Minería de Datos IIC2433

Reglas de asociación Vicente Domínguez

¿Qué veremos esta clase?

- Reglas de asociación

Reglas de asociación

Base de datos de transacciones

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

{huevos, mantequilla}

{huevos, mantequilla}

quizás para hacer tostadas...

{huevos, mantequilla}



quizás para hacer tostadas...

{huevos, harina}

{huevos, harina}

quizás para repostería...

{huevos, harina} azúcar

quizás para repostería...

TID Items 1 Pan, Coca Cola, Leche 2 Cerveza, Pan 3 Pan, Coca Cola, Pañales, Leche 4 Cerveza, Pan, Pañales, Leche 5 Coca Cola, Pañales, Leche

1. Listar **itemsets** posibles

- Pan}
- {Coca Cola}
- {Leche}
- {Pañales}
- {Cerveza}
- {Pan, Coca Cola}
- {Pan, Leche}
- {Pan, Pañales}
- {Pan, Cerveza}
- {Coca Cola, Leche}
- {Coca Cola, Pañales}
- {Coca Cola, Cerveza}
- {Leche, Pañales}
- {Leche, Cerveza}
- {Pañales, Cerveza}
- {Pan, Coca Cola, Leche}
- {Pan, Coca Cola, Pañales}
- {Pan, Coca Cola, Cerveza}
- {Pan, Leche, Pañales}
- {Pan, Leche, Cerveza}
- {Pan, Pañales, Cerveza}

- {Coca Cola, Leche, Pañales}
- {Coca Cola, Leche, Cerveza}
- {Coca Cola, Pañales, Cerveza}
- {Leche, Pañales, Cerveza}
- {Pan, Coca Cola, Leche, Pañales}
- {Pan, Coca Cola, Leche, Cerveza}
- {Pan, Coca Cola, Pañales, Cerveza}
- {Pan, Leche, Pañales, Cerveza}
- {Coca Cola, Leche, Pañales, Cerveza}
- {Pan, Coca Cola, Leche, Pañales, Cerveza}

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

2. Contar la frecuencia de cada itemset

- $\sigma(\{Pan\}) = 4$
- σ({Coca Cola, Pañales}) = 2
- σ({Pan, Leche, Cerveza}) = 1

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

3. Calcular el soporte

- s({Pan}) = \%
- s({Coca Cola, Pañales}) = %
- s({Pan, Leche, Cerveza}) = 1/5

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

4. Todos los itemsets cuyo soporte sea mayor a un **umbral** se consideran **frecuentes**

Por ejemplo, umbral = 0,6

- $s({Pan}) = \% = 0.8$
- s({Coca Cola, Pañales}) = % = 0,4
- s({Pan, Leche, Cerveza}) = \% = 0,2

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

4. Todos los itemsets cuyo soporte sea mayor a un **umbral** se consideran **frecuentes**

Por ejemplo, umbral = 0,6

- $s({Pan}) = \% = 0.8$ **frecuente**
- s({Coca Cola, Pañales}) = % = 0,4 no frecuente
- s({Pan, Leche, Cerveza}) = ½ = 0,2 **no frecuente**

Formalizando los conceptos:

Itemset

- Una colección de uno o más ítems
- Ejemplo: {Leche, pan, cerveza}

Contador del soporte (σ)

- Frecuencia de ocurrencia de un itemset
- \circ Ej. σ ({Leche, pan, cerveza}) = 1

Soporte

- Fracción de las transacciones que contiene un itemset
- Ej. s({Leche, pan, cerveza}) = ½

Itemset frecuente

 Un itemset cuyo soporte es mayor o igual a un determinado umbral. • s({Pañales, Leche}) = 3/5 = 0,6

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

¿Qué podemos decir sobre las siguientes reglas de asociación?

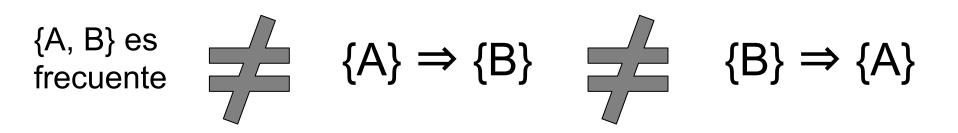
{Pañales} => {Leche} {Leche} => {Pañales}

• s({Pañales, Leche}) = 3/5 = 0,6

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

¿Qué podemos decir sobre las siguientes **reglas de asociación**?

