

FYS-STK3155/4155 Applied Data Analysis and Machine Learning - Project 3

Lotsberg, Bernhard Nornes
Nguyen, Anh-Nguyet Lise

<https://github.com/liseanh/FYS-STK4155-project3>

November - December 2019

Abstract

Whole page: 6.24123in Column: 3.01682in

1 Introduction

In popular culture the neural network is probably the most well known form of machine learning. In recent years many other statistical learning methods have proven themselves as well however. In this project we compare the performance of neural networks and gradient boosting in the case of binary classification. In addition to these, we also use the much simpler k-nearest neighbours method as a baseline for classification performance.

2 Data

The data set we will analyse in this project is the MAGIC Gamma Telescope data set retrieved from the UCI Machine Learning Repository, which was generated by a Monte Carlo (MC) program described by D. Heck et. al. [4] to simulate high energy gamma particle registration in a Cherenkov gamma telescope. The set consists of ten explanatory variables and a binary response variable `class` which specifies whether the measured photons resulted from a gamma par-

ticle (`class = g`) or a hadron (`class = h`). The entire data set consists of 19020 instances with no missing values, with outcome distribution as shown in Figure 1. The explanatory and response variables are defined as the following by the UCI Machine Learning Repository [2]:

1. `fLength`: continuous # major axis of ellipse [mm]
2. `fWidth`: continuous # minor axis of ellipse [mm]
3. `fSize`: continuous # 10-log of sum of content of all pixels [in #phot]
4. `fConc`: continuous # ratio of sum of two highest pixels over `fSize` [ratio]
5. `fConc1`: continuous # ratio of highest pixel over `fSize` [ratio]
6. `fAsym`: continuous # distance from highest pixel to center, projected onto major axis [mm]
7. `fM3Long`: continuous # 3rd root of third moment along major axis [mm]
8. `fM3Trans`: continuous # 3rd root of third moment along minor axis [mm]
9. `fAlpha`: continuous # angle of major axis with vector to origin [deg]
10. `fDist`: continuous # distance from origin

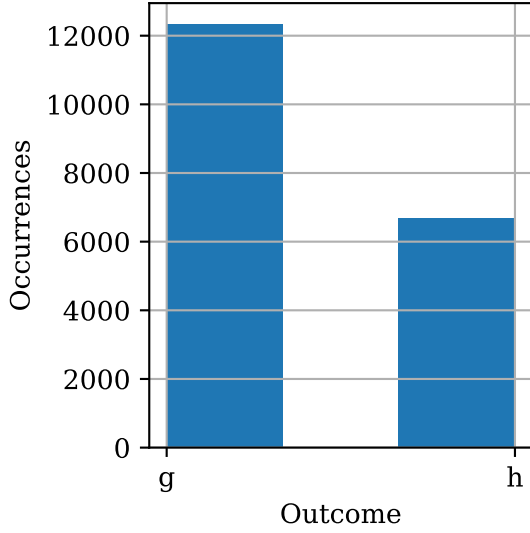


Figure 1: Frequencies of the outcomes g and h in the data set. The numbers of instances for the categories were 12332 and 6688 for g and h respectively.

- to center of ellipse [mm]
11. class: g, h # gamma (signal), hadron (background)
 - m

3 Methods

3.1 k-Nearest Neighbour (kNN)

The k -nearest neighbour method is a non-parametric method that consists of using the k observations in the training set in feature space closest to a point x to form a model \hat{y} . For regression and binary classification, the method considers a point x and the k closest points x_i in the neighbourhood $N_k(x)$ defined by k , such that

$$\hat{y} = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i, \quad (1)$$

i.e. the model output is the average of the k closest points to x .

As we are using this method for binary classification, we assign the classes such that

$$\hat{Y} = \begin{cases} g & \text{if } \hat{y} \leq 0.5 \\ h & \text{if } \hat{y} > 0.5 \end{cases} \quad [3]. \quad (2)$$

Alternatively, the observation can be assigned according to the majority class of its neighbours, such that the observation is classified as the most common class in its neighbourhood. This is the method used for multi-class classification.

To implement kNN and the hyperparameter optimisation in this project we use the Scikit-Learn library.

3.1.1 Hyperparameter tuning

The kNN method is a simple approach that only requires tuning of a single hyperparameter k . Using $k = 1$ results in a highly complex model with high variance, while using $k = N$, where N is the total number of samples in the training set, results in a simple, highly biased model. To find the optimal value for k in our case, we use 5-fold cross-validation on the training set, keeping 1/3rd of the total data as test set.

