

# H6751 Text and Web Mining

Zhao Rui

Course Webpage

# Course Instructors



**Zhao Rui**  
(Instructor)

rui.zhao@ntu.edu.sg



**Chen Zhenghua**  
(Instructor)

zhenghua.chen@ntu.edu.sg

Course web: <https://h6751.github.io/>

# Goals of this Course

**Learn how to analyse unstructured text data**

- Principles and concepts of text and web mining
- Various text mining techniques
  - Pre-processing, text categorization, document clustering, information extraction
- Practical text mining applications
  - Spam detection, sentiment analysis, knowledge graph

# Course Assessment

- Class Participation (5%)
- Assignments:
  - Kaggle Competition (8%):
  - A 90-minutes in-class assignment (12%)
- Group Project (25%)
  - Project Report (15%)
  - Final Presentation (10%)
- Final Exam(50%)

# Course Participation

- Class Participation (5%)

1. **Attending guest speakers' lectures:** In the semester, we have two invited speakers, who are making a great efforts to come lecture for us. We do not want them speaking to a empty room. Your attendance at lectures with guest speakers is expected! In addition, it will be a very awesome chance for networking! You will get 1% per speaker (total 2%) for attending.
2. **Attending two random lectures:** We are going to take attendance at two randomly-selected (non-guest) lectures in the quarter. Each is worth 1% (total 2%).
3. **Karma Point:** Any other act that improves the class, which instructors notice and deems worthy: 1%.

# Assignments

- For Kaggle, it will be a text classification problem
- For in-class Assignments, it will be code-based exam. Open-book and Open-Internet.
- Details will be updated before the release of these assignments.

| <b>Date</b>      | <b>Topic</b>                      |            |
|------------------|-----------------------------------|------------|
| Sat a.m<br>01/18 | Introduction to Text Mining       | <b>ZR</b>  |
| Sat a.m<br>02/01 | Pre-processing for Text Mining I  | <b>ZR</b>  |
| Sat p.m<br>02/01 | Pre-processing for Text Mining II | <b>ZR</b>  |
| Sat a.m<br>02/15 | Information Extraction            | <b>ZR</b>  |
| Sat p.m<br>02/15 | Text Categorization I             | <b>CZH</b> |
| Sat a.m<br>02/29 | Text Categorization II            | <b>CZH</b> |
| Sat p.m<br>02/29 | Document Clustering               | <b>CZH</b> |
| Sat a.m<br>03/21 | Sentiment Analysis                | <b>CZH</b> |
| Sat p.m<br>03/21 | Topic Modeling                    | <b>ZR</b>  |
| Sat a.m<br>04/04 | Neural Network                    | <b>ZR</b>  |
| Sat p.m<br>04/04 | Word Embeddings                   | <b>ZR</b>  |
| Sat a.m<br>04/18 | CNN and RNN                       | <b>ZR</b>  |
| Sat p.m<br>04/18 | Group Presentation                |            |



Let us Start

# Twitter in Chief

- The President not only Make America Great Again, but also Twitter
- He tweets 4178 per year and 11 to 12 per day



**Donald J. Trump** 

@realDonaldTrump

Follow

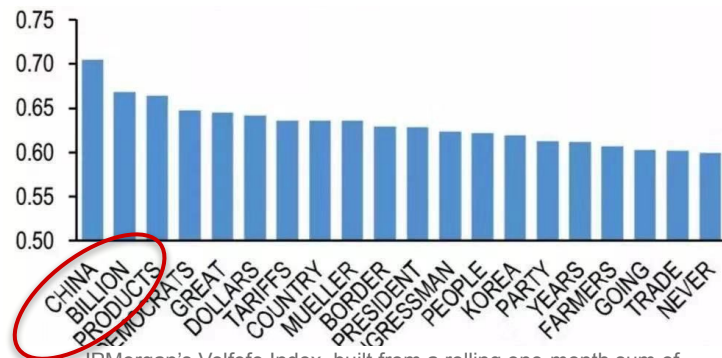


....place in TRADE, it's taking shape in Military Competition." Johnathan Ward, author and China expert. We are winning, and we will win. They should not have broken the deal we had with them. Happy Birthday China!

5:31 AM - 30 Sep 2019

# Volfe Index

- Quantify the market impact of Trump's tweets
- Supervised learning and Natural Language Processing techniques are used to spot “market-moving” tweets
- Volfe Index can explain moves in implied volatility



JPMorgan's Volfe Index, built from a rolling one-month sum of inferred probability that each tweet is market moving (Source: Bloomberg)

# What is More



**Donald J. Trump** ✓

@realDonaldTrump

Follow

Toyota Motor said will build a new plant in Baja, Mexico, to build Corolla cars for U.S. NO WAY! Build plant in U.S. or pay big border tax.

