



CALIFORNIA STATE UNIVERSITY  
**FULLERTON**™

---

## Applying the Association Rules Mining Technique to Identify Critical Graduation Pathway Courses

2014 CAIR Conference

*Afshin Karimi (akarimi@fullerton.edu)*  
*Ed Sullivan, Ph.D. (esullivan@calstate.edu)*

*November 19, 2014*

# Data Mining Techniques Used in Higher Education

- Prediction (and/or Classification)
- Clustering
- Relationship Mining

# Relationship Mining

- Goal is to discover relationships between variables with data set with large number of variables
- 4 types of Relationship Mining:
  - *Association Rules Mining*
  - *Sequential Pattern Mining*
  - *Correlation Data Mining*
  - *Causal Data Mining*

# Association Rules Mining

- Proposed by Agrawal et al in 1993
- If-then rules amongst variables
- Initially used for Market Basket Analysis
- Milk Purchase -> Cereal Purchase (5% support, 80% confidence)
  - 5% support: customers who buy both product (in any order) are 5% of all customers in the database
  - 80% confidence: 80% of those who buy milk also buy cereal
- If student takes courses A and B, she will take course C (not necessarily in that order)

# Association Rules Mining Examples

- Walmart study found young males buying beer on Friday afternoons also buy baby diapers
- Amazon recommending items based on your current browsing/buying selections as well as other customers' purchasing patterns
- Google search's auto-complete where after a word is typed in the search box, it suggests a follow-up associated search term

# The Apriori Algorithm

- The best known algorithm for Association Rules Mining
- The algorithm is a two step process:
  - Find frequent itemsets
  - Use frequent itemsets to generate rules

# Apriori algorithm, continued...

## **Step 1:** Finding frequent itemsets:

Iterative process starting with scanning the database to find frequent 1-itemsets (that meet min. support), then using a Join operation find larger frequent itemsets (through k-itemset)

## **Step 2:** Generating association rules:

Using the found frequent itemsets and minimum support and confidence, rules are established

## An example

- Transaction data
- Assume:

minsup = 30%

minconf = 80%

- An example **frequent itemset**:

{Chicken, Clothes, Milk} [sup = 3/7]

- Association rules** from the itemset:

Clothes → Milk, Chicken [sup = 3/7, conf = 3/3]

...

...

Clothes, Chicken → Milk, [sup = 3/7, conf = 3/3]

t1: Beef, Chicken, Milk  
t2: Beef, Cheese  
t3: Cheese, Boots  
t4: Beef, Chicken, Cheese  
t5: Beef, Chicken, Clothes, Cheese, Milk  
t6: Chicken, Clothes, Milk  
t7: Chicken, Milk, Clothes



# Input Data (Association Rules Mining)

Customer ID	Beer	Wine	Soda	Cheese	Soap	Apples	Ground Beef	Chips	Pasta Sauce	Gum	Wall Calendar	Ground Coffee	Postcard	Magazine	Mints	Oranges
1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	0
2	1	0	1	1	0	1	0	1	0	0	1	0	0	0	0	0
3	0	1	0	0	1	0	0	0	0	0	0	0	1	0	1	0
4	0	0	1	0	0	0	1	0	1	0	0	0	0	0	0	1
5	0	0	0	0	1	0	0	1	0	0	0	1	0	1	0	0
6	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
7	1	0	0	1	0	1	0	0	0	1	1	0	0	0	1	0
8	1	1	0	0	0	0	0	1	0	0	0	0	1	0	0	1
9	0	1	1	0	0	0	1	0	1	0	0	0	0	1	0	0
10	0	0	1	0	1	1	0	0	0	0	0	1	1	0	0	0
11	0	0	0	1	0	0	1	0	0	0	1	0	0	0	0	0
12	1	1	1	0	0	0	0	0	0	1	0	1	0	0	0	0
13	0	0	0	1	0	0	1	1	0	0	1	0	0	0	1	0
14	0	0	0	0	1	0	0	0	1	0	0	0	0	0	1	1
15	1	0	0	0	0	1	0	0	0	0	0	0	1	1	0	1
16	0	1	1	0	1	0	0	0	0	1	0	0	0	1	0	0
17	0	0	1	0	0	0	0	0	0	0	0	1	0	0	0	0
18	0	0	0	1	0	0	1	1	0	0	1	0	0	0	1	0
19	1	0	1	0	0	0	0	0	1	0	0	0	0	0	0	1
20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
21	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
22	0	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0
23	1	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
24	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
25	0	0	1	0	1	0	1	0	1	0	1	0	1	0	1	0
26	1	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1

# Problem with Association Rules Mining

- **Problem:** Algorithm discovers huge number of association rules (between one or more variables with one or more other variables), many of which are irrelevant
- **A Solution:** Use ‘interestingness’ measures to reduce the rule set

# Interestingness

- Objective Interestingness:
  - Support
  - Confidence
  - Cosine
  - Added value
  - Lift
- Subjective Interestingness:
  - Unexpectedness
  - Actionability

# Support

Let  $|X, Y|$  be the number of transactions that contain both X and Y

Support is the proportion of all transactions that contain both X and Y

$$\text{Sup}(X \rightarrow Y) = |X, Y| / n \quad \text{OR} \quad P(X, Y)$$

$$\text{Sup}(X \rightarrow Y) = \text{Sup}(Y \rightarrow X)$$

# Confidence

Let  $|X|$  the number of transactions that contain  $X$ .

