

Class Projects



Group Assignment

Group 1. Binary Classification of Duplicate Issues / PRs

Group 2. Classification of Issues

Group 3. Refactoring and SATD

Group 4. Prediction of who should refactor the code

Group 5. Detection of Performance Regression

Group 6. Classification of Commit Messages

Group 7. Identification of Extract Method

Group 8. Code Review Performance Analysis using Bug Report



About Projects

- Each topic has an empirical / exploratory component into it, then a predictive modeling.
 1. The first part helps students better understand their data and reveal any hidden dependencies, characteristics.
 2. The second part helps in using these characteristics as features to extract patterns that the machine can learn.



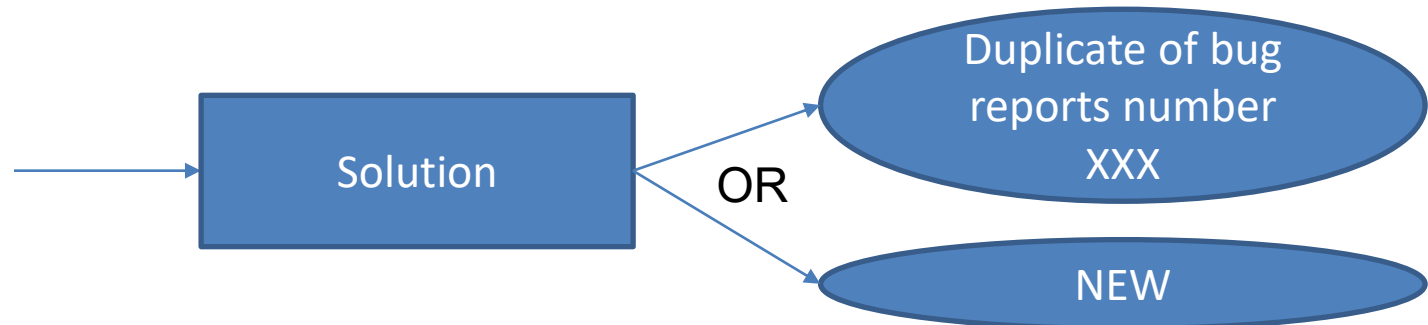
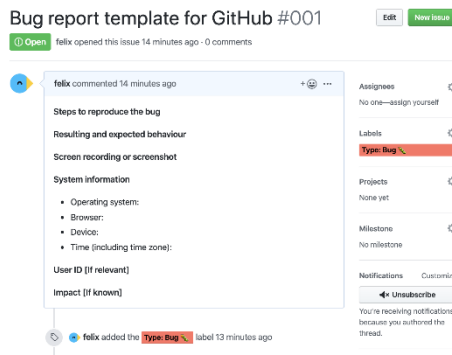
Binary Classification of Duplicate Issues / PRs

Problem?

- Many issue trackers receive bug reports of a similar bug that has been previously reported. Triaging these reports for a known bug is time consuming. Is there a way to help triagers discover whether a newly reported bug is already saved in the tracker, and being addressed by the team?

Solution

- Given a new issue tagged as “bug report”, compare it with all the open bug reports, and if it is highly similar to one of them, then flag it as potential duplicate.



