# GitHub User Analysis

Xin Guan | Ziqian Ge

## About GitHub

- a web-based hosting service of version control using Git. (Wikipedia)

- The world's **largest** community of developers to discover.

- 31 million users (October 2018)

- 100 million repositories (November 2018)

# Goal of Project

- Who are those developers

- Trends among developers

## Content of the Project

- Data Collection
- Data Cleaning
- User Analysis
- Repository Analysis
- Prediction

## Data Collection

- GitHub API
- Third party data from GHTorrent

```
[1] "user id: 206, user name: Jeff Smick"
[1] "user id: 217, user name: Tim Kersey"
[1] "user id: 228, user name: Alex Vollmer"
[1] "user id: 249, user name: Mike Vincent"
[1] "user id: 269, user name: Kamal Fariz Mahyuddin"
[1] "user id: 270, user name: Ben Bleything"
[1] "user id: 278, user name: Hemant Kumar"
[1] "user id: 279, user name: Patrick Ewing"
[1] "user id: 307, user name: Norio Shimizu"
[1] "user id: 325, user name: Edward Ocampo-Gooding"
[1] "user id: 347, user name: Tobias Lütke"
[1] "user id: 348, user name: James Tucker"
[1] "user id: 253, user name: Chris Anderson"
[1] "user id: 1017, user name: Jiang Jiang"
[1] "user id: 2621, user name: Brad Fitzpatrick"
[1] "user id: 3499, user name: Tatsuhiko Miyagawa"
[1] "user id: 4970, user name: Kang-min Liu"
[1] "user id: 5526, user name: Marcus Ramberg"
[1] "user id: 6545, user name: Lu Yibin"
[1] "user id: 8465, user name: Hironao OTSUBO"
[1] "user id: 11427, user name: Xin Liu"
[1] "user id: 14242, user name: Yuval Kogman"
[1] "user id: 14658, user name: Reeze Xia"
[1] "user id: 17814, user name: icyleaf"
character(0)
[1] "user id: 20723, user name: 唐鳳"
[1] "user id: 21084, user name: Tokuhiro Matsuno"
[1] "user id: 22623, user name: Victor Igumnov"
```

## Difficulty

- Large amount of data
- restricted access to API (5000 per hour)
- Not in the form of relational data base
- Uncleaned Data

## Data Cleaning

- Ignore users with 0 repositories, 0 followers and 0 following

- Ignore users without self-description (no location, company or bio)

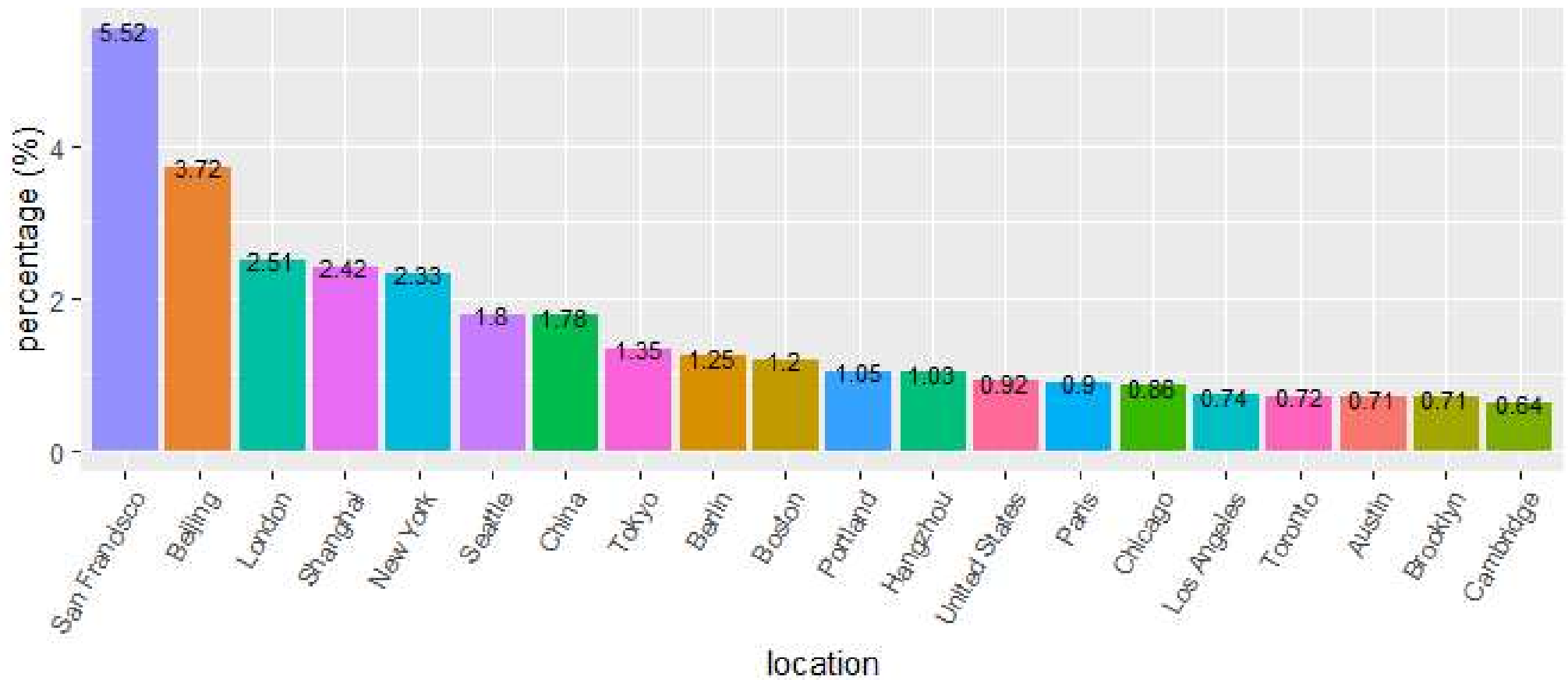- Ignore repos with 0 stars, 0 forks, 0 open issues, or 0 subscibers

