

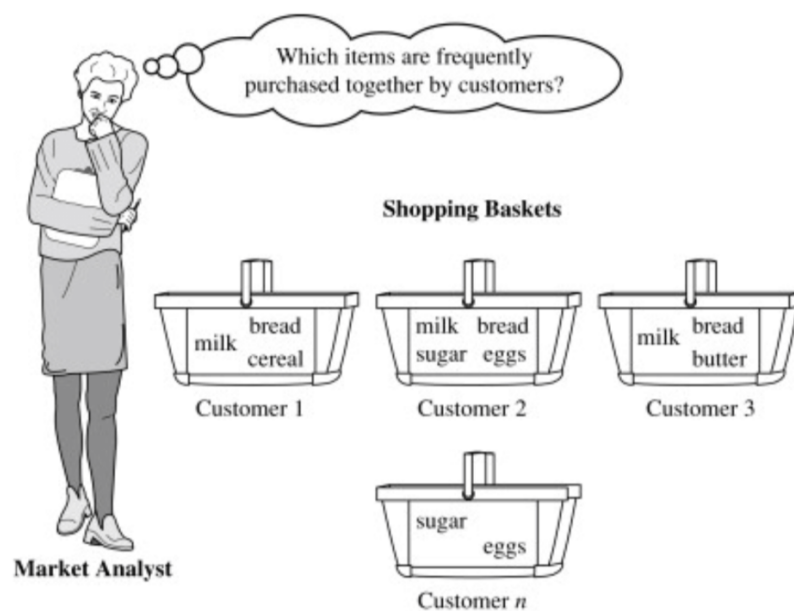
Probabilidad y Análisis de Datos

Daniel Fraiman

Maestría en Ciencia de Datos, Universidad de San Andrés

1 / 36

Basket market



2 / 36

Aplicaciones en espacios grandes discretos

Supongamos que estamos interesados en entender:

“qué y cómo compra una persona cuando va al supermercado”

- En el supermercado hay N items
- Una compra es: $\mathbf{X} = (X_1, X_2, \dots, X_N)$, con $X_i = \{0, 1\}$ (no compró, compró).

3 / 36

Aplicaciones en espacios discretos grandes discretos

Supongamos que estamos interesados en entender:

“qué y cómo compra una persona cuando va al supermercado”

Concretamente

Objetivo: Identificar productos que tiendan a comprarse de forma conjunta.

Con el fin de:

- situarlos en posiciones cercanas dentro de la tienda y maximizar la probabilidad de que los clientes compren.
- presentar nuevos combos de productos al consumidor de manera de aumentar las ventas.
- Si es una compra web, sugerir otros productos.

4 / 36

Aplicaciones en espacios discretos grandes discretos

Customers Who Bought This Item Also Bought



The screenshot shows five recommended books with their covers, titles, authors, ratings, and prices. A left arrow icon is visible on the far left.

Book Title	Author	Rating (Stars)	Count	Format	Price
Pattern Recognition and Machine Learning (Information Science and...)	Christopher Bishop	★★★★☆	115	Hardcover	\$60.76 Prime
Learning From Data	Yaser S. Abu-Mostafa	★★★★☆	88	Hardcover	
The Elements of Statistical Learning: Data Mining, Inference, and Prediction...	Trevor Hastie	★★★★☆	50	Hardcover	\$62.82 Prime
Probabilistic Graphical Models: Principles and Techniques (Adaptive...	Daphne Koller	★★★★☆	28	Hardcover	\$91.66 Prime
Foundations of Machine Learning (Adaptive Computation and...	Mehryar Mohri	★★★★☆	8	Hardcover	\$65.68 Prime

5/36

Reglas de Asociación e items frecuentes

El problema matemático es el siguiente:

Supongamos que sabemos que la persona tiene en su carrito de compras los productos A y B , ¿hay altas chances de que compre T si lo encuentra antes de llegar a la caja?

Llamemos S al evento que tiene las compras actuales.

$$S = A \cap B, \text{ ¿} S \rightarrow T?$$

¿Le recomendamos T ?

6/36

Reglas de Asociación e items frecuentes

Problema matemático:

Dado S , ¿cuál es el producto T que maximiza la probabilidad de compra? Llamemos $\Theta^[-S]$ al conjunto que tiene a todos los productos del supermercado menos los que ya fueron comprados S .

$$T = \underset{Y \in \Theta^{-S}}{\operatorname{argmax}} \mathbb{P}(Y|S)$$

- En realidad estaremos interesados en ordenar los productos según su probabilidad. Quizás ofrecemos los 5 primeros.

