

Information Modelling & Analysis – Project 2

Claudio Maggioni

Code Repository

The code and result files, part of this submission, can be found at:

- Repository: <https://github.com/infoMA2023/project-02-bug-prediction-maggicl>
- Commit ID: **TBD**

Data Pre-Processing

I use the sources of the CLOSURE repository that were already downloaded using the command:

```
defects4j checkout -p Closure -v 1f -w ./resources
```

and used the code in the following subfolder for the project:

```
./resources/defects4j-checkout-closure-1f/src/com/google/javascript/jscomp/
```

relative to the root folder of the repository. The resulting CSV of extracted, labelled feature vectors can be found in the repository at the path:

```
./metrics/feature_vectors_labeled.csv
```

relative to the root folder of the repository.

Unlabeled feature vectors can be computed by running the script `./extract_feature_vectors.py`. The resulting CSV of unlabeled feature vectors is located in `./metrics/feature_vectors.csv`.

Labels for feature vectors can be computed by running the script `./label_feature_vectors.py`.

Feature Vector Extraction

I extracted **291** feature vectors in total. Aggregate metrics about the extracted feature vectors, i.e. the distribution of the values of each code metric, can be found in Table 1.

Table 1: Distribution of values for each extracted code metric.

| Metric | Minimum | Average | Maximum |
|--------|---------|----------|---------|
| BCM | 0 | 13.4124 | 221 |
| CPX | 0 | 5.8247 | 96 |
| DCM | 0 | 4.8652 | 176.2 |
| EX | 0 | 0.1134 | 2 |
| FLD | 0 | 6.5773 | 167 |
| INT | 0 | 0.6667 | 3 |
| MTH | 0 | 11.6529 | 209 |
| NML | 0 | 13.5622 | 28 |
| RET | 0 | 3.6735 | 86 |
| RFC | 0 | 107.2710 | 882 |
| SZ | 0 | 18.9966 | 347 |
| WRD | 0 | 314.4740 | 3133 |

Feature Vector Labelling

After feature vectors are labeled, I determine that the dataset contains **75** buggy classes and **216** non-buggy classes.

Classifiers

In this section I explain how I define and perform training for each classifier.

Since the dataset has an unbalanced number of feature vectors of each class, in order to increase classification performance I upsample the dataset by performing sampling with replacement over the least frequent class until the number of feature vectors matches the most frequent class¹.

Other than for the **GaussianNB** (Naive Bayes) classifier, the classifiers chosen for the project offer to select hyperparameter values. In order to choose them, I perform a grid search over each classifier. The hyperparameter values I have considered in the grid search for each classifier are the following:

- For *DecisionTreeClassifier*:

| Parameter | Values |
|-----------|---------------|
| criterion | gini, entropy |
| splitter | best, random |

- For *SVC*:

| Parameter | Values |
|-----------|----------------------------|
| kernel | linear, poly, rbf, sigmoid |
| gamma | scale, auto |

- For *MLPClassifier*:

| Parameter | Values |
|--------------------|---|
| max_iter | 500000 |
| hidden_layer_sizes | [5, 10, 15, ..., 100], [15, 30, 45, 60, 75, 90] ² , [20, 40, 60, 80, 100] ³ |
| activation | identity, logistic, tanh, relu |
| solver | lbfgs, sgd, adam |
| learning_rate | constant, invscaling, adaptive |

Note that the [...]² denotes a cartesian product of the array with itself, and [...]³ denotes the cartesian product of [...]² with the array (i.e. [...]³ = [...]² × [...] = ([...] × [...]) × [...]).

Note also the high upper bound on iterations (500000). This is to allow convergence of the less optimal hyperparameter configurations and avoid **ConvergenceWarning** errors.

- For *RandomForestClassifier*:

| Parameter | Values |
|--------------|------------------------------|
| criterion | gini, entropy |
| max_features | sqrt, log2 |
| class_weight | balanced, balanced_subsample |

¹Upsampling due to unbalanced classes was suggested by *Michele Cattaneo*, who is attending this class.

The script `./train_classifiers.py`, according to the random seed 3735924759, performs upscaling of the dataset and the grid search training, by recording precision, accuracy, recall and the F1 score of each configuration of hyperparameters. These metrics are then collected and stored in `./models/models.csv`.