- Remove NA

- Outliers

**User Analysis**

- Who?
  1. Extract users with introduction
  2. Delete meaningless words
  3. Count the frequency of words
  4. Sort and Plot

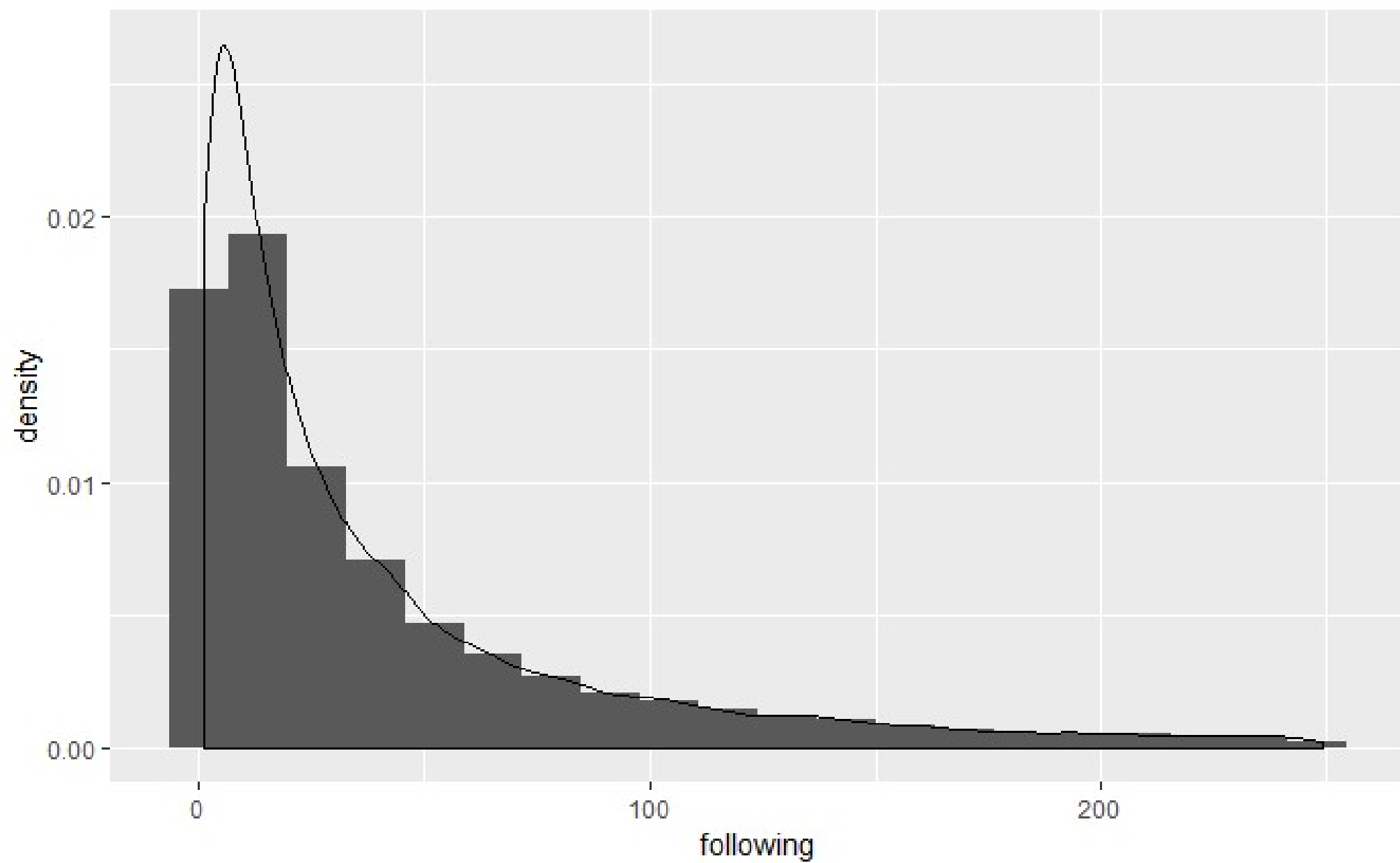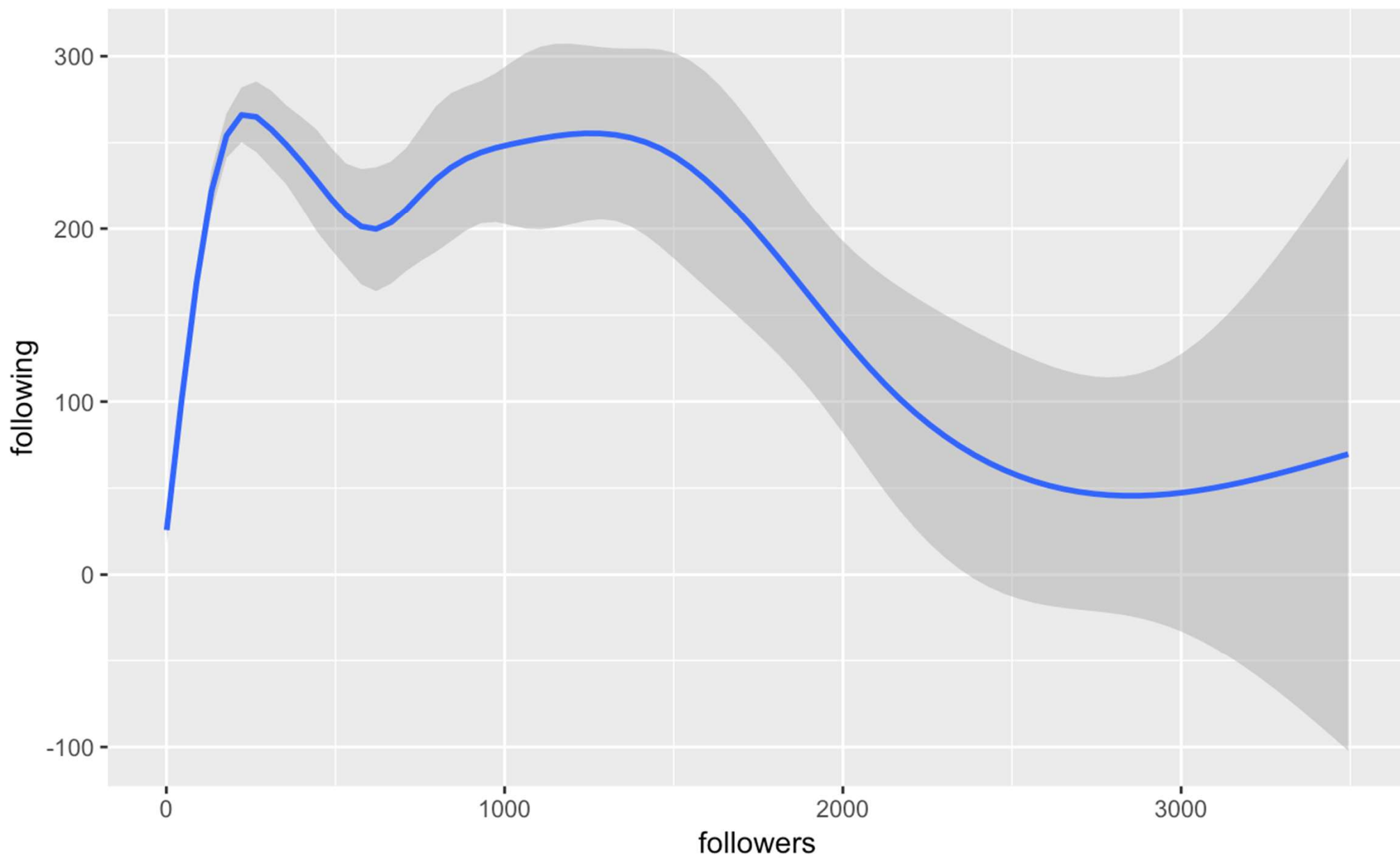# User Analysis

