# Text Analysis for Reddit in Big Data

Group 12 - Jahyeon Jee, Pi-Te Chen, Audrey Jung, Seunggyun Shin, Hyun Ji Lee

# Agenda

- Objectives
- Pipeline
- Performance Metrics
- Data Explanation
- Pre-processing and Cleaning
- Analysis
- Conclusions

# Objectives

- The project will utilize AWS-based workflow for data collection, preprocessing, and analysis using visualization.
- The project represents a convergence of technology and social science, bridging the gap between Big Data analytics and online community engagement.
- The objective is to develop a robust Big Data Pipeline to extract actionable insights from the subreddit's content.

# Objectives

The aim is to contribute meaningfully to the advancement of knowledge in employment dynamics and offer practical solutions to challenges faced by job seekers.
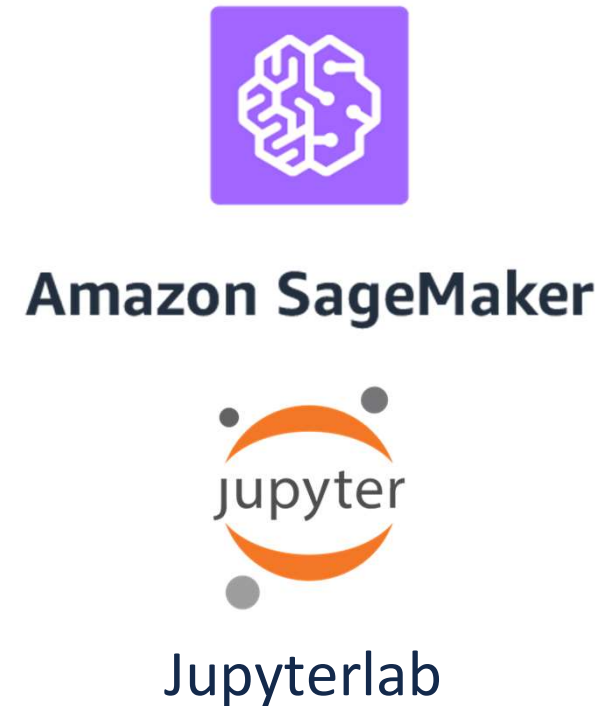
# Pipeline



**Data Extracting**       **Data Storing**       **Preprocessing, Cleaning, Analysis**

Jupyterlab

# Performance Metrics - Resources & Cost

## EC2 (2 weeks usage)

| Name ✎ ▽ | Instance ID | Instance state ▽ | Instance type ▽ | Platfor... ▽ |
|---|---|---|---|---|
| jjee2_ec2 | i-02ea7f5e56ea54621 | ⊘ Running ⊕ ⊖ | t2.micro | Linux/UNIX |

## S3 (2 weeks usage)

| Total storage | Object count | Average object size |
|---|---|---|
| 582.2 KB | 1 | 582.2 KB |

## SageMaker (2 weeks usage)

| Status | Notebook instance type | Platform identifier |
|---|---|---|
| ⊘ InService | ml.t2.medium | Amazon Linux 2, Jupyter Lab 3 (notebook-al2-v2) |

| Creation time | Elastic Inference | Minimum IMDS Version |
|---|---|---|
| Apr 20, 2024 23:36 UTC | - | 2 |

| Last updated | Volume Size | |
|---|---|---|
| Apr 25, 2024 04:22 UTC | 5GB EBS | |

Cost
- AWS: **Free**
- Reddit API: **Free**

# Data Explanation

**reddit**    **PREFERENCES**    options    apps    RSS feeds    friends    blocked    password/email    delete      peter7173 (1) | 📧 2 | 💬 | **preferences** | logout

## developed applications

**Project**
personal use script
X30HzbH-SUlMDphZqxM1Hg

API for project

edit      Developers: peter7173

**Project Peter**
personal use script
mt9zM-Scxh-IhkNx-V61xg

API for project

edit      Developers: peter7173

create another app...

BADM 558
Big Data Infrastructure

# Data Explanation

create application

By creating an app, you agree to Reddit's Developer Terms and Data Api Terms. You must also register to use the API.

