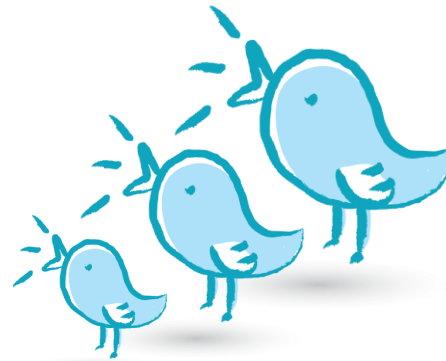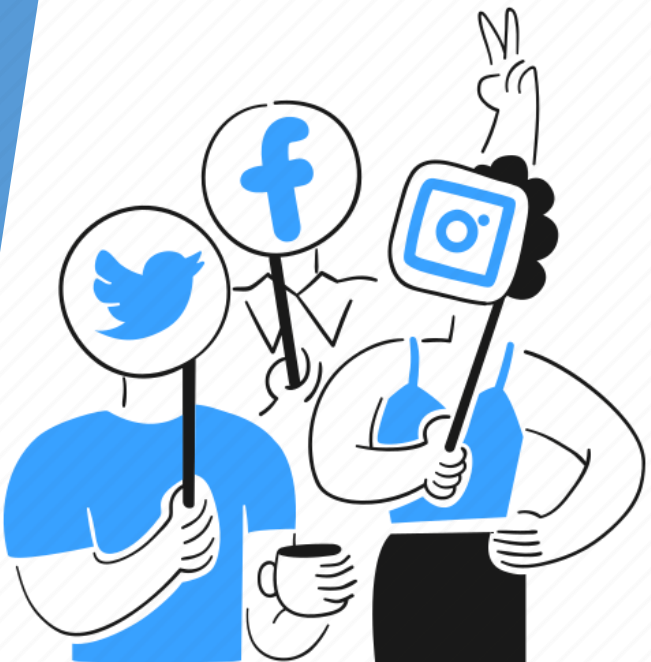# Building Predictive Models for Social Media Data

## Dr. Geetha Raju

**ML Engineer**

**Tamil Nadu e-Governance Agency, Chennai**

# Social Media and Data Analytics

- Social Media
  - Reinforce social bonds; manage social identities
  - Easy and open access
  - Social networking (Facebook, LinkedIn, Google+)
  - Microblogging (Twitter, Tumblr)
  - Photo sharing (Instagram, Snapchat, Pinterest)
  - Video sharing (YouTube, Facebook Live, Periscope, Vimeo)
- Practitioners, researchers and analyst – rich resources – social media data
- What they do?
  - Gather data
  - Find meaning / context
  - Derive insights that support decision making
  - Analyze / predict performance
- What is data in SM?
  - Post specific data - likes, reactions, comments, clicks, previews, etc.,
  - User specific data – name, DOB, followers, friends, etc.,
  - Network specific data – followers, following, friends, community / group, etc.,

What's Happening ?!?

