

# Project Requirements Document: Diabetes Prediction Model Development

## Objective

Develop a machine learning model to predict whether a patient has a chronic condition (`Chronic_flag`), leveraging a dataset of diabetes patients. The project involves data exploration, preprocessing, feature engineering, and model building, with the following considerations and requirements:

## Dataset Overview

### Columns Description

Column Name	Description	Notes
Unique_Identifier	Unique identifier for each patient.	Will not be used in the model; serves only for tracking.
Gender	Patient's gender (e.g., Male, Female).	Needs to be encoded into numeric format (e.g., 1 for Male, 0 for Female).
Religion	Patient's religious affiliation.	Suggested to drop, as it is unlikely to significantly impact the model's performance and doesn't have high correlation with output column.
Nationality	Patient's nationality.	Contains – for missing values; requires discussion on whether to fill values, drop rows, or exclude column specially because it doesn't have high correlation with output column. Also contain data at uppercase with lower case solve it in future by transform to one look before loading into data or better make it as choices for user not entered text data.
Avg_HBA1C Results	Average HbA1c results, indicating blood sugar control.	Contains 0 values that are incorrect and should be handled via imputation or removal. Range of values should fit in: Normal Range: 4% to 5.6% Prediabetes: 5.7% to 6.4% Diabetes: 6.5% or higher
HBA1C test Compliance	Indicates the patient's adherence to taking HbA1c tests.	Binary; no special notes at this stage.
D_Of_Birth	Patient's date of birth.	Needs transformation to calculate Age for inclusion in the model.
Diagnosis_Type	Type of diabetes diagnosis (e.g., Type I, Type II).	Contains only Type II in the dataset; suggested to drop as it provides no variability.
Chronic_flag	Target variable indicating chronic condition (1 = Yes, 0 = No).	Imbalanced dataset; oversampling (e.g., SMOTE) will be needed to address bias.

