# Identifying Key Health Indicators and Sociodemographic Factors Associated with Diabetes Risk

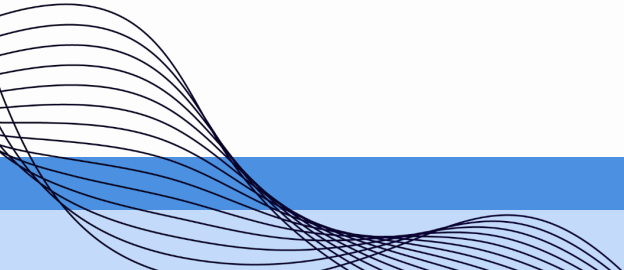**Alice Kang, Qidi An, Manwen Jia, Leon Li, Jiahao Chen**

# Background

**What is Diabetes?**
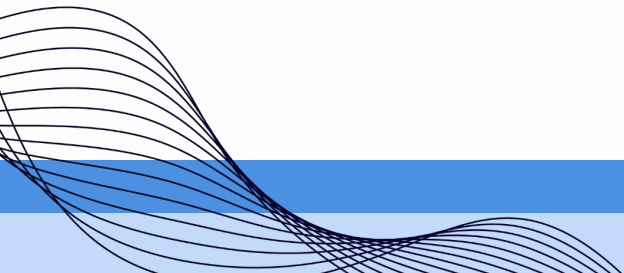
**Why Study Diabetes?**

- 1 in 10 adults in the U.S. has diabetes
- Rising global prevalence → 422M+ cases worldwide
- Early prediction = better prevention & management
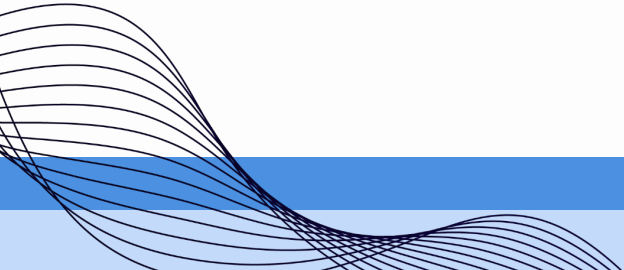
# Background

