

Comparison of Classification Techniques

SripriyaMuppidi and Bhatracharyulu N.Ch.

muppidisripriya@gmail.com and dwarakbhat@osmania.ac.in
Department of Statistics, University College of Science,
Osmania University, Hyderabad-7, TS, India.

Abstract

In this paper an attempt is made to compare the classification techniques, linear Discriminant, K-nearest Neighbourhood, Perceptron learning, Naïve Bayes Classifier, Logistic regression. The comparison made is with respect to their methods, merits, and demerits. The methods were implemented for a credit card bank data set and evaluated their accuracy and found that perceptron learning neural network algorithm have more accuracy among them.

Key words: Machine learning techniques, Credit Card, supervised learning.

Date of Submission: 13-02-2023

Date of Acceptance: 26-02-2023

I. INTRODUCTION

A classifier partitions the feature space into K disjoint subspaces, C_1, C_2, \dots, C_K , such that for a sample with expression profile $\mathbf{X} = (x_1, x_2, \dots, x_p) \in C_k$ the predicted class is k . Classifiers are built from a learning or training set: $L = (\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$; where $Y_i \in \{C_1, C_2, \dots, C_K\}$. Classifier D built from a learning set L : $D(\cdot, L): \mathbf{X} \rightarrow \{C_1, C_2, \dots, C_K\}$. Predicted class for observation \mathbf{X} : $D(\mathbf{X}, L) = C_k$ if \mathbf{X} is in C_k .

II. CLASSIFICATION TECHNIQUES

Several researchers contributed in developing different classification algorithms. In this section, the methods of popular classification techniques were presented.

METHOD-1: Fisher (1936) developed a discriminant function for classifying the different species of Iris setosa and Iris versicolor having more than two features. It is a multivariate statistical technique used to separate the two or more populations by constructing a boundary function between the populations.

Step-1: Let $T = \{(\mathbf{X}_i, Y_i), i = 1, 2, \dots, n\}$ be the sample of training data set, measured on feature space X in 'p' dimensions. Let X_0 be the sample observation to be classified to one of the class. Let $D = b'X$ be the discrimination function that classifies the population into two or more classes (where b' be the vector of discriminant score).

Step-2: The vector b' of discriminant scores can be obtained by maximizing $\lambda = b'Bb / b'Wb$ where $B = \sum (\mathbf{X}_i - \bar{X})(\mathbf{X}_i - \bar{X})'$ and $W = \sum (\mathbf{X}_{ij} - \bar{X}_i)(\mathbf{X}_{ij} - \bar{X}_i)'$.

Step-3: Obtain eigen values λ using the characteristic equation $|W^{-1}B - \lambda I| = 0$. and b is the eigen vector corresponding to λ . The resulting discriminant function is $D = b'X$.

Step-4: Let X_0 be an observation to be allocated to one of the populations. Allocate X_0 to i^{th} population if $|b'X_0 - b'X_i| < |b'X_0 - b'X_j|$.

METHOD-2: Cover and Hart (1967) developed Nearest Neighbours algorithm that searches the pattern for the k -training tuples that are closest to the unknown tuple. The closeness is defined in terms of distance metric. The detailed method is presented below.

Step-1: Let $T = \{(\mathbf{X}_i, Y_i), i = 1, 2, \dots, n\}$ be the trained data set measured on feature space in 'p' dimensions and 'm' be the number of labeled classes that sample data set T is classified. Let X_0 be the sample observation to be classified to one of the class. Let k is a parameter, the value of which will be determined by minimizing the cross-validation error

Step-2: Find the shortest distance between the sample observation X_0 to the 'k' closest neighbor observations in the data set T .

Step-3: Choose the class that has majority voting among those 'k' neighbors, Let 'k' is a parameter, the value of which will be determined by minimizing the cross-validation error later.

METHOD-3: Rosenblatt (1957) developed an artificial neuron model used to classify the sample observation, called perceptron learning and explained the concept behind it.

Step-1: Let $T = \{(X_i, Y_i), i = 1, 2, \dots, n\}$ be the training sample data set measured on feature space in 'p' dimensions. Assume Y be the vector of desired binary output response.

Step-2: Initially choose a random vector of weights $W^{(1)}$ in p-dimension and choose a learning rate constant η ($0 \leq \eta \leq 1$).

Step-3: For each input vector X_i , evaluate $O_i = \text{Sign}[W^{(i)}X_i]$.

