# INSTAGRAM DATA ANALYSIS

**¹Poonam Narkhede, ²Sameksha Moolya, ³Priti Molawade, ⁴Shrusti Chavan**

¹ Assistant Professor, Shivajirao S. Jondhale College of Engineering, Thane, India

²⁻⁴Student, Department of Computer Engineering, Shivajirao S. Jondhale College of Engineering, Thane, India,

*Abstract:* Because of the spread of the Internet, social platforms big data pools. there we can learn the trends, culture, and hot topics. This project focuses on analysing the data from Instagram. It shows the relationship of Instagram filter data with location and number of likes to give users filter suggestion on achieving more likes based on their location. It analyses the popular hashtags in different locations to show visual culture differences between the cities. We explore Kaggle, which provides us a dataset for our real-world project. we grab Instagram data analysis for post reach and for followers' information. And also, we have explored Power BI tool form visualization purpose which give use proper idea of analysis! It connects to various data sources, create interactive dashboards, and generate reports to gain insights for data. Power BI is used to transform raw data into meaningful visualizations that help organizations make data-driven decisions and analyse vast amounts of user-generated content, including images, captions, likes, comments, and more. We have also worked upon JUPYTER NOTEBOOK in VS CODE which give us an easy platform to work on and is a user friendly. In JUPYTER NOTEBOOK we work with many libraries like NumPy, pandas, Matplotlib, and Seaborn, making us to work easily and efficiently. Overall, in this project we have worked upon various analysis on the topic likeFollowers Analysis, Post reach which will have various sub analysis of the post reach and follower which will have mutual connection. How much likes, is an account is private or not a private/public and many more related to followers!

*Keywords-* Instagram, Analysis, Post, Followers, Post Reach, Comments, Private, Public, Business account.

## INTRODUCTION

The American business Meta Platforms is the owner of Instagram, a social media platform for sharing photos and videos. Users of the app are able to upload media that may be altered with filters and arranged using hashtags and geotagging. Public or followers who have already been approved can see posts. Users can view trending content, like photos, followother users to add their content to a personal feed, browse other users' content by tag and location, and browse other users' content by location and tag. In order to fit the width of the iPhone display at the time, Instagram's original design limited how content could be framed to a square (1:1) aspect ratio of 640 pixels. With an upgrade to 1080 pixels in 2015,this restriction was loosened. Additionally, it added messaging capabilities, the capacity to post multiple images or videos, and a feature called Stories that allowed users to post content to a feed in chronological order with each entry being viewable by others for a 24-hour period. Stories was similar to Snapchat, the platform's main rival. There are500 milliondaily users of Stories as of January 2019 [1]. The development of society at the present time is impossible without the use of modern information technologies, which over the past half century have dramatically changed our lives and openedup new market opportunities. The development of information technologies was so rapid that they have in less than the last 30 years overcome the same path as traditional production technologies, for example, metallurgical ones, over 300 years. The role of information technologies has become so important that they first turned the traditional economy into information, then into collaborative and, finally, into digital. Based on this chain of transformations, is it possible to consider that the digital economy is a simple sum of the traditional economy and information technologies? Probably not. The term "digital economy" arose a long time ago, namely in 1994, and is associated with the publication of the book "The Digital Economy [2].

Instagram is a popular social network platform that allows users to edit and upload photos and short videos using a mobile app. Users can add captions, hashtags, and location based geotags to make their posts searchable by other users within the app. Instagram is not only a social tool for individuals but also for businesses that use the platform to promote their brand and products. Companies with business accounts have access to free engagement and impression metrics to measure their performance. One of the most important metrics for businesses on Instagram is the reach of a post, whichrefers to the number of people who see the post. Higher reach leads to more engagement with the post, such as likes, comments, and shares, and ultimately more exposure for the business. However, the platform's algorithms determine which posts are shown to which users, making it challenging for businesses to achieve high reach. Understanding the

factors that influence the reach of a post on Instagram and being able to predict future reach is crucial for businesses looking to maximize their visibility on the platform. Recent research has focused on analysing the reach and interactionsof posts on Instagram and developing methods to improve reach [3].

Big Data Analytic

Big data analytics is the process of examining and analysing large and complex datasets to uncover valuable insights, patterns, trends, and information that can be used to make informed business decisions, optimize processes, and gain a competitive advantage. It involves the use of advanced technologies and techniques to process and interpret massive amounts of data from various sources, often in real-time or near-real-time, to drive data driven decision making [4].

Power BI

Power BI is a business intelligence tool by Microsoft that helps you analyse and visualize data from various sources. It simplifies the process of creating interactive reports and dashboards, allowing users to gain insights and make data-drivendecisions. Power BI connects to data, transforms it, and presents it through visually appealing charts and graphs, makingcomplex data more accessible to wide range of users [5].

## 1.1 Needs

•Instagram data analysis serves several important needs, both for businesses and individuals, by providing valuable insights into user behaviour, content performance, and trends. Here are some key needs for Instagram data analysis:

•Audience Understanding: Understanding your Instagram audience is crucial for businesses and content creators. Data analysis helps in identifying the demographics, interests, and behaviours of your followers, allowing you to tailor your content to their preferences and needs.

•Content Strategy: Analysing data can help determine which types of content (photos, videos, stories, reels, etc.) resonate best with your audience. It also helps you discover which captions, hashtags, and posting times are most effective.