# REASONS FOR USING THE INTERNET

PRIMARY REASONS WHY GLOBAL INTERNET USERS AGED 16 TO 64 USE THE INTERNET

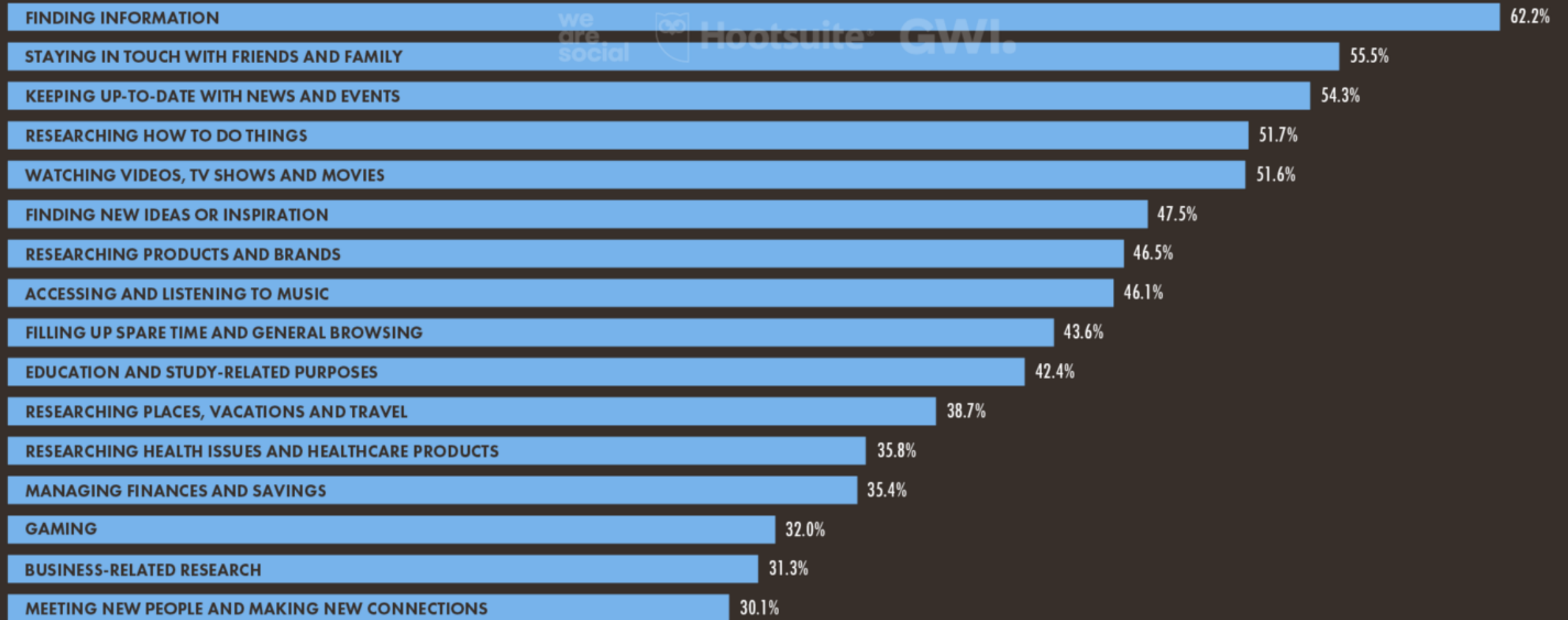| Reason | Percentage |
|---|---|
| FINDING INFORMATION | 62.2% |
| STAYING IN TOUCH WITH FRIENDS AND FAMILY | 55.5% |
| KEEPING UP-TO-DATE WITH NEWS AND EVENTS | 54.3% |
| RESEARCHING HOW TO DO THINGS | 51.7% |
| WATCHING VIDEOS, TV SHOWS AND MOVIES | 51.6% |
| FINDING NEW IDEAS OR INSPIRATION | 47.5% |
| RESEARCHING PRODUCTS AND BRANDS | 46.5% |
| ACCESSING AND LISTENING TO MUSIC | 46.1% |
| FILLING UP SPARE TIME AND GENERAL BROWSING | 43.6% |
| EDUCATION AND STUDY-RELATED PURPOSES | 42.4% |
| RESEARCHING PLACES, VACATIONS AND TRAVEL | 38.7% |
| RESEARCHING HEALTH ISSUES AND HEALTHCARE PRODUCTS | 35.8% |
| MANAGING FINANCES AND SAVINGS | 35.4% |
| GAMING | 32.0% |
| BUSINESS-RELATED RESEARCH | 31.3% |
| MEETING NEW PEOPLE AND MAKING NEW CONNECTIONS | 30.1% |

we are social

Hootsuite®

# A Minute on the Internet in 2021

Estimated amount of data created
on the internet in one minute

**NETFLIX** — 28,000 subscribers watching

2m views

695,000 stories shared

1.6m USD spent online

9,132 connections made

2m Swipes

69m messages sent

197.6m Emails sent

5,000 downloads

500 hours of content uploaded

60 Sec

Source: Lori Lewis via AllAccess

statista

# Twitter Statistics 2022

Twitter reached
**211 million**
daily active users in the third quarter of 2021

**79%**
of marketers will continue
investing in Twitter Spaces in 2022
according to HubSpot

Twitter stands
**15th**
in the list of the world's most
'active' social media platforms

**42% of all Twitter**
users have graduated college

**83% of all**
the world's leaders are on Twitter

**26% of US users**
check their Twitter account
several times a day

**42% of all registered**
Twitter users visit the platform daily

WHAT'S NEW

# Social Media and Data Analytics

▶ Challenges in social media data
  ▶ Time Sensitivity
  ▶ Short length
  ▶ Unstructured form

▶ More than 7 million web pages of text are being added to our collective repository, daily

▶ Processing speed
  ▶ 15,000- 250,000 pages an hour – TM software
  ▶ 60 pages for humans