- **Followers and Following**
  - **Clean Up data – NA / 0 cases**
  - **Outliers (93 % < 250)**
  - **Plot**

## User Analysis

- **Not Finished Task:**
  - Identify their company / institution / job
  - Relate users with their repositories

## Repo Analysis

- **Languages**
  - 1. Count repos that has changed in a day by language
  - 2. Compare popular language with history
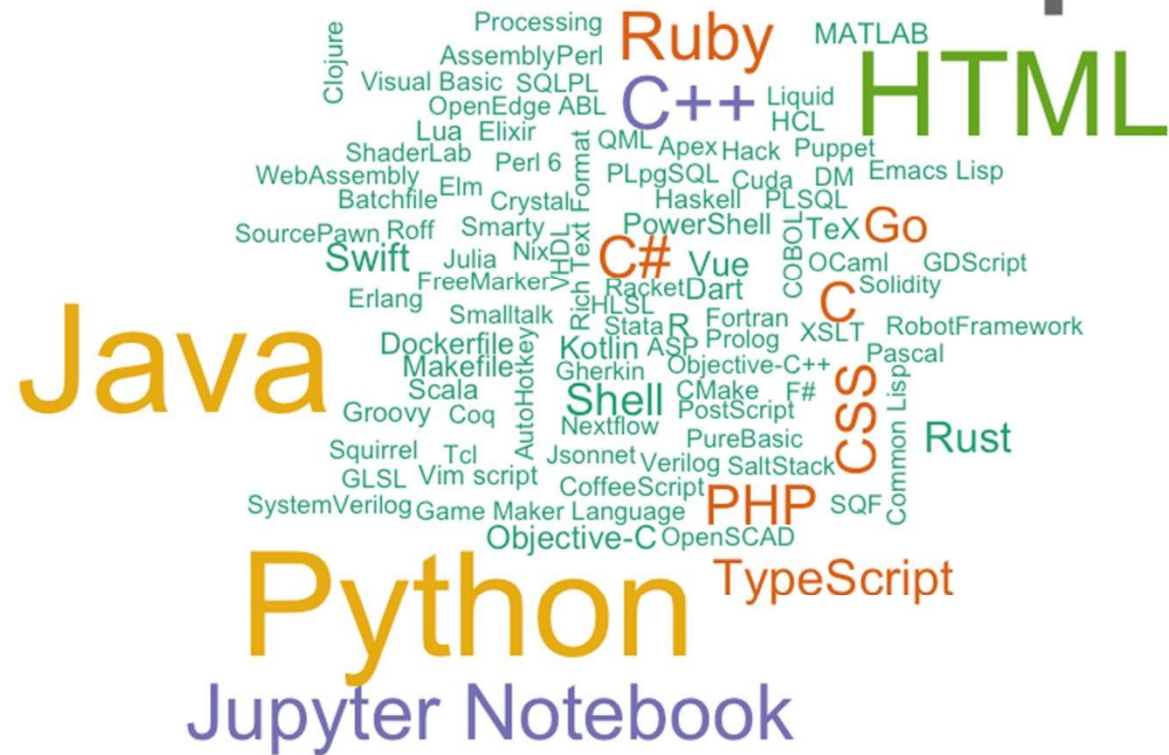- **Popularity Analysis**
