## Introduction:

In today's digitally interconnected world, understanding and managing customer sentiment is crucial for businesses striving to enhance customer satisfaction and optimize their processes. Natural Language Processing (NLP) offers a powerful toolkit for extracting insights from textual data, enabling organizations to delve deeper into customer feedback and predict the intensity of emotions expressed within it. This project aims to develop an intelligent system leveraging NLP techniques to predict the intensity of emotions, particularly focusing on happiness, anger, and sadness, expressed in text reviews.

The project encompasses several key phases, starting with data collection, where a diverse range of intensity-labelled text data will be gathered. Following data collection, the focus shifts to data preprocessing, where techniques such as cleaning, normalization, and transformation will be applied to ensure the data is suitable for subsequent analysis. Feature engineering will play a crucial role in extracting relevant features and identifying key variables that influence the intensity of emotions expressed in the text.

Model selection involves choosing appropriate machine learning algorithms capable of accurately classifying the intensity of emotions based on the processed textual data. Various algorithms will be explored and evaluated to determine the most effective approach for this specific task. Subsequently, the selected models will undergo rigorous training using the pre-processed data, aiming to optimize their performance.

Evaluation of the trained models will be conducted using appropriate metrics to assess their effectiveness in predicting emotion intensity. Hyperparameter tuning will be employed to fine-tune the models, enhancing their predictive accuracy and robustness. Finally, the trained model will be deployed in a production environment, enabling real-time predictions and proactive optimization of processes to improve overall customer satisfaction.

By developing an intelligent system capable of predicting emotion intensity in text reviews, businesses can gain valuable insights into customer sentiment, allowing them to make data-driven decisions to enhance customer experiences and drive organizational success.

## Data cleaning

In this project, we employed a combination of powerful tools and libraries, including **NLTK, regex, and string**, to effectively clean and preprocess the collected textual data. NLTK (Natural Language Toolkit) served as a fundamental component for various NLP tasks, providing robust functionality for tokenization, lemmatization, and stopwords removal. Additionally, regex (regular expressions) played a vital role in pattern matching and text manipulation, allowing us to efficiently handle tasks such as punctuation removal, URL extraction, and special character handling. Moreover. By leveraging these tools synergistically, we were able to preprocess the raw textual data effectively, ensuring that it was cleaned, standardized, and transformed into a format suitable for subsequent machine learning model training and analysis.

## Analysing data

In our project, we utilized word tokenization as a critical step in the data preprocessing pipeline. Word tokenization involves breaking down the text into individual words or tokens, enabling us to analyse and process each word independently. We leveraged NLTK's tokenization capabilities to split the text into meaningful units, taking into account punctuation, whitespace, and other linguistic features. This allowed us to effectively handle the textual data at a granular level, facilitating subsequent analysis and feature extraction.

Furthermore, data visualization played a crucial role in understanding the frequency and distribution of words within the text reviews. By visualizing word frequencies using techniques such as bar charts, word clouds, and histograms, we gained insights into which words were most commonly repeated across the dataset. This exploration helped us identify prominent themes, sentiments, and topics within the text reviews, providing valuable context for subsequent analysis. Moreover, visualizing word frequencies enabled us to detect any patterns or anomalies that could influence the intensity of emotions expressed in the text. Overall, word tokenization and data visualization served as essential techniques in our NLP workflow, facilitating a deeper understanding of the textual data and informing subsequent modelling decisions.

## Model Selection and Training

Firstly, we employed Count Vectorizer to convert text data into a matrix of token counts, representing the frequency of each word in the corpus. This approach allowed us to capture the frequency of individual words without considering their importance or rarity within the corpus. Additionally, we utilized the TF-IDF method, which considers not only the frequency of words but also their significance in the entire dataset by penalizing common words and amplifying rare ones.

Subsequently, we trained three different classification models on both the Count Vectorizer and TF-IDF representations of the data. The Multinomial Naive Bayes Classifier is a probabilistic model commonly used for text classification tasks, particularly when dealing with sparse data like text. Logistic Regression, on the other hand, is a linear model that estimates probabilities using a logistic function and is widely used in binary classification tasks. Lastly, Linear SVM is a powerful algorithm for text classification, known for its effectiveness in handling high-dimensional data and its ability to find optimal hyperplanes for separating different classes.

By experimenting with these combinations of text representation techniques and classification algorithms, we aimed to identify the most suitable approach for our task of predicting emotion intensity in text reviews. Through rigorous evaluation and comparison of model performance metrics, including accuracy, precision, recall, and F1-score, we sought to determine the optimal combination that would provide accurate and reliable predictions in real-world scenarios.

