

Apriori Algorithm for Sub-category Classification Analysis of Handwriting

Sung-Hyuk Cha
School of Computer Science and Information Systems
Pace University
Pleasantville NY 10570-9913
scha@pace.edu

Sargur N. Srihari
CEDAR
State University of New York at Buffalo
Buffalo, New York 14228-2567
srihari@cedar.buffalo.edu

Abstract

The sub-category classification problem is that of discriminating a pattern to all sub-categories. Not surprisingly, sub-category classification performance estimates are useful information to mine as many researchers are interested in any trend of pattern in specific sub-category. This paper presents a datamining technique to mine a database consisting of experimental and observational unit variables. Experimental unit variables are those attributes which make sub-categories of the entity, e.g., demographic data and observational unit variables are features observed to classify the entity, e.g., test results or handwriting styles, etc. Since there are an enormously large number of sub-categories based on the experimental unit variables, we apply the Apriori algorithm to select only sub-categories that have enough support among all possible ones in a given database. Those selected sub-categories are then discriminated using observational unit variables as input features to the Artificial Neural Network (ANN) classifier. The importance of this paper is twofold. First, we propose an algorithm that quickly selects all sub-categories that have enough both support and classification rate. Second, we successfully applied the proposed algorithm to the field of handwriting analysis. The task is to determine similarity of handwriting style of a specific group of people. Document examiners are interested in trends in the handwriting of specific groups, e.g., (i) does a male write differently from a female? (ii) can we tell the difference in handwriting of age group between 25 and 45 from others?, etc. Subgroups of white males in the age group 15-24 and white females in the age group 45-64 show 87 % correct classification performance.

1 Introduction

The problem of discrimination and classification is mainly concerned with differentiating between $g (g \geq 2)$

mutually exclusive populations and with classifying subjects or objects on the basis of multivariate observation [8, 9]. The mathematical problem of discrimination is to divide the observation space R^p into g mutually exclusive and exhaustive regions R_1, \dots, R_g , so that if an observation x falls in R_i , then it is classified as belonging to the i th population. The sub-category classification problem involves defining mutually exclusive populations according to their attributes and classifying an unseen instance into its sub-categories using multiple features. Borrowing terminologies from *design of experiments*, we call the attributes contributing defining the sub-categories as *experimental unit variables* and observed features as *observational unit variables*.

For example of classifying a person, possible attributes (experimental unit variables) are gender, ethnicity, age, etc, which is well known as *demographic data*. One can observe various features such as a person's appearances (observational unit variables) to classify him or her to the proper sub-categories. Consider a subject with three experimental unit variables where each variable has two values: $V_1 = \{v_{1,1}, v_{1,2}\}$, $V_2 = \{v_{2,1}, v_{2,2}\}$, $V_3 = \{v_{3,1}, v_{3,2}\}$. As shown in Figure 1, there are seven possible sub-category classification problems. Subsets of p_4, p_5, p_6 , and p_7 also form sub-category classification problem, e.g., $p_i \rightarrow \{(v_{1,1}, v_{2,2}, v_{3,1}) : (v_{1,1}, v_{2,2}, v_{3,2})\}$, etc. Hence, there are combinatorially large number of sub-category classification problems that one database can create. Enumerating and classifying them all is inefficient.

For this reason, we present an efficient technique to mine a database consisting of both experimental and observational unit variables. The information extracted from the technique is the classification rate of each sub-category. Not surprisingly, sub-category discrimination analysis is valuable information to be mined because many researchers are interested in any trend in specific subgroup.

In order to answer whether one can build a machine that can classify an unseen instance into its sub-category, each class (subgroup) must have a substantial number of instances for the sake of valid statistical inference. This Sine-

