

Association Analysis



Soroosh Nazem

Introduction

- One of the most important areas in data mining
- Talks about correlations between items on a set in transactions
- Application Areas:
 - Market Basket Analysis
 - Recommendation Systems
 - Customer Relationship Management (CRM)
 - Medical Diagnosis
 - Census Data
 -

Application:Market Basket

- Predicting of customers behaviors based on previous transactions.

- Example (based on table):

- Most customers buy “whole milk”.
- Customers with “curd” buy “whole milk” 100%.

Why Supermarket need this info?

- loyalty program management
- Location of items in supermarket
- promotions/discount management
-

TRANSACTION	ITEMS
T1	Tropical fruit, yogurt, coffee
T2	Whole milk
T3	Whole milk, butter, yogurt, rice
T4	Whole milk, cereals
T5	Citrus fruit, tropical fruit, whole milk, butter, curd, yogurt, flour, bottled water, dishes
T6	chicken, tropical fruit
T7	Root vegetables, other vegetables, whole mik, beverages, sugar
T8	Berries, yogurt
T9	Whole milk, curd, yogurt, pastry

Application: Recommender Systems

Linked in

Linkedin:

- Connection suggestion
- Job search suggestion

amazon.com

Amazon:

- Product recommendation based on previous searches/orders

NETFLIX

Netflix:

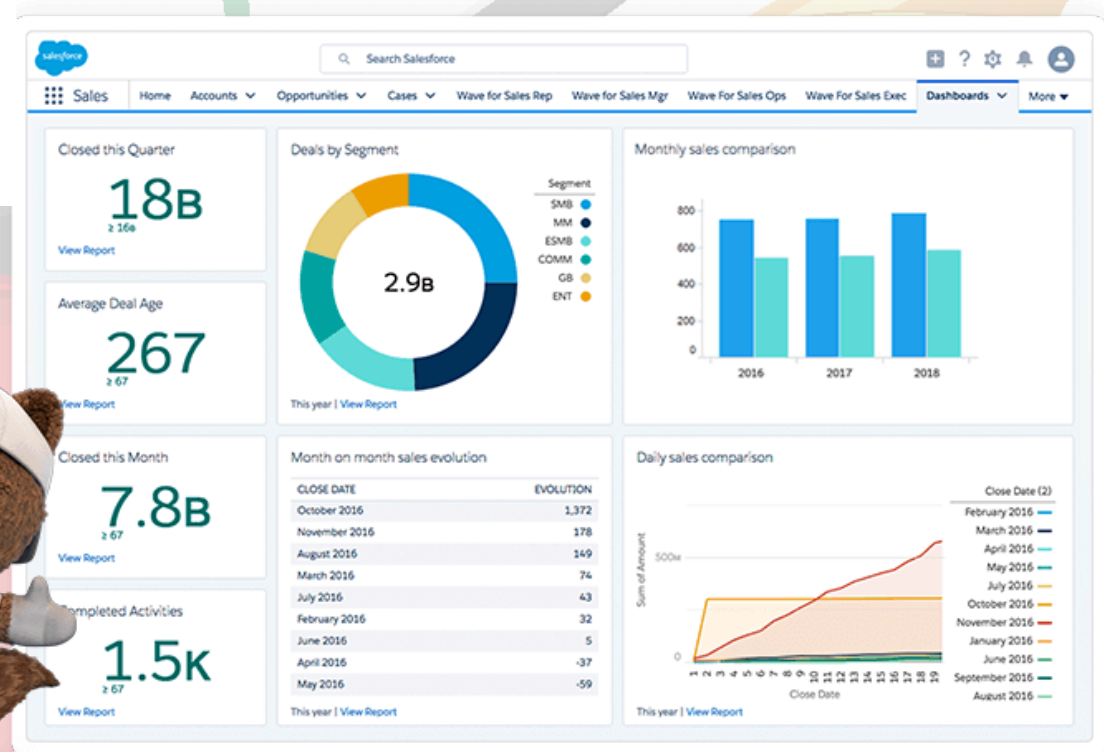
- Propose movies based on previous watched movies

Application: CRM

HubSpot



- What parameters have Increased the number of leads?
- What kinds of leads are more likely to be converted to a contract?
-



Definitions:

Consider a set of items $I=\{I_1, I_2, \dots, I_n\}$:

- A Transaction T_i is a items I_j where $1 \leq j \leq n$ and $T_i \subseteq I$
- $T=\{T_1, T_2, \dots\}$ a set of Transactions
- An “Association Rule (AR)” written in the form of $X \Rightarrow Y$ where $X \cap Y = \emptyset$ and $X, Y \subseteq I$.

