

# New York Times Article Abstract Analysis using Hadoop and NLTK

By: Lucas Nunno ([lnunno@cs.unm.edu](mailto:lnunno@cs.unm.edu))

# Part 1: Data Acquisition

# Data acquisition

- Used the python **requests** module.
  - Used the offset parameter to load new pages of abstracts and slept 1/8<sup>th</sup> of a second between each request to abide by the NYT API terms of use.
- Loaded JSON response into python dictionary and then exported as a single large JSON file containing all the articles and all metadata.  
(~40,000)
- In a separate script, I export this JSON data to a CSV file with the docIDs, URLs, and abstracts.
  - This is also where I check for **duplicates**. I have a set of URLs that the exporter has seen, if this URL is in this set the program prints a warning and does not export it.

## Part 2: Preprocessing and tf-idf

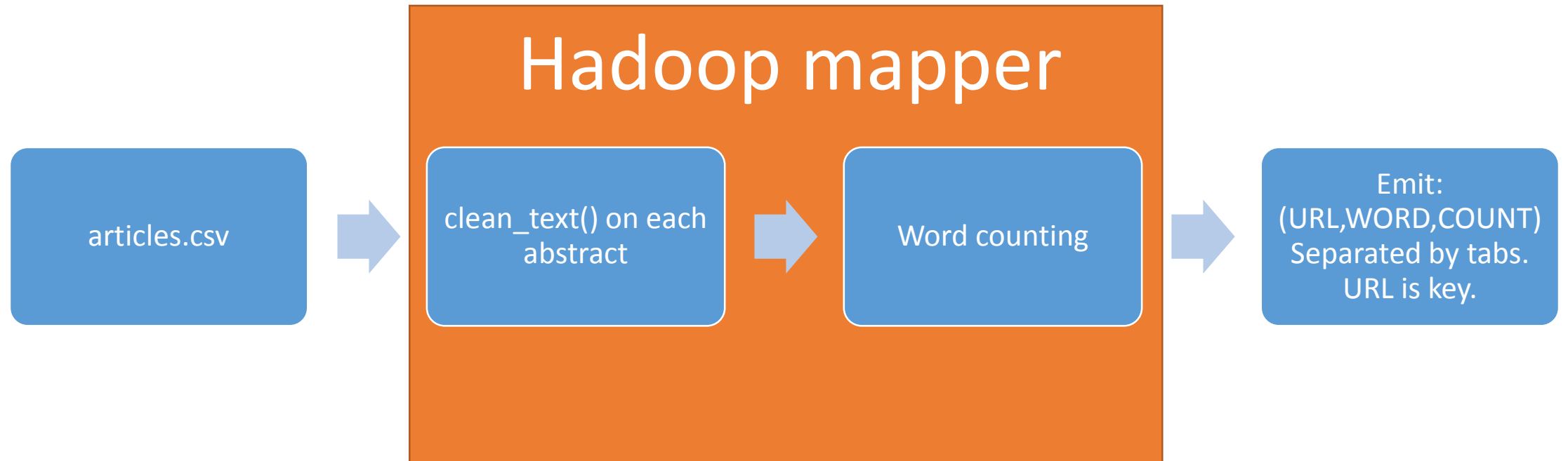
# Preprocessing

- Used the python natural language toolkit (NLTK) module for most of the preprocessing tasks. The algorithm is as follows:
  1. Convert text to lowercase.
  2. Remove punctuation and numbers.
    - Simple regex substitution:  
`remove_pattern = re.compile(r'^a-z\s')`
  3. Remove stopwords.
    - See: `nltk.corpus.stopwords`
  4. Stem all the remaining words.
    1. See: <http://www.nltk.org/api/nltk.stem.html#module-nltk.stem.porter>
  5. Output the cleaned abstract.

# tf-idf is broken up into a couple of stages

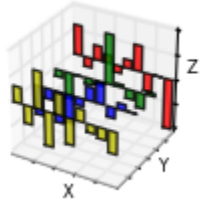
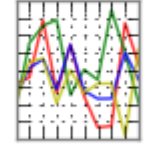
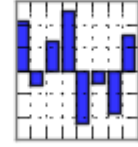
1. Term frequency map reduce job.
2. Document frequency and calculating the total number of documents. Done as a separate map reduce job.

