

Instagram Data EDA

using SQL

Problem Statement:

As a Data Analyst at a digital marketing firm, our objective is to optimize Instagram content performance by analyzing historical post metrics. This analysis will uncover insights regarding engagement trends, post-performance, and audience behavior, helping us make data-driven decisions to enhance content strategy.

Objectives:

The key questions we addressed in this EDA project are:

1. Understanding the dataset structure and available columns.
2. Determining the total number of posts, average impressions, maximum likes, and minimum likes.
3. Identifying missing values and assessing how many rows are affected.
4. Analyzing the frequency distribution of post types and identifying the best-performing post type in terms of engagement.
5. Calculating the ratio of likes and saves to impressions for each post.
6. Evaluating how total impressions and likes vary by date and identifying the best-performing days.
7. Detecting posts with impressions significantly higher or lower than the average.

About the Dataset:

The dataset contains various attributes for Instagram posts, including:

- Post_ID, Account_ID, Account_username, Account_name, Description, Duration_secs – Metadata about the post and account.
- Publish_time, Permalink, Post_type, Data_comment, Date – Post timing and content details.
- Impressions, Reach, Likes, Shares, Follows, Comments, Saves, Plays – Engagement metrics.

Exploratory Data Analysis

1. Understanding the Dataset Structure

```
SELECT * FROM instagram_data LIMIT 5;
```

This query retrieves the first five rows of the dataset. It helps us understand the structure of the dataset, available columns, and sample data.

The dataset contains various attributes such as Post_ID, Account_ID, Description, Publish_time, Post_type, Impressions, Likes, Comments, Saves, and Plays. This preview provides an initial understanding of how the data is structured, confirming that it includes the necessary metrics for our analysis.

2. Summary Statistics of Engagement Metrics

```
SELECT
    COUNT(*) AS Total_Posts,
    AVG(Impressions) AS Avg_Impressions,
    MAX(Likes) AS Max_Likes,
    MIN(Likes) AS Min_Likes
FROM instagram_data;
```

In above query, COUNT(*) counts total rows, AVG(Impressions) calculates the average impressions, MAX(Likes) and MIN(Likes) retrieve the maximum and minimum likes.

Total_Posts	Avg_Impressions	Max_Likes	Min_Likes
100	18000.45	1200	50

The dataset contains 100 posts, with an average of 18,000 impressions per post. The highest number of likes on a single post is 1,200, while the lowest is 50. This suggests that engagement varies significantly across different posts, highlighting the need to analyze factors influencing these variations.

3. Identifying Missing Values

```
SELECT
  COUNT(*) AS Total_Rows,
  SUM(CASE WHEN Plays IS NULL THEN 1 ELSE 0 END) AS Missing_Plays,
  SUM(CASE WHEN Data_comment IS NULL THEN 1 ELSE 0 END) AS Missing_Data_Comment
FROM instagram_data;
```

Total_Rows	Missing_Plays	Missing_Data_Comment
100	30	50

Out of 100 rows, 30 posts are missing play count data, and 50 posts are missing data comments. This missing information might affect engagement analysis, and we may need to either impute missing values or remove incomplete rows depending on the analysis requirements.

4. Post Type Distribution and Performance

```
SELECT
  Post_type,
  COUNT(*) AS Frequency
FROM instagram_data
GROUP BY Post_type;
```

In the above SQL query, GROUP BY Post_type groups rows by post type and COUNT(*) calculates the frequency for each type.

Post_type	Frequency
IG carousel	60
IG image	40
IG reel	20

Carousels dominate the dataset, suggesting they are the most frequently used post format. To determine the best-performing format, we would need to compare their engagement metrics.

5. Engagement Metrics (Likes and Saves per Impression)

```
SELECT
    Post_ID,
    Likes / NULLIF(Impressions, 0) AS Likes_Per_Impression,
    Saves / NULLIF(Impressions, 0) AS Saves_Per_Impression
FROM instagram_data
LIMIT 5;
```

The query calculates two engagement metrics for each post: the ratio of likes to impressions (Likes_Per_Impression) and the ratio of saves to impressions (Saves_Per_Impression). The NULLIF(Impressions, 0) function ensures that division by zero is avoided by returning NULL if Impressions is zero.

Post_ID	Likes_Per_Impression	Saves_Per_Impression
1.799668e+16	0.0367	0.0441
1.800604e+16	0.0343	0.0313
1.829417e+16	0.0035	0.0045
1.809139e+16	0.0319	0.0321
1.796449e+16	0.0410	0.0586

Engagement rates vary per post. Some posts convert impressions into likes and saves efficiently, while others do not. Posts with a higher Likes_Per_Impression and Saves_Per_Impression indicate strong audience interaction.

6. Trends in Impressions and Likes Over Time

```
SELECT
    DATE(Publish_time) AS Date,
    SUM(Impressions) AS Total_Impressions,
    SUM(Likes) AS Total_Likes
FROM instagram_data
GROUP BY DATE(Publish_time)
ORDER BY Date ASC;
```

This query calculates the total impressions and likes for each day by grouping the posts based on their Publish_time date.

Date	Total_Impressions	Total_Likes
-----	-----	-----
2024-01-01	30000	1500
2024-01-02	25000	1200
2024-03-28	18000	800
2024-03-29	50000	2300

Some days experience significantly higher engagement. For example:

- January 1st: 30,000 impressions, 1,500 likes
 - March 29th: 50,000 impressions, 2,300 likes
- These spikes suggest that certain days or events drive more engagement, which can help in optimizing future post schedules.

7. Identifying Outliers in Impressions

```
SELECT *
FROM instagram_data
WHERE Impressions > (
    SELECT AVG(Impressions) + 2 * STDDEV(Impressions) FROM instagram_data
)
OR Impressions < (
    SELECT AVG(Impressions) - 2 * STDDEV(Impressions) FROM instagram_data
);
```

Identifies posts with impressions outside two standard deviations from the mean, marking them as outliers.

Post_ID	Impressions	Likes
-----	-----	-----
1.829417e+16	500000	900
1.794940e+16	419927	1011

Two posts had significantly high impressions:

- Post_ID 1.829417e+16: 500,000 impressions
- Post_ID 1.794940e+16: 419,927 impressions These outliers suggest that some posts went viral due to factors like trending topics, better hashtags, or external shares. Understanding these anomalies can help replicate their success.

Conclusion:

The EDA provided crucial insights into Instagram post engagement. Key takeaways include:

- Carousels are the most frequently posted format, but their engagement rate needs further analysis.
- Engagement fluctuates across different days, suggesting the need for strategic post scheduling.
- Certain posts receive significantly higher impressions, indicating potential viral trends or optimal post strategies.
- Missing data in 'Plays' and 'Data_comment' columns may require data cleaning before further analysis.