

Reductionist View: A Priori Algorithm and Vector-Space Text Retrieval

Sargur Srihari

University at Buffalo

The State University of New York

A Priori Algorithm for Association Rule Learning

- Association rule is a representation for local patterns in data mining
- What is an Association Rule?
 - It is a probabilistic statement about the co-occurrence of certain events in the data base
 - Particularly applicable to sparse transaction data sets

Examples of Patterns and Rules

- Supermarket
 - 10 percent of customers buy wine and cheese
- Telecommunications
 - If alarms A and B occur within 30 seconds of each other, then alarm C occurs within 60 seconds with probability 0.5
- Weblog
 - If a person visits the CNN website there is a 60% chance person will visit the ABC News website in the same month

Form of Association Rule

- Assume all variables are binary
- Association Rule has the form:
If $A=1$ and $B=1$ then $C=1$ with probability p
where A, B, C are binary variables and $p = p(C=1|A=1, B=1)$
- Conditional probability p is the *accuracy* or *confidence* of the rule
- $p(A=1, B=1, C=1)$ is the *support*

Accuracy vs Support

If $A=1$ and $B=1$ then $C=1$ with probability $p = p(C=1|A=1, B=1)$
 $p(A=1, B=1, C=1)$ is the *support*

- Accuracy is a conditional probability
 - Given that A and B are present what is the probability that C is present
- Support is a joint probability
 - What is the probability that A, B and C are all present
- Example of three students in class

Goal of Association Rule Learning

If $A=1$ and $B=1$ then $C=1$ with probability $p = p(C=1|A=1, B=1)$
 $p(A=1, B=1, C=1)$ is the *support*

- Find *all* rules that satisfy the constraint that
 - Accuracy p is greater than threshold p_a
 - Support is greater than threshold p_s
- **Example:**
 - Find all rules that satisfy the constraint that accuracy greater than 0.8 and support greater than 0.05

Association Rules are Patterns in Data

If $A=1$ and $B=1$ then $C=1$ with probability $p = p(C=1|A=1, B=1)$
 $p(A=1, B=1, C=1)$ is the *support*

- They are a *weak* form of knowledge
 - They are summaries of co-occurrence patterns in data
- Rather than *strong* statements that characterize the population as a whole
- *If-then-else* here is inherently correlational and not causal

Origin of Association Rule Mining

- Applications involving “market-basket data”
- Data recorded in a database where each observation consists of an actual basket of items (such as grocery items)
- Association rules were invented to find simple patterns in such data in a computationally efficient manner

Basket Data

<i>Basket\Item</i>	A_1	A_2	A_3	A_4	A_5
t_1	1	0	0	0	0
t_2	1	1	1	1	0
t_3	1	0	1	0	1
t_4	0	0	1	0	0
t_4	0	1	1	1	0
t_5	1	1	1	0	0
t_6	1	0	1	1	0

For 5 items there will be $2^5 = 32$ different baskets
Set of baskets typically has a great deal of structure

Data matrix

- N rows (corresponding to baskets) and K columns (corresponding to items)
- N in the millions, K in tens of thousands
- Very sparse since typical basket contains few items

General Form of Association Rule

- Given a set of 0,1 valued variables A_1, \dots, A_K a rule would have the form

$$\left((A_{i_1} = 1) \wedge \dots \wedge (A_{i_k} = 1) \right) \Rightarrow A_{i_{k+1}} = 1$$

where $1 \leq i_j \leq K$ for all $j=1, \dots, k$

Subscripts allow for any combination of variables in rule

- Can be written more briefly as

$$(A_{i_1} \wedge \dots \wedge A_{i_k}) \Rightarrow A_{i_{k+1}}$$

- Pattern such as $(A_{i_1} = 1) \wedge \dots \wedge (A_{i_k} = 1)$
 - Is known as an *itemset*

Frequency of Itemsets

- A rule is an expression of the form $\theta \rightarrow \phi$
 - where θ is an itemset pattern
 - and ϕ is an itemset pattern consisting of a single conjunct
- Frequency of itemset
 - Given an itemset pattern θ
 - its frequency $fr(\theta)$ is the number of cases in the data that satisfy θ
- Frequency $fr(\theta \wedge \phi)$ is the support
- Accuracy of the rule

$$c(\theta \Rightarrow \varphi) = \frac{fr(\theta \wedge \varphi)}{fr(\theta)}$$

- Conditional probability that ϕ is true given that θ is true
- Frequent Sets
 - Given a frequency threshold s , all itemset patterns that are frequent

Example of Frequent Itemsets

Basket\Item	A_1	A_2	A_3	A_4	A_5
t_1	1	0	0	0	0
t_2	1	1	1	1	0
t_3	1	0	1	0	1
t_4	0	0	1	0	0
t_4	0	1	1	1	0
t_5	1	1	1	0	0
t_6	1	0	1	1	0

t_7	1	0	1	1	0
t_8	0	1	1	0	0
t_9	1	0	0	1	0
t_{10}	0	1	1	0	1

- Frequent sets for threshold 0.4 are:
 - $\{A_1\}, \{A_2\}, \{A_3\}, \{A_4\},$
 $\{A_1A_3\}, \{A_2A_3\}$
- Rule $A_1 \rightarrow A_3$ has accuracy $4/6=2/3$
- Rule $A_2 \rightarrow A_3$ has accuracy $5/5=1$

