

# Probabilidad y Análisis de Datos

Daniel Fraiman

Maestría en Ciencia de Datos, Universidad de San Andrés

1 / 35

## 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).

2 / 35

# 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.

3 / 35

# Aplicaciones en espacios discretos grandes discretos

Customers Who Bought This Item Also Bought



The screenshot displays a row of five book covers recommended for customers who purchased 'Pattern Recognition and Machine Learning' by Christopher M. Bishop. Each recommendation includes the book title, author, star rating, number of reviews, and price. The books are:

- Pattern Recognition and Machine Learning (Information Science and...)** by Christopher Bishop. 4.5 stars, 115 reviews. Hardcover, \$60.76. Prime.
- Learning From Data** by Yaser S. Abu-Mostafa. 4.5 stars, 88 reviews. Hardcover.
- The Elements of Statistical Learning: Data Mining, Inference, and Prediction...** by Trevor Hastie. 4.5 stars, 50 reviews. Hardcover, \$62.82. Prime.
- Probabilistic Graphical Models: Principles and Techniques (Adaptive...** by Daphne Koller. 4.5 stars, 28 reviews. Hardcover, \$91.66. Prime.
- Foundations of Machine Learning (Adaptive Computation and...** by Mehryar Mohri. 4.5 stars, 8 reviews. Hardcover, \$65.68. Prime.

4 / 35

## 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 conjunto que tiene las compras actuales.

$$S = \{A, B\}, ¿S \rightarrow T?$$

¿Le recomendamos  $T$ ?

5/35

## Reglas de Asociación e items frecuentes

Problema matemático:

Dado  $S$ , ¿cuál es el producto  $T$  que maximiza la probabilidad de compra? Llamemos  $\Theta$  al conjunto que tiene a todos los productos del supermercado ( $\Theta = \{A, B, \dots, W\}$ ).

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

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

Quiero el producto que maximiza la probabilidad de que lo compre las personas que van a buscar otro producto distinto, y terminan agregando esta compra.

6/35

# 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):

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

Cual es la probabilidad de llevar leche en polvo y pañales

los tickets, deben tener si o si ambos productos y pueden tener mas tambien, pero no pueden faltar la leche en polvo y pañales

7/35

# Reglas de Asociación e items frecuentes

## Dificultad en la estimación:

El conjunto  $\Theta$  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  $\mathbb{P}(Y|S) \ll 1$ .

tengo 2 alternativas, comprar o no comprar

## 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.

la probabilidad de comprar Y, cuando tenga S, donde S es lo que tengo en el carrito e Y es lo que quiero hacer comprar

producto que ya tengo en el changuito

producto que quiero que lleven

8/35

## 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.

vos podes tener coca y te ofrezco fernet, vale lo mismo que vos ya tengas el fernet y te ofrezca coca cola

La probabilidad de intercepcion de los productos, se lo llama SOPORTE

La confianza es cuanta “fe”, le tenes que cuando lleven el producto A, lleven tmb el producto B

9/35

## 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  $S$ .

10/35

### 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  $s$  (Soporte:  $\mathbb{P}(S \cap T) \geq s$ )
- Fijamos  $c$  (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  $c$ .

11/35

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

12/35

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

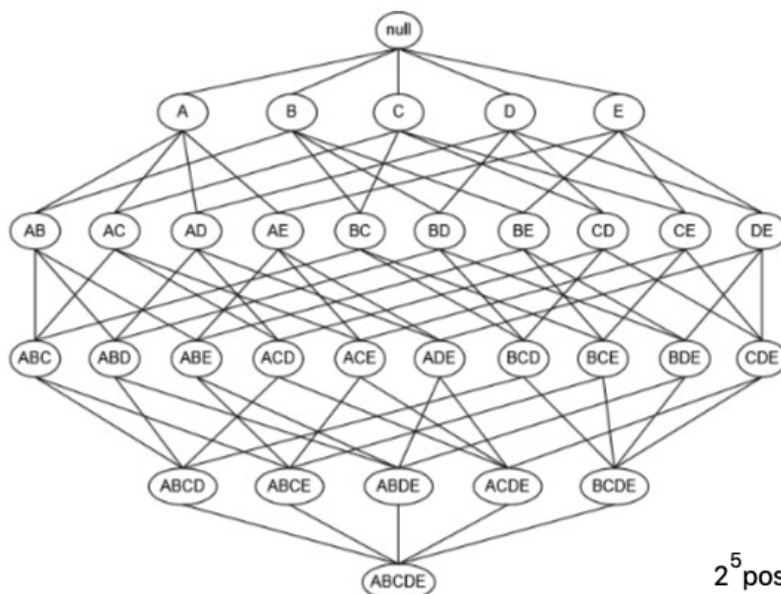
**Objetivo:** Encontrar todos los conjuntos  $Z$  (con  $|Z| < k$ ) que cumplen  $\mathbb{P}(Z) \geq s$ .

