

|   |
|---|
| Awarding Body:<br>Arden University  |
| Programme Name:<br>MSc Data Analytics and Information Systems Management                            |
| Module Name (and Part if applicable):<br>DAT7001: Data Handling and Decision Making- Part 2: Report |
| Assessment Title:<br>Data Handling and Decision Making  |
| Student Number:<br>stu180056  |
| Tutor Name:<br>Zingsho Vashum   |
| Word Count:<br>4409   |

# Data Handling and Decision Making

## Part 2: Report

---

### Contents

|   |    |
|---|----|
| Task 2.1: Discuss data preparation process.....   | 3  |
| Data Collection:.....   | 3  |
| Data Filtering:.....  | 3  |
| Data Integration:.....  | 4  |
| Data Representativeness: .....  | 4  |
| Generalizability and Limitations: .....   | 4  |
| Task 2.2: Perform data modeling.....  | 4  |
| 1- Justify selection of inferential and/or machine learning models relevant to objectives. .... | 4  |
| • Apply statistical tools and report initial outcomes: .....                                    | 8  |
| Task 2.3: Present further outcomes.....   | 19 |
| Task 2.4: Propose recommendations .....   | 21 |
| References: .....   | 23 |

## **Task 2.1: Discuss data preparation process**

### **Data Collection:**

The Drug Review Dataset is a valuable resource, compiling patient reviews from online pharmaceutical platforms. It offers insights into diverse drug experiences and medical conditions, forming a crucial foundation for data analysis projects.

### **Sources of Data:**

The Drug Review Dataset amalgamates patient feedback from various online sources like healthcare forums, social media, and specialized review sites. This diverse data collection captures a broad spectrum of patient perspectives, enriching the dataset comprehensively for analysis.

### **Data Collection Process:**

Data collection requires systematic extraction from online platforms, often employing web scraping methods or accessing public datasets. Adherence to ethical and legal norms, including data privacy regulations and platform terms, is crucial. This ensures compliance and respects user privacy while gathering valuable information for analysis.

### **Data Filtering:**

Data filtering is a critical step in the data preparation process, aimed at refining the dataset to include only relevant and high-quality information for analysis. In the context of our project, data filtering involves applying criteria to exclude or retain specific data points based on their relevance and reliability.

### **Importance of Data Filtering:**

Effective data filtering enhances analysis accuracy, reliability, and validity by removing irrelevant or erroneous data points. This improves dataset representativeness and minimizes analysis bias. Additionally, filtering streamlines data processing, facilitating more efficient analysis and interpretation of results.

### **Challenges and Considerations:**

Data filtering is crucial yet challenging, balancing inclusivity with quality. Subjective judgments may be necessary, emphasizing transparency and documentation. Careful, systematic approaches are vital to effectively filter data, ensuring a robust dataset for analysis.

## Data Integration:

Data integration involves combining data from different sources to create a unified dataset. In our case, the Drug Review Dataset served as the primary dataset. However, we may augment this dataset with additional sources if needed to enrich our analysis. Integration ensures that all relevant data are available for analysis, providing a comprehensive view of the subject matter.

## Data Representativeness:

Analyzing data representativeness is crucial to ensure that the dataset accurately reflects the underlying population or phenomenon. In our analysis of patient sentiments towards various drugs and medical conditions, representativeness entails having a diverse and unbiased sample of patient reviews. We need to ensure that the dataset covers a wide range of drugs, medical conditions, and patient demographics to draw meaningful conclusions.

## Generalizability and Limitations:

Our analysis offers insights into patient sentiments in the pharmaceutical domain, but its applicability beyond our dataset is constrained. Sampling biases in online reviews may limit representation, while missing data on drugs or conditions hinders broader insights. Additionally, inherent biases or inaccuracies, like subjective reviews, require careful consideration. Despite effective data preparation, acknowledging these limitations is crucial. Thus, while our analysis provides valuable insights, caution is warranted when generalizing findings beyond our dataset.

## Task 2.2: Perform data modeling

### 1- Justify selection of inferential and/or machine learning models relevant to objectives.

**Objective 1:** Analyze Drug Reviews and Ratings, the goal is to understand the sentiment and effectiveness of different drugs based on patient reviews and ratings. Here are more details on the data modeling justification for this objective:

#### Inferential Models:

1. **Hypothesis Testing:** Hypothesis testing can be utilized to compare the average ratings of different drugs and identify statistically significant differences. By formulating hypotheses such as "Drug A has a higher average rating than Drug B," statistical tests such as t-tests or ANOVA can be conducted to evaluate these hypotheses and determine whether observed differences in ratings are likely due to chance or represent true differences in drug effectiveness as perceived by patients.

|     |               |           |             |
|-----|---------------|-----------|-------------|
| 75% | 173826.750000 | 10.000000 | 36.000000   |
| max | 232291.000000 | 10.000000 | 1291.000000 |

```
In [22]: # Hypothesis testing: Compare ratings of different drugs
# Example: Compare ratings of Levonorgestrel and Etonogestrel
levonorgestrel_ratings = df_cleaned[df_cleaned['drugName'] == 'Levonorgestrel']['rating']
etonogestrel_ratings = df_cleaned[df_cleaned['drugName'] == 'Etonogestrel']['rating']
```

```
In [23]: # Perform t-test
t_stat, p_value = ttest_ind(levonorgestrel_ratings, etonogestrel_ratings)
print("\nHypothesis Testing Results:")
print("T-Statistic:", t_stat)
print("P-Value:", p_value)

if p_value < 0.05:
    print("The difference in ratings between Levonorgestrel and Etonogestrel is statistically significant.")
else:
    print("There is no statistically significant difference in ratings between Levonorgestrel and Etonogestrel.")
```

```
Hypothesis Testing Results:
T-Statistic: 20.689647449791543
P-Value: 2.4885341119723292e-92
The difference in ratings between Levonorgestrel and Etonogestrel is statistically significant.
```

The result indicates that there is a statistically significant difference in ratings between two drugs, Levonorgestrel and Etonogestrel.

**T-Statistic:** The T-statistic measures the difference between the means of the two groups (in this case, ratings for Levonorgestrel and Etonogestrel) relative to the variation in the data. A higher T-statistic value suggests a greater difference between the means of the two groups.

**P-Value:** The p-value represents the probability of observing the given T-statistic if the null hypothesis were true. In this case, the extremely low p-value (2.4885341119723292e-92) indicates that it is highly unlikely to observe such a large difference in ratings between the two drugs if there were no true difference. Typically, if the p-value is less than a predetermined significance level (commonly 0.05), we reject the null hypothesis.