Itemset frecuente y regla de asociación son conceptos diferentes



Formalizando más conceptos:

Regla de asociación

- Una expresión de la forma $X \Rightarrow Y$, donde X e Y son itemsets
- Ejemplo: {Leche, pañales} ⇒ {cerveza}

Soporte de una regla de asociación (s)

- Soporte de la unión del antecedente y el consecuente
- $\circ \quad \mathsf{s}(\mathsf{X} \Rightarrow \mathsf{Y}) = \mathsf{s}(\mathsf{X} \; \mathsf{U} \; \mathsf{Y})$
- \circ Ejemplo: s({Leche, pañales} \Rightarrow {cerveza}) = s({Leche, pañales, cerveza})

Confianza de una regla de asociación:

$$c(X \Rightarrow Y) = \frac{s(X \cup Y)}{s(X)}$$

Ejemplo: c({Leche, pañales} ⇒ {cerveza}) =

$$\frac{s(\{leche, paales, cerveza\})}{s(\{leche, paales\})} = \frac{\sigma(\{leche, paales, cerveza\})}{\sigma(\{leche, paales\})}$$

Tenemos un pequeño problema...

- Con 5 items, obtuvimos 2⁵ 1 posibles ítems
 - Es decir, 31

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

- {Pan}
- {Coca Cola}
- {Leche}
- {Pañales}
- {Cerveza}
- {Pan, Coca Cola}
- {Pan, Leche}
- {Pan, Pañales}
- {Pan, Cerveza}
- {Coca Cola, Leche}
- {Coca Cola, Pañales}
- {Coca Cola, Cerveza}
- {Leche, Pañales}
- {Leche, Cerveza}
- {Pañales, Cerveza}
- {Pan, Coca Cola, Leche}
- {Pan, Coca Cola, Pañales}
- {Pan, Coca Cola, Cerveza}
- {Pan, Leche, Pañales}
- {Pan, Leche, Cerveza}
- {Pan, Pañales, Cerveza}

- {Coca Cola, Leche, Pañales}
- {Coca Cola, Leche, Cerveza}
- {Coca Cola, Pañales, Cerveza}
- {Leche, Pañales, Cerveza}
- {Pan, Coca Cola, Leche, Pañales}
 - {Pan, Coca Cola, Leche, Cerveza}
- {Pan, Coca Cola, Pañales, Cerveza}
- {Pan, Leche, Pañales, Cerveza}
- {Coca Cola, Leche, Pañales, Cerveza}
- {Pan, Coca Cola, Leche, Pañales, Cerveza}

Tenemos un pequeño problema...

- Con 5 items, obtuvimos 2⁵ 1 posibles
 ítems
 - Es decir, 31
- Con n items, tenemos 2ⁿ 1 posibles ítems
- Una tienda suele tener varios ítems
 - Imaginemos n = 100

Tenemos un pequeño problema...

- Si n = 100
- 2¹⁰⁰ 1 posibles ítems
- 2^{100} 1 =

1267650600228229401496703205376

Un computador actual puede hacer ~ 3 millones de operaciones por segundo (3 GHz)...

Entonces demoraríamos aproximadamente

64403322675823264816693248 segundos en sólo

encontrar los itemsets posibles

es decir 2040860051243962497 años

Solución: Algoritmo Apriori

Principio de Monotonicidad:

Si un itemset es frecuente, entonces todos los subgrupos de éste también son frecuentes

- Si {pan, cerveza} es frecuente, entonces {pan} y {cerveza} deben ser frecuentes.
- Si {pan, cerveza, pañales}, entonces...

Regla inversa (anti-monotonía):

Si un itemset no es frecuente, entonces todos sus supersets deben también ser infrecuentes

- Si {pan} no es frecuente, entonces ningún conjunto que contenga panes será frecuente
- Si {pan, coca cola}, entonces...

Buscando Itemsets frecuentes

1. Definimos el valor del umbral Umbral = 0,6

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

Buscando Itemsets frecuentes

TID Items 1 Pan, Coca Cola, Leche 2 Cerveza, Pan 3 Pan, Coca Cola, Pañales, Leche 4 Cerveza, Pan, Pañales, Leche 5 Coca Cola, Pañales, Leche

2. Calculamos el soporte de los 1-itemsets (itemsets de tamaño 1)

- $\sigma(\{Pan\}) =$
- σ({Coca Cola}) =
- $\sigma(\{\text{Leche}\}) =$
- $\sigma(\{Cerveza\}) =$
- $\sigma(\{Pa\tilde{n}ales\}) =$

Buscando Itemsets frecuentes

TID Items 1 Pan, Coca Cola, Leche 2 Cerveza, Pan 3 Pan, Coca Cola, Pañales, Leche 4 Cerveza, Pan, Pañales, Leche 5 Coca Cola, Pañales, Leche