# How is this parallelized in Hadoop? – Term frequency



# Term frequency

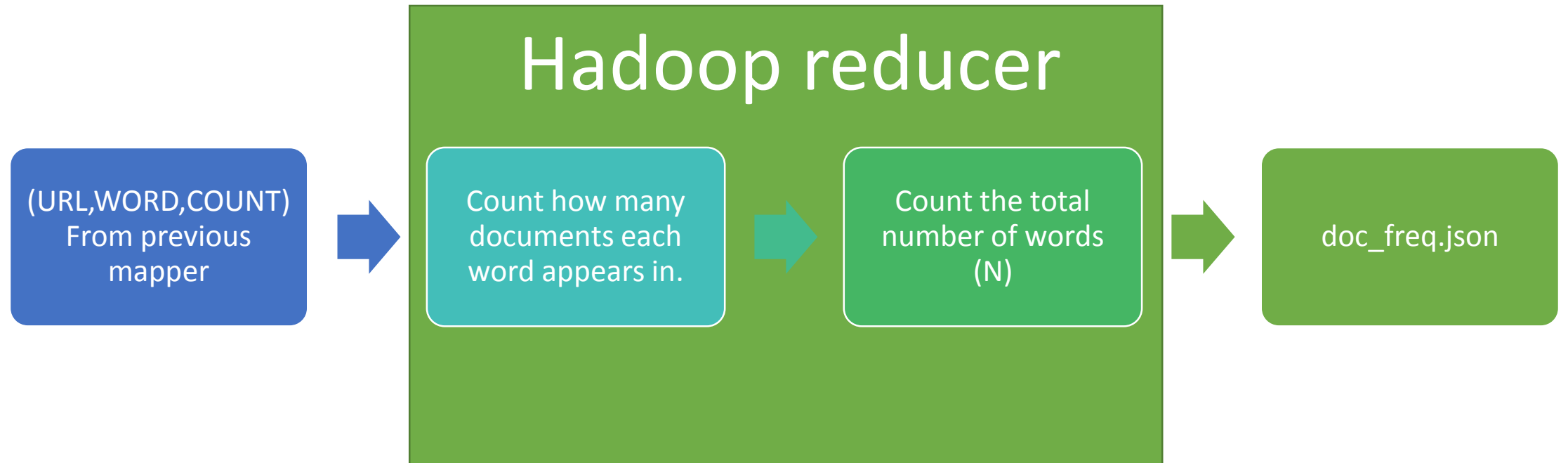
pandas  
 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$



- Dictionary is constructed of each document, the words it contains, and the frequencies of these words – this is all provided from the mapper.
- A pandas Series is constructed for all documents. This is sparse, because it specifies the entries which are non-zero.
- Augmented (normalized) term frequency vector is calculated.
  - $tf(t, d) = \frac{f(t, d)}{\sum f(t, d) \forall t \in d}$
  - This removes the bias for longer documents.



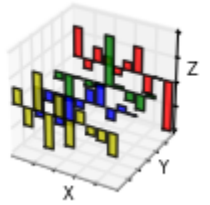
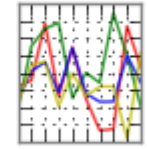
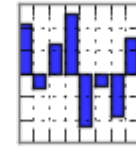
# How is this parallelized in Hadoop? – Doc frequency



# df and tf-idf

pandas

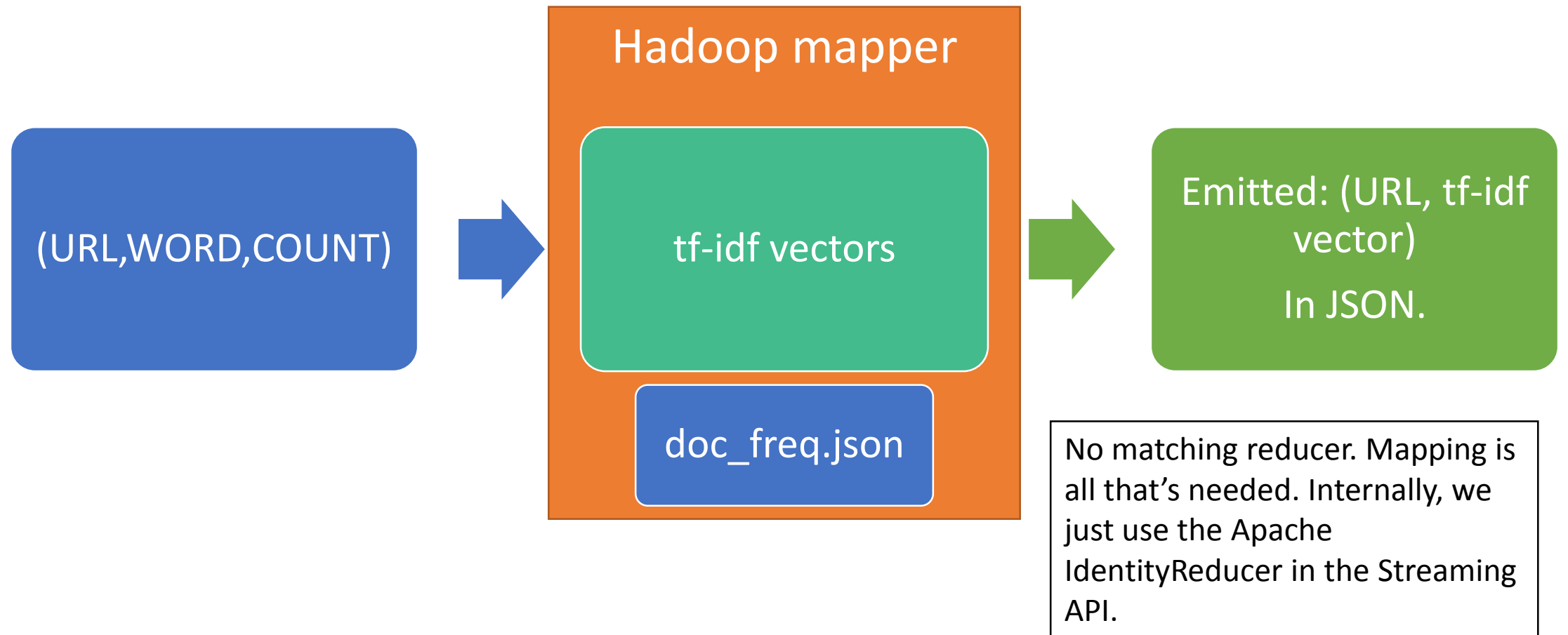
$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