•User Engagement: Knowing which posts receive the most likes, comments, shares, and saves can help you gauge audience engagement. You can also identify what content triggers conversations and discussions.

•Competitive Analysis: By monitoring the performance of competitors' content, you can gain insights into what works in your industry or niche. This can inform your own content strategy and help you stand out.

•Influencer Partnerships: Brands often collaborate with influencers on Instagram. Data analysis helps identify influencerswhose followers align with the brand's target audience, and it can also measure the effectiveness of influencer partnerships.

•Trend Identification: Data analysis reveals emerging trends, popular hashtags, and challenges on Instagram. Businessescan leverage these trends to create relevant content and connect with their audience.

•User Experience Improvement: By analysing user interactions, Instagram can enhance the platform's user experience. This includes improving algorithms for content recommendations and user interface design.

## 1.2 Applications

•Safety and Moderation: Employ data analysis to detect and filter out spam, hate speech, and inappropriate content toensure a safe and positive user experience on your account or platform.

•Content Recommendation: For social media platforms, use data analysis to recommend content to users based on theirpreferences and past interactions, enhancing the user experience.

•Research and Insights: Academics and researchers can use Instagram data analysis to study trends, behaviour, and socialphenomena, providing valuable insights and data for various fields of study.

•Personalization: Customize the Instagram user experience by providing personalized content recommendations, helpingusers discover content relevant to their interests.

•Fraud Detection and Security: Identify and mitigate fraudulent activities, such as fake accounts, spam, and cyber threats.

•User Engagement Campaigns: Design engagement campaigns and challenges that encourage user participation andinteraction.

## II. LITERATURE SURVEY

The study paper [1] Influencer Marketing on Instagram (2022): This section references an empirical research article oninfluencer marketing on Instagram. The research uses various methods, including content analysis, binomial regression, data collection, and categorization of influencers and advertising appeals is design by Taylor & Francis Group, LLC and it's disadvantage is Lack of Emotional Appeal, Potential for Boredom, Limited Creativity May Not Suit All Products, competition Based on Features.

The second paper [2] Research in the Instagram Context (2021): This part briefly mentions a research study in the Instagram context conducted in Hong Kong was designed by Hung Hom, Kowloon, Hong Kong disadvantage is requirement to follow to view content and the potential weaknesses in privacy settings.

The third paper [3] Social Media Visualization (2021): The last section discusses the development of a framework andtools for predictive analysis, personalization, and scalability in social media visualization and was designed by Sharmaand Jain It notes challenges related to complexity and cost, as well as data quality.

The fourth system [4] The big picture on Instagram research: Insights from a bibliometric analysis (2021) This article,published in 2021, is part of Elsevier's Web of Science database. It employs bibliometric analysis, citation analysis, collaboration analysis, and visualization techniques. The study's focus is on providing a comprehensive understandingof research related to Instagram by using quantitative insight and visual representation. The goal is to objectively assess the field of Instagram research and identify any research gaps.

## III. METHODOLOGY

The methodology for Instagram data analysis is a structured approach aimed at extracting valuable insights from the wealth of data generated on the Instagram platform. It begins by setting clear objectives, defining specific goals for the analysis, and then proceeds to collect data from various sources, encompassing user profiles, posts, comments, and more.Data preprocessing ensures that the collected data is clean and structured, while exploratory data analysis provides initialinsights into user behaviour and content performance. Following data transformation and feature engineering, the analysis phase involves the application of appropriate methods like sentiment analysis or user behaviour modelling to derive insights, and these findings are effectively communicated through data visualization and reporting. Privacy and ethical considerations are paramount throughout, and continuous monitoring and adaptation are crucial to keeping the analysis aligned with Instagram's evolving landscape. Ultimately, the insights generated through this methodology informdecision-making processes and content optimization strategies.

Incorporating these insights into decision-making and thoroughly documenting the analysis process are fundamental components of this methodology, facilitating knowledge sharing and promoting data-driven decision-making within the organization. This structured approach ensures that the analysis remains focused on objectives and ethical data usage, enabling organizations to harness the full potential of Instagram data for informed decision making and strategic planning.

## IV. EXISTING SYSTEM

An existing software system is any software application that is currently in use. It includes everything from newly released software to those that have existed for years. System analysis refers to the process of gathering data, interpretinginformation, identifying issues and using the results to recommend or develop possible system improvements. During this stage, companies may also evaluate future business needs and how improvements may answer them.
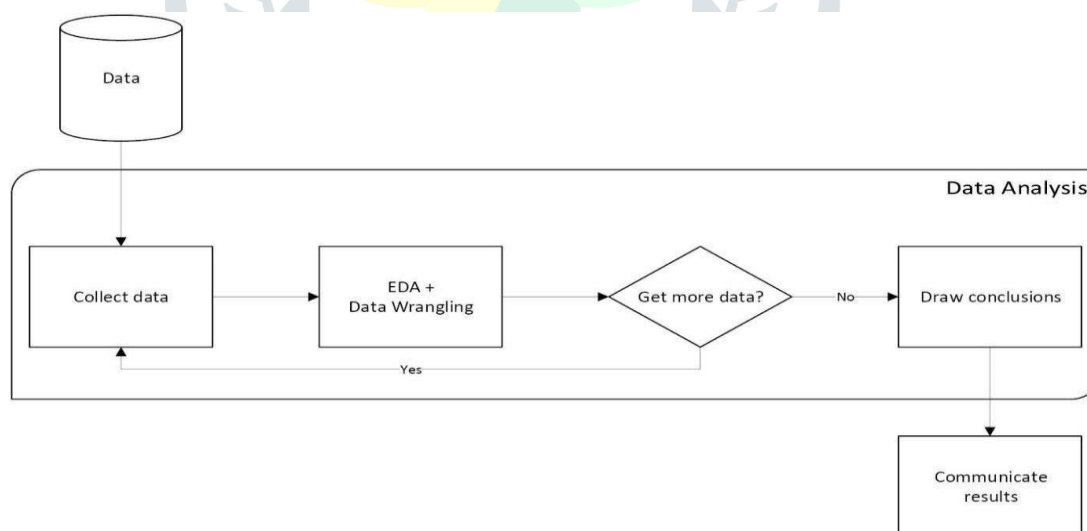
Figure 1. Architecture of existing system

