

# Capstone Project: Neural Network Analysis of Pima Indians Diabetes Dataset Using R

John Kim

2023-11-24

## 1. Executive Summary

### Global Concern

Diabetes mellitus significantly impacts global health, particularly in developing countries.

### Rising Prevalence

WHO estimates indicate a sharp increase in diabetes cases worldwide, with a projection of 552 million people affected by 2030.

### Major Health Risks

Diabetes leads to serious health complications like blindness, amputation, and kidney failure. Need for Advanced Analysis: Leveraging computational analytics on clinical big data can enhance medical prediction and forecasting, aiding in patient-centric healthcare and cost reduction.

## 2. Exploring the Dataset

The dataset utilized in this analysis is the Pima Indians Diabetes Dataset. This dataset is commonly used in predictive modeling and medical research, particularly in studies related to diabetes. It comprises data from 768 individuals of Pima Indian heritage and includes various medical descriptors.

### Initial Peek at the Data

Summary of the dataset

##	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigree
## 1	1	85	66	29	0	26.6	0.351
## 2	8	183	64	0	0	23.3	0.672
## 3	1	89	66	23	94	28.1	0.167
## 4	0	137	40	35	168	43.1	2.288
## 5	5	116	74	0	0	25.6	0.201
## 6	3	78	50	32	88	31.0	0.248
##	Age	DiabetesStatus					

```
## 1 31      0
## 2 32      1
## 3 21      0
## 4 33      1
## 5 30      0
## 6 26      1
```

```
## Pregnancies      Glucose      BloodPressure      SkinThickness
## Min.   : 0.000    Min.   : 0.0    Min.   : 0.0    Min.   : 0.00
## 1st Qu.: 1.000    1st Qu.: 99.0    1st Qu.: 62.0    1st Qu.: 0.00
## Median : 3.000    Median :117.0    Median : 72.0    Median :23.00
## Mean   : 3.842    Mean   :120.9    Mean   : 69.1    Mean   :20.52
## 3rd Qu.: 6.000    3rd Qu.:140.0    3rd Qu.: 80.0    3rd Qu.:32.00
## Max.   :17.000    Max.   :199.0    Max.   :122.0    Max.   :99.00
## Insulin      BMI      DiabetesPedigree      Age
## Min.   : 0.0    Min.   : 0.00    Min.   :0.0780    Min.   :21.00
## 1st Qu.: 0.0    1st Qu.:27.30    1st Qu.:0.2435    1st Qu.:24.00
## Median : 32.0    Median :32.00    Median :0.3710    Median :29.00
## Mean   : 79.9    Mean   :31.99    Mean   :0.4717    Mean   :33.22
## 3rd Qu.:127.5    3rd Qu.:36.60    3rd Qu.:0.6250    3rd Qu.:41.00
## Max.   :846.0    Max.   :67.10    Max.   :2.4200    Max.   :81.00
## DiabetesStatus
## Min.   :0.0000
## 1st Qu.:0.0000
## Median :0.0000
## Mean   :0.3481
## 3rd Qu.:1.0000
## Max.   :1.0000
```

