# Lab 2

**Team:**

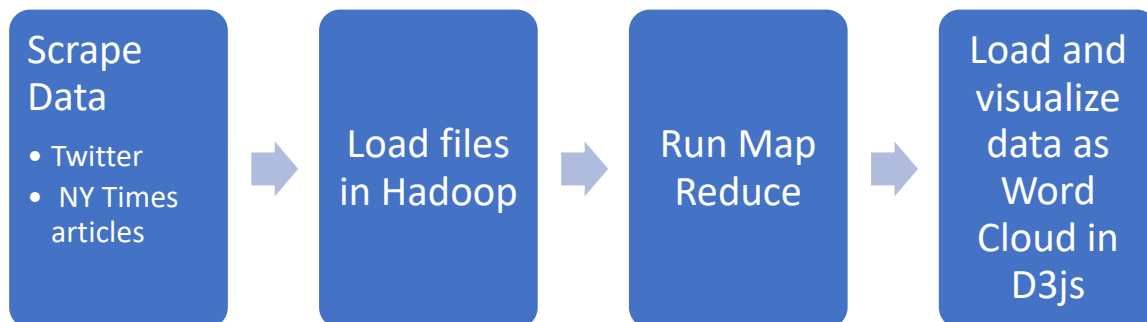| Name | UBIT Name | UBID |
|---|---|---|
| Krithika Krishnan | kk242 | 50169047 |
| Varad Tupe | varadsha | 50249001 |

**Topic used for analysis:**
Gun violence

**Data Source:**
1. New York Times articles
   - All articles details are fetch Article Search API provided by NY Times.
   - Using request and Beautiful Soup packages each article is scraped using the article urls procured from the article search API.
   - Each article is written to a separate file.
2. Tweets from Twitter
   - All tweets are fetched using Tweepy Package in Python.
   - UTF-8 encoding is used while saving the tweets in file.
   - Tweets collected each day are stored in separate txt files.

**Key words used in search:**

gun violence, gun control, school shooting, parkland, gun, gun, mass shooting, gun laws

**Data Flow**

| Scrape Data<br>• Twitter<br>• NY Times articles | → | Load files in Hadoop | → | Run Map Reduce | → | Load and visualize data as Word Cloud in D3js |
| --- | --- | --- | --- | --- | --- | --- |

**Algorithm**

**Top 100 words:**
Mapper:
1. Read input file.
2. Remove unwanted symbols and punctuation
3. Remove stop words.
4. For all words in doc
   Emit (word, 1)
Reducer:
1. Read input from Mapper.
2. Group the data by key.
3. For each key in data, put key and its count in dictionary.
4. Sorted the dictionary based on count descending.
5. For each key in dictionary
   Emit (key, count) [for first 100 keys]


**Top co-occurrence:**
Mapper:
1. Read input file.
2. Remove unwanted symbols and punctuation
3. Remove stop words.
4. For all words w1 in doc
   For all words w2 in Neighbours(w1)
   Emit (w1-w2, 1)
Reducer:
1. Read input from Mapper.
2. Group the data by key.
3. For each key in data, put key and its count in dictionary.
4. Sorted the dictionary based on count descending.
5. For each key in dictionary
   Emit (key, count) [for first 20 keys]

**Results:**
At first, data from a single day was fed into map reduce.
In the NY Times article, commonly occurring words included,
Students, shooting, gun, school, parkland,..
In the Twitter data, commonly occurring words included,
Students, shooting, sheriff, school, parkland,..
There were commonly occurring words between the two sets of data.

When the mapreduce was performed for the bigger dataset, collected over a week,
The commonly occurring words among the NY articles and Twitter data converged
with school, gun, shooting featuring among the top 20 words along with some interesting
words like vegan, democrats, female and education, along with the youtube.