The metrics for each classifier and each hyperparameter configuration in decreasing order of accuracy are reported in the following sections.

For each classifier, I then choose the hyperparameter configuration with highest accuracy.

Decision Tree (DT)

| critierion | splitter | precision | accuracy | recall | f1 |
|------------|----------|-----------|----------|----------|----------|
| gini | best | 0.788462 | 0.850575 | 0.953488 | 0.863158 |
| gini | random | 0.784314 | 0.83908 | 0.930233 | 0.851064 |
| entropy | random | 0.736842 | 0.816092 | 0.976744 | 0.84 |
| entropy | best | 0.745455 | 0.816092 | 0.953488 | 0.836735 |

Naive Bayes (NB)

| precision | accuracy | recall | f1 |
|-----------|----------|----------|----------|
| 0.8 | 0.678161 | 0.465116 | 0.588235 |

Support Vector Machine (SVP)

| gamma | kernel | precision | accuracy | recall | f1 |
|-------|---------|-----------|----------|----------|----------|
| scale | rbf | 0.717391 | 0.735632 | 0.767442 | 0.741573 |
| scale | linear | 0.75 | 0.735632 | 0.697674 | 0.722892 |
| auto | linear | 0.75 | 0.735632 | 0.697674 | 0.722892 |
| auto | rbf | 0.702128 | 0.724138 | 0.767442 | 0.733333 |
| scale | sigmoid | 0.647059 | 0.678161 | 0.767442 | 0.702128 |
| auto | sigmoid | 0.647059 | 0.678161 | 0.767442 | 0.702128 |
| auto | poly | 0.772727 | 0.643678 | 0.395349 | 0.523077 |
| scale | poly | 0.833333 | 0.597701 | 0.232558 | 0.363636 |

Multi-Layer Perceptron (MLP)

For sake of brevity, only the top 100 results by accuracy are shown.

| activation | hidden_layer_sizes | learning_rate | max_iter | solver | precision | accuracy | recall | f1 |
|------------|--------------------|---------------|----------|--------|-----------|----------|----------|----------|
| logistic | (60, 80, 100) | constant | 500000 | lbfgs | 0.895833 | 0.942529 | 1 | 0.945055 |
| logistic | (40, 80, 100) | adaptive | 500000 | lbfgs | 0.86 | 0.91954 | 1 | 0.924731 |
| tanh | (40, 80, 100) | invscaling | 500000 | adam | 0.86 | 0.91954 | 1 | 0.924731 |
| tanh | (60, 100, 80) | adaptive | 500000 | lbfgs | 0.86 | 0.91954 | 1 | 0.924731 |
| tanh | (100, 60, 20) | constant | 500000 | adam | 0.86 | 0.91954 | 1 | 0.924731 |
| tanh | (100, 80, 80) | constant | 500000 | adam | 0.86 | 0.91954 | 1 | 0.924731 |
| relu | (75, 30) | adaptive | 500000 | lbfgs | 0.86 | 0.91954 | 1 | 0.924731 |
| logistic | (20, 40, 60) | adaptive | 500000 | lbfgs | 0.875 | 0.91954 | 0.976744 | 0.923077 |
| logistic | (40, 60, 80) | adaptive | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| logistic | (80, 40, 20) | constant | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | 30 | invscaling | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | 60 | constant | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |

| activation | hidden_layer_sizes | learning_rate | max_iter | solver | precision | accuracy | recall | f1 |
|------------|--------------------|---------------|----------|--------|-----------|----------|----------|----------|
| tanh | 85 | adaptive | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (30, 30) | constant | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (45, 45) | adaptive | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (60, 60) | invscaling | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (75, 45) | invscaling | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (75, 75) | adaptive | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (90, 90) | invscaling | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (20, 40, 60) | invscaling | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (20, 100, 20) | invscaling | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (40, 20, 100) | constant | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (40, 80, 60) | invscaling | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (40, 80, 100) | adaptive | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (60, 20, 40) | adaptive | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (60, 60, 80) | constant | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (60, 80, 80) | adaptive | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (80, 20, 40) | adaptive | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (80, 40, 80) | constant | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| tanh | (80, 60, 60) | constant | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| relu | (20, 20, 80) | constant | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| relu | (20, 40, 100) | constant | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| relu | (20, 60, 20) | adaptive | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| relu | (20, 60, 100) | adaptive | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| relu | (20, 100, 20) | constant | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| relu | (20, 100, 40) | adaptive | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| relu | (40, 20, 80) | constant | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| relu | (40, 80, 60) | constant | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| relu | (60, 20, 100) | constant | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| relu | (80, 20, 60) | constant | 500000 | lbfgs | 0.843137 | 0.908046 | 1 | 0.914894 |
| relu | (80, 60, 20) | adaptive | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| relu | (100, 20, 60) | invscaling | 500000 | adam | 0.843137 | 0.908046 | 1 | 0.914894 |
| logistic | (20, 60, 80) | invscaling | 500000 | lbfgs | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| logistic | (60, 20, 20) | adaptive | 500000 | lbfgs | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| tanh | (15, 45) | constant | 500000 | lbfgs | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| tanh | (45, 90) | invscaling | 500000 | lbfgs | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| tanh | (90, 30) | constant | 500000 | lbfgs | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| tanh | (20, 80, 100) | invscaling | 500000 | adam | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| tanh | (20, 80, 100) | adaptive | 500000 | lbfgs | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| tanh | (40, 40, 40) | adaptive | 500000 | lbfgs | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| tanh | (40, 60, 100) | adaptive | 500000 | adam | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| tanh | (60, 80, 60) | constant | 500000 | lbfgs | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| tanh | (100, 40, 60) | invscaling | 500000 | lbfgs | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| tanh | (100, 80, 100) | adaptive | 500000 | adam | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| relu | (30, 30) | adaptive | 500000 | lbfgs | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| relu | (20, 20, 40) | adaptive | 500000 | lbfgs | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| relu | (20, 40, 40) | adaptive | 500000 | adam | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| relu | (40, 20, 100) | adaptive | 500000 | adam | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| relu | (60, 80, 20) | invscaling | 500000 | lbfgs | 0.857143 | 0.908046 | 0.976744 | 0.913043 |
| logistic | (40, 80, 60) | adaptive | 500000 | lbfgs | 0.87234 | 0.908046 | 0.953488 | 0.911111 |
| logistic | 35 | adaptive | 500000 | lbfgs | 0.826923 | 0.896552 | 1 | 0.905263 |
| logistic | (15, 60) | invscaling | 500000 | lbfgs | 0.826923 | 0.896552 | 1 | 0.905263 |
| logistic | (45, 45) | constant | 500000 | lbfgs | 0.826923 | 0.896552 | 1 | 0.905263 |
| logistic | (20, 20, 60) | adaptive | 500000 | lbfgs | 0.826923 | 0.896552 | 1 | 0.905263 |
| logistic | (60, 60, 80) | adaptive | 500000 | lbfgs | 0.826923 | 0.896552 | 1 | 0.905263 |
| logistic | (80, 40, 100) | invscaling | 500000 | lbfgs | 0.826923 | 0.896552 | 1 | 0.905263 |
| logistic | (100, 100, 100) | constant | 500000 | lbfgs | 0.826923 | 0.896552 | 1 | 0.905263 |

| activation | hidden_layer_sizes | learning_rate | max_iter | solver | precision | accuracy | recall | f1 |
|------------|--------------------|---------------|----------|--------|-----------|----------|--------|----------|
| tanh | 60 | constant | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (15, 15) | invscaling | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (15, 45) | adaptive | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (30, 30) | invscaling | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (30, 60) | constant | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (60, 90) | invscaling | 500000 | lbfgs | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (75, 15) | constant | 500000 | lbfgs | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (75, 45) | constant | 500000 | lbfgs | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (90, 15) | adaptive | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (90, 45) | invscaling | 500000 | lbfgs | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (20, 40, 20) | invscaling | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (20, 40, 40) | invscaling | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (20, 60, 20) | adaptive | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (20, 80, 60) | constant | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (20, 80, 80) | constant | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (20, 80, 100) | constant | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (40, 20, 60) | invscaling | 500000 | lbfgs | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (40, 60, 60) | constant | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (40, 60, 60) | invscaling | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (40, 80, 20) | adaptive | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (40, 100, 60) | constant | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (60, 40, 20) | constant | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (60, 40, 40) | invscaling | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (60, 40, 80) | constant | 500000 | lbfgs | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (60, 60, 20) | constant | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (60, 80, 60) | constant | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (60, 80, 80) | invscaling | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (60, 100, 20) | invscaling | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (60, 100, 40) | adaptive | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (60, 100, 60) | constant | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (60, 100, 60) | adaptive | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (60, 100, 80) | constant | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |
| tanh | (80, 40, 40) | constant | 500000 | adam | 0.826923 | 0.896552 | 1 | 0.905263 |