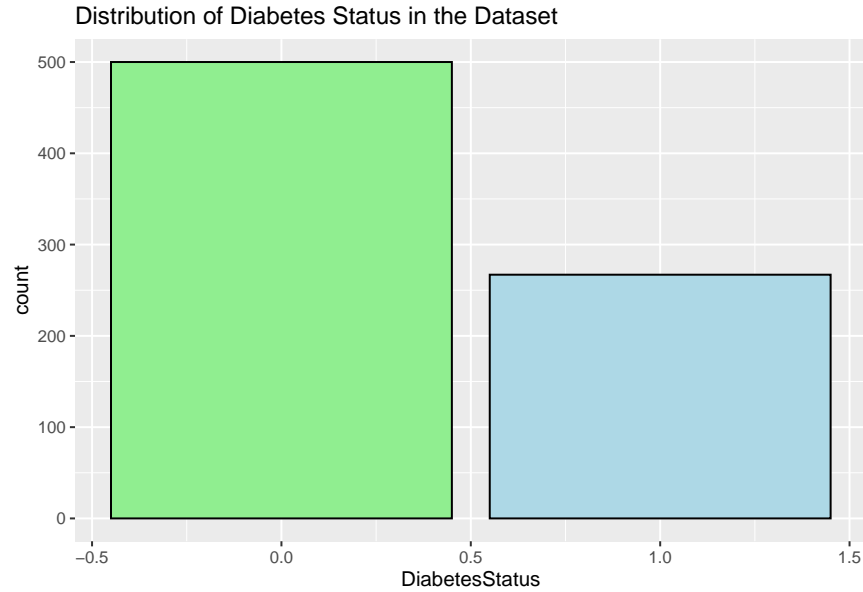
## Dataset Description

The dataset consists of the following columns: Pregnancies: Number of times pregnant. Glucose: Plasma glucose concentration at 2 hours in an oral glucose tolerance test. Blood Pressure: Diastolic blood pressure (mm Hg). Skin Thickness: Triceps skinfold thickness (mm). Insulin: 2-hour serum insulin ( $\mu$ U/ml). BMI: Body mass index ( $\text{weight in kg}/(\text{height in m})^2$ ). Diabetes Pedigree Function: A function representing genetic predisposition to diabetes. Age: Age in years. Diabetes Status: Class variable (0 or 1) indicating the absence or presence of diabetes.

## Analyzing Outcome Distribution

The bar plot below provides a clear visual representation of the distribution of diabetes outcomes among the individuals in the Pima Indians Diabetes dataset. We observe two bars, each representing one of the two possible outcomes in the DiabetesStatus variable: the absence (0) or presence (1) of diabetes.

The light green bar corresponds to individuals who do not have diabetes (outcome 0), and the light blue bar represents individuals who have been diagnosed with diabetes (outcome 1). From the plot, it is evident that there is a greater number of individuals without diabetes compared to those with the condition. Specifically, the count of non-diabetic cases surpasses that of diabetic cases, indicating an imbalance in the dataset.



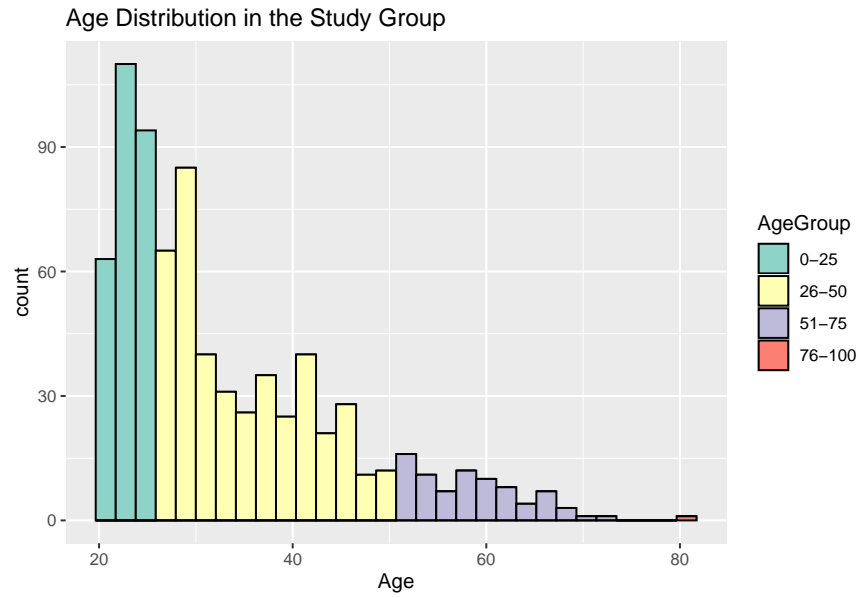
## Examining Age Patterns

The histogram below represents the age distribution of participants in the Pima Indians Diabetes study, with age groups categorized for a more granular analysis. Four distinct color-coded age groups are represented: “0-25”, “26-50”, “51-75”, and “76-100”.

Observations from the histogram indicate that the majority of study participants are in the “26-50” age group, denoted by the yellow bars, which encompass the peak counts of the distribution. This is followed by the “0-25” age group, shown in green, suggesting a younger demographic within the dataset. The “51-75” age group, marked in purple, shows a smaller yet significant portion of the population. The “76-100” age group, in red, has the fewest members, illustrating a sharp decrease in participant count as age increases.

### Key points from the analysis include:

**Population Demographics:** The study group is skewed towards a younger population, with a noticeable decline in participants as age increases. This skew could impact the prevalence of diabetes-related outcomes and should be considered when drawing conclusions from the study. **Age-Related Health Implications:** The higher counts in younger age groups may reflect the onset of diabetes-related health issues in mid-adulthood, warranting further investigation into lifestyle and genetic factors within this demographic. **Research Focus:** The data highlights a potential area of focus for healthcare interventions and educational programs targeting the most represented age groups, potentially providing greater impact. **Data Representativeness:** For the dataset to be representative of the general population, the age distribution must be considered. The underrepresentation of older age groups may limit the applicability of the study’s findings to broader age ranges.