**Dataset Snapshot**

- 100,000 individuals
- Clinical + demographic features
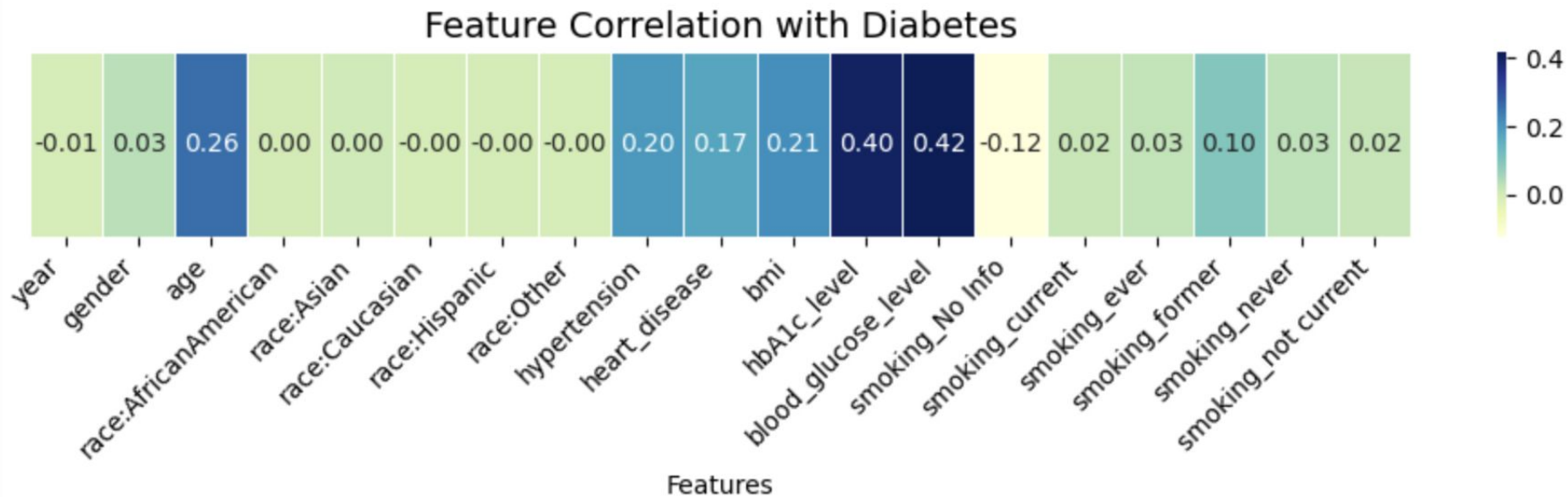- Source: *Diabetes Health & Demographics Dataset by ZIYA*
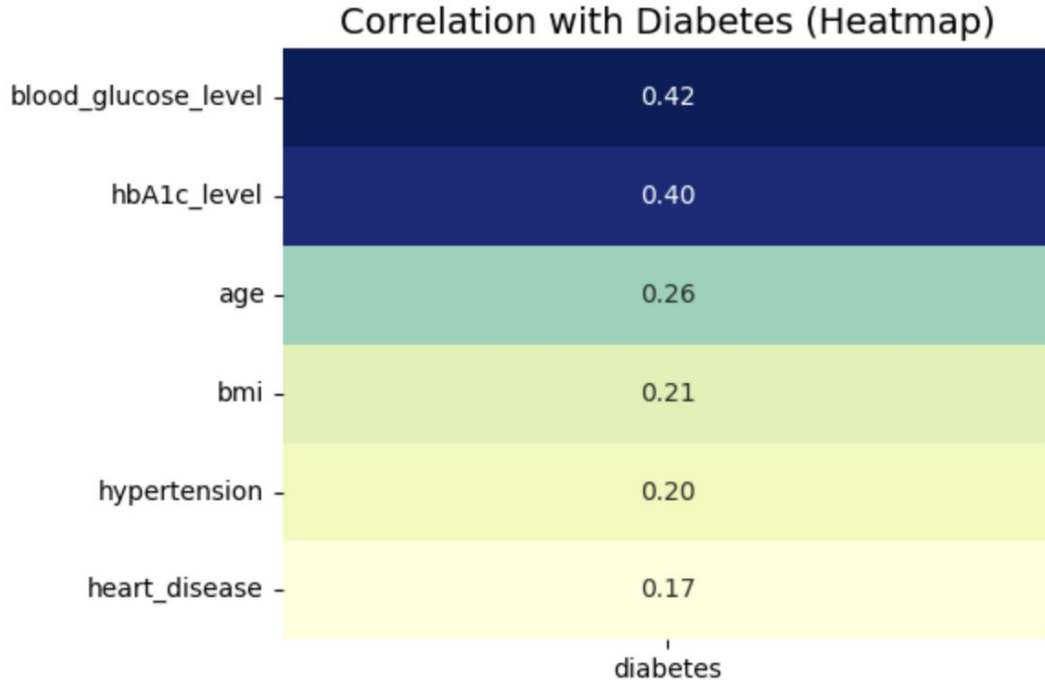
# Background

**Our Project Goal:**

- Predict diabetes status using health indicators
- Identify key risk factors from medical and demographic data

# Heatmap of correlation coefficients of each characteristic with diabetes



Feature Correlation with Diabetes

# Heatmap of Selected Feature characteristic



Correlation with Diabetes (Heatmap)

| | diabetes |
|---|---|
| blood_glucose_level | 0.42 |
| hbA1c_level | 0.40 |
| age | 0.26 |
| bmi | 0.21 |
| hypertension | 0.20 |
| heart_disease | 0.17 |

➔ **Blood Glucose Level**

➔ **HbA1c Level**

➔ **Age**

➔ **Bmi**

➔ **Hypertension**

➔ **Heart Disease**

# age

## bmi



# hbA1c_level

## blood_glucose_level

**hypertension**                    **heart_disease**

# Methods

➢ Rescaling : use standardscaler()
➢ Baseline Model Consideration:
   ○ 91.5% of individuals do not have diabetes, and only 8.5% do.
   ○ Predicting Probabilities of "no diabetes" gets 91.5% accuracy.



