# Duplicate bug report detection through machine learning techniques

Irving Muller Rodrigues
December 10, 2018

Prof. Daniel Aloise and Prof. Michel Dagenais

# Introduction

# Introduction

# Bug Tracking System

# Bug Tracking System

# Bug Tracking System



User
Tester
Developer

Report

Bug Report

- Incomplete Bugs
- Invalid Bugs
- Duplicate Bugs

# Bug Tracking System

# Bug Tracking System

### Triage Process

- Manual checking
- Time and money consuming
- Large user base project: Firefox ~300 new reports per day

# Objective



- Increase software quality and save resource
  - Decrease triage team overload
  - Avoid two or more developers fixing the same bug
  - Avoid to fix a bug already solved

# Duplicate bug report detection

- Detect whether a bug is duplicate or not
- Master set
  - Master report
  - Duplicate reports
  - Every report is in a master set
- Three approaches
  - Decision-making approach
  - Binary classification approach
  - Ranking approach

# Decision-making approach

- Pairs of bug reports (Training and Evaluation)
- Drawbacks
  - Too Easy
  - High probability to create easy non-duplicate pairs
  - Far from the real scenario
    - Compare new bug with a set of bugs in the dataset

# Binary classification approach

- Automatic prediction of the report as duplicate or not
  - General information extracted from the database and the new bug reports
- False negative can have a great impact
- Really difficult task

Bug Report $X$ → **System** → Prediction

# Ranking approach

- Recommend a similarity list
- A person check the list and label the report as duplicate or not
  - Decrease the decision time
- The most used approach in the literature
- Metric: Recall Rate
  - Rate of reports whose the lists have at least one bug report from the same master set

# Ranking approach

- Two methodologies: Deshmukh et al. 2017 and Sun et al. 2011
- Deshmukh et al. 2017
  - Training, validation and test datasets are randomly generated
  - Evaluation: similarity list are created using bug from the test dataset
  - Unrealistic scenario
  - It makes the problem easier
    - Decrease number of comparisons
    - Concept Drift mitigation
- Sun et al. 2011
  - Reports are sorted by creation date
  - Training, validation and test are generate by period of time
  - New bug report is compared with all previous bug reports
  - More realistic scenario

# Our Solution

- Ranking approach + Sun's Methodology
- Only textual data
  - Summary and description
- Baseline: TF-IDF
- Model: Word Embeddings + Convolution Neural Network

# TF-IDF

Document

Content

adapter creation
gets broken

| Term | Value |
|------|-------|
| adapter | $w_1$ |
| gets | $w_2$ |
| broken | $w_3$ |
| creation | $w_4$ |

# TF-IDF

Document

Content

adapter creation
gets broken

| Term | Value |
|------|-------|
| adapter | $w_1$ |
| gets | $w_2$ |
| broken | $w_3$ |
| creation | $w_4$ |

$w_4$ = Term Frequency x Inverse Document Frequency

# TF-IDF

Content

Document

adapter creation
gets broken

| Term | Value |
|------|-------|
| adapter | $w_1$ |
| gets | $w_2$ |
| broken | $w_3$ |
| creation | $w_4$ |

$w_4$ = Term Frequency x Inverse Document Frequency

# TF-IDF

Document

Content

adapter creation
gets broken

| Term | Value |
|------|-------|
| adapter | $w_1$ |
| gets | $w_2$ |
| broken | $w_3$ |
| creation | $w_4$ |

$w_4 = 1$      x Inverse Document Frequency

# TF-IDF

Document

Content

adapter creation
gets broken

| Term | Value |
|------|-------|
| adapter | $w_1$ |
| gets | $w_2$ |
| broken | $w_3$ |
| creation | $w_4$ |

$$w_4 = 1 \quad x \text{ Inverse Document Frequency}$$

$$\log\left(\frac{\text{Number of documents}}{\text{Document Frequency}}\right)$$

