

RESEARCH TOPIC :

# Application Of Machine Learning In Pregnancy Risk Prediction

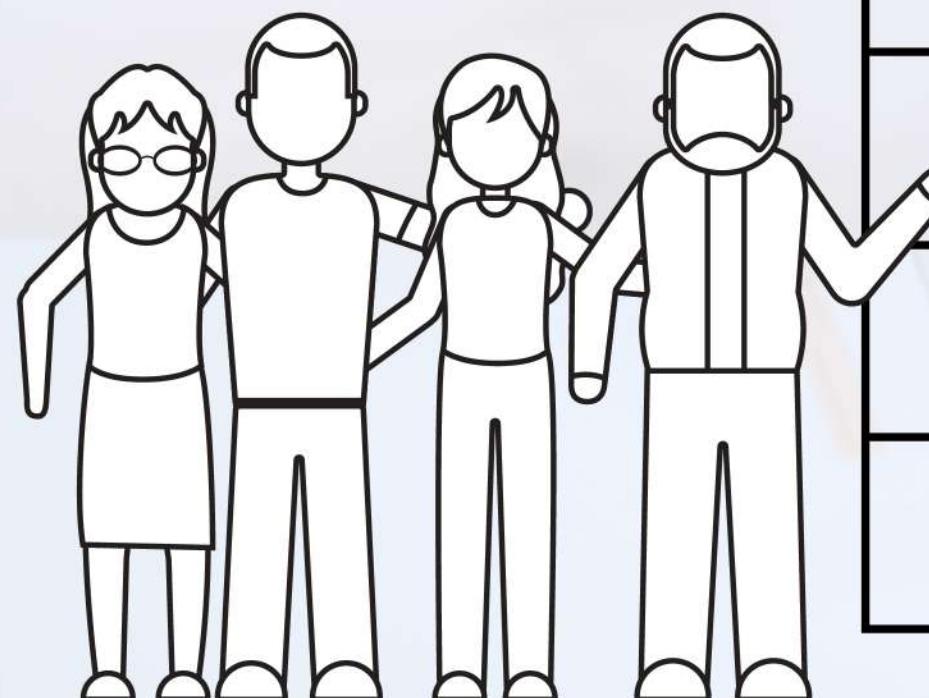
GROUP 6 - MET4

PhD. Cand. Huy Anh Nguyen

Dr. Emmanuel Lance Christopher VI M. Plan

# MEMBERS OF GROUP 6

STUDENT ID	FULL NAME
22080010	NGUYỄN QUỲNH ANH
22080041	NGUYỄN THỊ NGỌC HUYỀN
22080043	NGUYỄN DIỄN KHÁNH
22080044	NGÔ QUÝ KHOA
22080073	VŨ HÀ PHƯƠNG
22080076	NGUYỄN NHƯ QUỲNH
22080082	NGUYỄN ĐỨC THÀNH
22080088	PHẠM XUÂN TRUNG



# CONTENTS OF OUR PROJECT



INTRODUCTION

RESEARCH QUESTIONS

DESCRIPTIVE STATISTICS AND DATA EXPLORATION

ANALYTICS AND INSIGHT

ANALYTICAL MODEL

RECOMMENDATIONS TO STAKEHOLDERS

# Introduction

This study seeks to develop a predictive model that assesses risk factors influencing pregnancy.

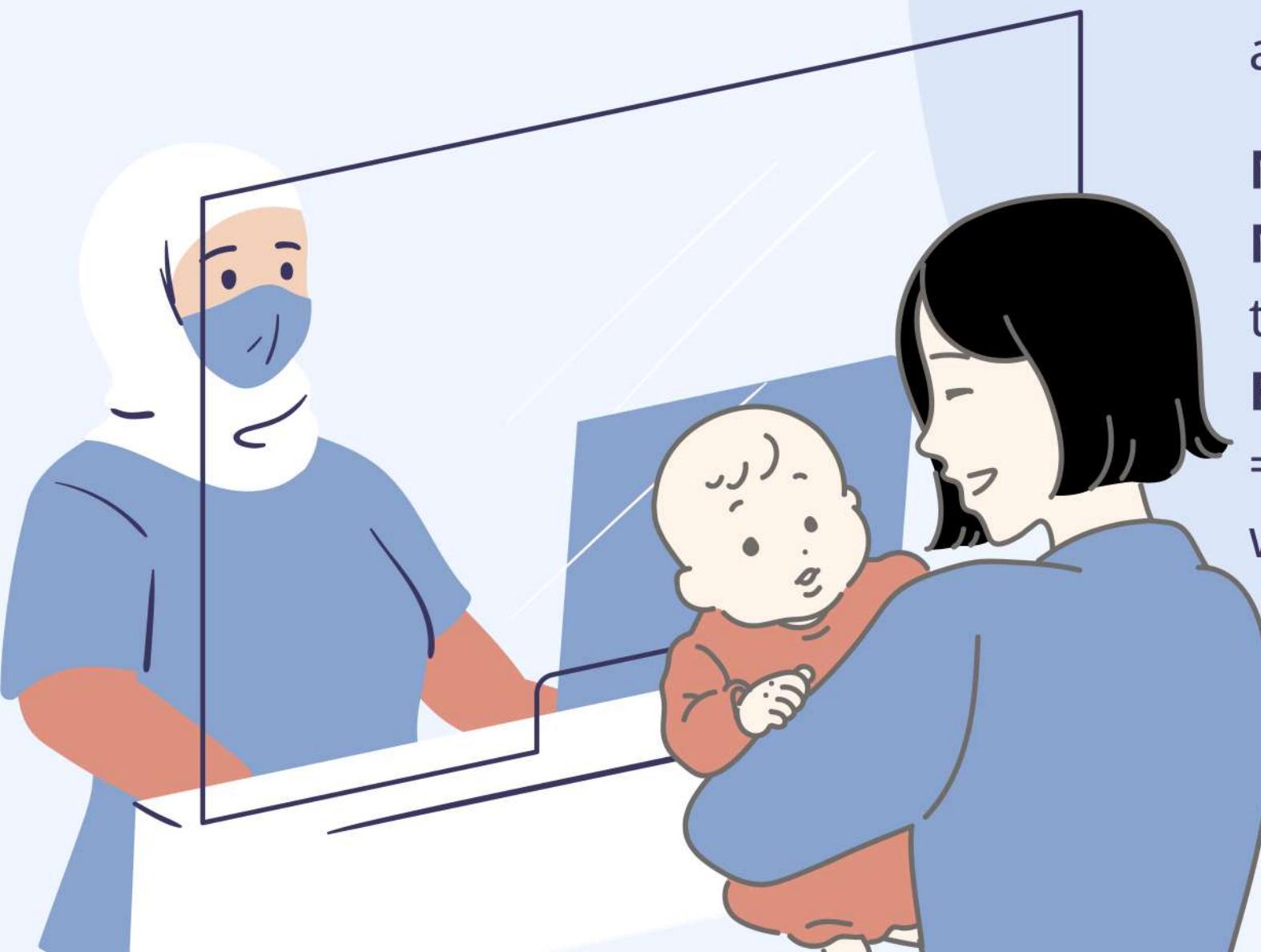
The dataset from the UCI Machine Learning Repository focuses on maternal health risks, collected from healthcare centers in rural Bangladesh using an IoT-based monitoring system. It is accessible on platforms like Kaggle but requires validation with trusted sources like WHO, UNICEF, or peer-reviewed research for accuracy.

**Number of records: 1,014 rows (cases).**

**Number of features:** 7 columns, including 6 input factors and 1 target variable.

**File format:** CSV file.

=> Each row represents a health assessment for a pregnant woman, covering physiological and lifestyle factors.



# Research Questions

**What are the factors that most strongly affect maternal health, and what is the degree of their impact?**

**Where can this predictive model be effectively applied (e.g., hospitals, clinics, rural areas, etc.)?**

**When does the model work most effectively during pregnancy, and in what stages of fetal development?**

**Who will be the primary beneficiaries of this model's implementation (e.g., obstetricians, pregnant women, unborn children, healthcare policymakers)?**

**Why is it important to develop and apply a predictive model to assess maternal health risks?**

**How will this predictive model be implemented to improve maternal healthcare outcomes?**

# Descriptive Statistics and Data Exploration

## Data Cleaning and Preprocessing



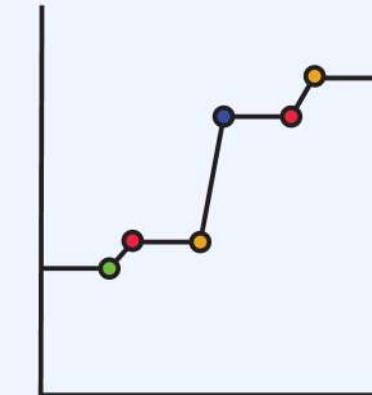
### Handling Duplicate Data

Some rows are completely duplicated.



### Detecting and Handling Outliers

Contains abnormal values  
eg: Heart Rate = 7 in 2 rows (outside healthy range: 60–100 bpm).



### Encoding the Target Variable

Low Risk → Assigned the value 0  
Mid Risk → Assigned the value 1  
High Risk → Assigned the value 2

# Descriptive Statistics and Data Exploration

## Data Cleaning and Preprocessing

	Age	SystolicBP	DiastolicBP	BS	BodyTemp	HeartRate
count	1014.000000	1014.000000	1014.000000	1014.000000	1014.000000	1014.000000
mean	29.871795	113.198225	76.460552	8.725986	98.665089	74.301775
std	13.474386	18.403913	13.885796	3.293532	1.371384	8.088702
min	10.000000	70.000000	49.000000	6.000000	98.000000	7.000000
25%	19.000000	100.000000	65.000000	6.900000	98.000000	70.000000
50%	26.000000	120.000000	80.000000	7.500000	98.000000	76.000000
75%	39.000000	120.000000	90.000000	8.000000	98.000000	80.000000
max	70.000000	160.000000	100.000000	19.000000	103.000000	90.000000

**Age Range:** Majority are within reproductive age (18–45 years); a few outliers.

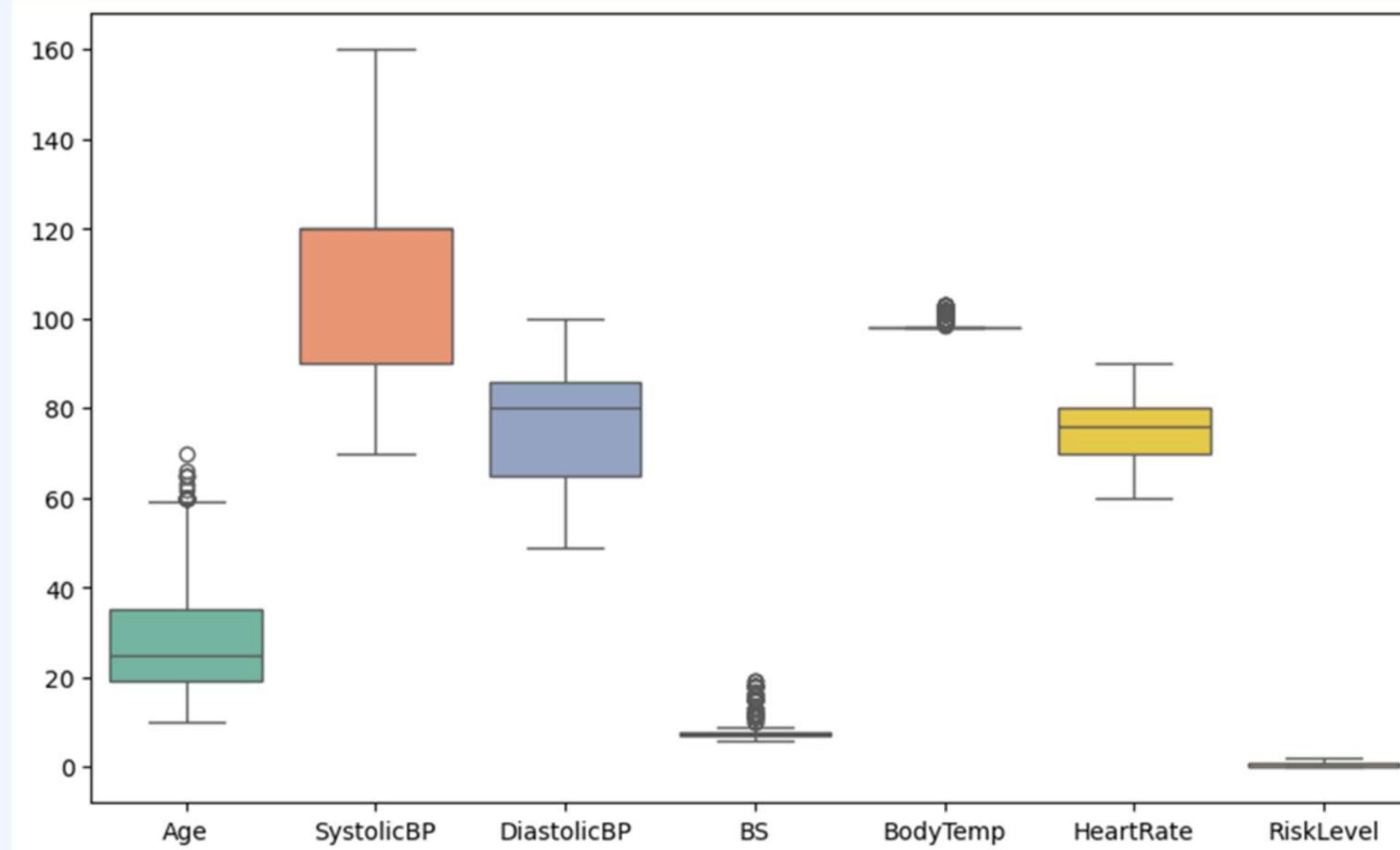
**Blood Pressure:** Generally healthy, but some high readings (e.g., SystolicBP 160 mmHg) suggest hypertension.