Random Forest (RF)

| criterion | class_weight | max_features | precision | accuracy | recall | f1 |
|-----------|--------------------|--------------|-----------|----------|----------|----------|
| gini | balanced | sqrt | 0.836735 | 0.885057 | 0.953488 | 0.891304 |
| entropy | balanced | sqrt | 0.807692 | 0.873563 | 0.976744 | 0.884211 |
| gini | balanced_subsample | sqrt | 0.807692 | 0.873563 | 0.976744 | 0.884211 |
| entropy | balanced_subsample | sqrt | 0.807692 | 0.873563 | 0.976744 | 0.884211 |
| gini | balanced | log2 | 0.82 | 0.873563 | 0.953488 | 0.88172 |
| entropy | balanced | log2 | 0.82 | 0.873563 | 0.953488 | 0.88172 |
| gini | balanced_subsample | log2 | 0.803922 | 0.862069 | 0.953488 | 0.87234 |
| entropy | balanced_subsample | log2 | 0.803922 | 0.862069 | 0.953488 | 0.87234 |

Evaluation

Output Distributions

Add a boxplot showing mean and standard deviation for **Precision** values on all 6 classifiers (5 trained + 1 biased)

Add a boxplot showing mean and standard deviation for **Recall** values on all 6 classifiers (5 trained + 1 biased)

Add a boxplot showing mean and standard deviation for **F1** values on all 6 classifiers (5 trained + 1 biased)

