

**UNIVERSITY OF  
WESTMINSTER** 

**7BUIS025W.2**

## **Web and Social Media Analytics**

**Assessment 02**

**May 2024**

**Module Leader:**

Philip Worrall

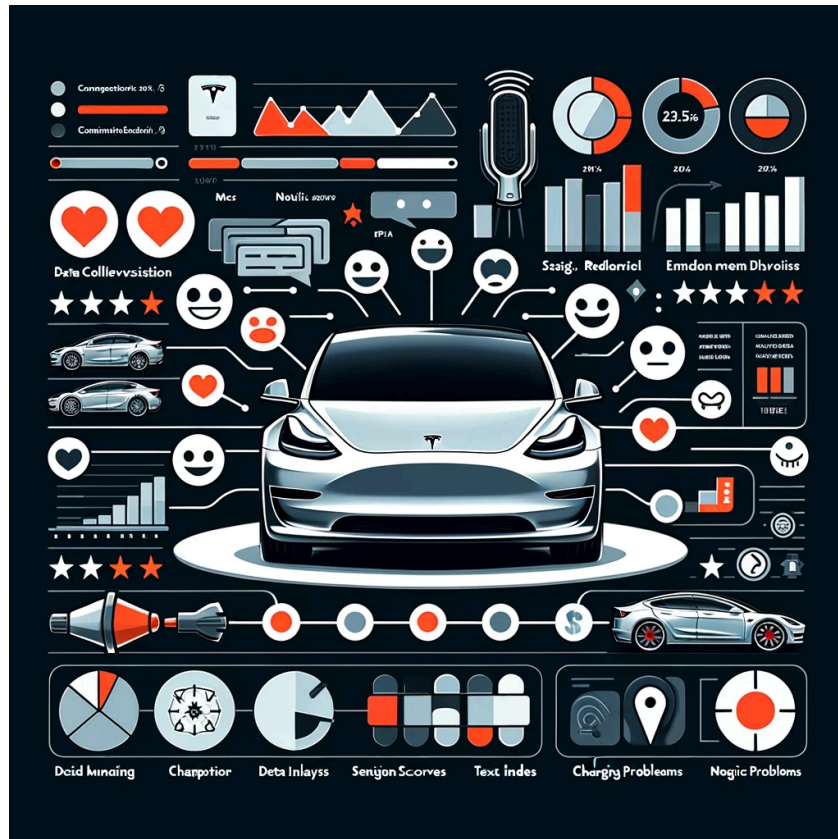
**Akram Naoufel Tabet**

W1979246

**Words count: 5179**

## Project Overview

In this project, we apply social media analytics to explore the experiences of electric vehicle (EV) owners. Using data from platforms like Reddit, YouTube, and others, we aim to uncover common themes, sentiments, and issues related to EV ownership. Our focus will be on the Tesla Model 3, a popular choice in the electric vehicle market. The analysis will involve data collection, preprocessing, exploratory data analysis, and text mining to draw meaningful insights.



## Introduction

Electric vehicles (EVs) have been increasingly adopted due to their environmental benefits and advancements in technology. Among these, the Tesla Model 3 stands out as one of the most popular fully electric vehicles. Launched in 2017, the Model 3 offers a balance of performance, range, and affordability, making it accessible to a broad audience. Industry trends indicate a growing shift towards sustainable transportation, supported by government incentives and advancements in battery technology. Current models available in the market include the Tesla Model 3 Standard Range Plus and the Tesla Model 3 Long Range, both of which cater to different consumer needs based on budget and range requirements.

## Data Collection and Preprocessing

### Keywords and Search Terms

To gather relevant social media data, I used keywords such as "Tesla Model 3," "electric vehicle tesla model 3 reddit," "EV tesla model 3 owners reviews," and "Tesla model 3 honest reviews."

### Data Collection Procedure

Data was collected from Reddit and YouTube. Using the ASYNCPRAW (Python Reddit API Wrapper) library, I retrieved comments from specific Reddit threads discussing the Tesla Model 3. For YouTube, I utilised the YouTube Data API to extract comments from relevant videos. The following steps outline the data collection process:

For the Reddit data collection, I used the `asynccpraw` library to interact with the Reddit API asynchronously. This approach allows us to efficiently fetch comments from multiple Reddit threads. The data collected spans from 2019 to 2024, with more data points in recent years, particularly in 2024.

### Code Explanation:

1. **API Setup:**
  - I initialised the `Reddit` instance with the provided API credentials (`API_ID` and `SECRET`).
2. **Fetching Comments:**
  - For each submission ID in `submission_ids`, I retrieved the submission and its comments.
  - Using `replace_more(limit=None)` ensures that all comments are fetched, including those beyond the initial limit.
  - For each comment, I extracted the comment ID, body, and the date of creation.
3. **Storing Data:**
  - Each comment's data, including the submission ID, comment ID, body, and date, was appended to a list and then converted into a DataFrame.

```
import asyncio
from datetime import datetime

SECRET = "rzHwncJbCMMBDONV3iSGFMStXLHmwQ"
API_ID = "zPDCzunJHwK4DzXTM1k4w"
submission_ids = ["1caeaog", "1agsjvfv", "1bd23yk", "14fnh1v", "up5qtz"]

async def fetch_comments(api_id, secret, submission_ids):
    reddit = asynccpraw.Reddit(
        client_id=api_id,
        client_secret=secret,
        user_agent="Comment Extraction"
    )

    comments_data = []

    for submission_id in submission_ids:
        try:
            submission = await reddit.submission(submission_id)
            comments = await submission.comments()
            await comments.replace_more(limit=None)
            all_comments = await comments.list()
            for comment in all_comments:
                try:
                    comment_date = datetime.fromtimestamp(comment.created_utc).strftime('%Y-%m-%d %H:%M:%S')
                    comments_data.append({
                        "submission_id": submission_id,
                        "comment_id": comment.id,
                        "comment_body": comment.body,
                        "comment_date": comment_date
                    })
                except AttributeError:
                    pass
            except Exception as e:
                print(f"Failed to process submission {submission_id}: {str(e)}")
                continue

    return pd.DataFrame(comments_data)

# fetch comments
df reddit = await fetch_comments(API_ID, SECRET, submission_ids)
```

Figure 01: Displays the code used for collecting data from Reddit

	submission_id	comment_id	comment_body	comment_date
0	1caeaog	i0rczzc	Putting aside, just for a moment, my feelings about Tesla's CEO, the Model 3 is a brilliant vehicle at getting folks from a-to-b.	2024-04-22 16:11:18
1	1caeaog	i0rpnmb	I like the new look. Signal stalk, speedometer hud and some build quality for the interior is what would entice me to buy one. \n\nThat throttle house video where they compared a new vs used tesla was eye opening	2024-04-22 17:23:12
2	1caeaog	i0s3zbo	Considering how he completely denegated the BMW i5 M50i. I'd say his review of the Model 3 was a compliment.	2024-04-22 18:44:36
3	1caeaog	i0rew7l	This is a really good deal for the money. If the tax credit ever comes back it'll be an even better deal	2024-04-22 16:22:07
4	1caeaog	i0bwnk3	It's crazy to me that you can't option a heads up display	2024-04-23 02:08:38
...	...	...	...	...
947	up5qzt	i924b6n	you know there are aftermarket suspension kits right ?	2022-05-18 10:51:18
948	up5qzt	i8n2mjl	Sounds like a business idea, a self contained side view mirror with blind spot indicator	2022-05-15 00:40:12
949	up5qzt	i8kkw63	It was an M50	2022-05-14 12:51:10
950	up5qzt	i92dddo	Good luck finding an after market that does both road and track use better than Porsche, BMW, or Audi. Nearly every single aftermarket either stiffens or raises the vehicle. I don't know of any brand making aftermarket adaptive suspension and claims a larger delta from stuff to comfy	2022-05-18 12:29:38

**Figure2:** Shows the resulting DataFrame with columns: **submission\_id**, **comment\_id**, **comment\_body**, and **comment\_date**.

### 1. YouTube Data Collection:

For YouTube, I used the YouTube Data API to extract comments from multiple videos on the Tesla channel. This approach ensured a broad collection of opinions and feedback from various videos, enhancing the diversity of the dataset.

#### Code Explanation:

##### 1. API Setup:

- I initialised the YouTube API client using the provided API key.

##### 2. Fetching Comments:

- For each video ID, I retrieved the top-level comments and their replies.
- I used pagination to handle videos with a large number of comments, fetching subsequent pages of comments using the **nextPageToken**.

