

Predicting Continuous Integration Build Outcome

Project Report (Draft)

Johannes Kästle
University of Alberta

Abstract

In order to improve developers productivity and increase pace, it is tried to predict the outcome of continuous integration builds. For that, 2,500,756 builds with 37 features were analyzed, using three different machine learning algorithms with over 700 configurations. Generally, all classifiers struggled with predicting failures, but had high sensitivity. In the end, the decision tree classifier could predict the outcome with an accuracy of 76% and performed significantly better than its competitors.

1 Introduction

In modern software engineering using continuous integration has evolved to a best practice. Its goal is it to archive higher productivity while developing and reduce the number of faults and bugs. Thanks to its interweaving with GitHub, TravisCI is commonly used for open source projects. Therefore, there lots of continuous integration data is available. It has been shown that broken builds delay a project significantly [1], hence the prediction if a build will fail and why can increase development pace. With this information it is possible to develop tools and techniques which can reduce number of failing builds.

In previous researches it was shown that there are correlations between build metrics and outcome [2], and that cascading classifiers work well for predicting the build outcome [3]. Another research tried to predict build outcome with decision trees, as well, using data from previous builds [4]. This work wants to try a simpler approach and see if it is possible to predict the build outcome only using the current's build information and its associated commits. For this, the TravisTorrent [5] dataset is analyzed in order to classify the build outcome. Then, three different algorithms, naive Bayes, neural networks and decision trees, are used to classify this dataset and it is evaluated if one of them is able to accurately learn this problem.

2 Background

TravisCI is a continuous integration tool to automated compile, build, integrate and test GitHub projects.

2.1 Dataset

For the MSR 2017 challenge the dataset TravisTorrent[5] was created. This study uses the version from January 11, 2017. It contains 3,702,595 samples. A detailed explanation of

all the features can be found on the homepage of TravisTorrent¹. Because not all features are useful for this study, the number of features were restricted to some features. A list of all 37 used features can be found under table 6. Only builds that were either successful, erroneous or failed were considered. Successful builds were considered as class 1, or true, while the other builds were marked as failed, i. e. class 0, or false. Aborted or canceled builds were ignored, therefore only 2,500,756 builds were considered in the final dataset. In the final dataset 41.6% of the builds did fail.

2.2 Algorithms

As a baseline an algorithm is used which labels all data to the most frequent class. Then three supervised classification algorithms are tried. The first one is Naive Bayes which assumes every feature to be Gaussian data. The second one is a neural net with one hidden layer and the logistic activation function. As solver for the weight optimization Adam[6], an optimized stochastic gradient descent, is used. The batch size is determined automatic, the maximum epochs are set to 200 and the dataset is shuffle between each epoch. As third algorithm a decision tree is used. A decision tree uses subsequent rules to classify the data. For example, the first node could be if the programming language is Java. If that is the case, the sample goes to the left child, otherwise to the right child. This is done, until the sample is classified at a leaf node. The criterion for splitting is set to the gini impurity, for performance it takes the best random split.

3 Methodology

The dataset is stored in a MySQL database. From there all samples are extracted and stored in a matrix. The build outcome is put into a vector y with the value 1 if the build was successful and 0 otherwise. The 'language' and 'previous build' feature are stored in a vector; for example the language 'java' is encoded into $[1, 0, 0]$, while 'javascript' is encoded to $[0, 1, 0]$. The date features are split up into month, day, daytime in seconds and weekday as integers. Since not every sample contains all features, the missing features are set to their respective mean, so it does not impact the classification.

The actual learning is divided into two parts. First, the best parameters for the neural network and the decision tree have to be learned, and afterwards the algorithms need to be compared. Because the dataset is big enough, the dataset is split up in two equal halves, one for every of the two tasks. The split and all following splits are done stratified. This complete split is used so that the parameter search does not bias the inter algorithm comparison.

The used scripts and results can be found in the authors GitHub repository².

3.1 Parameter

To look for the best parameters, the first half of the dataset is used for an exhaustive grid search. This cross validation learns a model for every possible combination of the

¹https://travis torrent.testroots.org/page_dataformat/

²<https://github.com/jodokae/cmpu551-mini-project/>

parameters. This is done on a K-Fold with $k = 5$. A list of all tried parameters can be found in tables 7 and 8. A detailed description of all possible parameters can be found in the documentation of sklearn for the decision tree³ and the neural network classifier⁴. Accuracy is used as scoring method.

3.2 Cross Algorithms

The other half of the data is used to compare the different algorithms. For that the dataset is split up into ten random sets, where the size of the test set is 20%. On this, every of the three and the baseline algorithm is run. For evaluation the confusion matrix is saved for every run. Then the means and standard deviation over the runs is computed for every algorithm.