From the above Figure 1 System analysis refers to the process of gathering data, interpreting information, identifying issues and using the results to recommend or develop possible system improvements. During this stage, companies mayalso evaluate future business needs and how improvements may answer them. Users can exchange photographs,

documents, user locations, and other content in addition to sending text and voice messages, making audio and video conversations and sharing messages. They are able to upload text, images, and other types of media that are shared openlyor, depending on the privacy settings, with any other users who have accepted to be their "friends." Additionally, users have direct access to one another, have the option to join groups with similar interests, and can subscribe to notificationsfor their friends' and their favourite sites' actions.

Users can post pictures and quick videos, subscribe to other users' feeds, and geotag pictures with the name of a place. Users have the option to make their accounts "private," which makes it necessary for them to consent to any new followerrequests. Users can share posted photographs to other social networking websites by linking their Instagram accounts toother platforms. The programme was updated with new and live filters, quick tilt-shift, high resolution images, optionalborders, one-click rotation, and a new icon. Photos were formerly limited to a square, 1:1 aspect ratio; however, the appalso supports portrait and widescreen aspect ratios. Previously, users could examine a map of a user's geotagged pictures[1].

## V. Classification

We used the `Instaload` Python module to collect data from Instagram. We scraped data for a set of users and saved it ina CSV file. The data included attributes such as username, number of followers, number of posts, likes, and time of posting of the last 10 posts. We collected a total of 1000 records for our analysis**.**
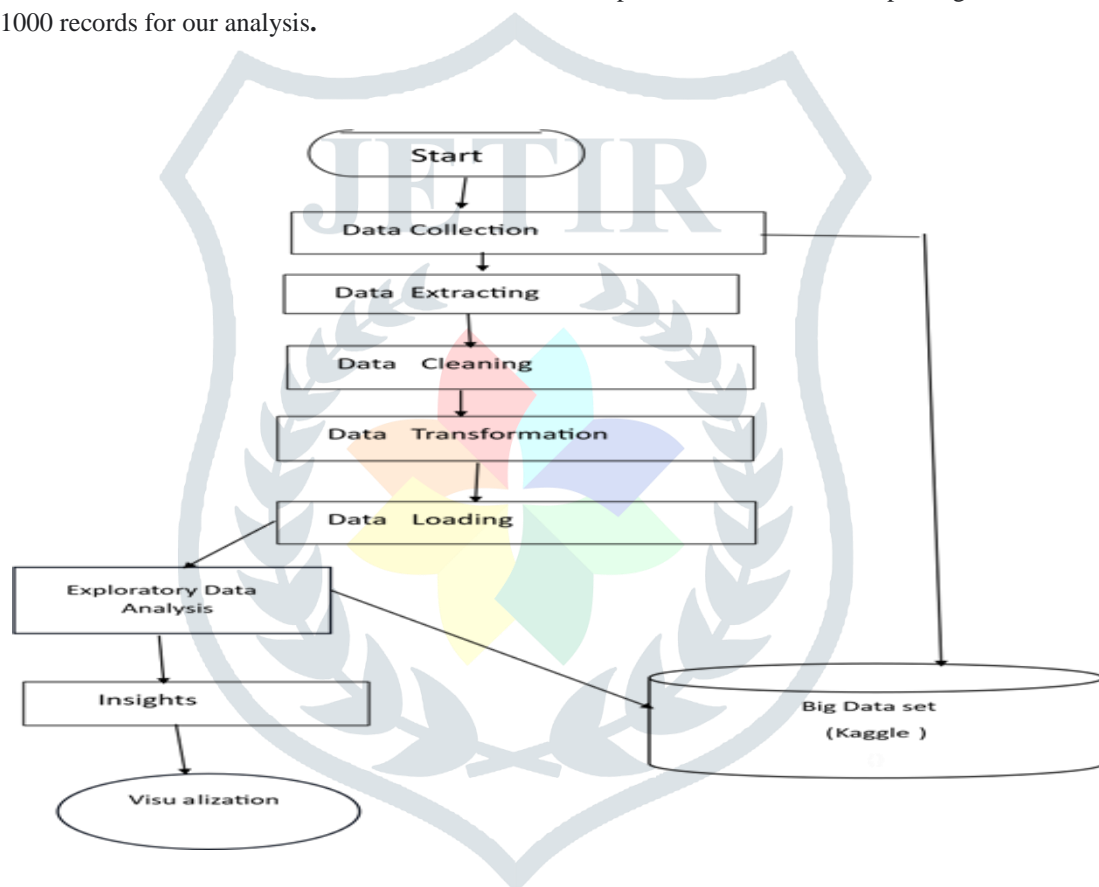


Figure 5.1. Architecture of Proposed System

**Data Preprocessing:** Before training the linear regression model, we pre-processed the data to ensure that it was in a suitable format for analysis.

