# Diabetes Prediction Documentation

Data Mining Final term Project
Prof. Yasser Abdullah
CS-634 104

-Vinit Walke(vnw3)

## Introduction

Diabetes prediction is a crucial healthcare task for identifying people at risk of developing the disease. This document outlines the process of using machine learning models for diabetes prediction. It covers data preparation, scaling, transformation, cross-validation, and three models: Naive Bayes, Random Forest, and Long Short-Term Memory (LSTM).

## Data Preparation

The diabetes dataset includes features like blood pressure, glucose levels, and body mass index (BMI), with a target variable indicating diabetes diagnosis. Before applying machine learning models, the data undergoes preprocessing:

1. **Data Cleaning:** Handling missing values, outliers, and errors.
2. **Feature Selection:** Choosing relevant features for predicting diabetes.
3. **Data Encoding:** Converting categorical variables to numerical format (if needed).
4. **Data Scaling:** Scaling numerical features to a similar range for equal contribution to the model.

**Data Wrangling and Transformation**

Python

```python
# Data Cleaning and Transformation
# - Handle missing values
# - Feature selection
# - Data encoding (if necessary)
# - Data scaling
```

**Cross-Validation**

Cross-validation is used to assess machine learning models' performance. It involves splitting the dataset into multiple subsets (folds), training the model on one subset, and testing it on the remaining data. This process is repeated, and performance metrics are averaged across all folds.

Python

```python
# Cross-validation for Naive Bayes
# Cross-validation for Random Forest
# Cross-validation for LSTM
```

**Model Implementation**

**Naive Bayes**

Naive Bayes is a probabilistic classifier based on Bayes' theorem, assuming independence between features.

Python

```python
# Naive Bayes Model
from sklearn.naive_bayes import GaussianNB
nb_model = GaussianNB()
```

**Random Forest**

Random Forest is an ensemble learning method that combines predictions from multiple decision trees.

Python

```python
# Random Forest Model
from sklearn.ensemble import RandomForestClassifier
rf_model = RandomForestClassifier(n_estimators=100,
random_state=42)
```

**LSTM (Long Short-Term Memory)**

LSTM is a type of recurrent neural network (RNN) architecture, particularly useful for sequence prediction tasks.

Python

```python
# LSTM Model
from keras.models import Sequential
from keras.layers import LSTM, Dense

# Reshape input data for LSTM
X_lstm = X_scaled.reshape((X_scaled.shape[0], 1,
X_scaled.shape[1]))

lstm_model = Sequential()
lstm_model.add(LSTM(50, activation='relu', input_shape=(1,
X_scaled.shape[1])))
lstm_model.add(Dense(1, activation='sigmoid'))
```

```
lstm_model.compile(optimizer='adam', loss='binary_crossentropy',
metrics=['accuracy'])
```

**Evaluation Metrics-**

Evaluating the performance of classification models is vital, and there are several metrics used for this purpose. Here's a breakdown of some common metrics applicable to the Naive Bayes, Random Forest, and LSTM models you mentioned in the context of diabetes prediction:

Accuracy:

- This is a basic metric representing the proportion of correct predictions made by the model. It's calculated as the number of correctly classified instances divided by the total number of instances.
- While seemingly straightforward, accuracy can be misleading, especially for imbalanced datasets where one class might significantly outnumber the others.

Precision:

- Precision focuses on the positive predictions. It tells you what proportion of instances labeled positive by the model are actually positive.
- In diabetes prediction, this would represent the ratio of people correctly identified as having diabetes out of all those predicted to have it.

Recall:

- Recall, also known as sensitivity, focuses on completeness. It tells you what proportion of actual positive cases were identified by the model.
- In diabetes prediction, this would represent the ratio of people with diabetes that the model correctly identified out of all the people who actually have diabetes.

F1-Score:

- F1-Score is a harmonic mean of precision and recall, combining their strengths. It provides a single metric that considers both how many relevant instances were retrieved (recall) and how many irrelevant instances were not retrieved (precision).
- A high F1-Score indicates a good balance between precision and recall.

ROC AUC (Area Under the Receiver Operating Characteristic Curve):

- ROC AUC is a performance measure for classification models at various classification thresholds. It depicts the model's ability to distinguish between positive and negative classes.
- A higher ROC AUC indicates better performance. It's particularly useful for imbalanced datasets where accuracy might not be a reliable metric.

Choosing the Right Metric:

The choice of metric depends on the specific problem and its priorities. Here's a general guideline:

- If correctly classifying all positive cases is crucial (e.g., identifying all diabetic patients), prioritize recall.
- If minimizing false positives is important (e.g., avoiding unnecessary medical tests), prioritize precision.

- If a balanced approach is desired, consider F1-Score or ROC AUC.

It's often recommended to use a combination of metrics to get a more comprehensive understanding of the model's performance.

**Conclusion**

Diabetes prediction is crucial in healthcare, and machine learning models offer valuable tools. By following the outlined process of data preparation, feature engineering, model implementation, and evaluation, we can make accurate predictions to identify individuals at risk of developing diabetes, enabling early interventions and improved healthcare management.

This documentation provides a high-level overview of the diabetes prediction process. Adjustments might be needed based on the specific dataset and prediction task requirements.

Sources-

Dataset-https://www.kaggle.com/datasets/iammustafatz/diabetes-prediction-dataset

Links-

https://medium.com/analytics-vidhya/introduction-to-long-short-term-memory-lstm-a8052cd0d4cd

https://medium.com/@kangeugine/long-short-term-memory-lstm-concept-cb3283934359

https://medium.com/analytics-vidhya/understanding-cross-validation-for-beginners-31e0c606ebe0