10:14 AM - 5 Jan 2017

27,560 Retweets 95,213 Likes



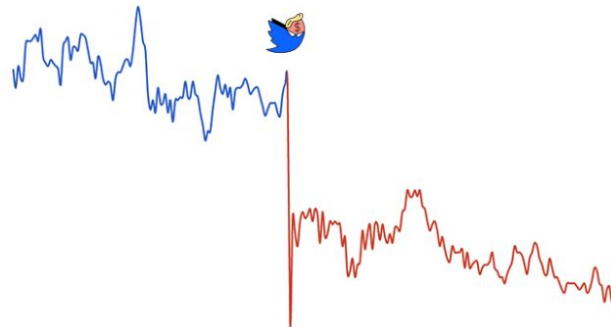
**Donald J. Trump** ✓

@realDonaldTrump

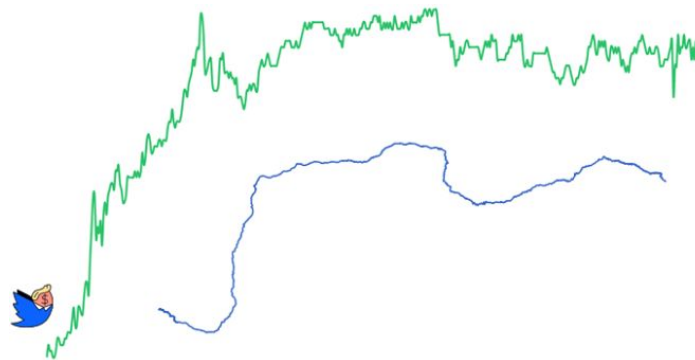
Thank you to Ford for scrapping a new plant in Mexico and creating 700 new jobs in the U.S. This is just the beginning - much more to follow

♥ 76.8K 9:19 PM - Jan 4, 2017

💬 22K people are talking about this



Toyota's NYSE:TM stock price on January 5th 2017



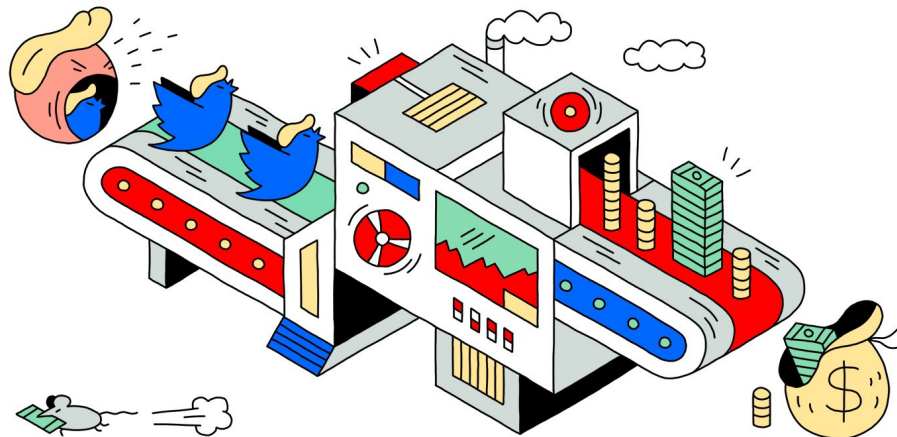
Ford's NYSE:F stock price on January 4th 2017 (Rio Grande for scale)

Source:

<https://medium.com/@maxbraun/this-machine-turns-trump-tweets-into-planned-parenthood-donations-4ece8301e722#yovbh4qc1/>

# Trump2Money

- 1 Open your laptop and write some code
2. Monitor Trump's twitter feed
3. Analyze the twitter  
If it mentions of any publicly traded stocks  
and compute its sentiment
  - a. Long it if the sentiment is positive
  - b. Short it if the sentiment is negative



Source: <https://github.com/maxbbraun/trump2cash>

# What is Machine Learning



**Mat Velloso**

@matvelloso

Follow



Difference between machine learning  
and AI:

