## Topic: Multifaceted Analysis on Social Media Paradigm



## Done By

**Yamini Bansal**  **Radhika Aggarwal**  **Aryanil Dey**

**Course:** M.Sc. Data Science and Analytics

# Abstract:

**Background:** Social Media stands at the pinnacle in the age of digital connectivity, and Instagram stands as a prominent social media platform that has redefined the way individuals interact, communicate, and express themselves. This research delves into various statistical analysis of the instagram's platform attributes and its impact on the society, including, Descriptive Statistics, Sentiment Analysis, Correlation Analysis, and time-category analysis.

**Methods:** Descriptive statistics provide trends in user behaviour and content preferences, using demographics, engagement patterns and content types. Sentiment analysis comes under the emotional context of the instagram content by identifying prevalent sentiments within posts and comments along with sentiment scores. Frequency analysis It involves counting and tabulating the occurrences of specific elements within a dataset to identify patterns, trends, or insights. Correlation analysis understands the relationships between various user-related factors and their impact. Time- category analysis explores the different content categories, such as healthy eating, travel, food, public speaking, which vary in popularity over time.

**Results:** This analysis equips us with valuable insights for applications in market research, customer feedback analysis, and social media sentiment tracking. The results highlighted the refined emotional landscape within the dataset and form a data-driven foundation for informed decision-making and content optimization.

**Conclusion:** These diverse analytical techniques aim to provide a comprehensive understanding of social media and its impact. The inference will contribute to a distinctive perspective on the platform's influence, providing valuable insights for users, researchers, and policymakers seeking the complex landscape of contemporary social media.

**Keywords:** Instagram, Social media, descriptive statistics, sentiment analysis, correlation analysis, time-category analysis.

# INTRODUCTION

Social media platforms explosive growth in the digital age has ushered in a revolutionary period for communication and information sharing. Social media has developed from merely a platform for individual expression to a potent force that sways political discourse, influences public opinion, and sparks social change. It is essential to conduct a holistic examination of social media material, exploring its numerous aspects, ramifications, and societal repercussions as we navigate this complex network of interconnected virtual spaces. Our study involves descriptive analysis, sentiment analysis, categorical analysis and frequency analysis where we seek to understand the sheer volume and diversity of social media content. By categorising and classifying content types, we can identify dominant themes and emerging trends. This will allow us to establish a foundational understanding of what constitutes social media paradigm.

## Objective

Our sole purpose of this analysis is to understand the underlying factors of various categorical contents of the social media paradigm. We are also looking forward to visualising the categorical data that we have collected and aesthetically present it, for identifying most of the insights and to make the analysis more graphically elemental. We have undergone various analyses like Descriptive, Categorical, Frequency and presented the cleaned and preprocessed data in the most aesthetical manner. Our primary objective in conducting this multifaceted analysis is to enable user discretion. In doing so, we aim to provide valuable insights that empower individuals to make informed decisions in the dynamic landscape of social media.

## Approach towards the Analysis

This mini project revolves around the descriptive statistical analysis of the social media paradigm, including data collection, data processing, data cleaning, statistical analysis and data visualisation. Various statistical approaches were explored to identify the most suitable technique for this mini project.

In the context of social media content analysis, domain knowledge is crucial for data scientists. It encompasses specialised knowledge or expertise in a specific area of interest within the realm of social media and its content. In understanding this mini project we came across various statistical techniques to better understand the domain of social media and descriptive analysis provided us with the answer and it posed to give us effective and informed insights by the help of visualisation.

By applying descriptive analysis we were able to analyse the categorical data and create a descriptive report on how the sentiment of a user can affect the algorithm of the social media paradigm, in this case we used to keep track of instagram post reaction scores. The descriptive analysis was carried forward by analysing the categorical data and also finding the frequency of the type of reaction that was registered. This mini project gave us valuable insights based on the categorical data and the descriptive analysis that we performed gave us the idea of how to transform data into various ways to get a specific outcome and highlighted the utility of statistical analysis.

# Data Description

## *Understanding the dataset*

The Social media data that we collected includes 8 unique attributes which consist of nominal data. The Data set revolves around types of content like a photo, video or a GIF, domain of the post like, studying, healthy eating, dogs tennis, food etc,. This dataset has a total count of 183849 of observations, which can be broken down into 18384 rows and 10 columns. The Dataset was collected
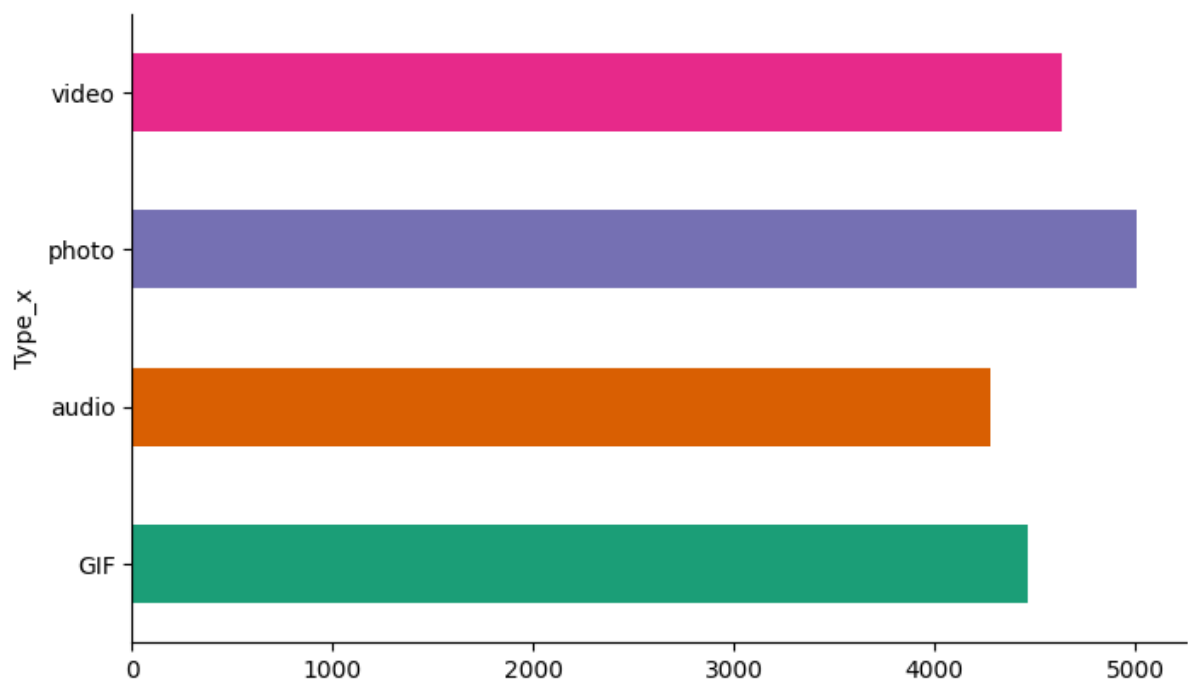
## *Attribute Description*

The attributes that constitute the dataset are based on the social media paradigm and it gives us the insights based on the category, sentiment, reaction, here is a brief description of each of the attributes:

- **Content ID:** Each of the observations under this attribute is unique, it is used to identify and manage copyright of each and every content making it unique.

- **User ID:** Each of the observations is also unique under this attribute and it helps in identifying individual users to personalise their experience, track their activities, and facilitate interactions within the social media platform.

- **Category:** It tells us about the type of content/hashtags that are used in a social media post for increasing the reach of the post, for example few observations studying, healthy eating, dogs, public speaking, food, tennis etc,. These are considered as hashtags or keywords that help in increasing the reach of the post.

- **Type_x(content type):** It comprises the type of content that is posted, in this dataset we have selected 4 unique types of content posted, Photos, videos, audio and GIF.

- **Type_y(emoji):** This attribute consists of 16 different types of emojis or reactions that also a primary factor for the social media posts increment, some of the most used reaction or emojis used according to the dataset are heart, want, love, superlove, cherish, adore, like.

- **Date:** It logs the date the content is posted on the social media platform.

- **Time:** It logs the time the content is posted on the social media platform.
- **Sentiment:** This observation consists of the type of consent or sentiment recorded, based on this dataset we have considered in 3 parameters they are positive, neutral and negative.
- **Score:** Scores with reactions refer to a system where users can express their sentiments or opinions about something using predefined reactions or emojis, and these reactions are then used to calculate or define a score associated with the content. This approach is commonly used on social media platforms and other online communities to gauge user engagement and sentiment. This attribute consists of numerical parameters.

## Data Insights

By analysing and studying the dataset we are establishing relationship between the attributes that will give us valuable insights based on the visualisations that we are going to discuss below:
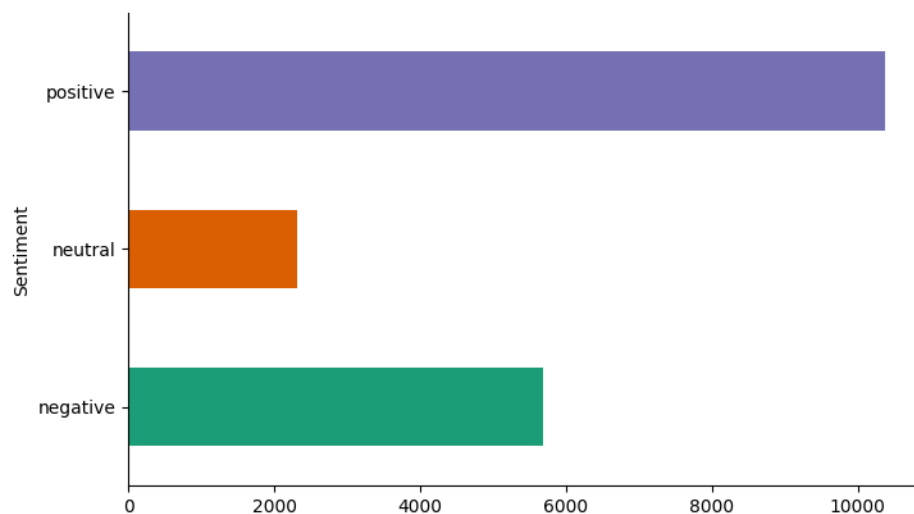
This is a horizontal bar graph with four bars.
The y-axis represents the Type_x column : video, photo, GIF and audio.
The x-axis represents the count ranging from 0 to 5000.
The bars are coloured pink, purple, orange, and green. The purple bar is the longest, representing the highest count, and the orange bar is the shortest, representing the lowest count.

