# ANALYSIS OF PUREGYM REVIEWS: DATA CLEANING, SENTIMENT ANALYSIS, AND TOPIC MODELING

# Contents

## Introduction

Customer feedback is an invaluable resource for businesses looking to enhance user experience and operational efficiency. This report analyzes customer reviews of PureGym collected from Google and Trustpilot. Through data cleaning, sentiment analysis, and topic modeling, this study aims to identify recurring issues, understand the emotional tone of negative reviews, and provide actionable recommendations for improvement.

## Data Collection and Cleaning

The initial dataset consisted of reviews from Google and Trustpilot. The Google dataset contained 23,250 reviews with seven features, while the Trustpilot dataset had 16,673 reviews with fifteen features. To ensure high-quality data, extensive cleaning procedures were applied. Reviews that contained empty comments or were written in a language other than English were removed. Additional preprocessing included eliminating stopwords, removing non-alphabetic characters, and tokenizing the text for easier analysis.

Following the cleaning process, the Google dataset was reduced to 13,898 data points, while the Trustpilot dataset remained unchanged at 16,673.

## Sentiment Analysis of Cleaned Data

Examining the Google dataset after cleaning revealed that the comments contained 13,492 unique words with a total word count of 241,042. The dataset encompassed 512 unique gym locations. Analyzing the most frequently used words revealed that positive reviews commonly featured terms like "clean," "great," "nice," and "love" (see Figure 1 below).

Figure 1 Google Review Word cloud

Negative reviews accounted for 2,785 data points, containing 8,414 unique words and a total of 81,163 words. A word cloud analysis revealed a shift in frequently used terms, with words like "dirty" and "broken" becoming more prominent (see Figure 2 below).



Figure 2 Negative Reviews word cloud

Similarly, in the Trustpilot dataset, the comments contained 12,859 unique words with a total of 277,004 words, spanning 374 unique gym locations. Like the Google reviews, positive comments featured words like "clean" and "great." The negative reviews accounted for 3,543 data points, consisting of 8,410 unique words and a total of 95,444 words. The word cloud analysis revealed common complaints, particularly concerning staff behavior and shower conditions.

## Emotion Analysis of Negative Reviews

To further understand customer dissatisfaction, emotion analysis was conducted using the **bert-base-uncased-emotion** model. Both datasets exhibited the same distribution of emotions, with the only difference being that sadness and joy swapped positions in ranking. Anger emerged as the dominant emotion across both datasets, with over 1,000 occurrences each.
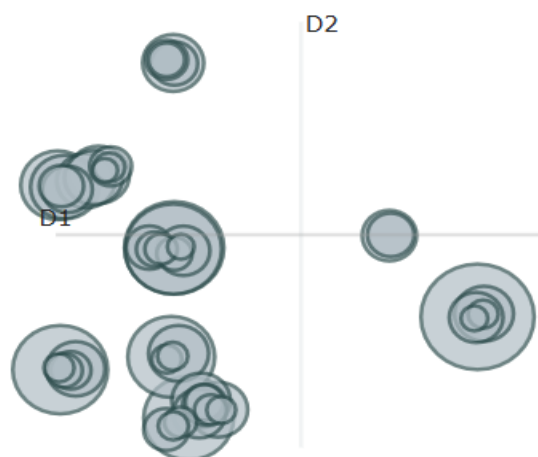
## Topic Modeling Insights

### BERTopic Analysis of Negative Reviews

Applying BERTopic to the negative reviews resulted in the classification of the dataset into 51 topics, including outliers. A total of 1,516 reviews were labeled as outliers, which was a significant proportion given that the dataset contained only 3,888 negative reviews.

An analysis of the top two topics, which contained 185 and 169 occurrences respectively, revealed that customers were particularly concerned about air conditioning problems and facility issues such as cold showers. The intertopic distance map illustrated the presence of eight distinct topics, with the remaining topics clustering around these core issues  (see figure 3 below)



Intertopic Distance Map

The top concerns identified from the analysis were related to air conditioning problems, rude staff, issues with gym classes, access problems involving PIN code entry, dirty toilets, parking fines, poor gym cleanliness, and excessively loud music. Among all the locations, London Stratford had the highest occurrence of anger-related reviews, followed by Leicester Walnut Street and London Enfield.

## Expanding Topic Modeling with Additional Data

To further refine the findings, the dataset was enhanced by merging the Google and Trustpilot reviews for locations common to both datasets. The focus remained on negative reviews to maintain consistency in sentiment analysis. Additional topics were created using the Phi model and subsequently re-analyzed using BERTopic.

Despite an increase in the number of topics to over 100, no new insights emerged. The primary issues remained the same, emphasizing concerns over staff behavior, air conditioning malfunctions, gym cleanliness, and problems related to gym classes.

## Gensim Topic Modeling

To validate the topic modeling results, Gensim was used to conduct an additional round of topic modeling with a maximum topic limit of ten. The results mirrored those obtained from BERTopic, reaffirming that the core issues remained unchanged. No new insights were gained from this approach.

## Recommendations

To address these recurring issues, actionable recommendations were developed using the Phi model. The dataset was split into two and processed separately to avoid computational errors, which led to some recommendations being repeated. These repetitions were excluded from the final list.

Overcrowding was identified as a major concern, and it is recommended that PureGym implement a reservation system to better manage capacity. Customer service training should be prioritized to ensure that staff interactions with members are professional and courteous. Regular maintenance of gym equipment is essential to prevent breakdowns and improve user experience. Hygiene standards should be reinforced by increasing cleaning frequency and enforcing stricter sanitation policies.

The analysis also highlighted the need for temperature control in gym showers, suggesting that adjustable thermostats be installed to maintain consistent water temperature. Membership fees should be reevaluated, with a consideration for introducing tiered membership options to cater to different customer needs. The gym atmosphere can also be enhanced by fostering a welcoming and motivating environment.

Technical issues related to the gym's app functionality were frequently mentioned, indicating a need for improvements to enhance usability. Outdoor facilities should be regularly maintained and upgraded to ensure their continued usability. Long wait times during peak hours can be addressed by improving the booking system and increasing staff presence.

Communication between the gym and its members should be improved, with clear and consistent updates provided regarding policy changes and facility updates. Changing rooms should be maintained to a higher standard to ensure cleanliness and comfort. Finally, specific locations such as Millhouses should be prioritized for facility upgrades to enhance overall user experience.

## Conclusion

Through data cleaning, sentiment analysis, and topic modeling, this study has identified key areas of concern for PureGym customers. Despite utilizing multiple modeling techniques, no additional insights emerged beyond the already identified pain points. The results consistently pointed to staff interactions, air conditioning issues, gym cleanliness, and access problems as the primary sources of customer dissatisfaction.