

**University of San Carlos**  
**DCISM**

# **Final Culminating Activity**

MATH 3109

**PREPARED BY:**  
Woodrow A. Fajardo

**SUBMITTED TO:**  
Prof. Elmer S. Poliquit

**MATH-3109**  
**NOVEMBER 2025**

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Project Objective	1
1.2	Literature Review	1
1.3	Raw Data	1
<b>2</b>	<b>Data Cleaning</b>	<b>1</b>
2.1	Converting Zeroes Values to NA	1
2.2	MICE Imputation	2
2.3	Imputation Validation	3
<b>3</b>	<b>Exploratory Analysis</b>	<b>7</b>
3.1	Correlation Matrix	7
3.2	Analyzing Linearity	7
<b>4</b>	<b>Modelling Logistic Regression</b>	<b>7</b>
4.1	Model Specification	7
4.2	Odds Ratio Interpretation	8
4.3	Model Performance	8
<b>5</b>	<b>Linear Regression (Predict Insulin)</b>	<b>8</b>
5.1	B-Spline Justification	9
5.2	Comparison Linear vs. B-Spline	10
5.3	Residual Diagnostics	11
<b>6</b>	<b>Discussions</b>	<b>12</b>
6.1	Validation of Literature Claims	12
6.2	Key Findings	12
6.3	Limitations (if any)	12
<b>7</b>	<b>Appendix</b>	<b>12</b>

Raw Data based from Kaggle:

## 1 Introduction

### 1.1 Project Objective

### 1.2 Literature Review

### 1.3 Raw Data

#### STATISTICAL MODELING

```
1 library(mice)
2 library(corrplot)
3 library(ggplot2)
4 library(splines)
5 library(performance)
6 library(see)
7 library(ggpubr)
8
9 setwd("~/School/MATH-3109/FINAL-1")
10 raw_data <- read.csv("raw_dataset.csv")
11 View(raw_data)
12 head(raw_data)
```

#### TERMINAL OUTPUT

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI
↪	DiabetesPedigreeFunction	Age	Outcome			
1	6	148	72	35	0	33.6
↪	0.627	50	1			
2	1	85	66	29	0	26.6
↪	0.351	31	0			
3	8	183	64	0	0	23.3
↪	0.672	32	1			
4	1	89	66	23	94	28.1
↪	0.167	21	0			
5	0	137	40	35	168	43.1
↪	2.288	33	1			
6	5	116	74	0	0	25.6
↪	0.201	30	0			

## 2 Data Cleaning

### 2.1 Converting Zeroes Values to NA

#### STATISTICAL MODELING

```
1 colSums(raw_data == 0)
```

## TERMINAL OUTPUT

Pregnancies	Glucose	BloodPressure	BMI	Outcome
↪ SkinThickness	Insulin			
↪ DiabetesPedigreeFunction		Age		
NA	NA		NA	
↪ NA	NA		NA	
↪ 0	0		500	

## STATISTICAL MODELING

```

1 col_w_zeroes <- c("Pregnancies", "Glucose", "BloodPressure", "SkinThickness",
  ↪ "Insulin", "BMI")
2 raw_data[col_w_zeroes] <- lapply(raw_data[col_w_zeroes], function(x) ifelse(x ==
  ↪ 0, NA, x))
3
4 summary(raw_data)
```