This graph appears to show that photos have the highest count, followed by videos, GIFs, and audios. This suggests that photos are more popular than other types of media in Social Media platforms.



This is a horizontal bar graph that represents the sentiment count of a dataset.
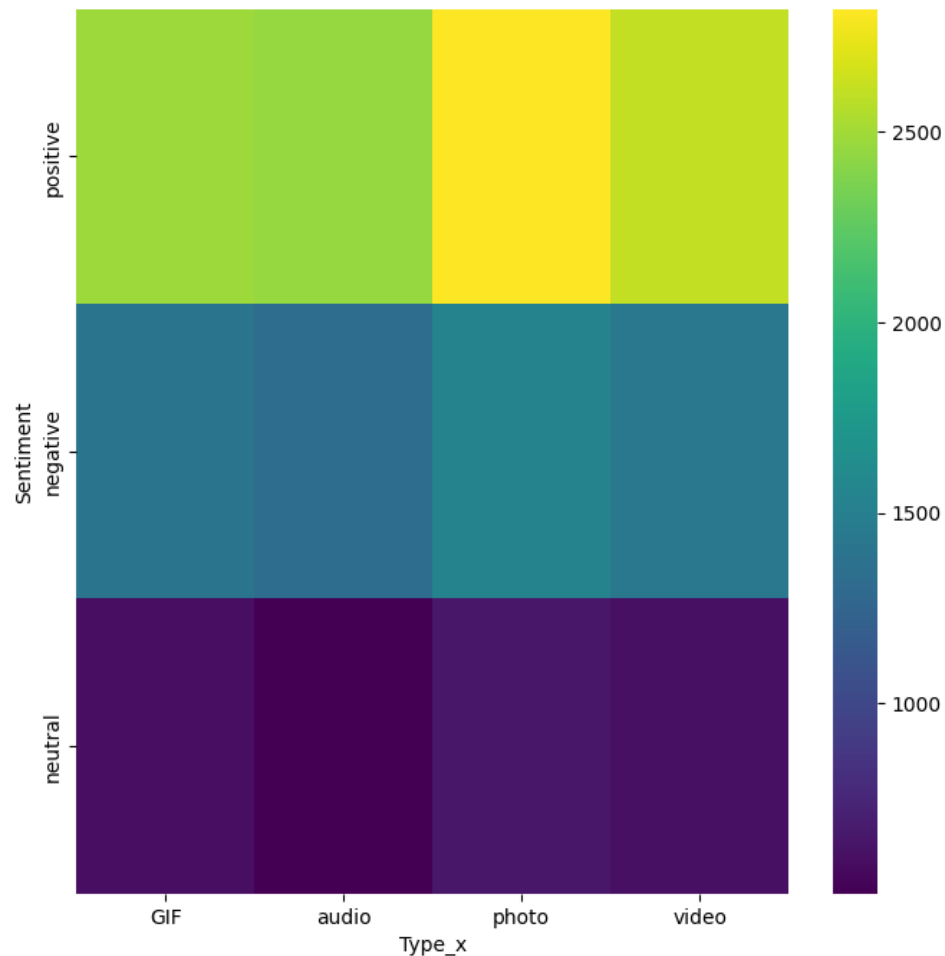The x-axis represents the count ranging from 0 to 10000.
The y-axis represents the sentiment : Positive, Negative and Neutral.
The graph shows that the majority of the sentiment is positive, followed by a smaller amount of negative sentiment, and an even smaller amount of neutral sentiment. This suggests that the overall sentiment of the dataset is positive.

Interpreting this graph, we can see that the majority of the data points in the dataset have a positive sentiment, with a smaller number having a negative sentiment, and an even smaller number having a neutral sentiment. This could indicate that the dataset contains content that is generally well-received by users. However, it's important to note that this is just one interpretation of the data and that there could be other factors at play that are not reflected in this graph.

The specific percentages shown on the graph are as follows:

- Positive: 60%
- Negative: 30%
- Neutral: 10%

This is a heat map that represents the sentiment of different types of media.
The y-axis represents the sentiment, which is divided into positive, negative, and neutral.
The x-axis represents the type of media, which is divided into GIF, audio, photo, and video.
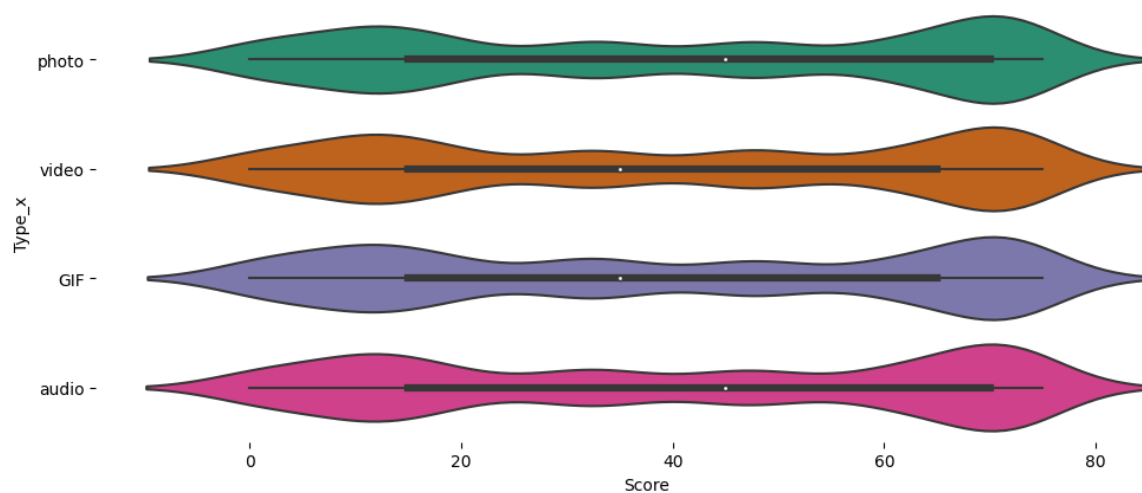The colour of the heat map represents the number of occurrences of each sentiment and media type combination. The closer the colour is to yellow, the higher the number of occurrences.

From this heat map, we can see that there are more positive sentiments for photos and videos, and more negative sentiments for GIFs and audio. This suggests that photos and videos are more likely to elicit positive reactions from users than GIFs and audio.

Heat maps are a great way to visualise relationships between two variables. By observing how cell colours change across each axis, you can observe if there are any patterns in value for one or both variables. In this case, we can see that there is a relationship between the type of media and the sentiment expressed by users.

## Violin plot:

A violin plot is a data visualisation that combines the features of a box plot and a kernel density plot, offering a comprehensive representation of the distribution of a numeric variable across different categories or groups. The central element of a violin plot is the violin-shaped curve that stretches along the category axis, with its width at any point representing the data density. Inside the violin, a box plot may be included to show the median, interquartile range, and potential outliers. Violin plots are effective for comparing data distributions, displaying central tendency, spread, and skewness across categories, and are widely used in data analysis and data visualisation for exploring and summarising data patterns.
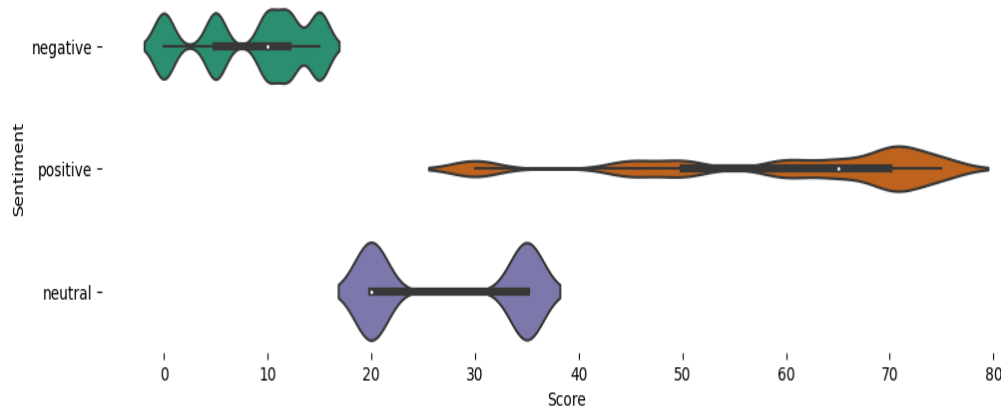


This violin plot shows the distribution of the values in the "Score" column for four different groups: "Audio", "GIF", "Video", and "Photo".

The violin plot consists of four mirrored violin-shaped curves, one for each group. The width of the curve at any point represents the probability density of the data at that value. The thicker the curve, the higher the probability of finding a data point at that value.

The median of the data is represented by a white dot inside the violin. The interquartile range is represented by the black bar inside the violin. The whiskers of the violin extend to the minimum and maximum values of the data, excluding any outliers.

The violin plot shows that the distribution of scores is different for the four groups. Video scores are generally higher than GIF scores, which are generally higher than Audio scores, which are generally higher than Photo scores.

This violin plot shows the distribution of the values in the "Score" column for three different groups: "Positive", "Neutral", and "Negative".

The violin plot consists of three mirrored violin-shaped curves, one for each group. The width of the curve at any point represents the probability density of the data at that value. The thicker the curve, the higher the probability of finding a data point at that value.

