**SIMAD UNIVERSITY**

**MASTER OF DATA SCIENCE**

**SEMESTER 5**

**APRIL 2024**

**SOCIAL MEDIA ANALYTICS AND CONSUMER BEHAVIORAL**

**COURSE PROJECT**

**"COMPREHENSIVE JOB MARKET ANALYSIS IN SOMALIA: TRENDS, PATTERNS, AND INSIGHTS FACEBOOK + WEB DATA"**

**ABDIQANI AHMED ABDULLE**

**LECTURER: DR. FUAD MIRE**

# INTRODUCTION

This report presents the findings of project titled "Comprehensive Job Market Analysis in Somalia: Trends, Patterns, and Insights." The initiative was undertaken as a part of the 'Social Media Analytics and Consumer Behavioral' course. It showcases the capabilities of Python in scraping social media and web data to gain a deeper understanding of job market dynamics. The project aimed to scrape job postings from the prominent Somali job site 'qaranjobs.so' alongside the analysis of corresponding engagement metrics from the website's Facebook page.

## Objective of the Project

The primary objective of this project was to unravel the intricacies of the Somali job market by analyzing various trends and patterns within the employment sector. By studying job titles, company hiring activities, location-based job availability, and social media engagement, we sought to gain insights into the demand for specific skills and professions in Somalia.

## Project Significance

The analysis stands significant as it synthesizes a vast array of job market data to unearth trends that are critical for job seekers, employers, and policymakers. It helps in understanding which sectors are booming with opportunities and garners attention on social platforms. The derived insights serve as a foundation for strategies aimed at improving job visibility and alignment with market needs, ultimately contributing to economic growth and employment rates.

## Project Scope

The project's scope entails an elaborate exploration of the job market through web scraping of 'qaranjobs.so' and analyzing the public response to job postings on Facebook. It assesses several facets such as the distribution of jobs across various industries, the volume of postings per month, and the day-to-day trends. Additionally, the study involves examining Facebook interactions to identify which job listings draw the most attention and engagement, thereby offering a window into the job market's responsiveness and public interest.

**Background**

The project is set against the backdrop of a burgeoning Somali job market, where digital platforms have become crucial in connecting employers with potential employees. Despite the growing use of online job portals and social media for recruitment, there has been a lack of comprehensive data analysis to guide job seekers and employers effectively.

# DATASET

Data was sourced from 'qaranjobs.so,' using a custom Python scraper tailored to extract job listing information including job titles, companies, locations, and posting dates. Additionally, a Python script utilizing the Facebook Cookies method was developed to gather engagement metrics for these postings. This included likes, comments, shares, and detailed reactions to each job post.

**Data Overview**

The final dataset consisted of a compilation of 4000 job records scraped from 'qaranjobs.so' and 5000 Facebook posts. Post-cleaning, an inner join operation on the 'url' column was performed to merge web-scraped data with Facebook engagement data. The merged dataset included columns such as 'Title,' 'Company,' 'Location,' 'Posted,' 'Url,' 'Created_time,' 'Facebook_Post,' 'Shares_Count,' 'Reactions_Count,' and 'Comments_Count.' Each entry provides a snapshot of a job posting and its social engagement, painting a holistic picture of the job market's online footprint.

Here is web scraped data:

| title | company | location | posted | url |
|---|---|---|---|---|
| Somali jobs: The Federal Republic of Somalia (... | Federal Republic of Somalia | Somalia | March 31, 2024 | https://qaranjobs.com/job/somali-jobs-the-fede... |
| Somali jobs at Save the Children (6 positions) | Save the Children | Somalia | March 30, 2024 | https://qaranjobs.com/job/somali-jobs-at-save-... |
| Somali jobs at UNICEF (3 Positions) | United Nations Children's Fund (UNICEF) | Somalia | March 24, 2024 | https://qaranjobs.com/?post_type=job_listing&p... |
| Jobs at Somali Red Crescent Society (SRCS) (4 ... | Ururka Bisha Cas (SRCS) | Somalia | March 23, 2024 | https://qaranjobs.com/?post_type=job_listing&p... |
| Finance Assistance | Amtel | Somalia | March 20, 2024 | https://qaranjobs.com/job/finance-assistance-a... |
| ... | ... | ... | ... | ... |
| Child Protection Information Management Officer | IRC | Awbare, Shedder and Kebribeyah, Somali Regiona... | July 27, 2022 | https://qaranjobs.com/job/child-protection-inf... |
| Terms of Reference (ToR) for Open tender for P... | Candlelight | Buhodle, Somaliland | July 27, 2022 | https://qaranjobs.com/job/terms-of-reference-t... |
| Saferty And Security Officer | Polish Humanitarian Action (PAH) | Mogadishu, Somalia | July 27, 2022 | https://qaranjobs.com/job/saferty-and-security... |
| Finance Officer | NRC | Kismayo, Somalia | July 27, 2022 | https://qaranjobs.com/job/finance-officer-kism... |
| Consultant International pour Appui technique ... | UNICEF | Djibouti | July 27, 2022 | https://qaranjobs.com/job/consultant-internati... |

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4320 entries, 0 to 4319
Data columns (total 6 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   title     4320 non-null   object
 1   company   4320 non-null   object
 2   location  4320 non-null   object
 3   posted    4320 non-null   object
 4   url       4320 non-null   object
 5   stem      4320 non-null   int64
dtypes: int64(1), object(5)
memory usage: 202.6+ KB
```

Here is Facebook data:

```
df.head()
```
Python

| | id | created_time | permalink_url | full_picture | shares | reactions | message | comments |
|---|---|---|---|---|---|---|---|---|
| | 1.905016e+29 | 2024-03-31 11:40:46+00:00 | https://www.facebook.com/100069752805511/posts... | https://scontent.fmgq1-2.fna.fbcdn.net/v/t39.3... | {'count': 3} | {'data': [], 'summary': {'total_count': 30}} | Fursado Shaqo:\nTSFP/MCHN/BSFP Registration Cl... | {'data': [{'id': '726961929638894_324291780293... |
| | 1.905016e+29 | 2024-03-31 09:02:41+00:00 | https://www.facebook.com/100069752805511/posts... | https://scontent.fmgq1-2.fna.fbcdn.net/v/t39.3... | NaN | {'data': [], 'summary': {'total_count': 26}} | 5 Fursado Shaqo Ah:\nJamhuuriyadda Federaalka ... | {'data': [{'id': '726894056312348_478021441447... |
| | 1.905016e+29 | 2024-03-31 08:15:21+00:00 | https://www.facebook.com/100069752805511/posts... | https://scontent.fmgq1-2.fna.fbcdn.net/v/t39.3... | {'count': 3} | {'data': [], 'summary': {'total_count': 21}} | Fursad Shaqo:\nDhamaan Dadwaynaha Ku Nool Deeg... | {'data': [{'id': '726875066314247_433778902458... |
| | 1.905016e+29 | 2024-03-31 06:36:19+00:00 | https://www.facebook.com/100069752805511/posts... | https://scontent.fmgq1-2.fna.fbcdn.net/v/t39.3... | {'count': 1} | {'data': [], 'summary': {'total_count': 16}} | 4 Fursado Shaqo Ah:\nSalaam Somali Bank Waxaa ... | {'data': [], 'summary': {'total_count': 0}} |
| | 1.905016e+29 | 2024-03-30 13:08:50+00:00 | https://www.facebook.com/100069752805511/posts... | https://scontent.fmgq1-2.fna.fbcdn.net/v/t45.1... | {'count': 9} | {'data': [], 'summary': {'total_count': 154}} | 6 Fursado Shaqo Ah:\nHay'adda Save the Childre... | {'data': [{'id': '726381423030278_114188169057... |

