# Diabetes Prediction

YueFeng Xue & Irina Lee

# Motivation

How can machine learning be used to help predict whether people have diabetes or prediabetes?

Symptoms being mild and developing over time

Take Years to diagnostic

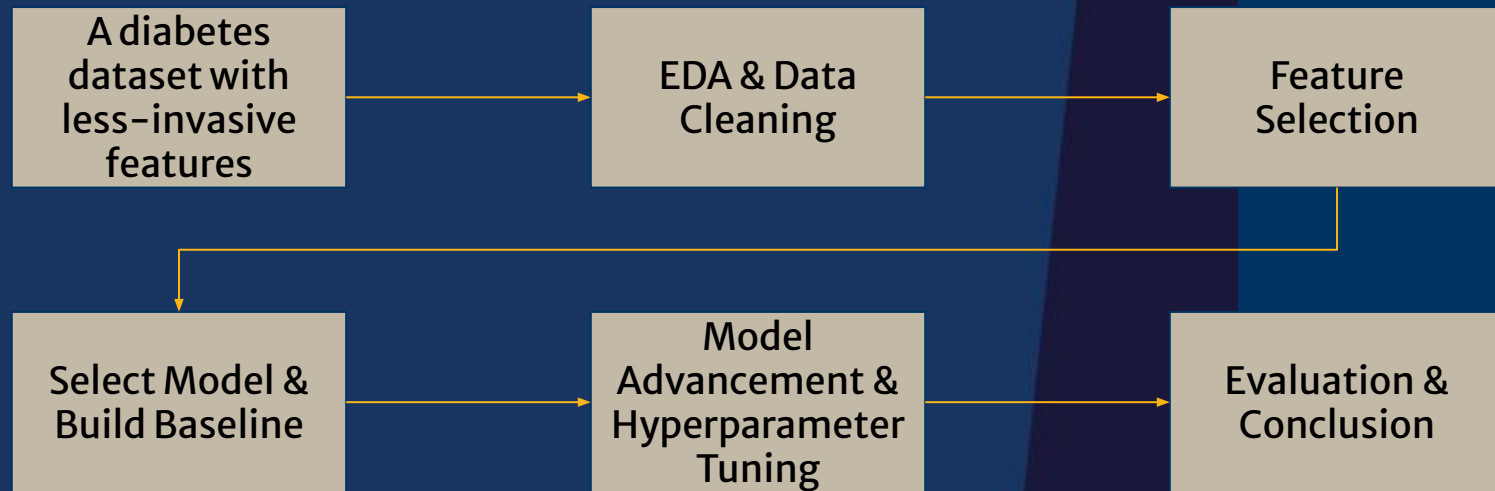#8 Ranked Cause of Death

Berkeley
UNIVERSITY OF CALIFORNIA

# Research Conducted in this Area

Promising Results in Diabetes Prediction Using Machine Learning

Increasing Interests in Predicting Using Non-Invasive Data

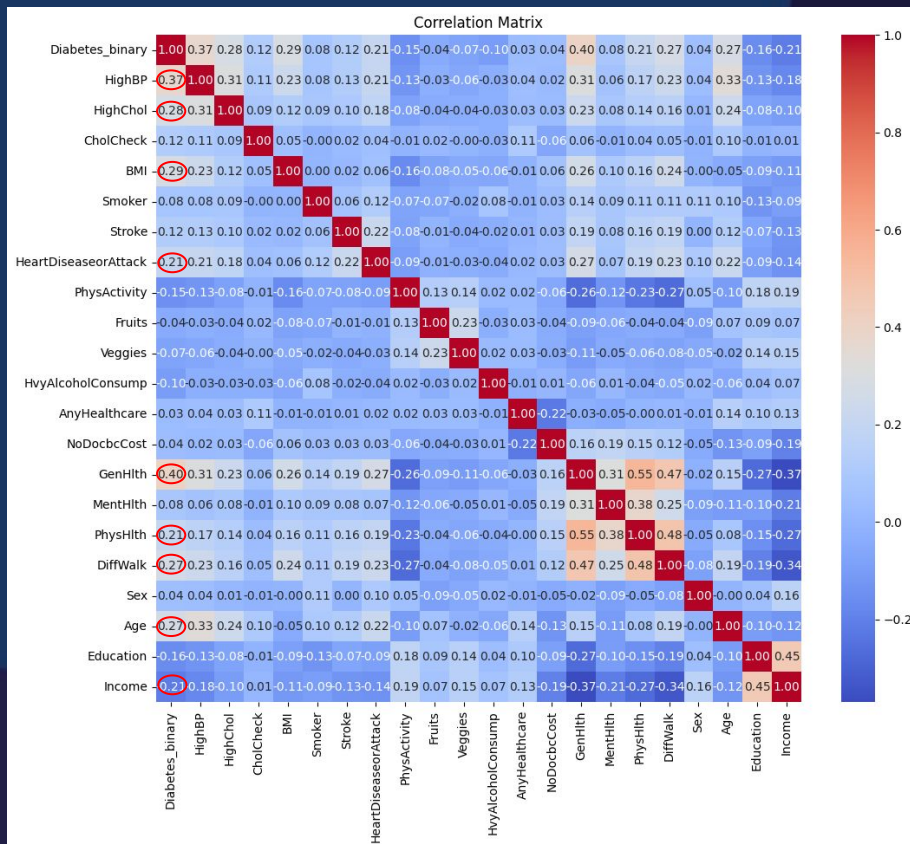Looking for Other Cost & Time Efficient Approach

