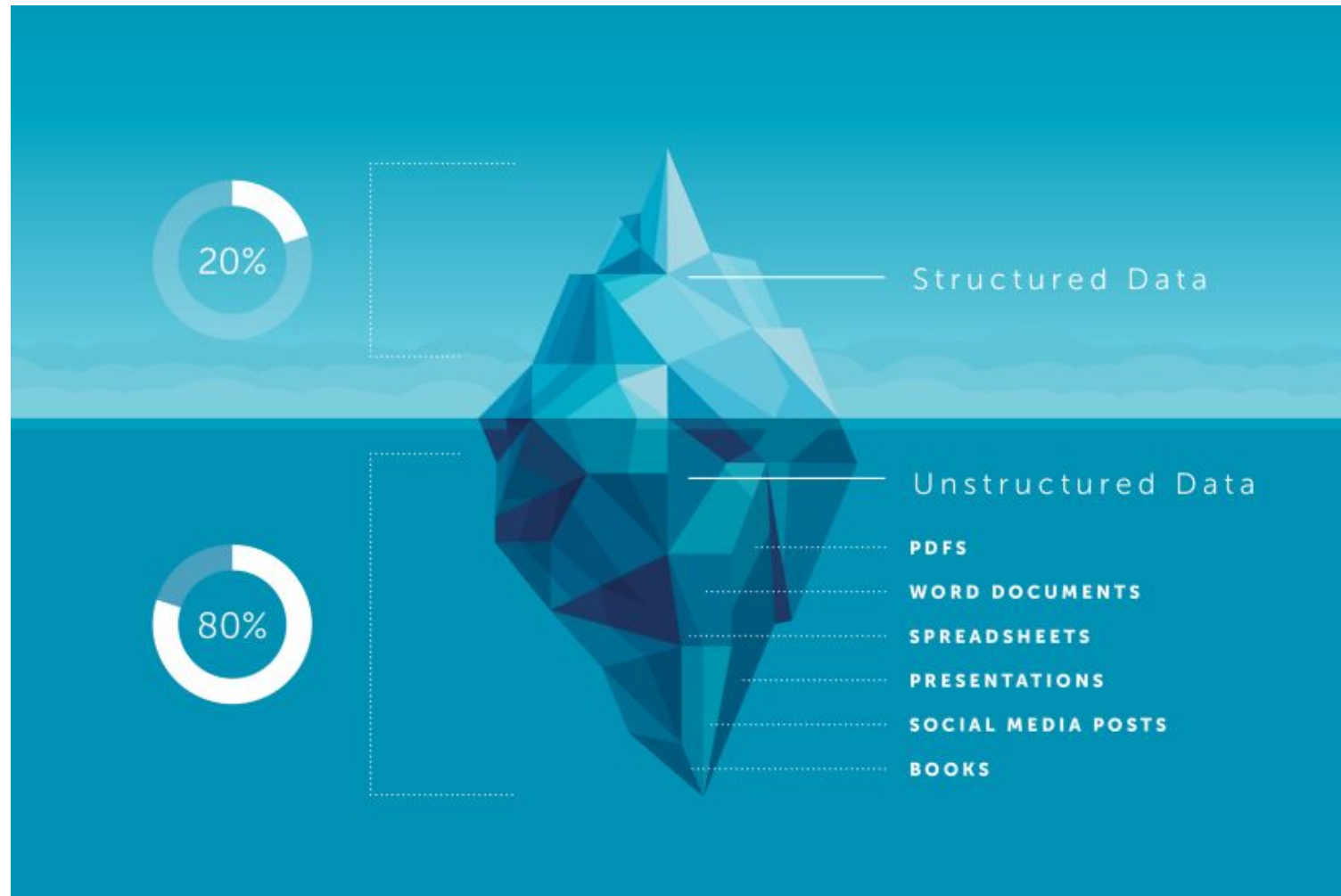# Modul 7: Text Mining

# Module Overview

**Topics**

- What is Text Mining

- Text Mining Process

- Text Classification

- Topic Modelling

**Activities**

- Coding Practice

- Orange

# Data In The World



20%