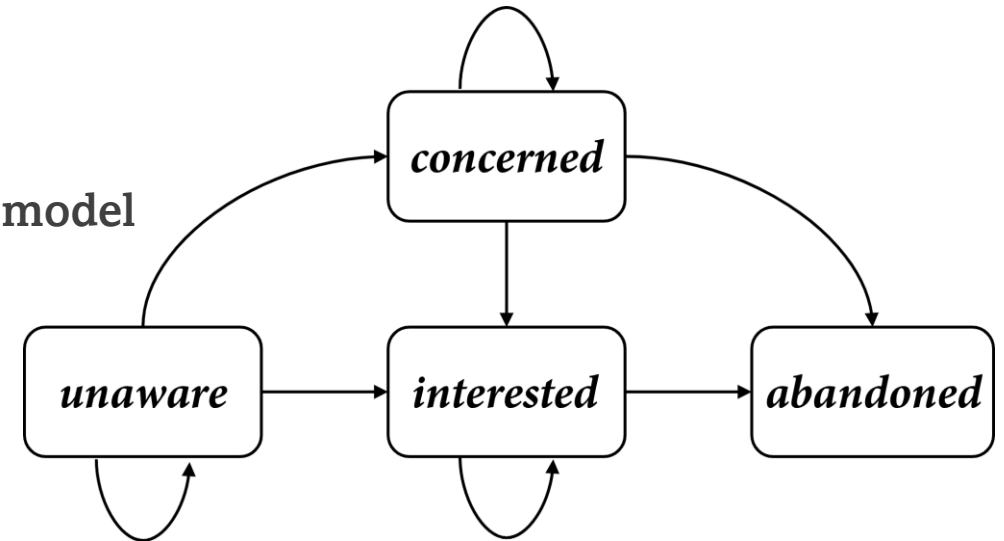# Real Life Vs Online Social Media



FAMILY TREE

SOCIAL MEDIA

# Information Spread and Flow in Social Media

▶ Information sharing

  ▶ Why - Pleasure and helps in gaining public attention

  ▶ What - Profile, Status, Location and shred content

  ▶ How - Noisy and unstructured

▶ Online social network information spreading (OSIS) model

▶ Information amplification

▶ Influenced by psychological and social factors

▶ False news

▶ How to analyze information flows?

  ▶ Subgraph constructions

  ▶ Activity and degree distributions

  ▶ Network analysis

# What ML / AL can do?



| Category | Types of Analytics | Questions Answered |
|---|---|---|
| Prescriptive | ▸ Optimization<br>▸ Randomized Testing | ▸ What is the best that can happen?<br>▸ What happens if we try this? |
| Predictive | ▸ Predictive modeling / forecasting<br>▸ Statistical modeling | ▸ What will happen next?<br>▸ What is making this happen? |
