

PROGRESS REPORT(PROJECT DIARY) :

Name : Hiren Bavaskar

Roll : 190050049

Project : Unscripted SOC 2021

DATE	PROGRESS
April Week 1, 2021	<p>Revised basics of Python and Git. Revised Numpy, Pandas and Matplotlib libraries.</p> <p>Solved the following exercises from Google developers course:</p> <ul style="list-style-type: none">• Numpy colab• Pandas colab
April Week 2, 2021	<p>Read 2 chapters of the prescribed ML textbook.</p> <p>Read about univariate regression, supervised learning, unsupervised learning, batch learning, online learning and exercises on scikit.</p> <p>Solved some exercises in python for ML in Jupyter notebook.</p>
April Week 3 - May Week 1	4th Semester Endsems
May Week 2, 2021	<p>I learned about NLP from Stanford lectures and other resources[Resource]. Got to know about word vectors, stop words, regex, word and sentence tokenization.</p> <p>Read about the nltk package in python and looked over some code snippets.</p>
June Week 1, 2021	<p>Started exploring Speech Recognition</p> <p>Learnt about ASR(Automatic Speech Recognition)</p> <p>There are 2 kinds of models:</p> <ul style="list-style-type: none">• Acoustic model - Deals with sound to phonetic representation• Language model - Deals with the language structure <p>Learnt use of hidden markov model (HMM) and Deep Neural networks (DNN) in ASR</p> <p>Started exploring DeepSpeech library for STT</p> <p>Learnt about WER (word error rate) and the use of</p>

	<p>jiwer library in python for comparing transcribed and original text</p> <p>The STT and WER calculation files can be found in the github repo.</p>
June Week 2- 3	<p>Couldn't do much because of unavailability of my laptop</p>
June Week 4, 2021	<p>Started huggingface crash course [link]</p> <p>Learnt about the use of transformers library and the pipeline module in various applications of NLP like:</p> <ul style="list-style-type: none"> • Classifying whole sentences • Classifying each word in a sentence • Generating text content • Extracting an answer from a text • Generating a new sentence from an input text <p>Pipeline object of transformers library is useful for executing NLP tasks on one or several lines of text. Pipeline by itself chooses a model appropriate for the task and does the required pre and post processing with the pre-trained model to make reasonable predictions. Eg: Sentiment analysis, classification based on labels(zero shot classification), text generation, filling masks(fill in the blanks).</p> <p>An ideal way to use transformers is to take a pre-trained model and then fine tune it according to our use.</p> <p>Learnt about encoder models, decoder models and sequence to sequence model(combination of the 2)</p> <p>Encoder: Takes the input and gives a vector representation of its features. Useful for understanding inputs like sentence classification.</p> <p>Decoder: Takes representation of features along with the input for generating output.</p> <p>Useful for generative tasks like text generation.</p> <p>Tokenizer -> Model -> Post processing</p> <p>Tokenizer analyzes the raw text and converts them into numbers for neural networks to understand</p> <p>Models can only process numbers so they take the output of tokenizers and do the further predictions.</p>

	<p>Implemented wav2vec model and transformers for speech to text in english using audio interval feature of librosa (to avoid crashing of file in colab)</p> <p>Colab file link</p> <p>[Number of hours spent this week: 11-12] (Studying course: 8-9 hrs, coding and references: 3-4 hrs)</p>
July Week 1, 2021	<p>Agenda for this week was to find good datasets to fine tune our model for speech to text and for documents summarization.</p> <p>https://www.openslr.org/12/ The above site contains about 300MB of data. The content involves .flac files of 2-10 seconds with their corresponding correct transcripts.</p> <p>https://www.kaggle.com/sunnysai12345/news-summary?select=news_summary.csv This site provides the text files (news articles) and provides their summarized text. The file is in the type of CSV files. It is about 50 MB of data which will be sufficient for fine tuning.</p> <p>https://colab.research.google.com/drive/1bewbhx9Ej4_iyeM9eBO_M_bPhHu48_Ap?usp=sharing This is the testing done on our voices. It's still gibberish so we are thinking of not using the wav2vec2 model ahead.</p> <p>https://colab.research.google.com/drive/13yuNiuNmYP9Bl-itd2t5cn5_QsV6ZSoe#scrollTo=ZHLHvkPy85Zy Further part of text summarization was taken upon by Harshit.</p> <p>[Number of hours spent this week: 9-10]</p>
July Week 2, 2021	<p>Final week before project submission</p> <p>Used google speech recognition library for speech to text for large audio files. However, it is running correctly only on US-English language.</p>

	<p>Tested on 52MB audio file(5 mins of audio) and the results were good enough</p> <p>Added summarizer python file which summarizes the transcription into any number of lines of summary the user wants.</p> <p>Added documentation of the project in github</p> <p>[Number of hours spent this week: 14-15 hrs]</p>
--	---