Example:

$I=\{\text{milk, sugar, salt, yogurt, cheese, butter, vegetable, bottled water, rice, fruit}\}$

$T_1=\{\text{milk, cheese, butter}\}$, $T_2=\{\text{milk, rice, salt}\}$, $T_3=\{\text{vegetable, fruit, rice}\}$,
 $T_4=\{\text{yogurt, salt, rice}\}$, $T_5=\{\text{milk, sugar, salt, butter, rice}\}$, $T_6=\{\text{vegetables, bottled water}\}$, $T_7=\{\text{vegetables, fruit}\}$

$AR1:\{\text{rice}\} \Rightarrow \{\text{salt}\}$, $AR2:\{\text{cheese, butter}\} \Rightarrow \{\text{milk}\}$, $AR2: \{\text{fruit}\} \Rightarrow \{\text{bottled water,}$

Metrics: SUPPORT

DEF

$SUPP(X \Rightarrow Y) =$
 $(\# \text{ Transaction that contain } X \text{ and } Y) / \# \text{ Total Transactions}$

EX

$SUPP(\{whole\ milk\} \Rightarrow \{yogurt\}) = 3/9$
 **$SUPP(\{tropical\ fruit\}$
 $\Rightarrow \{rice\}) = 0/9 = 0$**

NOTE

**$SUPP(\{yogurt\} \Rightarrow \{whole$
 $\{x\}) = 3/9$
 $SUPP(X \Rightarrow Y) =$
 $SUPP(Y \Rightarrow X)$
 $Range: [0, 1]$**

TRANSACTION	ITEMS
T1	Tropical fruit, yogurt, coffee
T2	Whole milk
T3	Whole milk, butter, yogurt, rice
T4	Whole milk, cereals
T5	Citrus fruit, tropical fruit, whole milk, butter, curd, yogurt, flour, bottled water, dishes
T6	chicken, tropical fruit
T7	Root vegetables, other vegetables, whole milk, beverages, sugar
T8	Berries, yogurt
T9	Whole milk, curd, yogurt, pastry

Metrics: CONFIDENCE

DEF

$$\text{CONF}(X \Rightarrow Y) = \text{SUPP}(X \Rightarrow Y) / \text{SUPP}(X)$$

EX

$$\text{CONF}(\{\text{whole milk}\} \Rightarrow \{\text{yogurt}\}) = 3/6 = 50\%$$

$$\text{CONF}(\{\text{tropical fruit}\} \Rightarrow \{\text{rice}\}) = 0/9 = 0$$

$$\text{CONF}(\{\text{yogurt}\} \Rightarrow \{\text{whole milk}\}) = 3/5 = 60\%$$

NOTE

$$\begin{aligned} \text{CONF}(X \Rightarrow Y) &\neq \text{CONF}(Y \Rightarrow X) \\ \text{Range: } &[0, 1] \end{aligned}$$

TRANSACTION	ITEMS
T1	Tropical fruit, yogurt, coffee
T2	Whole milk
T3	Whole milk, butter, yogurt, rice
T4	Whole milk, cereals
T5	Citrus fruit, tropical fruit, whole milk, butter, curd, yogurt, flour, bottled water, dishes
T6	chicken, tropical fruit
T7	Root vegetables, other vegetables, whole milk, beverages, sugar
T8	Berries, yogurt
T9	Whole milk, curd, yogurt, pastry

Metrics: LIFT

DEF

$$LIFT(X \Rightarrow Y) = CONF(X \Rightarrow Y) / SUPP(Y)$$

EX

$$LIFT(\{whole\ milk\} \Rightarrow \{yogurt\}) = (3/6) / (5/9) = 0.9$$

$$LIFT(\{tropical\ fruit\} \Rightarrow \{rice\}) = 0$$

$$LIFT(\{yogurt\} \Rightarrow \{whole\ milk\}) = (3/5) / (6/9) = 0.9$$

NOTE

$$LIFT(X \Rightarrow Y) \equiv LIFT(Y \Rightarrow X)$$

Range: $[0, \infty)$

TRANSACTION	ITEMS
T1	Tropical fruit, yogurt, coffee
T2	Whole milk
T3	Whole milk, butter, yogurt, rice
T4	Whole milk, cereals
T5	Citrus fruit, tropical fruit, whole milk, butter, curd, yogurt, flour, bottled water, dishes
T6	chicken, tropical fruit
T7	Root vegetables, other vegetables, whole milk, beverages, sugar
T8	Berries, yogurt
T9	Whole milk, curd, yogurt, pastry