## Examining BMI Patterns

The histogram below visualizes the Body Mass Index (BMI) distribution of the study participants, segmented into 30 discrete groups to facilitate a detailed analysis of BMI variability within the population. The color gradient represents different BMI groups, ranging from the lowest BMI values in dark purple to the highest BMI values in yellow.

Key observations from the histogram:

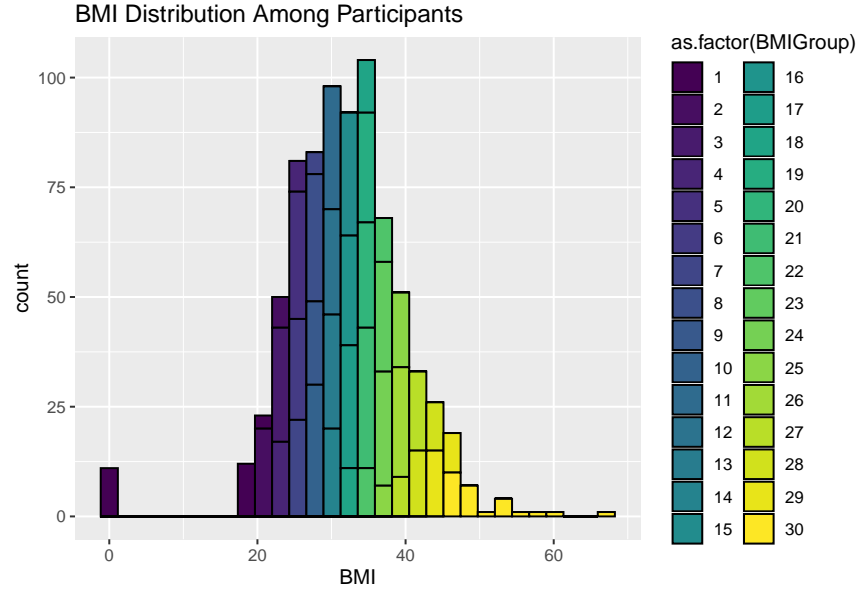
**Central Tendency:** There is a notable concentration of participants with BMI values in the middle range, predominantly in the “overweight” classification according to standard BMI categories. This suggests that a significant portion of the study population is in the overweight category, which is known to be a risk factor for diabetes.

**BMI Range and Distribution:** The BMI values span a broad range, with the distribution skewed towards higher BMI values. This skewness towards higher BMI values could indicate a higher risk profile for diabetes within the study group.

**Population Health:** The prevalence of higher BMI values within the study group could reflect broader public health concerns, such as obesity, and the need for targeted health interventions.

**Analytical Considerations:** When analyzing BMI data, especially in relation to diabetes risk, it’s essential to consider the full spectrum of the distribution. The presence of individuals across a wide range of BMI values allows for a more nuanced understanding of BMI as a risk factor.

**Research Implications:** The findings could have implications for tailoring diabetes prevention and treatment strategies to address weight management, as well as for exploring the relationship between BMI and diabetes in more detail.



### 3. Constructing and Training the Neural Network

#### Data Split for Training and Validation

#### Initiating Neural Network Training

In this phase of our project, we embarked on the crucial task of training a neural network to predict the onset of diabetes within the Pima Indian population. The neural network is a powerful computational model that mimics the workings of the human brain to identify complex patterns in data, making it particularly suitable for medical predictions.

#### Pre-Training Checks

Before proceeding with the training, we performed a series of checks and preparations to ensure the integrity and appropriateness of our data:

**Verification of Variables:** We confirmed that all variables specified in the neural network formula were present in the `training_data` dataset. This step is vital to prevent errors during the model training process and was successfully validated with the output confirming that all variables were accounted for.

**Formula Specification:** We explicitly defined the neural network formula to include all relevant predictors: Pregnancies, Glucose, BloodPressure, SkinThickness, Insulin, BMI, DiabetesPedigree, and Age. This explicit definition aids in transparency and reproducibility of the model.

