wrong,both could be correct or only one of them might be accurate.

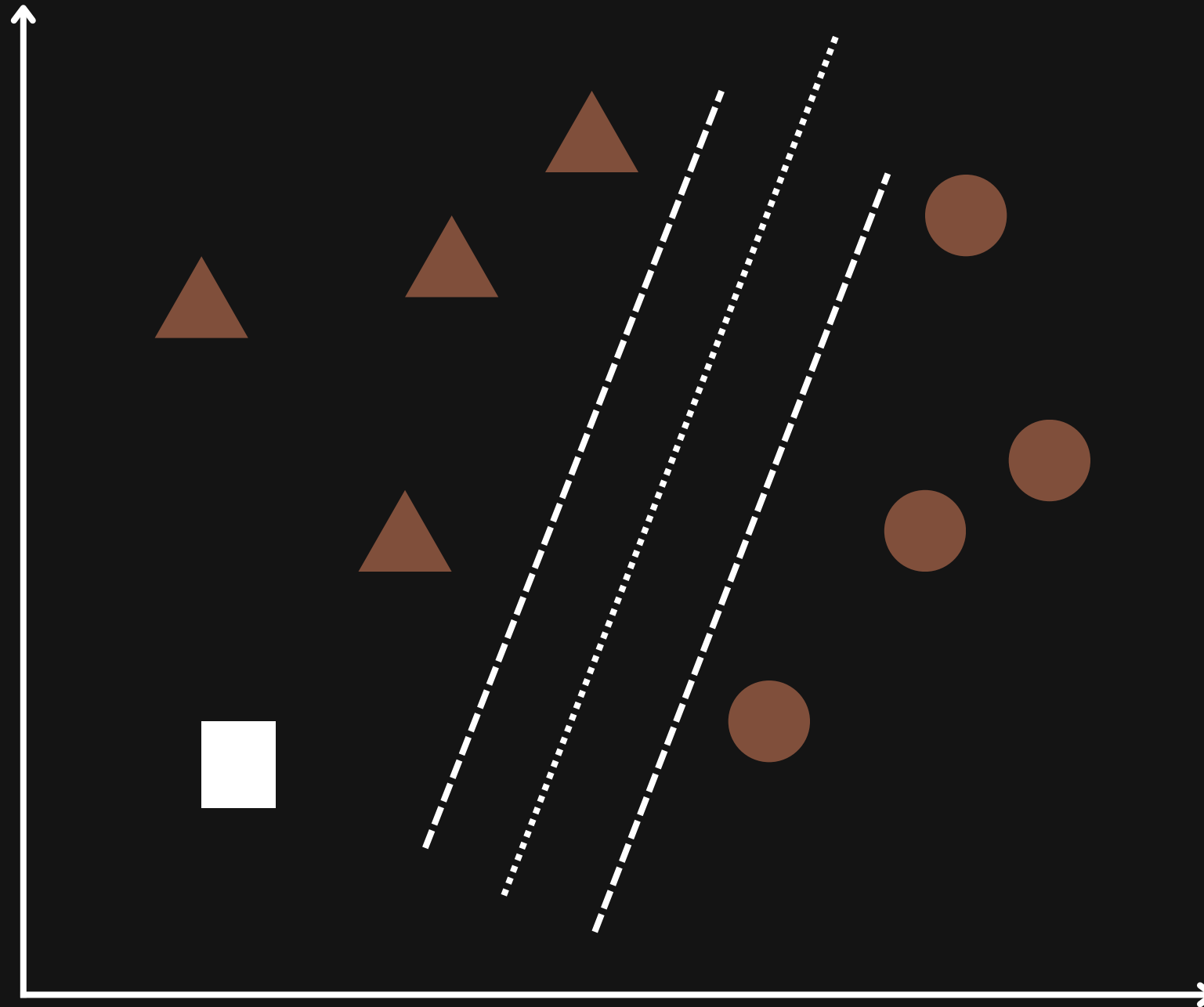Now u know the basics of KNN algo!!Well done

# KNN Algorithm

In simple Words KNN look for it's K nearest neighbours in by their  distances and then mark the class of the point according to the majority of the neighbours.

