

# Report for project 3

Jon Audun, Mikael Ravndal and Adam P W Sørensen

December 12, 2018

In this project we analysed some cooking data and did fairly well.

## 1 Introduction

The concept of machine learning have gained a hughe popularity boost over the last couple of years. And this is no wonder: the techniques have a wide range of applications and can be a major asset if you know when to use what.

When to use what is exactly what we're going to have a brief peek into in this project. We will evaluate the performance of three different methods for classification on a data set consisting of recipies labeled with the cuisines they belongs to. That is we are looking at a classification problem with multiple classes. The methods we will look at are logistic regression, support vector machines and random forests.

All the code we have implemented can be found at <https://github.com/???>. A print of the jupyter notebook is also attached at the end of this report.

## 2 Methods

First we need to establish some notation. In the following we assume that we are given  $N \in \mathbb{N}$  samples consisting of  $p \in \mathbb{N}$  features and one target. Further we let  $K \in \mathbb{N}$  be the number of classes in our classification problem. We then denote by  $x_{i,j}$  the  $j$ -th feature of the  $i$ -th sample and by  $t_i \in \{1, 2, \dots, K\}$  the target of the  $i$ -th sample.

## 2.1 Logistic Regression

Let  $\mathbf{x}_i := (1, x_{i,1}, x_{i,2}, \dots, x_{i,p})$  be the row-vector in  $\mathbb{R}^{p+1}$  consisting of a one followed by the features of sample  $i$ . Also let  $\boldsymbol{\beta}_k$  be a column-vector in  $\mathbb{R}^{p+1}$  for  $k \in \{1, 2, \dots, K-1\}$  and let  $\boldsymbol{\theta} := [\boldsymbol{\beta}_1 \ \boldsymbol{\beta}_2 \ \dots \ \boldsymbol{\beta}_{K-1}]$  be the matrix with the  $\boldsymbol{\beta}_k$ 's as it's columns. In the case of multinomial logistic regression our model is then given by

$$\begin{aligned} P_k(\mathbf{x}; \boldsymbol{\theta}) &= \frac{\exp(\mathbf{x}\boldsymbol{\beta}_k)}{1 + \sum_{l=1}^{K-1} \exp(\mathbf{x}\boldsymbol{\beta}_l)} & k \in \{1, 2, \dots, K-1\} \\ P_K(\mathbf{x}; \boldsymbol{\theta}) &= \frac{1}{1 + \sum_{l=1}^{K-1} \exp(\mathbf{x}\boldsymbol{\beta}_l)} \end{aligned} \quad (1)$$

where the  $\boldsymbol{\beta}_k$ 's are the coefficients we wish to estimate.  $P_k(\mathbf{x}; \boldsymbol{\theta})$  is then the probability that a sample with features  $\mathbf{x}$  belongs to category  $k$ .

The way we fit our model in the case of logistic regression deviates slightly from the usual procedure with the introduction of a loss function. We now instead define a function we wish to maximize, namely the log-likelihood function given by

$$L(\boldsymbol{\theta}) = \sum_{i=1}^N \sum_{k=1}^K \chi_{\{k\}}(t_i) \log(P_k(\mathbf{x}_i; \boldsymbol{\theta})) \quad (2)$$

However the difference doesn't go further than the fact that this function shouldn't be interpreted exactly as a loss function. Other than that we proceed as usual by minimizing the negative of this function in order to maximize the original function. To do this we first need to compute the derivative of this function with respect to  $\boldsymbol{\theta}$ . Some computations gives us

$$\frac{\partial L(\boldsymbol{\theta})}{\partial \beta_{k,j}} = \sum_{i=1}^N \chi_{\{k\}}(t_i) x_{i,j} (1 - P_k(\mathbf{x}_i; \boldsymbol{\theta})) \quad (3)$$

We are now all set to use the gradient method of your choice to minimize  $-L$  as a function of  $\boldsymbol{\theta}$ .

### 2.1.1 Gradient descent

The perhaps simplest approach to minimizing a real valued function  $f$  iteratively, is to set some starting point  $\mathbf{x}_0$ , compute the gradient  $\nabla f$  in  $\mathbf{x}_0$ , and move  $\eta > 0$  in the opposite direction of this gradient in order to obtain  $\mathbf{x}_1 = \mathbf{x}_0 - \eta \nabla f(\mathbf{x}_0)$ . Then one iterates by setting  $\mathbf{x}_{i+1} = \mathbf{x}_i - \eta \nabla f(\mathbf{x}_i)$  for  $i = 1, 2, \dots$ . This is the framework of gradient descent. For small enough  $\eta$  this implies  $f(\mathbf{x}_i) > f(\mathbf{x}_{i+1})$ . The problem with this approach is to find out what  $\eta$  is small enough. If we choose  $\eta$  too small, it might take days for our algorithm to converge. If we on the contrary choose  $\eta$  to big, it might not converge at all.

### 2.1.2 Newton-Raphson

The Newton-Raphson method is a kind of gradient descent method, except we choose the learning rate in a more clever way. We thus start out by choosing a random starting point  $\mathbf{x}_0$ , just as in the case of gradient descent. The update step however is a bit different:

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \frac{\partial^2 f(\mathbf{x}_i)}{\partial \mathbf{x}_i \partial \mathbf{x}_i^T} \frac{\partial f(\mathbf{x}_i)}{\partial \mathbf{x}_i} \quad (4)$$

The result is a method that converges way faster than normal gradient descent. The draw back is however the need to compute the Hessian of your objective function. For a thorough derivation of this method consult [?]