Structured Data

80%

Unstructured Data

PDFS
WORD DOCUMENTS
SPREADSHEETS
PRESENTATIONS
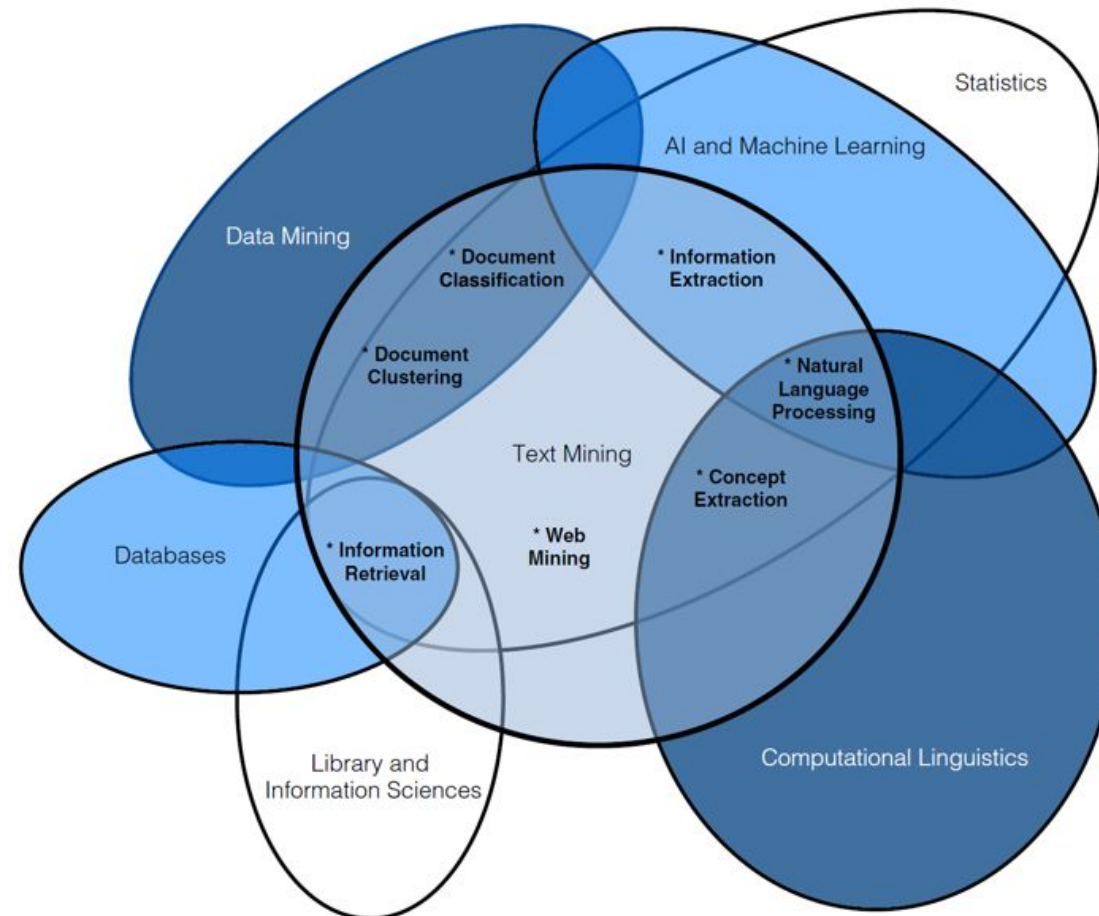SOCIAL MEDIA POSTS
BOOKS

# Text Mining

- Text mining is the process of extracting the implicit knowledge from textual data

