# LO1 – LO2

Examine data mining through social media tracking

# Objectives

- Introduction to Data analytics

- Collecting and extracting social media data

- Data analysis, visualization and exploration

# Objective of this session

After attending this session, you should be able

- Difference between structured and unstructured data
- Explain how data analysis is performed on a typical structured dataset
- List some of the techniques of quantitative data analysis

# Structured vs unstructured data

- Data can be represented in various types of structures, formats, and media



GDP per capita from world bank

SAT data from NYC open data (nycopendata.Socrata.com)

Blog of new York public library (https://www.nypl.org/blog)

# Multimedia data



Multimedia data (YouTube)

# Social media data

# Structured data

- Structure data refers to data with a high level of

- Organization, such as in relational databases and spreadsheets

- Depends on data model- a model of the data types and how they will be stored, processed and accessed

- Easily entered, stored, queried and analyzed

- Structured query language (SQL) is used for management of structured data (e.g. MySQL)

# Unstructured data

- Unstructured data means all things that cannot be classified and fit into one simple model

- Photos and graphics images

- Videos

- Streaming instrument data

- Webpages, emails, blog entries, wikis

- Pdf files, PowerPoint presentations, and word processing documents

# Structured and Unstructured features in social media data

- Structured data can be used for number-driven (quantitative) approaches

- Who, what, when, where, how and how many?

- Unstructured data for qualitative approaches:

- Why?

- Sentimental analysis



Positive Sentiment and Negative Sentiment Ticker Count

# Class Exercise

**1.1 Explain Various Kinds of Website/Social Media Data**

**Objective:** Understand different types of social media data available on websites.

**Activity 1: Brainstorming Session (20 mins)**

- "What social media data can we extract from websites?"
- "How is this data useful for businesses?"

**Discuss and categorize:**

- **Engagement Metrics:** Likes, shares, retweets, comments
- **User Data:** Demographics, location, follower count
- **Content Data:** Post text, hashtags, media type
- **Temporal Data:** Posting times, engagement patterns
- **Sentiment Data:** Positive/negative/neutral tone in comments

**Activity 2: Techniques to Analyze Social Media Data (10 mins)**

Submit pdf in the Dropbox '**Social Media Data Analysis**'

SASKATCHEWAN POLYTECHNIC | Tomorrow in the making

# Structure – Real-World Social Media Data

- **X/Twitter post → structured**: username, follower count, retweet count, likes, timestamp, verified badge; **unstructured**: tweet text, images, emojis

- **Instagram profile → structured**: follower count, following, number of posts, location tag; **unstructured**: bio text, stories, reels video

- **YouTube video page → structured**: view count, likes/dislikes, subscriber count, upload date, comments count; **unstructured**: video itself, title, description, comments text

- **TikTok video → structured**: likes, comments, shares, play count, duration, music name; **unstructured**: video, caption, hashtags, stitched/duet content

# Analyzing structured data: Descriptive statistics

- Descriptive statistics are used to describe the basic features of the data

- Summaries about the sample and the measures

  - Distribution (frequency table)

  - Central tendency (mean, median)

  - Dispersion (standard deviation)

- Simple graphics analysis

  - Bar chart, pie chart

  - Box plot



Number of retweets and favorites per tweet (log)
(http://www.martingrandjean.ch/)

| | KATY PERRY | BARACK OBAMA | HISTORYINPICS | REALTIMEWWII | GALLICABNF | UKWARCABINET |
|---|---|---|---|---|---|---|
| MAX | 46 968 | 19 280 | 15 957 | 615 | 90 | 14 |
| QUARTILE 3 | 23 493 | 2 647 | 3 652 | 136 | 25 | 2 |
| MEDIAN | 16 755 | 1 676 | 2 567 | 90 | 14 | 1 |
| QUARTILE 1 | 11 416 | 1 296 | 1 516 | 63 | 7 | 0 |
| MIN | 9 424 | 832 | 108 | 32 | 0 | 0 |

# Correlation

- In general, correlation refers to the extent to which two variables have a linear relationship with each other

- We will learn how to generate and test the correlation in python and R in upcoming session



Positive correlation between height and weight



Negative correlation between GPA and video game

# Regression

- Once a correlation is found, we do regression analysis to estimate the relationship among variables

- This relationship knowledge can then be used to predict one variable using others

# Visualization of Data

- Data analytics techniques could give us insightful values, but it may take the right visualization to convey the meaning for decision-making.
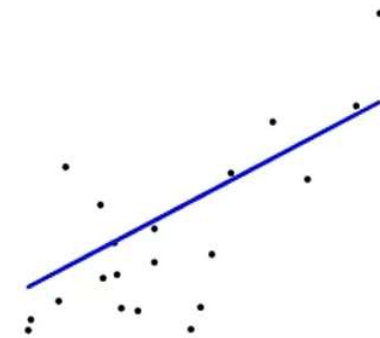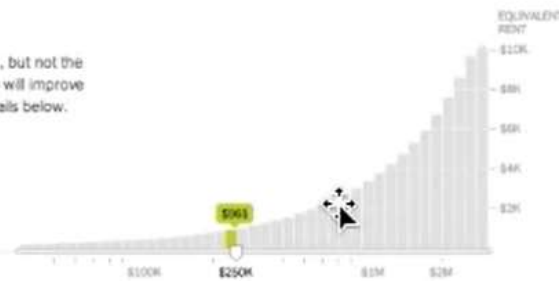
- Visualization helps understand the characteristics of data and provides insights from it