3.2 Multilayer Perceptron (MLP)

A multilayer perceptron (MLP) is a feed-forward neural network. It contains an input layer, one or more hidden layers, with tunable number of nodes and layers, and a final output layer. The information flows only in one direction from the input layer to the output layer. The input and output values of the nodes are determined by an activation function, in addition to the weights and biases of the nodes. For a more extensive explanation of how the different components of the MLP works, please refer to our previous work, *Project 2: Regression and Classification* [5].

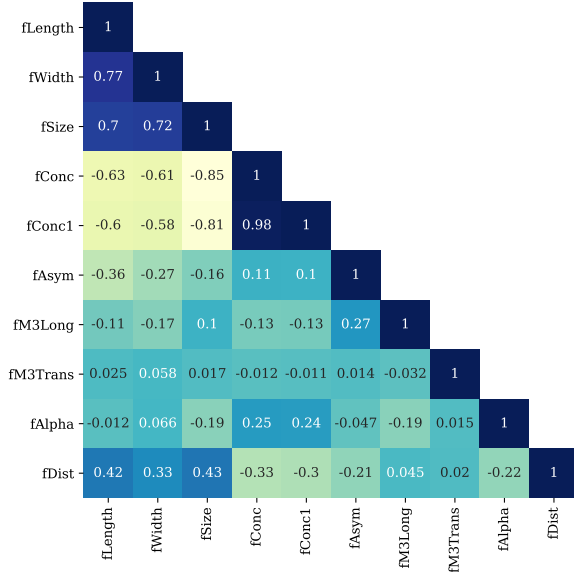


Figure 2: Correlation matrix of the features in the train set. Upper triangle excluded for readability.

We have chosen to use the Rectified Linear Unit (ReLU) function as our activation function $f(z)$ between the layers, which is given by

$$f(z) = \max(0, z), \quad (3)$$

and its gradient by

$$f'(z) = \begin{cases} 0, & z < 0 \\ 1, & z > 0 \end{cases}. \quad (4)$$

For the activation function of the final output layer, we use the softmax function to achieve the categorical outcomes,

$$f(z_i) = \frac{\exp(z_i)}{\sum \exp(z_j)} \quad (5)$$

For this project we will be using the TensorFlow library to implement the MLP with the addition of the Scikit-Optimize library to optimise the hyperparameters.

3.2.1 Hyperparameter tuning

In this project we choose to limit ourselves to two hidden layers for computational reasons,

and tune the amount of nodes in both of those layers.

To counteract the high variance of the model, we introduce regularisation using dropout. Dropout is tuned for each layer, and is the chance of any node in the given layer to be ignored at an iteration in training. This reduces variance by forcing the other nodes to be able to work more independently from each other, lessening the chance for overfitting.

We also tune the batch size. This leaves us with five hyperparameters to tune. A grid search would be too expensive to perform, so instead we perform a random search. Ideally we would employ cross validation on the train set evaluate the candidate hyperparameters, but for computational reasons we have decided to instead set aside 10% of the training set for validation. After finding the tuning parameters that maximise F_1 on the validation set, we use these parameters to fit the model on the entire training data including the validation data. This is known as a refit.

3.3 Gradient Boosting and Extreme Gradient Boosting (XGB)

Boosting is a learning technique based on the idea of combining several weak learners b_m into a final strong learner $f_{m_{\text{stop}}}$, where each f_m is improved from the previous iteration f_{m-1} through the step h_m ,

$$f_{m_{\text{stop}}} = \sum_{m=0}^{m_{\text{stop}}} h_m. \quad (6)$$

Gradient boosting is one such method. We consider a loss function $L(y, f(x))$, which is minimised by calculating the negative gradient of the loss function to obtain the minimisation direction. To prevent overfitting, the boosting step size, also referred to as learning rate, $\nu \in (0, 1)$ is introduced as a tuning parameter, with smaller values of ν resulting in larger shrinkage,

$$f_m(x) = f_{m-1}(x) + \nu h_m. \quad (7)$$

The gradient boosting algorithm goes as follows: [1]

1. Initialize the estimate $f_0(x)$
2. for $m = 1, \dots, m_{\text{stop}}$:
 - (a) Compute negative gradient vector

$$u_m = - \left. \frac{\partial L(y, f(x))}{\partial f(x)} \right|_{f(x)=\hat{f}_{m-1}(x)}$$

- (b) Fit the base learner to the negative gradient vector, $b_m(u_m, x)$
 - (c) Update the estimate

$$f_m(x) = f_{m-1}(x) + \nu b_m(u_m, x)$$

3. Compute final estimate,

$$\hat{f}_{m_{\text{stop}}}(x) = \sum_{m=0}^{m_{\text{stop}}-1} \nu b_m(u_m, x)$$

For this project, we have chosen to boost decision trees using a variant of gradient boosting called extreme gradient boosting (XGB).