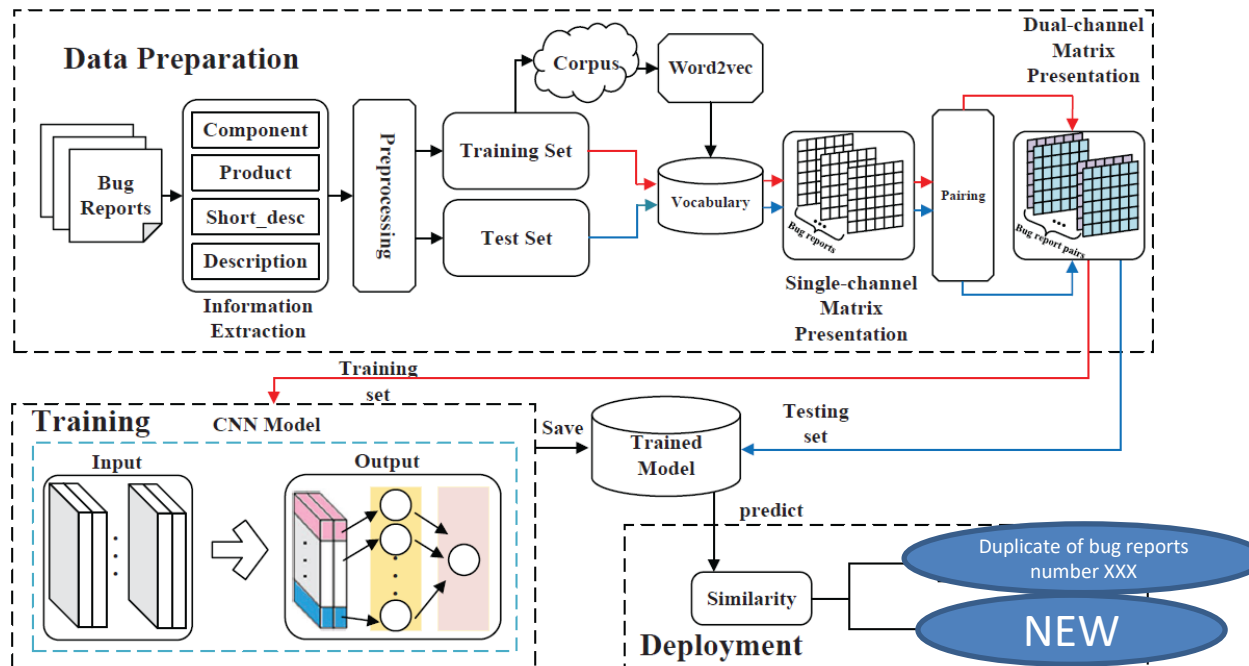
Binary Classification of Duplicate Issues / PRs

Empirical

- Collect information about duplicate issues and PRs. This will help in better understanding the feasibility of the duplicate detection process.

Recommendation

- Use appropriate solution to solve the problem. Example:



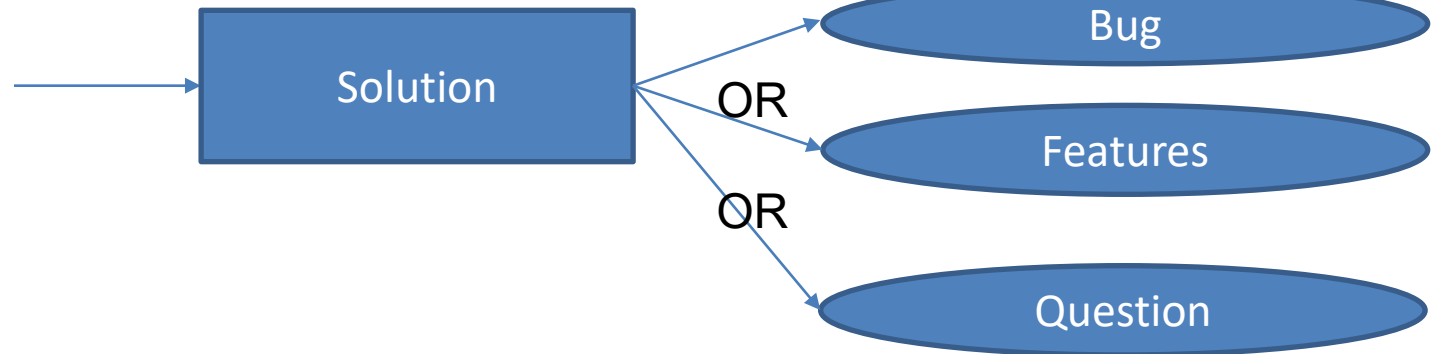
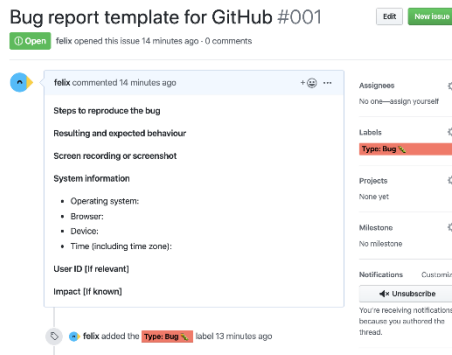
Classification of Issues

Problem?

