# News Categorization

Sarath Joseph
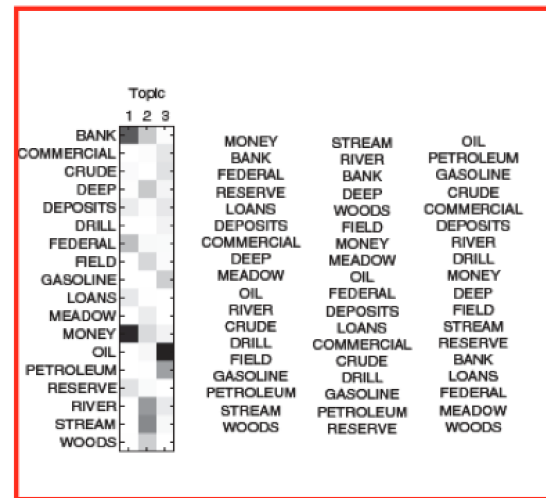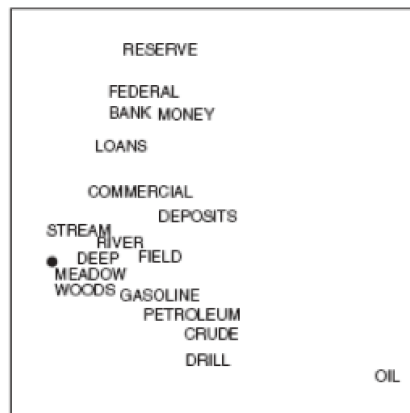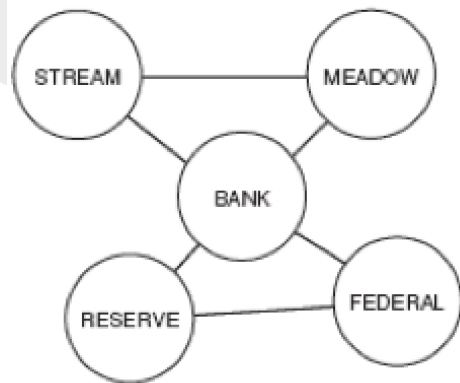
Saurabh Mehta

# Index

# Introduction

- as more information becomes available, it becomes more difficult to find and discover what we need

- we need a way to help us to organize and understand these vast amounts of information

- ex. google books

# Project Description

- We plan to develop a system to aggregate and classify news from diverse sources using machine learning algorithms and techniques.
- We develop the system from the point view of an online news startup.
- As a potential startup that works to aggregate news from multiple sources, it is essential to categorize news coming in from various sources.
- There are several existing ways to achieve this ranging from automatically sorting scraped articles into categories based on keywords or scraping individual parts of the website for categories to obtaining tag based news using News APIs.
- However, dynamic tagging of news data can be challenging. The strong advent of social media platforms as news sources can also be exploited to obtain trending information and for buzz monitoring and machine learning techniques can be employed to generate easily accessible news ranging a wide span of topics.

# Semantic Representation of Text

# Topic Modelling

Topic modelling provides methods for automatically organizing, understanding and summarizing large electronic archives.
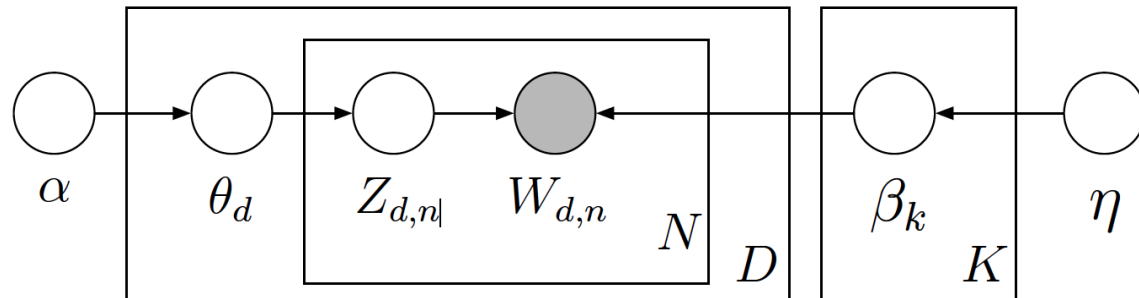
1. Discover the hidden themes that pervade the collection.
2. Annotate the documents according to those themes.
3. Use the annotations to organize and summarize the texts.

Documents are mixture of topics and a topic is a probability distribution over words.