- Document frequency vector is calculated in a reducer by counting the number of unique documents that each term appears in.
- We then multiply each normalized term frequency vector with the idf vector to calculate each document's tf-idf vector.
  - A special implementation detail is that we remove all nonzero entries and only output the information where the tf-idf score is zero. This saves a massive amount of space and time. Every tf-idf vector is very sparse.
- From these data structures we can calculate the inverse document frequency (idf) vector.

$$\text{idf}(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|}$$

# How is this parallelized in Hadoop?



# Part 3: Clustering and Visualization

So now that we have the tf-idf vectors, how do we cluster?



# Clustering methodology

- Apache Spark makes it (relatively) easy to do distributed clustering.
- The API for Spark allows for a text file to be stored in a distributed manner.
- We use the previously calculated tf-idf vectors and transform them into Apache Spark's own SparseVectors.
- Spark's MLlib module is then used for distributed K-means clustering.
  - This uses the SparseVectors that we created and can rapidly separate the vectors into clusters.
- The centroids of the K-means model were analyzed to determine the most discriminative terms per cluster.

# Why K-means?

- Aside from it being the only clustering algorithm available in Spark (which is a good reason), it is one of the most popular clustering algorithms available and used in production in many Big Data applications.
- The Mahout K-means algorithm was also considered, but Apache Spark offered an easier set-up and a Python API so it was used.

# K-means pros and cons

- Can easily tweak the k value.
- Can easily tweak the number of iterations and see if this improves the solution.
- Implementation is simple and easy to understand.
- Will find a local error minimum.
- Widely used.
- Must pick the appropriate K value.
- Have to choose how many iterations to perform.
- Will probably not find the optimal solution.
- Sensitive to outliers.

# How many clusters?

- If you look at the sections in the New York Times (right), there are around 20 primary sections of the paper.
- We initially started out at 20 clusters and then quantitatively analyzed how the number of clusters affected the results.

## Categories

### NEWS

#### [NYT Front Page](#)

#### [International](#)

- [Africa](#)
- [Americas](#)
- [Asia Pacific](#)
- [Middle East](#)
- [Europe](#)

#### [National](#)

- [Education](#)

#### [New York Region](#)

#### [Politics](#)

- [Campaigns](#)
- [Congressional Guide](#)
- [Supreme Court Guide](#)
- [Governor Guide](#)
- [White House Guide](#)
- [Politics Navigator](#)

#### [Business](#)

- [Media & Advertising](#)
- [World Business](#)
- [Your Money](#)
- [Markets](#)
- [Company Research](#)
- [Mutual Funds](#)

#### [Technology](#)

- [E-Business](#)
- [Circuits](#)

#### [Science](#)

- [Earth Science](#)
- [Life Science](#)
- [Physical Science](#)
- [Social Science](#)
- [Space](#)

#### [Health](#)

- [Aging](#)
- [Anatomy](#)
- [Children](#)
- [Fitness](#)
- [Genetics](#)
- [Men](#)
- [Nutrition](#)
- [Policy](#)
- [Psychology](#)
- [Women](#)

#### [Arts](#)

- [Art & Design](#)
- [Dance](#)
- [Music](#)
- [Television](#)
- [Theater](#)

#### [Weather](#)

- [N.Y.C. Metro](#)
- [U.S. Regions](#)
- [International](#)
- [Travel Forecast](#)

#### [Sports](#)

- [Baseball](#)
- [Basketball, College](#)
- [Basketball, Pro](#)
- [Football, College](#)
- [Football, Pro](#)
- [Golf](#)
- [Hockey](#)
- [Other Sports](#)
- [Soccer](#)
- [Tennis](#)

#### [Obituaries](#)

#### [Editorials](#)

#### [Letters](#)

#### [Op-Ed](#)

#### [Corrections](#)

### FEATURES

#### [Automobiles](#)

- [Auto Classifieds](#)

#### [Books](#)

#### [Travel](#)

#### [New York Today](#)

#### [College Times](#)

- [Students](#)
- [Faculty](#)



# Analysis on number of clusters

- Each cluster contains a vector of the 20 highest discriminative terms for that cluster.
- Notice how the average is maximized at 14 clusters.
- Due to these results, we will be using the 14 clusters found to have the highest average discriminative terms.