```
df.info()
```
✓ 0.0s

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5075 entries, 0 to 5074
Data columns (total 9 columns):
 #   Column         Non-Null Count  Dtype
---  ------         --------------  -----
 0   id             5075 non-null   float64
 1   created_time   5075 non-null   datetime64[ns, UTC]
 2   permalink_url  5075 non-null   object
 3   full_picture   5074 non-null   object
 4   shares         2987 non-null   object
 5   reactions      5075 non-null   object
 6   message        5075 non-null   object
 7   comments       5075 non-null   object
 8   attachments    1 non-null      object
dtypes: datetime64[ns, UTC](1), float64(1), object(7)
memory usage: 357.0+ KB
```

# Data Cleaning and Integration

In the preliminary phase, the data cleaning process was initiated by aligning the 'url' column across both the web-scraped job postings and Facebook engagement datasets. This step was crucial to ensure that the subsequent merging operation would be seamless and error-free. Both datasets were then merged based on the 'url' column, which serves as the unique identifier linking the job posting to its social engagement data.

Following the merge, the dataset was pruned to include only pertinent columns necessary for analysis. These columns encapsulated crucial information such as the job title, hiring company, job location, posting date, URL, and corresponding Facebook data including post creation time, content, share count, total reaction count, and total comment count. This curated selection of columns allowed for a focused examination of the interplay between job postings and their social resonance.

The merged dataset then underwent a renaming process to standardize column names, thereby improving readability and accessibility of the data for analysis:

- 'title' was renamed to 'Title'

- 'company' to 'Company'

- 'location' to 'Location'

- 'posted' to 'Posted'

- 'url' to 'Url'

- 'created_time' to 'Created_time'

- 'message' to 'Facebook_Post'

- 'shares.count' to 'Shares_Count'

- 'reactions.summary.total_count' to 'Reactions_Count'

- 'comments.summary.total_count' to 'Comments_Count'

To address the issue of missing engagement metrics, which could potentially skew the analysis, NaN values within the engagement metric columns were replaced with zero. This decision was made to reflect a conservative estimate of engagement levels, assuming that missing data signified an absence of interaction.

Subsequent to the cleaning operations, the dataset was further refined by removing duplicate entries, ensuring that each record in the dataset represented a unique job posting and its associated engagement metrics.

Lastly, to maintain the integrity of the text data, parentheses and any enclosed text were excised from the dataset. This step was essential to eliminate extraneous information that could distort text analysis algorithms or sentiment evaluation.

Post-cleanup, the 'df_final' DataFrame consisted of 2765 entries, each containing 10 well-defined columns, with all null values addressed and textual data sanitized for analysis. This meticulous approach to data cleaning established a robust foundation for the subsequent analytical phases of the project.

Here is how the Post-cleanup Final dataframe look like:

```python
df_final.head()
```
Python