- Because the implicit knowledge which is the output of text mining does not exist in the given storage, it should be distinguished from the information which is retrieved from the storage.
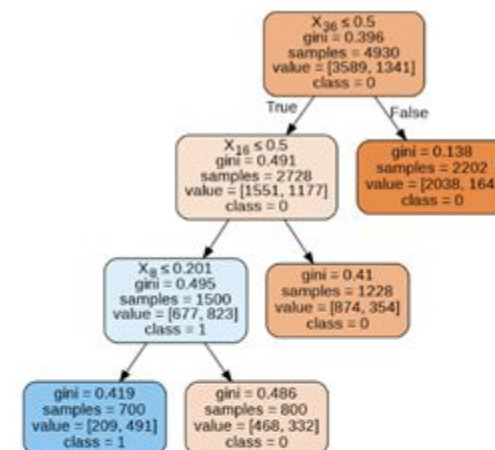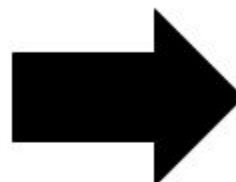
# Text Mining

# Text Mining Method

| Retrieval | Obtaining information resources relevant to an information need from a collection of information resources |
|---|---|
| Extraction | Extracting structured information from unstructured and/or semi-structured machine-readable documents. |
| Summarization | process of automatically generating a compressed version of a specific text that holds valuable information for the end-user. |
| Categorization | Assigned Document to a predefined set of topics depending upon their content. |
| Cluster | identify intrinsic structures in textual information and organize them into relevant subgroups or 'clusters' for further analysis. |

# Text Mining Application

| Manufacturers | Government | Financial Institutions |
|---|---|---|
| • Identify root causes of product issues quicker<br>• Identify trends in market segments<br>• Understand competitors' products | • Identify fraud<br>• Understand public sentiments about unmet needs<br>• Find emerging concerns that can shape policy | • Use contact center transcriptions understand customers<br>• Identify money laundering or other fraudulent situations |
| **Retail**<br>• Identify profitable customers and understand the reasons for their loyalty<br>• Manage the brand on social media | **Legal**<br>• Identify topics and keywords in discovery documents<br>• Find patterns in defendant's communications | **Healthcare**<br>• Find similar patterns in doctor's reports<br>• Use social media to detect disease outbreaks earlier<br>• Identify patterns in patient claims data |
| **Telecommunications**<br>• Prevent customer churn<br>• Suggest up-sell/cross-sell opportunities by understanding customer comments | **Life Sciences**<br>• Identify adverse events in medicines or vaccines<br>• Recommend appropriate research materials | **Insurance**<br>• Identify fraudulent claims<br>• Track competitive intelligence<br>• Manage the brand on social media |

zencos

Sumber: https://www.zencos.com/blog/text-mining-examples-advanced-analytics/

# My Previous Project

- Customer Satisfaction Analysis in E-commerce

| No. | Customer Chat | Sentiment Class |
|-----|---------------|-----------------|
| 1. | I'm having a look now Thank you so much for your help You were amazing | Positive |
| 2. | Will you be getting a sample pack of this design in? | Neutra |
| 3. | This is not addressing again at all how bad of experience i have had the delays in my printing and approval the service horrific | Negative |



Neutral 67%
Negative 8%
Positive 25%

# My Previous Project

- Customer Satisfaction Analysis in E-commerce

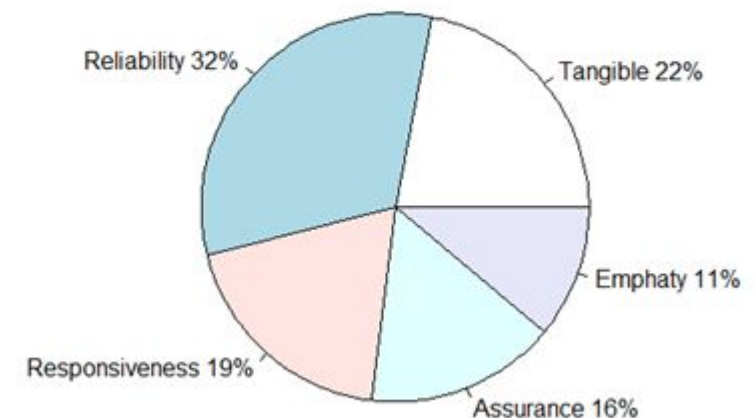| No. | Negative Customer Chat | Service Quality Dimension |
|---|---|---|
| 1. | Flow on this web is confusing, I can't find what I'm looking for | Tangible |
| 2. | I tried to contact Jamie to discuss this design problem, but he gave no response. | Responsiveness |
| 3. | This is not addressing again at all how bad of experience I have had the delays in my printing and approval the service horrific | Reliability |