Distribution of Diabetes Cases

# Methods

Stratified Sampling for Data Splitting

- ➢ **Training Set : 80%**

- ➢ **Validation Set : 10%**

- ➢ **Test Set: 10%**

- ➢ **Stratify by diabetes proportion**

```
Train set diabetes proportion: 0.0850
Validation set diabetes proportion: 0.0850
Test set diabetes proportion: 0.0850
```

# Methods

**Model Selection & Tuning**

- Tested models: **Logistic Regression**, **Decision Tree**, **SVM**, **KNN**

- Each model was tuned twice:

  - Once with **accuracy** as the scoring metric
    i. measure the overall correctness of predictions

  - Once with **recall**, which is crucial for disease detection
    i. focuses specifically on finding all positive instances

# Methods

| Model | Accuracy score (optimizing for accuracy) | Accuracy score (optimizing for recall) |
|---|---|---|
| **Logistic** | **0.96**19 | **0.96**23**(Finalized Model)** |
| **Decision Tree** | **0.97**31 | **0.95**70 |
| **KNN** | **0.96**77 | **0.96**98 |
| **SVM** | **0.96**85 | **0.96**85 |

# Overfitting Curve

## Logistic Model
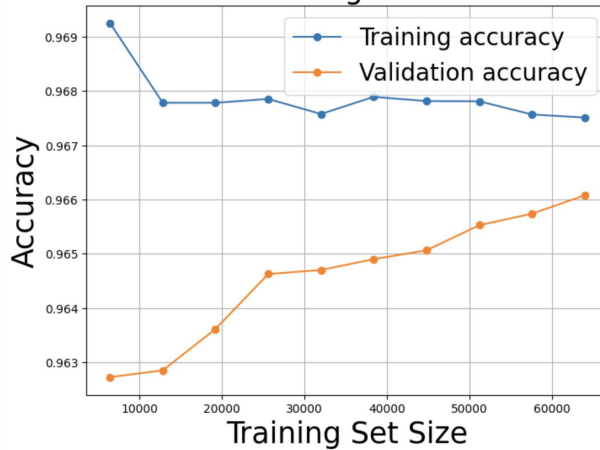
### Learning Curve



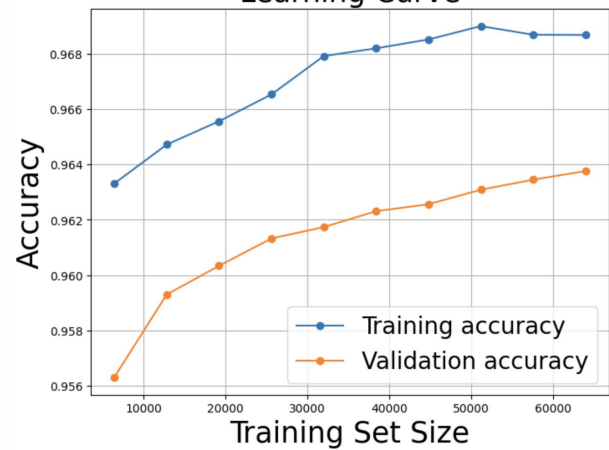## Decision Tree Model

### Learning Curve



## SVM Model

### Learning Curve



## KNN Model

### Learning Curve

# Results

❖ Our final Model is **Logistic Regression Model.**

$$P(\text{diabetes} = 1 \mid X) = \frac{1}{1 + \exp\left(5.005 - \sum_i w_i x_i\right)}.$$

$$\sum_i w_i x_i = 2.53 \cdot \text{hbA1c\_level} + 1.35 \cdot \text{blood\_glucose\_level} + 1.04 \cdot \text{age} + 0.61 \cdot \text{bmi}$$

$$+ 0.20 \cdot \text{hypertension} + 0.16 \cdot \text{heart\_disease}$$

❖ Hyperparameters We Choose    ❖ Test Data Accuracy ≈ **95.8%** (> 91.5%)

- c= **10** **(low regularization)**
- penalty= **'l1'** **(Lasso)**
- solver= **'liblinear'** **('liblinear' optimization algorithm)**