3.3.1 Hyperparameter tuning

The XGBoost package has many hyperparameters added on top of the ordinary gradient boosting algorithm.

The learning rate ν

As with the neural network model, we also do not use cross validation here. We use the same train-validation split as for the neural network. We also employ a randomised search instead of a grid search.

3.4 Model evaluation

For classification problems, one can visualise the performance of a model using the so called confusion matrix. In our case it is defined as.

$$\begin{aligned} A &= \begin{pmatrix} \sum(\hat{y} = g, y = g) & \sum(\hat{y} = h, y = g) \\ \sum(\hat{y} = g, y = h) & \sum(\hat{y} = h, y = h) \end{pmatrix} \\ &= \begin{pmatrix} g_g & g_h \\ h_g & h_h \end{pmatrix}. \end{aligned} \quad (8)$$

A perfect model would then only have elements on the diagonal.

A common way to evaluate classification models is the F_1 -score. This is a metric quantifying the amount of diagonal and off-diagonal elements of the confusion matrix and is defined as

$$F_1 = \frac{2g_g}{2g_g + g_h + h_g} \quad (9)$$

The F_1 score can vary between 0 and 1, 1 representing a perfect fit with no off-diagonal elements. To us the best performing model is the model which maximises F_1 and minimises g_h on the test set, as false negatives are worse than false positives in this case because they lead to loss of important data.

4 Results

Figure 2 shows the correlation of the features in the training set. Some of the features are highly correlated to each other.

Table 3 shows the train and test F_1 scores for the three models employed in this project. XGBoost seems to outperform the neural network by a small margin for this split of the data. Both the neural network and XGBoost clearly outperform kNN. This is also evident in the confusion matrices shown in Figure 3, 4 and 5. Both the number of false positives and false negatives are lowest for XGBoost.

For kNN, the tuning parameter k was found to be best at $k = 5$. This is also evident in Figure 3. The estimated hyperparameters for MLP and XGBoost are stated in Table 1 and 2.

Table 1: Estimated optimal hyperparameters found using randomised search for the multilayer perceptron model.

Hyperparameter	Optimal value
Batch size	109
Epochs	55
Dropout, first layer	0.13
Dropout, second layer	0.095
Nodes, first layer	17
Nodes, second layer	4

Table 2: Estimated optimal hyperparameters found using randomised search for the extreme gradient boosted tree model.

Hyperparameter	Optimal value
Learning rate	0.048
Epochs	125
Maximum tree depth	9
L1 shrinkage parameter	0.0074
L2 shrinkage parameter	0.0035
Minimum child weight	3

Table 3: F_1 score for the training and test set for our k-nearest neighbour (kNN), multi-layer perceptron (MLP) and extreme gradient boost (XGB) models .

F_1 score	kNN	MLP	XGB
Training set	0.88	0.87	0.95
Test set	0.83	0.87	0.88

5 Discussion

6 Conclusion

Acknowledgements

We want to thank the University of California, Irving for providing us with the data used in this project.

References

- [1] Riccardo De Bin. STK-IN4300 Statistical Learning Methods in Data Science: Lecture 10. https://www.uio.no/studier/emner/matnat/math/STK-IN4300/h19/slides/lecture_10_ho.pdf, 2019. Retrieved: 09-12-2019.
- [2] Dheeru Dua and Casey Graff. UCI Machine Learning Repository. <http://archive.ics.uci.edu/ml>, 2017.
- [3] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The Elements of Statistical Learning*, volume 27. Springer Series in Statistics, 2009.
- [4] Dieter Heck, G Schatz, J Knapp, T Thouw, and JN Capdevielle. CORSIKA: a Monte Carlo code to simulate extensive air showers. Technical report, 1998.
- [5] Bernhard Nornes Lotsberg and Anh-Nguyet Lise Nguyen. *Project 2: Classification and Regression*. FYS-STK3155/4155

