

15 Introduction to Data Mining

15.1 Introduction to principle methods

15.2 Mining association rule

see also: A. Kemper, Chap. 17.4, Kifer et al.: chap 17.7 ff

15.1 Introduction



"**Discovery** of useful, possibly unexpected **patterns in data**"
J. Ullman

"Data mining is the **process of discovering** meaningful new **correlations, patterns and trends** by sifting through **large amounts of data** stored in repositories, using pattern recognition technologies as well as statistical and mathematical techniques."

The Gartner Group

015-DBS-DMining-2

Introduction



Example

- Large amount of data

Assumption: 1 phone call per day per adults in Germany
⇒ 50 Mill / day, **15 Billion phone calls / year**

- Find "**hidden knowledge**" – e.g. correlations between attributes

"Phone calls in the evening are longer than during the day"

... a DWH task!

"Young people in the age group 20-25 with a new (< 6 month) mobile phone contract living in big cities make phone calls in the evening longer than on average".

Might be very interesting in order to optimize tariffs.

015-DBS-DMining-3

Examples



More examples

Personalized recommendations

Recommend products to customers which 'similar' customers like (have purchased, ...)

Data: all purchases, all shopping carts, click streams

Degree of credit worthiness

Analyze potential borrower: income? profession? where living?...

Data: borrowers in the past and their *known* worthiness

Estimate scholastic success

Predict success of studies based on grades, social background, ...

Data: historic data from Campus MGM ;)

Typical: prediction of the future for individuals based on history of others

015-DBS-DMining-4

Methods



All kinds of **statistical techniques**

- simple counting
- estimation of probabilities
- finding parameters of distributions,

Challenge for DB technology:

scalable algorithms for **very large data sets**

New challenges: **Real Time Data Mining** on streams

015-DBS-DMining-5

(Un) Supervised learning



Principle methods

(1) Learn a model (supervised learning)
from a large set of retrospective data)

Very different models, e.g.

Association rules, parameters of probability distributions, decision trees, classes of "homogeneous" objects

(2) Cluster data (unsupervised learning)
group data according to some similarity criterion

⇒ **Machine Learning**

015-DBS-DMining-6

Classification Freie Universität Berlin

Methods of supervised Learning

Classification: the general problem
 A **training set** of classified records is available.
 Infer a model, which predicts the class of a new record

Predict attribute value $x=c$ according to model F :
 $F(a_1, \dots, a_n) = c$

May sometimes be written as **classification rule** :

$(age < 40) \wedge (sex = 'm') \wedge (make = 'Golf GTI') \wedge (hp > 150)$
 $\Rightarrow (risk = 'high')$

015-DBS-DMining-7

Classification (2) Freie Universität Berlin

$k=2$ classes, **features** learned in training phase.

New objects assigned according to their features.

Classification by decision tree model – which has to be learned.
 Model: the decision tree

015-DBS-DMining-8

Classification (3) Freie Universität Berlin

Naïve Bayes Classification

Suppose you are going to meet a person named *Gerit Robben* at the Amsterdam airport. Unfortunately you do not know if *Gerit* is a male or female first name...

Would like to know:
 $p(\text{female} | \text{Gerit})$

Probability, that a male is named Gerit

$$p(\text{female} | \text{Gerit}) = \frac{p(\text{Gerit} | \text{female}) p(\text{female})}{p(\text{Gerit})}$$

Probability of being named Gerit

Probability of being a female

Probabilistic Modeling

015-DBS-DMining-9

Classification (4) Freie Universität Berlin

Support vector Machine (SVM)

Find a hyperplane which separates classes in an optimal way

Optimization problem
 Model: ~ geometry

015-DBS-DMining-10

Classification (5) Freie Universität Berlin

k Nearest Neighbors

Calculation of similarity of probe to its (classified!) neighbors.
 Model: the classes

015-DBS-DMining-11

Unsupervised Learning (1) Freie Universität Berlin

Unsupervised learning methods
Clustering

Group homogenous data into a cluster according to some similarity measure
 e.g. find subsets of customers with similar shopping patterns

hierarchical / taxonomic

needed: **similarity measure**

015-DBS-DMining-12

Similarity is ...



... hard to define. "We know it when we see it"

form E. Keogh

015-DBS-DMining-13

Unsupervised Learning(2)

Association rules

Market basket analysis:

customer transaction data: `tid, time, {articles}`

Find rules $X \Rightarrow Y$, with particular confidence

e.g.