If it is written in Python, it's probably  
machine learning

If it is written in PowerPoint, it's  
probably AI

5:25 PM - 22 Nov 2018

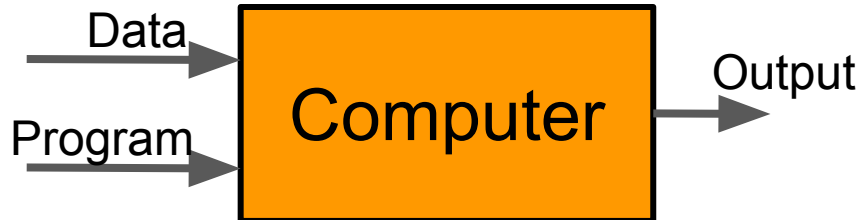
8,541 Retweets 23,778 Likes



# Python Programming

```
In [1]: a = 3  
b = 1  
q = 3*a + 2*b  
print('result is {}'.format(a + b))
```

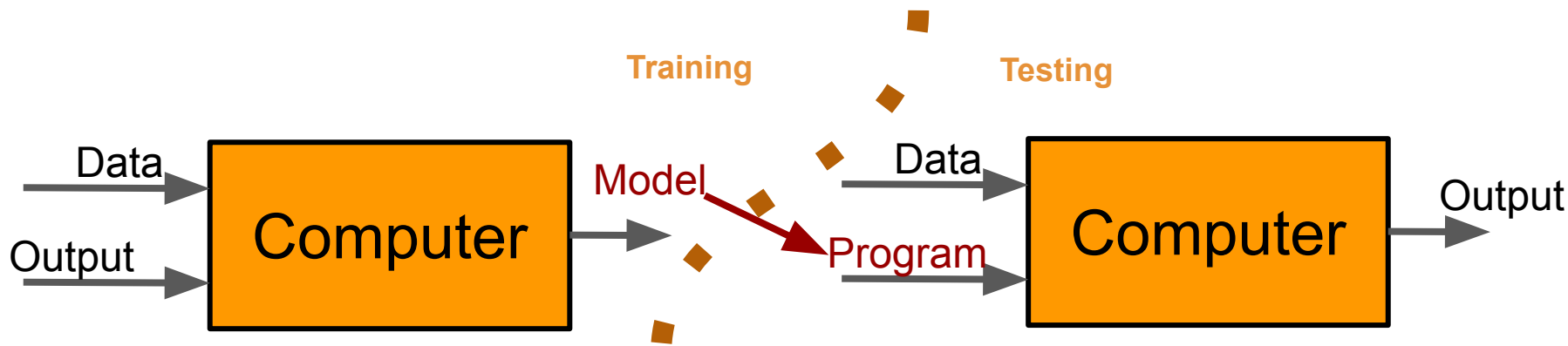
result is 4





# Machine Learning

```
] from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
#create an object of KNN
neigh = KNeighborsClassifier(n_neighbors=3)
#train the algorithm on training data and predict using the testing data
pred = neigh.fit(data_train, target_train).predict(data_test)
```



# Definition of Machine Learning

“A computer program is said to learn from **experience E** with respect to some class of **tasks T** and performance **measure P**, if its performance at tasks in **T**, as measured by **P**, improves with experience **E**”



**Tom Mitchell**

**T**, **P**, **E** are three basic elements to define a complete machine learning tasks

# AlphaGo

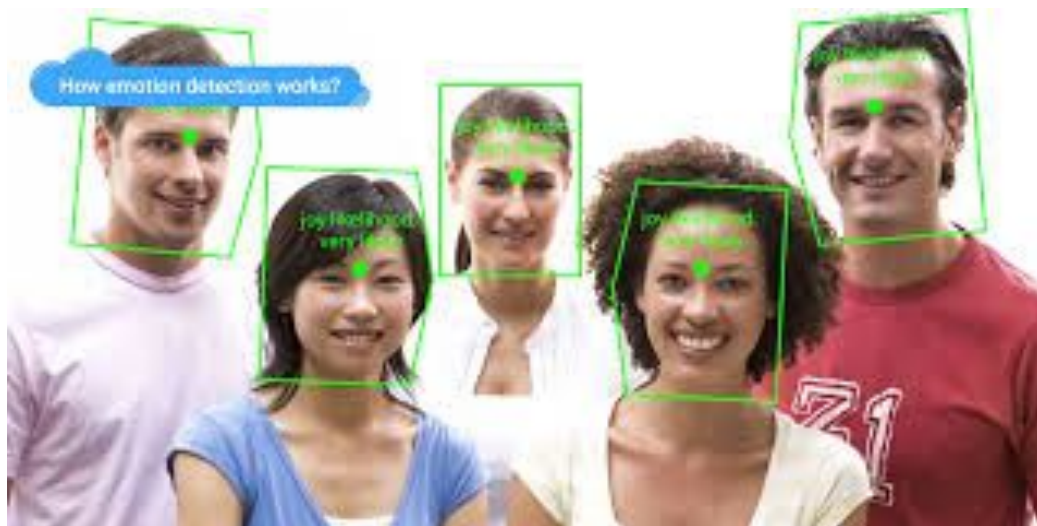


**T**: Play Go Games

**P**: Win rates of all matches

**E**: Match Experiences with many go players or itself

# Face Recognition



**T**: Identify or verify human faces

**P**: Accuracy that human faces are detected

**E**: Dataset of labelled human faces

# More E

- For machine learning algorithms, E is **data**.
- When **data is text (unstructured data type)**, we then have text mining.

