TAWI

A MORINGA SCHOOL CHATBOT









Business Understanding



Data Understanding



Data Preparation



EDA



Modelling



Recommendations







Business Problem

Chatbot Objectives

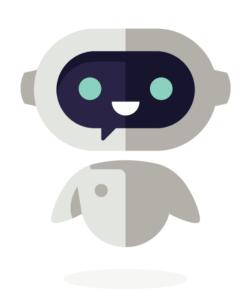


BUSINESS PROBLEM

A chatbot is a program whose aim is to simulate conversation with human users.

Research suggests that by the end of 2023, more than 75% of customer queries will be resolved by chatbots.

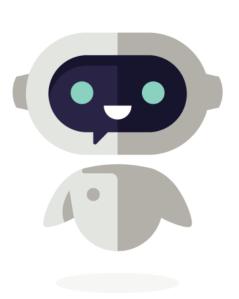
Moringa is one of the best tech schools in Kenya and as it continues to grow, so will the demand for its services. An efficient way to fill this gap would be to employ the use of a chatbot.





OBJECTIVES

- Provide fast, 24/7 Service
- •Improve Customer Experience
- Provide Access to Information
- Offer Technical Support
- Data Collection



link_scraper.py contains a function designed to go through the Moringa School website and extract hyperlinks present on its pages

DATA UNDERSTANDING

web_scraper.py contains a function that utilizes the links obtained from (O1) to scrape text content from the corresponding web pages

scraped_links,json
contains a list of unique URLs
obtained during the link scraping
process.

moringa_text_corpus.json stores
the text content scraped from
the corresponding URLs,
organized in a dictionary
structure



DATA PREPARATION

The goal is to create an intent list, which helps the chatbot understand the user's goals. The intent list is stored in "final_intents.json" with the following structure:

"tag": The category of the intent. Eg "software_engineering_ course information"

"questions": A list of questions.

Eg

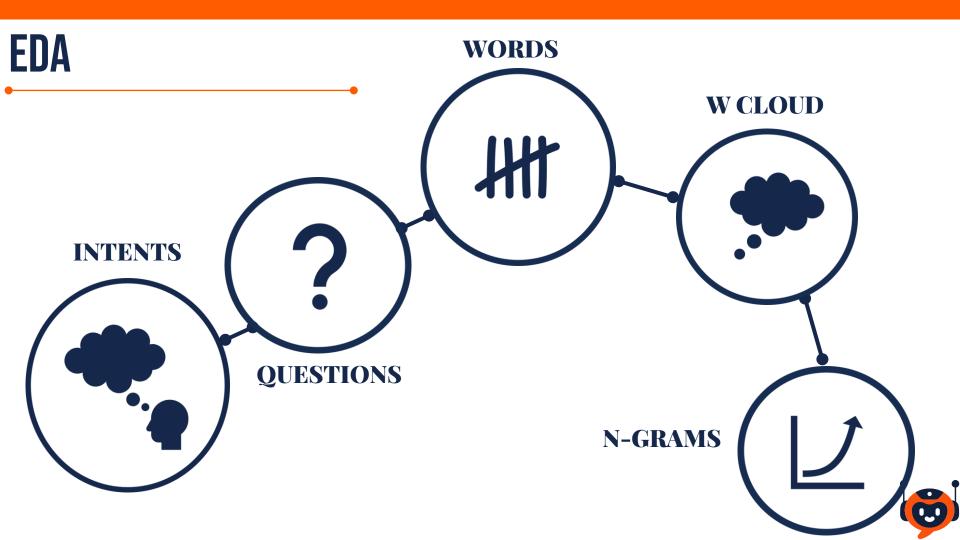
"What topics does the software engineering course cover?"

"responses": A list of responses

Eg

"The software engineering course covers HTML and CSS".

