

## COURSE PRESENTATION:

### *Discriminative Frequent Pattern Analysis for Effective Classification*

Presenter: Han Liang

1

## Outline

- Motivation
- Introduction
- Academic Background
- Methodologies
- Experimental Study
- Contributions

2

## Motivation

- ❖ *Frequent patterns are potentially useful in many classification tasks, such as association rule-based classification, text mining, and protein structure prediction.*
- ❖ *Frequent patterns can accurately reflect underlying semantics among items (attribute-value pairs).*

3

## Introduction

- ❖ *This paper investigates the connections between the support of a pattern and its information gain (a discriminative measure), and develops a method to set the minimum support in pattern mining. It also proposes a pattern selection algorithm. Finally, the generated frequent patterns can be used for building high quality classifiers.*
- ❖ *Experiments on UCI data sets indicate that the frequent pattern-based classification framework can achieve high classification accuracy and good scalability.*

4

# Classification – A Two-Step Process

*Step1: Classifier Construction: learning the underlying class probability distributions.*

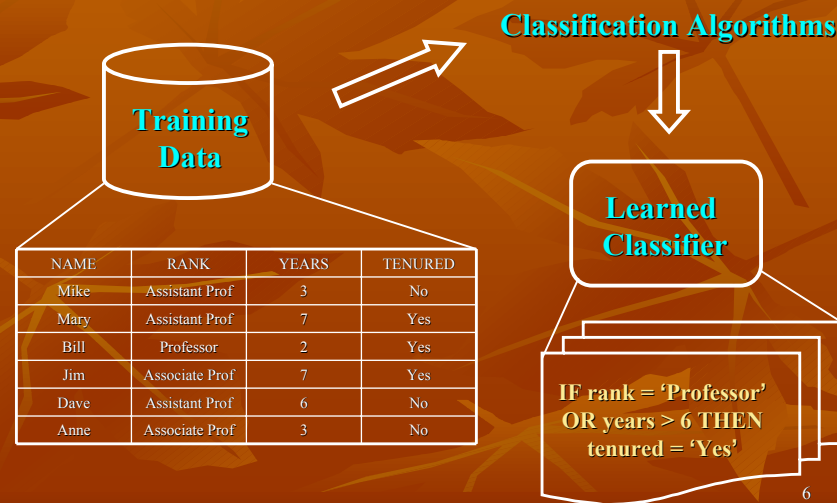
- > The set of instances used for building classifiers is called *training* data set.
- > The learned classifier can be represented as classification rules, decision trees, or mathematical formulae (e.g. Bayesian rules).

*Step2: Classifier Usage: classifying unlabeled instances.*

- > Estimate accuracy of the classifier.
- > Accuracy rate is the percentage of test instances which are correctly classified by the classifier.
- > If the accuracy is acceptable, use the classifier to classify instances whose class labels are not known.

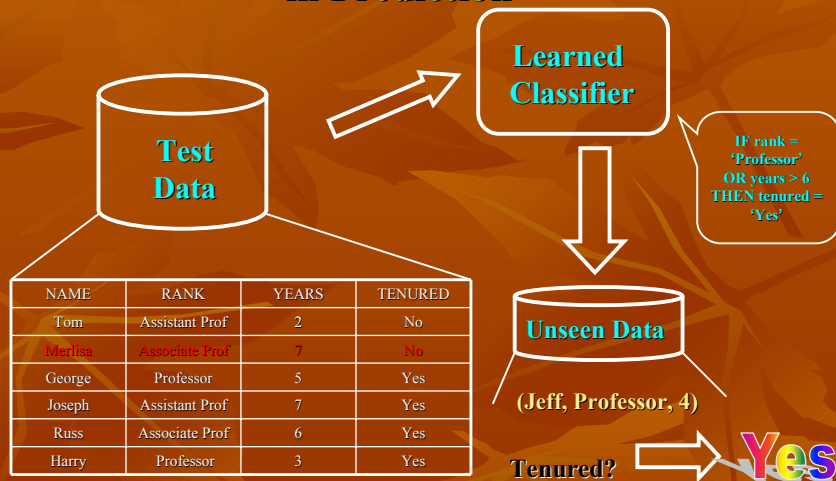
5

## Classification Process I – Classifier Construction



6

## Classification Process II – Use the Classifier in Prediction



7

## Association-Rule based Classification (ARC)

*Classification Rule Mining: discovering a small set of classification rules that forms an accurate classifier.*