7/36

Reglas de Asociación e items frecuentes

Objetivo:

Estimar $\mathbb{P}(Y|S)$ para todo $Y \in \Theta$.

Recordemos:

$$\mathbb{P}(Y|S) = \frac{\mathbb{P}(Y \cap S)}{\mathbb{P}(S)} \text{ con } \mathbb{P}(S) > 0.$$

Estimación:

Basado en la historia de compras (tickets):

- $\mathbb{P}(S) \approx \frac{\#\{\text{tickets con } S\}}{\#\{\text{tickets}\}}$
- $\mathbb{P}(Y \cap S) \approx \frac{\#\{\text{tickets con } Y \text{ y } S\}}{\#\{\text{tickets}\}}$

8/36

Reglas de Asociación e items frecuentes

Dificultad en la estimación:

El conjunto Θ es muy grande (N), por lo tanto el Espacio de Probabilidad de compras será gigante (2^N), no tendremos una muestra (historia de compras) suficientemente grande y entonces las estimaciones tendrán mucho error (o varianza).

- Pero independiente de lo anterior, no nos interesan realmente los productos Y con $\mathbb{P}(Y|S) \ll 1$.

Nuevo planteo

Vamos a pedir:

- $\mathbb{P}(S \cap T) \geq s$, con s algún valor prefijado.
- $\mathbb{P}(T|S) \geq c$, con c algún valor prefijado.

9/36

Reglas de Asociación e items frecuentes

Nuevo planteo

Vamos a pedir:

- $\mathbb{P}(S \cap T) \geq s$, con s algún valor prefijado.
 - $\mathbb{P}(T|S) \geq c$, con c algún valor prefijado.
-
- $\text{Soporte}(S \rightarrow T) = \text{Soporte}(T \rightarrow S) := \mathbb{P}(S \cap T) \geq s$, con s algún valor prefijado.
 - $\text{Confianza}(S \rightarrow T) := \mathbb{P}(T|S) \geq c$, con c algún valor prefijado.

10/36

Reglas de Asociación e items frecuentes

- $\text{Soporte}(S \rightarrow T) := \mathbb{P}(S \cap T) \geq s$, con s algún valor prefijado.
- $\text{Confianza}(S \rightarrow T) := \mathbb{P}(T|S) \geq c$, con c algún valor prefijado.

- Soporte = “cuántas ventas $S \cap T$ espero tener”
- Confianza = qué confianza (chances) tengo en que compren el producto T recomendado cuando tienen S .

11 / 36

Reglas de Asociación e items frecuentes

En la práctica:

Ponemos un límite al tamaño de S y T . Por ejemplo $|S| < 3$ y $|T| = 1$, Si compró salchichas y pan ($|S| = 2$), ¿qué le ofrezco? ¿Y si compró ojotas y protector solar?

- Fijamos s ($\text{Soporte} := \mathbb{P}(S \cap T) \geq s$)
- Fijamos c ($\text{Confianza} := \mathbb{P}(T|S) \geq c$)
- Y damos todas las recomendaciones (T) para todos los conjuntos S compatibles con tener un soporte y una confianza mayor a s y c .