**Prototype wordcloud from NY Times article**
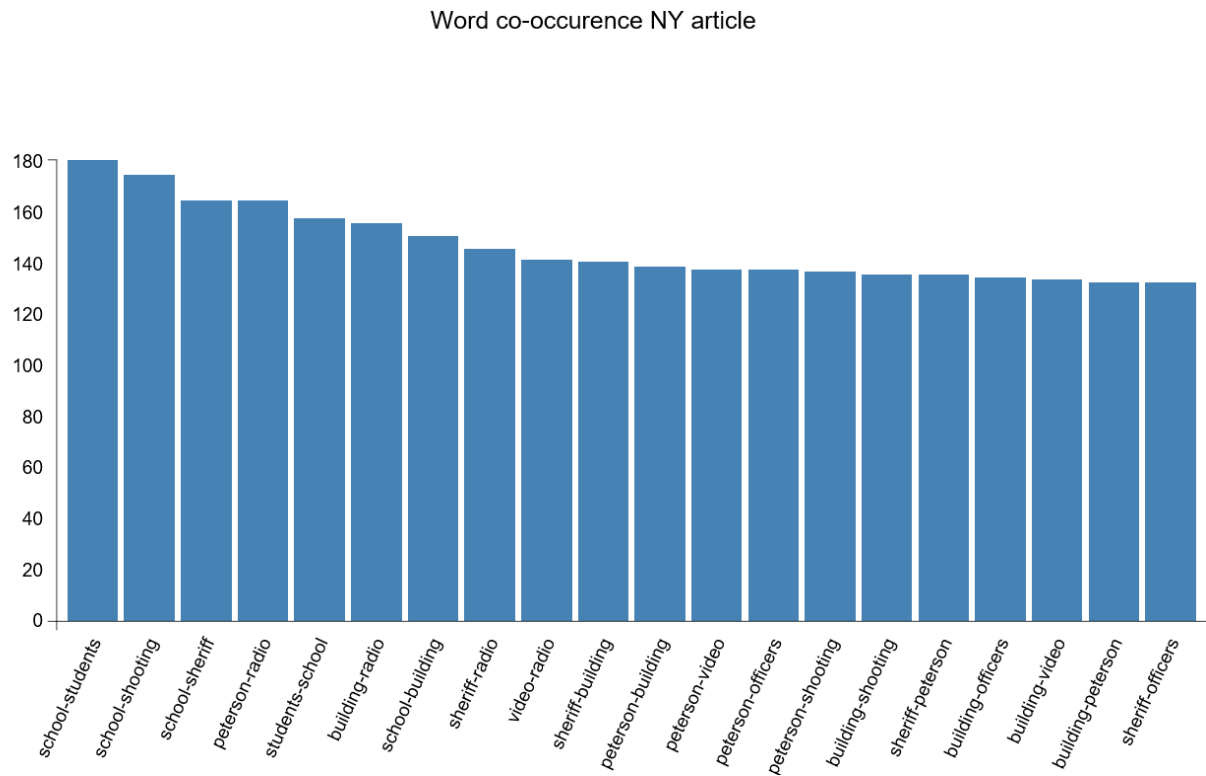
**Prototype wordcloud from Twitter**

group iwillredpillu nprobably nyoutube woman draconian
white aghdam dantdm youtubeshooter won warnings
lives conservative male people nasim attack
failed handgun activist mentally left year
ncame kill ill amp libs pretty mlk
disease liberals nof female shooter
make realjack quiet shooting gun hq laws give
crazy bruno nthey nmaybe damn time registered
falls death yesterday vegan ve promised talk
bully ar san iheartmindy
facts control education media obama greatest
rfk muslim aa
youtube animal hollywood didn
nicsanely guns killed foreign politicize doesn cwhite purchased
cure stop common
called niranian chrisloesch california legally strictest
speech america understand dem notice

**NY Times article collected over a week**

response law chicago ve schools shot deputy part
including movie life florida rally month
state marjory baltimore hogg
week peterson time place
twitter saturday building gunman student lives
voice sheriff douglas daily shooting students stoneman shootings
parents people sunday city teachers
cruz school research police times change
washington gun make killed israel
york march violence college fla times high called
wednesday radio listen parkland video guns
stop officers stalin friends office health
show control national inside left
find news county group year home university laws

**Twitter data collected over a week**

concerts nutjob
family florida today wear legally registered stricter
students april guns campus public candidate
barnitt john yesterday female hq realcandaceo gunned tragedy
anniversary day youtube hrs school amp
vegan run violence don month
marks democrats gun laws columbine marjory
media narrative shooting march angels
control people high muslim
nmaybe chicago chicag stoneman
julian nra libs parkland state liberals
issue orange year education pretty paul didn Instagram
work town years douglas hall
time wi quiet honor edelman
liberal shooter survivor shoots
ryan gop kids mass make

**Word co-occurrence:**
**NY Data**

Word co-occurence NY article



**Twitter Data**

Word co-occurence of Twitter data

**Discussion of Word pair results:**

In the graphs above, the actual count of the word co-occuring are plotted as bar plots with the count on the y axis and the corresponding word pairs plotted in the x axis.
- There seems to be repeated word pairing in the NY and twitter data.
- School-students, school-shooting, students-school featured in among the top 20 word co-occurences for the NY Times articles.
- Yesterday-shooting, Liberals-shooting, gun-control, are among the featured in among the top 20 word co-occurences for the Twitter data.

There seems to be a co-occurrence of shooting, students on the most co-occuring word pairs.

**Folder Structure:**
- lab2
  - dataFetch
    - Python notebooks for fetching NY & Twitter data
  - Hadoop
    - Input
      - TwitterData
      - protoTwitterData
      - NewsData
      - protoNewsData
    - Output
      - Twitter & NY wordcount for top 100
      - Twitter & NY wordpair for top 20 co-occurences
    - README.txt
  - results
    - visualization
      - index.html (to view wordcloud)
      - barchart
        - barchartTwitter.html
        - barchartNY.html

**Citation:**
**Michael Noll for mapReduce**
**And d3.js from Jason Davies d3.js wordcloud template**