

Github Repository Classifier

Rami Aly¹, Andre Schurat²

¹ University of Hamburg

² Technical University of Dortmund

January 13, 2017

1 Abstract

Contents

1	Abstract	2
2	Selecting features	4
3	Gathering selected features from Github	4
4	Removing irrelevant information from selected features	4
5	Building the Prediction Model	4
5.1	Choosing a prediction Model	4
6	Training Set	5
7	Optimizing our Neural Network	5
8	Validation of created Classifier	5
9	Extensions	5
10	Source Code and used external Libraries	6

2 Selecting features

3 Gathering selected features from Github

4 Removing irrelevant information from selected features

5 Building the Prediction Model

5.1 Choosing a prediction Model

Of course there are many different approaches to the problem. A static algorithm to classify repositories is rather impractical because the parameters of our classify function would be strongly influenced by our interpretation of weights of the features for each class. However we quickly noticed that the complexity is very high, so that a normal algorithm must limit the aspects which can be considered. Above all the problem is non-linear and through the static analysis we would lose the possibility to freely improve or change the classifier. As the software and use-case market of Github rises the possible need of further classes could arise.

Hence to ensure a classifier who is as dynamic and as extensible as possible we choose to use some form of machine learning. The problem which needs to be solved by the Prediction Model is a classification problem: We have a fixed number of values for selected features as input and as an output the class to which the values fit the most. As a result of this fact it was pretty clear to us that a supervised learning method would be optimal.

In the next step we thought about the pro- and contra arguments of non-parametric and parametric learning. For example Gaussian-Process-Models could be used in principle, as one does not need to specify a fixed number of parameters and therefore be non-parametric. The main problem with Gaussian-Process-Models is that they scale rather poorly with a complexity of $O(n^3)$ [1]. Moreover if we keep the huge dimension of repository datasets in mind and as such the possible complexity of the classifier function, our choice will lead us to a parametric neural Network.

6 Training Set

7 Optimizing our Neural Network

8 Validation of created Classifier

9 Extensions

References

- [1] Gaussian Processes 2006, David P. Williams <http://people.ee.duke.edu/~lcarin/David1.27.06.pdf>

10 Source Code and used external Libraries

Source Code Github: <https://github.com/Crigges/InformatiCup2017>