

# Clasificación y agrupamiento

Métodos para formar colecciones de objetos.

## Contenidos

---

- ▶ Agrupamiento
  - ▶ Descripción
  - ▶ Ejercicio con K-means
  
- ▶ Clasificación



## Agrupamiento (*clustering*)

- ▶ Consiste en acomodar elementos en grupos.
  - ▶ Elementos en el mismo grupo deben ser más similares entre ellos
    - ▶ ... y menos similares con respecto a otros grupos.
- ▶ Los grupos **no** se conocen de antemano.
- ▶ En los enfoques clásicos, la información se representa con vectores de características
  - ▶ ej. coordenadas en un plano cartesiano, vectores binarios, pesos reales, etc.

▶ 3

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## Datos

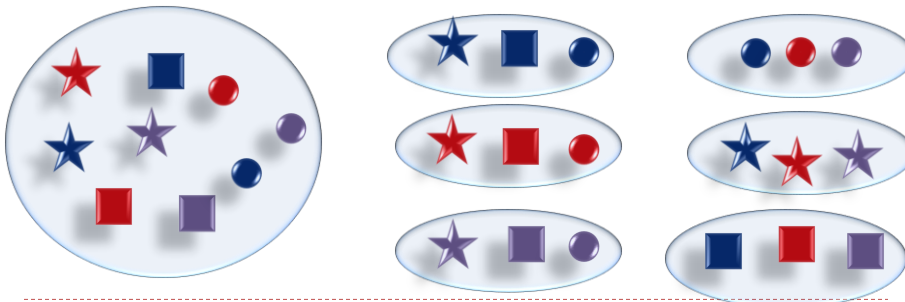
- ▶ Los datos pueden representar lo que Uds. quieran.
  - ▶ Películas
  - ▶ Páginas Web (documentos)
  - ▶ Superhéroes
  - ▶ Proteínas
  - ▶ Canciones
  - ▶ Personas
  - ▶ El límite es la creatividad.

▶ 4

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## Teorema del patito feo

- ▶ Necesitamos un criterio para agrupar
  - ▶ i.e., nunca tenemos grupos libres de sesgos.
- ▶ Existen diferentes métricas y enfoques.



▶ 5

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## Distancia y similitud

- ▶ Similitud = parecido entre un par de elementos
- ▶ Distancia = disimilitud entre un par de elementos
- ▶ Métricas comunes
  - ▶ Índice Jaccard → Conjuntos
  - ▶ Distancia euclidiana → Espacios geométricos
  - ▶ Coseno → Vectores de pesos

▶ 6

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## Tipos de clusters de acuerdo a contención de elementos



▶ 7

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## Algoritmo K-means

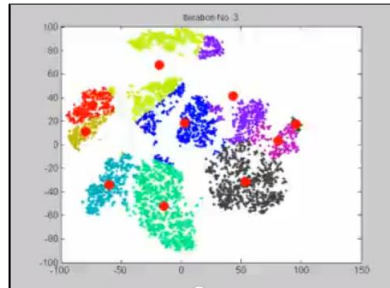
- ▶ Algoritmo clásico de agrupamiento
  - ▶ Particional
- ▶ Inspiración para muchos otros algoritmos
- ▶ Introducido en los 70's
- ▶ Español: K-medias

▶ 8

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## K-means

- ▶ <http://www.youtube.com/watch?v=74rv4snLI70&feature=endscreen&NR=1>

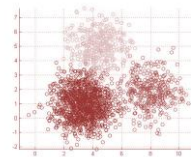


▶ 9

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## K-means: ¿Cómo funciona?

- ▶ Dado un conjunto de datos (vectores)
  - ▶ en N dimensiones
- ▶ Colocar **k** centroides.
  - ▶ al azar
- ▶ Asignar cada vector al centroide más cercano.
- ▶ Con los grupos formados, calcular nuevos centroides.
  - ▶ Que queden en el centro del grupo recién formado...
- ▶ Repetir esto hasta que el *error* sea menor a un umbral.