##### 3. Storing Data:

- Each comment's data, including the comment text and date, was appended to a list and then converted into a DataFrame.

```
from googleapiclient.discovery import build
import pandas as pd
from datetime import datetime

api_key = 'AIzaSyAFx9pF2In8eFDPag_0l6B0G6OGfZCcBn8'

def video_comments(video_ids):
    youtube = build('youtube', 'v3', developerKey=api_key)
    all_comments = []

    for video_id in video_ids:
        try:
            video_response = youtube.commentThreads().list(
                part='snippet,replies',
                videoId=video_id,
                maxResults=100,
                textFormat='plainText'
            ).execute()

            while video_response:
                for item in video_response.get('items', []):
                    topLevelComment = item['snippet']['topLevelComment']['snippet']
                    comment_text = topLevelComment.get('textDisplay', '')
                    comment_date = topLevelComment.get('publishedAt', '')

                    if len(comment_text) > 10 and comment_date:
                        comment_date = datetime.strptime(comment_date, '%Y-%m-%dT%H:%M:%SZ')
                        all_comments.append({
                            "video_id": video_id,
                            "comment": comment_text,
                            "comment_date": comment_date
                        })

                # Process replies if they exist
                if item.get('replies'):
                    replies = item['replies']['comments']
                    for reply in replies:
                        reply_text = reply['snippet'].get('textDisplay', '')
                        reply_date = reply['snippet'].get('publishedAt', '')
```

**Figure 03:**Displays the code used for collecting data from YouTube.



```

dem = demoji.findall(text)

for item in dem.keys():

    text = text.replace(item, '')

return text

def is_english(text):
    try:
        return detect(text) == 'en'
    except:
        return False

df['is_english'] = df['comment'].apply(is_english)
df = df[df['is_english']].drop(columns='is_english')

def clean_text(text, remove_emojis):
    text = text.lower()
    if remove_emojis:
        text = remove_em(text)
    else:
        text = demojize(text)
    text = re.sub(r'https?://\S+|www\.\S+', '{link}', text)
    text = re.sub(r'<a\s+href="[^"]*">([<]*)</a>', '{link}', text)
    text = re.sub(r'@[w-.\_]+', ' ', text)
    text = re.sub(r'&\d+;', ' ', text)
    text = text.replace('\n', ' ')
    text = re.sub(r'^\w\s\''|_', ' ', text)
    text = re.sub(r'@\w+', ' ', text)
    text = re.sub(r'""', '', text)
    text = re.sub(r'isn\'t', "is not", text)
    text = re.sub(r"isn't", "is not", text)
    text = re.sub(r"it's", "it is", text)
    text = re.sub(r"it's", "it is", text)
    text = re.sub(r"wasn't", "was not", text)
    text = re.sub(r"wasn't", "was not", text)
    text = re.sub(r"didn't", "did not", text)
    text = re.sub(r"didn't", "did not", text)
    text = re.sub(r"\re", " are", text)

```

**Figure 05:** Shows the code used for cleaning process

```

from nltk.tag import pos_tag
from nltk.tokenize import word_tokenize
from nltk.corpus import wordnet
from nltk.stem import WordNetLemmatizer

nltk.download('wordnet')
nltk.download('averaged_perceptron_tagger')
nltk.download('omw-1.4')
nltk.download('punkt')
nltk.download('stopwords')

# Initialize the lemmatizer
lemmatizer = WordNetLemmatizer()

def get_wordnet_pos(word):
    """Map POS tag to first character lemmatize() accepts"""
    tag = pos_tag([word])[0][1][0].upper()
    tag_dict = {
        'J': wordnet.ADJ,
        'N': wordnet.NOUN,
        'V': wordnet.VERB,
        'R': wordnet.ADV
    }
    return tag_dict.get(tag, wordnet.NOUN)

def lemmatize_text(text):
    tokens = word_tokenize(text)
    lemmatized_tokens = [lemmatizer.lemmatize(token, get_wordnet_pos(token)) for token in tokens]
    return ' '.join(lemmatized_tokens)

```

**Figure 06:** Shows the code used for removing stop words and lemma

## Code Explanation

### 1. Import Libraries:

- I import necessary libraries for text processing, including `nltk`, `re`, and `langdetect`.

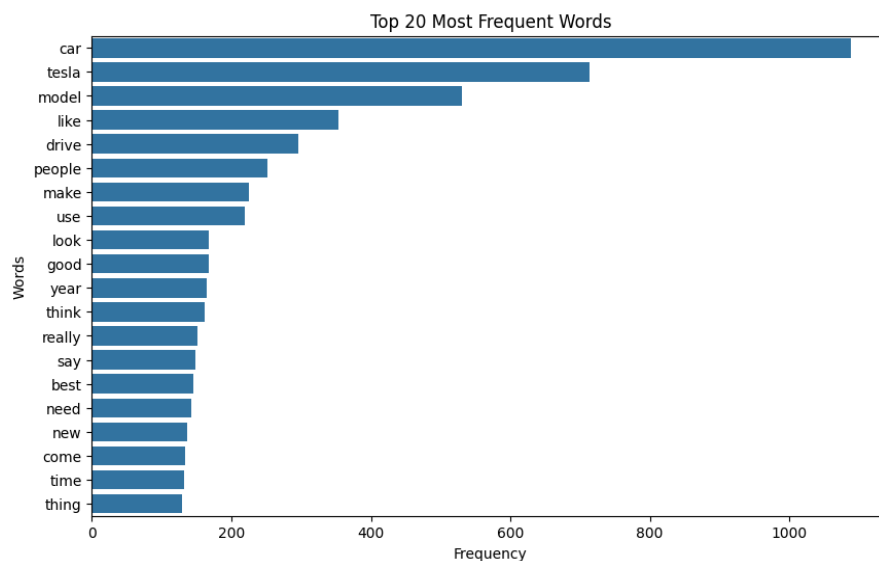
## Introduction to EDA

EDA helps in summarising the main characteristics of the data, often visualising them to identify patterns, anomalies, and insights that are not immediately obvious. In this analysis, I focused on:

1. **Word Frequency Analysis:** Identifying the most frequently occurring words in the comments to understand the primary topics of discussion.
2. **Sentiment Analysis Over Time:** Analysing the sentiment trends over the years to understand how user opinions have evolved.
3. **N-gram Analysis:** Identifying common bigrams and trigrams to capture more contextual insights from the comments.

## Most Frequent Words

The first step in our EDA was to analyse the most frequently occurring words in the dataset. This analysis helps in identifying the primary topics and concerns discussed by the users.



**Figure 08: Top 20 Most Frequent Words**

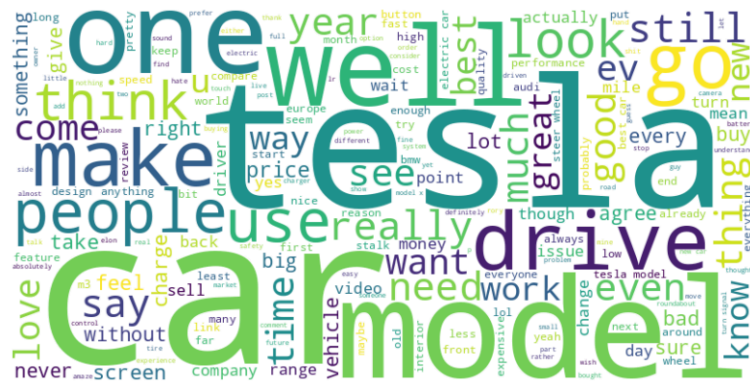
The bar chart above displays the top 20 most frequently used words in the comments.

- **Key Findings:**
  - The word "car" is the most frequent, indicating that the discussions are generally centred around the vehicle itself.
  - "Tesla" and "model" are also among the top words, highlighting the focus on the Tesla brand and the Model 3 specifically.
  - Words like "like", "drive", "people", "make", and "use" suggest common topics related to user experiences, preferences, and functionality.
  - Positive words such as "good", "best", and "better" indicate favourable opinions, whereas words like "need" and "time" may suggest areas for improvement or common concerns.

## Word Cloud Visualization

The word cloud visualisation provides an intuitive representation of the most common terms found in the comments about the Tesla Model 3. Larger words indicate higher frequency, giving us a quick insight into the primary topics and sentiments expressed by users. In this word cloud, terms like "Tesla," "car," "model," and "drive" stand out prominently, suggesting that discussions are heavily centred around the vehicle itself and its driving experience. Other notable terms such as "one," "make," "people," and "time" hint at more detailed aspects of user experiences and opinions. This visualisation helps us identify key themes and topics at a glance, setting the stage for more detailed text analysis.



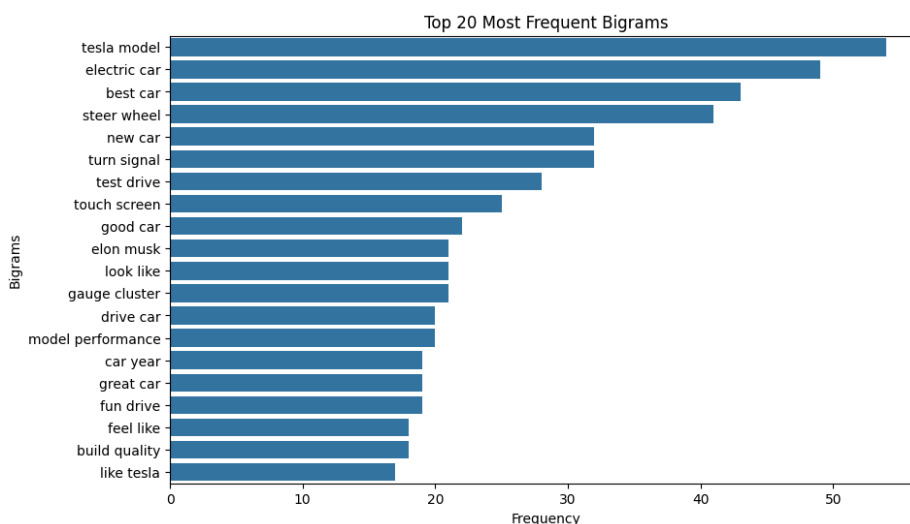


## Bigrams Analysis

The bigrams analysis provides deeper insights into the common phrases used by Tesla Model 3 owners, highlighting specific areas of interest and concern. The top bigram, "Tesla Model," signifies the frequent mention of the specific vehicle model, emphasising its centrality in discussions.

