

Minería de datos **Práctica 1:Clustering** **knn-means**

Jose Ignacio Sánchez
Josu Rodríguez

4 de octubre de 2014

ÍNDICE DE CONTENIDO

1. Introducción	1
2. Recursos	1
3. Clasificación NO-supervisada o <i>Clustering</i>	1
3.1. Clustering <i>k-means</i>	1
4. Diseño	1
4.1. Algoritmo en pseudocódigo	3
5. Implementación	4
5.0.1. Problemas encontrados	4
5.0.2. Soluciones adoptadas	4
6. Validación del <i>software</i>	4
6.1. Diseño del banco de pruebas	4
7. Análisis de resultados	4
7.1. Modificando inicializaciones	4
7.2. Modificando distancia Minkowski	4
7.3. Criterios de convergencia	4
7.3.1. Número fijo de iteraciones	4
7.3.2. Disimilitud entre <i>codebooks</i>	4
7.4. Distintas métricas	4
7.4.1. Manhattan	4
7.4.2. Euclídea	4
7.4.3. Minkowski	4
8. Clasificación supervisada respecto de	4
9. Conclusiones	4
10. Valoración subjetiva	4

ÍNDICE DE TABLAS

ÍNDICE DE FIGURAS

1.	Esquema de dependencias del sistema	2
----	---	---

1. Introducción

El objetivo principal de esta práctica es obtener la capacidad de formular un algoritmo de aprendizaje automático de clasificación **No-Supervisada**. Por otra parte, se trabajarán la capacidad de sintetizar una técnica de aprendizaje automático no-supervisado, conocer su coste computacional así como sus limitaciones de representación y de inteligibilidad

2. Recursos

- PC con aplicación Weka.
- Bibliografía.
- Librerías de Weka.
- Manual de Weka.
- Guía de la práctica.
- Ficheros para los datos de la práctica: [food.arff](#), [colon.arff](#).
- Otros ficheros que no están en formato *.arff*:
 - En formato *.txt*: [ClusterData.atributos.txt](#) (este fichero si tiene la clase asociada para evaluar la calidad del *clustering* en [ClusterData.clase.txt](#)).
 - En formato *.csv* [bank-data.csv](#)clustering

3. Clasificación NO-supervisada o *Clustering*

(Definición) 3.1 Se considera **clasificación no-supervisada** cuando el conjunto de entrenamiento no están las instancias etiquetadas con el valor de la clase. Es un experimento exploratorio, que trata de agrupar las instancias en grupos definidos por similitud entre las características de las instancias que pertenecen al mismo grupo y disimilitud entre las que pertenecen a grupos distintos. Técnicamente estos grupos son llamados *Clusters*.

3.1. Clustering *k-means*

4. Diseño

Estructuramos la ejecución del algoritmo en fases¹, las cuales se definen a continuación.

Primera fase: carga de datos y configuración

Inicialmente se encarga de cargar el fichero en una estructura de datos adecuada para el cálculo del algoritmo. Además se carga la configuración establecida por el usuario, es decir que tipo de inicialización para el *codebook*, número de clusters, distancia a utilizar. . .

Segunda fase: Preproceso de datos

En el preproceso se normaliza el espacio, haciendo uso de la función estadística *zscore*, es decir a cada *Feature* de la instancia se le resta su media y dividiendo el resultado por la varianza de dicha *Feature* se consigue que todos los valores de los atributos se encuentren dentro del intervalo $[-1,1]$. Con esto conseguimos evitar la mayor influencia de algunos atributos debido a que se encuentren dentro de un dominio de valores mayor. Por ejemplo si vamos a agrupar viviendas y utilizamos como características el número de habitaciones y el precio, seguramente el número de habitaciones no será mucho mayor de 3 y en cambio el precio puede ser mayor incluso de 300000 €.[1]

Tercera Fase: Algoritmo K-means

En esta fase se implementa el algoritmo **K-means**.

1. En primer lugar inicia los *centroides* con el criterio establecido por el usuario, o la matriz de bits de pertenencias.
2. Recorre las instancias del conjunto y calcula la distancia a cada uno de los *codeword* actualizando la matriz de bits de pertenencia, el valor del bit es uno si es el centroide más cercano a la instancia.
3. Se calcula de nuevo el vector promedio para cada cluster.
4. Iterar los pasos dos y tres hasta converger.

Cuarta Fase: Evaluación

En esta fase se tratará de automatizar la evaluación del algoritmo frente a los datos obtenidos con distintas ejecuciones, variando los parámetros o incluso con los resultados obtenidos con el modelo de **K-means** que ya implementa el API de **weka**.

