



UANL®



FCFM

CLUSTERING

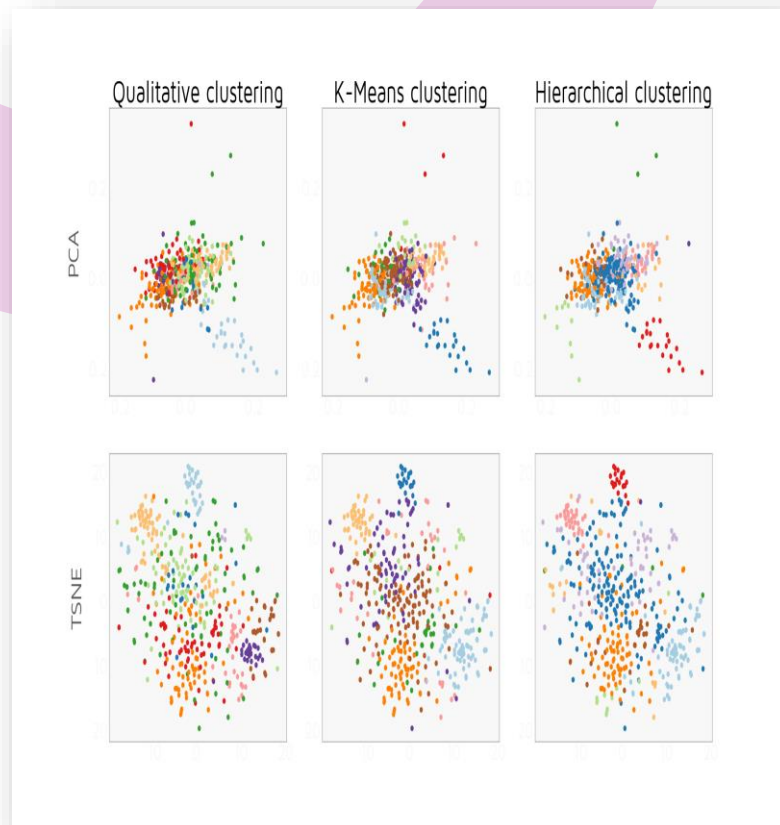
DATA MINING

Anahí Alemán Alvarado	1821952
Ricardo Zarek Sánchez Olivares	1795134
Juan Pablo Nasser Benavides	1753367
Eduardo Almaguer Alanís	1741322
Oscar Saúl Vega Macias	1626997

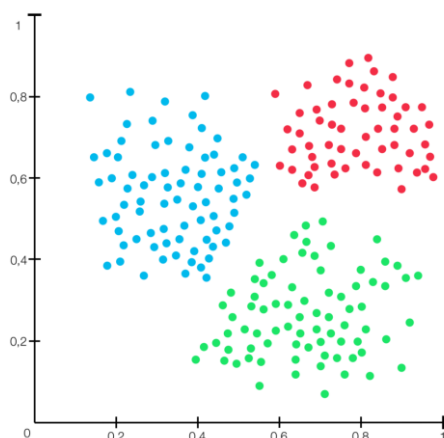
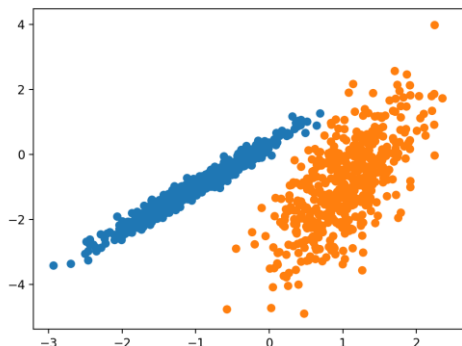
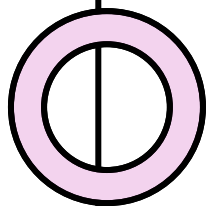


¿Qué es?

También conocido como agrupamiento, es una de las técnicas de minería de datos, el proceso consiste en la división de los datos en grupos de objetos similares.



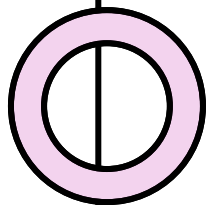
Las técnicas de Clustering son las que utilizando algoritmos matemáticos se encargan de agrupar objetos. Usando la información que brindan las variables que pertenecen a cada objeto se mide la similitud entre los mismos, y una vez hecho esto se colocan en clases que son muy similares internamente y a la vez diferente entre los miembros de las diferentes clases.



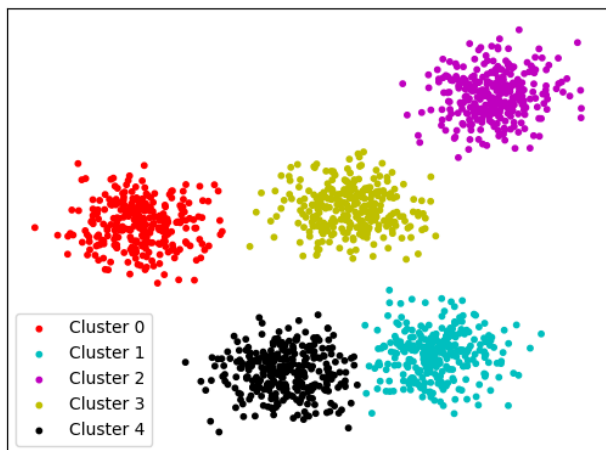
Un **cluster** es una colección de objetos de datos. Similares entre sí dentro del mismo grupo. Disimilar a los objetos en otros grupos.



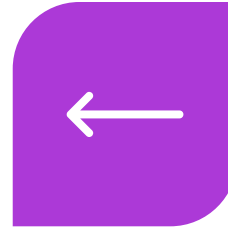
Análisis de cluster: dado un conjunto de puntos de datos tratar de entender su estructura. Encuentra similitudes entre los datos de acuerdo con las características encontradas en los datos. Es un aprendizaje no supervisado ya que no hay clases predefinidas.



