

Language Translation Project

Machine Learning Programming - PROG8245

Tessa Nejla Ayvazoglu &
Satyam Patel 15/04/2024



Agenda



Introduction

Project Overview

Web Scraping and Dataset Preparation

Dataset Loading, Cleaning, and Annotation

Text Preprocessing

Feature Extraction

Model Building

Model Training and Evaluation

Deployment and Interface

Conclusion

Introduction

01

Language translation involves converting text or speech from one language to another while maintaining meaning and context. It's essential for cross-cultural communication and understanding in today's interconnected world.

02

Machine learning has transformed translation by enabling computers to learn from data and improve accuracy over time. It allows for more natural translations, scalability, and faster speeds compared to traditional methods.

03

This project aims to explore the application of machine learning, particularly neural networks, in language translation. We'll train models using large datasets to achieve accurate and fluent translations. Our goal is to evaluate different neural network architectures and optimization techniques to enhance translation quality and contribute to advancing translation technology.

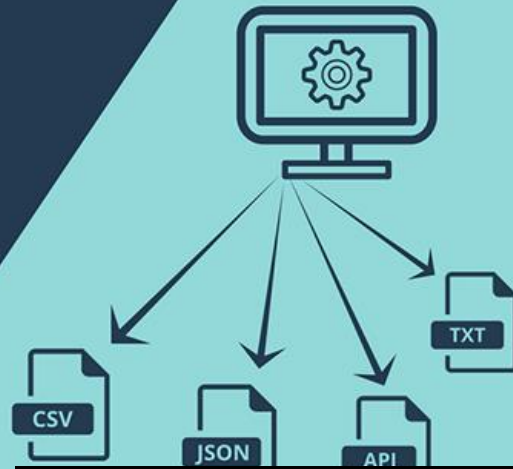
Project Overview

Objectives

- To scrape French-English sentence pairs and prepare a dataset.
- To preprocess and prepare the data for training a translation model.
- To train and evaluate a machine translation model.
- To deploy the model in a simple web interface for real-time translation.

What is Web Scraping?

The beauty of web scraping is that you can fetch images, URLs, email addresses, phone numbers and other text data like key-value pairs and paragraphs.



This Photo by Unknown author is licensed under CC BY-SA-NC.

Web Scraping and Dataset Preparation

Web Scraping and Dataset Preparation

Objective: Gather large datasets for training translation model.

Methods: Utilize web scraping with Python, BeautifulSoup, and requests.

Ethical Considerations: Comply with website terms of service.

Data Integrity: Normalize, remove duplicates, correct misalignments.

Text Processing: Tokenization, stemming, lemmatization with NLTK.

Challenges: Dynamic content, rate limiting, ensuring data quality.

Conclusion: Effective web scraping and dataset preparation are crucial for training accurate translation models, requiring attention to ethical considerations and robust data processing techniques.

Dataset Loading, Cleaning, and Annotation

OBJECTIVE: PREPARE THE DATASET FOR TRAINING THE TRANSLATION MODEL.

LOADING: IMPORT DATA INTO THE PROJECT ENVIRONMENT.

CLEANING: REMOVE IRRELEVANT INFORMATION, HANDLE MISSING VALUES.

ANNOTATION: LABEL DATA FOR SUPERVISED LEARNING IF NECESSARY.

QUALITY CHECK: ENSURE DATA INTEGRITY AND CONSISTENCY.

CONCLUSION: PROPER DATASET LOADING, CLEANING, AND ANNOTATION ARE ESSENTIAL FOR BUILDING RELIABLE TRANSLATION MODELS, LAYING THE FOUNDATION FOR ACCURATE TRAINING AND EVALUATION.

Sentiment Annotation

Using a pre-trained sentiment analysis model from Hugging Face, we annotate each English sentence with a sentiment label (positive, negative, neutral) based on the emotional context of the sentences, which is critical in translation for maintaining the original tone.