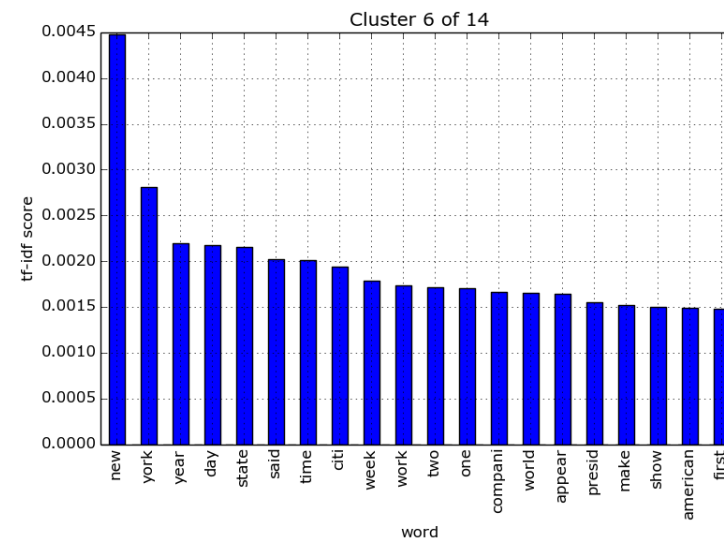
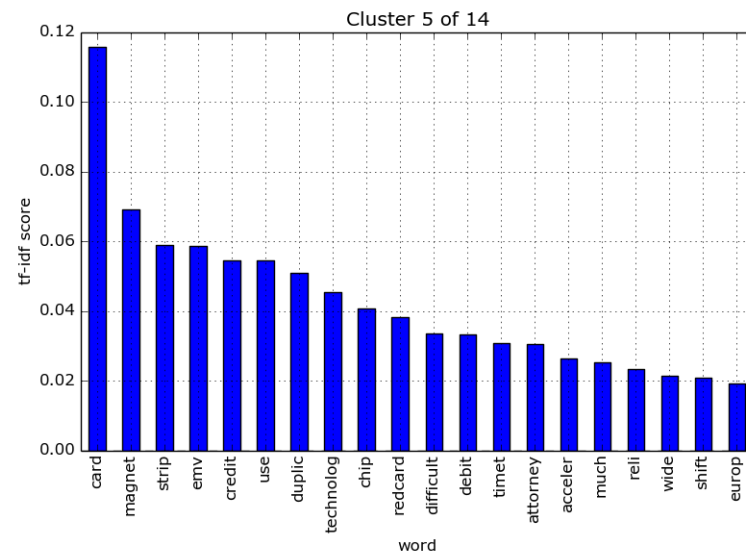
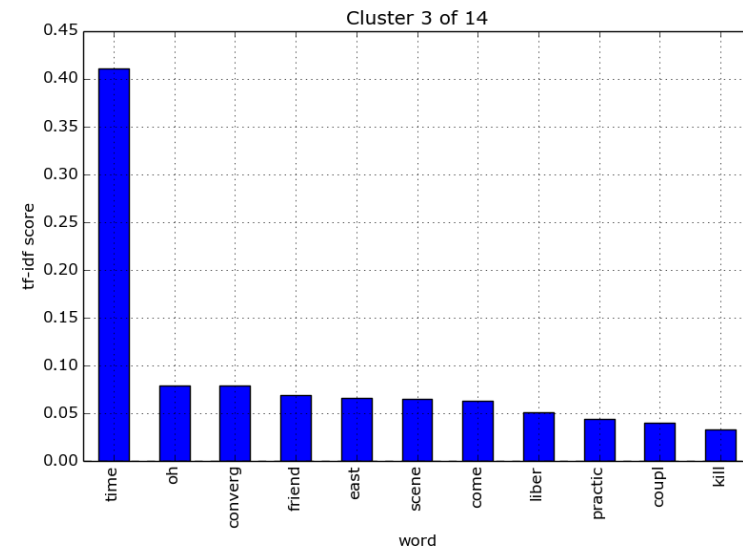
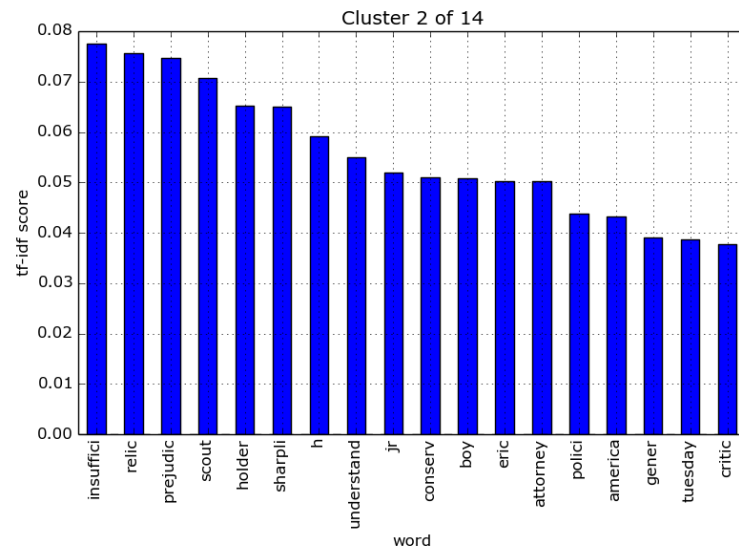
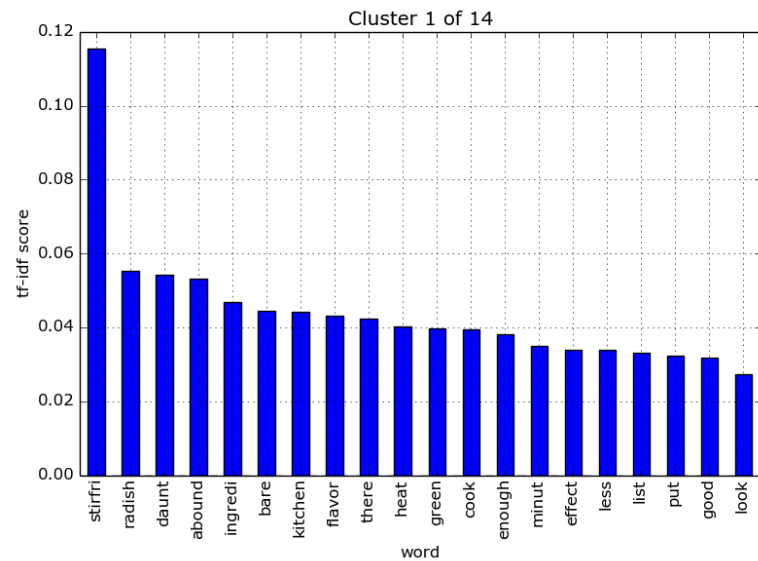
# Things to consider when clustering

