

TECHNISCHE UNIVERSITÄT BERLIN

AIM-3: SCALABLE DATA ANALYSIS & DATA MINING

---

# **Unsupervised Clustering of Hacker News Stories**

---

*Authors:*

Bram LEENDERS & Marc ROMEYN

bcleenders@gmail.com & marc.romeyn@gmail.com

July 14, 2015

# 1 INTRODUCTION

## 1.1 WHAT IS HACKERNEWS

Hacker News is a social news site: it aggregates news by allowing users to submit stories. Interesting submissions can be upvoted by other users and all submissions are ranked by popularity.

The content on Hacker News is mostly related to science, in particular computer science. The guidelines for what content can be posted are very broad; the guidelines specify on topic as “*anything that gratifies one’s intellectual curiosity*”<sup>1</sup>.

Over the last eight years Hacker News has experienced rapid growth, resulting in a daily 2.6 million pageviews and 3.5 unique visitors per month<sup>2</sup>. One of the reasons suggested for this popularity is Hacker News’ similarity to how Reddis used to be<sup>3</sup>: user-submitted content with a very minimalistic, terminal-like interface.

## 1.2 RESEARCH QUESTION

The world changed a lot over the last eight years, especially in the field of computer science. The userbase of Hacker News and their interests have changed as well. In this research, we want to see how they changed; what new . Informally phrased, we want to find trends in the popularity of several topics over the last eight years. This lead to the following research question:

RESEARCH QUESTION: how did the popularity of topics on Hacker News change over time?

This question depends on two other questions, since we have not yet specified what we mean by popular nor what topics we mean exactly. These subquestions are:

SUBQUESTION 1: how does one quantify the popularity of a topic?

SUBQUESTION 2: what topics does the Hacker News content consist of?

## 1.3 OUTLINE

We will address the first subquestion by ranking posts and users over the entire Hacker News history. This will help characterize the dataset and show how various ways of measurent are often very similar and reveal some flaws in several ways of measuring popularity.

Dividing the content in topics is a more academically challenging task. We have used two methods to divide the articles into categories: latent Dirichlet allocation and a combination of Word2Vec and k-means. Both methods are briefly explained in section 4.

---

<sup>1</sup><https://news.ycombinator.com/newsguidelines.html>

<sup>2</sup><https://news.ycombinator.com/item?id=9219581>

<sup>3</sup><http://techcrunch.com/2013/05/18/the-evolution-of-hacker-news/>

## 2 DATASET

As explained in the previous section, our analysis is on the set of all stories posted on Hacker News. It is a large set of news articles, blog posts, essays, tutorials and other types of textual media. Nearly all of it is English, although there are a few other languages used as well.

Let us first describe what exactly we refer to when we use the word story. To do this, we have to describe the various types of content on Hacker News, which can be categorized in these four categories:

- Stories: the majority of submissions are stories. A story can be either a link to another webpage or a relatively short text by the submitter.
- Jobs: companies sponsored by YCombinator (the seed investor behind Hacker News) can post job offers on Hacker News. The percentage of Jobs is very small: well under one percent of the total volume.
- Polls: users can submit multiple choice questions for other users to answer.
- Comments: the three types mentioned above can receive comments by other users.

Since this research focuses only on the stories, we have left out the other three types. The dataset used in our research contains all stories between February 19th 2007 (the date Hacker News was launched<sup>4</sup>) and June 10th 2015 (the day we ran our crawler). This is a time span of 3033 days, during which a total of over 1.5 million stories were submitted.

**+ 10GB comments for dictionary training data!**

### 2.1 DATA RETRIEVAL

To crawl all stories, we used the official Hacker News API<sup>5</sup>. This API returns some basic data about the story, such as the submitter, title, points (upvotes), a (possibly empty) story text and a (possibly empty) url. Stories generally either have a story text or a url, most only have a url.

For the stories with a non-empty story text field, fetching the story from the API is enough. For other stories, we also fetched the content the URL links to. An important remark here is that some urls (especially the old ones) have become invalid over the course of time.

To crawl all stories (including the linked content), we used a homemade crawler that fetches the content, strips out meaningless text and html and saves the result in chunks of 1 day's worth of content. The code for this crawler is publicly available on GitHub<sup>6</sup>.