**Service Quality Dimention - Negative**

Reliability 32%
Tangible 22%
Emphaty 11%
Assurance 16%
Responsiveness 19%

# My Previous Project

- Job Opportunities for Business School Graduates

# My Previous Project

- Twitter COVID-19 Detection

# Text Mining Process

Text Data Collection → Text Preprocessing → Analyzing → Visualizing → Interpretation

# Text Data Collection

Some examples of text data sources:

- Social Media (Twitter, Facebook, WhatsApp, Etc. )

- Website

- News

- Review

- Customer Chat

- Book

- Etc.

# Text Preprocessing

- Text Preprocessing is a process of transforms text into a more digestible form so that machine learning algorithms can perform better.

- One of the Important Task

- Why? **Garbage in Garbage Out**

# Common Task in Text Preprocessing

- Text Cleaning

- Feature Extraction

# Text Cleaning



siam 🔞 @mimiamiamia95 · 3m
O gitu caranya banya followers bikin GA ipon, pas udah banyak yang pollow GAnya gajadi diumumin deh atau bilangnya udah ada yang menang hhhhhh Bisa saja **klean** ini ngibulnya 😊😊

| Text Cleaning | Result |
| --- | --- |
| **Tokenization** | O gitu caranya banya followers bikin GA ipon, pas sudah banyak yang pollow GAnya gajadi diumumin deh atau bilangnya sudah ada yang menang hhhhhh Bisa saja klean ini ngibulnya |
| **Slang word** | O gitu caranya banyak followers bikin GA iphone, pas sudah banyak yang follow GAnya tidak jadi diumumin deh atau bilangnya sudah ada yang menang hhhhhh Bisa saja kalian ini ngibulnya |
| **Stemming** | O gitu cara banyak followers bikin GA iphone, pas sudah banyak yang follow GA tidak jadi umum deh atau bilang sudah ada yang menang hhhhhh Bisa saja kalian ini ngibul |
| **Lemmatization** | O gitu cara banyak followers buat GA iphone, pas sudah banyak yang follow GA tidak jadi umum deh atau sebut sudah ada yang menang hhhhhh Bisa saja kalian ini tipu |
| **Stop word** | cara banyak followers GA iphone, banyak follow GA tidak umum sebut menang kalian tipu. |

# Feature Extraction

- Machine Learning only process numeric input



"Suka banget pake OK-JEK, drivernya ramah dan lucu"

Input is the form of sentence (Text)

ML-based Classifier

Corpus

Positive

# Numerical Representation of Text

D = Indonesia adalah negara terindah di dunia.

**Tokenize**

("Indonesia", "adalah", "negara", "terindah" , "di", "dunia")

$W1$   $W2$   ...   $Wn$

$X1$   $X2$   ...   $Xn$

$W = Word\ (Text)$

$X = Some\ numeric\ encoding\ of\ Word$

D = X1, X2, X3, ... Xn

# Feature Extraction

**Popular Technique:**

- One-Hot Encoding

- Frequency-Based: (TF-IDF, Co-occurrence)

- Prediction-Based: (Lda2Vec)

# One Hot Encoding

| Tweet |
| --- |
| Direkomendasikan banget nih |
| Drivernya bau |
| Ramah banget |
| Gak bakal naik lagi |

| All Word |
| --- |
| Direkomendasikan |
| Banget |
| Nih |
| Drivernya |
| Bau |
| Ramah |
| Gak |
| Bakal |
| Naik |
| Lagi |

# One Hot Encoding

| All word | Direkomendasikan banget nih | Drivernya bau | Ramah banget |
|---|---|---|---|
| direkomendasikan | 1 | 0 | 0 |
| banget | 1 | 0 | 1 |
| nih | 1 | 0 | 0 |
| drivernya | 0 | 1 | 0 |
| bau | 0 | 1 | 0 |
| ramah | 0 | 0 | 1 |
| gak | 0 | 0 | 0 |
| bakal | 0 | 0 | 0 |
| naik | 0 | 0 | 0 |
| lagi | 0 | 0 | 0 |

# One Hot Encoding

**Advantages**:

- Simple

**Disadvantages:**

- Large vocab – Enormous Feature Vector.