	V ₁₁		V ₁₂	
	V ₃₁	V ₃₂	V ₃₁	V ₃₂
V ₂₁				
V ₂₂				

- $p_1 \rightarrow \{v_{1,1} : v_{1,2}\}$
 $p_2 \rightarrow \{v_{2,1} : v_{2,2}\}$
 $p_3 \rightarrow \{v_{3,1} : v_{3,2}\}$
 $p_4 \rightarrow \{(v_{1,1}, v_{2,1}) : (v_{1,1}, v_{2,2}) : (v_{1,2}, v_{2,1}) : (v_{1,2}, v_{2,2})\}$
 $p_5 \rightarrow \{(v_{1,1}, v_{3,1}) : (v_{1,1}, v_{3,2}) : (v_{1,2}, v_{3,1}) : (v_{1,2}, v_{3,2})\}$
 $p_6 \rightarrow \{(v_{2,1}, v_{3,1}) : (v_{2,1}, v_{3,2}) : (v_{2,2}, v_{3,1}) : (v_{2,2}, v_{3,2})\}$
 $p_7 \rightarrow \{(v_{1,1}, v_{2,1}, v_{3,1}) : (v_{1,1}, v_{2,1}, v_{3,2}) : (v_{1,1}, v_{2,2}, v_{3,1}) : (v_{1,1}, v_{2,2}, v_{3,2}) : (v_{2,1}, v_{2,2}, v_{3,1}) : (v_{2,1}, v_{2,2}, v_{3,2}) : (v_{2,2}, v_{2,2}, v_{3,1}) : (v_{2,2}, v_{2,2}, v_{3,2})\}$

Figure 1. Sample sub-category classification problems

qua-non is called *support*. We apply the *Apriori* algorithm to select all sub-categories that have enough support among all possible ones in a given database. Those selected sub-categories are then discriminated using the *Artificial Neural Network (ANN)* classifier. Finally, the performance measures for each selected sub-category problem are reported as final outputs. We use the *Artificial Neural Network (ANN)* because it is equivalent to multivariate statistical analysis. There is a wealth of literature regarding a close relationship between neural networks and the techniques of statistical analysis, especially multivariate statistical analysis, which involves many variables [7, 9].

Apriori Algorithm was originally designed for the purpose of efficient association rule mining by *Agrawal et al* [3, 2]. The concept of association rules was introduced in 1993 [1] and many researchers have endeavored to improve the performance of algorithms that discover the association rules in large datasets. The *Apriori algorithm* is an efficient association discovery algorithm that filters item sets by incorporating item constraints (support). We apply this filtration by support approach to mine the subgroup classification information in a large database with one table with experimental units (writer information) and the other table with observational units (document image features).

As an illustrative example, we consider a CEDAR letter database [6] consisting of 1,000 writer data with six writer attributes and features extracted from a handwriting sample to determine similarity of a specific group of people. Document examiners are interested in any trend of handwriting

in specific group, e.g., *i*) does male write differently from female? *ii*) can we tell the difference in handwriting of age group between 25 and 45 from others?, etc.

2 Database

The CEDAR letter database is a database consisting of 3,000 handwritten document images written by 1,000 subjects representative of the US population by stratification and proportional allocation [6]. Each individual provided three handwriting samples, written with black ballpoint pens on plain white sheets by copying text that has all alphabets in every position of a word, as well as numerals. The database was originally created for the purpose of the *writer identification* study [5, 4].

We stratified our database along six experimental unit variables: gender (G) {male, female}, handedness (H) {right, left}, age (A) {under 15, 15 ~ 24, 25 ~ 44, 45 ~ 64, 65 ~ 84, over 85}, ethnicity (E) {white, black, hispanic, asian and pacific islander, American Indian, Eskimo, Aleut}, highest level of education (D) {below highschool graduate, above}, and place of schooling (S) {USA, Foreign}. We can analyze the association between different combinations of these variables. Studying the association between only one variable is called *1-constraint* subgroup analysis (male vs. female, or black vs. white vs. hispanic), studying the association between pairs of variables is called *2-constraint* subgroup analysis (white male vs. white female), and so on. The size of 2-constraint subgroup is the multiplication of each size of 1-constraint subgroup. Figure 2 shows some of the possible combinations. As mentioned earlier, there can be a combinatorially large number of subgroup classification problems one can analyze. Since our database was not stratified heavily across each variable, it is necessary to identify those subgroups that have enough support or coverage.

There seem to be conflicting views about whether or not group-specific handwriting features can be attributed to the sex, age, ethnicity, or handedness of writers. Correlations between these groups and handwriting features are dealt with in [10]. For instance, while tremors caused by aging may have a bearing on handwriting, the direction of horizontal strokes, amount of pressure exerted on up-strokes or down-strokes, consistency in letter slopes, slope of writing, direction of curves may be affected by handedness. Ethnicity may also affect handwriting. Hispanic writers have a tendency to ornateness in the formation of capital letters, slope of writing in France, United Kingdom and India tends to be vertical or even slightly backhand and is clearly forehand in Germany.