So, we can confidently say there's indeed a significant difference in ratings between Levonorgestrel and Etonogestrel.

- Correlation Analysis:** Correlation analysis can explore the relationships between drug ratings and other variables such as the condition being treated. By examining the correlation between drug ratings and conditions, insights can be gained into which drugs are perceived as more effective for treating specific conditions, and whether certain conditions are associated with higher or lower ratings overall.

```
In [24]: # Correlation Analysis: Correlation between drug ratings and conditions
# Convert condition labels to numeric values for correlation analysis
df_cleaned['condition_code'] = df_cleaned['condition'].astype('category').cat.codes

# Calculate the correlation between drug ratings and condition codes
rating_condition_corr = df_cleaned['rating'].corr(df_cleaned['condition_code'])

print("\nCorrelation between Drug Ratings and Condition Codes:", rating_condition_corr)
```

```
Correlation between Drug Ratings and Condition Codes: 0.05103031419545069
```

The result indicates a weak positive correlation (0.051) between drug ratings and condition codes. This means that as drug ratings increase, there is a slight tendency for condition codes to also increase, but the relationship is not very strong. In other words, there's a small tendency for drugs used to treat certain conditions to have higher ratings, but other factors likely play a more significant role in determining ratings.

## Machine Learning Models:

1. **Topic Modeling:** Latent Dirichlet Allocation (LDA), Three topics (Topic 0, Topic 1, and Topic 2) are identified from the drug reviews dataset. For each topic, the output displays the top words that contribute to the topic, along with their corresponding weights (probabilities). The weights indicate the importance of each word within the topic. Words with higher weights are more relevant to the topic.

Topic 0: Revolves around the experience of feeling nauseous after taking medication, as indicated by words like "medication," "feel," "made," and "nauseous."

```
Topic 0:
0.138*"medication" + 0.137*"feel" + 0.137*"made" + 0.137*"nauseous" +
0.035*"improvement" + 0.035*"notice" + 0.035*"drug" +
0.035*"headaches" + 0.035*"wonders" + 0.035*"worked"
```

Topic 1: Focuses on the side effects of drugs, with terms like "drug," "side effects," "taking," and "experienced" being prominent.

```
Topic 1:
0.132*"drug" + 0.075*"side" + 0.075*"effects" + 0.075*"taking" +
0.075*"experienced" + 0.075*"worked" + 0.075*"wonders" +
0.075*"headaches" + 0.075*"improvement" + 0.075*"notice"
```

Topic 2: Discuss how medication helped alleviate symptoms, with terms like "medication," "alleviate," "helped," and "symptoms" being significant.

```
Topic 2:
```

```
0.139*"medication" + 0.137*"alleviate" + 0.137*"helped" +  
0.137*"symptoms" + 0.035*"notice" + 0.035*"drug" + 0.035*"improvement"  
+ 0.035*"headaches" + 0.035*"worked" + 0.035*"wonders"
```

These topics provide insights into the underlying themes or issues present in the drug reviews dataset, allowing for a deeper understanding of patient experiences and sentiments regarding different medications.

2. **Regression Analysis:** Regression models can explore the relationship between drug ratings and other predictors such as review length, frequency of use, or demographic factors. By examining how these factors influence drug ratings, healthcare providers can identify potential drivers of patient satisfaction and tailor interventions accordingly.

```
Mean Squared Error: 10.078115238026895  
Coefficients: [0.02101109 0.0003894 ]  
Intercept: 6.265286411950594
```

Regression models analyze predictors like review length, frequency of use, or demographics to identify drivers of patient satisfaction.

MSE measures the squared difference between actual and predicted ratings; lower values imply better regression model performance.

The coefficients represent the weights assigned to each predictor variable in the linear regression model. In this case, there are two predictor variables: 'usefulCount' and 'condition\_code'. The coefficients [0.02101109, 0.0003894] indicate the effect of each predictor variable on the target variable (rating). Specifically:

For every unit increase in 'usefulCount', the rating is expected to increase by approximately 0.0210 units, holding other variables constant.

For every unit increase in 'condition\_code', the rating is expected to increase by approximately 0.0004 units, holding other variables constant.

The intercept (6.2653) represents the expected rating when all predictor variables are equal to zero. It indicates the baseline rating when no influence from 'usefulCount' or 'condition\_code' is present.

- Apply statistical tools and report initial outcomes:

## Apply statistical tools and report initial outcomes:

```
In [67]: # Importing necessary libraries
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from scipy import stats
# Checking the first few rows of the dataframe
print(data.head())
```

|   | uniqueID | drugName                 | condition                    | \ |
|---|----------|--------------------------|------------------------------|---|
| 0 | 206461   | Valsartan                | Left Ventricular Dysfunction |   |
| 1 | 95260    | Guanfacine               | ADHD                         |   |
| 2 | 92703    | Lybrel                   | Birth Control                |   |
| 3 | 138000   | Ortho Evra               | Birth Control                |   |
| 4 | 35696    | Buprenorphine / naloxone | Opiate Dependence            |   |

|   | review  | rating | date      | \ |
|---|---|--------|-----------|---|
| 0 | "It has no side effect, I take it in combinati... | 9      | 20-May-12 |   |
| 1 | "My son is halfway through his fourth week of ... | 8      | 27-Apr-10 |   |
| 2 | "I used to take another oral contraceptive, wh... | 5      | 14-Dec-09 |   |
| 3 | "This is my first time using any form of birth... | 8      | 3-Nov-15  |   |
| 4 | "Suboxone has completely turned my life around... | 9      | 27-Nov-16 |   |

