



Universidad Simón Bolívar
Departamento de computación
Inteligencia artificial CI5437

Proyecto 3

Clustering

Profesora :
Carolina Martínez

Grupo:
Carlos Farinha - 09-10270
Javier López - 11-10552
Nabil Márquez - 11-10683

16 de marzo del 2017

Introducción

K-means es un método de agrupamiento, que tiene como objetivo la partición de un conjunto de n observaciones en k grupos en el que cada observación pertenece al grupo cuyo valor medio es más cercano. Es un método muy utilizado en minería de datos debido a su capacidad de ubicar patrones dentro de un espacio definido de observaciones.

En el siguiente informe se presenta una aplicación del algoritmo k-means y su representación gráfica en el conjunto de datos “Iris Data Set” el cual cataloga 3 categorizaciones de una planta de Iris, el objetivo es que el algoritmo pueda generar la definición de a qué conjunto pertenece una observación en base a longitud y ancho del sépalo y pétalo de la flor, para categorizar a qué especie pertenece.

Implementación

Librerías y Funciones Externas Usadas

- *PIL - Image* : Generador de Imágenes jpg
- *Math - random* : Generador de valores pseudo-aleatorios
- *os* : Creador de directorios para el guardado de los gráficos generados
- *numpy* : Control de estructuras en el algoritmo k-means
- *matplotlib - pyplot* : Generador de Scatterplots de los conjuntos

Algoritmo K-means

Para la implementación del algoritmo k-means se creó el módulo *k-means.py* el cual contiene las siguientes funciones definidas:

- *puntos_del_cluster*: Clasifica los puntos según el centroide más cercano a la observación
- *centrar*: Calcula el punto medio de las observaciones clasificadas
- *iteracion*: Realiza una iteración de la clasificación de los puntos siempre que el cambio en el centroide varíe.

Programa principal

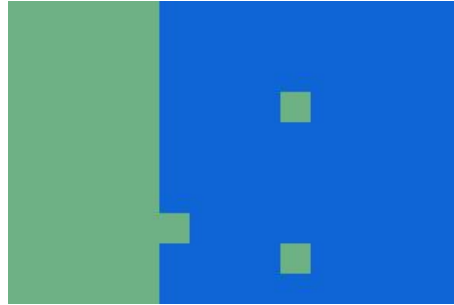
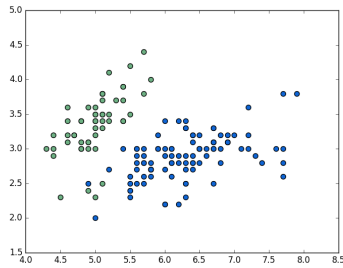
El archivo principal está ubicado como *main.py* y el cual lee los datos del archivo principal *bezdekIris.data* realiza una iteración continua centrando los datos y definiendo los nuevos centroides del conjunto en base a un *k* definido en la función, las iteraciones continúan formalmente hasta lograr una convergencia, una vez realizada la convergencia se comprime la información en un Scatter plot donde se toman 2 de las variables para presentar la información (longitud y ancho del sépalo) además genera una imagen comprimida de los resultados de cada uno de las observaciones, ambos gráficos son guardados en las carpetas “/Imágenes/Ejercicio_a” o “/Imágenes/Ejercicio_b” respectivamente dependiendo del enunciado.

Presentación de los resultados

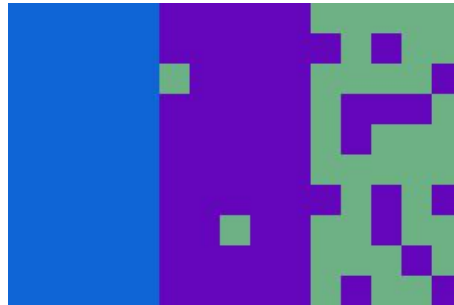
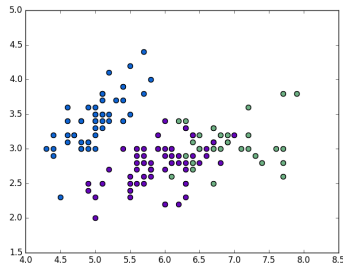
Para la presentación de los resultados se escogieron las mejores corridas del algoritmo en base a k a fin de poder presentar la información mas veraz del mismo.