The algorithm used for the extraction of meaningful content is GoOse<sup>7</sup>. It uses heuristics to rank the importance and relevance of html elements on a webpage. For example: if it detects a `<div>...</div>` block with a lot of words inside, then that is likely to be important. If,

---

<sup>4</sup><https://news.ycombinator.com/hackernews.html>

<sup>5</sup><https://hn.algolia.com/api>

<sup>6</sup><https://github.com/bcleenders/AIM/tree/master/crawler>

<sup>7</sup><https://github.com/advancedlogic/GoOse>

on the other hand, it finds an html block `<button>Login</button>`, then it will remove the block for it is probably not a relevant part of the text of the page. The resulting dataset is about 4.4 GB in size.

### 3 INITIAL ANALYSIS

We have now established some semantics of the data gathered for our research. Before diving into the data analysis, we will first quantify the data by providing the reader with some statistics.

#### 3.1 OVERALL ACTIVITY

Statistics	
Stories	1544261
Submitters	165126
Upvotes	16668848
Comments	7383865

An interesting remark, is that January 6th 2014 did not have any submissions. The reason for this was a long downtime of the Hacker News website <sup>8</sup>.

#### 3.2 UPVOTE DISTRIBUTION

If users like a story, they can express their approval by upvoting it. These upvotes are then used to rank the stories by popularity and calculate the (currently) most popular submissions. Since the formula for determining a posts “score” strongly favours newly submitted stories, even stories with few upvotes can spike to the frontpage.

To enter the frontpage (the top 30), a story often needs at least ten upvotes, stories that stay on the frontpage for a longer period of time will often have over a hundred upvotes. The distribution of upvotes per story can be seen in figure 3.1. Stories placed in buckets of size 100 based on the number of upvotes they received. The graph plots the size of each bucket. Note that these graphs are scaled logarithmically: in the graph showing all posts (figure 3.1a), 97.5% of the stories are in the first bucket. That means only 2.5% of the stories get over 100 upvotes.

The second graph (3.1b) is a zoomed in version with a bucket size of 1, and shows the distribution for points with 30 upvotes or less, a set of over 1.4 million stories.

There are only six stories with over 2100 upvotes. For completeness’ sake, these are the stories whose popularity is quite literally “off the charts”:

---

<sup>8</sup><https://twitter.com/hackernewsunion/status/420068968464789505> - “Hacker News is DOWN, but your chances of getting into YC if you know how to scale a plain text website are UP”

