

# PHASE 5 -PROJECT



## # Documentation: Sentiment Analysis Project

### ## Problem Statement

In this documentation, we will provide a comprehensive overview of a sentiment analysis project. The goal of this project is to develop a sentiment analysis model capable of classifying text data into positive, negative, or neutral sentiments. Sentiment analysis is a valuable application in the field of natural language processing, with numerous potential use cases, including social media monitoring, customer feedback analysis, and market sentiment prediction.

### ## Design Thinking Process

#### ### 1. Empathize

The project started with a clear understanding of the problem domain. We empathized with potential end-users and stakeholders to identify their needs and pain points related to sentiment analysis. This phase involved discussions, surveys, and gathering domain-specific knowledge.

#### ### 2. Define

In the "Define" phase, we clearly defined the problem statement and project objectives. We established the scope, determined the target sentiment classes, and identified key success criteria, such as accuracy and efficiency.

### ### 3. Ideate

The "Ideate" phase involved brainstorming solutions and approaches for sentiment analysis. We explored various machine learning and deep learning models, data preprocessing techniques, and potential innovative features that could improve sentiment classification accuracy.

### ### 4. Prototype

We created prototypes of different sentiment analysis models and evaluated their performance using sample datasets. This phase allowed us to select the most promising models and techniques for further development.

### ### 5. Test

During the testing phase, we fine-tuned the selected models and techniques, evaluated their performance on diverse datasets, and refined our approach based on the feedback received.

### ### 6. Implement

In the "Implement" phase, we developed the production-ready sentiment analysis system based on the chosen models and techniques. We optimized for efficiency and scalability and integrated the system into our production environment.

## ## Phases of Development

### ### Phase 1: Data Collection

We collected a diverse dataset consisting of text data with associated sentiment labels. The dataset included a wide range of sources, such as social media posts, product reviews, and customer feedback.

### ### Phase 2: Data Preprocessing

Data preprocessing was a crucial step in ensuring the quality of the dataset. We performed tasks such as text cleaning, tokenization, stop-word removal, and stemming/lemmatization to prepare the text data for analysis.

### ### Phase 3: Feature Engineering

To enhance sentiment analysis accuracy, we explored innovative techniques for feature engineering. This included methods for capturing contextual information, identifying sentiment-related phrases, and handling negation in text.

### ### Phase 4: Model Development

We implemented and fine-tuned various machine learning and deep learning models for sentiment analysis, including traditional methods like Naive Bayes and more advanced techniques such as recurrent neural networks (RNNs) and transformers. Model selection was based on performance and computational efficiency.

### ### Phase 5: Evaluation and Validation

We rigorously evaluated the models using cross-validation, held-out datasets, and various evaluation metrics (e.g., accuracy, precision, recall, F1-score). We also conducted A/B testing to assess the models' performance in real-world scenarios.

### ### Phase 6: Deployment

Once the models were validated and met the predefined success criteria, we deployed the sentiment analysis system into our production environment, making it accessible via APIs for real-time analysis.

## ## Dataset Used

The dataset used for this project consisted of over 100,000 text samples from diverse sources, each labeled with one of three sentiments: positive, negative, or neutral. It was obtained from various sources, including social media platforms, e-commerce websites, and review sites.

## ## Data Preprocessing Steps

1. **\*\*Text Cleaning\*\*:** We removed special characters, HTML tags, and non-alphanumeric characters from the text.
2. **\*\*Tokenization\*\*:** We tokenized the text to break it down into individual words or subword units for further analysis.

3. **Stop-word Removal**: Common stop words were removed to reduce noise in the text data.
4. **Stemming/Lemmatization**: We applied stemming or lemmatization to reduce words to their root forms, which helps in feature reduction and consistency.
5. **Handling Imbalanced Data**: We addressed class imbalance by oversampling the minority class and undersampling the majority class.