The bigram "electric car" appears frequently, reflecting the general discussions about electric vehicles, their benefits, and challenges. This indicates that users often compare Tesla Model 3 with other electric cars or discuss the electric vehicle market in general. Bigrams such as "best car" and "new car" suggest that many users consider the Tesla Model 3 to be among the best or are discussing its new features and updates. Mentions of specific features like "steer wheel," "test drive," and "touch screen" highlight the importance of these components in user discussions. This could indicate areas where users have strong opinions, either positive or negative, about the vehicle's design and functionality.

This analysis helps us understand the prevalent themes and issues in more context compared to single-word frequency analysis.



## Trigrams Analysis

The trigrams analysis goes a step further by capturing three-word phrases, providing more context to the discussions. Trigrams like "turn signal stalk" and "touch screen drive" suggest specific features that are frequently talked about, possibly indicating areas where users

have had notable experiences or issues. The appearance of "anonymized redact link" and "mass delete anonymized" indicates concerns or discussions around data privacy and content moderation, which might be relevant in forums where personal data security is a topic of concern."Tesla come India" highlights geographic interest, suggesting that there is significant discussion about Tesla's market expansion into India. This could reflect both excitement and concerns about the vehicle's availability and performance in new markets.Trigrams such as "best sell car" and "brand new car" reinforce the notion that users often discuss the Tesla Model 3 in the context of it being a top-selling or newly updated vehicle, further underscoring its popularity and relevance in the electric vehicle market.

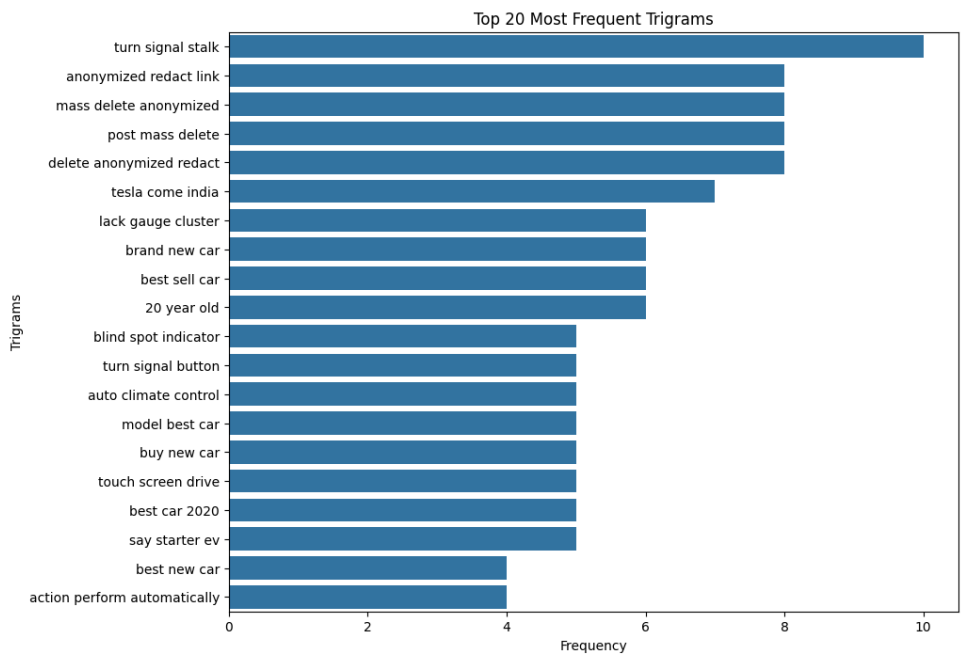
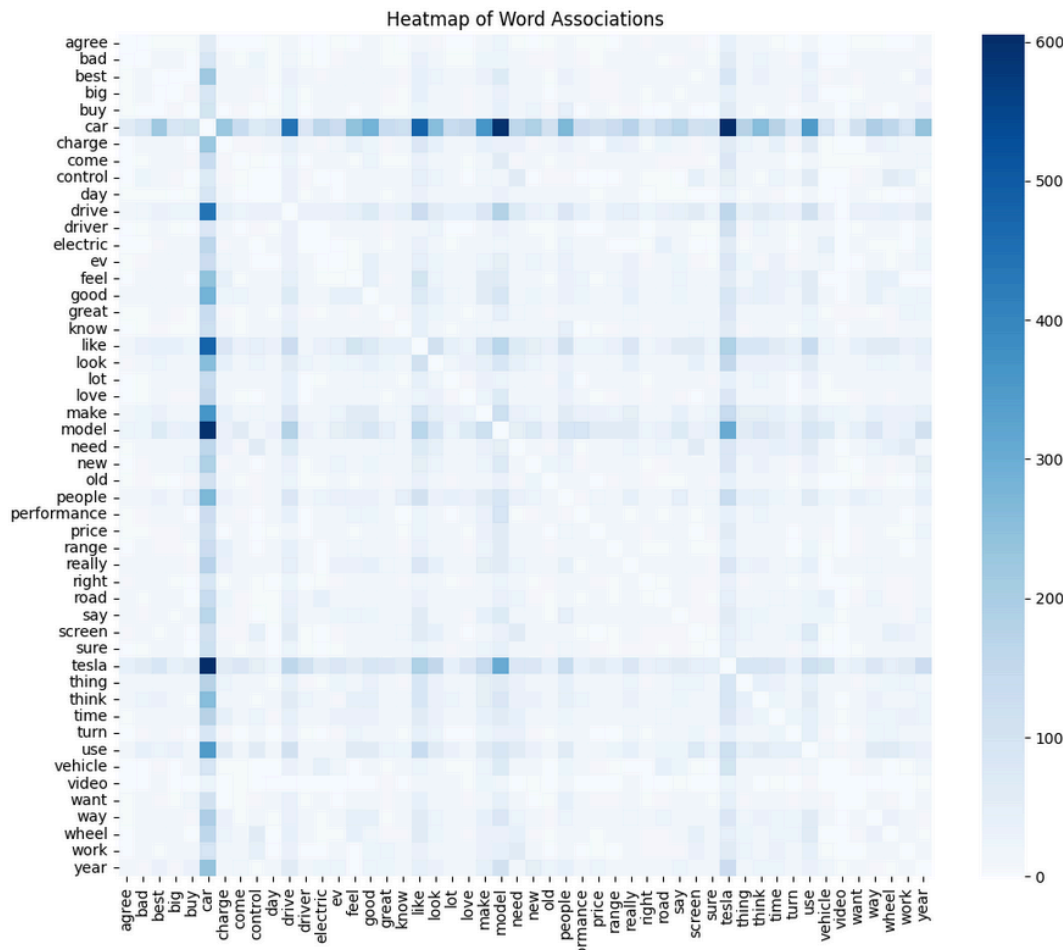


Figure 10: Top 20 most frequent trigrams

Heatmap of Word Associations



**Figure 11: Heatmap of Word Associations**

The heatmap visualises the associations between the most frequently occurring words in the comments. Darker shades indicate stronger associations between pairs of words. For example, "car" and "Tesla" are strongly associated, reflecting discussions specifically about Tesla cars. Similarly, words like "charge" and "time" show a strong association, likely highlighting user concerns about charging duration and infrastructure. The heatmap also reveals associations like "drive" and "feel," suggesting that users frequently discuss their driving experiences and feelings about the car's performance.

### Comment Volume Over Time

The line chart represents the volume of comments over time, highlighting peaks and troughs in user engagement. Significant spikes in the number of comments, such as those in early 2019 and late 2023, may correspond to major events or announcements related to the Tesla Model 3, such as new model releases or updates. Periods with lower comment volumes might indicate times of less activity or fewer notable events. This analysis of comment volume over time provides insights into the periods of heightened interest and engagement within the Tesla community, helping to identify when the vehicle was most discussed and why.

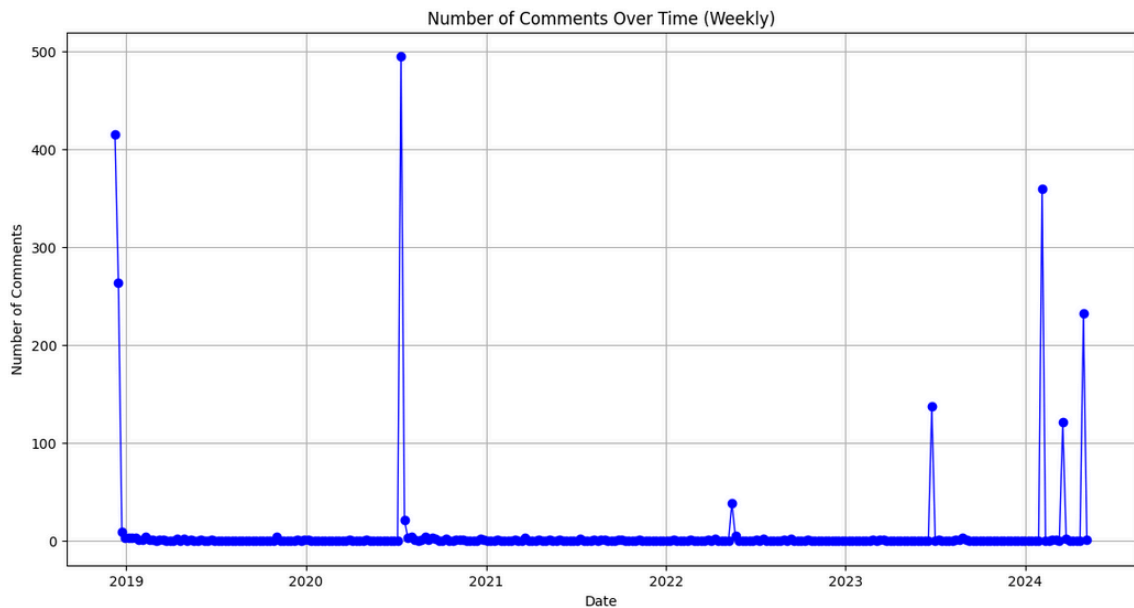


Figure 15: Number of Comments Over Time

## Text Mining Analysis

In the text mining section of this project, I applied various techniques to extract meaningful insights from the collected comments about the Tesla Model 3. My approach included sentiment analysis to gauge the overall sentiment of the comments, LDA topic modelling to identify prevalent discussion themes, and sentence transformers to detect issues and negative comments based on semantic similarity. This comprehensive analysis helps me to understand the sentiments, identify key topics, and pinpoint specific issues experienced by Tesla Model 3 owners.