*Text mining is not only limited to machine learning approaches, since we can also hand-craft rules (old days).*

# Text Mining

# What is Text Mining

- Is finding **interesting regularities** in large **textual** dataset.
  - Where **interesting** means non-trivial, hidden, previously unknown and potentially useful.
  - E.g., extract **relations** between all of the entities.
  - E.g., **NTU** is in **Singapore**.
- Is finding semantic and abstract information from the surface form of text data:
  - E.g., predict sentiment towards products
- The International Data Corporation estimated that approximately **80%** of the data in an organization is **text-based**.
- Text mining is also called **text analytics**.

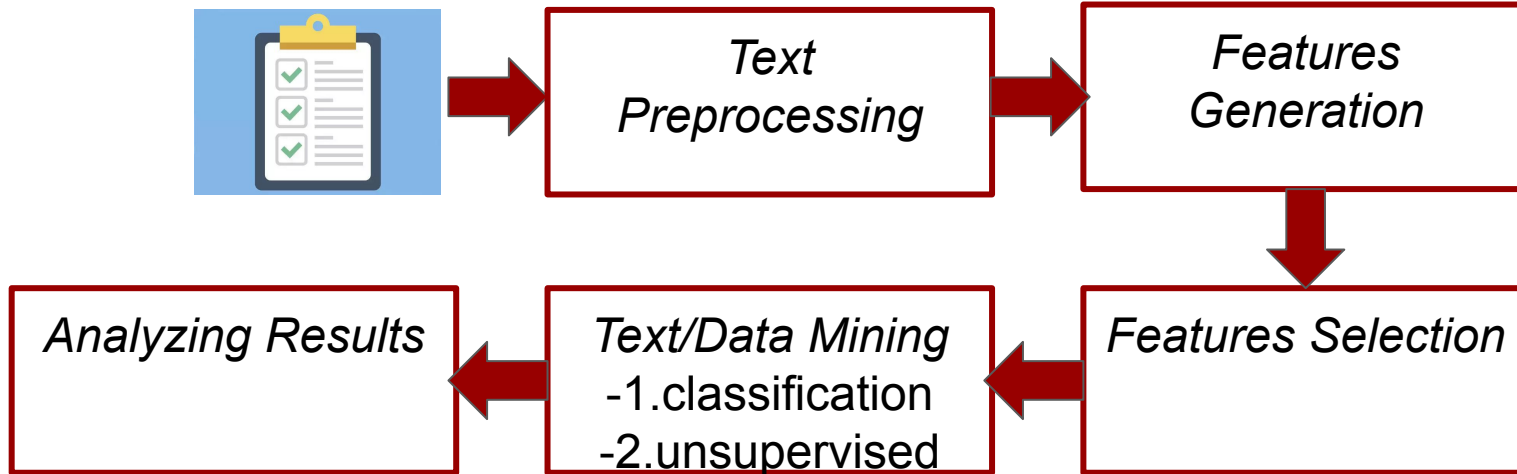
# Which Topics are related to Text Mining

- Data Mining
- Machine Learning
- **Natural Language Processing**
  - Computational Linguistics
- **Information Retrieval**
  - Search & full-text indexing
- **Knowledge Management**
  - Knowledge Representation and Reasoning
  - Used in Question & Answering Systems



# Text Mining Process Flow

- A typical text mining project involves 5 steps



# Unstructured Data: Text

# Structured Data

- Structured Data
  - Machine learning/predictive algorithms need fixed-length vectors as inputs
  - Structured data is easily to be handled/prepared by our humans
  - Can be represented by columns and rows.
  - Each row is a data sample. Each column is attribute/feature.
- A toy task: predict the position of the basketball player



# Structured Data for Toy Example

- Structured: just like the excel file or csv

| Player   | Height (inches) | Weight (pounds) | Position |
|----------|-----------------|-----------------|----------|
| Player 1 | 76              | 225             | C        |
| Player 2 | 75              | 195             | PG       |
| Player 3 | 72              | 180             | SF       |
| Player 4 | 82              | 231             | PF       |

**Features** (points to Height and Weight columns)

**Labels** (points to Position column)

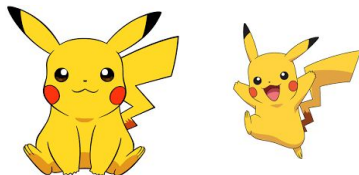
**Feature Values** (points to the numerical values in the Height and Weight columns)

**Data Sample** (points to the entire row for Player 4)

# Unstructured

- The original data can not be stored in an “table”
- More abstract, more fuzzy, and more high-dimensionality

