

# Capstone Project - 4

**Unsupervised – ML – Zomato Restaurant  
Clustering and Sentiment Analysis**

**Sethupathy M**

# Contents

- 1) What is Zomato?
- 2) How Zomato works?
- 1. 3) Zomato funds and stats
- 2. 4) Zomato business model
- 3. 5) Addressing the problem
- 4. 6) Feature summary
- 5. 7) EDA
- 6. 8) Text Pre-processing
- 7. 9) Clustering
- 10) KMeans Clustering
- 11) DBSCAN Clustering
- 12) Sentiment Analysis
- 13) Conclusion

# What is **zomato**?

Zomato is one of the most comprehensive and user-friendly apps where people can search for nearby restaurants and cafe's, order food online, and get it delivered at their doorstep in no time. Moreover, you can also get accurate information about restaurants as it provides menus, reviews, and ratings. Based on that, users can place orders and enjoy lip-smacking food at their homes.

Zomato was founded by **Deepinder Goyal** and **Pankaj Chaddah**, two **Delhi IIT** graduates, in **2008**. Till November 2010, Zomato was known as “**Foodiebay**.” Once they saw their colleagues who were seeking menus of different restaurants to order food. That’s when the idea took birth, and they thought of converting these manual menus into a digital format. In the year 2012, Zomato had spread its wings across the globe and started to list out the number of restaurants in the market.

# How **zomato** works ?

The main work of Zomato is to suggest local and nearby restaurants to users and receive orders from them. Users can place orders from their favorite restaurant based on ratings and reviews shared by previous customers.

**Step 1:** From the dedicated app solution or website, users can explore various restaurants and order meals.

**Step 2:** Particular restaurant owners receive an order request and start preparing a meal.

**Step 3:** Once the food is ready to dispatch, it will be handed over to delivery providers.

**Step 4:** Delivery providers deliver the meal to the customer's preferred location.

**Step 5:** From the given payment options, customers can make payments and share reviews based on their experience.





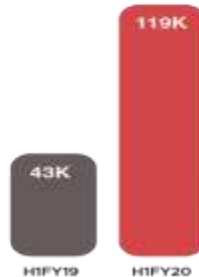
## Funds and Stats

Zomato received a total number of 909.6 million from different investors. Their funding was from Private Equity in 2020. Info Edge is a leading investor of Zomato. Other than that, Ant Financial, Delivery Hero, Shunwei Capital, Vy Capital, and many others are the investors of Zomato who have contributed their major stack to make Zomato popular worldwide. Now let's have a look at some interesting figures about Zomato's growth.

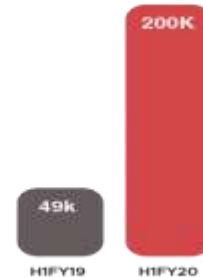
**Average Monthly  
Transacting Users**  
↑ 211%



**Average Monthly  
Active Restaurants**  
↑ 177%



**Average Monthly  
Active Delivery Partners**  
↑ 308%



**zomato**

# Business Model

Zomato's Business Model is aimed at providing quality food services, information related to restaurants, their menus and user reviews. The Business Model of Zomato consists of providing food delivery services, information, user reviews and menu's of partner restaurants. It has created a revolution in industries doing food business by including different restaurants and facilitating people to look for restaurants more conveniently.

## Zomato Business Model




# Addressing the Problem

The Project focuses on Customers and Company, and to analyze the sentiments of the reviews given by the customer in the data and made some useful conclusion in the form of Visualizations. Also, clustering the zomato restaurants into different segments. The data is visualized as it becomes easy to analyze data at instant. The Analysis also solve some of the business cases that can directly help the customers finding the Best restaurant in their locality and for the company to grow up and work on the fields they are currently lagging in.

The Project contains : Exploratory data analysis, Clustering, Sentiment Analysis which could help the customers to choose best restaurants and the restaurant owners to improve the restaurants in various aspects.

# Features Summary

## Zomato Restaurant Names and Metadata

- **Name** – Name of the Restaurants.
- **Links** – Links of the Restaurants.
- **Cost** – Average cost of the meal in Restaurants.
- **Collections** - The Collections in Zomato features popular restaurants across specific themes and trends at a particular location.
- **Cuisine** – A cuisine is specific set of cooking traditions and practices, often associated with a specific culture or region.
- **Timings** – Opening and closing time of the Restaurants.



# Features Summary (continued)

## Zomato Restaurant Reviews

- **Restaurants** – Name of the Restaurants.
- **Reviewer** – Name of the Reviewer.
- **Review** – Experience of the reviewer in the restaurant expressed in words.
- **Rating** - Experience of the reviewer in the restaurant expressed in numbers.
- **Metadata** – A set of data that describes and gives information about reviewers.
- **Time** – Time at which the reviews and ratings are given.
- **Pictures** – Number of pictures taken by the reviewer.

# NULL value treatment

In the 'Zomato Restaurant Names and Metadata' dataset 'Collections' feature contains more than 50 % of its values as NULL. So 'Collections' feature is dropped.

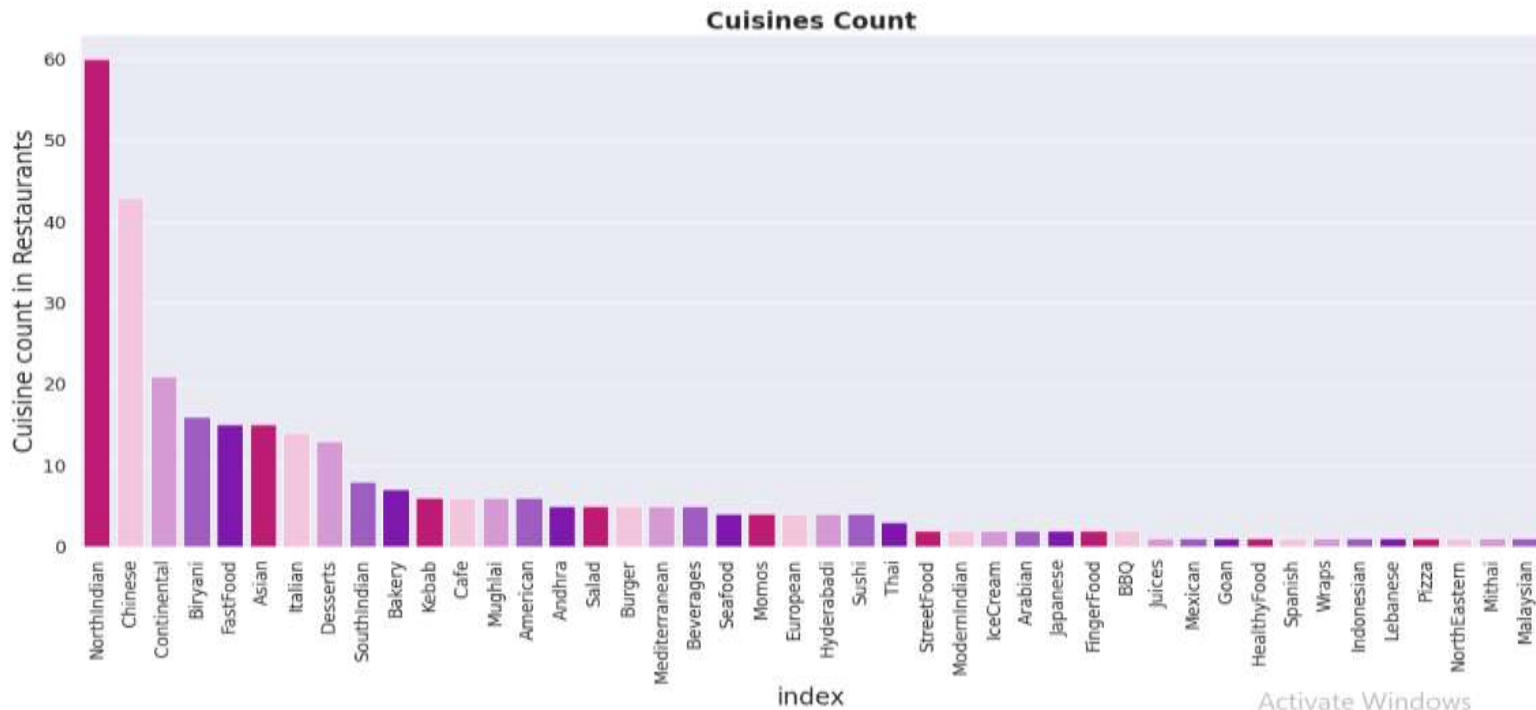
```
# Null values
# Percentage of null values for each features
# print(list(df.columns))

for col in list(df.columns):
    if ((df[col].isnull().sum())/(len(df[col]))*100) > 50:
        # print((df[col].isnull().sum())/(len(df[col]))*100)
        print('Feature with more than 50% of the observations are NULL values:',col)

# As in the feature 'Collections' more than 50% of the observations are NULL values, so feature 'Collections' is dropped
```

# Exploratory Data Analysis

On **Zomato Restaurant Names and Metadata** dataset.