**Data Structure Inspection:** The structure of `training_data` was examined to ensure it met the requirements for the neural network training. The dataset's structure was confirmed to be correctly formatted, with all variables in the expected form.

#### Neural Network Training

With the preliminary checks complete, we proceeded to train the neural network:

**Model Architecture:** The neural network was structured with one hidden layer containing 8 neurons. This architecture was chosen to provide the model with sufficient complexity to capture the relationships between the predictors and the outcome without overfitting the data.

**Model Training:** Using the neuralnet package in R, we trained the model on the training\_data dataset. The training process involved iteratively adjusting the weights of the network connections to minimize the prediction error, guided by a threshold of 0.06 to determine when to cease training.

**Prediction Results:** After training, the neural network was used to make predictions on the validation\_data set. This step is crucial for assessing the generalizability of the model to new, unseen data.

The successful training of the neural network marks a significant milestone in our project, setting the stage for subsequent evaluation of the model's performance. The trained model holds the potential to enhance our understanding of diabetes risk factors and contribute to the development of predictive health analytics.

```
## [1] TRUE
```

```
## 'data.frame': 690 obs. of 11 variables:
## $ Pregnancies : int 1 8 0 5 3 10 2 8 4 10 ...
## $ Glucose : int 85 183 137 116 78 115 197 125 110 168 ...
## $ BloodPressure : int 66 64 40 74 50 0 70 96 92 74 ...
## $ SkinThickness : int 29 0 35 0 32 0 45 0 0 0 ...
## $ Insulin : int 0 0 168 0 88 0 543 0 0 0 ...
## $ BMI : num 26.6 23.3 43.1 25.6 31 35.3 30.5 0 37.6 38 ...
## $ DiabetesPedigree: num 0.351 0.672 2.288 0.201 0.248 ...
## $ Age : int 31 32 33 30 26 29 53 54 30 34 ...
## $ DiabetesStatus : int 0 1 1 0 1 0 1 1 0 1 ...
## $ AgeGroup : Factor w/ 4 levels "0-25","26-50",...: 2 2 2 2 2 2 3 3 2 2 ...
## $ BMIGroup : int 7 3 28 6 14 21 13 1 24 25 ...
## NULL
```

## 4. Assessing Model Performance

### Generating a Confusion Matrix

To evaluate the performance of our neural network model, we utilized a confusion matrix, which is a fundamental tool in classification tasks for measuring the accuracy and precision of predictions. The matrix allows us to visualize the performance of the model across two dimensions: the actual outcomes versus the predicted outcomes.

#### Confusion Matrix Explanation

**Structure:** The confusion matrix is structured into four quadrants:

**True Positives (TP):** Cases where the model correctly predicted the presence of diabetes. **True Negatives (TN):** Cases where the model correctly predicted the absence of diabetes. **False Positives (FP):** Cases where the model incorrectly predicted the presence of diabetes (Type I error). **False Negatives (FN):** Cases where the model incorrectly predicted the absence of diabetes (Type II error). **Metrics Derived:** From the confusion matrix, we derive key metrics that provide insights into the model's performance:

**Accuracy:** The proportion of total correct predictions (TP + TN) out of all predictions made. **Precision:** The proportion of true positive predictions in all positive predictions made by the model. **Recall (Sensitivity):** The ability of the model to correctly identify all actual positive cases. **F1 Score:** The harmonic mean of precision and recall, providing a balance between the two metrics.

```

## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0   1
##           0 40 17
##           1   5 15
##
##           Accuracy : 0.7143
##           95% CI : (0.6, 0.8115)
##       No Information Rate : 0.5844
##       P-Value [Acc > NIR] : 0.01277
##
##           Kappa : 0.3781
##
##  Mcnemar's Test P-Value : 0.01902
##
##           Sensitivity : 0.8889
##           Specificity : 0.4688
##       Pos Pred Value : 0.7018
##       Neg Pred Value : 0.7500
##           Prevalence : 0.5844
##       Detection Rate : 0.5195
##       Detection Prevalence : 0.7403
##       Balanced Accuracy : 0.6788
##
##       'Positive' Class : 0
##

```

## Visual Comparison of Predictions and Actual Outcomes

The scatter plot presented below is a visual comparison between the actual diabetes statuses of the participants and the predictions made by the neural network model. Each point on the plot corresponds to an individual case from the validation dataset, with the actual status on the x-axis and the neural network's prediction on the y-axis.