**Data Cleaning:** We removed records with missing or invalid data.

**Data Transformation:** We transformed categorical variables such as username into numerical values.

**Kaggle Dataset:** A Kaggle dataset is a collection of structured data available on Kaggle, a popular platform for data science and machine learning competitions. These datasets cover a wide range of topics and are provided by the Kaggle community, researchers, and organizations. They serve as valuable resources for data analysis, model training, and research. Kaggle datasets can include data in various formats, such as CSV, JSON, or SQL, and are often accompanied

by descriptions, context, and potential research questions. Users on Kaggle can explore and download these datasets fortheir data science projects and machine learning experiments.

**Exploratory Data Analysis:** Exploratory Data Analysis, is a critical process in data analysis that involves examining and visualizing data to understand its key characteristics, uncover patterns, identify anomalies, and generate initial insights. EDA helps analysts and data scientists gain a preliminary understanding of the data's distribution, relationships, and potential issues before more in-depth analysis or modelling. It often includes tasks such as data summary statistics, data visualization, and basic statistical analysis to inform subsequent data-driven decisions and strategies.

**Data visualization:**

Data visualization is the practice of representing data and information graphically through charts, graphs, and othervisual elements. It helps make complex data more understandable and accessible by presenting it in a visual format, allowing trends, patterns, and insights to be quickly identified. Data visualization is widely used in data analysis, reporting, and decision-making to convey information in a clear and intuitive manner, making it a powerful tool for bothprofessionals and non-experts to comprehend and interpret data. Data visualization means turning numbers and facts intopictures or graphs. It's like using pictures to show information, making it easier for people to understand and see patternsin the data. This helps us make better decisions and tell stories with data, whether it's about sales, weather, or anything else.

## VI.    System Design

With reference to Fig 6.1 our system design takes a data as a input from a Kaggle dataset then next step performed is data preprocessing in this step data is resized, rescale, noise and duplicates is removed from input data set. After the datapreposing, we categorized the data in different category like user accounts, Followers, Post Reach, Likes, Comments, Locations, Gender and many more in this. After the data categorising, we will visualize the data that we have extracted from dataset. We did the visualization from the power bi dashboard tool and got the expected results. Visualization is theprocess of representing data or information using visual elements like charts, graphs, maps, or images to make complexconcepts or patterns more understandable and accessible to humans. It helps in gaining insights, spotting trends, and communicating                information                effectively                through                visual                representations.
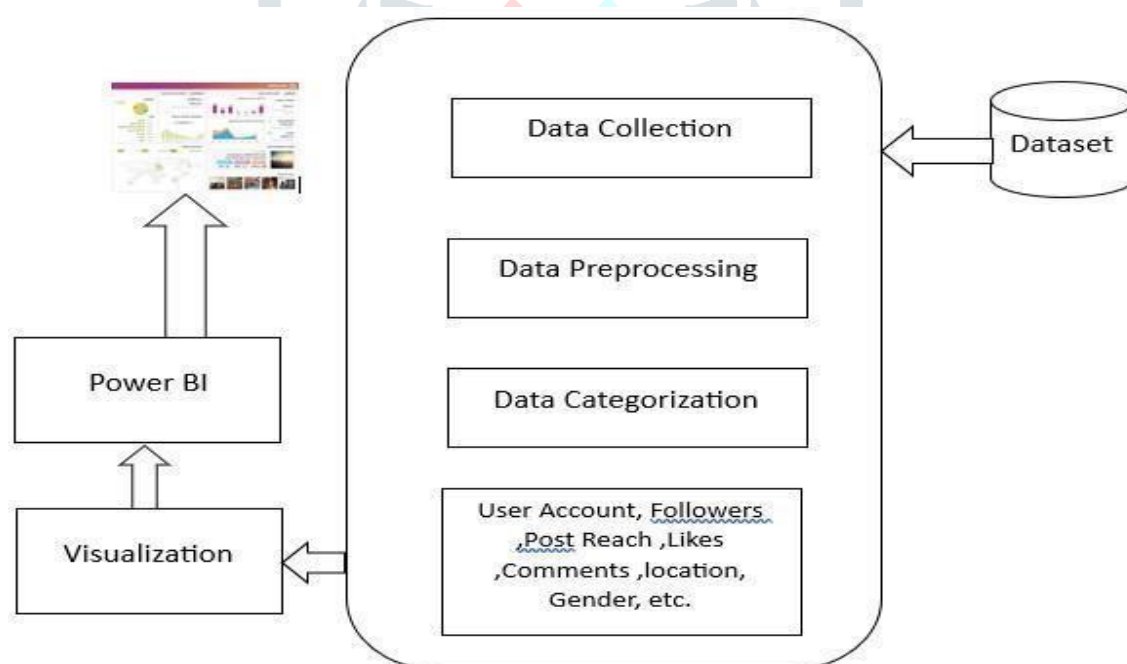


Figure 6.1.a. System Design

Our system design takes a data as a input from a Kaggle dataset then next step performed is data preprocessing in this step data is resized, rescale, noise and duplicates is removed from input data set. After the data preposing, we categorizedthe data in different category like user accounts, Followers, Post Reach, Likes, Comments, Locations, Gender and many

more in this. After the data categorising, we will visualize the data that we have extracted from dataset. We did the visualization from the power bi dashboard tool and got the expected results.