Aplicaciones



Métodos de agrupación



ASIGNACIÓN
JERÁRQUICA
FRENTE A PUNTO



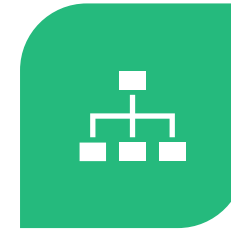
DATOS
NUMÉRICOS Y/O
SIMBÓLICOS



DETERMINÍSTICA
VS.
PROBABILÍSTICA



EXCLUSIVO VS.
SUPERPUESTO



JERÁRQUICO VS.
PLANO.



DE ARRIBA A
ABAJO Y DE
ABAJO A ARRIBA

The background features a white canvas with several decorative elements. On the left, there are two horizontal rows of black zigzag lines. A large, light pink curved shape sweeps across the left side. In the top right corner, there is a pink circle with a horizontal line through its center, resembling a binder ring. At the bottom right, there are four parallel black diagonal lines. The central text is enclosed in a white rectangular box with a thin black border and a pink drop shadow.

ALGORITMOS DE CLUSTERING



Simple K-Means

Este algoritmo debe definir el número de clusters que se desean obtener, así se convierte en un algoritmo voraz para particionar. Pasos:

1

Se determina la cantidad de clusters en los que se quiere agrupar la información, en este caso las simulaciones.

2

Se asume de forma aleatoria los centros por cada clusters. Una vez encontrados los primeros centroides el algoritmo hará los tres pasos siguientes:

- Determina las coordenadas del centroide.
- Determina la distancia de cada objeto a los centroides.
- Agrupa los objetos basados en la menor distancia.

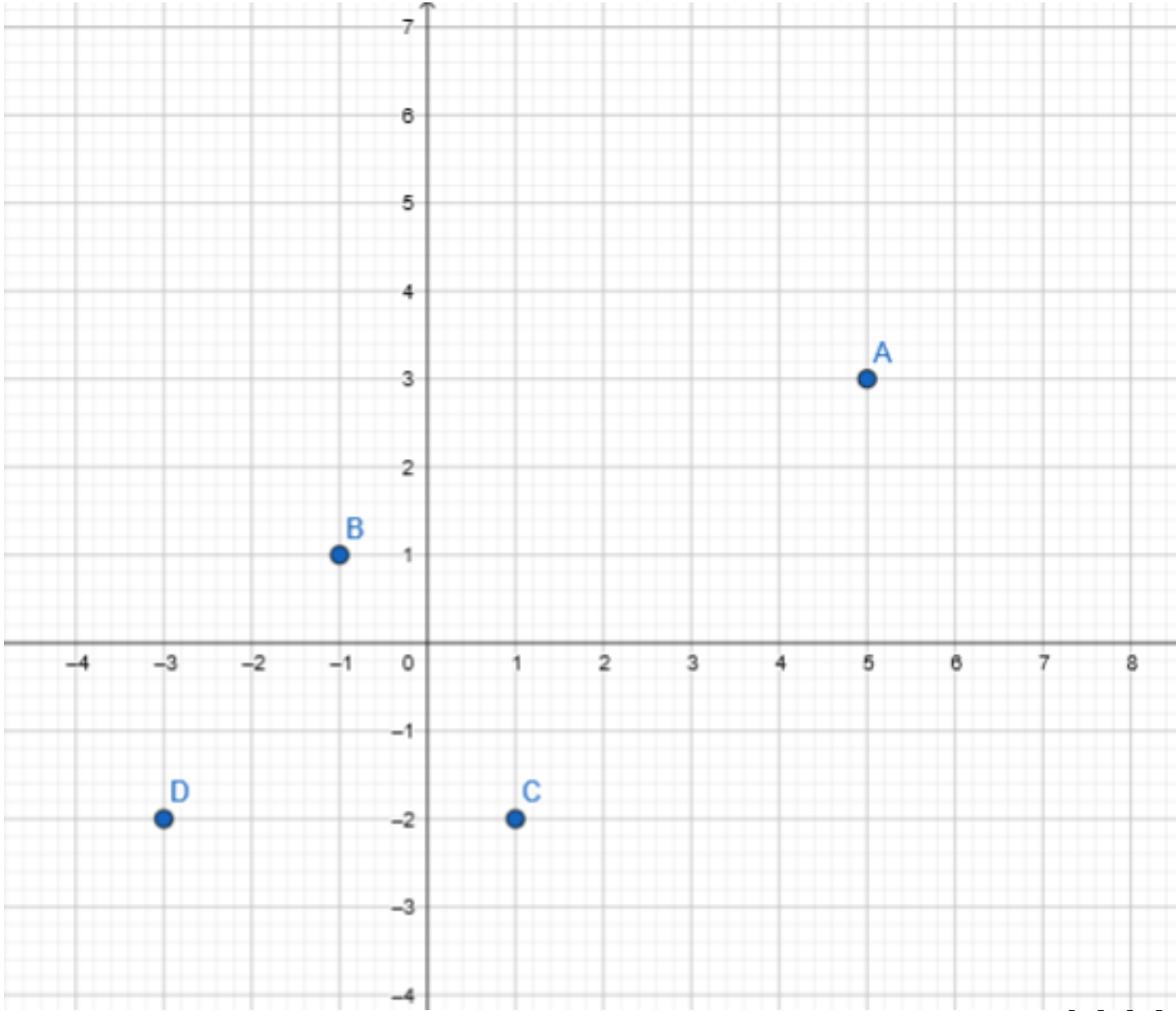
3

Finalmente quedarán agrupados por clusters, los grupos de simulaciones según la cantidad de clusters que el investigador definió en el momento de ejecutar el algoritmo.