### 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.

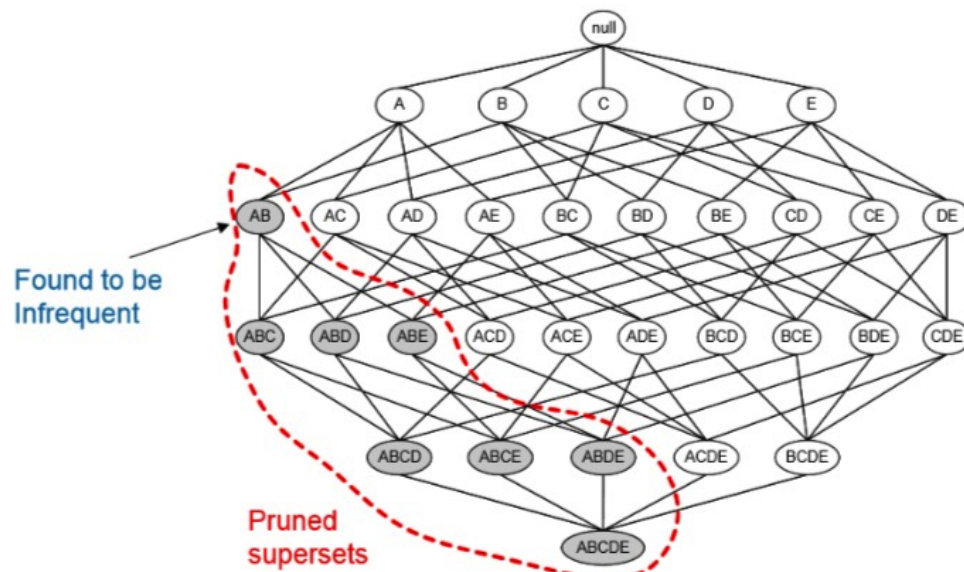
13 / 35

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



14 / 35

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



15/35

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

Ej :En cuantos tickets de todos los que tengo compraron Pan → en 4 tickets

Item	Count
Bread	4
Coke	2
Milk	4
Beer	3
Diaper	4
Eggs	1

Soporte  $s=3$  ventas

que grupos podria armar con los productos que sobrevivieron, y luego, vuelvo a los que tienen menos de 3 ventas

Itemset	Count
{Bread,Milk}	3
{Bread,Beer}	2
{Bread,Diaper}	3
{Milk,Beer}	2
{Milk,Diaper}	3
{Beer,Diaper}	3

Itemset	Count
{Bread,Milk,Diaper}	3

Me voy a quedar con los productos que se vendieron mas de 3 veces, esa es la funcion del soporte

16/35



## Observación:

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

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

### Proof

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.

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

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

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

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

17/35

## REGLAS DE ASOCIACIÓN

18/35

## 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  $T$ )

Confianza  $\rightarrow$  Probabilidad de que compres  $T$ , dado que llevaste  $S$

19/35

## Reglas de Asociación

Supongamos que el conjunto de items  $\{A, B, C\}$  es uno de los que tiene <sup>SOPORTE</sup> 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  $c$ .