▶ 10

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## Recalculando centroides

- ▶ Obtener el promedio por cada dimensión del vector.

- ▶ Ejemplo (3D):

<b>Ci = (5, 7, 10)</b>	$Ci_x = \frac{pa_x + pb_x + pc_x + pd_x}{\text{cant. puntos en } Ci}$	$Ci_y = \frac{pa_y + pb_y + pc_y + pd_y}{\text{cant. puntos en } Ci}$
pa = (1, 2, 5)	$= \frac{1+3+8+4}{4}$	$= \frac{2+4+12+17}{4}$
pb = (3, 4, 6)	$= \frac{16}{4}$	$= \frac{35}{4}$
pc = (8, 12, 20)	$= \frac{16}{4}$	$= \frac{35}{4}$
pd = (4, 17, 5)	$= 4$	$= 8.75$

▶ 11

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## Recalculando centroides

<b>Ci = (5, 7, 10)</b>	$Ci_z = \frac{pa_z + pb_z + pc_z + pd_z}{\text{cant. puntos en } Ci}$
pa = (1, 2, 5)	$= \frac{5+6+20+5}{4}$
pb = (3, 4, 6)	$= \frac{36}{4}$
pc = (8, 12, 20)	$= 9$
pd = (4, 17, 5)	

Ci = ~~(5, 7, 10)~~

Ci = (4, 8.75, 9)

▶ 12

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## Ejemplo

Puntos (datos)	Centroides
P1=(1,7)	C1=(3,5)
P2=(2,6)	C2=(12,13)
P3=(4,8)	
P4=(10,19)	
P5=(11,15)	

Cálculo de distancia a cada centroide (para asignar al más cercano)

Punto	Distancia euclídana a C1 (3, 5)	Distancia C2 (12, 13)
(1,7)	2.8	12.5
(2,6)	1.4	12.2
(4,8)	3.2	9.4
(10,19)	15.7	6.3
(11,15)	12.8	2.2

Ejemplo de cálculo distancia euclidiana

$$\text{dist}(P1, C1) = \sqrt{(1-3)^2 + (7-5)^2}$$

$$\text{dist}(P1, C1) = \sqrt{4+4} = \sqrt{8}$$

$$\text{dist}(P1, C1) = 2.8$$

Centroides recalculados

Asignación de puntos a grupos

Grupo 1 (C1)	Grupo 2 (C2)
(1,7)	(10, 19)
(2,6)	(11, 15)
(4,8)	

$$C1(x) = (1+2+4)/3 = 7/3 = 2.3$$

$$C1(y) = (7+6+8)/3 = 21/3 = 7$$

$$C1 = (2.3, 7)$$

$$C2(x) = (10+11)/2 = 21/2 = 10.5$$

$$C2(y) = (19+15)/2 = 34/2 = 17$$

$$C2 = (10.5, 17)$$

► I3

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## Ejemplo

Puntos (datos)
P1=(1,7)
P2=(2,6)
P3=(4,8)
P4=(10,19)
P5=(11,15)

Centroides
C1=(3,5)
C2=(12,13)

Ejemplo de cálculo distancia euclidiana

$$\text{dist}(P1, C1) = \sqrt{(1-3)^2 + (7-5)^2}$$

$$\text{dist}(P1, C1) = \sqrt{4+4} = \sqrt{8}$$

$$\text{dist}(P1, C1) = 2.8$$

Centroides recalculados

$$C1(x) = (1+2+4)/3 = 7/3 = 2.3$$

$$C1(y) = (7+6+8)/3 = 21/3 = 7$$

$$C1 = (2.3, 7)$$

$$C2(x) = (10+11)/2 = 21/2 = 10.5$$

$$C2(y) = (19+15)/2 = 34/2 = 17$$

$$C2 = (10.5, 17)$$

► I4

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

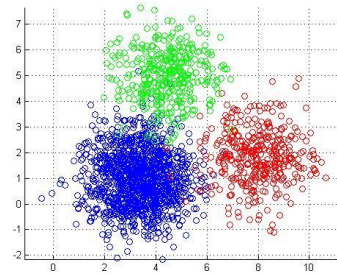
## K-means: ejercicio

### ► Puntos:

- $p_1 = (1, 1)$ ,
- $p_2 = (2, 4)$ ,
- $p_3 = (3, 2)$ ,
- $p_4 = (7, 2)$ ,
- $p_5 = (8, 3)$

