Brendan Lim

CS4395.001

Jan 28, 2023

## Assignment 0: Getting Started

Natural Language Processing, or NLP for short, is the study of making computers able to compute human language through the application and combination of various approaches. Computation of human language is usual split into two: understanding what is being input and generating a response based on a query. Natural Language Understanding is where machines must understand what the author wants/feels based on the input passage and Natural Language Generation creates a passage of text based on an input. These two disciplines are tied very closely together, because to understand what kind of passage to generate, a model must first understand what is being said to them. NLP has many different applications like sentiment analysis. where the intent behind a passage of text is determined. and email spam detection.

NLP is important for artificial intelligence development as the ability to understand and mimic human speech/writing is a great way for humans to think that the machine itself is intelligent. This is because language is one indicator of intelligence and humans have the best linguistic capabilities currently known to us. The most famous example of NLP and AI being closely linked is the Turing Test, where a human cannot tell if the responses to a question have been generated by a human or a computer.

## Three Main Approaches to NLP

There are three approaches to NLP that have developed as computing power has exponentially increased and as the field has matured with computer science. The first is rules-based, which utilizes grammar rules used in language as well as computing theory to mimic language. This approach has been used for basic applications such as grammar checkers and chatbots. However, rules-based approaches tend to be too static and languages are more than just their rules. Languages and their user evolve over time and commonly break their own rules. Additionally, different languages will have different rules and cultures behind them, so a model for English cannot be easily translated for Chinese.

The second approach is using statistics and probabilities. The idea behind this approach is that the next word can be "predicted" based on the previous sequence of words and a corpus of text to determine the probability of the next word. It is possible to train smaller machine learning models using this type of approach, where the model optimizes the probabilities for the next word. Well-known applications of this approach include the next word predictor for your phone or search engine. An issue with this approach is that it is too short sighted; the phrase or sentence structure may make sense, but the overall paragraph will talk about the same thing in circles or run off into a tangent.

The third approach, deep learning, seems to be an extension of the statistical approach. The difference is that deep learning uses much larger and more complex machine learning models combined with massive datasets. This is recently possible because of our ability to store petabytes of data in large data centers. Deep learning approaches have shown themselves to be able to create large bodies of text that makes sense and elocute a central point. The main problem with this approach is it requires massive amounts of data storage and computation power to train a deep learning model, which may not be feasible for smaller projects.

## My Interest in NLP

As I've studied different subjects over the years, I've always had an interest in language and computation, so NLP seemed like a natural progression for me (although writing essays is another story). Starting from my high school Latin classes, I've always been intrigued by the complex structures of language juxtaposed by the fact that it is so simple that it becomes second nature to humans. Additionally, when I was first getting into machine learning, I did a sentiment analysis project to predict the positivity/negativity of a passage of text. While I was able to do the project by following a tutorial, I did not understand any of the NLP techniques and reasoning that went into making the project. After this course I plan on going back to that tutorial and see if I can understand it.