

Diabetes Risk Prediction

Empowering Early Detection and Management Within Children



Alex Gerwer

Overview

This groundbreaking project confronts the escalating global health crisis of diabetes by harnessing the power of advanced machine learning. We have developed sophisticated predictive models specifically designed for the early identification of diabetes risk, with a crucial focus on safeguarding the health of younger generations. By intelligently analyzing readily accessible health indicators—such as age, BMI, blood pressure, lifestyle choices, and relevant medical history—our approach provides a powerful tool for proactive healthcare. This isn't just about prediction; it's about empowerment. Our rigorously validated models enable timely, personalized interventions and foster effective preventative strategies long before the disease manifests fully. This initiative represents a significant leap forward, aiming to shift the paradigm from reactive treatment to proactive prevention, ultimately securing better long-term health outcomes and enhancing the well-being of our youth.

Methodology: From Vision to Validation

Our project progressed through a rigorous, iterative methodology, translating strategic planning into demonstrable results. Each phase built systematically on the previous, driven by empirical evidence and analytical refinement.

Laying a Robust Data Foundation: We began by evaluating three distinct dataset variations provided. After careful analysis comparing ternary vs. binary target encodings and dataset completeness, we strategically selected the most comprehensive file (binary target, full dataset) as the optimal foundation. Following this selection, meticulous data preparation prioritized integrity, employing sophisticated KNN imputation to address outliers without discarding valuable information, and establishing crucial performance baselines.

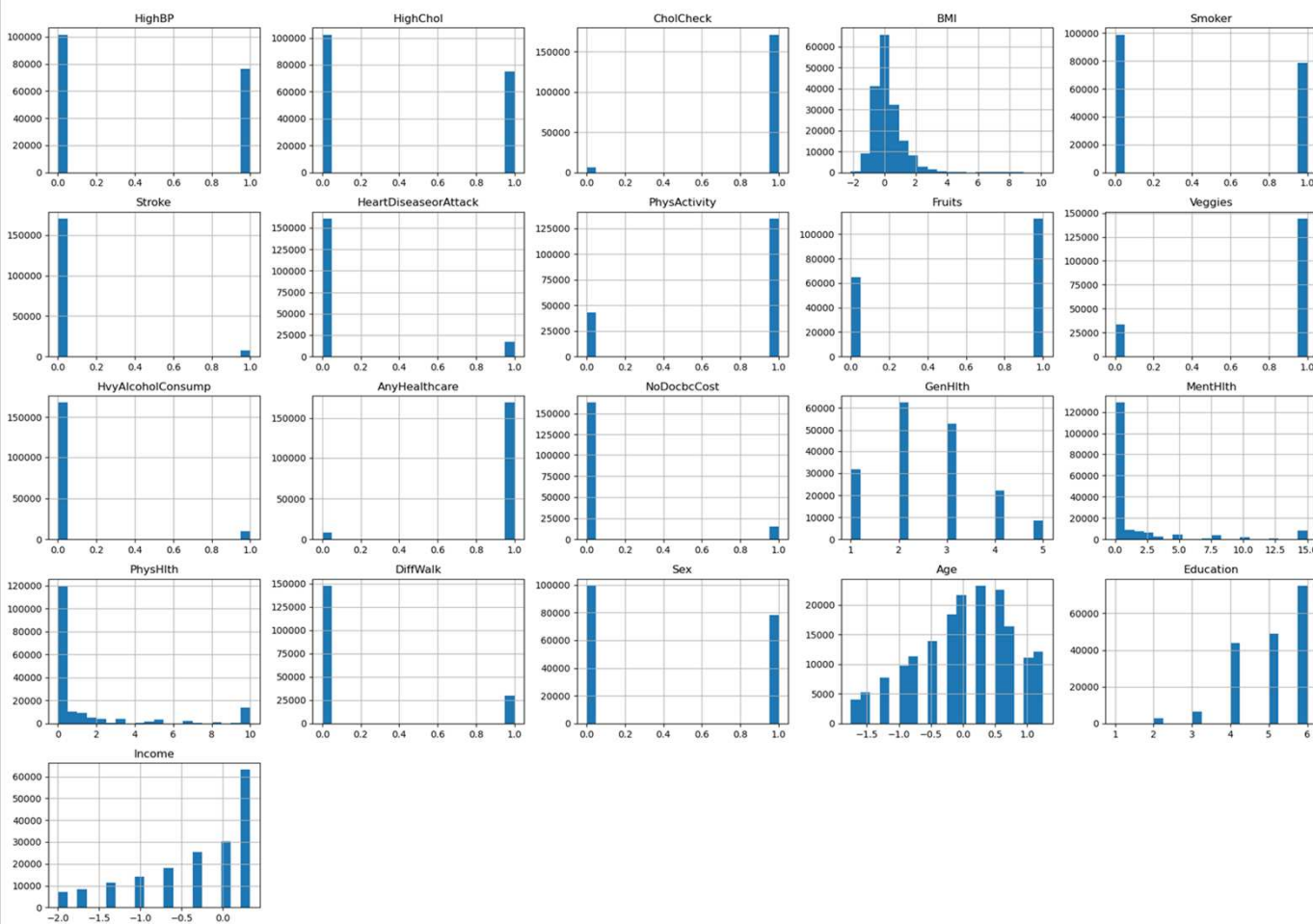
Harnessing EDA for Strategic Refinement: Exploratory Data Analysis was central, moving beyond visualization to actively guide strategy using the chosen dataset. Insights from EDA directly informed targeted feature engineering and a disciplined feature selection process, including rigorous VIF analysis to mitigate multicollinearity and enhance model efficiency.