Dependencias

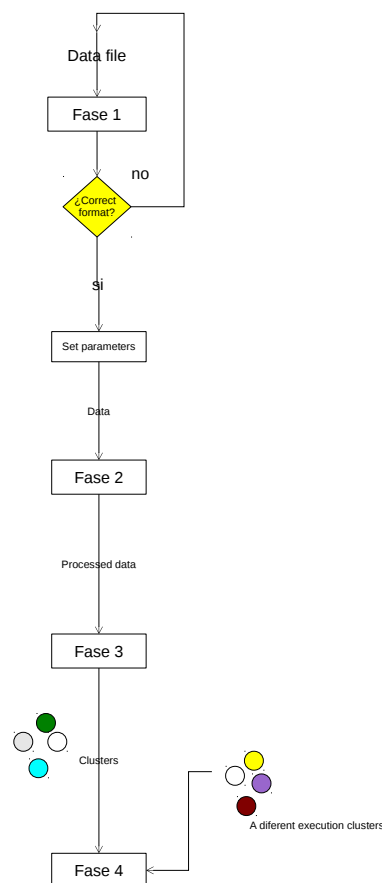


Figura 1: Dependencias del sistema

4.1. Algoritmo en pseudocódigo

```
1  Let k be the number of clusters to partition the data set
2  Let  $X = \{x_1, x_2, \dots, x_n\}$  be the data set to be analyzed
3  Let  $M = \{m_1, m_2, \dots, m_k\}$  be the code-book associated to the clusters
4  Let  $\text{dist}(a, b)$  be the desired distance metric
5  Let  $B = \{B_{11}, B_{12}, \dots, B_{nk}\}$  be the temporary pertenece bit matrix
6
7  Ensure:  $C = \{C_1, C_2, \dots, C_k\}$  set of clusterized instances
8
9  Begin:
10
11  for each  $m_j$ 
12     $m_j = \text{randomsample}(X)$ 
13  end
14
15  for each  $x_n$ 
16     $C_j = \text{argmin dist}(x_i, m_j)$  // para cada j, como la formulita del argmin de siempre. en vez del
    igual tengo que mirar algo que simbolice aadir
17  end
18
19  while changes in M do
20    for each  $m_j$ 
21       $m_{j\text{new}} = \text{calculatecentroid}(C_j)$ 
22      if  $m_{j\text{new}} == m_j$ 
23        changes = false
24      else
25        changes = true
26      end
27       $m_j = m_{j\text{new}}$ 
28    end
29
30    for each  $x_n$ 
31      //  $C_j = \text{argmin dist}(x_i, m_j)$  // lo mismo que arriba, mirar a hacerlo con frmulas de latex despues
      de que lo repases posibilidad 1, sin matriz de pertenencia
32
33       $B_{nj} = 1$  sii  $\text{argmin dist}(x_n, m_j) = m_j$  (para cada  $m_j$ ) si no  $B_{nj} = 0$ 
34    end
35
36    for each  $B_{ij}$ 
37      if  $B_{ij} == 1$ 
38         $C_j = x_i$  (aadir  $x_i$  a  $C_j$ )
39      end
40    end
41  end
42
43  return  $C = \{C_1, C_2, \dots, C_k\}$ 
44 end
```

5. Implementación

5.0.1. Problemas encontrados

5.0.2. Soluciones adoptadas

6. Validación del *software*

6.1. Diseño del banco de pruebas

7. Análisis de resultados

7.1. Modificando inicializaciones

7.2. Modificando distancia Minkowski

7.3. Criterios de convergencia

7.3.1. Número fijo de iteraciones

7.3.2. Disimilitud entre *codebooks*

7.4. Distintas métricas

7.4.1. Manhattan

7.4.2. Euclídea

7.4.3. Minkowski

8. Clasificación supervisada respecto de

9. Conclusiones

- Breve descripción de las motivaciones para llevar a cabo técnicas de clustering.
- Conclusiones a la vista de los resultados más relevantes.
- Conclusiones generales.(Análisis de fortalezas del sw y reflexiones sobre la tarea.
- Análisis de puntos débiles y propuestas de mejoras.

10. Valoración subjetiva

1. ¿Has alcanzado los objetivos que se plantean?
2. ¿Te ha resultado de utilidad la tarea planteada?
3. ¿Qué dificultades has encontrado?Valora el grado de dificultad de la tarea.
4. ¿Cuánto tiempo has trabajado en esta tarea? Desglosado:

Coste temporal	
Diseño de software	1
Implementación de software	
Tiempo trabajando con Weka	
Búsqueda bibliográfica	
Informe	

5. Sugerencias para mejorar la tarea. Sugerencias para que se consiga despertar mayor interés y motivación en los alumnos.
6. Críticas(construktivas).

Referencias

- [1] Andrew Ng. Machine learning, 2014.