

An Approach to Detecting Duplicate Bug Reports using Natural Language and Execution Information

Xiaoyin Wang¹, Lu Zhang^{1*}, Tao Xie^{2*}, John Anvik³ and Jiasu Sun¹

¹Key laboratory of High Confidence Software Technologies, Ministry of Education, Institute of Software, EECS, Peking University, Beijing, 100871, P. R. China, {wangxy06, zhanglu, sjs}@sei.pku.edu.cn

²Department of Computer Science, North Carolina State University, Raleigh, NC 27695, USA
xie@csc.ncsu.edu

³Department of Computer Science, University of Victoria
PO Box 3055, STN CSC
Victoria, BC, Canada, V8W 3P6
janvik@cs.uvic.ca

ABSTRACT

Categories and Subject Descriptors
Software Engineering

General Terms

Keywords

1. INTRODUCTION

*

2.1 Browser-Closing Bug

Bug-260331 *Bug-239223*

Bug-260331

Bug-239223

Bug-260331: After closing Firefox, the process is still running. Cannot reopen Firefox after that, unless the previous process is killed manually

Bug-239223: (Ghostproc) – [Meta] firefox.exe doesn't always exit after closing all windows; session-specific data retained

•

•

•

•

2.2 Document-Contain-No-Data Bug

Bug-244372 *Bug-219232*

Bug-244372: "Document contains no data" message on continuation page of NY Times article

Bug-219232: random "The Document contains no data." Alerts

Bug-244372

Bug-219232

Bug-219232

2. MOTIVATING EXAMPLES

4. THE PROPOSED APPROACH

2.3 Motivation

3. BACKGROUND

$$w_i = tf_i \times idf_i$$

tf_i term frequency
 idf_i inverse document frequency

$$idf_i = \log(D_{sum} / D_{w_i})$$

$$Sim = \frac{\sum_{i=1}^n w_{i1} w_{i2}}{\sqrt{\sum_{i=1}^n w_{i1}^2 \times \sum_{i=1}^n w_{i2}^2}}$$