### 6.1. Why Random Forest?

The Random Forest classifier is a powerful ensemble learning method comprising multiple decision trees trained on different subsets of the dataset. Rather than relying on a single decision tree, it aggregates the predictions from each tree and uses a voting mechanism to determine the final output. By leveraging the collective wisdom of multiple trees, it typically improves predictive accuracy compared to individual trees. In contrast, the Passive Aggressive model, while effective in certain scenarios, may not match the accuracy achieved by Random Forest due to its different underlying approach and characteristics.
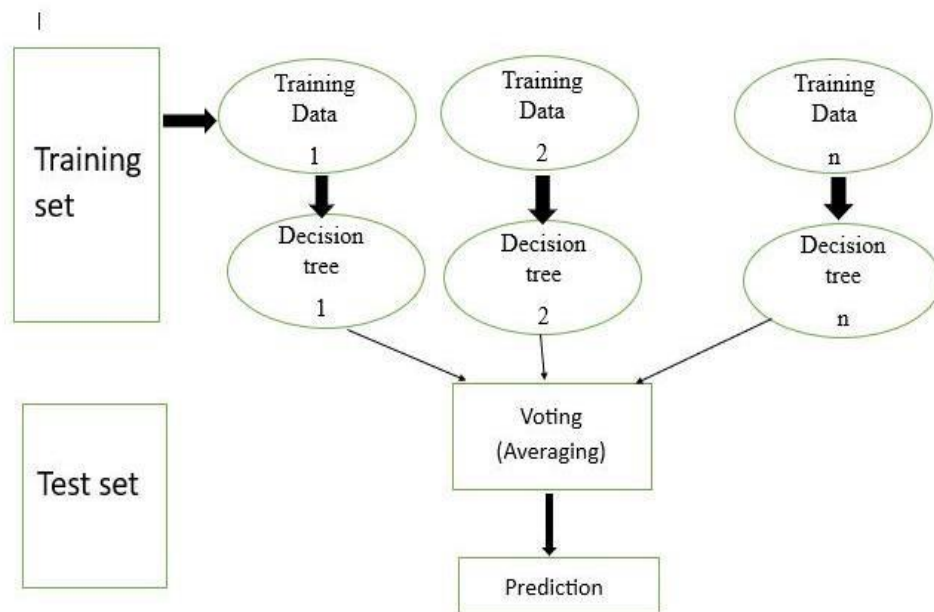


Figure 5.1.b. Random Forest

Prediction is essentially about estimating or guessing the value of something when you don't have that information readily available. It's like making an educated guess about what might happen or what something might be based on what you already know. For example, if you're trying to predict the best treatment for a certain illness, you're essentially trying to figure out which treatment is most likely to work based on factors like the patient's symptoms, medical history, and other relevant information. So, prediction involves using available data to make informed guesses or estimates about unknown or future outcomes.

### 6.2. Passive Aggressive Regressor:

Passive: If the prediction is correct, keep the model and do not make any changes. i.e., the data in the example is not enough to cause any changes in the model.

Aggressive: If the prediction is incorrect, make changes to the model. i.e., some change to the model may correct it. The Passive Aggressive Regressor is a type of algorithm used to predict numerical values based on given data. It's handy when dealing with data that keeps coming in overtime and needs to be analysed as it arrives, without delay.

### 6.3. Why not Passive Aggressive Regressor?

If the data exhibits complex patterns or relationships that the Passive Aggressive Regressor struggles to capture, or if the dataset is relatively small and doesn't require real-time processing, other algorithms like Random Forest may be more appropriate.

### 6.4. Hyperparameter Tuning:

Hyperparameter tuning involves finding the optimal values for the parameters of a machine learning model. ForRandom Forest Classifier, common hyperparameters to tune include:

n_estimators: Number of trees in the forest.max_depth: Maximum depth of each tree.

min_samples_split: Minimum number of samples required to split a node. min_samples_leaf: Minimum number of samples required at each leaf node.

max_features: Number of features to consider when looking for the best split.bootstrap: Whether bootstrap samples are used when building trees.

