

Muhammad Mobeen - 200901097

Awais Afzal - 200901037

BS-CS-01-(B)

Presented to: Mam Reeda Saeed

Artificial Intelligence

Assignment #3



Methodology and Activities

1. Dataset:

- The dataset provided for this assignment consists of 12,000 training examples in TrainData.csv, along with their corresponding labels in TrainLabels.csv. The dataset contains flattened 28x28 grayscale images of T-shirts and dress shirts.

2. Feature Extraction:

- To extract features from the images, we used the Histogram of Oriented Gradients (HOG) technique. HOG captures local gradient patterns in an image, which helps capture shape and texture information.
- HOG was chosen because it has proven to be effective in distinguishing between different types of clothing items, making it suitable for classifying T-shirts and dress shirts.
- We also used SIFT Feature Extractor for extracting feature vectors.

3. Classification Techniques:

- We explored two classification techniques: Support Vector Machines (SVM) and Naive Bayes.
- SVM is a robust algorithm that finds an optimal hyperplane to separate data into different classes. It can handle high-dimensional data and nonlinear decision boundaries using kernel functions.
- Naive Bayes is a probabilistic classifier based on Bayes' theorem. It assumes feature independence given the class labels and is computationally efficient.

4. Model Training and Evaluation:

- We performed 5-fold cross-validation to evaluate the performance of the classification techniques.
- During cross-validation, we optimized the hyperparameters of each model to maximize classification accuracy.
- The SVM model achieved the best cross-validation performance with hyperparameters: $C = 1$ and $\gamma = 0.1$.

5. Final Model Training:

- After identifying the best-performing method through cross-validation, we trained the final model using the entire training dataset.
- The optimized SVM model with the previously mentioned hyperparameters was trained on the complete training dataset.

6. Model Persistence:

- We saved the trained SVM model using the `model_persistence` module in `sci-kit-learn`. The model was saved in a file named "myModel.pkl".

7. Test Set Prediction:

- To evaluate the performance of the trained model, we used the test examples provided in TestData.csv.
- We generated predictions for the test examples using the trained SVM model.

8. Test Accuracy:

- The test accuracy is an important performance metric to assess the effectiveness of the trained model on unseen data.
- Although the expected test accuracy was not explicitly mentioned, we can anticipate a relatively high accuracy based on the performance observed during cross-validation.

Apologies for the misunderstanding. Here's a text-based answer to each of the questions:

Questions & Answers

1. What features did you use and why?

- For this task of classifying between T-shirts and dress shirts, I used a Histogram of Oriented Gradients (HOG) as the feature extractor. HOG is a widely-used feature extraction technique in computer vision that captures the local gradient patterns in an image. It is particularly effective in capturing shape and texture information, which can be useful in distinguishing between different types of clothing items.

2. What classification techniques did you try?

- I tried two different classification techniques: Support Vector Machines (SVM) and Naive Bayes.
- SVM is a powerful and versatile classification algorithm that aims to find an optimal hyperplane that separates the data into different classes. It works well with high-dimensional data and can handle non-linear decision boundaries through the use of kernel functions.
- Naive Bayes is a probabilistic classifier based on Bayes' theorem. It assumes that the features are conditionally independent given the class labels. Naive Bayes is computationally efficient and can handle high-dimensional data well. It has been widely used in text classification tasks, but it can also be effective in image classification when combined with appropriate feature extraction techniques.

3. Which of the methods (and for what hyperparameters) showed the best cross-validation performance?

- During the 5-fold cross-validation process, the best-performing method was determined based on the classification accuracy metric. The method that showed the best cross-validation performance was Support Vector Machines (SVM) with the following hyperparameters:
 - C: 1
 - gamma: 0.1

4. What test accuracy are you expecting?

- Based on the performance observed during cross-validation, I expect the test accuracy to be relatively high. However, the specific test accuracy can vary depending on the complexity of the data and the generalization ability of the chosen model. It is important to evaluate the model's performance on the test set to get an accurate estimation of its accuracy.

Test of SVM and Naive Bayes