# Exploratory Data Analysis (continued)

On **Zomato Restaurant Names and Metadata** dataset.



### Cost of Restaurants

Restaurant Name	Cost (INR)
Collage - Pyatt Hyderabad	2650
Feast - Sheraton Hyderabad Hotel	2500
10 Downing Street	1900
mathan's kitchen - Holiday Inn Express & Suites	1850
Caracoli - Sheraton Hyderabad Hotel	1750
Zoya - Sheraton Hyderabad Hotel	1700
Republic Of Noodles - Lemon Tree Hotel	1650
Mazzo - Marriott Executive Apartments	1600
Barbecue Nation	1550
Alegra - Bar & Kitchen	1550
BD-Jobs	1500
The Indi Grill	1450
AB's - Absolute Barbecues	1450
Komalose - Holiday Inn Express & Suites	1450
The Fisherman's Wharf	1450
SKHRV	1400
Marvecharu - Bar & Kitchen	1350
East India Company	1300
Mustang Terrace Lounge	1250
Over The Moon Brew Company	1200
Prism Club & Kitchen	1150
Diner's Pavilion	1150
The Arabain Food Court	1150
Yum Yum Tree - The Arabain Food Court	1150
JB's - Buddies, Bar & Barbecue	1100
Urban Asia - Kitchen & Bar	1050
Khaani Saab	1050
La La Land - Bar & Kitchen	1000
Chinese Pavilion	1000
Pasta House	1000
Hyperlocal	1000
Qwm Mon Mon	950
Club Rogue	900
Dine O China	850
Gai Panahi D	800
Shah Ghouse Hotel & Restaurant	800
Beyond Flavours	800
Shanghae Chef 2	800
Al Saba Restaurant	750
Absolute Sizzlers	750
Olive Garden	700
Shal Bhoj	700
Banana Leaf Multicuisine Restaurant	700
Royal Spry Restaurant	700
Green Bawarchi Restaurant	700
Café Eclat	700
Marqalla Food Company	700
Squeeze @ The Lime	650
The Chocolate Room	600
Karachi Cafe	600
American Wild Wings	600
Driven Cafe	600
Befrouz Brynari	600
Rabos	600
Hyderabad Davaas	600
Hyderabad Chels	600
Zing's Northeast Kitchen	600
Northern Kitchen	600
Dunkin' Donuts	600
The Foodie Monster Kitchen	600
Tandron Food Works	600
Madhura Vraa	600
Blynz	600
Udipi's Uppahar	600
Angaara Counts 3	600
Shree Sanjivuh Dhabba Family Restaurant	600
Kritunga Restaurant	600
Hirechi Bawarchi Food Zone	600
WFC	600
GO's	600
Karachi Bawarchi	600
13 Dhabba	600
Being Hungry	600
Pakwaan Gharid	600
WCCY	600
Sandarija's Chaats & Munchies	600
Hotel Zora H-Fri	600
Dominio's Pizza	600
Desi Bytes	600
Creem Stone	600
The Old Madras Baking Company	600
Shah Ghouse SpQ Shaleema	600
Which Pessae	600
Assam Meal Box	600
Hunger Maggi Point	600
Momos Delight	600
Sweet Basket	600
Mohammeda Shaleema	600

# Text Preprocessing

On **Zomato Restaurant Names and Metadata** dataset.

## Stemming

It is a technique used to extract the base form of the words by removing affixes from them. It is just like cutting down the branches of a tree to its stems.

## Lemmatization

Lemmatization is a method responsible for grouping different inflected forms of words into the root form, having the same meaning.

```
def stem_words(words):  
    """Stem words in list of tokenized words"""  
    stemmer = LancasterStemmer()  
    stems = []  
    for word in words:  
        stem = stemmer.stem(word)  
        stems.append(stem)  
    return stems  
  
def lemmatize_verbs(words):  
    """Lemmatize verbs in list of tokenized words"""  
    lemmatizer = WordNetLemmatizer()  
    lemmas = []  
    for word in words:  
        lemma = lemmatizer.lemmatize(word, pos='v')  
        lemmas.append(lemma)  
    return lemmas  
  
def normalize(words):  
    words = remove_non_ascii(words)  
    words = to_lowercase(words)  
    words = remove_punctuation(words)  
    words = replace_numbers(words)  
    return words
```

# Text Preprocessing (continued)

On **Zomato Restaurant Names and Metadata** dataset.

- **Non-ASCII** characters are those that are not encoded in ASCII, such as Unicode, EBCDIC..
- Converting Uppercase letters to **Lowercase**.
- **Removing punctuations**.

```
def remove_non_ascii(words):
    """Remove non-ASCII characters from list of tokenized words"""
    new_words = []
    for word in words:
        new_word = unicodedata.normalize('NFKD', word).encode('ascii', 'ignore').decode('utf-8', 'ignore')
        new_words.append(new_word)
    return new_words

def to_lowercase(words):
    """Convert all characters to lowercase from list of tokenized words"""
    new_words = []
    for word in words:
        new_word = word.lower()
        new_words.append(new_word)
    return new_words

def remove_punctuation(words):
    """Remove punctuation from list of tokenized words"""
    new_words = []
    for word in words:
        new_word = re.sub(r'[^\w\s]', '', word)
        if new_word != '':
            new_words.append(new_word)
    return new_words

def replace_numbers(words):
    """Replace all integer occurrences in list of tokenized words with textual representation"""
    p = inflect.engine()
    new_words = []
    for word in words:
        if word.isdigit():
            new_word = p.number_to_words(word)
            new_words.append(new_word)
        else:
            new_words.append(word)
    return new_words
```

# Text Preprocessing (continued)

On **Zomato Restaurant Names and Metadata** dataset.

## TFIDF vectorizer

Term frequency-inverse document frequency is a text vectorizer that transforms the text into a usable vector. It combines 2 concepts, Term Frequency (TF) and Document Frequency (DF). The term frequency is the number of occurrences of a specific term in a document.

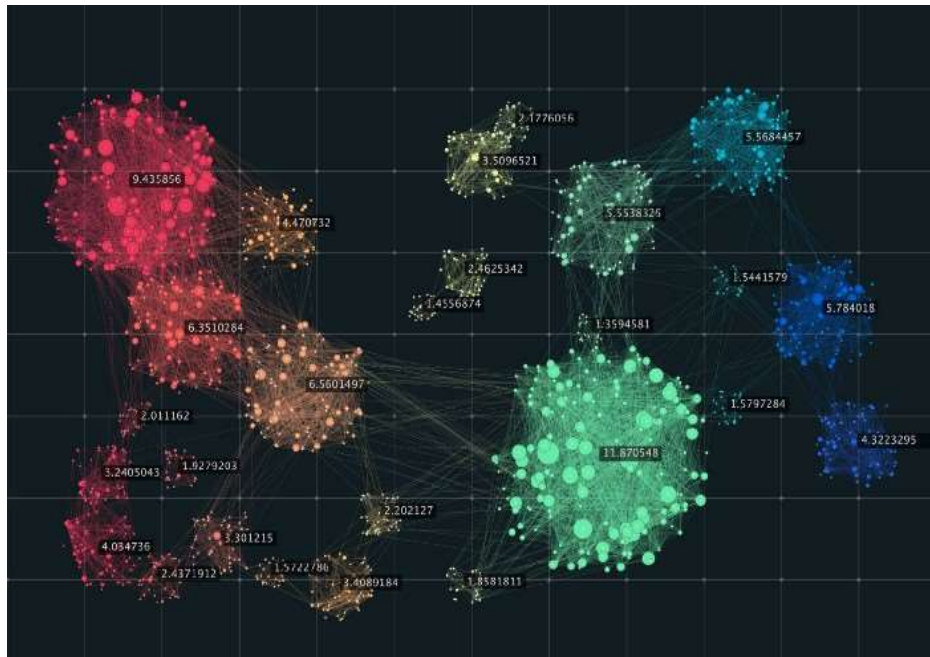
```
# Tfidf vectorizer
vectorizer = TfidfVectorizer(stop_words= 'english')
X = vectorizer.fit_transform(df['Cuisines'])
```



# Clustering

On **Zomato Restaurant Names and Metadata** dataset.