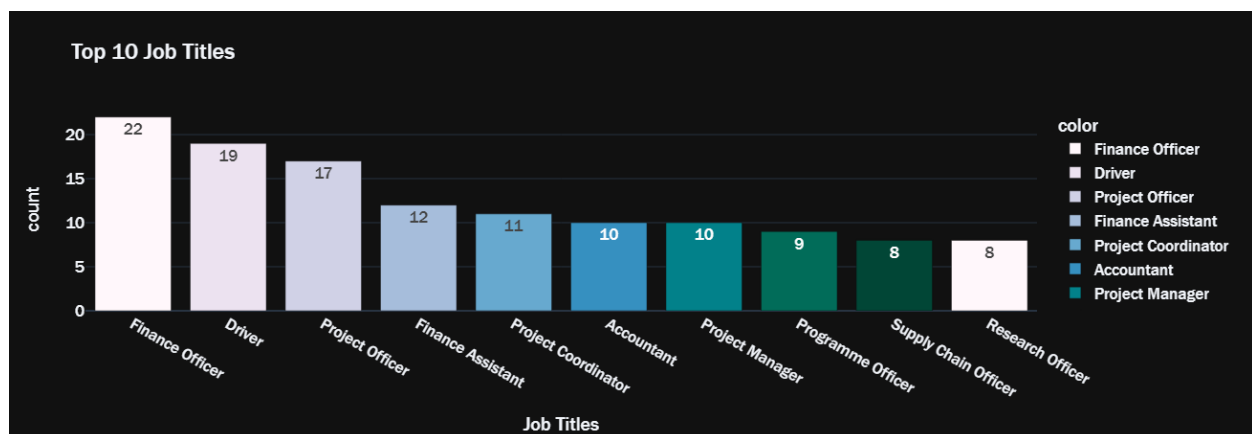| | Title | Company | Location | Posted | Url | Created_time | Facebook_Post | Shares_Count | Reactions_Count | Comments_Count |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Somali jobs: The Federal Republic of Somalia (... | Federal Republic of Somalia | Somalia | March 31, 2024 | https://qaranjobs.com/job/somali-jobs-the-fede... | 2024-03-31 09:02:41+00:00 | **5 Fursado Shaqo Ah:**\nJamhuuriyadda Federaalka ... | 0.0 | 28 | 1 |
| 1 | Somali jobs at Save the Children (6 positions) | Save the Children | Somalia | March 30, 2024 | https://qaranjobs.com/job/somali-jobs-at-save-... | 2024-03-30 13:08:50+00:00 | 6 Fursado Shaqo Ah:\nHay'adda Save the Childre... | 9.0 | 161 | 16 |
| 2 | Finance Assistance | Amtel | Somalia | March 20, 2024 | https://qaranjobs.com/job/finance-assistance-a... | 2024-03-23 13:16:28+00:00 | 2 Fursado Shaqo Ah:\nShirkadda Isgaarsiinta AM... | 5.0 | 59 | 14 |
| 4 | Financial Management Specialist, Baidoa, Somalia | Southwest State of Somalia, Municipality of Ba... | Baidoa, Somalia | March 31, 2024 | https://qaranjobs.com/job/financial-management... | 2024-03-31 08:15:21+00:00 | Fursad Shaqo:\nDhamaan Dadwaynaha Ku Nool Deeg... | 3.0 | 22 | 3 |
| 5 | CHIEF SECURITY OFFICER, Mogadishu, Somalia | United Nations Support Office in Somalia (UNSOS) | Mogadishu, Somalia | March 30, 2024 | https://qaranjobs.com/job/chief-security-offic... | 2024-03-30 12:08:34+00:00 | Fursad Shaqo:\nCHIEF SECURITY OFFICER, Mogadis... | 0.0 | 0 | 0 |

```python
df_final.info()
```
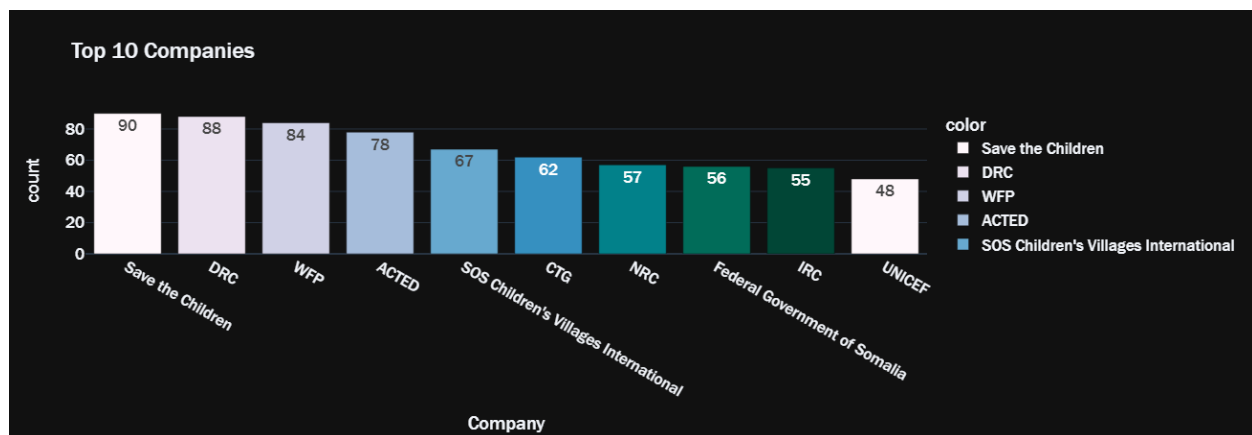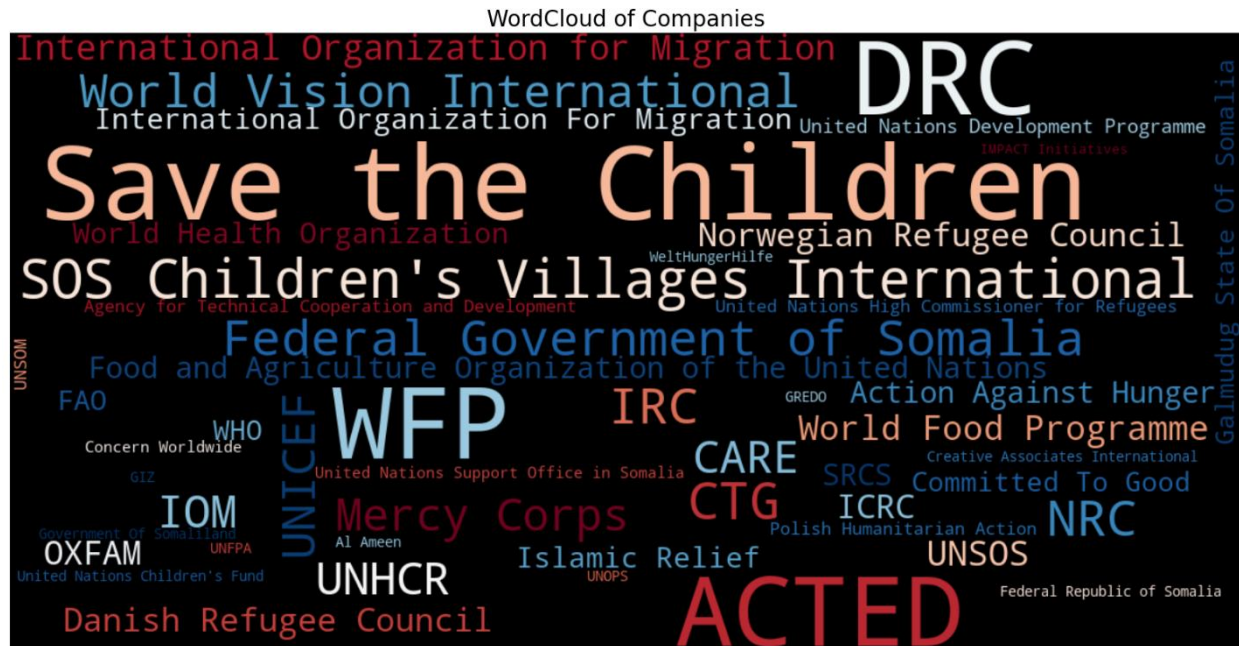3]