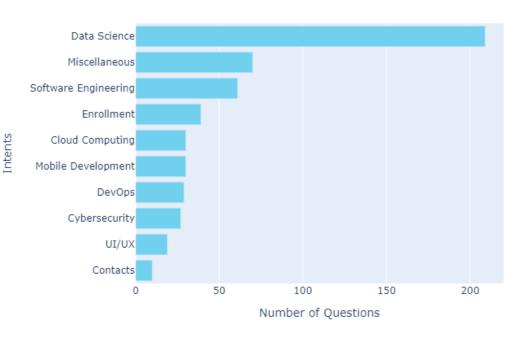




INTENTS



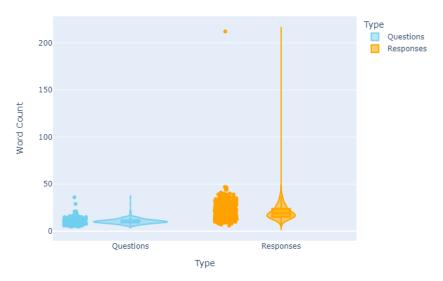
Intents by Number of Questions



- The dataset comprises 10 unique intents, with an average of approximately 52.40 questions per intent.
- The median number of questions per intent is 30, indicating a balanced distribution.
- Notably, the "Contacts" intent stands out with the fewest questions, totaling 10, suggesting it represents a category with relatively fewer questions compared to others in the dataset.



Distribution of Word Counts in Questions and Responses

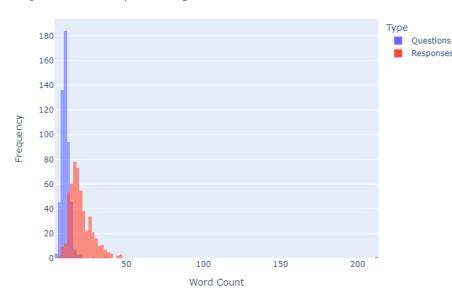


This is a plot showing the distribution of word counts in questions and responses.

Questions typically have fewer words, while responses vary more in length.

6

Question and Response Lengths



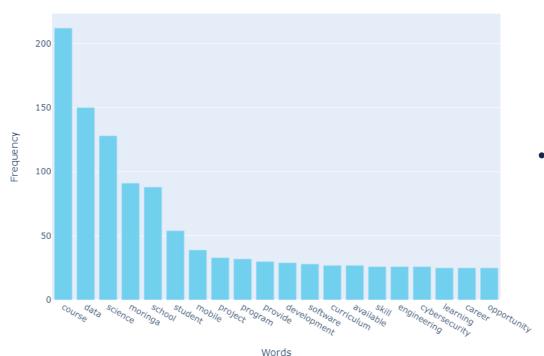
On average, questions are relatively short (around 10 words), with a moderate level of variability.

Responses are longer on average (around 21 words), with greater variability in their lengths.





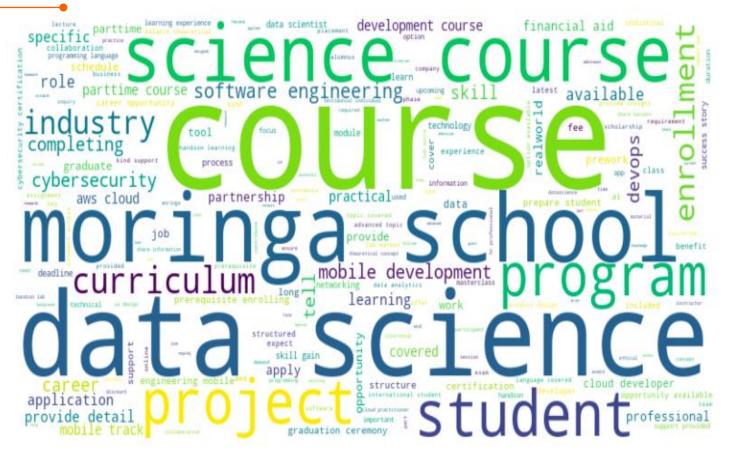
Top 20 Most Common Words in Questions



- The analysis highlights key themes in
 the dataset, emphasizing courses like data science, the school Moringa, and others like education
- 'course,'
 'data,' 'science,' 'moringa,' 'student,'
 'develop,' and 'learn' indicates a focus
 on data science education with a
 learner-centric approach.

