

DBA5106 Group Project

Predictive Modeling to Identify Diabetic Patients

A0231956Y / Gino Martelli Tiu
A0231905L / Susan Koruthu
A0231857Y / Widya Gani Salim
A0231930N / Xhoni Shollaj



Meeting Agenda

- **Introduction**
 - Problem Statement
 - Proposed Solution and Metrics
- **Data Analysis**
 - Datasets Overview
 - Exploratory Data Analysis (EDA)
- **Model Selection**
 - Model Performance
 - Results and Expected Benefits
- **Implementation**
 - Business Implication and Implementation
 - Next Steps

The true costs of diabetes

7th

Leading Cause
of Death

Annual Cost
US\$327
Billion

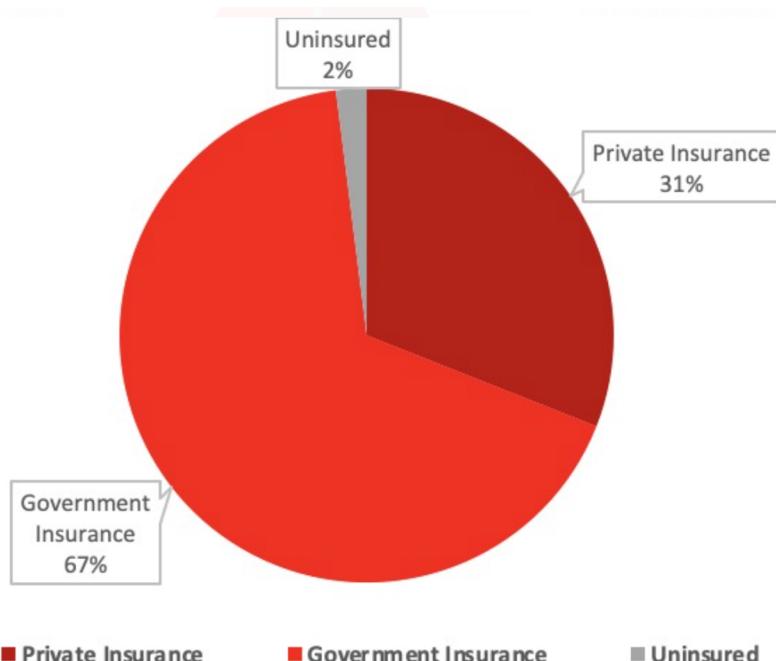
1 in 4

Healthcare Dollar
Expenditure

Who covers these costs? We do.

A costly problem

Distribution of Diabetes Treatment Claims Coverage



Annual cost (per individual)
Up to \$18,652

Existing classification method
Average premium charged
detects only around **30%** of
diabetics cases.

Annual loss (per individual)
\$11,182

What can we do to manage this risk?

Our proposed solution

A classification model with 2 key objectives:

- (1) Compute the probability of diabetes diagnosis given lifestyle & medical test variables
- (2) Identify high risk individuals for monitoring and intervention

Medical and personal data are readily available to us, let's take advantage.