Code smells | Comments | Unit tests

```
import os      `os` imported but unused
from transformers import pipeline    Import "transformers" could not be resolved
```

Function to annotate sentiment using a pre-trained model

Type Hints | Code Smells | Generate Docstring | Unit tests

```
def annotate_sentiment(data):
```

```
    classifier = pipeline('sentiment-analysis', model="distilbert/distilbert-base-uncased-finetuned-sst-2-english")    "distilbert": Unknown word.
```

Split the data into batches

```
batch_size = 200 # Adjust based on your GPU memory
```

```
batches = [data[i:i + batch_size] for i in range(0, len(data), batch_size)]
```

```
sentiments = []
```

```
for batch in batches:
```

```
    batch_predictions = classifier(batch)
```

```
    batch_sentiments = [prediction['label'] for prediction in batch_predictions]
```

```
    sentiments.extend(batch_sentiments)
```

```
return sentiments
```

Annotate the dataset

```
sentiments = annotate_sentiment(dataframe['English'].tolist())    "dataframe": Unknown word.
```

```
dataframe['Sentiment'] = sentiments    "dataframe": Unknown word.
```

```
print(dataframe.head())    "dataframe": Unknown word.
```

Dataset Loading, Cleaning, and Annotation

Text Preprocessing

Objective: Prepare textual data for machine learning (ML) models using natural language processing (NLP) techniques.

Key Steps:

- **Tokenization:** Split text into individual words or tokens.
- **Normalization:** Standardize text by converting to lowercase and removing punctuation.
- **Stopword Removal:** Eliminate common words (e.g., "the," "is") that add little semantic value.
- **Stemming/Lemmatization:** Reduce words to their base or root form to capture semantic similarity.
- **Entity Recognition:** Identify and label named entities (e.g., names, locations) for better understanding.

Benefits:

- Improves model performance by reducing noise and irrelevant information.
- Enhances interpretability and generalization of ML models.

Challenges:

- Handling domain-specific terms or slang.
- Balancing between preserving meaning and reducing dimensionality.

Conclusion: Text preprocessing is a critical step in NLP ML pipelines, enabling models to effectively learn from textual data and make accurate predictions.

This section covers the essential text preprocessing steps which include:

- **Tokenization:** Splitting text into meaningful elements called tokens.
- **Lowercasing:** Normalizing all text to lowercase to ensure uniformity.
- **Removing Stop Words:** Eliminating common words that may not be useful in the analysis.
- **Lemmatization:** Reducing words to their base or root form.

Effective preprocessing is crucial for reducing noise and improving the quality of feature extraction.

Code smells | Comments | Unit tests

```
import nltk      Import "nltk" could not be resolved
nltk.download('punkt')      "punkt": Unknown word.
nltk.download('stopwords')   "stopwords": Unknown word.
nltk.download('wordnet')     "wordnet": Unknown word.
from nltk.corpus import stopwords      Module level import not at top of cell
from nltk.stem import WordNetLemmatizer      Module level import not at top of cell
from nltk.tokenize import word_tokenize      Module level import not at top of cell
```

Text preprocessing function

Type Hints | Code Smells | Generate Docstring | Unit tests

```
def preprocess_text(text):
    tokens = word_tokenize(text.lower())
    tokens = [token for token in tokens if token.isalpha() and token not in stopwords.words('english')]
    lemmatizer = WordNetLemmatizer()      "lemmatizer": Unknown word.
    lemmatized = [lemmatizer.lemmatize(token) for token in tokens]      "lemmatized": Unknown word.
    return ' '.join(lemmatized)      "lemmatized": Unknown word.
```

Apply preprocessing

```
dataframe['English_clean'] = dataframe['English'].apply(preprocess_text)      "dataframe": Unknown word.
dataframe['French_clean'] = dataframe['French'].apply(preprocess_text)      "dataframe": Unknown word.
print(dataframe.head())      "dataframe": Unknown word.
dataframe.to_csv('./cleaned_data.txt', sep='\t', index=False)
```


Text Preprocessing

Tokenization

Lowercasing

Stop words removal

Lemmatization

Discuss how
preprocessing
improves the quality
of feature extraction



What is TF-IDF?

TF-IDF, short for Term Frequency-Inverse Document Frequency, is a technique used in natural language processing (NLP) to transform text data into a numerical format suitable for machine learning models.