Berkeley
UNIVERSITY OF CALIFORNIA

# Our Plan

| A diabetes dataset with less-invasive features | → | EDA & Data Cleaning | → | Feature Selection |
|---|---|---|---|---|

| Select Model & Build Baseline | → | Model Advancement & Hyperparameter Tuning | → | Evaluation & Conclusion |
|---|---|---|---|---|

# Data

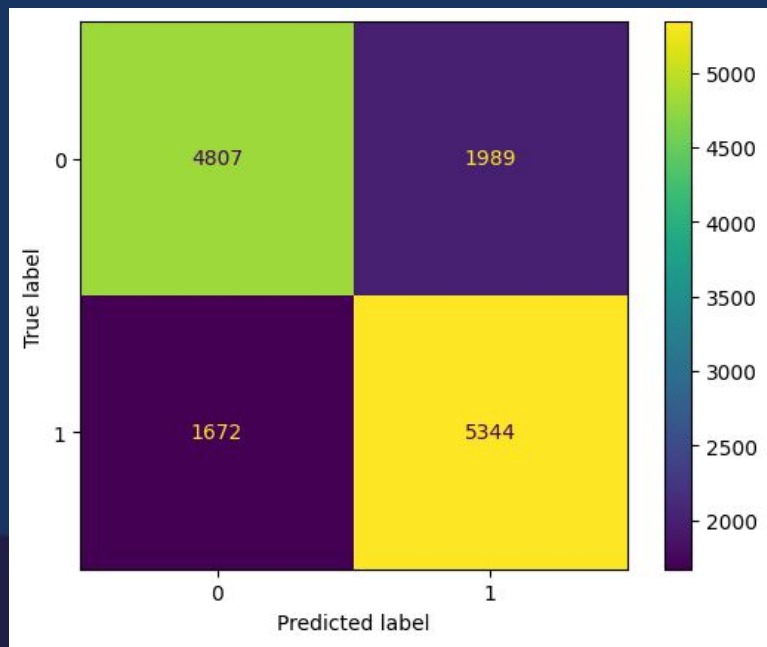| Category | Description |
|---|---|
| Method of Data Collection | Health-related telephone survey that is collected annually by the CDC in 2015 |
| Number of Features | 22 |
| Outcome Feature | Diabetes_binary<br>(0: no diabetes; 1: prediabetes or diabetes) |
| Outcome Feature Balance | Yes, 50/50 Split |
| Missing Values | 0 |
| Total Rows | 70692 |
| Duplicated Rows | 1635 |

