

A stylized sunburst graphic in shades of purple and blue, located in the top-left corner of the slide.

# Minería de patrones frecuentes y reglas de asociación

Máster Online en Ciencia de Datos

## Dr. José Raúl Romero

Profesor Titular de la Universidad de Córdoba y Doctor en Ingeniería Informática por la Universidad de Málaga. Sus líneas actuales de trabajo se centran en la democratización de la ciencia de datos (*Automated ML* y *Explainable Artificial Intelligence*), aprendizaje automático evolutivo y analítica de software (aplicación de aprendizaje y optimización a la mejora del proceso de desarrollo de software).

Miembro del Consejo de Administración de la *European Association for Data Science*, e investigador senior del Instituto de Investigación Andaluz de *Data Science and Computational Intelligence*.

Director del **Máster Online en Ciencia de Datos** de la Universidad de Córdoba.



A stylized sunburst graphic in shades of purple and blue, located in the top-left corner of the slide. It features a semi-circle on the left with several rays extending outwards to the right.

# Generación de reglas de asociación

# Recordemos...

- Habitualmente, el proceso de **minería de reglas de asociación** se divide en **dos subproblemas** (Han and Kamber, 2001):
  - **Encontrar los itemsets frecuentes** cuyas ocurrencias exceden un umbral de soporte mínimo predefinido
  - **Derivar reglas de asociación** a partir de aquellos itemsets frecuentes (considerando un umbral para la confianza, si procede)
- Estos subproblemas se resuelven iterativamente hasta que no aparezca más nuevas reglas (**condición de cierre**)

# Reglas de asociación

Una implicación de la forma  $\mathbf{X} \rightarrow \mathbf{Y}$ ,  
donde  $\mathbf{X}$  e  $\mathbf{Y}$  son *itemsets*.

Ejemplo:

$\{\text{Leche, Pañales}\} \rightarrow \{\text{Cerveza}\}$

## Medidas de evaluación de una regla:

### Soporte (s)

Fracción de transacciones que  
contienen los *itemsets*  $\mathbf{X}$  e  $\mathbf{Y}$

### Confianza (c)

Mide con qué frecuencia los ítems  
de  $\mathbf{Y}$  aparecen en las transacciones  
que contienen  $\mathbf{X}$

TID	Items
1	Pan, Leche
2	Pan, Pañales, Cerveza, Huevos
3	Leche, Pañales, Cerveza, Coca-cola
4	Pan, Leche, Pañales, Cerveza
5	Pan, Leche, Pañales, Coca-cola

Ejemplo:  $\{\text{Leche, Pañales}\} \Rightarrow \text{Cerveza}$

$$s = \frac{\sigma(\text{Leche, Pañales, Cerveza})}{|T|} = \frac{2}{5} = 0.4$$

$$c = \frac{\sigma(\text{Leche, Pañales, Cerveza})}{\sigma(\text{Leche, Pañales})} = \frac{2}{3} = 0.67$$

# Generación de reglas de asociación

Dado un *itemset* frecuente  $L$ , encontrar todos los subconjuntos no vacíos  $f \subset L$  tal que  $f \rightarrow L - f$  satisfaciendo el **umbral de confianza**

- Si  $\{A,B,C,D\}$  es un *itemset* frecuente, las reglas candidatas son:

$ABC \rightarrow D, \quad ABD \rightarrow C, \quad ACD \rightarrow B, \quad BCD \rightarrow A,$   
 $A \rightarrow BCD, \quad B \rightarrow ACD, \quad C \rightarrow ABD, \quad D \rightarrow ABC,$   
 $AB \rightarrow CD, \quad AC \rightarrow BD, \quad AD \rightarrow BC, \quad BC \rightarrow AD,$   
 $BD \rightarrow AC, \quad CD \rightarrow AB,$

Si  $|L| = k$ , hay  $2k - 2$  reglas de asociación candidatas (ignorando  $L \rightarrow \emptyset$  y  $\emptyset \rightarrow L$ )

## ¿Cómo generar eficientemente reglas a partir de los *itemsets* frecuentes?

- En general, **la confianza no tiene la propiedad anti-monótona**

$c(ABC \rightarrow D)$  puede ser mayor o menor que  $c(AB \rightarrow D)$

- Pero la confianza de las **reglas generadas a partir del mismo *itemset* sí tienen la propiedad anti-monótona**

e.g.,  $L = \{A, B, C, D\}$ :

$$c(ABC \rightarrow D) \geq c(AB \rightarrow CD) \geq c(A \rightarrow BCD)$$

- La confianza es anti-monótona con respecto al número de ítems del consecuente de la regla



# Generación de reglas de asociación

Algoritmo de generación de reglas



# Generación de reglas

Función de generación de reglas de asociación

**AssocRules()**

forall large datasets  $l_k, k > 2$

$GenRules(l_k, l_k);$

Generar todas las reglas válidas  $a \rightarrow (l_k - a)$ , para todas  $a \subset a_m$