3.3 Statistics

The algorithms are compared using Welch’s t-test, because the variances of the algorithms are different. The baseline algorithm has, as expected, a variance of 0, whereas in test runs the neural networks had generally a high standard deviation. As comparing metrics accuracy, sensitivity and specificity are used, in which the latter ones are equivalent to the AUC of the ROC curve. Those three metrics are used to give a general sense of how good the algorithms are performing and then in detail how good they are at predicting successful and failing builds separately.

The two tailed t-test is done pairwise for every algorithm with $H_0 : \mu_0 = \mu_1$ and $\alpha = 0.05$ to produce a ranking between the three (plus baseline) algorithm.

4 Results

The methodology was applied to the full dataset. The following results were found:

4.1 Parameters

When running on the first 50% of the dataset, 243 different parameter combinations were fitted on five folds for the neural nets. This resulted into 1215 different learnings. With the highest accuracy of 70.78%, the parameters which the highest number of hidden units and the smallest alpha was chosen. The learning rate was with 0.001 the median value, as well as the beta_1 parameter, whereas beta_2 was given the highest value.

For the decision tree classifier, 486 different parameter sets were tried. On five folds, which were not necessarily the same as for the neural net, this resulted into 2430 different decision trees. The best parameters found were the the maximum depth of 100, the split and leaf ratios all at $1e - 5$ and the maximum leaf nodes as 10,000. As maximum features the logarithm was taken. With this settings, an accuracy of 77.63% was achieved.

³<http://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html>

⁴http://scikit-learn.org/stable/modules/generated/sklearn.neural_network.MLPClassifier.html

Table 4: t-test sensitivity (t,p)

	Baseline	Naive Bayes	Decision Tree	Neural Network
Baseline	(0, 1)	(21.29, 5.22×10^{-9})	(17.04, 3.71×10^{-8})	(4.445, 0.00161)
Naive Bayes	(-21.29, 5.22×10^{-9})	(0, 1)	(11.63, 3.84×10^{-7})	(-0.2063, 0.841)
Decision Tree	(-17.04, 3.71×10^{-8})	(-11.63, 3.84×10^{-7})	(0, 1)	(-8.229, 1.73×10^{-7})
Neural Network	(-4.445, 0.00161)	(0.2063, 0.841)	(8.229, 1.73×10^{-7})	(0, 1)

Table 5: t-test specificity (t,p)

	Baseline	Naive Bayes	Decision Tree	Neural Network
Baseline	(0, 1)	(-25.65, 1×10^{-9})	(-16.05, 62.7×10^{-8})	(-4.785, 0.000995)
Naive Bayes	(25.65, 1×10^{-9})	(0, 1)	(-14.59, 1.35×10^{-7})	(-0.4998, 0.629)
Decision Tree	(16.05, 6.27×10^{-8})	(14.59, 1.35×10^{-7})	(0, 1)	(13.69, 2.96×10^{-8})
Neural Network	(4.785, 0.000995)	(0.4998, 0.629)	(-13.69, 2.96×10^{-8})	(0, 1)

From table 2 it can be seen, that the algorithms don't seem to differ that much in terms of accuracy, apart from the decision tree which is ahead with almost 6 points. This is probably the case because the decision tree is the only algorithm which has an acceptable specificity rate. Although 31% is high in comparison, it is still a low value.

To confirm the feeling, that the decision tree is the best algorithm in this case, a two tailed t-test was executed for all three metrics pairwise between all algorithms. The null hypothesis, that the algorithms have the same mean, are therefore equal, is tested with $\alpha = 0.05$. The full results can be seen in tables 3 - 5. Since the baseline algorithm achieves a full score in sensitivity and zero score at specificity by its nature, the t-test for those metrics is only used for comparison of the algorithms under study.

According to accuracy (table 3) the Naive Bayes classifier performs worse than the baseline, the neural net performs equally good and only the decision tree can achieve a higher accuracy than the mean predictor. Neural net and Naive Bayes are not statistically different in direct comparison. In terms of sensitivity (table 4) the tree performs the worst. Neural net and Bayes do not differ again, this does not change for specificity (table 5), as well. The decision excels clearly.

With the results of the t-tests, the decision tree classifier algorithm is considered as winner over the others. Despite losing in sensitivity, the difference is only 4%, whereas it wins in classifying failures 10 times more accurate than the other two algorithms. This is confirmed by the t-test of accuracy. The other two algorithms are not worth using, since they do not perform better than the baseline algorithm.

5 Threats to validity

Not all possibilities could be checked in this project which influence the strength of the results. First of all, because of lack of time (one run of the complete methodology takes over ten hours on a standard laptop), the complete project, i. e. find parameters with cross validation and compare algorithms on different splits, was only done once and not on separate splits of the complete dataset. This was also necessary, because different runs

would produce different optimal parameters which would then harden the conclusion of the comparison. Because the dataset is quite large, this issue is considered to be acceptable.

