



Credit risk evaluation by using nearest subspace method

Arekatla Madhava Reddy, Arekatla Jaganmohan Reddy, Shaik Guntur Mahabub
Subhani, Mr. Merugu Anand Kumar

Assistant Professor^{1,4}, Associate Professor^{2,3}

amreddy2008@gmail.com¹, jagan.arekatla@gmail.com²,
subhanimehandi@gmail.com³, meruguanand502@gmail.com⁴

Department of CSE, A.M. Reddy Memorial College of Engineering and Technology,
Petlurivaripalem, Narasaraopet, Andhra Pradesh

Abstract

In this study, we use a classification strategy called the closest subspace approach to assess credit risk. Identifying "good" and "bad" creditors via credit risk assessment is a common categorization challenge. There has been a lot of talk lately about using machine learning techniques like support vector machine (SVM) to assess credit risk. Yet there is plenty No tried-and-true pattern recognition or AI-based classification techniques for use in assessing creditworthiness exist. This work proposes using the closest subspace classification technique, a robust approach to facial recognition, in the context of credit scoring. When evaluating creditworthiness, the nearest subspace credit evaluation method uses the subspaces spanned by creditors in the same class to extend the training set, with the Euclidean distance between a test creditor and the subspace serving as the similarity measure for classification.

1. Introduction

In the realm of financial risk management, credit risk appraisal analysis is a subject of intense interest. It is a common classification challenge to sort "good" creditors from "bad" ones. Numerous data mining approaches, including logit analysis [1], probity analysis [2], ANN [3], ANN [4], genetic algorithms [5, 6], and genetic programming [7, 8], have been applied to the task of assessing credit risk in recent years. Support vector machine (SVM) [8-14], genetic algorithm (GA) [5], multiple criterion linear programming (MCLP) [6] [7], etc. Even though there are a growing number of learning approaches being used to credit assessment, certain very efficient classification strategies from the fields of pattern recognition and artificial intelligence have yet to be explored. In this research, we evaluated credit risks using a closest subspace classification approach [15, 16, and 17]. In closest subspace classification, the subspaces covered by the training samples of each class are used to represent the training set, and the samples of interest are assigned to the class that contains the nearest subspace. The issues of facial

recognition have been effectively used to this method of categorization [16, 17]. For facial image data, [16] proposes the closest feature subspace (NFS) technique, which uses feature extraction and then uses arbitrary k ($k > 3$) feature training samples to

Span subspaces in each class; the resulting subspaces are used as extensions for the training set. As shown to be optimal for classification in [17], we additionally employ subspaces covered by all samples inside each class to represent each class. We discover that existing credit data are often low dimensional data; hence the feature reduction approach is not employed in this article for credit assessment. In other words, we utilize all the samples from each class to generate subspaces, and a test sample will be a member of the class represented by the subspace to which it most closely corresponds. This approach to credit scoring is known as nearest subspace (NS) credit scoring. The NS credit evaluation approach outperforms the SVM method and the 1-NN method on a U.S. credit dataset. The rest of the paper is structured as follows: The closest subspace method is described in Section 2. Experiments on a credit assessment dataset are provided in Section 3. The last portion is the summary.

2. nearest Subspace Algorithm

The NS approach's central notion is to increase the representational capacity of class prototypes using subspace. This effectively offers an endless number of prototype points, which may explain more variations in the prototypical form than the original samples. The query vector is projected to the class's space in order to determine the distance between the two. Samples from this group occupy a subspace. To get the projection point, choose the linear combination that comes closest to answering the question. Classification is based on how far away the projection point is from the query point. The closest subspace credit assessment approach takes into account all creditors in the same class and utilizes a linear combination (subspace) of them to approximate the various versions of creditors. This means that the original set of creditors used for training may now be combined indefinitely by linear means. Each set of class creditors will include a credit record that is quite close to the one we need for our test creditor. It's possible that the closest credit history isn't from a single lender, but rather the linear average of all lenders in the same category. Finally, the closest creditor record to the test creditor is considered to be of the same class as the test creditor. Linearly combining creditors yields just the subspace occupied by creditors, which is the set of linear combinations. Thus, a creditor is virtually subdivided into the closest subspace of creditors using the NS approach. More potential creditors are generated by NS than by the standard 1-NN or k -NN classifiers. This increases the potential of the pool of established creditors. Here, we provide the NS approach to credit assessment and explain the measure of a test creditor to a subspace.

