



# Diabetes Detection

Using Machine Learning to Make Health Predictions

August 31, 2022

Nashra Khan | Jason Noble | Marie Sanon | Will Wright



# Agenda

- Why this Dataset?
- Visualizing the target and features variables
- Data Cleaning
- Machine Learning Algorithms
- Tuning the Models
- The BEST Model
- Web App Demo



# Why this dataset?

Diabetes is among the most common chronic diseases in the United States impacting millions of Americans each year.

This dataset is the result of a Behavioral Risk Factor Surveillance System (BRFSS) telephone survey that is collected annually by the CDC and was made available on Kaggle.

The CDC estimates that 1 in 5 diabetics, and roughly 8 in 10 pre-diabetics are unaware of their risk. Early diagnosis can lead to lifestyle changes and more effective treatment.

Using this dataset, a predictive model seeks to assist in early detection.

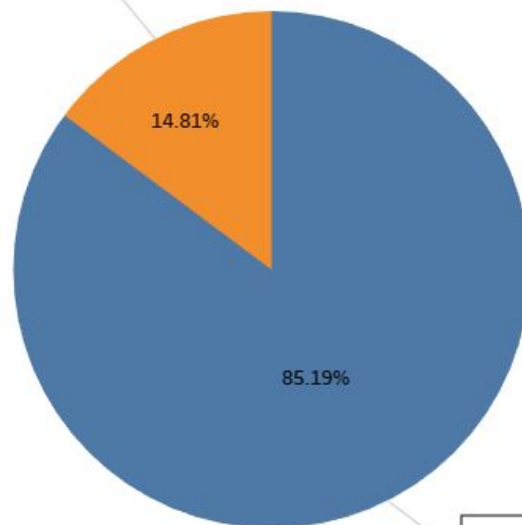
Source: [https://www.kaggle.com/datasets/alexteboul/diabetes-health-indicators-dataset?select=diabetes\\_binary\\_health\\_indicators\\_BRFSS2015.csv](https://www.kaggle.com/datasets/alexteboul/diabetes-health-indicators-dataset?select=diabetes_binary_health_indicators_BRFSS2015.csv)