## Testing Analysis

The **ML Algorithms** used for prediction are listed as follows:

**Building models using different classifiers (Count vectorizer):**

Model 1: **Multinomial Naive Bayes Classifier** - Accuracy **70%**
Model 2: **Linear SVM** - Accuracy **82%**
Model 3: **Logistic Regression** - Accuracy **90%**

**Building models using different classifiers (TF-IDF vectorizer):**

Model 1: **Multinomial Naive Bayes Classifier** - Accuracy **48%**
Model 2: **Linear SVM** - Accuracy **54%**
Model 3: **Logistic Regression** - Accuracy **48%**

During the testing phase of our project, where we evaluated the performance of various machine learning models on   data, we discovered that Logistic Regression, combined with Count Vectorization, exhibited the highest accuracy of 90%. This finding suggests that the logistic regression model trained on the counts of words in the text data performed exceptionally well in predicting the intensity of emotions expressed in reviews.

The accuracy metric, which measures the proportion of correctly classified instances out of the total instances, serves as a valuable indicator of model performance. In our case, achieving an accuracy of 86% indicates that the logistic regression model with count vectorization effectively captured the underlying patterns and relationships in the text data, enabling accurate predictions of emotion intensity.

This result underscores the effectiveness of logistic regression in handling text classification tasks and highlights the utility of count vectorization as a text representation technique. By leveraging these methods, our model demonstrates promising capabilities in accurately predicting emotion intensity in text reviews, thereby empowering businesses to gain valuable insights into customer sentiment and enhance overall customer satisfaction.

# Prediction of intensity from sentences

|  | input_text | predicted_emotion |
|---|---|---|
| 0 | I am so angry at you!!!!! | Anger |
| 1 | you ve hit a new low with a danger of blm fascist slogan please stop it before too late stop | Anger |
| 2 | I love my doggg | Happy |
| 3 | I think i'm gonna be sick :'â€'( | Happy |
| 4 | I hate you so much | Happy |
| 5 | @TheTombert i was watching Harpers Island, lol... there was no vodka involved | Happy |
| 6 | sometimes i wish things could go back to the way they were the beginning of last summer | Sad |
| 7 | it's your 18th birthday finally!!! yippeeeee | Happy |
| 8 | oh no he is hospitalised!!! | Anger |

# Future Prospects

The successful implementation of the project opens up several promising future prospects and avenues for further development:

**Fine-tuning and Optimization**: Although achieving an 86% accuracy rate is commendable, there's always room for improvement. Future efforts could focus on fine-tuning model parameters, exploring different feature engineering techniques, and experimenting with advanced algorithms to enhance predictive performance further.

**Integration with Feedback Systems**: Integrating the developed model into existing feedback systems or customer relationship management (CRM) platforms could provide real-time insights into customer sentiment. This integration would enable businesses to promptly address issues, prioritize actions, and tailor responses to individual customers based on their expressed emotions.

**Sentiment Analysis in Multimodal Data**: Extending the project to handle multimodal data, including text, images, and audio, could provide a more comprehensive understanding of customer sentiment. By incorporating advanced techniques such as multimodal sentiment analysis, businesses can leverage diverse sources of customer feedback to gain deeper insights and make more informed decisions.

**Personalization and Recommendation Systems**: Leveraging the predictive capabilities of the model, businesses can develop personalized recommendation systems tailored to individual customer preferences and emotions. By analyzing historical data and real-time interactions, these systems can offer personalized product recommendations, content suggestions, and marketing messages, thereby enhancing customer engagement and satisfaction.

**Deployment in Various Industries**: The project's framework and methodology are applicable across a wide range of industries, including e-commerce, hospitality, healthcare, and finance. Tailoring the model to specific industry domains and use cases can unlock opportunities for improving customer experiences, optimizing business processes, and driving competitive advantage.

**Ethical Considerations and Bias Mitigation**: As with any AI-driven system, addressing ethical considerations and mitigating potential biases is paramount. Future research could focus on developing methods to ensure fairness, transparency, and accountability in the model's predictions, thereby promoting trust and confidence among users and stakeholders.

Overall, the future prospects of the project are promising, with opportunities to further advance the state-of-the-art in customer sentiment analysis, enhance decision-making processes, and ultimately, elevate the quality of customer experiences across various domains and industries.