# Connecting the Dots Between News Articles

Guide : Amitabha Mukerjee
Ankit Modi (10104, amodi@)
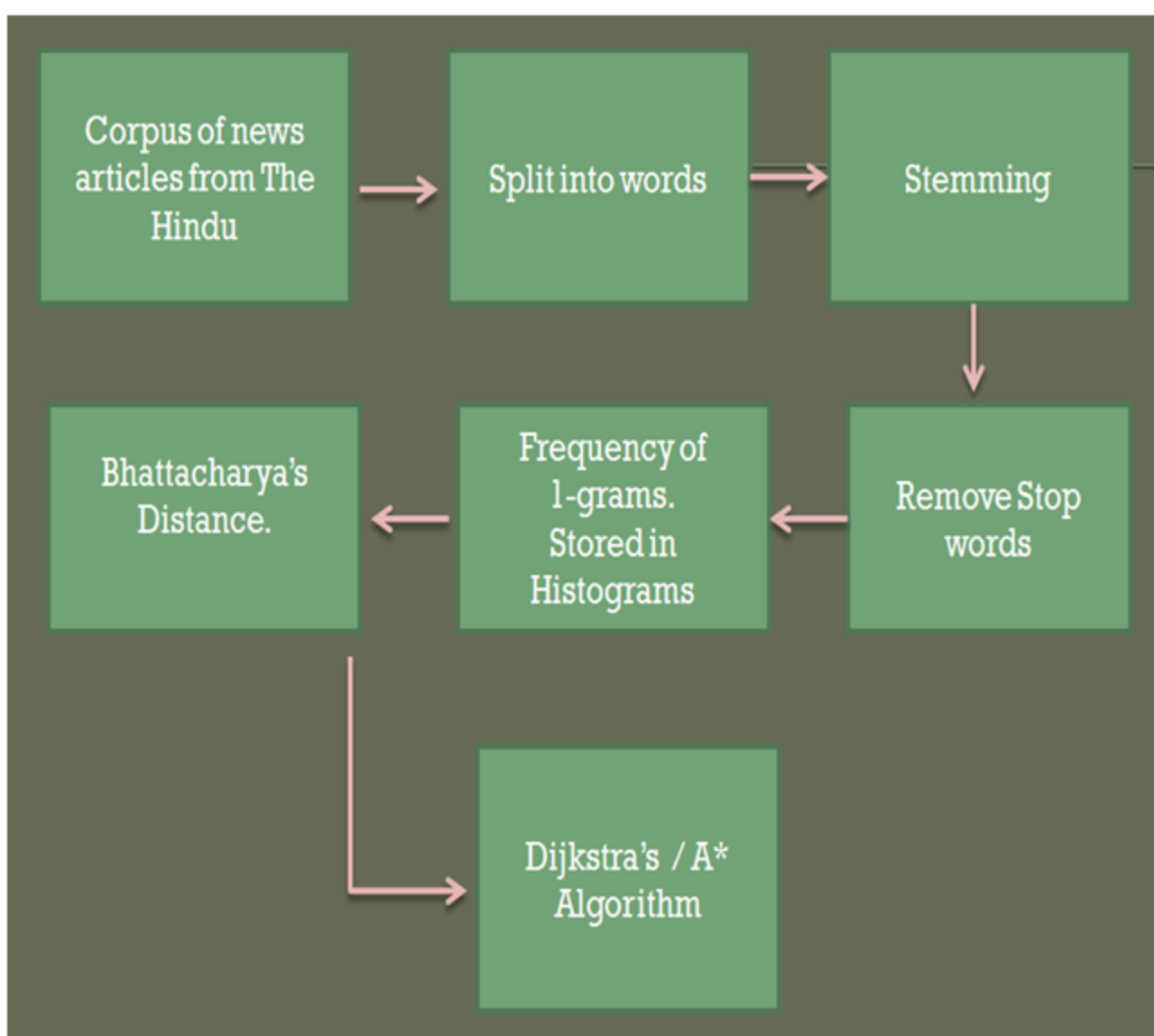Chirag Gupta (10212, gchirag@)