- Many issue trackers receive various issues in need to be tagged. This manual process can be automated if we can leverage previously manually labeled issues.

Solution

- Given a new issue, can we automatically assign a tag?



Classification of Issues

Empirical

- Collect information about all issues tags (given a dataset). Analyze the distribution of each tag to have a sense of the data balance. Then choose appropriate solution to recommend a TAG.

Recommendation

- Use appropriate solution to solve the problem. Example:



*fast*Text

5,862 Open ✓ 68,774 Closed		Author ▾	Labels ▾	Projects ▾	Milestones ▾	Assignee ▾	Sort ▾
1	Avoid shadowed variable declarations	#80311 opened 20 days ago by aeschli	feature-request under-discussion				3
1	scripts/test.sh --debug crashes in await	#80310 opened 20 days ago by bpasero	engineering upstream				3
1	The warning about deleting file under version control is too strong / unnecessary	#80309 opened 20 days ago by borekb	feature-request file-explorer				7
1	Web Smoke Test: Run against produced build	#80308 opened 20 days ago by bpasero	broken-build debt		September 2019		
1	Proposal for selecting language mode not using file extension	#80303 opened 20 days ago by TonyGravagno					2
1	[Feature Request] snippet support RANDOM	#80297 opened 20 days ago by seognil	feature-request help wanted snippets				1
1	Extensions can potentially block autocompletion from other extensions	#80295 opened 20 days ago by pavelvlavov	needs more info suggest under-discussion				4
1	textDocument/codeAction request has end character set as 1.7976931348623157E+308	#80288 opened 21 days ago by rjmholt	bug		October 2019		2
1	Change application icon to indicate notifications are present	#80286 opened 21 days ago by SteveBenz	workbench-notifications workbench-os-integration				4
1	Comment buttons/actions too far away on a wide screen	#80276 opened 21 days ago by evangrayk	comments ux				
1	File not syncing between duplicate workspaces	#80275 opened 21 days ago by karrtikr					1
1	Extension activationEvent onLanguage:plaintext always fires, even if no plaintext files are open or a non-plaintext file is opened	#80268 opened 21 days ago by valentin					9
1	Debug with arguments without launch.json	#80267 opened 21 days ago by anatologyourakov	debug feature-request needs more info				4

GitHub issue tracker of Microsoft's vscode



Refactoring and SATD

Problem?

- To what extent refactoring helps in reducing technical debt? We want to investigate whether refactoring activity correlates with the removal of technical debt.

Solution

- Given a set of SATD comments removals and their associated refactorings, use the statistical analysis to identify whether there is a correlation between the existence of refactoring and the removal of SATD comment. Also, explore what type of refactoring is popular in removing SATD comments. Then, design and implement a solution to learn from them to identify, for an SATD to be removed what is the appropriate refactoring type to apply



The screenshot shows a commit message and a code change. The commit message is "Rename getter function to hopefully be a bit more descriptive..." by user "master" (bekvon) committed on Aug 22, 2012. The code change is in the file "src/com/bekvon/bukkit/residence/protection/FlagPermissions.java". The change shows a refactoring of a getter function. The original code (lines 41-43) is: `public static void removeMaterialFromUseFlag(Material mat) { matUseFlagList.remove(mat); }`. The new code (lines 44-47) is: `public static EnumMap<Material,String> getMaterialToUseList(){ public static EnumMap<Material,String> getMaterialUseFlagList(){ return (EnumMap<Material, String>) matUseFlagList; }`. The change is highlighted in green.