# ***Tableau Visualizations***

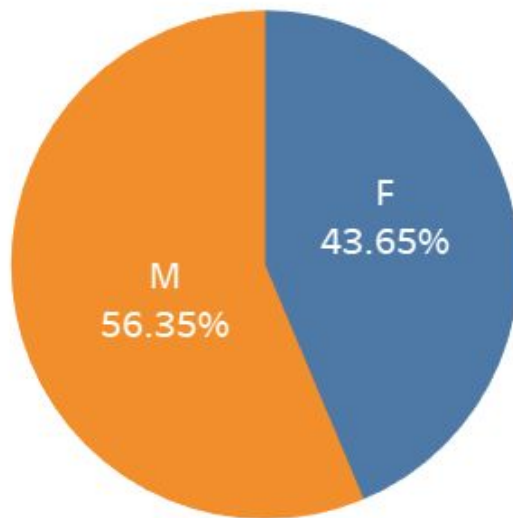
## Target Variable Distribution

Prediabetic or Has Diabetes  
# of Respondents: 33,531

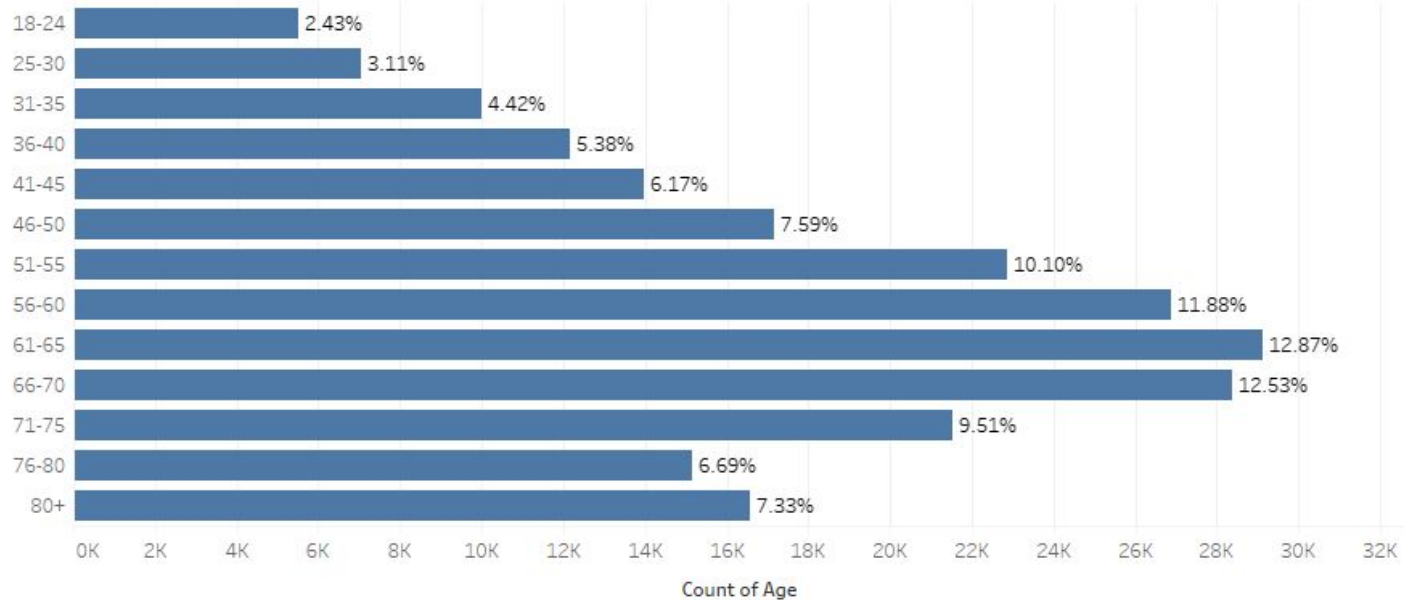


No Diabetes  
# of Respondents: 192,811

Sex Distribution of Respondents

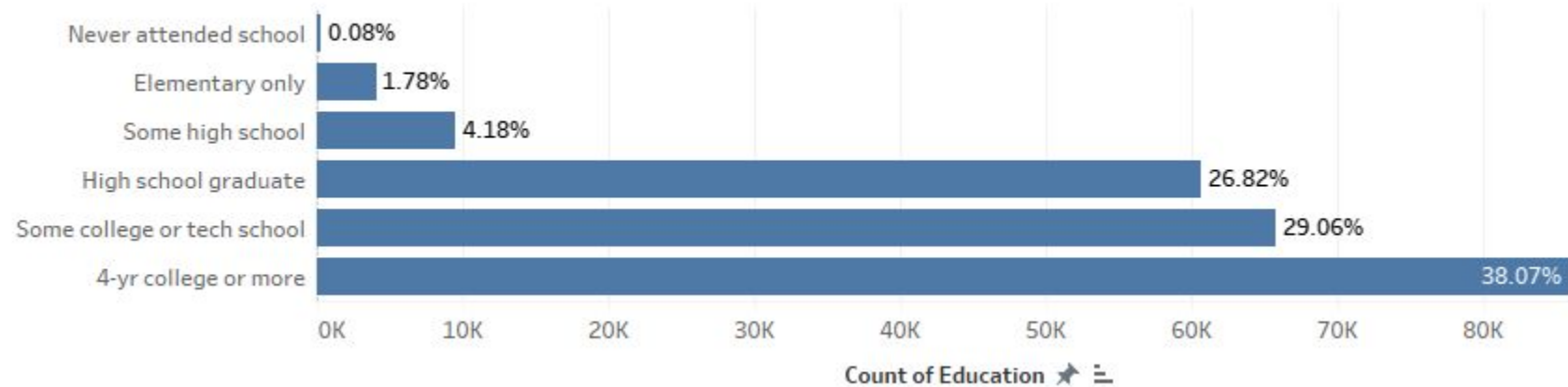


### Age Distribution of Respondents



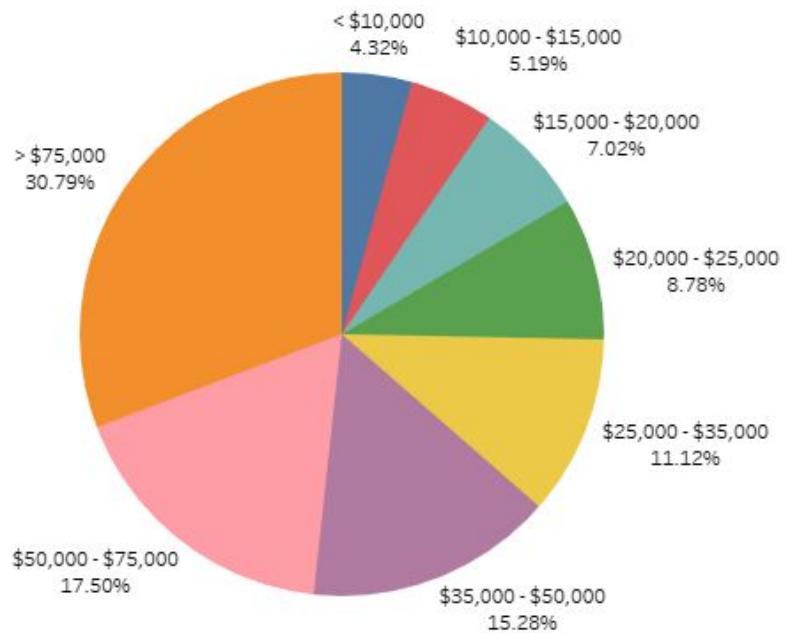


## Education Attainment of Respondents

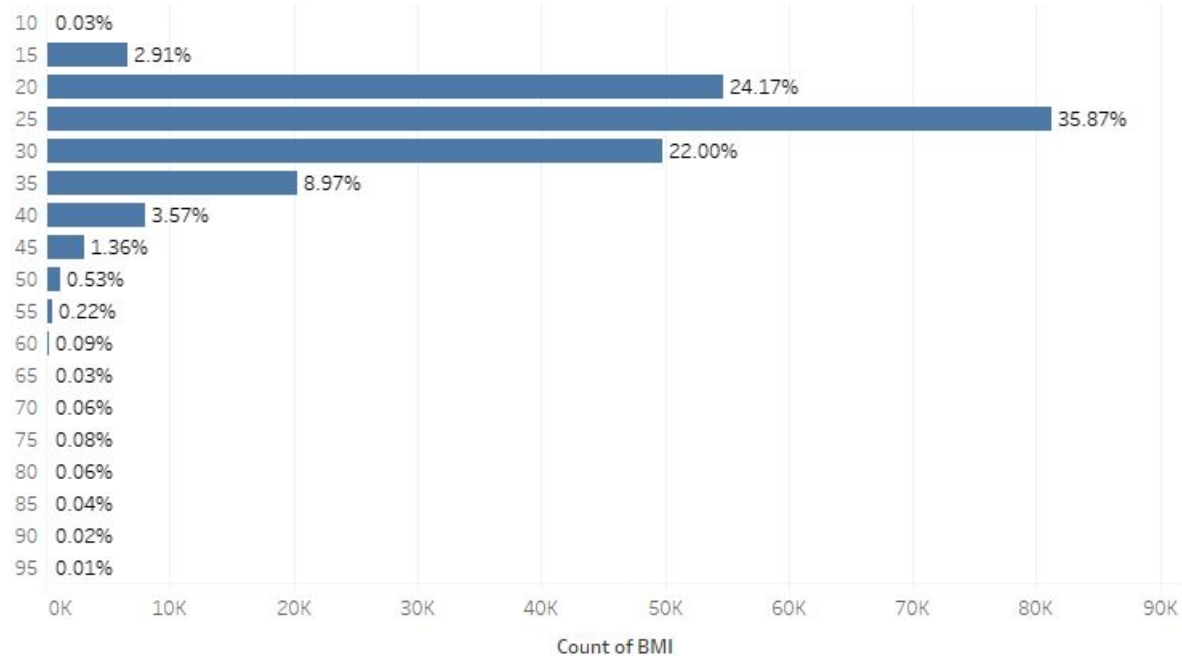




## Income Distribution of Respondents



### BMI Distribution of Respondents





## Data cleaning steps

- Check for nulls and remove rows that include nulls
  - Our dataset did not have null values
- Check for and remove duplicate entries
  - 24,206 duplicate entries were removed (including target variable)
  - 1,556 duplicate entries were removed (excluding target variable)
- 226,342 entries remained after data cleaning



## Algorithms used

- Logistic Regression
- Random Forest Classifier
- AdaBoost Classifier

## Libraries used

- Tableau: visualizations
- Pandas: data manipulation
- Sklearn: ML algorithms
- Pickle / Streamlit: app



# Tuning the models

## Logistic Regression

- Evaluated different values for max\_iter

## Random Forest Classifier

- Evaluated different values for n\_estimators and bootstrap setting

## AdaBoost Classifier

- Evaluated different values for n\_estimators and learning\_rate