Systematic Model Development & Optimization: Building on this solid groundwork, we systematically explored and evaluated a diverse range of machine learning algorithms. This comparative approach, coupled with iterative optimization, ensured the identification of models offering the highest predictive accuracy and reliability for identifying diabetes risk within the selected data context.

This disciplined, data-centric journey, starting with deliberate dataset selection, culminated in the development of demonstrably effective and impactful predictive models, showcasing our commitment to delivering robust, actionable healthcare insights.

Raw Data Set:

Histograms of Features



Features Name	Description
High Blood Pressure (BP)	0 = No High BP, 1 = High BP
High Cholesterol (Chol)	0 = No High Cholesterol, 1 = High Cholesterol
Cholesterol Check	0 = No, 1 = Yes
Body Mass Index (BMI)	Numerical value representing weight adjusted for height
Smoker	0 = No (has not smoked 100+ cigarettes in lifetime), 1 = Yes
Stroke	0 = No, 1 = Yes (has had a stroke)
Heart Disease or Attack	0 = No (no coronary heart disease or myocardial infarction), 1 = Yes
Physical Activity	0 = No physical activity in past 30 days (excluding work), 1 = Yes
Fruit Consumption	0 = No, 1 = Consumes fruits 1 or more times per day
Vegetable Consumption	0 = No, 1 = Consumes vegetables 1 or more times per day
Heavy Alcohol Consumption	0 = No, 1 = Meets criteria for heavy drinking
Health Insurance Coverage	0 = No, 1 = Has any health insurance/prepaid plan
Difficulty Affording Healthcare	0 = No, 1 = Faced difficulty affording healthcare in the past year
General Health Perception	Scale 1-5 (1 = Excellent, 5 = Poor)
Mental Health Days	Number of days in the past 30 with poor mental health (scale 1-30)
Physical Health Days	Number of days in the past 30 with poor physical health (scale 1-30)
Difficulty Walking	0 = No, 1 = Has serious difficulty walking/climbing stairs
Sex	0 = Female, 1 = Male
Age Category	13-level category (18-29 to 80+)

Data Preparation

- **Outlier Management (Focus: Data Preservation; Version 1):**

- Evaluated various methods (Capping, Removal, Imputation).
- Selected K-Nearest Neighbors (KNN) Imputation: Replace outliers with values based on the average of the nearest neighbors of the outlier data to address outliers in numerical features (like BMI, Age).
- Contextual Handling: Ordinal features (e.g., GenHlth, Education) were excluded from numerical outlier techniques as IQR/KNN are inappropriate for their scale.

- **Feature Selection and Engineering (Version 2):**

- Used pair plots with Kernel Density Estimate plots off-diagonal and histograms on diagonal for feature selection and engineering.

- **Feature Density Handling (Version 3):**

- Used violin plots to further understand feature distributions (modality and densities) and identify outliers.
- Indicated need for the use of One Hot Encoding.