Our proposed performance metric

Expected Value of Diabetes Expenses Claims =

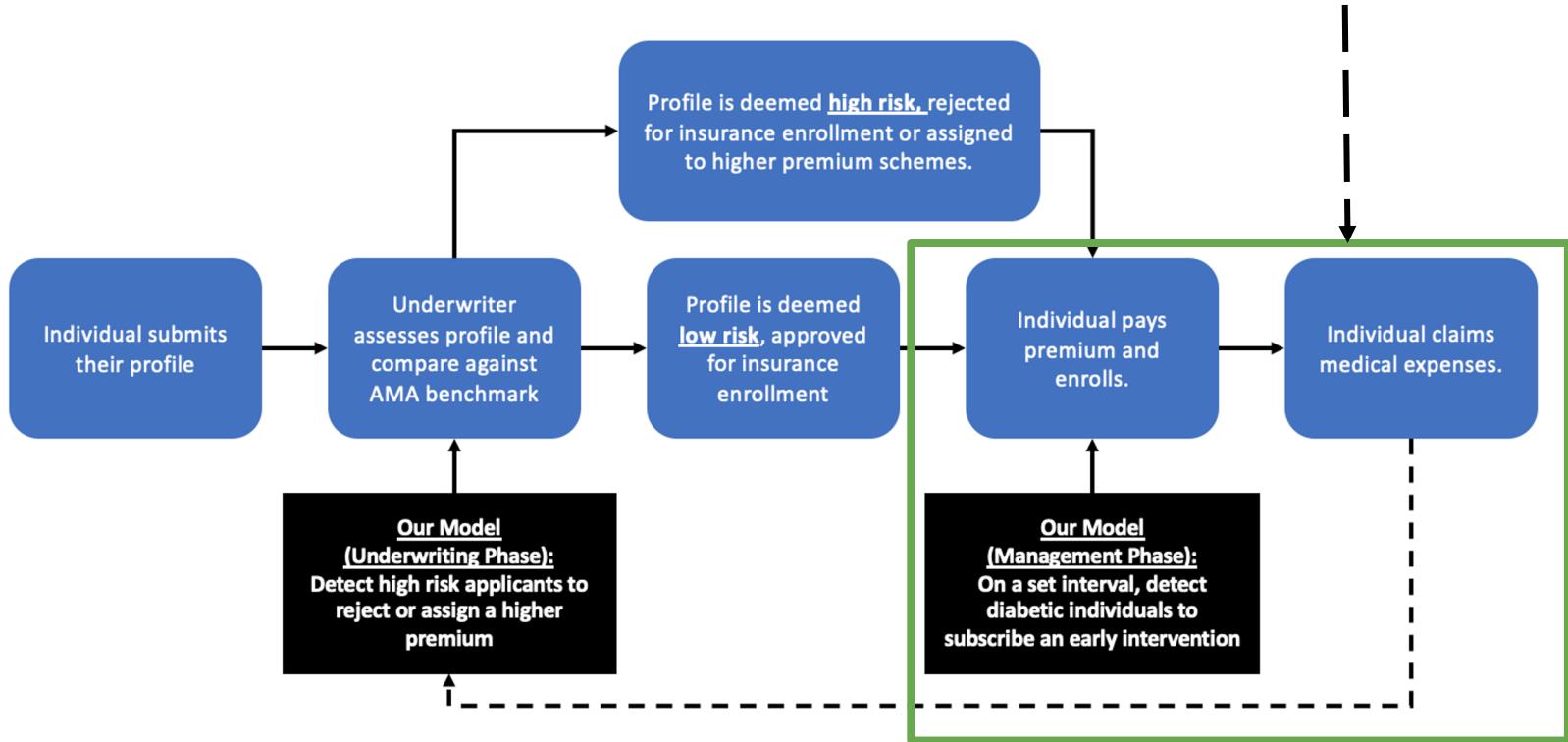
$$\begin{aligned} & P(\text{Diabetic}) \times [P(\text{Predicted Diabetic} | \text{Diabetic}) \times \$16,752 + P(\text{Predicted NonDiabetic} | \text{Diabetic}) \times \$18,652] \\ & + P(\text{NonDiabetic}) \times [P(\text{Predicted NonDiabetic} | \text{NonDiabetic}) \times \$0 + P(\text{Predicted Diabetic} | \text{NonDiabetic}) \times \$200] \end{aligned}$$

Probability	Case	Cost	Remark
$P(\text{Predicted Diabetic} \text{Diabetic})$	True Positive	\$16,572	Base cost of diabetes treatment.
$P(\text{Predicted NonDiabetic} \text{Diabetic})$	False Negative	\$18,652	Base cost + cost associated with complications due to the lack of early treatment (\$1,900).
$P(\text{Predicted NonDiabetic} \text{NonDiabetic})$	True Negative	\$0	No cost.
$P(\text{Predicted Diabetic} \text{NonDiabetic})$	False Positive	\$200	Cost for diabetes screening test.