Step-4: Update the weight vector using the relation

$$W^{(i+1)} = W^{(i)} + \eta * (Y_i - O_i) \cdot X_i$$

Step-5: Repeat the steps 3 and 4 'n' times (until all training samples are over)

METHOD-4: Cortes and Vapnik (1995) presented a training algorithm that maximizes the margin between the training patterns and the decision boundary in support vector machines, which is used to classify the given observation to one of the label classes as follows. Let x be the feature space. Let $y_i \in \{-1, 1\}$ denote its class label for binary response.

Step 1: Let $X = (X_1, X_2, \dots, X_p)$ be the feature vector and Y be the binary response variable. Let $T = \{(X_i, Y_i), i = 1, 2, \dots, n\}$ be the trained data set measured on feature space in 'p' dimensions and 'm' be the number of labeled classes that sample data set T is classified. Let X_0 be the sample observation to be classified to one of the classes.

Step-2: The decision boundary of a linear classifier can be expressed as $W \cdot X + b = 0$, where W and b are parameters of the model. The margin is given by two parallel hyperplanes that are separated by the maximum possible distance $2/\|w\|$ with no data points inside the margin. Then the hyperplanes based on training data set can be constructed as $W \cdot X_i + b \geq 1$ for X_i in the +1 class and $W \cdot X_i + b \leq -1$ for X_i in the -1 class.

Step-3: Consider the optimization model, Minimize: $Z = \frac{\|W\|^2}{2}$ Subject to the constraints $y_i(w \cdot \phi + b) \geq 1$. The Lagrangian function for the optimization problem is

$$L = \frac{1}{2} W \cdot W - \sum_{i=1}^N \alpha_i [y_i (W \cdot X_i + b) - 1]$$

Where the parameters λ_i are called the Lagrange multipliers. To minimize the Lagrangian function, we can obtain the Karush-Kuhn-Tucker conditions as

$$\frac{\partial L}{\partial W} = W - \sum_{i=1}^N \alpha_i y_i x_i = 0 \text{ and } \frac{\partial L}{\partial b} = \sum_{i=1}^N \alpha_i y_i = 0$$

$$\lambda_i \geq 0, \text{ and } \lambda_i [y_i (W \cdot X_i + b) - 1] = 0$$

Step 4: The dual formulation is maximizing over the Karush-Kuhn-Tucker multipliers $\lambda_i = (\lambda_1, \lambda_2, \dots, \lambda_N)$ the function:

$$L_D = \sum_{i=1}^l \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j x_i x_j$$

Step 5: The decision boundary can be expressed as follows

$$\sum_{i=1}^N \lambda_i y_i x_i x + b = 0$$

Step-6: $f(X_0) = \text{Evaluate sign}(W \cdot X_0 + b) = \text{sign}(\sum_{i=1}^N \alpha_i y_i x_i \cdot x_0 + b)$. X_0 can be classified as positive class if $f(x_0) = +1$, and negative class if $f(x_0) = -1$,

METHOD-5: Naive Bayes classifier algorithm is used to classify the given sample observation based on data set is presented below.

Step-1: Let $X = (X_1, X_2, \dots, X_p)$ be the feature vector and Y be the binary response variable. Let $T = \{(X_i, Y_i), i = 1, 2, \dots, n\}$ be the trained data set measured on feature space in 'p' dimensions and 'm' be the number of labeled classes C_1, C_2, \dots, C_m that sample data set T is classified. Let X_0 be the sample observation to be classified to one of the class.

Step-2: evaluate the Prior probabilities $P(C_i)$ and likelihood probabilities $P(C_i/X)$ based on the training data set. Evaluate the posterior probability

$$P(X/C_i) = \frac{P(C_i/X) \cdot P(C_i)}{P(C_i)} \text{ for each class } C_i.$$

Step-3: The object or instance 'x' is allocated to class C_i iff

$$P(X/C_i) \cdot P(C_i) > P(X/C_j) \cdot P(C_j) \text{ for } 1 \leq j \leq m, j \neq i.$$

the predicted class label is the class C_i for which $P(X/C_i) \cdot P(C_i)$ is the maximum.

METHOD-6: The Logistic regression model can be used as a classifier is presented below.

Step-1: Let $X = (X_1, X_2, \dots, X_p)$ be the feature vector and Y be the binary response variable, Let $T = \{(X_i, Y_i), i = 1, 2, \dots, n\}$ be the trained data set measured on feature space in 'p' dimensions and 'm' be the number of labeled classes that sample data set T is classified. Let X_0 be the sample observation to be classified to one of the classes.