## Overall

- Looked at results using dataset with non-binary and non-binned features scaled
- Looked at results using unscaled dataset

```
: 1 def model_tester(model, X, y):
2     X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42)
3     clf = model.fit(X_train, y_train)
4     y_pred = clf.predict(X_test)
5     print(classification_report(y_test, y_pred))
6     print(f'Training Score: {clf.score(X_train, y_train)}')
7     print(f'Testing Score: {clf.score(X_test, y_test)}')

: 1 # Look at different Logistic Regression models and find better performing for further tuning
2 model_tester(LogisticRegression(random_state=42), X, y)
3 model_tester(LogisticRegression(random_state=42, max_iter=500), X, y)
4 model_tester(LogisticRegression(random_state=42, max_iter=1000), X, y)
5 model_tester(LogisticRegression(random_state=42, max_iter=10000), X, y)
```



# The Best Model - AdaBoost Classifier

- Unscaled dataset
- Hyperparameters
  - `n_estimators = 1000`
  - `Learning_rate = 0.1`
- Results
  - Training Score: 0.85691
  - Testing Score: 0.85848
  - Recall: 0.97
- Future Work
  - Use a balanced dataset. Our target variable was unbalanced towards the non-diabetic result (85-15).
  - See if XGBoost performs better than AdaBoost.

# Web App made with Streamlit

## Diabetes Predictor



If you are 18 or older, please answer the following survey questions:

Do you have High Blood Pressure?

Yes

Do you have High Cholesterol?

Yes



**Questions?**