## Motivation

- Tackling Information Overload
- Navigate Between Topics
- Seeing the bigger picture
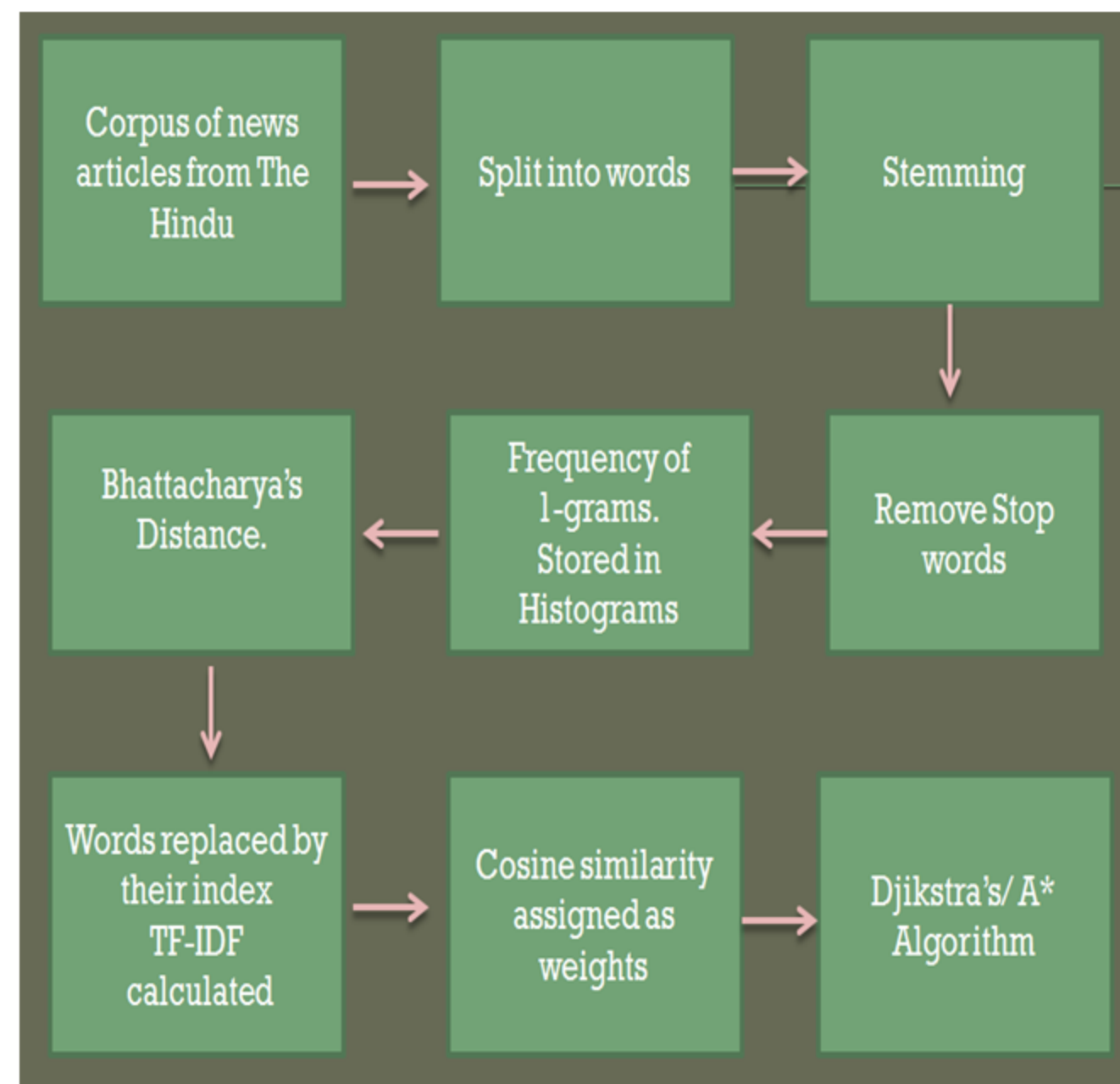- News Browsing : Primary use of Internet

## Related work

- Dafna Shahaf and Prof. Carlos Guestrin used Bipartite Graph and notion of influence (2010)

- Hossain, Gresock, Edmond, helm, potts and ramakrishnan used *Soergel Distance* and *A\* algorithm* (2012)

- Shahaf, Guestrin and Horvitz used m-coherence to compare similarity between articles (2012)

- R. Nallapati, A. Feng, F. Peng, and J. Allan. Event threading within news topics. In *CIKM '04, 2004*

- Q. Mei and C. Zhai. Discovering evolutionary theme patterns from text: an exploration of temporal text mining. In *KDD '05, 2005*

## Approach 1



Corpus of news articles from The Hindu → Split into words → Stemming → Remove Stop words → Frequency of 1-grams. Stored in Histograms → Bhattacharya's Distance. → Dijkstra's / A\* Algorithm

## Approach2



Corpus of news articles from The Hindu → Split into words → Stemming → Remove Stop words → Frequency of 1-grams. Stored in Histograms → Bhattacharya's Distance. → Words replaced by their index TF-IDF calculated → Cosine similarity assigned as weights → Dijkstra's/ A\* Algorithm

## Coherence

We evaluate our results using *coherence* as a parameter. We used two different methods to calculate *coherence*.