## TERMINAL OUTPUT

Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin
↪ BMI	DiabetesPedigreeFunction	Age	Outcome	
Min. : 1.000	Min. : 44.0	Min. : 24.00	Min. : 7.00	Min. : 14.00
↪ Min. : 18.20	Min. : 0.0780	Min. : 21.00	Min. : 0.000	
1st Qu.: 2.000	1st Qu.: 99.0	1st Qu.: 64.00	1st Qu.: 22.00	1st Qu.: 76.25
↪ 1st Qu.: 27.50	1st Qu.: 0.2437	1st Qu.: 24.00	1st Qu.: 0.000	
Median : 4.000	Median : 117.0	Median : 72.00	Median : 29.00	Median : 125.00
↪ Median : 32.30	Median : 0.3725	Median : 29.00	Median : 0.000	
Mean : 4.495	Mean : 121.7	Mean : 72.41	Mean : 29.15	Mean : 155.55
↪ Mean : 32.46	Mean : 0.4719	Mean : 33.24	Mean : 0.349	
3rd Qu.: 7.000	3rd Qu.: 141.0	3rd Qu.: 80.00	3rd Qu.: 36.00	3rd Qu.: 190.00
↪ 3rd Qu.: 36.60	3rd Qu.: 0.6262	3rd Qu.: 41.00	3rd Qu.: 1.000	
Max. : 17.000	Max. : 199.0	Max. : 122.00	Max. : 99.00	Max. : 846.00
↪ Max. : 67.10	Max. : 2.4200	Max. : 81.00	Max. : 1.000	
NA's : 111	NA's : 5	NA's : 35	NA's : 227	NA's : 374
↪ NA's : 11				

## 2.2 MICE Imputation

## STATISTICAL MODELING

```

1 init <- mice(raw_data, maxit = 0)
2 methods <- init$method #$$
3 methods[col_w_zeroes] <- "pmm"
4 imputed_data <- mice(raw_data, method = methods, m = 5, maxit = 5, seed = 123)
5 final_data <- complete(imputed_data, 1)
6 head(final_data)
```

View full R script output here.