Solution

Rename



Refactoring and SATD

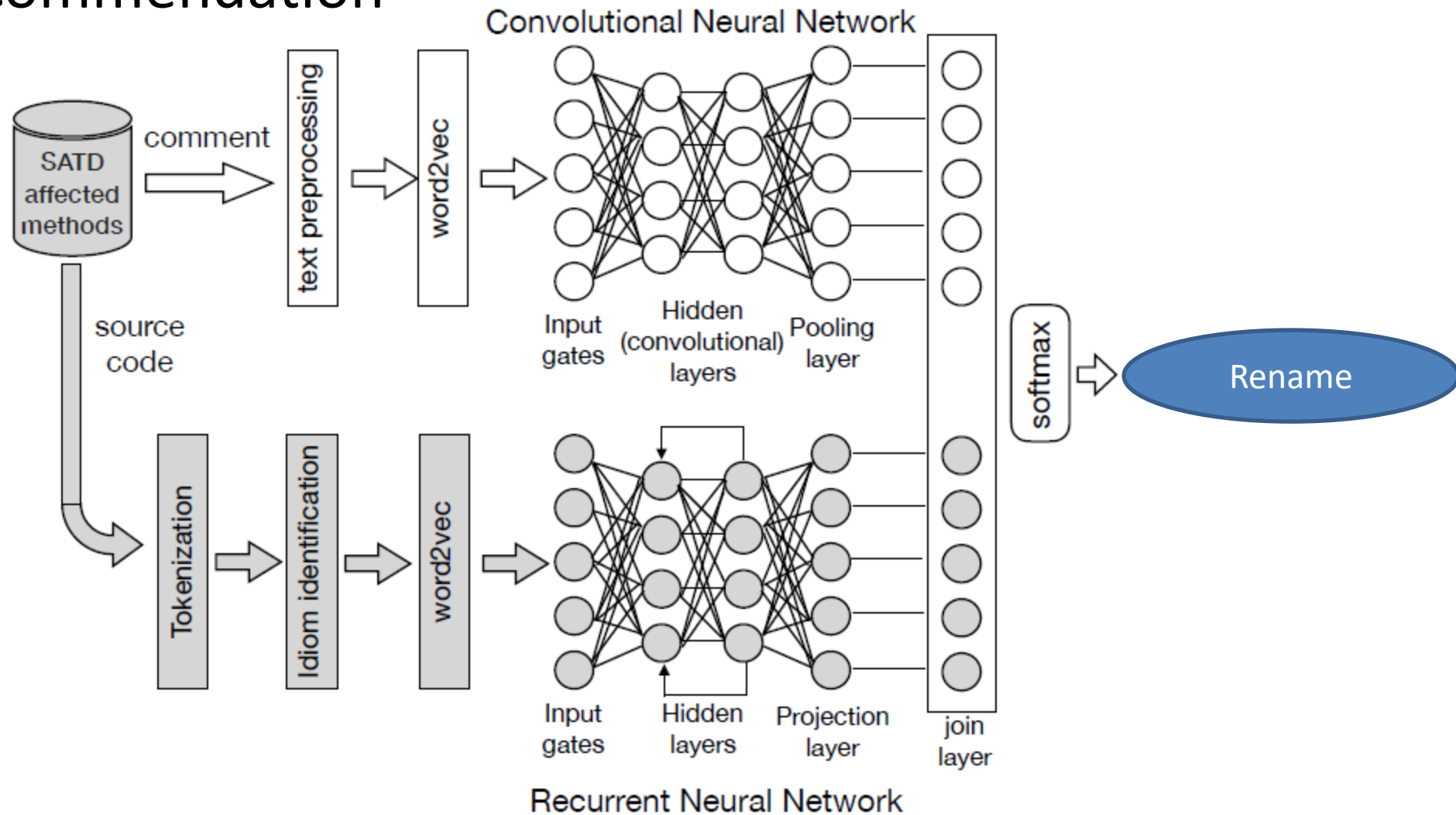
Empirical

- Collect information about colocation of refactoring and SATD removal (same commit).
- Analyze this subset of commits to see if the distribution of refactoring in SATD removal is uniform compared to the one in other mainstream commits.
- Analyze the distribution of SATD removal refactoring per type (frequency).



Refactoring and SATD

Recommendation




Prediction of who should refactor the code

Problem?

- Given a code change, can we identify who made it? Can we identify who is more suitable to make a given change in the source code?