**Blood Sugar (BS):** Elevated averages, with a maximum of 19 mmol/L, indicating possible gestational diabetes.

**Heart Rate:** Corrected values (70–90 bpm) align with normal resting rates for adults.

# Descriptive Statistics and Data Exploration

## Observation



### Age:

- Primarily younger participants (19-35), left-skewed
- Outliers at 60+ age with a max 70

### Systolic Blood Pressure:

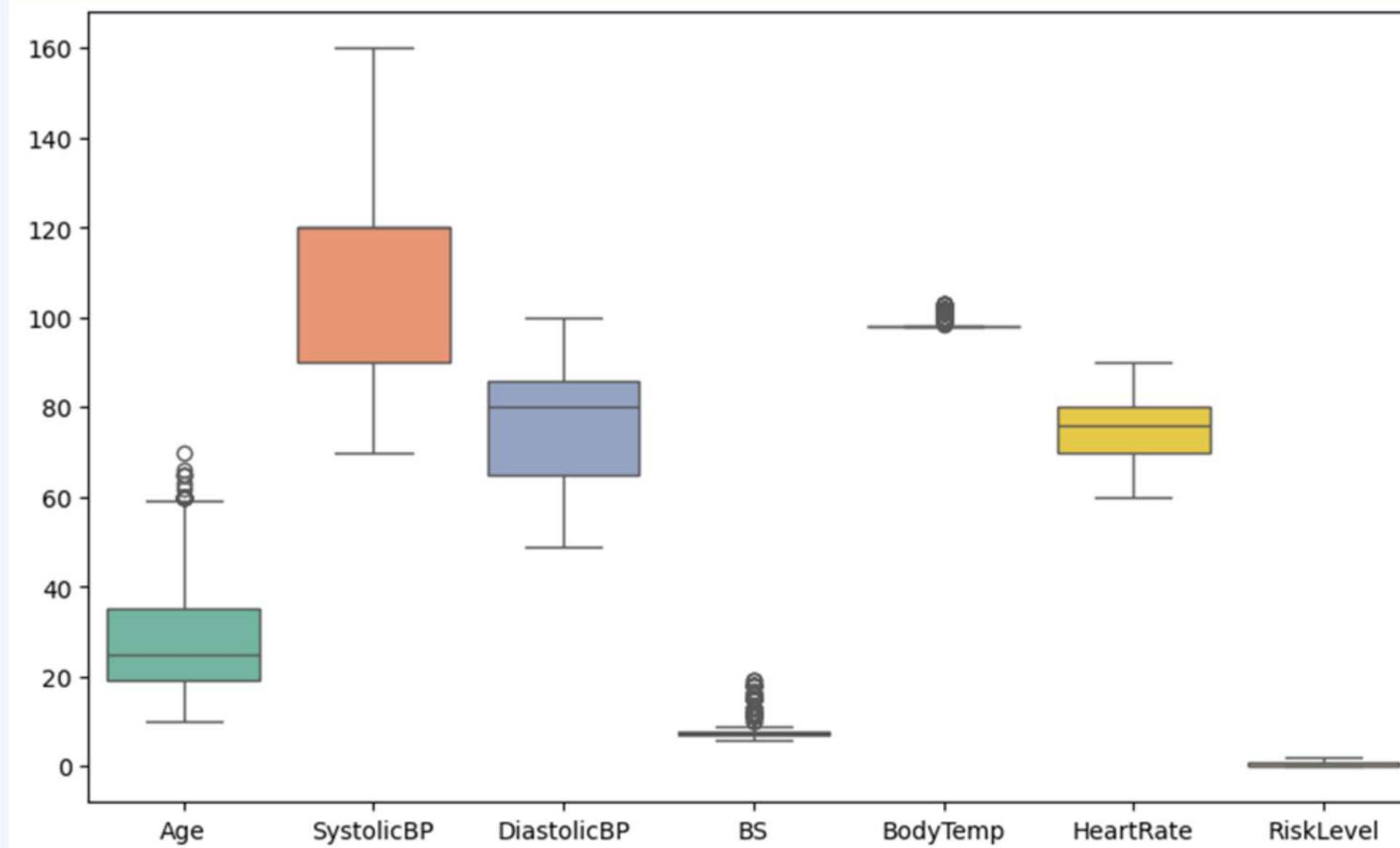
- Concentrated 90-120 mmHg, left-skewed
- No outliers observed up to max 160mmHg

### Diastolic Blood Pressure:

- Evenly distributed 65-86 mmHg, median 80 mmHg
- No outliers up to max 100mmHg

# Descriptive Statistics and Data Exploration

## Observation



### Blood Glucose:

- Outliers for values over 10 mmol/L, max 19 mmol/L
- Significant number of elevated levels

### Body Temperature:

- Mostly around 98°F
- Outliers observed from 98.4°F+, max temp of 103°F

### Heart Rate:

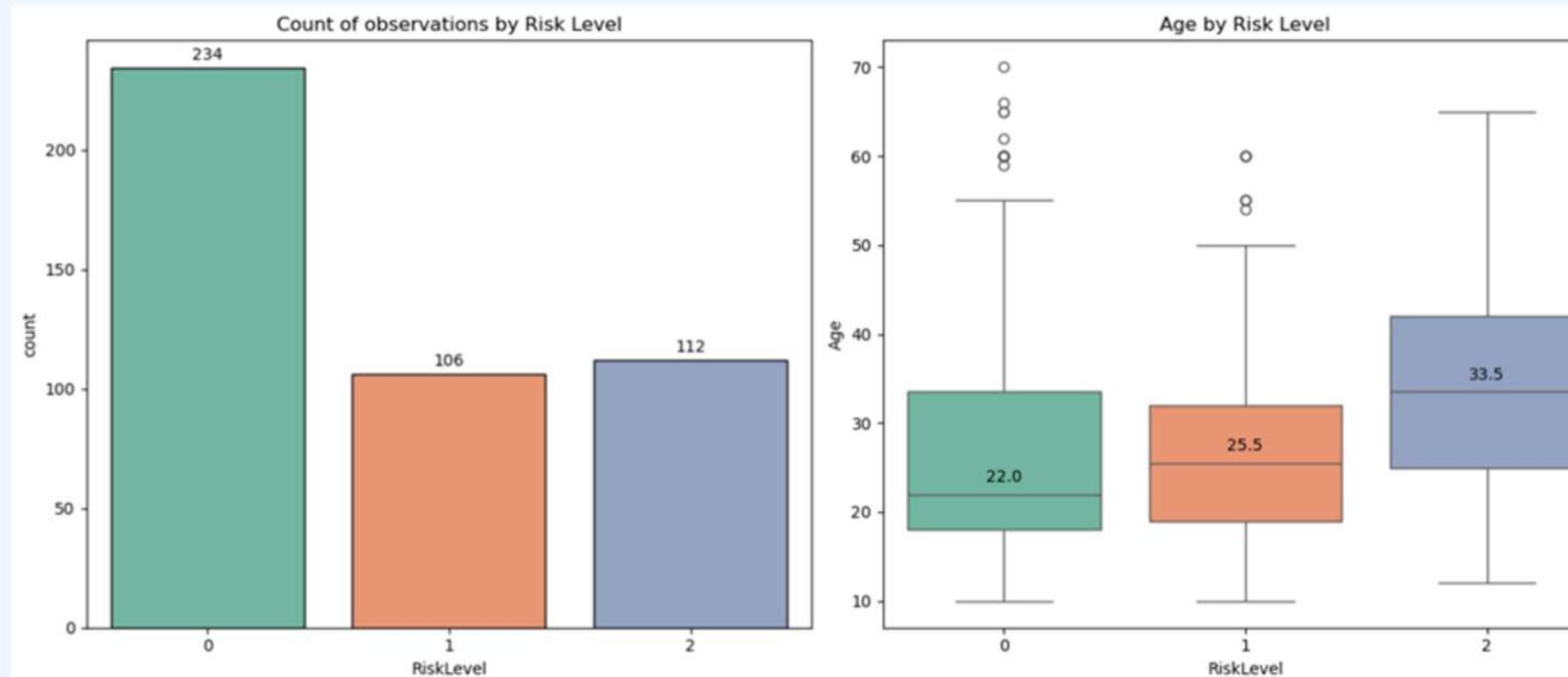
- Evenly distributed 70-80 bpm, median 76 bpm
- No outliers up to max 90 bpm

### Risk Level:

- Distributed 0-1, no outliers observed

# Analytics and Insights

## AGE



### Sample Distribution by Risk Level:

- Low Risk: 234 samples (largest group).
- Medium Risk: 106 samples.
- High Risk: 112 samples.

### Age and Risk Correlation:

- Average age increases with risk level:
  - Low Risk: ~22 years.
  - High Risk: ~33.5 years.

### Outliers:

Present across all risk levels, especially for Low Risk (0) and High Risk (2), representing individuals significantly older or younger than the group average.

# Analytics and Insights

## BLOOD GLUCOSE LEVEL



- No clear linear relationship between blood sugar levels and age.
- Data points are evenly distributed, showing high variability across all ages.

=> This suggests that blood sugar levels are influenced by factors other than age, with no consistent trend observed.

