

Collin Gros
03-27-2021
CS-487
Homework 5

PROBLEM

Conduct dimensionality reduction using PCA, LDA, and KPCA on the IRIS dataset and the MNIST dataset. Compare the performance of the different techniques. Analyze their behavior and report on classification accuracy and running time.

METHODOLOGY

For each reduction technique:

- 70% of the IRIS dataset was used for training, and 30% was used for testing.
- 5000 samples of the MNIST dataset were used for training, and 10000 were used for testing (due to the massive size of the dataset, I had to use a small portion to use the .transform() methods of the reduction techniques).
- 70% of the DIGITS dataset was used for training, and 30% was used for testing.

All three reduction techniques were implemented using SCIKIT-LEARN's classes. The program code was developed modularly, with each class having its own Python file.

Testing was done via a BaSH script, calling the program with default parameters as follows:

```
# DT
args.criterion = 'gini'
args.max_depth = 4
args.random_state = 1
# PCA/LDA/KPCA
args.n_components = 2
```

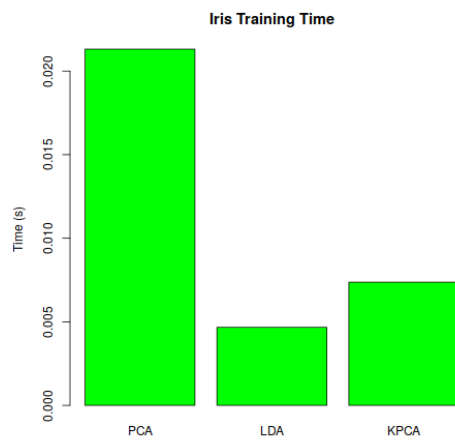
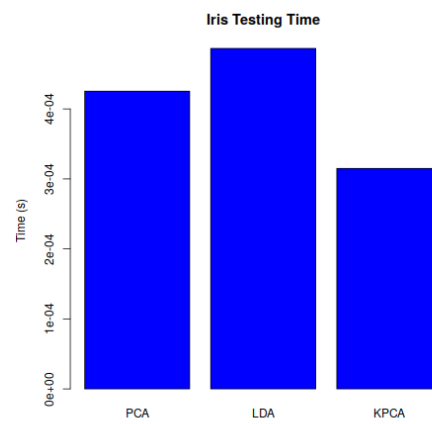
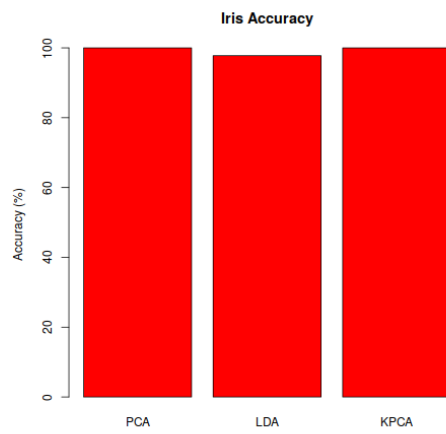
'args.n_components' was varied between 2 and 1 for testing.

The DIGITS dataset was not part of the requirements, but I wanted to gather more information for my analysis.