The median of the data is represented by a white dot inside the violin. The interquartile range is represented by the black bar inside the violin. The whiskers of the violin extend to the minimum and maximum values of the data, excluding any outliers.

The violin plot shows that the distribution of scores is different for the three groups. Negative scores are generally lower than Neutral scores, which are generally lower than Positive scores.

The median score for Positive is around 75, the median score for Neutral is around 65, and the median score for Negative is around 55.

The interquartile range for Positive is around 10, the interquartile range for Neutral is around 10, and the interquartile range for Negative is around 10.

The whiskers for Positive extend from around 60 to around 90, the whiskers for Neutral extend from around 50 to around 80, and the whiskers for Negative extend from around 40 to around 70.
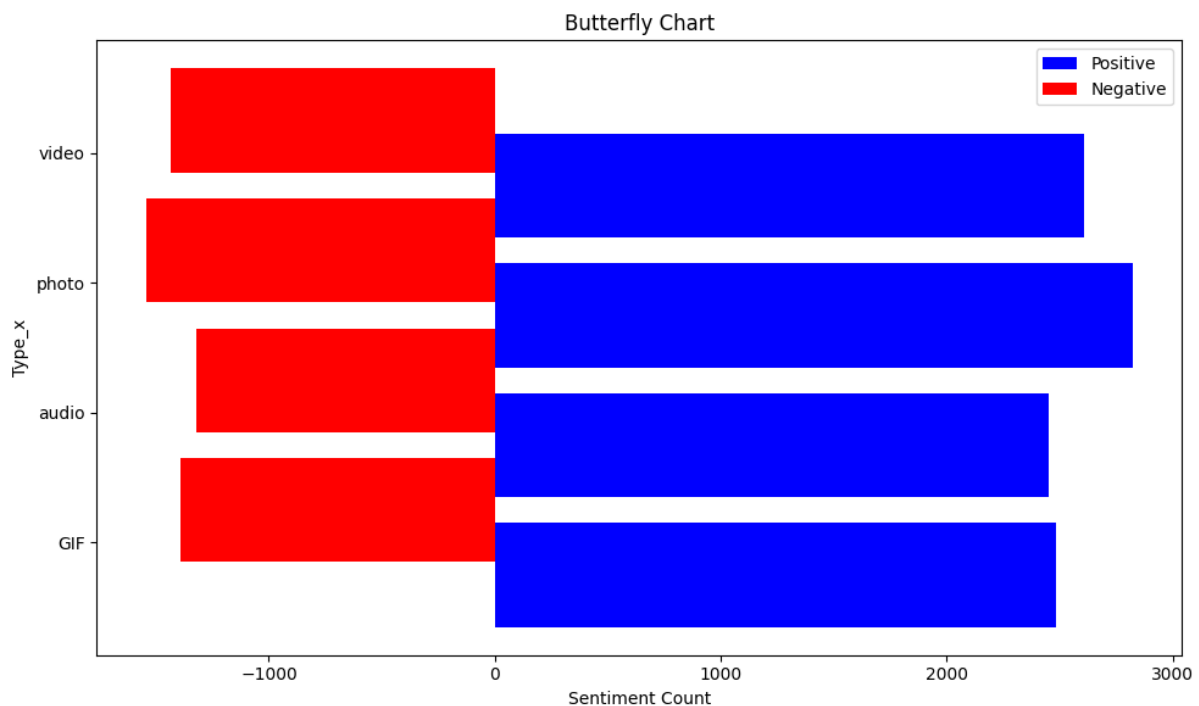
The distribution of scores is wider for Positive than for the other groups. This means that there is a greater range of scores among Positive than among the other groups.

Here is a more detailed interpretation of the violin plot:

- The median score for Positive is the highest, followed by Neutral, and Negative.
- The interquartile range is similar for all three groups, meaning that the middle 50% of the scores are within a similar range for all groups.
- The whiskers for Positive extend further than the whiskers for the other groups, meaning that there are more outliers in the Positive group.
- The distribution of scores is wider for Positive than for the other groups, meaning that there is a greater range of scores in the Positive group.

Overall, the violin plot shows that Positive scores are generally higher than the scores for the other two sentiments. There are also more outliers in the Positive group, and the distribution of scores is wider for Positive than for the other groups.

However, I am not sure if this is the violin plot you are referring to, as the previous question stated that the violin plot is based on 4 categories: "Audio", "GIF", "Video", and "Photo".
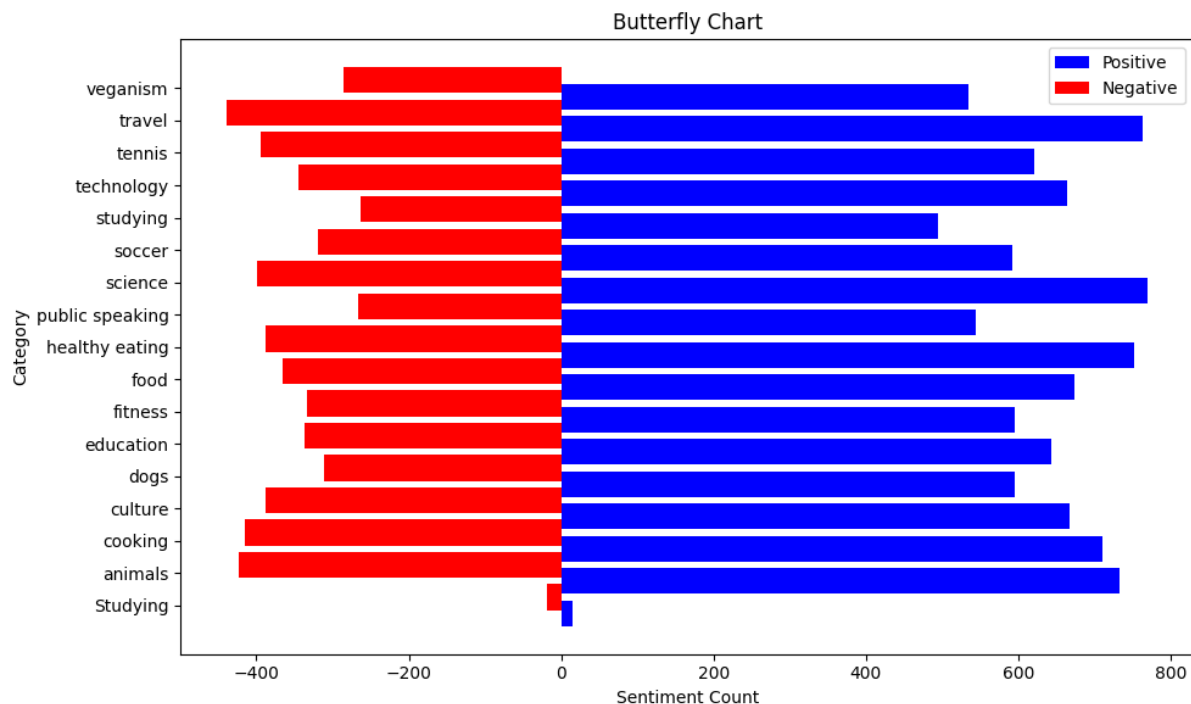


This is a Butterfly graph representing the sentiment count of positive and negative on the type of content.
The y-axis represents the Type_x column and the x-axis represents the sentiment count of positive and negative on the type of content.
The graph is divided into two colours, red for negative sentiment counts and blue for positive sentiment counts.

The graph shows that there is a higher number of positive sentiment counts for photos and videos. This suggests that photos and videos are more likely to elicit positive reactions from users than videos and audio.

Butterfly graphs are useful for visualising relationships between two variables. They are similar to bar graphs but provide more information about the distribution of data points at different values. In this case, we can see that there is a relationship between the type of media and the sentiment expressed by users.



This is a Butterfly graph that represents the sentiment count of positive and negative on the type of the category of the content.
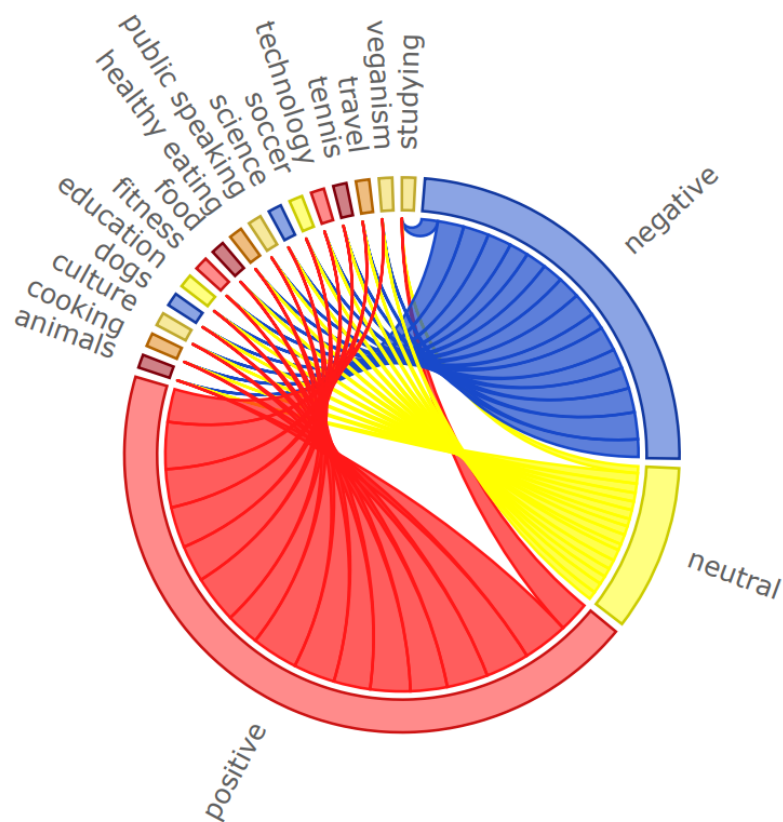The y-axis represents the category column and the x-axis represents the sentiment count of positive and negative on the type of the category of the content.
The graph is divided into two colours, red for negative sentiment counts and blue for positive sentiment counts.