# Analytics and Insights

## BLOOD PRESSURE

### Blood Pressure Patterns:

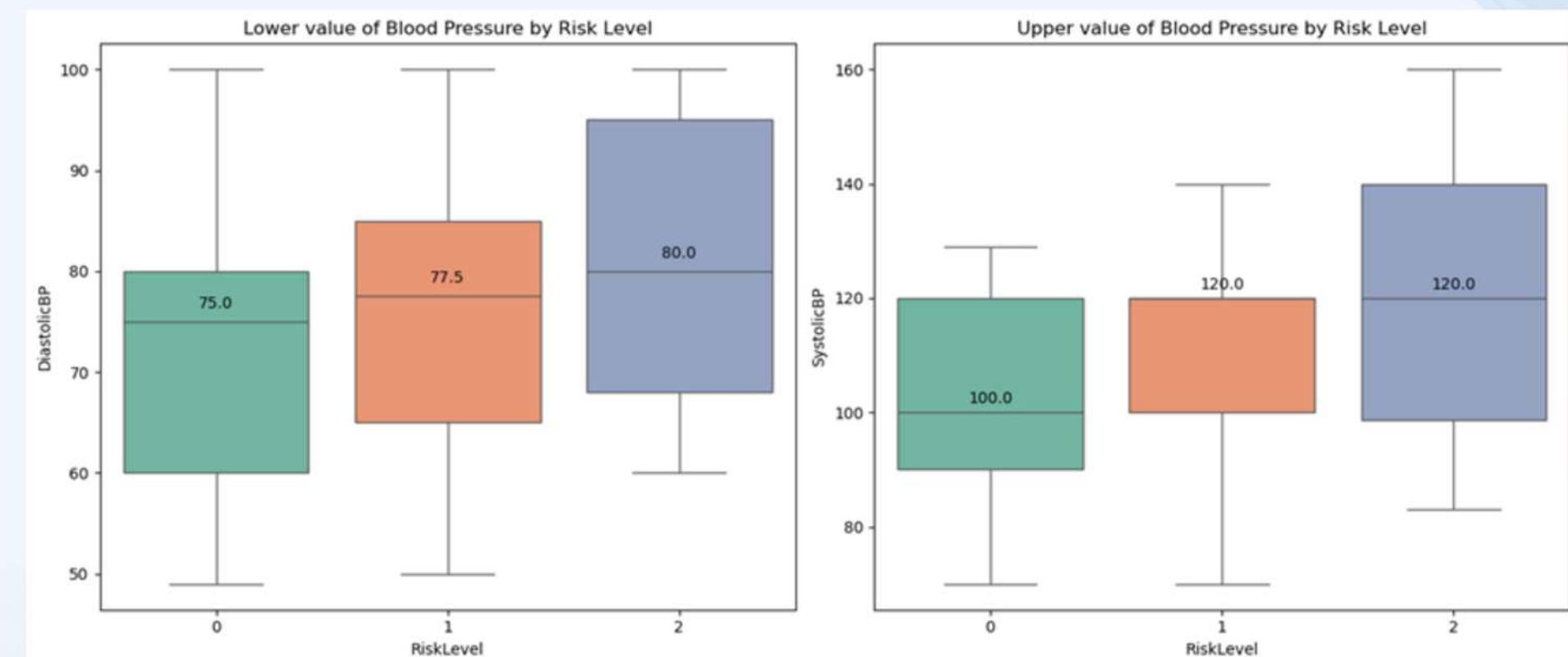
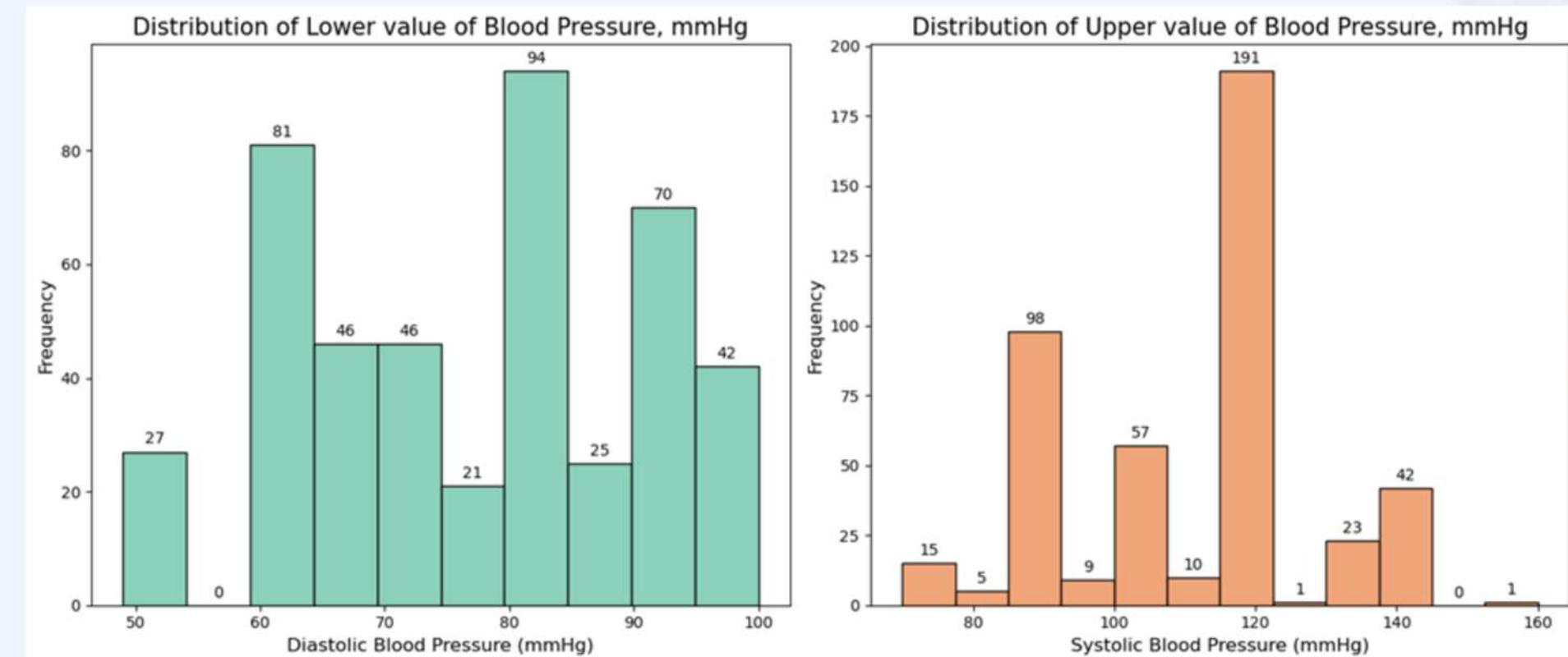
- Diastolic: 70–90 mmHg, peak at 80 mmHg.
- Systolic: 100–120 mmHg, sharp peak at 120 mmHg.

### Risk Correlation:

- Higher risk levels have increased median BP and broader ranges.
- High-risk groups show higher maximum BP values.

### Clinical Insight:

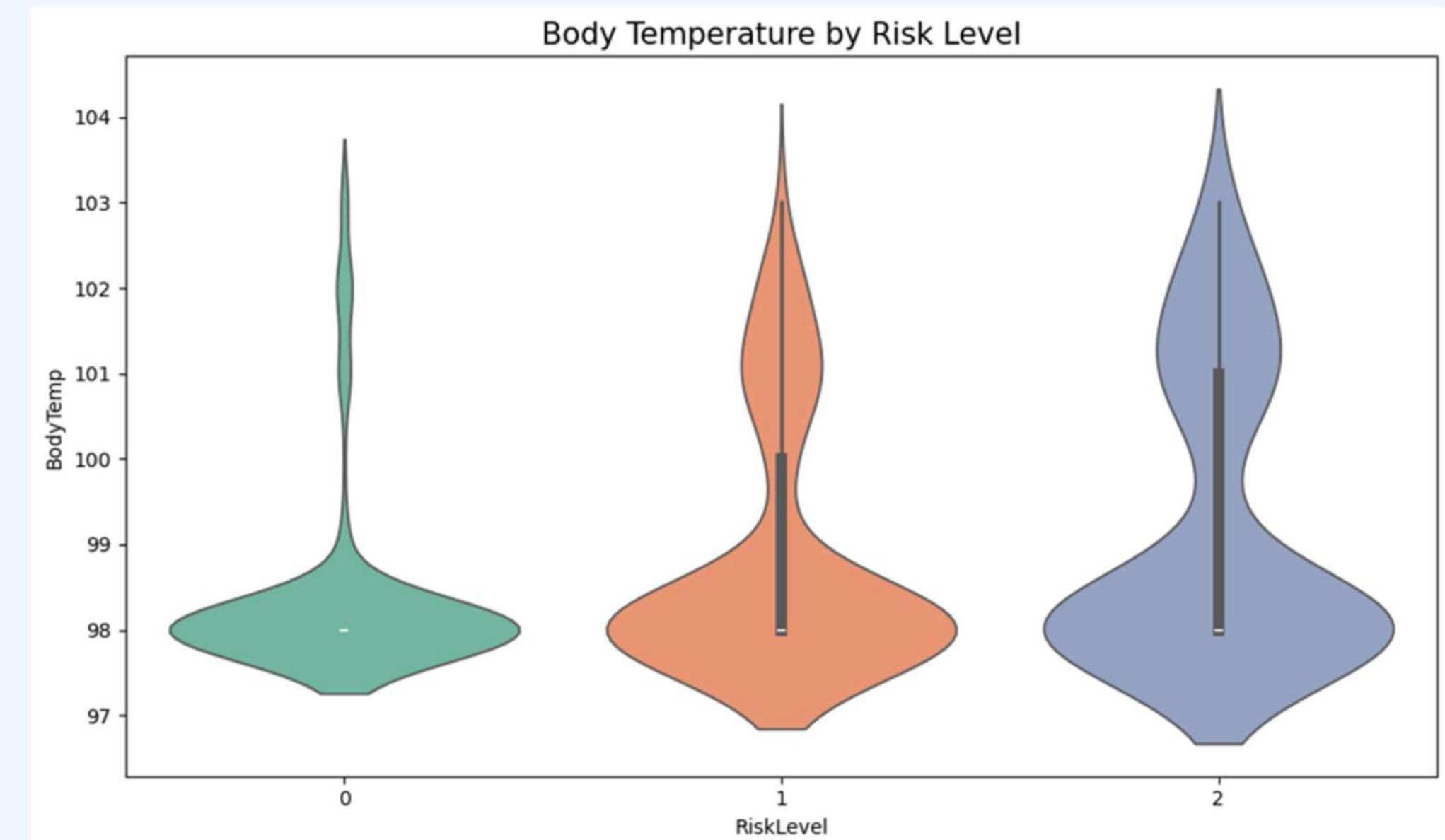
- Elevated BP correlates with higher pregnancy risk.
- Safe ranges: Systolic: 90–120 mmHg, Diastolic: 60–80 mmHg.



# Analytics and Insights

## BODY TEMPERATURE

- Low Risk (0): Narrowly distributed around 98°F, indicating stability.
  - Mid Risk (1): Slightly broader range, with cases exceeding 101°F, signaling potential risks.
  - High Risk (2): Widest variability, including extremes above 103°F, strongly linked to severe maternal risks.
- => Elevated body temperature is a significant indicator of pregnancy-related risks, emphasizing its importance in risk prediction.



# Analytics and Insights

## LOOKING CORRELATION

- Evaluates relationships between health factors (e.g., Age, SystolicBP, DiastolicBP, BS, BodyTemp, HeartRate) and RiskLevel.
- Correlation values range from 0 (weak) to 1 (strong).

### Notable Correlations:

- Strong Correlation: Blood Sugar (BS) with RiskLevel (0.55).

### Moderate Correlations:

- SystolicBP (0.33) and DiastolicBP (0.25) with RiskLevel.
- Weak Correlation: Age (0.18) and BodyTemp (0.26) with RiskLevel.

=> Blood Sugar significantly influences RiskLevel, while Age and BodyTemp have minimal impact. This highlights the importance of prioritizing key factors in risk modeling.