Step-2: Assume the functional relationship between the variables is

$$Y = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}$$

Estimate the vector of parameters $\beta = (\beta_0 \beta_1 \beta_2 \dots \beta_p)$ using least square method.

Step 3: Now optimize the β values by using $\beta_j = \beta_j + \alpha * (y - p) * x_j$ where α is the learning rate, p is the prediction of the model and x_j is the input value

Step 4: X_0 can be classified based on the decision rule: 0 if $p < 0.5$; and 1 if $p > 0.5$

III. COMPARISON OF VARIOUS CLASSIFICATION TECHNIQUES

In this section an attempt is made to compare the various classification methods considered in section 2.

Method	Merits	Demerits
Discriminant Model	<ol style="list-style-type: none"> It is a parametric supervised learning classifier, finds the vector which maximizes the separation between classes and builds a model. It uses the mean values of the classes and maximizes the distance between them. It minimizes the variance in the dataset by reducing the number of features. It can be used for both binary and multi class problems. 	<ol style="list-style-type: none"> It makes assumption about feature vector follows normal. It is based on the variance-covariance matrix of the data set.
K-Nearest Neighbor	<ol style="list-style-type: none"> It is a non-parametric supervised learning, but doesn't build any mathematical model but classifies the given sample point to one of the labeled classes based on 'k' nearest neighbors with majority voting. It uses the distance metric, Euclidean, Manhattan, Minkowski etc. It is a lazy learning and the accuracy is used to find the number of nearest neighbors. Root means square error is used to find the number of nearest neighbors. 	<ol style="list-style-type: none"> Choosing the value of k is difficult task. Computational cost is high for high dimensional data. It doesn't perform well on imbalanced dataset.
Perceptron	<ol style="list-style-type: none"> A Perceptron maps an m-dimensional input vector, onto a n-dimensional output vector. A distinct feature of a Perceptron is that the weights are not pre-calculated as in a McCulloch-Pitts model of neuron but are adjusted by an iterative process. Activation function is used to convert the input vector into a useful output. It doesn't make assumptions about input feature vector. A Multilayer perceptron can be used to solve complex non-linear problems and as the number of hidden layers increased the accuracy also increased. 	<ol style="list-style-type: none"> Training a multilayer perceptron is usually time consuming. It is difficult to predict how much the dependent variable affects each independent variable.
Support Vector Machine	<ol style="list-style-type: none"> It is a supervised optimization learning technique with a quadratic objective function and linear constraints. It uses nonlinear programming approach for its solvation. (Through Lagrange function) It constructs the optimal hyper-plane that classifies the original feature space into two sets by minimizing the squared norm of the separating hyper-plane. Number of support vectors will be identified. The hyperplane is chosen in such a way that it maximizes the margin, which is the distance between the hyperplane and support vectors. 	<ol style="list-style-type: none"> It is computationally expensive for large datasets due to quadratic optimization problem. Choice of kernel greatly affect the performance of SVM. Tuning the hyper parameters is difficult.
Naïve bayes Classifier	<ol style="list-style-type: none"> It is a supervised nonparametric probabilistic classifier works based on Bayes theorem. It calculates posterior probabilities based on the prior information. It assumes the independence of feature vector. It doesn't build any mathematical model rather than posterior probabilities. 	<ol style="list-style-type: none"> It assumes all the features are independent but it fails sometimes. Not suitable for continuous feature vectors.
Logistic	<ol style="list-style-type: none"> It is a non-parametric supervised learning classification technique. It constructs the linear decision boundary. Sigmoid function $f(x) = \frac{1}{1+e^{-x}}$ is used to convert the 	<ol style="list-style-type: none"> It assumes the linearity between dependent and independent variables. Due to linear decision boundary non-linear problems cannot be solved.

<p>predictions into probabilities.</p> <p>4. It provides the difference in the percentage of dependent variable and provides the rank of individual variable according to its importance.</p> <p>5. It is used to classify the low dimensional data having nonlinear boundaries,</p>	<p>3. Number of observations should be greater than the number of feature vectors otherwise it leads to overfitting.</p>
--	--

EXAMPLE 3.1: A data set of 13,444 customers of a national bank with a feature vector $\mathbf{X} = (X_1, X_2, \dots, X_{12})$ of 12 attributes under study with a categorical response variable Y of the status of credit-card is considered. For the classification, 80% (10,755) of the data is used for training the model and 20% (2,689) of the data set is used for testing the model and its accuracy.