- > The data set  $E$  is represented by a set of items (or attribute-value pairs)  $I = \{a_1, \dots, a_n\}$  and a set of class memberships  $C = \{c_1, \dots, c_m\}$ .
- > Classification Rule:  $X \Rightarrow Y$ , where  $X$  is the body and  $Y$  is the head.
  - >  $X$  is a set of items (a sub set of  $I$ , denoted as  $X \subset I$ ).
  - >  $Y$  is a class membership item.
- > Confidence of a classification rule:  $\text{conf} = S(X \cup Y) / S(X)$ .
  - > Support  $S(X)$ : the number of training instances that satisfy  $X$ .

8

## ARC-II: Mining - Apriori

Generate all the classification rules with support and confidence larger than predefined values.

- > Divide training data set into several subsets; one subset for each class membership.
- > For each subset, with the help of *on-the-shelf* rule mining algorithms (e.g. Apriori), mines all item sets above the minimum support, and call them *frequent* item sets.
- > Output rules by dividing frequent item sets in rule body (attribute-value pairs) and head (one class label).
- > Check if the confidence of a rule is above the minimum confidence.
- > Merge rules from each sub set, and sort rules according to their confidences.



9

## ARC-III: Rule Pruning

Prune the classification rules with the goal of improving accuracy.

- > Simple Strategy:
  - > Bound the number of rules.
- > Pessimistic error-rate based pruning.
  - > For a rule, if we remove a single item from the rule body and the new rule decreases in error rate, we will prune this rule.
- > Data set coverage approach.
  - > If a rule can classify at least one instance correctly, we will put it into the resulting classifier.
  - > Delete all covered instances from training data set.



10

## ARC-IV: Classification

Use the resulting classification rules to classify unseen instances.

- > Input:
  - > Pruned, sorted list of classification rules .
- > Two different approaches:
  - > Majority vote
  - > Use the first rule that is applicable to the unseen instance for classification.



11

## The Framework of Frequent-Pattern based Classification (FPC)

- ❖ Discriminative Power vs. Information Gain
- ❖ Pattern Generation
- ❖ Pattern Selection
- ❖ Classifier Construction

12

## Discriminative Power vs. Information Gain - I

The discriminative power of a pattern is evaluated by its information gain.

- Pattern based Information Gain:
  - Data set  $S$  has  $S_i$  training instances that belong to class  $C_i$ . Thus,  $S$  is divided into several subsets, denoted as  $S = \{S_1 \dots S_i \dots S_m\}$ .
  - Pattern  $X$  divides  $S$  into two subsets: the group where pattern  $X$  is applicable and the group where pattern  $X$  is rejected. (binary splitting)
  - The information gain of pattern  $X$  is calculated via:

$$\text{Gain}(X) = I(S_1, S_2, \dots, S_m) - E(X)$$

where  $I(S_1, S_2, \dots, S_m)$  is represented by:

$$I(S_1, S_2, \dots, S_m) = - \sum_{i=1}^m \frac{S_i}{S} \log_2 \frac{S_i}{S}$$

and  $E(X)$  is computed by:

$$E(X) = \sum_{j \in \{x, \bar{x}\}} \frac{|S_j|}{|S|} I(S_j)$$

## Discriminative Power vs. Information Gain - II

The discriminative power of a pattern is evaluated by its information gain.

- Information Gain is related to pattern support and pattern confidence:
  - To simplify the analysis, assume pattern  $X \in \{0,1\}$  and  $C = \{0,1\}$ . Let  $P(x=1) = \theta$ ,  $P(c=1) = p$  and  $P(x=1|c=1) = q$ .
  - Then,

$$E(X) = \sum_{j \in \{x, \bar{x}\}} \frac{|S_j|}{|S|} I(S_j)$$

can be instantiated as:

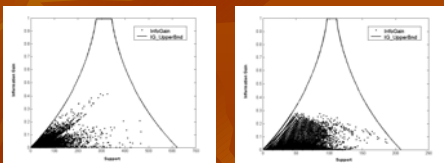
$$\begin{aligned} E(X) &= - \sum_{x \in \{0,1\}} P(x) \sum_{c \in \{0,1\}} P(c|x) \log P(c|x) \\ &= -\theta q \log q - \theta(1-q) \log(1-q) + (\theta q - p) \log \frac{p - \theta q}{1 - \theta} \\ &\quad + (\theta(1-q) - (1-p)) \log \frac{(1-p) - \theta(1-q)}{1 - \theta} \end{aligned}$$

where  $\theta$  and  $q$  are actually the support and confidence of pattern  $X$ .