Explaining TF-IDF



Term Frequency (TF):

Term Frequency measures the frequency of a term (word) within a document. It reflects how often a word appears in a document relative to the total number of words in that document.



Inverse Document Frequency (IDF):

IDF highlights the importance of a term across a collection of documents (corpus). It emphasizes terms that are rare across all documents in the corpus. IDF is calculated by taking the logarithm of the total number of documents divided by the number of documents containing the term.

Discussing Word Importance with TF-IDF



Highlighting Word Importance:

TF-IDF assigns higher weights to words that are frequent within a document but rare across the corpus.

Words with high TF-IDF scores are considered important features as they carry significant meaning within the context of a document.

Discussing the Role of Feature Extraction



Feature Extraction in NLP:

Feature extraction is the process of converting raw data into a format suitable for machine learning models.

In the context of text data, feature extraction involves transforming words or sentences into numerical features that ML algorithms can understand and process.



Importance of TF-IDF in Feature Extraction:

TF-IDF is a crucial feature extraction technique in NLP ML pipelines.

It captures the importance of words within individual documents and across the entire corpus, enabling ML models to effectively learn from textual data.

Feature Extraction

Exploring Feature Extraction

In this section, we apply the TF-IDF technique to transform the text into a numeric format that machine learning models can understand. TF-IDF highlights the importance of each word in the context of its document, balancing out the frequency of occurrence across the dataset.

[+ Code](#) [+ Markdown](#)

Code smells | Comments | Unit tests

```
from sklearn.feature_extraction.text import TfidfVectorizer
```

Import "sklearn.feature_extraction.text" could not be resolved from source

```
# TF-IDF Vectorizer
vectorizer = TfidfVectorizer()
features = vectorizer.fit_transform(dataframe['English_clean'])
print(features.shape)
```

Python

(232736, 13668)

Prepare the Environment

Code smells | Comments | Unit tests

```
# Ensure that TensorFlow is running with GPU support
```

```
import tensorflow as tf
```

Import "tensorflow" could not be resolved

```
print("Num GPUs Available: ", len(tf.config.list_physical_devices('GPU')))
if tf.test.gpu_device_name():
    print('Default GPU Device: {}'.format(tf.test.gpu_device_name()))
else:
    print("Please install GPU version of TF")
```

Python

Feature Extraction

Model Building

The translation model is structured upon a sequence-to-sequence architecture, comprising an encoder and a decoder. This framework allows for the processing of input sequences, in this case, French text, and the generation of output sequences, which are the English translations.

LSTM (Long Short-Term Memory) layers play a fundamental role in processing sequences within the model. These layers were specifically chosen due to their capability to effectively handle long sequences and maintain contextual information over time. By incorporating LSTM layers, the model can capture the sequential dependencies inherent in natural language data, facilitating accurate translation.

In the translation process from French to English, the encoder phase initially encodes the input French sequence using the LSTM layers. This encoding generates a fixed-length context vector that encapsulates the input information. Subsequently, the decoder utilizes this context vector to generate the English translation, token by token. At each decoding step, the decoder LSTM layer processes the previously generated token alongside the context vector to predict the subsequent token in the output sequence.

Layer (type)	Output Shape	Param #	Connected to
input_layer (InputLayer)	(None, None)	0	-
input_layer_1 (InputLayer)	(None, None)	0	-
embedding (Embedding)	(None, None, 256)	2,560,256	input_layer[0][0]
not_equal (NotEqual)	(None, None)	0	input_layer[0][0]
embedding_1 (Embedding)	(None, None, 256)	2,560,256	input_layer_1[0]...
lstm (LSTM)	[(None, 256), (None, 256), (None, 256)]	525,312	embedding[0][0], not_equal[0][0]
lstm_1 (LSTM)	[(None, None, 256), (None, 256), (None, 256)]	525,312	embedding_1[0][0... lstm[0][1], lstm[0][2]
dense (Dense)	(None, None, 10001)	2,570,257	lstm_1[0][0]

Total params: 8,741,393 (33.35 MB)

Trainable params: 8,741,393 (33.35 MB)

Model Training and Evaluation

The model was trained using GPU acceleration to expedite the training process. By leveraging Google Cloud's T4 GPUs, computational tasks were significantly accelerated, allowing for faster convergence and reduced training times. GPU acceleration enabled the model to process large volumes of data efficiently, thereby enhancing its learning capabilities and improving overall training performance.

