



Reddit posts - A Data Analysis

Joshua Sung, Data Scientist



Guiding questions

- What characteristics of a post on Reddit are most predictive of the overall interaction on a thread (as measured by number of comments)?
- Of those, which of them predict a post to be above or below the median number of comments the best?

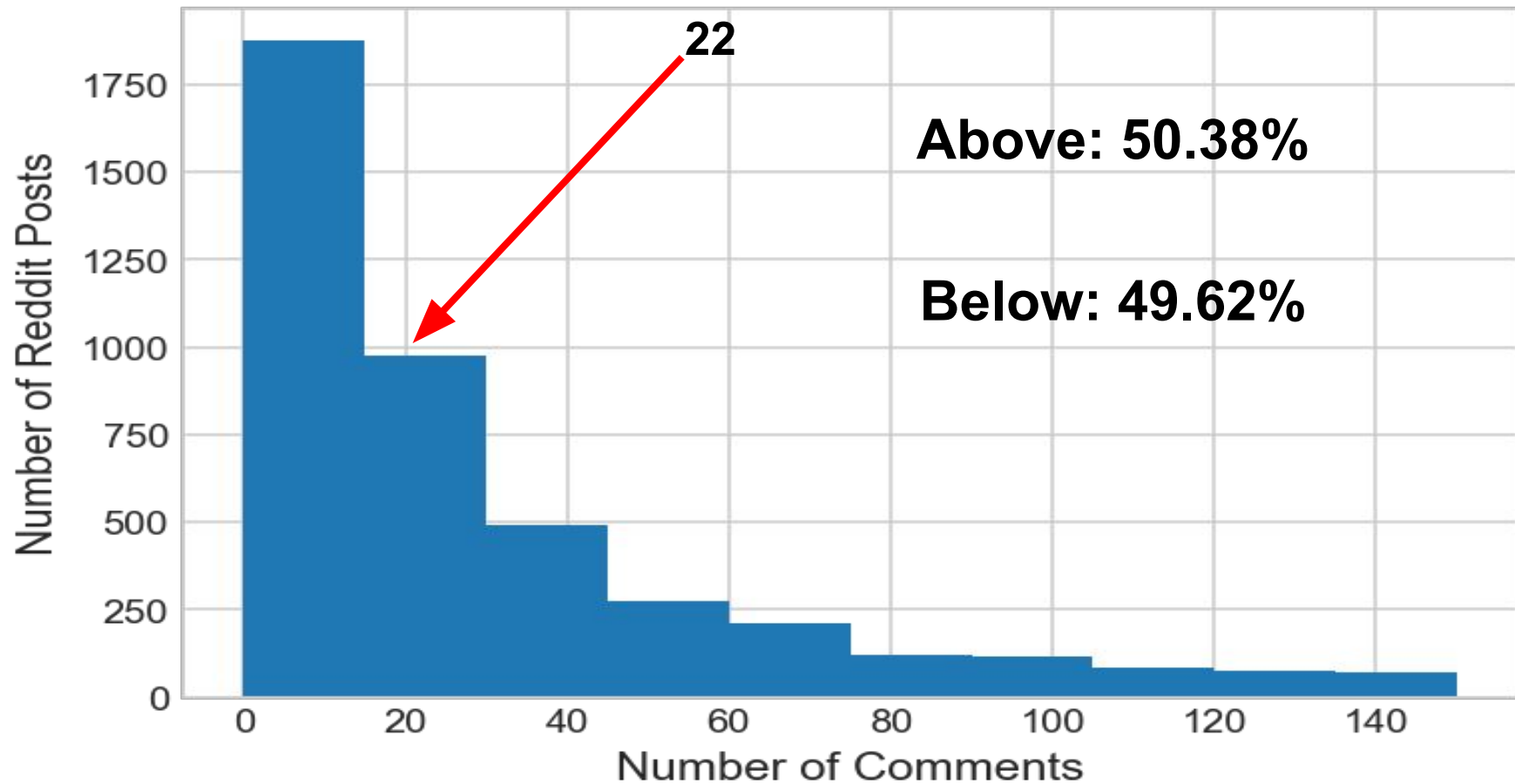
What's in a Reddit post?