Metrics: LEVERAGE

DEF

$$LEV(X \Rightarrow Y) = SUPP(X \Rightarrow Y) - SUPP(X) \cdot SUPP(Y)$$

EX

$$LEV(\{whole\ milk\} \Rightarrow \{yogurt\}) = -1/27$$

$$LEV(\{tropical\ fruit\} \Rightarrow \{rice\}) = 0 - (3/9)(1/9) = -1/27$$

$$LEV(\{yogurt\} \Rightarrow \{whole\ milk\}) = -1/27$$

$$LEV(\{curd\} \Rightarrow \{yogurt\}) = 2/9 - (2/9)(5/9) = 8/81$$

NOTE

$$LEV(X \Rightarrow Y) = LEV(Y \Rightarrow X)$$
$$Range: [-1, 1]$$

TRANSACTION	ITEMS
T1	Tropical fruit, yogurt, coffee
T2	Whole milk
T3	Whole milk, butter, yogurt, rice
T4	Whole milk, cereals
T5	Citrus fruit, tropical fruit, whole milk, butter, curd, yogurt, flour, bottled water, dishes
T6	chicken, tropical fruit
T7	Root vegetables, other vegetables, whole milk, beverages, sugar
T8	Berries, yogurt
T9	Whole milk, curd, yogurt, pastry

Metrics: CONVICTION

DEF

$$\text{CONV}(X \Rightarrow Y) = \frac{(1 - \text{SUPP}(Y))}{(1 - \text{CONF}(X \Rightarrow Y))}$$

EX

$$\begin{aligned}\text{CONV}(\{\text{whole milk}\} \Rightarrow \{\text{yogurt}\}) &= 8/9 \approx 0.89 \\ \text{CONV}(\{\text{tropical fruit}\} \Rightarrow \{\text{rice}\}) &= 8/9 / (1-0) \approx 0.89 \\ \text{CONV}(\{\text{yogurt}\} \Rightarrow \{\text{whole milk}\}) &= 5/6 \approx 0.83 \\ \text{CONV}(\{\text{curd}\} \Rightarrow \{\text{yogurt}\}) &= 4/9 / (1-1) = \infty\end{aligned}$$

NOTE

$$\begin{aligned}\text{CONV}(X \Rightarrow Y) &\neq \text{CONV}(Y \Rightarrow X) \\ \text{Range: } &[0, \infty)\end{aligned}$$

TRANSACTION	ITEMS
T1	Tropical fruit, yogurt, coffee
T2	Whole milk
T3	Whole milk, butter, yogurt, rice
T4	Whole milk, cereals
T5	Citrus fruit, tropical fruit, whole milk, butter, curd, yogurt, flour, bottled water, dishes
T6	chicken, tropical fruit
T7	Root vegetables, other vegetables, whole milk, beverages, sugar
T8	Berries, yogurt
T9	Whole milk, curd, yogurt, pastry

FREQUENT ITEMSET

DEF

An item set X where $SUPP(X) \geq \text{minsup}$ is called a frequent itemset

EX

$\text{minsup} = \frac{1}{3}$

frequent itemsets:

$SUPP(\{\text{whole milk}\}) = \frac{2}{3}$

$SUPP(\{\text{yogurt}\}) = \frac{5}{9}$

$SUPP(\{\text{yogurt, whole milk}\}) = \frac{1}{3}$

$SUPP(\{\text{tropical fruit}\}) = \frac{1}{3}$

TRANSACTION	ITEMS
T1	Tropical fruit, yogurt, coffee
T2	Whole milk
T3	Whole milk, butter, yogurt, rice
T4	Whole milk, cereals
T5	Citrus fruit, tropical fruit, whole milk, butter, curd, yogurt, flour, bottled water, dishes
T6	chicken, tropical fruit
T7	Root vegetables, other vegetables, whole milk, beverages, sugar
T8	Berries, yogurt
T9	Whole milk, curd, yogurt, pastry

FREQUENT ITEMSET

Q

How to find frequent item sets?

EX

Algorithms

- *Apriori*
- *FP-growth*
- *ECLAT*

Different methods, same results

TRANSACTION	ITEMS
T1	Tropical fruit, yogurt, coffee
T2	Whole milk
T3	Whole milk, butter, yogurt, rice
T4	Whole milk, cereals
T5	Citrus fruit, tropical fruit, whole milk, butter, curd, yogurt, flour, bottled water, dishes
T6	chicken, tropical fruit
T7	Root vegetables, other vegetables, whole milk, beverages, sugar
T8	Berries, yogurt
T9	Whole milk, curd, yogurt, pastry