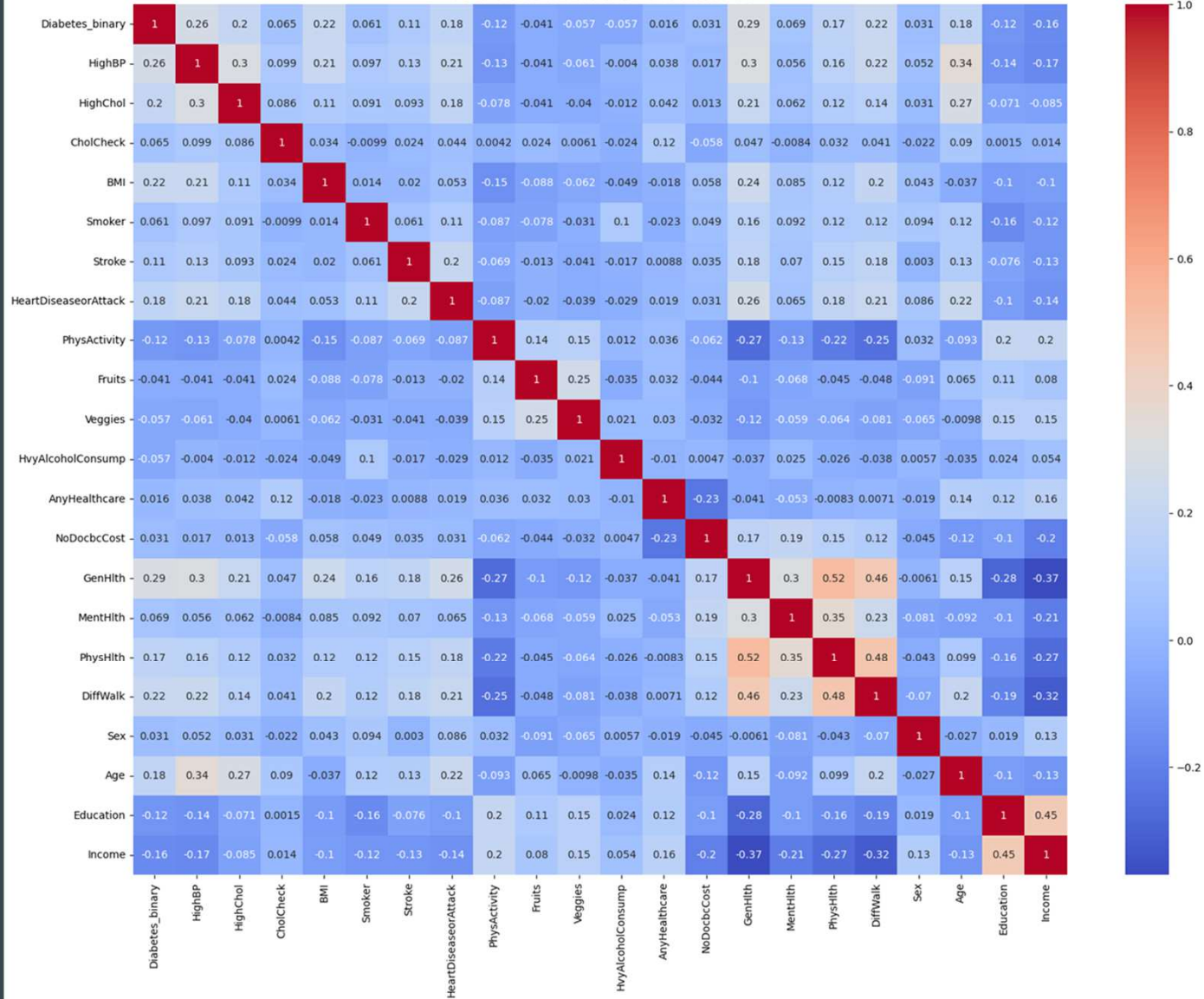
- **Feature Scaling (Focus: Consistency & Compatibility):**

- Experimented with RobustScaler vs. StandardScaler (Version 4).
- Adopted StandardScaler for feature normalization.

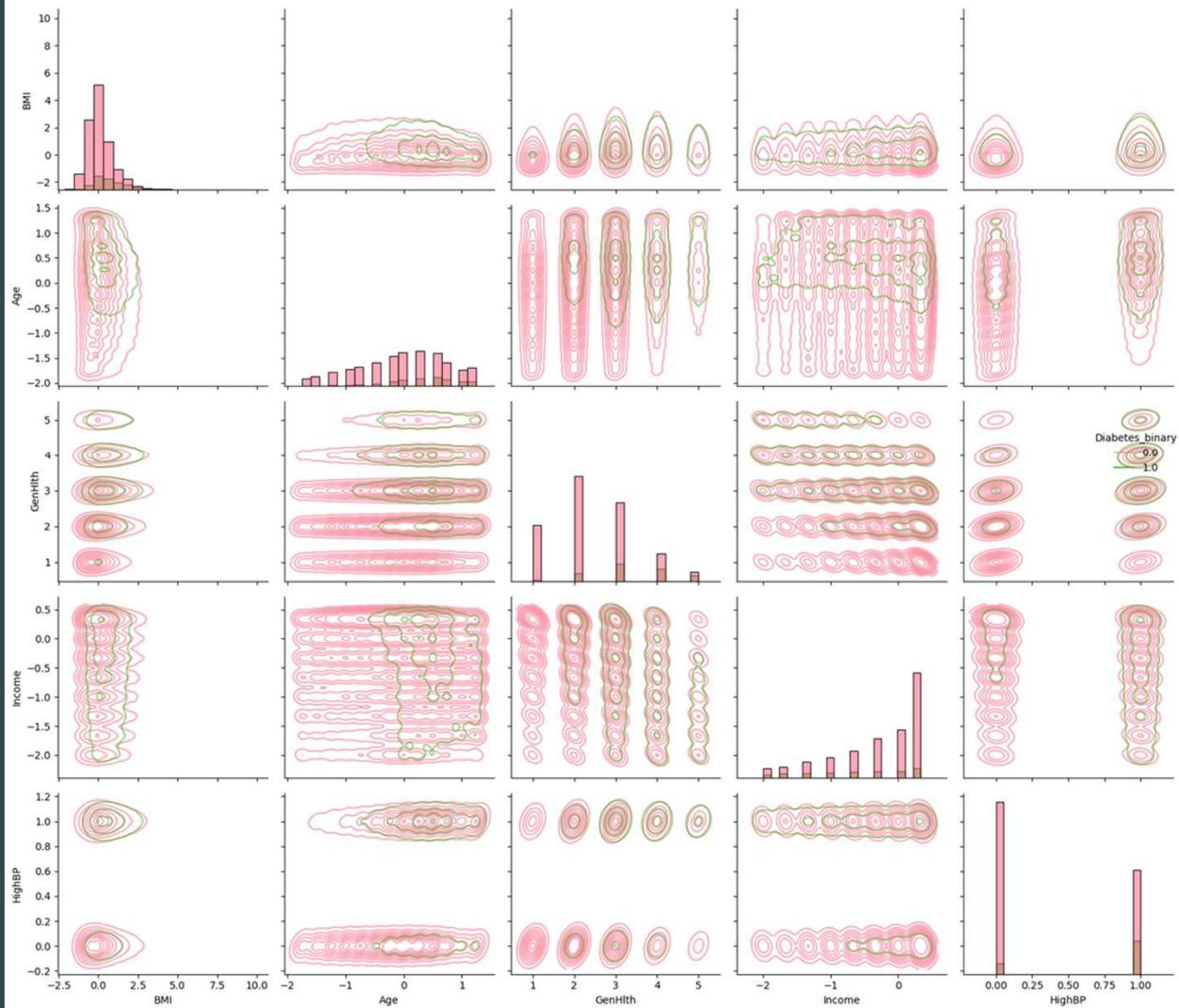
- **Multicollinearity Reduction (Focus: Model Stability):**