### VII. Experimental Scenario

Table 7.1. Comparison of Accuracies

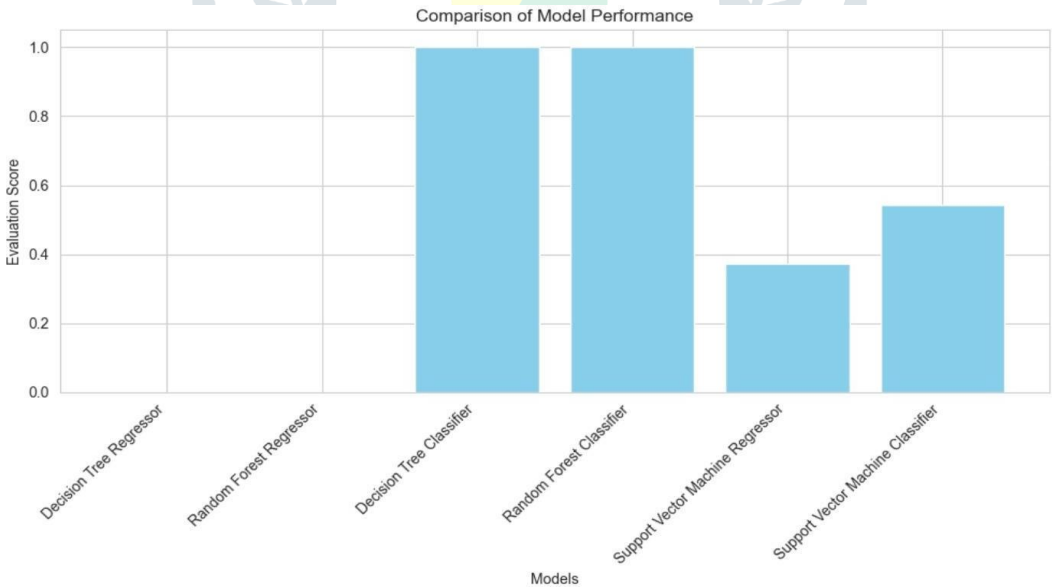| Model | Accuracy(%) |
|---|---|
| Decision Tree Regressor | 0.0 |
| Random Forest Regressor | 0.0013 |
| PassiveAggressiveRegressor | 65.33 |
| Decision Tree Classifier | 100.0 |
| RandomForestClassifier | 100.0 |
| Support Vector MachineRegressor | 37.38 |
| Support Vector MachineClassifier | 54.16 |



Figure 7.1. Graphical representation of accuracies

With the reference with figure 7.1 we used various model on same dataset and performed evaluation model for eachand every. We used Decision Tree ,Random Forest, Support Vector Machine Regressor.

## VI. RESULTS

```
Evaluation Scores:
Decision Tree Regressor: 0.0
Random Forest Regressor: 0.0013333333333333346
Decision Tree Classifier: 100.0
Random Forest Classifier: 100.0
Support Vector Machine Regressor: 37.38740613582891
Support Vector Machine Classifier: 54.166666666666664
```
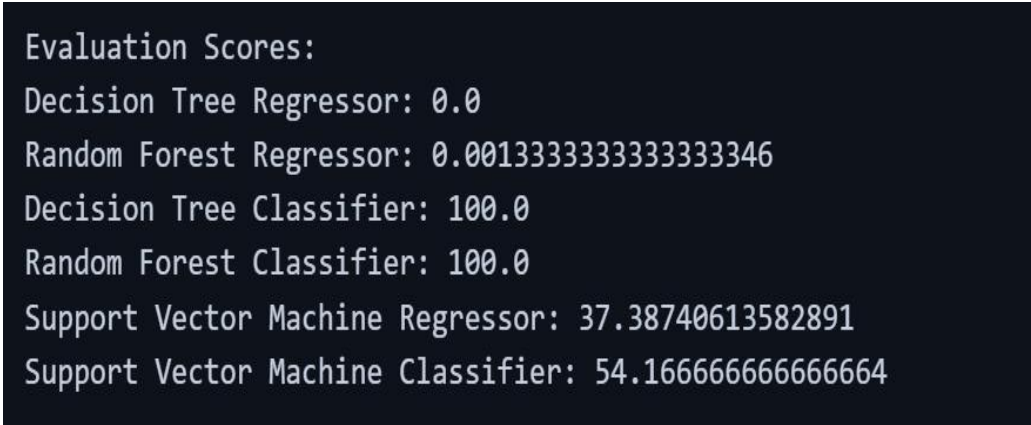
Figure 6.1  various Algorithm scores

Above the training and evaluation of each model on the identical training dataset, we analyse their performance using Decision Tree Regression ,Random Forest Regression ,Decision Tree Classifier ,Random Classifier , SVM (Support Vector Machine Regressor) , Support Vector Machine Classifier (CVMC) This assessment aids in determiningthe most suitable model is Decision Tree Classifier and Random Forest Classifier for our prediction task and dataset, taking into account factors such as predictive accuracy, model complexity, and computational efficiency. It gives 100% accuracy.
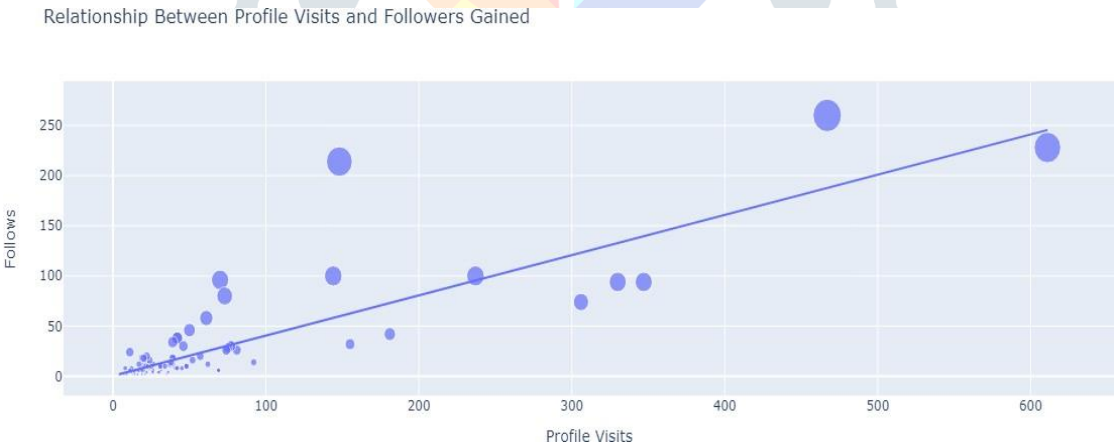