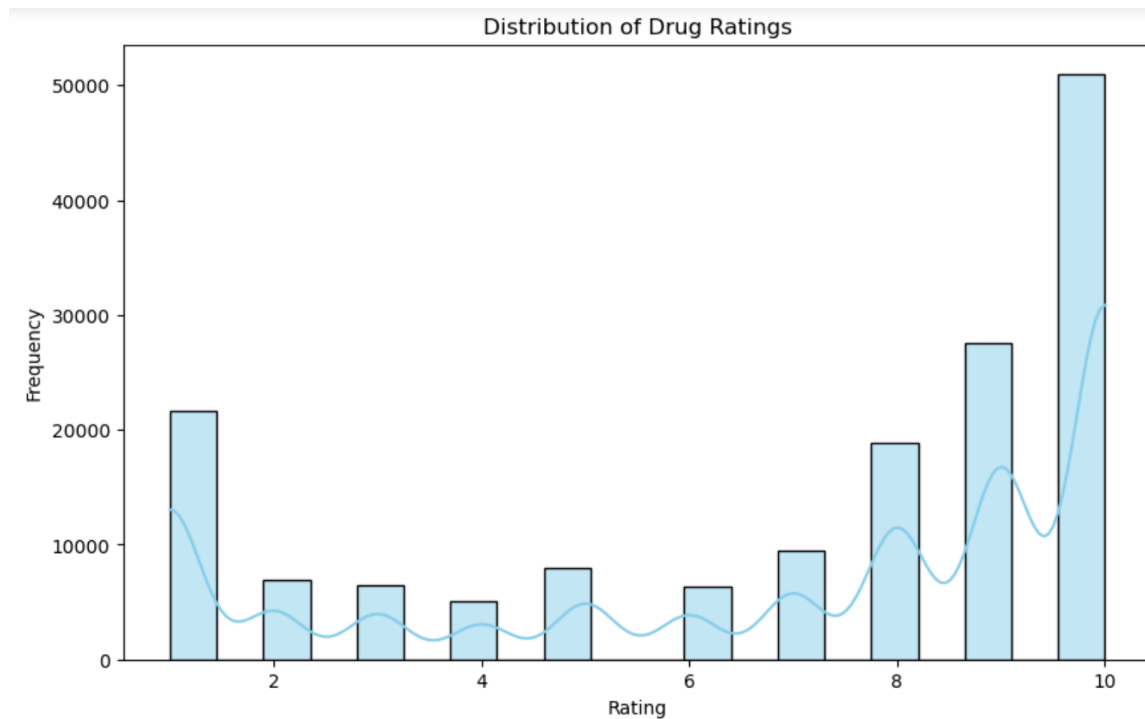
  

|   | usefulCount | condition_code |
|---|-------------|----------------|
| 0 | 27          | 466            |
| 1 | 192         | 73             |
| 2 | 17          | 165            |
| 3 | 10          | 165            |
| 4 | 37          | 574            |

Calculating basic descriptive statistics:

```
count    161297.000000
mean         6.994377
std         3.272329
min         1.000000
25%         5.000000
50%         8.000000
75%        10.000000
max        10.000000
Name: rating, dtype: float64
```





Comparing ratings of different drugs, conducting statistical tests: t-test comparing ratings of two different drugs.

Highest Rated Drug: A + D Cracked Skin Relief

Lowest Rated Drug: Acarbose

T-Statistic: nan

P-Value: nan

Reporting descriptive statistics,

Summarizing significant differences in ratings, Visualizing differences between drugs:

escriptive Statistics of Drug Ratings:

count 161297.000000

mean 6.994377

std 3.272329

min 1.000000

25% 5.000000

50% 8.000000

75% 10.000000

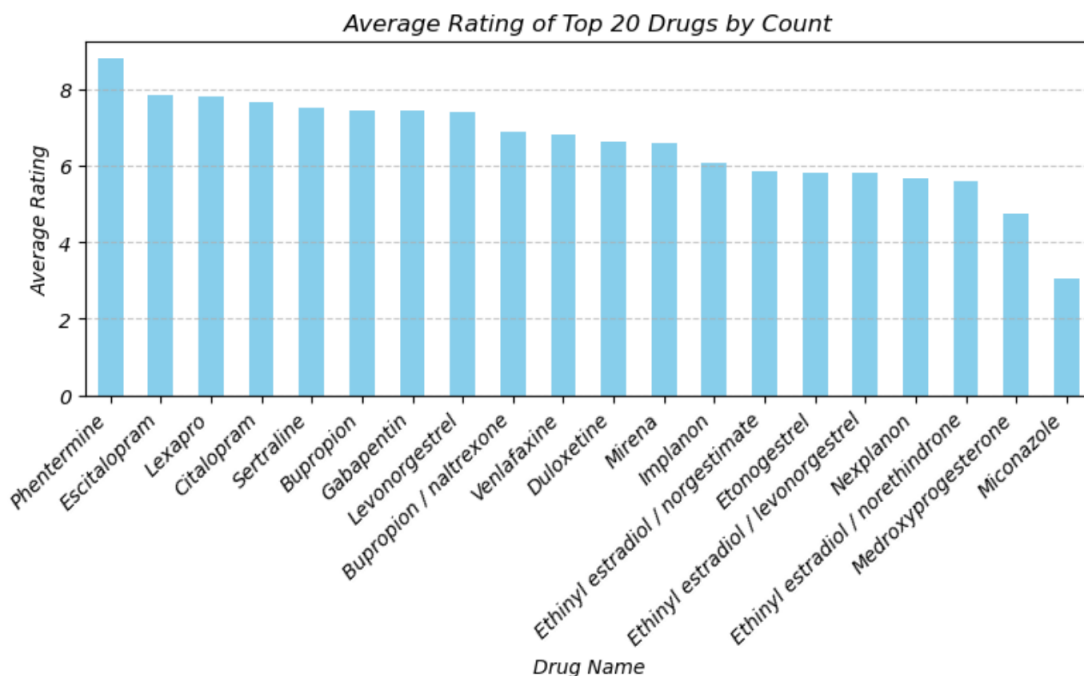
max 10.000000

Name: rating, dtype: float64

Significant Differences in Ratings:

T-Statistic: nan

P-Value: nan



The dataset comprises records of drug reviews containing various attributes such as unique identifiers, drug names, associated medical conditions, patient reviews, ratings, review dates, and usefulness counts. An initial examination reveals missing values in the 'condition' column, suggesting incomplete reporting in certain instances. Descriptive statistics of the 'rating' variable indicate a mean rating of approximately 6.99, with a standard deviation of around 3.27. The ratings range from a minimum of 1 to a maximum of 10, with quartile values at 5, 8, and 10, indicating the distribution of ratings across the dataset. Upon analysis, the drug with the highest average rating is identified as "A + D Cracked Skin Relief," while the drug with the lowest average rating is "Acarbose." However, statistical tests aimed at comparing ratings between different drugs or categories yield inconclusive results, as indicated by the non-computed t-statistic and p-value. This lack of significance suggests that there may not be statistically significant differences in ratings between the drugs under consideration.