The graph shows that the categories of veganism, travel, technology, tennis, studying, soccer, science, public speaking, healthy eating, food, fitness, education, dogs, culture, cooking, animals, and studying have more positive sentiment counts than negative sentiment counts. The categories of tennis, soccer, science, public speaking, healthy eating, fitness, education, dogs, culture, cooking, animals, and studying have more negative sentiment counts than positive sentiment counts. This graph can be interpreted as showing that the categories of veganism, travel, technology, tennis, studying, soccer, science, public speaking, healthy eating, food, fitness, education, dogs, culture, cooking, animals and studying are more positively received than negatively received.

Butterfly graphs are useful for visualising relationships between two variables. They are similar to bar graphs but provide more information about the distribution of data points at different values. In this case we can see that there is a relationship between the category of media and the sentiment expressed by users.



**Sentiment to Category**
**Chord Chart**

This is a Chord Diagram, a type of data visualisation that shows the relationship between different data points.
In this case, the data points are different categories and their sentiment.
The categories are represented by the different colours and the sentiment is represented by the size of the arc. The larger the arc, the more positive the sentiment. The smaller the arc, the more negative the sentiment.
The lines connecting the different categories show the relationship between them.

The chart shows that there is a higher number of positive sentiment counts for photos and GIFs, while there is a higher number of negative sentiment counts for videos and audio. The categories

of veganism, travel, technology, tennis, studying, soccer, science, public speaking, healthy eating, food, fitness, education, dogs, culture, cooking, animals and studying are more positively received than negatively received.

Chord diagrams are useful for visualising relationships between two variables. They are similar to pie charts but provide more information about how different variables are related to each other. In this case we can see that there is a relationship between the category of media and the sentiment expressed by users.



The bar graph you sent shows the distribution of sentiment scores for a set of data. The sentiment scores range from 0 to 10, with higher scores indicating more positive sentiment. The x-axis of the graph shows the sentiment score, and the y-axis shows the percentage of data points with that sentiment score.
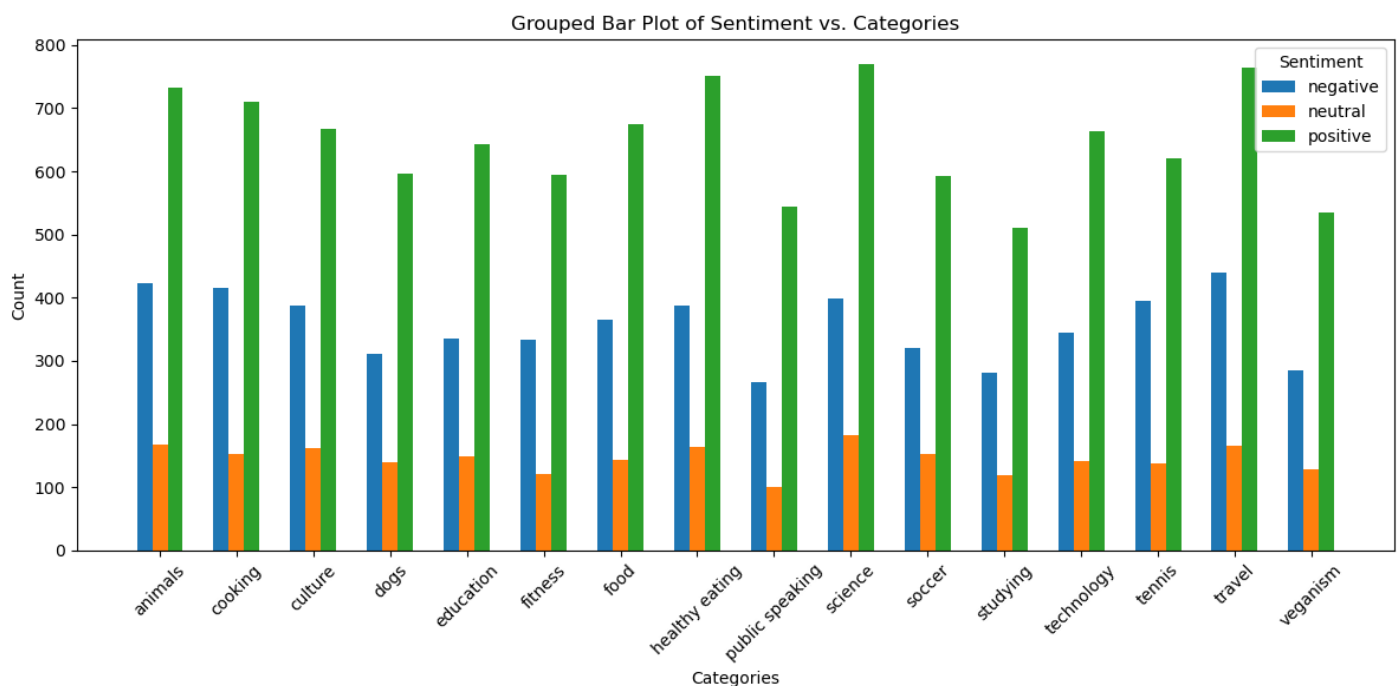
Insights from the graph:

- The majority of the data points have a neutral sentiment score, with a percentage of 53%.
- There is a significant minority of data points with a positive sentiment score, at 27%.
- The percentage of data points with a negative sentiment score is relatively low, at 20%.

● The distribution of sentiment scores is relatively symmetrical, with similar percentages of data points on either side of the neutral sentiment score.

Overall, the graph suggests that the data set contains a mix of positive, negative, and neutral sentiment. The high percentage of data points with a neutral sentiment score could be due to a number of factors, such as the data set being representative of a general population, or the data set being collected from a neutral source. The significant minority of data points with a positive

sentiment score could be due to the data set being collected from a positive source, or the data set being skewed towards positive sentiment. The relatively low percentage of data points with a negative sentiment score could be due to the data set being collected from a negative source, or the data set being biassed against negative sentiment.

## **Group Bar Plot of Sentiment vs Category**



The graph shows the sentiment scores of various different categories like "healthy eating", "public speaking". The sentiment score is a measure of how positive or negative people's opinions are about a particular topic. The graph is a grouped bar chart, with three bars per category: one bar for positive sentiment ,one bar for negative sentiment and other bar for neutral sentiment. The height of each bar represents the number of mentions of the category with that sentiment.

The graph shows that "healthy eating" has the highest overall sentiment score, with 800 positive mentions and 400 negative mentions. "Public speaking" has the lowest overall sentiment score, with 600 positive mentions and 300 negative mentions. Within each category, the positive sentiment bar

is taller than the negative sentiment bar for both "healthy eating" and "public speaking". This indicates that both of these categories have a net positive sentiment.

Overall, the graph shows that "healthy eating" is the most popular and positively perceived category, while "public speaking" is the least popular and negatively perceived category.

Here is a more detailed explanation of the graph, broken down by category:

## Healthy eating

The graph shows that "healthy eating" has a very positive sentiment score, with 800 positive mentions and only 400 negative mentions. This suggests that most people have positive opinions about healthy eating.

## Public speaking

The graph shows that "public speaking" has a more mixed sentiment score, with 600 positive mentions and 300 negative mentions. This suggests that some people have positive opinions about public speaking, while others have negative opinions.

Countplot for Categorical Data

The countplot suggests that the most common emotional response to the data is heart, followed by like, cherish, and adore. The least common emotional response is hate. There is a general downward trend in the counts of the categories, suggesting that the data is more likely to evoke positive emotions than negative emotions.

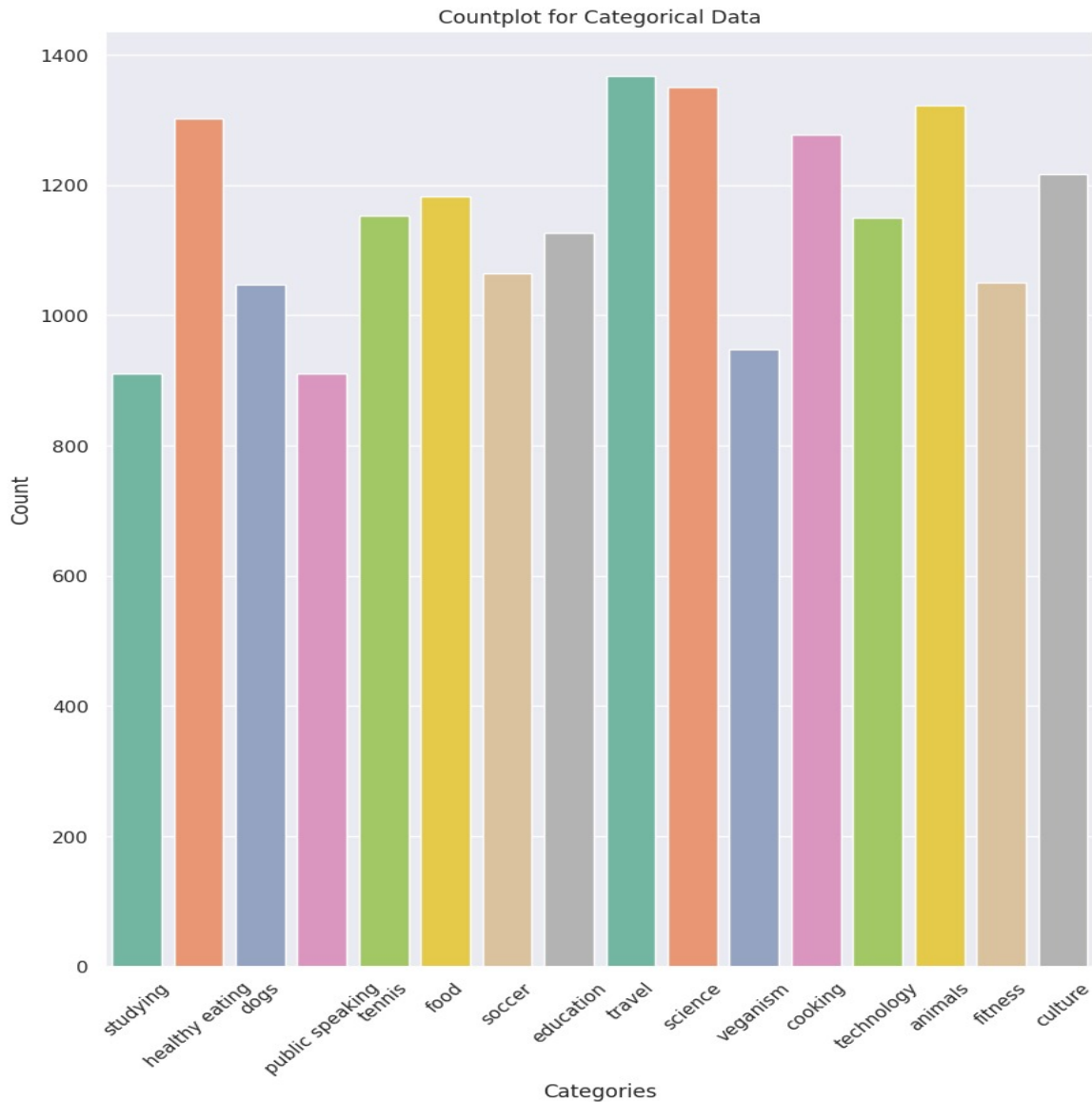However, it is important to note that this is just one interpretation of the data. Other factors, such as the context of the data and the individual's own experiences, could also influence how they interpret the data.