Solution

- Given a set of code changes and their associated authors, design and implement a solution to learn from them to identify, for a new code change the appropriate developer who can perform it



```
Commit Message
Rename getter function to hopefully be a bit more descriptive...
master
bekvon committed on Aug 22, 2012

Commit Code Change snippet
src/com/bekvon/bukkit/residence/protection/FlagPermissions.java
@@ -41,7 +41,7 @@ public static void addMaterialToUseFlag(Material mat, String flag) {
41 41     public static void removeMaterialFromUseFlag(Material mat) {
42 42         matUseFlagList.remove(mat);
43 43     }
44 -    public static EnumMap<Material, String> getMaterialToUseList(){
44 +    public static EnumMap<Material, String> getMaterialUseFlagList(){
45 45         return (EnumMap<Material, String>) matUseFlagList;
46 46     }
47 47     public static void addFlag(String flag) {
```

Solution

XYZ

bekvon

ABC

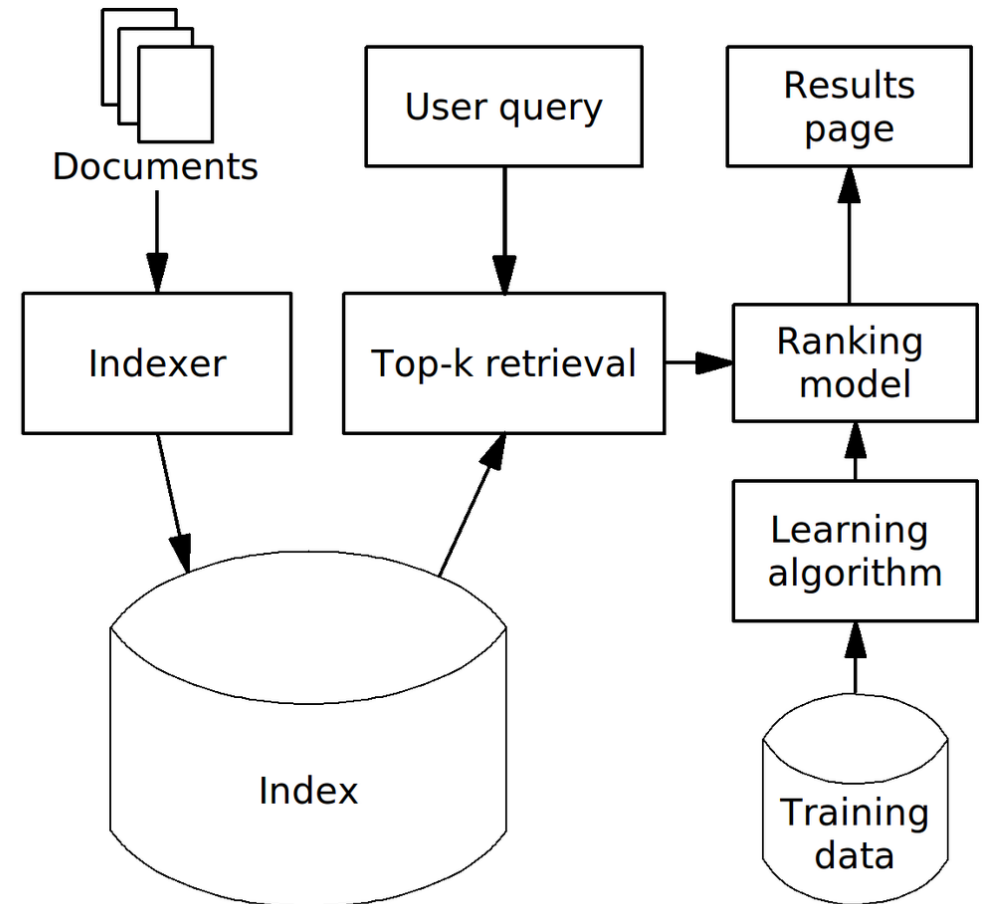
Prediction of who should refactor the code

Empirical

- Analyze the distribution of change per authors (developers) and the uniformity of its distribution.

Recommendation

- Tackle it as a ranking problem?



