

# 09 - Association Analysis

Frequent Itemsets, Market Basket Analysis, and Association Rule Mining

SYS 6018 | Fall 2020

09-association.pdf

## Contents

<b>1</b>	<b>Reading</b>	<b>3</b>
<b>2</b>	<b>Case: Instacart</b>	<b>3</b>
<b>3</b>	<b>Association Analysis Motivation</b>	<b>3</b>
3.1	Market-Basket Analysis . . . . .	3
3.2	Other Applications . . . . .	3
<b>4</b>	<b>Market Basket Model</b>	<b>6</b>
4.1	Data . . . . .	6
4.2	Definitions . . . . .	6
4.3	Practice . . . . .	7
<b>5</b>	<b>Apriori Algorithm</b>	<b>8</b>
5.1	Introduction . . . . .	8
5.2	Apriori Approach . . . . .	8
5.3	Targeted Items . . . . .	10
<b>6</b>	<b>Example: Instacart</b>	<b>11</b>
6.1	Instacart Case . . . . .	11
6.2	Load Data . . . . .	11
6.3	Clean Data . . . . .	12
6.4	arules package . . . . .	12
6.5	Find Frequent Itemsets . . . . .	12
6.6	Find Association Rules . . . . .	12
6.7	Plot Association Rules . . . . .	13
6.8	Target specific itemsets . . . . .	13
<b>7</b>	<b>More Interesting Rules</b>	<b>14</b>
7.1	Measures of Interestingness . . . . .	14
7.2	Lift . . . . .	14
7.3	Other Measures . . . . .	16
7.4	Adding additional interesting measures in arules package . . . . .	17

<b>8</b>	<b>Statistical Issues</b>	<b>18</b>
8.1	Repeatable Patterns . . . . .	18
8.2	Probability Estimation and Statistical Significance . . . . .	18
<b>9</b>	<b>Extensions</b>	<b>21</b>
9.1	Non-Market Basket Format . . . . .	21
9.2	Supervised Learning Problem . . . . .	22
9.3	Research Opportunities . . . . .	23

---

## 1 Reading

- MMDS 6.1 and 6.2
- ITDM 5.1-5.3; 5.7-5.8
- ESL 14.2
- R package `arules`
- R package `arulesViz`

## 2 Case: Instacart

[Instacart Case](#)

## 3 Association Analysis Motivation

### 3.1 Market-Basket Analysis

- A grocery store records the items purchased in a set of transactions. For example,

TID	Items
1	{Bread, Milk}
2	{Bread, Diapers, Beer, Eggs}
3	{Milk, Diapers, Beer, Cola}
4	{Bread, Milk, Diapers, Beer}
5	{Bread, Milk, Diapers, Cola}

- If they can discover items **frequently** purchased together, they may be able to increase business.
  - If {Diapers, Beer} are frequently purchased together, the store could put diapers on sale and increase the price of beer.
  - If {Bread, Peanut Butter} are frequently purchased together, they may put some Peanut Butter on the shelf next to bread.
  - In on-line retail, you will often see “People who buy {X,Y} also tend to buy {Z}”.
- **Association Rules** are rules that describe potential relationships *of interest* between itemsets.
  - E.g., Given a set of transactions, find *rules that will predict* the occurrence of an item based on the occurrences of other items in the transaction
  - For example, the rule {Diapers}  $\rightarrow$  {Beer} suggests that customers who buy Diapers tend to purchase Beer more than *expected*.
  - Of course, we need to discuss metrics that will permit us to quantitatively evaluate the discovered rules.
  - **Co-occurrence doesn't necessitate causality!**

### 3.2 Other Applications

- Market Basket
  - Items: Products for sale
  - Baskets: Sets of products purchased in single trip
  - Needs rule to occur frequently to make worthwhile

TID	Items
1	{ Bread, Milk }
2	{ Bread, Diapers, Beer, Eggs }
3	{ Milk, Diapers, Beer, Cola }
4	{ Bread, Milk, Diapers, Beer }
5	{ Bread, Milk, Diapers, Cola }

- Text Mining (1)
  - Items: Words
  - Baskets: Books/Articles/Webpages/News
  - Frequent itemsets indicate words that are connected

News Article	Words
1: NY Times	{ Trump, Tweet, Fake }
2: Wash Post	{ Collusion, Russian, Trump }

- Text Mining (2)
  - Items: Documents
  - Baskets: Sentences/Code
  - Frequent itemsets could indicate plagiarism

