# Automatic identification of bot accounts in *Open-Source projects*

Miguel Ángel Fernández Sánchez

MSc on Data Science
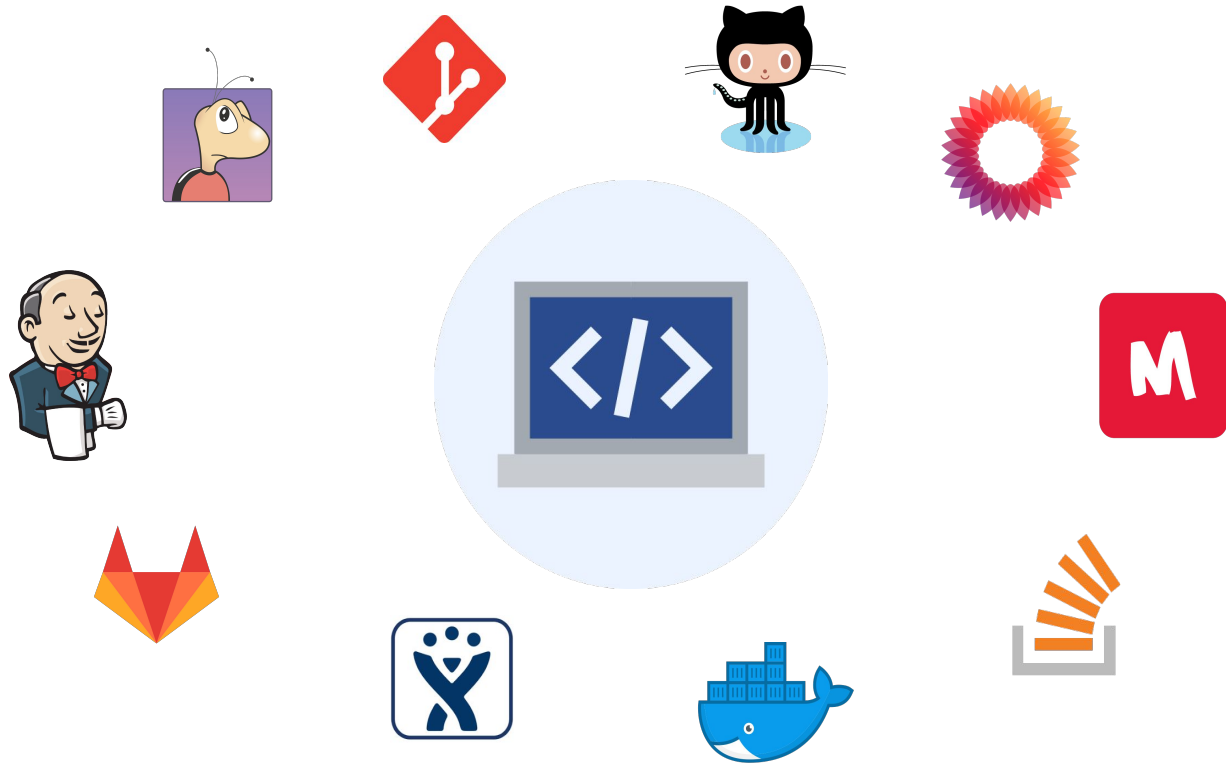
Tutor: José Felipe Ortega Soto, PhD.

April 20th, 2023

Universidad Rey Juan Carlos

# /motivation

# /motivation

**Accounts**

| | |
|---|---|
| git | jane.doe@urjc.es |
| | janedoe@libresoft.com |
| GitHub | janedoe |
| stack**overflow** | doe.jane |
| Mattermost | janedoe |

**Organizations**

Universidad
Rey Juan Carlos

**From:** September 1st, 2015
**To:** January 31st, 2017

GSyC
LibreSoft
we study libre software

**From:** February 1st, 2017
**To:** Present day

3

# /motivation

Merge these profiles

👤 Identity A

👤 Identity B

Add this organization

Universidad Rey Juan Carlos

Mark this profile as **bot**

Name: autobot
Bot: Yes

# /problem

Commits from humans

Commits from bots

Comparison: Number of commits over time

* Number of commits from January, 2008 to September, 2021

● By Bots  ● By Humans

# /problem

# /proposal

# /choosing the community to analyse

# /details from a Git commit



**[schema] Support filtering individuals by last updated date**

```
The accepted formats are controlled by regular expressions
matching two patterns:
* A comparison operator (>, >=, <, <=) and a date
(e.g. `>=YYYY-MM-DDTHH:MM:SSZ`).
* A range operator (..) between two dates
(e.g. `YYYY-MM-DDTHH:MM:SSZ..YYYY-MM-DDTHH:MM:SSZ`)

The accepted date format is ISO 8601, YYYY-MM-DDTHH:MM:SSZ,
also accepting microseconds and time zone offset
(YYYY-MM-DDTHH:MM:SS.ms+HH:HH).

Signed-off-by: Miguel Ángel Fernández <████████@████████.com>
```