# Hyperparameter tuning:

- GridSearchCV
- 3-fold cross-validation
- Scoring metric: "Recall"

- Tuning parameters: C=[0.01,0.1,10] penalty=['l1', 'l2'] solver=['liblinear', 'sage']

```
Fitting 3 folds for each of 16 candidates, totalling 48 fits
Logistic Best Params: {'C': 10, 'penalty': 'l1', 'solver': 'liblinear'}
Logistic Best CV Accuracy: 0.62
Logistic Validation Accuracy: 0.96
```
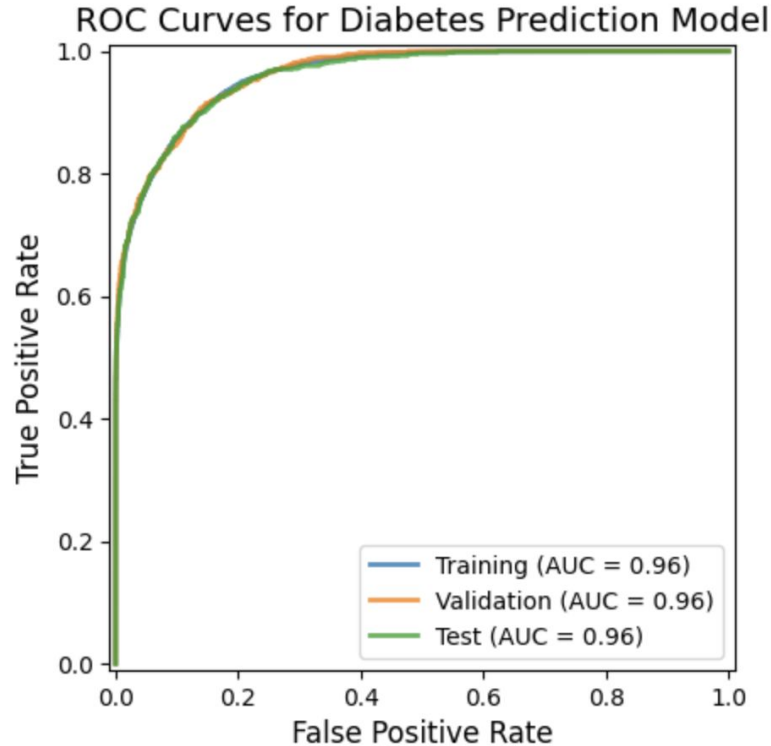
## Test Data Confusion Matrix:

|  | Predicted 0(Negative) | Predicted 1(Positive) |
|---|---|---|
| Actual 0(No Diabetes) | 9057 | 93 |
| Actual 1(With Diabetes) | 327 | 523 |