# Analytics and Insights

## HEART RATE

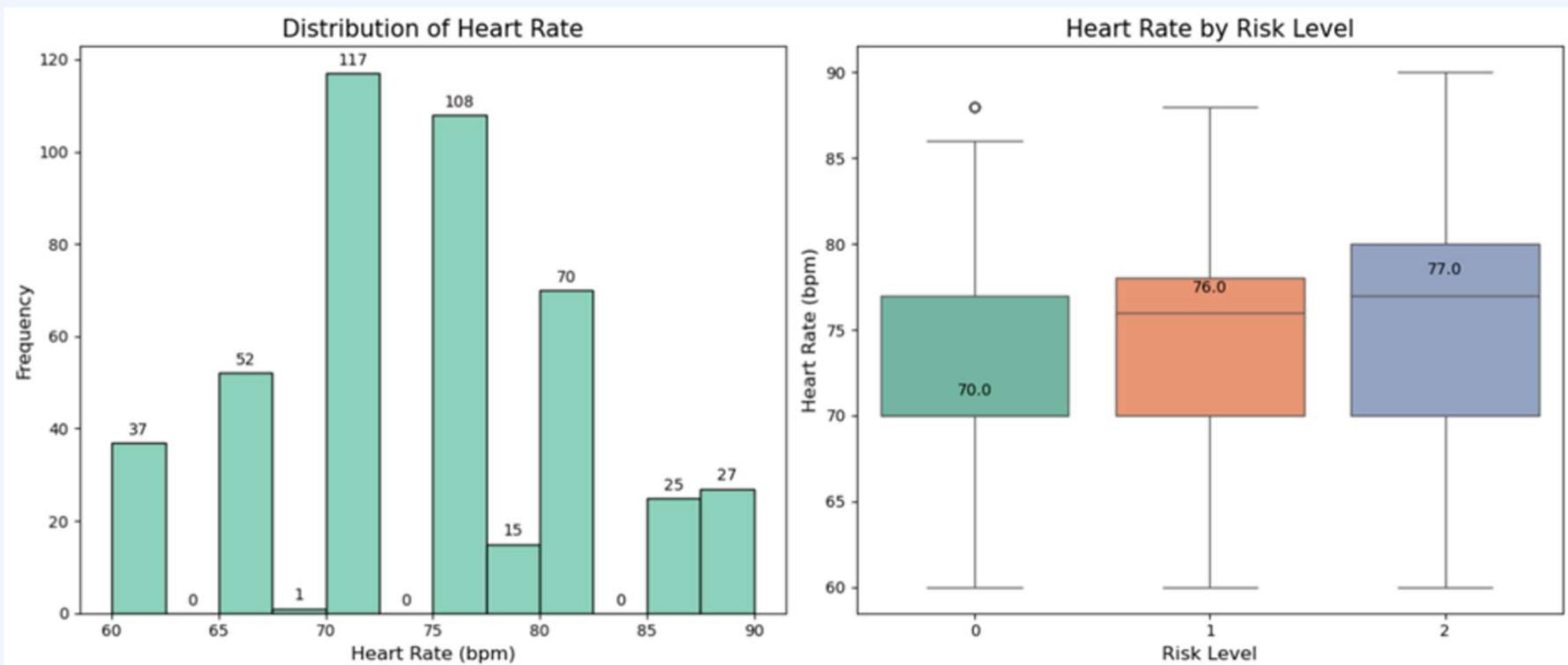
### Heart Rate Distribution:

- Bell-shaped curve centered around 70 bpm, normal for adults.

### Median Heart Rates by Risk Level:

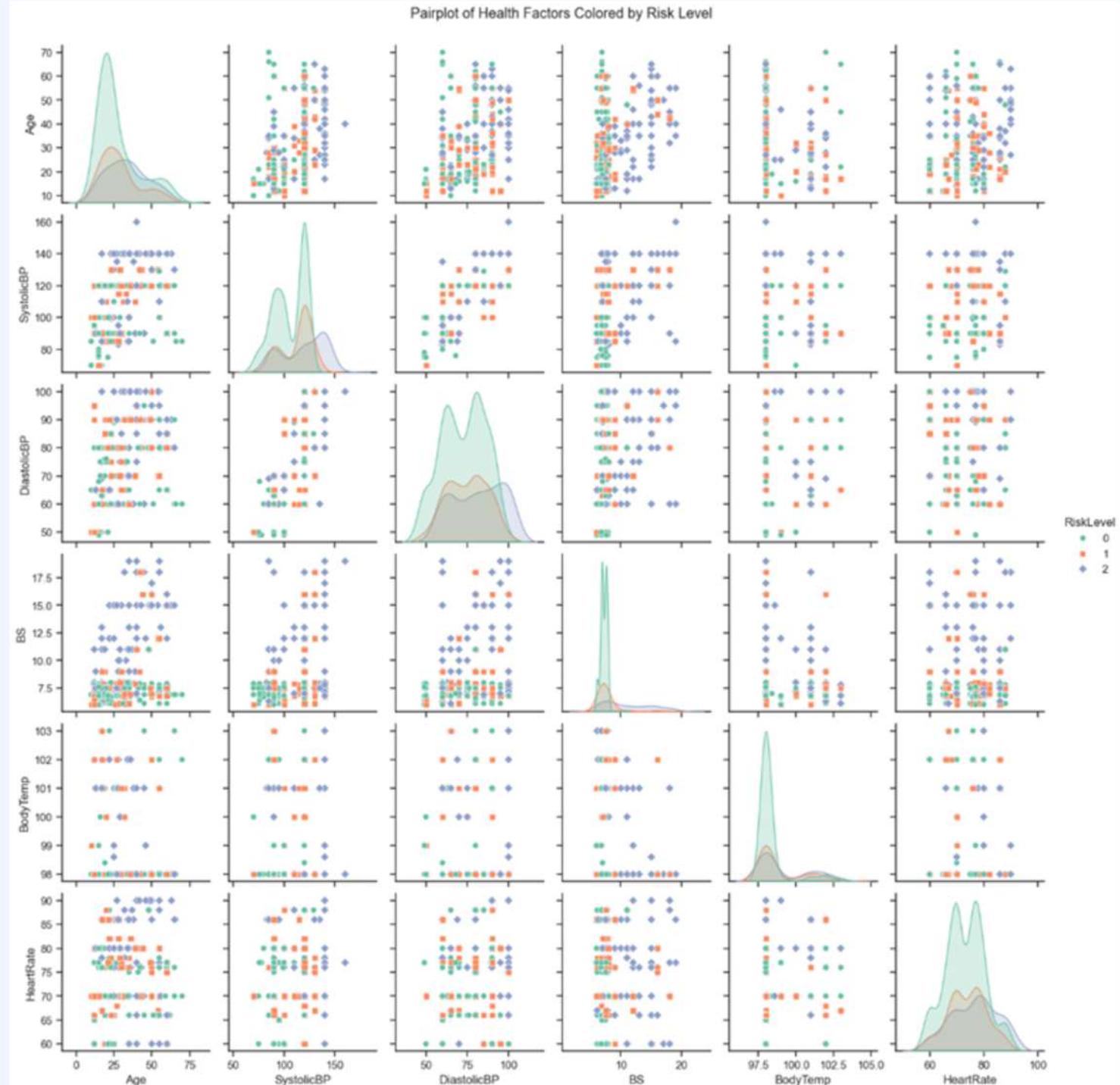
- Low Risk: 70 bpm.
- Mid Risk: 76 bpm.
- High Risk: 77 bpm.

=> Heart rate increases with risk level, highlighting its significance in assessing pregnancy-related risks.



# Analytics and Insights

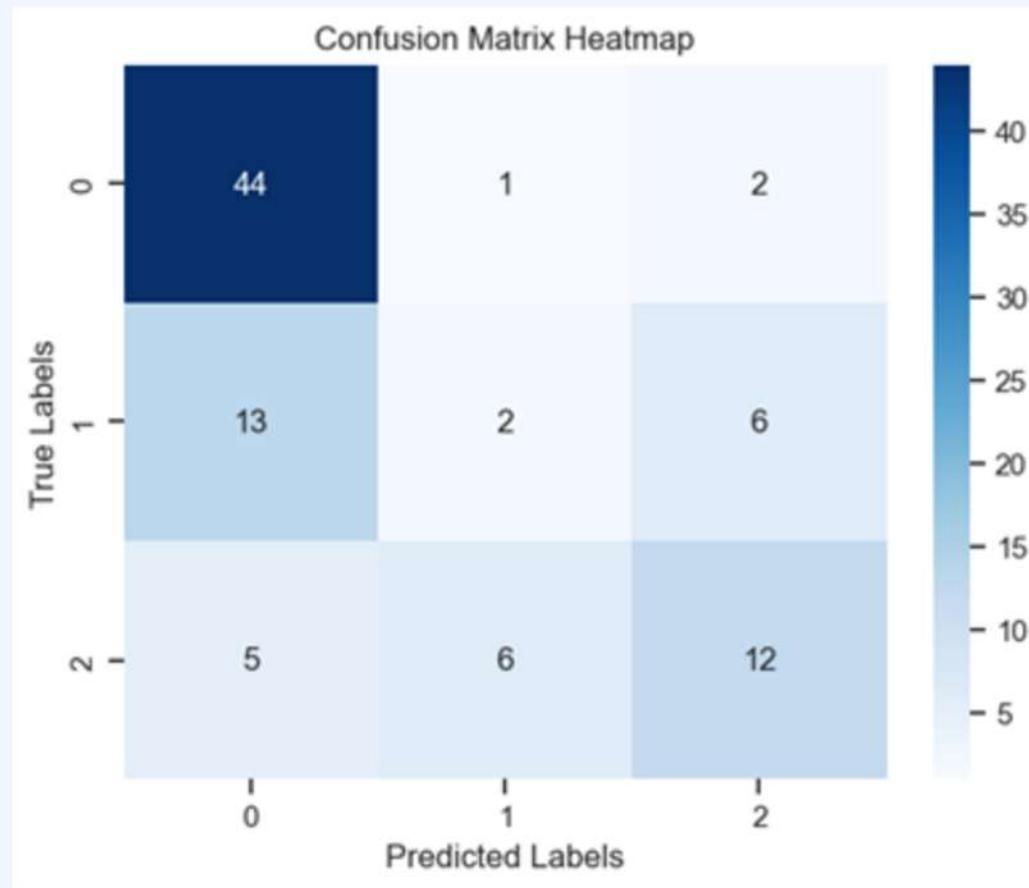
## LOOKING CORRELATION



- Visualizes relationships between health factors (e.g., systolic/diastolic blood pressure, age, body temperature, heart rate) and risk levels.
- Systolic & Diastolic BP: Strong positive correlation—when systolic increases, diastolic also increases.
- No clear correlation observed between blood pressure, heart rate, and body temperature.  
=> Different health factors contribute uniquely to risk assessment, highlighting the complexity of pregnancy-related risk prediction.

# Analytical Model

## LOGISTIC REGRESSION



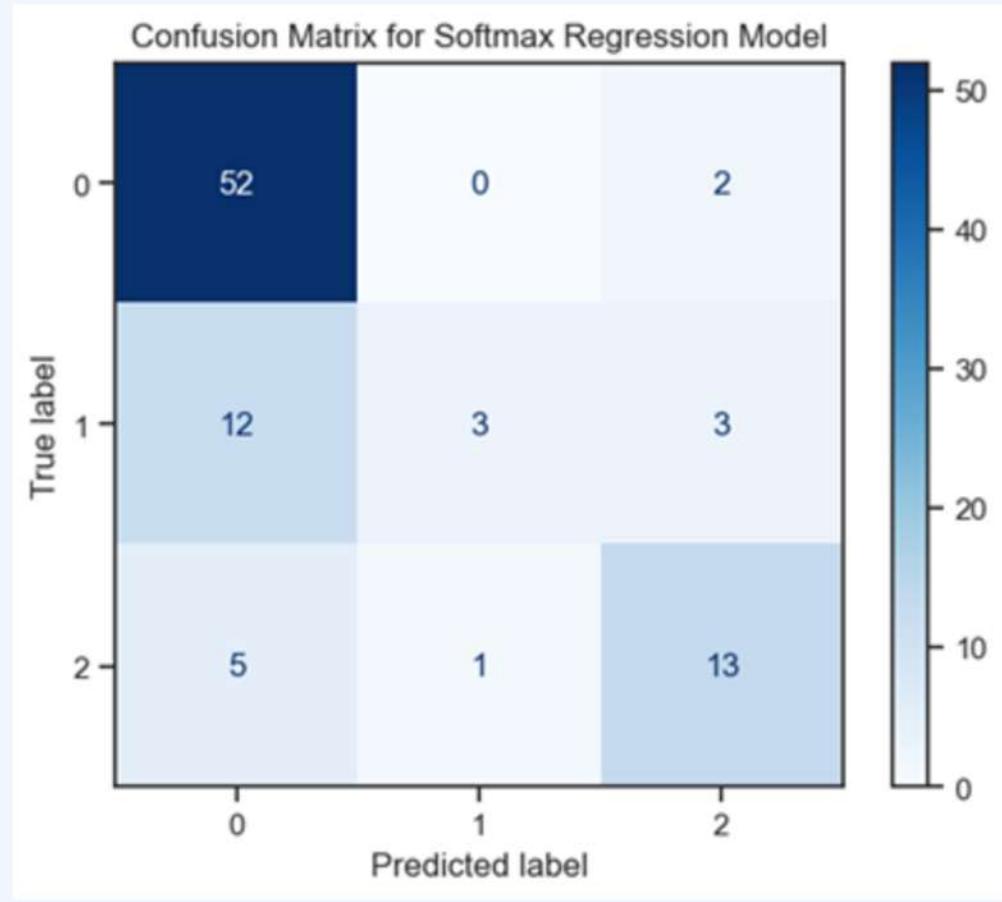
Classification Report:					
	precision	recall	f1-score	support	
0	0.71	0.94	0.81	47	
1	0.22	0.10	0.13	21	
2	0.60	0.52	0.56	23	
accuracy			0.64	91	
macro avg	0.51	0.52	0.50	91	
weighted avg	0.57	0.64	0.59	91	

