# CIS 6930: Trustworthy Machine Learning
## Project Proposal: Decrypting job title classification model

Tanvi Jain
*(Point of Contact)*
tjain@ufl.edu

Nimish Bajaj
nimishbajaj@ufl.edu

September 21, 2021

## 1  Introduction

Determine whether explainable/interpretable ML techniques like LIME and others can provide useful insights into the cause of the unfairness.

A model's accuracy is not always enough to state whether it will perform well in the wild because the accuracy highly depends on the data that was used to train and test the model. Explaining the results of a model, and identifying what leads to a particular classification can help provide insights into the model. Once you add interpretability to the model, its result can be easily understood by a domain expert can be verified for correctness. This greatly enhances the trustworthiness of the model.

We are proposing to build a text classification model and explaining the results using LIME, and validate if LIME can provide stable and meaningful features.

## 2  Background and Related Work

Model interpretation techniques have been in place for the last couple of years, but they are not always able to provide stable inferences for a given problem. Work done on inference healthcare model[3], Perturbing Inputs for Fragile Interpretations in Deep Natural Language Processing [5] and many more have a common objective, i.e. to have robust interpretations for determining the trustworthiness of the model.

For our project, we are planning to decrypt the Job classification model on Resume corpus. Our work will include categorizing resumes using CNN and RNN. We are going to train Neural Network and use the inference technique 'LIME'. Since the data contains technical words which are sometimes outside of the English vocabulary, we are going to tune word embedding for the problem domain.

Existing work has explored a CNN for training a classification model, as part of the project we will be using an RNN to learn the temporal features. The data contains technical words which are sometimes outside of the English vocabulary, we are going to tune word embedding for the problem domain by utilizing the training dataset and other datasets from the same domain.

## 3  Proposed Approach & Plan

Following is the approach we are going to take:
1. Fetch data and take samples as per the distribution of data for train, test, and validation
2. Clean and vectorize the data for Neural Network (NN) approach
3. Build a NN architecture to learn classification labels
4. Train, test, and validate the neural network
5. Use LIME for model inferencing

6. Manually validate LIME output

For the test and experimentation, we are currently foreseeing that the local machine would be sufficient for carrying out the work. If needed we will use the hyper-gator resources that are availed to us as part of this project.



(a) Example Input



(b) Model prediction

Figure 1: Model output description [3]

# 4 Timeline

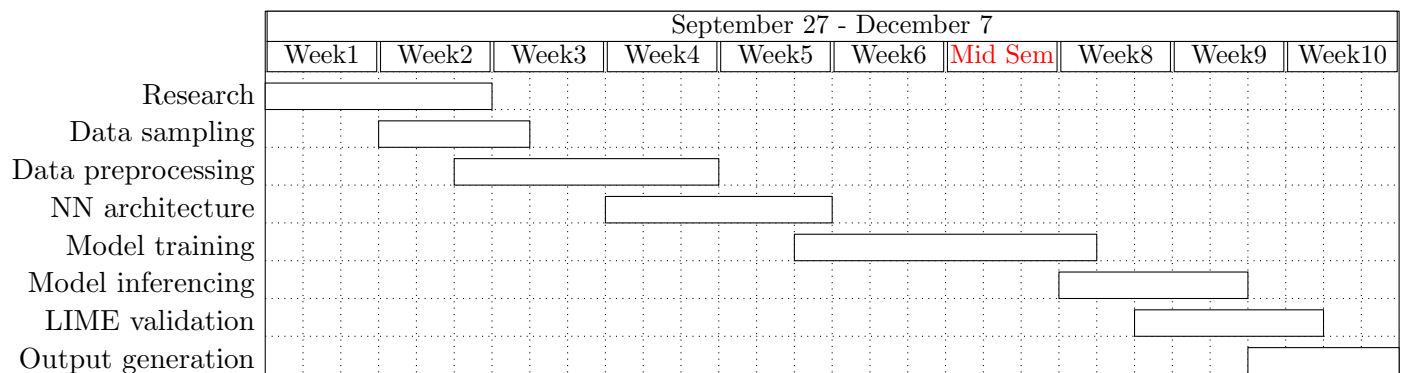| | September 27 - December 7 | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Week1 | Week2 | Week3 | Week4 | Week5 | Week6 | Mid Sem | Week8 | Week9 | Week10 |
| Research | | | | | | | | | | |
| Data sampling | | | | | | | | | | |
| Data preprocessing | | | | | | | | | | |
| NN architecture | | | | | | | | | | |
| Model training | | | | | | | | | | |
| Model inferencing | | | | | | | | | | |
| LIME validation | | | | | | | | | | |
| Output generation | | | | | | | | | | |

Figure 2: Gantt Chart

# References

[1] Damien Garreau and Ulrike von Luxburg. Explaining the explainer: A first theoretical analysis of lime. 2020.

[2] Kameni Florentin Flambeau Jiechieu and Norbert Tsopze. Skills prediction based on multi-label resume classification using cnn with model predictions explanation. *Neural Computing and Applications*, 33(10):5069–5087, May 2021.

[3] Khansa Rasheed, Adnan Qayyum, Mohammed Ghaly, Ala Al-Fuqaha, Adeel Razi, and Junaid Qadir. Explainable, Trustworthy, and Ethical Machine Learning for Healthcare: A Survey. 4 2021.

[4] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the predictions of any classifier. 2016.

[5] Sanchit Sinha, Hanjie Chen, Arshdeep Sekhon, Yangfeng Ji, and Yanjun Qi. Perturbing inputs for fragile interpretations in deep natural language processing, 2021.

[6] Tim Zimmermann, Leo Kotschenreuther, and Karsten Schmidt. Data-driven hr - resume analysis based on natural language processing and machine learning, 2016.