12 / 36

ALGORITMO A PRIORI: $\mathbb{P}(S \cap T) \geq s$

13/36

Algoritmo a Priori: $\mathbb{P}(S \cap T) =: \mathbb{P}(Z) \geq s$

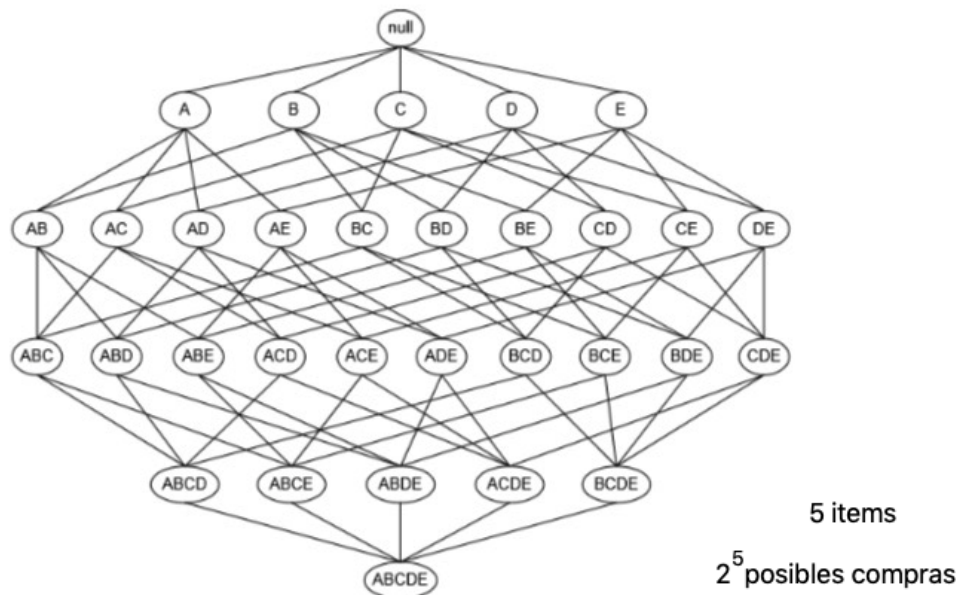
Objetivo: Encontrar todos los conjuntos Z (con $|Z| < k$) que cumplen $\mathbb{P}(Z) \geq s$.

Algoritmo A Priori

1. Generar una lista con todos los conjuntos Z de tamaño 1.
2. $k = 1$.
3. Podar (prune) de la lista los candidatos Z de tamaño k que $\mathbb{P}(Z) < s$.
4. Generar todos los conjuntos Z' de tamaño $k + 1$ que tienen como subconjunto a los elementos de la lista.
5. $k = k + 1$ y go to 3.