# Features of KNN

Lazy Learning

Instance-based learning

Non -Parametric:

# How do u find the optimal value of K?

Choosing the K value is the hardest part of building a KNN model. As a rule of thumb we can choose K values as sqrt(n) Where 'n' is the no of datapoints . But The best way to do it is make a graph between loss vs K and choosing the correct K value for the implementation of the model.

# How do we measure the distance to figure out the closeness of the objects?

$$\text{Minkowski Distance} = \left( \sum_{i=1}^{n} \left| x_i - y_i \right|^{p} \right)^{1/p}$$

$p = 2$

**Euclidean distance :**

$$d(x,y) = \sqrt{\sum_{i=1}^{n}(y_i - x_i)^2}$$

**Manhattan distance:**

$p = 1$

$$\text{Manhattan Distance} = d(x,y) = \left(\sum_{i=1}^{m}|x_i - y_i|\right)$$

# Limitation of KNN

Time Complexity　　Space Complexity

# Naive Bayes Algo

Suppose we want to classify emails as spams or normal messages.We will do this on the basis of presence of certain words,namely dear,friend,lunch and money.

We collect all the messages we have and we go message by message and check if a particular word is present.

We will try to figure out the probability of  the message being a nrml mssge given it has the words-dear and friend.Then we will do the same for spams and will figure out which is larger.

We can say that

P(message being a spam given it consists of dear,friend)

*P(message consisting of dear,friend)

=

P(message being a spam and consists of dear,friend)

=

P(message is a spam)

*P(message consisting of dear,friend given it is a spam)

P(message being a normal mssge given it consists of dear,friend)
*P(message consisting of dear,friend)

=
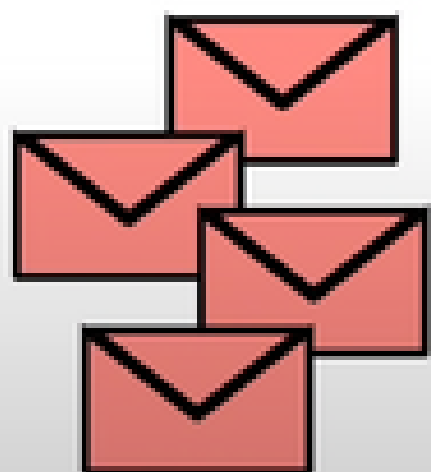
P(message being a nrml mssge and consists of dear,friend)

=

P(message is a nrml mssge)
*P(message consisting of dear,friend given it is a nrml mssge)

We start with histograms of all the words in the normal messages…

…and all of the words in the spam.

p( Dear | N ) = 0.47
p( Friend | N ) = 0.29
p( Lunch | N ) = 0.18
p( Money | N ) = 0.06

p( N ) = 0.67

Then we made an initial guess about the probability of seeing a **normal message**.

p( Dear | S ) = 0.29
p( Friend | S ) = 0.14
p( Lunch | S ) = 0.00
p( Money | S ) = 0.57

# Bayes Theorem

P (A and B happening)= P(A happening provided B happened)
* P(B happening)

P (A and B happening)= P(B happening provided A happened)
* P(A happening)

P(A happening provided B happened) * P(B happening)
=
P(B happening provided A happened) * P(A happening)

# Thus,

$$P(A|B) * P(B) = P(B|A) * P(A)$$

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

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

So far the feature that we considered(presence of a particular word) can only take  2 possible values-yes or no This is called mutinomial bayesian classification.
What if we had a feature that took continous values?

# Now let us add one more feature into consideration while classifying an email-No of words

| No of words | Type |
| --- | --- |
| 35 | SPAM |
| 40 | SPAM |
| 70 | NORMAL |
| 76 | NORMAL |
| 56 | SPAM |

Suppose we get a new message which has the words Dear(20)money(43)

Now we need to include one more factor into calculations P(message has 63 words given it is a nrml mssge) and P(message has 63 words given it is a spam)

That is,

P(mssge is a spam)*P(mssge consists of 63 words given it is a spam)
*P(mssge consists of word dear(20)money(43)

And

P(mssge is a nrml)*P(mssge consists of 63 words given it is nrml)
*P(mssge consists of word dear(20)money(43) given it is nrml)

| No of words | Type |
| --- | --- |
| 35 | SPAM |
| 40 | SPAM |
| 70 | NORMAL |
| 76 | NORMAL |
| 56 | SPAM |

But we see that there is no data available for letters that are 63 words long.