# Discover topics from a corpus

| human | evolution | disease | computer |
|-------|-----------|---------|----------|
| genome | evolutionary | host | models |
| dna | species | bacteria | information |
| genetic | organisms | diseases | data |
| genes | life | resistance | computers |
| sequence | origin| | bacterial | system |
| gene | biology | new | network |
| molecular | groups | strains | systems |
| sequencing | phylogenetic | control | model |
| map | living | infectious | parallel |
| information | diversity | malaria | methods |
| genetics | group | parasite | networks |
| mapping | new | parasites | software |
| project | two | united | new |
| sequences | common | tuberculosis | simulations |

# LDA

# Tools

- Mallet (Java)
- Flask (Python Microframework)
- MongoDb

Python NewsPaper API

Pymongo for MongoDb

# Roadmap

- Extraction of News articles using Python NewsPaper API (Text and metadata) and storing in Mongodb

- Feeding the article text into Mallet
- Parse the Mallet output using Python
- Create new collections in Mongodb with Parsed output
- Creating a Web service for topic based access
- Labelling of Topics
- Creating a client for consumption of news web-service

# Architecture

# Topic Results

Topic 0 -  Middle East

Topic 1 -  International

Topic 2 -  Justice

Topic 3 -  Showbiz

Topic 4 -  Life

Topic 5 -  Stories

Topic 6 -  Travel

Topic 7 -  Politics

Topic 8 -  Education

Topic 9-   The UK and India

Topic 10-  Religion

Topic 11-  Food

Topic 12-  Health

Topic 14-  Trends

Topic 15-  Economy

Topic 16-  Sports

Topic 18-  Theatre and Art

Topic 19 - Housing and Design

Topic 21-  Business

Topic 22 - China

Topic 23-  Technology

Topic 24 - Opinion

Topic 25 - Books and Authors

Topic 26 - Spanish

Topic 27-  Flight and Travel

Topic 28 - US Crime

Topic 29-  Ebola

**Miscellaneous**

Topic 13 - Spam

Topic 17 -  Decisions

Topic 20 - Blogs

# Evaluation

40 articles from the categorization results for
5 topics were observed and the
precision was noted

# Evaluation (contd)

Topic 0
Middle
East


30/40
Precision ~
0.75

# Evaluation (contd)

## Topic 8 Education

34/40

Precision ~ 0.86

# Evaluation (contd)

## Topic 18

### Theatre & Art

35/40

Precision ~ 0.88

# Evaluation (contd)

## Topic 16
Sports

40/40

Precision ~ 1.00

# Evaluation (contd)

Topic 12

Health

39/40

Precision ~ 0.98



### THE 7 BEST STRENGTH EXERCISES YOU'RE NOT DOING

(Life by DailyBurn) -- Every exercise in your strength program has a purpose -- to help you build strength and muscle, burn fat, and improve your fitness. While there's a time and a place for nearly any exercise under the right circumstance, some movements are simply more effective than others. And it should be no surprise that the ones that build a foundation for skills that you'll use in real life will be the most beneficial for improving your fitness and quality of life.DailyBurn: 5 CrossFit workouts that will kick your buttSo how does a lifter ensure they're making all the right moves? If you've plateaued or just aren't seeing the results you're banking on, it's time to get back to basics with these seven moves. From increased strength, better core stability, greater athleticism, and improved overall health, these key exercises need to find their way into your routine.Squats are an exercise many people struggle to perform safely and effectively. Luckily, the goblet squat is a great...

Read more



### HELLO, GREEN MAN

A few days after I wrote about conditions that can mimic dementia, reader Sue Murray emailed me from Westchester County. Her subject line: Have you heard of Charles Bonnet Syndrome?I hadnt, and until about six months ago, neither had Ms. Murray.Her mother Elizabeth, who is 91, has glaucoma and macular degeneration, and has been gradually losing her vision, Ms. Murray explained. So at first, her family was excited when Elizabeth seemed to be seeing things more clearly. Maybe, they thought, her vision was returning.But the things she was seeing patterns and colors, strangers, a green man werent there. She insisted that there were people in the cellar, people on the porch, people in the house, Ms. Murray said. Shed point and say, Dont you see them? And shed get mad when we didnt.Elizabeth and her husband Victor, 95, live in Connecticut, in a house they bought 50 years ago. For a while, the Green Man, as Elizabeth began calling him, seemed to have moved in, too. Shed start hiding things...

# Future Work

- Developing a workflow for

  - Assignment of new articles to topics

  - Reindexing articles

- Exploring relationships between topics

- Live data

# Bibliography

- D. Blei , A. Ng, M. Jordan, **Latent dirichlet allocation** , *Journal of machine Learning research* , 2003

- H. Wallach, **Topic modeling: beyond bag-of-words** , *International conference on Machine learning* , 2006

- Hayes, P., Knecht, L., & Cellio, M. (1988). **A news story categorization system**. *Proceedings of the Second Conference on Applied Natural Language Processing (ANLC '88), 9-17.* http://dx.doi.org/10.3115/974235.974238

- OuYang, L. (n.d.). Newspaper: Article scraping & curation. Retrieved from Lucasao website: http://newspaper.readthedocs.org/en/latest/

Thank you !