Accuracy Score: 0.64

- The logistic regression model's accuracy is 0.64, which is not too low but not good enough for evaluation
  - Other metrics like precision, recall, F1-score, and confusion matrix are needed
  - The model performs best on class 0, poorly on class 1, and sub-optimally on class 2
- => The model has limitations and may not be the most suitable choice for this dataset

# Analytical Model

## SOFTMAX REGRESSION



Classification Report:

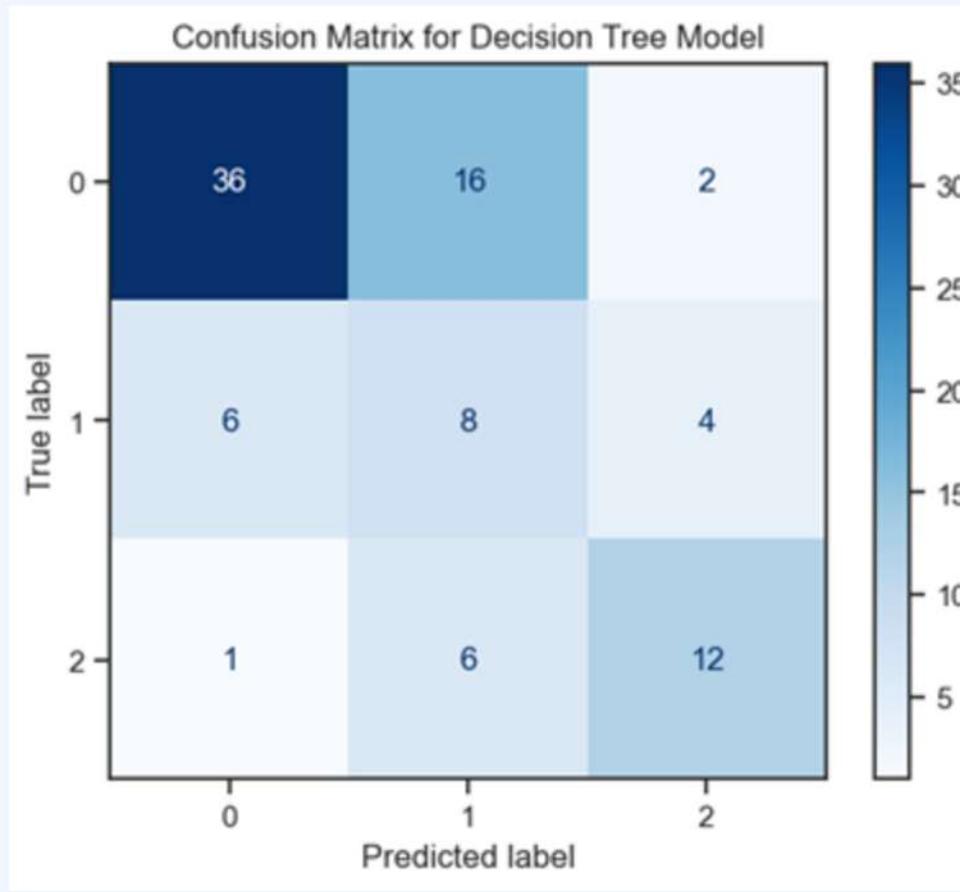
	precision	recall	f1-score	support
0	0.75	0.96	0.85	54
1	0.75	0.17	0.27	18
2	0.72	0.68	0.70	19
accuracy			0.75	91
macro avg	0.74	0.60	0.61	91
weighted avg	0.75	0.75	0.70	91

Accuracy Score: 0.7472527472527473

- Accuracy: 75%
- Class 0 (Low Risk): Recall 0.96, F1-score 0.85 (Best performance)
- Class 1 (Mid Risk): Recall 0.17, F1-score 0.27 (Poor performance)
- Class 2 (High Risk): F1-score 0.70 (Moderate performance, occasional misclassifications)
- Recommendations: Address class imbalance, optimize features.

# Analytical Model

## DECISION TREE



Classification Report:

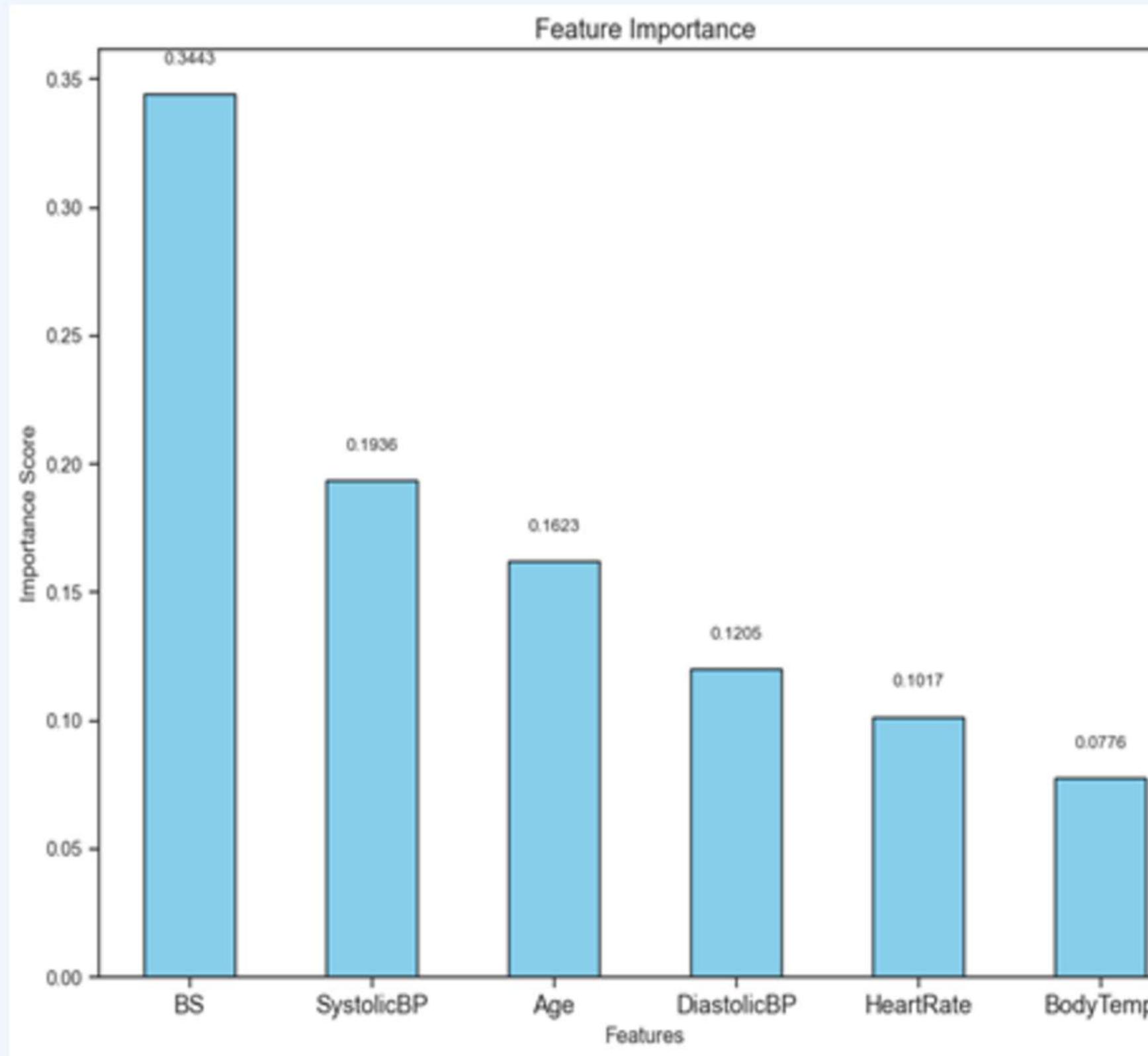
	precision	recall	f1-score	support
0	0.84	0.67	0.74	54
1	0.27	0.44	0.33	18
2	0.67	0.63	0.65	19
accuracy			0.62	91
macro avg	0.59	0.58	0.57	91
weighted avg	0.69	0.62	0.64	91

Accuracy Score: 0.6153846153846154

- The decision tree model has a moderate accuracy of 0.62, which is not fully indicative of its performance.
- Accuracy alone is not enough to judge the model's effectiveness. Other metrics are needed.
  - Performs best for class 0 (high precision and F1-score)
  - Struggles with class 1 (low precision and recall)
  - Class 2 performance is moderate and can still be improved

# Analytical Model

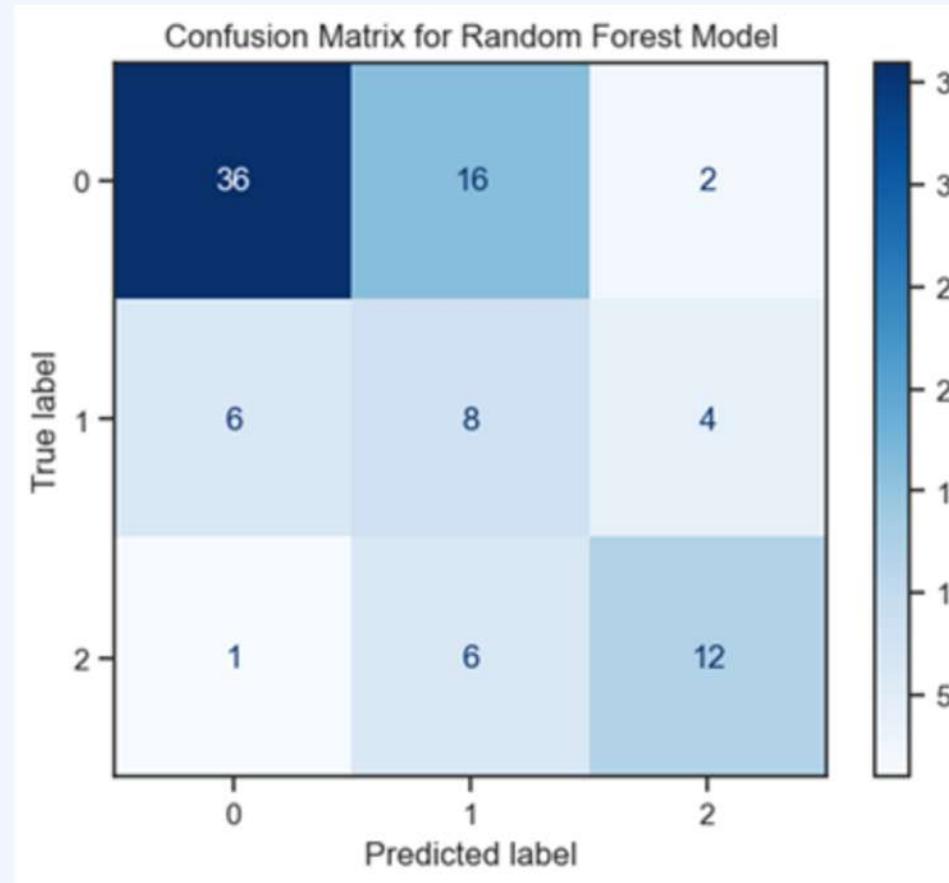
## DECISION TREE



The Important Features chart shows that the model has made good use of special important features such as BS (Blood Sugar), Systolic Blood Pressure (SBP) and Maternal Age.