### ► Centroides

- $k=2$
- $c_1 = (3, 5)$
- $c_2 = (9, 1)$



► 15

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## K-means: ejercicio

- Utiliza la distancia euclidiana como métrica de *disimilitud*.
- Realiza dos iteraciones del algoritmo.

► 16

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento



## K-means: OJO

---

- ▶ No siempre es fácil colocar los centroides iniciales.
  - ▶ También es difícil saber cuántos poner.
- ▶ Susceptible a intrusos (*outliers*)

---

▶ 17

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## Otros tipos de clustering

---

- ▶ Basado en grafos
- ▶ Basado en densidad
- ▶ Kernels
- ▶ Co-clustering

---

▶ 18

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## ¿Cómo evaluar un agrupamiento?

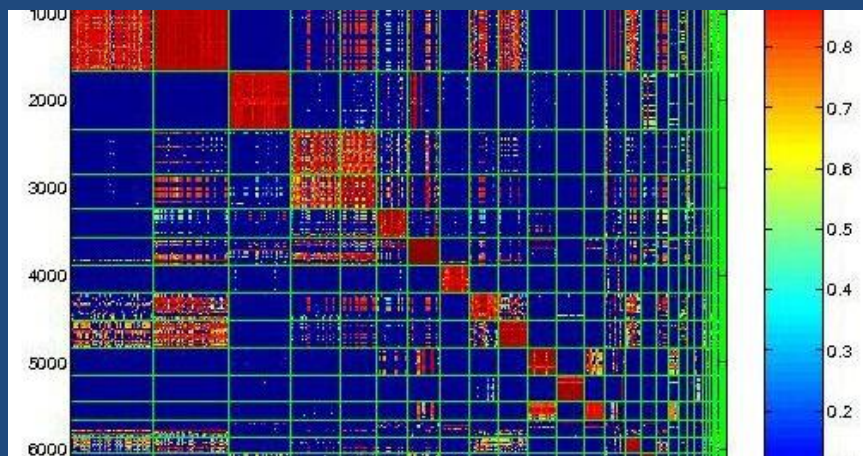
- ▶ Visualmente
  - ▶ Matrices de similitud
- ▶ Precisión y recuerdo (*precision and recall*)
  - ▶ Correctitud y completez
  - ▶ Medidas integradoras: F

▶ 19

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

### Matriz de similitud

Fuente: <http://blogoutthemaps.blogspot.com/2011/01/similarity-matrix.html>



▶ 20

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## Clasificación

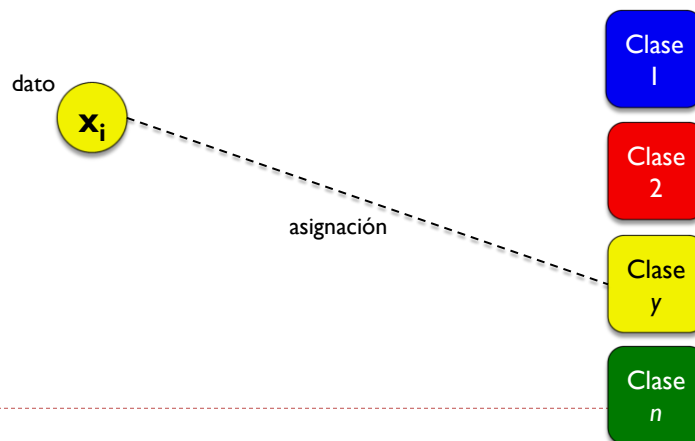
- ▶ Igual que agrupamiento, pero los grupos (llamados clases o categorías) ya se encuentran pre-establecidos.

