

# Evaluation Report: Predicting Heart Attack Likelihood

Name: Khushi Naik

## 1. Introduction

The models explored include:

1. **Logistic Regression** (Baseline Model)
2. **Deep Neural Network (DNN) with batch size = 16**
3. **Deep Neural Network (DNN) with batch size = 8**
4. **Ensemble of Feedforward Network (FFN) and Convolutional Neural Network (CNN) with batch size = 8**
5. **Ensemble of FFN and CNN with batch size = 16**
6. **Ensemble of FFN and CNN (Expanded Architecture) with batch size = 32**

Each model's performance is assessed based on accuracy, recall, and AUC-ROC scores. Additionally, this report discusses data preprocessing strategies, training behavior, and the rationale behind each methodological choice.

## 2. Data Preprocessing

### 2.1 Data Cleaning

- **Blood Pressure Feature Engineering:** The dataset originally contained a single column for blood pressure recorded as systolic/diastolic. This was decomposed into two independent variables to enhance interpretability and facilitate differential analysis of systolic and diastolic pressure effects.
  - **Rationale:** Given that systolic and diastolic blood pressures are physiologically distinct and indicative of separate cardiovascular risk factors, isolating them allows the model to infer their independent contributions.
- **Categorical Variable Encoding:**

- **Binary Features:** Label encoding was applied to binary categorical variables such as gender and urban/rural classification.
- **Multiclass Features:** One-hot encoding was employed for variables with multiple categories (e.g., regional classification, stress level) to mitigate the risk of misleading ordinal assumptions.
- **Rationale:** Label encoding preserves a natural binary distinction, while one-hot encoding prevents the model from incorrectly assuming hierarchical relationships among categorical values.

## 2.2 Missing Value Imputation

- Missing values were imputed using domain-specific heuristics:
  - **Numerical Features:** Mean imputation was applied where data was missing at random.
  - **Categorical Features:** Mode imputation was utilized for categorical data to maintain categorical integrity without introducing spurious variance.
  - **Rationale:** Imputation prevents the loss of valuable training instances while maintaining statistical consistency within feature distributions (Little & Rubin, 2002).

## 2.3 Addressing Class Imbalance

- **Synthetic Minority Over-sampling Technique (SMOTE)** was implemented to counteract class imbalance by synthesizing new instances of the minority class.
  - **Rationale:** Unbalanced classes can lead to biased learning where the model disproportionately favors the majority class. SMOTE ensures equitable learning across both classes (Chawla et al., 2002).

## 2.4 Feature Scaling

- Both StandardScaler and RobustScaler were evaluated.
- **RobustScaler** was applied to standardize numerical features, including age, cholesterol, BMI, and heart rate.
  - **Rationale:**
    - RobustScaler is resilient to outliers, unlike traditional normalization methods.
    - Standardized feature magnitudes expedite convergence in gradient-based optimization.
    - Prevents dominance of features with larger numerical ranges.

# 3. Model Architectures

## 3.1 Logistic Regression (Baseline Model)

- A simple linear model used for binary classification.
- **Accuracy: 87.37%**
- **Limitations:** Lacks the ability to model complex, non-linear relationships.

### 3.2 Deep Neural Networks (DNNs)

- Implemented as a **4-layer fully connected network** with:
  - **ReLU activation** for hidden layers
  - **Sigmoid activation** for binary classification
  - **Batch Normalization** for stabilizing training
  - **Dropout layers** to mitigate overfitting

#### Model Variants:

- **DNN (Batch size = 16)** → Balanced generalization and efficiency (**Accuracy: 87.25%**)
- **DNN (Glorot Uniform Initialization, Batch size = 8)** → Smoother weight initialization, improved convergence (**Accuracy: 87.35%**)
- **DNN (He Normal Initialization, Batch size = 8)** → Stable convergence and better learning dynamics (**Accuracy: 87.38%**)

### 3.3 Ensembles of FFN and CNN

- FFNs capture global feature interactions, while CNNs identify local patterns within structured medical data.

