# Predicting the bug fixing likelihood

Florian Spychiger

University of Zurich

December 9, 2018

# Roadmap

# Background Information

The annual cost of software bugs is estimated at \$59.5 billion[1]. For the Eclipse project, there are thousands of bugs reported. An efficent bug-triaging can help developers to focus their resources and thus, save companies a lot of money.

---

[1]P Bhattacharya and I Neamtiu, "Fine--grained incremental learning and multi--feature tossing graphs to improve bug triaging", Software Maintenance (ICSM) 2010 (ieeexplore.ieee.org)

# Problem Formulation & Goal

### Problem

Bug-triaging is an important, but labor-intensive process if done manually.

### Goal

Train a bug-triaging machine, which predicts whether a bug is likely to be fixed.

## Raw data

The Eclipse data set can be found at
https://github.com/ansymo/msr2013-bug_dataset.
The raw data set consists of 12 tables:

| Eclipse Bug Data Set | |
| --- | --- |
| reports | priority |
| assigned_to | product |
| bug_status | resolution |
| cc[2] | severity |
| component | short_desc |
| op_sys | version |

---

[2]The data has been newly formated with Excel VBA.

## Data preselection

After a visual exploratory analysis, four datasets were excluded:

| Eclipse Bug Data Set* | |
|---|---|
| reports | priority ❶ |
| assigned_to | product |
| bug_status | resolution |
| cc | severity ❷ |
| component | short_desc ❸ |
| op_sys | version ❹ |

*Duplicates are excluded.

❶ Priority is set by the assignee, but as we want to help them triaging the bug, we exlude it.

❷ Severity is currently set by the triaging team.

❸ The descriptions are hard to encode.

❹ The version dataset is quite messy and sometimes it is not clear which version is being referred to.
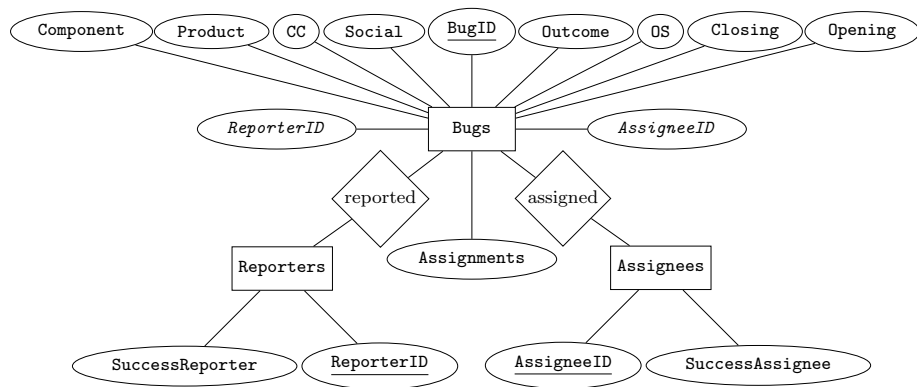
# Data model



Figure: ER model of data used.

## Feature Creation

From the data model, the feature matrix $X$ is constructed with:

$x_1 = $ OpeningTime (Open - Close ) $\qquad$ [discrete]

$x_2 = $ Assignments Nr. of assignees) $\qquad$ [discrete]

$x_3 = $ CCs (Nr. of interested parties) $\qquad$ [discrete]

$x_4 = $ Product (Affected product) $\qquad$ [discrete]

$x_5 = $ OS (Major OS) $\qquad$ [discrete]

$x_6 = $ SuccessAssignee (Success rate of Assignee) $\qquad$ [proportion]

$x_7 = $ SuccessReporter (Success rate of Reporter) $\qquad$ [proportion]

$x_8 = $ Component (The affected subcomponent) $\qquad$ [discrete]

$x_9 = $ Social (Past bug collaborations) $\qquad$ [binary]

$x_{10} = $ Equal (Reporter equals Assignee) $\qquad$ [binary]