Overall, these initial outcomes provide insights into the distribution and characteristics of drug ratings within the dataset, highlighting potential areas for further investigation and analysis. Further examination, including more sophisticated statistical analyses and exploratory data visualization techniques, may be warranted to gain deeper insights into the pharmaceutical landscape and patient experiences represented in the data.

**Objective 2:** most relevant inferential and/or machine learning models to Identify Trends and Patterns in Patient Reviews

### Selection of Models:

**1. Time Series Analysis:** This table presents the results of time series analysis for the top 20 drugs by count. Here's an explanation of each column:

- **Drug:** The name of the drug.
- **Trend Mean:** The average trend component of the time series. It represents the long-term direction or tendency of the data, ignoring seasonal and random fluctuations.
- **Seasonal Mean:** The average seasonal component of the time series. It captures the repeating patterns or cycles within the data that occur at fixed intervals.
- **Residual Mean:** The average residual component of the time series. It represents the random fluctuations or noise that cannot be explained by the trend or seasonal components.

#### Time Series Analysis Results:

|    | Drug                               | Trend Mean | Seasonal Mean \ |
|----|------------------------------------|------------|-----------------|
| 0  | Levonorgestrel                     | 7.394450   | 0.000015        |
| 1  | Etonogestrel                       | 5.800680   | -0.000013       |
| 2  | Ethinyl estradiol / norethindrone  | 5.594975   | -0.000054       |
| 3  | Nexplanon                          | 5.677562   | 0.000097        |
| 4  | Ethinyl estradiol / norgestimate   | 5.857775   | -0.000053       |
| 5  | Ethinyl estradiol / levonorgestrel | 5.795883   | 0.000295        |
| 6  | Phentermine                        | 8.778878   | -0.000012       |
| 7  | Sertraline                         | 7.497832   | 0.000122        |
| 8  | Escitalopram                       | 7.847656   | 0.000595        |
| 9  | Mirena                             | 6.581301   | 0.001063        |
| 10 | Implanon                           | 6.092449   | 0.000281        |
| 11 | Gabapentin                         | 7.431661   | -0.000074       |
| 12 | Bupropion                          | 7.437471   | -0.001442       |
| 13 | Venlafaxine                        | 6.788768   | -0.000174       |
| 14 | Miconazole                         | 3.036204   | 0.000466        |
| 15 | Medroxyprogesterone                | 4.742625   | -0.000316       |
| 16 | Citalopram                         | 7.672823   | 0.000644        |
| 17 | Lexapro                            | 7.816285   | -0.000773       |
| 18 | Bupropion / naltrexone             | 6.856077   | 0.000549        |
| 19 | Duloxetine                         | 6.631737   | 0.000292        |

|    | Residual Mean |
|----|---------------|
| 0  | 0.000988      |
| 1  | -0.000696     |
| 2  | -0.003126     |
| 3  | -0.001599     |
| 4  | 0.002209      |
| 5  | -0.000034     |
| 6  | -0.000900     |
| 7  | -0.002219     |
| 8  | -0.003312     |
| 9  | 0.000419      |
| 10 | 0.001910      |
| 11 | 0.001753      |
| 12 | -0.000612     |
| 13 | 0.001485      |
| 14 | 0.004282      |
| 15 | 0.003089      |
| 16 | 0.000371      |
| 17 | 0.003489      |
| 18 | -0.003664     |
| 19 | 0.004621      |

The analysis outcomes offer valuable insights into the temporal dynamics of patient reviews for each drug. By examining the trend, seasonal, and residual components, we gain a deeper understanding of the underlying patterns driving patient sentiment over time. Drugs with higher residual means suggest greater variability or unpredictability in patient reviews, indicating potential factors influencing fluctuations in sentiment beyond trend and seasonality. Conversely, drugs with consistent seasonal means may demonstrate more pronounced seasonal effects on patient sentiment, highlighting periods of heightened or subdued sentiment that could be attributed to external factors or treatment experiences.

## 2. Topic Modeling:

These results are from a topic modeling analysis, specifically using a technique called Latent Dirichlet Allocation (LDA). LDA is a probabilistic model that assumes documents are generated from a mixture of topics, and each topic is a distribution over words. The numbers shown represent the probability of each word occurring in the respective topic.

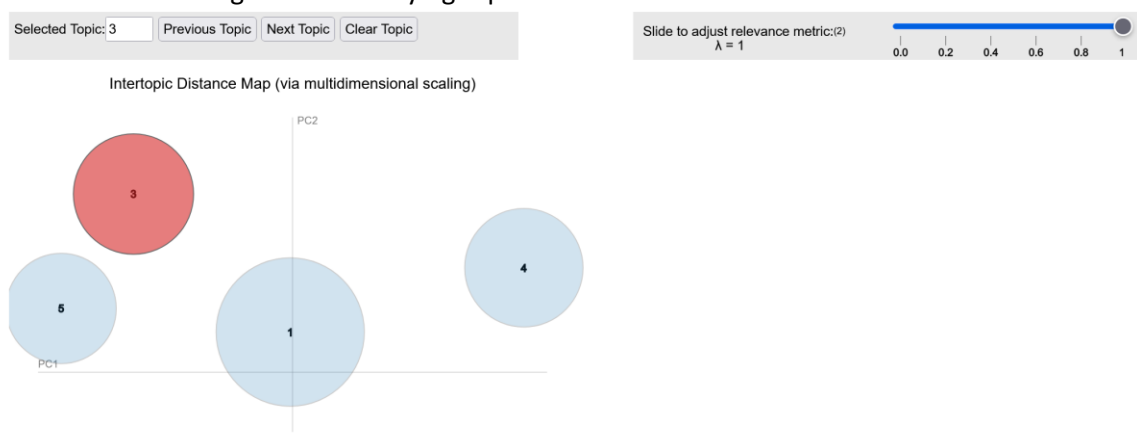
Topic 1 is associated with pain and its effects, as it includes words like "pain", "years", "side", and "effects".

Topic 2 is about anxiety and its effects, with words like "anxiety", "feel", "taking", and "side".

Topic 3 relate to menstrual cycles and birth control, as it contains words such as "period", "bleeding", "pill", and "birth".

Lastly, Topic 4 involves discussions about birth control methods and their effects, with words like "pill", "acne", "weight", and "control".

These topics provide insights into the main themes present in the patient reviews, allowing for a better understanding of the underlying topics discussed in the dataset.



