

Text Classification Using LangChain: Enhancing Natural Language Understanding

Principal Investigator:

Alina Hasan, Student Research Assistant, ahqgm@umsystem.edu

Objectives:

The primary objective of this work is to use LangChain to create a reliable text classification model. This model will consistently and accurately categorize text across various data types by employing cutting-edge techniques.

We seek to enhance the precision and effectiveness of natural language understanding by improving our model's comprehension of context, sentiment, and other language nuances. This will enable faster information processing and more accurate text classification.

We plan to develop an open-source, scalable system for text classification that can operate on various technological platforms, including smaller devices and cloud systems, and handle large volumes of data efficiently.

Goal:

Our key aim is to ensure that our text categorization model is highly accurate and efficient, performing at the forefront of current technology and excelling in tests and benchmarks. By offering comprehensive guidance and assistance to developers, we aim to facilitate the smooth integration of our model into their projects across various contexts. Benefit the scholarly and open-source communities by disseminating our findings and resources, thus promoting collaboration and further advancements in natural language processing. To promote cooperation and further advancements in the field of natural language processing, we will share our research, attend conferences, and make our code publicly available.

Cross-cutting Focus Areas:

1. Advancing AI Methods

Our project focuses on enhancing text classification by integrating innovative techniques. Utilizing deep learning architectures designed for natural language processing (NLP), hybrid models like BERT and GPT-3, we aim to improve the accuracy of text classification and handle complex text material with greater ease.

2. Creating Open-Source Foundation Models

We are dedicated to developing robust and adaptable text classification models, which we will share openly on platforms like Hugging Face. These models will cater to various applications, ensuring users can easily adapt and customize the models.

3. Empowering Use of Scientific Data

We aim to enhance the processing of scientific literature by classifying numerous scholarly articles and identifying patterns across various fields. We will develop interfaces and visualization tools to make these findings accessible.

4. Privacy-Preserving Methods

Data security and privacy protection are paramount in our project. We will use methods like federated learning and differential privacy to safeguard sensitive data and ensure compliance with ethical standards.

Domain-Specific Focus Areas:

1. AI Safety, Reliability, Security, and Privacy

We are developing AI models that prioritize safety and reliability, capable of handling diverse and adversarial text data. To address privacy concerns, we'll use advanced techniques such as secure multi-party computation and homomorphic encryption to process sensitive data securely.

2. Health Outcomes and Cancer Treatment

Our text classification model will analyze medical literature, focusing on cancer research. By categorizing and summarizing relevant studies, we aim to uncover insights that can drive new treatments and improve patient outcomes. This model will help healthcare professionals and researchers stay updated with the latest advancements in cancer treatment.

3. Environmental and Climate Challenges

Our model will process and analyze environmental data from research papers, reports, and policy documents to address issues like climate change and biodiversity loss. By organizing and interpreting large volumes of information, our model aims to enhance climate science research and offer actionable insights to policymakers, scientists, and the public.

Model Analysis:

1. GPT-J

Developer: EleutherAI

Architecture: Transformer-based, similar to GPT-3.

Size: 6 billion parameters.

Use Cases: Text generation, language modeling, completion tasks.

Performance: Offers higher quality text generation compared to GPT-Neo due to its larger parameter size.

Accessibility: Can be run locally without the need for API keys.

Pros for Text Classification:

Large Parameter Size: Helps capture nuanced language patterns, which can improve the model's ability to understand and classify text.

Local Execution: Accessible without the need for API keys, making it cost-effective for local deployments.

Cons for Text Classification:

Primary Use Case: Initially designed for text generation, requiring adaptation and fine-tuning for classification tasks.

Training Effort: May involve additional work to adapt from its primary use case to text classification, including extensive fine-tuning.

2. GPT-NeoX

Developer: EleutherAI

Architecture: Advanced transformer architecture.

Size: 20 billion parameters.

Use Cases: Complex text generation, large-scale language tasks.

Performance: Superior contextual understanding and high-quality text generation.

Accessibility: Requires significant computational resources and specialized hardware for local execution.

Pros for Text Classification:

Contextual Understanding: Offers excellent performance in understanding context, which is beneficial for complex classification tasks.

High-Quality Output: Capable of generating high-quality text, which can translate to more nuanced and accurate classification.

Cons for Text Classification:

Computational Requirements: High resource needs may limit accessibility, making it less practical for all users.

Adaptation Effort: Like GPT-J, GPT-NeoX requires adaptation from text generation to classification tasks, which can be resource-intensive.

3. DialoGPT

Developer: Microsoft

Architecture: Transformer-based, fine-tuned version of GPT-2 for dialogue.

Size: Small (117 million parameters), medium (345 million parameters), large (762 million parameters).

Use Cases: Conversational AI, dialogue generation.

Performance: Specializes in generating coherent and contextually relevant dialogue.

Pros for Text Classification:

Dialogue Specialization: Well-suited for tasks involving conversational data, such as customer support dialogues.