Sentences/Code	Items
1: Logistic Regression is ...	{ Doc 1, Doc 4, Doc 5 }
2: ggplot(aes(x=size, y =weight))+ geom_...	{ Doc 1, Doc 5 }

- Health Mining
  - Items: Drugs and Side-Effects
  - Baskets: Patients
  - Discover combinations of drugs that result in particular side-effects
    - \* { Drug 1, Drug 4, Drug 6 } → { Psychosis }
  - Note that absence of drug/side-effect could also be important
    - \* { Drug 1, Drug 4, Not Drug 7 } → { Psychosis }

Patient	Items
1	{ Drug 1, Drug 3, Rash, Weight-Loss }
2	{ Drug 2, Drug 3, Headaches }

- Transportation
  - Items: Features collected after vehicle crash
  - Baskets: The set of vehicle crashes
  - Discover dangerous/safe combinations

CrashID	Items
1	{Wet, Night, Rural}
2	{Alcohol, Night, Lights-off}
3	{Day, Flat tire, Elderly}

- Survey Data

- Items: Potential Responses to questionnaire
- Baskets: Set of responses by respondent
- {Married, Own house}  $\rightarrow$  {income  $\in$  [60K, 100K]}

ID	Items
1	{Sex = <i>Male</i> , Age = [20,25], Type of home = <i>House</i> , ... }
2	{Sex = <i>Female</i> , Age = [30,35], Type of home = <i>Apt</i> , ... }
3	{Sex = <i>Male</i> , Age = [65+], Type of home = <i>House</i> , ... }

## 4 Market Basket Model

### 4.1 Data

There are two equivalent data formats:

- **Transaction Lists**

TID	Items
1	{Bread, Milk}
2	{Bread, Diapers, Beer, Eggs}
3	{Milk, Diapers, Beer, Cola}
4	{Bread, Milk, Diapers, Beer}
5	{Bread, Milk, Diapers, Cola}

- **Binary Incident Table/Matrix**

TID	Bread	Milk	Diapers	Beer	Eggs	Cola
1	1	1	-	-	-	-
2	1	-	1	1	1	-
3	-	1	1	1	-	1
4	1	1	1	1	-	-
5	1	1	1	-	-	1

### 4.2 Definitions

- $K$  possible items
- $N$  total transactions
- Let  $I$  be an itemset
  - E.g.,  $I = \{\text{Cola, Bread, Fruit}\}$ ,  $I = \{\text{Bread}\}$
  - There are  $2^K - 1$  possible itemsets
- The **support** of  $I$  is the fractions of transactions that contain  $I$ 
  - $S(I) = (\# \text{ of transactions containing } I) / N$
- The support is an estimate of the probability of getting  $I$  in a future basket.
  - $S(I) \hat{=} \Pr(I \subseteq T)$ , where  $T$  is a future basket
- A **Frequent Itemset** is one that has support greater than or equal to the **support threshold (minsup)**,  $s$ .
- **Association Rule** is an if-then expression about the content of the baskets
  - The expression  $I \rightarrow J$  means that if a basket contains  $I$ , then it is *more likely* to also contain  $J$ .
  - $I$  and  $J$  should be *disjoint*, e.g.,  $I \cap J = \emptyset$
- The **support** of an association rule is the frequency in which both itemsets are in the same basket

- $S(I \rightarrow J) = (\# \text{ of transactions containing } I \text{ and } J)/N$
  - $S(I \rightarrow J) = S(I, J) \triangleq \Pr(I \subseteq T, J \subseteq T)$ , where  $T$  is a future basket.
  - The **confidence** of a rule is an estimate of the conditional probability of  $J$  being in a basket given the basket already contains  $I$ 
    - $C(I \rightarrow J) = \frac{S(I, J)}{S(I)} \triangleq \frac{\Pr(I \subseteq T, J \subseteq T)}{\Pr(I \subseteq T)} = P(J \subseteq T | I \subseteq T)$
  - A **High Confidence** rule is one that has confidence greater than or equal to the **confidence threshold** (**minconf**),  $c$ .
  - The classic **Association Rule Discovery** task is to find all rules that have *support*  $\geq s$  and *confidence*  $\geq c$
1. Find all *frequent* itemsets,  $I$ 
    2. For every subset  $A \subset I$ , generate the rule  $A \rightarrow I \setminus A$
    3. It is easy to find the *largest* confidence set derived from  $I$ . Do you see how?