### Sentiment Analysis

Sentiment analysis is a crucial part of understanding user opinions and emotions expressed in the comments. I used the VADER (Valence Aware Dictionary and sEntiment Reasoner) sentiment analysis tool from the NLTK library, which is particularly effective for social media text. VADER provides a compound sentiment score for each comment, which I used to classify the comments into negative, neutral, and positive categories.

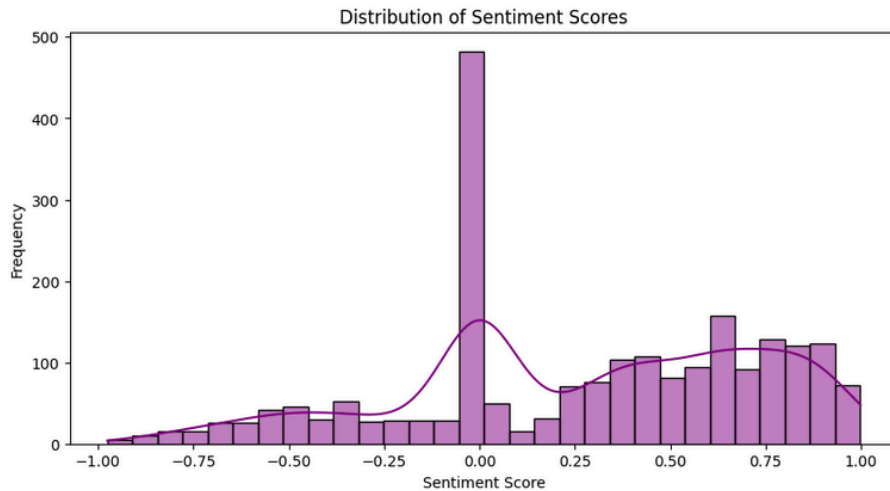
```
from nltk.sentiment.vader import SentimentIntensityAnalyzer
nltk.download('vader_lexicon')

# Initialize sentiment analyzer
sia = SentimentIntensityAnalyzer()

# Apply sentiment analysis
df['sentiment'] = df['cleaned_comment_with_emotes'].apply(lambda x: sia.polarity_scores(x)['compound'])
```

cleaned_comment	processed_comment	processed_comment_with_emotes	sentiment
for 95 of car buyers fun to drive is not in th...	95 car buyer fun drive top 5 priority fuel eco...	95 car buyer fun drive top 5 priority fuel eco...	0.8271
god i wish that were me	god wish	god wish	0.5859
the packaging is poor interior space is crampe...	packaging poor interior space cramped storage ...	packaging poor interior space cramped storage ...	-0.7971
shame i do not have 30 000 lying around	shame 30 000 lie around	shame 30 000 lie around	-0.4767
if they are minor parts they are probably not ...	minor part probably put car service	minor part probably put car service	0.0000

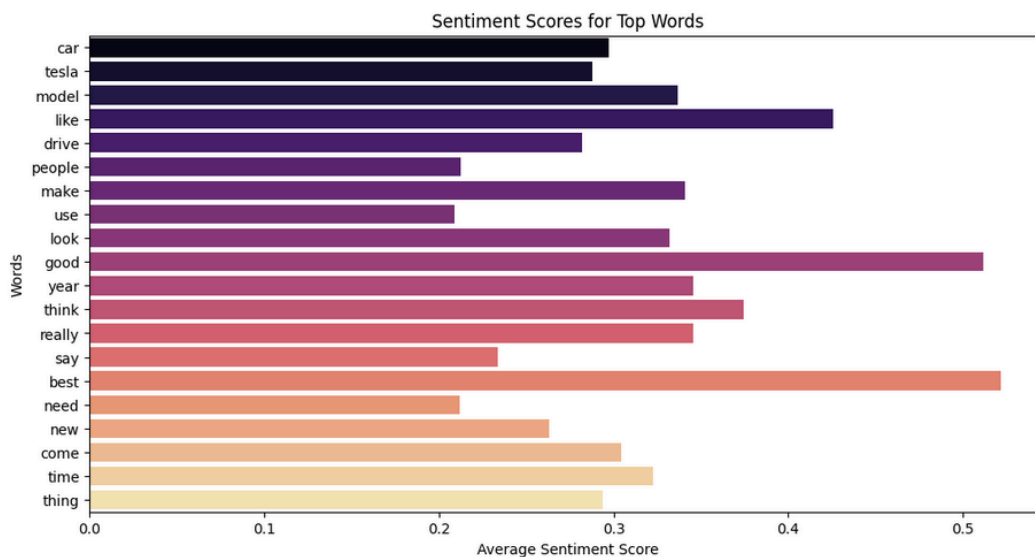
The result of the sentiment analysis includes a new column in the DataFrame called **sentiment**, which categorises each comment as either negative, neutral, or positive based on the compound sentiment score. Here is a sample of the resulting DataFrame.



**Figure 16:** Distribution of Sentiment Scores

This histogram displays the distribution of sentiment scores for the comments. The sentiment scores, calculated using the VADER sentiment analysis tool, range from -1 (most negative) to +1 (most positive). The distribution is somewhat skewed towards the positive end, indicating that a majority of comments express positive sentiments about the Tesla Model 3. The high frequency of comments around a sentiment score of 0 suggests that many comments are neutral, offering balanced or factual statements without strong emotional tones. The presence of negative sentiment scores highlights areas where users have expressed dissatisfaction or concerns.

- **Positive Sentiment:** The dominance of positive sentiment scores suggests that Tesla Model 3 owners generally have favourable opinions about their vehicles.
- **Neutral Comments:** A significant number of neutral comments might indicate factual discussions or balanced reviews, possibly focusing on technical aspects or performance details.
- **Negative Sentiment:** The negative sentiment scores highlight specific issues or areas of dissatisfaction that could be explored further to identify common pain points among owners.



**Figure 17:** Sentiment Scores for Top Words

This bar chart presents the average sentiment scores for the top words mentioned in the comments. Words like "car," "Tesla," and "model" are associated with positive sentiment scores, reinforcing the overall positive perception of the Tesla Model 3. Words such as "drive" and "people" also have positive sentiments, suggesting that driving experiences and community interactions are well-regarded. The term "need," although not strongly negative, indicates areas where users feel improvements or additional features are necessary.

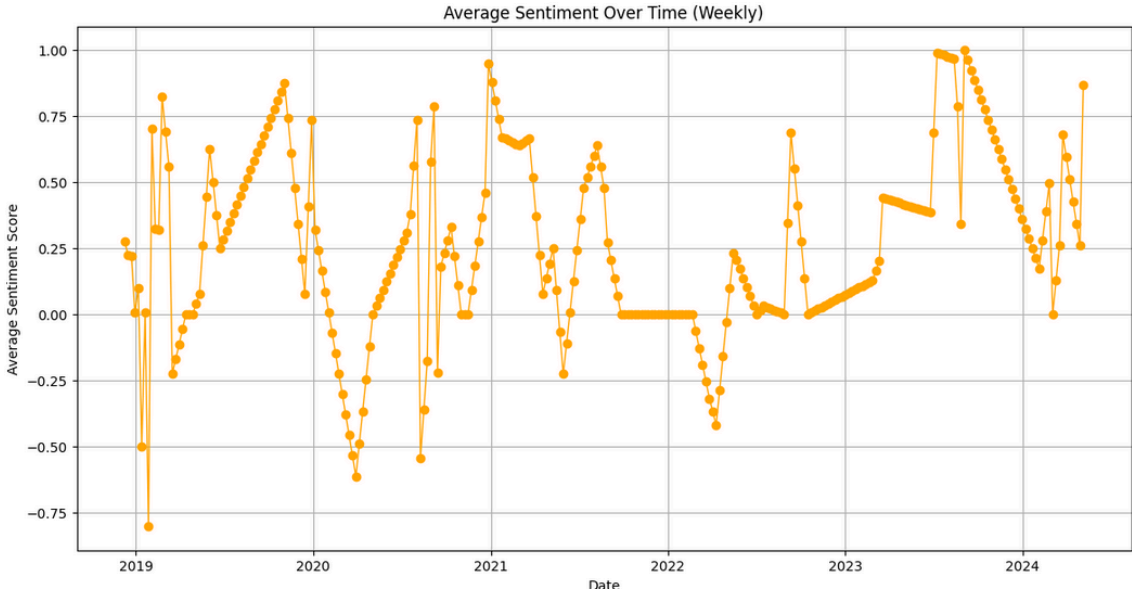


Figure 18: Average Sentiment Over Time

The line chart shows the average sentiment score of comments over time, from 2019 to 2024. This analysis reveals how user sentiment has fluctuated over the years. Notable peaks in sentiment, such as in mid-2020 and late 2023, suggest periods of particularly positive user experiences, which could coincide with new model releases, significant software updates, or positive media coverage. The troughs, particularly in early 2021, might correspond to issues like product recalls or negative press.