Parameter Options: Availability of different model sizes allows for flexibility based on resource constraints.

Cons for Text Classification:

Primary Focus: Primarily designed for conversational contexts, which may limit its effectiveness for general text classification tasks.

Fine-Tuning: Adapting it for broader text classification tasks requires fine-tuning and may not be as effective as models specifically designed for classification.

Technical Approach:

```
▶ # Install the required packages
!pip install langchain
!pip install -U langchain-openai
!pip install langchain_community
!pip install transformers
!pip install openai
!pip install openai langchain
!pip install --upgrade openai
!pip install openai==1.0.0
```

```
▶ from langchain.prompts import PromptTemplate
from langchain.chains import LLMChain
from langchain.memory import ConversationBufferMemory
from langchain.llms import OpenAI
from transformers import pipeline

# Define a prompt template
prompt_template = PromptTemplate(
    input_variables=["question"],
    template="You are a helpful assistant. Answer the following question: {question}"
)

# Create a memory object to store conversation history
memory = ConversationBufferMemory()

# Initialize the language model (using OpenAI as an example)
llm = OpenAI(api_key="sk-proj-SNpaeGRgRYKFRBY3dtBT3BlbkFJUjputpe2f24SGXeFikED")

# Create a chain with the prompt and memory
chain = LLMChain(
    prompt=prompt_template,
    llm=llm,
    memory=memory
)

# Initialize the text generation pipeline
text_generation_pipeline = pipeline("text-generation", model="microsoft/DialogPT-medium")

# Example conversation history
conversation_history = "Human: How is the weather?\nAI:"
#How is the weather?
#How do you like flowers?
#Do you like cookies?
#How do you feel about the current state of the economy?
#How do you feel about the new software update?
#How do you feel about the new social media platform you joined?
#How do you feel about your fitness routine?
#Do you like sunset or sunrise?
#How was your day?
#how are you doing?

# Generate a response using the pipeline
response = text_generation_pipeline(
    conversation_history,
    max_length=100,
    truncation=True,
    pad_token_id=text_generation_pipeline.tokenizer.eos_token_id
)

# Extract the generated response
generated_text = response[0]['generated_text'].strip()
print("DialogPT Response:", generated_text)

# Initialize the sentiment analysis pipeline
sentiment_pipeline = pipeline("sentiment-analysis")

# Perform sentiment analysis on the generated text
sentiment = sentiment_pipeline(generated_text)

# Print the sentiment analysis result
print("Sentiment Analysis:", sentiment)
```

Output

```
➡ No model was supplied, defaulted to distilbert/distilbert-base-uncased-finetuned-sst-2-english and revision af0f99b (https://huggingface.co/distilbert/distilbert-base-uncased-finetuned-sst-2-english)
Using a pipeline without specifying a model name and revision in production is not recommended.
DialogPT Response: Human: How is the weather?
AI: It's cold.
Sentiment Analysis: [{'label': 'NEGATIVE', 'score': 0.9992030262947083}]
```

Conclusion:

This research aims to develop an advanced text classification model using LangChain and several state-of-the-art models like GPT-J, GPT-NeoX, and DialoGPT. Our goal is to create a model that not only understands text more accurately but also works efficiently across different platforms, whether on smaller devices or cloud systems.

We've taken a detailed approach to training these models, preparing data carefully, and building scalable systems. Our focus on privacy and open-source contributions highlights our commitment to ethical practices and collaboration within the research community.

We believe this project will significantly advance the field of text classification and natural language understanding, benefiting areas such as healthcare, environmental science, and conversational AI. By making our work publicly available, we hope to encourage further research and innovation. We appreciate the support of our collaborators and the developers of the models we used. Their groundbreaking work has been crucial to our project.

Acknowledgments:

I would like to extend my heartfelt thanks to my faculty advisor, Yusuf Sarkar, for his invaluable help and guidance throughout this project. His support has been crucial to the success of this research.

References:

- EleutherAI. (2021). GPT-J: A 6 Billion Parameter Language Model. Available at [EleutherAI's GitHub repository](<https://github.com/EleutherAI/gpt-j>).
- EleutherAI. (2022). GPT-NeoX: A Large-Scale Transformer Language Model. Available at [EleutherAI's GitHub repository](<https://github.com/EleutherAI/gpt-neox>).
- Microsoft. (2021). DialoGPT: Large-Scale Generative Pretrained Transformer for Dialogue. Available at [Microsoft's GitHub repository](<https://github.com/microsoft/DialoGPT>).
- Google. (2020). T5 (Text-To-Text Transfer Transformer). Available at [Google's Research Blog](<https://ai.googleblog.com/2020/02/exploring-transfer-learning-with-t5.html>).
- Facebook AI. (2020). RoBERTa: A Robustly Optimized BERT Pretraining Approach. Available at [Facebook AI's GitHub repository](<https://github.com/pytorch/fairseq/tree/main/examples/roberta>).