### 4.3 Practice

#### Your Turn #1

TID	Items
0001	{a, d, e}
0024	{a, b, c, e}
0012	{a, b, d, e}
0031	{a, c, d, e}
0015	{b, c, e}
0022	{b, d, e}
0029	{c, d}
0040	{a, b, c}
0033	{a, d, e}
0038	{a, b, e}

1. What are the value of  $K$  and  $N$ ?
2. What is the *support* of  $\{e\}$ ,  $\{b, d\}$ , and  $\{b, d, e\}$ ?
3. What is the *confidence* of  $\{b, d\} \rightarrow \{e\}$  and  $\{e\} \rightarrow \{b, d\}$ ? Is confidence symmetric?
4. Think of an approach/algorithm to find *all Association Rules* with  $s = 0.4$  and  $c \geq .70$ .

## 5 Apriori Algorithm

### 5.1 Introduction

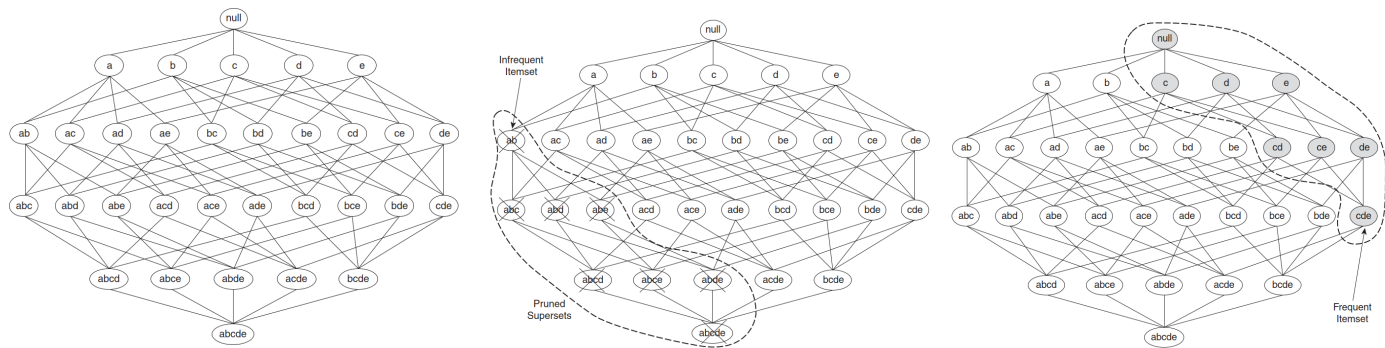
- **Goal:** Discover associations between itemsets.
  - Interested in discovering *frequent itemsets* - those that appear in  $\geq s$  baskets.
  - Interested in discovering *high confidence rules* - those with confidence  $\geq c$  scores.
- **Problem:** There are  $2^K - 1$  possible itemsets and  $3^K - 2^{K+1} + 1$  possible rules.
  - Wal-mart has 100K items
  - 10B webpages
  - Memory and computing time issues
- **Solution:** The **Apriori** principle
  - *If an itemset is frequent, then all of its subsets must also be frequent.*
  - Equivalently, *If an itemset is not frequent, then none of its unions will be frequent.*
- Thus, if item  $i$  does not appear in at least  $s$  baskets (i.e., not frequent), then no itemsets that include item  $i$  will appear in at least  $s$  baskets.

### 5.2 Apriori Approach

- Because of the size of association data, traditional approaches have focused on finding the rules that have support  $\geq s$  and confidence  $\geq c$ .
- This can be efficiently approached in a two step process
  1. Find all frequent itemsets
  2. Find rules, using the frequent itemsets, with high confidence



