

University of Waterloo
ECE 657A: Data and Knowledge Modeling and Analysis
Winter 2019

Assignment 1: Data Cleaning and Dimensionality Reduction

Due: February 16th, 2019 11:59pm

Overview

Assignment Type: done in groups of up to three students.

Hand in: One report (PDF) or python notebook per group, via the LEARN dropbox. Also submit the code / scripts needed to reproduce your work. (If you are submitting by PDF, if you don't know \LaTeX you should try to use it, it's good practice and it will make the project report easier)

Objective: To study how to apply some of the methods discussed in class on three datasets. The emphasis is on analysis and presentation of results not on code implemented or used. You can use libraries available in python, R or any other programs available to you. You need to mention explicitly the source with proper references.

Data sets

Available on LEARN, if you aren't registered yet they will be placed in the owncloud dropbox. The datasets are in csv format (.csv), if you are using python you can load them using the Scipy method `loadmat`.

Dataset A : This is a time-series dataset which is collected from a set of motion sensors for wearable activity recognition. The data is given in time order, with 19,000 samples and 81 features. Some missing values are denoted by Not Available (NA) and also some outliers are present. (note: The negative values are not outliers) This data is used to illustrate the data cleaning and preprocessing techniques. (File: [DataA.csv](#))

Dataset B : Handwritten digits of 0, 1, 2, 3, and 4 (5 classes). This dataset contains 2066 samples with 784 features corresponding to a 28 x 28 gray-scale (0-255) image of the digit, arranged in column-wise. This data is used to illustrate the difference between feature extraction methods. (File: [DataB.csv](#))

Questions

I. Data Cleaning and Preprocessing (for dataset A)

1. Detect any problems that need to be fixed in dataset A. Report such problems.
2. Fix the detected problems using some of the methods discussed in class.
3. Normalize the data using min-max and z-score normalization. Plot histograms of feature 9 and 24; compare and comment on the differences before and after normalization.

II. Feature Extraction (for dataset B)

1. Use PCA as a dimensionality reduction technique to the data, compute the eigenvectors and eigenvalues.
2. Plot a 2 dimensional representation of the data points based on the first and second principal components. Explain the results versus the known classes (display data points of each class with a different color).
3. Repeat step 2 for the 5th and 6th components. Comment on the result.
4. Use the Naive Bayes classifier to classify 8 sets of dimensionality reduced data (using the first 2, 4, 10, 30, 60, 200, 500, and all 784 PCA components). Plot the classification error for the 8 sets against the retained variance of each case.
5. As the class labels are already known, you can use the Linear Discriminant Analysis (LDA) to reduce the dimensionality, plot the data points using the first 2 LDA components (display data points of each class with a different color). Explain the results obtained in terms of the known classes. Compare with the results obtained by using PCA.

III. Nonlinear Dimensionality Reduction (for dataset B)

Apply the nonlinear dimensionality reduction methods Locally Linear Embedding (LLE) and ISOMAP to the dataset B, set the number of nearest neighbours to be 5, the projected low dimension to be 4.

1. Apply LLE to the images of digit '3' only. Visualize the original images by plotting the images corresponding to those instances on 2-D representations of the data based on the first and second components of LLE, see Figure for an example of what this looks like for random location of images on of the number 1-3. Describe qualitatively what kind of variations is captured.
2. Repeat step 1 using the ISOMAP method. Comment on the result. Does ISOMAP do better in some way? Are the patterns being found globally based or locally based?
3. Use the Naive Bayes classifier to classify the dataset based on the projected 4-dimension representations of the LLE and ISOMAP. Train your classifier by randomly selected 70% of data, and test with remained 30%. Retrain for multiple iterations (using different random partitions of the data) and use the average accuracy of multiple runs for your analysis. Justify why your number of iterations was sufficient. Based on the average accuracies compare their performance with PCA and LDA. Discuss the result.

Deliverables

For submitting your assignment please consider the following notes:

- Submit all of your work as one compressed file (.zip, .rar) named as Gx.zip or Gx.rar where "x" indicate your group number. (You will be able to see your group number on LEARN, if you have any question please contact: ifadakar@uwaterloo.ca)

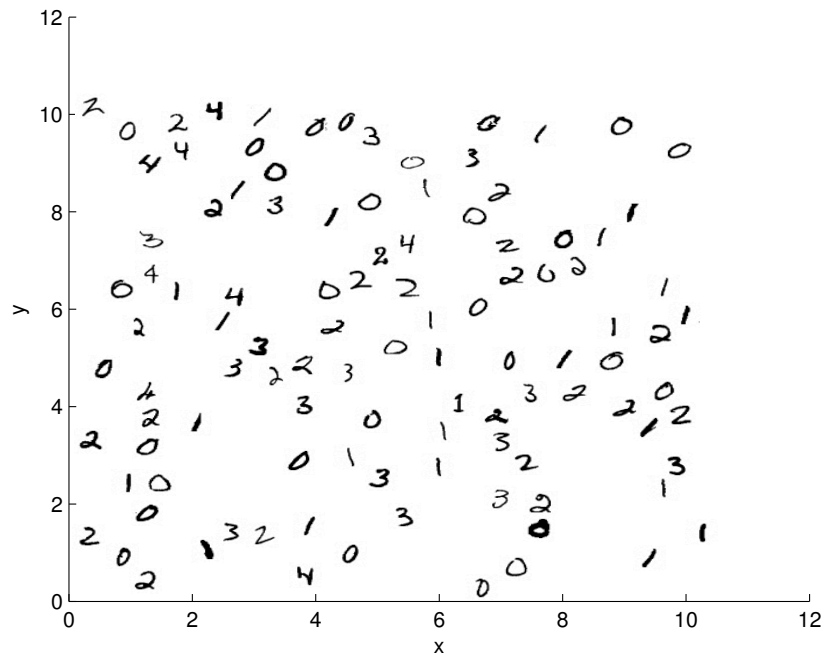


Figure 1: Using plotImages to plot pictures in a sample 2D space; Output of the sample code.

- Your compressed file should have all code, images, etc in addition to your report's document.
- Write a technical document as your report and submit its PDF format included it in your compressed file.
- Your report (.pdf file) should have the name and student number of all members of your group at the beginning and separated sections for the answer of each part of each question .
- Late submissions (up to 3 days) are accepted with penalty of 10% per day.
- All code should be clearly written and commented and be runnable on another system with just the data set files beside the code in the same folder.
- Do not upload the data set files.
- One member of each group should upload the report to your group's dropbox on Learn. Each member does not need to submit same version. The last version submitted will be the one which is graded.