- Unordered- Lost All Context.

- Lost Freq. Distribution (Binary).

- Do not capture any semantic information.

# Frequency-based Embeddings

**Count:** Capture how often a word occurs in a document

**TF-IDF:** Capture how often a word occurs in a document as well as the entire corpus

**Co-occurrence:** Similar word will occur together and will have similar context

# Count

Count, capture how often a word occurs in a document. Ex Count or the Frequency

**Advantages**:

- Simple

**Disadvantages:**

- Large vocab – Enormous Feature Vector.

- Lost All Context.

- Do not capture any semantic information.

# TF-IDF

TF-IDF: Capture how often a word occurs in a document as well as the entire corpus

$$Xi = tf(Wi) \times idf(wi)$$

*tf = (Number of repetitions of word in a document) / (# of words in a document)*

*idf = Log[(# Number of documents) / (Number of documents containing the word)] and*

Frequently in a single document
**Might be important**

Frequently in the corpus
**Probably a common word**

# TF-IDF

**Advantages**:

- Capture Frequency and Relevance
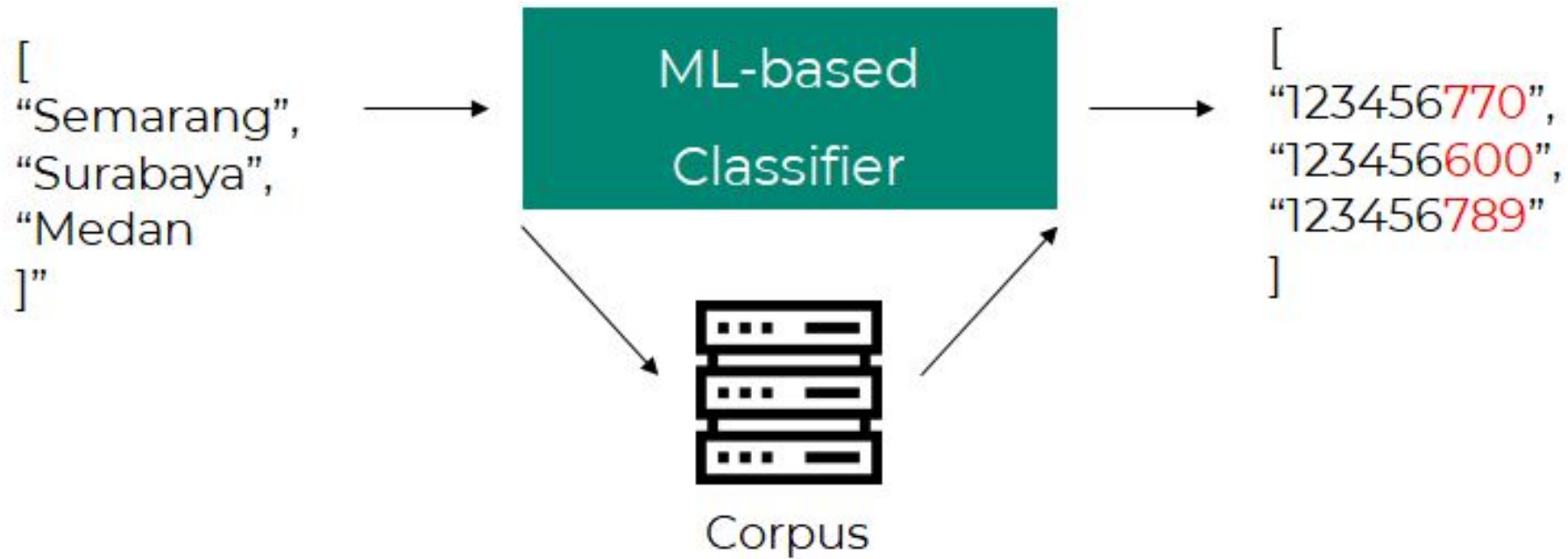
**Disadvantages:**

- Context still not captured.