### 5.2.1 Finding Frequent Itemsets

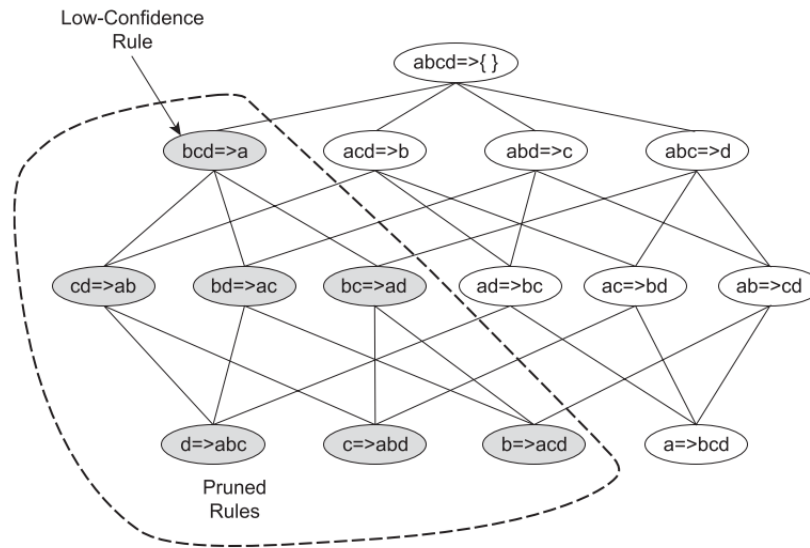


### 5.2.2 Finding High Confidence Rules

#### Your Turn #2

A frequent itemset consists of  $K = 4$  items {a,b,c,d}.

1. How many association rules (of any size) can be formed from data with  $K = 4$  items?
2. How many association rules can be formed from *this itemset* (exactly  $K = 4$  items)?
3. Suppose the *confidence* for rule  $\{b,c,d\} \rightarrow \{a\}$  is low (i.e.,  $C(\{b, c, d\} \rightarrow \{a\}) < c$ ).
  - a. Do we need to calculate  $C(\{d\} \rightarrow \{a, b, c\})$ ?
  - b. Do we need to calculate  $C(\{a, b\} \rightarrow \{c, d\})$ ?



### 5.3 Targeted Items

- Suppose we are especially interested in one item (e.g., Honey).
- If we want to find an association itemset that *increases the probability* that targeted item is in a basket (i.e., item of interest on the rhs), then it may help to convert the problem into a *supervised learning* classification problem.
  - Same “association is not causation” caveat applies
- If we want to see how the target item associates with other items (i.e., target item on lhs), then we don’t need to run a full search, but only build up the rhs items sequentially.

## 6 Example: Instacart

### 6.1 Instacart Case

See the [Instacart Case](#) for details.

### 6.2 Load Data

The R Code *instacart.R* (see course webpage) contains a detailed example.

1. The data have already been downloaded and extracted to the course website. The two relevant datasets can be load into R using `readr::read_csv()`

```
1 library(readr)
2 #-- Load the "orders_products__train.csv" (notice two __) and "products.csv"
3 data.dir = "https://mdporter.github.io/SYS6018/data/"
4 orders = read_csv(file.path(data.dir, "order_products__train.csv"))
5 products = read_csv(file.path(data.dir, "products.csv"))

#> # A tibble: 1,384,617 x 4
#>   order_id product_id add_to_cart_order reordered
#>   <dbl>      <dbl>         <dbl>      <dbl>
#> 1         1         49302             1         1
#> 2         1         11109             2         1
#> 3         1         10246             3         0
#> 4         1         49683             4         0
#> 5         1         43633             5         1
#> 6         1         13176             6         0
#> 7         1         47209             7         0
#> 8         1         22035             8         1
#> 9         36         39612             1         0
#> 10        36         19660             2         1
#> # ... with 1,384,607 more rows
#> # A tibble: 49,688 x 4
#>   product_id product_name      aisle_id department_id
#>   <dbl> <chr>          <dbl>      <dbl>
#> 1         1 Chocolate Sandwich Cookies      61         19
#> 2         2 All-Seasons Salt             104         13
#> 3         3 Robust Golden Unsweetened Oolong Tea      94          7
#> 4         4 Smart Ones Classic Favorites Mini Rigatoni~      38          1
#> 5         5 Green Chile Anytime Sauce                5         13
#> 6         6 Dry Nose Oil                     11         11
#> 7         7 Pure Coconut Water With Orange           98          7
#> 8         8 Cut Russet Potatoes Steam N' Mash        116          1
#> 9         9 Light Strawberry Blueberry Yogurt        120         16
#> 10        10 Sparkling Orange Juice & Prickly Pear Beve~  115          7
#> # ... with 49,678 more rows
```