▶ 21

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

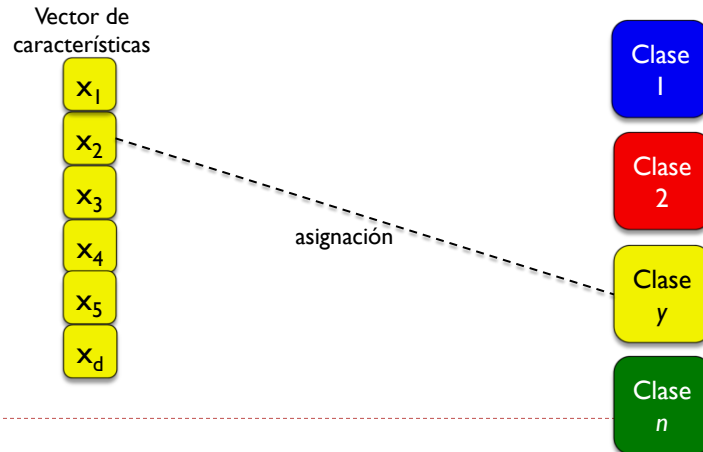
## Clasificación

- ▶ También conocida como categorización.

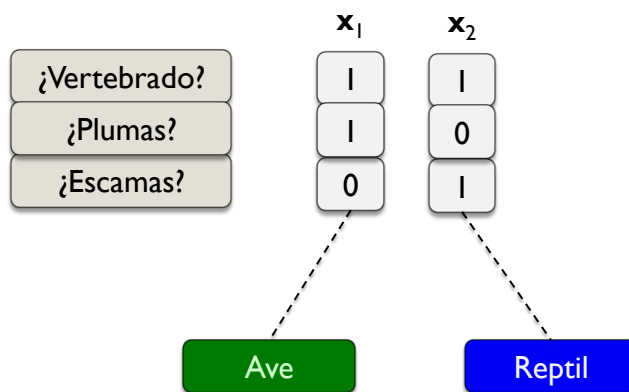


## Clasificación

- También conocida como categorización.



## Ejemplo con vectores binarios



## Clasificación: Análisis de sentimiento

- ▶ Una aplicación reciente es el análisis de sentimiento para documentos (opiniones).
  - ▶ Tres clases: Positivo, negativo, neutro.
  - ▶ Ejemplo **positivo**: “La película estuvo muy divertida.”
  - ▶ Ejemplo **negativo**: “El nuevo dispositivo es poco funcional.”
  - ▶ Ejemplo **neutro**: “El motor viene con dos sensores.”

▶ 25

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## Resumen

- ▶ Agrupamiento = colocar datos en grupos de acuerdo a su parecido.
- ▶ Tipos de grupos
  - ▶ Sin traslape
  - ▶ Con traslape
- ▶ K-medias (k-means)
  - ▶ Asignar el dato al centroide más cercano. Al final, centroides.



▶ 26

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## Resumen

- ▶ Agrupamiento = colocar datos en grupos de acuerdo a su parecido.
- ▶ Tipos de grupos
  - ▶ Sin traslape
  - ▶ Con traslape
- ▶ K-medias (k-means)
  - ▶ Asignar el dato al centroide más cercano. Al final, recalculan los centroides.



▶ 27

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento

## Referencias

- ▶ Tan et al. Introduction to Data Mining. Addison Wesley, EUA, 2006.
- ▶ A Tutorial on Clustering Algorithms.  
[http://home.dei.polimi.it/matteucc/Clustering/tutorial\\_html/index.html](http://home.dei.polimi.it/matteucc/Clustering/tutorial_html/index.html)

▶ 28

Programación de Sistemas Adaptativos:  
Clasificación y agrupamiento