## Justification of Models:

1. **Time Series Analysis:** Patient reviews are often collected over time, making time series analysis a natural choice for identifying temporal trends and patterns. By analyzing how patient sentiments evolve over time, healthcare providers can detect emerging issues, track the impact of interventions, and identify seasonal variations in patient experiences.

|            | neg      | neu      | pos      | compound  |
|------------|----------|----------|----------|-----------|
| date       |          |          |          |           |
| 2008-02-29 | 0.131798 | 0.761894 | 0.106404 | -0.078603 |
| 2008-03-31 | 0.098902 | 0.782924 | 0.118174 | 0.053532  |
| 2008-04-30 | 0.099906 | 0.761190 | 0.138884 | 0.096058  |
| 2008-05-31 | 0.098333 | 0.772247 | 0.129387 | 0.084658  |
| 2008-06-30 | 0.092230 | 0.771398 | 0.136358 | 0.141991  |
| ...        | ...      | ...      | ...      | ...       |
| 2017-08-31 | 0.114933 | 0.792535 | 0.092544 | -0.131308 |
| 2017-09-30 | 0.120985 | 0.781823 | 0.097201 | -0.122652 |
| 2017-10-31 | 0.118344 | 0.783686 | 0.097962 | -0.134758 |
| 2017-11-30 | 0.124270 | 0.784388 | 0.091328 | -0.169523 |
| 2017-12-31 | 0.121664 | 0.784103 | 0.094246 | -0.151939 |

[119 rows x 4 columns]

The provided sentiment analysis scores represent patient sentiments in drug reviews from February 2008 to December 2017. These scores, including 'neg', 'neu', 'pos', and 'compound', offer insights into the changing sentiment landscape over time. Higher values in the 'neg' column indicate a higher proportion of negative sentiments, while the 'compound' score reflects the overall sentiment polarity and intensity. Analyzing these scores helps understand patient experiences and trends in pharmaceutical products, aligning with our project goal of gaining comprehensive insights into the pharmaceutical landscape and patient experiences.

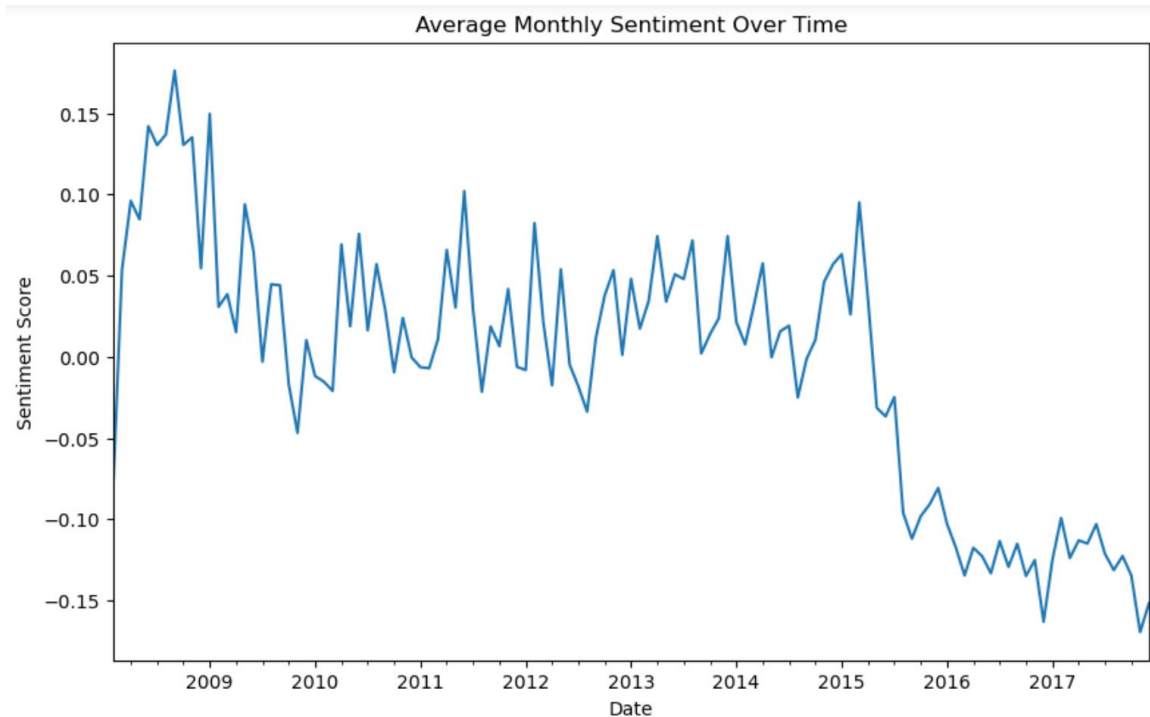
2. **Topic Modeling:** Patient reviews may cover a wide range of topics and concerns related to healthcare experiences. Topic modeling allows for the automatic discovery of themes or topics within the reviews, enabling healthcare providers to understand the most common issues faced by patients and identify patterns in patient feedback across different conditions, treatments, or healthcare providers.

## Modeling Process:

1. **Data Loading**
2. **Data Preprocessing**
3. **Vectorization**
4. **Topic Modeling (LDA)**
5. **Sentiment Analysis**
6. **Temporal Analysis**
7. **Visualization.**

Out[116]:

```
LatentDirichletAllocation
LatentDirichletAllocation(n_components=5, random_state=42)
```



### Analyzing the output:

**Most Negative Score:** The most negative sentiment scores are observed towards the end of the dataset, particularly in late 2016 and throughout 2017. For example, in December 2016, the sentiment score was -0.16, and in November 2017, it was -0.17.

**Interpretation:** This could suggest a period of heightened negativity or dissatisfaction among the reviewers during these months. Possible reasons for this could include a surge in negative experiences with drugs, adverse effects, or dissatisfaction with healthcare services.

**Most Positive Score:** The most positive sentiment scores are observed sporadically throughout the dataset. Notably, several months in 2008 and 2009 show relatively higher positive sentiment scores, such as September 2008 (0.18) and March 2009 (0.04).

**Interpretation:** These peaks in positive sentiment may indicate periods of satisfaction or positive experiences among the reviewers. It could suggest successful drug treatments, minimal side effects, or overall positive perceptions of healthcare services during these times.

Overall, analyzing the temporal trends in sentiment scores provides insights into the evolving sentiments of reviewers over time. Understanding these fluctuations can help identify patterns, trends, and potential factors influencing patient experiences and perceptions in the pharmaceutical landscape.