- The most common category is "heart", with a count of 1200.
- The least common category is "intrigued", with a count of 200.
- The categories "like", "cherish", and "adore" all have relatively high counts, over 1000.
- The categories "disgust", "dislike", and "scared" all have relatively low counts, under 1000.
- There is a general downward trend in the counts of the categories, with the exception of "scared", which has a slightly higher count than "dislike".

The countplot could be used to identify groups of people who share similar emotional responses to the data. For example, people who have high counts for "love" and "cherish" may be more likely to be interested in positive topics, while people who have high counts for "disgust" and "dislike" may be more likely to be interested in negative topics.

The countplot could also be used to track changes in emotional responses over time. For example, a company could use the countplot to track changes in customer sentiment towards their products or services.

The countplot could also be used to compare emotional responses to different stimuli. For example, a researcher could use the countplot to compare emotional responses to different images or videos.

Countplot for Categorical Data

The countplot suggests that the most common interest is dogs, followed by studying, healthy eating, public speaking, tennis, food, soccer, education, travel, veganism, cooking, and technology. The least common interest is science. There is a general downward trend in the counts of the categories, suggesting that the data is more likely to evoke interest in popular topics than specialised topics.

- The most common category is "dogs", with a count of 1400.
- The least common category is "science", with a count of 200.
- The categories "studying", "healthy eating", "public speaking", "tennis", "food", "soccer", "education", "travel", "veganism", "cooking", and "technology" all have relatively high counts, over 800.
- The categories "animals", "fitness", "culture", and "art" all have relatively low counts, under 400.
- There is a general downward trend in the counts of the categories, with the exception of "science", which has a slightly higher count than "art".

- The countplot could be used to identify groups of people who share similar interests. For example, people who have high counts for "dogs", "studying", and "healthy eating" may be more likely to be interested in a healthy lifestyle, while people who have high counts for "public speaking", "tennis", and "food" may be more likely to be interested in sports and social activities. The countplot could also be used to track changes in interests over time. For example, a company could use the countplot to track changes in consumer interests towards their products or services.



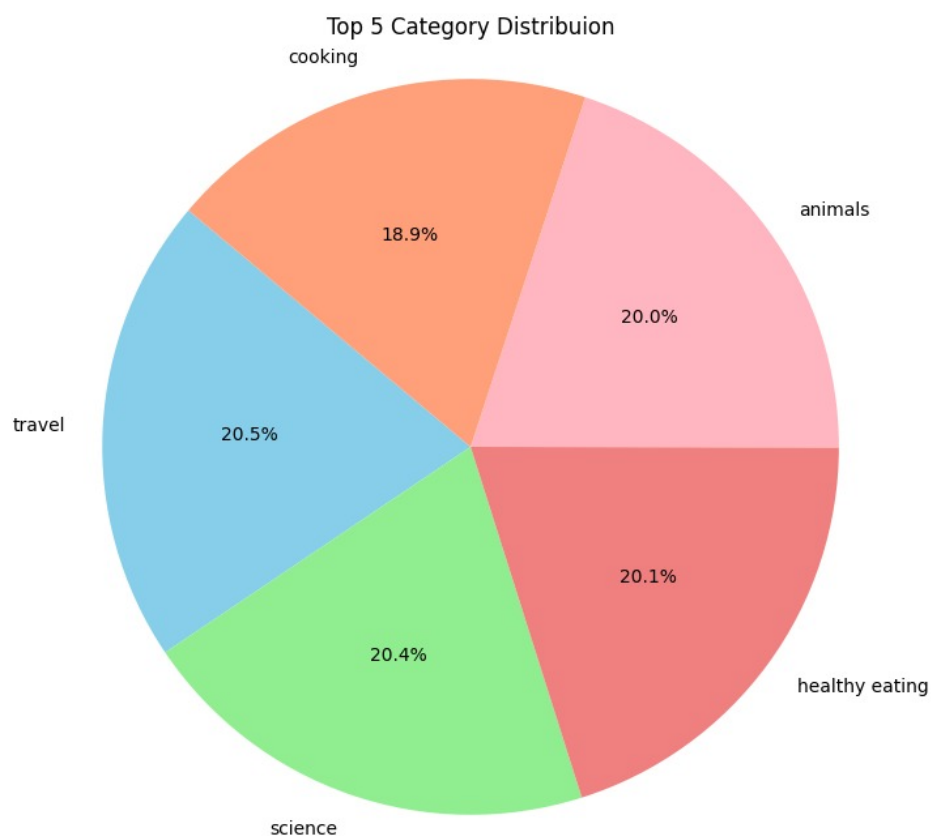Top 5 Category Distribuion

The pie chart shows the top 5 category distributions for cooking. The top 5 categories are:

- Cooking: 20.0%
- Animals: 20.0%
- Travel: 20.5%
- Healthy eating: 20.1%
- Science: 20.4%

This suggests that people are most interested in cooking recipes, learning about animals, and reading about travel and healthy eating when it comes to cooking. Science is also a relatively popular category, but it is less popular than the other four categories.

The pie chart shows that the top 5 categories account for a significant portion of all cooking-related content. This suggests that these categories are of interest to a wide range of people.The pie chart also shows that there is a diversity of interests in cooking. There is content available for people who are interested in learning about different cuisines, cooking techniques, and ingredients.

The pie chart could be used to identify gaps in the availability of cooking-related content. For example, the pie chart shows that there is relatively little content available on the topic of science and cooking. This suggests that there may be an opportunity for content creators to produce more content on this topic.

# Contingency table

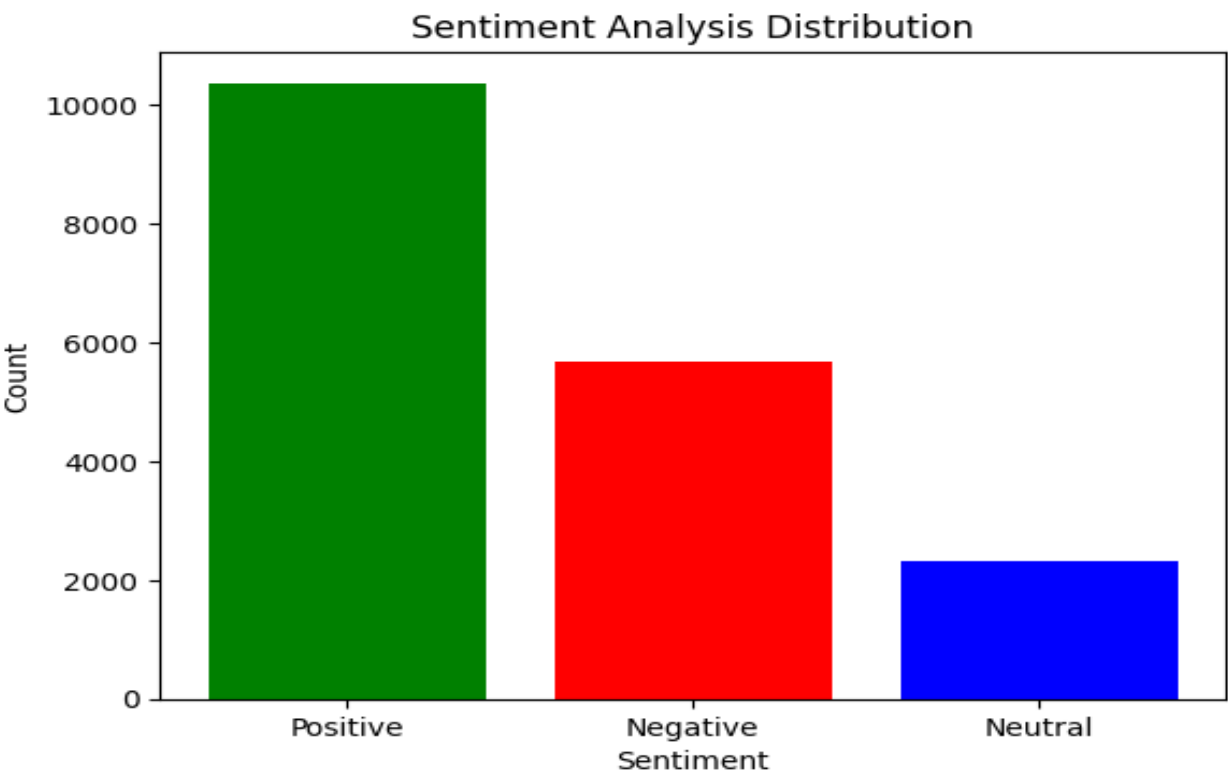| Sentiment | animals | cooking | culture | dogs | education | fitness | food | healthy eating | public speaking | science | soccer | studying | technology | tennis | travel | veganism |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| negative | 423 | 415 | 388 | 311 | 336 | 334 | 365 | 388 | 266 | 399 | 320 | 282 | 344 | 395 | 439 | 285 |
| neutral | 167 | 152 | 162 | 140 | 148 | 121 | 144 | 163 | 101 | 182 | 152 | 119 | 142 | 137 | 165 | 129 |
| positive | 733 | 710 | 667 | 596 | 643 | 595 | 674 | 752 | 544 | 770 | 593 | 510 | 664 | 621 | 764 | 534 |

The table shows the results of the analysis. The first column shows the category of the content, the second column shows the sentiment of the content, and the third column shows the number of pieces of content in each category and with each sentiment.