# Analytical Model

## RANDOM FOREST



Classification Report:					
	precision	recall	f1-score	support	
0	0.83	0.80	0.81	54	
1	0.26	0.28	0.27	18	
2	0.65	0.68	0.67	19	
accuracy			0.67	91	
macro avg	0.58	0.59	0.58	91	
weighted avg	0.68	0.67	0.67	91	

Accuracy Score: 0.6703296703296703

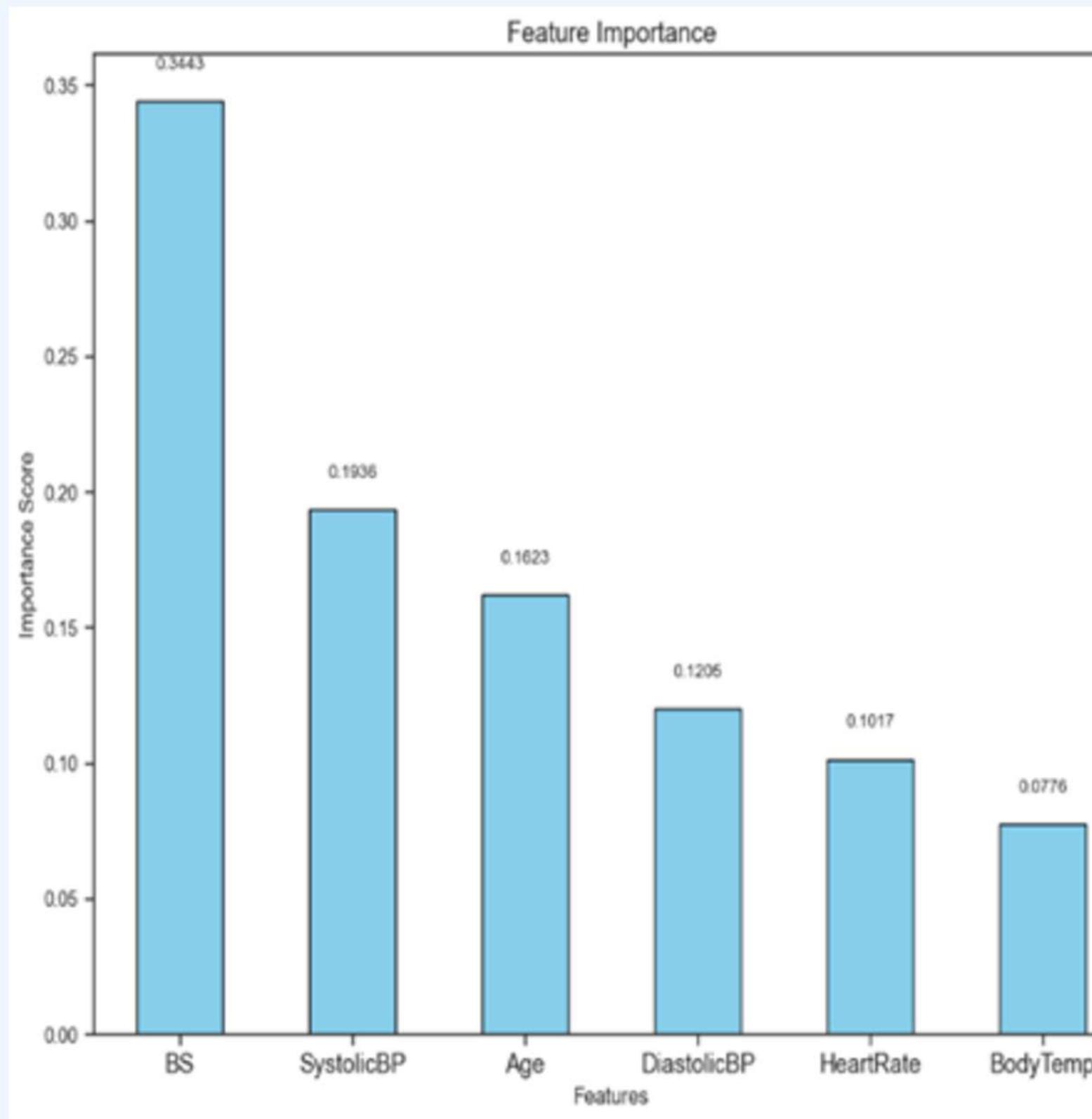
The Random Forest model has an accuracy of 67%

- Best for Class 0 (Low Risk) with high precision (0.83) and recall (0.80)
- Struggles with Class 1 (Mid Risk), showing low precision (0.26) and recall (0.28)
- Class 2 (High Risk) has moderate performance with an F1-score of 0.67

=> While good for low-risk, the model needs optimization to better classify mid and high-risk cases for balanced predictions across all classes

# Analytical Model

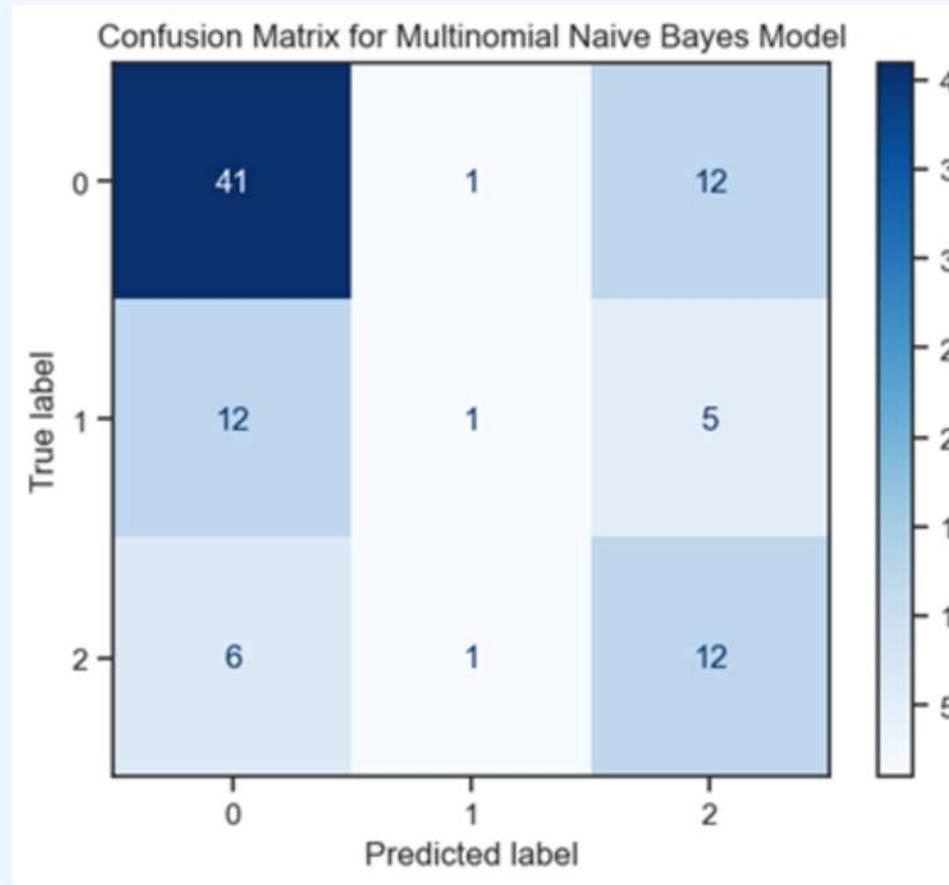
## RANDOM FOREST



The Feature Importance chart also shows that the model has effectively utilized important features, especially health-related factors such as BS and Systolic BP, which clearly reflects the model's logic and practicality in prediction.

# Analytical Model

## MULTINOMIAL NAIVE BAYES



Classification Report:					
	precision	recall	f1-score	support	
0	0.69	0.76	0.73	54	
1	0.33	0.06	0.10	18	
2	0.41	0.63	0.50	19	
accuracy			0.59	91	
macro avg	0.48	0.48	0.44	91	
weighted avg	0.56	0.59	0.55	91	

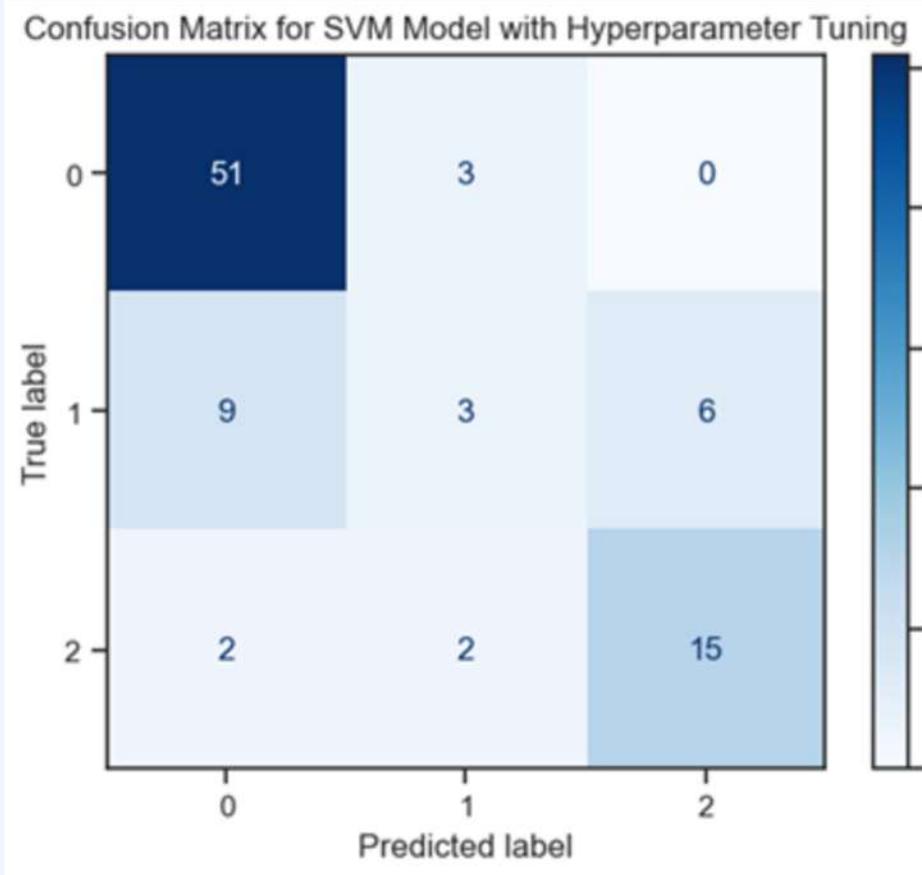
Accuracy Score: 0.5934065934065934

The MNB model has a accuracy of around 0.59, indicating limited classification performance

- Class 1 is barely recognized (recall - 0.06, F1 score - 0.10)
  - The difference between macro F1 (0.42) and weighted F1 (0.53) => The model performs better on majority classes than minority classes
- => The MNB model is imbalanced and struggles to classify minority classes effectively.

# Analytical Model

## SUPPORT VECTOR MACHINE



Classification Report:

	precision	recall	f1-score	support
0	0.82	0.94	0.88	54
1	0.38	0.17	0.23	18
2	0.71	0.79	0.75	19
accuracy			0.76	91
macro avg	0.64	0.63	0.62	91
weighted avg	0.71	0.76	0.72	91

Accuracy Score: 0.7582417582417582

The SVM model shows improved data separation with an accuracy of 0.76.