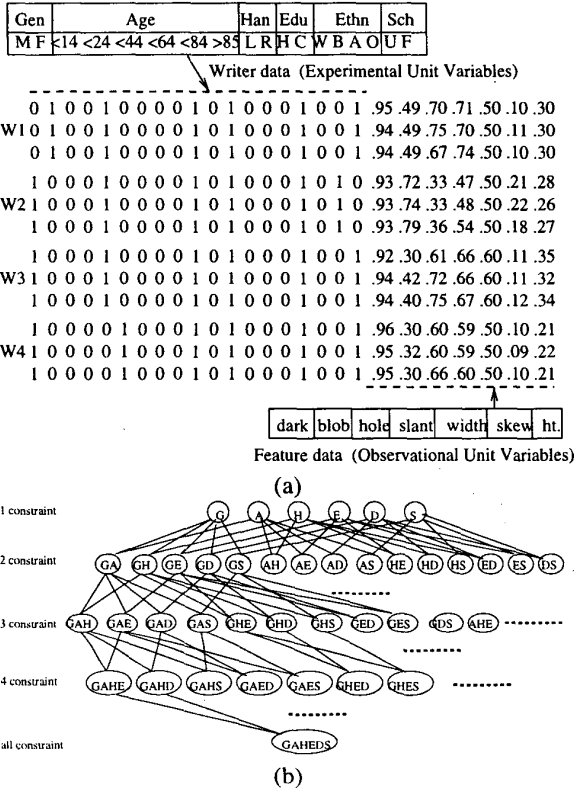


Figure 2. (a) Sample Entries of CEDAR letter database and (b) List of sub-categories where G, A, H, E, D, and S correspond to Gender, Age, Handedness, Ethnicity, Degree of education, and place of Schooling, respectively.

3 Algorithm

The algorithm has two filtering stages with two user defined threshold values. Figure 3 illustrates the algorithm figuratively. To find the subgroup classification problem

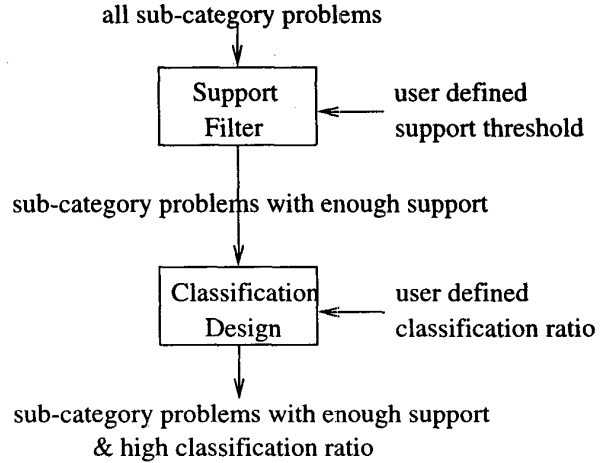


Figure 3. Apriori algorithm overview.

lem that is statistically valid (having enough supporting instances), we show the efficient *Apriori* algorithm without aggregation. In this algorithm, the subgroup with aggregation such as {white vs. (black \cup hispanic)} is not considered. First, for each attribute value of every variable, count the occurrences to find whether the sum exceeds the user defined minimum support. If so, add it to the *1-constraint* output list: {(Male: Female), (15 ~ 24 : 25 ~ 44 : 45 ~ 64), (white: black: hispanic)}. Second, from the *1-constraint* output list, generate the *2-constraint* output list. They are {(Male, 15 ~ 24 : Female, 15 ~ 24 : Male, 25 ~ 44 : Female, 25 ~ 44), (Male, white: Female, white), (white, 15 ~ 24 : white, 25 ~ 44)}. Next, from the *2-constraint* output list, generate the *3-constraint* output list, it is {(white, male, 15 ~ 24 : white female, 15 ~ 24)}. Repeat generating the higher constraint list until it is empty or it reaches the all-constraint list. It is efficient because we do not consider all possible combinations but generate the higher constraint list from only those elements in the lower constraint list.