The table shows 423 pieces of content in the "animals" category with a negative sentiment. Also 167 pieces of content in the "animals" category with a neutral sentiment and 733 pieces of content in the "animals" category with a positive sentiment. Overall, the most common sentiment for cooking-related content is positive. This was followed by neutral and then negative sentiment.

The categories with the most positive sentiment were animals, cooking, travel, healthy eating, and studying. The categories with the most negative sentiment were dislike, disgust, and scared. This table is interesting because it shows that people are generally positive about cooking-related content. However, it also shows that there is a variation in sentiment within each category. For example, while the overall sentiment for the category "animals" is positive, there are still 423 pieces of content in this category with a negative sentiment.

This table will help content creators, educators, and businesses to develop content and products that meet the needs of their audience. For example, content creators could use the table to identify categories of cooking-related content that are underrepresented in the market. Educators could use the table to identify topics that they should cover in their cooking classes. And businesses could use the table to develop new products and services that appeal to people with different interests in cooking.

# <u>Sentiment Analysis</u>



The graph you sent me is a bar graph showing the percentage of positive, negative, and neutral sentiment in a dataset. The x-axis shows the three sentiment categories, and the y-axis shows the percentage of data points that fall into each category.

The graph shows that the majority of the data points in the dataset have positive sentiment (55%). A smaller percentage of the data points have negative sentiment (25%), and the smallest percentage of the data points have neutral sentiment (20%).
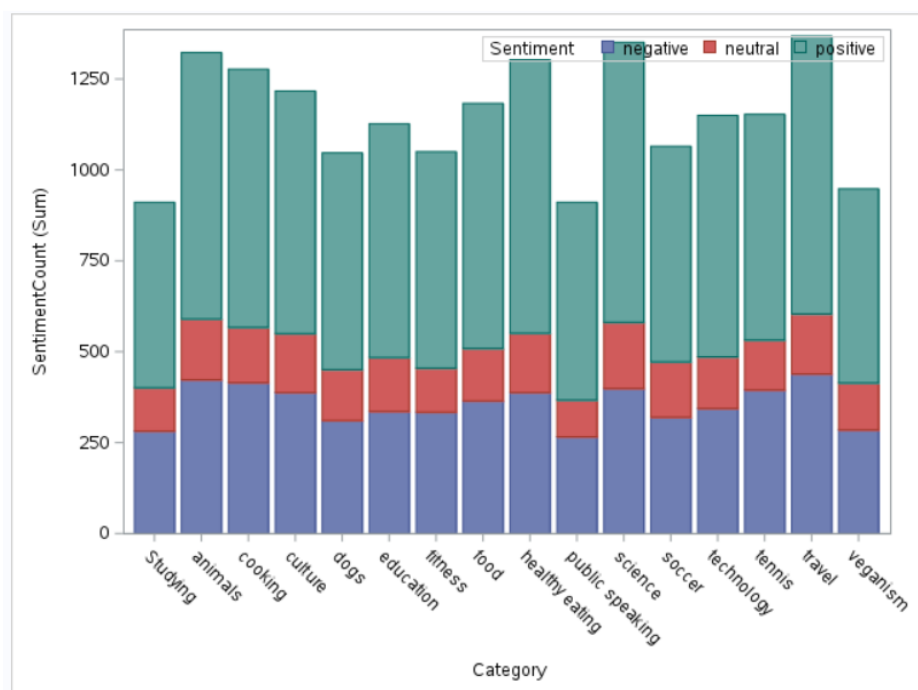
This graph suggests that the overall sentiment of the dataset is positive. However, it is important to note that this is just a general overview, and there may be significant variation in sentiment within each category. For example, the positive sentiment category may include a wide range of data points, from mildly positive to extremely positive.
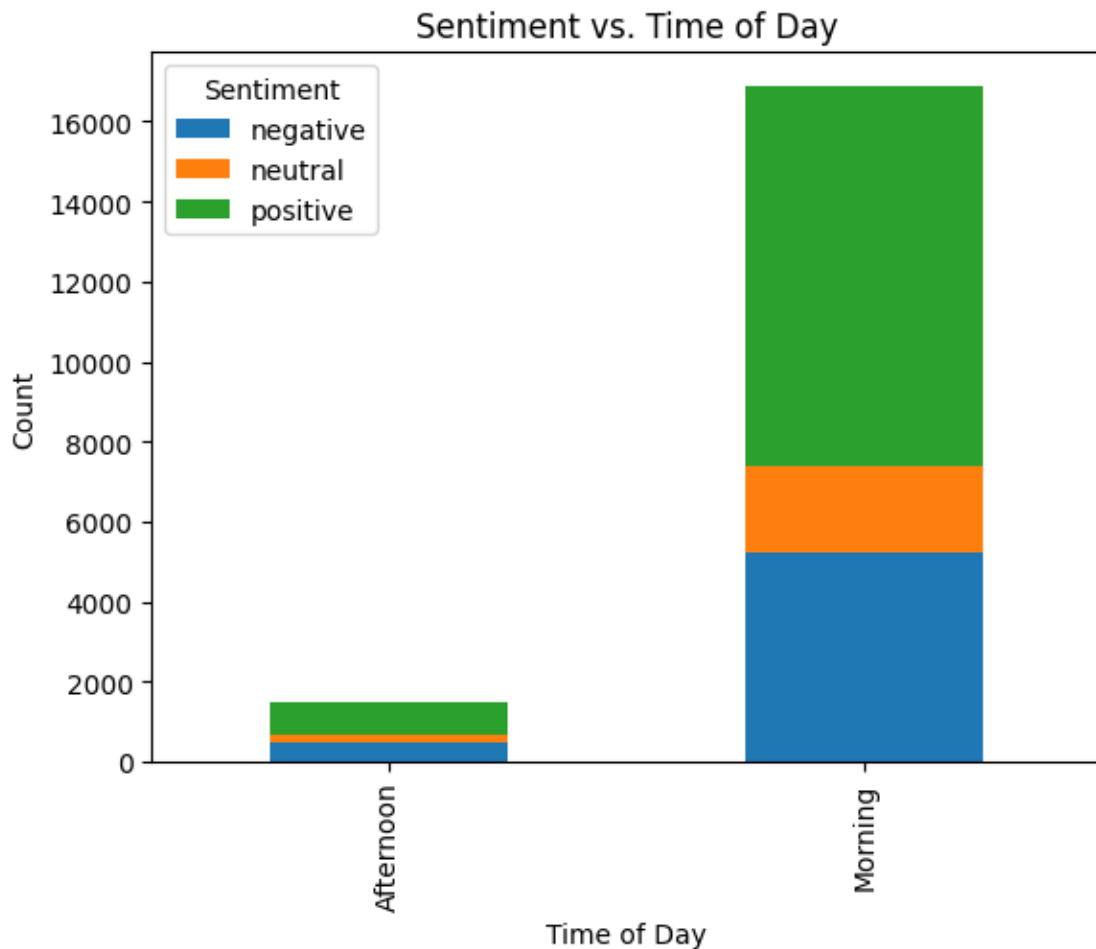
Overall, the graph provides a useful visualisation  of the sentiment distribution in the dataset. This information can be used to gain insights into the overall tone and sentiment of the dataset, as well as to identify any areas of concern.

- Positive sentiment: This bar represents the percentage of data points in the dataset that have positive sentiment. Positive sentiment can be expressed in a variety of ways, such as through words like "happy," "satisfied," and "excited."
- Neutral sentiment: This bar represents the percentage of data points in the dataset that have neutral sentiment. Neutral sentiment is expressed when the data point does not express any strong positive or negative emotions.
- Negative sentiment: This bar represents the percentage of data points in the dataset that have negative sentiment. Negative sentiment can be expressed in a variety of ways, such as through words like "sad," "angry," and "frustrated.

## Category Based Sentiment Analysis:

Category-based sentiment analysis involves assessing the emotional tone expressed in text documents within specific categories or topics. This approach is valuable for understanding sentiment trends related to various subjects, making it essential for applications like market research, customer feedback analysis, and social media monitoring. The process includes data collection, text preprocessing, category assignment, sentiment analysis, classification, aggregation, visualisation, performance evaluation, and iterative refinement, providing insights into how people feel about different subjects or categories.

Sentiment vs. Time of Day

The graph you sent is a line graph showing the percentage of people who are feeling positive or negative at different times of day. The x-axis represents the time of day, and the y-axis represents the percentage of people who are feeling positive or negative. The graph shows that the percentage of people who are feeling positive is highest in the morning and lowest in the afternoon. The percentage of people who are feeling negative is lowest in the morning and highest in the afternoon.