- **Coherence1**$(d_1, ..., d_n) = _{i=1...n-1} \min \Sigma_w 1(w \in d_i \cap d_{i+1})$
This is the minimal transition score .
Ref : Dafna Shahaf and Prof. Carlos Guestrin : Connecting the dots between news articles. *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD) 2010*.

- **Coherence2**$(d_1, ..., d_n) = _{i=1...n-1} \min \{ \text{cosine-similarity } (d_i, d_{i+1}) \}$
This measure guarantees a value between 0 and 1, and hence, gives the percentage similarity between documents.

## Tools And Resources used

- **Gensim:** www.radimrehurek.com/gensim
- **NetworkX:** http://networkx.github.io/documentation/latest/index.html
- **matplotlib:** http://matplotlib.org/

## Data Set

We have created our own corpus of 1000 articles downloaded from The Hindu spread over around 40 topics. Some of the topics are Nitish Katara murder case, Jessica Lal murder case, US presidential elections, India's world cup triumph etc.
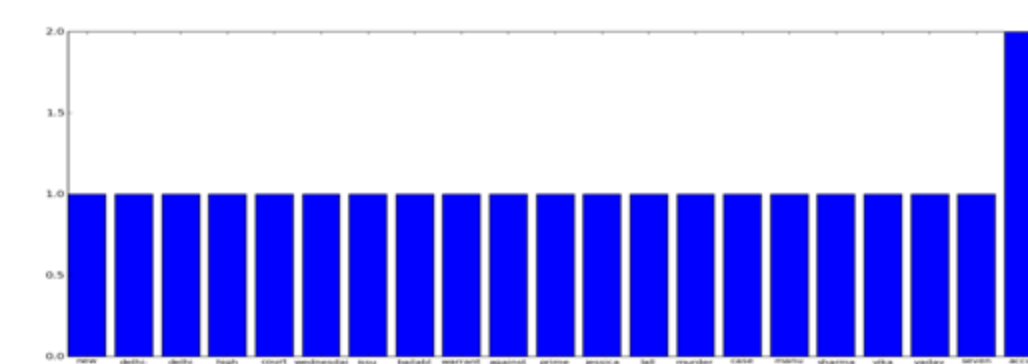We expect it to be very useful for future projects related to Data-Mining. Following is a snapshot of our data set