**Images**



**Audio**



**Video**

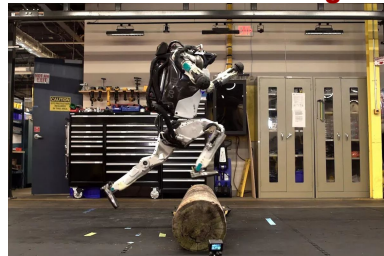


**Text**

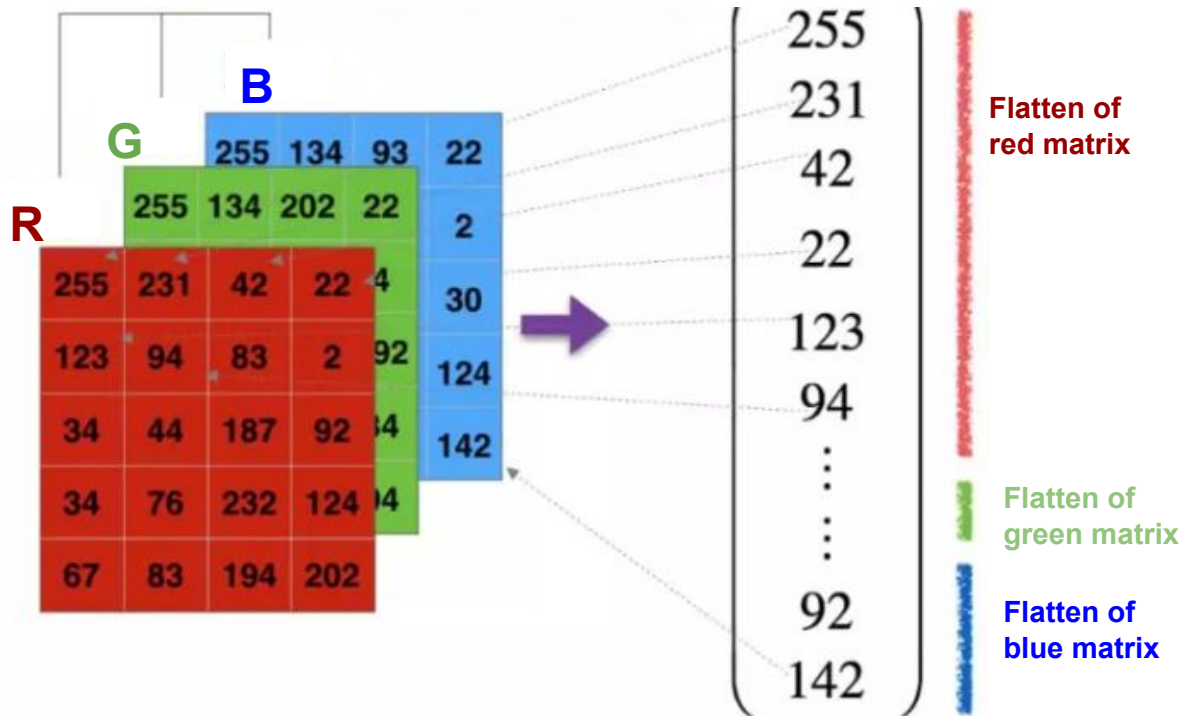
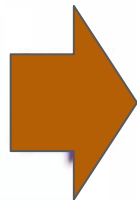
**Content**

This module provides students a deep overview of various advanced machine learning techniques applied to business analytics tasks. The focus of this course will be the key and intuitive idea behind machine learning models and hands-on examples instead of theoretical analysis. The tentative topics include machine learning pipeline, unsupervised learning, structure learning, Bayesian learning, deep learning and generative models. The programming languages used will be Python.

**Environment around agent**



# For Images



# For Text

- One of the main themes supporting text mining is **the transformation of text into numerical data**.
- Although the initial presentation is document format, the data move into a classical data-mining encoding (from unstructured to structured).
  - Each data is a vector
  - The length of the vector should be fixed
- Each row represents a document and each column a word.

|                          |
|--------------------------|
| The cat and the dog play |
| The cat is on the mat    |

*corpus*

|  |
|--|
| and, the, cat, dog,<br>play, on, mat, is |
|--|

*vocab.*

|   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
| 1 | 2 | 1 | 1 | 1 | 0 | 0 |
| 1 | 2 | 0 | 0 | 1 | 1 | 1 |