20/35

## Reglas de Asociación

### Confianza

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

La confianza de la “Regla S-T” es igual

21 / 35

## 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

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

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

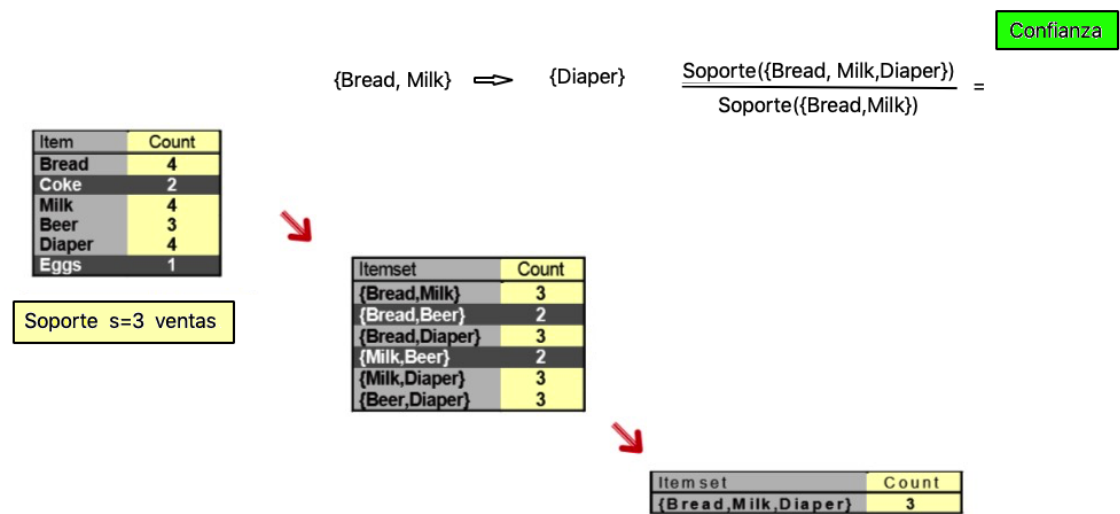
$T \cup S = \{A, B, C, D\}$ .

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

En programacion, para decirte que quieren las 2 cosas, utilizan UNION, en vez de intercepcion para decir que es la suman de las 2 cosas

