

Identification of Digits from Sign Language Images

Aprendizagem Automática – Professor: Pétia Georgieva

José Santos, 98279
DETI
Universidade de Aveiro

Henrique Sousa, 98324
DETI
Universidade de Aveiro

Abstract—The purpose of this work is to implement and compare machine learning models capable of identifying digits from sign language images. In this paper, we tried to obtain a good result with the models using a dataset provided by Kaggle. Some changes are discussed, based on the work of others that positively affected our work.

Index Terms—Sign language recognition, Digit recognition, Machine learning

I. INTRODUCTION

In recent years, there has been growing interest in using machine learning to develop computer vision systems capable of recognizing sign language gestures. Such systems could be used to improve communication between hearing and non-hearing individuals, as well as to facilitate the development of new technologies for the deaf and hard-of-hearing community.

In this paper, we present a novel approach to digit recognition from sign language images using machine learning. We explore several different models, including neural networks, support vector machines, and decision trees, and compare their performance on a dataset of sign language images. We also investigate the impact of some preprocessing techniques.

II. STATE OF THE ART

Over the past few years, there have been several studies and projects focused on recognition from sign language images using machine learning techniques. Many of these approaches have utilized deep learning methods such as convolutional neural networks (CNNs), which have been shown to be effective in image recognition tasks.

TODO: Add some references to the state of the art

III. DATASET ANALYSIS

For the development of this work, we used a dataset provided by Kaggle that contains 2062 images of sign language digits. The dataset is well balanced, with a minimum of 204 examples for a label and a maximum of 208, ensuring that every label has a similar number of examples (Figure 1). The dataset includes data for the digits 0 to 9, resulting in 10 labels in total as we can see on the figure 2.

The images are provided in a .npy format, but we found it easier to work with the raw images to manipulate them and apply preprocessing techniques to improve the model's performance. The dataset's size and balanced distribution make it an ideal choice for training and evaluating machine learning models for digit recognition from sign language images.

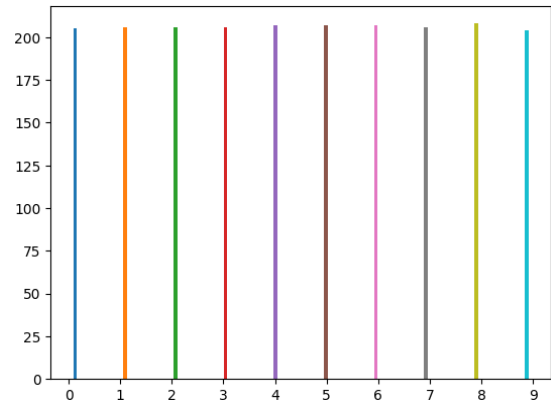


Fig. 1. Dataset balanced distribution

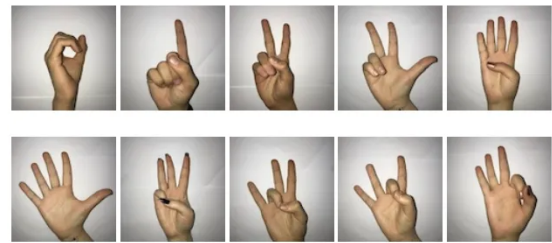


Fig. 2. Dataset label examples

IV. DATASET PREPROCESSING

A. Image Augmentation

The first thing we did to the images was to augment them to create more data for model training. We used *imgaug.augmenters* to perform the augmentation. The following changes were made to the images:

- **Rotate** - Rotate the image by a random angle between -20 and 20 degrees.
- **Gaussian noise** - Add gaussian noise with standard deviation of 0 to 0.05×255
- **Gamma contrast** - Change the contrast of the image by a random factor between 0.5 and 1.5.

For every image in the dataset, we created 5 new images with the above changes. Due to the small size of the dataset, we decided to perform the augmentation offline, before splitting the dataset into training and test data.

B. Image Preprocessing

We also resized the images to 50x50 pixels (down from the original size of 64x64), as we found that this size was sufficient for the models to achieve good results. Then, we converted the images to grayscale and flattened them, as we found that this improved the performance of the models.

Finally, we split the dataset into training and test data. We used an 80/20 split, with 80% of the data being used for training and 20% for testing. Resulting in 9897 total images for the training data and 2475 for the test data.

V. MODELS

In order to find the best training model, we tested several different models, including neural networks, support vector machines, and decision trees. We also tested different preprocessing techniques to see if they had any impact on the model's performance as discussed before. The two evaluation metrics used to measure the performance of the models are Accuracy and F1 Score. Accuracy is the proportion of correct predictions over the total number of predictions, while F1 Score considers both precision and recall of the model's predictions.

First, we tried a lot of models with their default parameters to understand which ones had value and were worth to be improved. These were the first results:

Model	Accuracy	F1 Score
Logistic Regression	0.750	0.749
Decision Tree Classifier	0.631	0.632
Random Forest Classifier	0.876	0.876
Naive Bayes	0.502	0.506
Support Vector Machines	0.888	0.888
Neural Networks (Multilayer Perceptron Classifier)	0.092	0.015

TABLE I
MODELS' PERFORMANCE WITH DEFAULT PARAMETERS

Among all the models, Support Vector Machines has the highest accuracy and F1 Score, both of which are 0.888. Random Forest Classifier also shows a good performance with accuracy and F1 Score of 0.876. Logistic Regression comes in third place with accuracy and F1 Score of 0.750 and 0.749, respectively. The other models have lower performance, so for the next step, we will skip the Decision Tree Classifier and Naive Bayes. We will still try to find ideal parameters for the MLP Classifier because the model ended up predicting the same label for each example. This is probably because the model was not trained properly, so we will give it a chance and try to improve it.