	Q1	Q2
A	5	3
B	-1	1
C	1	-2
D	-3	-2

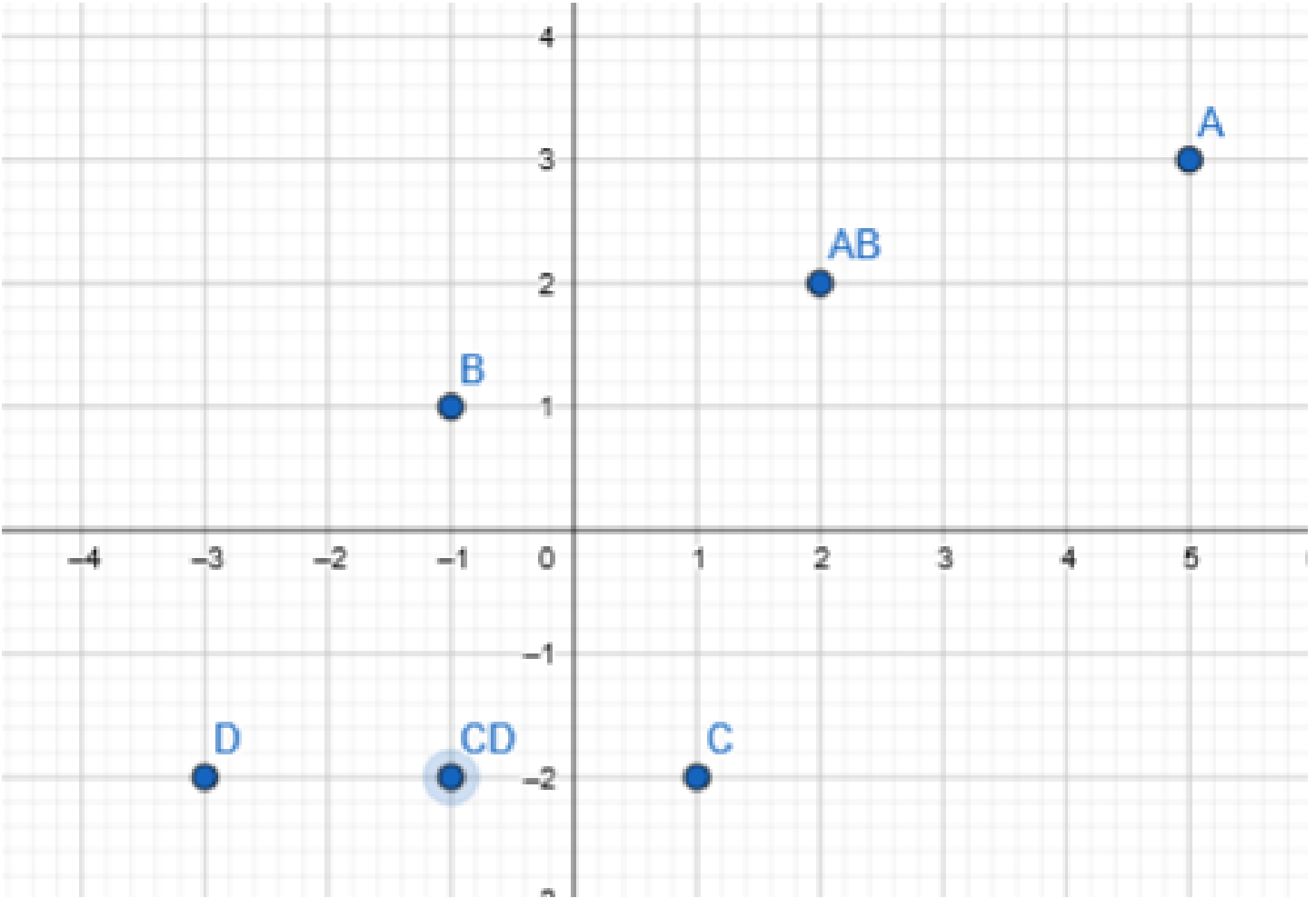




DEFINIMOS EL VALOR DE K

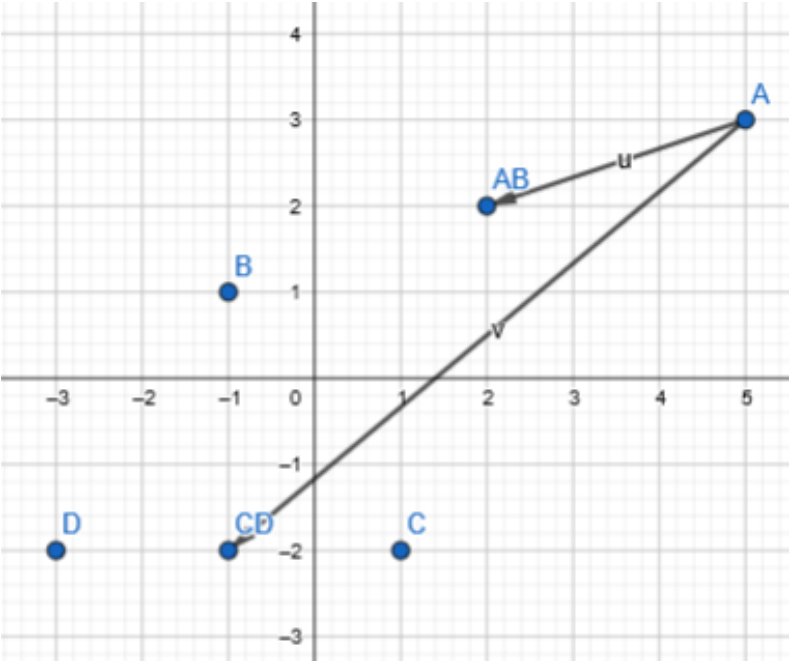
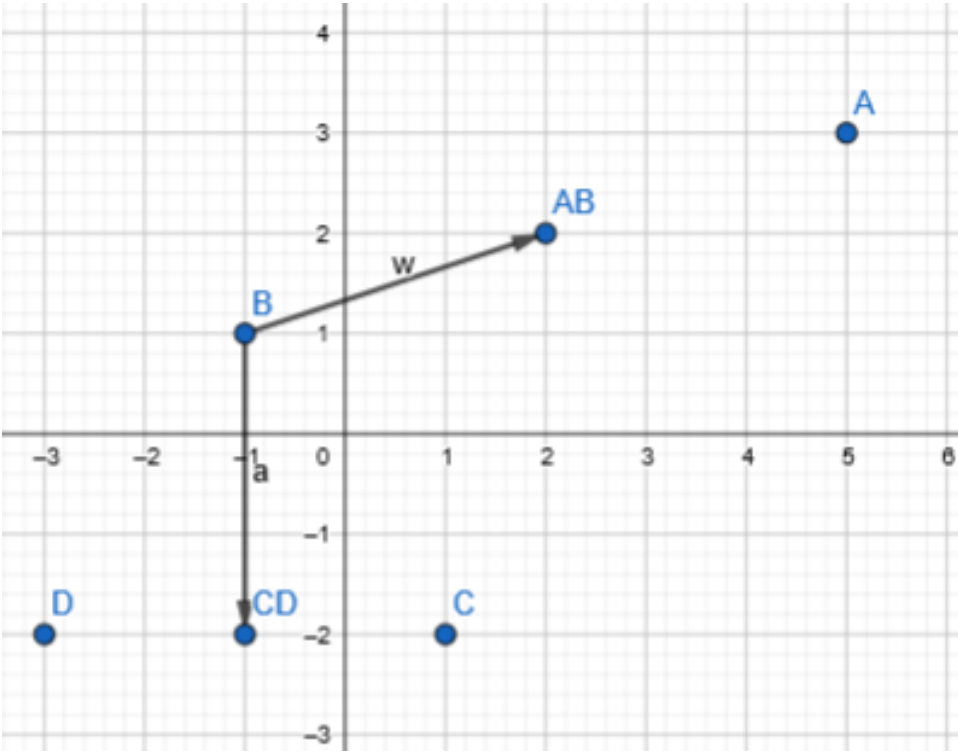
K	2