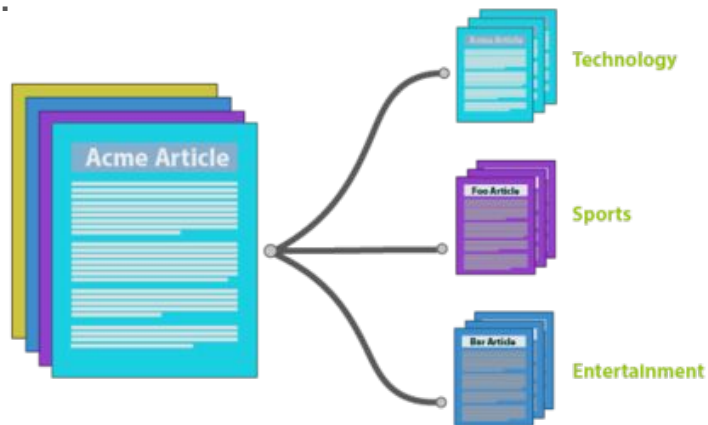
*countVec*

# Text Mining Applications



# Applications

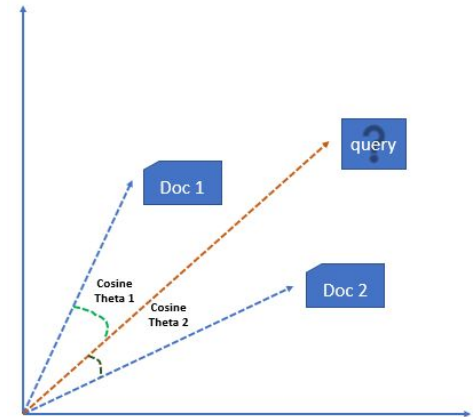
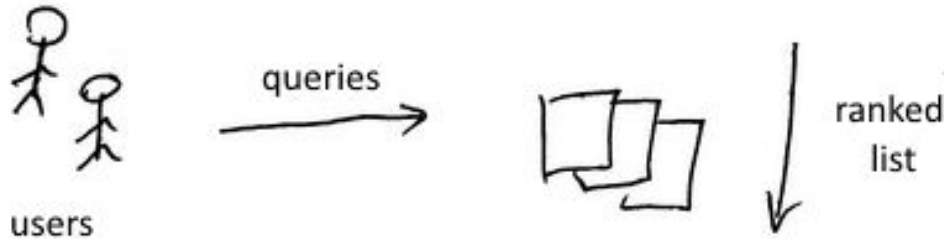
- **Document Classification:** given a sample of documents and correct answers (text categories) for each document, the objective is to find the correct answers for new documents.



- Assign topic into each document/piece of text
- Email spam detection (binary classification) or new topic categorization (multiple classification)

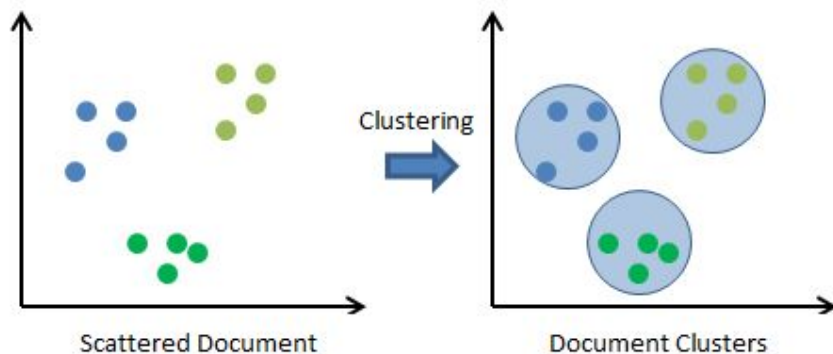
# Applications

- **Information Retrieval** is the science of searching for documents or information in documents.
  - The input document is matched to all documents, retrieving the best-matched documents.
  - A basic concept for IR is **measuring similarity**: a comparison is made between two documents, measuring how similar the documents are.
  - Similarity can be computed after documents have been encoded as vectors



# Applications

- **Document Clustering** is used when we have a collection of **documents with no known structure or no predefined categories**.
  - E.g., email complaints by users are clustered, and can learn about the categories and types of complaints.
- Because there are many ways to cluster documents, it is not quite as powerful as assigning answers(i.e., known correct labels) to documents.



- An example of Document Clustering: consider the comments made by the patients about the best thing they liked about the hospital.
- Because there are many ways to cluster documents, it is not quite as powerful as assigning answers(i.e., known correct labels) to documents.

