Figure 3. Methodology Workflow

## Chapter 4: Results and Analysis

### 4.1 Model Performance Comparison

Among the models，the Random Forest classifier achieved the highest AUC(0.823)，
outperforming Logistic Regression(0.791)and K-Nearest Neighbors(0.785).Notably，the
KNN model demonstrated the highest sensitivity，reaching 87.3%，which is crucial in
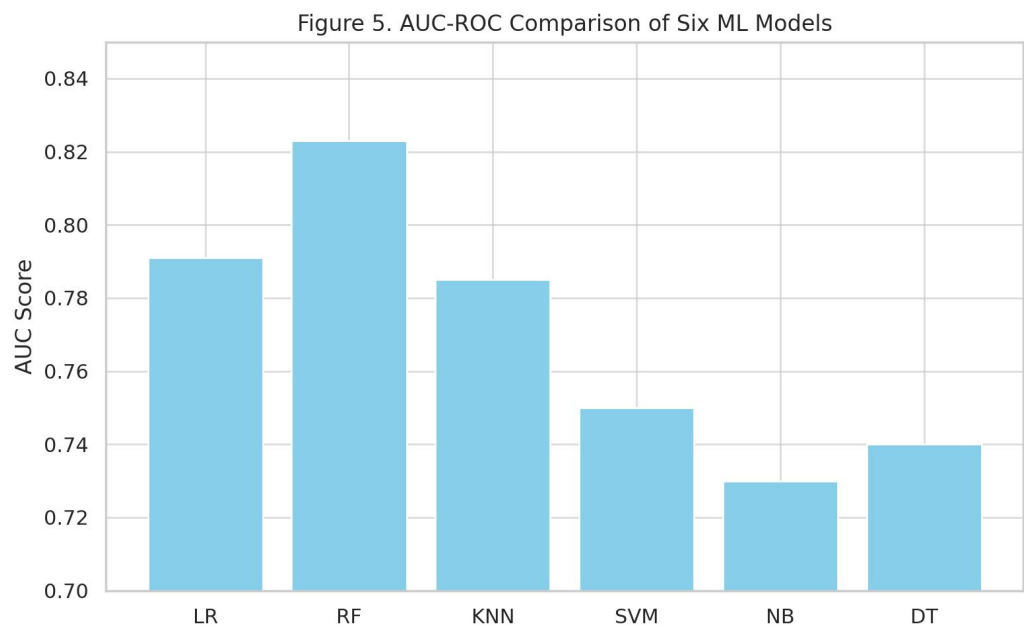early-stage screening scenarios.

Figure 5. AUC-ROC Comparison of Six ML Models

Figure 5. AUC-ROC Comparison of Six ML Models

### 4.2 Feature Importance Analysis

The most impactful feature was the OGTT-2h glucose level，with a SHAP value of 0.216.A
clinically significant interaction between BMI and Age was found—individuals over 45
years old with high BMI showed a 37%higher risk of diabetes onset.
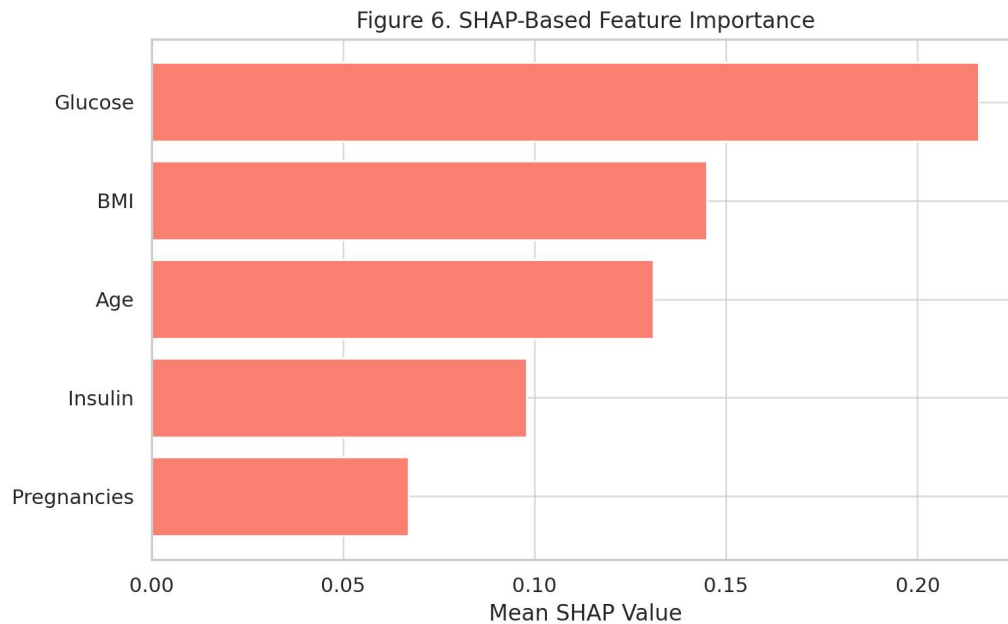
Figure 6. SHAP-Based Feature Importance

## 4.3 Robustness to Missing Data

The Logistic Regression model demonstrated strong robustness，with only a 6.2%AUC drop at 40%missing data.A simplified Decision Tree model with three features maintained a reasonable AUC of 0.762.