Ejercicio a

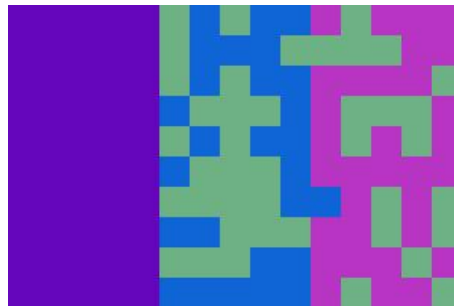
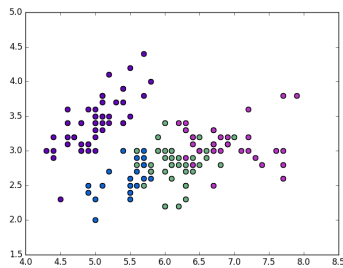
- $k = 2$



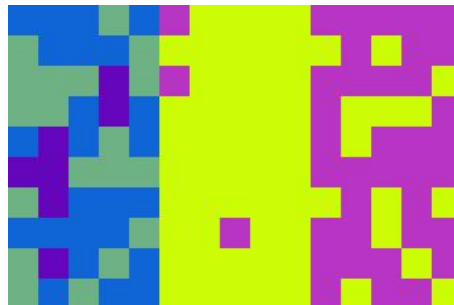
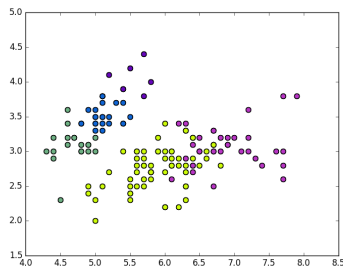
- $k = 3$



- $k = 4$



- $k = 5$



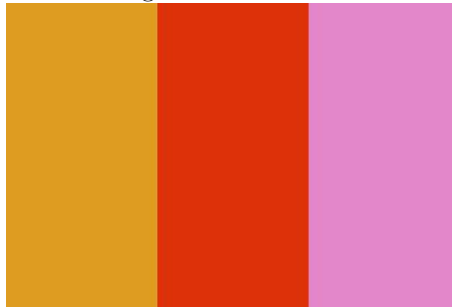
Se puede apreciar que para las corridas con $k = \{2,3\}$ el algoritmo puede clasificar con un alto grado de eficiencia el cluster al cual pertenece cada observación.

- $k = 2$ Para el caso de $k = 2$ el algoritmo genera un solo cluster que contiene 2 clasificaciones en conjunto a pesar de ello usualmente entre 3 a 5 de las 150 observaciones es catalogado de manera incorrecta teniendo en cuenta que el conjunto “Iris versicolor” e “Iris virginia” pertenecen al mismo cluster. Generando así un margen de error de 2% a 3% en la clasificación
- $k = 3$ Para el caso de $k = 3$ el algoritmo genera un solo cluster que contiene 2 clasificaciones en conjunto a pesar de ello usualmente entre 10 a 16 de las 150 observaciones es catalogado de manera incorrecta teniendo en cuenta que el conjunto “Iris setosa” es catalogado perfectamente en el 100% de los casos pero debido a la similitud de las observaciones de “Iris versicolor” e “Iris virginia” el algoritmo tiende a categorizarlas de manera errónea, se deben realizar múltiples corridas del algoritmo para poder tener en cuenta cual realmente es la clasificación de los últimos 2 conjuntos. Generando así un margen de error de 6% a 10% en la clasificación.

Ejercicio b

Originalmente el conjunto a estudiar se puede presentar en una imagen comprimida. En este caso cada pixel en la imagen original es un ejemplo, que será asignado en la imagen comprimida al cluster k_i . K será el número de colores en la imagen comprimida.