Subspace Distance

In the NS approach, the distribution estimate of a class is the subspace covered by training samples, and the similarity measure for classification is the Euclidean distance between a test sample and the subspace. A subspace of the set S is the set that contains all linear combinations of samples in S , given that S is a set from the class SRd, 1, 2,, k .

$$F(S) = \sum_{i=1}^k \alpha_i x_i, \quad (1)$$

An exact expression for the distance between a query x and its corresponding subspace of S may be expressed as follows:

$$\begin{aligned} d^2(x, F(S)) &= \min_{y \in F(S)} \|x - y\|^2 \\ &= \min_{\alpha} \left\| x - \sum_{i=1}^k \alpha_i x_i \right\|^2 \\ &= \min_{\alpha} \left[(x \cdot x) - 2 \sum_{i=1}^k \alpha_i (x \cdot x_i) + \sum_{i,j=1}^k \alpha_i \alpha_j (x_i \cdot x_j) \right] \end{aligned} \quad (2)$$

$$d^2(x, F(S)) = \min_a (x^T x - 2x^T Xa + a^T X^T Xa) \quad (3)$$

The solution to Eq. (3) is a straightforward computation of an unconstrained optimum problem:

$$a = (X^T X)^+ X^T x \quad (4)$$

or

$$a = (X^T X + \alpha I)^{-1} X^T x \quad (5)$$

Where I is the identity Matrix for k and $(X^T X)^+$ is the pseudo-inverse of $X^T X$; 0. Once the projection y is calculated, the coefficient may be used to express the optimal linear combination of samples in this subspace.

$$y = \sum_{i=1}^k \alpha_i x_i \quad (6)$$

The Euclidean distance $d(x, y)$ is the distance from x to F(S), where y is the closest point to x in subspace F(S). We can then calculate $d^2(x, F(S))$ by y as follows:

$$d^2(x, F(S)) = \|x - y\|^2 \quad (7)$$

Nearest Subspace Algorithm

When training on many classes simultaneously, the it class's training set looks like $(1, 2, 1, S, S, S, S)$. Training sets for various categories have contributed the following subspaces: We calculate the distance between any two points for any given query. D Si R L () xRd F S1 F (SL) x and the subspaces of all classes:

$$d^2(x, F(S_1)) = \min_{y_1 \in F(S_1)} \|x - y_1\|^2,$$

....,

$$d^2(x, F(S_l)) = \min_{y_l \in F(S_l)} \|x - y_l\|^2.$$

Classifying x into the closest neighbor subspace, we use $d(x, F(S_i))$ as the similarity between $d(x, F(S_i))$ and the it class. In other words, it is a member of the set, $j \in \{1, 2, \dots, l\}$. $x_j \in S_j$ $\arg \min_j d^2(x, F(S_j))$, [1, 2], S x Here we outline the big picture of the NS approach for creditor assessment using a made-up example creditor x. The first step is to calculate the best weights for each category of creditors.

$$a = (X^T X)^+ X^T x \text{ or } a = (X^T X + \alpha I)^{-1} X^T x$$

Step 2: Determine the optimal linear combination of each category's creditor records.