| Diagnostic | ▸ Data exploration<br>▸ Intuitive visuals | ▸ Why did this happen?<br>▸ What insights can I gain? |
| Descriptive | ▸ Alerts<br>▸ Query / drill down<br>▸ Ad Hoc reports / scorecards<br>▸ Standard reports | ▸ What actions are needed?<br>▸ What is the problem?<br>▸ How many, often, where?<br>▸ What happened? |

Capability →

# Scope of Predictive analytics on Social Media Data

- Predictive systems
  - Stock Market Predictions
  - Spam Detection
- Content based Recommender Systems
  - News, Movie, Product reviews / suggestions.
- Social media and Text Analytics
  - Sentiment Analysis – Product / Brand / Topic
  - Topic / Keyword / Phrase Identification
  - Cybercrimes and Cyberbullying
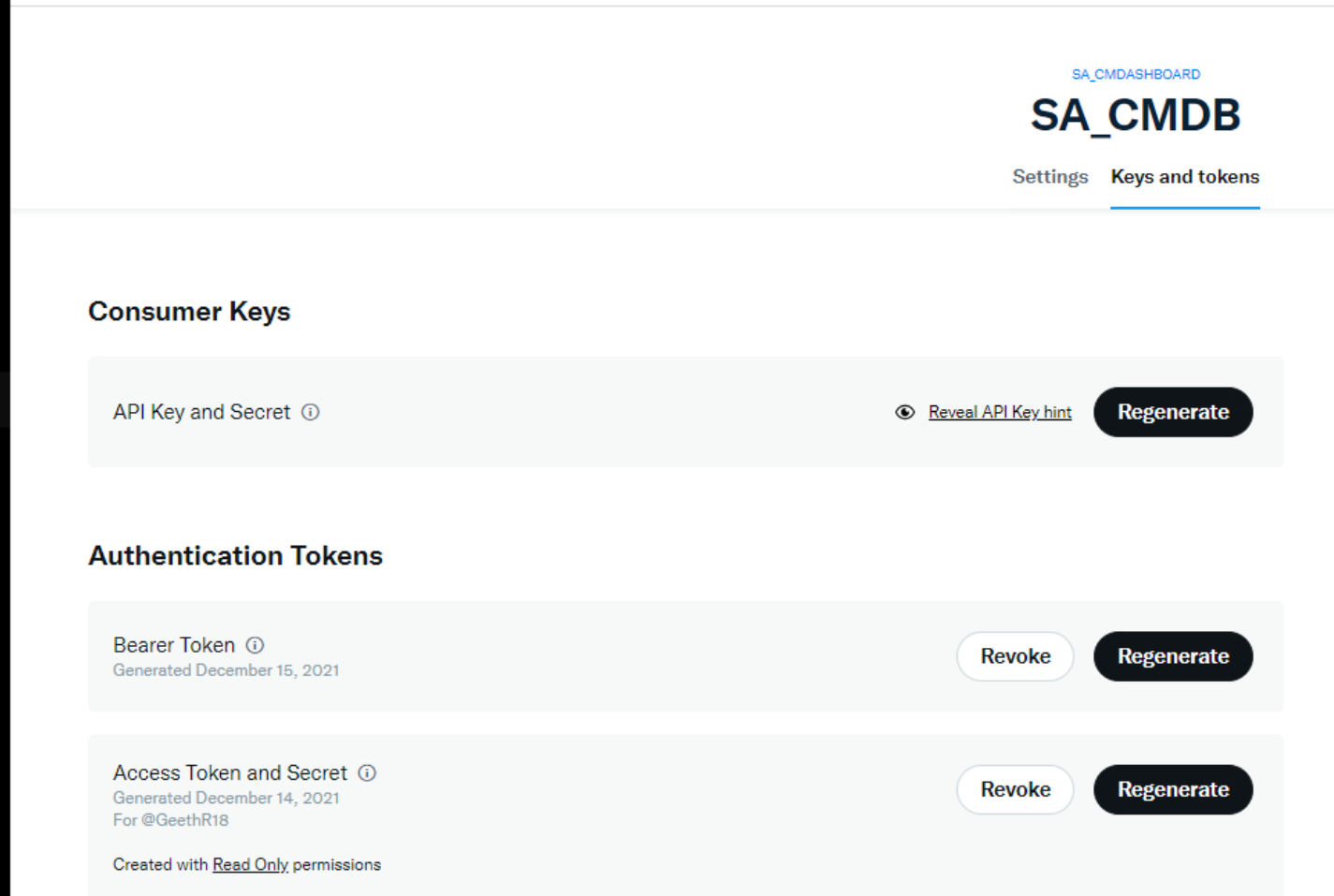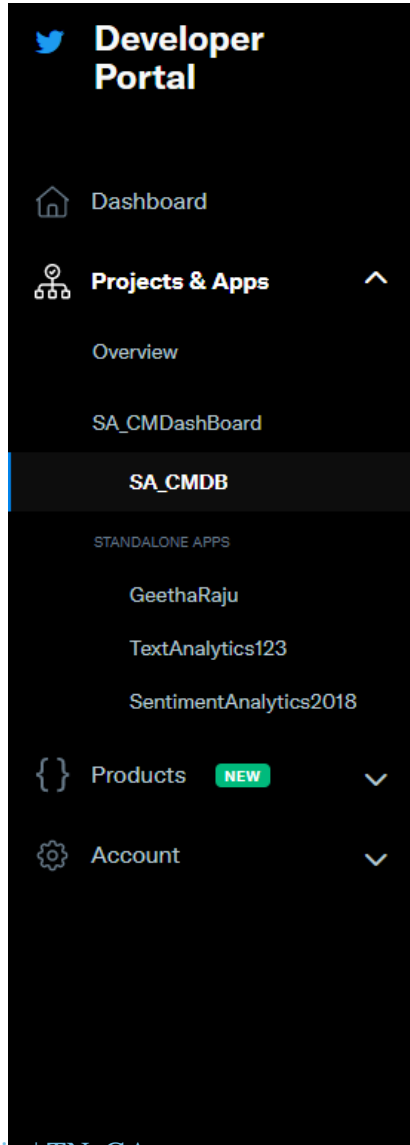- Linguistic Rules and Machine Learning Analysis