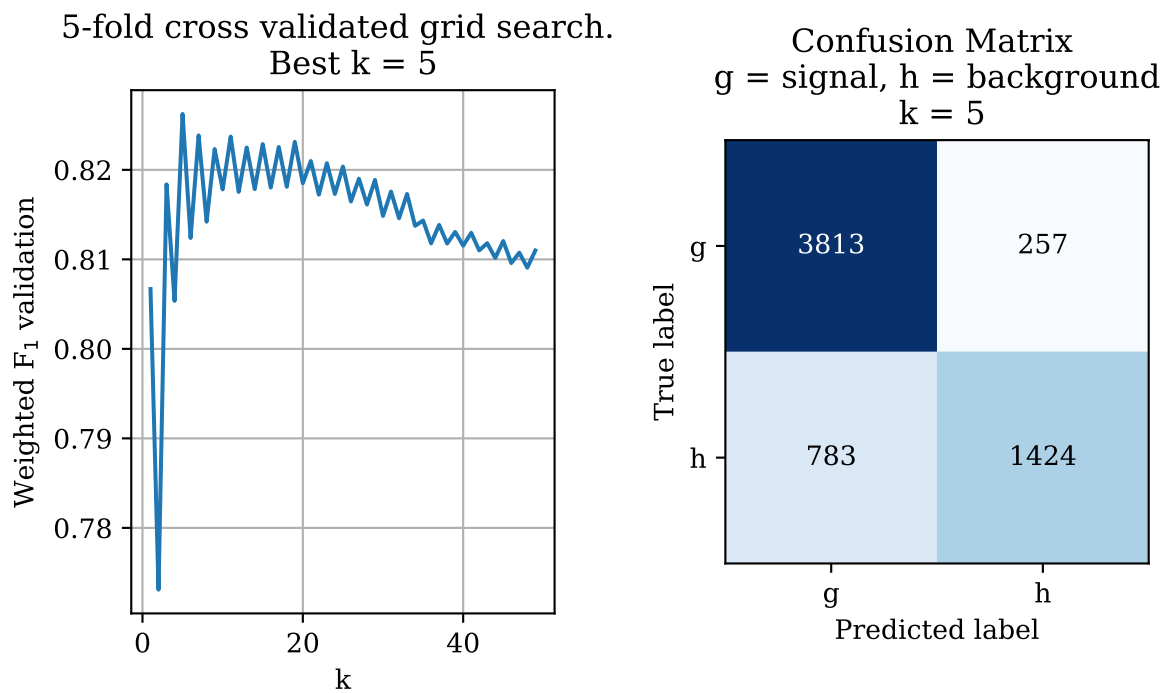


Figure 3: Results from tuned kNN using cross validation. The confusion matrix was found using the test set.

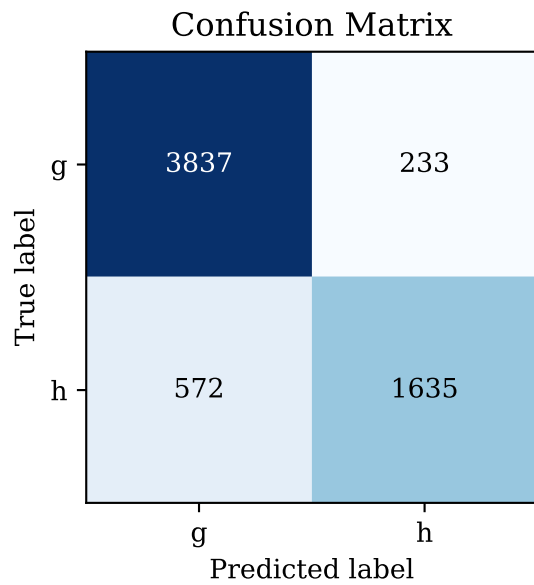


Figure 4: Confusion matrix of the neural network model applied to the test set.

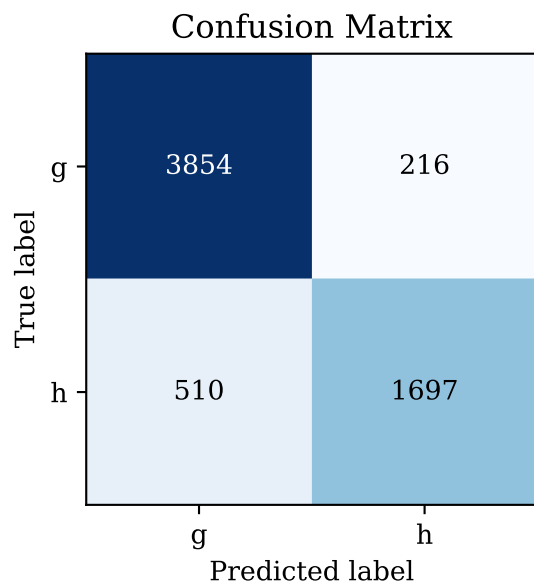


Figure 5: Confusion matrix of the gradient boosted model applied to the test set.