- Initial correlation analysis considered (Version 5).
- Primarily used Variance Inflation Factor (VIF) analysis for diagnosis.
- Performed an iterative feature removal process based on high VIF scores (Versions 5, onwards).

Correlation Matrix of Features



Pair Plots with KDE Off-Diagonal and Histograms Diagonal



Feature Selection and Feature Engineering Using Pair Plots

Diagonal:

- **BMI:** Shows a right-skewed distribution indicating higher BMI values are more common in diabetic individuals. There's a clear visual separation.
- **Age:** Also shows a right-skewed distribution the separation isn't as pronounced as BMI, but shows distinct age groups.
- **GenHlth:** The non-diabetic group is heavily concentrated in the "Excellent" and "Very Good" categories (lower numerical values), while the diabetic group has a much higher proportion in "Fair" and "Poor" categories (higher numerical values). This feature shows strong discriminatory power.
- **Income:** Similar to GenHlth, Non-diabetic individuals are more prevalent in higher income brackets, and diabetic individuals are more represented in lower income brackets. This feature also seems informative.
- **HighBP:** Non-diabetic group is heavily skewed towards HighBP=0 (no high blood pressure), while the diabetic group has a significant proportion with HighBP=1 (high blood pressure). This is a strong indicator.

Off-Diagonal:

- **BMI vs. Age:** The KDE plot shows that the highest density for non-diabetic individuals is at lower BMI and younger ages (bottom left), indicating that the combination of higher BMI and older age significantly increases diabetes likelihood.
- **BMI vs. GenHlth:** Non-diabetic individuals are concentrated at lower BMI and better general health (lower GenHlth numerical values), confirming the combined influence of these factors.
- **GenHlth vs. HighBP:** Strong separation. Excellent/Very Good health is strongly associated with no HighBP, while Fair/Poor health is strongly associated with HighBP, especially for the diabetic group.

Logistic Regression (Versions 6-9):

The Logistic Regression model served as an interpretable baseline. Initial optimization via GridSearchCV (Versions 6-7) showed limited gains, suggesting inherent limitations of the linear model for this dataset. Feature importances were directly extracted from coefficients (Version 8), and further analysis using Calibration Curves and Odds Ratios (Version 9) provided deeper insights into its probabilistic predictions and feature influence, though its predictive ceiling was apparent (AUC-ROC ~0.824).

To make the GridSearchCV optimization for Logistic Regression as thorough as practically possible within the code, we expanded the `param_grid_lr` to explore a wider and finer range of hyperparameter values. The focus was on the `C` parameter (regularization strength) and also included different penalty options for LogisticRegression with the 'liblinear' solver.

Hyperparameter tuning was performed using GridSearchCV and, while thorough, GridSearchCV didn't significantly boost the metrics for Logistic Regression. This suggests that further optimization of the Logistic Regression model itself might be limited, and the issue might lie elsewhere.

Random Forest (Version 10):

The Random Forest model was introduced to capture non-linearities. Optimization using RandomizedSearchCV (balancing thoroughness and efficiency) provided improved performance over the baseline Logistic Regression (AUC-ROC ~0.803), demonstrating the value of tree-based ensembles.

Efficient Random Forest Optimization: The core change is using RandomizedSearchCV instead of GridSearchCV for the Random Forest. Instead of trying every combination of hyperparameters (like GridSearchCV), RandomizedSearchCV samples a specified number of combinations (n_iter) randomly from the provided distributions. This is much faster, especially with large search spaces.

- **Balanced n_iter:** The n_iter parameter is the primary control over the trade-off between search thoroughness and computation time. 20 iterations is a good starting point; you can adjust it based on your available time and the specific problem.
- **Strategic Parameter Grid:** Even with randomized search, a well-chosen param_grid_rf is important. The ranges provided cover the most important hyperparameters and typical values that work well.
- **Reduced CV Folds:** Using cv=5 for the Random Forest further speeds up the process without sacrificing too much accuracy in the performance estimate.