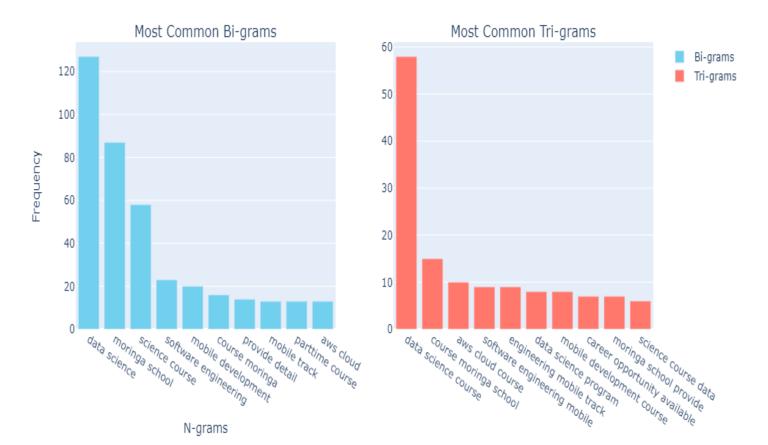
WORD CLOUD





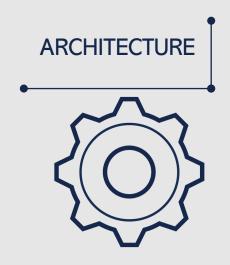
N-GRAMS

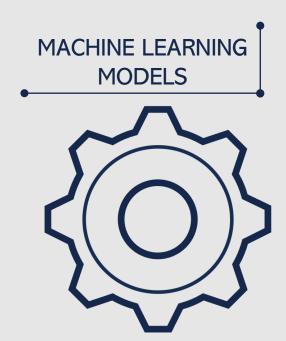
Most Common N-grams

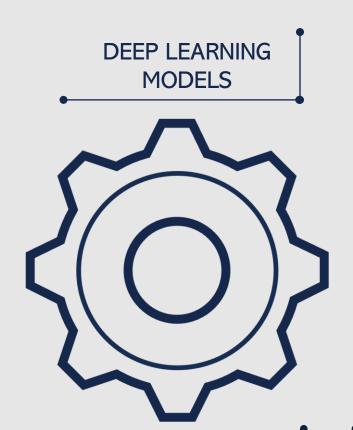




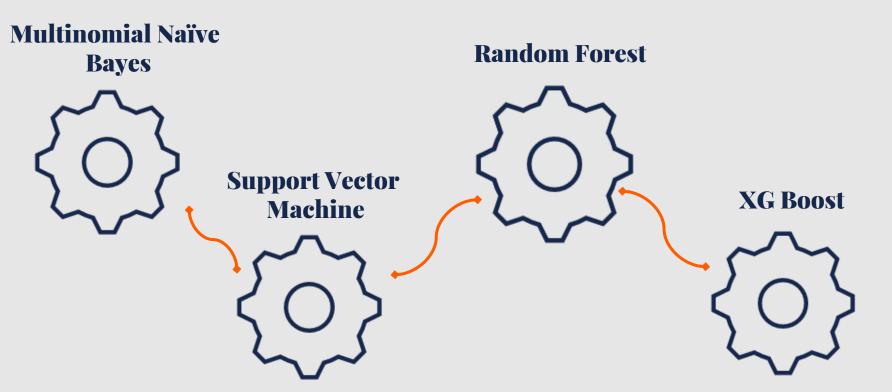
MODELLING







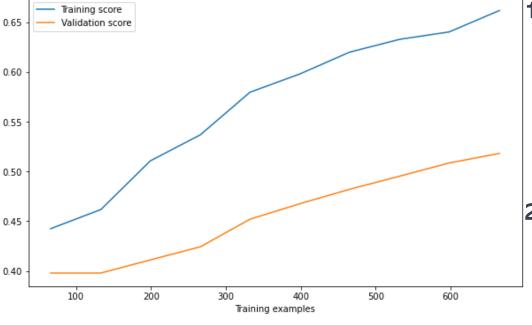
MACHINE LEARNING MODELS





MULTINOMIAL NAÏVE BAYES





1. Multinomial Naive Bayes:

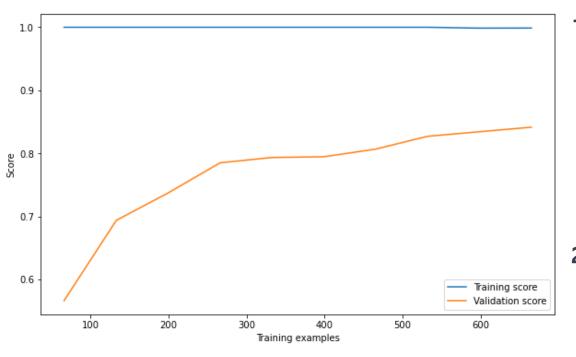
- 1. Accuracy: 78%
- 2. Notable variations in precision and recall among classes.
- 3. Adjustments suggested for classes with low precision or recall.