Summary Statistics

In addition to the sentiment analysis, I computed summary statistics to provide an overview of the data. This includes the total number of comments, the average sentiment score, and identifying the most positive and most negative comments.

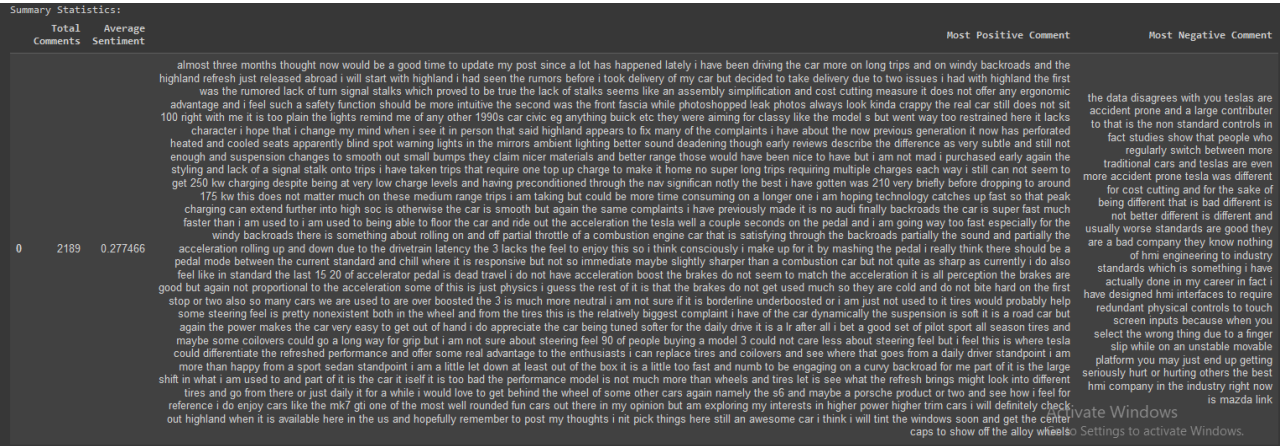


Figure 19: represent the extremes of user sentiment

LDA Topic Modeling

LDA (Latent Dirichlet Allocation) is a topic modelling technique used to uncover hidden themes or topics within a collection of documents. In our analysis, I applied LDA to identify the main discussion themes in the comments about the Tesla Model 3. Due to space constraints in

this report, I will focus on the LDA modelling of bigrams, which provides richer context by considering pairs of words together. For a more comprehensive analysis, including unigrams and trigrams, please refer to the associated Jupyter notebook.

### Code for LDA Topic Modeling

I used the `gensim` library to perform LDA topic modelling. The process involves several steps:

1. **Preprocessing:** This includes tokenization, removing stop words, and creating bigrams.
2. **Dictionary and Corpus Creation:** I created a dictionary from the processed text and converted the text into a bag-of-words corpus.
3. **LDA Model Training:** Using the `gensim.models.LdaModel`, I train the LDA model on our corpus.

Here is the code used for LDA topic modelling on bigrams:

```
import gensim
from gensim import corpora

# Create a list of bigrams for each document
bigrams_list = bigram_vectorizer.get_feature_names_out()

# Convert the bigrams matrix to a list of lists (required format for LDA)
bigrams_docs = bigram_matrix.toarray()
bigrams_docs_list = [[bigrams_list[i] for i in range(len(bigrams_list)) if doc[i] > 0] for doc in bigrams_docs]

# Create a dictionary representation of the documents
dictionary = corpora.Dictionary(bigrams_docs_list)

# Filter out extremes to limit the number of features
dictionary.filter_extremes(no_below=1, no_above=0.5)

# Create the Bag-of-Words model
corpus = [dictionary.doc2bow(doc) for doc in bigrams_docs_list]

# Apply LDA
lda_model = gensim.models.LdaModel(corpus, num_topics=5, id2word=dictionary, passes=10)

topics = lda_model.print_topics(num_words=5)
for topic in topics:
    print(topic)
```

**Figure 20:** Shows the code for LDA Topic Modeling

**Topic 2 Analysis:** The bigrams identified in Topic 2 reveal several key themes discussed by Tesla Model 3 owners. The bigram "good deal" suggests that many users perceive the car as offering good value for money, which is a positive indicator of customer satisfaction regarding pricing. On the other hand, the presence of "remove stalk" highlights a specific design feature that some users might find problematic or unnecessary, indicating an area where Tesla could consider making adjustments based on user feedback. The term "sh\*\* car" represents a starkly negative sentiment, suggesting that there are users who have had particularly bad experiences with their vehicles. This kind of feedback, although harsh, is crucial for identifying severe issues that need immediate attention. Additionally, "Tesla come India" reflects a strong interest in Tesla's market expansion, suggesting a demand for the vehicle in new geographic regions. Other bigrams such as "lack gauge cluster" and "auto climate control" indicate specific features that users frequently discuss, whether positively or negatively. The mention of "best sell car" points to a recognition of the Tesla Model 3's success in the market, while "blind spot indicator" and "touch screen drive" highlight important safety and usability features that are of significant interest to the community.

- **Value Perception:** The bigram "good deal" suggests that many users see the Tesla Model 3 as offering good value, which is a positive sign for the company's pricing strategy.
- **Design and Features:** Terms like "remove stalk" and "lack gauge cluster" indicate areas where users have strong opinions about the design and features of the car, highlighting opportunities for Tesla to refine these aspects.
- **Negative Feedback:** The appearance of "sh\*\* car" underscores the importance of addressing negative experiences to improve overall customer satisfaction.
- **Market Expansion:** "Tesla come India" reflects a strong interest in expanding Tesla's market presence to new regions, suggesting a potential area for growth.



- **Safety and Usability:** Features like "blind spot indicator" and "touch screen drive" are critical to users, emphasising the importance of these aspects in user satisfaction.

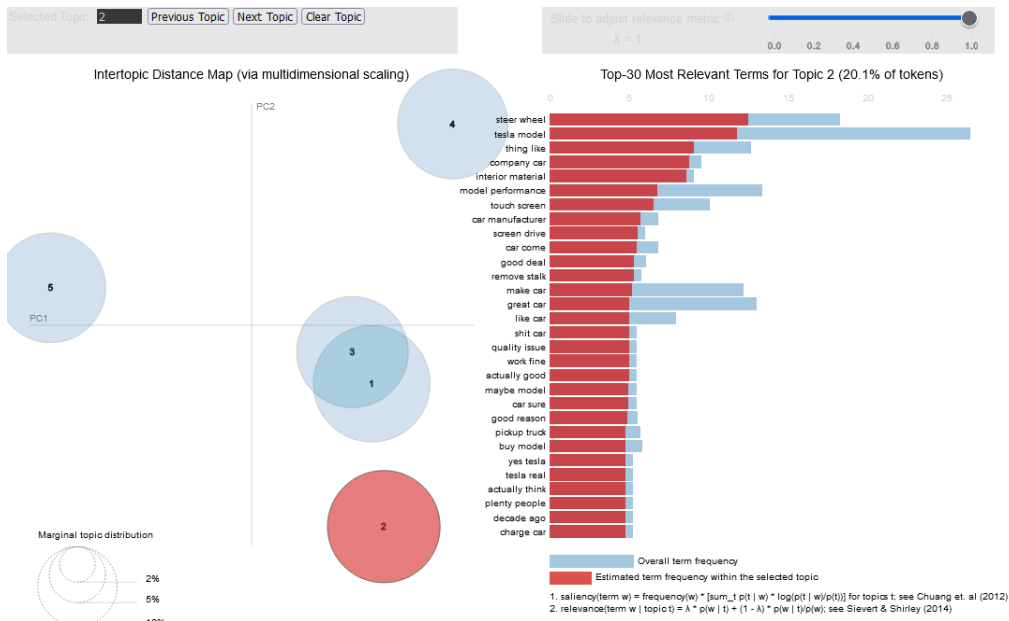


Figure 21: Shows LDA topic 2 with most relevant terms

## Identifying Themes and Issues in Negative Comments

In this section, we focus on identifying themes and issues related to the experiences of electric vehicle (EV) owners by analysing negative comments. By examining the top bigrams and trigrams in these comments and applying LDA topic modelling and sentence transformers based on their semantics, we can uncover specific areas of concern. The following figures present the results of our analysis.