- **Precision** = 523/(523+93) ≈ **0.849**
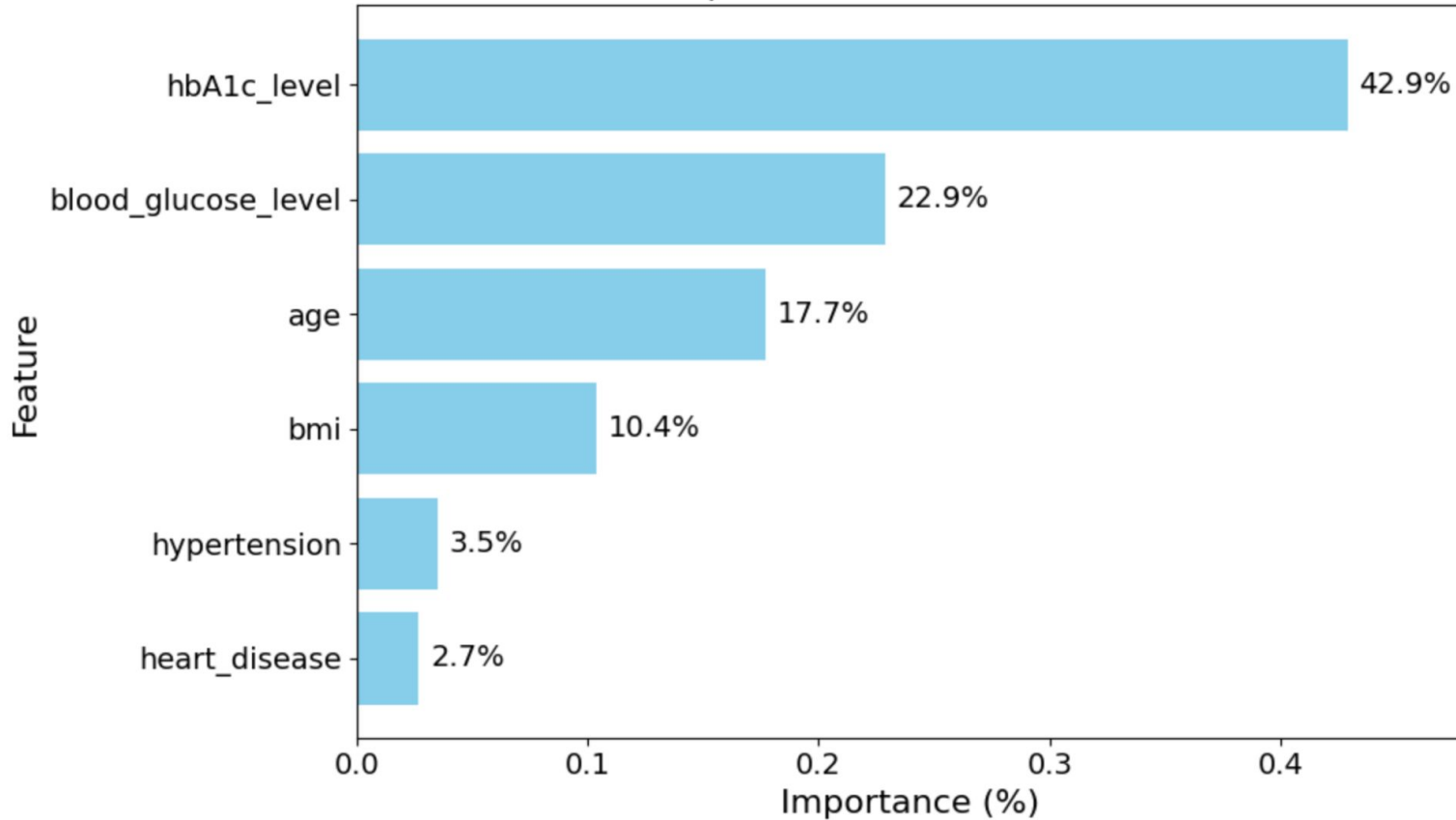High Confidence in Predicting Diabetes

- **Recall** = 519/(519+327) ≈ **0.613**
Identified Patients: 61.3%      Missed Patients: 38.7%

# ROC Curve Demonstrates Robust Generalization



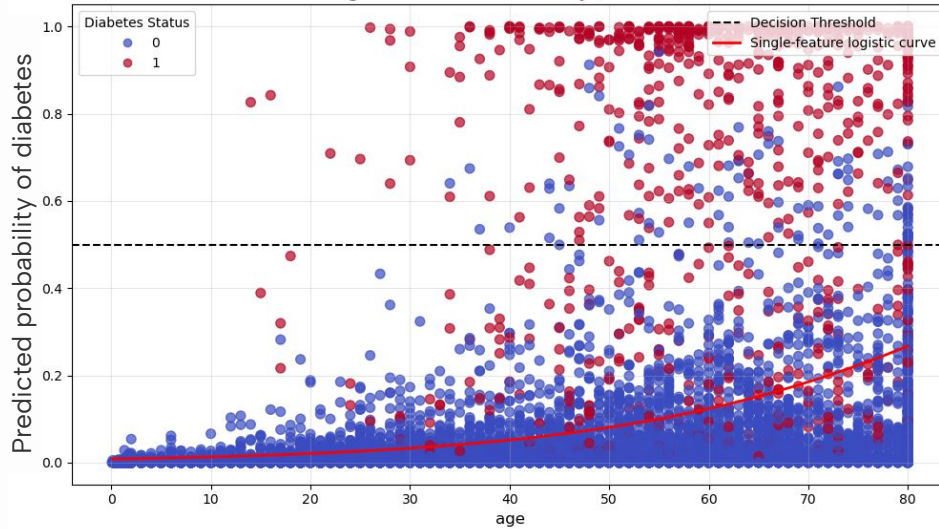ROC Curves for Diabetes Prediction Model

- **Consistent AUC** shows the model generalizes equally well to unseen data.