We cannot expect that either,since no of words is a continous (approximately) quantity and we cannot have data which includes all the possible values that the "no of words" feature can take.

Note that the mean value for the no of words in spams is   43.67

Note that the standard deviation for the no of words in spams  is   10.96

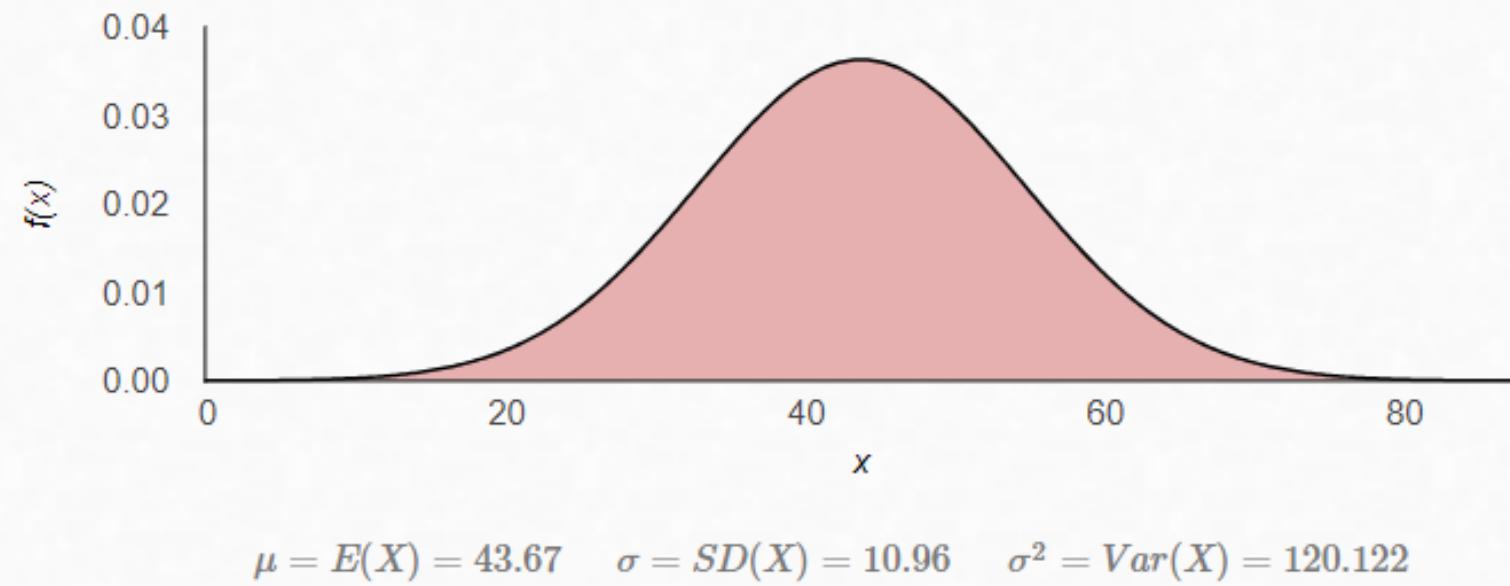Note that the mean value for the no of words in normal msg is 73