Gradient Boosting (Versions 9 [initial], 11-12):

This is a machine learning technique for regression and classification problems, which produces a prediction model in the form of an ensemble of weak prediction models, typically decision trees. This model emerged as a top performer. Initial exploration (Version 9 baseline) showed promise. Systematic optimization using GridSearchCV (Version 11) and further refinement incorporating subsample, max_features, and crucially, early stopping (Version 12, detailed in Notes), proved highly effective. Version 11 achieved the highest AUC-ROC (0.828541) and F1-Score (0.837105), indicating a strong ability to discriminate and balance precision/recall. Version 12 maintained strong performance while potentially improving generalization and efficiency through early stopping. The optimization process took around two hours, with very little improvement.

Base Model

Accuracy: 0.8675102491327656

Confusion Matrix:

```
[[42795  944]
 [ 5778 1219]]
```

Classification Report:

	precision	recall	f1-score	support
0.0	0.88	0.98	0.93	43739
1.0	0.56	0.17	0.27	6997
accuracy			0.87	50736
macro avg	0.72	0.58	0.60	50736
weighted avg	0.84	0.87	0.84	50736

Optimized Model

Accuracy: 0.8678650268054242

Confusion Matrix:

```
[[42817  922]
 [ 5782 1215]]
```

Classification Report:

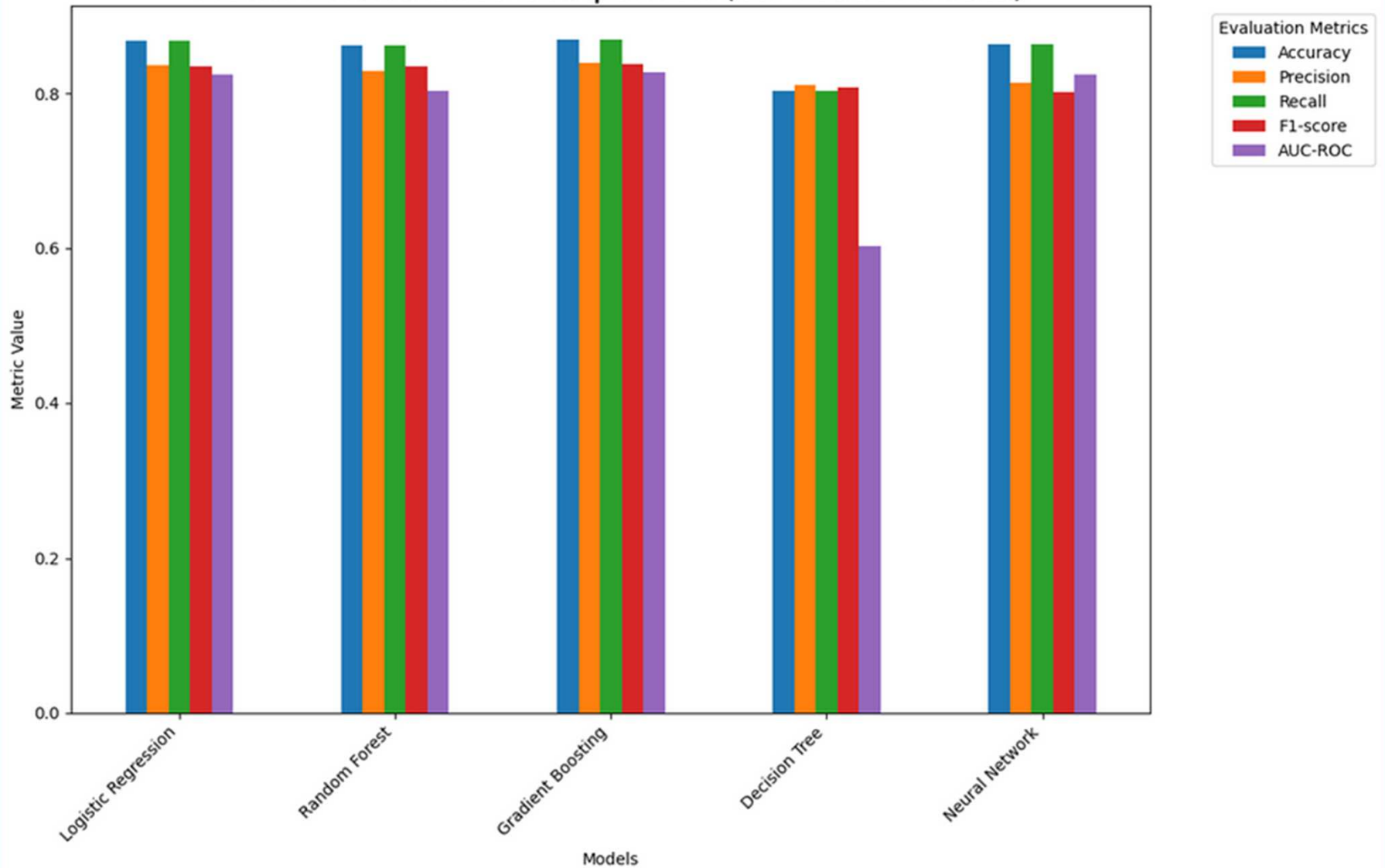
	precision	recall	f1-score	support
0.0	0.88	0.98	0.93	43739
1.0	0.57	0.17	0.27	6997
accuracy			0.87	50736
macro avg	0.72	0.58	0.60	50736
weighted avg	0.84	0.87	0.84	50736