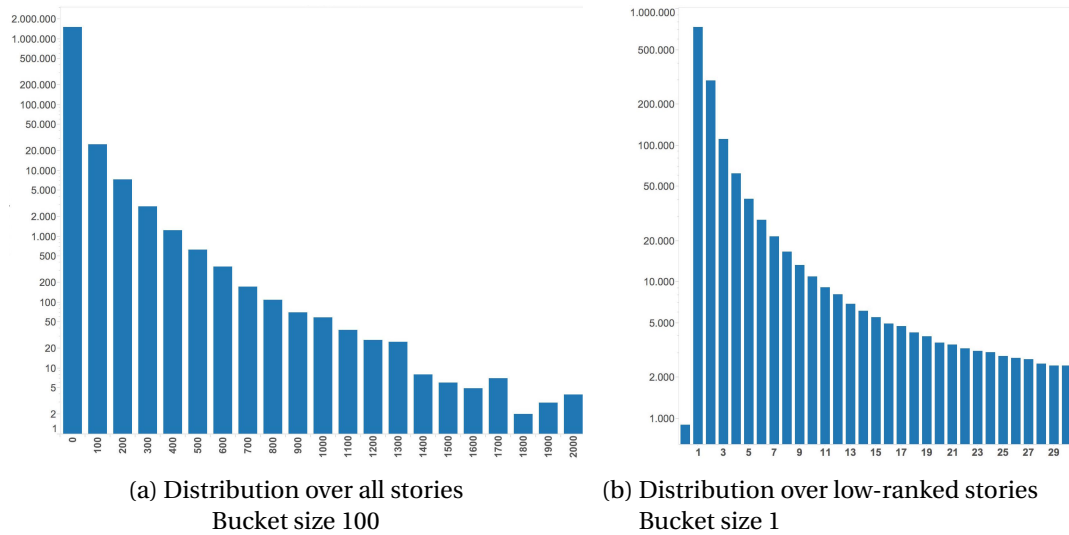


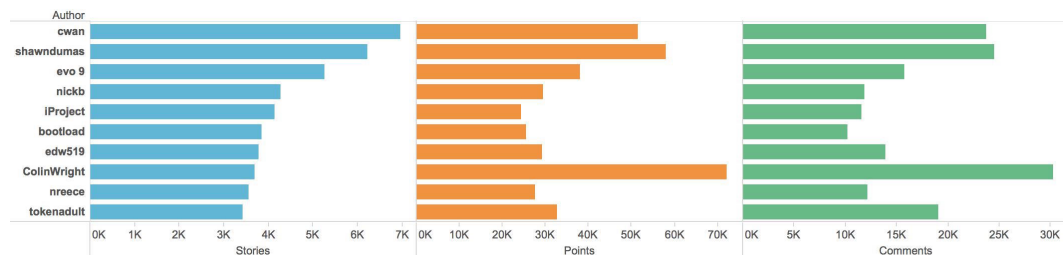
Figure 3.1: Distributions of upvotes over posts (log scale)

Most upvoted stories	
Title	Points
Steve Jobs has passed away	4271
Tim Cook Speaks Up	3086
2048	2732
Don't Fly During Ramadan	2617
Hyperloop	2549
Microsoft takes NET open source and cross platform	2376

### 3.3 TOP USERS

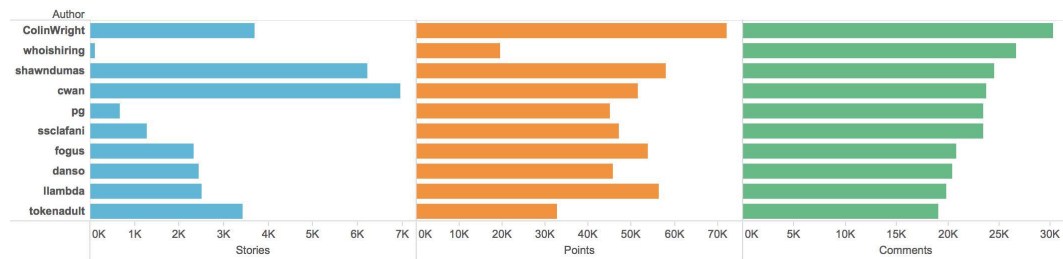
In figure 3.2, we show the statistics for some of the most active submitters. These statistics illustrate that Hacker News has a very active group of users. Together, this top 10 submitted 45,190 stories, which equals an average of 1.5 stories per user per day over the last 7 years.

Figure 3.2: Top 10 submitters by submitted stories



If we take the top 10 not by the number of stories a user has submitted but by the number of comments his or her stories have received, we get the rankings in figure 3.3.

Figure 3.3: Top 10 submitters by received comments



The interesting second place in this top 10 is user “\_whoishiring“. This user has almost no posts (a mere 114) but received 26,649 comments. Experienced Hacker News readers may have already expected this, for this user starts a monthly thread where companies can post job offers and others can respond to these. As such, the user only posts one story per month but its stories are discussed very actively.

### 3.4 TOP DOMAINS

As already stated in the description of the dataset, many of the stories are links to external sites. To provide some insights into which sites attract a lot of attention from the Hacker News community, we made three top ten lists that rank the sites by the same criteria as we did in the previous section: by number of stories, by points and by number of comments.

Most popular domains		
By Number of Stories	By Points	By Comments
techcrunch.com (27.711)	github.com (400.373)	techcrunch.com (172.854)
github.com (26.596)	techcrunch.com (365.207)	nytimes.com (153.471)
youtube.com (21.977)	nytimes.com (287.234)	github.com (127.812)
nytimes.com (18.125)	arstechnica.com (176.633)	arstechnica.com (80.666)
medium.com (14.172)	wired.com (161.698)	wired.com (75.406)
arstechnica.com (12.657)	medium.com (132.468)	washingtonpost.com (56.257)
wired.com (10.867)	bbc.co.uk (98.031)	medium.com (53.924)
bbc.co.uk (8.118)	washingtonpost.com (97.537)	bbc.co.uk (51.828)
en.wikipedia.org (7.058)	youtube.com (96.339)	theatlantic.com (41.530)
businessinsider.com (6.877)	theatlantic.com (77.628)	online.wsj.com (36.729)

These rankings provide some insights into what types of news are popular. The big geeky news sites (Techcrunch, Ars Technica and Wired) are present and are about as popular as the big newspapers (NY Times, BBC, Washington Post).

The high ranking of github.com (a code hosting site, *not* a news site) can be explained by the large number of open source projects hosted on GitHub that submit links to new versions and press releases on Hacker News. Some examples of these projects are Facebook’s React framework, Twitter’s Bootstrap, SQLite and io.js.