1. *Friendliness of the doctor and staff*
2. *Service at the eye clinic was fast.*
3. *The doctor and other people were very, very friendly.*
4. *Waiting time has been excellent and staff has been very helpful.*
5. *The way the treatment was done.*
6. *No hassles in scheduling an appointment.*
7. *Speed of the service.*
8. *The way I was treated and my results.*
9. *No waiting time, results were returned fast, and great treatment.*

Table 1.2: Clustering Results from Text Mining

| Cluster No. | Comment | Key Words                          |
|-------------|---------|------------------------------------|
| 1           | 1, 3, 4 | doctor, staff, friendly, helpful   |
| 2           | 5, 6, 8 | treatment, results, time, schedule |
| 3           | 2, 7    | service, clinic, fast              |

# Applications

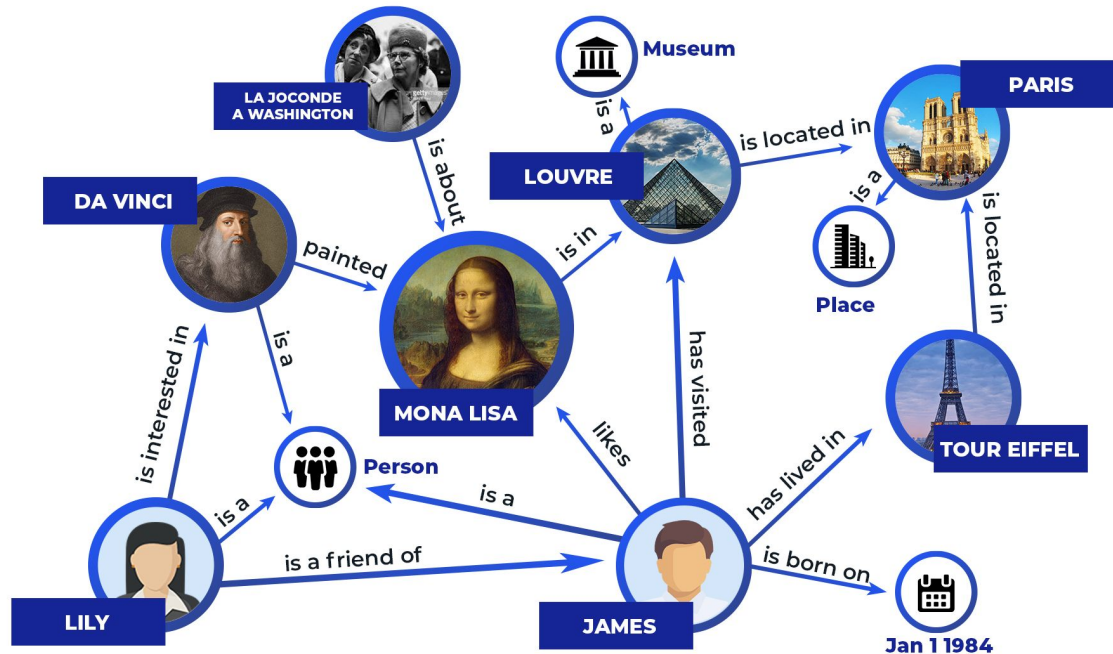
- **Text Summarization**

- Task: the task is to produce shorter, summary version of an original document.
- Two main approaches to the problem:
  - Extraction-based: output consists from topmost text units
  - Abstraction-based: perform semantic analysis, representing the meaning and generating the text satisfying length restriction.

# Applications

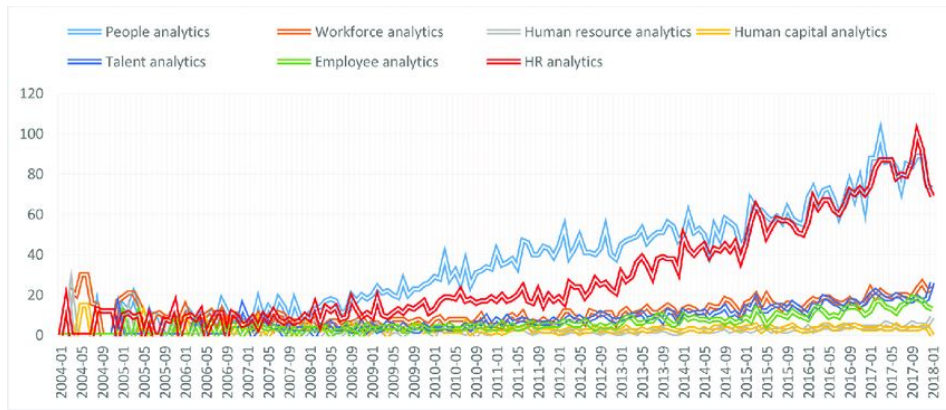
- **Knowledge Management**

- Knowledge Graph: nodes(entities) and edge (relationship between entities)



# Applications

- **Trend Analysis:** Given a set of documents with a time stamp, text mining can be used to identify trends of different topics that exist in the text.
- **Examples**
  - Tracking the trends in research from scientific literature
  - Summarizing events from news articles.
- **Google Trends** provides a facility to identify the trends in various topics over a period of time.
  - **Topic: Text Analytics**



