

Classificação de *issues* do Github relacionadas a Segurança

Aprendizagem de Máquina

Bruno Gonçalves de Oliveira

`bruno.mphx2@gmail.com`

Diogo Cezar Teixeira Batista

`diogocezar@utfpr.br`

Universidade Federal do Paraná - UFPR

Curitiba - 2020

Agenda

- 1 Introdução
- 2 Obtenção dos Dados
- 3 Pré-processamento
- 4 Extração de Características
- 5 Resultados
- 6 Conclusão

- Gerenciamento e manutenção de arquivos: desafio;
 - backups não realizados;
 - sobrescrita de arquivos;
 - difícil manutenabilidade em times;
- Diferentes soluções no mercado: *CVS*, *Subversion*, *TFS*, *Mercurial*.

- *Git + GitHub = OpenSource.*
- *GitHub* que é uma plataforma para versionamento, gerenciamento e colaboração de projetos, que utiliza o *Git* como base. [CS20]

- Dentre outras ferramentas temos o controle de *Issues*:
 - documentar possíveis *bugs*, melhorias, ou novas *features* para os projetos

Exemplo Issues

Filters ▾

🔍 is:issue is:open

🏷 Labels 33

📅 Milestones 0

New issue

🔔 332 Open ✓ 8,850 Closed

Author ▾	Label ▾	Projects ▾	Milestones ▾	Assignee ▾	Sort ▾
🔔 Add feature to catch errors in a production build for compatibility with external error tracking	feature request				💬 2
#11666 opened 5 days ago by ssmulders					
🔔 transition-group with flex parent causes removed items to fly	transition				💬 1
#11654 opened 12 days ago by turbosheep44					
🔔 slot is reused with v-if/v-else	bug	has workaround			💬 5
#11652 opened 12 days ago by jiankafei					
🔔 `vnode.data.on` is empty for custom components					
#11635 opened 19 days ago by posva					
🔔 TemplateRender.renderScripts breaks when there are no preloaded files (empty array)			🔗 1		💬 1
#11612 opened 26 days ago by john-ko					

Figura: Exemplos de Issues do Projeto Vue.js

- Eventualmente, as *issues* podem estar relacionadas a tópicos de segurança.
- Quando consideradas críticas, podem ser analisadas por outros especialistas;
- Como identificar quais *issues* que são relacionadas com segurança?
- Como classificar estas *issues* para que especialistas possam analisar os códigos?

- Criação de uma ferramenta que utilize técnicas de **aprendizagem de máquina** para o desenvolvimento de um classificador que consiga analisar as palavras contidas nas mensagens das *issues* de um dado projeto, e classificar se esta *issue* está ou não relacionada no contexto de segurança da informação.

Obtenção dos Dados

- Utilizou-se o *github-csv-tools*¹ que possibilita a exportação dos dados de um repositório do *GitHub*, salvando as informações em um arquivo no formato CSV.
- Dados tratados para um CSV com 2 colunas:

```
security,PushObserver can be used to push ↵  
serverinitiated HTTP/2 requests into an ↵  
OkResponseCache...  
not,Handle LOCKED in conversions.Motivation...
```

Código 1: CSV Exemplo com Base de Dados

¹<https://github.com/gavinr/github-csv-tools>

- Base de Testes: *issues* do projeto *Wildfly*²;
- Base de Treinamento: *issues* dos projetos: *okhttp*³, *jgit*⁴ e *couchbase*⁵
- Os dados de treinamento possuem **199** entradas, enquanto que para a base de teste foram utilizadas **211** entradas.

²<https://github.com/wildfly/wildfly>

³<https://github.com/square/okhttp>

⁴<https://github.com/eclipse/jgit>

⁵<https://github.com/couchbase>

- Completar...

- Completar...

- Completar...

- Completar...

- Completar...



Scott Chacon and Ben Straub, *Pro git*, Apress, 2020.