Another runtime related problem is that not all possible parameters were considered. For example, this work restricts itself to the gini impurity for the decision tree. Hence, it cannot be generalized that decision trees are strictly better than neural networks for this dataset, but only under the tested parameters. The terms 'better' and 'winner' are only regarding accuracy, sensitivity and specificity and do not consider runtime or other metrics at all. Because this work only wants to compare the possibility of predicting this type of data at all, runtime was considered a less important aspect.

Since the dataset only is a small snapshot of all CI projects out there, especially restricted only to three programming languages and only open source projects, it cannot be generally said that build outcomes are predictable. But it is believed, that the considered features do not differ significantly between closed and open source projects. Furthermore, the dataset contained missing data for some features, which may impact the outcome of the analysis. This was mitigated due to preprocessing of the data. In future, it might be interesting to see if the accuracy of the prediction can be improved with feature selection.

6 Conclusion

In order to predict build outcomes of continuous integration builds, it was found, that a decision tree classifier can predict the outcomes with an accuracy of 76%. Although it has high variance in its learning, it still outperforms Naive Bayes and neural networks with one hidden layer, by far. Of course not all possible parameters were tested, but it gives a general idea direction should be looked at. The biggest problem is optimizing the prediction of failing builds, which have only a specificity of 31%.

References

- [1] N. Kerzazi, F. Khomh, and B. Adams, "Why do automated builds break? an empirical study," in *2014 IEEE International Conference on Software Maintenance and Evolution*, pp. 41–50, Sept 2014.
- [2] M. R. Islam and M. F. Zibran, "Insights into continuous integration build failures," in *2017 IEEE/ACM 14th International Conference on Mining Software Repositories (MSR)*, pp. 467–470, May 2017.
- [3] A. Ni and M. Li, "Cost-effective build outcome prediction using cascaded classifiers," in *2017 IEEE/ACM 14th International Conference on Mining Software Repositories (MSR)*, pp. 455–458, May 2017.
- [4] A. E. Hassan and K. Zhang, "Using decision trees to predict the certification result of a build," in *Proceedings of the 21st IEEE/ACM International Conference on Automated Software Engineering, ASE '06*, (Washington, DC, USA), pp. 189–198, IEEE Computer Society, 2006.
- [5] M. Beller, G. Gousios, and A. Zaidman, "Travistorrent: Synthesizing travis ci and github for full-stack research on continuous integration," in *Proceedings of the 14th working conference on mining software repositories*, 2017.

- [6] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *CoRR*, vol. abs/1412.6980, 2014.

A Features

Table 6: Feature list

feature	type
gh_lang	[java, javascript, ruby]
gh_num_commits_in_push	int
git_prev_commit_resolution_status	[build_found, merge_found, no_previous_build]
gh_team_size	int
git_num_all_built_commits	int
gh_num_issue_comments	int
gh_num_commit_comments	int
gh_num_pr_comments	int
git_diff_src_churn	int
git_diff_test_churn	int
gh_diff_files_added	int
gh_diff_files_deleted	int
gh_diff_files_modified	int
gh_diff_tests_added	int
gh_diff_tests_deleted	int
gh_diff_src_files	int
gh_diff_doc_files	int
gh_diff_other_files	int
gh_num_commits_on_files_touched	int
gh_sloc	int
gh_test_lines_per_kloc	double
gh_test_cases_per_kloc	double
gh_asserts_cases_per_kloc	double
gh_by_core_team_member	double
gh_description_complexity	int
gh_pushed_at	date
gh_build_started_at	date
tr_duration	int
tr_log_bool_tests_ran	double
tr_log_bool_tests_failed	double
tr_log_num_tests_ok	int
tr_log_num_tests_failed	int
tr_log_num_tests_run	int
tr_log_num_tests_skipped	int
tr_log_tests_failed	int
tr_log_test_duration	double
tr_log_build_duration	double
tr_status	[errored, failed, passed]

B Parameter

Table 7: Parameters for Decision Tree Classifier

parameter	values	description
max_depth	[1, 10, 100]	maximum depth of the tree
min_samples_split	[1e-5, 1e-4, 1e-3]	minimum fraction of samples needed on both child nodes
min_samples_leaf	[1e-5, 1e-4, 1e-3]	minimum fraction of samples per leaf
min_weight_fraction_leaf	[1e-5, 1e-4, 1e-3]	minimum weighted fraction of samples per leaf
max_features	['sqrt', 'log2']	maximum amount of features considered per split, either $\sqrt{ features }$ or $\log_2 features $
max_leaf_nodes	[100, 1000, 10000]	maximum number of leaves in tree

Table 8: Parameters for Neural Network Classifier

parameter	values	description
hidden_layer_sizes	[(10,), (50,), (100,)]	number of hidden nodes
alpha	[1e-4, 1e-3, 1e-2]	l2 regularization parameter
learning_rate_init	[1e-3, 1e-2, 1e-1]	step size
beta_1	[0.1, 0.5, 0.9]	adam specific parameter
beta_2	[0.1, 0.5, 0.9]	adam specific parameter