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2765 entries, 0 to 2919
Data columns (total 10 columns):
 #   Column           Non-Null Count   Dtype
---  ------           --------------   -----
 0   Title            2765 non-null    object
 1   Company          2765 non-null    object
 2   Location         2765 non-null    object
 3   Posted           2765 non-null    object
 4   Url              2765 non-null    object
 5   Created_time     2765 non-null    object
 6   Facebook_Post    2765 non-null    object
 7   Shares_Count     2765 non-null    float64
 8   Reactions_Count  2765 non-null    int64
 9   Comments_Count   2765 non-null    int64
dtypes: float64(1), int64(2), object(7)
memory usage: 237.6+ KB
```

# Exploratory Data Analysis

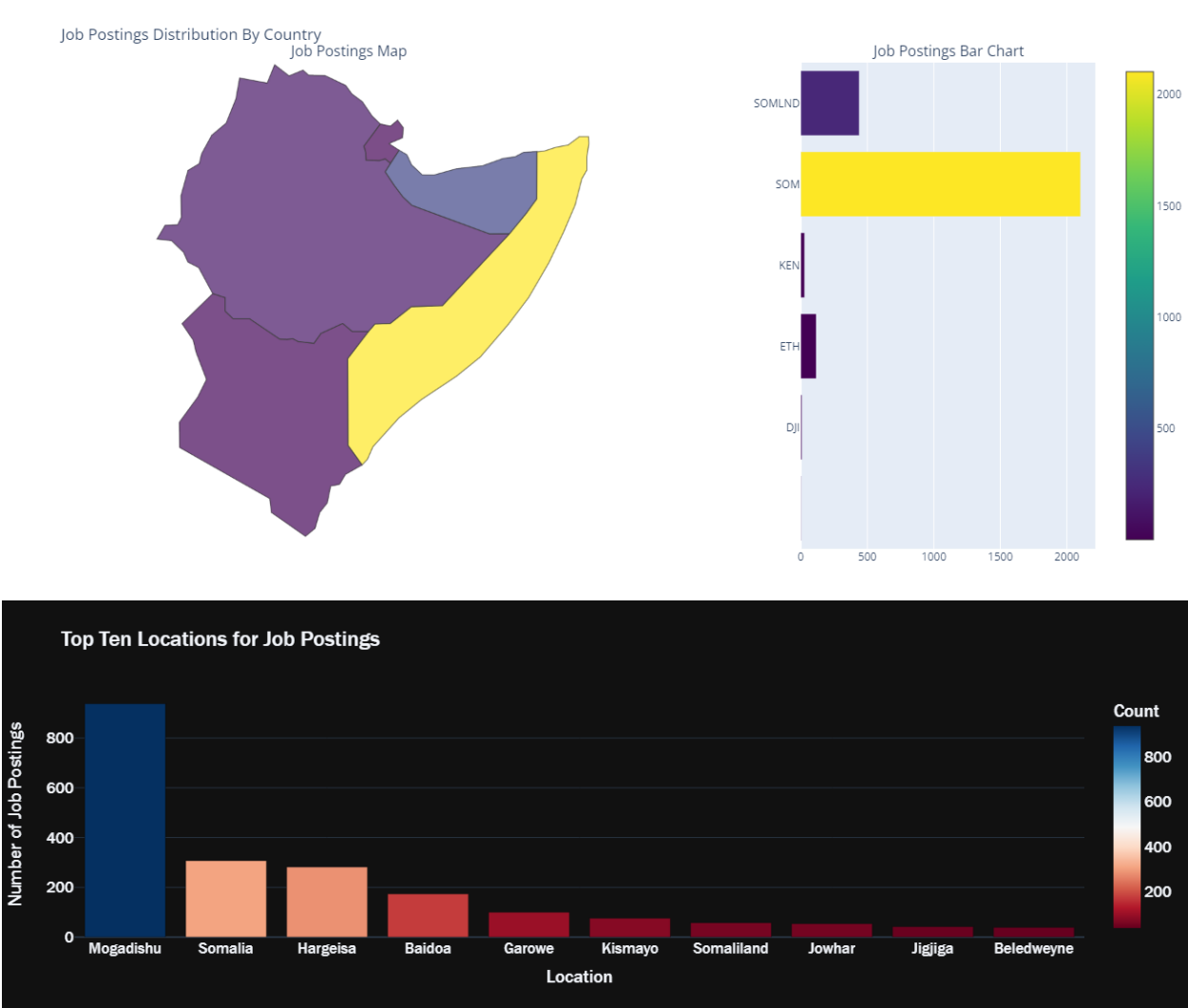**Lets visualize top 10 job titles to see which comes on top:**





The WordCloud and "Top 10 Job Titles" charts above reveals a job market geared towards finance and project management roles. 'Finance Officer' leads with 22 postings, emphasizing a significant demand in financial management. Drivers and project-related roles are also prominent, highlighting operational and administrative priorities in the sector. Supply chain and research roles, while less frequent, indicate niche areas of demand.

This WordCloud and the bar charts above "Top 10 Companies" illustrates that Save the Children, DRC, and WFP are leading employers, collectively contributing to a significant portion of job offerings. This trend highlights a robust presence of international NGOs in the job market, possibly reflecting a focus on humanitarian and developmental work. The presence of governmental roles suggests a balance of public sector opportunities amidst the NGO-dominated landscape.

# Let's visualize in what CITY/COUNTRY were the top jobs posted

Job Postings Distribution By Country

Job Postings Map



Job Postings Bar Chart



**Top Ten Locations for Job Postings**



**Job Postings Distribution by Country** Somalia dominates the job market, with a significant drop in postings from neighboring countries, indicating a centralized job market.
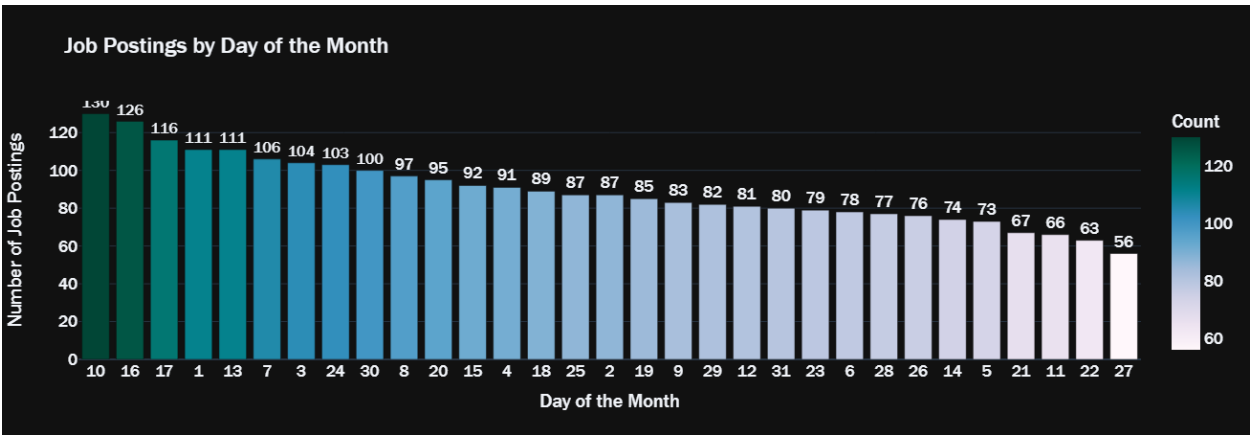
**Top Ten Locations for Job Postings**

**:** Mogadishu leads by a wide margin, followed by Hargeisa and Baidoa, which could be due to their status as economic hubs attracting a larger workforce.

**Let's visualize on most day(s) of the month and months were jobs posted?**



The Monthly Job Posting graph indicates a significant number of job postings in November 2022 and a noticeable dip by February 2024. March 2023 saw a notable uptick in job listings.



The Job Postings by Day of the Month chart above reveals the 10th as the peak day for job postings, with a gradual decrease towards the end of the month, suggesting a front-loaded posting trend within each month.
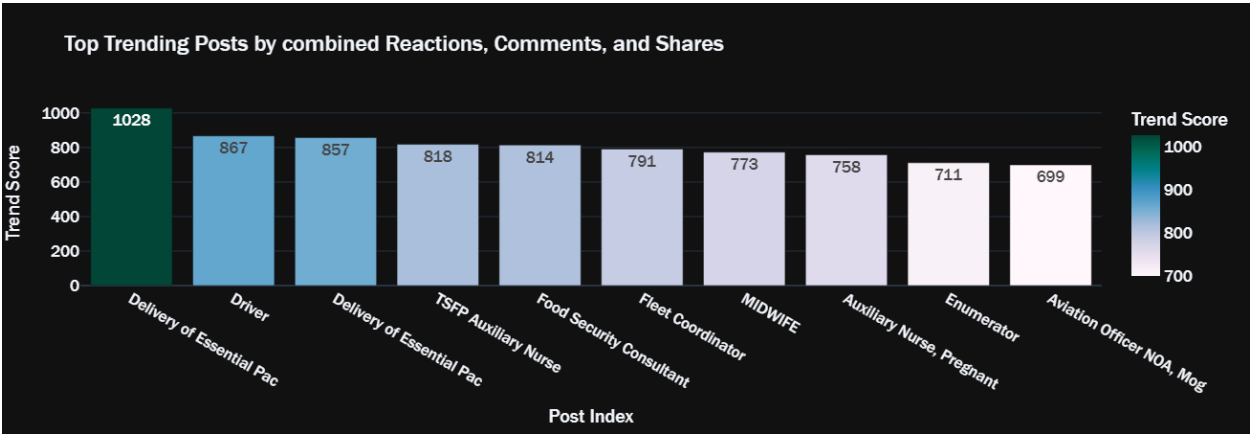
The pie chart above on STEM vs Non-STEM job distribution starkly highlights the dominance of Non-STEM jobs, constituting 98.4% of the total postings, while STEM-related jobs represent only a slender slice at 1.59%. This disparity underscores a heavy inclination towards Non-STEM opportunities in the job market.
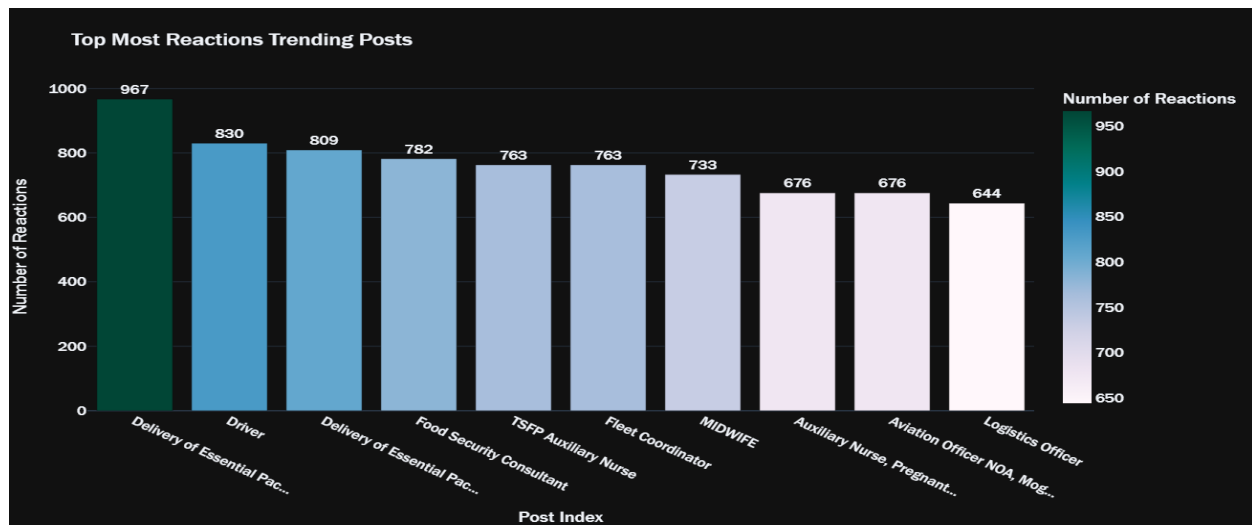