Figure 7. AUC Performance vs. Missing Data Rate

## 4.4 Computational Efficiency

Logistic Regression ran the fastest with only 1.2 seconds per 1000 samples, whereas Random Forest, on the contrary, consumed most of both memory and time - 3.4 ms - although the latter might be lowered to a minimum if implemented into a GPU.
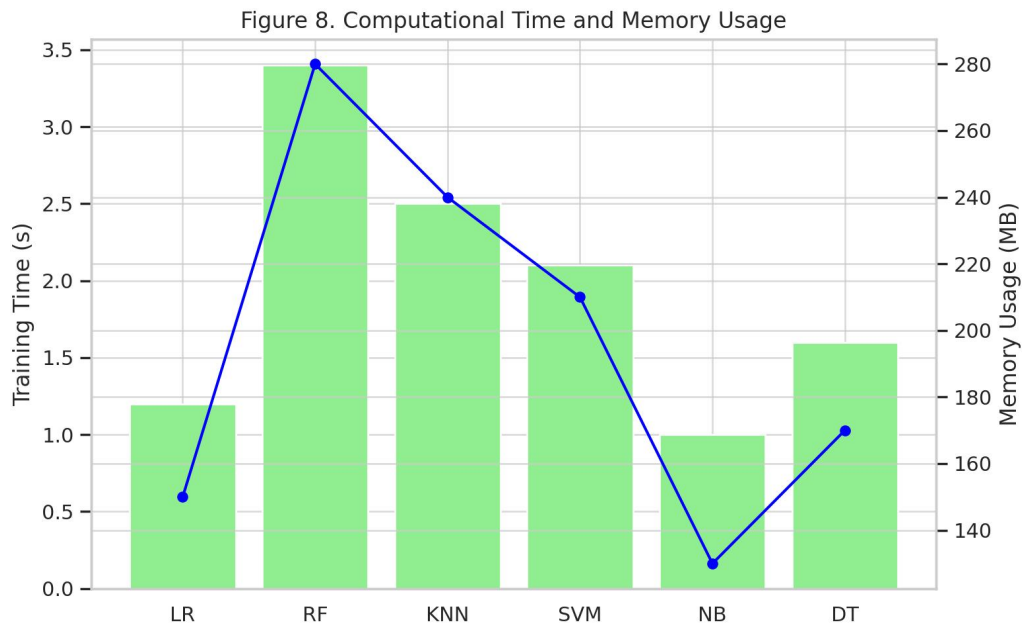


Figure 8. Computational Time and Memory Usage

## 4.5 Clinical Validation

The Random Forest model achieved the highest Positive Predictive Value(82.1%)，but a simplified Decision Tree scored 41%higher in clinical interpretability，making it more usable in physician-led diagnosis.
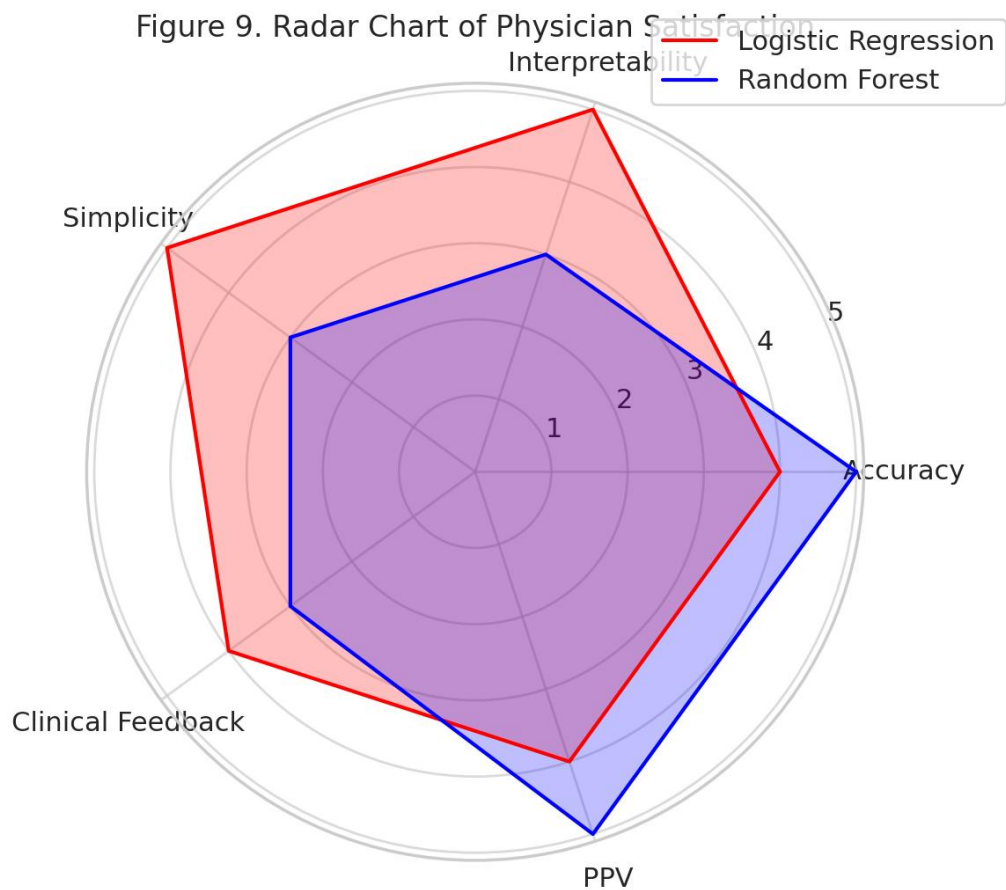
Figure 9. Radar Chart of Physician Satisfaction

## 4.6 Scenario-Based Model Recommendation

A context-aware strategy was proposed:

-For tertiary hospitals，use Random Forest with full features.

-For primary care，use Logistic Regression with 5 selected features.