

To,

IITD-AIA Foundation of Smart Manufacturing

Subject: **Weekly Progress Report for Week 2**

Dear sir,

Following is the required progress report to the best of my knowledge considering relevant topics to be covered.

What's happening this week:

- Deep Reinforcement Learning
- Deep Learning Limitations and new Frontiers
- Python Libraries like Textblob
- Text cleaning and Preprocessing techniques
- Data extraction techniques from website
- Normalization and Stemming
- Word Embedding or vectorization
- Word2Vec

My Understanding of INTP23-ML-01: Chatbot for FSM

Scope:

This project involves developing an interactive conversational interface that can assist users by providing instant responses to their queries, guiding them through the website, and resolving common issues. The chatbot should be equipped with natural language processing capabilities, allowing it to understand and respond to user input accurately. It can offer student support, retrieve information from the website, provide student with the information he wants, guide him to the website of fsm.

Solution:

This kinda chatbot can be made using natural language processing (NLP) where it involves designing and implementing an intelligent conversational bot that can understand and respond to user queries in a natural and helpful manner.

Approach:

It involves collecting all the data from the scratch of the website and then designing the conversational flow and user interface. NLP and ML models will be applied to make the bot work with different languages.

## **Weekly Progress:**

### **June 12:**

Covered topics like Deep reinforcement learning.

- Reinforcement learning is a general-purpose framework for decision making which comes with states and action which turns in reasoning.
- Deep reinforcement learning defines the useful state space, action space, and reward by learning the data and gets us the insights from the this process.
- DRL framework consists of an agent, an environment, and a learning process. The agent interacts with the environment by observing the current state, taking actions, and receiving rewards. It then uses this information to update its policy or value function using techniques such as deep Q-networks (DQN).

### **June 13:**

Learnt about limitations of deep learning and it's new frontiers. And also python libraries used for NLP.

- Universal Approximation theorem: A single layer is enough to make a arbitrary precision or and function. that layer can be exponential big. Deep learning limitations :
- Deep Learning Limitations:
  1. The system is very fond of data.
  2. Easily fool by adversarial examples(Forcefully given inputs).
  3. poor at predicting uncertainty(he is not sure about his own probability).
- Python Libraries:
  1. TextBlob: It provides a simple API for common NLP tasks such as tokenization, part-of-speech tagging, noun phrase extraction, sentiment analysis, and more. It is built on the top of NLTK library.
  2. Gensim: It is an open-source Python library for topic modeling and document similarity analysis. Its applications include like topic modeling, document similarity analysis, text classification, and recommendation systems.

### **June 14:**

Studied some of the concepts of text cleaning and its pre-processing. And also done its practical implementation on a dataset.

- Studied text cleaning and pre-processing methods like tokenization , Noise Entities Removal, Data Visualization for Text Data, Parts of Speech (POS) Tagging.
- Tokenization are of 3 types White-space Tokenization, Regular Expression Tokenization, Word and sentence tokenization and this used to transform the data into smaller tokens for further processing.
- Noise Entities Removal are used to Remove unwanted Punctuation marks and stopped words.
- For visualizing the data we use techniques like frequency graphs,frequency distribution and word clouds.
- Parts of speech (POS) tagging is for identifying various kinds of speech and used in sematic analysi. POS tags are the basis of the lemmatization process for converting a word to its base form (lemma).
- This all the process is done on a red wine dataset which can be viewed on my github.

**June 15:**

Studied techniques of extracting the data from a website.

- Studied how to extract data from the website.
- Learned and implement web scrapping on a web page product.
- Studied how we can extract data from API in an web page.

**June 16:**

Studied Text cleaning techniques like Normalization, Stemming and Lemmatization with its practical implementation on red wine dataset. also chunking too.

- Normalization is the process of converting a token into its base form. It makes the text more uniform and easier to work with. Some of the normalization techniques are lower-casing, removing punctuation, removing special character, handling nos and lemmatization or stemming.
- Stemming refers to the process of removing the suffixes and prefixes of the word to obtain the root word.
- Lemmatization is an organized & step by step procedure of obtaining the root form of the word, as it makes use of vocabulary (dictionary importance of words) and morphological analysis (word structure and grammar relations).
- Difference between both of them is lemmatization converts the word to its meaningful base form, whereas stemming simply chops off the ends of a word using heuristics that often leading to incorrect meanings and spelling errors.
- In Chunking, we try to extract meaningful phrases from unstructured data in the form of text by tokenizing. It works on top of Part of Speech(POS) tagging.

**June 17:**

Studied Word embedding/ Vectorization, Term frequency-inverse document frequency with practical implementation.

- As it is difficult to process the raw form of the text so we require numerical numbers as inputs to perform any sort of task. For that we use Word Embedding or Vectorization which converts texts into number vectors and there may be different numerical representations of the same text.
- Vector representation of a one-hot encoded vector represents in the form of 1, and 0 where 1 stands for the position where the word exists and 0 everywhere else.
- In the N-Gram method, a document term matrix is generated, and each cell represents the count.
- Term frequency-inverse document frequency ( TF-IDF) gives a measure that takes the importance of a word into consideration depending on how frequently it occurs in a document and a corpus. Here Term frequency denotes the frequency(TF) of a word in a document and Inverse document frequency (IDF)is the logarithmic ratio of no. of total documents to no. of a document with a particular word.

**June 18:**

Studied word2vec model.

- Word2Vec model is used for Word representations in Vector Space and is a neural network model that attempts to explain the word embeddings based on a text corpus.
- Word2Vec model composed of two preprocessing modules or techniques i.e., Continuous Bag of Words (CBOW), Skip-Gram.
- The aim of the CBOW model is to predict a target word in its neighborhood, using all words. To predict the target word, this model uses the sum of the background vectors. This model works in both the cases single context word and multiple context words.
- The skip-gram model is the exact opposite of the CBOW model. In this case, the target word is given as the input, the hidden layer remains the same, and the output layer of the neural network is replicated multiple times to accommodate the chosen number of context words.
- Skip-Gram model works well with a small amount of the training datasets, and can better represent rare words or phrases. However, the CBOW model is observed to train faster than Skip-Gram, and can better represent more frequent words which mean gives slightly better accuracy for the frequent words.
- Done sentiment analysis on twitter's positive and negative comments.

**GANTT CHART FOR INTP23-ML-01**