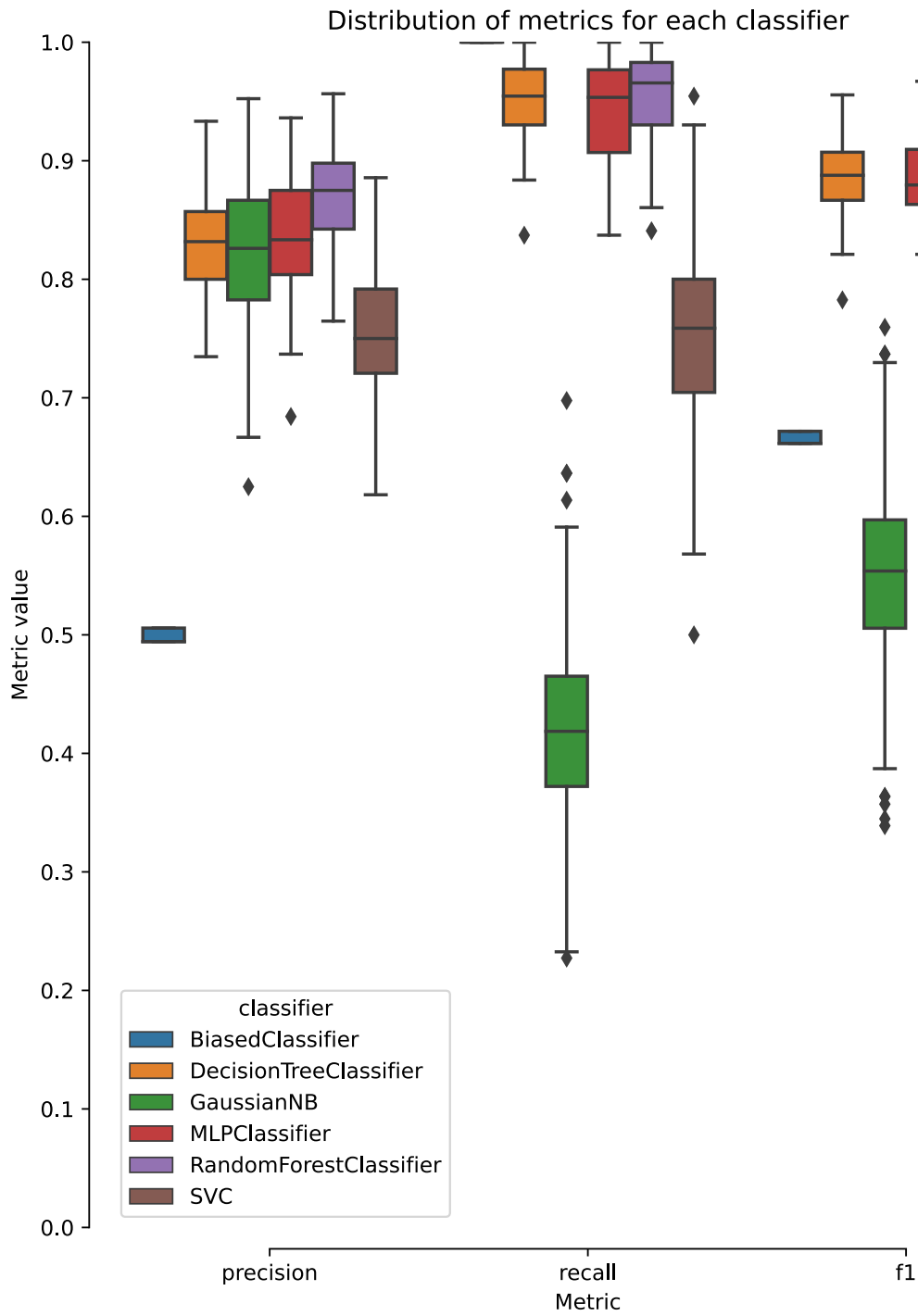


Figure 1: Precision, Recall and F1 score distribution for each classifier for 20-times cross validation.

Comparison and Significance

For every combination of two classifiers and every performance metric (precision, recall, f1) compare which algorithm performs better, by how much, and report the corresponding p-value in the following subsections:

Table 11: Pairwise Wilcoxon test for precision for each combination of classifiers.

| | DecisionTreeClassifier | GaussianNB | MLPClassifier | RandomForestClassifier | SVC |
|------------------------|------------------------|------------|---------------|------------------------|--------|
| BiasedClassifier | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| DecisionTreeClassifier | – | 0.0893 | 0.4012 | 0.0000 | 0.0000 |
| GaussianNB | – | – | 0.0348 | 0.0000 | 0.0000 |
| MLPClassifier | – | – | – | 0.0000 | 0.0000 |
| RandomForestClassifier | – | – | – | – | 0.0000 |

Table 12: Pairwise Wilcoxon test for recall for each combination of classifiers.

| | DecisionTreeClassifier | GaussianNB | MLPClassifier | RandomForestClassifier | SVC |
|------------------------|------------------------|------------|---------------|------------------------|--------|
| BiasedClassifier | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| DecisionTreeClassifier | – | 0.0000 | 0.0118 | 0.3276 | 0.0000 |
| GaussianNB | – | – | 0.0000 | 0.0000 | 0.0000 |
| MLPClassifier | – | – | – | 0.0001 | 0.0000 |
| RandomForestClassifier | – | – | – | – | 0.0000 |

Table 13: Pairwise Wilcoxon test for the F1 score metric for each combination of classifiers.

| | DecisionTreeClassifier | GaussianNB | MLPClassifier | RandomForestClassifier | SVC |
|------------------------|------------------------|------------|---------------|------------------------|--------|
| BiasedClassifier | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| DecisionTreeClassifier | – | 0.0000 | 0.4711 | 0.0000 | 0.0000 |
| GaussianNB | – | – | 0.0000 | 0.0000 | 0.0000 |
| MLPClassifier | – | – | – | 0.0000 | 0.0000 |
| RandomForestClassifier | – | – | – | – | 0.0000 |

F1 Values

-
- ...

Precision

(same as for F1 above)

Recall

(same as for F1 above)

Practical Usefulness

Discuss the practical usefulness of the obtained classifiers in a realistic bug prediction scenario (1 paragraph).