	AB	CD
Q1	2	-1
Q2	2	-2



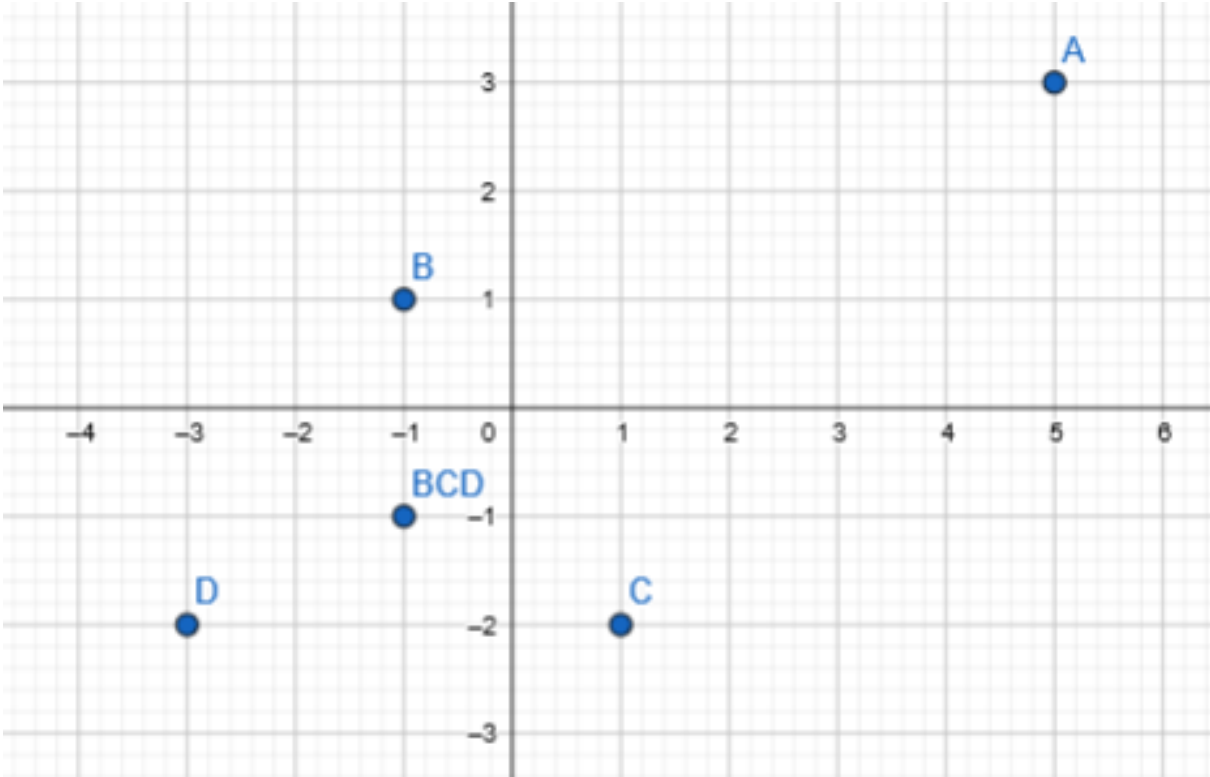


	DISTANCIA
A,(AB)	10
A,(CD)	61
B,(AB)	10
B,(CD)	9



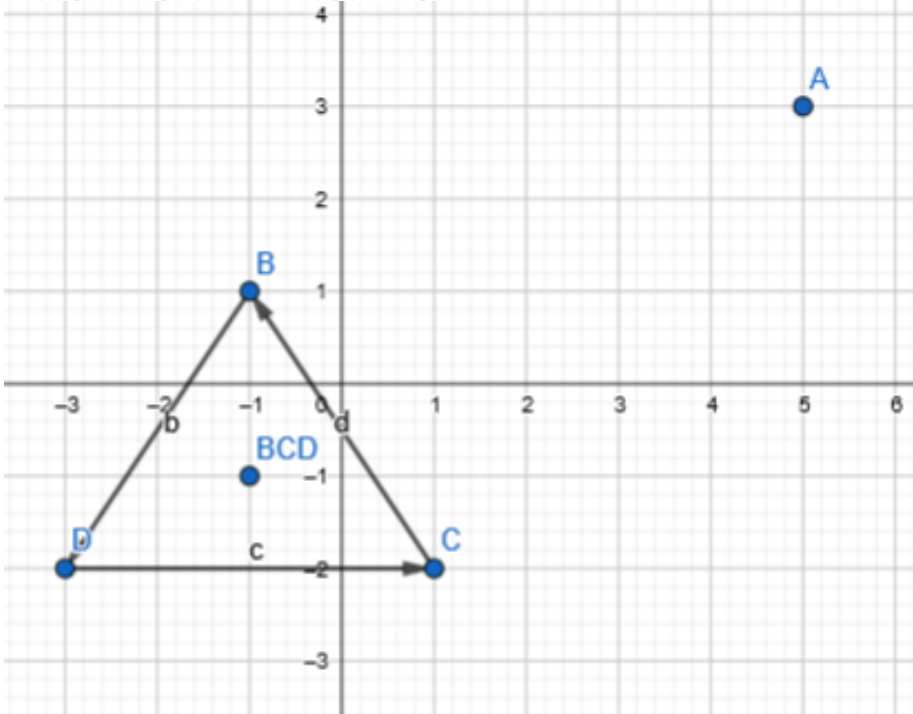
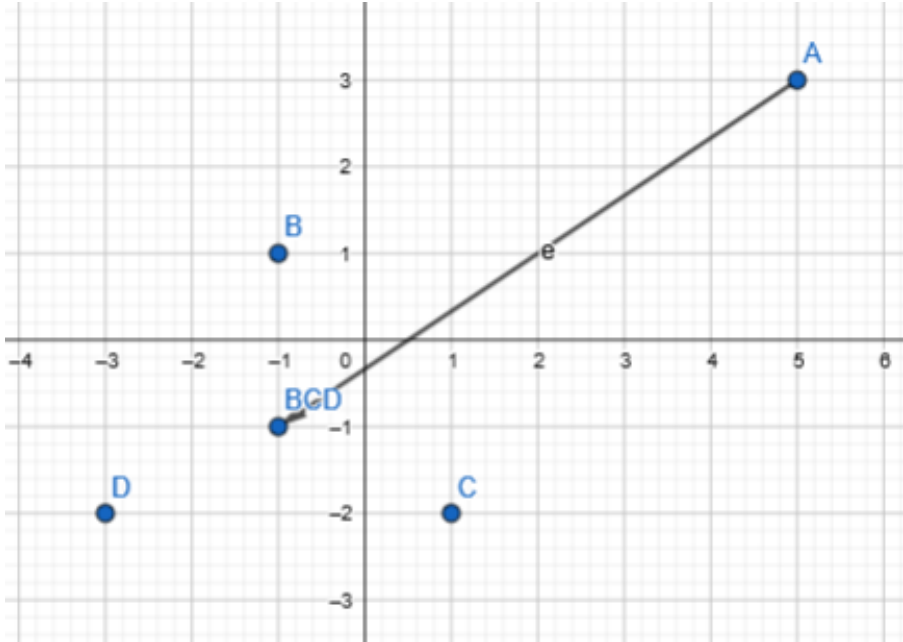