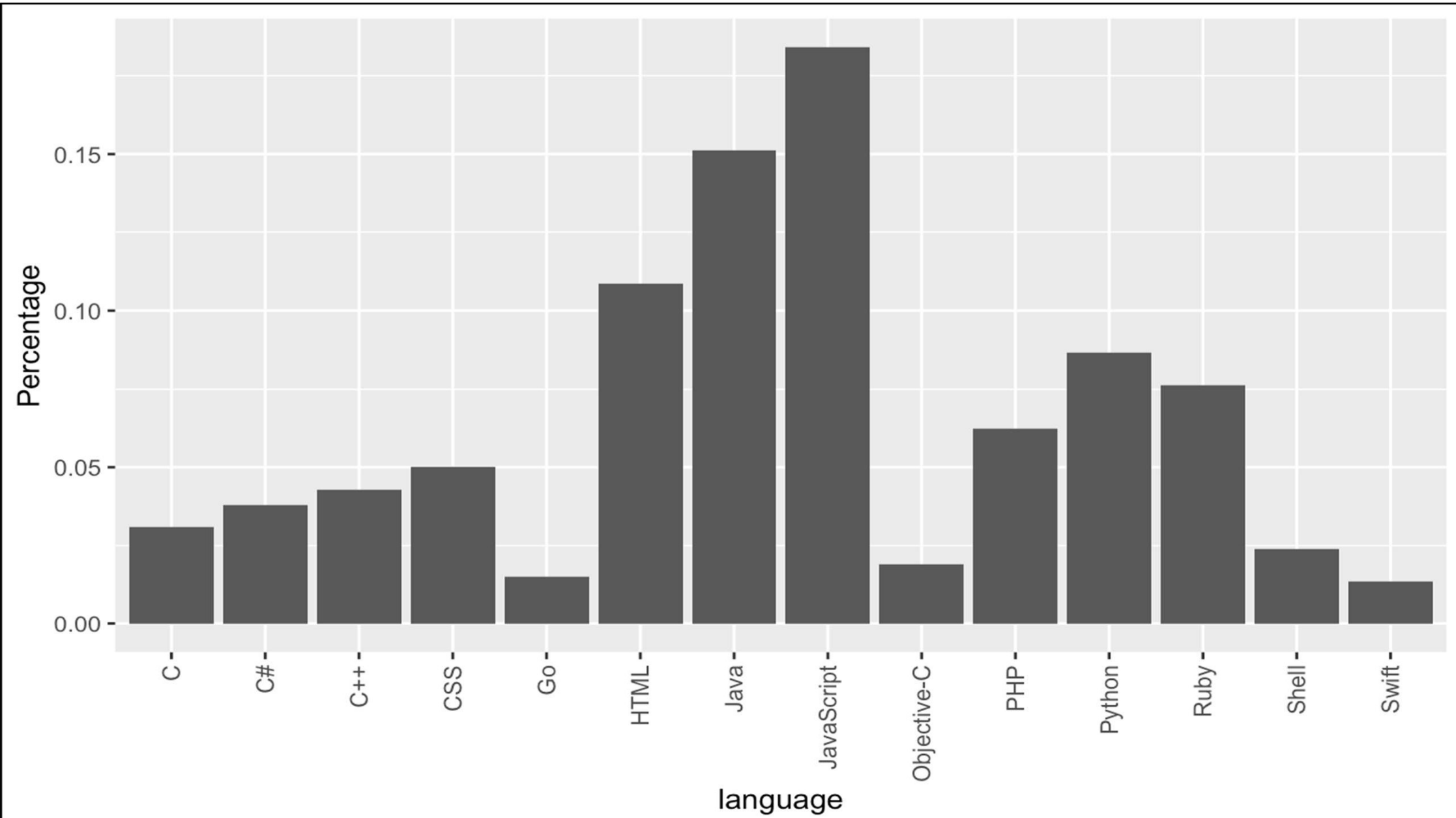  - Open issues, forks, stars, subscribers

## Repo Analysis

- **To be finished**
  - Relate with users, i.e. repo owner and contributors
  - Relate with commits

## Prediction

- Based on properties of the owner of a repository, and his/her other repositories, predict the popularity of the repo.
  - Evaluate the popularity by a popularity index
  - Build a model on properties of a user
  - Predict using regression