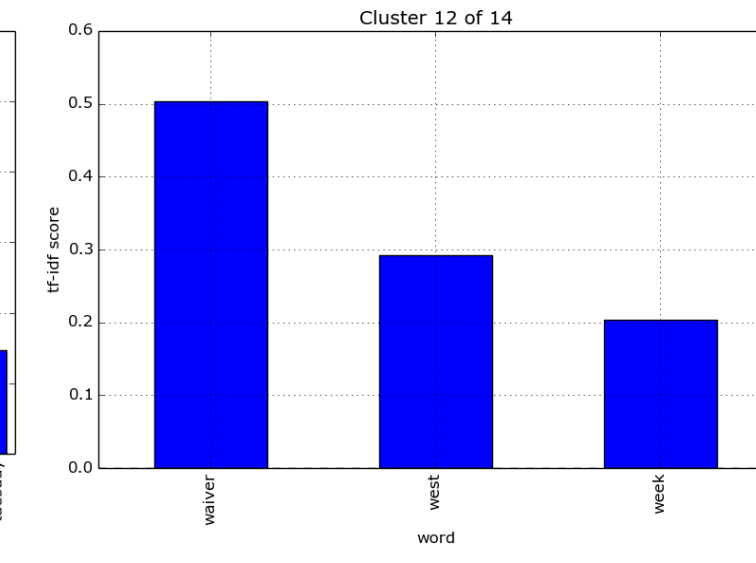
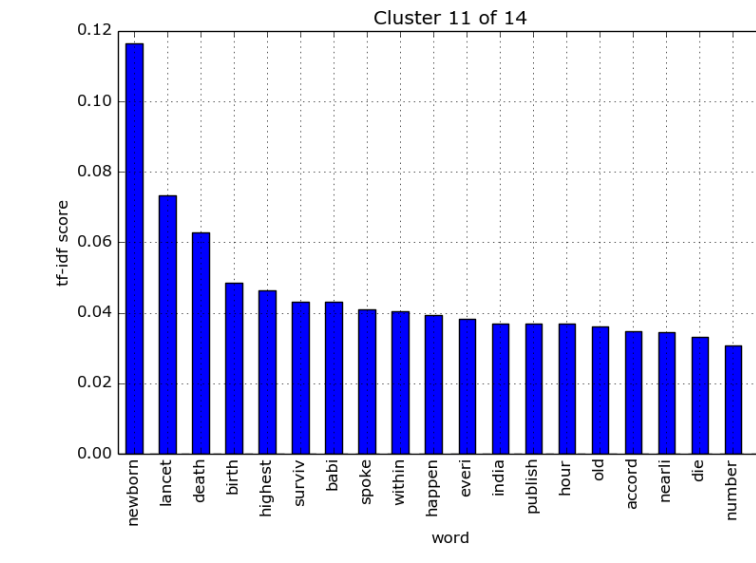
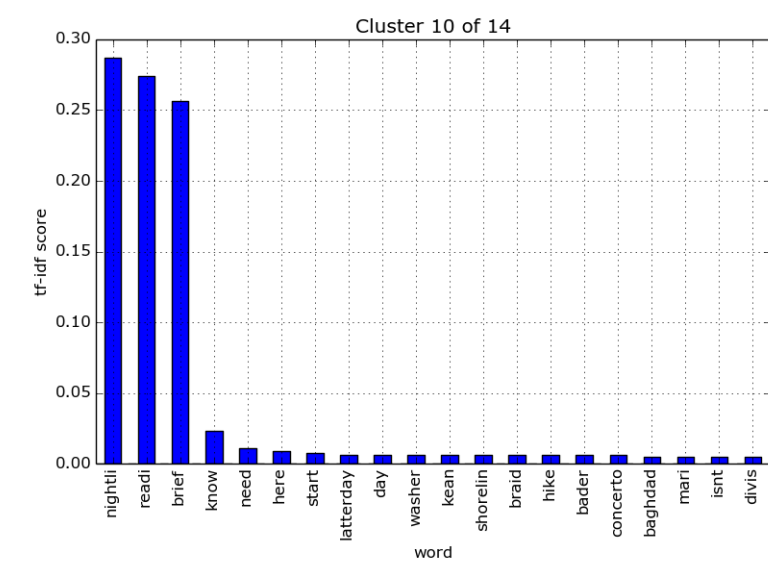
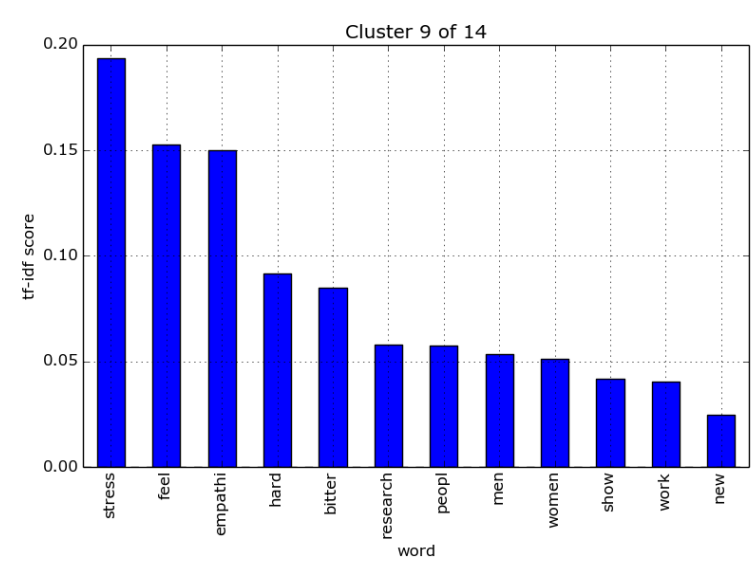
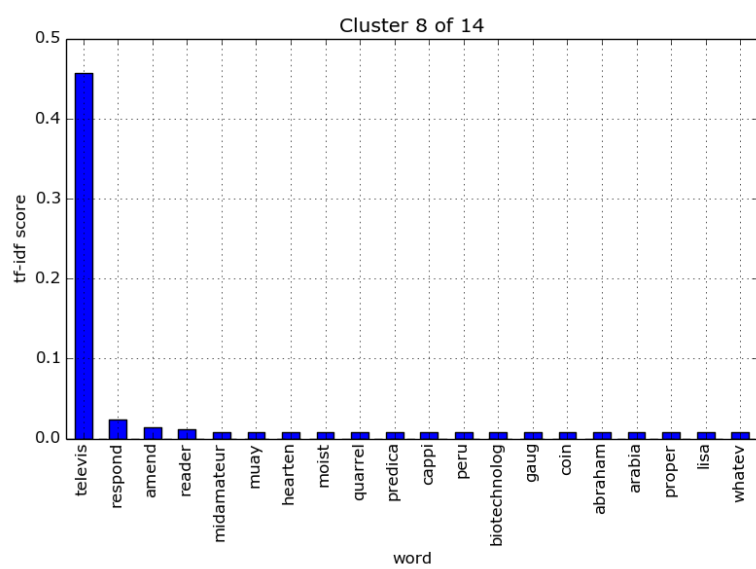
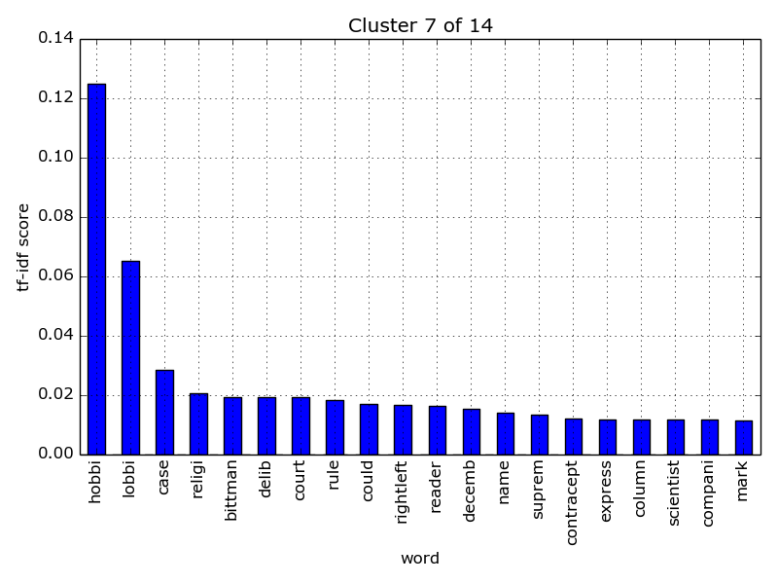
- Our sample size is relatively small (~40k) and not all sections may be represented sufficiently.
- Sampling methodology.
  - It is highly dependent on which articles were being produced the most at the time of collection. We did not discriminate on the types of articles found.
  - This is a combination of the limitations of the New York Times API and the methodology employed in this study.
- This may help to explain some of the results. Some clusters seem to be much more informative than others.