- **High TPR/FPR trade-off** allows tuning threshold for either fewer false alarms or fewer misses.

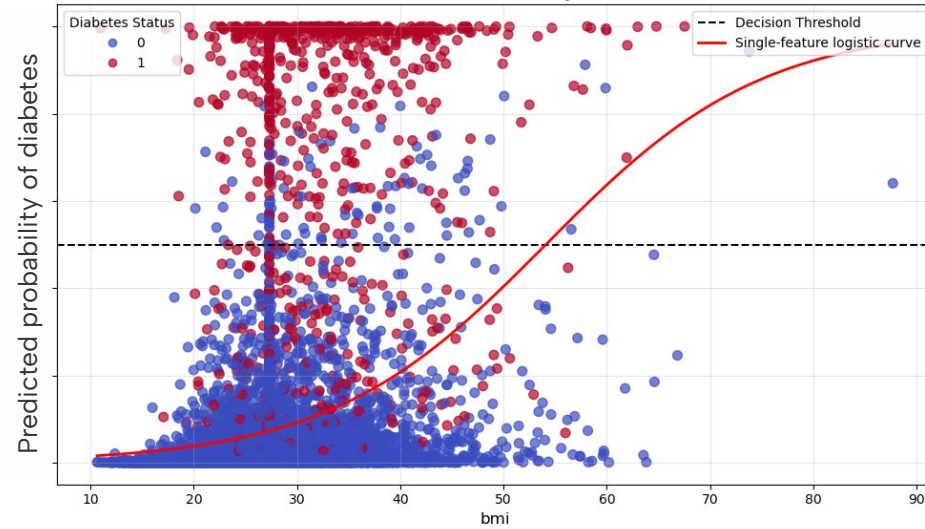Feature Importance for Diabetes Prediction

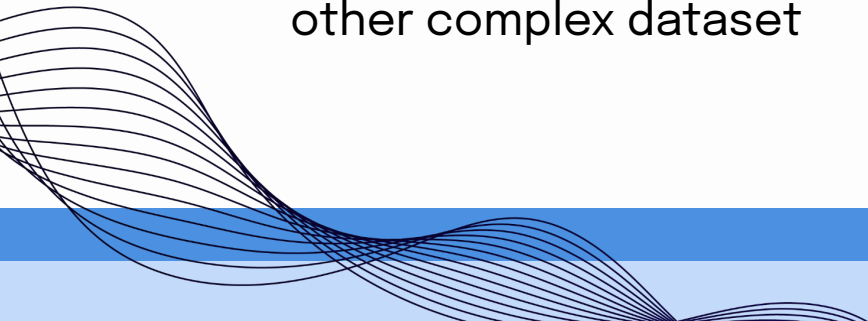# Visualize our logistic curves



Age vs Diabetes (Test set)

BMI vs Diabetes (Test set)

- As age/BMI increases, the possibility of having diabetes will increase.
- Test data points are well-separated

# Future Directions

- Interaction Terms
  - such as HbA1c * age, BMI * hypertension
- Include Other Potential Variables to Increase the Accuracy:
  - other medical conditions
  - family medical history
- Make new models to predict Type I, Type II, disease severity for other complex dataset

# Thanks!

**Do you have any questions?  :)**

# Reference List

- World Health Organization (WHO). (2016). Global report on diabetes. Geneva: World Health Organization. https://www.who.int/publications/i/item/9789241565257

- Centers for Disease Control and Prevention (CDC). (2023). National Diabetes Statistics Report. https://www.cdc.gov/diabetes/data/statistics-report/index.html

- ZIYA. (2023). Diabetes Clinical Dataset (100K Rows). Available on Kaggle: https://www.kaggle.com/datasets/ziya07/diabetes-clinical-dataset100k-rows