2. Join orders and products to get data frame of *order\_id*, *product\_id*, and *product\_name*.
  - Here using `dplyr` functions.

```
1 library(dplyr)
2 Y = left_join(orders %>% select(order_id, product_id),
3               products %>% select(product_id, product_name),
4               by = "product_id") %>%
5   arrange(order_id)

#> # A tibble: 1,384,617 x 3
```

```
#>   order_id product_id product_name
#>   <dbl>      <dbl> <chr>
#> 1         1         49302 Bulgarian Yogurt
#> 2         1         11109 Organic 4% Milk Fat Whole Milk Cottage Cheese
#> 3         1         10246 Organic Celery Hearts
#> 4         1         49683 Cucumber Kirby
#> 5         1         43633 Lightly Smoked Sardines in Olive Oil
#> 6         1         13176 Bag of Organic Bananas
#> 7         1         47209 Organic Hass Avocado
#> 8         1         22035 Organic Whole String Cheese
#> 9         36         39612 Grated Pecorino Romano Cheese
#> 10        36         19660 Spring Water
#> # ... with 1,384,607 more rows
```

### 6.3 Clean Data

- The data should have a one-to-one mapping between *product\_name* and *product\_id*
- Transactions should not contain duplicate items
- May need to clean item names or ids

### 6.4 arules package

- The R package *arules* has functions for performing association analysis
  - More details in the vignette <https://cran.r-project.org/web/packages/arules/vignettes/arules.pdf>
- The main function is `apriori()`
  - the arguments *parameter*, *appearance*, and *control* are used to run certain analyses
- The data needs to be converted to *transaction* class (essentially a sparse binary matrix)
  - This is done by using the function `as(tList, "transactions")`, where *tList* is a **transaction list**
  - This means we need to convert our transaction data frame into a transaction *list*

### 6.5 Find Frequent Itemsets

- Frequent itemsets can be found by running the function

```
1 apriori(trans, parameter = list(support = .01, target="frequent"))
```

- *trans* is the transactions object
- support threshold is set to .01
- target="frequent" indicates we only want frequency, not rules
- The `parameter=` argument can be used focus on certain classes of itemsets
  - E.g., support, minlen, maxlen

### 6.6 Find Association Rules

- Association rules can be found by running the function

```
1 apriori(trans, parameter = list(support = .001, confidence=.50,
2                               target="rules"))
```

- *trans* is the transactions object

- support threshold is set to  $s = .01$
- confidence threshold is set to  $c = .50$
- target="rules" indicates we only want rules

## 6.7 Plot Association Rules

- the `arulesViz` R package provides some plotting functionality
  - More details in the vignette <https://cran.r-project.org/web/packages/arulesViz/vignettes/arulesViz.pdf>

## 6.8 Target specific itemsets

- Specific itemsets can be targeted by setting the `appearance =` parameter

## 7 More Interesting Rules

### 7.1 Measures of Interestingness

The utility of a particular association rule is based on the needs of the business/research (e.g., sell more products, increase profit). When in a purely exploratory mode, the **interestingness** of a rule determines its value. So far, we have only examined *support* and *confidence*, but there are many other measures of interestingness.

#### 7.1.1 Think About It

##### Your Turn #3

For each of the following questions, provide an example of an association rule from the market basket domain that satisfies the following conditions. Also, describe whether such rules are subjectively interesting.

1. A rule that has high support and high confidence.
2. A rule that has reasonably high support but low confidence.
3. A rule that has low support and low confidence.
4. A rule that has low support and high confidence.

### 7.2 Lift

- Why would the rule Vodka → Caviar be interesting, even if it has low support?
- In general, an interesting rule is one that occurs more (or less) than *expected*.
- Support and Confidence are based on an expected **uniform** distribution.
- **Lift**, on the other hand is based on statistical independence.
  - Notice that lift is **not** directional

$$L(I, J) = \frac{C(I \rightarrow J)}{S(J)}$$

$$= L(J, I)$$

- If  $I$  and  $J$  are independent, then  $\Pr(I, J) = \Pr(I) \Pr(J)$ .
  - Lift  $> 1$  indicates positive association
  - Lift  $< 1$  indicates negative association ( $I$  and  $J$  *inhibit* each other)
- Because its a ratio, *lift* can be large even when the *support* is low and small when the *confidence* is large.

#### Your Turn #4

Consider the following beverage preferences from 1000 people:

	Coffee	No Coffee
Tea	150	50
No Tea	650	150

1. What is the *support* of  $\{\text{Tea}\} \rightarrow \{\text{Coffee}\}$ ?
2. What is the *confidence* of  $\{\text{Tea}\} \rightarrow \{\text{Coffee}\}$ ?
3. What is the *lift* of  $\{\text{Tea}\} \rightarrow \{\text{Coffee}\}$ ?
4. Interpret.

