

ELEN-0062: Introduction to Machine Learning

Project 3 - Human Activity Prediction

Maxime Goffart
180521

Olivier Joris
182113

Academic year 2021 - 2022

Contents

1	Introduction	2
2	Data pre-processing	2
2.1	Filling missing values	2
2.2	Feature selection	2
2.3	Implementation	3
3	Studied models	3
3.1	K-nearest neighbors	3
3.1.1	1-NN	3
3.1.2	25-NN	3
3.1.3	55-NN and 49-NN	3
3.1.4	Multiple K-nearest neighbors models	4
3.1.5	Mutiple K-nearest neighbors models and samples modification	5
3.1.6	Conclusion on K-nearest neighbors	5

1 Introduction

In this report, we will present and justify the different techniques we used for the third project of the course. The goal is to provide the reasoning behind our choices and to present the results obtained on Kaggle.

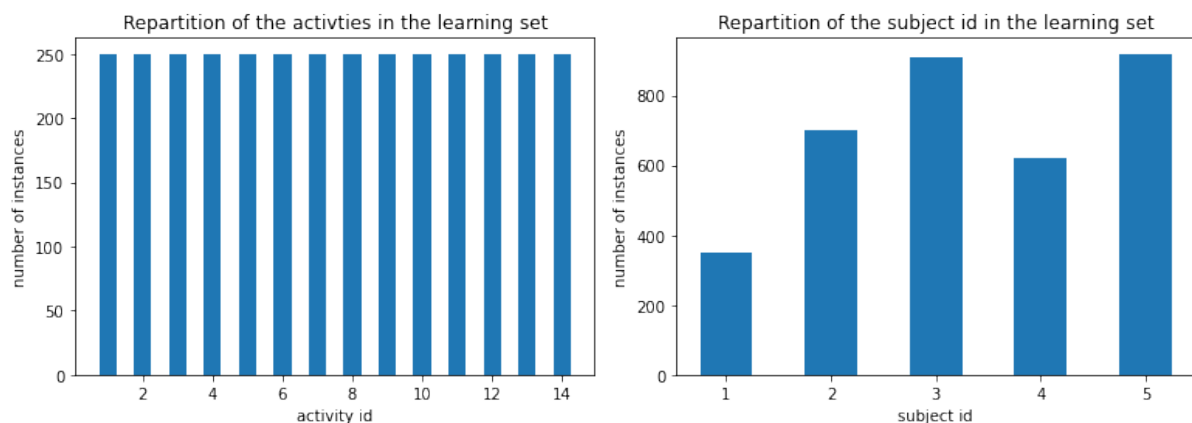
2 Data pre-processing

We started by analyzing the data in order to pre-process them. Indeed, pre-processing is important to have good results when doing classification. This section includes the observations we made about the dataset and different approaches we tried in order to solve the problems linked to the dataset.

2.1 Filling missing values

The learning set contained a certain number of missing values represented by the `-999999.99` value. We thus started by looking at several methods proposed by the scikit-learn library¹ to replace these missing values. We tried several of these methods and finally decided to use the `KNNImputer` which completes missing values using k-nearest neighbors : each missing value is imputed using the mean value from the nearest neighbors found in the training set.

2.2 Feature selection



We observed that the learning set seems to be balanced. We can see on the figure at the top right that the activity ID's are well distributed: they are each composed of 250 instances of the subjects in the training set. However, when looking at the subject ID's of the learning set, we observed that some subjects have been less used to perform measures. This can be observed on the figure at the top right. This can lead to overfit: all subjects have not the same behavior when performing an activity and the sensors measures are correlated to one subject.

These observations and the fact that the dataset was really huge ($3500 * 31 * 512$ values), motivated us to perform features selection. We decided to perform feature selection using `SelectFromModel`². It is performing ranking based on one estimator and then remove the less ranked features recursively. We chose the `ExtraTreeClassifier`³ as estimator for this method.

¹<https://scikit-learn.org/stable/modules/impute.html>

²https://scikit-learn.org/stable/modules/feature_selection.html

³<https://www.geeksforgeeks.org/ml-extra-tree-classifier-for-feature-selection/>

It is similar to a random tree classifier and only differ in the manner it is constructed which lead to the creation of multiple de-correlated decision trees. This seems us to be a good model to rank the features because we have seen in the theoretical course that decision trees are nice to discriminate features but are too much related to the dataset. It is why the way the Extremely Randomized Trees are constructed reduce the variance link to the dataset and thus reduce the chances of overfitting.