In contrast, STEM stands for Science, Technology, Engineering, and Mathematics. It refers to fields of study and sectors that focus on those disciplines. Careers in STEM typically involve research, innovation, and the application of scientific and mathematical principles to solve complex problems. Examples include software development, engineering roles, scientific research, and data analysis.

Non-STEM fields encompass all other areas that do not fall under the science, technology, engineering, and mathematics umbrella. This can include humanities, arts, social sciences, education, healthcare outside of medical research, business, law, and many other disciplines. These careers might focus more on human interaction, creativity, subjective analysis, and communication.
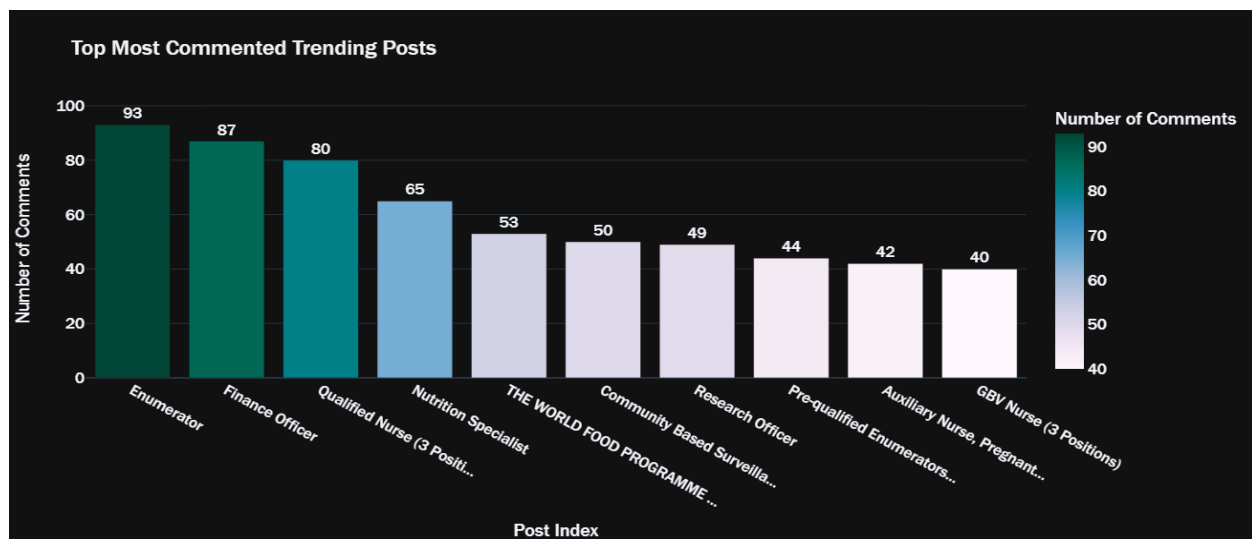
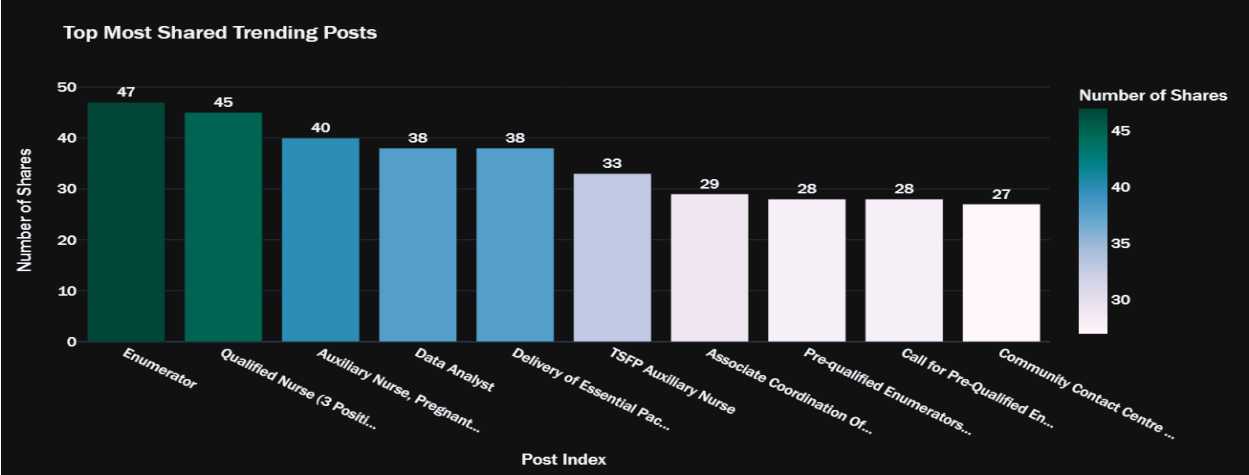**Let's analyze and visualize the most trending posts?**



The "Delivery of Essential Package of Health Service in Hiran and Middle Regions" job post significantly outperforms others in combined reactions, comments, and shares, indicating high engagement and interest. Posts about "Health", "Driver" and "Food Security Consultant" also show substantial interaction, suggesting these job titles resonate with the audience. These trending posts could reflect urgent hiring needs or particularly appealing opportunities.

**Top Most Reactions Trending Posts**

The chart above showcases the job posts that garnered the most reactions, indicating user interest and engagement. The 'Delivery of Essential package of Health Service in Hiran and Middle Regions.' post leads significantly, suggesting content relevance or urgency. Posts related to driving, food security, and healthcare also stimulated a strong response, reflecting their importance to the audience.

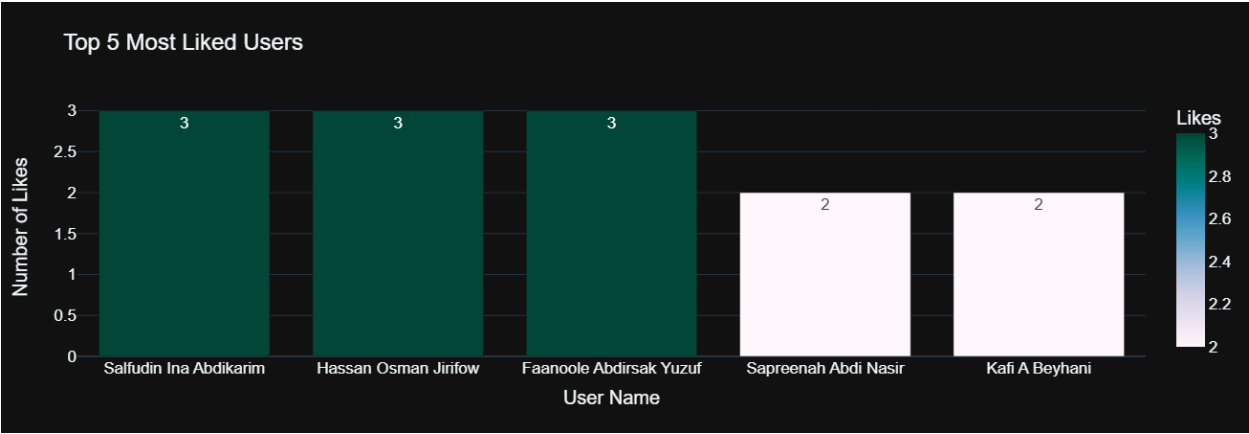

**Top Most Commented Trending Posts**

Enumerator and finance-related job posts sparked the most discussion, indicative of high interest or potential ambiguity needing clarification. Health sector jobs, like qualified nurse positions, also elicited significant commentary, possibly reflecting concerns or inquiries about qualifications and job specifics.
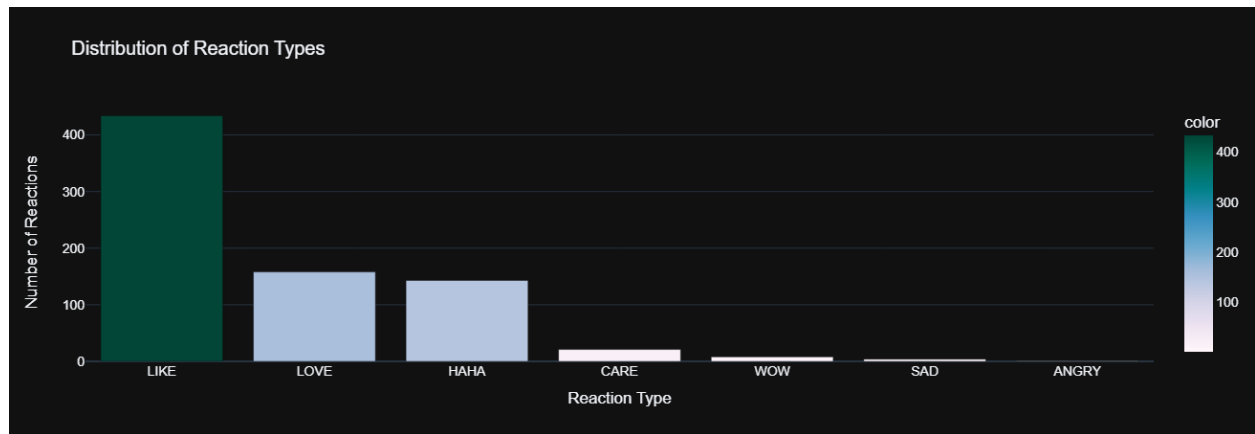
Enumerator roles led in shares, suggesting high community interest or a perceived urgency to disseminate these opportunities. Roles in nursing and data analysis were also prominently shared, indicating a community-driven effort to highlight these sectors, perhaps due to their critical importance or demand in the job market.

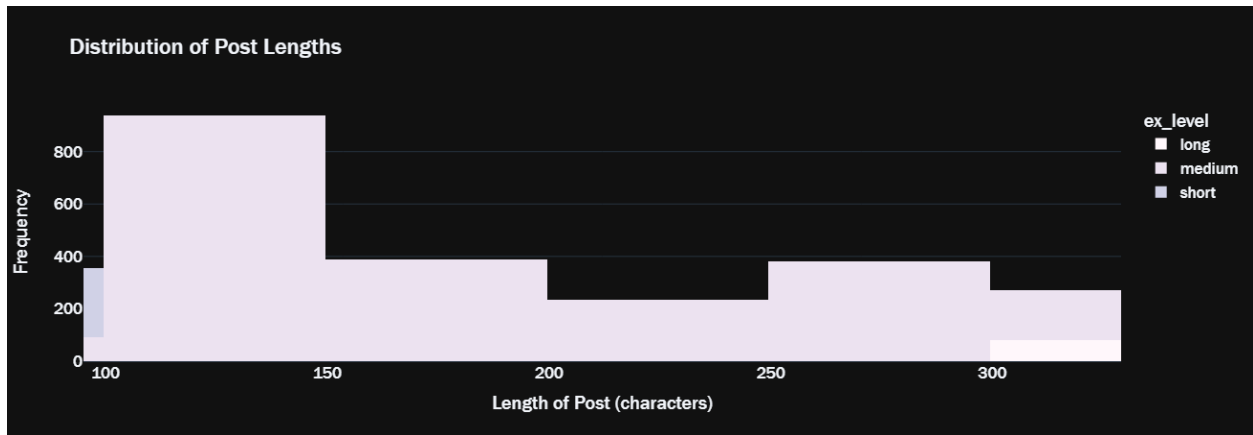## Let's analyze and visualize the top 5 users by number of likes?



"Top 5 Most Liked Users" exhibits a narrow distribution of likes among users. Three users share the highest count, each receiving three likes, suggesting a small, engaged community or possibly a focused discussion topic that garnered specific attention. Two other users follow closely, highlighting their contributions to the platform's discussions.
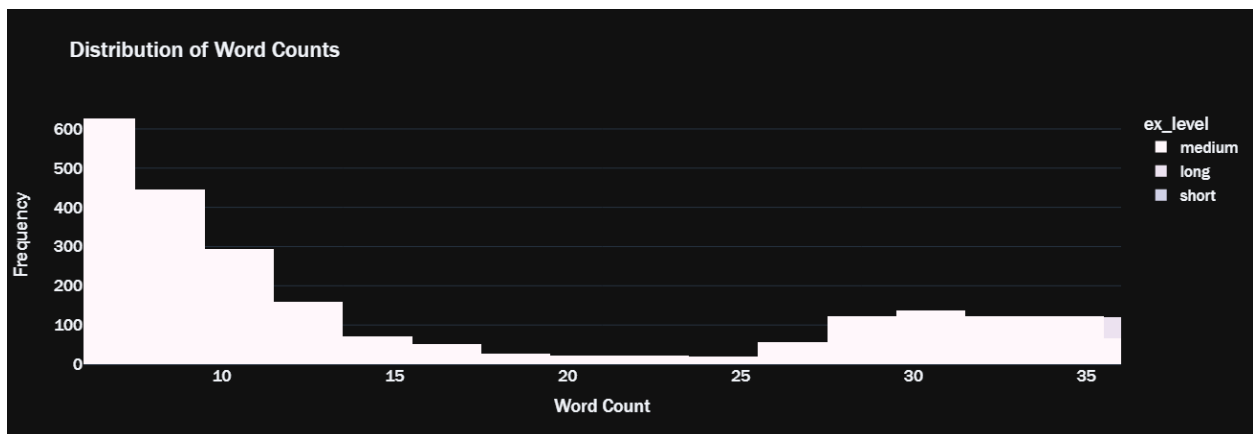