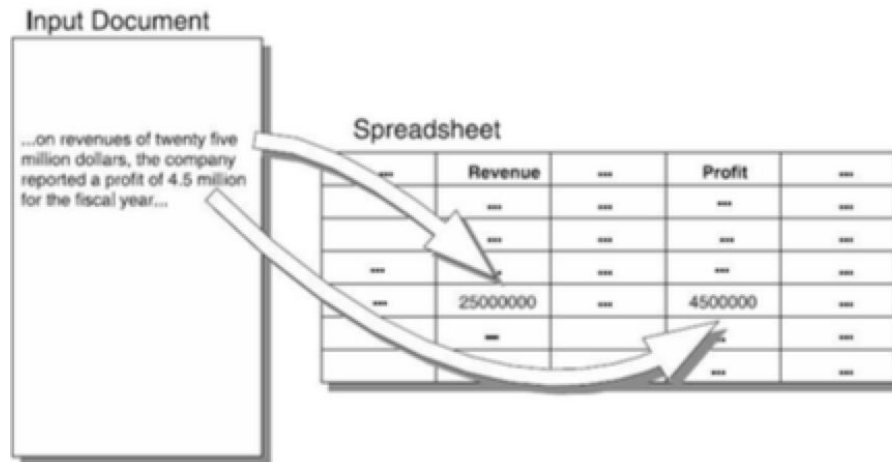
# Applications

- **Information Extraction**

- Take an unstructured document and automatically turn them into structured format
- In the structured format, the columns are not just words but higher-level concepts that are found by the information extraction process.
  - E.g., people, organization, places, addresses, dates.



Figure 1: An example of NER application on an example text

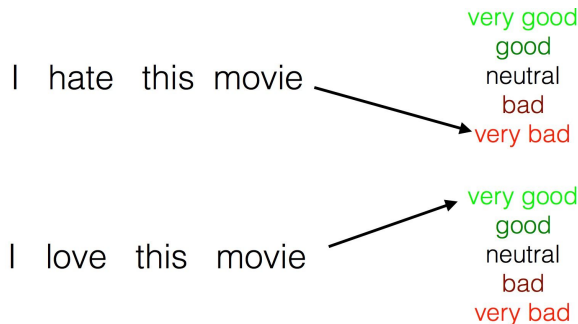




# Applications

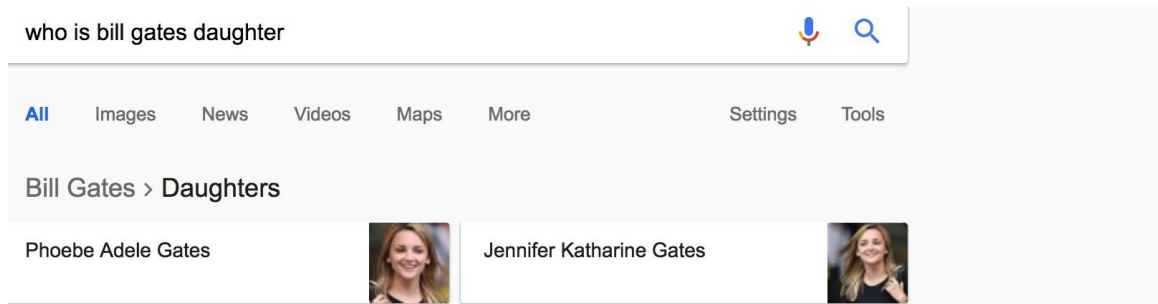
- **Sentiment Analysis**

- A type of subjective analysis which analyzes sentiment in a given textual unit with the objective of understanding the sentiment polarities (i.e. positive, negative, or neutral) of the opinions toward various aspects of a subject.
- It is also called as opinion mining.
- Importance of social media and online opinions
  - Online shoppers are influenced by product reviews and are willing to pay more for products highly rated by other consumers.
  - Users are more influenced by reviews of fellow consumers rather than those generated by professionals.



# Applications

- Question Answering



- Visual Question Answering

Is the umbrella upside down?  
yes                      no



How many children are in the bed?  
2                                      1



# Why Text Mining is Tough?

- Many ways to represent similar concepts
  - E.g., space ship, flying saucer, and UFO
- “Countless” combinations of subtle, abstract relationships among concepts
  - E.g., relationship between drugs and diseases
- High dimensionality
  - Tens of hundreds of thousands of features
- Data Variation
  - We have ImageNet, while we do not have such huge labelled volume text data
- Ambiguity of Language
  - Word level: bank
  - Sentence level: I heard his cell phone in my office

# Text mining/NLP is really hard