The projection of the test creditor x in subspace F (Si) may be found using the weights for the it class creditor set $I = 1, 2, k$.

$$y_i = \sum_{l=1}^k \alpha_l x_l.$$

Third, we calculate how far away creditor x is from each subspace. A class subspace F (Si) is said to be "distant" from x if and only if

$$d^2(x, F(S_i)) = \|x - y_i\|^2$$

3. Credit Evaluation Experiments

Evaluating credit risk to determine which debtors are "good" and which are "bad" is a common example of a classification problem [18] [19]. In this study, we evaluate credit risks using the closest subspace approach. We compare NS with SVM using a linear kernel and an RBF kernel ($k = \exp(-0.5(x-y)/2)$) on a real-world dataset to assess its usefulness in creditor evaluation. Information about credit in the United States. The experimental credit card dataset comes from a prominent U.S. financial institution. It's 6,000 strong and has 66 calculated fields. Only 960 of the 6,000 records are considered to be in "good" shape, while the remaining 5040 are in bankruptcy [18]. In our tests, we will assess the classifiers based on three different accuracies: good accuracy, bad accuracy, and total accuracy.

$$\text{"Good" Accuracy} = \frac{\text{number of correctly classified "Good" samples in test set}}{\text{number of "Good" samples in test set}},$$

$$\text{"Bad" Accuracy} = \frac{\text{number of correctly classified "Bad" samples in test set}}{\text{number of "Bad" samples in test set}},$$

$$\text{Total Accuracy} = \frac{\text{number of correct classification in test set}}{\text{number of samples in test set}}.$$

Where "Good" accuracy and "Bad" accuracy quantify the classifiers' ability to distinguish between "Good" and "Bad" users. Precision of classification for the dangerous class must be increased to an acceptable quality in the actual world for the unique goals of preventing credit fraud, without unduly detracting from the precision of classification for the safe class. Categories for other categories. Thus, increasing the precision of the "Bad" category is one of the primary goals of credit scoring randomly selecting p ($p=10, 20, 30, 100$) samples from each class as training set and the remainder for test; we conduct experiments on each dataset. Each classifier test is conducted 20 times, and the average results are reported. The Mat lab 7.0 environment is used for all of our investigations. Optimal techniques in Mat lab are used to address the convex quadratic programming issue underlying support vector machines. Tables 1 and 2 provide "Bad" and "Good" classifier comparisons, whereas Table 3 displays overall accuracy comparisons. Both the SVM and KASNP use the same RBF kernel parameters ($=10000$) and the same SVM penalty constant ($C=$).

Method comparisons for "bad" accuracy (percent) on a US dataset Table 1

Number of training data per class	"Bad" accuracy (%) comparisons on USA dataset			
	1-NN	Linear SVM	RBF SVM	NS
10	66.67 %	59.68%	64.51%	58.75 %
20	63.37 %	66.23%	65.43%	65.03 %
30	64.38 %	65.25%	64.48%	72.13 %
40	63.77 %	63.97%	65.34%	76.83 %
50	65.78 %	65.21%	66.20%	76.32 %
60	64.82 %	65.82%	66.01%	74.32 %
70	65.52 %	65.89%	68.31%	75.22 %
80	65.46 %	67.37%	69.99%	72.14 %
90	65.60 %	66.94%	70.62%	71.99 %
100	64.83 %	66.32%	70.69%	71.62 %

Table 2: Comparisons of "Good" accuracy (percent) among techniques on a US dataset

Number of training data per class	"good" accuracy (%) comparisons on USA dataset			
	1-NN	Linear SVM	RBF SVM	NS
10	56.48 %	60.97%	61.50%	66.36 %
20	59.40 %	66.60%	66.82%	69.89 %
30	59.83 %	65.03%	65.64%	67.23 %
40	62.41 %	67.12%	67.62%	64.43 %
50	61.55 %	66.46%	66.62%	65.81 %
60	62.45 %	66.46%	67.84%	67.33 %
70	62.65 %	66.65%	67.49%	68.73 %
80	62.15 %	66.65%	66.35%	71.24 %
90	62.54 %	66.67%	67.74%	69.70 %
100	63.23 %	67.02%	67.97%	68.44 %

Table 3: Total accuracy (%) comparisons between approaches on a dataset from the United States