**Distribution of Reaction Types**

The "Distribution of Reaction Types" graph shows that "Like" is by far the most common reaction, indicating general positive engagement. "Love" and "Haha" have moderate occurrences, suggesting specific appreciation and humor, while "Care," "Wow," "Sad," and "Angry" are less frequent, depicting a smaller range of emotions expressed by the audience.

# Lets analyze and visualize the distribution of post length and word counts?
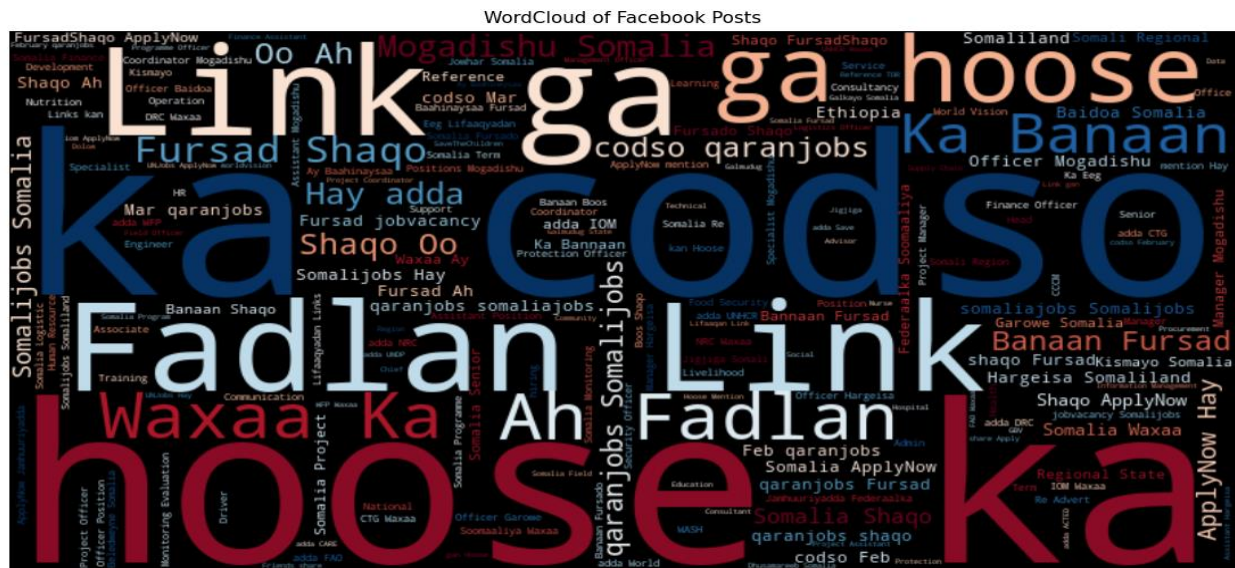


1. **Distribution of Post Lengths**: This histogram showcases the frequency of post lengths in characters. The majority of posts are relatively short, with lengths ranging around 100 characters. There's a significant drop as the character count increases, with a small resurgence for posts that are approximately 300 characters long. This suggests that most job postings are succinct, possibly only including essential information like the job title and a brief description. The presence of longer posts could indicate more detailed job descriptions or additional information pertinent to the application process.



2. **Distribution of Word Counts**: In the second histogram, we observe the distribution of word counts per post. Similar to the post length in characters, there's a high frequency of posts with fewer words, peaking around 10 words per post. The frequency diminishes as the word count increases, with fewer posts containing 20 to 35 words. This graph further

corroborates the trend towards brevity in job postings on the platform. It could be inferred that job posters prefer concise descriptions that are likely easy to read and quick to comprehend for potential applicants.

**Let's visualize the prevalent words in posts using WordCloud?**



WordCloud of Facebook Posts

The word cloud provide a visual representation of the most frequently occurring words in Facebook posts and comments. The first word cloud shows frequent use of words like "Shaqo" (job), "Fursad" (opportunity), and "qaranjobs," highlighting a focus on employment and the QaranJobs platform.

WordCloud of Comments

In this word cloud, personal names dominate, suggesting active community engagement with individuals often personally addressed or mentioned in the comments.

# Model Training

In the final phase of the project, we concentrate on data preparation and model training. We amalgamate datasets from different sources to create a comprehensive corpus for sentiment analysis. Specifically, we collate four separate datasets, labeled **df1** to **df4**, each comprising various segments of scraped social media labeled data.

A crucial step in data pre-processing involves standardizing labels to ensure uniformity across datasets. This is achieved by defining a mapping that associates various label formats with standard categories of sentiment, namely 'Negative-Xumaan', 'Positive-Wanaag', and 'Neutral-Dhexdhexaad'. Applying this mapping corrects inconsistencies and homogenizes sentiment labels across all dataframes.

Subsequent to label normalization, we conduct a sentiment distribution analysis. The **value_counts** method reveals the sentiment distribution across the datasets, with 'Negative-Xumaan' leading, followed closely by 'Neutral-Dhexdhexaad' and 'Positive-Wanaag'. This distribution suggests a balanced dataset, which is crucial for unbiased model training.