Notably, GitHub's competitors (e.g. BitBucket, GitLab and Beanstalk) do not show up in these rankings. This demonstrates GitHub's overpowering popularity in the code hosting market, at least in the open source community.

## 4 TOPIC DETECTION

It's data reduction and clustering, really...

Say that we did two types of unsupervised clustering: k-means and LDA. K-means

### 4.1 LATENT DIRICHLET ALLOCATION

Let's start by formally define the terms used in Latent Dirichlet Allocation (LDA):

- **Word:** the basic unit of discrete data, defined to be an item from a vocabulary indexed by  $\{1, \dots, V\}$ . We represent words using unit-basis vectors that have a single component equal to one and all other components equal to zero. Thus, using superscripts to denote components, the  $v$ th word in the vocabulary is represented by a  $V$ -vector  $w$  such that  $w^v = 1$  and  $w^u = 0$  for  $u \neq v$ .
- **Document:** a sequence of  $N$  words denoted by  $\mathbf{w} = (w_1, w_2, \dots, w_N)$ , where  $w_n$  is the  $n$ th word in the sequence.
- **Corpus:** collection of  $M$  documents denoted by  $\mathbf{D} = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M\}$ .

LDA is a generative probabilistic model of a corpus. The basic idea is that documents are represented as random mixtures over latent topics, where each topic is characterized by a distribution over words.[?]

-

### 4.2 WORD2VEC & K-MEANS

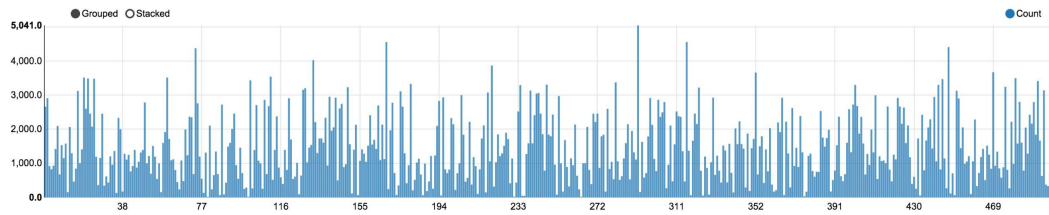
Explain two steps and how awesome it worked

### 4.3 DISCUSS IMPLEMENTATION

## 5 RESULTS

Make some awesome plots from the topics over time, rank by (upvotes | comments | number of articles published)

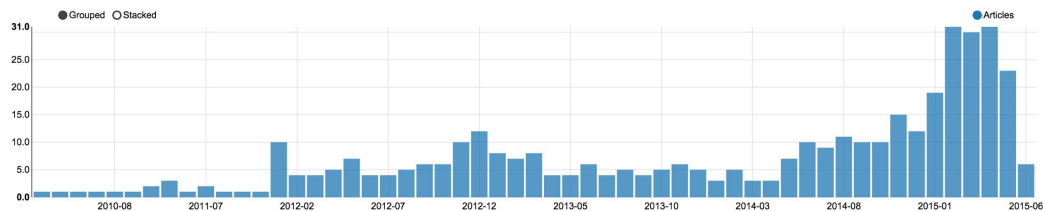
Figure 5.1: Topic group sizes



## 5.1 POPULARITY PLOTS

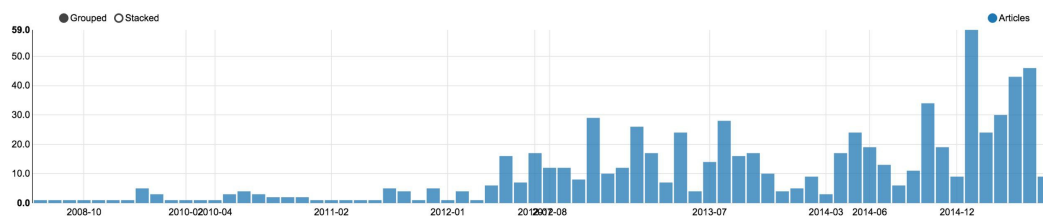
In this section, we will show some awesome plots of how the popularity of various topics changed over time.

Figure 5.2: Popularity of Docker, CoreOS, etcd, containers, OpenStack



We start the overview with a trend plot (figure 5.2) of a new and upcoming technique: containerization. Docker and CoreOS were released in March resp. October 2014 and before that, containerization was already used as a term for separating Linux processes. The trends show how Docker started gaining some real traction within half a year and is really booming in the first half of 2015.