- Title
- Source (i.e. URL)
- Author
- Upvotes
- Subreddit
- Time posted

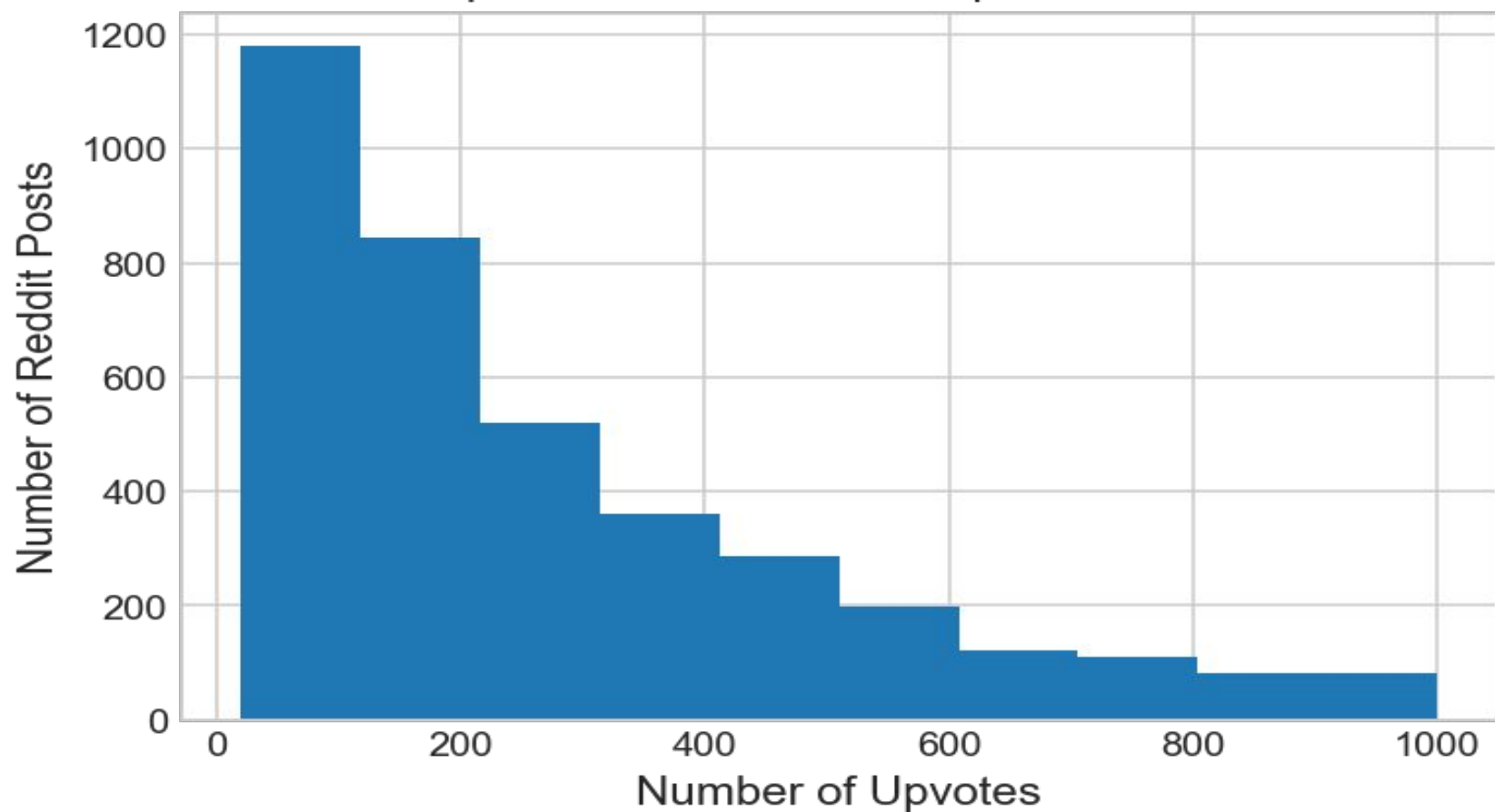
What did I collect?