```
ml_data_df.head()
✓ 0.0s                                                                                      Python
```

| | Qoraalka | Qeybta |
|---|---|---|
| 0 | Qaadka Kenya ee Deni u diiday Puntland miyuu ahaa ujeedka socdaalka XSM Dulaal Qaad miyeynu Dooranay Belo ku ciyaadhay | Neutral-Dhexdhexaad |
| 1 | Hassan Sheikh ayaa ka qaatay Aduun garaya 5 Milion Dalka Kenya si uu Qaadka igu Fasaxo | Neutral-Dhexdhexaad |
| 2 | Garoonka Loftus Road Stadium Garoonkaan waxaa ku ciyaarta kooxda sanadkii lasoo dhaafi loo celiyay horyaalka xeerka 3xaad ee Ingiriiska Kooxda Queen Park Rangers Gobalka Shepherd s Bush dalka Ingiriiska La furay sanadii 1904 Qaadka taageero 18 439 | Neutral-Dhexdhexaad |
| 3 | Xaqbuu u leeyahay inuu qayilo qofkaasi illeen wuu ina iibsan hadii kale e Nin murqaamay oo Libaax ka sameeyay Caleentii qaadka ee uu cunayay kkkk X Saalax Sheekh Dayax | Positive-Wanaag |
| 4 | Qaadka ay Itoobiya inoo soo dhoofiso | Neutral-Dhexdhexaad |