Figure 6.2 .Relationship between profile visits and followers gained

Figure 6.2 illustrates the correlation between profile visits and the number of followers gained. The graph depicts how changes in profile visits relate to the increase or decrease in followers over a given period. This visual representation helps to understand the relationship between these two variables, providing insights into potential patterns or trends thatmay exist. By analysing this relationship, one can gain valuable insights into factors influencing follower growth on a profile.
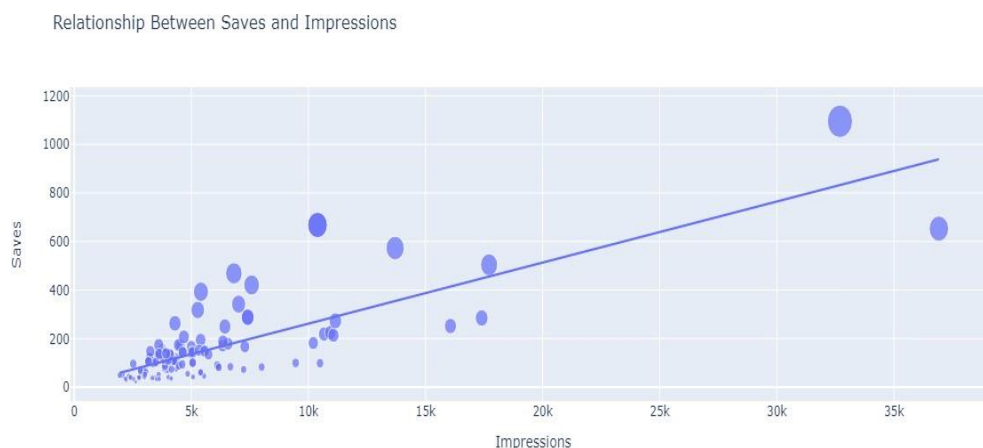
Figure 6.3 Relationship between saves and impressions

Figure 6.3 displays the connection between the number of saves and impressions. Through this graph, we can observe the relationship between these two variables and how they impact each other. By examining this relationship, we can gain insights into the effectiveness of content in terms of engagement and reach. This visualization aids in understanding how the number of saves influences the impressions a piece of content receives, providing valuable information for content optimization and strategy planning.



Figure 6.4 Relationship between shares and impressions

In Figure 6.4, we explore the relationship between shares and impressions. This visual representation offers insights into how the number of shares impacts the overall impressions of content. By examining this relationship, we can discern patterns and trends that shed light on the effectiveness of content in generating engagement and visibility. Understanding this correlation is crucial for optimizing content strategies to enhance reach and audience engagement.

Figure 6.5 most common used hashtags words

In Figure 6.5 These words are often used to describe the process, findings, and interpretations of data analysis. data analysis, machine learning, data, visualization, patterns, trend, corelation.



| Impressions | 1.000000 |
| From Explore | 0.893607 |
| Follows | 0.889363 |
| Likes | 0.849835 |
| From Home | 0.844698 |
| Saves | 0.779231 |
| Profile Visits | 0.760981 |
| Shares | 0.634675 |
| From Other | 0.592960 |
| From Hashtags | 0.560760 |
| Comments | -0.028524 |

Name: Impressions, dtype: float64

Figure 6.6 scores for each

With reference from Figure 6.6 score for each features have been calculated from where the user reached As mentionedin the above figure we had calculated mode for all the columns mentions in the datasets like from explore, follows, likes,from home, saves, shares, profile visits, from hashtags, from other, comments.



| | username | mutual_connecting | total_following | total_followers | name | private_account | business_account | recently_joined | total_likes | total_comments | total_posts | follow_back |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 7olga7777777 | 0 | 107 | 69 | Ольга Петришина | 0 | 0 | 0 | 61 | 7 | 2 | 1 |
| 1 | aadamkathst | 2 | 212 | 40 | aadam kathst | 0 | 0 | 0 | 43 | 0 | 12 | 1 |
| 2 | aadish__khan | 2 | 116 | 65 | aadish__khan | 0 | 0 | 0 | 184 | 6 | 18 | 0 |
| 3 | aarzoo6466 | 2 | 222 | 40 | Aarzoo Khan | 0 | 1 | 1 | 3 | 0 | 1 | 0 |
| 4 | abbasyadavad | 0 | 25 | 17 | Abbas Yadavad | 0 | 0 | 0 | 3 | 0 | 2 | 0 |

Figure 6.7 Databases table of users

With reference fig 6.7 we analyse from the data set user name, Mutual connecting, Total followers, Name, Private account, Business account, recently joined Instagram ,total likes, total comments , total post.

```
Total private account are: 2
Total business accounts are: 11
```
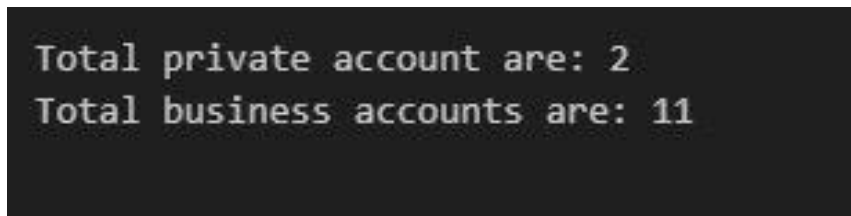
Figure 6.8 Count of private and public account