**Clustering** – It is the task of dividing the population or data points into a number of groups such that data points in the same groups are more similar to other data points in the same group than those in other groups.

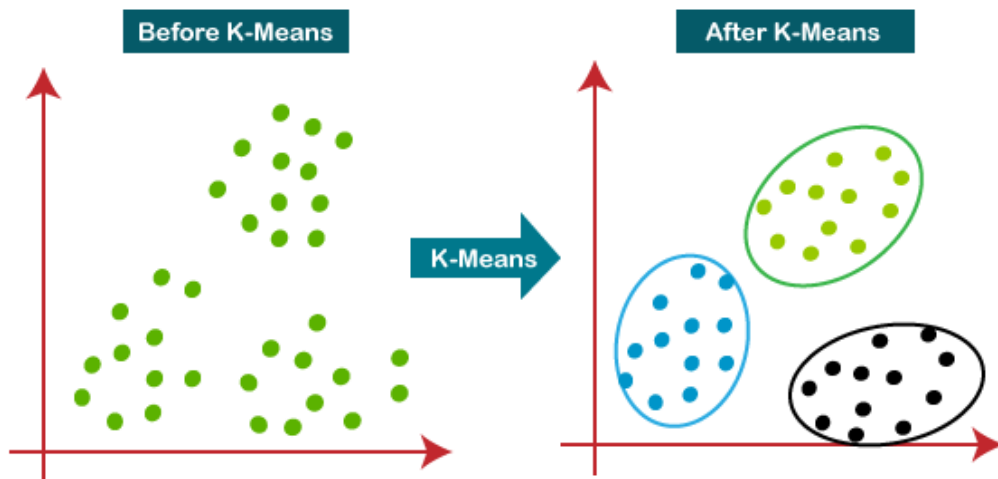


# Clustering (continued)

On **Zomato Restaurant Names and Metadata** dataset.

## K Means Clustering

The K-means clustering algorithm computes centroids and repeats until the optimal centroid is found. It is presumptively known how many clusters there are. It is also known as the flat clustering algorithm. The number of clusters found from data by the method is denoted by the letter 'K' in K-means.

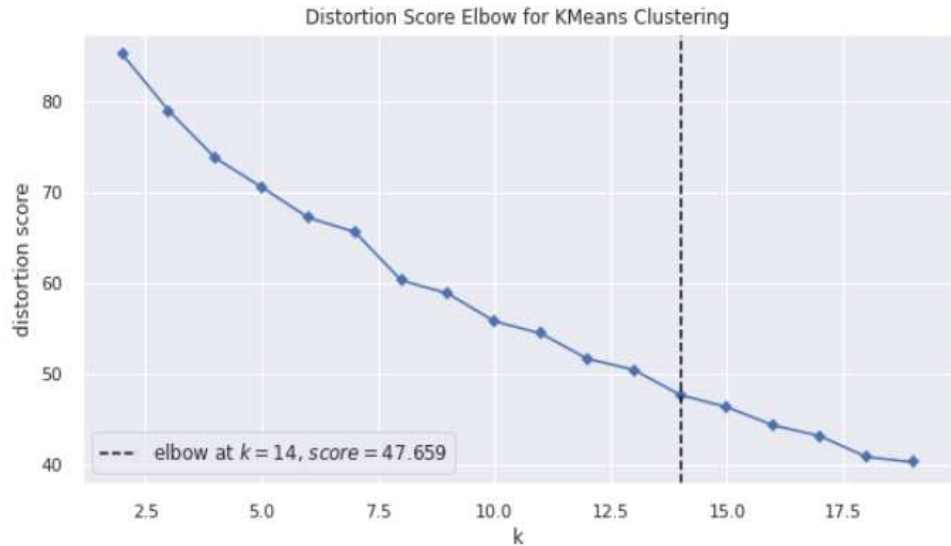


# Clustering (continued)

On **Zomato Restaurant Names and Metadata** dataset.

## Finding appropriate 'K' value

```
# Function to find appropriate 'K' value
def KElbowvisualizer(metric):
    model = KMeans(init="k-means++",max_iter=300,random_state=0)
    plt.figure(figsize=(10,5))
    sns.set(font_scale = 1)
    visualizer = KElbowVisualizer(model, k=(2,20),metric= metric,
    # plt.title(fontweight='bold')
    # # Fit the data to the visualize
    visualizer.fit(X)
    visualizer.poof()
```



# Clustering (continued)

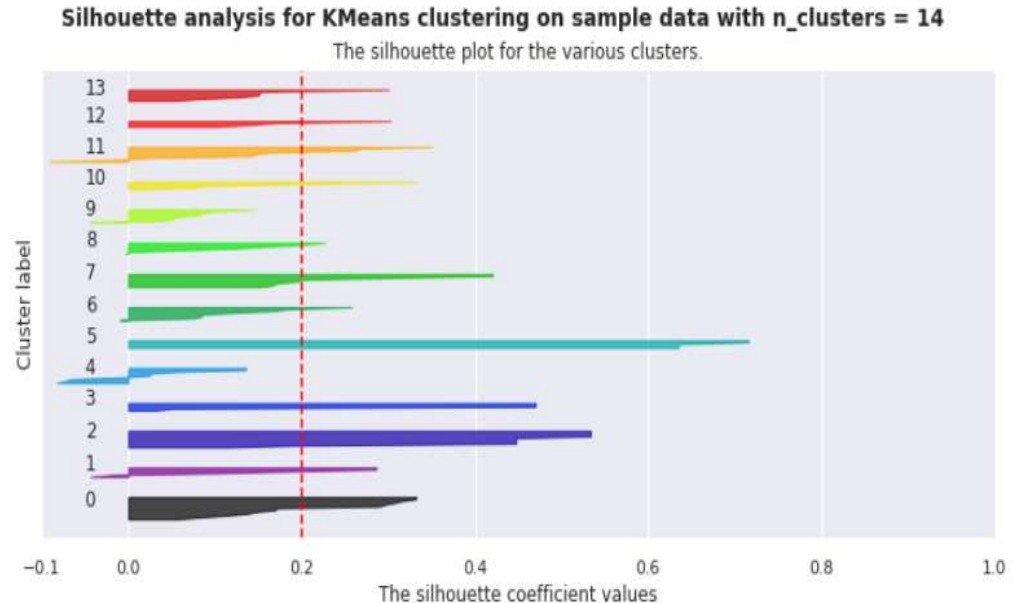
## On Zomato Restaurant Names and Metadata dataset.

- **Working of K-Means Algorithm**
  - The following stages will help us understand how the K-Means clustering technique works-
  - **Step 1:** First, we need to provide the number of clusters,  $K$ , that need to be generated by this algorithm.
  - **Step 2:** Next, choose  $K$  data points at random and assign each to a cluster. Briefly, categorize the data based on the number of data points.
  - **Step 3:** The cluster centroids will now be computed.
  - **Step 4:** Iterate the steps below until we find the ideal centroid, which is the assigning of data points to clusters that do not vary.