2. Calculamos el soporte de los 1-itemsets (itemsets de tamaño 1)

•
$$\sigma(\{Pan\}) = 4$$

•
$$\sigma(\{\text{Coca Cola}\}) = 3$$

•
$$\sigma(\{\text{Leche}\}) = 4$$

•
$$\sigma(\{Cerveza\}) = 2$$

Buscando Itemsets frecuentes

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

3. Seleccionamos los datasets frecuentes

- s({Pan}) =
- s({Coca Cola}) =
- s({Leche}) =
- s({Cerveza}) =
- s({Pañales}) =

C₁ = {{Cerveza}, {Coca Cola}, {Leche}, {Pan}, {Pañales}}

Ordenar alfabéticamente

L₁ = {{Coca Cola}, {Leche}, {Pan},{Pañales}}

Buscando Itemsets frecuentes

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

3. Seleccionamos los datasets frecuentes

- $s(\{Pan\}) = 0.8$
- s({Coca Cola}) = 0.6
- $s(\{Leche\}) = 0.8$
- $s(\{Cerveza\}) = 0.4$
- s({Pañales}) = 0.6

C₁ = {{Cerveza}, {Coca Cola}, {Leche}, {Pan}, {Pañales}}

Ordenar alfabéticamente

L₁ = {{Coca Cola}, {Leche}, {Pan},{Pañales}}

Buscando Itemsets frecuentes

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

operación join.

L₁ = { {Coca Cola}, {Leche}, {Pan}, {Pañales}}

$$C_2 = L_1 \bowtie L_1$$

Operación Join (⋈)

Buscando Itemsets frecuentes

Conj 1[∞] Conj 2 = { {A, B, C, D}, {E, F, G, H}}

Buscando Itemsets frecuentes

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

K = 2
Construimos 2-itemsets (itemsets de tamaño 2) candidatos) a partir de los 1-itemsets frecuentes, utilizando la operación join.

$$L_1 = L_1 = L_1$$

$$C_2 = L_1 \bowtie L_1$$

Buscando Itemsets frecuentes

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

K = 2
Construimos 2-itemsets (itemsets de tamaño 2) candidatos) a partir de los 1-itemsets frecuentes, utilizando la operación join.

$$C_2 = L_1 \bowtie L_1$$

Buscando Itemsets frecuentes

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

```
K = 2
```

- s({Coca Cola, Leche) =
- s({Coca Cola,Pan}) =
- s({Coca Cola, Pañales}) =
- s({Leche, Pan) =
- s({Leche, Pañales}) =
- s({Pan, Pañales}) =

Buscando Itemsets frecuentes

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

```
K = 2
```

- s({Coca Cola, Leche) =
- s({Coca Cola,Pan}) =
- s({Coca Cola, Pañales}) =
- s({Leche, Pan) =
- s({Leche, Pañales}) =
- s({Pan, Pañales}) =

Buscando Itemsets frecuentes

K = 3

TID	Items
1	Pan, Coca Cola, Leche
2	Cerveza, Pan
3	Pan, Coca Cola, Pañales, Leche
4	Cerveza, Pan, Pañales, Leche
5	Coca Cola, Pañales, Leche

En la pizarra

Buscando reglas de asociación

Para cada itemset frecuente sacar sus subconjuntos y calcular la **confianza** entre ellas.

Recordando:

Regla de asociación

- \circ Una expresión de la forma X \Rightarrow Y, donde X e Y son itemsets
- Ejemplo: {Leche, pañales} ⇒ {cerveza}

Soporte de una regla de asociación (s)

- Soporte de la unión del antecedente y el consecuente
- $\circ \quad \mathsf{s}(\mathsf{X} \Rightarrow \mathsf{Y}) = \mathsf{s}(\mathsf{X} \; \mathsf{U} \; \mathsf{Y})$
- \circ Ejemplo: s({Leche, pañales} \Rightarrow {cerveza}) = s({Leche, pañales, cerveza})

Confianza de una regla de asociación:

$$c(X \Rightarrow Y) = \frac{s(X \cup Y)}{s(X)}$$

Ejemplo: c({Leche, pañales} ⇒ {cerveza}) =

$$\frac{s(\{leche, paales, cerveza\})}{s(\{leche, paales\})} = \frac{\sigma(\{leche, paales, cerveza\})}{\sigma(\{leche, paales\})}$$

Formalizando más conceptos:

Lift

Mide qué tan correlacionados están X e Y

$$lift(X \Rightarrow Y) = \frac{c(X \Rightarrow Y)}{s(Y)} = \frac{s(X \cup Y)}{s(X)s(Y)}$$

$$Lift = \frac{c(X \to Y)}{s(Y)} = \begin{cases} < 1 \\ = 1 \\ > 1 \end{cases}$$

Negativamente correlacionadas Independientes Positivamente correlacionadas

Ejercicio

TID	Items
T1	11, 12, 15
T2	12, 14
Т3	12, 13
T4	11, 12, 14
T5	11, 13
T6	12, 13
T7	11, 13
T8	11, 12, 13, 15
Т9	11, 12, 13