```
ml_data_df.shape
✓ 0.0s
```

(20385, 2)

```python
ml_data_df.info()
```
✓  0.0s

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20385 entries, 0 to 20384
Data columns (total 2 columns):
 #   Column   Non-Null Count  Dtype
---  ------   --------------  -----
 0   Qoraalka  20382 non-null  object
 1   Qeybta    20385 non-null  object
dtypes: object(2)
memory usage: 318.6+ KB
```

```python
ml_data_df['Qeybta'].value_counts()
```

```
Negative-Xumaan       7035
Neutral-Dhexdhexaad   6811
Positive-Wanaag       6539
Name: Qeybta, dtype: int64
```

## Model Evaluation

We deploy cross-validation to assess the model's robustness and generalizability. The 5-fold cross-validation yields an average accuracy of approximately 70.13%, indicating consistent model performance across different subsets of the data.

The final test accuracy stands at 71.53%, with the model demonstrating higher precision in predicting the majority class 'Negative-Xumaan'. The recall for 'Positive-Wanaag' is relatively low, suggesting the model is more conservative in predicting positive sentiments. The F1-scores further elucidate this, with a significantly higher score for the 'Negative-Xumaan' class compared to 'Positive-Wanaag'.

```
CV Accuracy (5-fold): [0.71098266 0.70872766 0.69348055 0.69873817 0.69453207]
Mean CV Accuracy: 0.7012922205405931


Test Accuracy: 0.7153069806560135
[[1558   71]
 [ 606  143]]
              precision    recall  f1-score   support

           0       0.72      0.96      0.82      1629
           1       0.67      0.19      0.30       749

    accuracy                           0.72      2378
   macro avg       0.69      0.57      0.56      2378
weighted avg       0.70      0.72      0.66      2378
```
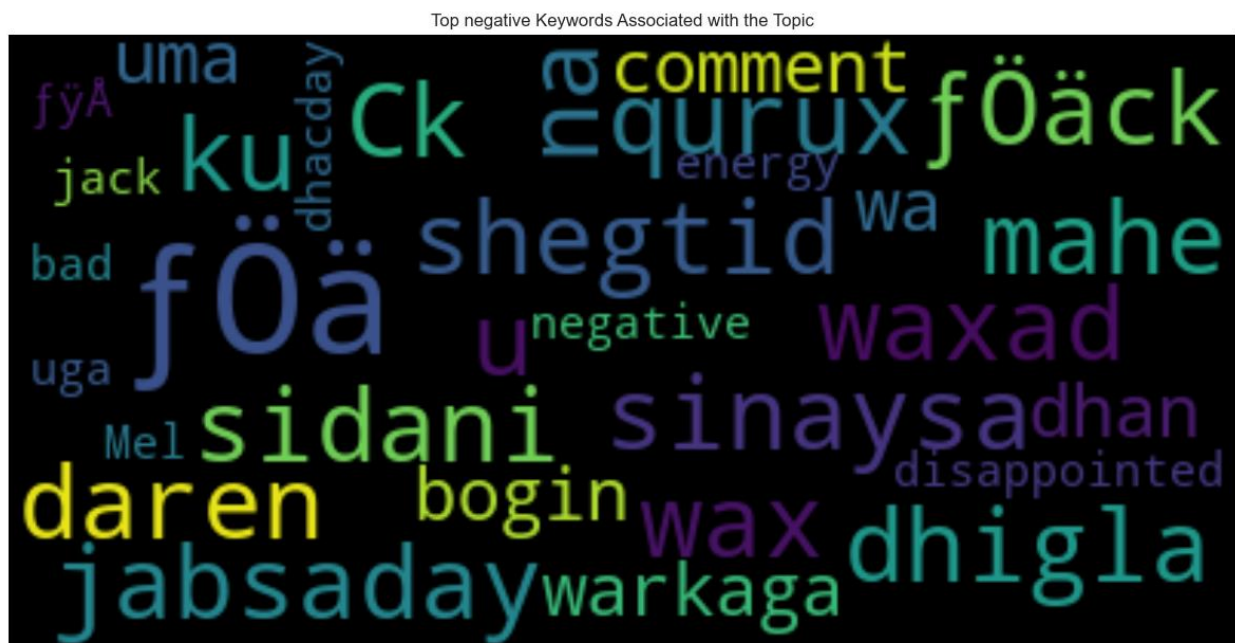
**- Let's visualize the sentiment percentages (positive, negative, neutral)**



Data label Distribution.

**Let's visualize the top positive keywords associated with the data**



**Let's visualize the top negative keywords associated with the data**



Top negative Keywords Associated with the Topic

# CONCLUSION

The culmination of the 'Comprehensive Job Market Analysis in Somalia' project underscores the intricate dance between employment opportunities and community engagement within the digital ecosystem. With a dataset amalgamating social media engagement metrics and job postings, the project casts a spotlight on the prevalent job market trends and patterns that signify economic vitality and social responsiveness.

The analysis embarked on a journey through a treasure trove of data, encompassing 4,000 job postings from Qaranjobs.so and a substantive 5,000 posts from Facebook, reflecting the interplay between online presence and job market dynamics. Through meticulous data cleaning and integration processes, a robust dataset was curated to embody the pulse of the job market in Somalia.

Our model, built on the back of Logistic Regression, demonstrated a commendable predictive capability, registering a mean cross-validation accuracy of approximately 70% and a testing accuracy that slightly surpasses this mark. This attests to the model's efficacy in gleaning insights from the data while being mindful of overfitting.

Diving into sentiment analysis, the project has painted a tableau of public sentiment, with posts ranging across the spectrum from positive to neutral, and negative. Each sentiment echoed the community's voice, offering a quantitative testament to the collective hopes, neutrality, and challenges faced by job seekers.

In conclusion, the project has not only provided a quantitative assessment of the job market but also qualitatively enriched our understanding of the relationship between job availability and community engagement. It has bridged the data analysis with human experience, where numbers met narratives. The insights gained hold the potential to inform policymakers and businesses alike, shaping strategies that align with the evolving job landscape. The project stands as a beacon, guiding towards a data-driven approach in understanding and optimizing the job market for economic growth and social betterment.

# END