# Clustering (continued)

## On Zomato Restaurant Names and Metadata dataset. K Means Clustering Model Validation

For `n_clusters = 10` The average silhouette\_score is : 0.1825  
For `n_clusters = 11` The average silhouette\_score is : 0.1747  
For `n_clusters = 12` The average silhouette\_score is : 0.1937  
For `n_clusters = 13` The average silhouette\_score is : 0.1866  
For `n_clusters = 14` The average silhouette\_score is : 0.1998  
For `n_clusters = 15` The average silhouette\_score is : 0.1945

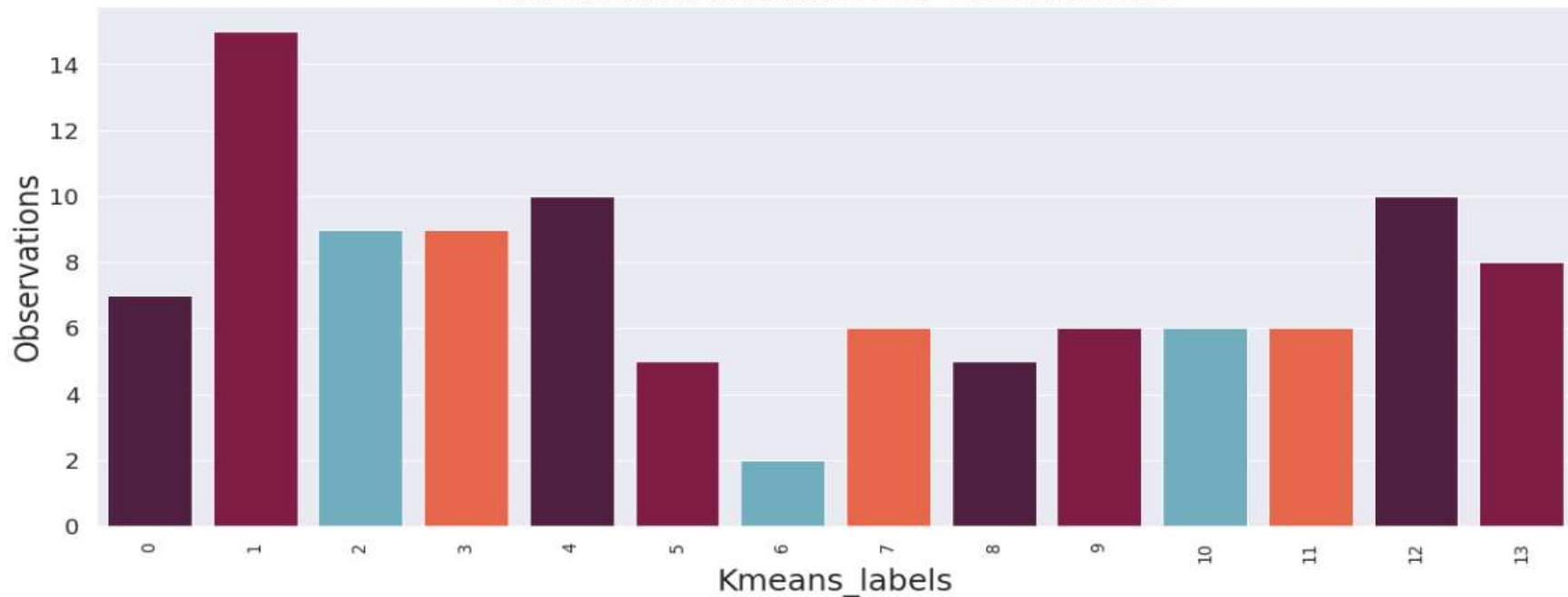


# Clustering (continued)

On Zomato Restaurant Names and Metadata dataset.

## K Means Clustering

Labels with n number of observations

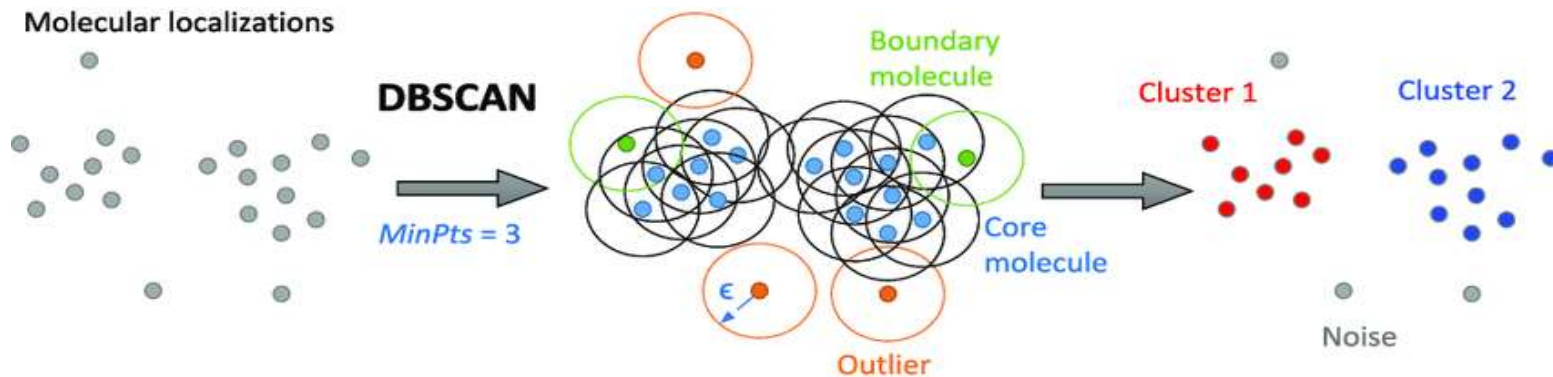


# Clustering (continued)

**On Zomato Restaurant Names and Metadata dataset.**

## DBSCAN Clustering

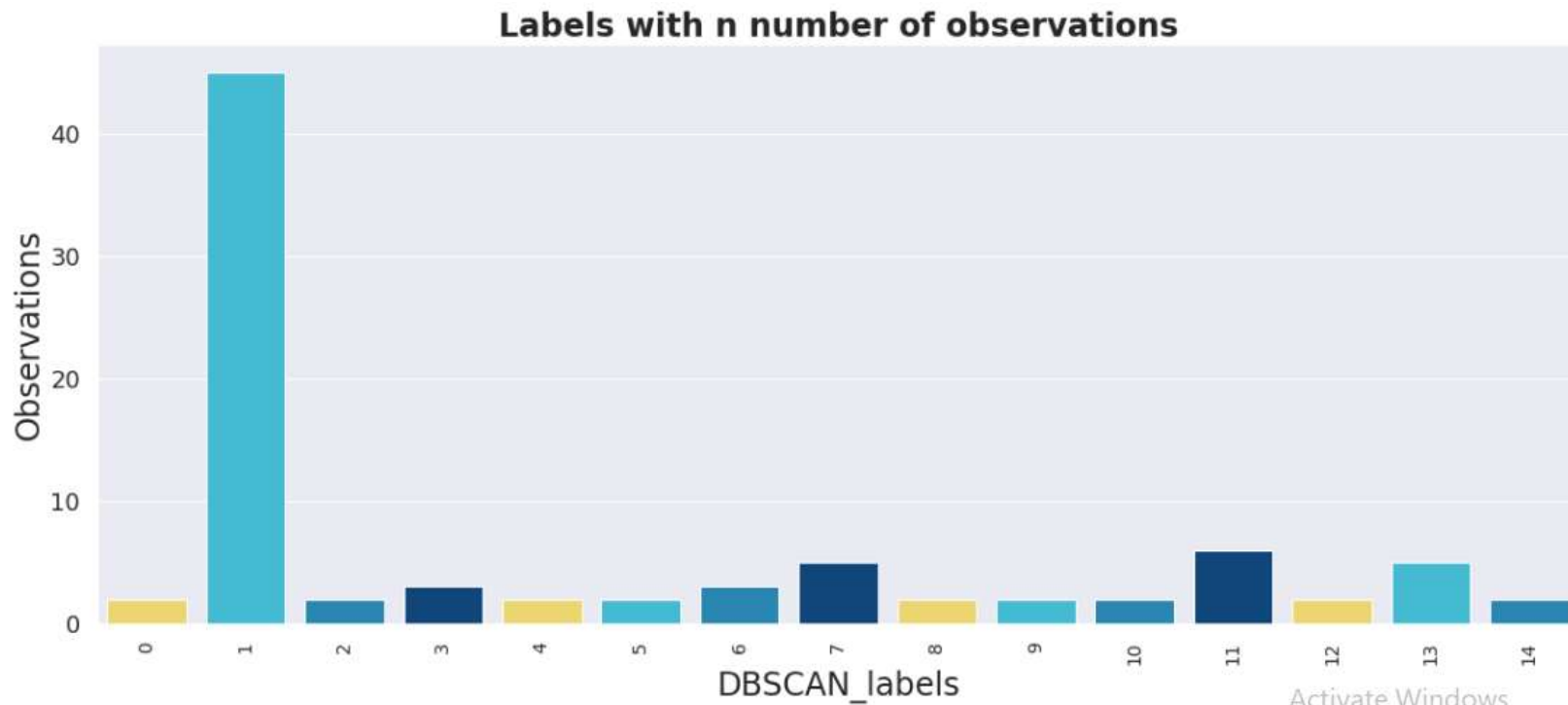
DBSCAN is a density-based clustering algorithm that works on the assumption that clusters are dense regions in space separated by regions of lower density. It groups 'densely grouped' data points into a single cluster. It can identify clusters in large spatial datasets by looking at the local density of the data points. The most exciting feature of DBSCAN clustering is that it is robust to outliers.