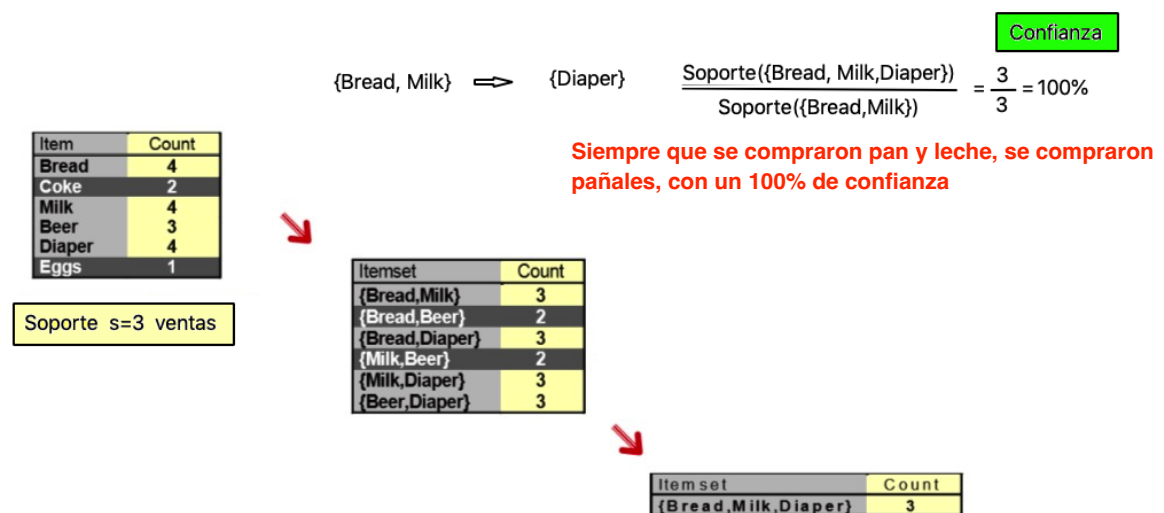
22 / 35

# Reglas de Asociación



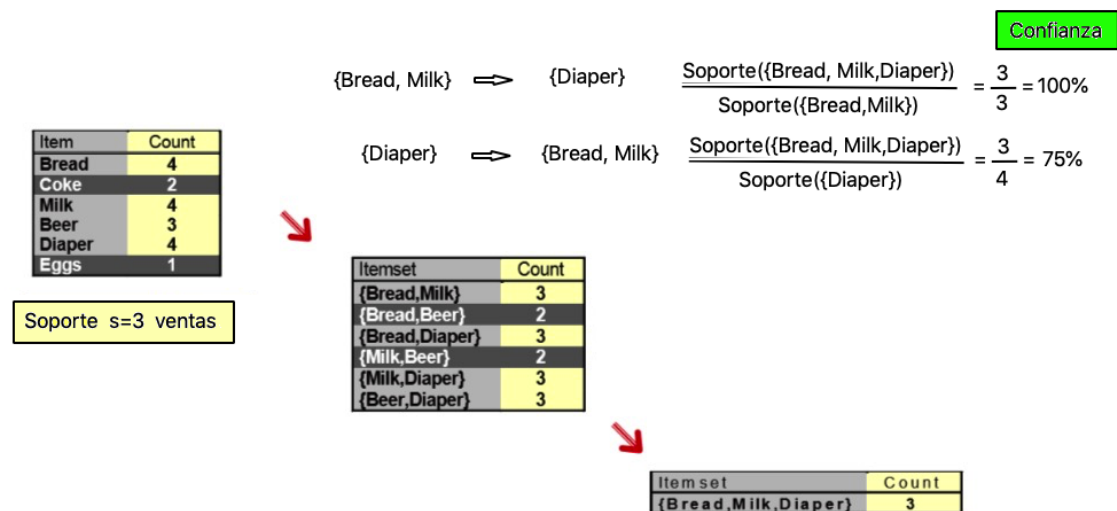
23 / 35

# Reglas de Asociación



24 / 35

# Reglas de Asociación



25 / 35

# 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 arbol de reglas de decisión.

26 / 35

## Reglas de Asociación con número de items fijos

Al fijar el 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:

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

Propiedad:

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

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

27 / 35

## 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

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

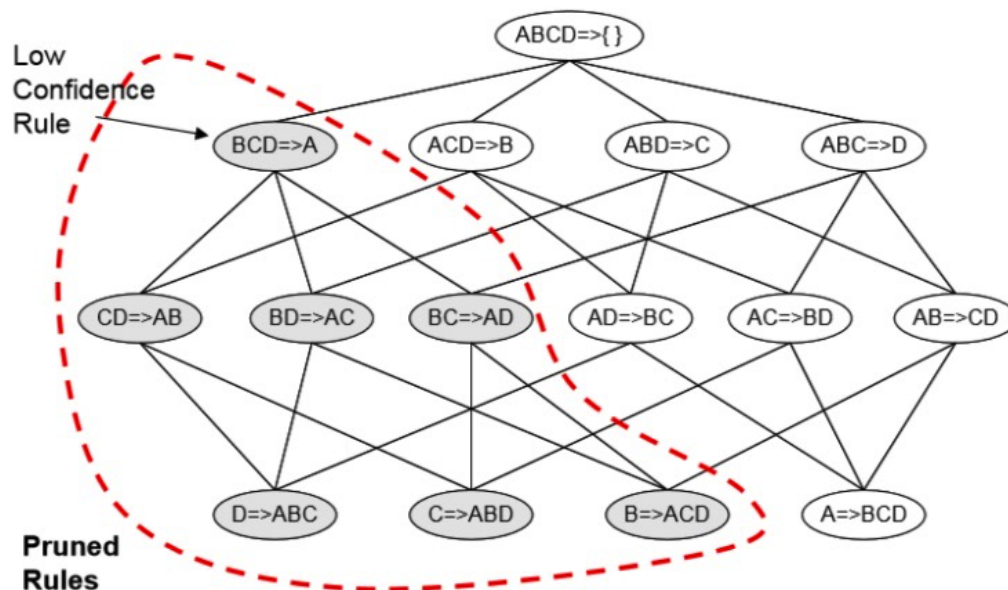
Demostración.

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



28 / 35

## Pruning de Reglas de Asociación



29 / 35

ALTERNATIVAS A LA CONFIANZA: LIFT,  
LEVERAGE

30 / 35

## 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.

independiente

### Lift:

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

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

31 / 35

## 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.

32 / 35



## Otras medidas

Tambien se debe mirar el soporte, la cantidad de ventas de S, el producto que lleva originalmente mi cliente, y que luego yo le voy a sugerir el producto T

Coverage:

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

Leverage:

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

$$\text{Leverage}(S \rightarrow T) = \text{Support}(S \cup T) - \text{Support}(S)\text{Support}(T)$$

Es lo mismo que lift, pero en vez de ser mayor a 1, menor a 1, o 0, puede ser positivo, negativo o 0.

33 / 35

## Reglas de Asociación con

arules → “reglas de asociacion”

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)
```

34 / 35

- 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.