

Truth or DeGPTion: Evaluating Lie Detection Capabilities of GPT-3.5 through Fine-Tuning on Personal Opinions, Autobiographical Memories, and Intentions

Tanner Graves

tanneraaron.graves@studenti.unipd.it

Marco Uderzo

marco.uderzo@studenti.unipd.it

Francesco Vo

francesco.vo@studenti.unipd.it

Mehran Faraji

mehran.faraji@studenti.unipd.it

Claudio Palmeri

claudio.palmeri@studenti.unipd.it

Abstract

This paper aims at evaluating the capabilities of GPT3.5 in the task of Lie Detection. This is done through the fine-tuning of GPT-3.5 on three English-language datasets encompassing personal opinions, autobiographical memories, and future intentions. Fine-tuning of LLMs consists in adapting a pre-trained language model to a specific task by further training the model on task-specific data, thereby enhancing its ability to generate contextually relevant and coherent text in line with the desired task objectives. In our investigation, the objective is to discern and classify instances of truth or deception.

1. Introduction

Multiple papers consistently show that the capability of humans to discern truth from deception is at chance level, there is a growing interest in employing Machine Learning methods, especially based on the Transformer Model, to more accurately predict the truthfulness of a statement. Indeed, the inherent pattern recognition capability of ML Models allows them to pick up subtle cues that humans just seem to miss. In this paper, we will use OpenAI’s GPT-3.5 Large Language Model (LLM), performing benchmarks on the performance of the base model, and then on a GPT-3.5 model specifically fine-tuned on the Opinion Dataset (Deceptive Opinions), Memory Dataset (Hippocampus) and Intention Dataset.

2. Methods

2.1. Dataset Preprocessing

The dataset was processed according to ‘Scenario 3’ of Locante et al.[1]. In Scenario 3, aggregation was performed

on the three train and test sets from Scenario 1 (opinions, memories, intentions), followed by fine-tuning the model on the aggregated sets. This scenario evaluates the model’s ability to classify truthful and deceptive statements across various contexts.

Regarding the implementation of this process, the single datasets were merged into a single one after being cleaned. Then we divided it into training, validation and test set. It is worth noting that we did not use all the data on the initial dataset, this is because the training of a LLM was made simpler by using a smaller training dataset. Finally, the datasets were formatted into JSON to align with the expected input format of the OpenAI API.

2.2. GPT-3.5 Fine-Tuning

To manage the potential high computational costs, the model was trained on a subset of the dataset. The model was trained utilizing the OpenAI API, and its performance was assessed through testing and comparison with GPT-3.5. Further experimentation was conducted to assess the impact of engineering the system prompt on overall performance.

Specifically, we noticed that the baseline GPT-3.5 prefers giving verbose or indecisive answers. Verbose answers, that actually classify a statement as genuine or deceptive can be classified easily. Nonetheless, the model decides not to give a definitive answer when it thinks it does not have enough information to classify the statement. The following example shows this behaviour.

User: "Each and every abortion is essentially a tragedy. The potential mother will suffer unforeseen consequences. Society as a whole will be deprived of the potential it could have received from the new life."

Baseline GPT-3.5: "There is no objective truth to the statement as it expresses subjective opinions and beliefs about abortion. It cannot be definitively classified as 'True' or 'False'."

To address this issue it was necessary to engineer a system prompt that discourages this behaviour and adequately explains the task. This prompt is provided to the model at every example query, so instructions should be concise to minimize any token overhead that leads to increased cost of training and queries.

System Prompt to Fine-Tuned GPT-3.5:

"You are an expert capable of discerning truthful from deceptive opinions based on speech patterns. Definitively classify the following statement as 'True' or 'False', based on the likelihood the statement represents a genuinely held belief or a deception."

This issue is avoided in the fine-tuned models as the training process rewards our expected behaviour and output format.

3. Results

4. Discussion

5. Code Availability

The datasets used and all the code used for this project is available at the following GitHub Repository.

6. References

List and number all bibliographical references in 9-point Times, single-spaced, at the end of your paper. When referenced in the text, enclose the citation number in square brackets, for example [1]. Where appropriate, include the name(s) of editors of referenced books.

References

- [1] Loconte et al. Verbal lie detection using large language models. 2023.

7. Appendix

7.1. Python Code

This should be full width, 1 column