**name** [_____]

◉ web app       A web based application
○ installed app   An app intended for installation, such as on a mobile phone
○ script         Script for personal use. Will only have access to the developers accounts

**description**
[_____]

**about url** [_____]

**redirect uri** [_____]

✓ 我不是機器人   reCAPTCHA
隱私權 - 條款

[create app]   application created

# Data Explanation

# Data Explanation

```
>>> subreddit = reddit.subreddit('jobs')
>>> submissions = []
>>> for submission in subreddit.hot(limit=1000):
...     submission_info = {
...         "title": submission.title,
...         "score": submission.score,
...         "upvotes": submission.ups,
...         "downvotes": submission.downs,
...         "num_comments": submission.num_comments,
...         "submission_id": submission.id,
...         "url": submission.url,
...         "created_utc": submission.created_utc,
...         "Content": submission.selftext  # Add the content of the post
...     }
...     # Append the submission_info dictionary to the submissions list
...     submissions.append(submission_info)
...
```

# Data Explanation

Transform the data to CSV file

```
>>> import pandas as pd
>>> df = pd.DataFrame(submissions)
>>> df.to_csv(csv_file, index=False)
```

# Data Explanation

| | title | score | upvotes | downvotes | num_comments | submission_id | url | created_utc | Content |
|---|---|---|---|---|---|---|---|---|---|
| 9 | I got put on PIP…is it too late to tell my bos... | 88 | 88 | 0 | 90 | 1c6zq86 | https://www.reddit.com/r/jobs/comments/1c6zq86... | 1.713435e+09 | I started working for a nonprofit a year ago. ... |
| 568 | Companies need to learn to build from the bott... | 0 | 0 | 0 | 0 | 1c5b2q7 | https://www.reddit.com/r/jobs/comments/1c5b2q7... | 1.713256e+09 | Look I get it, everybody is lazy and wants to ... |
| 216 | Scams used to believable | 2 | 2 | 0 | 2 | 1c6vls2 | https://i.redd.it/wqec23xec6vc1.jpeg | 1.713418e+09 | Copywriter/research assistant/ typer/ and edit... |
| 320 | Haven't heard back from a job after calling th... | 1 | 1 | 0 | 1 | 1c6gg6j | https://www.reddit.com/r/jobs/comments/1c6gg6j... | 1.713376e+09 | I had an interview on Tuesday and a working in... |
| 265 | Is it time to move on? | 3 | 3 | 0 | 2 | 1c6fs7i | https://www.reddit.com/r/jobs/comments/1c6fs7i... | 1.713375e+09 | I recently had a job interview on the 8th, see... |

# Pre-processing & Cleaning

**Reddit Post**: "This is the weekly success and disappointment Megathread for the week. Please post all of your successes and disappointments for this week, including job offers and other victories, as well as any venting of frustration, in this thread, and this thread only. Thanks!"

**Remove Punctuations**
(re package: RegEx)
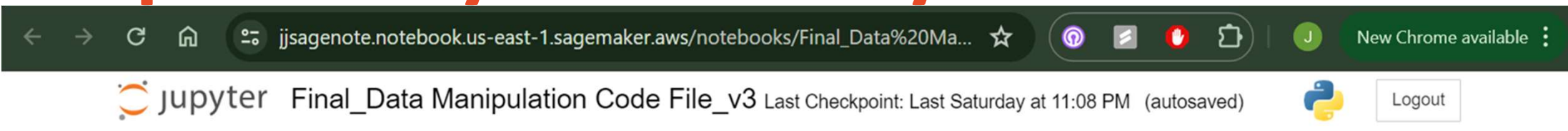
this is the weekly success and disappointment megathread for the week please post all of your successes and disappointments for this week including job offers and other victories as well as any venting of frustration in this thread and this thread only thanks
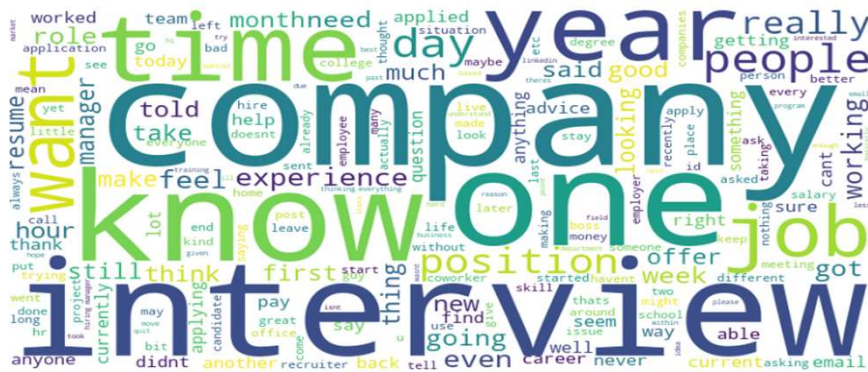
**Stop Word Dictionary**

(nltk package: Natural Language Toolkit)

weekly success disappointment megathread week please post successes disappointments week including job offers victories well venting frustration thread thread thanks