# Clustering (continued)

On Zomato Restaurant Names and Metadata dataset.

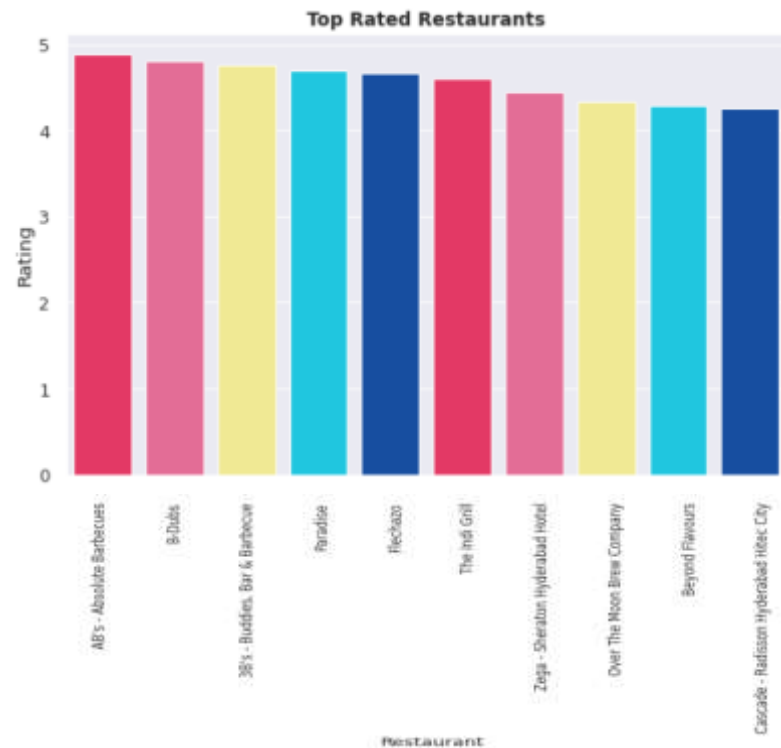
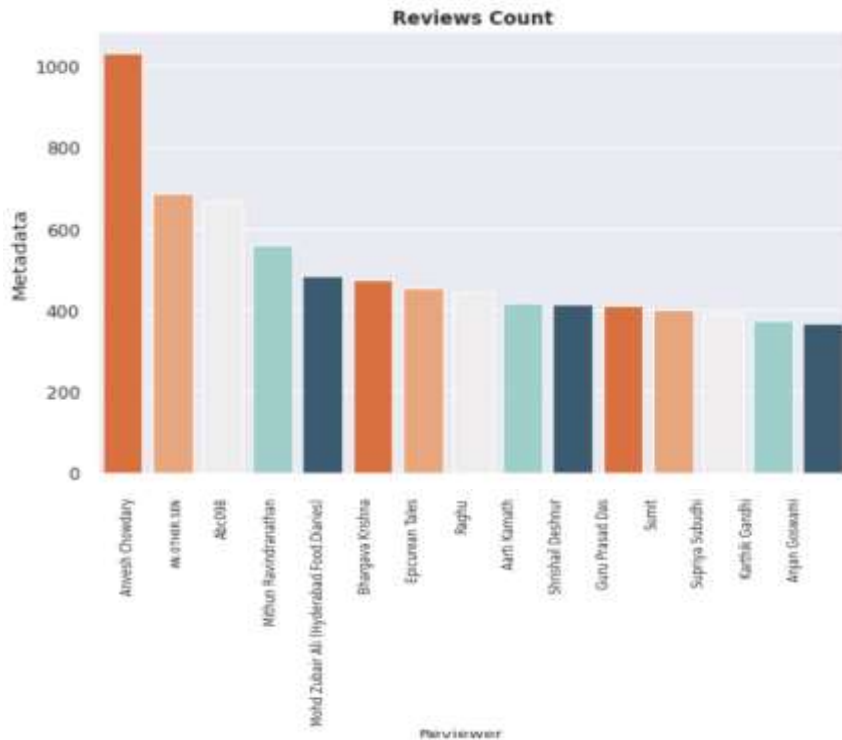
## DBSCAN Clustering





# Exploratory Data Analysis

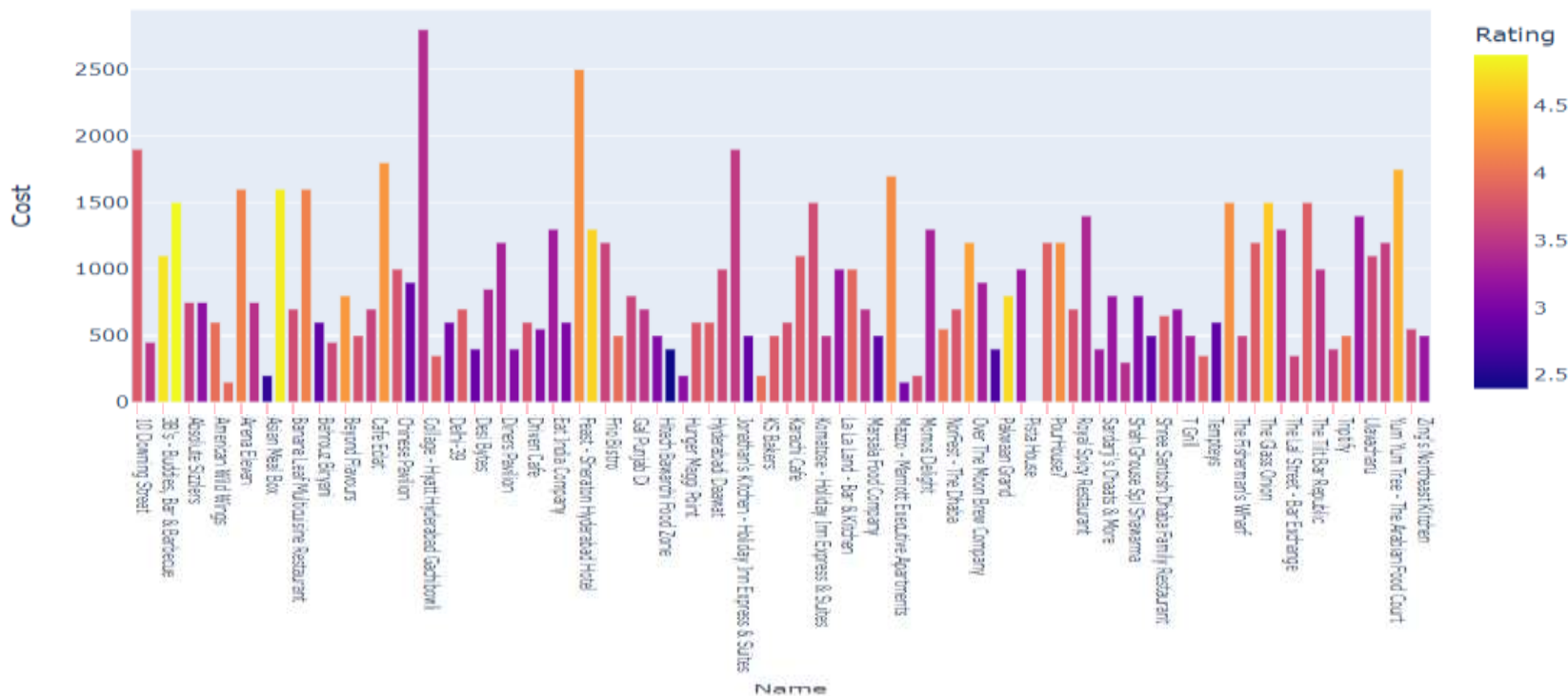
On Zomato Restaurant Reviews dataset.