Note that the standard deviation for the no of words in normal mssge  is 3

Note that the mean value for the no of words in spams is   43.67

Note that the standard deviation for the no of words in spams  is   10.96

Note that the mean value for the no of words in normal msg is 73

Note that the standard deviation for the no of words in normal mssge  is 3

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$
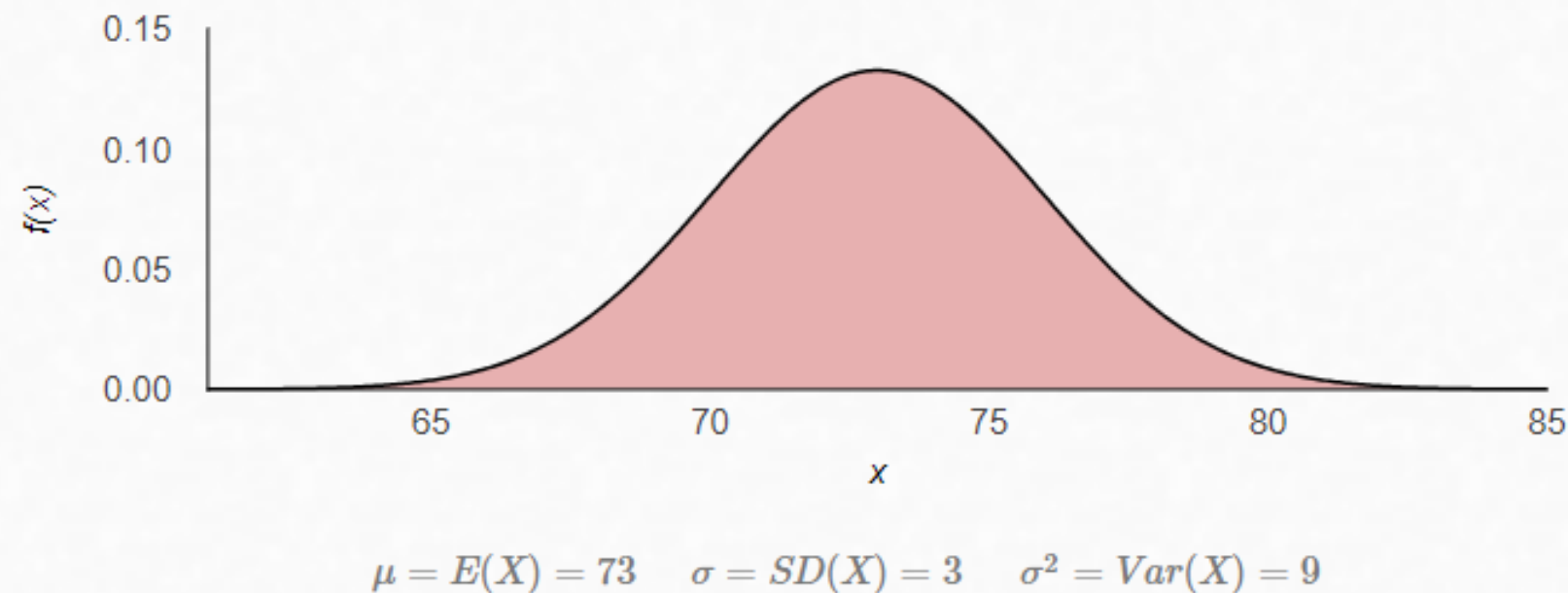
$\sigma$    = standard deviation