### 2.1.3 Stochastic gradient descent

To avoid overfitting when training neural networks, training on randomly picked batches of the training data for each iteration is a good regularization method. If we let  $b$  be our batch size, the first step of this method consists of dividing the training data into  $N/b$  batches. One does this by picking from the  $N$  data-points without replacement. The second step is then to perform a gradient descent as described in section 2.1.1 with a new batch for each iteration.

## Support Vector Machines

### 2.2

To call a Support Vector Machine classifier with no kernel we use the following Python code:

```
svm_clf = svm.LinearSVC(C = 0.1)
svm_clf.fit(x_train, y_train)
preds = svm_clf.fit(x_test)
print('We did this good: %.4f' % accuracy_score(y_test, preds))
```

### 2.3 Decision trees

Wikipedias overview on decision trees:

A decision tree is a flowchart-like structure in which each internal node represents a “test” on an attribute (e.g. whether a coin flip comes up heads or tails), each branch represents the outcome of the test, and each leaf node represents a class label (decision taken after computing all attributes). The paths from root to leaf represent classification rules.

A decision tree typically uses all of the features as its features. The result of this is then that you get a classifier which predicts well on the training data, but falls short on testing data. This is because it fits its classifier perfectly to the data which it has seen. So a lone decision tree with all features will almost always be overfitted.

### 2.3.1 Rather a forest, then just a tree

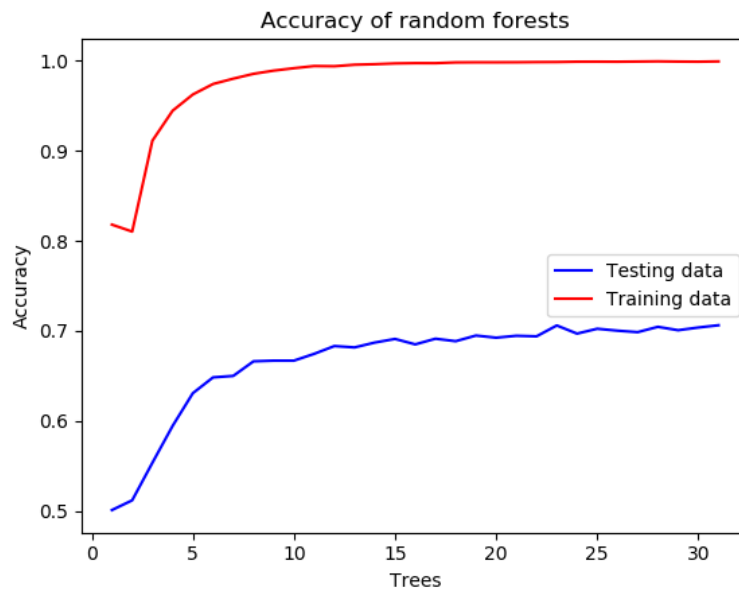
This is why random forests are more popular. A random forest can usually generalise much better than what a lone tree can.

A random forest is constructed by choosing how many decision trees you want in the forest and then train every tree on your data. But the construction of these trees is done differently than a lone decision tree.

Each tree just use a subset of the features, typically  $features = \sqrt{all - features}$ . This makes it so each tree functions a bit differently and will then predict different things. The final prediction then becomes which of the classes which gets voted most for when every tree predicts their own.

### 2.3.2 Difference in number of trees

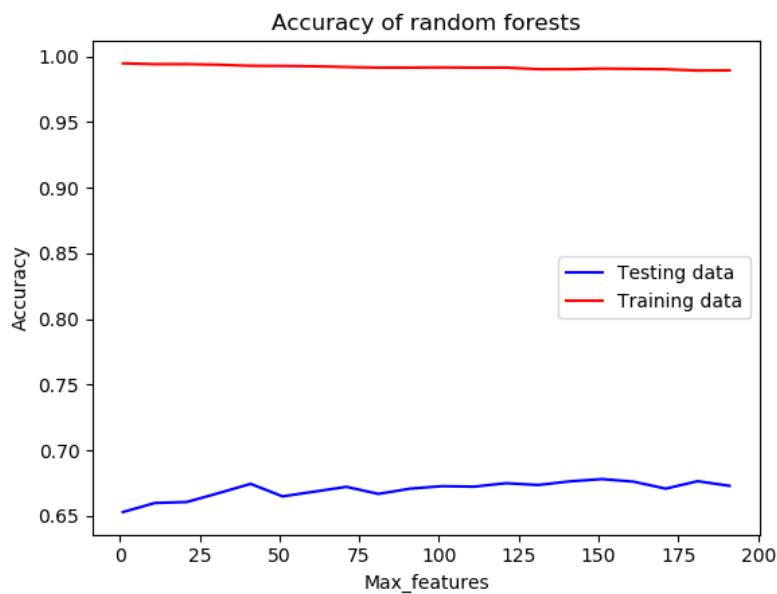
The number of trees in the forest is the main hyper parameter we have used to tweak this classifier. But it seems like the rule of thumb was the more, the merrier. Here is a plot of how well the random forest did:



So more trees provides a better accuracy. But the classifier starts to have a decent accuracy after 10 trees.

### 2.3.3 Difference in number of max\_features

Here is a plot on how well the accuracy is using different numbers as max\_features:



The number of trees used is 10.

The default in the scikitlearn classifier is using the root of the number of all the features as `max_features`. Since we have approximately 2000 features, the root becomes 44.72. That's why I computed different accuracies around this number.

The plot tells us that it doesn't seem to matter much what we use as `max_features`. The optimum in the plot is around 40 which gives us just more reason to stick with the default of using the root of all the features as `max_features`.

### **3 Model Selection and Verification**

### **4 Results**

### **5 Conclusion**

We did good.

### **References**