A. Hyperparameter Tuning & Cross-Validation

Hyperparameter tuning is the process of finding the best values for the parameters that result in the best performance of the model on the given dataset. The purpose of cross-validation is to estimate how well a machine learning model will generalize to new data. By evaluating the model on multiple subsets of the data, it is less likely that the model's performance is biased towards a particular subset of the data.

Cross-validation is commonly used in machine learning to evaluate the performance of different models or to compare the performance of different hyperparameters. To implement these methods and get the best model performance, we created a function that given a model, the parameters to test and the dataset, it will try the different combinations of parameters and return the best model and its performance. The function also performs cross-validation on the training data to evaluate the model's performance. The following code shows the function's implementation:

```
def hyperparameters(model, params, X, y):  
    model = GridSearchCV(model, params,  
                          scoring="accuracy")  
    model.fit(X, y)  
    print("Best Parameters :")  
    print(model.best_params)
```

These function is going to be applied in the models with different parameters to find the best model for each one.

B. Support Vector Machines

Support Vector Machines (SVM) proved to be the best-performing model for the classification task, with an initial accuracy and F1 Score of 88.8% using default parameters. In an effort to improve the performance of the SVM model, we experimented with different values for the C, kernel, degree, and gamma parameters, as shown in Table V-B.

Parameter	Values	Best Value
C	[0.001, 0.01, 0.1, 1, 10, 100, 1000]	100
kernel	[linear, poly, rbf, sigmoid]	rbf
degree	[2, 3, 4, 5]	2
gamma	[scale, auto]	scale

TABLE II
SVM PARAMETERS

After running several iterations, we were able to identify the optimal parameter values, which resulted in a significant improvement in performance. The best-performing parameters were:

- **C** = 100
- **kernel** = rbf
- **degree** = 2
- **gamma** = scale

The SVM model with these parameters achieved an accuracy and F1 Score of 95%, which represents a significant improvement over the default parameters. We can better visualize the performance of the SVM model by looking at the confusion matrix, as shown in Figure 3.

As we can see in the confusion matrix, the SVM model can correctly classify most of the images. Some classes are misclassified. This is probably because the images of these classes are very similar to each other, so it is not surprising that the model has a hard time classifying them correctly.

Looking at the precision scores, which measure the proportion of true positives among all samples predicted as positive, we can see that the model performed well for most classes,

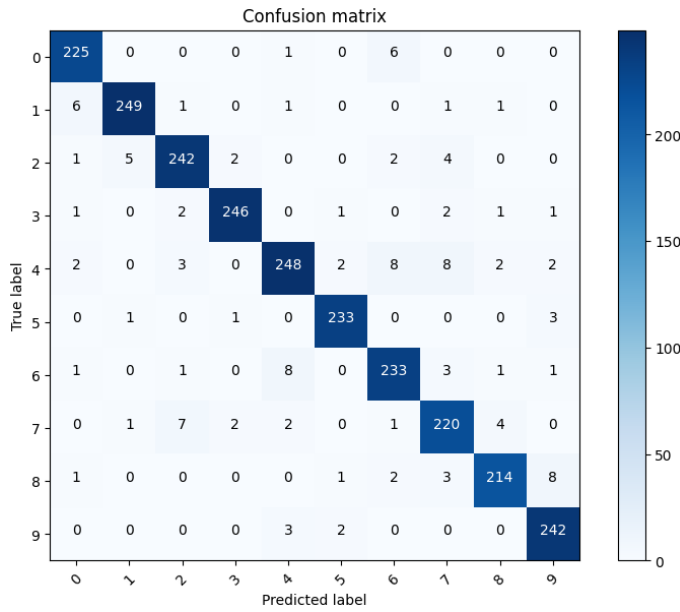


Fig. 3. SVM Confusion Matrix

Label	Precision	Recall	F1 Score	Occurrences
0	0.95	0.97	0.96	232
1	0.97	0.96	0.97	259
2	0.95	0.95	0.95	256
3	0.98	0.97	0.97	254
4	0.94	0.90	0.92	275
5	0.97	0.98	0.98	238
6	0.92	0.94	0.93	248
7	0.91	0.93	0.92	237
8	0.96	0.93	0.95	229
9	0.94	0.98	0.96	247

TABLE III
SVM CLASSIFICATION REPORT

with scores ranging from 0.91 to 0.98. This indicates that when the model predicted a certain class, it was usually correct.

The recall scores, which measure the proportion of true positives among all actual positive samples, are also high for most classes, ranging from 0.90 to 0.98. This suggests that the model was able to identify most instances of each class correctly.

In addition to high precision and recall scores, the F1-score, which balances these two metrics, also indicates strong performance for most classes, with scores ranging from 0.92 to 0.98. This suggests that the model achieved a good trade-off between precision and recall, and is able to correctly classify samples while minimizing false positives and false negatives. Overall, these results suggest that the classification model was effective at distinguishing between the different classes in the dataset, with high levels of precision, recall, and F1-score for most classes.

VI. CONCLUSION

VII. REFERENCES

REFERENCES

- [1] Akanksha Telagamsetty, Sign Language Digits Classification <https://medium.com/analytics-vidhya/sign-language-classification-64fe8ad0fc2c>