# Visualization notes

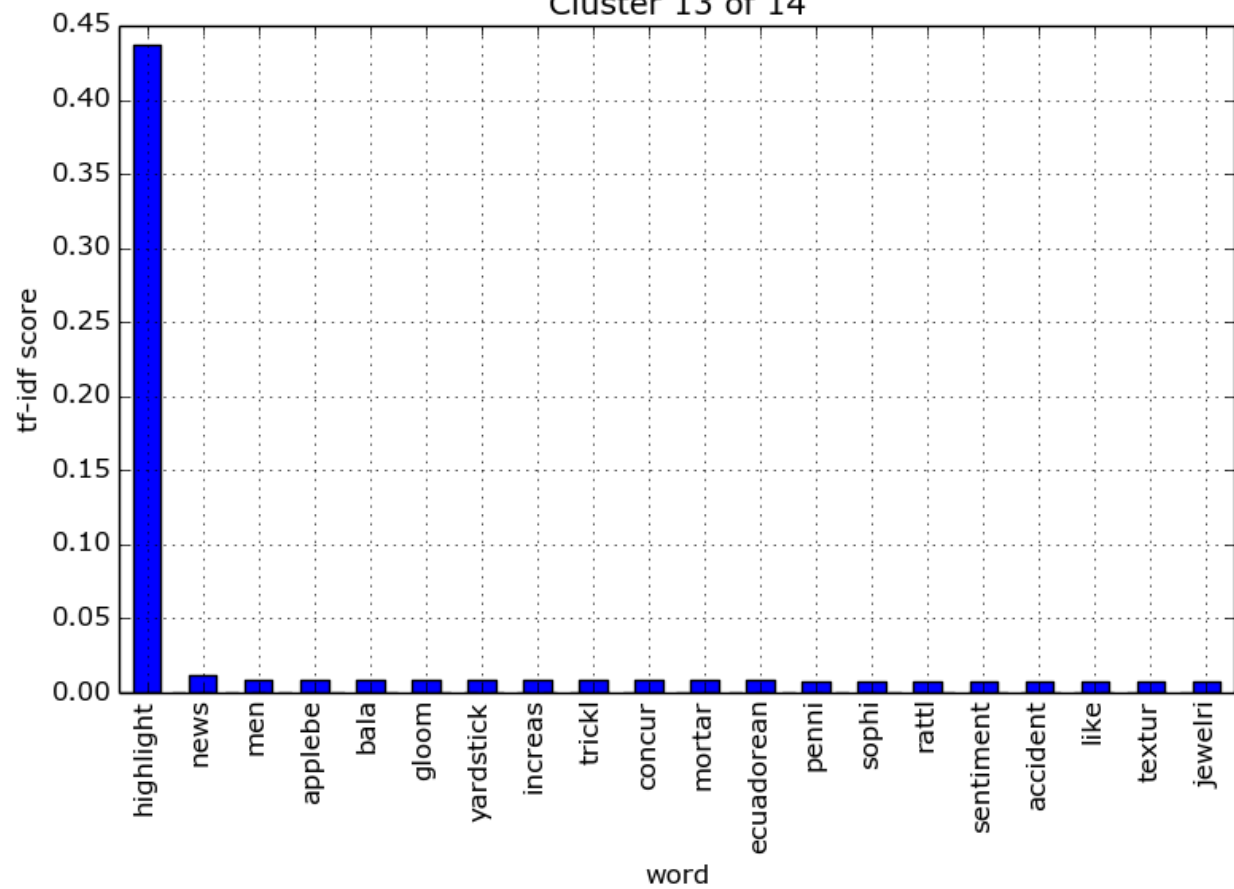
- The top terms were selected based on their tf-idf value for that cluster. This was simply `row.nlargest(NUM_TERMS)` for each DataFrame row.
- Some clusters had less than 20 significant terms, the terms with zero tf-idf scores were eliminated, therefore some of the plots show less than 20 terms per cluster.
- This is an effect of the clustering algorithm used.