- What is the lift for more than two items (e.g.,  $L(I, J, L)$ )?

### 7.3 Other Measures

There are [many other measures](#). Here are a few:

- **Piatetsky-Shapiro (PS)/Leverage**: difference of observed and expected.

$$PS(I, J) = S(I, J) - S(I)S(J)$$

- **Correlation ( $\phi$ )**: normalized version of **PS**.

$$\phi(I, J) = \frac{PS(I, J)}{\sqrt{S(I)(1 - S(I))S(J)(1 - S(J))}}$$

- **Added Value (AV)**: difference between conditional and unconditional.

$$\begin{aligned} AV(I \rightarrow J) &= C(I \rightarrow J) - S(J) \\ &= PS(I, J)/S(I) \\ &\cong \Pr(J|I) - \Pr(J) \end{aligned}$$

- **Conditional Entropy  $H(J | I)$** : amount of info needed to describe distribution of  $J$  given  $I$ .

$$H(J|I) = H(I, J) - H(I)$$

- **Mutual Information (MI)**: type of Kullback-Leibler divergence where independence is used for alternative distribution.

$$MI(I, J) = H(J) - H(J|I)$$

- Michael Hahsler, author of the `arules` package, details [many other interest measures](#)

#### Your Turn #5

Consider the following results from 1000 people:

	Honey	No Honey
Tea	100	100
No Tea	20	780

1. What is the *support* of  $\{\text{Tea}\} \rightarrow \{\text{Honey}\}$ ?
2. What is the *confidence* of  $\{\text{Tea}\} \rightarrow \{\text{Honey}\}$ ?
3. What is the *lift* of  $\{\text{Tea}\} \rightarrow \{\text{Honey}\}$ ?
4. What is the *Added Value (AV)* of  $\{\text{Tea}\} \rightarrow \{\text{Honey}\}$ ?
5. What is the *Piatetsky-Shapiro (PS)* score for  $\{\text{Tea}\} \rightarrow \{\text{Honey}\}$ ?
6. Interpret.



## 7.4 Adding additional interesting measures in **arules** package

The `arules` package provides the function `interestMeasure()` that will calculate a wide range of metrics.

- For example, *lift* can be calculated for all frequent itemsets

## 8 Statistical Issues

### 8.1 Repeatable Patterns

- **Co-occurrence does not mean causality**
- The patterns/associations of interest are those that are *repeatable*, i.e., the patterns will appear in future transactions.
- If future transactions will come from a different distribution than the observed data, proceed with caution!
- Rule metrics are constructed from *counts* e.g.,  $N(I)$ ,  $N(J)$ ,  $N(I, J)$ .
  - If the transactions are *independent*, then  $N(A) \sim \text{Bino}(p_A, N)$ .
    - \* What is a good estimate of  $p_A$ ?
    - \* What are the confidence intervals for  $p_A$ ?
  - What if the transactions are **not** independent (same person makes multiple transactions)?

#### Your Turn #6

Suppose we have  $N = 10\,000$  transactions and found two itemsets  $I$  and  $J$  with  $N(I, J) = 500$ ,  $N(I) = 1000$ ,  $N(J) = 4000$ .

1. What is a good point estimate of  $\Pr(I, J)$ ?
2. What are the confidence intervals for  $\Pr(I, J)$  (pick the confidence level)?
3. What is a good point estimate of  $P(J|I)$ ?
4. Confidence interval for  $P(J|I)$ ?
5. What should we do if the transactions are *not* independent? E.g., the  $N$  transactions are made by  $M = 100$  customers.

### 8.2 Probability Estimation and Statistical Significance

- **Sampling:** If metrics are based on counts/probability estimates, then sampling can help alleviate the computational burden while still discovering the associations

- Much in the same way as polling is used to estimate the percentage of people who will vote for a particular candidate.
- Note: Should you sample items or transactions, or both?
- **Multiple Testing:** Statistical significance and confidence intervals are approaches to help us avoid *spurious associations*. One challenge with association analysis is that we are searching through a large data set and are likely to find (empirically) strong associations that are due to chance alone.
  - We need to be especially careful about rules/associations that are based on small counts as they will have the most variation (especially metrics based on ratios).
- **Simpson's Paradox:** See ITDM 5.7.3

