

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

Assignment 2: Data Cleaning and Dimensionality Reduction

Due: March 25th, 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 L^AT_EX you should try to use it, it's good practice and it will make the project report easier)

Objective: To gain experience on the use of classification methods. The emphasis is on analysis and presentation of results not on code implemented or used. You can use libraries available in Python, R which are available to you. You need to mention explicitly the source of any other references used.

Data sets (available on the UW 'LEARN' system)

Data set A: The data is related with direct marketing campaigns of a Portuguese banking institution. The marketing campaigns were based on phone calls. Often, more than one contact to the same client was required, in order to access if the product (bank term deposit) would be ('yes') or not ('no') subscribed. Use dataset bank-additional.csv for this assignment. The classification goal is to predict if the client will subscribe (yes/no) a term deposit (variable y). You can also check the dataset at the following link:

<https://archive.ics.uci.edu/ml/datasets/bank+marketing>

Data set B: This data is the splice junctions on DNA sequences. The given dataset includes 2200 samples with 57 features. It is a binary class problem. The class labels are either +1 or -1, which is given in the last column. Parameter selection and classification tasks are conducted on this dataset.

Questions

Question I:[40] Revisiting HW4 Bank Classification with New Tools (for dataset A)

1. Load a sample dataset and perform some basic data preprocessing to fill out "unknowns", outliers or other invalid data. Explain what preprocessing was performed and why. Also, change categorical data into numerical features using `pandas.get_dummies` [5].
2. Divide data into train and test portions, justify your split decision [5].

3. Apply classification using Decision Trees (DT), Random Forests (RF) and Neural Networks (NN) and run using standard libraries in your language or choice. Indicate the classification properties (example: depth of tree, size of neural network, ensembles for RF) you have chosen and justify. Briefly describe the algorithms employed by the libraries you are using (example ID3, C4.5, C5.0 and CART for DT). Make sure to run the classification on 10 features or more [10].
4. Create a few plots of your model on the test data, two of the data dimensions at a time, indicating the predicted elements of each class using different colors or shapes. You may need to try plotting various pairs of dimensions to see which provide some interesting result. Be sure to label your axis and legend. Why is separation better on some plots than others [10]?
5. Produce a table with the true/false positive/negative metrics as well as accuracies. Compare the values using bar charts [5]. HINT: `classification_report` from `sklearn.metrics`
6. Provide a short explanation of the results you have shown and what it means. Which classification method performed better? Why? Contrast performance with classification from the previous homework and comment on the difference, if any [5].
7. **For Fun/Bonus:** attempt at least one method to tackle the discrepancy in the size of the classes (imbalanced data) [+5].

Question 2:[60] Parameter Selection and Classification (for dataset B)

Classify dataset B using four classifiers: k-NN, Support Vector Machine (with RBF kernel), Random Forests and simple Neural Networks (MLPs). The objective is to experiment with parameter selection in training classifiers and to compare the performance of these well-known classification methods.

1. Preprocess the given data using the Z-score normalization on the data. Justify the choice of Z-score normalization here, as opposed to min-max normalization. Why do you need normalization in general? Justify why you would normally split the test and training set randomly. What is the distribution of the +1,-1 classes in the dataset?[5]
2. Parameter Selection:
 - (a) For k-NN you need to evaluate the best value **k** to use. Using 5-fold cross validation on the training set evaluate k-NN on the values **k**=[1, 3, 5, 7, , 31]. The following link can be helpful:

https://scikit-learn.org/stable/modules/cross_validation.html

 Plot a figure that shows the relationship between the accuracy and the parameter **k**. Report the best **k** in terms of classification accuracy. Explain why you didnt evaluate directly on the test set [7.5].
 - (b) For the RBF kernel SVM, there are two parameters to be decided: the soft margin penalty term "**c**" and the kernel width parameter "**sigma**". Again use 5-fold cross validation on the training set to select the parameter "**c**" from the set [0.1, 0.5, 1, 2, 5, 10, 20, 50] and select the parameter "**sigma**" from the set [0.01, 0.05, 0.1, 0.5, 1, 2, 5, 10]. Report the best parameters in terms of classification accuracy including plotting the ROC curves [7.5].

3. Train (at least) six classifiers and report the results:
 - (a) Classify the test set using k-NN, SVM, Random Forests and Neural Networks. Use the chosen parameters from the parameter selection process in question 2 for k-NN and SVM. For the next two classifiers use the default setups listed at the end for Random Forests and Neural Networks [15].
 - (b) For the fifth and sixth classifiers, you should explore the parameters of the Random Forests and Neural Network models to devise your own classifier instance that does better than the other methods. For example, you could consider a deeper neural network with multiple layers, use different optimization/solver algorithms, you could modify the Random Forests using different parameter settings for depth and number of trees or enable boosting. Play around with options and choose a setting for RFs and NNs that performs better [10].
 - (c) Repeat each classification method 20 times by varying the split of the training-test set as in question 2-2. Report the average and standard deviation of classification performance on the test set regarding: accuracy, precision, recall, and F- Measure. Also report the training time and classification time of all the methods. Explain why the classification was repeated 20 times [5].
4. Comment on the obtained results, what are the benefits and weaknesses of each method on this dataset. How could this analysis help to make the choice of the right method to use for a dataset of this type in the future? [5]
5. If you had to remove 1 features from the dataset, which feature would you select to remove from the dataset and why? What would have happened if you did classification on two dimensions only? [5]

Notes

1. For classification methods of k-NN, Naive Bayes, SVM, and Decision trees and Random Forests (bagged decision trees), Python has implemented functions in Sklearn:
 - http://scikit-learn.org/stable/modules/neural_networks_supervised.html
 - <http://scikit-learn.org/stable/modules/tree.html>
 - <http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>
 - <http://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>
 - <http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html>
2. For simple Neural Networks in python, the sklearn package MLPClassifier can be helpful(section 1.17.2):

http://scikit-learn.org/stable/modules/neural_networks_supervised.html

3. Late submissions (up to 3 days) are accepted with penalty of 10% per day.

Default Parameters:

For easier comparison of results please try to use default parameters for the first four classifiers of question 2. For all other parts you can select your own parameters but you must justify your choices. Also, for tuning hyperparameters you can use `GridSearchCV`