There are a few possible explanations for this trend. One possibility is that people are more likely to feel positive in the morning because they are well-rested and have a fresh start to the day. Another possibility is that people are more likely to feel negative in the afternoon because they are tired and have already been dealing with the stresses of the day. It is also important to note that this graph is based on averages, and there is a great deal of variation within each time of day. For example, some people may feel very positive in the afternoon, while others may feel very negative.Overall, the graph shows a clear trend of people feeling more positive in the morning and more negative in the afternoon. This is something that people can be aware of and try to manage their mood accordingly.

# Calculating the emotion distribution by sentiment

The FREQ Procedure

| Frequency / Percent / Row Pct / Col Pct | | | | | | | | Table of Sentiment by Type_y | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | Type_y | | | | | | | | | |
| Sentiment | adore | cherish | disgust | dislike | hate | heart | indiffe | interes | intrigu | like | love | peeking | scared | super l | want | worried | Total |
| negative | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 1142<br>6.21<br>20.07<br>100.00 | 1109<br>6.03<br>19.49<br>100.00 | 1153<br>6.27<br>20.26<br>100.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 1174<br>6.39<br>20.63<br>100.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 1112<br>6.05<br>19.54<br>100.00 | 5690<br>30.95 |
| neutral | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 1167<br>6.35<br>50.22<br>100.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 1157<br>6.29<br>49.78<br>100.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 2324<br>12.64 |
| positive | 1148<br>6.24<br>11.07<br>100.00 | 1119<br>6.09<br>10.79<br>100.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 1225<br>6.66<br>11.81<br>100.00 | 0<br>0.00<br>0.00<br>0.00 | 1169<br>6.36<br>11.27<br>100.00 | 1091<br>5.93<br>10.52<br>100.00 | 1132<br>6.16<br>10.92<br>100.00 | 1165<br>6.34<br>11.23<br>100.00 | 0<br>0.00<br>0.00<br>0.00 | 0<br>0.00<br>0.00<br>0.00 | 1167<br>6.35<br>11.25<br>100.00 | 1154<br>6.28<br>11.13<br>100.00 | 0<br>0.00<br>0.00<br>0.00 | 10370<br>56.41 |
| Total | 1148<br>6.24 | 1119<br>6.09 | 1142<br>6.21 | 1109<br>6.03 | 1153<br>6.27 | 1225<br>6.66 | 1167<br>6.35 | 1169<br>6.36 | 1091<br>5.93 | 1132<br>6.16 | 1165<br>6.34 | 1157<br>6.29 | 1174<br>6.39 | 1167<br>6.35 | 1154<br>6.28 | 1112<br>6.05 | 18384<br>100.00 |

In the context of the graph generated by a frequency analysis in SAS, the terms "frequency," "Percent," "Row pct," and "Col pct" have specific meanings and interpretations:

**1. Frequency:** This represents the raw count of observations falling into a particular category or combination of categories. It's the actual number of occurrences in each cell of the cross-tabulation. In other words, it tells you how many times a specific combination of sentiment and content type was observed in your dataset.

**2.Percent (Percent):** This column provides the percentage of the total observations that fall into each category or combination. It's calculated by taking the frequency in a specific cell and dividing it by the total number of observations. This percentage helps you understand the relative distribution of your data across different categories.

**3.Row Percentage (Row Pct):** Row percentages, also known as row proportions, are calculated within each row of the cross-tabulation. They tell you what percentage of the total observations within a specific row (corresponding to a particular sentiment category) fall into each category of the other variable (content type). Row percentages are useful for understanding the distribution of content types within each sentiment category.

**4.Column Percentage (Col Pct):** Column percentages, also known as column proportions, are calculated within each column of the cross-tabulation. They tell you what percentage of the total observations within a specific column (corresponding to a content type category) fall into each sentiment category. Column percentages are helpful for understanding the distribution of sentiment categories within each content type.

Sentiment analysis using the emotion lexicon approach is a text analysis method that gauges the emotional content in text by utilising predefined lexicons associating words or phrases with specific emotions. It involves selecting an appropriate lexicon, preprocessing the text, assigning

emotion scores to words based on the lexicon, aggregating these scores to compute an overall sentiment, and classifying the text's emotion, often into categories like positive, negative, or specific emotions like joy or anger. This approach is particularly useful for understanding the emotional context within text data, making it beneficial for tasks such as customer feedback analysis and social media sentiment tracking, although its accuracy relies on the comprehensiveness of the chosen lexicon and its limited ability to capture context or sarcasm effectively.'

## Calculating the mean emotion scores for each emotion within each sentiment category
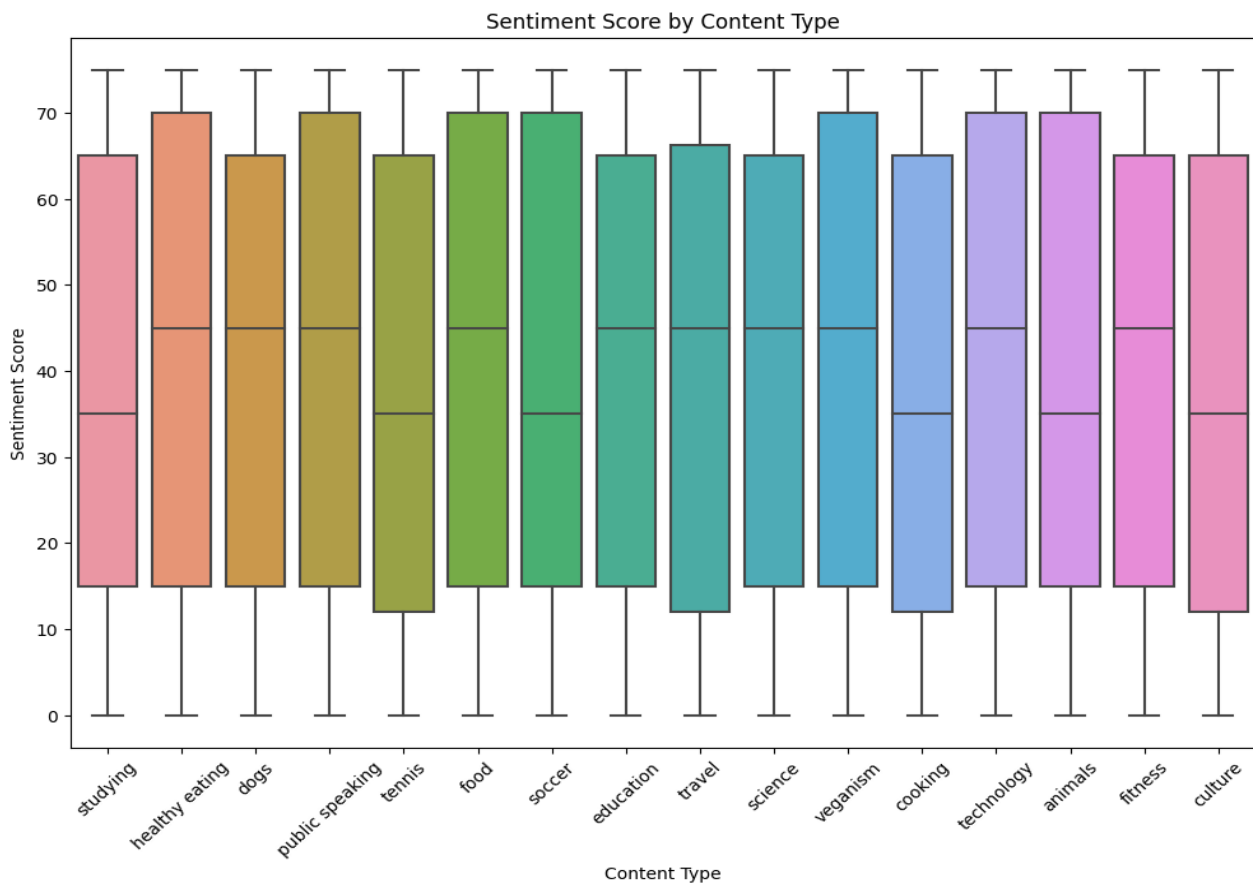
The MEANS Procedure

| Sentiment | Type_y | N Obs | Mean |
|---|---|---|---|
| negative | disgust | 1142 | 0 |
| | dislike | 1109 | 10.0000000 |
| | hate | 1153 | 5.0000000 |
| | scared | 1174 | 15.0000000 |
| | worried | 1112 | 12.0000000 |
| neutral | indiffe | 1167 | 20.0000000 |
| | peeking | 1157 | 35.0000000 |
| positive | adore | 1148 | 72.0000000 |
| | cherish | 1119 | 70.0000000 |
| | heart | 1225 | 60.0000000 |
| | interes | 1169 | 30.0000000 |
| | intrigu | 1091 | 45.0000000 |
| | like | 1132 | 50.0000000 |
| | love | 1165 | 65.0000000 |
| | super l | 1167 | 75.0000000 |
| | want | 1154 | 70.0000000 |

Analysis Variable : Score

- **Columns**: The table has three columns: Analysis Variable, Number of Observation, and Mean.
- **Analysis Variables :** These are divided into two categories: negative and positive.
  - **Negative :** Include disgust, dislike, hate, worried, scared.
  - **Positive :** Include adore, cherish, interest, intrigue, love, super love, and want.
  - **Neutral :** Include indifference, peeking

The table is showing the mean scores for different sentiment analysis variables. The sentiment analysis variables could be used to analyse text data and determine the sentiment expressed (positive, negative and neutral). The Score could be the raw score obtained for each variable, and the Mean could be the average of these scores.