Berkeley
UNIVERSITY OF CALIFORNIA

# Correlation Matrix



Correlation Matrix

# Baseline Model:
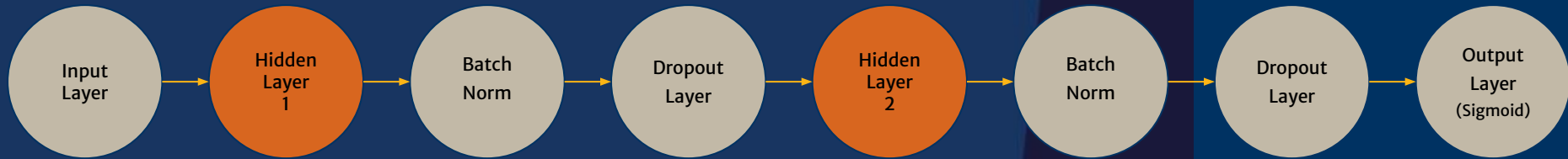# Logistic Regression



Training Accuracy: 0.7401

Validation Accuracy: 0.7479

Testing Accuracy: 0.7349

Test Set - Diabetic & Pre-Diabetic Accuracy: 0.7617

Test Set - Non-Diabetic Accuracy: 0.7073

Berkeley
UNIVERSITY OF CALIFORNIA

# Improved Model
## Neural Network

# Experiment

## Learning Rate: 0.001; Batch Size: 64; Epoch: 200

| Index | Hidden Size | Activation | Optimizer | Dropout Rate | Parameters | Training Accuracy | Validation Accuracy |
|-------|-------------|------------|-----------|--------------|------------|-------------------|---------------------|
| 1 | [64, 32] | relu | SGD | [0.5, 0.3] | 3139 | 0.7402 | 0.7508 |
| 2 | [64, 32] | relu | Adam | [0.5, 0.3] | 9029 | 0.7438 | 0.7503 |
| 3 | [64, 32] | tanh | SGD | [0.5, 0.3] | 3139 | 0.7405 | 0.7500 |
| 4 | [64, 32] | tanh | Adam | [0.5, 0.3] | 9029 | 0.7426 | 0.7503 |
| 5 | [128, 64] | relu | SGD | [0.5, 0.3] | 10371 | 0.7420 | 0.7501 |
| 6 | [128, 64] | relu | Adam | [0.5, 0.3] | 30341 | 0.7447 | 0.7514 |
| 7 | [128, 64] | tanh | SGD | [0.5, 0.3] | 10371 | 0.7407 | 0.7495 |
| 8 | [128, 64] | tanh | Adam | [0.5, 0.3] | 30341 | 0.7444 | 0.7506 |

Berkeley
UNIVERSITY OF CALIFORNIA

# Experiment

## Learning Rate: 0.001; Batch Size: 64; Epoch: 200

| Index | Hidden Size | Activation | Optimizer | Dropout Rate | Parameters | Training Accuracy | Validation Accuracy |
|-------|-------------|------------|-----------|--------------|------------|-------------------|---------------------|
| 9 | [128, 64] | relu | Adam | [0.3, 0.3] | 30341 | 0.7443 | 0.7506 |
| 10 | [128, 64] | relu | Adam | [0.4, 0.4] | 30341 | 0.7450 | 0.7516 |
| 11 | [128, 64] | relu | Adam | [0.5, 0.5] | 30341 | 0.7447 | 0.7511 |
| 12 | [128, 64] | relu | Adam | [0.4, 0.3] | 30341 | 0.7456 | 0.7517 |
| 13 | [128, 64] | relu | Adam | [0.3, 0.4] | 30341 | 0.7446 | 0.7524 |
| 6 | [128, 64] | relu | Adam | [0.5, 0.3] | 30341 | 0.7442 | 0.7514 |
| 14 | [128, 64] | relu | Adam | [0.3, 0.5] | 30341 | 0.7449 | 0.7514 |
| 15 | [128, 64] | relu | Adam | [0.4, 0.5] | 30341 | 0.7449 | 0.7520 |
| 16 | [128, 64] | relu | Adam | [0.5, 0.4] | 30341 | 0.7447 | 0.7521 |

# Experiment

## Hidden Size: [128, 64]; Activation: relu; Optimizer: Adam; Dropout: [0.3, 0.4]

| Index | Batch Size | Learning Rate | Epoch | Parameters | Training Accuracy | Validation Accuracy |
|-------|-----------|---------------|-------|------------|-------------------|---------------------|
| 17 | 128 | 0.001 | 200 | 30341 | 0.7442 | 0.7501 |
| 13 | 64 | 0.001 | 200 | 30341 | 0.7446 | 0.7524 |
| 18 | 32 | 0.001 | 200 | 30341 | 0.7447 | 0.7505 |
| 19 | 128 | 0.0001 | 200 | 30341 | 0.7467 | 0.7508 |
| 20 | 64 | 0.0001 | 200 | 30341 | 0.7450 | 0.7507 |
| 21 | 32 | 0.0001 | 200 | 30341 | 0.7458 | 0.7500 |
| 22 | 128 | 0.01 | 200 | 30341 | 0.7438 | 0.7530 |
| 23 | 64 | 0.01 | 200 | 30341 | 0.7436 | 0.7500 |
| 24 | 32 | 0.01 | 200 | 30341 | 0.7435 | 0.7501 |