Confidence is the proportion of transactions that contain  $Y$  amongst the ones that contain  $X$ .

$$\text{conf}(X \rightarrow Y) = |X, Y| / |X| \quad \text{OR} \quad P(X, Y) / P(X)$$

$$\text{conf}(X \rightarrow Y) \neq \text{conf}(Y \rightarrow X)$$

# Cosine

(borrowing from cosine of angle between two vectors...)

$$\text{Cosine } (X \rightarrow Y) = |X, Y| / \sqrt{|X| \cdot |Y|}$$

- The closer cosine (X → Y) is to 1, the more transactions containing item X also contain Y
- The closer cosine (X → Y) is to 0, the more transactions contain item X without containing Y
- Cosine is a symmetric measure:  $\text{cosine}(X \rightarrow Y) = \text{cosine}(Y \rightarrow X)$

# Lift

$$\text{lift}(X \rightarrow Y) = \text{conf}(X \rightarrow Y) / P(Y)$$

If  $P(X, Y) = P(X) \cdot P(Y)$ , lift is 1. This is the worst case (occurrence of X and occurrence of Y in the same transactions are independent events)

# Subjective Interestingness

Subjective Interestingness is application domain- specific. Two such measures are:

- **Unexpectedness:** Grocery chain already knows about (Beer -> Chips) association rule, but not about the (Beer -> Diapers) association rule.
- **Actionability:** Rules that offer strategic information on which user can act on.



# Association Rules Example

- Transfer Student Success Project in the Mihaylo College of Business & Economics
- Identify the gateway courses that prevent MCBE transfer students from timely graduation

# Association Rules Example Continued...

## MCBE Transfer Students Success:

- Examine CBE courses that new transfer students take AND fail during 1<sup>st</sup> term at Fullerton
- Find all Association Rules between all the variables (course failures) and a new variable that represents graduation in 4 years or less
- Use interestingness measures to focus on the relevant associations

# Association Rules Example Continued...

## Input File Format

- Rows: fall 08 & 09 new transfer MCBE students who took at least one MCBE course during their 1<sup>st</sup> term (1,807 students)
- Columns: MCBE courses above students took during their 1<sup>st</sup> term PLUS Graduation variable that indicates if student graduated in 4 years or less (43 columns)
- Values:
  - 1: failed the course in 1<sup>st</sup> term (grade of C- thru F, including WU)
  - 0: passed the course in 1<sup>st</sup> term (grade of C or above) *OR* didn't take course in 1<sup>st</sup> term

# Example Input File

cwid	NotGraduated	ACCT201A	ACCT201B	ACCT301A	ACCT302	ACCT307	ACCT364	BUAD201	BUAD210	BUAD301	ECON201	ECON202	ECON310	ECON315	ECON320	ECO
8906	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0797	1	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
8860	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
9249	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3102	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1375	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1853	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
7215	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0532	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
6824	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6247	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6475	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
8455	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6961	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
9981	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4777	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
7802	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5741	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6024	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
6881	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
046	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
174	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
874	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
963	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
592	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5482	1	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0
6836	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
896	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
987	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
473	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
347	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
982	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5975	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2369	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
5081	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
7762	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
135	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0



# Association Rules Example Continued...

- Algorithm finds large number of rules between one or more variables with one or more (other) variables
- Here we focus on association rules between different course variables and graduation variable:  
(X -> Grad in 4 Yrs) where X is any of the 42 CBE courses.
- Furthermore, narrow the list by using Support & Confidence measures

# RapidMiner 5 Software Demo

The screenshot displays the RapidMiner 5 software interface. The main workspace shows a workflow diagram with the following components:

- Retrieve**: A purple box with a database icon and an 'out' port.
- Numerical to ...**: A pink box with a 3D bar chart icon. It has 'exa' and 'ori' ports on the left and 'exa' and 'ori' ports on the right.
- FP-Growth**: A green box with a shopping cart icon. It has 'exa' and 'fre' ports on the left and 'exa' and 'fre' ports on the right.
- Create Associ...**: A green box with a shopping cart icon. It has 'ite' and 'rul' ports on the left and 'ite' and 'rul' ports on the right.
- Read Excel**: A grey box with a document icon, currently highlighted with an orange border. It has 'fil' and 'out' ports.

The workflow is connected as follows: Retrieve connects to Numerical to ...; Numerical to ... connects to FP-Growth; FP-Growth connects to Create Associ...; and Read Excel is positioned below the main flow.

At the bottom, a log window shows the following messages:

- 2014 9:05:04 AM CONFIG: Loading perspectives.
- 2014 9:05:04 AM CONFIG: Ignoring update check. Last update check was on 3/11/14 2:26 PM.
- 2014 9:05:47 AM INFO: Reading example set...
- 2014 9:06:21 AM INFO: Reading example set...
- 2014 9:06:34 AM INFO: No filename given for result file, using stdout for logging results!
- 2014 9:06:34 AM INFO: Process //Local Repository/processes/AssociationRules1 starts
- 2014 9:06:34 AM INFO: Loading initial data.
- 2014 9:06:34 AM INFO: Saving results.
- 2014 9:06:34 AM INFO: Process //Local Repository/processes/AssociationRules1 finished successfully after 0 s
- 2014 9:07:29 AM INFO: No filename given for result file, using stdout for logging results!

# Association Rules Example Continued...

Results:

- Top 3 identified gateway courses are all 200 level courses (lower division core courses) that new transfer students take AND fail
- Graduation variable not really the 'target' variable

# Future Work/Summary

- Further study of the identified gateway courses
- If order of events is important, use Sequential Mining method instead (not covered in this presentation)
- No need to have intimate knowledge of the algorithm used. Just need to compile model's input data file



# Questions/Comments?

Contact: [akarimi@fullerton.edu](mailto:akarimi@fullerton.edu)