### 8.2.1 Simulation

Here is some simulated transaction data. All items are generated independently and all transactions are independent. Ten items (1-10) have probability of .25 of being included in any transaction, ten items (11-20) have probability of .10 of being included in any transaction, and the remaining items have probability of .01.

```

1  #-- Settings
2  K = 2000 # total number of items
3  N = 50000 # number of transactions
4  p = c(rep(.25, 10), rep(.10, 10), rep(.01, K-20)) # probabilities
5
6  #-- function to generate transactions
7  get_transaction <- function(p) {
8    K = length(p)
9    (1:K)[runif(K) <= p]
10 }
11
12 #-- Simulate N transactions
13 X = replicate(N, get_transaction(p))
14
15 #-- Run apriori
16 library(arules)
17 trans = as(X, "transactions")
18
19 apriori(trans, parameter=list(support=100/N, confidence=.20, minlen=2),
20        control=list(verbose=FALSE)) %>%
21  sort(by="lift") %>% head(n=5) %>% inspect()
22
23 #>      lhs      rhs support confidence coverage lift count
24 #> [1] {5,10,15} => {1} 0.00208 0.3549      0.00586 1.408 104
25 #> [2] {2,3,16}  => {8} 0.00204 0.3290      0.00620 1.323 102
26 #> [3] {1643}    => {2} 0.00322 0.3286      0.00980 1.321 161
27 #> [4] {3,7,16}  => {9} 0.00222 0.3265      0.00680 1.308 111
28 #> [5] {1,10,15} => {5} 0.00208 0.3220      0.00646 1.301 104

```

- There are no “real” patterns in these data, but we can find many apparent patterns
- Notice the relatively large values of *lift* (expected lift is 1)
- Notice from which groups the items come from (high, medium, low probabilities). Why do we get these patterns?
- **Statistical Testing Approaches:**
  - Due to multiple comparisons, the exact distribution of a test statistic is often elusive.
  - However, it may not be too difficult estimate the distribution of a test statistic using simulation.

- **Independence:** randomly permute each column (in the transaction matrix) independently. However, this may change the distribution of the transaction length (number of items in a transaction). A method called *swap randomization* can help.
- Bootstrap the transactions
- Cross-Validation / Hold-out analysis

## 9 Extensions

### 9.1 Non-Market Basket Format

- Association analysis can be performed on non-market basket data.
- It only requires a data set of *binary* features
  - Categorical Data is dummy coded
  - Numerical data is binned

#### 9.1.1 Categorical Data

- A *categorical variable* is one that can one of  $k$  possible values
  - *nomial* is unordered (e.g., *blue*, *green*, *brown*)
  - *ordinal* is ordered (e.g., *Great*, *Good*, *Average*, *Below*, *Poor*)
- Dummy coding (*one-hot-encoding*) is the way to convert a categorical variable to a set of *binary* variables
  - A  $k$  level variable will create  $k$  binary variables
  - This is different than the usual dummy coding for regression, where  $k - 1$  levels are generated
  - This full  $k$  level coding is referred to as *one-hot-encoding* in the machine learning literature

```

1  #-- categorical vector with k=3 levels
2  x.cat = c('a', 'b', 'a', 'c', 'b')
3
4  #-- Convert to matrix with 3 columns
5  model.matrix(~x.cat-1)
6  #>   x.cata x.catb x.catc
7  #> 1      1      0      0
8  #> 2      0      1      0
9  #> 3      1      0      0
10 #> 4      0      0      1
11 #> 5      0      1      0
12 #> attr(,"assign")
13 #> [1] 1 1 1
14 #> attr(,"contrasts")
15 #> attr(,"contrasts")$x.cat
16 #> [1] "contr.treatment"
```

- May need to be creative in making interesting levels.
  - *Eyes = blue* vs. *Eyes ≠ blue*

#### 9.1.2 Numeric Data

- Numeric data must be converted into binary format
- For example,  $age \leq 25 = TRUE$
- Numeric data could also be discretized or binned. This converts it into categorical data which is then dummy encoded.
  - E.g., *age 14-17*, *age 18-24*, ..., *age 65+*