Proposed metric measures cost saving as compared to current classification method.

Integration with our operations

Trial run should focus on our existing customer first.

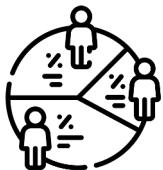


Our existing dataset

4
Tables

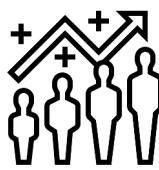
10,175
Observations

83
Original Features



Demographics

1. Gender
 2. Age Group
 3. Race
 4. Education
 5. Household income
- ...



Physical Exam

1. Body Mass Index
2. Blood pressure
3. Height
4. Weight
5. Grip strength

...



Blood Chemistry

1. Glycohemoglobin
2. Potassium
3. Lead
4. Cadmium
5. Insulin

...



Diabetic (y/n)

Identifying key features

Notable Key Features Engineered & Relabeled:

BMI Group

Age Group

Socioeconomic Status

Marital Status

High Risk

High Risk

(based on AMA Guidelines)



Body Mass Index (BMI) ≥ 25



Glycohemoglobin (%) ≥ 5.7

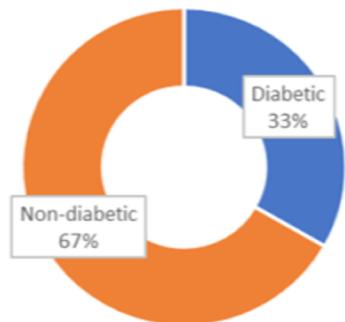


Oral Glucose Tolerance Test Level (mg/dL) $\geq 140^{**}$

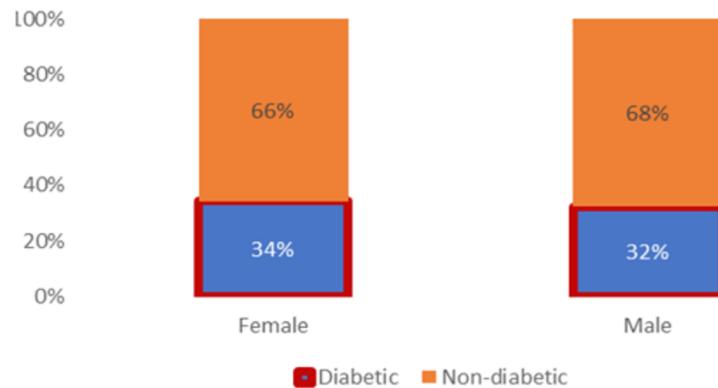
Exploratory data analysis

Focus is to find relevant features that could affect predisposition to diabetes.