Browse files

master (#401)

0.8.1 ... 0.8.0-rc.1

**mafesan** committed on Dec 3, 2020        1 parent d0af9c6 | commit 66c3b16

Showing **2 changed files** with **372 additions** and **2 deletions**.        Split | Unified

> 84 ▮▮▮▮▮ sortinghat/core/schema.py

> 290 ▮▮▮▮▯ tests/test_schema.py

**Author**

**Modified files**

**Added lines**

**Removed lines**

**Commit message**

**Commit identifier**

**Commit date**

# /choosing fields using *GQM methodology*

**Author data**

- author_uuid
- author_name
- author_bot

**Author data, from commit**

- commit_name
- author_date

**Commit data**

- hash
- message
- files
- lines_added
- lines_removed
- utc_commit
- repo_name

- - - - - - - - - - - - - - - - - - - - - - - - - - - - -

- commit_date_weekday
- time_to_commit_hours

# /experiments phase

Pre-processing of the data

Clasification models

**05**

**01**

**02**

**03**

**04**

Validation of results

Choosing the classifier

Hyperparameter adjustment

# /exploratory data analysis

| Commits | | Commits | | Humans |
|---------|---|---------|---|--------|

Commits

1.916.010

Commits

1.895.701

Humans

3.717

Filter

Authors

16.284

Authors

3.747

Commits ≥ 10

Bots

30

Repositories

2.876

# /input dataset

**User 1 - Commits**
- Hash 1.1
- Hash 1.2
- Hash 1.M

**File 1**

**User 2 - Commits**
- Hash 2.1
- Hash 2.2
- Hash 2.M

**File 2**

...        ...

**User 1 - Commits**
- Hash N.1
- Hash N.2
- Hash N.M

**File N**

```
{
    "git__hash": "601dbd92f4bfb99e1af8ce4dda0...",
    "git__lines_added": 28,
    "git__lines_removed": 15,
    "git__files": 3,
    "git__utc_commit": 1328767481000,
    "git__grimoire_creation_date": 1328751851000,
    "git__commit_date_weekday": 3,
    "git__commit_name": "Jane Doe",
    "git__message": "Initial commit",
    "git__time_to_commit_hours": 19.65999984741211,
    "git__repo_name": "https://github.com/example-repo.git",
    "author_uuid": "001c1fc408eccb58365b8...",
    "author_bot": false,
    "author_name": "Jane Doe",
    "author_date": 1328751851000
},
{

    ...

},
{

    ...

}
```

**Author 1**
- Metric 1.1
- Metric 1.2
- Metric 1.M

**Author 1**

**Author 2**
- Metric 2.1
- Metric 2.2
- Metric 2.M

**Author 2**

...        ...