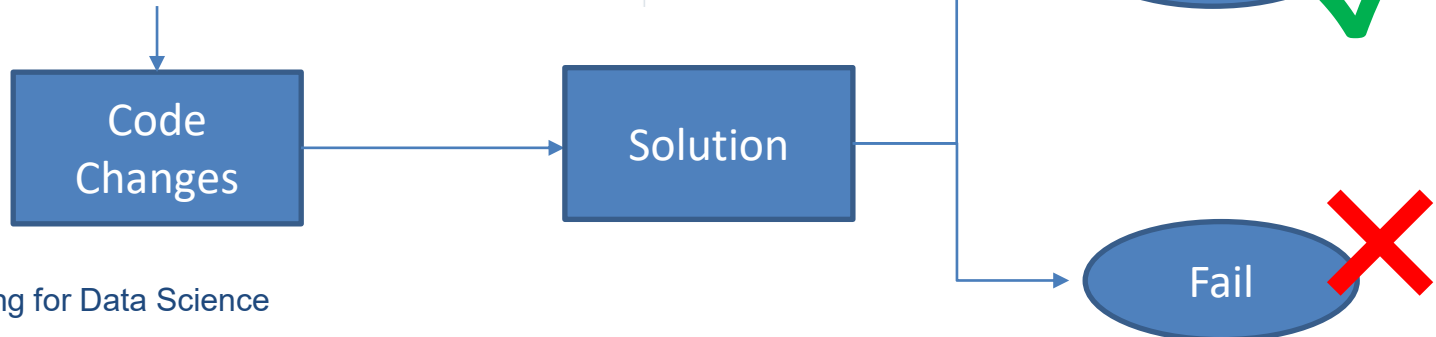
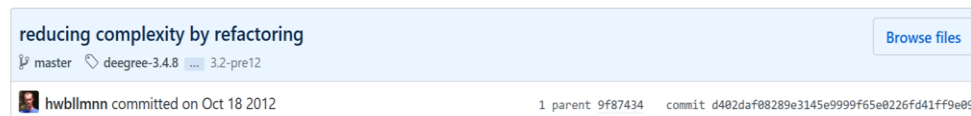
Prediction of performance regression

Problem?

- Given a code change, can we identify whether it can introduce performance regression?

Solution

- Given a set of code changes and their associated performance test results, design and implement a solution to learn from them to identify, for a new code change, whether it can trigger performance regression