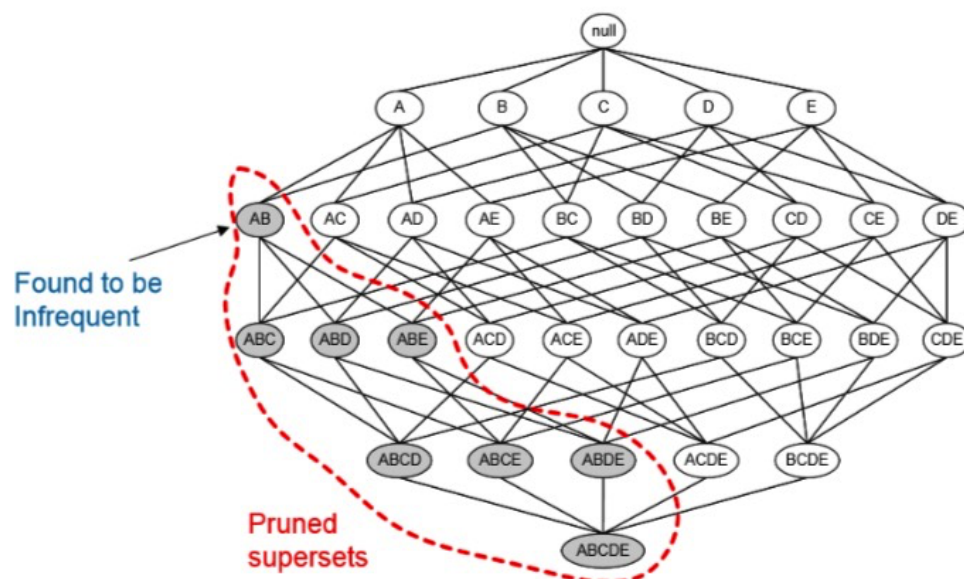
14/36

Algoritmo a Priori: $\mathbb{P}(Z) \geq s$



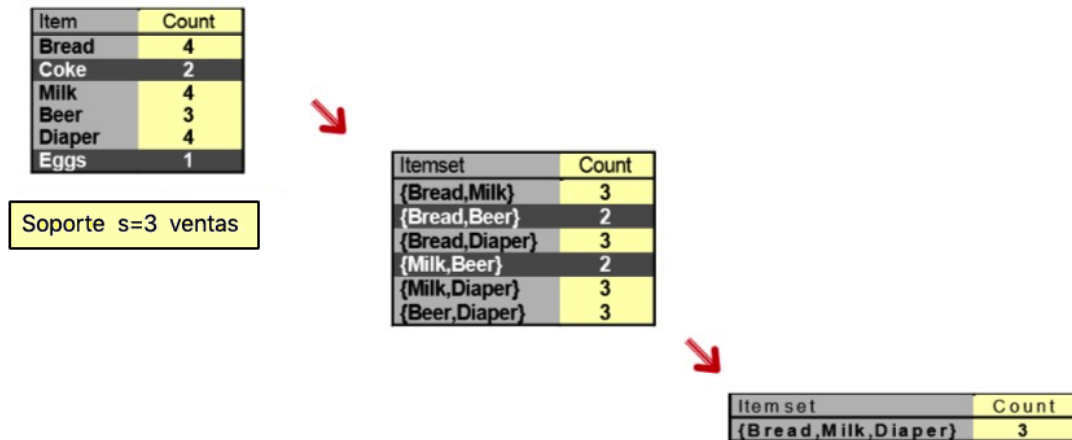
15/36

Algoritmo a Priori: $\mathbb{P}(Z) \geq s$



16/36

Algoritmo a Priori: $\mathbb{P}(Z) \geq s$



17/36

Observación:

Propiedad:

Sea Z un conjunto de k items ($|Z| = k$) que tiene soporte c .

- Cualquier subconjunto, Z' , no vacío de Z tiene soporte $\geq c$.

Demostración.

Supongamos que Z corresponde a los productos de las primeras k coordenadas del vector de compras.

$$s = \mathbb{P}(Z) = \mathbb{P}((X_1, X_2, \dots, X_k) = (1, 1, \dots, 1))$$

Z' es un subconjunto de Z , por ejemplo los dos primeros productos.

$$\mathbb{P}(Z') = \mathbb{P}((X_1, X_2) = (1, 1))$$

$$= \sum_{j_1, j_2, \dots, j_{k-2} \in \{0,1\}} \mathbb{P}((X_1, X_2, X_3, \dots, X_k) = (1, 1, j_1, j_2, \dots, j_{k-2}))$$

$$= \mathbb{P}((X_1, X_2, \dots, X_k) = (1, 1, \dots, 1)) + \text{probabilidades}$$

$$= s + \text{probabilidades} \geq s$$

18/36

REGLAS DE ASOCIACIÓN

19/36

Reglas de Asociación

Una vez que tenemos nuestro listado con los distintos Z con soporte $\geq s$.

- 1 Particionamos cada Z . $S \cup T = Z$ con $S \cap T = \emptyset$.
- 2 Calculamos la confianza de la regla $S \rightarrow T$

(si compraste S te recomiendo T)

20/36

Reglas de Asociación

Supongamos que el conjunto de items $\{A, B, C\}$ es uno de los que tiene confianza $\geq s$

- ① Particionamos $\{A, B, C\}: S \rightarrow T$
 - $\{A\} \rightarrow \{B, C\}, \{B\} \rightarrow \{A, C\}, \{C\} \rightarrow \{A, B\}$
 $\{A, B\} \rightarrow \{C\}, \{A, C\} \rightarrow \{B\}, \{B, C\} \rightarrow \{A\}$
- ② Calculamos la confianza de cada una de estas reglas de asociación $S \rightarrow T$.
- ③ Presentamos las que tienen confianza mayor a c .

21 / 36

Reglas de Asociación

Confianza

$$\text{Confianza}(S \rightarrow T) = \mathbb{P}(T|S) = \frac{\mathbb{P}(T \cap S)}{\mathbb{P}(S)} = \frac{\text{Soporte}(S \cup T)}{\text{Soporte}(S)}. \text{ (raro, ¿no?)}$$

22 / 36

Reglas de Asociación

Eventos

T = compra los artículos C y D . $\mathbf{X} = (X_1, X_2, 1, 1, X_5, \dots, X_N)$

S = compra el artículos A y B . $\mathbf{X} = (1, 1, X_3, X_4, X_5, \dots, X_N)$

$T \cap S$ = compra los A, B, C , y D . $\mathbf{X} = (1, 1, 1, 1, X_5, \dots, X_N)$

$$\text{confianza}(S \rightarrow T) = \mathbb{P}(T|S) = \frac{\mathbb{P}(T \cap S)}{\mathbb{P}(S)}$$

Conjuntos de artículos

$T = \{C, D\}$.