# Conclusion - Result
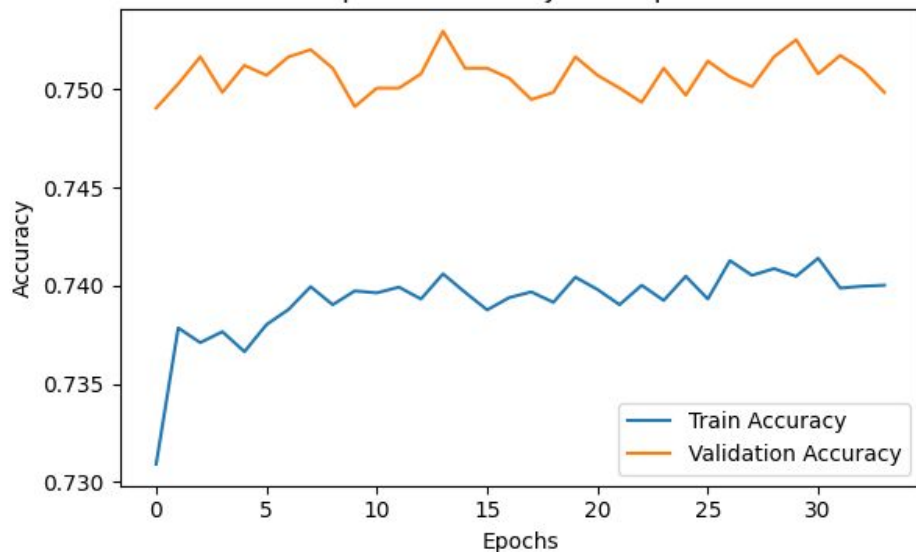


Training Accuracy: 0.7438

Validation Accuracy: 0.7530

Testing Accuracy: 0.418

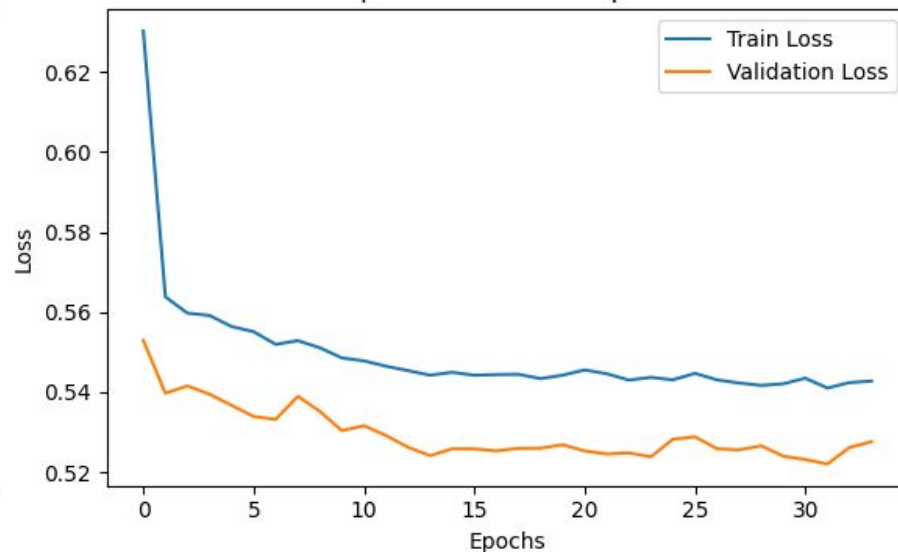Test Set - Diabetic & Pre-Diabetic Accuracy: 0.8033

Test Set - Non-Diabetic Accuracy: 0.6783

# Conclusion - Result

# Conclusion - Avenue for Future

Considering more advanced model, for example: XGBoost, etc.

Exploring better hyperparameter tuning approaches, such as grid search

Looking for better datasets with potential for more feature engineering

Berkeley
UNIVERSITY OF CALIFORNIA

# Link to Github Repository

https://github.com/yuefengxue/mids-w207-final-YueFeng-Xue-Irina-Lee



Berkeley
UNIVERSITY OF CALIFORNIA

# Contribution

**YueFeng Xue:** Data preparation, data processing, EDA, modeling, evaluation & conclusion, design and draft PowerPoint, PowerPoint preparation, code clean-up and organize.

**Irina Lee:** Topic research and suggestion, data research and making a decision on what data we are using, modeling, evaluation & conclusion, code review, PowerPoint review and preparation.

Berkeley
UNIVERSITY OF CALIFORNIA

# Thank You!

Berkeley
UNIVERSITY OF CALIFORNIA