	A	BCD
Q1	5	-1
Q2	3	-1





	DISTANCIA
A,(BCD)	52
B,(BCD)	4
C,(BCD)	5
D,(BCD)	5

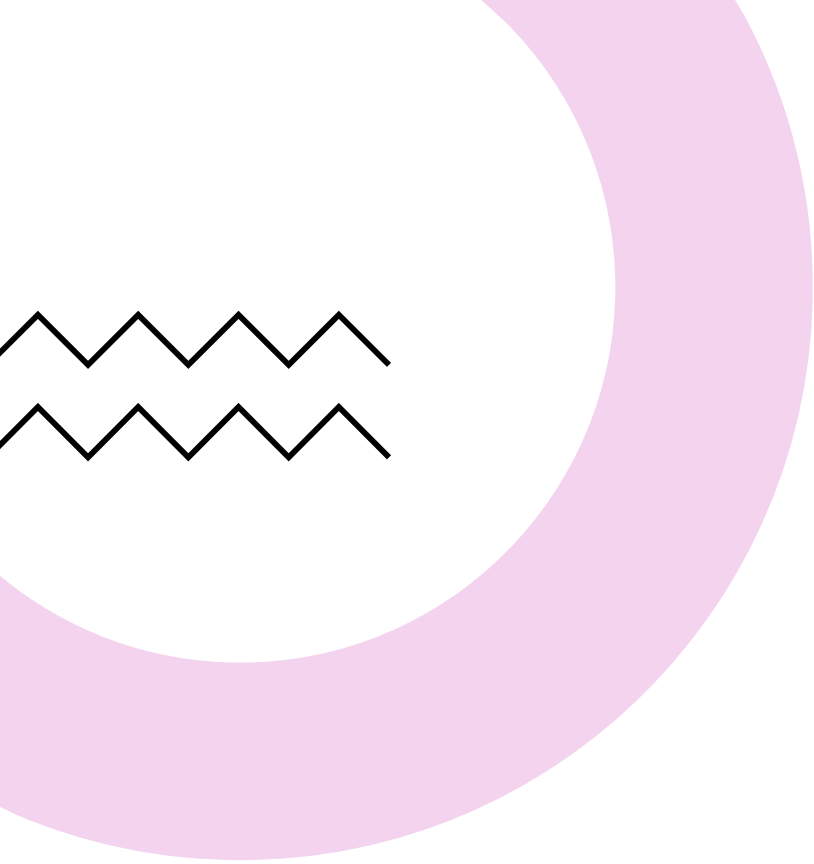


○ X-Means

Este algoritmo es una variante mejorada del **K-Means**.