#### 9.1.3 Categorical and Numeric Data in **arules** package

- The **arules** package has functionality to automatically handle categorical and numeric data

```

1 library(arules)
2
3 ##-- Load data
4 data("IncomeESL") # load income data from arules package
5 summary(IncomeESL) # summary of data
6
7 #>      income      sex      marital status      age
8 #> [0,10) :1745   male  :4075   married      :3334   14-17: 878
9 #> [50,75):1308   female:4918   cohabitation: 668   18-24:2129
10 #> [30,40):1110                                divorced    : 875   25-34:2249
11 #> [40,50): 969                                widowed     : 302   35-44:1615
12 #> 75+      : 884                                single      :3654   45-54: 922
13 #> [20,25): 813                                NA's        : 160   55-64: 640
14 #> (Other):2164                                65+       : 560
15
16 #>      education      occupation      years in bay area
17 #> grade <9          : 264   professional/managerial:2820   <1 : 270
18 #> grades 9-11       :1046   student              :1489   1-3 :1042
19 #> high school graduate:2041   clerical/service     :1062   4-6 : 686
20 #> college (1-3 years):3066   sales                 : 770   7-10: 900
21 #> college graduate  :1524   laborer               : 767   >10 :5182
22 #> graduate study    : 966   (Other)               :1949   NA's: 913
23 #> NA's              : 86    NA's                  : 136
24
25 #>      dual incomes  number in household  number of children
26 #> not married:5438   2      :2664          0      :5724
27 #> yes              :2211   3      :1670          1      :1506
28 #> no               :1344   1      :1620          2      :1148
29 #>                  :      4      :1526          3      : 412
30 #>                  :      5      : 686          4      : 117
31 #>                  (Other): 452          5      : 46
32 #>                  NA's   : 375          (Other): 40
33
34 #>      householder status      type of home      ethnic classification
35 #> own                        :3256   house      :5073   white      :5811
36 #> rent                       :3670   condominium: 655   hispanic:1231
37 #> live with parents/family:1827   apartment  :2373   black      : 910
38 #> NA's                       : 240   mobile Home: 151   asian      : 477
39 #>                            :      other      : 384   other      : 225
40 #>                            NA's      : 357   (Other)    : 271
41 #>                            NA's      : 68
42
43 #>      language in home
44 #> english:7794
45 #> spanish: 579
46 #> other   : 261
47 #> NA's    : 359
48
49 #>
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```

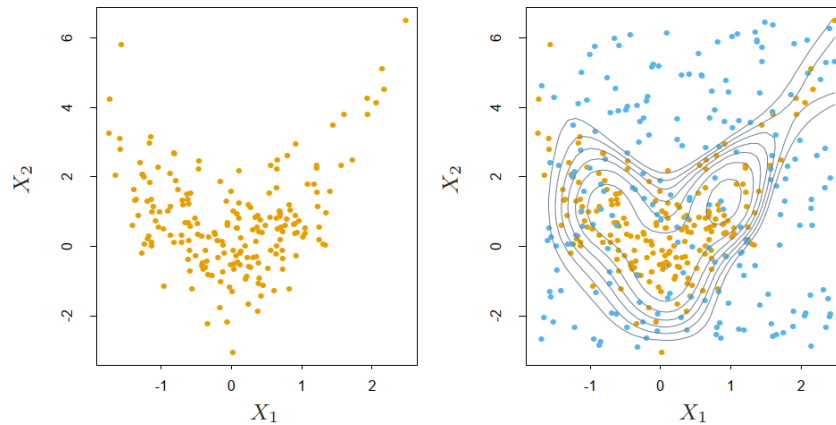
```

1
2 ##-- convert to transactions object
3 trans = as(IncomeESL, "transactions")

```

## 9.2 Supervised Learning Problem

- The ESL text (ESL 14.2.4 - 14.2.7) outlines an interesting approach for finding frequent itemsets
- The idea is to simulate additional data under an appropriate null distribution (e.g., independence, uniform), combine the simulated data with the observed data, and assign a *label* indicating if the observation comes from the observed or simulated data



- Now any supervised learning procedure can be used to find regions of deviation from the null
  - That is, regions where there are more actual observations than simulated observations
- This concept can be incorporated into other tasks, like density estimation and clustering

### 9.3 Research Opportunities

- How can “repeatable patterns” be confirmed? E.g., statistical significance.
- How to incorporate resampling (bootstrap, cross-validation)?
- How to handle structure in items?
  - e.g., *small oranges*, *large oranges*, *bag of oranges* are all oranges
- Detect temporal changes in transaction data
- Find outliers and anomalies in transaction data
- Transaction clustering
- Weighted transactions