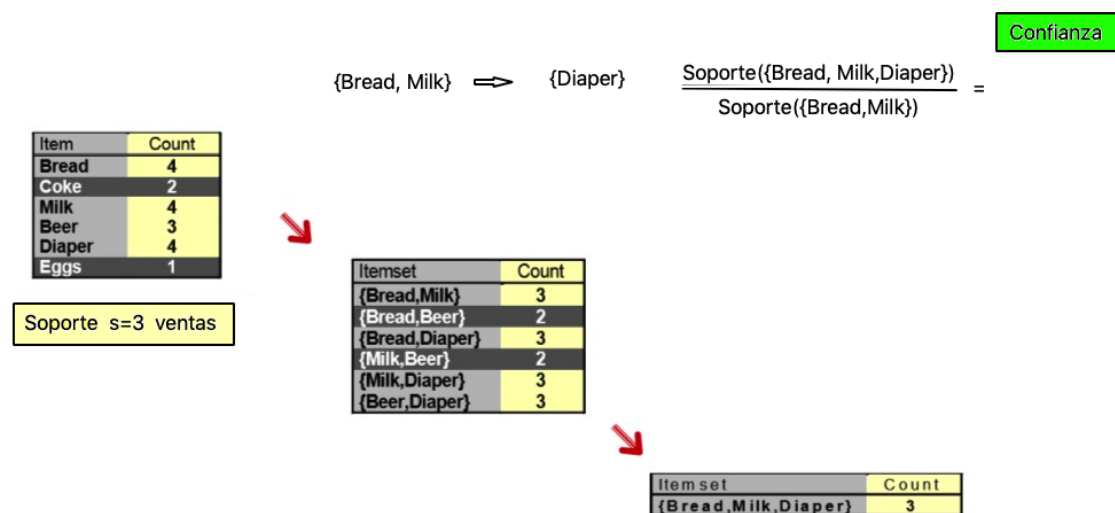
$S = \{A, B\}$.

$T \cup S = \{A, B, C, D\}$.

$$\text{Confianza}(S \rightarrow T) = \frac{\text{Soporte}(S \cup T)}{\text{Soporte}(S)}$$