2.GridSearchCV Multinomial Naive Bayes:

- 1. Overall accuracy: 73.33%
- 2. Notable variations in precision and recall, emphasizing the need for adjustments.

SUPPORT VECTOR MACHINE





1.TF-IDF Vectorizer and SVM:

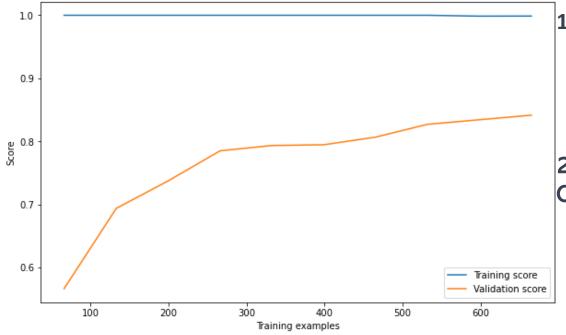
- 1. Accuracy: 81.9%
- 2. Varied precision and recall values across classes.
- 3. Consideration of hyperparameter tuning for enhanced performance.

2.Cross-validation with TF-IDF + SVM:

- 1. Average accuracy: 80.2%
- 2. Challenges observed in precision and recall for certain classes.

RANDOM FOREST





1. Random Forest Classifier:

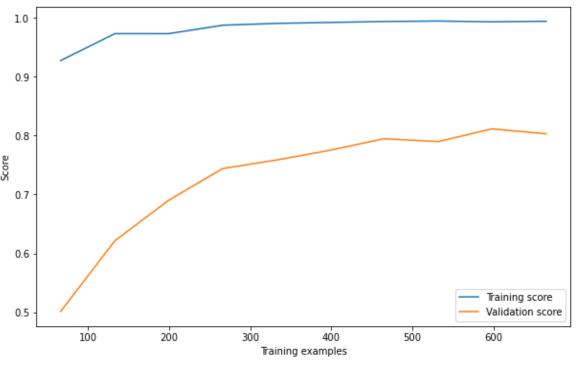
- 1. Accuracy: 81%
- 2. Challenges in precision and recall for specific classes.

2.GridSearchCV for Random Forest Classifier:

- 1. Average cross-validation accuracy: 81.87%
- 2. Class-specific insights provided for refinement.

XG BOOST



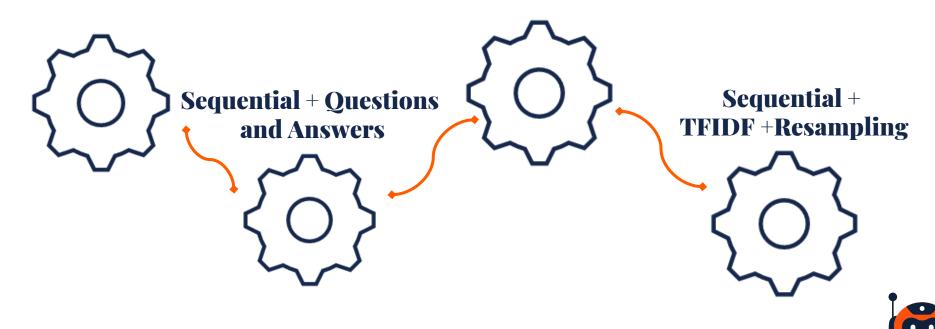


- 1. Accuracy: 82.86%
- 2. Varied performance across classes, suggesting further tuning or additional data.