- Data scraped from Reddit.com
 - Titles
 - Number of comments
 - Time posted
 - Number of upvotes
 - Subreddit

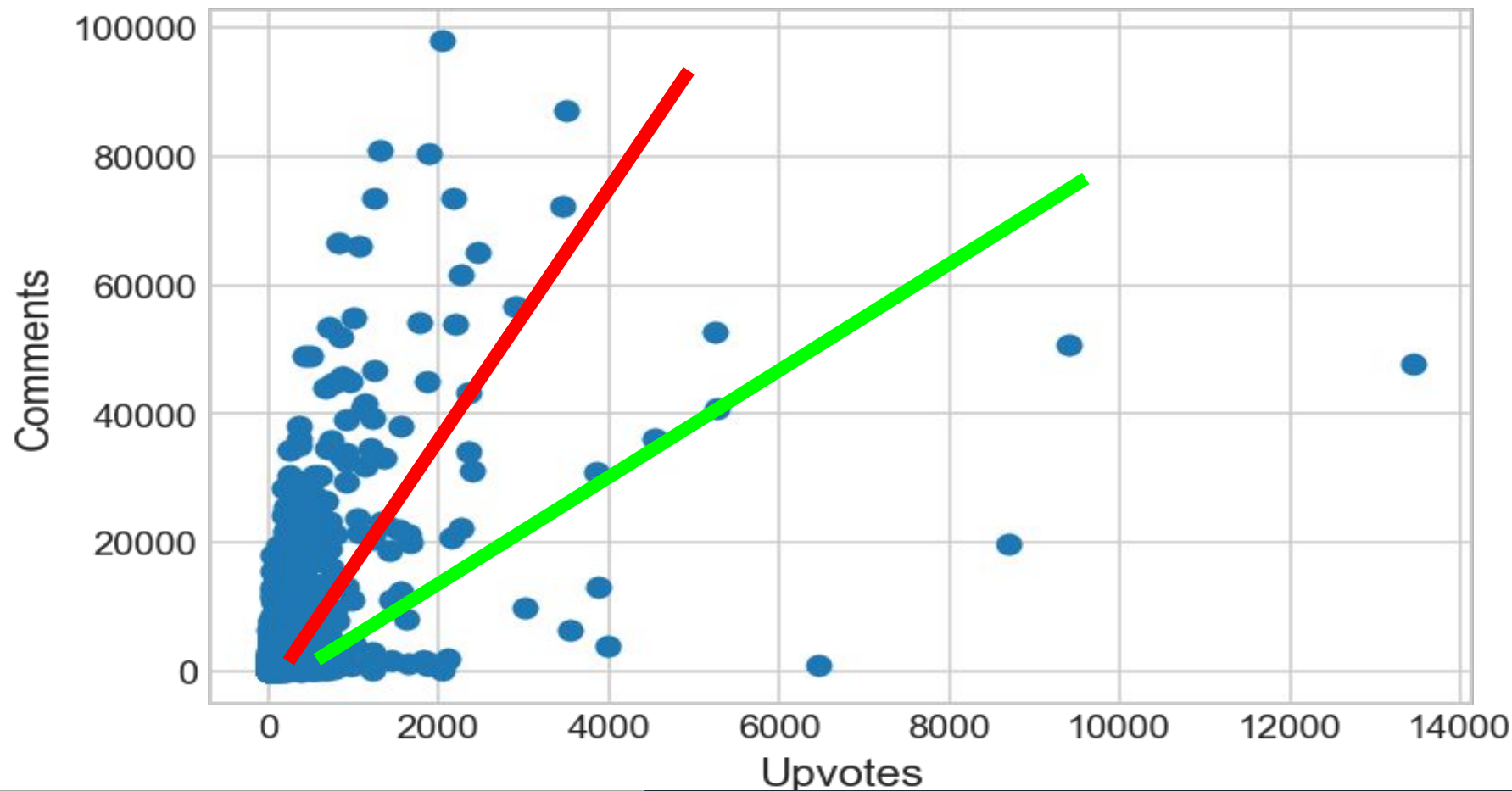
About half are above and below the median number of comments!