Su ventaja fundamental está en haber solucionado una de las mayores deficiencias presentadas en K-Means, el hecho de tener que seleccionar a priori el número de clusters que se deseen obtener, a **X-Means** se le define un límite inferior **K-min** (número mínimo de clusters) y un límite superior **K-Max** (número máximo de clusters) y este algoritmo es capaz de obtener en ese rango el número óptimo de clusters, dando de esta manera más flexibilidad al usuario.





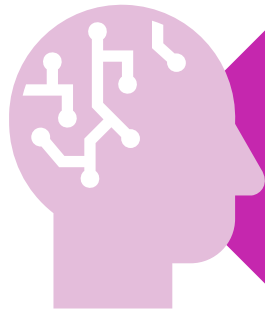
Cobweb

Pertenece a la familia de algoritmos **jerárquicos**. Se caracteriza por la utilización de aprendizaje incremental, esto quiere decir, que realiza las agrupaciones instancia a instancia. Durante la ejecución del algoritmo se forma un árbol (**árbol de clasificación**) donde las hojas representan los segmentos y el nodo raíz engloba por completo el conjunto de datos.

Además, en el algoritmo también hay que tener en cuenta dos parámetros muy importantes:



Acuity: es un parámetro muy necesario, pues la utilidad de categoría está basada en la estimación de la media y la desviación estándar del valor de un atributo para un nodo en particular.



Cut-off: este parámetro es usado para evitar el crecimiento descontrolado de la cantidad de segmentos. Indica el grado de mejoría que se debe producir en la utilidad de categoría





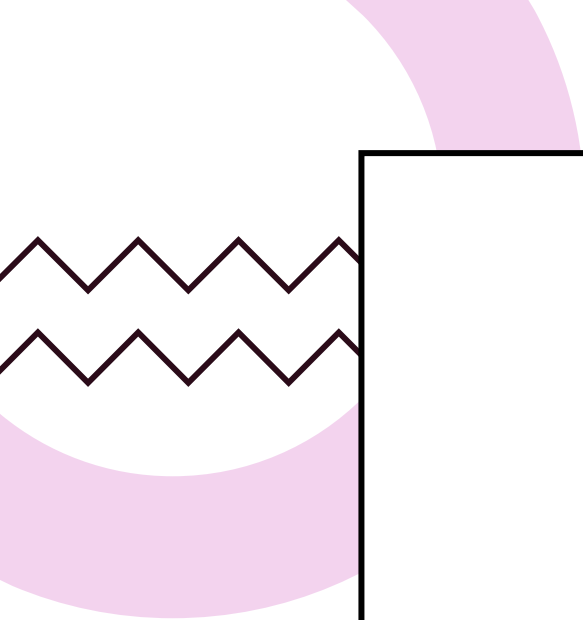
EM

Este algoritmo pertenece a una familia de modelos que se conocen como Finite Mixture Models, los cuales se pueden utilizar para segmentar conjuntos de datos. Está clasificado como un método de particionado y recolocación, o sea, **Clustering Probabilístico**. Se trata de obtener la FDP (Función de Densidad de Probabilidad) desconocida a la que pertenecen el conjunto completo de datos.

El algoritmo EM, procede en dos pasos que se repiten de forma iterativa:

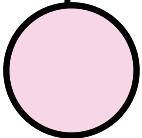
- Expectation: Utiliza los valores de los parámetros, iniciales o proporcionados por el paso Maximization, obteniendo diferentes formas de la FDP buscada.
- Maximization: Obtiene nuevos valores de los parámetros a partir de los datos proporcionados por el paso anterior.

Finalmente se obtendrá un conjunto de clusters que agrupan el conjunto de proyectos original. Cada uno de estos cluster estará definido por los parámetros de una distribución

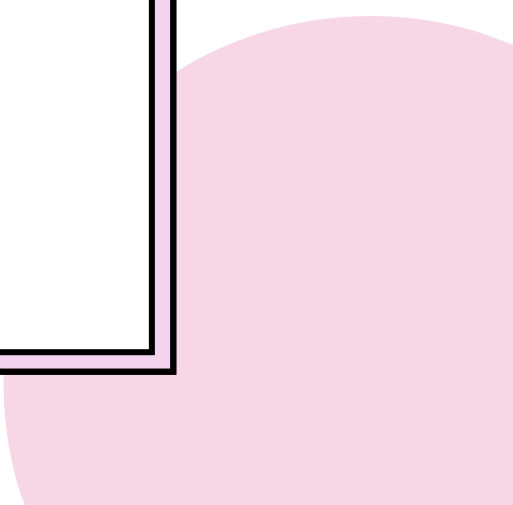


EJEMPLO: ALGORITMO K-MEDIAS

ES UN ALGORITMO DE APRENDIZAJE NO
SUPERVISADO.



SOLO SE LE DARÁ UN CONJUNTO DE VALORES Y EL
ALGORITMO SERÁ CAPAZ DE CLASIFICAR ESOS
VALORES EN CLASES, ¿EN CUANTAS?, A VECES ES
UN PROBLEMA SABER CUANTAS, A VECES NO



- Se parte de 2 vectores

```
In [85]: #X = [[3,5],[1,4],[1,6],[2,6],[1,5],[6,8],[6,6],[6,7],[5,6],[6,7],[7,1],[8,2],[9,1],[8,2],[9,3],[9,2],[8,3], ...]
v1=[0, 3, 1, 1, 2, 1, 6, 6, 5, 6, 7, 5, 6, 7, 8, 10, 9, 8, 9, 9, 9, 11, 13, 9, 15, 14, 13, 12, 14, 12]
v2=[5, 5, 4, 6, 6, 5, 8, 6, 7, 7, 7, 6, 7, 1, 2, , 1, 2, 3, 2, 3, 12, 11, 10, 13, 12, 13, 10, 10, 11]
```

- Los vectores se convierten a vectores tipo numpy.

```
x1 = np.array(v1)
x2 = np.array(v2)
```

- Se unen esos dos valores en vectores de dos elementos.

```
X = np.array(list(zip(x1, x2))).reshape(len(x1), 2)
print(X)
```

```
print(X)
```