$\mu$    = mean

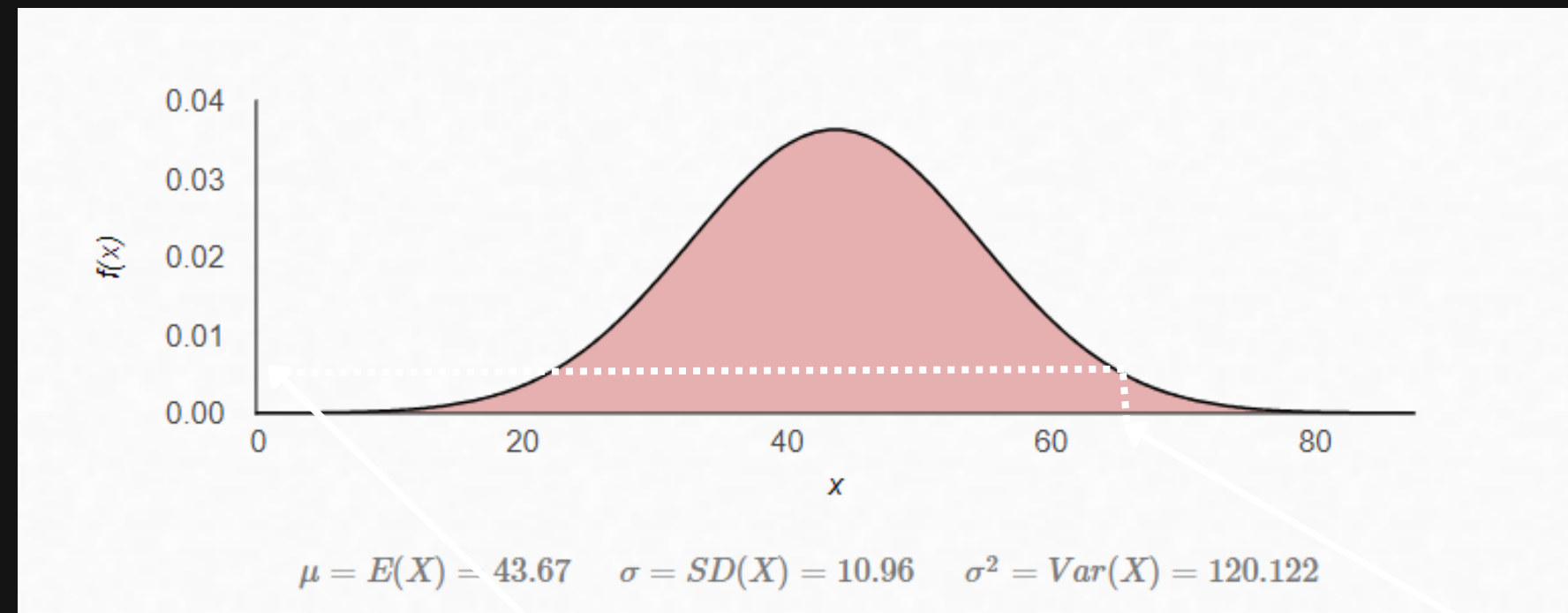If we plug these values into this formula and  plot the respective graphs,

Spam



$\mu = E(X) = 43.67 \quad \sigma = SD(X) = 10.96 \quad \sigma^2 = Var(X) = 120.122$

The value in the y axis represents the probability that message has x words given it is a spam.

Nrml msge



$\mu = E(X) = 73 \quad \sigma = SD(X) = 3 \quad \sigma^2 = Var(X) = 9$

Suppose we need the
P(mssge has 63 words|spam)

$\mu = E(X) = 43.67 \quad \sigma = SD(X) = 10.96 \quad \sigma^2 = Var(X) = 120.122$

63

P(mssge has 63 words|spam)

**Lunch Money Money Money Money**

p( Dear | N ) = 0.47
p( Friend | N ) = 0.29
p( Lunch | N ) = 0.18
p( Money | N ) = 0.06

p( N ) = 0.67

...and that means we will always classify the messages with **Lunch** in them as normal, no matter how how many times we see the word **Money**.

$$p( N ) \times p( \textbf{Lunch} | N ) \times p( \textbf{Money} | N )^4 = 0.000002$$

$$p( S ) \times p( \textbf{Lunch} | S ) \times p( \textbf{Money} | S )^4 = 0$$

p( Dear | S ) = 0.29
p( Friend | S ) = 0.14
**p( Lunch | S ) = 0.00**
p( Money | S ) = 0.57

