



DSCI 551 Project Proposal

Early Stage Diabetes Risk Prediction





Project Title

Early Stage Diabetes Risk Prediction

Topic

Healthcare

Potential Dataset

[https://archive.ics.uci.edu/ml/datasets/Early+stage+diabetes+risk+prediction+dataset.](https://archive.ics.uci.edu/ml/datasets/Early+stage+diabetes+risk+prediction+dataset)



Motivation

1

Diabetes is a chronic (long lasting) health condition that affects how your body turns food into energy. As there isn't a cure yet for diabetes, diagnosing it at an early stage is essential to be clinically treated.

2

Due to the presence of a relatively long asymptomatic phase, early detection is a difficult task and hence we feel that if data science is used for this purpose it would be of immense social value.

3

Since diabetes is linked to the lifestyle of a person, based on a simple questionnaire we can predict whether a person is diabetic or not. This questionnaire includes critical questions like what is the age and gender of the person, whether the person is experiencing excessive urination or not, or is he or she experiencing excessive thirst or not and so on which help us predict whether the person is diabetic or not.



Project Overview

Data: For this project we will need a labelled dataset which is open sourced. While looking to satisfy these requirements we found a dataset where data was collected using direct questionnaires and diagnosis results from the patients in the Sylhet Diabetes Hospital in Sylhet, Bangladesh. The dataset contains the signs and symptoms of newly diabetic or would be a diabetic patient.

Machine Learning Model: We plan on using this publicly available dataset which has a list of carefully chosen 16 questions along with the corresponding diagnosis results to fit a carefully chosen and designed supervised machine learning classification algorithm like the random forest, xgboost (subject to change if a better alternative is found).

Web App (Browser): To complement this we will also be developing a web based User Interface (UI) where the user will be able to answer those 16 questions himself and see what our model predicts corresponding to his inputs (whether he is diabetic or not). Apart from this, the UI will also have insightful data visualizations where we explore the most common features associated with diabetic risk.



Team

Praveen Iyer | MS Applied Data Science | Email ID - praveeni@usc.edu | USC ID - 5549510111

Tejas Sujit Bharambe | MS Applied Data Science | Email ID - tbharamb@usc.edu | USC ID - 2160849114

Project Plan, Milestones & Team Members' Responsibilities

[illegible]



Relevant Papers

- Likelihood Prediction of Diabetes at Early Stage Using Data Mining Techniques
 - [\[Web Link\]](#) - Computer Vision and Machine Intelligence in Medical Image Analysis. Springer, Singapore, 2020. 113-125
 - Authors and affiliations:
 - M. M. Faniqul IslamEmail
 - Rahatara Ferdousi
 - Sadikur Rahman
 - Humayra Yasmin Bushra
- https://link.springer.com/chapter/10.1007/978-981-16-5655-2_41



Thank you.