# TF-IDF

Document

Content

adapter creation
gets broken

| Term | Value |
|------|-------|
| adapter | $w_1$ |
| gets | $w_2$ |
| broken | $w_3$ |
| creation | $w_4$ |

$w_4 = 1$    x Inverse Document Frequency

$$\log\left(\frac{10}{8}\right)$$

# TF-IDF

Document

Content

adapter creation
gets broken

| Term | Value |
|---|---|
| adapter | $w_1$ |
| gets | $w_2$ |
| broken | $w_3$ |
| creation | 0.09 |

# Represent word as vector

- **Word Embedding**
  - Dense vectors with real numbers
  - More compact representation
  - Semantic and syntactic information

| Word | Representation |
|------|----------------|
| adapter | [0.5, 0.6] |
| broken | [0.3, 0.2] |
| gets | [0.1, 0.7] |
| creation | [0.6, 0.3] |

# Convolution Neural Network for NLP

Input

|          |     |     |
|----------|-----|-----|
| adapter  | 0.5 | 0.6 |
| creation | 0.6 | 0.3 |
| gets     | 0.1 | 0.7 |
| broken   | 0.6 | 0.3 |

Filters

Filter 1

| 1. | 2. |
|----|----|
| 5. | 1. |
| 2. | 3. |

Filter 2

| 2. | 5. |
|----|----|
| 1. | 1. |
| 1. | 1. |

# Convolution Neural Network for NLP

# Convolution Neural Network for NLP

Input

| | | |
|---|---|---|
| adapter | 0.5 | 0.6 |
| creation | 0.6 | 0.3 |
| gets | 0.1 | 0.7 |
| broken | 0.6 | 0.3 |

Filter 1

$\odot$

| 1. | 2. |
|---|---|
| 5. | 1. |
| 2. | 3. |

$= \text{sum} \,(\,$

| 0.5 | 1.2 |
|---|---|
| 3.0 | 0.3 |
| 0.2 | 2.1 |

$\,) =$ | 7.3 |

# Convolution Neural Network for NLP

Input

|          |     |     |
|----------|-----|-----|
| adapter  | 0.5 | 0.6 |
| creation | 0.6 | 0.3 |
| gets     | 0.1 | 0.7 |
| broken   | 0.6 | 0.3 |

$\odot$

Filter 1

| 1. | 2. |
|----|----|
| 5. | 1. |
| 2. | 3. |

$= \text{sum} ($

| 0.6 | 1.6 |
|-----|-----|
| 0.5 | 0.7 |
| 1.2 | 0.9 |

$) = $ 5.5

# Convolution Neural Network for NLP

# Convolution Neural Network for NLP

# Our Deep Learning Model

- Encoder
  - Represent the report as vector

# Our Deep Learning Model



P(D)

Output Layer

Hidden Layer

Hidden Layer

$v^1$ $\bullet\bullet\bullet\bullet\bullet$     $v^2$ $\bullet\bullet\bullet\bullet\bullet$     $|v^1 - v^2|$     $v^1 \odot v^2$

Encoder     Encoder

Bug Report 18042     Bug Report 137861

# Our Deep Learning Model



**Cross Entropy**

$y \times \log(P(D)) + (1 - y) \log(1 - P(D))$

$P(D)$

Output Layer

Hidden Layer

Hidden Layer

$v^1$ ●●●●●  $v^2$ ●●●●●  $|v^1 - v^2|$  $v^1 \odot v^2$

Encoder  Encoder

Bug Report 18042  Bug Report 137861

# Preliminar Results

| Model | Top-5 | Top-10 | Top-15 | Top-20 |
|-------|-------|--------|--------|--------|
| TF-IDF | 44.80% | 51.27% | 54.97% | 57.88% |
| DL Model | 37.11% | 43.95% | 48.61% | 52.03% |

# Our Deep Learning Model

- Challenge:
  - Generate relevant non-duplicate pairs (negative) can be difficult
  - Most non-duplicate pairs are easy
  - ~ $n^2$ different combinations
  - n = 174,002 ⇨ $n^2 \cong 30 \times 10^9$
- Solution: Random subsample negative examples each epoch
  - Constraint: loss has to be greater than 0
  - Keep rate between positive and negative examples