## Exploratory Data Analysis (continued)

## On Merged dataset.

### Restaurant Cost vs Rating



# Sentiment Analysis

## Vader Model

Sentiment analysis is a text analysis method that detects polarity (e.g. a positive or negative opinion) within the text, whether a whole document, paragraph, sentence, or clause. Sentiment analysis aims to measure the attitude, sentiments, evaluations, attitudes, and emotions of a speaker/writer based on the computational treatment of subjectivity in a text.

## Sentiment Analysis



**Positive**



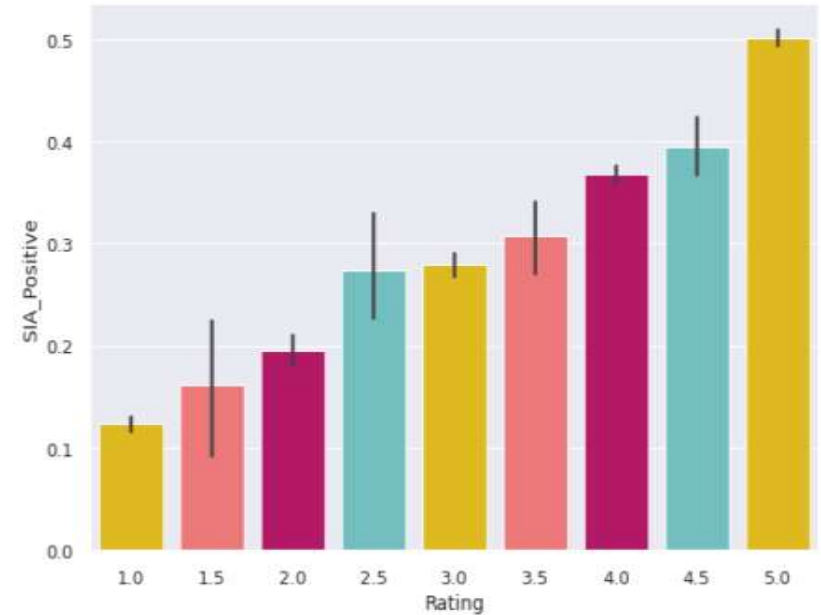
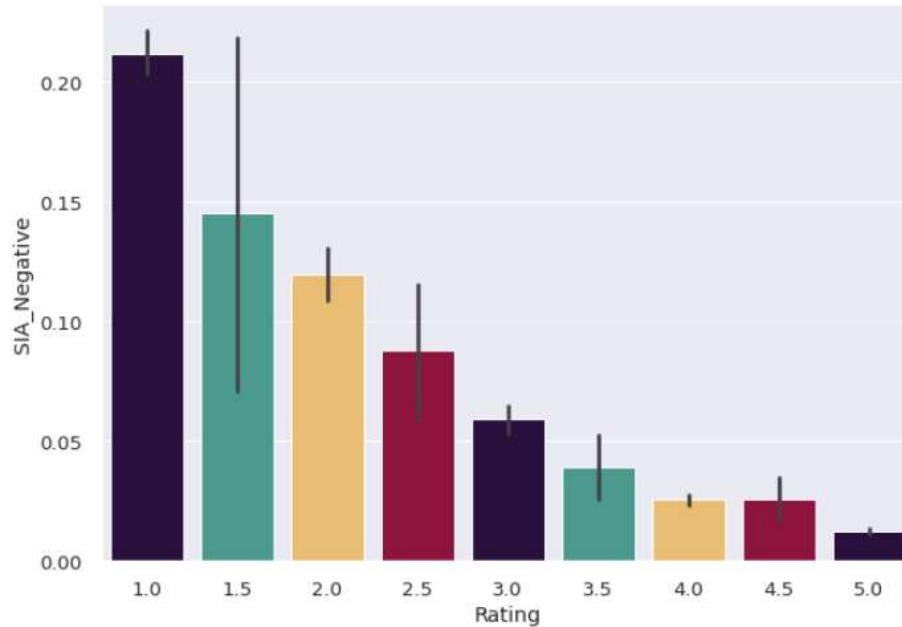
**Negative**



**Neutral**

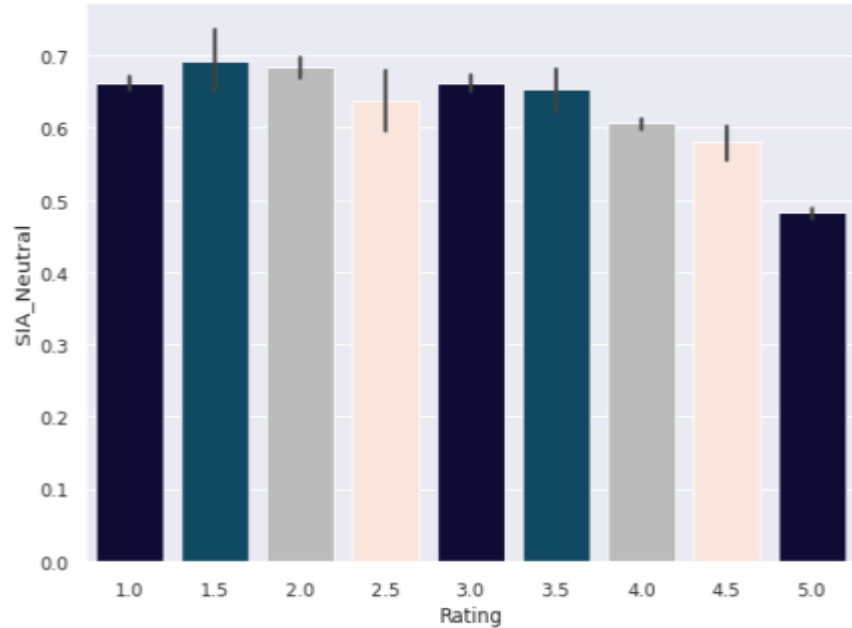
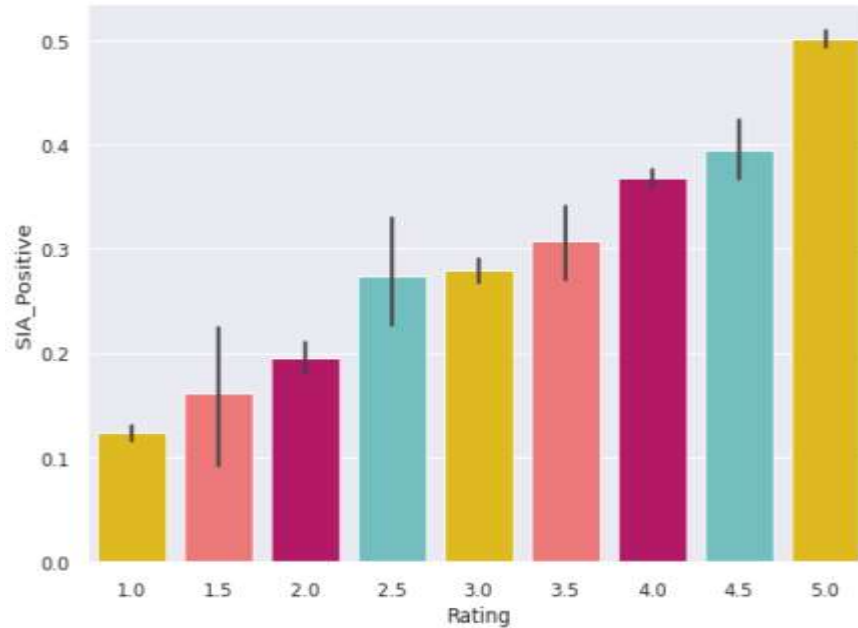
# Sentiment Analysis (continued)

