

Patterns for making sense of bug data / Patterns for cleaning and analyzing bug data

Rodrigo Souza Christina Chavez

Roberto Bittencourt
Software Engineering Labs
Department of Computer Science - IM
Universidade Federal da Bahia (UFBA), Brazil
rodrigors@dcc.ufba.br, flach@dcc.ufba.br, roberto@uefs.br

[illegible]

Index Terms—TODO; TODO; TODO; TODO

I. INTRODUCTION

Someone said that data analysis is easy, the hard part is cleaning the data. This is true also of bug databases. You deal with data that may be incomplete, inconsistent, or just plain wrong [cite Bird, Aranda etc.]. Without proper guidance, it is easy to get lost. We hope to provide some hints for those adventurous enough to analyse bug data, by suggesting some filtering, and also general analyses that support more specific analyses.

Conventions: lines beginning with > denote R code.

II. DATA

The R snippets refer to data extracted from Bugzilla databases. Other databases may have different schemas, so the scripts may need to be adapted.

Tables: `bugs_activity` (we'll call it `changes`), `bugs`, `longdescs` (we'll call it `comments`) Let's assume that we have imported such tables as a R data frame.

```
> head(events)
```

III. PATTERNS

The format we use is:

Write a short story referencing the patterns, that shows how one typically uses them

IV. PATTERN CLUSTER: GATHER

V. GRAB THE RELEASE DATES

Problem

How to discover the release dates?

Context

releases are important events. They help to determine phases in the software lifecycle. The activity in a project that is near release date may be different from the activity in regular periods. For example, in Eclipse/Platform, verifications are much more frequent just before a release.

Solution

project websites do not usually keep info about all the previous releases. Fetch previous versions of the sites using The Internet Archive, archive.org. Other solution is look for tags in the version control system.

VI. PATTERNS CLUSTER: CLEAN

VII. LOOK OUT FOR MASS EDITS

Problem

Discover changes to bug reports that were the result of a mass edit.

Context

Sometimes lots of bug reports are updated almost at the same time. These mass edits are often a result of a cleanup due to process changes, or phase change. They should not be regarded as actual effort, because they may influence statistic analysis. For example, when many bugs are marked as VERIFIED in a short amount of time, it usually does not mean that someone actually verified that the resolution for the bug was adequate, it just means that the bug database was cleaned up.

Solution

At least two solutions can be applied to this problem. The first one is more exploratory and visual, and the second one is quantitative.

(1) Compute the accumulated number of changes over time and then plot this data as a simple line chart, with time in the horizontal axis. The line is always increasing, but periods with many changes will stand out as steep ascents. Examine your vertical axis to assess whether such ascents represent a large number of changes (e.g.,

thousands). If this is the case, then it is likely that the changes were caused by a mass edit.

```
> # Sort changes by time
> #sorted.changes <- changes[order(changes$time), ]
> #cum.changes <- 1:nrow(sorted.changes)
> #plot(cum.changes ~ sorted.changes$time)
```

(2) Choose a type of change you are interested

Construct a contingency table to count how many changes

(3) Perform a hash of the comments (this is a optimization to avoid working with the full comments). Compute a contingency table to count, for each hash, the number of occurrences per day, per developer. You can sort the hashes by

(please note that there is a chance of hash collision, but it is unlikely, so such event, if it occurs, is probably sporadic and will not affect the analyses)

```
> #library(digest)
> #comments$hash <- sapply(comments$text, digest) # or select md5() --SQL
> #comments$date <- as.Date(comments$time)
> #
> #tab <- count(comments, c("hash", "date"))
> #tab <- tab[order(tab$freq, decreasing=TRUE), ]
> #head(tab)
```

Discussion

Examples

[2]

Related Patterns

Be sure to Choose a SuitableThreshold. It is also a good idea to ReadTheFineMessages

VIII. MEET THE TEAMS / KNOW YOUR SUBJECTS

Problem

How to discover developers' roles?

Context

Some developers are triagers, other are "fixers", there are also QA experts. When studying human factors, it is important to know that developers are not all equal and have different roles.