## STATISTICAL MODELING

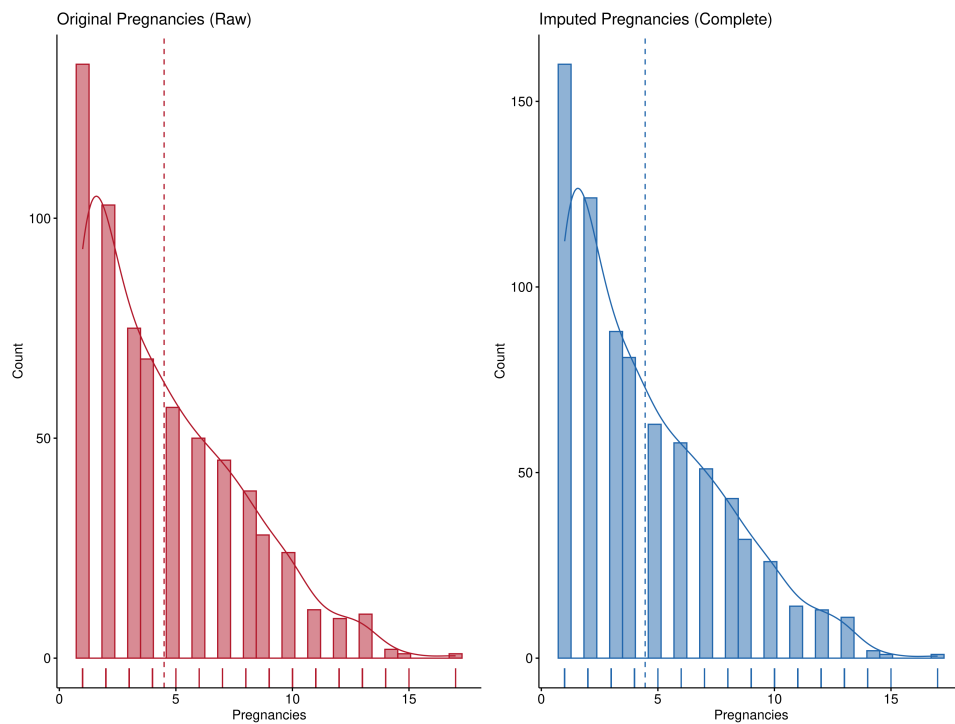
```
1 anyNA(final_data)
```

## TERMINAL OUTPUT

```
anyNA(final_data)
```

## 2.3 Imputation Validation

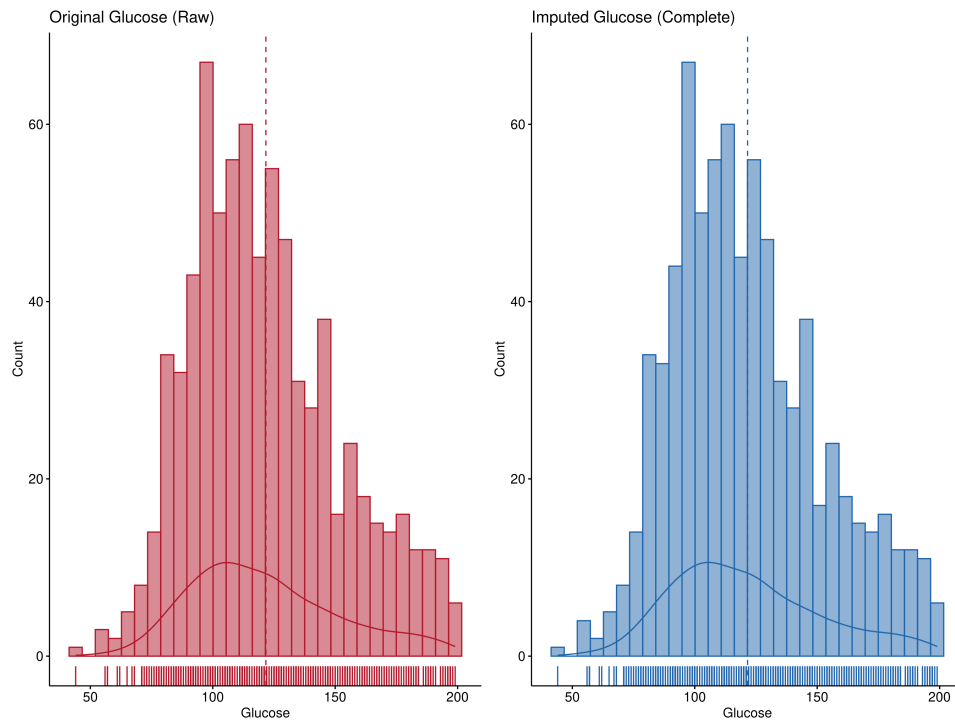
Refer to Appendix Section for the full R script.



**FIGURE 1**

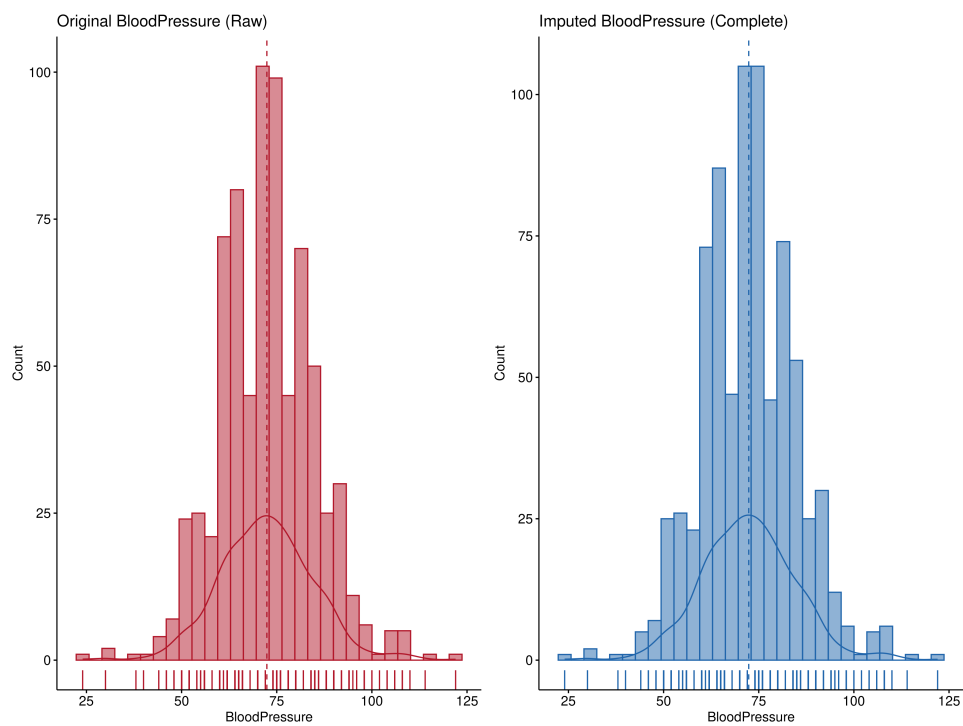
A standard centered figure.