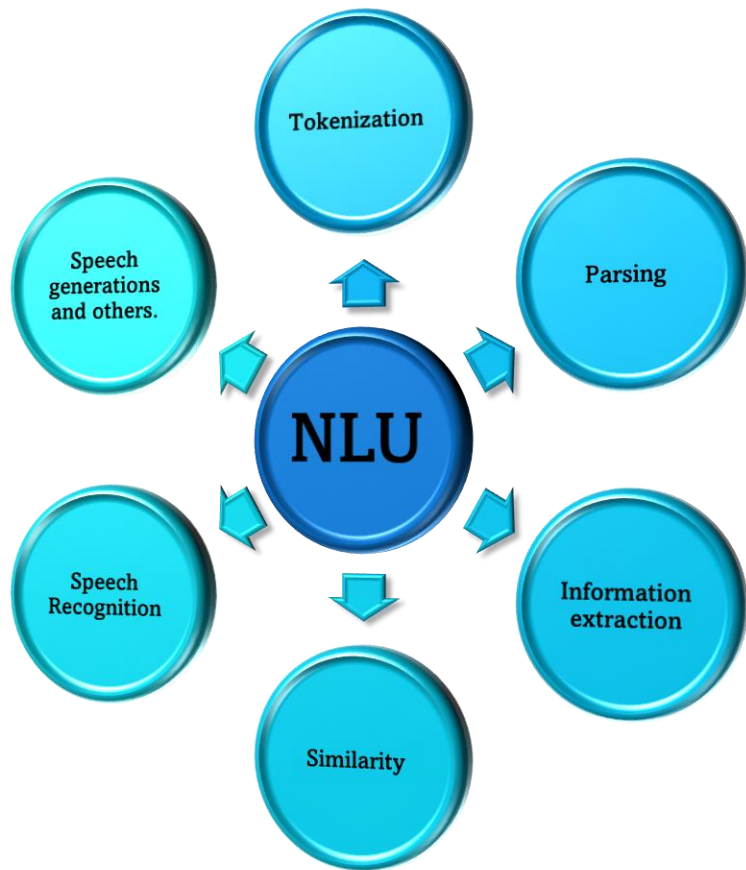# So, Lets Start with Data Preparation

# Data Collection Pipeline – Twitter Developer Account

# Text Data Features
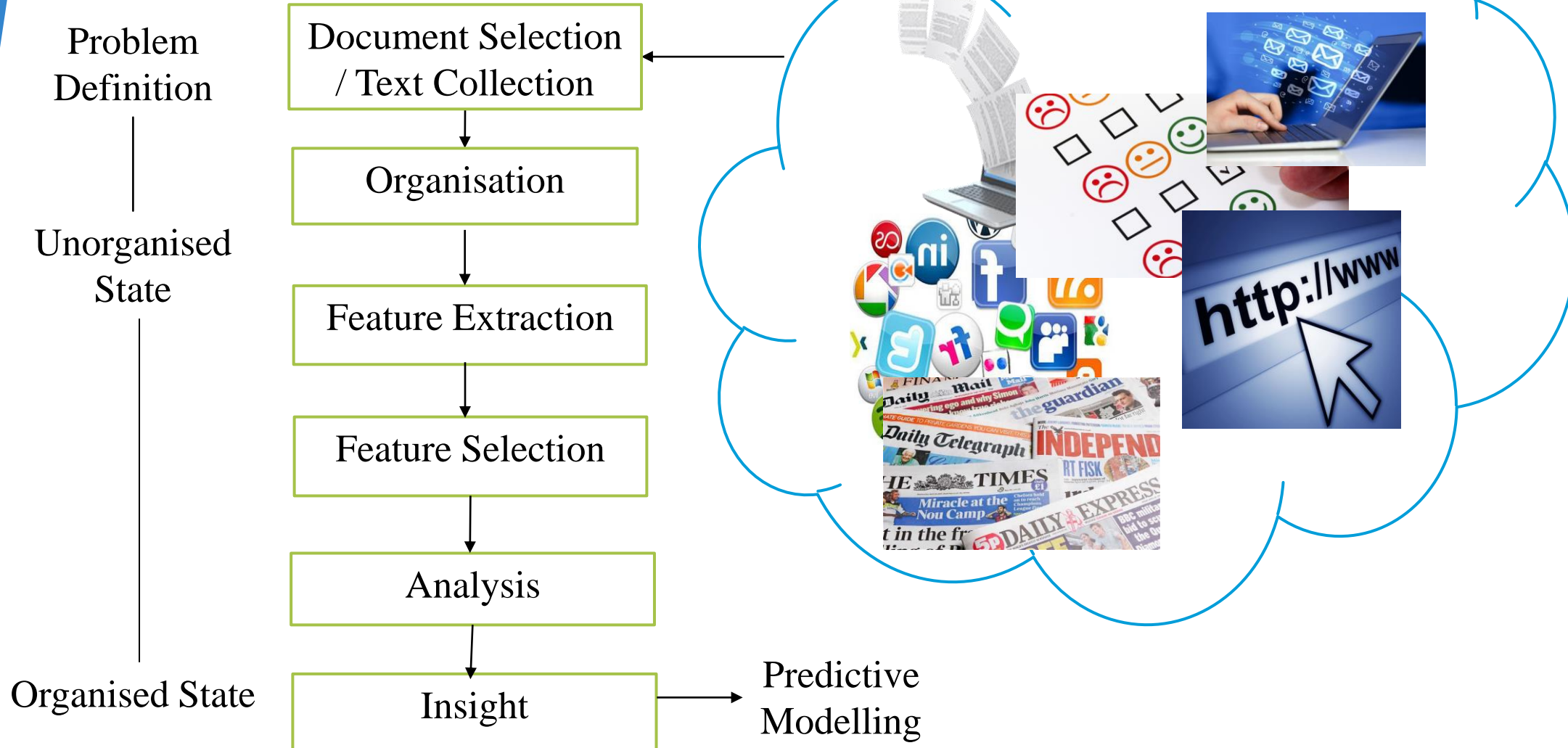
▶ Generate structure where there is no structure

▶ Challenges

    ▶ unstructured textual form (80% - 90% of world's data)

    ▶ large text data

    ▶ high dimension but sparse

    ▶ word / phrase types in various languages

    ▶ complex relationship between concepts in text

    ▶ word ambiguity and context sensitivity

    ▶ noisy data

# NLU Vs NLP

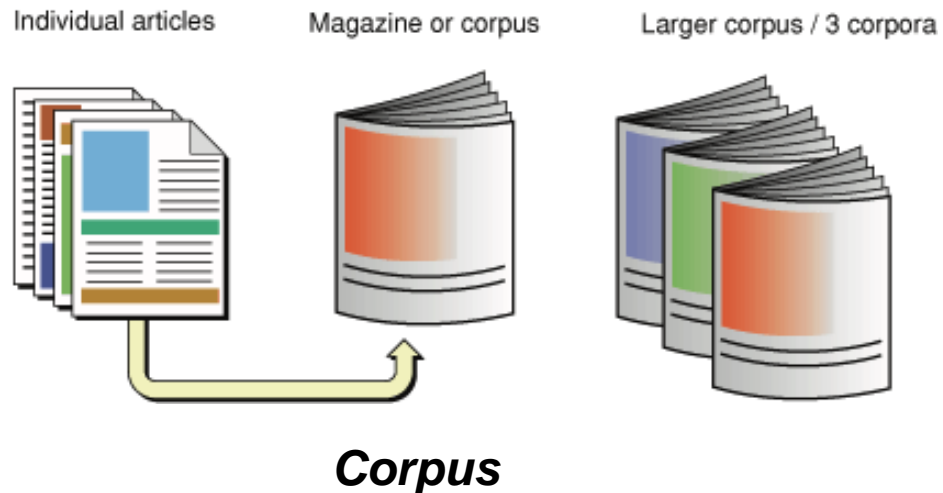# Overview of building a Text based Predictive Model