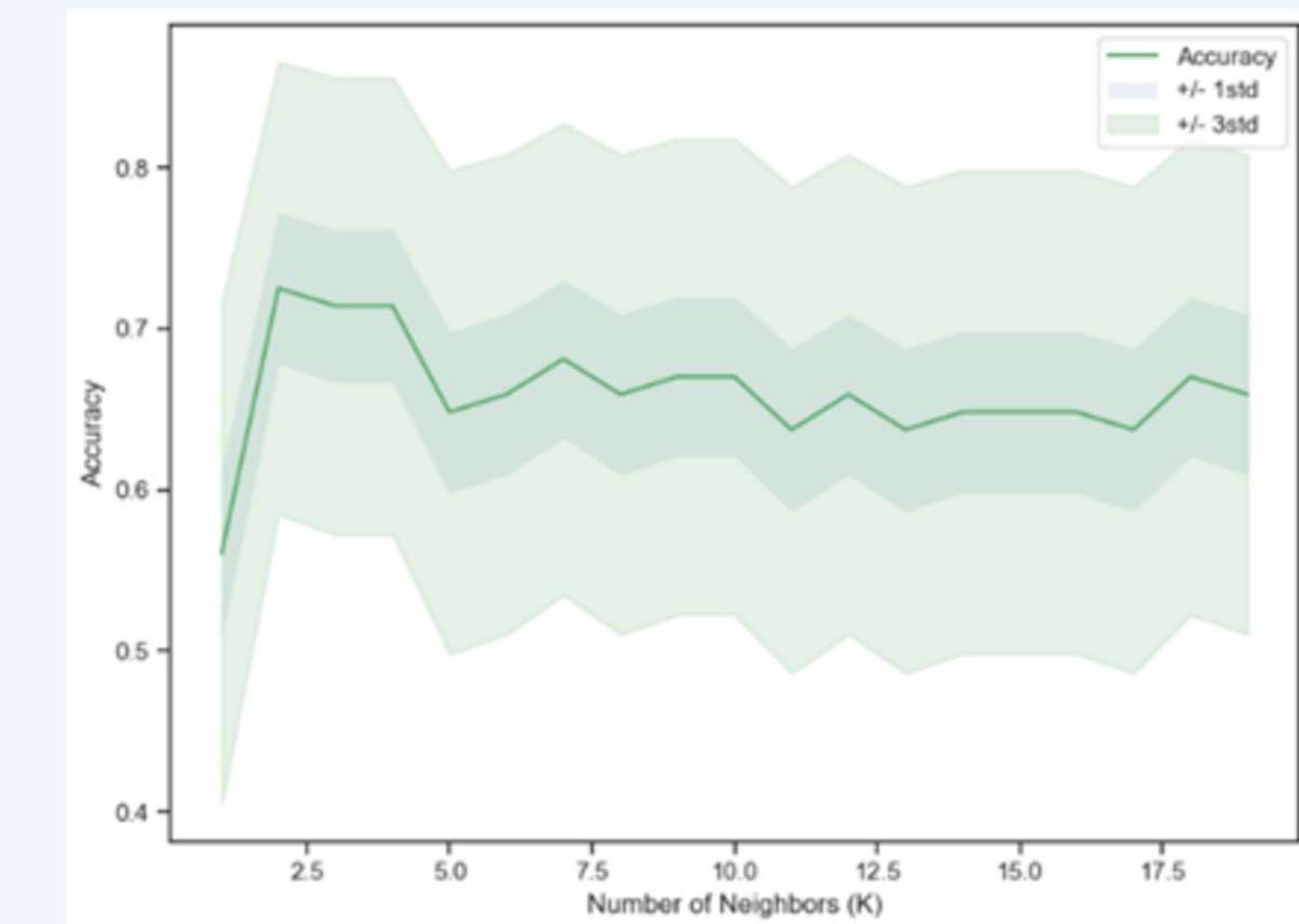
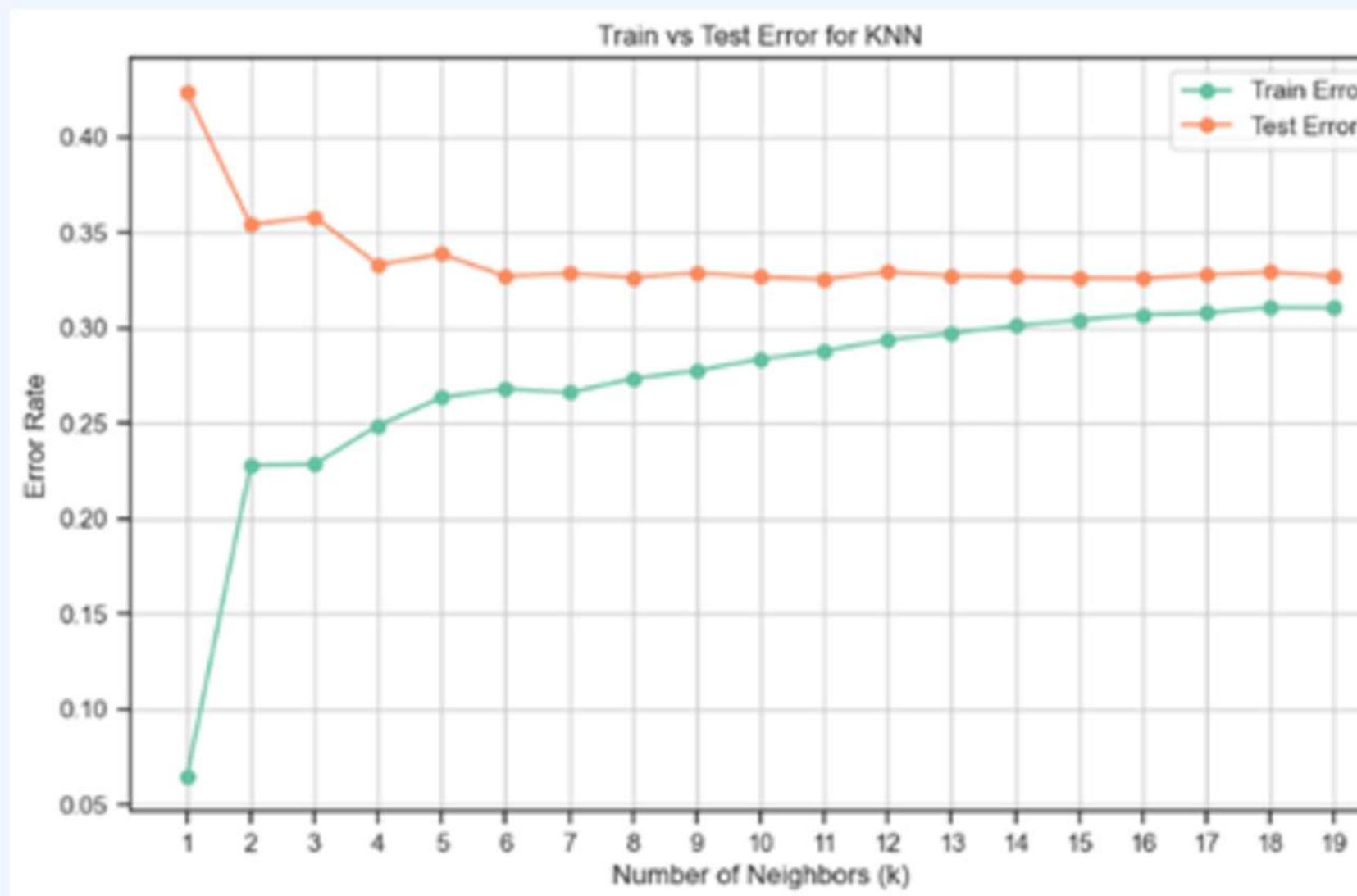
- Class 0 was recognized almost perfectly, showing the high sensitivity
- Class 1 still faces significant challenges ( $F1=0.23$ ), indicating difficulty in distinguishing it from other classes

=> The problem may lie in the data or the characteristics of Class 1, it is necessary to balance the class distribution for better optimization.

# Analytical Model

## KNN

### SELECT OPTIMAL K



- Test error decreases as 'k' increases from 1 to around 9-11
  - Beyond k=11, accuracy improvement is insignificant, possibly plateauing or decreasing
  - k=10 or k=11 provides a good balance, avoiding overfitting (small k) and underfitting (large k), while ensuring stability
- => Two methods will be applied to validate and choose the optimal 'k' for highest accuracy

# Analytical Model

## KNN

### Select the optimal K by Using GridSearchCV

Best Parameters: {'n\_neighbors': 18, 'p': 2, 'weights': 'uniform'}

Test Accuracy of Best KNN: 0.6703

- The KNN model was optimized using GridSearchCV by testing different values of K, distance metric (p), and weighting scheme
  - Weights='uniform' means that all neighbors are treated equally, regardless of distance
- => The KNN model with K=18, p=2, and weights='uniform' achieved the best accuracy on the test set after cross-validation

### Select the optimal K by Test Error

Best k: 11

Test Error for Best k: 0.3254

Evaluation for Best k:

Train Accuracy: 0.7036

Test Accuracy: 0.7363

- K was chosen by identifying the value that resulted in the lowest test error
- Optimizing on test error ensures good model performance on unseen data
- After retraining with the selected K=11, the KNN model exhibited high and stable performance

After applying both methods (GridSearchCV and Test Error analysis) to select the optimal K, we decided to choose K=11 as the optimal value due to the higher test accuracy it provides

# Analytical Model

## KNN

### TRAINING MODEL KNN

Confusion Matrix on Test Set:

```
[[52  2  0]
 [12  3  3]
 [ 4  3 12]]
```

Classification Report on Test Set:

	precision	recall	f1-score	support
0	0.76	0.96	0.85	54
1	0.38	0.17	0.23	18
2	0.80	0.63	0.71	19
accuracy			0.74	91
macro avg	0.65	0.59	0.60	91
weighted avg	0.69	0.74	0.70	91

```
[[43 10  1]
 [10  5  3]
 [ 5  3 11]]
```

	precision	recall	f1-score	support
0	0.74	0.80	0.77	54
1	0.28	0.28	0.28	18
2	0.73	0.58	0.65	19
accuracy			0.65	91
macro avg	0.58	0.55	0.56	91
weighted avg	0.65	0.65	0.65	91

The initial KNN model with weights="uniform" performs well in classifying Class 0 and Class 2, with Class 0 achieving high recall. However, it struggles with Class 1 (low recall), resulting in an overall accuracy of 74%, indicating decent performance overall.

Using weights='distance', the model prioritizes closer neighbors but does not significantly improve Class 1 classification. Overall accuracy drops to 65%, with Class 0 performance slightly decreasing while Class 1 sees some improvement, enhancing class generalization.

# Conclusion

**Random Forest and KNN are the most effective for predicting pregnancy risk levels.**

## Performance

- Random Forest achieves 67% accuracy, comparable to other models.
- Excels in classifying low, mid, and high-risk levels with reduced bias.

1

## Strengths:

- Captures complex relationships among key features (Blood Sugar, SystolicBP, Age).
- Feature importance analysis enhances interpretability.

2

## Applicability:

- Reliable and practical for real-world healthcare scenarios.

3

NO	QUESTION	ANSWER
1	<b>WHAT</b>	Key factors include systolic blood pressure, blood sugar, age, and heart rate. High blood pressure and elevated blood sugar increase risks of preeclampsia and gestational diabetes, while maternal age under 18 or over 35 raises complications. Abnormal heart rates reflect cardiovascular stress, impacting both mother and fetus.
2	<b>WHERE</b>	The model can be applied in hospitals, clinics, and rural areas. In hospitals and clinics, it aids in early risk detection, while in rural areas, integration with IoT devices provides real-time monitoring and support for underserved populations.
3	<b>WHEN</b>	The model is most effective in the first and second trimesters, allowing proactive monitoring and intervention. However, it remains valuable throughout pregnancy for continuous risk assessment.
4	<b>WHO</b>	Primary beneficiaries include pregnant women, who receive early risk detection, and healthcare providers, who gain tools for better decision-making. Unborn children benefit through healthier pregnancies, while policymakers can use insights to improve healthcare systems.
5	<b>WHY</b>	A predictive model identifies high-risk pregnancies early, enabling timely interventions and reducing complications. It supports resource allocation, improves healthcare quality, and prevents maternal and fetal mortality.
6	<b>HOW</b>	The model will be integrated with IoT devices for real-time monitoring and embedded in hospital systems for seamless use. Training programs and community outreach will ensure accessibility in rural and underserved areas, improving maternal healthcare outcomes.