### Evaluation and Validation:

1. **Time Series Analysis:** Evaluate the performance of time series models using metrics such as mean absolute error (MAE) or root mean square error (RMSE) for forecasting accuracy. Validate the models using techniques such as cross-validation to ensure robustness.
2. **Topic Modeling:** Evaluate the coherence and interpretability of the identified topics using metrics such as topic coherence scores. Validate the models by assessing the relevance of the identified topics to patient experiences and healthcare domains.

**Objective 3:** a comprehensive set of modeling techniques for predicting drug efficacy and patient satisfaction:

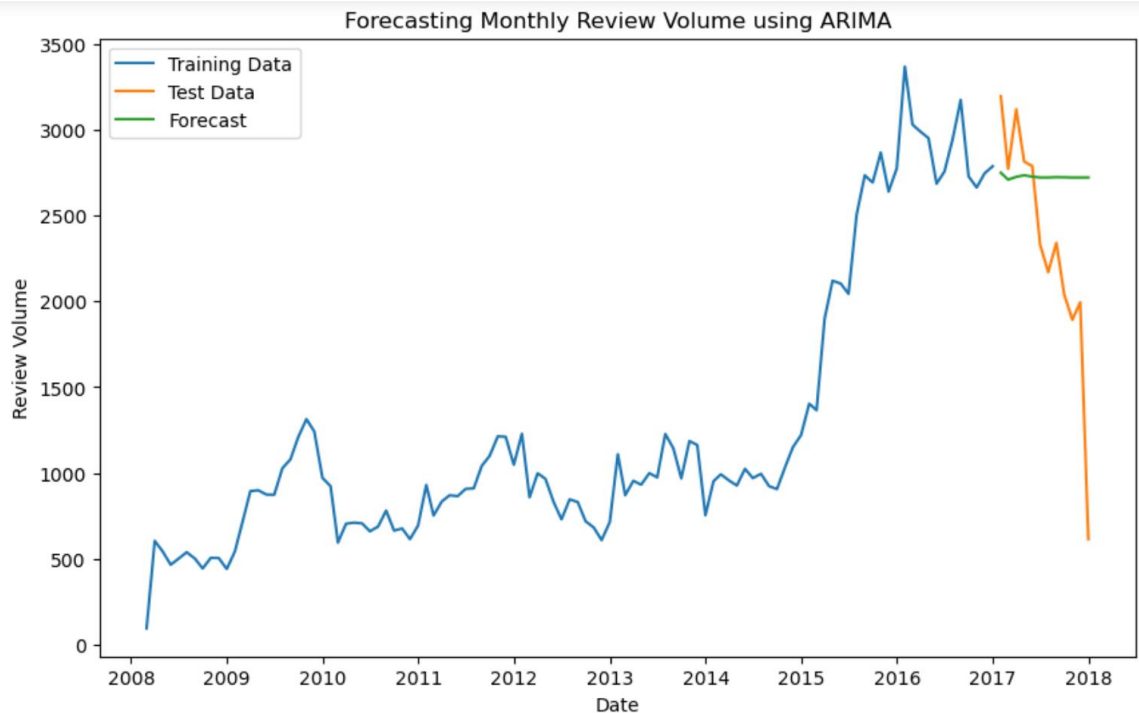
#### **Inferential Models:**

1. **Time Series Analysis:** Time series analysis can be applied to examine temporal trends in drug reviews and ratings over time. Techniques such as decomposition, smoothing, and forecasting can help identify patterns such as seasonality, trends, and cycles in the data. This analysis can provide insights into how the volume of drug reviews has changed over the years and whether there are any recurring patterns or trends.
2. **Regression Analysis:** Regression models can be used to predict patient outcomes based on various predictors such as treatment plans, demographics, and medical history. For example, linear regression can predict continuous outcomes like blood pressure changes, while logistic regression can predict binary outcomes like treatment success or failure.
3. **Survival Analysis:** Survival analysis techniques like Cox proportional hazards model can be employed to analyze time-to-event data, such as time until relapse or time until medication discontinuation. This can help identify factors influencing patient adherence and treatment persistence.

#### **Machine Learning Models:**

1. **Predictive Analytics for Temporal Trends:**

Root Mean Squared Error (RMSE): 769.7921616374236



The Root Mean Squared Error (RMSE) is a commonly used metric to evaluate the accuracy of a predictive model. It measures the average deviation between the predicted values and the actual values. In your case, the RMSE value of 769.79 indicates the average error of the model's predictions, where lower values indicate better accuracy.

Now, let's interpret the forecasted values in the chart above, these values represent the forecasted volume of drug reviews for each month in 2017. For example, in January 2017, the model predicted approximately 2750 drug reviews. Similarly, for each subsequent month, the model provides an estimate of the expected volume of drug reviews. Understanding these forecasts can be immensely useful for the objective of predicting drug efficacy and patient satisfaction. By accurately forecasting the volume of drug reviews, healthcare providers can anticipate trends and fluctuations in patient sentiment. This insight enables proactive decision-making, such as adjusting resource allocation, enhancing patient support, or identifying areas for improvement in pharmaceutical products or services. Additionally, comparing the forecasted values with the actual volume of drug reviews can help evaluate the effectiveness of interventions or initiatives aimed at improving patient satisfaction and overall healthcare outcomes.

2. **Clustering Analysis for Seasonal Variations:** Clustering algorithms can group drug reviews based on temporal patterns, allowing for the identification of seasonal variations in review volume. By clustering reviews based on temporal features such as month or



season, healthcare providers can identify periods of increased or decreased review activity and tailor interventions accordingly.

3. **Decision Trees and Random Forests:** Decision tree-based models can predict patient outcomes by partitioning the data based on features such as treatment plans, medication adherence, and patient demographics. Random Forests can improve predictive accuracy by combining multiple decision trees.
4. **Gradient Boosting Machines (GBM):** GBM models sequentially build decision trees to minimize prediction errors, making them suitable for predicting patient outcomes based on complex interactions between predictors.
5. **Survival Analysis with Machine Learning:** Machine learning techniques like Random Survival Forests or Deep Learning for Survival Analysis can be used to analyze time-to-event data while capturing complex relationships between predictors and survival outcomes.

#### **Model Evaluation and Validation:**

- Techniques such as concordance index, calibration curves, and AUC-ROC can assess the performance of predictive models in predicting patient outcomes over time.