With reference to Fig 6.8 we have analyse various account of different users account types . In this we have categorizedtwo part that is private and public account and analyse the type with total calculation of account types in the above we have calculated total in which private accounts is 2 and business accounts are 11 in total.
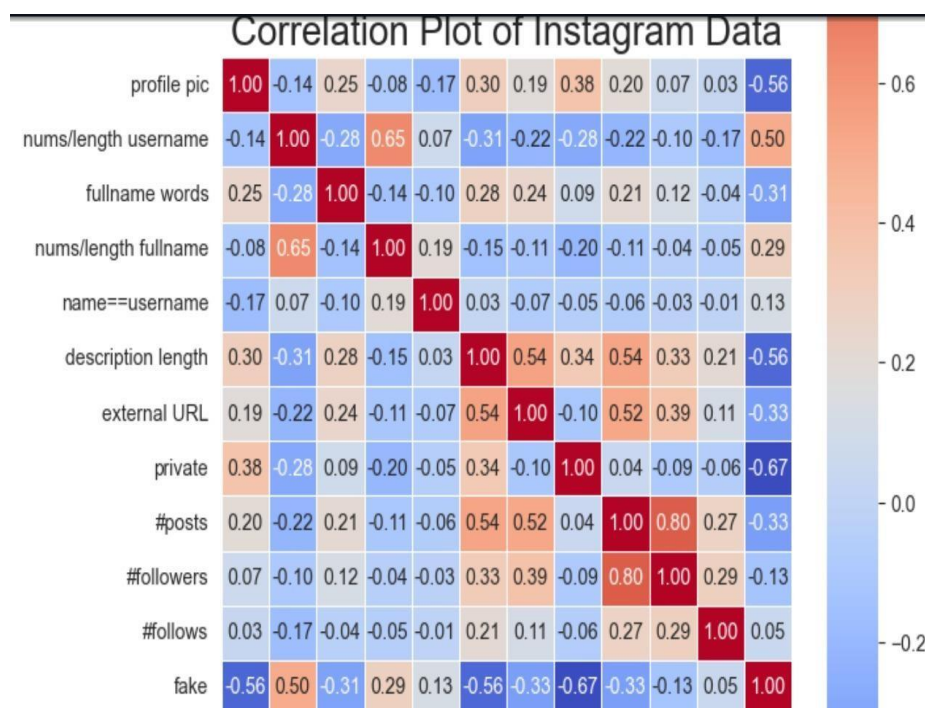


Figure 6.9 Correlation plot of Instagram data

With reference to Fig 6.9 A heatmap is a popular visualization tool for depicting the correlation matrix visually. In aheatmap: Every cell represents the correlation coefficient between two variables.
The strength of correlation is depicted using a colour gradient: warmer hues like red signify positive correlations, whilecooler tones like blue denote negative correlations. Neutral colours like white indicate no correlation.
Additionally, you can choose to include correlation coefficient values within each cell to aid interpretation.

## VI. CONCLUSION

In the analysis of Instagram data using Big Data tools and techniques has provided valuable insights into user behaviour, content trends, and platform performance. This comprehensive analysis has the potential to benefit various stakeholders, including marketers, businesses, and the platform itself. By leveraging the vast amount of data generated on Instagram, we can make informed decisions, optimize marketing strategies, enhance user experiences, and ultimately,drive engagement and growth on the platform. As Big Data technology continues to evolve, the opportunities for deeper and more meaningful analysis of Instagram data will only expand, offering even greater potential for insights and improvements in the future ..

**REFERENCES**

1. Dr C K Gomathy-Assistant Professor , "BIG DATA ANALYTICS IN INSTAGRAM", India Nov-2022 (IJSREM).

2. Taylor & Francis Group, LLC, " Influencer Marketing on Instagram": Empirical Research on Social MediaEngagement with Sponsored Posts 2022.

3. Hung Hom, Kowloon, Hong Kong, "Research in the Instagram Context": Approaches & Methods 2021.

4. Sharma and Jain, "Social Media Visualization" 2021.

5. Dr. Daryl D. Green, Dr. Richard Martinez, Amalan Kadjar, Lauran Evenson, Lisa MacManus, Stepanie Birbeck in a"world of social media: A case study Analysis of Instagram" 2020.

6. Elsevier Ltd, "The big picture on Instagram research": Insights from a bibliometric analysis 2021.

7. Sivaraj Stieglitz et al., " Social Media Analytics" 2020.

8. Dr.C K Gomathy, Article: "A Study on the recent Advancements in Online Surveying , International Journal ofEmerging technologies" ( JETIR ) Volume 5 | Issue 11 | ISSN : 2349-5162, P.No:327-331, Nov-2018.

9. Dr.C.K.Gomathy,P.Sarvani Divya jyothsna,M.Srimayi, Article: A study on the Mobile Application Advancements in Anti-Ragging, SSRG International Journal of Computer Science And Engineering(SSRG-IJCSE)-Volume 6 issue 3,March 2019.

10. Dr.C K Gomathy, THE LOAN PREDICTION USING MACHINE LEARNING, International Research Journal of Engineering and Technology (IRJET), Volume: 08 Issue: 10 | e-ISSN: 2395-0056, p-ISSN: 2395-0072, Impact Factor value: 7.529, Available at www.irjet.net Oct 2021.