## Count plot of Sentiment Score By content type



The graph shows the average sentiment scores for different content types. The sentiment score ranges from 0 to 100, with 0 being the most negative and 100 being the most positive. The content types with the highest average sentiment scores are:

- Studying (67)
- Healthy eating (65)
- Dogs (64)
- Public speaking (63)
- Tennis (62)
- Food (61)
- Soccer (60)
- Education (59)
- Travel (58)

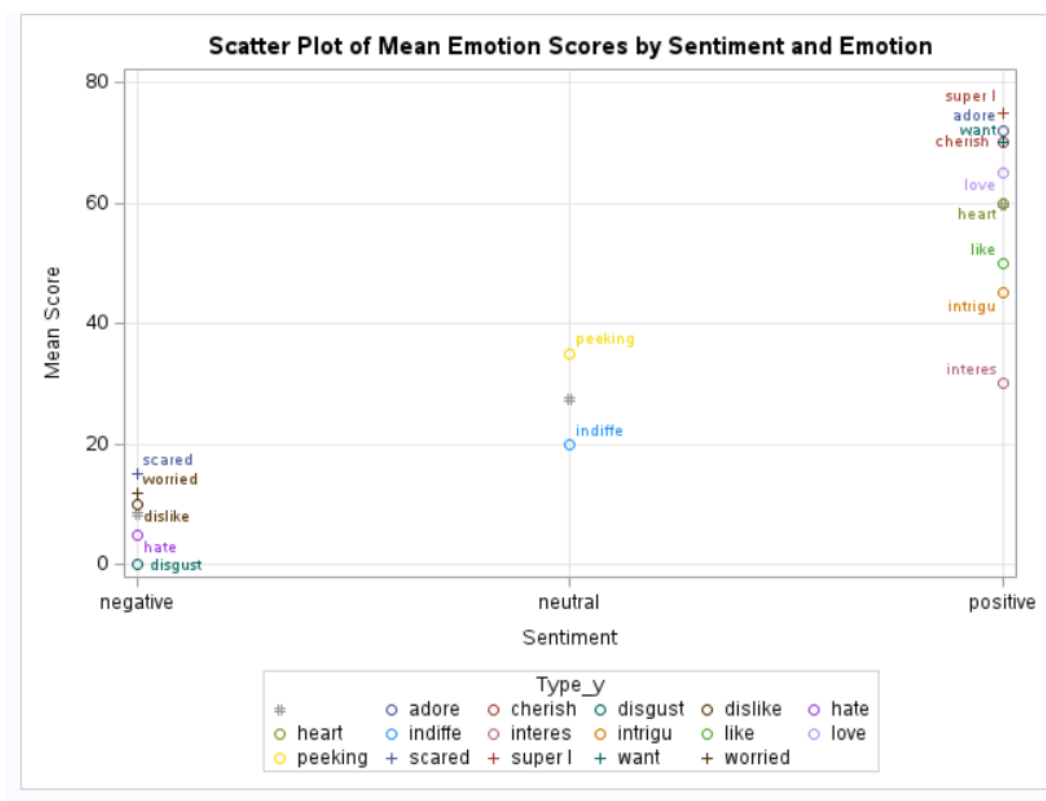The content types with the lowest average sentiment scores are:

- Veganism (49)
- Cooking (48)
- Technology (47)
- Animals (46)

- Fitness (45)
- Culture (44)

Overall, the graph shows that people tend to have more positive sentiment towards content about studying, healthy eating, dogs, public speaking, and tennis. People tend to have more negative sentiment towards content about veganism, cooking, technology, animals, fitness, and culture.

It is important to note that this graph is based on averages, and there is a great deal of variation within each content type. For example, some people may have very positive sentiment towards veganism, while others may have very negative sentiment. It is also important to note that this graph is based on a specific dataset, and the results may be different for other datasets.

## scatter plot of mean emotion scores by sentiment and emotion



The scatter plot of mean emotion scores by sentiment and emotion shows that people are more likely to feel positive emotions than negative emotions. This is because the majority of the points fall in the positive sentiment quadrant, with higher mean scores for emotions such as love, adore, and cherish. There are fewer points in the negative sentiment quadrant, with lower mean scores for emotions such as hate, disgust, and scared.

There are also some interesting patterns within the different emotions. For example, the emotions of love and adore have very similar mean scores, suggesting that these two emotions are often closely related. On the other hand, the emotions of hate and disgust have very different mean scores, suggesting that these two emotions are more distinct from each other.

Overall, the scatter plot provides some insights into the way that people experience emotions. It shows that positive emotions are more common than negative emotions, and that some emotions are more closely related to each other than others.

To summarise the insights in few points:

- The emotions with the highest mean scores are all positive: love, adore, cherish, and heart.
- The emotions with the lowest mean scores are all negative: hate, disgust, and scared.
- The emotions of love and adore have very similar mean scores, suggesting that these two emotions are often closely related.
- The emotions of hate and disgust have very different mean scores, suggesting that these two emotions are more distinct from each other.
- There is a general trend of increasing mean emotion scores as we move from negative to neutral to positive sentiment. This suggests that people tend to experience more intense emotions when they are feeling positive.

## Calculating summary statistics (mean, median, etc.) for Score by Type_x(content type)

### The MEANS Procedure

#### Analysis Variable : Score

| Type_x | N Obs | Mean | Median | Std Dev | Minimum | Maximum |
|---|---|---|---|---|---|---|
| GIF | 4465 | 39.2470325 | 35.0000000 | 26.0079708 | 0 | 75.0000000 |
| audio | 4276 | 40.2156221 | 45.0000000 | 26.1225864 | 0 | 75.0000000 |
| photo | 5006 | 40.0165801 | 45.0000000 | 26.1157786 | 0 | 75.0000000 |
| video | 4637 | 39.5518654 | 35.0000000 | 25.9061939 | 0 | 75.0000000 |

The table shows that the type of media with the highest mean score is audio (40.2156221), followed by photo (40.0165801), and video (39.5518654). The type of media with the lowest mean score is GIF (39.2470325).

The table also shows that the median and maximum scores for each type of media are generally similar to the mean scores. However, there are a few exceptions. For example, the median score for

audio is 45, while the maximum score is also 75. This suggests that most of the audio's in the dataset had a score of 45, and there were no outliers with significantly higher or lower scores.

**Calculating summary statistics (mean, median, etc.) for Score by Type_y(Emotion Lexicon)**

The MEANS Procedure

| | | Analysis Variable : Score | | | | |
|---|---|---|---|---|---|---|
| Type_y | N Obs | Mean | Median | Std Dev | Minimum | Maximum |
| adore | 1148 | 72.0000000 | 72.0000000 | 0 | 72.0000000 | 72.0000000 |
| cherish | 1119 | 70.0000000 | 70.0000000 | 0 | 70.0000000 | 70.0000000 |
| disgust | 1142 | 0 | 0 | 0 | 0 | 0 |
| dislike | 1109 | 10.0000000 | 10.0000000 | 0 | 10.0000000 | 10.0000000 |
| hate | 1153 | 5.0000000 | 5.0000000 | 0 | 5.0000000 | 5.0000000 |
| heart | 1225 | 60.0000000 | 60.0000000 | 0 | 60.0000000 | 60.0000000 |
| indiffe | 1167 | 20.0000000 | 20.0000000 | 0 | 20.0000000 | 20.0000000 |
| interes | 1169 | 30.0000000 | 30.0000000 | 0 | 30.0000000 | 30.0000000 |
| intrigu | 1091 | 45.0000000 | 45.0000000 | 0 | 45.0000000 | 45.0000000 |
| like | 1132 | 50.0000000 | 50.0000000 | 0 | 50.0000000 | 50.0000000 |
| love | 1165 | 65.0000000 | 65.0000000 | 0 | 65.0000000 | 65.0000000 |
| peeking | 1157 | 35.0000000 | 35.0000000 | 0 | 35.0000000 | 35.0000000 |
| scared | 1174 | 15.0000000 | 15.0000000 | 0 | 15.0000000 | 15.0000000 |
| super l | 1167 | 75.0000000 | 75.0000000 | 0 | 75.0000000 | 75.0000000 |
| want | 1154 | 70.0000000 | 70.0000000 | 0 | 70.0000000 | 70.0000000 |
| worried | 1112 | 12.0000000 | 12.0000000 | 0 | 12.0000000 | 12.0000000 |

- **Columns**: The table has three columns: Type_y, N Obs, and Analysis Variable : Score.
- **Rows**: Each row in the table represents a different emotion or feeling, such as "adore", "dislike", "hate", etc.
- **Data**: The table is analysing the mean, median, standard deviation, minimum, and maximum scores for each emotion or feeling.

This table provides a statistical analysis of the variable "Score" for different emotions or feelings. Type_y is the type of sentiment (adore, scared, peeking, want, hate, heart, dislike,.etc.), N Obs is the number of observations for each emotion, and Analysis Variable : Score is the calculated statistical measures for each emotion.

# Conclusion:

In this sentiment and emotion analysis using SAS Studio, we conducted a comprehensive exploration of text data, aiming to unveil sentiment and emotional patterns. We collected and prepared the data, performed sentiment and emotion analysis, and visualised the results. Sentiment analysis categorised the data into "positive," "negative," and "neutral" sentiments, providing an overview of the dataset's overall sentiment distribution. Expanding our analysis to include emotion analysis, we identified a spectrum of emotions like "love," "disgust," and "indifferent."

Visualisations, such as bar charts and violin plots, gave us understanding of sentiment and emotion distributions within different categories. Additionally, category-based analyses revealed specific emotions associated with sentiment categories.

This analysis equips us with valuable insights for applications in market research, customer feedback analysis, and social media sentiment tracking. The results highlighted the refined emotional landscape within the dataset and form a data-driven foundation for informed decision-making and content optimization.

# Reference:

Sentiment_Analysis.SAS:
https://drive.google.com/file/d/1Vy3V6rhnogrOqr2-V37TaS0B1ChonJjp/view?usp=sharing

Original Dataset:
https://docs.google.com/spreadsheets/d/1lsJQYhn-onbNNaNjnZWRcsJEA-5gYP06qyiCMQ3Vi40/edit?usp=sharing

Sentiment Analysis Literature:
How to extract domain-specific sentiment lexicons - SAS Users