Following training, the model was evaluated on unseen data to assess its translation accuracy and performance. This evaluation aimed to validate the model's generalization capabilities and its ability to produce accurate translations on real-world data. Unseen data sets were carefully selected to represent diverse linguistic patterns and usage scenarios, ensuring robust evaluation across different language domains.

The evaluation process employed metrics such as BLEU scores to quantify translation quality against a set of standard translations. BLEU scores provide a measure of similarity between the model's translations and human-generated references, offering insights into the model's accuracy and fluency. Additionally, other accuracy metrics, including precision, recall, and F1-score, were computed to comprehensively assess the model's performance across various linguistic dimensions.

Evaluation

Evaluate the model on unseen data

Code smells | Comments | [Unit tests](#)

```
# Fill NaN values with a placeholder string and ensure all entries are strings
test_df['French_clean'] = test_df['French_clean'].fillna('').astype(str)    "fillna":
test_df['English_clean'] = test_df['English_clean'].fillna('').astype(str)    "fillna":
```

```
# Now prepare the sequences
```

```
test_french_seq = pad_sequences(french_tokenizer.texts_to_sequences(test_df['French_clean']),
                               maxlen=train_french_seq.shape[1], padding='post')
test_english_seq = pad_sequences(english_tokenizer.texts_to_sequences(test_df['English_clean']),
                                maxlen=train_english_seq.shape[1], padding='post')
```

```
# Prepare the decoder input for the test data
```

```
test_decoder_input = np.roll(test_english_seq, shift=1, axis=1)
```

```
test_decoder_input[:, 0] = 0 # Assuming '0' is used for the 'start of sequence' token
```

```
# Evaluate the model on the prepared test data
```

```
results = model.evaluate([test_french_seq, test_decoder_input], test_english_seq)
```

```
print("Test Loss:", results[0])
```

```
print("Test Accuracy:", results[1])
```

5/1455 ————— 57s 39ms/step - accuracy: 0.7631 - loss: 2.4667

t Loss: 2.452239751815796

t Accuracy: 0.7626456618309021

Deployment and Interface

A simple graphical user interface (GUI) was developed using Tkinter to facilitate user interaction with the translation model. Tkinter, a standard Python library for creating GUI applications, was chosen for its ease of use and cross-platform compatibility. The GUI provided a user-friendly interface for inputting French text and receiving real-time translations.

The GUI allowed users to input French text into a designated text entry field. Upon input, the translation model processed the text and generated real-time translations, which were displayed in a separate output field within the interface. Users could easily view the translated text and interact with the model without the need for complex commands or technical expertise.

The developed GUI and translation model have practical applications across various domains, including language learning, communication, and international business. Users can leverage the model to translate French text into English accurately and efficiently, enabling seamless communication between individuals who speak different languages. Moreover, the user-friendly interface enhances accessibility, making the translation model accessible to a wider audience with varying levels of technical proficiency.

Conclusion



Project Summary:

Data Handling: Streamlined data collection, cleaning, and preprocessing to train the model.

Model Architecture: Employed a sequence-to-sequence framework with LSTM networks for effective handling of long sequences.

Training & Evaluation: Leveraged GPU acceleration for efficient training and utilized BLEU scores for rigorous performance evaluation.

Deployment: Developed a user-friendly GUI using Tkinter for real-time French to English translations.



Achievements:

Successfully built and deployed a robust translation model.

Enhanced user interaction through a simple yet functional graphical interface.



Challenges Overcome:

Addressed significant hurdles in data preprocessing, model optimization, and GUI scalability.

Innovated solutions to maintain high translation accuracy under operational constraints.



Future Directions:

Expand Linguistic Coverage: Adapt the model to include more languages.

Enhance Model Accuracy: Integrate advanced neural network architectures like transformers.

Improve User Experience: Upgrade the GUI with additional features like auto-detection of input language and audio output capabilities.



Q&A

OPEN THE FLOOR FOR QUESTIONS
FROM THE AUDIENCE