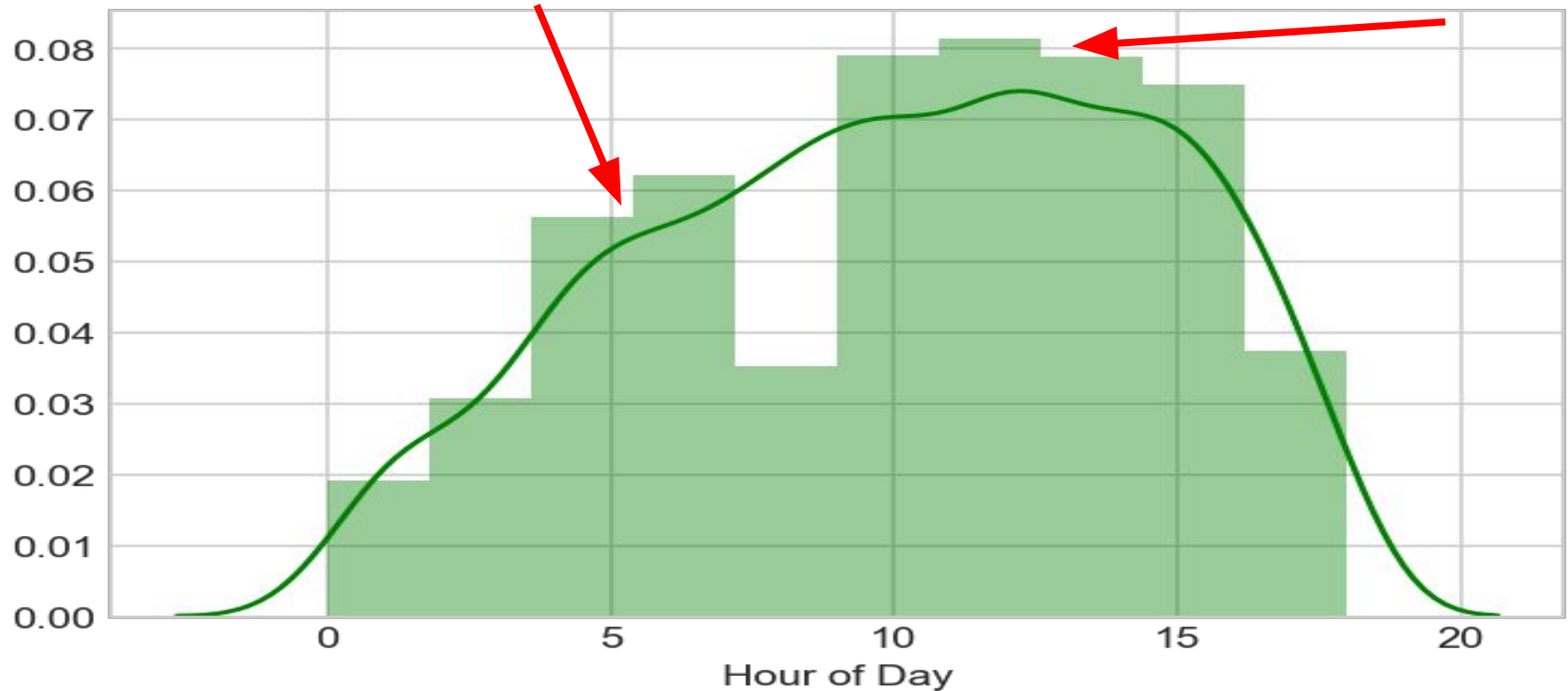
The distribution of upvotes has a similar shape to the comments distribution



Upvotes and Comments have a positive relationship!



Most people post on Reddit during the work hours and in the early morning from 4am to 7am.



Findings

- Posting on Reddit during 'work day hours' and 'morning hours' along with the number of upvotes → 70% accuracy
- Subreddits only and words in a title not so much → 51% accuracy

Conclusion

- Collected data at a specific point in time.
 - Future data collection: over a longer period of time, different times of day, scrape less and more frequent
- Distinguishing between posts based on author, the source?
- Post during the work day, early in the morning