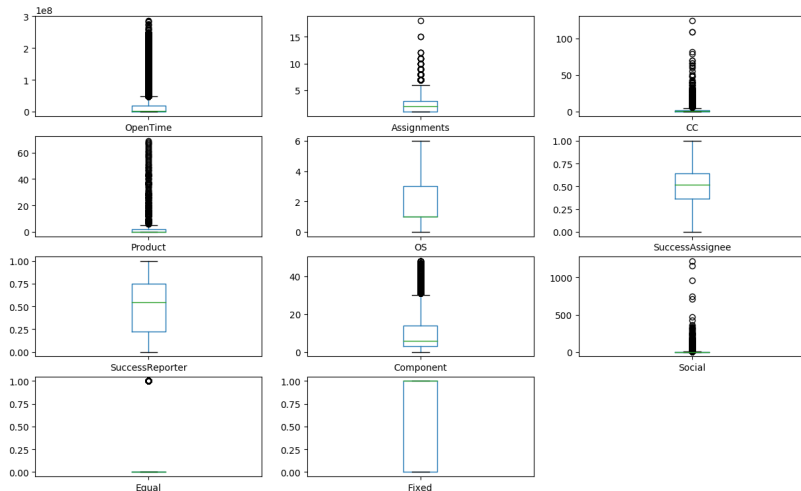# Univariate Analysis



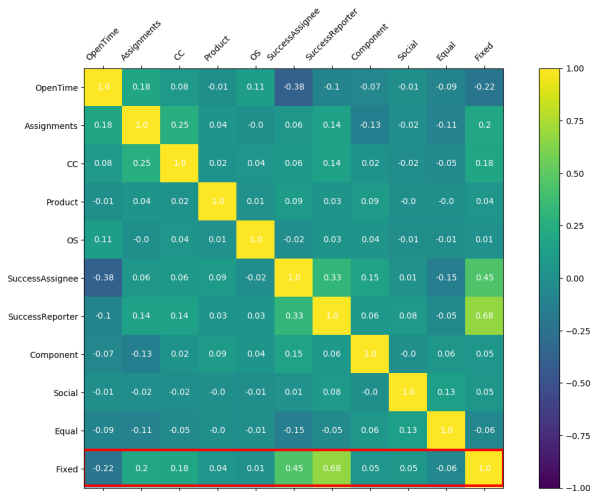Figure: ER model of data used.

# Correlation Analysis



Figure: ER model of data used.

## Models

We consider 6 models:

1. Naive Bayes
2. Logistic Regression
3. Random Forest
4. Boosting Classifier
5. Support Vector Machine
6. Neural Network

We split the data set into a training (50%), a cross-validation (25%) and a test (25%) set. The training set is used to train the models and we calibrate the parametes on the cross-validation set. The final accuracy is caculated on the test set.

## Accuracy

We achieve the following accuracies on the test set:

| | |
|---|---|
| Naive Bayes | 82.8098% |
| Logistic Regression | 84.9409% |
| Random Forest | 86.1529% |
| Boosting Classifier | 85.4661% |
| Support Vector Machine | 85.9105% |
| Neural Network[3] | 86.1125% |

---

[3]Results are not exactly reproducible, as some randomness with GPU usage cannot be avoided.
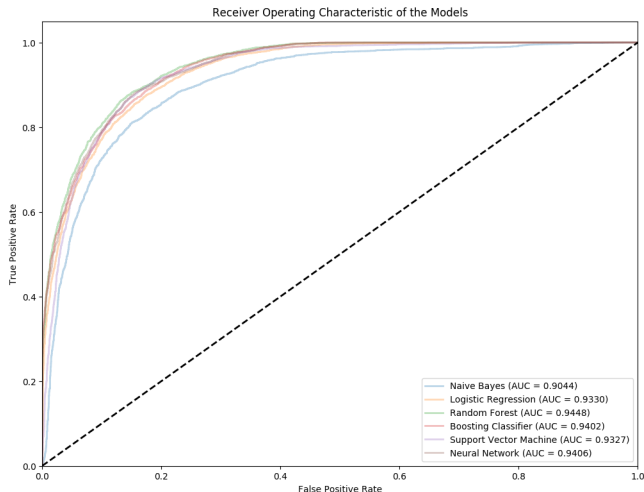
# ROC-Curves



Figure: ER model of data used.

# Thank you!

The code of the project can be found at

https://github.com/Speaker90/BusinessAnalytics_RPIcase