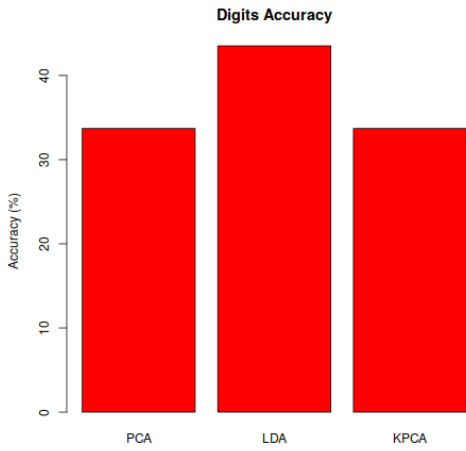
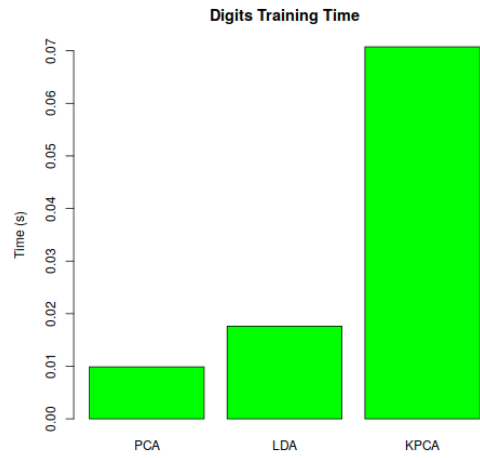
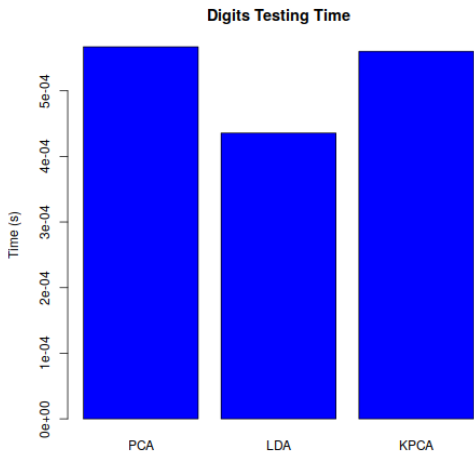
RESULTS