4 Experimental Results

Once we identify those subgroup classification problems that have enough support (over 200 samples), we train an ANN (Artificial Neural Network) to evaluate the accuracy with which we can classify these subgroups as shown in

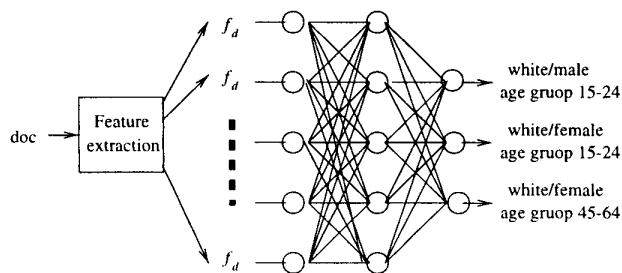


Figure 4. Artificial Neural Network classifier for writer subgroup classification problem.

Figure 4. The samples of interest were divided into 4 groups - one for training, one for validation, and two test sets. An ANN was trained using the feature values of samples. Examples of 1-constraint subgroups and their classification rates are Male vs. Female (70.2%), black vs. white (64.5%), left vs. right (59.5%), below vs. above highschool graduate (61%). The 2-constraint subgroup of white males, and white females was studied and 68% (type-I error: 24%, type-II error: 40.5%) performance was obtained. The 3-constraint subgroup of white males in the age group 15-24, and white females in the age group 45-64 was studied. The performance on the two tests sets was 83% (type-I error: 21%, type-II error: 12%) and 87% (type-I error: 14%, type-II error: 12%). It is observed that the more constraint, the higher classification rate.

5 Conclusion

In this paper, we presented a datamining technique for sub-category classification problem. *A priori* algorithm is applied to filter the sub-category problems with insufficient support. We considered a database consisting of writer data and features obtained from a handwriting sample, statistically representative of the US population, to determine similarity of a specific group of people. Higher classification rate is achieved with higher constraint sub-categories. In this paper, the subgroup with aggregation was not considered. The algorithm can be altered to handle the subgroup with aggregation.

Acknowledgments

This research has been possible funded by National Institute of Justice (NIJ) in response to the solicitation entitled *Forensic Document Examination Validation Studies: Award Number 1999-IJ-CX-K010* [11]. Also, we would like to thank Eugenia Smith for helping us build the database.

References

- [1] Rakesh Agrawal, Tomasz Imielinski, and Arun Swami. Mining association rules between sets of items in large databases. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, volume 22, pages 207–216, June 1993.
- [2] Rakesh Agrawal, Heikki Mannila, Ramakrishnan Srikant, Hannu Toivonen, and Inkeri Verkamo. Fast discovery of association rules. In Usama M. Fayyad, Gregory Piatetsky-Shapiro, Padhraic Smyth, and Ramasamy Uthurusamy, editors, *Advances in Knowledge Discovery and Data Mining*, chapter 12. AAAI Press, 1996.
- [3] Rakesh Agrawal and Ramakrishnan Srikant. Fast algorithms for mining association rules. In *Proceedings of the 20th Int'l Conference on Very Large Databases*, volume 2, pages 478–499, September 1994.
- [4] Sung-Hyuk Cha and Sargur N. Srihari. Multiple feature integration for writer verification. In *Proceedings of 7th IWFHR 2000*, pages 333–342, September 2000.
- [5] Sung-Hyuk Cha and Sargur N. Srihari. Writer identification: Statistical analysis and dichotomizer. In *Proceedings of SS&SPR 2000 LNCS - Advances in Pattern Recognition*, volume 1876, pages 123–132. Springer-Verlag, September 2000.
- [6] Sung-Hyuk Cha and Sargur N. Srihari. Handwritten document image database construction and retrieval system. In *Proceedings of SPIE, Document Recognition and Retrieval*, volume 4307, pages 13–21, January 2001.
- [7] Vladimir Cherkassky, Jerome H. Friedman, and Harry Wechsler. *From Statistics to Neural Networks: Theory and Pattern Recognition Applications*. Springer, nato asi edition, 1994.
- [8] Sung C. Choi and Ervin Y. Rodin. *Statistical Methods of Discrimination and Classification, Advances in Theory and Applications*. Pergamon Press, 1986.
- [9] Richard O. Duda, Peter E. Hart, and David G. Stork. *Pattern Classification*. John Wiley & Sons, Inc., 2nd edition, 2000.
- [10] Roy A. Huber and A. M. Headrick. *Handwriting Identification: Facts and Fundamentals*. CRC Press LLC, 1999.
- [11] Jeremy Travis. Forensic document examination validation studies. Solicitation: <http://ncjrs.org/pdffiles/sl297.pdf>, October 1998.