## Vader Model



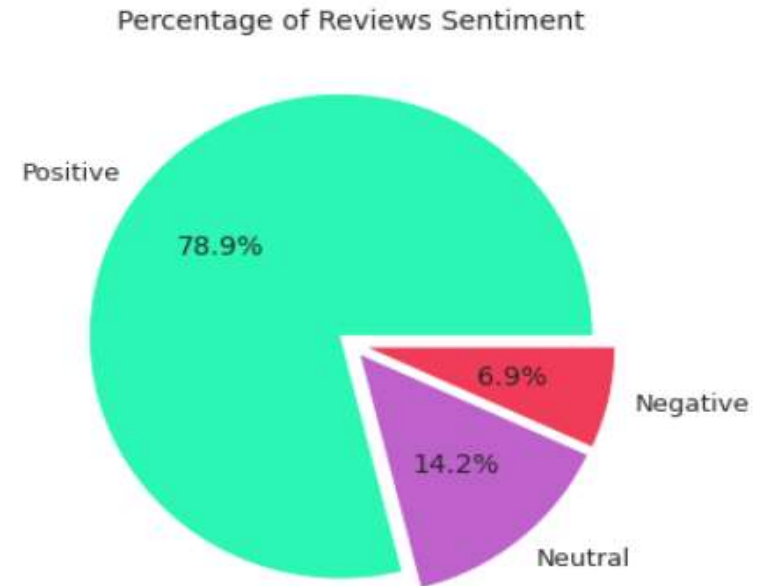
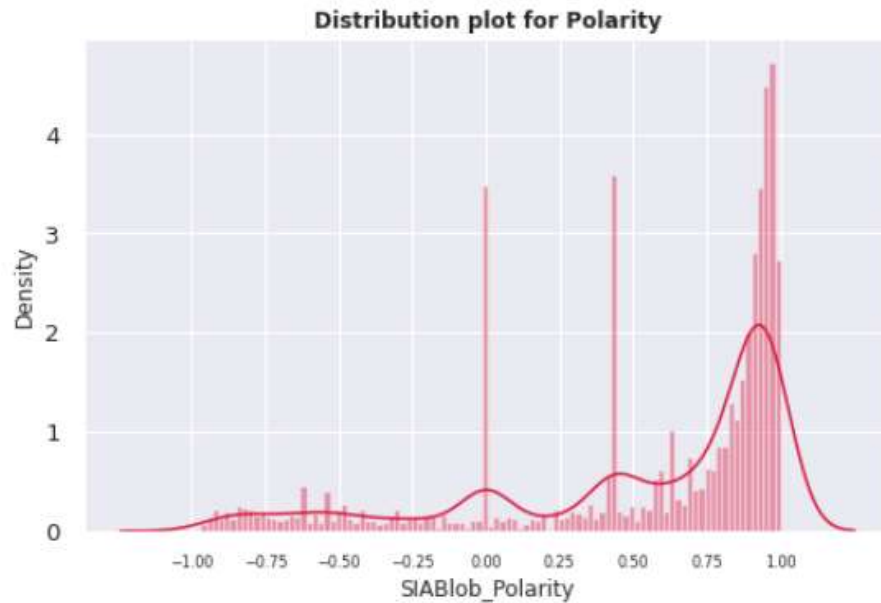
# Sentiment Analysis (continued)

## Vader Model



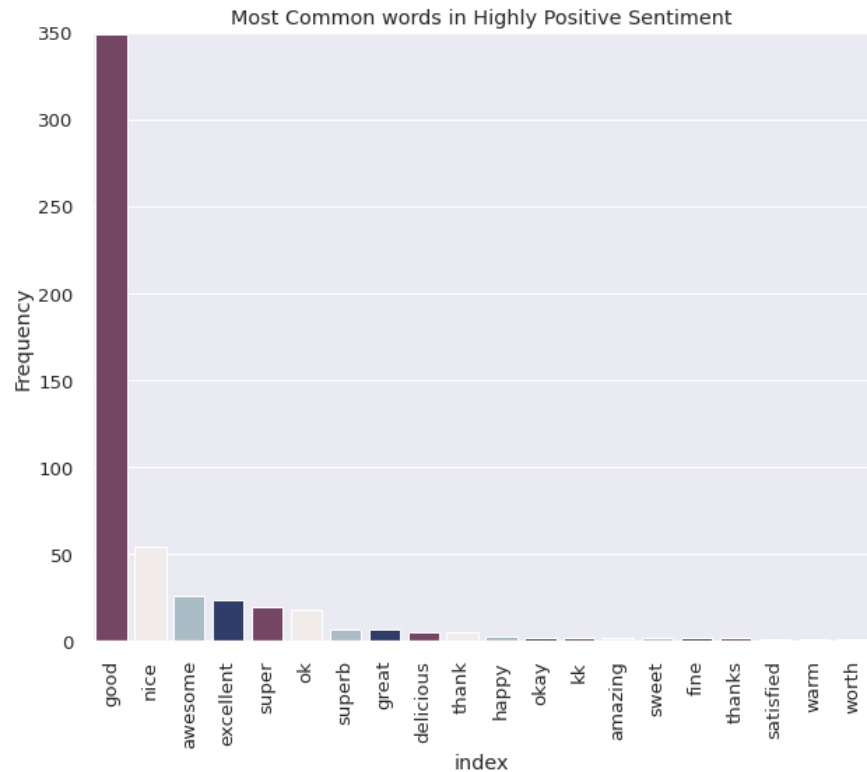
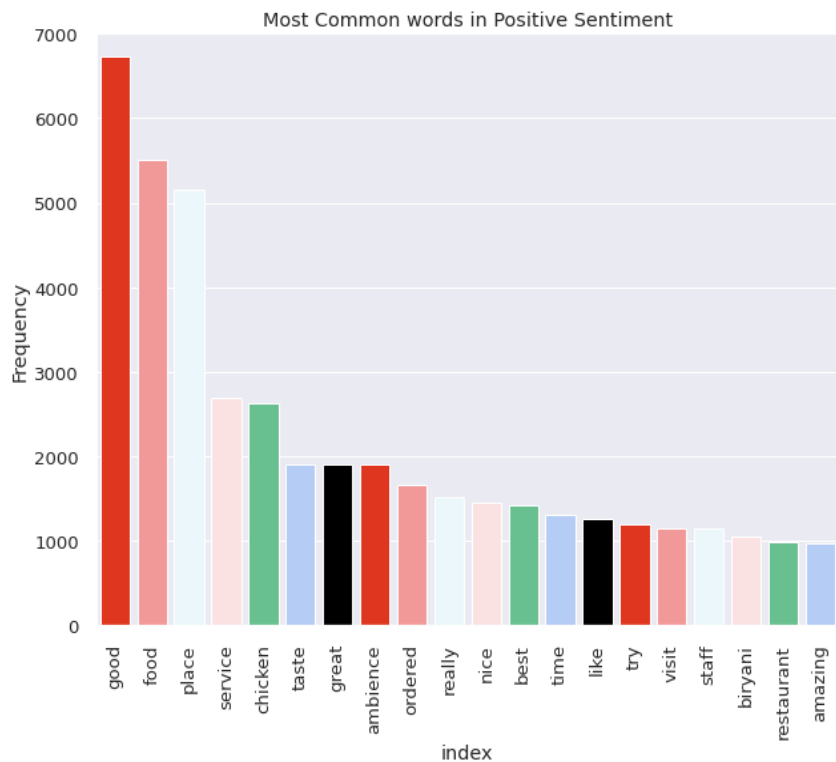
# Sentiment Analysis (continued)