Variable	Variable description	Label
Y	Credit card status	Credit card status
X ₁	Age	Age
X ₂	Months living at current address	ACADMOS
X ₃	1+No of dependents	ADEPCNT
X ₄	No. of Major Derogatory reports	MAJORDRG
X ₅	No. of Minor Derogatory reports	MINORDRG
X ₆	Own-rent	Own-Rent
X ₇	Income	Income
X ₈	Self-employed	Self-employed
X ₉	Income divided by number of dependents	INCPER
X ₁₀	Ratio of monthly credit card expenditure to yearly income	EXP_INC
X ₁₁	Average monthly credit card expenditure	SPEND
X ₁₂	Log of Spending	LOGSPEND

Sample data Set:

Y	X ₁	X ₂	X ₃	X ₄	X ₅	X ₆	X ₇	X ₈	X ₉	X ₁₀	X ₁₁	X ₁₂
...
...
1	27.66667	16	1	0	0	0	1650	0	23900	0.0953847	157.384808	5.0586938
0	33.75	18	0	0	0	1	1833.333	1	31000	3.87E-04	17.5041667	3.7810954
1	25.91667	54	0	1	1	1	1918	0	23016	0.169264	324.64833	5.7827425
1	38.58333	24	1	0	0	1	5000	0	30000	0.0202707	101.353331	4.6186127
0	42.41667	2	3	0	0	0	2916.667	0	8750	3.43E-04	69.3808347	4.6803215
1	43.25	118	4	0	1	1	3333.333	0	9000	2.67E-04	0.8888889	-0.117783
0	42	36	2	0	0	1	1250	0	6600	6.06E-04	8	0.7731899
1	40.58333	36	0	0	0	1	3000	0	42000	0.042216	126.647852	4.8414104
1	28.83333	26	0	0	0	1	2000	0	24500	0.0387955	77.5910243	4.3514518
1	26.33333	24	0	0	0	1	866.6667	0	18720	0.1779535	154.22639	5.0384216
...
...

Comparison with respect to the existence of the model, confusion matrix, accuracy, mis-classification rate evaluated using R- programming are summarized in the following table.

Table 3.4: A comparison w.r.t Model, Confusion matrix, Accuracy, Mis-classification rate

Classification Method	Classification Model	Confusion matrix	Accuracy	Rate of Mis-classification
Linear Discriminant	Linear Model	$\begin{pmatrix} 287 & 103 \\ 291 & 2008 \end{pmatrix}$	0.8535	0.1465
K-Nearest Neighbour	No model	$\begin{pmatrix} 348 & 90 \\ 242 & 2009 \end{pmatrix}$	0.8765	0.1235
Perceptron Learning	No model	$\begin{pmatrix} 562 & 6 \\ 8 & 2113 \end{pmatrix}$	0.9948	0.0052
Support Vector Machine	Linear model	$\begin{pmatrix} 576 & 18 \\ 18 & 2077 \end{pmatrix}$	0.9866	0.0134
Naïve Bayes Classifier	Probabilities	$\begin{pmatrix} 572 & 108 \\ 6 & 2003 \end{pmatrix}$	0.9576	0.0424
Logistic Regression	Logistic model	$\begin{pmatrix} 580 & 12 \\ 11 & 2086 \end{pmatrix}$	0.9914	0.0086

IV. CONCLUSION

It can observe that some classification methods provide mathematical linear functions to classify the given sample observations. It can be observed that for the chosen data set perceptron neural network learning technique provides 99.48% accuracy for the testing data set and logistic regression provides 99.14% accuracy.

REFERENCES

- [1]. Cortes, C. and Vapnik, V. (1995): "Support vector networks", Machine Learning, Vol 20, pp 273-297.
- [2]. Cover, T.M.and Hart, P. E. (1967): "Nearest neighbor pattern classification", IEEE Transactions on Information Theory, Vol 13 (1): pp 21-27.
- [3]. Fisher R.A. (1936): "The Use of Multiple Measurements in Taxonomic Problems" (PDF). *Annals of Eugenics*. Vol 7(2), pp 179-188.
- [4]. Rosenblatt, F. (1957): The Perceptron, a Perceiving and Recognizing Automaton Project Para. Cornell Aeronautical Laboratory. Technical report.
- [5]. Zurada J.M. (1994): Introduction to Artificial Neural Systems, Jaico Publishing house.
- [6]. Johnson R. A. and Wichern D.W. (2012): Applied Multivariate Statistical Analysis, PHI, Eastern Economic Edition, 6th edition.
- [7]. Draper N.R and Smith H. (1998): Applied Regression Analysis, 3rd Edition, John Wiley Publications.