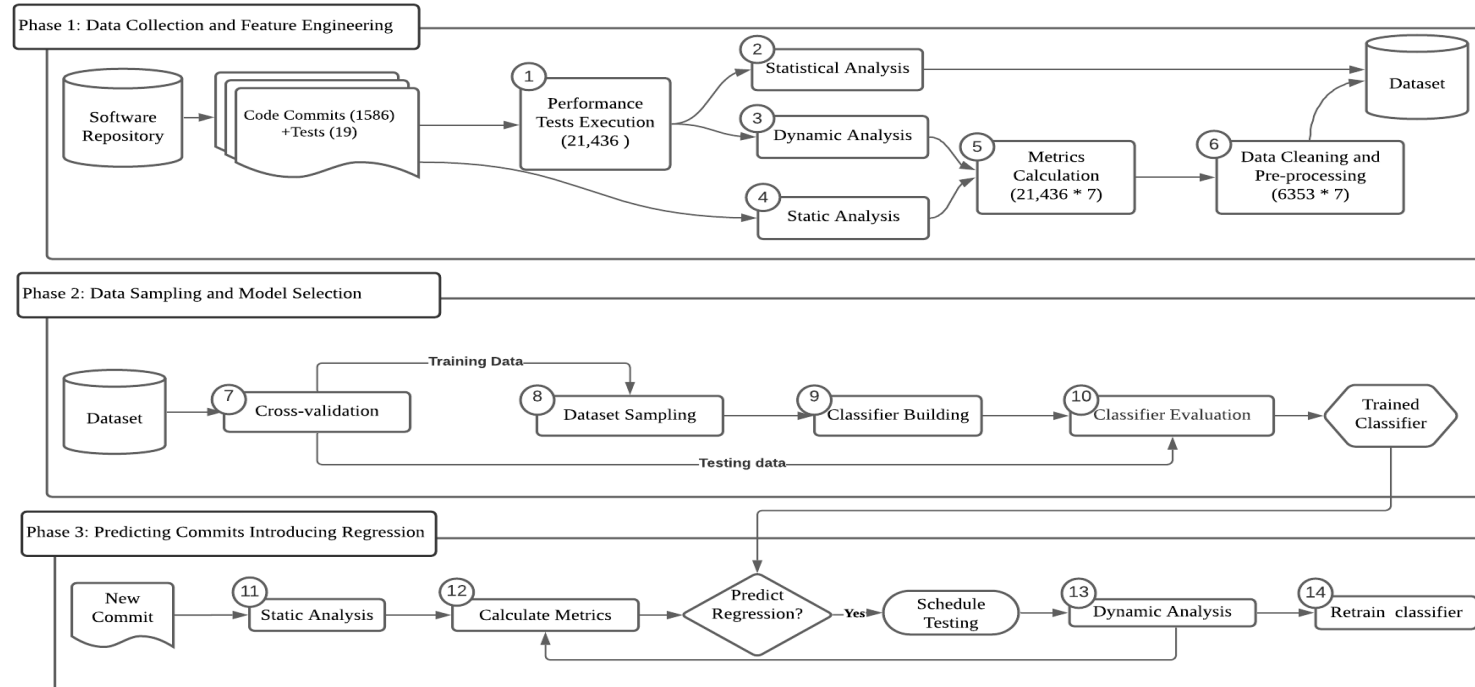
Prediction of performance regression

Empirical

- Analyze the distribution of problematic changes vs. unproblematic ones. Try to find any different patterns between them.

Recommendation

- Tackle it as a binary classification problem?



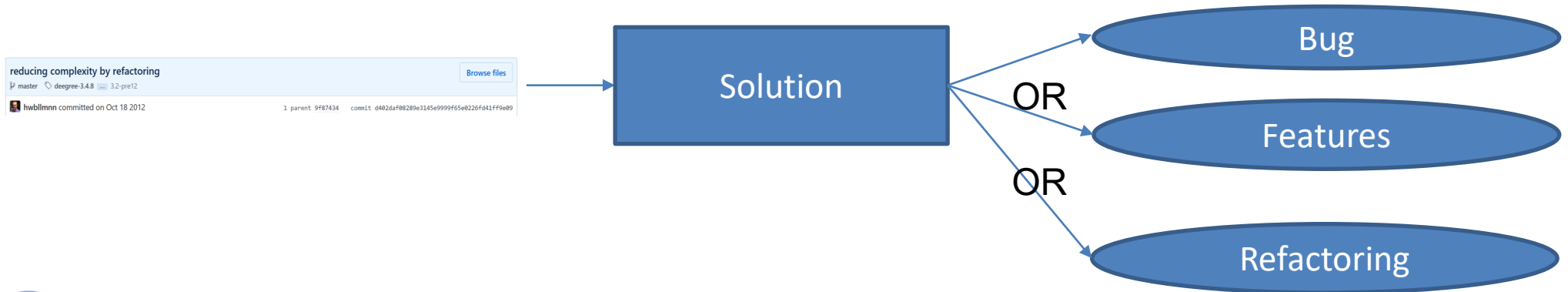
Classification of Commit Messages

Problem?

- Developers document their code changes via commit messages. Can we predict what type of code change was performed given its inline documentation

Solution

- Given a new commit message, can we automatically assign a tag?



Classification of Commit Messages

Empirical

- Similarly to issues, it is important to understand how developers document a given code change. This will help identify some textual patterns for a given class.

Recommendation

- Use appropriate solutions to solve the problem. Example:



*fast*Text

Identification of Extract Method

Problem?

- Given two (or more) methods with shared instructions (duplicate code), can we identify whether we can extract it into a separate method and call it inside the previous ones instead?

Solution

- Given a set of extract method refactorings, can we learn some patterns that would help us make a proper decision about whether the duplicate code is should be extracted or not?

```
public void SimpleExample1()
{
    using (var connection = GetConnection())
    {
        DoSomething(connection);
    }
}

public void SimpleExample2()
{
    using (var connection = GetConnection())
    {
        DoMoreStuff(connection);
    }
}
```

Solution

Extract ✓

Keep ✗

Identification of Extract Method

Empirical

- Looks at what patterns are relevant for developers to trigger an extract method refactoring (instruction complexity? Length?)

Recommendation

- Tackle it as a binary classification problem?



Code Review Performance Analysis using Bug Report

Introduction

- Nowadays, many companies are focusing more on code quality and getting less bugs on deployed code. To achieve this, a code review should be done thoroughly. This project will help to visualize the performance of code reviews by analyzing the bugs reported. When any bug is reported, we will map the bug with code review instructions, to show which code review instruction was missed during the review of the code.

Code Review Performance Analysis using Bug Report

Empirical

- Extraction of code review guidelines and their corresponding potential violations (bug reports).

Recommendation

- Tackle it as a multi-class classification problem?