# Preliminar Results

| Model | Top-5 | Top-10 | Top-15 | Top-20 |
|---|---|---|---|---|
| TF-IDF | 44.80% | 51.27% | 54.97% | 57.88% |
| DL Model | 37.11% | 43.95% | 48.61% | 52.03% |
| DL Model - subsampling by epoch | 44.02% | 51.03% | 55.49% | 58.43% |

# Preliminar Results

| Model | Top-5 | Top-10 | Top-15 | Top-20 |
|---|---|---|---|---|
| TF-IDF | 44.80% | 51.27% | 54.97% | 57.88% |
| DL Model | 37.11% | 43.95% | 48.61% | 52.03% |
| DL Model - subsampling by epoch | 44.02% | 51.03% | 55.49% | 58.43% |

6.40%

# Future Work

- Bottleneck: select negative pairs
  - Try different approaches
- Encoder receives information from the first bug
  - Attention
- Combine different information sources
  - Categorical information, stack trace, tracing
- Use our solution to help our partners
  - Partner data

# Thank you for your attention!
## Questions?

Irving Muller Rodrigues
irving.muller-rodrigues@polymtl.ca

# References

- Deshmukh, J., M, A. K., Podder, S., Sengupta, S., & Dubash, N. (2017). Towards Accurate Duplicate Bug Retrieval Using Deep Learning Techniques. 2017 IEEE International Conference on Software Maintenance and Evolution (ICSME), 115–124. http://doi.org/10.1109/ICSME.2017.69
- Lazar, A., Ritchey, S., & Sharif, B. (2014). Generating duplicate bug datasets. Proceedings of the 11th Working Conference on Mining Software Repositories - MSR 2014, 392–395. http://doi.org/10.1145/2597073.2597128
- Sabor, K. K., Hamou-Lhadj, A., & Larsson, A. (2017). DURFEX: A feature extraction technique for efficient detection of duplicate bug reports. Proceedings - 2017 IEEE International Conference on Software Quality, Reliability and Security, QRS 2017, 240–250. http://doi.org/10.1109/QRS.2017.35

# References

- Anh Tuan Nguyen, Tung Thanh Nguyen, Tien N Nguyen, David Lo, and Chengnian Sun. Duplicate bug report detection with a combination of information retrieval and topic modeling. In Automated Software Engineering (ASE), 2012 Proceedings of the 27th IEEE/ACM International Conference on, pages 70–79. IEEE, 2012.
- Klaus Greff, Rupesh Kumar Srivastava, Jan Koutník, Bas R. Steunebrink, Jürgen Schmidhuber. LSTM: A Search Space Odyssey. CoRR abs/1503.04069 (2015)
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In Advances in neural information processing systems (pp. 3111-3119).

# References

- Kim, Yoon. "Convolutional Neural Networks for Sentence Classification." Paper presented at the meeting of the Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL, 2014.
- C. Sun, D. Lo, S. Khoo and J. Jiang, "Towards more accurate retrieval of duplicate bug reports," 2011 26th IEEE/ACM International Conference on Automated Software Engineering (ASE 2011), Lawrence, KS, 2011, pp. 253-262.

# Represent word as vector

- **One hot encoding**
  - Binary Vectors
  - Vector Size = Vocabulary Size
  - Curse of Dimensionality

| Word | Representation |
|---|---|
| adapter | [1,0,0,0] |
| broken | [0,1,0,0] |
| gets | [0,0,1,0] |
| creation | [0,0,0,1] |

# TF-IDF

Document

Content

adapter creation
gets broken

| Term | Value |
|------|-------|
| adapter | $w_1$ |
| gets | $w_2$ |
| broken | $w_3$ |
| creation | $w_4$ |

$w_4$ = Term Frequency x Inverse Document Frequency

$$\log\left(\frac{\text{Number of documents}}{\text{Document Frequency}}\right)$$