## Discriminative Power vs. Information Gain - III

The discriminative power of a pattern is evaluated by its information gain.

- Information Gain is related to pattern support and pattern confidence:
  - Given a dataset with a fixed class probability distribution,  $I(S_1, S_2, \dots, S_m)$  is a constant part.
  - $E(X)$  is a concave function, it reaches its lower bound w.r.t.  $q$ , for fixed  $p$  and  $\theta$ .
  - After mathematical analysis, we draw the following two conclusions that:
    - The discriminative power of a low-support pattern is poor. It will harm classification accuracy due to over-fitting.
    - The discriminative power of a very high-support pattern is also weak. It is useless for improving classification accuracy.
  - Experiments on UCI datasets.



Austral

Sonar

The X axis represents the support of a pattern and the Y axis represents the information gain. We can clearly see that both low-support and very high-support patterns have small values of information gain.

## Pattern Selection Algorithm MMRFS

- Relevance: A relevance measure  $S$  is a function mapping a pattern  $X$  to a real value such that  $S(X)$  is the relevance w.r.t. the class label.
  - Information gain can be used as a relevance measure.
  - A pattern can be selected if it is relevant to the class label measured by IG.
- Redundancy: A redundancy measure  $R$  is a function mapping two patterns  $X$  and  $Z$  to a real value such that  $R(X, Z)$  is the redundancy value between them.
  - Mapping function: 
$$R(X, Z) = \frac{P(X, Z)}{P(X) + P(Z) - P(X, Z)} \min(S(X), S(Z))$$
  - A pattern can be chosen if it contains very low redundancy to the patterns already selected.

# Pattern Selection Algorithm MMRFS-II

## Algorithm 1 Feature Selection Algorithm MMRFS

Input: Frequent patterns  $\mathcal{F}$ , Coverage threshold  $\delta$ , Relevance  $S$ , Redundancy  $R$   
 Output: A selected pattern set  $\mathcal{F}_s$

- 1: Let  $\alpha$  be the most relevant pattern;
- 2:  $\mathcal{F}_s = \{\alpha\}$ ;
- 3: **while** (true)
- 4: Find a pattern  $\beta$  such that the gain  $g(\beta)$  is the maximum among the set of patterns in  $\mathcal{F} - \mathcal{F}_s$ ;
- 5: If  $\beta$  can correctly cover at least one instance
- 6:  $\mathcal{F}_s = \mathcal{F}_s \cup \{\beta\}$ ;
- 7:  $\mathcal{F} = \mathcal{F} - \{\beta\}$ ;
- 8: If all instances are covered  $\delta$  times or  $\mathcal{F} = \emptyset$
- 9: **break**;
- 10: **return**  $\mathcal{F}_s$

- ✦ The MMRFS algorithm searches over the pattern space in a greedy way.
- ✦ In the beginning, a pattern with highest relevance value (information gain value) is selected. Then the algorithm incrementally selects more patterns from  $\mathcal{F}$ .
- ✦ A pattern is selected if it has the maximum estimated gain among the remaining patterns  $\mathcal{F} - \mathcal{F}_s$ . The estimated gain is calculated by:
 
$$g(\alpha) = S(\alpha) - \max_{\beta \in \mathcal{F}_s} R(\alpha, \beta)$$
- ✦ The coverage parameter  $\delta$  is set to ensure that each training instance is covered at least  $\delta$  times by the selected patterns. In this way, the number of patterns selected is automatically determined.

17

# Experimental Study

- ✦ Basic learning classifiers – used to classify unseen instances.
  - > C4.5 and SVM
- ✦ For each dataset, a set of frequent patterns  $F$  is generated.
  - > A basic classifier will be built using all patterns in  $F$ . We call it Pat\_All.
  - > MMRFS is applied on  $F$  and a basic classifier is built using a set of selected features  $F_s$ . We call the resulting classifier Pat\_FS.
  - > For comparisons, basic classifiers which are built on single features are also tested. Item\_All represents the basic classifier which is built on all single features, and Item\_FS built on a set of selected single ones.
- ✦ Classification Accuracy.
- ✦ All experimental results are obtained by use of ten-fold cross validation.