Number of training data per class	Total accuracy (%) comparisons on USA dataset			
	1-NN	Linear SVM	RBF SVM	NS
10	58.10%	60.77%	61.97%	63.15%
20	60.03%	66.55%	66.60%	69.13%
30	60.54%	63.83%	65.46%	67.99%
40	62.62%	67.81%	67.27%	66.36%
50	62.20%	66.27%	66.55%	67.43%
60	62.81%	66.44%	67.56%	68.40%
70	63.08%	66.54%	67.61%	69.72%
80	62.65%	67.39%	66.90%	71.38%
90	62.99%	65.13%	68.17%	70.04%
100	63.46%	66.92%	68.37%	68.91%

We found that the NS approach performed better than the other methods we tried for classifying risks. By comparing the findings in Tables 1 and 3, we see that the NS approach maintains a higher level of accuracy for all three measures ('bad,' 'good,' and 'Total'). One, the closest subspace (NS) technique is more effective than other classifiers in identifying "Bad" customers. It is clear from Table 1 shows that when the number of training samples for each class is more than 30, the NS technique consistently achieves classification accuracy of 70% or above, whereas other methods often fall short of this threshold. Table 2 compares NS and RBF SVM for recognizing "Good" customers, showing that the former achieves the maximum accuracy at p=10,20,30,70,80,90,100, while the latter achieves the highest accuracy at p=40,50,60. (3) In a nutshell, NS technique is superior to 1-NN and SVMs (see Table 3).

Based on our experiments with the U.S. credit dataset, we find that the NS approach is competitive with 1-NN and SVMs when it comes to classifying borrowers.

4. Conclusion

In this study, we provide a new approach to credit score classification: the closest subspace technique. To approximate the potential versions of creditors, the closest subspace credit assessment approach employs a linear combination (subspace) of all creditors belonging to the same class. In order to get the best estimate for a test creditor, the NS technique uses closest subspace, and then subdivides the class of test debtors into closest subspace classes. The NS achieves satisfactory results when used to assess creditworthiness in a real-world U.S. credit card dataset.

References

1. S. Scholes, *Discuss. Faraday Soc. No. 50 (1970) 222.*
2. O.V. Mazurin and E.A. Porai-Koshits (eds.), *Phase Separation in Glass, North-Holland, Amsterdam, 1984.*
3. Y. Dimitriev and E. Kashchieva, *J.Mater. Sci. 10 (1975) 1419.*
4. D.L. Eaton, *Porous Glass Support Material, US Patent No. 3 904 422 (1975).*
1. Wiginton, J. C. *A note on the comparison of logit and discriminant models of consumer credit behaviour. Journal of Financial Quantitative Analysis, 15 (1980), 757-770.*

2. Grablowsky, B. J., & Talley, W. K. Probit and discriminant functions for classifying credit applicants: A comparison. *Journal of Economic Business*, 33 (1981), 254-261.
3. Malhotra, R., & Malhotra, D. K. Evaluating consumer loans using neural networks. *Omega*, 31, 83-96.
4. Smalz, R., & Conrad, M. (1994). Combining evolution with credit apportionment: A new learning algorithm for neural nets. *Neural Networks*, 7 (2003), 341-351.
5. Varetto, F. Genetic algorithms applications in the analysis of insolvency risk. *Journal of Banking and Finance*, 22 (1998), 1421-1439.
6. Shi, Y., Peng, Y., Xu, W., & Tang, X. Data Mining via Multiple Criteria Linear Programming: Applications in Credit Card Portfolio Management, *International Journal of Information Technology and Decision Making*, Vol. 1 (2002), 131-151.
7. Shi, Y., Wise, M., Luo, M., & Yu, L. Data mining in credit card portfolio management: a multiple criteria decision making approach, in Koksalan, M. and Zionts, S. eds., *Multiple Criteria Decision Making in the New Millennium*, Springer, Berlin, (2001), 427-436.
8. Vapnik, V. N. *The Nature of Statistical Learning Theory*. New York: Springer-Verlag, (1995).
9. Van Gestel, T., Baesens, B., Garcia, J., & Van Dijke, P.: A support vector approach to credit scoring. *Bank en Financierwezen* 2 (2003): 73-82.
10. Bellotti, T., & Crook, J. Support vector machines for credit scoring and discovery of significant features. *Expert Systems with Applications*, 36(2), (2009) 3302-3308.