# Co-occurrence

- **Co-occurrence:** The number of times two words w1 and w2 have occurred together in a context window

- **Context Window:** A window centered around a word, which includes a certain number of neighboring words.

*Similar word will occur together and will have similar context*

# Prediction-Based Embedding



[
"Semarang",
"Surabaya",
"Medan
]"

ML-based Classifier

Corpus

[
"123456770",
"123456600",
"123456789"
]

- **Text classification** also known as text tagging or text categorization is the process of categorizing text into organized groups. Some of the most common examples and use cases for automatic text classification is Sentiment Analysis.
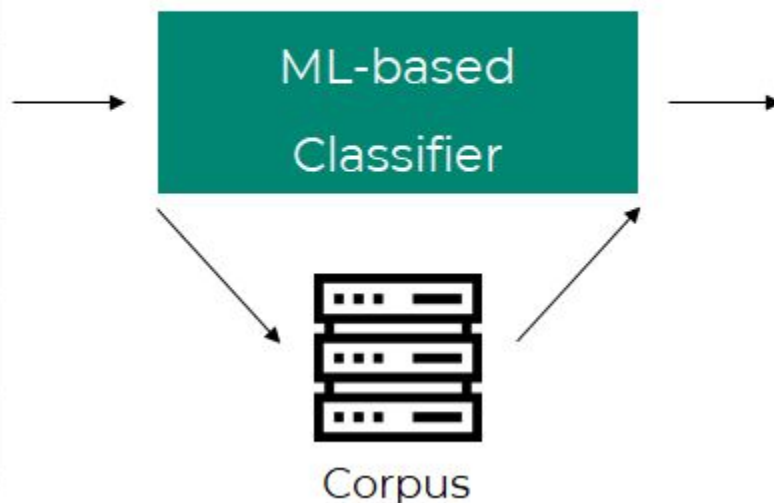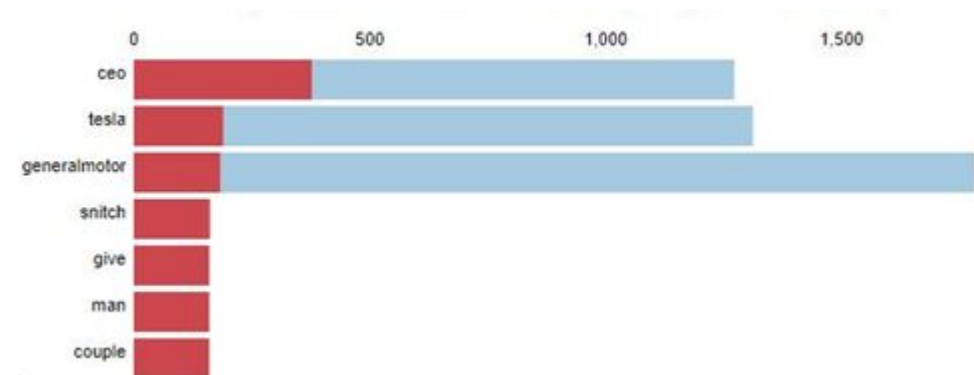
# Sentiment Analysis

| Tweets | Sentiment |
|---|---|
| kartu yg sdh sukses diregistrasi ttp di reject saat bayar via grabpay?mslh baru timbul semingguan terakhir | Negative |
| halo Grab kenapa ya lokasi ko ga kedetect dari tadi? sy mo order ga bisa | Negative |
| pengemudi payah..ambil penumpang dr bandara sll cancel! | Negative |
| apa pembatalan sama dengan order yang tidak di ambil? | Negative |
| klo mau ganti no rek cimb kemana ya? troubleshootnya lambat menanggapinya | Negative |
| SECEPAT datangnya GRAB\ud83d\02 Buruan pake !! Radiovenusmks | Neutral |
| min mau top up grabpay di alfamart sistemnya eror tuh | Negative |