Problem Definition

Unorganised State

Organised State

Document Selection / Text Collection

↓

Organisation

↓

Feature Extraction

↓

Feature Selection

↓

Analysis

↓

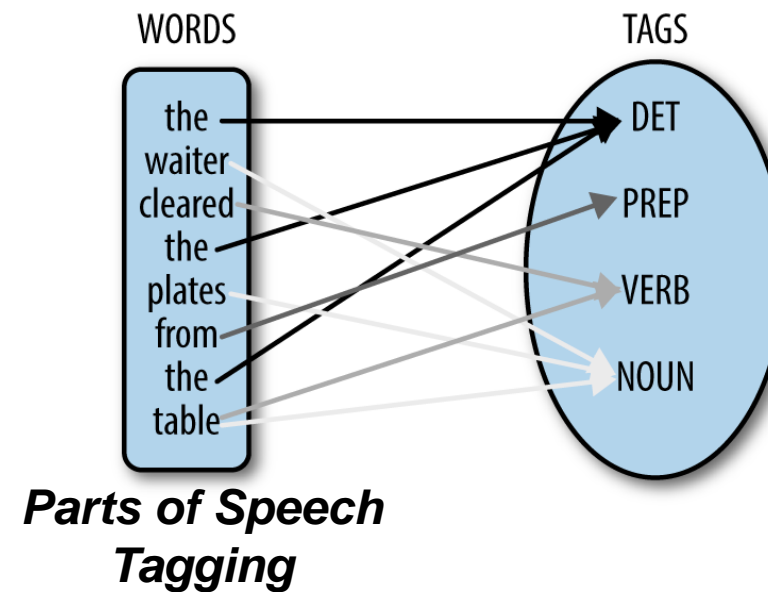Insight → Predictive Modelling

# Text Analytics - Terminologies


Document


Tokens / Terms

Stopwords



Corpus

Parts of Speech Tagging

# Textual Features – Extraction Methods

▶ Bag-of-words

▶ n-grams

▶ Scoring words – counts, frequencies, binary

  ▶ Histogram

  ▶ Document Term Matrix

  ▶ Term Frequency – Inverse Document Frequency

▶ Word hashing

▶ Word embeddings

# Document Term Matrix

| | Term 1 | Term 2 | . | . | Term n |
|---|---|---|---|---|---|
| Document 1 | | | | | |
| Document 2 | | | | | |
| . | | | | | |
| . | | | | | |
| Document m | | | | | |

# Term Frequency – Inverse Document Frequency (TF-IDF)

- Importance / Significance of a word.

- $TF = \dfrac{Number\ of\ times\ Term\ t\ occurs}{Total\ number\ of\ terms}$

- $IDF = Log\left(\dfrac{Total\ number\ of\ docs}{Number\ of\ docs\ containing\ Term\ t}\right)$

- $TF - IDF = TF * IDF$

# Image Feature

▶ **Matrix of numbers**

▶ **Size of this matrix depends on the number of pixels of the input image.**

▶ **Pixel Values**

   ▶ **Intensity and Brightness - how bright that pixel is?**

   ▶ **what color it should be?**

   ▶ **Important shape / objects / edges**

# Audio Features

## Abstraction Level

High – Understood and enjoyed by humans
- instrumentation, key, chords, melody, harmony, rhythm, genre, mood, etc

Mid – Perceived by humans
- pitch, beat-related descriptors, note onsets, fluctuation patterns, s, etc

Low - statistical features which sense to the machine, but not to humans
- amplitude envelope, energy, spectral centroid, spectral flux, zero-crossing rate,

## Temporal scope

Instantaneous
- Range in milliseconds – 10ms

Segment-level
- Wider range

Global
- Aggregate feature for whole word / sentence

## Signal domains

Time zone
- waveforms of the raw audio.