Those buying sauce, meat and spaghetti
buy red wine with high confidence

(whatever that may be)

Naiv algorithm: count how many spaghetti buyers
also bought red wine

015-DBS-DMining-14

Mining Process

The Data mining process

1. Data gathering, joining, reformatting

e.g. Oracle: max 1000 attributes \Rightarrow transform into
"transactional format": (id, attr_name, value)

2. Data cleansing

- eliminate outliers
- check correctness based on domain specific heuristics
- check values in case of redundancy, ...

3. Build model (training phase).

4. Apply to [new] data

015-DBS-DMining-15

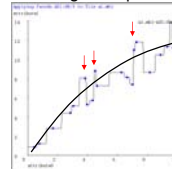
Errors

Training set error

Check with records of training set if predicted
value equals known value in record

Overfitting: Results are influenced by unimportant
details, fit perfect to training set, but spoiled by noise.

Training data: perfect fit, but..



015-DBS-DMining-16

15.2 Association rules

Applications

Discover patterns of co-occurring products in a shopping
basket

"How many customers buy a printer together
with a PC"

"Shopping basket" is just a metaphor...

015-DBS-DMining-17

Association rules: examples

Baskets = **Web pages**; items = **words**.

Unusual words appearing together in a large number of documents, e.g.,
"British Petrol" and "Gulf" may indicate an interesting relationship.

Baskets = **Documents**, items = **sentences**

If number of co-occurrences of sentences is high: suspicion of
plagiarism.

Common term for baskets, documents, or whatever:

transaction

Only slightly related to DB transactions in some, not all
applications.

015-DBS-DMining-18

Scale

Problem scale very different:

Kaufhof: ~ 100,000 products,
~ $n \cdot 10^9$ purchase transactions

Words in Web pages: ~ 10^6
Web pages: 10^{11} ?

015-DBS-DMining-19

Association rules

The abstract problem

Given a **set of objects** (items) $I = \{i_1, \dots, i_n\}$
and a **bag T** (multiset) of **non-empty subsets t** of I,
the transactions.

Find **association rules** $R_j: \{i_{j_1}, \dots, i_{j_k}\} \rightarrow \{i_{j_{(k+1)}}, \dots, i_{j_{(k+l)}}\}$
Semantics: if $i_{j_1}, \dots, i_{j_k} \in \mathbf{t}, \mathbf{t} \in \mathbf{T}$, then $\{i_{j_{(k+1)}}, \dots, i_{j_{(k+l)}}\} \subseteq \mathbf{t}$
with some predefined probability.

But: **Probability not sufficient** - $\{a, b, c\}$ occurs once in \mathbf{T} ,
and $\{a, b\}$ is a subset of exactly two subsets of I:
 \Rightarrow Probability of $c \in \mathbf{t} \mid \{a, b\} \subseteq \mathbf{t}$ is $1/2$

015-DBS-DMining-20

Association rules

Measures:

Def: support ($A \rightarrow B$) = $P(A, B)$, $A, B \subseteq I$
probability that A and B co-occur in the data set $\mathbf{t} \in \mathbf{T}$
e.g. 0.05 if 5% of all customers bought a printer and a PC

Def.: confidence ($A \rightarrow B$) = $P(B \mid A)$
fraction of transaction \mathbf{t} containing B if \mathbf{t} contains A,
e.g. 0.8: 4 of 5 bought also printer if they bought a PC

Find all rules $r: A \rightarrow B$ with
support(r) \geq **minSupport** and
confidence(r) \geq **minConfidence**

015-DBS-DMining-21

Counting

Base task is counting!

$r: A \rightarrow B$ has support support(r) if there are
support(r) * $|\mathbf{T}|$ subsets \mathbf{t} and $A \cup B \subseteq \mathbf{t}$

Example

$t_1 = \{m, c, b\}$ $t_2 = \{m, p, j\}$
 $t_3 = \{m, b\}$ $t_4 = \{c, j\}$
 $t_5 = \{m, p, b\}$ $t_6 = \{m, c, b, j\}$
 $t_7 = \{c, b, j\}$ $t_8 = \{b, c\}$

Association rule $r: \{m, b\} \rightarrow \{c\}$.

support(r) = $P(\{m, b, c\} \subseteq \mathbf{t}) = 1/4$

confidence(r) = $P(\{c\} \mid \{m, b\}) = 2/4$

015-DBS-DMining-22

Counting

(2) Given a frequent item set $F = \{i_1, \dots, i_k\}$,
then for each $L \subseteq F$
 $r: L \rightarrow F \setminus L$
is a rule with
confidence(r) = support F / support L

since

confidence($L \rightarrow F \setminus L$) = $P(F \setminus L \mid L)$
= $P(F \subseteq \mathbf{t} \wedge L \subseteq \mathbf{t}) / P(L \subseteq \mathbf{t}) = P(F) / P(L)$
= support(F) / support(L)

Task: **Find frequent item sets**
... a trivial counting task??