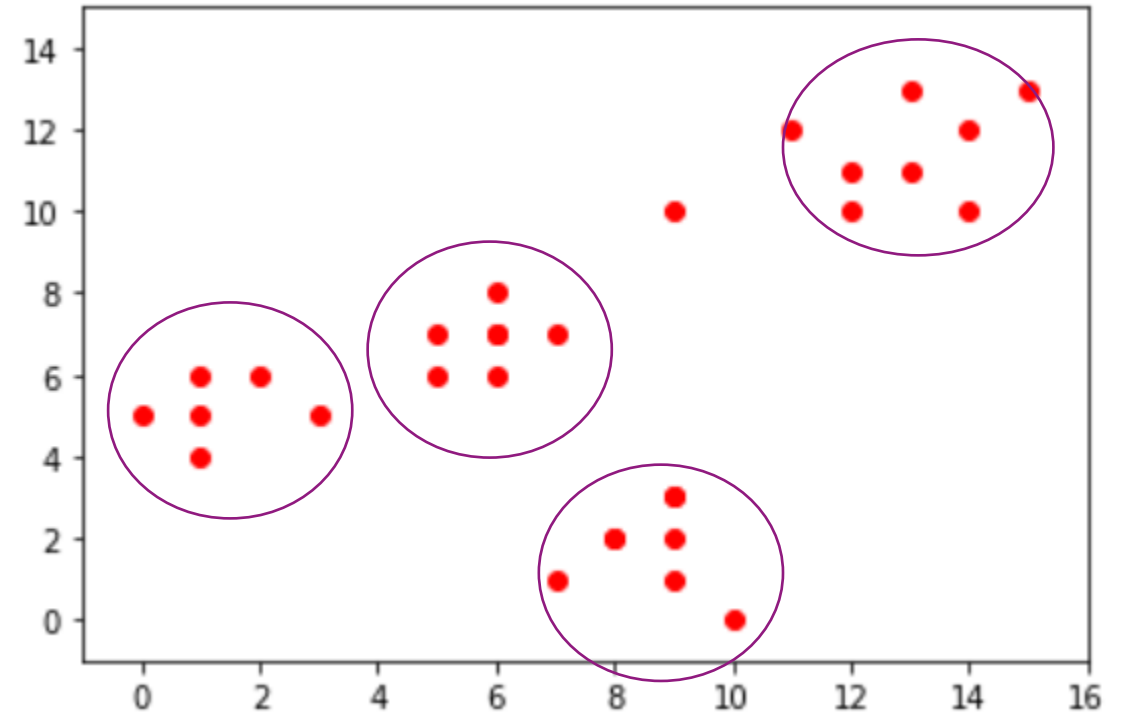
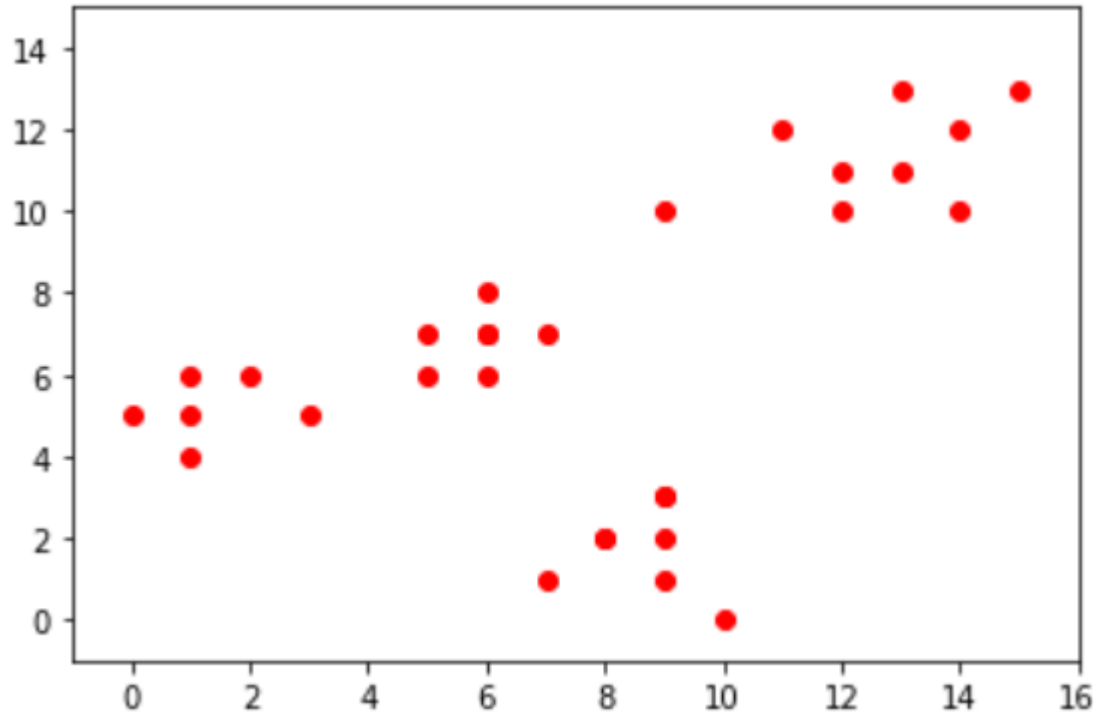
```
[[ 0  5] [ 7  7] [ 9 10]
 [ 3  5] [ 5  6] [15 13]
 [ 1  4] [ 6  7] [14 12]
 [ 1  6] [ 7  1] [13 13]
 [ 2  6] [ 8  2] [12 10]
 [ 1  5] [10  0] [14 10]
 [ 6  8] [ 9  1] [12 11]]
[ 6  6] [ 8  2]
[ 5  7] [ 9  3]
[ 6  7] [ 9  2]
[ 7  7] [ 9  3]
[ 5  6] [11 12]
[ 6  7] [13 11]
```



- Podemos visualizar los datos de la siguiente manera:

```
import matplotlib.pyplot as plt
plt.plot(v1, v2, 'ro')
plt.axis([-1, 16, -1, 15]) #Eje x: de -1 a 16; Eje Y: de -1 a 15
plt.show()
```

Al ver esta grafica podemos notar al menos 4 grupos diferentes.