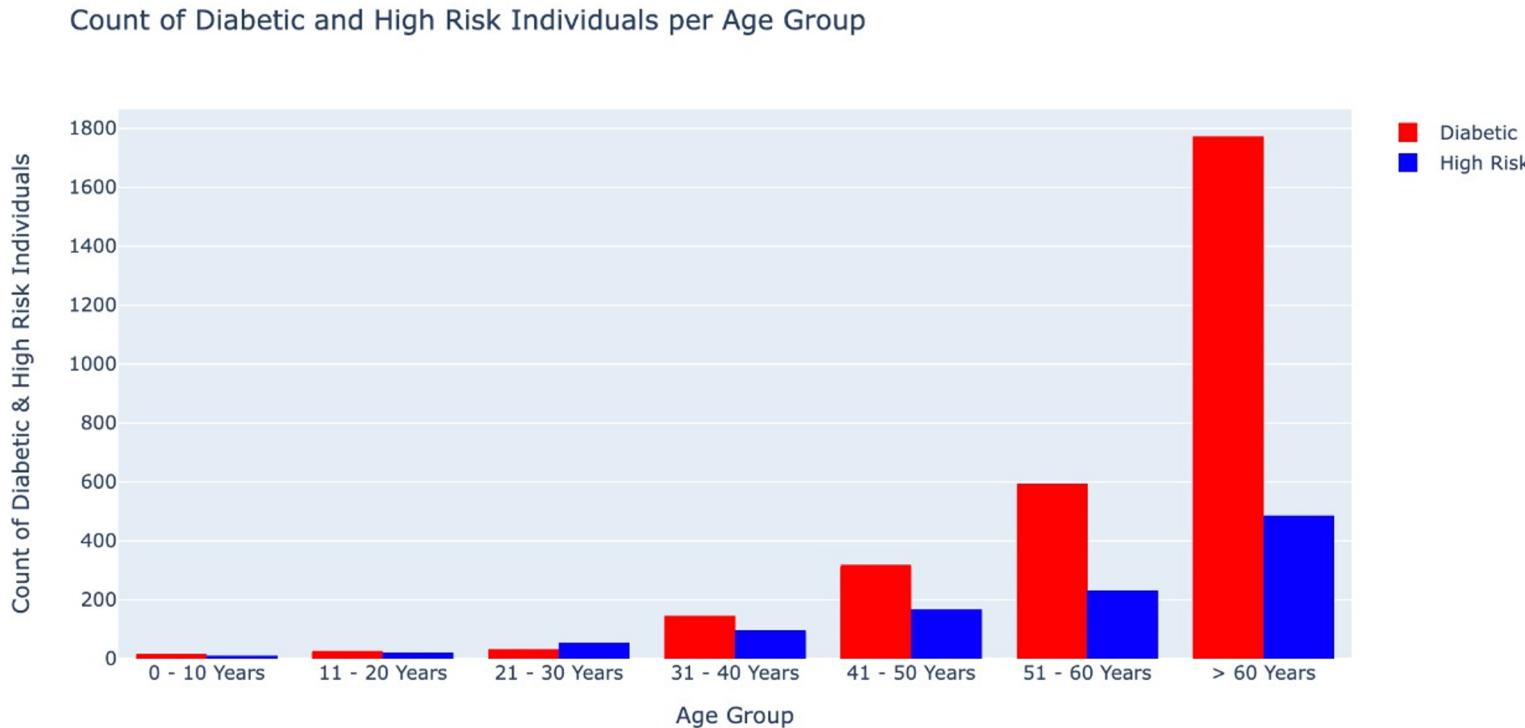
Diabetic : Non-diabetic Ratio



% Diabetics per Gender Segment