Neural Network (Versions 12 [initial], 13):

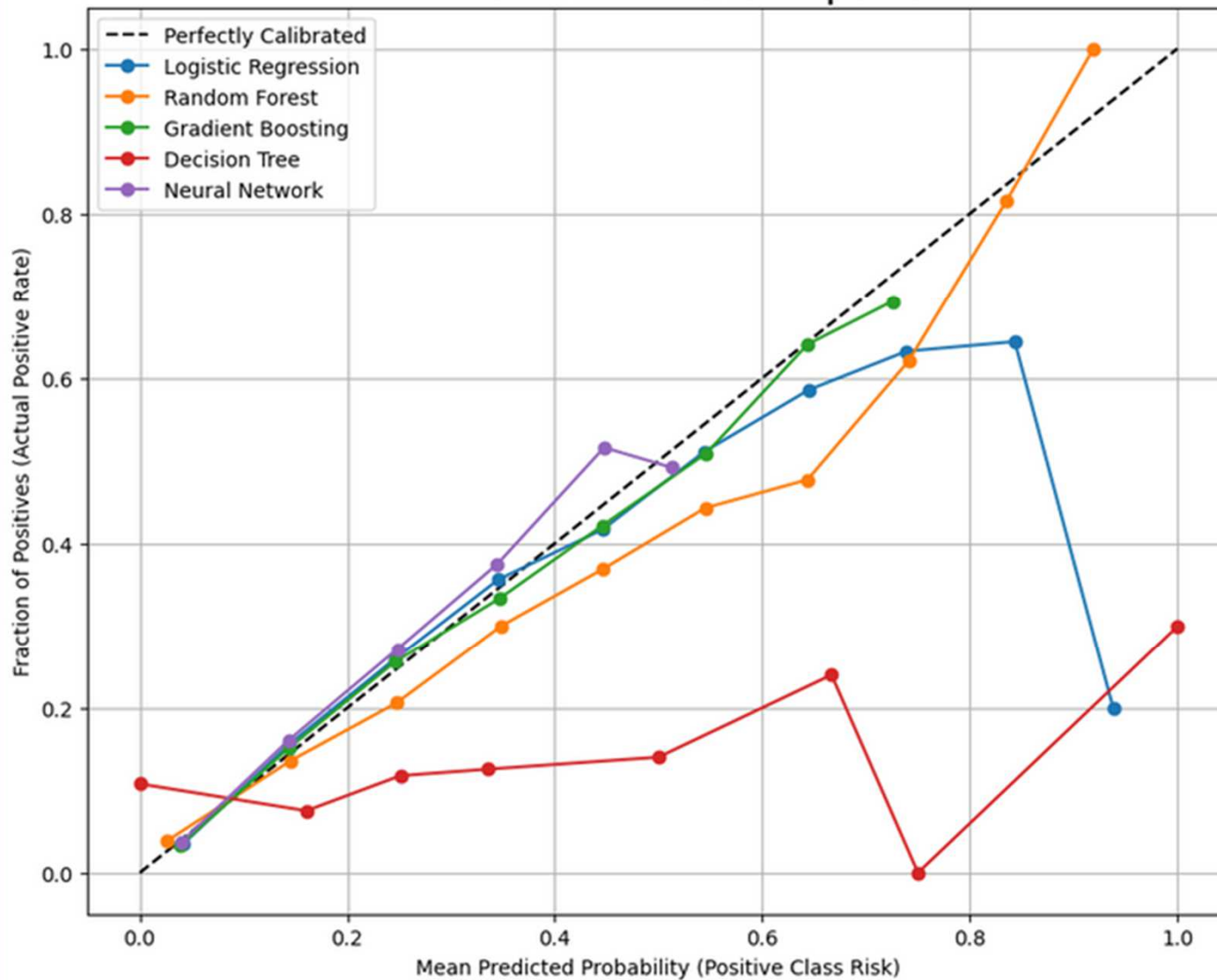
The Neural Network model allowed for the exploration of the potential of deep learning. Optimization using Keras Tuner (RandomSearch) with L2 regularization and early stopping (Version 13, detailed in Notes) was implemented. While achieving respectable results (AUC-ROC ~0.8246), it did not surpass the optimized Gradient Boosting models in this instance, potentially requiring further architectural tuning or larger data scales to unlock its full potential relative to the highly effective boosting methods.