## Results

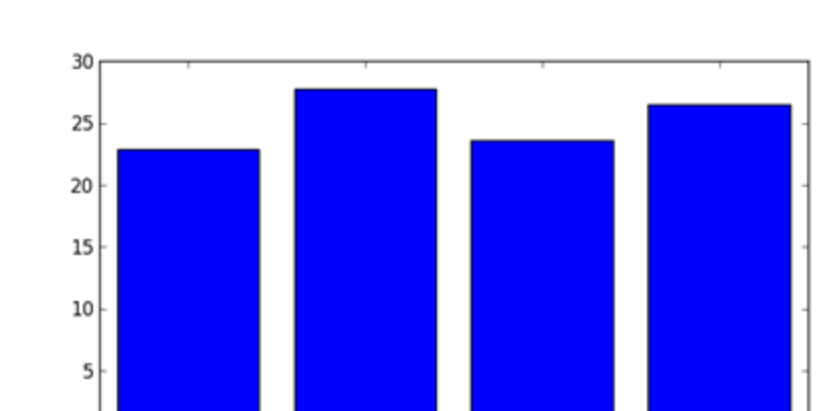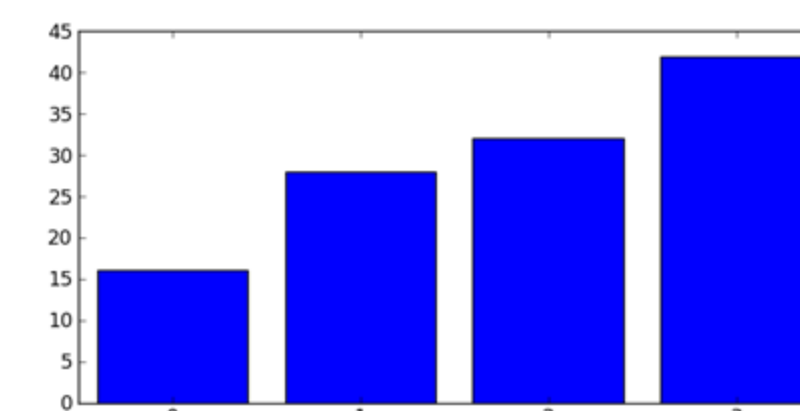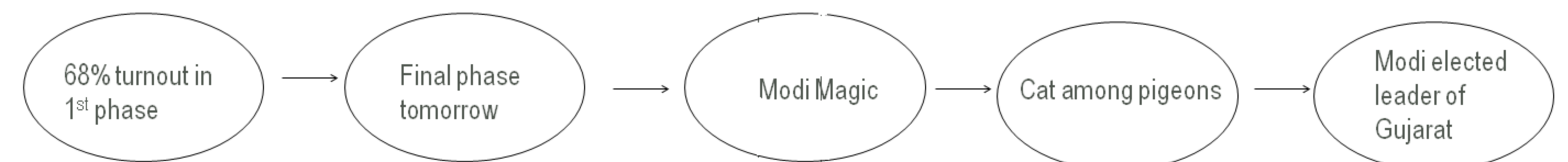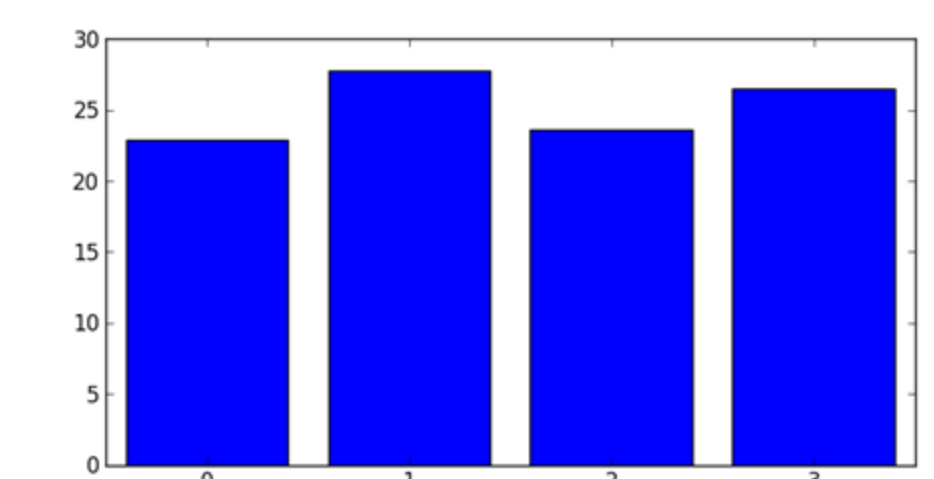| Weight | Algorithm | Chain | Coherence1 | Coherence2 |
|---|---|---|---|---|
| Bhatt. Dist | Dijkstra | 123-117-105-113-111 | 11 (105-113) | 16.05% (105-113) |
| Cosine-Similarity | Dijkstra | 123-103-102-110-111 | 16 (110-111) | 22.88% (110-111) |
| Bhatt. Dist | A\* | 123-118-103-117-111 | 18 (117-111) | 19.24% (118-103) |
| Cosine-similarity | A\* | 123-103-117-110-111 | 16 (110-111) | 20.45% (117-110) |

## Comparison between Word-count and TF-IDF



unigram



TF-IDF

## Comparison between *Coherence1* and *Coherence2*



68% turnout in 1st phase → Final phase tomorrow → Modi Magic → Cat among pigeons → Modi elected leader of Gujarat





## Comparison using *coherence2* between *Bhattacharya's Distance* and *Cosine Similarity*





## Conclusion

- TF-IDF is better than uniram
- *Coherence2* is a better evaluating parameter than *cohorence1*
- Cosine similarity using TF-IDF gives more coherent chain than Bhattacaraya's distance.

## Future Work

- Notion of *m-coherence* can be used instead of *coherence* for better evaluation of results.
Ref: Dafna Shahaf , Prof. Carlos Guestrin and Eric Horvitz : Trains of thought-Generating information maps. *International World Wide Web Conference (WWW), 2012*
- Using different corpuses (other than News articles) may help in important scientific discoveries.