Lect5: Tasks other than classification

Tokyo Tech. Intro. to Comp. & Data Lecture week5

- Intro. to the frequent item set mining.
 Intro. to the clustering.
- Frequent item set mining

 (and association rule mining as its application).
- 2. Clustering.
- 3. On Exercise #5.

* Some of the slide materials (in particular, green ones) are from the slides of the authors of the textbook and their group at the University of Waikato.

One of the earliest data mining examples investigated as a basic step for the association rule mining.



R. Agrawal, T. Imielinski, and A. N. Swami. Mining association rules between sets of items in large databases. in *Proc. <u>SIGMOD Conference</u>* 1993, pp. 207-216 (1993).

Famous database conference started from 1975!!

3

1. Frequent item set mining

glossaries

transaction database = a set of transactions. transaction = a set of (usually, a sequence of) *item*s recorded as one record in a database.



item set = a set of items in general.

frequency of an item set I = # of transactions that contain I. \uparrow also called support

frequent item set = an item set whose frequency is greater than or equal to a specified min. freq. θ .

a task of enumerating *all* frequent item sets from a given database \mathcal{D} and a min. freq. parameter θ .

Example: Weath	er database \rightarrow	Outlook	c Temp	Humidity	windy	Play		
$\theta = 2.$		Sunny	Hot	High	False	No		
One-item sets	Two-item sets		Three-item se	ts	Four-item sets			
Outlook = Sunny (5)	Outlook = Sunny		Outlook = Sur	nny	Outlook = Sunny			
	Temperature = Hot (2	2)	Temperature =	= Hot	Temperature = Hot			
		Humidity = High (2)			Humidity = High			
					Play = No (2)		
Temperature = Cool (4)		Outlook = Sur	nny	Outlook = F	Rainy			
	Humidity = High (3)		Humidity = Hi	gh	Temperature = Mild			
			Windy = False	: (2)	Windy = False			
					Play = Yes	(2)		
		Dainv	Mild	Normal	False	Yes		

In total: 12 one-item sets, 47 two-item sets
 sets, 39 three-item sets, 6 four-item sets
 and 0 five-item sets (with minimum support False Yes of two)

1.1. Algorithms

Item set mining
 Algorithms

There are quite good number of algorithms for enumerating frequent item sets. Here are two well-known approaches.

Apriori methods:

The first apriori algorithm was proposed independently by Agrawal-Srikant and Mannila-Toivonen-Verkamo.

R. Agrawal and R. Srikant
Fast algorithms for mining association rules in large databases. *in Proc. VLDB* 1994, pp. 487-499 (1994)
H. Mannila, H. Toivonen, and A. I. Verkamo.
Efficient algorithms for discovering association rules. *in Proc. KDD Workshop* 1994: pp. 181-192 (1994).

Backtrack methods:

Algorithms known as *LCM* by Uno et al. are typical examples.

See http://research.nii.ac.jp/~uno/code/lcm.html

1. Item set mining

1.2. Assoc. rule

Frequent item set mining Association rule mining: Application

Association rule mining is to derive a relation among items (in general, attribute values) with a certain "significance".

Outlook	Temp	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	Hiah	True	No

If windy = false and play = no then outlook = sunny and humidity = high

customer ID	purchase record										
12C3321	K. beer	K. chips									
18A2130	K. bread	M. milk	apple	orange							
15A2210	beef	pork ham	egg	• • •							
15B1213											

If (buying) diapers & chips then (buying) beer

Item set mining
 Assoc. rule

1.2. Association rule mining: Application

glossaries

- Support: number of instances p correctly
- Confidence: number of correct proportion of all instances the r
- Example: 4 cool days with norr
 - If temperature = cool then humidit Over
 - \Rightarrow Support = 4, confidence = 100% Sunn
- ✤ Normally: minimum support an ^{Sup} pre-specified (e.g. 58 rules with ^{Rai} and confidence ≥ 95% for weat ^{Sup}

utlook	Temp	Humidity	Windy	Play
unny	Hot	High	False	No
unny	Hot	High	True	No
vercast	Hot	High	False	Yes
ainy	Mild	High	False	Yes
ainy	Cool	Normal	False	Yes
ainy	Cool	Normal	True	No
vercast	Cool	Normal	True	Yes
unny	Mild	High	False	No
unny	Cool	Normal	False	Yes
ainy	Mild	Normal	False	Yes
unny	Mild	Normal	True	Yes
vercast	Mild	High	True	Yes
vercast	Hot	Normal	False	Yes
ainy	Mild	High	True	No

1 Frequent item act mining	Outlook	Temp	Humidity	Windy	Play
i. riequent item set mining	Sunny	Hot	High	False	No
1.2 Association rule mining: Λ	Sunny	Hot	High	True	No
1.2. Association fulle mining. A	Overcast	Hot	High	False	Yes
	Rainy	Mild	High	False	Yes
	Rainy	Cool	Normal	False	Yes
	Rainy	Cool	Normal	True	No
↔ Once all item sets with minimum support	Overcast	Cool	Normal	True	Yes
	Sunny	Mild	High	False	No
have been generated, we can turn them	Sunny	Cool	Normal	False	Yes
into rules	Rainy	Mild	Normal	False	Yes
• Eveneplei	Sunny	Mild	Normal	True	Yes
↔ Example:	Overcast	Mild	High	True	Yes
Humidity = Normal, Windy = False, Play = Yes	Overcast	Hot	Normal	False	Yes
	Rainy	Mild	High	True	No
Seven (2 ^{N-1}) potential rules:					
If Humidity = Normal and Windy = False then Play = Yes	4/4				
If Humidity = Normal and Play = Yes then Windy = False	4/6				
If Windy = False and Play = Yes then Humidity = Normal	4/6				
If Humidity = Normal then Windy = False and Play = Yes	4/7				
If Windy = False then Humidity = Normal and Play = Yes	4/8				
If Play = Yes then Humidity = Normal and Windy = False	4/9				
If True then Humidity = Normal and Windy = False and Play = Yes	4/12				