2.3 Implementation

The implementation of this part can be found in the `preprocessing.ipynb` notebook or in the `preprocessing.py` python script. We have performed some tests measuring the scores and it is a bit increasing the scores of our models.

3 Studied models

3.1 K-nearest neighbors

First, we decided to use the K-nearest neighbors method because it was the one provided by the pedagogical team with the assignment. Yet, the reasons that made us trying multiple versions of this technique are the facts that it is easy to interpret and to use.

3.1.1 1-NN

Our first submission was the one provided with assignment⁴ that was generated with the 1-nearest neighbor model⁵.

This submission yields a score of 0.52 on Kaggle.

3.1.2 25-NN

Based on 3.1.1, we decided to increase the value of K of the K-nearest neighbors method to 25. The motivation behind the choice of this value for K was theoretical. Since the dataset is huge and can contain errors in the measured values, we thought that increasing the value of K would increase the precision of the model because each prediction will not depend on a single neighbor that could have been misclassified.

This technique is implemented in the file `knn_basic_25.py`. This submission yields a score of 0.54 on Kaggle.

3.1.3 55-NN and 49-NN

We decided to keep using the K-NN method to study how well it could perform on the assignment. We decided to study the accuracy of the KNN technique by using a 10-fold cross validation technique with the negative log loss function as a scoring measure. We obtained the following graph of accuracies depending on the value of K:

⁴File `example_submission.csv`.

⁵Files `toy_script.py` or `knn_basic_1.py`.

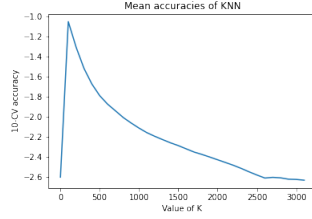


Figure 1: 10-CV accuracy of KNN depending on value of K

We could see that there is a peak. By using a divide and conquer technique, we found that the peak was obtained for values of K around 50. Thus, we decided to study the accuracy in a small range of K values around $K = 50$. We obtained the following graph:

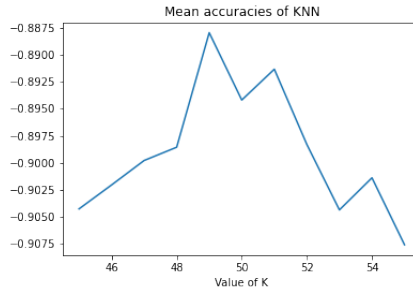


Figure 2: 10-CV accuracy of KNN depending on value of K

Based on the divide and conquer approach performed, we decided to use 55 neighbors. The choice of this value was motivated by the fact that we observed a peak around this value. This submission yields a score of 0.53142 on Kaggle and the implementation of it is inside the file `knn_basic_55.py`.

Then, based on the graph of figure 2, we decided to use 49 neighbors. The choice of this value was motivated by the fact that it yields the highest score on the learning set by using a 10-fold cross validation strategy with the negative log loss function as a scoring measure. This submission yields a score of 0.52571 on Kaggle which is less than the one for 55-NN. Yet, we expected a higher score based on the graph of figure 2. The implementation is available inside the file `knn_basic_49.py`.

Based on the two previous results, we came to the conclusion that using "vanilla" K-nearest neighbors was not sufficient. Thus, we decided to find new techniques explained in the following sections.

3.1.4 Multiple K-nearest neighbors models

We took a look back at the slides provided with the assignment and we focused on the features (sensors) of the datasets. We noticed that the features are not in the same units. Thus, we thought about using multiple 1-nearest neighbor models with one 1-NN model per feature. This idea was motivated by the fact that we thought we would obtain a better accuracy by using one 1-NN per feature because measurements in different units would not be mixed. We implemented this approach in the file `knn_split_1.py`. This yields a score of 0.52 on Kaggle which is the same as the toy script provided with the assignment.

3.1.5 Mutiple K-nearest neighbors models and samples modification

After the disappointing result obtained in 3.1.4, we decided to take a better a look at the datasets. We noticed that some samples were vectors full of -999999.99. We though that these measurements would badly influenced the models. Thus, we decided to replace each sample that was a vector full of -999999.99 by a vector full of 0.

We implemented this approach in the file `knn_splitted_1_filtered.py` and the score obtained on Kaggle was 0.53142.

3.1.6 Conclusion on K-nearest neighbors

After all the different tries perform using K-nearest neighbors models, we came to the conclusion that "vanilla" K-NN were not sufficient for this project.

The codes for the different plots are available inside the Jupyter notebooks `knn_basic.ipynb` and `knn_splitted.ipynb`.