- Discovery of new phenomena

- Sense-making of what data delivers to people

- Communication method between data analysts, decision-makers, service providers, etc

# Renting vs Buying (from NYTimes)

# Histogram

- A graphical representation of the distribution of numerical data

- Give a rough sense of the density of the underlying distribution of the data



When are people working? (from NPR)

# Bar chart

- Present grouped data with rectangular bard with lengths proportional to the values that they represent.



Interactive bar chart of NYC street trees

# Time plot

- Displays values against time

- Helps understand trends

# summary

- Social media data = mix of structured (numbers, counts, timestamps) and unstructured (text, images, video)

- Structured data → perfect for fast, quantitative analysis ("how many?", "when?", "who?")

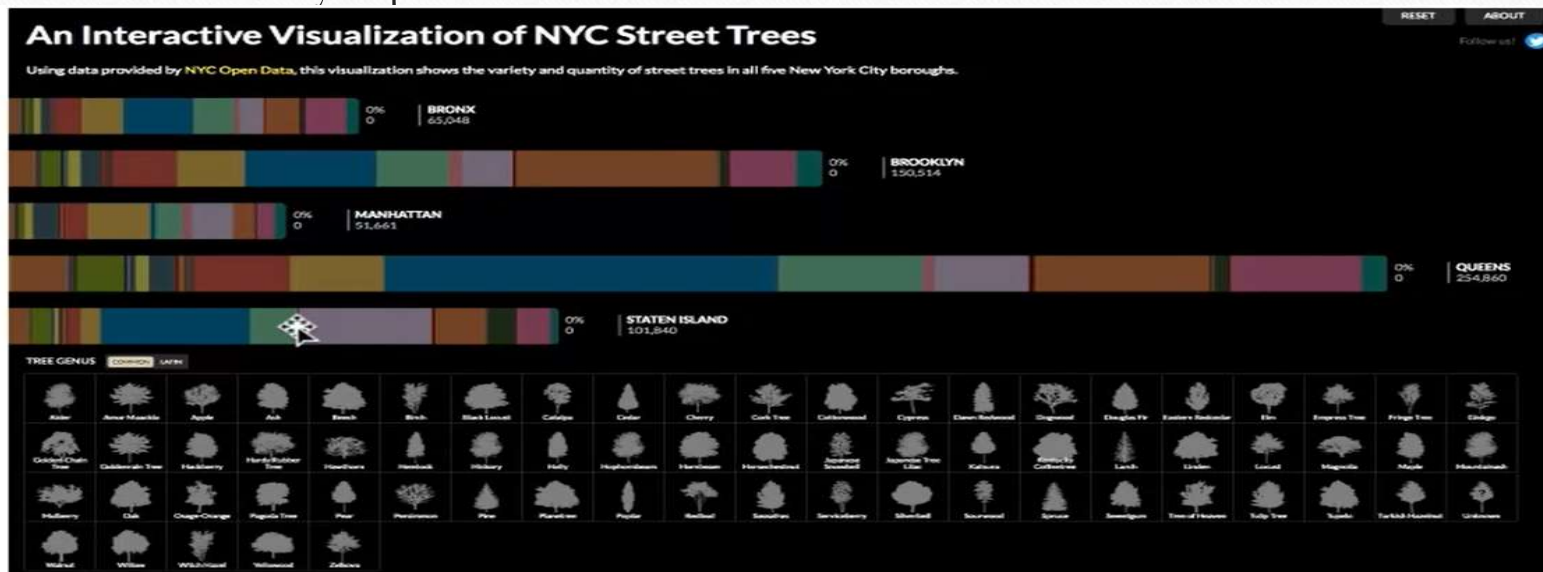- Unstructured data → needs NLP, sentiment analysis, etc. (we'll cover later)

- Focused only on the structured part → easier & immediately actionable

- **Descriptive Statistics**

    Central Tendency → Mean, Median, Mode

    Dispersion → Range, Standard Deviation, Variance

    Distribution → Frequency tables, Histograms

    Relationships → Correlation → Regression (prediction)

- **Visualization Toolbox**

    Histogram → see distribution shape

    Bar Chart → compare categories

    Time Plot / Line Chart → spot trends over time

    Pie Chart → show proportions

    Scatter plots, Box plots, Heatmaps

Himanshu Patel, Instructor
Saskatchewan Polytechnic
email: patelh@saskpolytech.ca
Mining building, Saskatoon