2. Clustering

- Clustering techniques apply when there is no class to be predicted
- Aim: divide instances into "natural" groups

There are several ways to represent clusters, and several ways to measure the "appropriateness" of clusters.

Representing clusters

- 1. Give a cluster label to each instance: disjoint sets.
- 2. By giving sets of instances: may have some overlaps.

3. Probabilistic assignment

	1	2	3	
a	0.4	0.1	0.5	
b	0.1	0.8	0.1	
с	0.3	0.3	0.4	
d	0.1	0.1	0.8	
е	0.4	0.2	0.4	
f	0.1	0.4	0.5	
g	0.7	0.2	0.1	
h	0.5	0.4	0.1	





NB: dendron is the Greek word for tree

Clustering Two major algorithms

We explain two major algorithms for clustering. Both are designed w.r.t. a certain way to measure the "appropriateness" of clusters. Here we mainly discuss these measures.

k-means: An algorithm for separating instances to k disjoint sets (for a given k) so that the total distance from each center becomes small. (\leftarrow usually, hard to get the smallest one)

- Simplest case: one numeric attribute
 - Distance is the difference between the two attribute values involved (or a function thereof)
- Several numeric attributes: normally, Euclidean distance is used and attributes are normalized
- Nominal attributes: distance is set to 1 if values are different, 0 if they are equal

2. Clustering

2.1. Two algorithms

Clustering Two major algorithms

heuristics

k-means: An algorithm to separating instances to k disjoint sets (for a given k) so that the total distance from each center becomes small. (\leftarrow usually, hard to get the smallest one)

- To cluster data into k groups: (k is predefined)
- 1. Choose k cluster centers
 - □ e.g. at random
- 2. Assign instances to clustersD based on distance to cluster centers
- 3. Compute *centroids* of clusters
- 4. Go to step 1
 - until convergence

well... simplified very much!



2. Clustering

Clustering Two major algorithms

EM-algorithm: Probabilistic version of k-means that tries to get "most likely" clusters for a given dataset.

under a certain probabilistic assumption

Most typically, we assume that instances are generated randomly under a mixture of (several) normal distributions.

da A	ta 51 43	B	62 47	B	64 51	AB	48 64	AB	39 62	A	51 48	model
B A A A A A	62 64 45 42 46 45 45	A B A A A A	52 64 51 65 48 49 46	A B A B A A	52 62 49 48 62 43 40	A B B B A	51 63 43 65 66 65 46	B A B A B A	64 52 63 64 48 64 48	B A A A	64 42 48 41	
												$\mu_{\rm A}$ =50, $\sigma_{\rm A}$ =5, $\rho_{\rm A}$ =0.6 $\mu_{\rm B}$ =65, $\sigma_{\rm B}$ =2, $\rho_{\rm B}$ =0.4

2. Clustering

2.1. Two algorithms

2. Clustering

EM-algorithm: Probabilistic version of k-means 2.1. Two algorithms that tries to get "most likely" clusters for a given dataset.

A B A A A A	51 43 62 64 45 42 46 45	B A B A B A A	62 47 52 64 51 65 48 49	B A B A B A A	64 51 52 62 49 48 62 43	A B A B B B B	48 64 51 63 43 65 66 65	A B A B A B A B	39 62 64 52 63 64 48 64	A B A A A	51 48 64 42 48 41	model A B
A	45 45	A	49 46	A	43 40	В А	65 46	В А	64 48			30 40 50 60 70
												$\mu_{\rm A}$ =50, $\sigma_{\rm A}$ =5, $\rho_{\rm A}$ =0.6 $\mu_{\rm B}$ =65, $\sigma_{\rm B}$ =2, $\rho_{\rm B}$ =0.4

Probability that instance x belongs to cluster A:

 $\Pr[A \mid x] = \frac{\Pr[x \mid A] \Pr[A]}{\Pr[x]} = \frac{f(x; \mu_A, \sigma_A) p_A}{\Pr[x]} \quad \text{with} \quad f(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma}} e^{\frac{(x-\mu)^2}{2\sigma^2}}$

★ Likelihood of an instance given the clusters: Pr[x] = ∑_j Pr[x|cluster_j] • Pr[cluster_j] = ∑_j Pr[x & cluster_j]
★ Log-likelihood of n instances in the training set: log∏_i Pr[x_i] = ∑_i log Pr[x_i] ← x₁, ..., x_n

2. Clustering

2.2. # of clusters

2. Clustering2.2. How to determine # of clusters?

A general question on clustering is to a way to determine the number k of clusters. In fact, this is the topic of our Ex.#5. Please do some literature study and try two or three (or even more) ways.

Note that the important point is a way to evaluate the "appropriateness" of k because if we have a good way for measuring the appropriateness of k, then we would be able to design a search algorithm for k.