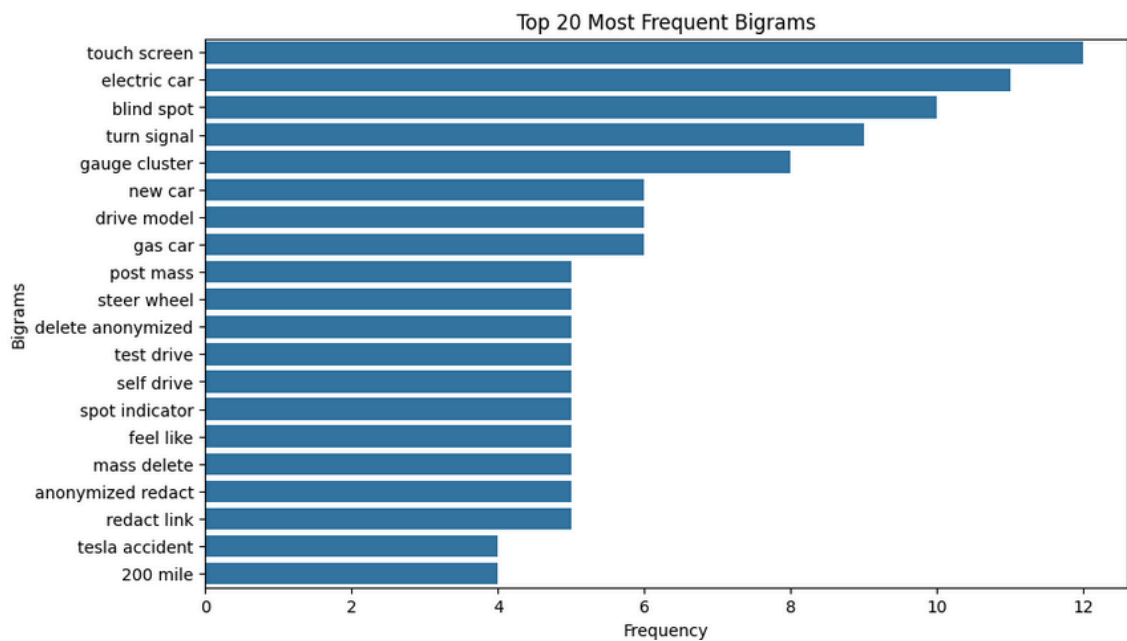
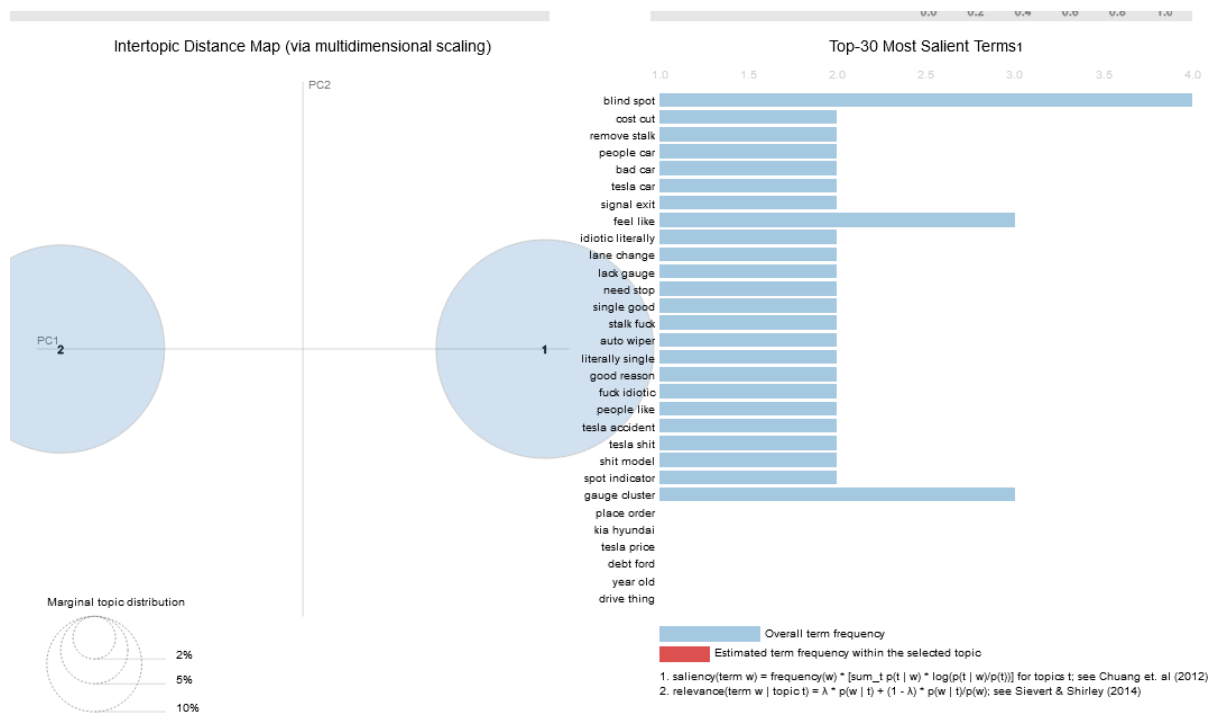


Figure 22: Top 20 Most Frequent Bigrams in Negative Comments



This bar chart shows the top 20 most frequent bigrams in negative comments. The most frequent bigrams include "post mass," "delete anonymized," and "anonymized redact," which suggest issues related to content moderation or data privacy. Other significant bigrams such as "blind spot," "turn signal," and "gauge cluster" indicate specific vehicle features that users are dissatisfied with. The bigram "remove stalk" suggests that many users find the stalk feature problematic. Additionally, "feel like" and "sh\*\* car" express strong negative sentiments about the vehicle's performance or overall experience.

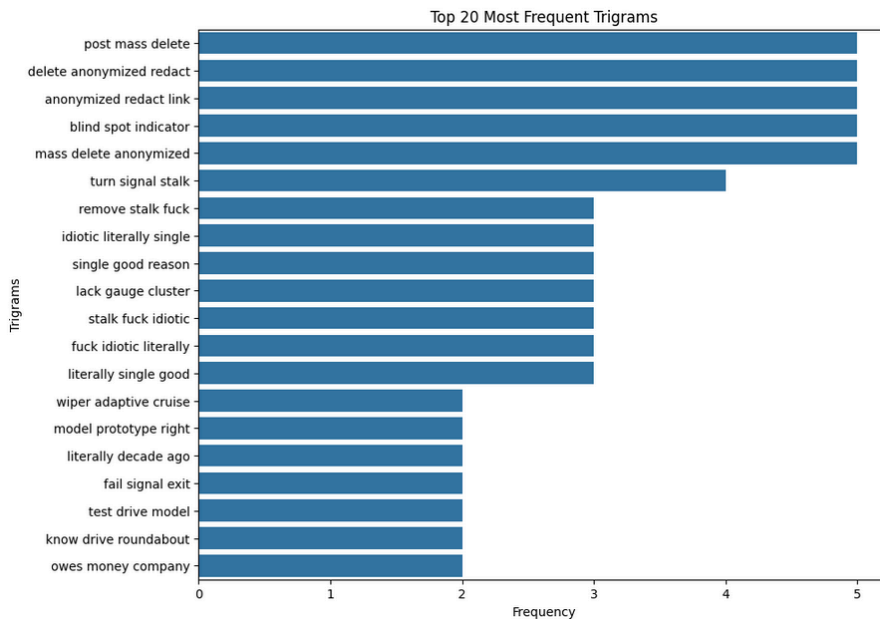
- **Content Moderation/Privacy Issues:** Frequent mentions of "post mass," "delete anonymized," and "anonymized redact" indicate concerns about content moderation and data privacy.
- **Vehicle Features:** Issues with features like "blind spot," "turn signal," and "gauge cluster" suggest areas where the vehicle's design or functionality could be improved.
- **Overall Dissatisfaction:** Bigrams like "feel like" and "sh\*\* car" highlight general dissatisfaction with the vehicle, pointing to a need for addressing broader performance or quality issues



**Figure 23: LDA Topic Modeling of Bigrams in Negative Comments**

The LDA topic modelling visualisation shows the most relevant bigrams for the identified topics. For instance, the bigrams "good deal," "remove stalk," and "sh\*\* car" highlight specific areas of user feedback. "Good deal" might indicate a paradox where users feel the car is priced well but still have complaints about specific features or overall performance. "Remove stalk" and "sh\*\* car" are direct indicators of dissatisfaction with certain design elements and the overall quality of the vehicle.

- **Pricing Paradox:** The mention of "good deal" alongside negative bigrams suggests that while users may find the car affordable, there are still significant concerns that need to be addressed.
- **Design and Quality Issues:** Bigrams like "remove stalk" and "sh\*\* car" clearly point to dissatisfaction with the vehicle's design and overall quality, indicating areas for potential improvement



**Figure 24:** Top 20 Most Frequent Trigrams in Negative Comments

This bar chart displays the top 20 most frequent trigrams in negative comments. Trigrams like "post mass delete," "delete anonymized redact," and "anonymized redact link" reinforce the content moderation and privacy concerns highlighted by the bigram analysis. Additionally, trigrams such as "turn signal stalk" and "remove stalk fu\*\*" indicate strong negative reactions to specific design features. Other trigrams like "idiotic literally single" and "fu\*\* idiotic literally" suggest a high level of frustration among some users.

- **Content Moderation/Privacy Issues:** The repetition of privacy-related trigrams confirms that this is a significant area of concern.
- **Design Feature Dissatisfaction:** Trigrams such as "turn signal stalk" and "remove stalk fu\*\*" highlight specific design features that users find particularly problematic.
- **High Frustration Levels:** Trigrams like "idiotic literally single" and "fu\*\* idiotic literally" show that some users express their dissatisfaction in very strong terms, indicating high levels of frustration that need to be addressed.