**ML-based Classifier**

**Corpus**

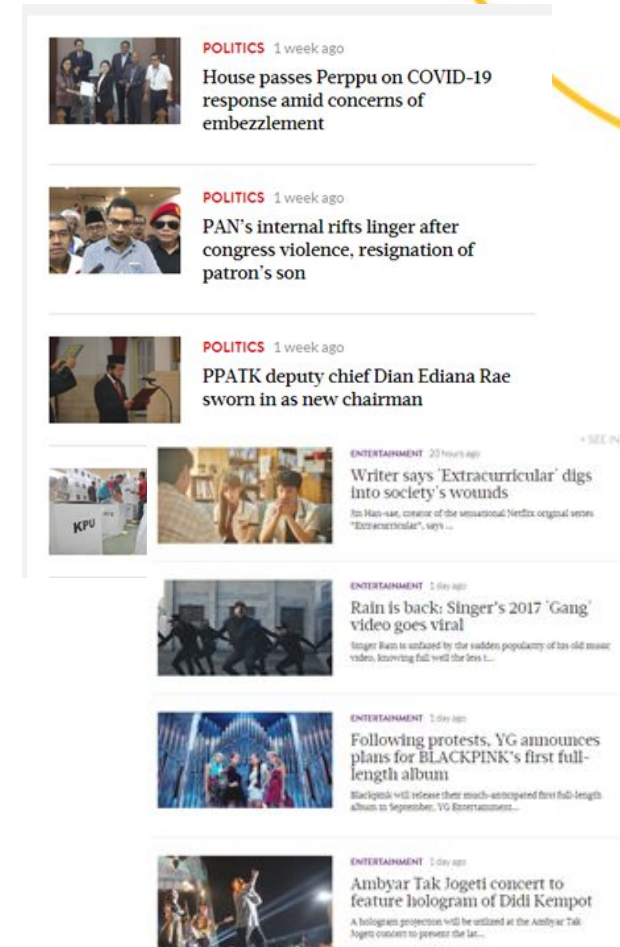| Tweets | Sentiment | Predicted Sentiment |
|---|---|---|
| tidak bisa dipakai sama sekali saldonya. Mohon di kroscek ya | | Negative |
| aplikasi grab kok gbsa ya? Keluar sendiri pas d klik | | Negative |
| halo kenapa applikasi grab saya selalu ketutup sendiri stlh order grabfood ya? Saya jd tdk bisa hub drivernya | | Negative |
| I can't access Grabchat. May I know why? | | Negative |
| kalo pesen grabfood..kdg di app open..tp aslinya close. driver nya suka kayak bete gt ditelfon seolah olah kt yg salah ya haha | | Negative |
| halo dari kemarin kalo jam segini mau pesen grab selalu \failed to contact the Grab server\ padahal kalo pagi ga ada masalah | | Negative |

# Analyzing: Topic Modelling

- Topic model is a type of statistical model for discovering the abstract "topics" that occur in a collection of documents

- Popular technique:

    - Latent Semantic Analysis (LSA)

    - Probabilistic Latent Semantic Analysis (PLSA)