Source: Generated by LaTeX



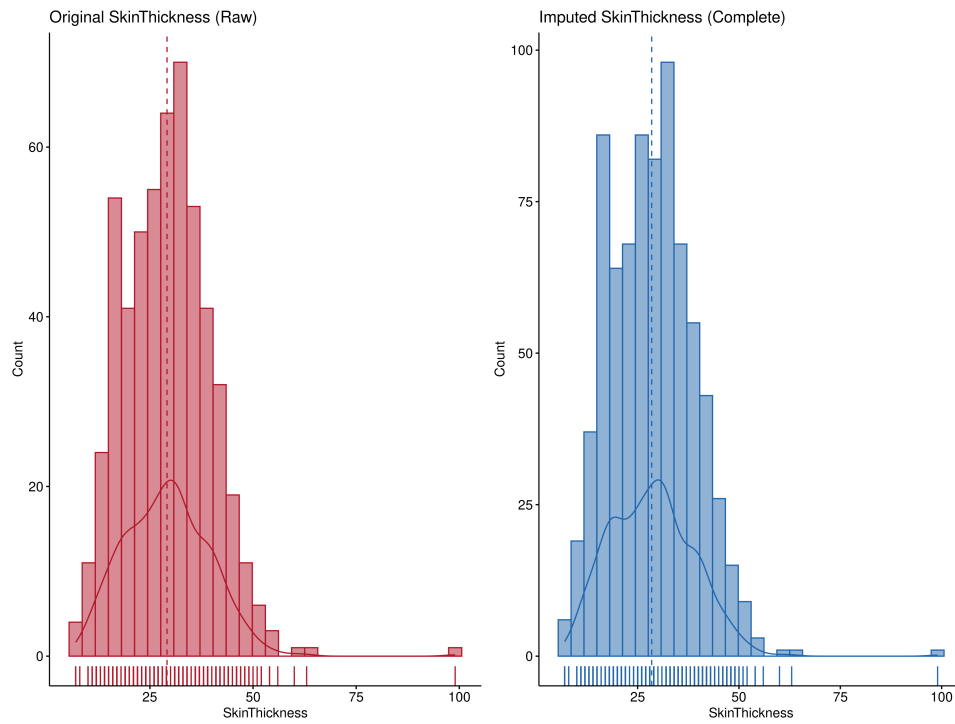
**FIGURE 2**  
A standard centered figure.