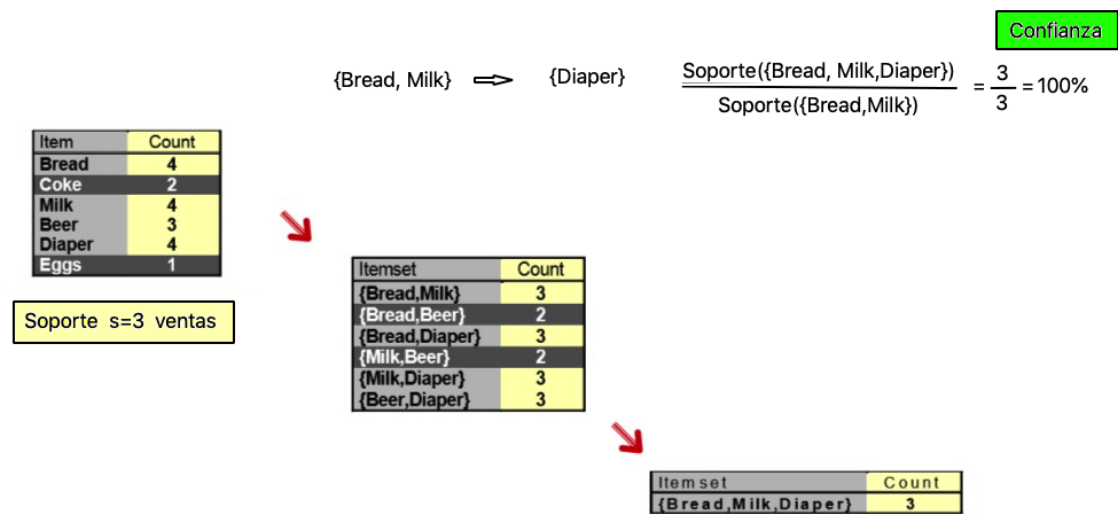
23 / 36

Reglas de Asociación



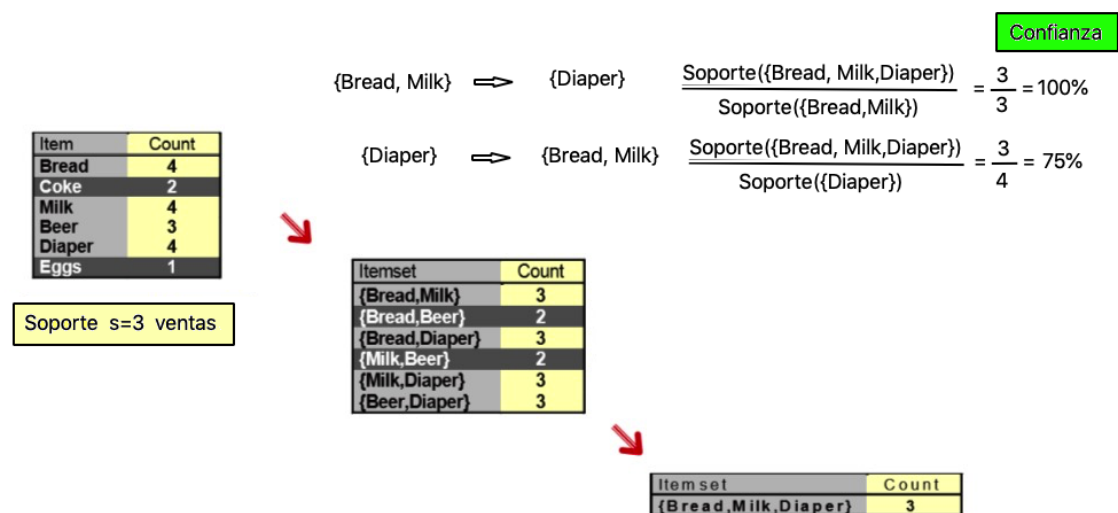
24 / 36

Reglas de Asociación



25 / 36

Reglas de Asociación



26 / 36

Reglas de Asociación con número de items fijos

Una vez que tenemos nuestro listado con los distintos Z con soporte $\geq s$.

- ① Particionamos cada Z . $S \cup T = Z$ con $S \cap T = \emptyset$.
- ② Calculamos la confianza de la regla $S \rightarrow T$

(si compraste S te recomiendo T)

Al fijar en número de items en la regla de asociación en $|Z|$. Donde $S \cup T = Z$ podemos podar el árbol de reglas de decisión.

27 / 36

Reglas de Asociación con número de items fijos

Al fijar en número de items en la regla de asociación en $|Z|$. Donde $S \cup T = Z$ podemos podar el árbol de reglas de decisión.

Supongamos que $Z = \{A, B, C, D\}$. Las posibles reglas son:

- $\{A\} \rightleftharpoons \{B, C, D\}, \{B\} \rightleftharpoons \{A, C, D\}, \{C\} \rightleftharpoons \{A, B, D\}, \{D\} \rightleftharpoons \{A, B, C\}$
- $\{A, B\} \rightleftharpoons \{C, D\}, \{A, C\} \rightleftharpoons \{B, D\}, \{A, D\} \rightleftharpoons \{B, C\}$

Propiedad:

Sea $S_1 \cup T_1 = Z$ y $S_2 \cup T_2 = Z$ con $\dim(S_1) > \dim(S_2)$. Se cumple

$$\text{Confianza}(S_1 \rightarrow T_1) \geq \text{Confianza}(S_2 \rightarrow T_2)$$

28 / 36

Reglas de Asociación con número de items fijos

Propiedad:

Sea $S_1 \cup T_1 = Z$ y $S_2 \cup T_2 = Z$ con $\dim(S_1) > \dim(S_2)$. Se cumple

$$\text{Confianza}(S_1 \rightarrow T_1) \geq \text{Confianza}(S_2 \rightarrow T_2)$$