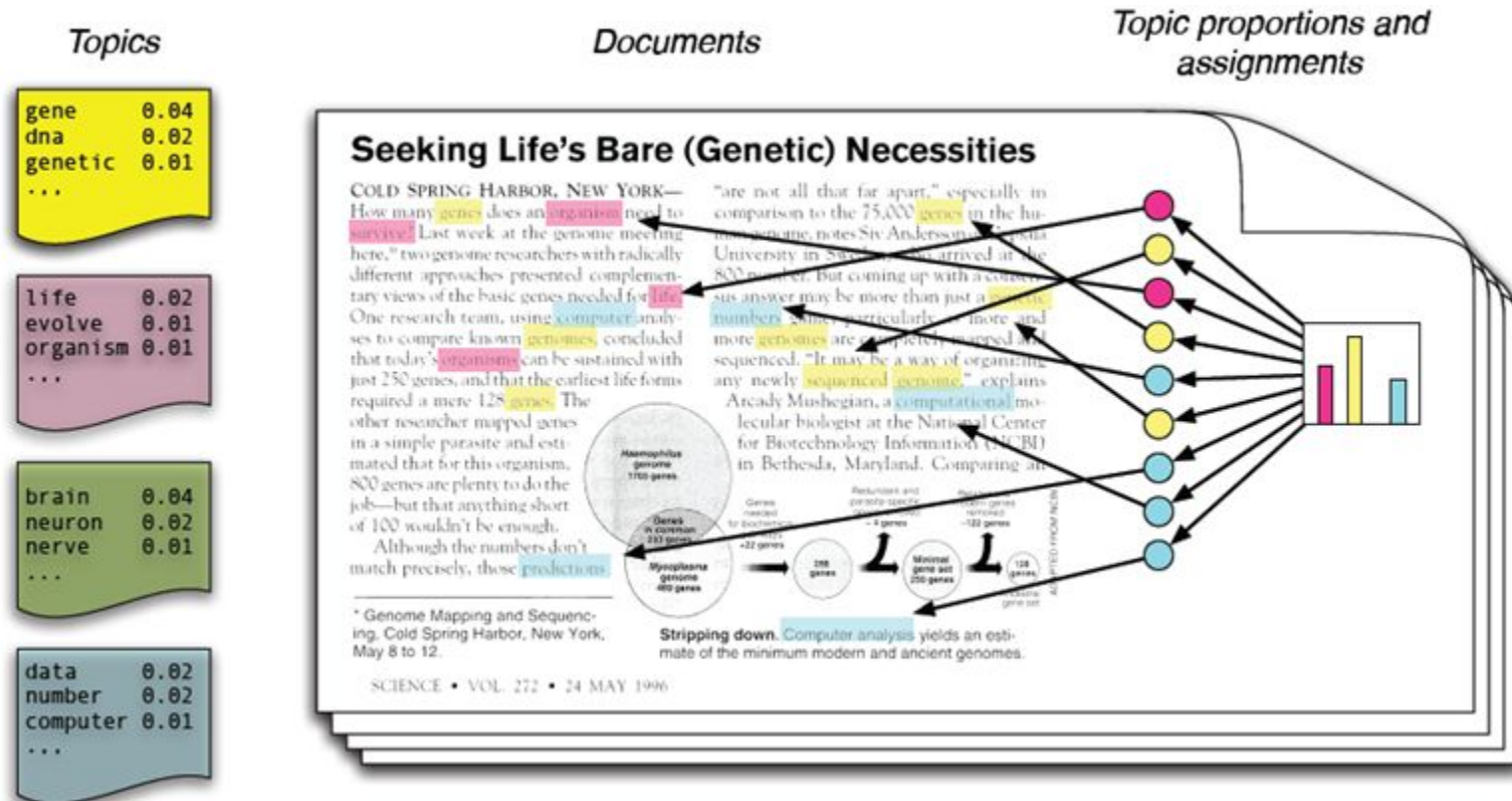    - Latent Dirichlet Allocation (LDA)

    - Lda2vec

# Latent Dirichlet Allocation (LDA)

- **Every document is a mixture of topics.**

  - For example, in a two-topic model we could say "Document 1 is 90% topic A and 10% topic B, while Document 2 is 30% topic A and 70% topic B."

- **Every topic is a mixture of words.**

  - For example, we could imagine a two-topic model of Indonesian news, with one topic for "politics" and one for "entertainment."

  - The most common words in the politics topic might be "President", "Congress", and "government", while the entertainment topic may be made up of words such as "movies", "television", and "actor".

  - Importantly, words can be shared between topics; a word like "budget" might appear in both equally.

Sumber: Probabilistic Topic Models By David M. Blei

# Assignment

- Build a model for detecting hate speech using dataset

  https://github.com/ialfina/id-hatespeech-detection

- Use Google Collab or Jupyter Notebook

- Post it on Github