p( S ) = 0.33

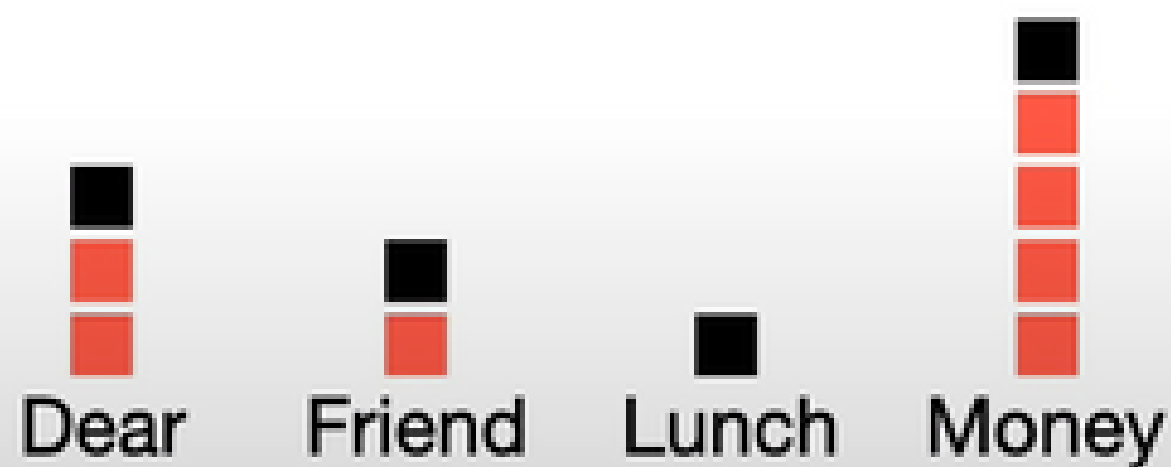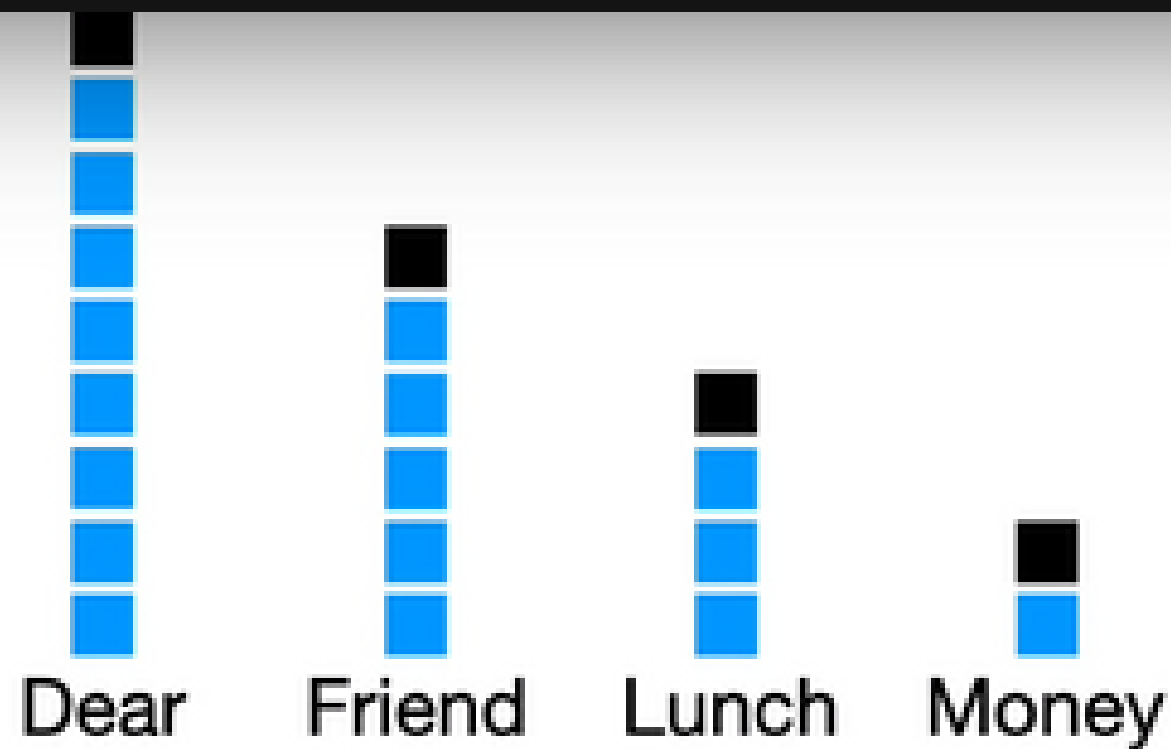p( N ) = 0.67

p( S ) = 0.33

To avoid this problem we add one more count to each feature,so that none of the Probabilities  vanishes.

This is called "smoothing"

Now let us consider two sentences with the same words

1)" John , you do take care of your family ,hey?"

2)"Hey John ,do you take care of your family?"

Note that 2 sentences constructed with same words

can appear in different contexts.

Since Naive Bayes does not care about the order of the words

It ignores the semantics and context of the sentence.

This is why it is "naive".