Exploratory data analysis

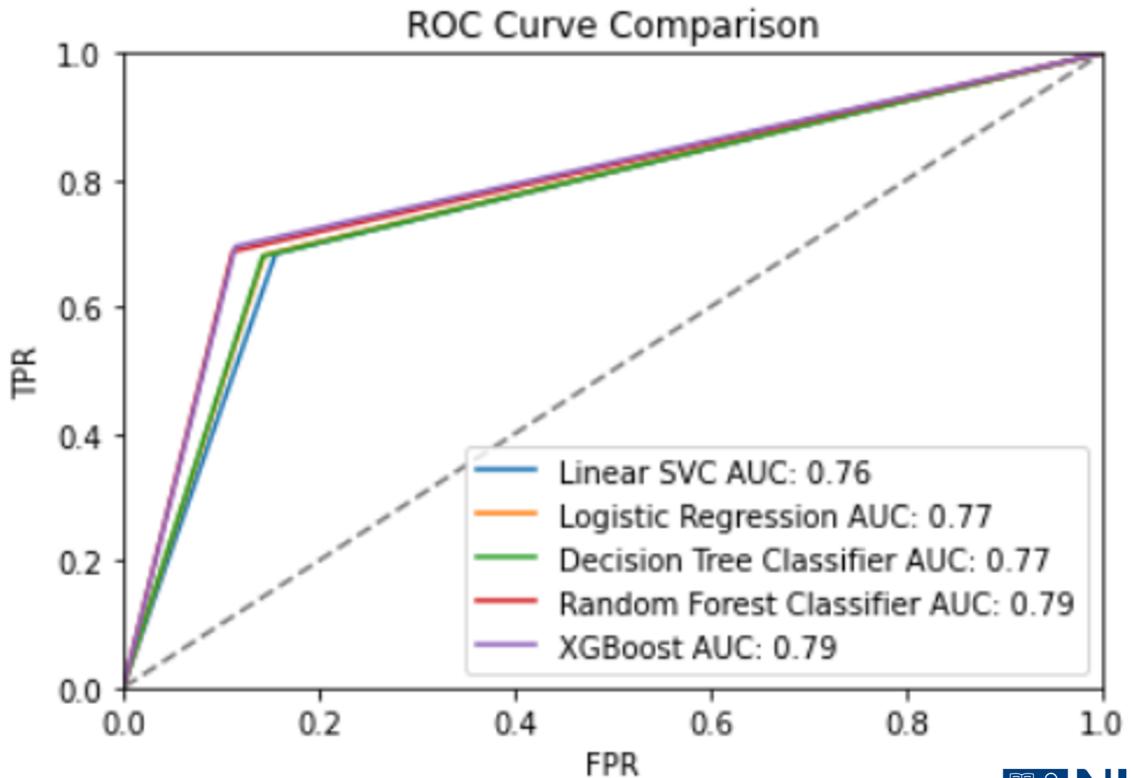


How can we use these insights to manage our risks?