Frequency zone
- Frequency component

Time-Frequency zone
- Combination of Time and Frequency

## Time domain

Amplitude

Root mean square energy

Zero crossing rate

# Exploratory Data Analytics

▶ Examine the data distribution

▶ Handling missing values of the dataset

▶ Handling the outliers

▶ Removing duplicate data

▶ Encoding the categorical variables
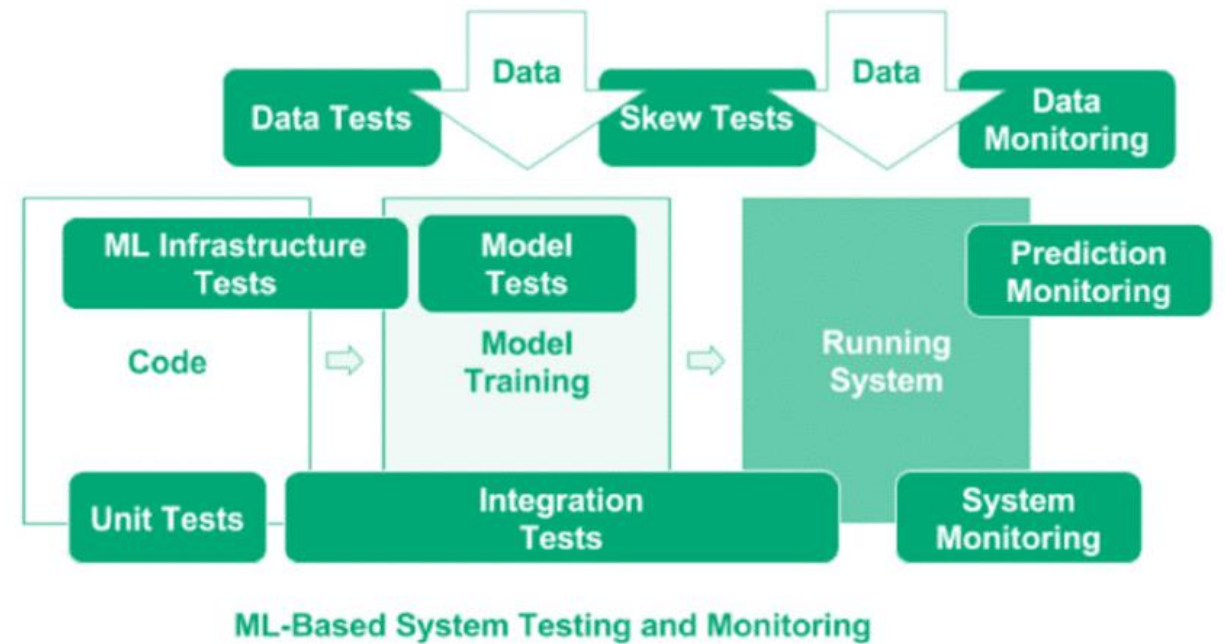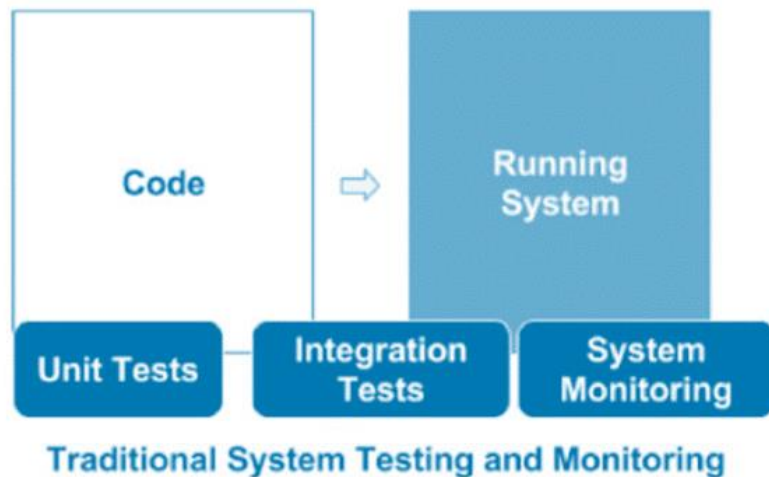
▶ Normalizing and Scaling

# Predictive Modeling

▶ **Statistical technique using machine learning and data mining**

▶ **Requires historical and existing data**

▶ **Forecast future based on historical and existing data**

▶ **Types of Predictive Modelling**

    ▶ **Classification Model**

    ▶ **Clustering Model**

    ▶ **Time Series Model**

    ▶ **Forecast Model**

    ▶ **Language Model**

Gradient boosting Model

Generalized Linear Models

K Means

Regression

Random Forest

# Traditional Vs ML System



Traditional System Testing and Monitoring

ML-Based System Testing and Monitoring

# ML-Ops – System Pipeline

▶ **Major challenges in ML Models and SM Data**

   ▶ Tracability – What ideas have been tried? What are successful?

   ▶ Reproducibility – How to reproduce successful ideas?

▶ **ML Lifecyle**

   ▶ Manage resources, data, code, time, and quality to meet objectives.