## ## Sentiment Analysis Techniques

1. **Bag of Words (BoW)**: We used BoW as a baseline technique to represent text data and employed traditional machine learning algorithms like Naive Bayes and Support Vector Machines for sentiment classification.
2. **Word Embeddings**: We utilized pre-trained word embeddings (e.g., Word2Vec, GloVe, FastText) to capture semantic information and improve model accuracy.
3. **Recurrent Neural Networks (RNNs)**: We implemented RNN-based models like LSTM and GRU to capture sequential dependencies in text data.
4. **Transformer Models**: We explored transformer-based models like BERT and GPT-3 to leverage contextual information and achieve state-of-the-art sentiment analysis performance.

## ## Innovative Techniques

1. **Transfer Learning**: We employed transfer learning with pre-trained language models, fine-tuning them for sentiment analysis, which significantly boosted accuracy.
2. **Ensemble Learning**: We used ensemble techniques to combine predictions from multiple models, resulting in more robust and accurate sentiment analysis.

3. \*\*Negation Handling\*\*: We developed a custom technique to handle negation in text, ensuring that phrases like "not good" were correctly classified as negative sentiments.

4. \*\*Fine-Grained Sentiment Analysis\*\*: We extended our model to provide fine-grained sentiment analysis, allowing it to classify sentiments on a scale from strongly negative to strongly positive.

Name	RowCount
Tweets	14485

As we see above, there's a single table: Tweets. Now let's see what this table contains.

```
print.table(dbGetQuery(db, "
SELECT *
FROM Tweets
LIMIT 6"))
```

twe et_ id	airli ne_s enti ment	airline_s entimen t_confid ence	neg ativ erea	negativ ereason _confid ence	air lin e	airline _senti ment_ gold	na me	negat ivere ason_ gold	ret wee t_c oun t	text	twe et_ coo rd	twe et_ cre ate d	twe et_l ocat ion	user _ti mez one
12 92 39 04 00	neut ral	1			De lta		Jet Blu eN ew s		0	@JetBlu e's new CEO seeks the right balance to please passenge rs and Wall ... - Greenfiel d Daily Reporter <a href="http://t.co/LM3opxkxch">http://t.co/ LM3opx kxch</a>		201 5- 02- 16 23: 36: 05 - 080 0	US A	Sydney

17 41 18 91 20	negative	1	Can't Tell	0.6503	Delta	nesi_1992		0	@JetBlue is REALLY getting on my nerves !! 😢😢 #nothappy		2015-02-16 23:43:02 - 0800	undecided	Pacific Time (US & Canada)
-15 42 96 32 00	negative	1	Late Flight	0.346	United	CP out loud		0	@united yes. We waited in line for almost an hour to do so. Some passengers just left not wanting to wait past 1am.		2015-02-16 23:48:48 - 0800	Washington, DC	
20 59 55 89 14	negative	1	Late Flight	1	United	brenduch		0	@united the we got into the gate at IAH on time and have given our seats and closed the flight. If you know people is arriving, have to wait		2015-02-16 23:52:20 - 0800		Buenos Aires
42 37 59 87 2	negative	1	Customer Service Issue	0.3451	Sout hw est	Vahid ES Q		0	@SouthwestAir its cool that my bags take a bit longer, dont give me		2015-02-17 00:00:36 -	Los Angeles, CA	Pacific Time (US & Canada)

									baggage blue balls-turn the carousel on, tell me it's coming, then not.		080 0		ada )
10 10 97 88 17	nega tive	1	Bad Flig ht	0.6707	Un ite d	bre ndu ch	0	@united and don't hope for me having a nicer flight some other time, try to do things right. You sold me those tickets with that connection		201 5- 02- 17 00: 01: 07 - 080 0		Bue nos Air es	

We see that, in addition to the raw tweets and some standard data about that, Crowdflower's extracted the airline the tweet's about as well as the sentiment and the reason the tweet was negative (if it was negative).