14 Cluster Plots

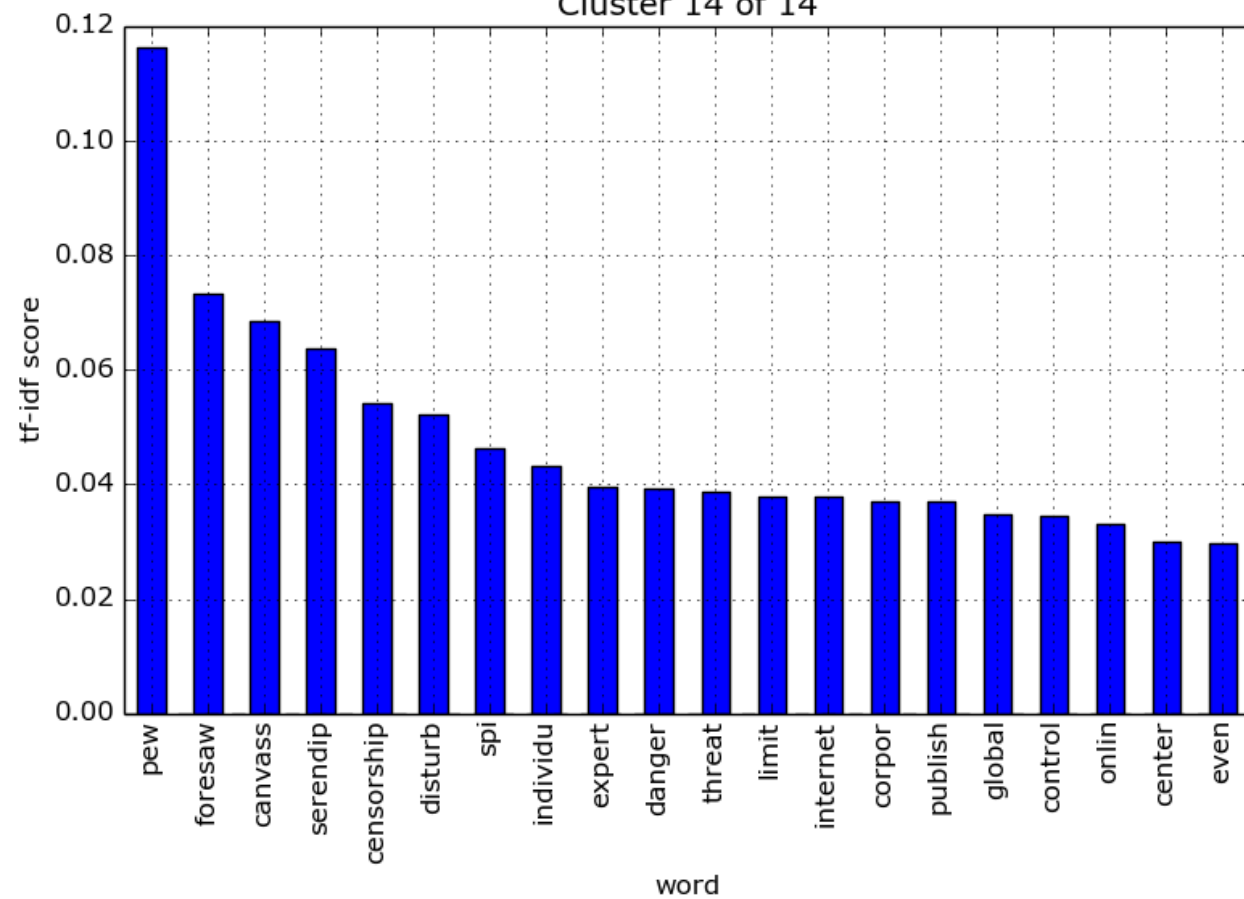




Cluster 13 of 14



Cluster 14 of 14



# 14 Clusters Selected Word Clouds

These were created by hand at <http://www.wordle.net/compose> using the tf-idf scores as weights. Words with higher weights are larger.

Most interesting included, some left out for brevity.



Cooking cluster?



A word cloud of terms related to the Eric Star case. The words are arranged in a roughly rectangular shape, with some words oriented vertically and others horizontally. The colors of the words range from dark blue to brown. The words include:

- scout
- gener
- tuesday
- critic
- boy
- conserv
- prejudice
- eric
- america
- sharp
- relic
- police
- attorney
- understand
- insufficient
- holder

east converg oh  
practic coupl come  
time friend scene  
liber kill

emv  
timet  
reli  
chip  
duplic  
debit  
magnet  
wide  
credit  
shift  
difficult  
strip  
use  
europ  
acceler  
attorney  
much  
redcard  
technology

Local news?



hobby  
lobbi  
case  
rightleft  
delib  
column  
could  
decemb  
bittman  
express  
mark  
reader  
suprem  
scientist  
religi  
contracept  
rule  
court  
name  
compani

respond  
televi  
s

A word cloud of various terms related to television and media, arranged in a circular pattern around the word 'televi'. The words are in different colors and orientations, including: respond, televi, s, reader, capi, proper, direction, amend, value, media, clip, press, function, and la.

televi

reader  
capi  
proper  
direction  
amend  
value  
media  
clip  
press  
function  
la

Lifestyle & health cluster?





# briefly readi

nightli

know

A word cloud with the word 'newborn' as the largest and most central element. Other prominent words include 'lancet', 'death', 'birth', 'happen', 'publish', 'surviv', 'highest', 'india', 'old', 'number', 'nearli', 'spoke', 'tuesday', 'accord', 'die', 'everi', 'babi', 'within', and 'hour'. The words are arranged in a circular pattern around the central 'newborn' and are colored in various shades of green, yellow, orange, and red.

newborn  
lancet  
death  
birth  
happen  
publish  
surviv  
highest  
india  
old  
number  
nearli  
spoke  
tuesday  
accord  
die  
everi  
babi  
within  
hour

A word cloud featuring various terms related to internet censorship and surveillance. The words are arranged in a roughly circular shape, with some words being significantly larger than others. The colors of the words range from dark purple to bright yellow. The words include: censorship, canvass, pew, foresaw, serendip, disturb, limit, expert, individu, control, corpor, spi, danger, publish, threat, internet, global, onlin, even, center, and limit.

# Thank you.

Check out the source on Github: <https://github.com/lnunno/big-data-nyt-tf-idf>

