

Handwritten Digit Classification Using Support Vector Machine Algorithm

Dr.R.V.Ramana Chary¹, Vasundhara Rao², S.V.S.Harshitha³

¹Professor and Associate Head, Department of Information Technology, B.V. Raju Institute of Technology (UGC Autonomous) Vishnupur, Narsapur, Medak (Dist.) - 502313

^{2,3}Department of Information Technology, B.V. Raju Institute of Technology (UGC Autonomous) Vishnupur, Narsapur, Medak (Dist.) - 502313

Abstract- The process of converting handwritten digits into digital format can be challenging because of many variations that can occur due to variation in size and orientation. They differ from person to person. There is also a possibility where a person can write a single digit in various styles. To a human, it is a very difficult task to determine the digits. This is where we can use machine learning approach to solve this problem.

Machine Learning is an invaluable tool when used to learn from large datasets, identify patterns and relationships and make decisions. One such machine learning algorithm is Support Vector Machine Algorithm. SVM is a classification algorithm which classifies data into categories based on patterns. In the proposed work, we developed a model using SVM which should correctly classify the handwritten digits from 0-9 based on the pixel values. We are using MNIST Dataset for classification.

Index Terms—Handwritten Digits, Machine Learning, MNIST Dataset, Support Vector Machine.

I. INTRODUCTION

Machine Learning is defined as the practice of teaching computers to learn from data using programming. It is done by using algorithms and statistical models to enable the computers to learn from data and make predictions based on the learning. It is rapidly evolving, driven by technological advancements. It is used in various fields such as health and finance. As more data becomes available and algorithms become more advanced, we can expect to see growth in areas such as natural language processing, robotics and autonomous systems.

One of the common application of machine learning is digit classification. The aim is to correctly identify handwritten digits based on their pixel values. Digit

Classification is used in various fields such as computer vision, handwriting recognition, and image processing. It can also be used in banking and financial institutions for automatic check processing and digitizing documents. Thus, we can say that digit classification is a powerful tool and can be used for myriad of fields.

Support Vector Machine or SVM Algorithm is a supervised machine learning algorithm used for classification as well as regression tasks. The primary idea behind SVM algorithm is to find the hyperplane that separates different classes. The margin is the distance between the hyperplane and nearest data points of each class. The objective of SVM is to find a hyperplane that maximizes the margin. The data points near the hyperplane are called support vectors which play a prominent role in determining the position and alignment of the hyperplane.

Kernels in SVM are the mathematical functions which take data and transform into the required form. Different kernels are used based on the type of data. When the data is linearly separable, we use linear kernel. Linearly separable data means a single line can be used to separate the data. When the data is not linearly separable i.e., we cannot find a line that separates the classes. In digit classification, the data is linearly inseparable. Therefore, we use non-linear kernel named radial basis function (rbf) kernel for handwritten digit classification.

Handwritten digit classification is a task whose goal is to identify handwritten images using computer algorithms. We use MNIST database for the classification.

The MNIST database (Modified National Institute of Standards and Technology database) is a collection of

handwritten digits. It is used as a benchmark dataset in machine learning. It contains thousands of 28 pixels by 28 pixels grayscale images of the handwritten digits. Each digit is labeled with corresponding numbers from 0-9. The dataset is in csv (comma separated values) file format. This dataset is used to train the machine learning model using support vector machine algorithm.