**Continuous Integration**
- Checkout code
- Complete task
- Validate against code base
- Perform unit testing
- Remerge code

**Continuous Training**
- Monitor
- Measure
- Retrain
- Serve

**Continuous Delivery**
- Build
- Test
- Release

# Predictive Modeling – Case Study

Predictive Modeling
Emotion Modeling

# Identifiable Entities in Social Media Data

| Entities | Description |
| --- | --- |
| ORG | referring to institutions, organizations, companies, agencies etc., |
| GPE | denoting countries, states, and cities |
| PERSON | representing people names and fictional characters |
| DATE | symbolizing relative or absolute dates or time periods |
| TIME | pertaining to time periods that are shorter than a day |
| NORP | mapping to nationalities, political and religious groups or communities |
| LOC | signifying any non-GPE entities like water bodies, mountains, etc., |
| PRODUCT | Characterizing things, objects, food, vehicles and other non-service entities |
| EVENT | for data indicating any disasters like battles, hurricanes, earthquakes, wars, sports and well-established happenings |
| PERCENT | implying any values represented in percentage format (%) |

# Text Data Preprocessing



Raw Tweet

1. Convert tweet to lower case

2. Remove '@' in user mentions

3. Remove '#' in hashtags

4. Remove punctuations

5. Replace links with 'URL'

6. Remove Expressions

7. Remove leading and lagging blank spaces

8. Remove stop words

9. Remove Non-English Words

10. Replace date / time values with "at_datetime"

11. Replace phone numbers with "at_tele"

Pre-processed Tweet

# Entity – Emotion Modeling

▶ Identification of entities through NER – Sequence labelling task

▶ Understanding the emotions through emotion mining mechanism

▶ Build features based on discovered entities and emotion

▶ Study 'Correlations' by joint probability between adjacent features

▶ Interdependency modelling by statistical analysis on univariate and multi-variate entity correlation

▶ Rank based correlation – Spearman's correlation coefficient

$$\rho = 1 - \frac{6 \sum_{i=1}^{n} (u_i - v_i)^2}{n(n^2 - 1)}$$

# Predictions – Insights on Citizen opinions

# Predictions – Insights on Citizen opinions



| | Share of Voice |
|---|---|
| Opening/using a medical savings account | 31% |
| Cannot afford cost of healthcare | 24% |
| Tax increases | 19% |
| Premium increases | 16% |
| Pay out -of-pocket | 10% |

Share of Voice

# Linear Regression model
## How strong is your relationship?

# Predictable Network Features in Social Media



| Feature | Value |
|---|---|
| Days since last communication | −0.762 |
| Days since first communication | 0.755 |
| Intimacy × Structural | 0.4 |
| Wall words exchanged | 0.299 |
| Mean strength of mutual friends | 0.257 |
| Educational difference | −0.223 |
| Structural × Structural | 0.195 |
| Reciprocal Serv. × Reciprocal Serv. | −0.19 |
| Participant-initiated wall posts | 0.146 |
| Inbox thread depth | −0.137 |
| Participant's number of friends | −0.136 |
| Inbox positive emotion words | 0.135 |
| Social Distance × Structural | 0.13 |
| Participant's number of apps | −0.122 |
| Wall intimacy words | 0.111 |

# Audio – Debater - An automated Debating System

Where are we now...?!

# Google Cloud Platform

| CATEGORY | PRODUCT | FEATURES |
|---|---|---|
| **Vertex AI** | **Vertex AI**<br><br>Unified platform to help you build, deploy and scale more AI models. | ✅ Prepare and store your datasets<br>✅ Access the ML tools that power Google<br>✅ Experiment and deploy more models, faster<br>✅ Manage your models with confidence |
| **Sight** | **AutoML Image**<br><br>Derive insights from object detection and image classification, in the cloud or at the edge. Try it now. | ✅ Use REST and RPC APIs<br>✅ Detect objects, where they are, and how many<br>✅ Classify images using custom labels<br>✅ Deploy ML models at the edge |
| | **AutoML Video**<br><br>Enable powerful content discovery and engaging video experiences. Try it now. | ✅ Annotate video using custom labels<br>✅ Streaming video analysis<br>✅ Shot change detection<br>✅ Object detection and tracking |
| **Language** | **AutoML Text**<br><br>Reveal the structure and meaning of text through machine learning. Try it now. | ✅ Integrated REST API<br>✅ Custom entity extraction<br>✅ Custom sentiment analysis<br>✅ Large dataset support |
| | **AutoML Translation**<br><br>Dynamically detect and translate between languages. Try it now. | ✅ Integrated REST and gRPC APIs<br>✅ Supports 50 language pairs<br>✅ Translate with custom models |
| **Structured data** | **AutoML Tabular**<br><br>Automatically build and deploy state-of-the-art machine learning models on structured data. Try it now. | ✅ Handles wide range of tabular data primitives<br>✅ Easy to build models<br>✅ Easy to deploy and scale models |

# Microsoft Azure

**Azure Cognitive Services**

Speech

Language

Vision

Decision

Improve customer experiences with **Cognitive Service for Speech**

**Speech to Text**

Transcribe audible speech into readable, searchable text.

**Text to Speech**

Convert text to lifelike speech for more natural interfaces.

**Speech Translation**

Integrate real-time speech translation into your apps.

**Speaker Recognition**

Identify and verify the people speaking based on audio.

# Microsoft Azure

# Amazon Web Services – AI / ML Services

https://aws.amazon.com/machine-learning/