18

# Experimental Study-II

**Table 1. Accuracy by SVM on Frequent Combined Features vs. Single Features**

Data	Single Feature			Freq. Pattern	
	Item_All	Item_FS	Item_RBF	Pat_All	Pat_FS
anneal	<b>99.78</b>	<b>99.78</b>	99.11	99.33	99.67
austral	85.01	85.50	85.01	81.79	<b>91.14</b>
auto	83.25	84.21	78.80	74.97	<b>90.79</b>
breast	97.46	97.46	96.98	96.83	<b>97.78</b>
cleve	84.81	84.81	85.80	78.55	<b>95.04</b>
diabetes	74.41	74.41	74.55	77.73	<b>78.31</b>
glass	75.19	75.19	74.78	79.91	<b>81.32</b>
heart	84.81	84.81	84.07	82.22	<b>88.15</b>
hepatic	84.50	89.04	85.83	81.29	<b>96.83</b>
horse	83.70	84.79	82.36	82.35	<b>92.30</b>
iono	93.15	94.30	92.61	89.17	<b>95.44</b>
iris	94.00	<b>96.00</b>	94.00	95.33	<b>96.00</b>
labor	89.99	91.67	91.67	94.99	<b>95.00</b>
lymph	81.00	81.62	84.29	83.67	<b>96.67</b>
pima	74.56	74.56	76.15	76.43	<b>77.16</b>
sonar	82.71	86.55	82.71	84.60	<b>90.86</b>
vehicle	70.43	72.93	72.14	73.33	<b>76.34</b>
wine	98.33	99.44	98.33	98.30	<b>100</b>
zoo	97.09	97.09	95.09	94.18	<b>99.00</b>

- ✦ Table 1 shows the results by SVM.
- ✦ Pat\_FS has significant improvement over Item\_All and Item\_FS. This conclusion indicates that:
  - > the discriminative power of some frequent patterns is higher than that of single features.
- ✦ The performance of Pat\_All is much worse than that of Pat\_FS. That confirms that redundant and non-discriminative patterns let classifiers over-fit the data and decrease the classification accuracy.
- ✦ The experimental results by C4.5 is similar to SVM's.

19

# Experimental Study-III: Scalability Test

**Table 3. Accuracy & Time on Chess Data**

min_sup	#Patterns	Time (s)	SVM (%)	C4.5 (%)
1	N/A	N/A	N/A	N/A
2000	68,967	44.703	92.52	97.59
2200	28,358	19.938	91.68	97.84
2500	6,837	2.906	91.68	97.62
2800	1,031	0.469	91.84	97.37
3000	136	0.063	91.90	97.06

**Table 4. Accuracy & Time on Waveform Data**

min_sup	#Patterns	Time (s)	SVM (%)	C4.5 (%)
1	9,468,109	N/A	N/A	N/A
80	26,576	176.485	92.40	88.35
100	15,316	90.406	92.19	87.29
150	5,408	23.610	91.53	88.80
200	2,481	8.234	91.22	87.32

**Table 5. Accuracy & Time on Letter Recognition Data**

min_sup	#Patterns	Time (s)	SVM (%)	C4.5 (%)
1	5,147,030	N/A	N/A	N/A
3000	3,246	200.406	79.86	77.08
3500	2,078	103.797	80.21	77.28
4000	1,429	61.047	79.57	77.32
4500	962	35.235	79.51	77.42

- ✦ Scalability tests are performed to show the frequent pattern-based framework is very scalable with good classification accuracy.
- ✦ Three large UCI data sets are chosen.
- ✦ In each table, experiments are conducted by varying min\_sup. #Patterns gives the number of frequent patterns. Time gives the sum of pattern mining and pattern selection time.
- ✦ min\_sup = 1 is used to enumerate all feature combinations. Pattern selection fails with such a large number of patterns.
- ✦ In contrast, the frequent pattern-based framework is very efficient and achieves good accuracy within a wide range of minimum support thresholds.

20

## Contributions

- ❖ This paper propose a framework of frequent pattern-based classification. By analyzing the relations between pattern support and its discriminative power, the paper shows that frequent patterns are very useful for classification.
- ❖ Frequent pattern-based classification can use the state-of-the-art frequent pattern mining algorithm for pattern generation, thus achieving good scalability.
- ❖ An effective and efficient pattern selection algorithm is proposed to select a set of frequent and discriminative patterns for classification.

*Thanks!*

*Any Question?*