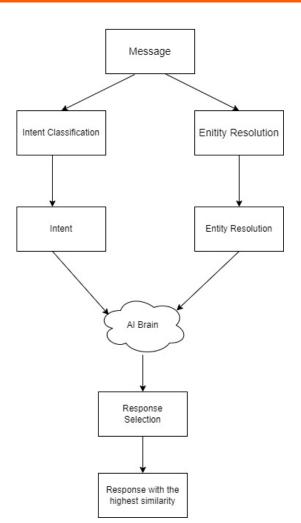
DEEP LEARNING MODELS

Base – Sequential Model

Sequential + **Regularization**



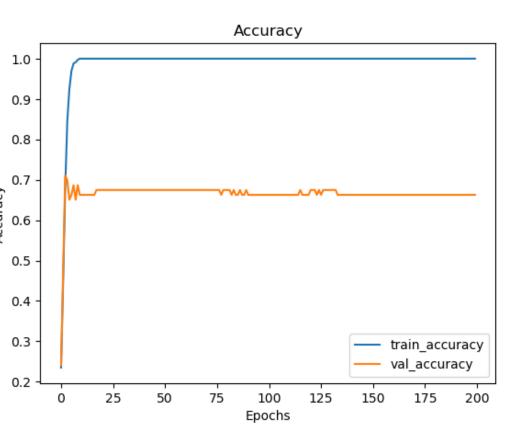
ARCHITECTURE





BASE - SEQUENTIAL



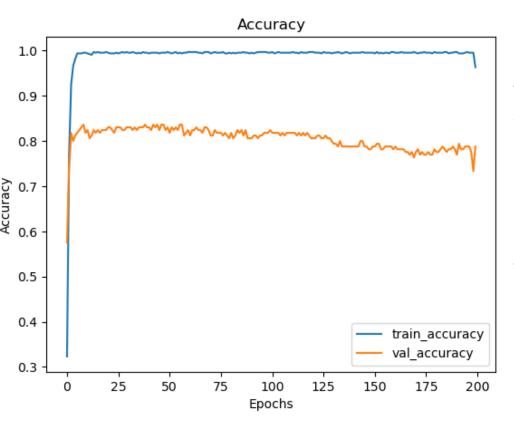


The base model demonstrates a high training accuracy of approximately 98%, suggesting effective learning from the training data.

However, a significant performance gap is observed, with the validation accuracy plateauing at 30-40%, raising concerns about overfitting

SEQUENTIAL + QNA



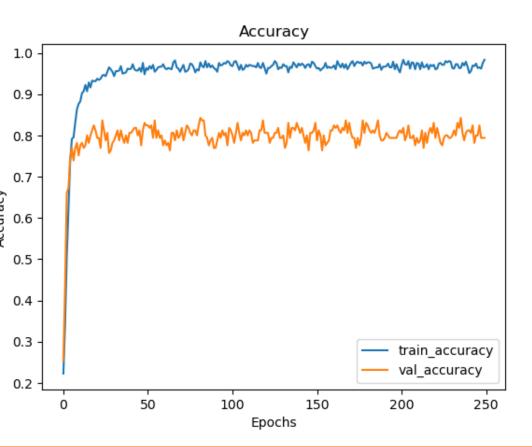


The model exhibits notable improvement, with training accuracy at approximately 99.42% and validation accuracy at around 58.78%.

Despite this progress, there is potential for further enhancement in the validation accuracy.

SEQUENTIAL + REGULARIZATION

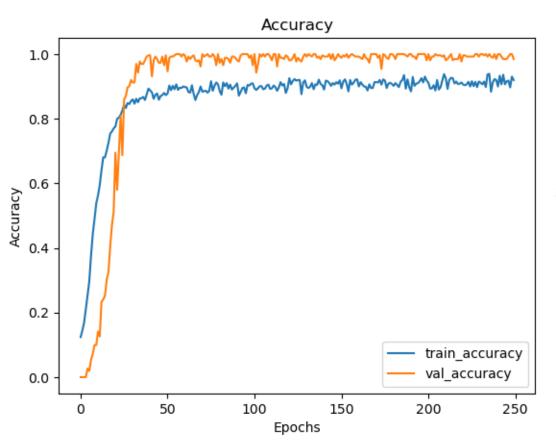




The gap between training and validation scores is not drastic, indicating that the regularization and dropout are having a positive effect.

SEQUENTIAL + TFIDF, RESAMPLING





The model achieved high accuracy on the training set, with the training accuracy steadily increasing and reaching nearly 1.0

Further steps may include using more advanced models

CONCLUSION & RECOMMENDATIONS

1.Intent Classification:

The intent classification model achieved good performance, accurately categorizing user queries into predefined intent categories. This model can effectively route user queries to appropriate response handlers or workflows.

2. Response Generation:

The response generation system, utilizing cosine similarity with TF-IDF vectors, showed promising results in mapping user queries to appropriate responses. However, there is room for improvement, particularly in handling context and generating diverse responses.

3. Continuous Learning and Adaptation:

Implement mechanisms for continuous learning and adaptation, allowing the chatbot to evolve over time based on user feedback, new data, and changing requirements.

Monitor performance metrics and user satisfaction to measure the effectiveness

Interact with TAWI at: https://tawi-chatbot.onrender.com/

THANKYOU!