## Vader Model



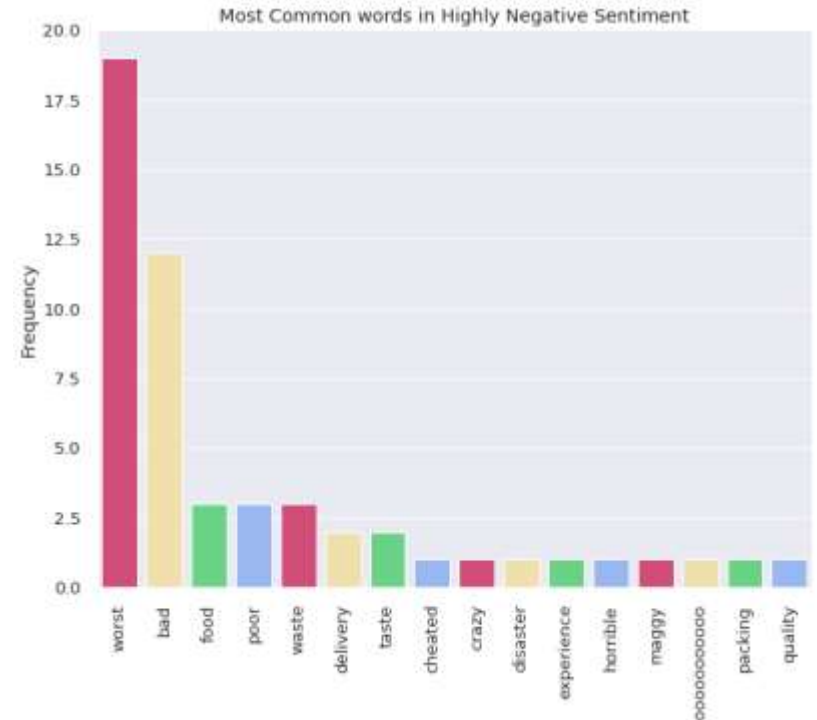
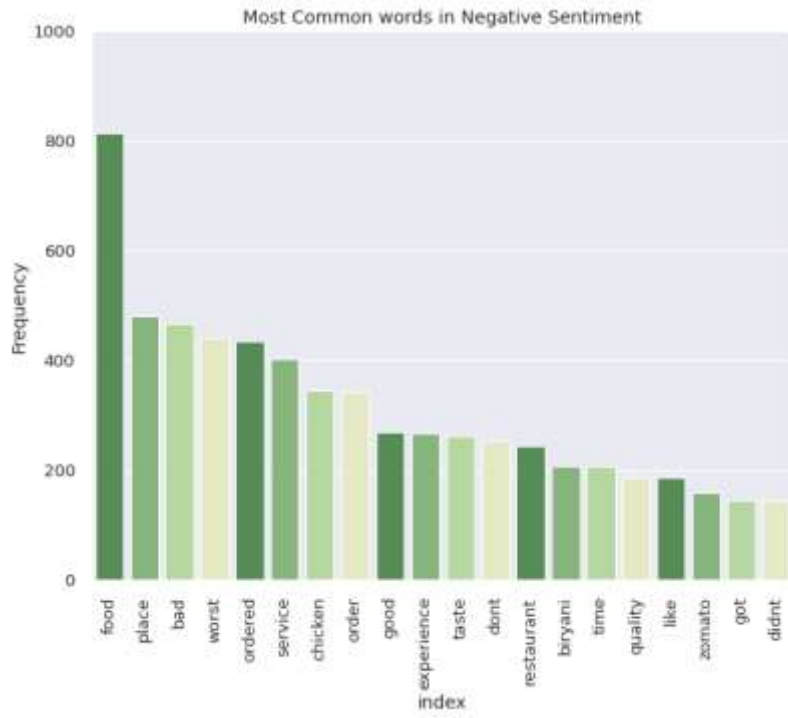
# Sentiment Analysis (continued)

## Vader Model



# Sentiment Analysis (continued)

## Vader Model





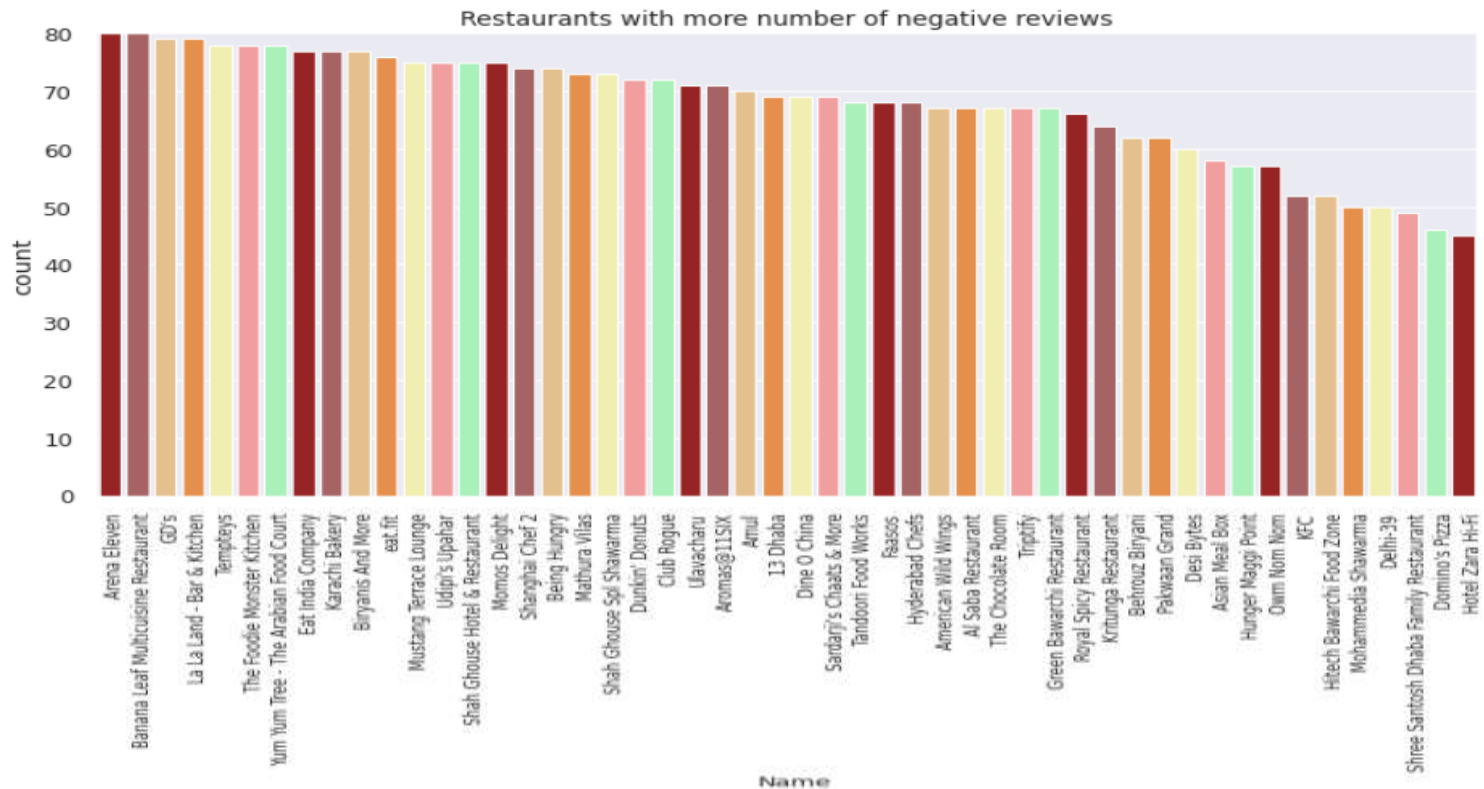
# Sentiment Analysis (continued)

## Vader Model



# Sentiment Analysis (continued)

## Vader Model



# Conclusion

- **North Indian** cuisine is most common cuisine found in the restaurants.
- **Collage - Hyatt Hyderabad Gachibowli** is most expensive restaurant.
- **Amul** and **Mohammedia Shawarma** are the most affordable restaurants.
- The Restaurants are clustered on cuisines into **15** clusters by using **KMeans** clustering algorithm with the **Silhouette score** of **0.195**.
- **DBSCAN** algorithm is also used to cluster the restaurants into **15** clusters and also helps us to detect the outliers with the **Silhouette score** of **0.107**.
- **Anvesh Chowdary** has given the most number of reviews.
- **AB's – Absolute Barbecues** is the top rated restaurant.
- Almost **79** percent of the observations have **Positive** sentiment and **14** and **7** percent of the observations have **Neutral** and **Negative** sentiments respectively

# Conclusion

- **Good** is the most common word in the Highly positive sentiment.
- **Worst** is the most common word in the Highly negative sentiment.
- **AB's – Absolute Barbecues, The Indi grill** and **B-Dubs** are the restaurants with the most number of positive reviews.
- **Arena Eleven** and **Banana leaf Multicuisine** restaurant are the restaurants with the most number of negative reviews.
- **Udipi's Upahar** is the most **affordable** restaurant with the **best** rating.
- **Feast - Sheraton Hyderabad Hotel** is the most **expensive** restaurant with the **best** rating.
- **Asian Meal Box** is the most **affordable** restaurant with the **worst** rating.
- **Club Rogue** is the **expensive** restaurant with **worst** rating.

**Thank you**