*Source: Generated by LaTeX*



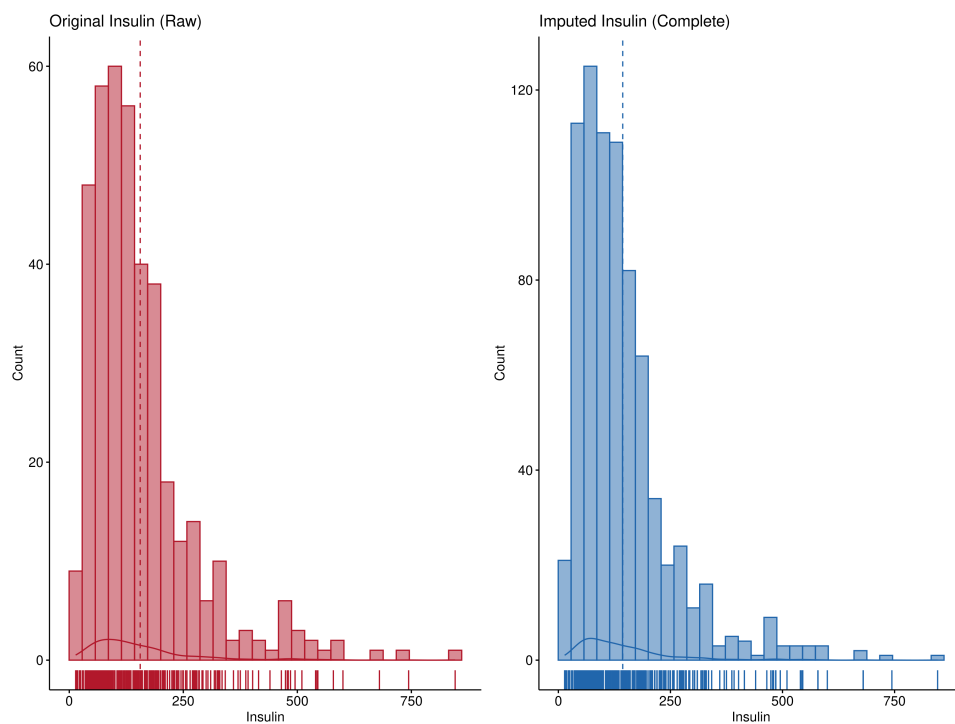
**FIGURE 3**  
A standard centered figure.

*Source: Generated by LaTeX*



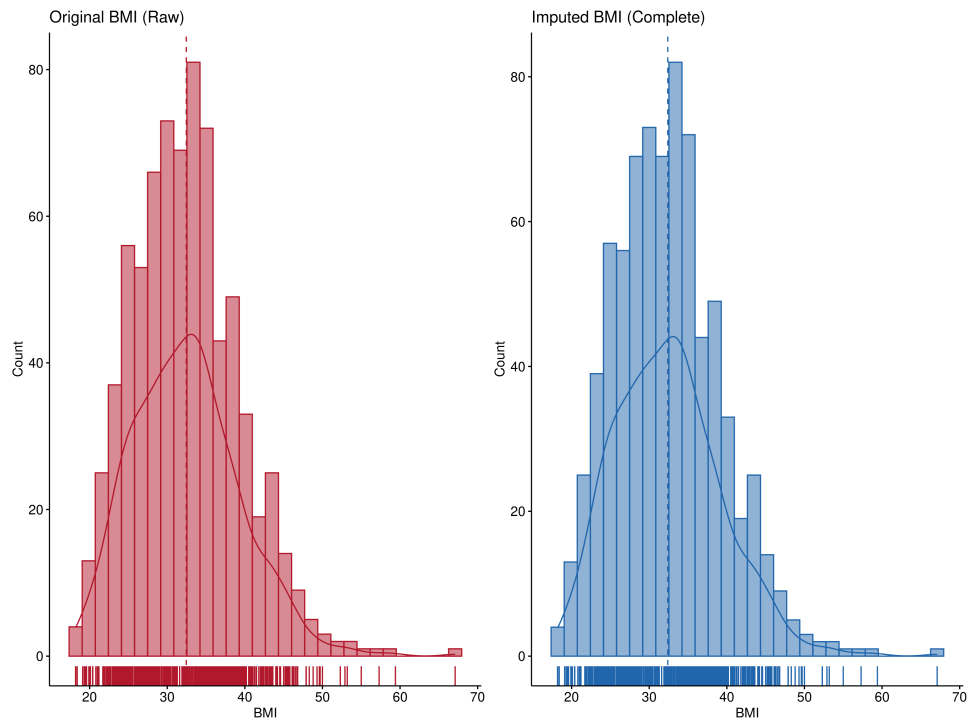
**FIGURE 4**  
A standard centered figure.