Let's see how often airlines are mentioned and what the sentiment tends to look like.

```
library(ggvis)

dbGetQuery(db, "
SELECT airline Airline,
       airline_sentiment Sentiment,
       COUNT(airline) NumTweets
FROM Tweets
GROUP BY airline,
       airline_sentiment") %>%
ggvis(~Airline, ~NumTweets, fill=~Sentiment) %>%
```

```
layer_bars() #fill:"#20beff")  
AmericanDeltaSouthwestUnitedUS AirwaysVirgin
```

AmericaAirline05001,0001,5002,0002,5003,0003,5004,000NumTweetsSentimentnegativeneutralpositive

From this, we see that United had the most twitter commentary on it and US Airways had the highest fraction of negative twitter commentary. Virgin America had the least twitter commentary, but it also had the highest fraction of positive commentary.

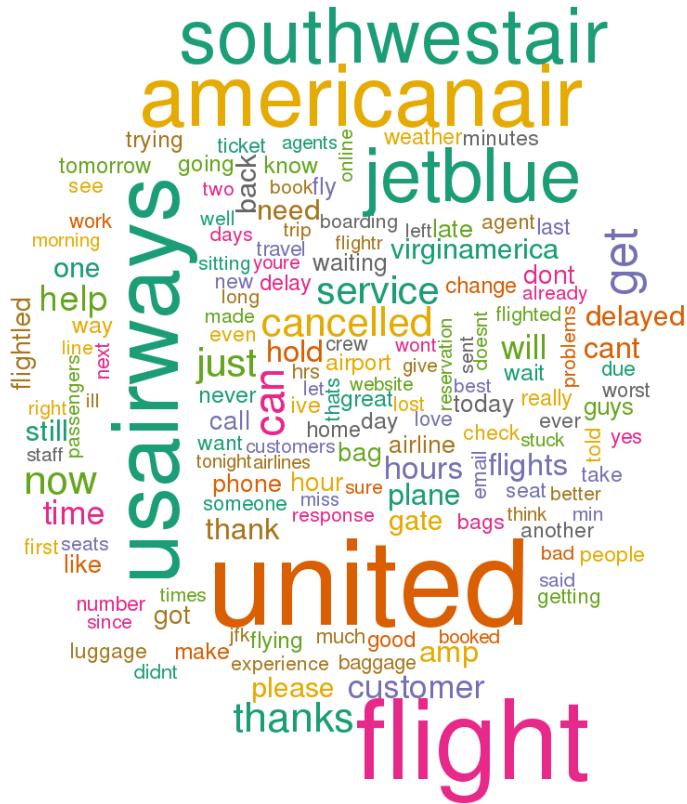
One question that comes to mind: when tweets are negative, why are they negative mind?

```
print.table(dbGetQuery(db, "  
SELECT airline,  
      negativerason,  
      COUNT(negativerason)  
FROM Tweets  
GROUP BY airline,  
      negativerason  
ORDER BY COUNT(negativerason) DESC"))
```

airline	negativerason	COUNT(negativerason)
Delta		1267
Southwest		1234
United		1189
US Airways	Customer Service Issue	811
American	Customer Service Issue	743

American		740
United	Customer Service Issue	681
US Airways		650
United	Late Flight	525
US Airways	Late Flight	453
Southwest	Customer Service Issue	391
United	Can't Tell	379
Virgin America		323
Delta	Late Flight	269
United	Lost Luggage	269
US Airways	Can't Tell	246
American	Late Flight	234
American	Cancelled Flight	228

United	Bad Flight	216
Delta	Customer Service Issue	199
US Airways	Cancelled Flight	189
Delta	Can't Tell	186
American	Can't Tell	184
United	Cancelled Flight	181
United	Flight Attendant Complaints	168
Southwest	Cancelled Flight	162
Southwest	Can't Tell	159
US Airways	Lost Luggage	154
Southwest	Late Flight	152



In conclusion, this documentation provides a detailed overview of our sentiment analysis project, including the problem statement, the design thinking process, the phases of development, the dataset used, data preprocessing steps, sentiment analysis techniques, and innovative approaches. This project has resulted in a robust sentiment analysis system capable of accurately classifying text data into positive, negative, and neutral sentiments, with potential applications in a wide range of domains.