- Formalmente la imagen original debería ser de la siguiente forma:

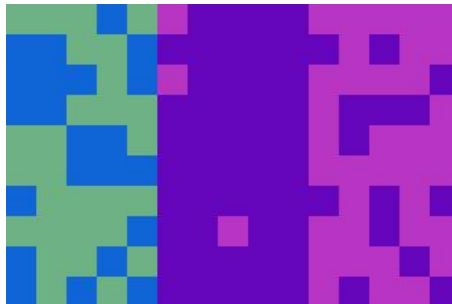


Se evalúa el algoritmo para los siguientes conjuntos:

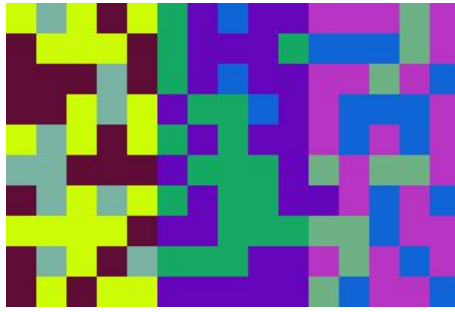
- $k = 2$



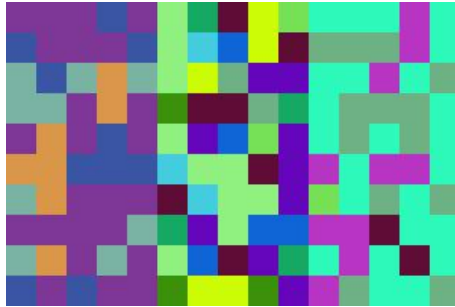
- $k = 4$



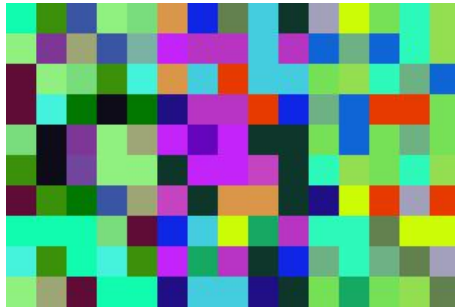
- $k = 8$



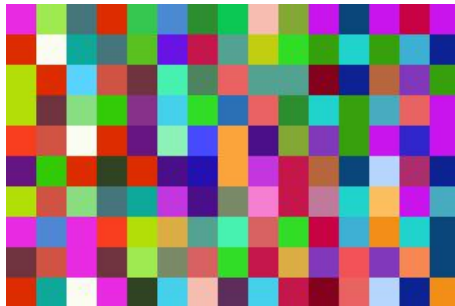
- $k = 16$



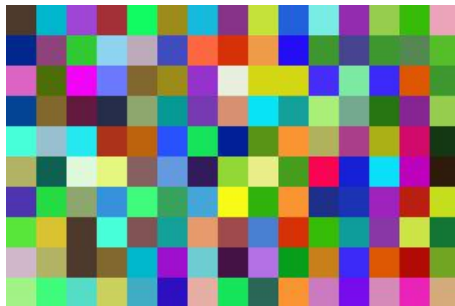
- $k = 32$



- $k = 64$



- $k = 128$



Conclusiones

Se pudo observar al realizar varias corridas del algoritmo que los resultados de k-means pueden variar drasticamente dependiendo de donde se inicialize los centroides de los clusters.

Mediante un Diagrama de Voronoi se podrian generar los espacios a los cuales cada conjunto se encuentra definido permitiendo predecir en que conjunto puede ser clasificado un punto.

En el caso de conjuntos que se encuentren muy proximos los casos bordes de dichos conjuntos pueden variar a la hora de categorizar los resultados.

Cuando se realiza una corrida con una gran cantidad de centroides los mismos convergiran a un solo valor en las observaciones creando overfitting d ela data.

En cualquiera de los casos superiores a $k=4$ los cluster que se generan son subconjuntos de los clusters cuando $k=3$ esto debido a que la tendencia predictiva del algoritmo converge hacia sonas similares definidas.