Solution

Analyze their activity in the bug tracking system, specially the status changes. Count the relative number of activities. Plot distribution of ratios. After discovering a group of developers with a specific role, see if the group is responsible for most of that activity in the project.

Related Patterns

Be sure to Choose a Suitable Threshold.

IX. CHOOSE A SUITABLE THRESHOLD

X. OLD WINE TASTES BETTER

(aka Let them mature / give 'em some time)

Problem

how to avoid analysing bugs that are not accurate

Context

it may be the case that a bug is later discovered to be invalid, or a fix is reopened, but you may not know this

and consider the bug as valid or fixed, because you do not have enough history.

Solution

] Discard newer bugs, that do not have enough history. Plot the distribution of bug durations (from open to last activity, or from open to a specific activity that you are interested as an outcome variable). (in my case, I discovered that more than 90% are reopened within 1 year, so it is probably safe to discard bugs that are newer than 1 year).

Related Patterns

Be sure to Choose a Suitable Threshold.

XI. DON'T BE CONFUSED

Problem

confusion variables affect outcomes.

Context

or select md5() --SQL

Solution

use mosaic plots. use correlation matrices. Consequences: mosaic plots are easy to read with two variables, barely readable with 3 or 4, but start to become confusing with 5 or more.

XII. COUNT 'EM ALL

Problem

How to know if your variables are associated?

Solution

count the cases. Built a contingency table. After that you can use Fisher's exact test to assess association.

Related Patterns

[Classify] your data (make sure your [data is classified])

XIII. PATTERN CLUSTER: ANALYZE

XIV. READ THE FINE MESSAGES

(aka Read the Fine Comments)

Problem

How to know if a specific status change really represents what you think it represents?

Context

a specific status change may have different meanings on different projects. For example, when a bug is marked as FIXED, does it mean that a patch was submitted? That the fix is available in the version control system?

Solution

don't be afraid and read the comments of a few cases in which that status change occurs. Specially comments of bugs that are exceptional in one way or another.

Related Patterns

it is more productive to read comments that are related. You may want to [Classify] your data before .

XV. KNOW YOUR PROJECT

Problem

Context

when you do not know your project, you run the risk of making wrong assumptions about the data.

Solution

read developer documentation on the website. You may find guidelines, well defined roles and processes. Seek material in developer conferences.

Examples

eclipsecon, the talk “eclipse way”. NetBeans docs talk about quality engineers.

XVI. KNOW THEIR ALIASES/NAMES

Problem

how to discover the multiple identities of a developer?

Solution

use the methods proposed by Bird et al. in Mining Email Social

XVII. EAVESDROPPING

(aka Be a regular reader/listener)

Problem

how to detect what is not explicit in the bug metadata?

Solution

do a text search. Maybe use regular patterns.

Examples

verification techniques reported by developers.

Related Patterns

read the fine comments.

XVIII. THIS CASE IS CLASSIFIED

(aka Classify your Cases / Make your cases classified / Separate the wheat from the chaff)

Problem

How to avoid analysing cases that you are not interested in?

Context

Solution

create scripts to automate the classification and enrich your data with derived observations.

Related Patterns

after classifying the cases, [Read the Fine Comments] of each class and verify if the classification is adequate

XIX. TRIANGULATE / CHECK

Check your inferences by reading a sample of the data.

XX. EXPLORE

Read some data to get the idiosyncrasies of each project. See Read the Fine Comments, Know Your Project, Eavesdropping.

REFERENCES

- [1] C. Bird, A. Bachmann, E. Aune, J. Duffy, A. Bernstein, V. Filkov, and P. Devanbu, “Fair and balanced?: bias in bug-fix datasets,” in *European Soft. Eng. Conf. and Symposium on the Foundations of Soft. Eng.*, ser. ESEC/FSE '09. ACM, 2009.
- [2] R. Souza and C. Chavez, “Characterizing verification of bug fixes in two open source ides.” in *MSR*, M. Lanza, M. D. Pent, and T. Xi, Eds. IEEE, 2012, pp. 70–73. [Online]. Available: <http://dblp.uni-trier.de/db/conf/msr/msr2012.html#SouzaC12>