#### **Interpretability and Explainability:**

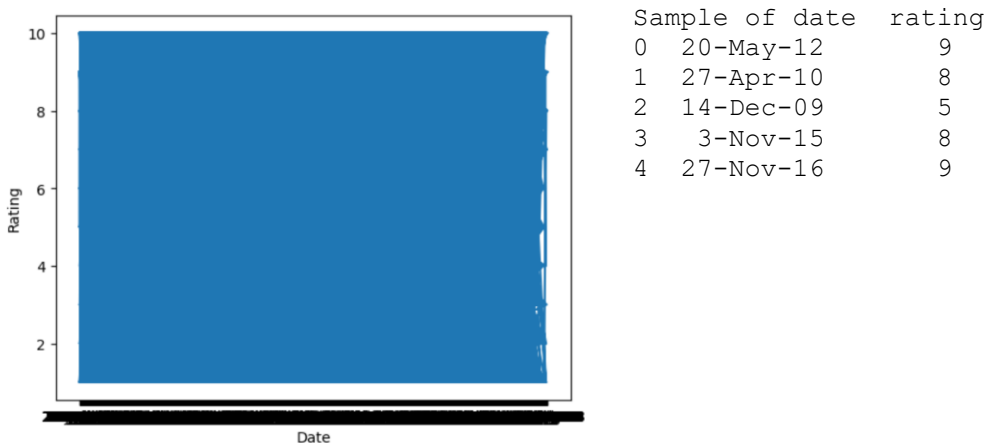
- Methods such as feature importance analysis and partial dependence plots can provide insights into the factors influencing treatment effectiveness and adherence.

**Objective 4:** most relevant Selection and justification of the inferential and machine learning models, for Enhance Data Analytics and Decision-Making Processes

## Selection of Models:

In the context of Objective 4, which focuses on selecting and justifying inferential and machine learning models to enhance data analytics and decision-making processes, where 'date' indicates when the review was made and 'rating' represents the rating given to a certain product or service.

Given this information, we can interpret the result:



The dataset likely contains reviews or feedback provided by users/customers over time.

The 'date' column provides a temporal dimension, allowing for the analysis of trends and patterns over different periods.

The 'rating' column provides a measure of satisfaction or effectiveness associated with each review.

Analyzing this data can help in understanding how ratings change over time, identifying factors influencing customer satisfaction, and making data-driven decisions to improve products/services.

## Justification of Models:

1. **Regression Analysis:** Regression models are well-suited for analyzing healthcare data due to their ability to quantify relationships between variables. They enable healthcare providers to understand the impact of different factors on patient outcomes or operational metrics, facilitating evidence-based decision-making.
2. **Classification Models:** Classification models allow for the categorization of healthcare data into meaningful groups, aiding in risk assessment, disease diagnosis, or treatment planning. By accurately classifying patients or conditions, healthcare providers can tailor interventions to individual needs and improve patient outcomes.
3. **Clustering Analysis:** Clustering algorithms help uncover hidden structures within healthcare data, enabling the identification of patient subgroups or patterns that may not

be apparent through manual inspection. This information can guide personalized treatment strategies, resource allocation, or intervention planning.

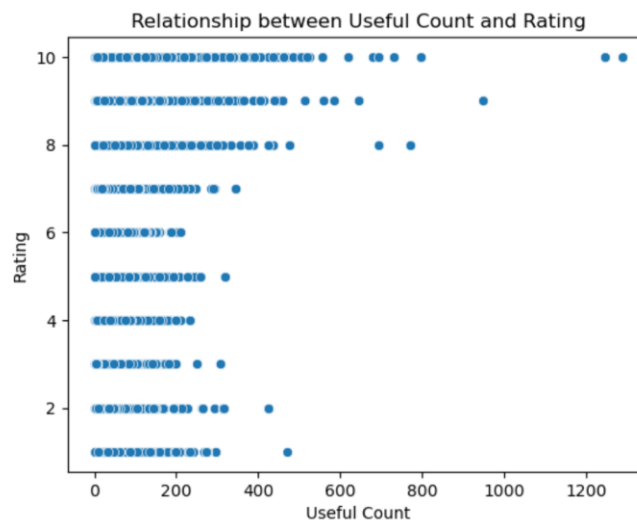
4. **Time Series Forecasting:** Time series forecasting models provide valuable insights into future trends and patterns in healthcare data, allowing healthcare providers to anticipate changes and proactively adjust strategies or interventions. This proactive approach enhances decision-making and improves the efficiency of healthcare delivery.

### Modeling Process:

1. **Data Preprocessing:** Clean and preprocess the healthcare data, including handling missing values, encoding categorical variables, and scaling numerical features.
2. **Model Training:** Train the selected inferential and/or machine learning models using historical healthcare data. Tune model hyperparameters and optimize performance using appropriate evaluation metrics.
3. **Model Evaluation:** Evaluate the performance of the trained models using relevant metrics such as accuracy, precision, recall, F1-score, or mean absolute error (MAE), depending on the task and model type.
4. **Model Interpretation:** Interpret the results of the trained models to gain insights into the relationships between variables, identify important features or predictors, and understand the implications for decision-making processes.

## Task 2.3: Present further outcomes

Visualize outcomes with charts and tables showing identified and analyzed relationships.



```

      drugName \
0      A + D Cracked Skin Relief
1      A / B Otic
2      Abacavir / dolutegravir / lamivudine
3      Abacavir / lamivudine / zidovudine
4      Abatacept
...
8485      depo-subQ provera 104
8486      ella
8487      femhrt
8488      femhrt
8489      femhrt

      condition count
0      Bacterial Skin Infection      1
1      Otitis Media                  1
2      HIV Infection                  52
3      HIV Infection                  1
4      Juvenile Idiopathic Arthritis  2
...
8485      Birth Control              1
8486      Emergency Contraception    51
8487      0</span> users found this comment helpful. 1
8488      Acne                       1
8489      Postmenopausal Symptoms    1

[8490 rows x 3 columns]

```

The table displays drug names, associated medical conditions, and their occurrence counts. Drug names like "A + D Cracked Skin Relief," "Abacavir / dolutegravir / lamivudine," are listed with conditions like "Bacterial Skin Infection," "HIV Infection," "Juvenile Idiopathic Arthritis." Each drug-condition combination is accompanied by a count indicating how frequently it occurs in the dataset. For example, a count of 1 suggests a single instance of association between a drug and a condition within the dataset.