**Tokenization**
**Lemmatization**
**Vectorization**

**Visualization (Word Cloud)**
**Sentiment Analysis**
**Trend Analysis**

**Post Date:** 1713490000 (Timestamp)

**Post Date:** 2024-04-19 01:23:41

# Exploratory Data Analysis

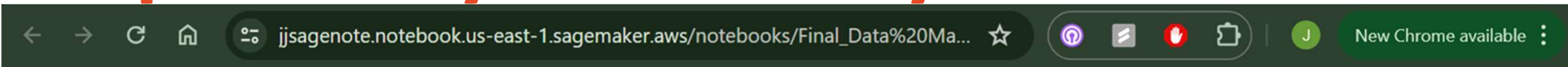jupyter   Final_Data Manipulation Code File_v3   Last Checkpoint: Last Saturday at 11:08 PM   (autosaved)   Logout

## Word Cloud



General words:
company, interview, know



Refined words:
Interview, time, position

# Exploratory Data Analysis



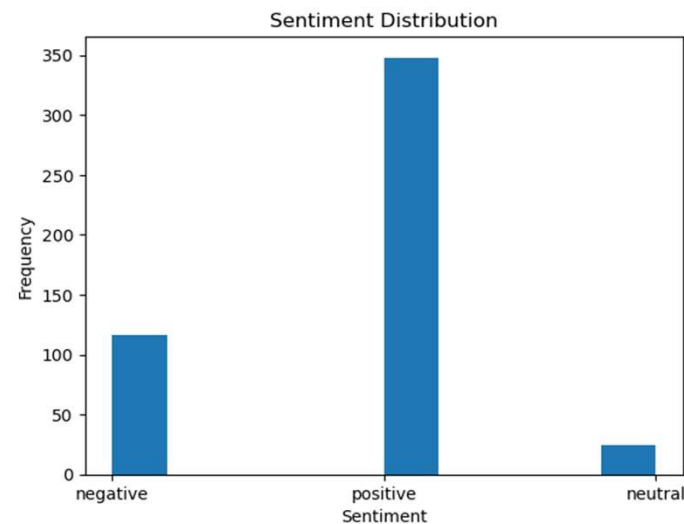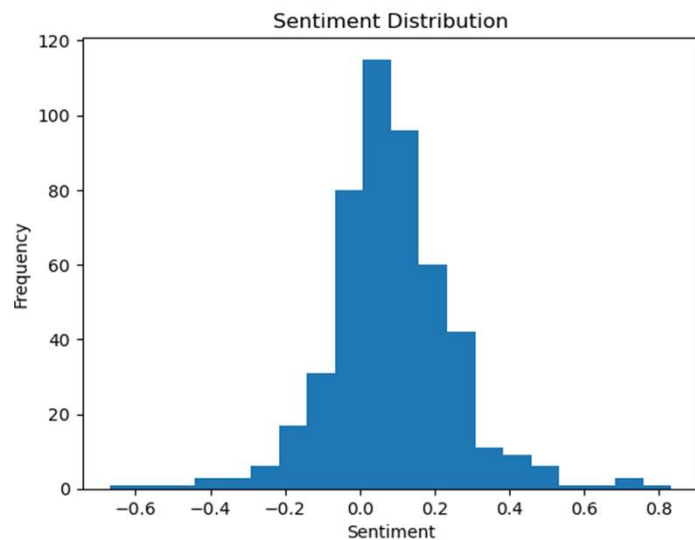**Sentiment Analysis: Overview**
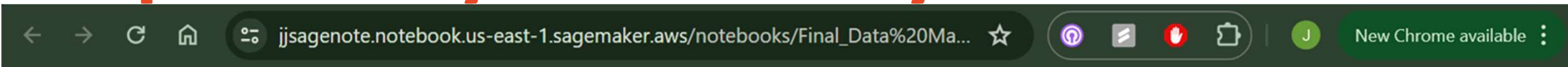