## Approach for Identifying Issues Using Sentence Transformers

In my analysis, I utilised Sentence Transformers to identify specific issues in the comments related to the Tesla Model 3. The aim was to leverage the power of semantic similarity to accurately categorise comments into predefined issues. This approach ensures that I capture the nuanced meanings and contexts of the comments, which is often lost with simple keyword matching.

### Why Sentence Transformers?

Sentence Transformers, especially models like 'all-MiniLM-L6-v2', are designed to generate dense vector representations (embeddings) for sentences, preserving their semantic meaning. Unlike traditional methods that rely on surface-level text similarities, Sentence Transformers capture the contextual meaning, making them ideal for our task. This allows us to:

1. **Capture Nuances:** Understand subtle differences and similarities in meaning.
2. **Handle Synonyms:** Recognize different words/phrases that convey the same meaning.
3. **Contextual Understanding:** Take into account the surrounding context to better understand the sentiment and topic of a comment.

### Identifying Issues: Step-by-Step Explanation

1. **Define Issues:** I start by defining a set of potential issues that might be mentioned in the comments. Each issue is described with a detailed description that captures various aspects of that issue.
2. **Generate Embeddings:** Using the Sentence Transformer model, I generate embeddings for both the issue descriptions and the comments. These embeddings are dense vector representations that capture the semantic meaning of the text.

3. **Compute Similarity:** I compute the cosine similarity between the comment embeddings and the issue embeddings. Cosine similarity measures the angle between two vectors, which in this context indicates how similar the meanings of the two texts are.
4. **Assign Issues:** For each comment, we find the issue with the highest similarity score and assign that issue to the comment. This way, each comment is categorised into one of the predefined issues based on its semantic similarity.

Here is the code that implements the above steps:

```
from sentence_transformers import SentenceTransformer, util

model = SentenceTransformer('all-MiniLM-L6-v2')

# Combine the issue descriptions into a list
issue_descriptions = list(issues.values())
issue_names = list(issues.keys())

# Generate embeddings for the issue descriptions
issue_embeddings = model.encode(issue_descriptions)

# Generate embeddings for the negative comments
negative_comments_text = negative_comments['cleaned_comment_with_emotes'].tolist()
comment_embeddings = model.encode(negative_comments_text)

# Compute the cosine similarity
similarity_matrix = util.pytorch_cos_sim(comment_embeddings, issue_embeddings)

# Find the index of the most similar issue for each comment
issue_indices = np.argmax(similarity_matrix, axis=1)

# Map the indices to issue names
assigned_issues = [issue_names[index] for index in issue_indices]

# Add the assigned issues to the DataFrame
negative_comments['assigned_issue'] = assigned_issues
```

**Figure 25:** Code implementation for the sentence transformer approach

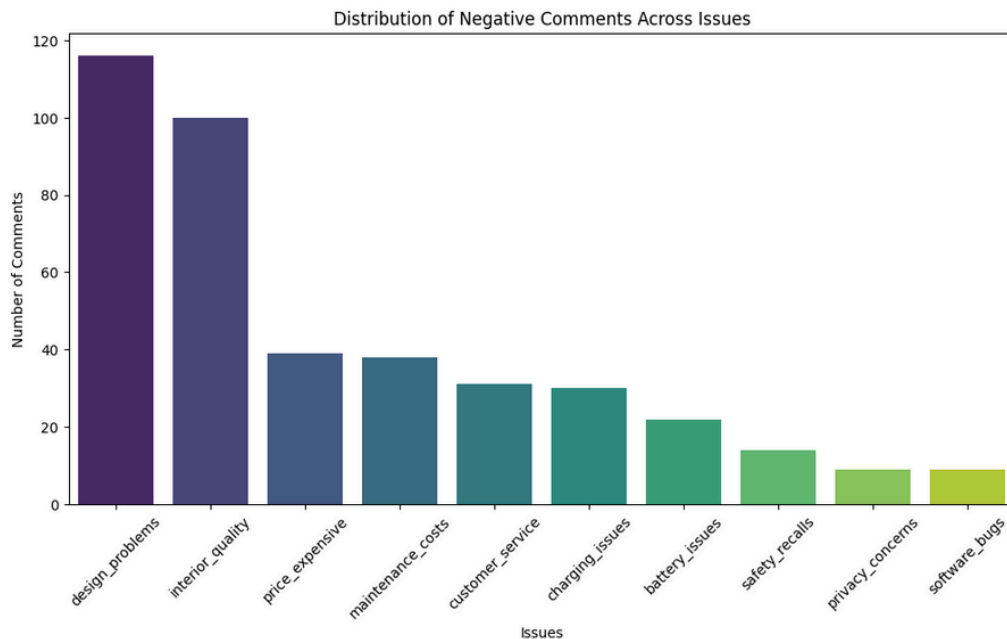
	cleaned_comment_with_emotes	assigned_issue
2	the packaging is poor interior space is cramped and storage is insufficient for the size and despite being a hatch it is less overall spacious than a model 3 that platform has been out since 2018 and it is overdue for a major revamp or replacement to stay competitive solved by the new ev dedicated platforms coming out you still do not have a working app with more than 4 functions cue the anecdotal b b but mine works the hardware was eol before it even launched and you haven't had a noteworthy feature addition since carplay the center infotainment screen has plasma tv sized bezels low resolution and can not even display two full features at the same time it still relies on a remote the size and style of an 80 is tv remote since the app does not work and finally the interior quality is nowhere near as good as you hype it up to be considering it about on par with the average hyundai kia the only reason why quality is always the first thing you polestar clones spurt out when it comes to your car is because it is the only attribute it really has it is mediocre otherwise as a 40k ev it is good as a 25k used ev here in the states it is fantastic as a new buy it is a no buy the ardent defenders and upvoters of this car on this sub has to be because they are so deeply upside down on their cars and or they are early stockholders way down it is very evident the 4 needs to replace this car entirely now cue the real fanboys worse than tesla by far in every sense of the word downvoting this	design_problems
4	shame i do not have 30 000 lying around	price_expensive
6	model 3 is a benchmark ev elon hate virus does not change that	battery_issues
8	never thought i would ever buy a new car then made the mistake of test driving a model s a few years later i found myself standing in line hours before the king of prussia mall opened to reserve a car i hadn't even seen yet	interior_quality
9	you can handwave all of this stuff all you want but the reality is the poorly performing driver assists safety features and the unintuitive dangerous cost cutting for basic io functions are all contributing factors as to why tesla is crash into things much more than any other brand by a large margin link link cars that crash often and are expensive to repair are always going to have very high insurance premiums and unfortunately tesla does not make this sort of safety or serviceability a priority in their designs	safety_recalls
...	...	...
2501	supercharged v8 for life screw nature	battery_issues
2504	the design looks outdated and lacks the modern touch	design_problems
2505	charging the car takes forever and it is a huge inconvenience	charging_issues
2506	the software is full of bugs and crashes frequently	software_bugs
2508	maintenance costs are ridiculously high	maintenance_costs

The table shows that each comment is assigned to an issue such as "design\_problems," "charging\_issues," "software\_bugs," and "maintenance\_costs." This categorisation helps in understanding the main areas of concern among Tesla Model 3 owners.

**Note:** To clarify, I chose to employ cleaned comments with emotes(which have stop words and not lemmatized) rather than processed comments(lemmatized,non-stop words) for sentiment analysis and semantic similarity due to their higher accuracy. However, for other

tasks such as LDA and identifying the most frequent words, I utilised the processed comments column.

### Distribution of Negative Comments Across Issues

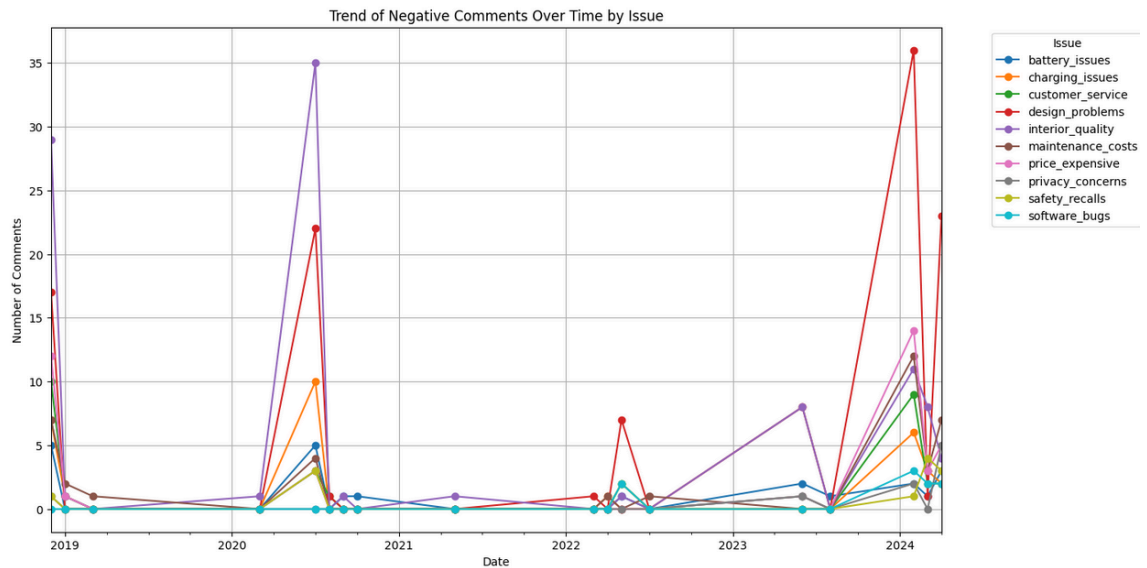


**Figure 27:** Shows the distribution of negative comments across issues

The bar chart above illustrates the distribution of negative comments across different issues. The most frequent issues are:

1. **Design and Interior Quality:** The majority of negative feedback is directed towards the design and interior quality of the Tesla Model 3, suggesting that these are critical areas for improvement.
2. **Cost-Related Issues:** High pricing and maintenance costs are significant concerns, indicating that customers feel the cost of ownership is too high.
3. **Customer Service:** Issues with customer service point to the need for Tesla to enhance its support and responsiveness to improve customer satisfaction.