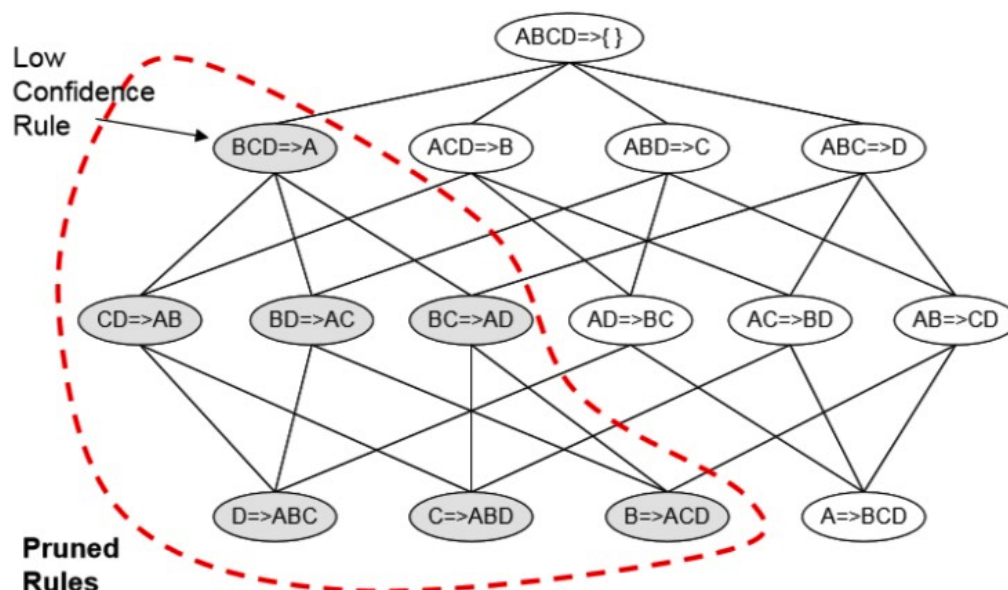
Demostración.

$\text{Confianza}(S_k \rightarrow T_k) = \frac{\text{Soporte}(Z)}{\text{Soporte}(S_k)}$ y como
 $\text{Soporte}(S_k) \geq \text{Soporte}(S_{k-1})$



29/36

Pruning de Reglas de Asociación



30/36

ALTERNATIVAS A LA CONFIANZA: LIFT, LEVERAGE

31 / 36

Confianza vs Lift

Confianza:

$$\text{Confianza}(S \rightarrow T) = \mathbb{P}(T|S)$$

Problema: ¿Qué pasa si T se compra casi siempre?

- Entonces $\mathbb{P}(T|S)$ probablemente sea alto.
- Peor aún si S y T son indep ($\mathbb{P}(T|S) = \mathbb{P}(T)$) S no predice nada.

Lift:

$$\text{Lift}(S \rightarrow T) = \frac{\mathbb{P}(T|S)}{\mathbb{P}(T)} = \frac{\text{Confianza}(S \rightarrow T)}{\text{Soporte}(T)}$$

- S y T son indep $\leftrightarrow \text{Lift}(S \rightarrow T) = 1$.
- Recomendamos $S \rightarrow T$ cuando $\text{Lift} > 1$.

32 / 36

Confianza vs Lift

Comparación

- Confianza: medida que dice las chances de que se compre T cuando compraste S .
- Lift: medida que compara el grado de dependencia entre S y T . Mide cuán buena es la regla respecto al azar.

33 / 36

Otras medidas

Coverage:

$$Coverage(S \rightarrow T) = \mathbb{P}(S) = Soporte(S) = \frac{Soporte(S \rightarrow T)}{Confianza(S \rightarrow T)} = \frac{\mathbb{P}(S \cap T)}{\mathbb{P}(T|S)}$$

Leverage:

$$Leverage(S \rightarrow T) = \mathbb{P}(S \cap T) - \mathbb{P}(S) \mathbb{P}(T)$$

$$Leverage(S \rightarrow T) = soporte(S \cup T) - soporte(S)soporte(T)$$

Added Value:

$$AD(S \rightarrow T) = \mathbb{P}(T|S) - \mathbb{P}(T) = confianza(S \rightarrow T) - soporte(T)$$

34 / 36

Paquete *arules*

```
> transacciones = read.transactions(file = "datos_groceries.csv",  
format = "single", sep = ";", header = TRUE, cols = c("id_compra",  
"item"), rm.duplicates = TRUE)  
> soporte=0.1; confianza=0.7  
> reglas=apriori(data = transacciones, parameter = list(Support =  
soporte, confidence = confianza)) #, minlen = 3, maxlen = 5  
> inspect(reglas)
```

35 / 36

Resumen

- Confianza: medida que dice las chances de que se compre T cuando compraste S .
- Lift: medida que compara el grado de dependencia entre S y T . Mide que tan buena es la regla respecto al azar.
- Soporte de la regla: indica el impacto en término de ventas totales.

36 / 36