**Method**
- Barplot with sentiment score/type as x and frequency as y

**Analysis**
- Almost normal distribution
- More positive reactions

# Exploratory Data Analysis



## Bivariate Analysis: Score vs. Sentiment

**Method**
- Boxplot and Heatmap

**Analysis**
- Weak pattern of worsening the sentiment by the score
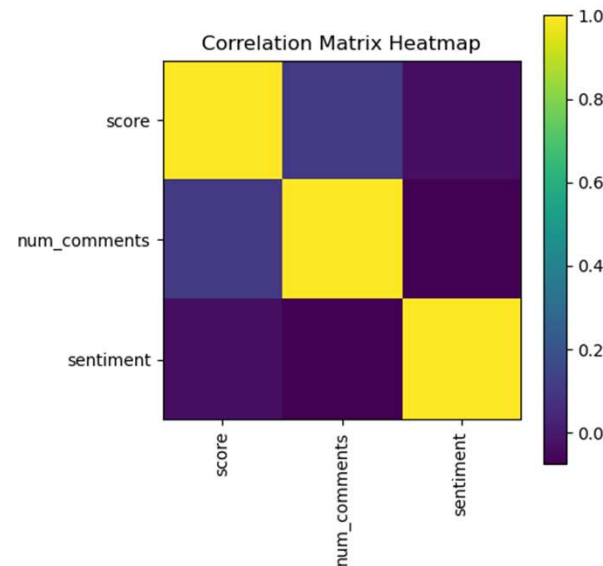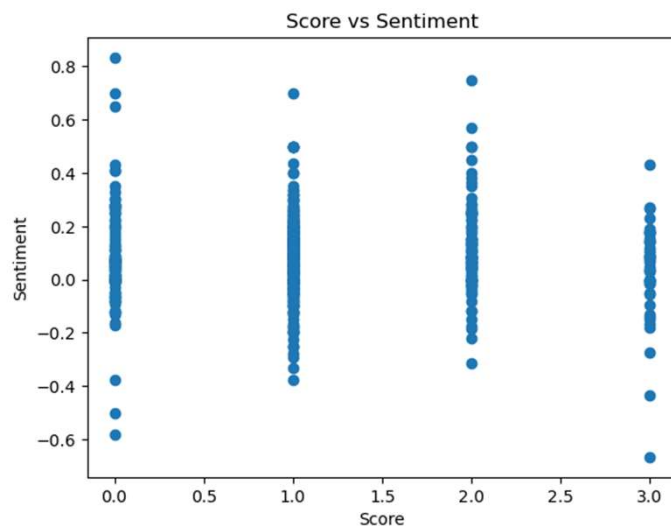- Weak correlation between score and number of comments

# Exploratory Data Analysis

jupyter  Final_Data Manipulation Code File_v3 Last Checkpoint: Last Saturday at 11:08 PM  (autosaved)    Logout
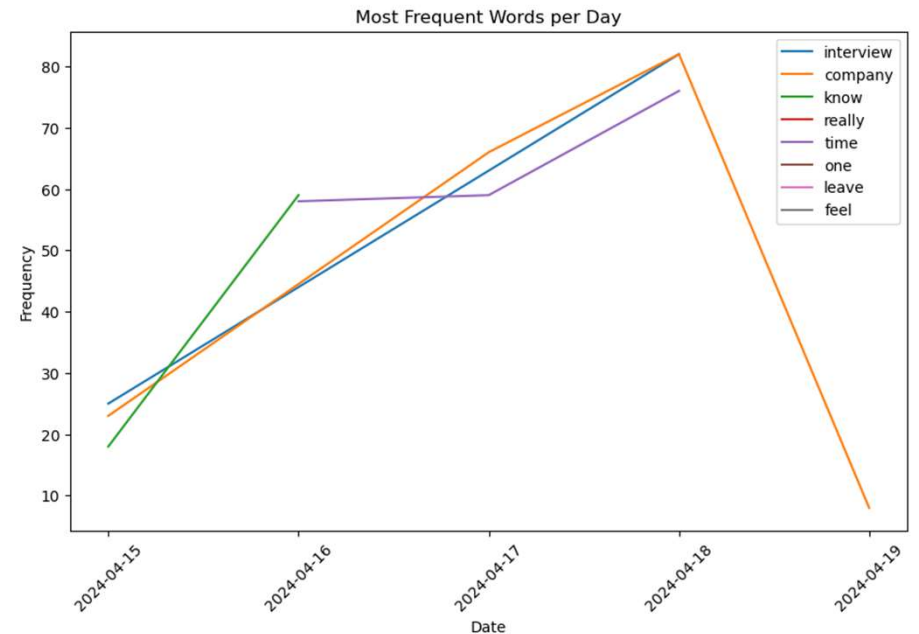
## Trend Analysis: Frequent Words by day

**Method**
- Plotting the word frequency data with the date as x, frequency as y with legend of words

**Analysis**
- Overlapping words for several days
- Similar trend with the total number of posts



Most Frequent Words per Day

Legend: interview, company, know, really, time, one, leave, feel

BADM 558
Big Data Infrastructure

# Conclusions

### Summary

By harnessing a range of tools and services on the AWS platform, including the Reddit API, we delved into the intricacies of big data infrastructure and established efficient pipelines. Our analysis provided valuable insights into the prevailing trends in web-generated content related to job searches, underscoring its immediate relevance and utility for our objectives.

### Innovation

Advanced Data Collection Methods: API usage is efficient and real-time
Natural Language Processing (NLP): Word frequency, Sentiment analysis, and the Trend
Scalability and Flexibility: Through AWS platform, data usage can be easily extended

### Future Potential

- More detailed and long term analysis
- Diverse opinions