**GenRules( $l_k, a_m$ )**

$A = \{(m-1) - \text{itemsets } a_{m-1} \mid a_{m-1} \subset a_m\}$

**forall**  $a_{m-1} \in A$  **begin**

$conf = \frac{\text{support}(l_k)}{\text{support}(a_{m-1})}$

**if** ( $conf \geq \text{minconf}$ ) **then**

output – rule  $a_{m-1} \rightarrow (l_k - a_{m-1})$

**if** ( $m - 1 > 1$ ) **then**

**GenRules**( $l_k, a_{m-1}$ )

Con confianza  $conf$   
Con support =  $\text{support}(l_k)$

Generar reglas con subconjuntos de  $a_{m-1}$  como antecedente

# Ejemplo (recuperamos: lección Apriori)

TID	Lista de Items
T100	I1, I2, I5
T101	I2, I4
T102	I2, I3
T103	I1, I2, I4
T104	I1, I3
T105	I2, I3
T106	I1, I3
T107	I1, I2, I3, I5
T108	I1, I2, I3

- La base de datos **D** contiene 9 transacciones
- Se supone que el soporte mínimo requerido es 2 ( $\text{min\_sup} = 2/9 = 22\%$ )
- Se establece un **umbral mínimo de confianza del 70%**
- Primero, debemos averiguar los itemsets frecuentes utilizando el algoritmo correspondiente
- Después, se generan las reglas de asociación según los umbrales mínimos de soporte y confianza

# Ejemplo

- **Procedimiento:**

- Para cada itemset frecuente  $I$ , se generan todos los subconjuntos de  $I$  no vacíos
- Para cada subconjunto no vacío  $a$  de  $I$ , se considera la regla  $a \rightarrow (I - a)$  si su confianza supera el umbral mínimo:

$$\text{support\_count}(I) / \text{support\_count}(s) \geq \text{min\_conf}$$

- Siguiendo con el ejemplo...

$L = \{\{I1\}, \{I2\}, \{I3\}, \{I4\}, \{I5\}, \{I1,I2\}, \{I1,I3\}, \{I1,I5\}, \{I2,I3\}, \{I2,I4\}, \{I2,I5\}, \{I1,I2,I3\}, \{I1,I2,I5\}\}.$

- Tomamos  $I = \{I1,I2,I5\}$  y **min\_conf** = 70%
- Todos sus subconjuntos no vacíos son  $\{I1,I2\}, \{I1,I5\}, \{I2,I5\}, \{I1\}, \{I2\}, \{I5\}$

# Ejemplo

A cada regla resultante se le calcula su confianza:

- R1:  $I1 \wedge I2 \rightarrow I5$ 
  - Confianza =  $\text{supp}\{I1, I2, I5\} / \text{supp}\{I1, I2\} = 2 / 4 = 50\%$
  - R1 es rechazada
- R2:  $I1 \wedge I5 \rightarrow I2$ 
  - Confianza =  $\text{supp}\{I1, I2, I5\} / \text{supp}\{I1, I5\} = 2 / 2 = 100\%$
  - R2 es seleccionada
- R3:  $I2 \wedge I5 \rightarrow I1$ 
  - Confianza =  $\text{supp}\{I1, I2, I5\} / \text{supp}\{I2, I5\} = 2/2 = 100\%$
  - R3 es seleccionada

# Ejemplo

- R4:  $I1 \rightarrow I2 \wedge I5$ 
  - Confianza =  $\text{supp}\{I1, I2, I5\} / \text{supp}\{I1\} = 2 / 6 = 33\%$ 
    - R4 es rechazada
- R5:  $I2 \rightarrow I1 \wedge I5$ 
  - Confianza =  $\text{supp}\{I1, I2, I5\} / \{I2\} = 2 / 7 = 29\%$ 
    - R5 es rechazada
- R6:  $I5 \rightarrow I1 \wedge I2$ 
  - Confidence =  $\text{sc}\{I1, I2, I5\} / \{I5\} = 2/2 = 100\%$ 
    - R6 es seleccionada

R2, R3 y R6 son reglas de asociación exactas

## Ejemplo 2

$L = \{\{1\}, \{2\}, \{3\}, \{5\}, \{1,3\}, \{2,3\}, \{2,5\}, \{3,5\}, \{2,3,5\}\}$

Fila	1	2	3	4	5
1	x		x	x	
2		x	x		x
3	x	x	x		x
4		x			x

$\text{min\_supp} = 2$

$\text{min\_conf} = 0.75$

### Reglas de asociación:

$\{1\} \rightarrow \{3\} : 1$

$\{2\} \rightarrow \{3\} : 0.67$

$\{2\} \rightarrow \{5\} : 1$

$\{3\} \rightarrow \{5\} : 0.67$

$\{2,3\} \rightarrow \{5\} : 1$

$\{3,5\} \rightarrow \{2\} : 1$

$\{3\} \rightarrow \{1\} : 0.67$

$\{3\} \rightarrow \{2\} : 0.67$

$\{5\} \rightarrow \{2\} : 1$

$\{5\} \rightarrow \{3\} : 0.67$

$\{2,5\} \rightarrow \{3\} : 0.67$

A stylized sunburst graphic in shades of purple and blue, located in the top-left corner of the slide.

# Generación de reglas de asociación

Uso de la herramienta Weka

A stylized sunburst graphic in shades of purple and blue, located in the top-left corner of the slide. It features a semi-circle on the left with several rays extending outwards to the right.

¡Gracias!

**UCO**  
ONLINE

A horizontal bar at the bottom of the slide consisting of three equal-width rectangles of yellow, red, and blue, representing the colors of the Spanish flag.