- **keras-tuner Integration:** The core addition is the use of keras-tuner's RandomSearch to optimize the Neural Network's hyperparameters. This is much more efficient than manually trying different combinations.
- **build_nn_model(hp) Function:** This function defines the model architecture, but parameterized by the hp object (hyperparameter space) provided by keras-tuner. This is how keras-tuner explores different configurations.
- **L2 Regularization:** Added `kernel_regularizer=l2(0.01)` to the Dense layers to help prevent overfitting.
- **Output Layer:** `activation='softmax'` is used for multi-class classification. The number of units is `y_train.shape[1]`, which corresponds to the number of classes
- **Compilation:** `loss='categorical_crossentropy'` is used for multi-class classification.
- **RandomSearch**

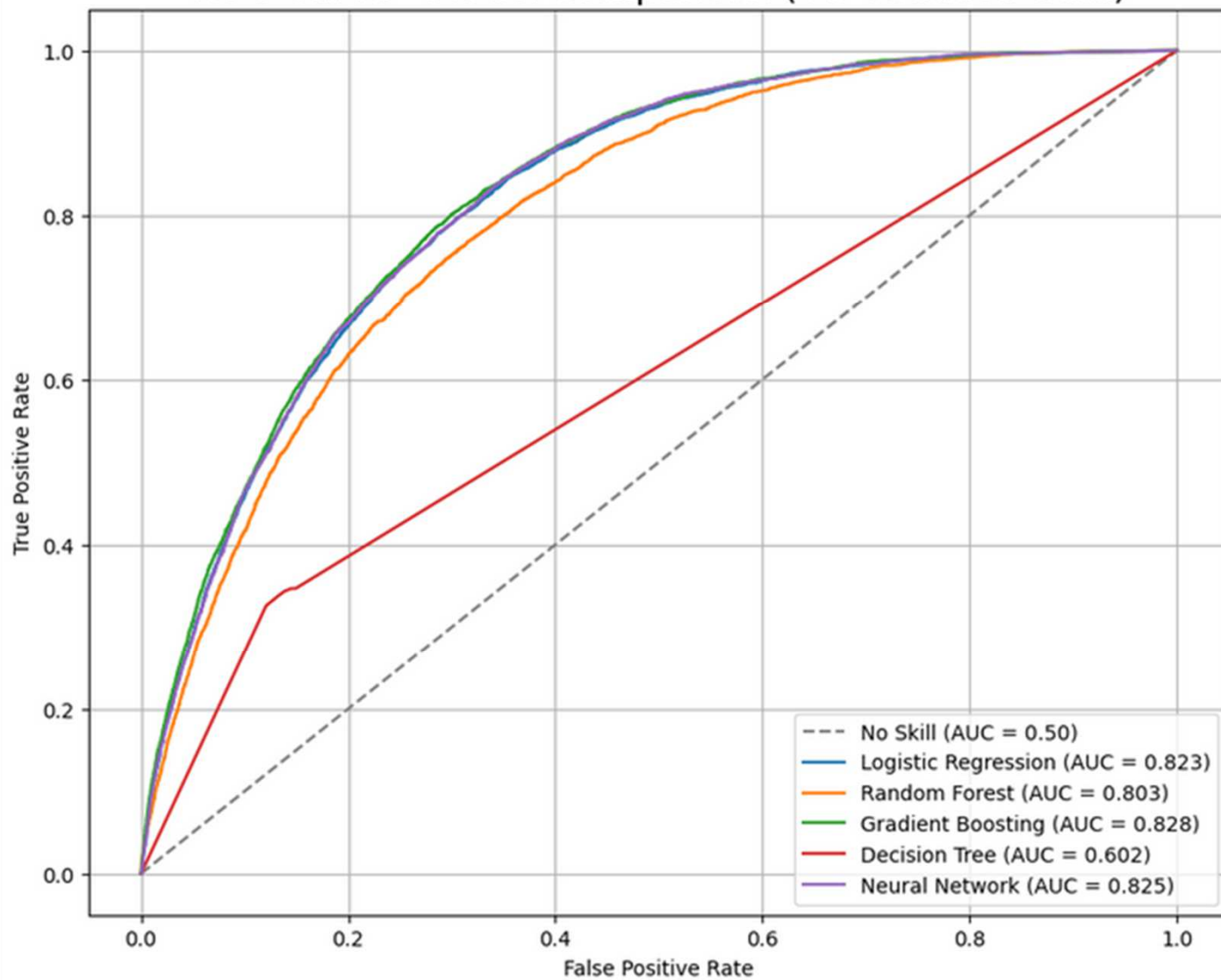
Model Performance Comparison (VIF-reduced data)