- Para nuestro algoritmo el numero de clases que logramos apreciar son 4, por lo que hacemos $k=4$.
- Llamamos a la función K-Means y se le da como parámetros el numero de clases.
- El `fit(x)` lo que hace es entrenar el sistema con los valores utilizados.

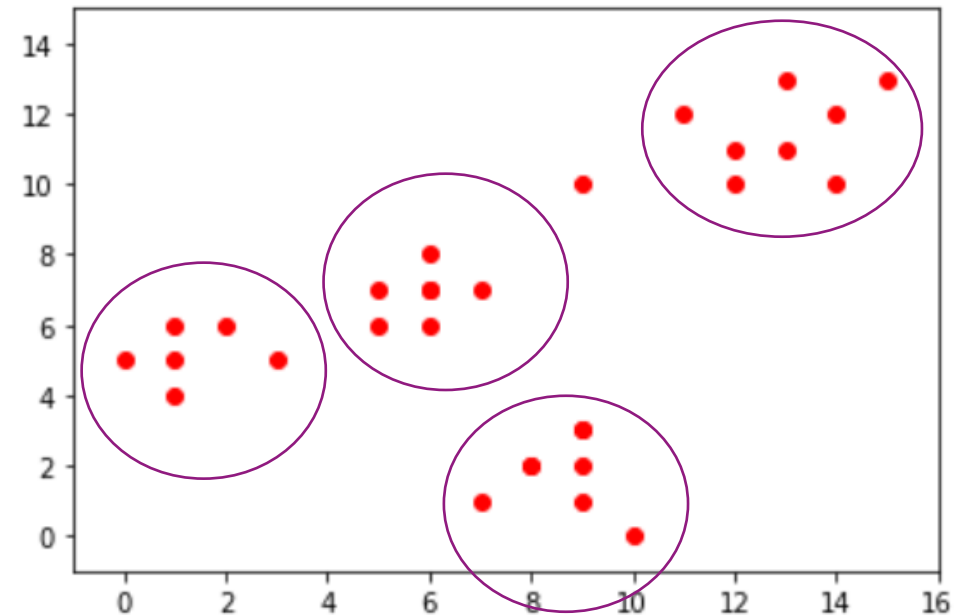
```
#Algoritmo K-Means  
K = 4  
kmeans_model = KMeans(n_clusters=K).fit(X)
```



- Si se quiere ver cuantas clases ha encontrado se puede hacer de la siguiente manera:

```
for i, l in enumerate(kmeans_model.labels_):  
    print("(x1,x2) -> Clase")  
    print("({0},{1}) ->:{2}".format(x1[i], x2[i], l))
```

```
(x1,x2) -> Clase (x1,x2) -> Clase (7,1) ->:0 (x1,x2) -> Clase  
(0,5) ->:2 (6,8) ->:3 (x1,x2) -> Clase (11,12) ->:1  
(x1,x2) -> Clase (x1,x2) -> Clase (8,2) ->:0 (x1,x2) -> Clase  
(3,5) ->:2 (6,6) ->:3 (x1,x2) -> Clase (13,11) ->:1  
(x1,x2) -> Clase (x1,x2) -> Clase (10,0) ->:0 (x1,x2) -> Clase  
(1,4) ->:2 (5,7) ->:3 (x1,x2) -> Clase (9,10) ->:1  
(x1,x2) -> Clase (x1,x2) -> Clase (9,1) ->:0 (x1,x2) -> Clase  
(1,6) ->:2 (6,7) ->:3 (x1,x2) -> Clase (15,13) ->:1  
(x1,x2) -> Clase (x1,x2) -> Clase (8,2) ->:0 (x1,x2) -> Clase  
(2,6) ->:2 (7,7) ->:3 (x1,x2) -> Clase (14,12) ->:1  
(x1,x2) -> Clase (x1,x2) -> Clase (9,3) ->:0 (x1,x2) -> Clase  
(1,5) ->:2 (5,6) ->:3 (x1,x2) -> Clase (13,13) ->:1  
 (x1,x2) -> Clase (9,2) ->:0 (x1,x2) -> Clase  
 (6,7) ->:3 (x1,x2) -> Clase (12,10) ->:1  
 (9,3) ->:0 (x1,x2) -> Clase  
 (14,10) ->:1  
 (x1,x2) -> Clase  
 (12,11) ->:1
```



- Para mostrar clases de cada valor:



```
predicciones = kmeans_model.predict(X)
print(predicciones)
```

```
[2 2 2 2 2 2 3 3 3 3 3 3 3 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1]
```

- Para valorar un único valor.

```
test=[[10,15]]
prediccion = kmeans_model.predict(test)
print("Prediccion {0}->{1}".format(test,prediccion))
```

```
Prediccion [[10, 15]]->[1]
```

```
x1 = np.array([2, 5, 8, 12])
x2 = np.array([5, 5, 4, 14])

X = np.array(list(zip(x1, x2))).reshape(len(x1), 2)
prediccion = kmeans_model.fit_predict(X)
print(prediccion)
```

```
[3 1 2 0]
```





- Se puede hacer que el sistema aprenda nuevos valores y que los intente clasificar.

```
x1 = np.array([2, 5, 8, 12])
x2 = np.array([5, 5, 4, 14])

X = np.array(list(zip(x1, x2))).reshape(len(x1), 2)
prediccion = kmeans_model.fit_predict(X)
print(prediccion)
```

```
[3 1 2 0]
```

- Se pueden representar los diferentes centros de la siguiente manera.

```
print(kmeans_model.cluster_centers_)
```

```
[[12. 14.]
 [ 5.  5.]
 [ 8.  4.]
 [ 2.  5.]]
```