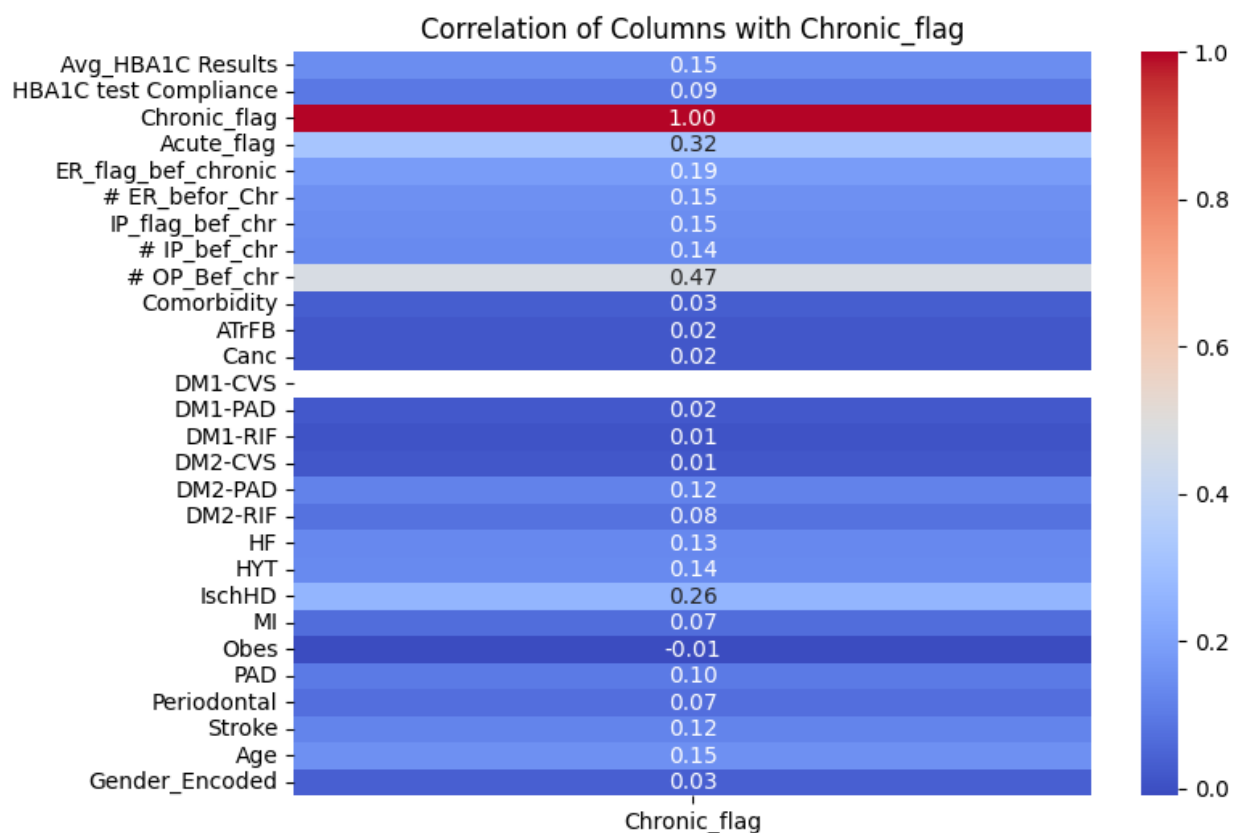
<b>Acute_flag</b>	Indicates acute conditions related to diabetes (1 = Yes, 0 = No).	Binary; needs scaling to match other features.
<b>ER_flag_bef_chronic</b>	Indicates whether the patient visited the ER before chronic diagnosis.	Binary; needs scaling.
<b># ER_befor_Ch</b>	Number of ER visits before chronic diagnosis.	Requires outlier handling and potential scaling.
<b>IP_flag_bef_chr</b>	Indicates whether the patient had inpatient admissions before chronic diagnosis.	Binary; needs scaling.
<b># IP_bef_chr</b>	Number of inpatient admissions before chronic diagnosis.	Requires outlier handling and potential scaling.
<b># OP_Bef_chr</b>	Number of outpatient visits before chronic diagnosis.	Requires outlier handling and potential scaling.
<b>Comorbidity</b>	Other conditions present alongside diabetes.	Requires encoding if categorical.
<b>ATrFB</b>	Presence of atrial fibrillation.	Binary; needs scaling.
<b>Canc</b>	Presence of cancer.	Binary; needs scaling.
<b>DM1-CVS / DM2-CVS</b>	Cardiovascular conditions associated with diabetes (Type I/II).	Binary; needs scaling.
<b>DM1-PAD / DM2-PAD</b>	Peripheral Artery Disease associated with diabetes (Type I/II).	Binary; needs scaling.
<b>DM1-RIF / DM2-RIF</b>	Renal insufficiency associated with diabetes (Type I/II).	Binary; needs scaling.
<b>HF</b>	Presence of heart failure.	Binary; needs scaling.
<b>HYT</b>	Presence of hypertension.	Binary; needs scaling.
<b>IschHD</b>	Presence of ischemic heart disease.	Binary; needs scaling.
<b>MI</b>	Presence of myocardial infarction (heart attack).	Binary; needs scaling.
<b>Obes</b>	Presence of obesity.	Binary; needs scaling.
<b>PAD</b>	Presence of peripheral artery disease.	Binary; needs scaling.

Periodontal	Presence of gum disease.	Binary; needs scaling.
Stroke	History of stroke.	Binary; needs scaling.

## Notes and Actionable Points

### 1. Correlation:

- The correlation heatmap shows no strong correlation between `Chronic_flag` and most other columns, including `Avg_HBA1C_Results`, which seems unexpected. Overall, there are no highly correlated variables with `Chronic_flag`.



### 2. Missing Values:

- Nationality:** Contains - for missing values. Discuss with the data team to decide on:
  - Obtaining correct values if available.
  - Dropping rows with missing values.
  - Excluding the column entirely.
- Avg\_HBA1C\_Results:** Contains 0 values, which are invalid. Potential solutions:
  - Use ML-based imputation methods based on other patient attributes.
  - Drop rows with 0 values if correct values are unavailable.
- Religion:** Likely not impactful for prediction; recommend removing the column.

3. **Imbalanced Target Variable:**

- `Chronic_flag` has a majority of 0 and a minority of 1. Address imbalance using oversampling methods like **SMOTE** to avoid bias.

4. **Feature Engineering:**

- **Age:** Calculate from `D_Of_Birth` and include it as a feature.
- **Binary Columns:** Normalize to a common range (e.g., 0–1) to ensure consistency.

5. **Redundant Features:**

- **Diagnosis\_Type:** Contains only one value (`Type II`) across all rows. Exclude from modeling as it offers no predictive power.

6. **Data Encoding:**

- Categorical columns (`Gender`, `Comorbidity`, etc.) must be encoded into numeric representations.

7. **Scaling:**

- Ensure all numerical data is scaled (e.g., using **StandardScaler**) to align feature magnitudes.

8. **Outlier Detection:**

- Columns like `# ER_befor_chr`, `# IP_bef_chr`, and `# OP_Bef_chr` may contain outliers. Apply robust scaling or capping methods.