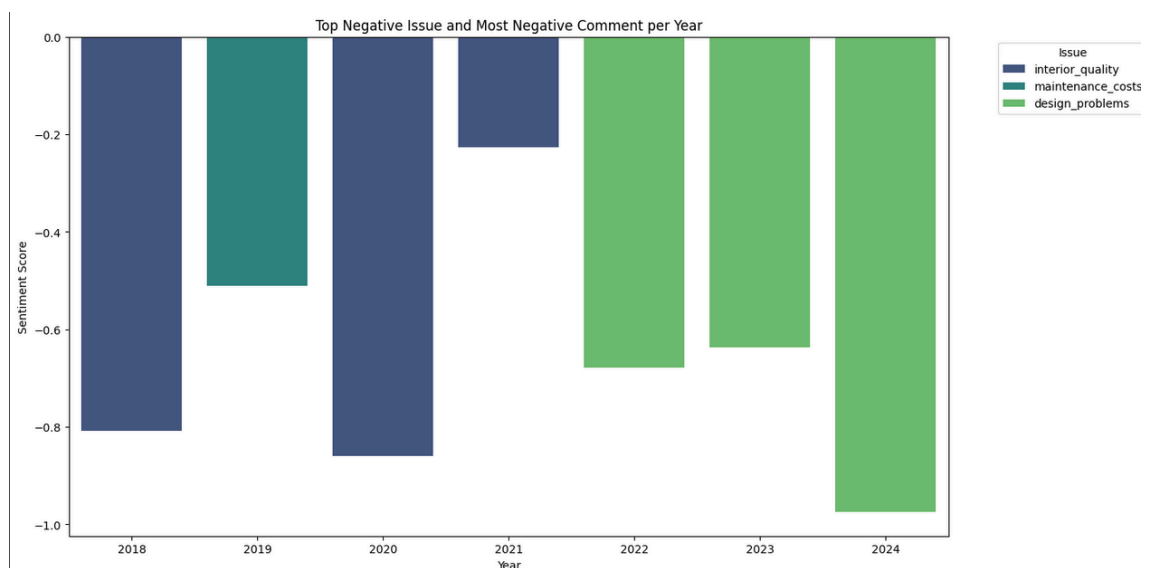
By addressing these key issues, Tesla can focus on improving the areas that matter most to their customers, potentially leading to increased satisfaction and loyalty.



**Figure 28: Trend of Negative Comments Over Time by Issue**

This chart visualises the trend of negative comments over time, categorised by different issues such as design problems, interior quality, maintenance costs, and others.

1. **Design Problems:** This issue shows significant spikes in certain periods, particularly around 2020 and 2024. This suggests that users were particularly vocal about design flaws during these times, potentially due to model updates or feature changes.
2. **Interior Quality:** Comments on interior quality also see peaks, indicating recurring concerns about the materials and build quality of the Tesla Model 3's interior.
3. **Maintenance Costs:** The issue of high maintenance costs appears consistently, reflecting ongoing concerns about the cost of ownership.
4. **Charging Issues:** There are notable spikes in comments about charging issues, especially towards the end of the timeline, suggesting recent dissatisfaction with charging infrastructure or technology.
5. **Customer Service:** Comments about customer service show steady mentions, indicating a persistent issue that Tesla might need to address to improve customer satisfaction.



**Figure 29: Top Negative Issue per Year**

This graph shows the top negative issue for each year, providing a clear view of the predominant concerns over time.

1. **Interior Quality (2018, 2020):** Early comments heavily focused on interior quality, indicating initial customer dissatisfaction with the materials and build quality.
2. **Maintenance Costs (2019):** Comments on maintenance costs became prominent in 2019, reflecting concerns about the cost of repairs and servicing.
3. **Design Problems (2022-2024):** From 2022 onwards, design problems became the dominant issue, with users consistently highlighting flaws in the vehicle's design, such as the turn signal stalk and overall aesthetics.

The most negative comments for each year provide detailed feedback on specific issues:

- **Interior Quality:** "The packaging is poor interior space is cramped and storage is insufficient for the size and despite being a hatch it is less overall spacious than a model 3 that platform has been out since 2018 and it is overdue..."
- **Design Problems:** "The data disagrees with you teslas are accident prone and a large contributor to that is the non-standard controls in fact studies show that people who regularly switch between more traditional cars and teslas are even more accident-prone..."

## Conclusion

### Summary of Main Findings

My comprehensive analysis of social media comments related to the Tesla Model 3 has revealed several critical insights into the experiences and sentiments of electric vehicle owners. By leveraging data from Reddit and YouTube, I collected a diverse dataset spanning from 2019 to 2024, ensuring a broad temporal coverage of user feedback. The key findings from my analysis are as follows:

1. **Prevalent Issues:**
  - **Design Problems:** The most frequently mentioned issue, with users expressing dissatisfaction with specific features like the turn signal stalk and the overall aesthetic appeal of the car.
  - **Interior Quality:** A recurring concern, with many comments highlighting the poor build quality and use of cheap materials in the interior.
  - **Maintenance Costs:** High maintenance costs and expensive repairs were significant pain points for many owners.
  - **Charging Issues:** Users reported long charging times and inconvenience due to inadequate charging infrastructure.
  - **Content Moderation/Privacy Issues:** Concerns about data privacy and content moderation, with users expressing unease about how their data is handled and the transparency of Tesla's data policies.
  - **Customer Service:** There were numerous complaints about unresponsive and unhelpful customer service.
2. **Positive Feedback:**
  - **Performance:** Many comments praised the performance of the Tesla Model 3, highlighting its acceleration, handling, and overall driving experience.
  - **Innovation:** Users appreciated the innovative features, such as Autopilot and the advanced infotainment system, which set the Tesla Model 3 apart from other vehicles.
  - **Environmental Impact:** Numerous comments mentioned the positive environmental impact of driving an electric vehicle, with users feeling good about reducing their carbon footprint.
  - **Community and Support:** Some comments highlighted the strong community of Tesla owners and the support and camaraderie they found among fellow owners.
3. **Sentiment Analysis:**
  - The sentiment scores showed a significant number of neutral comments, with a balanced mix of positive and negative feedback.

- Positive comments often praised the performance and innovative features of the Tesla Model 3, while negative comments focused on the issues mentioned above.
4. **Temporal Trends:**
- Specific issues like design problems and interior quality saw spikes in certain periods, indicating possible triggers such as model updates or feature changes.
  - The trend of negative comments over time highlighted persistent concerns, particularly in recent years, with design problems being the most dominant issue.

### Critical Evaluation

The analysis provided a detailed understanding of the sentiments and issues faced by Tesla Model 3 owners. However, there are several limitations and areas for improvement:

1. **Data Imbalance:** The dataset was not evenly distributed across the years, with more comments in 2024 than in previous years. This could skew the analysis towards more recent sentiments.
2. **Source Limitations:** The data was collected from Reddit and YouTube, which may not fully represent the broader population of Tesla Model 3 owners. Including data from other platforms could provide a more comprehensive view.
3. **Subjectivity in Comments:** The analysis relied on user-generated content, which is inherently subjective and may not always accurately reflect the actual issues with the vehicle.

### Suggested Actions

Based on my findings, I recommend the following actions to improve the experience for Tesla Model 3 owners:

1. **Design Improvements:**
  - Address specific design issues highlighted by users, such as the turn signal stalk and other aesthetic concerns.
  - Conduct user testing and gather feedback during the design phase to ensure new features meet customer expectations.
2. **Enhance Interior Quality:**
  - Invest in higher quality materials and improve the build quality of the interior to address concerns about cheap and uncomfortable interiors.
  - Consider offering customizable interior options to meet diverse customer preferences.
3. **Reduce Maintenance Costs:**
  - Introduce more affordable maintenance plans and transparent pricing for repairs and services.
  - Provide detailed maintenance guidelines to help owners manage their vehicles better and potentially reduce costs.
4. **Improve Charging Infrastructure:**
  - Expand the charging network and reduce charging times to alleviate the inconvenience faced by users.
  - Explore partnerships with other companies to provide more charging options and better coverage.
5. **Customer Service Enhancement:**
  - Invest in training and resources to improve the responsiveness and effectiveness of customer service.
  - Implement a robust feedback system to continuously monitor and improve customer support.
6. **Address Software and Privacy Concerns:**
  - Ensure regular updates to the software to fix bugs and enhance the user experience.
  - Maintain transparency about data privacy policies and ensure users' data is handled securely.
  - Provide clear information on content moderation practices and privacy measures to alleviate user concerns.

By taking these actions, Tesla can address the critical concerns of their customers, leading to increased satisfaction and loyalty. This proactive approach will not only improve the ownership experience but also strengthen Tesla's reputation as a leader in the electric vehicle market.