Association Rule Algorithm tuple

1. **Task** = description: associations between variables
2. **Structure** = probabilistic “association rules” (patterns)
3. **Score Function** = Threshold on accuracy and support
4. **Search Method** = Systematic search (breadth first with pruning)
5. **Data Management Technique** = multiple linear scans

Score Function

If $A=1$ and $B=1$ then $C=1$ with probability $p = p(C=1|A=1, B=1)$
 $p(A=1, B=1, C=1)$ is the *support*

1. Score function is a binary function (defined in 2)
Two thresholds:
 - p_s is a lower bound on the support for the rule
e.g., $p_s = 0.1$ want only rules that cover at least 10% of the data
 - p_a is a lower bound on the accuracy of the rule
e.g., $p_a = 0.9$ want only rules that are 90% accurate
2. A pattern gets a score of 1 if it satisfies both threshold conditions and a score of 0 otherwise
3. Goal is to find *all* rules (patterns) with a score of 1

Search Problem

- Searching for all rules is formidable problem
- Exponential number of association rules
 - $O(K2^{K-1})$ for binary variables if we limit ourselves to rules with positive propositions (e.g., $A=1$) in left- and right- hand sides
- Taking advantage of nature of score function can reduce run-time

Reducing Average Search Run-Time

If $A=1$ and $B=1$ then $C=1$ with probability $p = p(C=1|A=1, B=1)$
 $p(A=1, B=1, C=1)$ is the *support*

- **Observation:** If either $p(A=1) \leq p_s$ or $p(B=1) \leq p_s$ then $p(A=1, B=1) \leq p_s$
- First find all events (such as $A=1$) that have probability greater than p_s . This is a *frequent set*.
- Consider all possible pairs of these frequent events to be candidate frequent sets of size 2

Frequent Sets

- Going from frequent sets of size $k-1$ to frequent sets of size k , we can
 - prune any sets of size k that contain a subset of $k-1$ items that are not frequent
- E.g.,
 - If we had frequent sets $\{A=1, B=1\}$ and $\{B=1, C=1\}$ they can be combined to get $k=3$ set $\{A=1, B=1, C=1\}$
 - However, if $\{A=1, B=1\}$ is not frequent then $\{A=1, B=1, C=1\}$ is not frequent either and it could be safely pruned
- Pruning can take place without searching the data directly
- This is the “a priori” property

A priori Algorithm Operation

- Given a pruned list of candidate frequent sets of size k
 - Algorithm performs another linear scan of the database to determine which of these sets are in fact frequent
- Confirmed frequent sets of size k are combined to generate possible frequent sets containing $k+1$ events followed by another pruning etc
 - Cardinality of largest frequent set is quite small (relative to n) for large support values
- Algorithm makes one last pass through data set to determine which subset combination of frequent sets also satisfy the accuracy threshold

Summary: Association Rule Algorithms

- Search and Data Management are most critical components
- Use a systematic breadth-first general-to-specific search method that tries to minimize number of linear scans through the database
- Unlike machine learning algorithms for rule-based representations, they are designed to operate on very large data sets relatively efficiently
- Papers tend to emphasize computational efficiency rather than interpretation of the rules produced

Vector Space Algorithms for Text Retrieval

- Retrieval by content
- Query object and a large database of objects
- Find k objects in database that are similar to query

Text Retrieval Algorithm

- How is similarity defined?
- Text documents are of different length and structure
- Key idea:
 - Reduce all documents to a uniform vector representation as follows:
 - Let t_1, \dots, t_p be p terms (words, phrases, etc)
 - These are the variables or columns in data matrix

Vector Space Representation of Documents

- A document (a row in data matrix) is represented by a vector of length p
- Where the i^{th} component contains the count of how often term t_i appears in the document
- In practice, can have a very large data matrix
 - n in millions, p in tens of thousands
 - Sparse matrix
 - Instead of a very large $n \times p$ matrix, store a list for each term t_i of all documents containing the term

Similarity of Documents

- Similarity distance is a function of the angle between two vectors in p -space
- Angle measures similarity in *term space* and factors out any differences arising from fact that large documents have many occurrences of a word than small documents
- Works well -- many variations on this theme

Text Retrieval Algorithm tuple

1. Task = retrieval of k most similar documents in a database relative to a given query
2. Representation = vector of term occurrences
3. Score function = angle between two vectors
4. Search method = various techniques
5. Data Management Technique = various fast indexing strategies

Variations of TR Components

- In defining score function, we can specify similarity metrics more general than angle function
- In specifying search method, various heuristic techniques possible
 - Real time search since algorithm has to retrieve patterns in real time for a user (unlike other data mining algorithms meant for off-line searching for optimal parameters and model structures)

Text Retrieval Variations

- In searching legal documents, absence of particular terms might be significant – reflect this in score function
- Another context, down-weight the fact that certain terms are missing in two documents relative what they have in common