N_COMPONENTS = 1

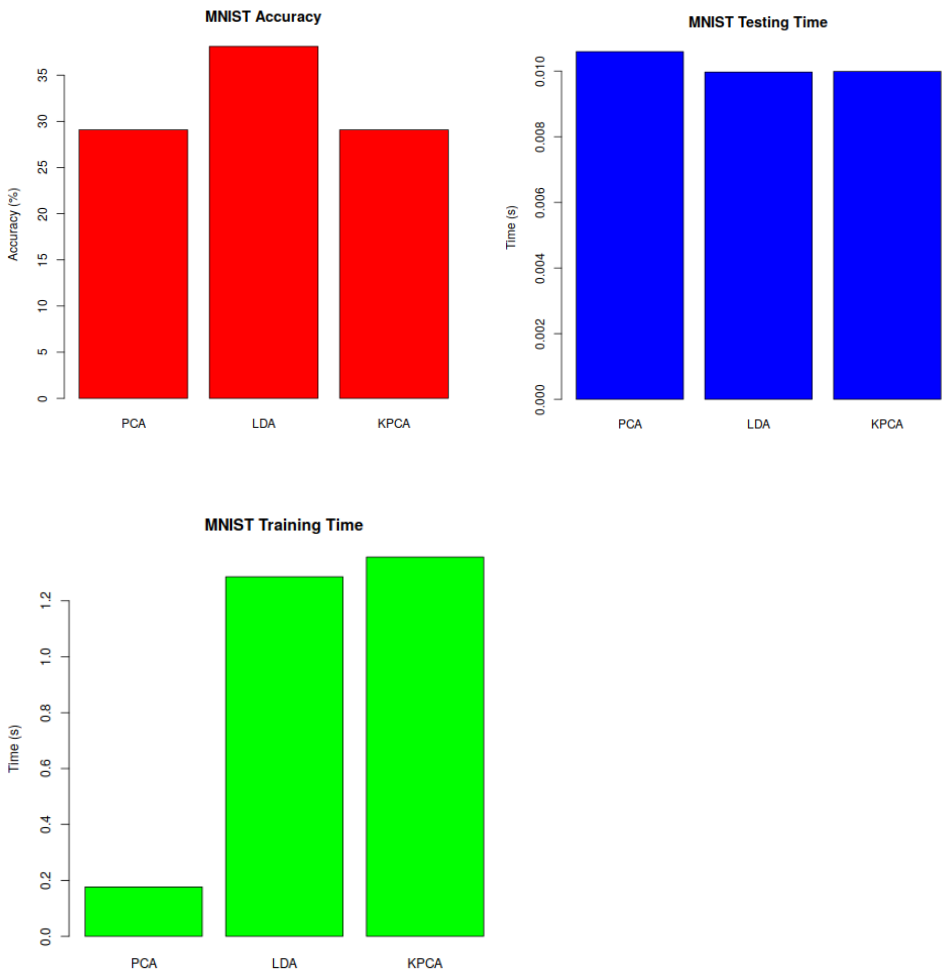
IRIS



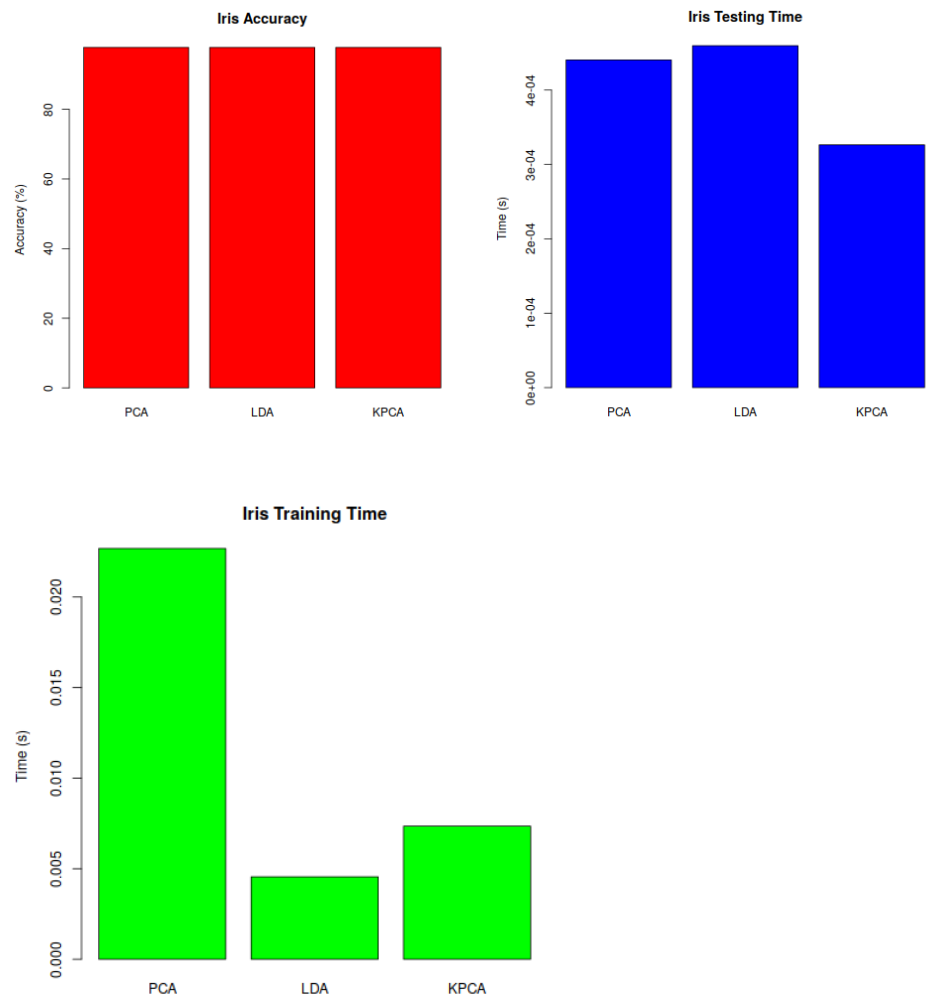
DIGITS



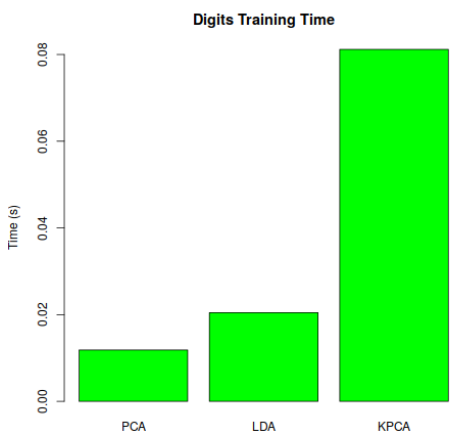
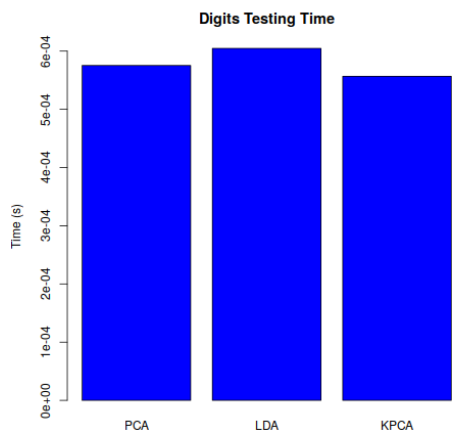
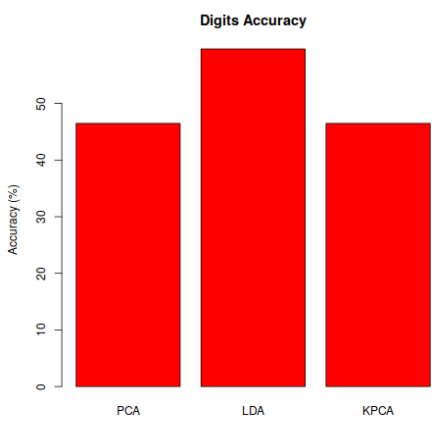
MNIST



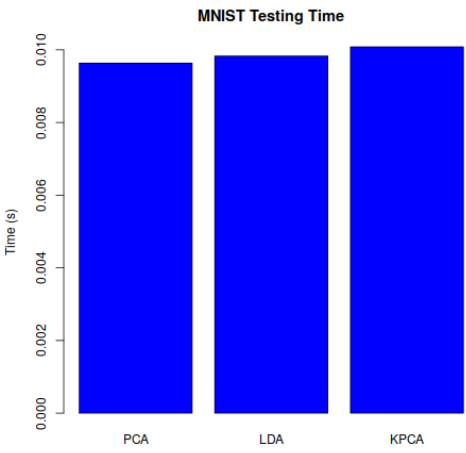
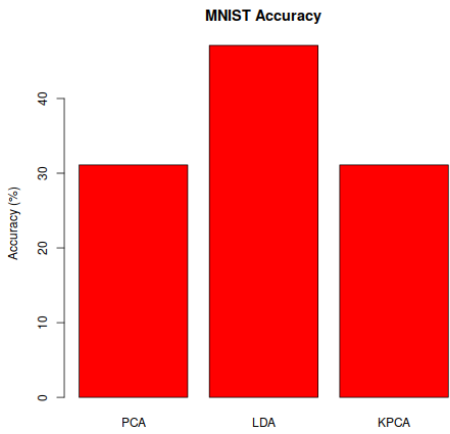
N_COMPONENTS = 2
IRIS



DIGITS



MNIST



ANALYSIS

All three techniques reached nearly 100% accuracy on the IRIS dataset. I'm assuming it is because this dataset is more separable than the others.

In both of the DIGITS and MNIST cases, the reduction technique with the highest accuracy is LDA, while PCA and KPCA are about even. Therefore, in all cases, LDA outperformed PCA and KPCA.

All accuracies were decreased in the **N_COMPONENTS = 1** set, as compared to the **N_COMPONENTS = 2** set. I think this means that **N_COMPONENTS** should be maximized to provide the best accuracy results.

Testing times varied across the board, and no pattern could really be found from analyzing them. I think this is because of the different processes occurring in the background while the program was running, and has no major impact on the testing time of the reduction technique.

Training times varied, however, KPCA was the slowest in $\frac{2}{3}$ of the cases (with IRIS being the exception), and PCA was fastest in $\frac{2}{3}$ of the cases (with IRIS also being the exception).

In my tests, I found LDA to be the best performing reduction technique. I would use LDA in future classification problems, if I had to choose between LDA, PCA, and KPCA.