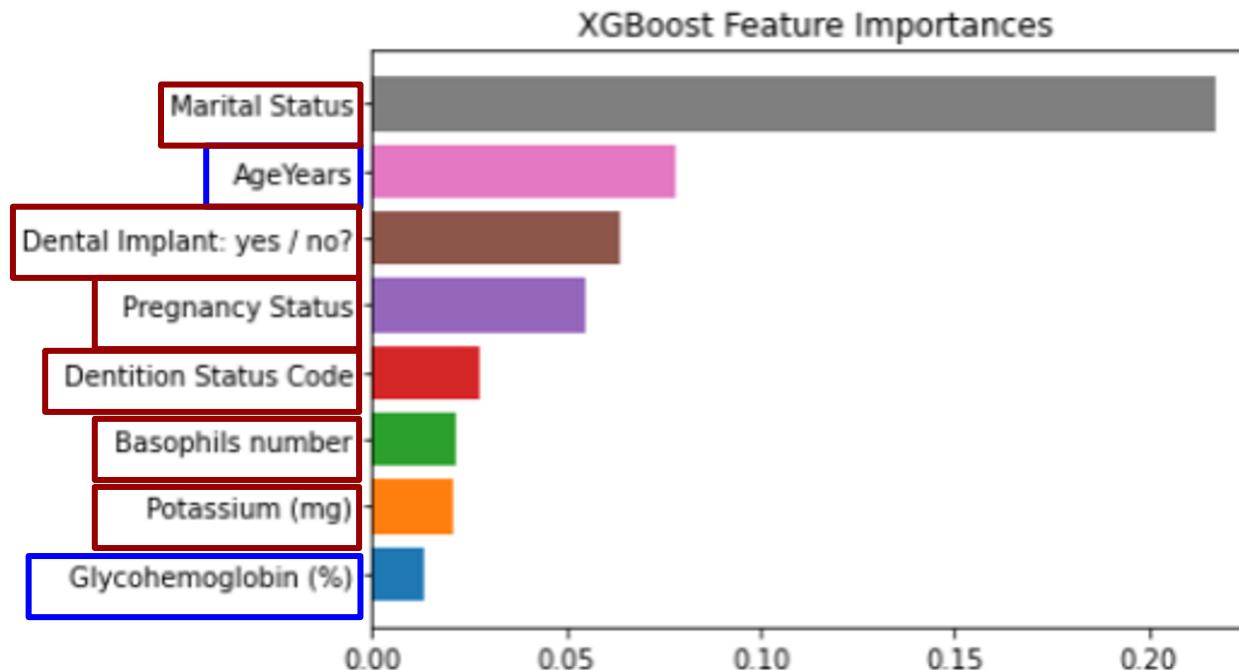
Selected model based on performance

5
Models
Compared

XG Boost
performed best



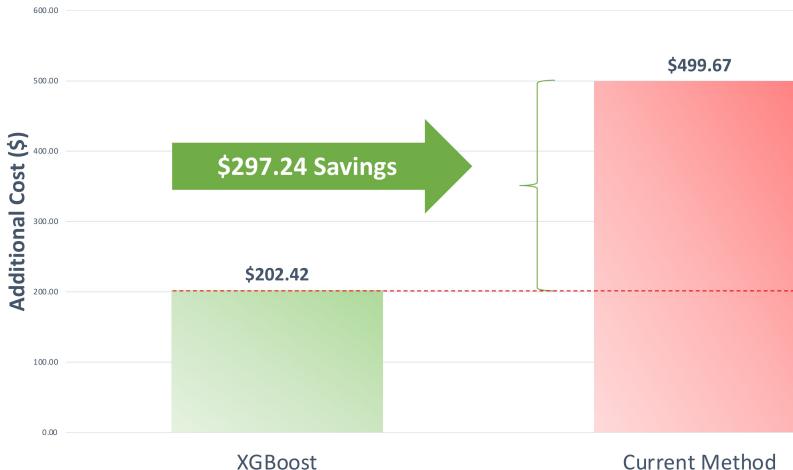
Interesting insights from top features



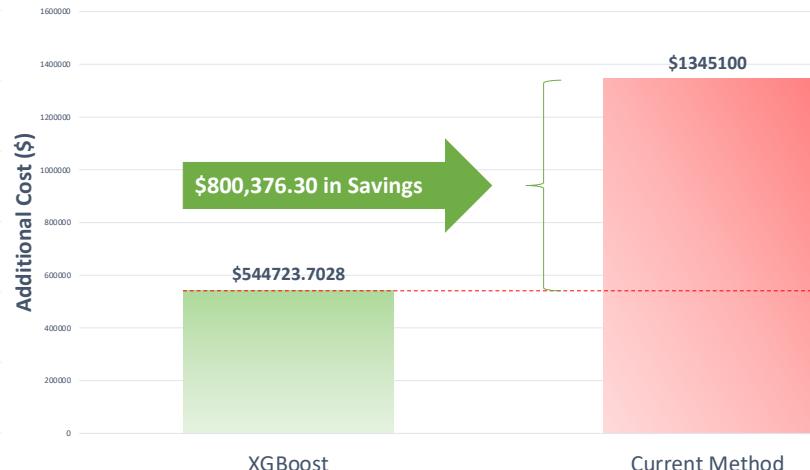
Extracted model features improve on the current detection capability and allow early intervention and better risk management.