4.1 Calculating NL-S

4.2 Calculating E-S

■

■

■

4.3 Retrieving Potential Target Bug Reports

4.3.1 Basic Heuristic

$$SIM_{combined} = f(SIM_{nlp}, SIM_{exe})$$

$$SIM_{combined} = \frac{SIM_{nlp} + SIM_{exe}}{2}$$

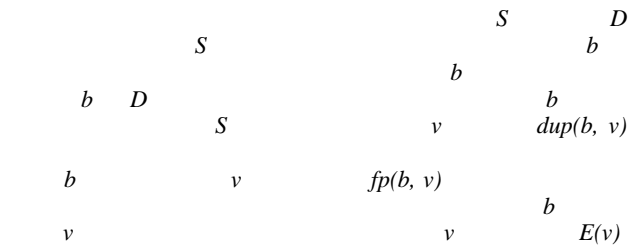
SIM_{exe} f

v v

$$SIM_{combined} = \frac{SIM_{nlp} + SIM_{exe}}{2}$$

4.3.3 Determining Credibility Thresholds

4.3.2 Classification-Based Heuristic



$$E(v) = \sum_{b \in D} dup(b, v) - fp(b, v)$$

4.4 Presenting Potential Target Bug Reports

■

5. EXPERIMENT

-

$$\text{recall rate} = \frac{N_{\text{recalled}}}{N_{\text{total}}}$$

-

5.1 Experimental Setup

-

5.2 Calibration and Evaluation on Eclipse

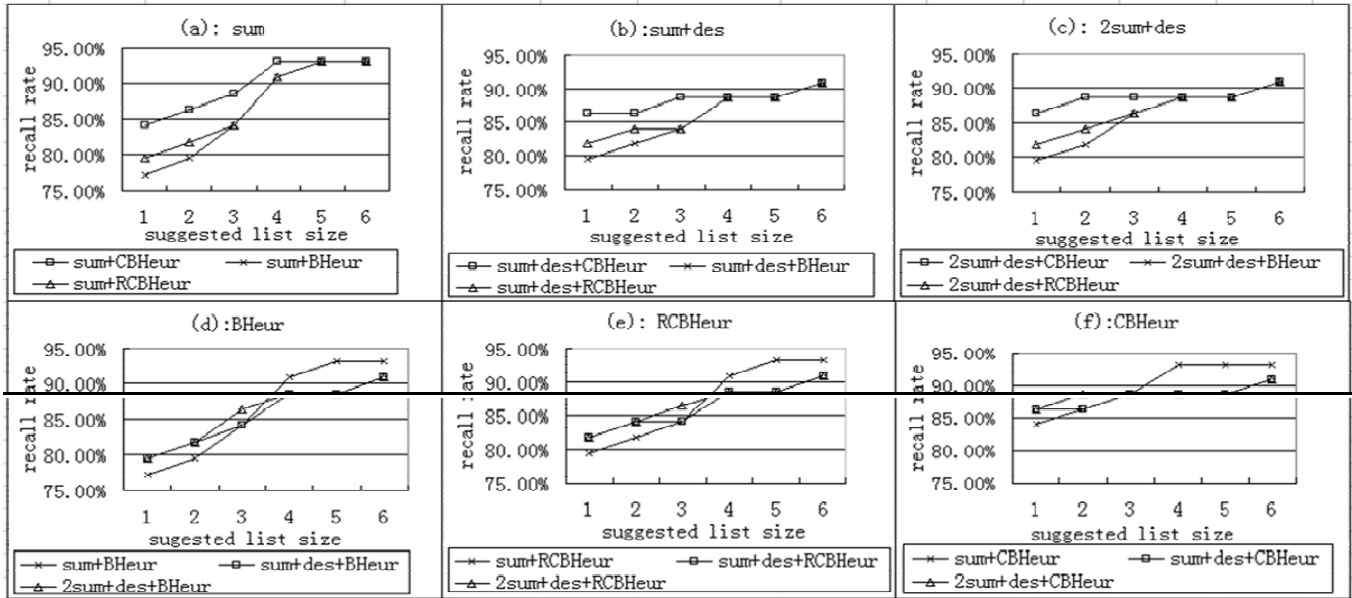


Figure 1: Recall rates using different parameters in Eclipse

5.3 Evaluation on Firefox

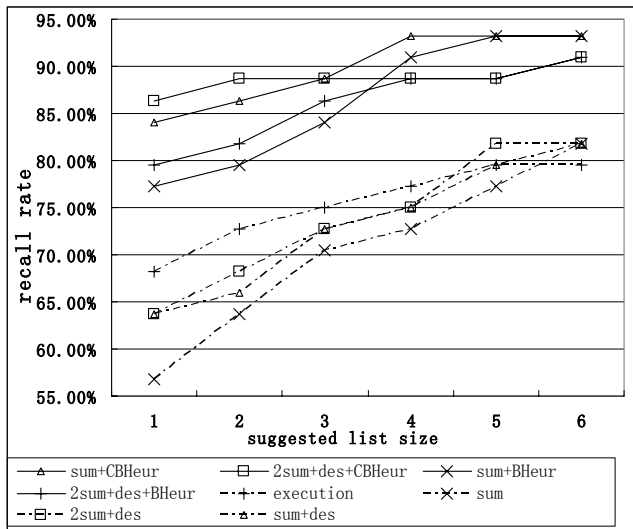


Figure 2: Recall rates using different similarities in Eclipse

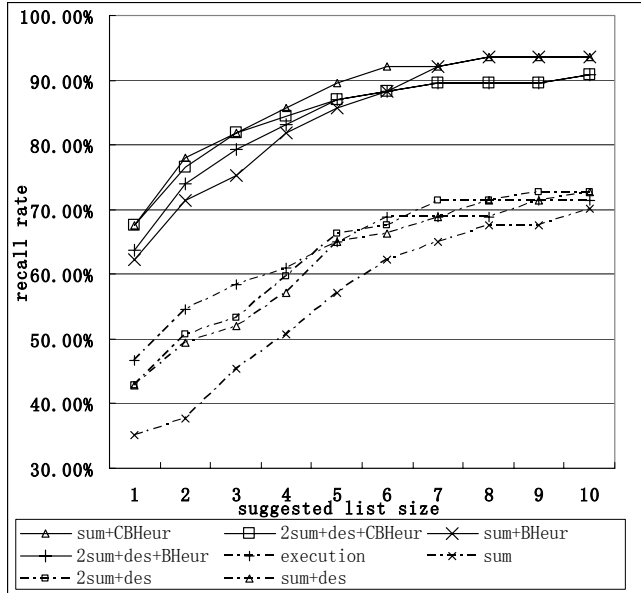


Figure 3: Recall rates using different similarities in Firefox

6. DISCUSSION

6.1 Costs of Using Execution Information

idf

idf

5.4 Threats to Validity

n
 n m
 m

6.2 Coping with Existing Bug Repositories

7. RELATED WORK

$$AvgSim_{exe} i = \frac{1}{n} \sum_{x \in R} Sim_{exe} i x$$

R i
 n R
 i

6.3 Evaluation Criteria for Duplicate-Bug-Report Detection

10. REFERENCES

In Proc. ICSE.

*In Proc. of OOPSLA Workshop on Eclipse
Technology eXchange (ETX)*

In Proc. ICSE

In Proc. SEKE

IEEE TSE

In Proc. ISSRE

8. FUTURE WORK

Assisted Detection of Duplicate Bug Reports

In Proc. ISSTA

*In
Proc. of IEEE Conf. on Visual Language and Human-Centric
Computing (VL/HCC)*

In Proc. PLDI

In Proc. ICSM,

9. CONCLUSION

ACM TOSEM

In Proc. ICSE

*Journal of the American So-
ciety for Information Science*

Proc. ICSE

In

In Proc. ICSE,

Acknowledgments

In Proc. MSR,