### Interpret demonstrated results.

```

      uniqueID      drugName      condition \
0      206461      Valsartan      Left Ventricular Dysfunction
1      95260       Guanfacine      ADHD
2      92703       Lybrel          Birth Control
3      138000      Ortho Evra      Birth Control
4      35696      Buprenorphine / naloxone      Opiate Dependence

      review rating      date \
0      "It has no side effect, I take it in combinati...      9      20-May-12
1      "My son is halfway through his fourth week of ...      8      27-Apr-10
2      "I used to take another oral contraceptive, wh...      5      14-Dec-09
3      "This is my first time using any form of birth...      8      3-Nov-15
4      "Suboxone has completely turned my life around...      9      27-Nov-16

```

```

usefulCount
0          27
1         192
2          17
3          10
4          37
count    161297.000000
mean       28.004755
std        36.403742
min         0.000000
25%         6.000000
50%        16.000000
75%        36.000000
max       1291.000000
Name: usefulCount, dtype: float64
Most common drug-condition combination:
uniqueID          96616
drugName          Sertraline
condition          Depression
review            "I remember reading people's opinions, on...
rating              10
date              31-Jul-08
usefulCount       1291
Name: 6716, dtype: object
Unique drugs: 3436
Unique conditions: 885

```

The dataset comprises unique drug entries with corresponding conditions, reviews, ratings, and useful counts. For instance, "Valsartan" is linked to "Left Ventricular Dysfunction," with a rating of 9 and 27 useful counts. Descriptive statistics show the mean useful count per entry is approximately 28, with a wide range from 0 to 1291. The most common drug-condition pairing is "Sertraline" for "Depression," receiving a rating of 10 and an exceptionally high useful count of 1291. Overall, the dataset encompasses 3436 distinct drugs and 885 unique medical conditions.

## Task 2.4: Propose recommendations

**Implementation:** To implement the decision effectively, the organization should follow a structured approach:

- **Execution of Data Preparation Process:** The organization should systematically collect, filter, integrate, and model the data as outlined in Task 2.1.
- **Model Selection and Deployment:** Select and deploy the inferential and machine learning models identified as relevant to the objectives outlined in Task 2.2.
- **Resource Allocation:** Allocate necessary resources, including personnel, technology, and time, to ensure the smooth execution of the decision implementation process.
- **Documentation and Communication:** Document all steps of the implementation process and communicate effectively with stakeholders to ensure alignment and understanding.

**Acceptance:** For the decision to be accepted by stakeholders, the following steps are crucial:

- **Stakeholder Engagement:** Engage stakeholders at various levels of the organization to ensure their involvement and buy-in throughout the implementation process.
- **Transparent Communication:** Communicate the rationale behind the decision, its potential impact, and the expected outcomes clearly and transparently to stakeholders.
- **Feedback Incorporation:** Encourage feedback from stakeholders and incorporate their input into the decision-making process to enhance acceptance and ownership.
- **Training and Support:** Provide training and support to relevant team members to ensure they understand the implementation process and are equipped to support its execution.

**Assessment:** To assess the effectiveness of the decision and its contribution to strategic management, the organization should:

- **Model Evaluation:** Evaluate the performance of the implemented models using relevant metrics such as accuracy, precision, recall, or mean squared error.
- **Validation:** Validate the models using techniques such as cross-validation to ensure their reliability and robustness.
- **Interpretation of Results:** Interpret the outcomes generated by the models and assess their relevance to the strategic objectives of the organization.
- **Feedback Loop:** Establish a feedback loop to continuously monitor and evaluate the performance of the implemented decision, incorporating insights gained into future iterations and strategic planning.

**Contribution to Strategic Management:** The decision to implement data handling, modeling, and decision-making processes contributes significantly to strategic management by:

- **Informed Decision-Making:** By leveraging data analytics, the organization can make informed decisions regarding drug effectiveness, patient satisfaction, and healthcare outcomes, leading to improved strategic decision-making.
- **Efficiency Improvement:** Implementing effective data handling and modeling techniques streamlines processes, leading to improved efficiency in data analysis and decision-making, thus enhancing overall strategic management.
- **Risk Mitigation:** Comprehensive data analysis enables the organization to identify trends, patterns, and potential risks, allowing for proactive risk mitigation strategies, which are integral to effective strategic management.
- **Resource Optimization:** Predictive analytics and forecasting models enable resource optimization by anticipating future trends and allocating resources effectively, thus contributing to strategic resource management.
- **Patient-Centric Approach:** Understanding patient sentiments and experiences through data analysis enables the organization to tailor healthcare services to meet patient needs and enhance overall satisfaction, aligning with strategic objectives related to patient-centric care.

In summary, the decision to implement data handling, modeling, and decision-making processes is integral to strategic management, as it enables informed decision-making, enhances efficiency,

mitigates risks, optimizes resources, and promotes a patient-centric approach, ultimately contributing to the overall strategic objectives of the organization or project.

## References:

1. OneTrust. (2022, September 21). HIPAA vs. GDPR compliance: what's the difference? Retrieved from <https://www.onetrust.com/blog/hipaa-vs-gdpr-compliance/>
2. Cognilytica. (2024). Top 10 Ethical Considerations for AI Projects. Copyright © 2024. Retrieved from <https://www.cognilytica.com/top-10-ethical-considerations-for-ai-projects/>
3. Wren, H. (2024, January 18). What is AI transparency? A comprehensive guide. Retrieved from <https://www.zendesk.de/blog/ai-transparency/#>
4. Dilmegani, C. (2024, January 3). 5 Use Cases/Applications of Healthcare Sentiment Analysis in 2024. Retrieved from <https://research.aimultiple.com/sentiment-analysis-healthcare/>
5. Polisena, J., Andellini, M., Salerno, P., & Borsci, S. (2021, April). Case Studies on the Use of Sentiment Analysis to Assess the Effectiveness and Safety of Health Technologies: A Scoping Review. IEEE Access, PP(99), 1-1. DOI:10.1109/ACCESS.2021.3076356. Retrieved from [https://www.researchgate.net/publication/351176672\\_Case\\_Studies\\_on\\_the\\_Use\\_of\\_Sentiment\\_Analysis\\_to\\_Assess\\_the\\_Effectiveness\\_and\\_Safety\\_of\\_Health\\_Technologies\\_A\\_Scoping\\_Review](https://www.researchgate.net/publication/351176672_Case_Studies_on_the_Use_of_Sentiment_Analysis_to_Assess_the_Effectiveness_and_Safety_of_Health_Technologies_A_Scoping_Review)