Calibration Curves Comparison



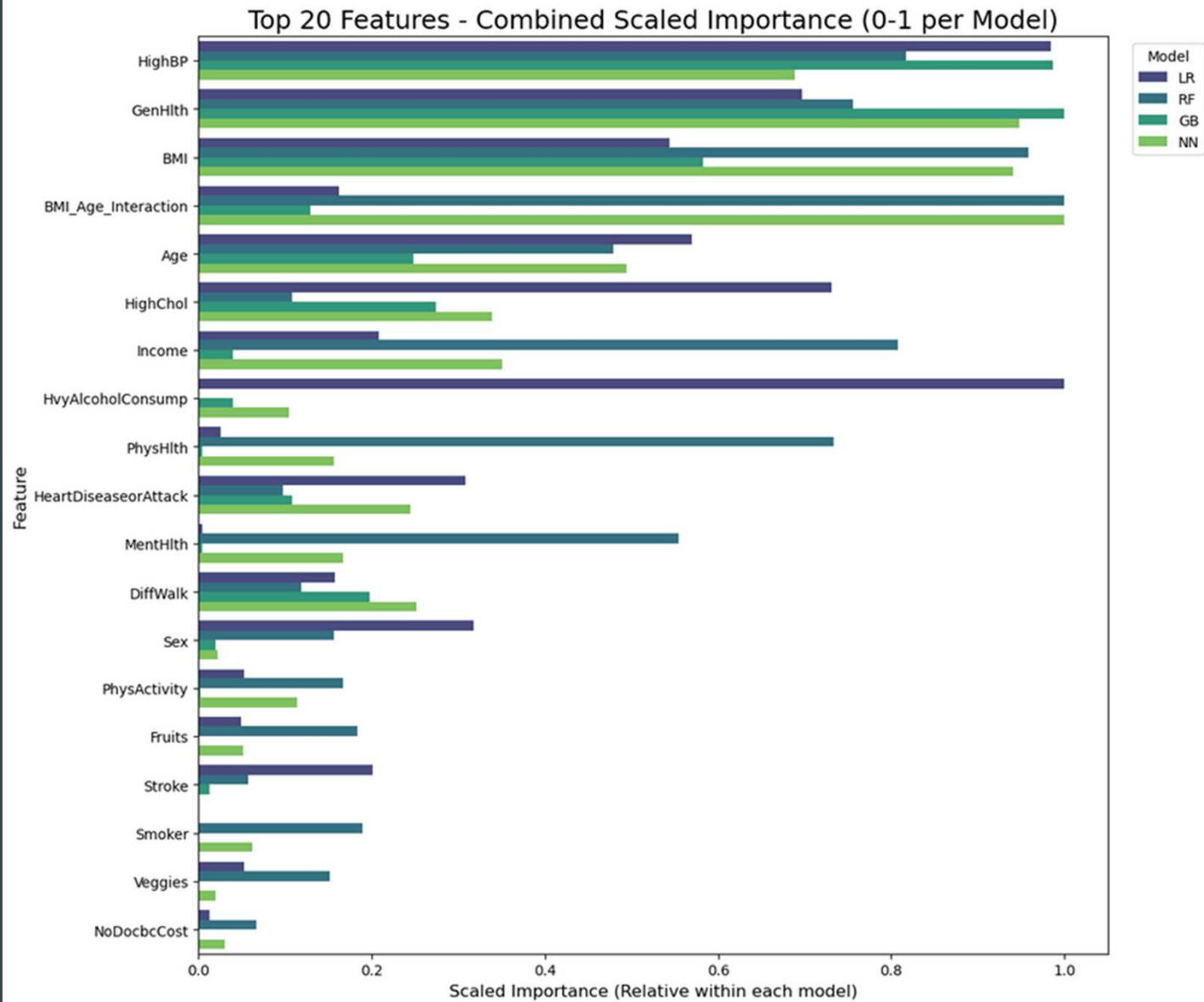
ROC Curve - Model Comparison (VIF-reduced data)



Model Performance Profile (Radar Chart)



	Accuracy	Precision	Recall	F1-score	AUC-ROC
Logistic Regression	0.867445	0.835481	0.867445	0.834772	0.823493
Random Forest	0.861505	0.828025	0.861505	0.834254	0.802718
Gradient Boosting	0.869021	0.839101	0.869021	0.837746	0.827785
Decision Tree	0.803532	0.810982	0.803532	0.807134	0.602289
Neural Network	0.863476	0.813205	0.863476	0.801695	0.824584



Conclusions

Top Performer

- Gradient Boosting Was the Top Performer: Across various evaluation metrics (AUC-ROC, F1-Score, Accuracy, Precision, Recall), the optimized Gradient Boosting model consistently demonstrated the strongest predictive performance (AUC \approx 0.828, F1 \approx 0.838).

Other Model Performances

- Logistic Regression provided a solid, interpretable baseline (AUC \approx 0.823) but was slightly outperformed by non-linear models.
- Neural Networks achieved competitive performance (AUC \approx 0.825), comparable to Logistic Regression, but didn't surpass Gradient Boosting, potentially needing more data or complex architecture.
- Random Forest showed decent performance but lagged behind GB, NN, and LR in AUC (\approx 0.803).
- A single Decision Tree performed poorly (AUC \approx 0.602), highlighting the benefit of ensemble methods like RF and GB.

Model Calibration

- Gradient Boosting, Logistic Regression, and Neural Network models demonstrated relatively good calibration, meaning their predicted probabilities align reasonably well with actual observed risk.

Key Risk Factors Identified

- Feature importance analysis consistently highlighted High Blood Pressure (HighBP), General Health perception (GenHlth), BMI, Age, and High Cholesterol (HighChol) as significant predictors across multiple models. Income and potentially engineered features (like BMI_Age_Interaction) also showed relevance.

Questions?