**Author N**
- Metric N.1
- Metric N.2
- Metric N.M

**Author N**

# /pre-processing of the data and feature transformation

# /feature generation

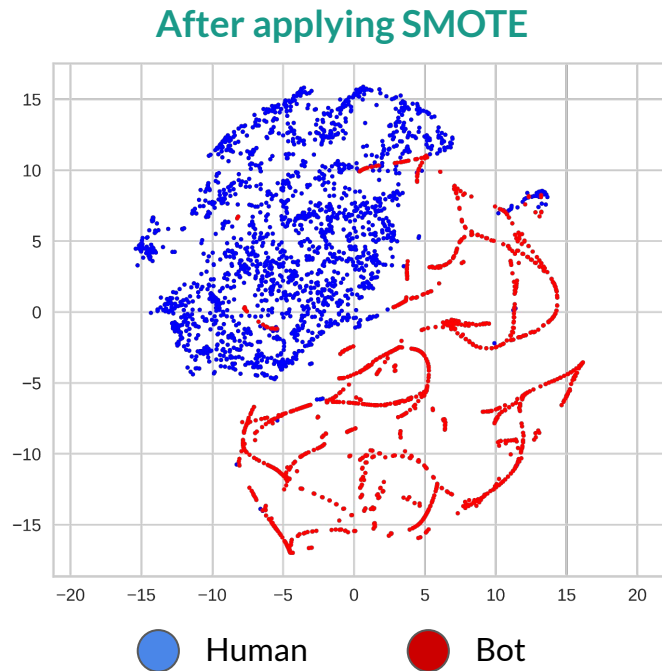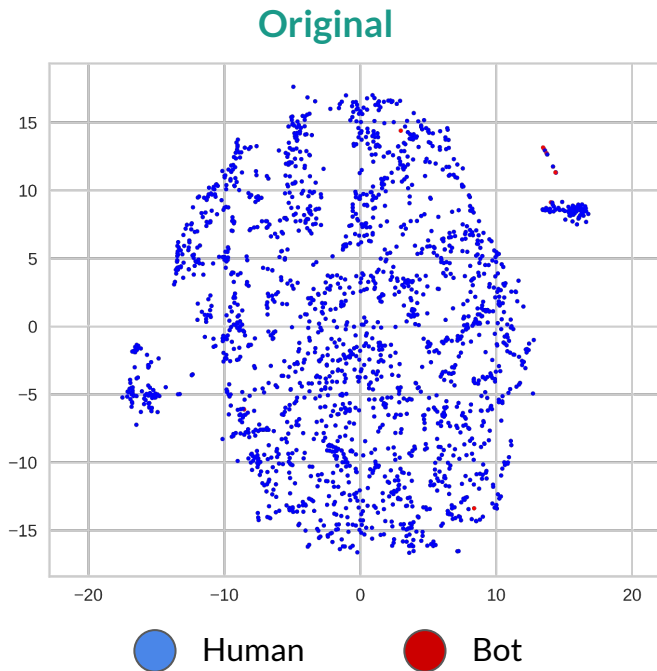| Level | Weight | Heuristic terms |
|---|---|---|
| 1 | 60 | bot, dependency, fix, integration, merge |
| 2 | 30 | auto, build, commit, copy, issue, release, request, review, sync, template, tool, travis |
| 3 | 10 | cd, ci, code, patrol, pr, pull |

$$\mathrm{Ts(term)} = 60\mathrm{Nl}_1 + 30\mathrm{Nl}_2 + 10\mathrm{Nl}_3$$

# /unbalanced data: SMOTE

# /class distribution using t-SNE



**Original**

**After applying SMOTE**

# /classifiers: evaluation

| | Predicted | |
|---|---|---|
| **Real** | TN | FP |
| | FN | TP |

- *Precision*
- *Recall*
- *Fβ-score*

| Classification model |
|---|
| Gaussian Naive-Bayes |
| Complement Naive-Bayes |
| LinearSVC |
| KNN |
| Decision Tree |
| Random Forest |
| XGBoost |

# /classifiers: Random Forest

| | |
|---|---|
| **Number of estimators** | 300 |
| **Split criterion** | *Gini* impurity |
| **Maximum depth** | 4 levels |



Bots       Humans

# /classifiers: test

| Predicted | | |
|---|---|---|
| | Human | Bot |
| **Real** Human | 826 | 3 |
| Bot | 1 | 6 |

| Model name | Precision | Recall | Fβ-score |
|---|---|---|---|
| Random Forest | 0.667 | 0.857 | 0.811 |

**+ 5-fold cross-validation**

# /classifiers: validation

| Real | | Predicted | |
|---|---|---|---|
| | | Human | Bot |
| | Human | 493 | 6 |
| | Bot | 1 | 3 |

| Model name | Precision | Recall | Fβ-score |
|---|---|---|---|
| Random Forest | 0.333 | 0.75 | 0.6 |

# /classifiers: most relevant features



Feature Importances of 12 Features using RandomForestClassifier

- terms_score
- git_log_iqr_len_words_commit_message
- git_log_median files
- git_log_num_commits
- git_log_num_weekend_commits
- git_log_num_signed_commits
- git_log_iqr_files
- git_log_median_lines_removed
- git_log_num_repos
- git_log_median_len_commit_message
- git_log_num_merge_commits
- git_log_median_lines_added

Relative importance

# /future work

## Improvement of the classification model

- Other projects / Other communities
- Text metrics / Digital footprints
- Summary of the user / *Concept Learning*
- Mixed accounts / Multi-class classification

## Integration with SortingHat/GrimoireLab

- Classification report
- Recommendation engine
- Feedback to the classification model after the user's decision

/end

Questions

**[mafesan.github.io/Memoria-TFM](mafesan.github.io/Memoria-TFM)**