

Machine Learning - SS 2021

Exercise 1: Classification

Federico Ambrogi, e1449911@student.tuwien.ac.at

Adam Höfler, e11847620@student.tuwien.ac.at

Matteo Panzeri 12039996@student.tuwien.ac.at

TU Wien

April 21, 2021

Contents

1 Data Sets Description	1
2 Data Exploration and Pre-processing	2
3 Holdout and Cross Validation	2
4 Performance Tests	2

Introduction

In this document we describe the results of the implementation of three algorithms to solve classification problems on four different dataset.

1 Data Sets Description

Here we briefly introduce our four datasets.

Drug consumption data set This is the data set taken from exercise 0, so here we only provide basic details. The primary task of our implementation is to classify the likeliness of the frequency of usage of a certain type of drug (from never used to current usage), given personal aspects such as gender, ethnicity, education background etc. and psychological traits.

Asteroids data set The asteroids data set is designed to perform a binary classification task. It consists of a total of **XXX** features, mostly astrophysical parameters of asteroids such as the distance from the Earth, the size of the major axis of the orbits, their mass etc. as well metadata regarding their name. The target feature for classification is a string variable that maps to boolean values "True" in the case the asteroid represents a concrete hazard for the Earth i.e. with high risk of impact, and "False" otherwise.

Blah

Blah

2 Data Exploration and Pre-processing

In this section we describe the necessary preliminary steps to import and prepare the data to make them suitable for the learning algorithms. The steps are handled by the script *data_preparation.py* which includes dedicated function to pre-process each data set.

3 Holdout and Cross Validation

Here we describe briefly the techniques of "holdout" and "cross validation" that are used to evaluate a model. The **holdout** method is essentially based on the splitting on the input data set into two subset, one used for training the model, and one used for testing the model, for example in 80% – 20% proportion, although there is not fix recipe for this split. Once the model is trained, the evaluation can be performed on the test data set, and this is therefore possible to check if the prediction of the model match the data. One big issue of this model is that the training strongly depend on the splitting of the initial data set, for example if the characteristic of the training data set are not representative of the whole data set.

A more powerful method is the **cross-validation** or "k-fold cross validation". Form the full dataset, a test is held out for final evaluation, but the validation set is no longer needed. For this, the training set is split into k sets, so that the model is trained using $k - 1$ folds as training data, and the remaining one is used for the validation as a test set (to compute the interesting metric of the model). Once we obtain such k number of metrics, the final result is the average of these parameters, obtained for each iteration of the cross-validation on each distinct fold.

4 Performance Tests

Confusion Matrix For each classifier, we produce a confusion matrix where each entry i, j corresponds to the number of observations in group i , but predicted to be in group j . We chose to normalize the entries according to the sum of each row. In case of binary classification, the matrix reduces to the number of true negatives (TN), false positives (FP), false negatives (FN) and true positives (TP).

Examples can be see in Fig. 1.

We remind here the definition of the metric parameters we will used to quantify the performance of our classifiers i.e. precision (P), recall (R), accuracy (A) and specificity S :

$$P = \frac{TP}{TP + FP} \quad R = \frac{TP}{TP + FN} \quad A = \frac{TP + TN}{all} \quad S = \frac{TN}{TN + FP} \quad (1)$$

It is straightforward to calculate these parameters for binary classification tasks out of the confusion matrix. In case of multiple labels, we need to calculate these parameter for each class, given that:

1. TP s are the values in the diagonal;
2. FN s for a certain class are the sum of values in the corresponding row excluding the TP ;
3. FP s for a certain class are the sum of values in the corresponding column excluding the TP ;
4. TN s for a certain class are the sum of all rows and columns, excluding the class's column and row.

Another convenient metric, particularly because it is calculated directly from the proper *scikit-learn* function, is called *f1 - score*, which is defined as the harmonic mean of the precision and recall:

$$f1 - score = 2 \times \frac{PR}{P + R} \quad (2)$$

Describe micro averaging and macro averaging, will we report both ?

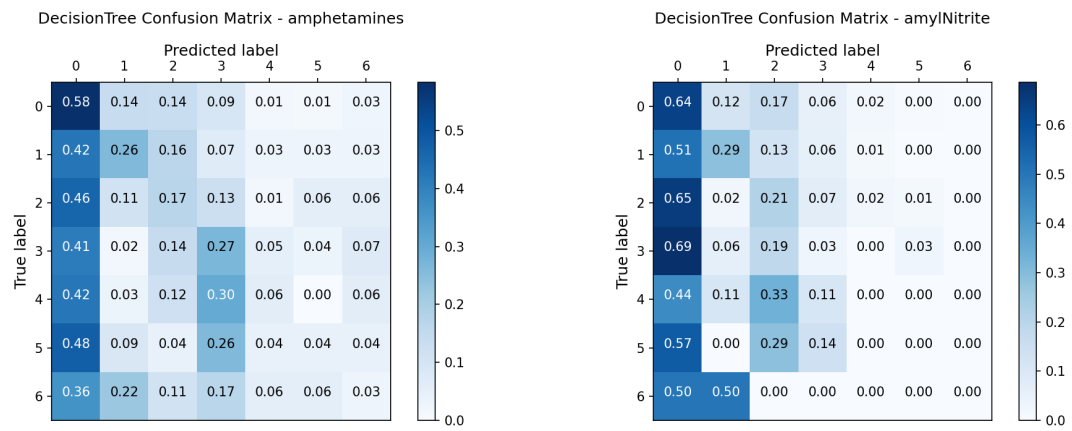


Figure 1: Example confusion matrices for XXX