015-DBS-DMining-23

Counting...

Not so easy....

Suppose 1000 items.

Naïve approach: count all possible subsets...how many?

$2^{1000} - 1$

Wanted:

Clever strategies for finding frequent item sets (fis)
with 1, 2, 3, ... items from a **large number of transactions**.

015-DBS-DMining-24

A Priori approach

Use a priori knowledge about frequent items.

Monotony property:

each subset of a frequent item set is frequent,
 i.e. if $X \subseteq I$ is a frequent item set,
 i.e. $\text{support}(X) > \text{minSupport}$
 then $\forall X' \subseteq X: \text{support}(X') > \text{minSupport}$

Idea for **algorithm** to determine all Frequent Item Sets (FIS):

- (1) Find FIS $\{F: |F| = m\}$ having $|F| > \text{minSupport}$
- (2) Generate candidates $F' = F \cup \{i\}$, $i \in I$ and $\text{support}(\{i\}) > \text{minSupport}$.
- (3) Check each candidate with $m+1$ elements if frequent.

A Priori Algorithm

```

for all items p {
  if p occurs more than minSupport make
  frequent item set with one element:  $F_1^p = \{p\}$ 
  k = 1
  repeat {
    for each  $F_k$  with k elements generate candidates
     $F_{k+1}$  with k+1 elements and  $F_k \subseteq F_{k+1}$ .
    check in database, which candidates occur at least
    minSupport times; (sequential scan of DB)
    k = k+1
  }
  until no new frequent item set found
    
```

A Priori algorithm for finding associations

TransID	Product
111	printer
111	paper
111	PC
111	toner
222	PC
222	scanner
333	printer
333	paper
333	toner
444	printer
444	PC
555	printer
555	paper
555	PC
555	scanner
555	toner

minSupport := 3/5; minCount := minSupport * |T| = 3
 confidence = 0.7

Select product, count(*) from Transactions
 group by product
 having count(*) > minSupport;

product, count(*)	k=1
printer 4	
paper 3	
PC 4	
toner 3	
scanner 2 X	

5 transactions only!

Transactions		k=2	
TransID	Product	FI-candidat	count
111	printer	{printer, paper}	3
111	paper	{printer, PC}	3
111	PC	{printer, scanner}	
111	toner	{printer, toner}	3
222	PC	{paper, PC}	2 X
222	scanner	{paper, toner}	3
333	printer	{paper, scanner}	
333	paper	{PC, scanner}	
333	toner	{PC, toner}	2 X
444	printer	{scanner, toner}	
444	PC	{printer, paper, PC}	2 X
555	printer	{printer, paper, toner}	3
555	paper	{printer, PC, toner}	2 X
555	PC	{paper, PC, toner}	2 X
555	scanner		
555	toner		

Priori-Algorithm

k=3

Generate association rules

Given: set of FI of frequent items

```

for each FI with support >= minSupport:
  { for each subset L ⊂ FI
    define rule R : L → FI \ L
    confidence (R) = support FI / support L
    if confidence(R) >= minConfidence: keep L
  }
    
```

Example:

FI = {printer, paper, toner}

SupportCount = 3

Rule: {printer} ⇒ {paper, toner},

Confidence = SupportCount({printer, paper, toner}) / SupportCount({printer})
 = (3/5) / (4/5)
 = ¾ = 0.75

Increase of confidence

Increase of left hand side (i.e. decrease of right hand side) of a rule increases confidence:

$L \subseteq L^+$, $R^c \subseteq R$ and $F = L \cup R = L^+ \cup R^c$
 $\Rightarrow \text{confidence}(L \rightarrow R) \leq \text{confidence}(L^+ \rightarrow R^c)$

Rule: {printer} ⇒ {paper, toner}

confidence = support({printer, paper, toner}) / support({printer})
 = (3/5) / (4/5)
 = ¾ = 75%

Rule: {printer,paper} ⇒ {toner}

confidence = S({printer, paper, Toner}) / S({printer,paper})
 = (3/5) / (3/5)
 = 1 = 100%

Summary data mining



Important **statistical technique**

Basis algorithms from machine learning

Many different methods and algorithms

Supervised versus **unsupervised**
learning

Efficient implementation on **very large data sets**
essential

Enormous **commercial interest**
(business transactions, web logs, ...)

015-DBS-DMining-31