# Actionable Recommendations

## IDEAL AGE

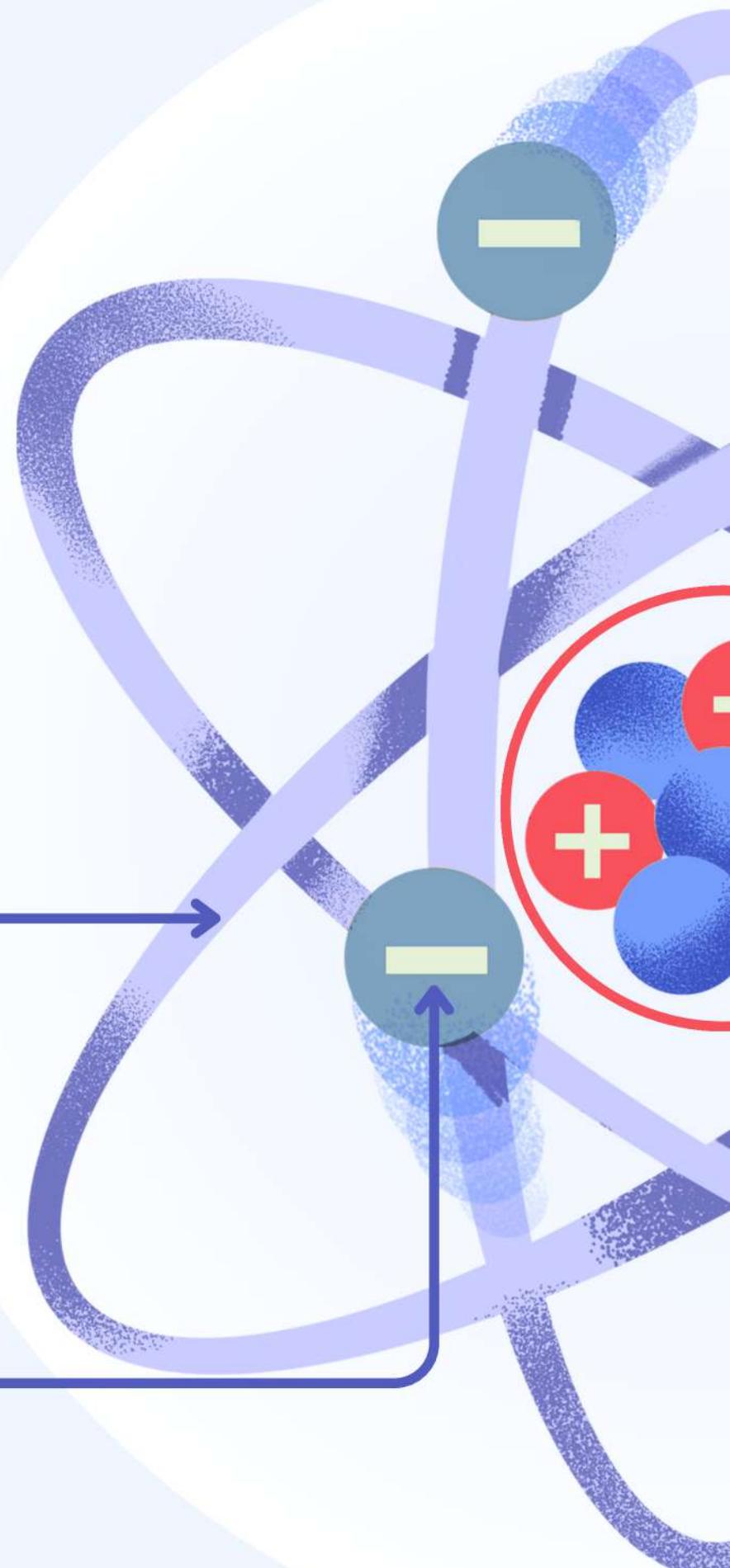
Ideal age: 18–35 years for the lowest risk of complications  
Risks increase significantly for women under 18 or over 35

## SYSTOLIC BLOOD PRESSURE

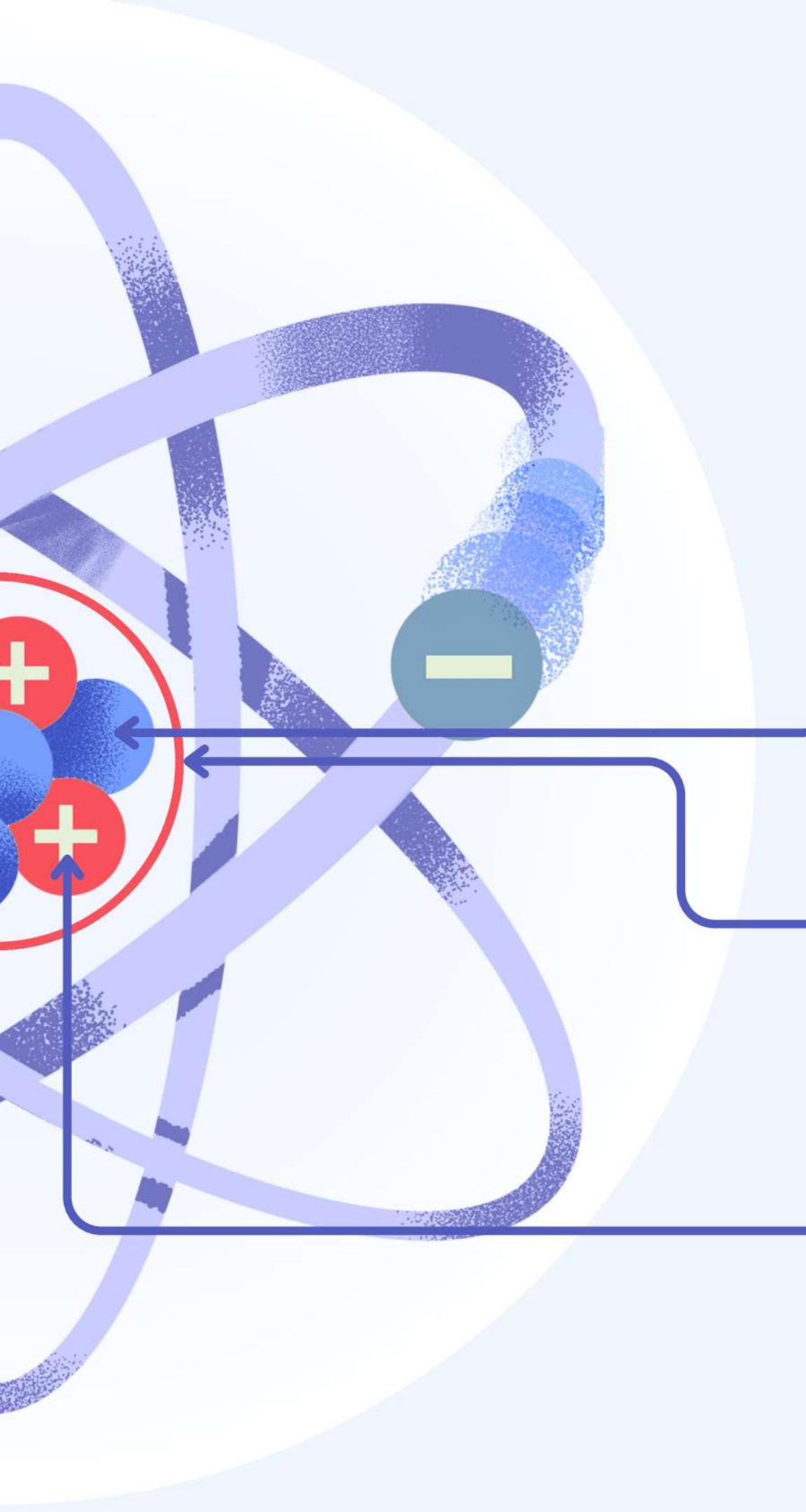
Maintain between 90–120 mmHg to avoid risks like preeclampsia ( $\geq 140$  mmHg) or reduced blood supply ( $< 90$  mmHg).

## DIASTOLIC BLOOD PRESSURE

60–80 mmHg to prevent gestational hypertension ( $\geq 90$  mmHg) or fetal growth issues ( $< 60$  mmHg).



# Actionable Recommendations



## BLOOD SUGAR

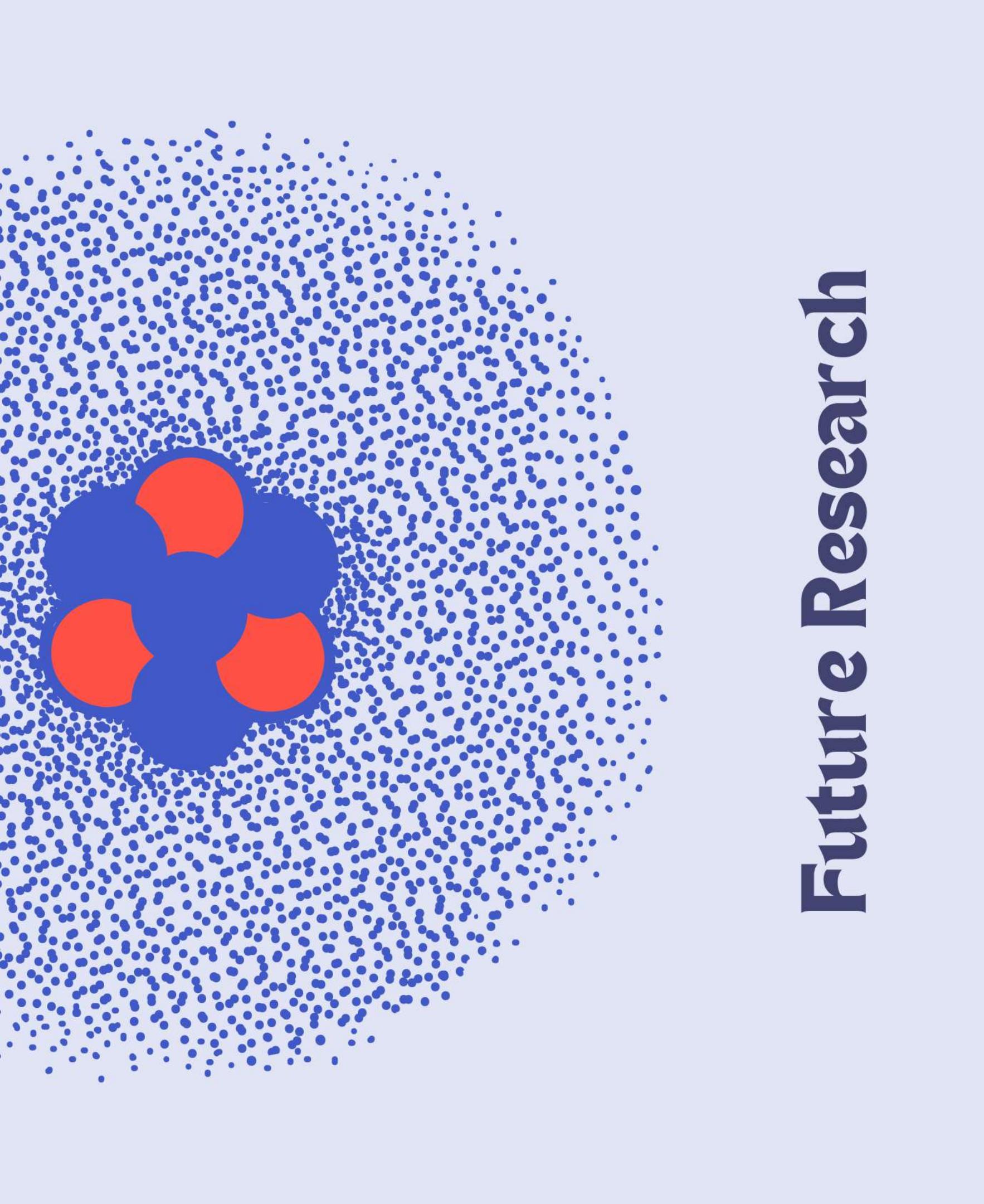
High levels ( $\geq 7.0$  mmol/L): Risk of gestational diabetes;  
low levels ( $< 3.5$  mmol/L): Hypoglycemia risks.

## HEART RATE

Normal range: 60–90 bpm to avoid stress or cardiovascular issues

## CONCLUSION

Regular monitoring and maintaining these indicators within the safe ranges reduce complications and ensure better maternal and fetal health outcomes. Seek medical attention for any deviations.



# Future Research

- 01 Exploration of Advanced Techniques**
- 02 Optimization of the Current Model**
- 03 Addressing Class Imbalance**
- 04 Dataset Expansion**
- 05 Incorporation of Longitudinal Data**
- 06 Real-world Application**

# Thank you for your attention