#### Tested Variants:

- **Ensemble (Batch Size = 8)** → Accuracy: **87.38%**
- **Ensemble (Batch Size = 16)** → Accuracy: **87.38%**
- **Wider Ensemble (Batch Size = 32)** → Best generalization with **87.44% Accuracy**

## 4. Justification for Neural Network Selection

- **Deep Neural Networks (DNNs)** were selected due to their ability to model complex, non-linear relationships in medical data, capturing subtle interactions between features that traditional models like logistic regression fail to recognize.
- **Batch Normalization** was included to stabilize learning by reducing internal covariate shift, leading to improved convergence speed.
- **Dropout layers** were integrated to prevent overfitting, particularly critical in medical datasets where data scarcity is common.
- **ReLU activation** was preferred for hidden layers due to its non-linearity and computational efficiency, promoting sparse activation and reducing vanishing gradient issues.

- **He Normal and Glorot Uniform Initializations** were used to ensure better weight initialization, aiding in stable and efficient model convergence.
- **CNN-based ensemble models** were chosen for their ability to extract spatial and local patterns, enhancing predictive performance through feature augmentation.

## 5. Model Training & Performance Analysis

### 5.1 Training Curve Analysis

- **Early Epochs:** A steep decline in loss indicates effective feature learning.
- **Plateau Phase (~10 epochs):** Models stabilized, suggesting diminishing returns with additional training.
- **Validation Loss Drop:** Learning rate adjustments improved model convergence.
- **Overfitting Risks:** Smaller batch sizes (8) introduced variability but improved generalization, whereas batch size 32 exhibited signs of overfitting.

### 5.2 Model Performance Comparison

Model	Accuracy	Precision	Recall	AUC-ROC
Logistic Regression	87.37%	100%	74%	0.87
DNN (Batch=16)	87.25%	100%	74%	0.87
DNN (Glorot Init)	87.35%	100%	73.9%	0.87
DNN (He Normal Init)	87.38%	100%	74%	0.87
FFN+CNN Ensemble (Batch=8)	87.38%	100%	73.9%	0.87
FFN+CNN Ensemble (Batch=16)	87.38%	100%	74.1%	0.87
FFN+CNN Ensemble (Batch=32)	<b>87.44%</b>	<b>100%</b>	<b>74.3%</b>	<b>0.87</b>

### 5.3 Observations

- **Ensemble models performed best** by leveraging FFN's high-level abstractions and CNN's pattern recognition capabilities.
- **Batch size 32 showed slight overfitting**, though it achieved the highest validation accuracy.
- **DNNs with He Normal Initialization** exhibited the most stable convergence among neural network models.

## 6. Conclusion

- **Logistic Regression, though interpretable, is not optimal** for capturing complex feature interactions.
- **Deep Neural Networks significantly enhance predictive power** by modeling intricate relationships.
- **Ensemble models (FFN+CNN) provide the highest accuracy (87.44%)**, benefiting from complementary feature extraction techniques.
- **Careful selection of batch size and weight initialization is critical** for optimizing model performance.

## 7. Recommendations

- **For interpretability:** Logistic Regression is preferable.
- **For best predictive accuracy:** FFN+CNN ensemble with batch size 32.
- **For balanced performance and training efficiency:** DNN with batch size 16.

## References

Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research*, 16, 321-357.

Little, R. J. A., & Rubin, D. B. (2002). *Statistical Analysis with Missing Data*. Wiley.

Pedregosa, F., et al. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830.

Snoek, J., Larochelle, H., & Adams, R. P. (2012). Practical Bayesian Optimization of Machine Learning Algorithms. *Advances in Neural Information Processing Systems*, 25.

Vaswani, A., et al. (2017). Attention is All You Need. *Advances in Neural Information Processing Systems*, 30.