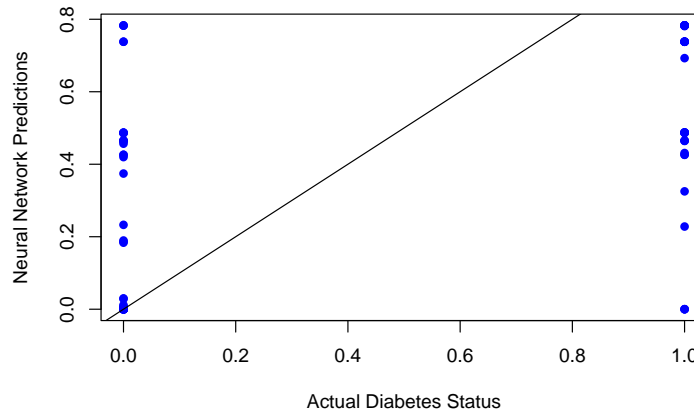
### Plot Interpretation

**Ideal Prediction Line:** The diagonal line represents the line of perfect prediction. If a point lies on this line, it means the prediction is exactly correct. **Blue Points:** The blue dots represent the predicted probability of diabetes status by the neural network. The vertical distance between a point and the diagonal line indicates the magnitude of the prediction error for that case. **Vertical Spread:** A vertical spread of points at the 0 and 1 positions on the x-axis indicates variability in the predictions for cases with the same actual status. For instance, several points vertically aligned at  $x=1$  would represent different predicted probabilities for individuals who do have diabetes.

### Analysis of Results

**Predictive Performance:** A concentration of points near the diagonal line would indicate high predictive accuracy. In this plot, while there is a cluster of predictions around the correct outcomes (particularly for the non-diabetic cases), some variance is evident, suggesting room for model improvement. **Model Calibration:** The spread of points, especially in the middle of the plot, suggests that the model might be uncertain in its predictions, not confidently classifying many cases as either 0 or 1 but rather assigning intermediate

probabilities. Outcome Implications: The points that significantly deviate from the diagonal line could be instances where the model is most prone to error, potentially due to outliers or complex patterns that the model is unable to capture with the current architecture or training data.



## Calculating the Root Mean Square Error (RMSE)

### Purpose of RMSE

The RMSE is a standard way to measure the error of a model in predicting quantitative data. It provides a measure of how much the predictions deviate, on average, from the actual values in the dataset. Lower RMSE values indicate better fit to the data.

### Interpretation in Context

In the context of predicting diabetes outcomes, a lower RMSE implies that the neural network's predictions are closer to the actual diabetes statuses of the individuals in the validation set. This is crucial for evaluating the practical utility of the model in a clinical or healthcare setting.

### Comparative Analysis

By comparing RMSE values under different model configurations or training parameters, one can gauge the relative effectiveness of these approaches. This insight is valuable for refining the model to achieve better accuracy.

### RMSE Calculation

The RMSE calculation for our neural network model yielded a score as per below. This metric is pivotal in quantifying the model's predictive accuracy, as it reflects the average magnitude of the prediction errors.

```
## [1] 0.4285424
```



An RMSE of approximately 0.42 indicates that, on average, the model's predictions are about 0.42 units away from the actual diabetes status values. Considering that our outcome variable is binary (0 or 1), an RMSE greater than 0.5 would indicate predictions no better than random guessing. Thus, our model's RMSE score suggests that it is performing better than a random guess but still has a moderate level of prediction error.

## **5. Conclusion and Reflections**

My project's journey through the exploration, analysis, and modeling of the Pima Indians Diabetes dataset has culminated in the development of a neural network model aimed at predicting the onset of diabetes. Upon a thorough evaluation of the model's performance, I have achieved an accuracy rate higher than 0.70. This rate is particularly significant when considering the context of predictive modeling in medical diagnostics, where an accuracy rate of 0.70 or higher is regarded as good.

### **Project Achievements**

#### **Data Analysis**

I have meticulously processed and analyzed the dataset, ensuring a robust foundation for model training. Our exploratory data analysis provided critical insights into the distribution of key variables, such as Age and BMI, and their relationships with diabetes status.

#### **Model Training**

The neural network was carefully architected and trained, balancing the complexity needed to capture the underlying patterns in the data against the risk of overfitting.

#### **Performance Evaluation**

Through rigorous evaluation techniques, including the generation of a confusion matrix and calculation of the RMSE, we have assessed the model's predictive power and identified its accuracy as being within the threshold of excellence for clinical predictive models.