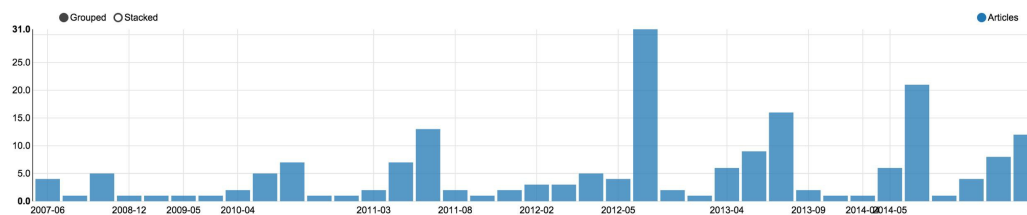
Figure 5.3: Popularity of Elon Musk, SpaceX, Tesla, Hyperloop



**Something about Elon's awesome projects :)**



Figure 5.4: Popularity of keynote, Apple, live, @scale, conference

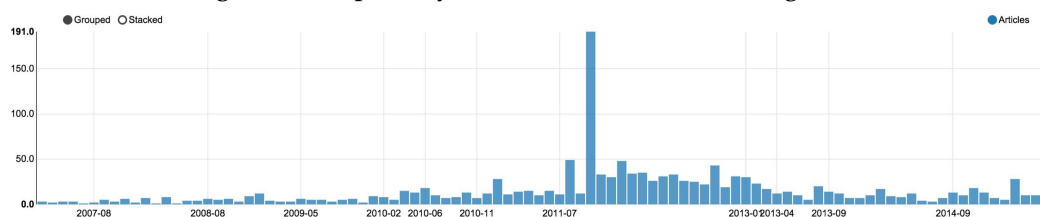


In the popularity charts of keynote-related topics (figure 5.4), one can see patterns of three months of increasing popularity. The highest point is always the month in which an Apple keynote took place; the two months before usually have rumours about what new products will be introduced.

The big spike halfway 2012 marks the announcement of Siri, Apple's artificial intelligence personal assistant.

Unfortunately, Hacker News started after the release of the iPhone (January 2007; one month too early), so we cannot see the data for that release.

Figure 5.5: Popularity of Steve Jobs, Wozniak, imagineers



On the topic of Apple: not only it's products spark the interest of the Hacker News community; it's former CEO also had a big impact in the community. We show figure ?? to indicate the enormous impact the death of Steve Jobs had. He shares this topic group with his co-founder Steve Wozniak, but the news of his death (October 5th, 2011) dwarfs all other events in Hacker News history.

Figure 5.6: Popularity of Raspberry Pi, Kindleberry

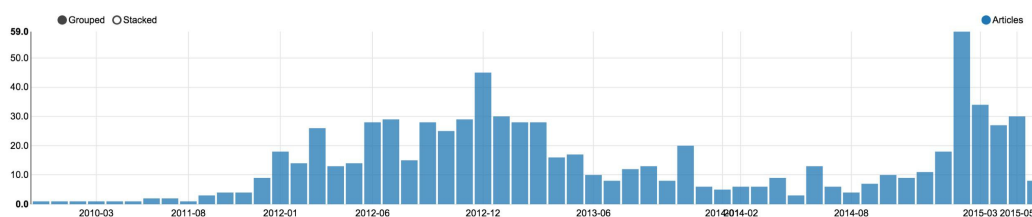
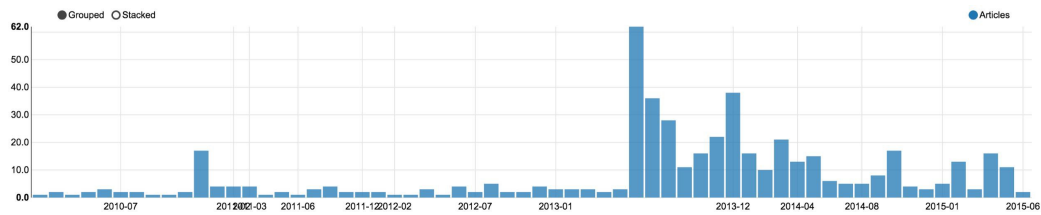


Figure 5.7: Popularity of Edward Snowden, Wistleblower, leaks, Reddit



## 6 CONCLUSION

We're awesome!