II. IMPLEMENTATION

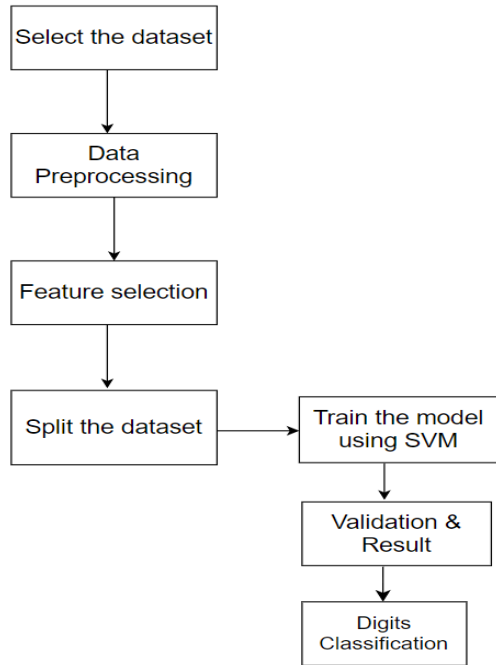


Figure 1: Steps in SVM Classification

A. Obtain the dataset

To implement any machine learning algorithm, we need data. For this, we need to obtain appropriate dataset. For the handwritten digit classification, we are using MNIST dataset. Here, we are using this dataset in the form of csv (comma separated values) file which consists of labels from 0-9 and corresponding pixel values.

```

df=pd.read_csv("/content/digit_svm.csv")
[4] df
   label  pixel0  pixel1  pixel2  pixel3  pixel4  pixel5  pixel6  pixel7  pixel8
0      1      0      0      0      0      0      0      0      0      0
1      0      0      0      0      0      0      0      0      0      0
2      1      0      0      0      0      0      0      0      0      0
3      4      0      0      0      0      0      0      0      0      0
4      0      0      0      0      0      0      0      0      0      0
...     ...     ...     ...     ...     ...     ...     ...     ...     ...
41995  0      0      0      0      0      0      0      0      0      0
  
```

Figure 2: Importing the dataset in Google Colaboratory

B. Data Preprocessing

In this step, we use various steps to prepare the data before applying machine learning algorithms. It includes removing errors and inconsistencies, combining data from various sources etc. With this process, we can improve the accuracy of our model and get better insights.

C. Feature Selection

It is a process of selecting a smaller set of important features in the dataset, which are essential for the task. In the dataset we used, we have pixel values as the features and hence all the features are considered important.

D. Splitting the dataset

It is an important step in machine learning. The idea is to divide the dataset into a training set and a validation set. The training set is used to train the machine learning model. We take the training set as 80% of the entire dataset.

Validation dataset is the subset of the original dataset. It is used to assess the performance of the machine learning model. We take the validating dataset as 20% of the entire dataset. While there are various ways of splitting, we used a random split. In this method, the data is randomly divided into train and test sets

```

[12] from sklearn.model_selection import train_test_split
[13] x_train,x_test,y_train,y_test= train_test_split(x,y,test_size=0.2,random_state=0)
[14] from sklearn.svm import SVC
  
```

Figure 3: Splitting dataset using Scikit-learn

E. Train the model

We use the training dataset to train the machine learning model. In digit classification, the data is nonlinear. Therefore, we use rbf kernel in SVM. It can be imported from Scikit-learn library.

```

[14] from sklearn.svm import SVC
[15] clf=SVC(kernel='rbf') #default kernel is rbf
[16] clf.fit(x_train,y_train) #fit the model
SVC()
  
```

Figure 4: Training model using SVM with rbf kernel

F. Validating the model

Using the validation dataset, we assess the performance of the model. Based on the result, the model can be adjusted.

G. Accuracy Results

For the model, we can predict the accuracy of the model using accuracy score function from scikit-learn library. It defines how many exact matches did we get i.e., how many correct values did it predict corresponding to the true values. We also used confusion matrix to predict how many true positives and true negatives were there.

```

SVM(rbf).ipynb
File Edit View Insert Runtime Tools Help All changes saved
+ Code + Text
[19] from sklearn.metrics import accuracy_score
accuracy_score(y_test, pred_y)
0.9747619047619047

from sklearn.metrics import confusion_matrix
confusion_matrix(y_test, pred_y)
array([[801, 0, 1, 0, 1, 3, 3, 0, 4, 0],
       [0, 953, 3, 2, 0, 0, 0, 1, 2, 0],
       [2, 1, 843, 2, 1, 0, 0, 5, 6, 0],
       [1, 1, 11, 821, 0, 12, 2, 4, 9, 2],
       [1, 3, 2, 0, 309, 1, 4, 2, 0, 14],
       [2, 0, 0, 4, 0, 743, 6, 0, 1, 0],
       [2, 0, 0, 0, 2, 5, 831, 0, 1, 0],
       [0, 2, 6, 1, 6, 2, 0, 873, 3, 0],
       [2, 3, 0, 3, 3, 5, 4, 0, 747, 1],
       [5, 2, 2, 4, 9, 1, 0, 11, 2, 770]])
    
```

Figure 5: Accuracy score and Confusion matrix

III. RESULTS & CONCLUSION

This model is efficient and has an accuracy of 97%. Therefore, we can say that rbf kernel is suitable for digit classification. Also, this model will make the classification much easier as it reduces human efforts. The values in the diagonal of the confusion matrix are the values which are predicted correctly i.e., predicted value is equal to the true value.

IV. FUTURE WORK

The model was implemented to only classify the images. Further we can add modules that predicts the type of digits when an image is given as input using advanced techniques in machine Learning.

REFERENCE

[1] Hafeez Ahamed, Syed Md Ishraq Alam. "SVM Based Real Time Hand-Written Digit Recognition". International Conference on Engineering Research and Education School of Applied sciences & technology, at SUST, Sylhet (January 2019).

[2] <https://scikit-learn.org/stable/modules/svm.html>