Our results and their expected benefit

Additional Cost Per Individual



Additional Costs for Test Sample of 2692 Individuals

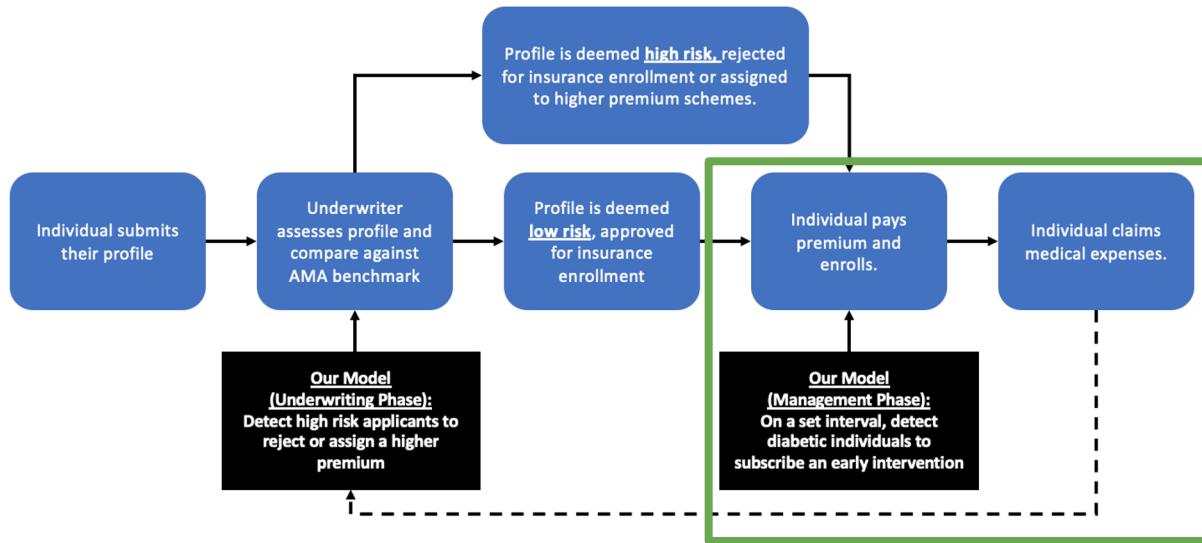


Potential cost saving is 59.49% from early and improved detection vs current industry standard...

Implementation & Next Steps

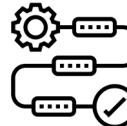
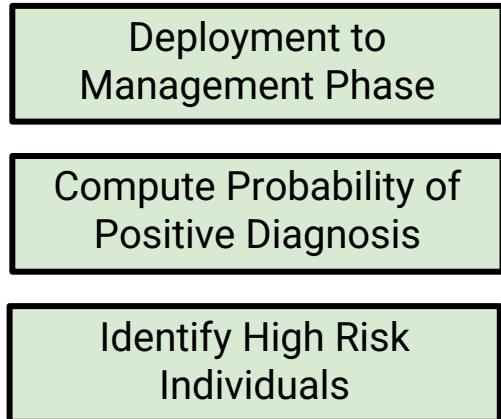
Three steps to up to 59.49% cost-saving:

- Model deployment in management phase.
- Incentive or penalty scheme for those detected as diabetic.
- Track cost-savings as a performance metrics.



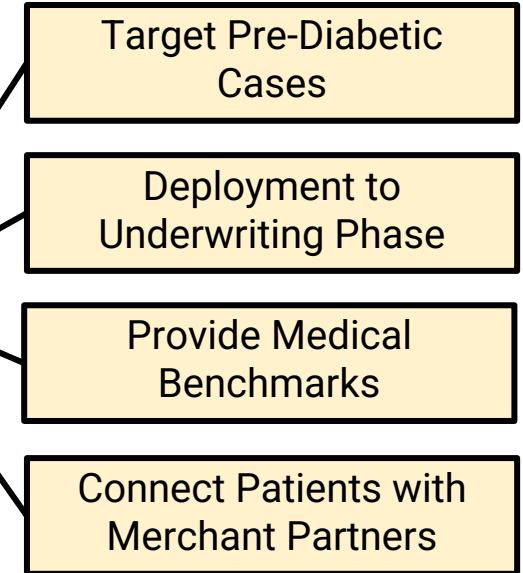
Future Considerations

Phase 1 (Now)



Performance & Business Impact Reviews

Phase 2



The model is scalable and its cost-saving impact can be easily measured.

Thank You

Appendix

References:

Dataset: <https://www.cdc.gov/nchs/nhanes/>

Model Libraries: https://scikit-learn.org/stable/user_guide.html

Facts on Diabetes(1): <https://www.who.int/news-room/fact-sheets/detail/diabetes>

Diabetes costs (2): <https://www.diabetes.org/resources/statistics/cost-diabetes>

CDCP Statistics (3): <https://www.cdc.gov/diabetes/library/features/diabetes-stat-report.html>

NUS MSBA 5106 Lectures