*Source: Generated by LaTeX*



**FIGURE 5**  
A standard centered figure.

*Source: Generated by LaTeX*



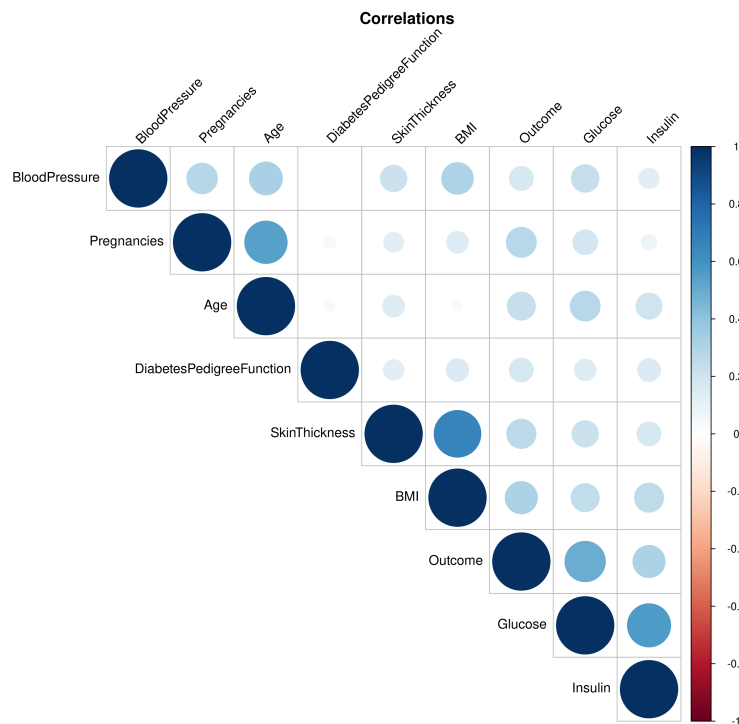
**FIGURE 6**  
A standard centered figure.

*Source: Generated by LaTeX*



### 3 Exploratory Analysis

#### 3.1 Correlation Matrix



**FIGURE 7**

A standard centered figure.

Source: Generated by LaTeX

#### STATISTICAL MODELING

```
1 cor_matrix <- cor(final_data, use = "complete.obs")
2 corrpplot(cor_matrix, method = "circle", type = "upper", order = "hclust",
3           tl.col = "black", tl.srt = 45, title = "Correlations", mar=c(0,0,1,0))
```

#### 3.2 Analyzing Linearity

### 4 Modelling Logistic Regression

#### 4.1 Model Specification

#### STATISTICAL MODELING

```
1 set.seed(123)
2 sample_index <- sample(1:nrow(final_data), 0.8 * nrow(final_data))
3 train_data <- final_data[sample_index, ]
4 test_data <- final_data[-sample_index, ]
5
6 log_model <- glm(Outcome ~ Glucose + BMI + Age + Pregnancies +
7                 ↪ DiabetesPedigreeFunction,
8                 data = train_data, family= "binomial")
9 summary(log_model)
```

## TERMINAL OUTPUT

```
glm(formula = Outcome ~ Glucose + BMI + Age + Pregnancies +
     ↪ DiabetesPedigreeFunction,
     family = "binomial", data = train_data)

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept)    -9.200623   0.809675 -11.363 < 2e-16 ***
Glucose         0.037413   0.004099   9.128 < 2e-16 ***
BMI             0.079814   0.016719   4.774 1.81e-06 ***
Age            0.009594   0.010003   0.959 0.337506
Pregnancies     0.130485   0.037653   3.465 0.000529 ***
DiabetesPedigreeFunction 0.659625   0.321580   2.051 0.040247 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 796.42  on 613  degrees of freedom
Residual deviance: 568.82  on 608  degrees of freedom
AIC: 580.82

Number of Fisher Scoring iterations: 5
```

## 4.2 Odds Ratio Interpretation

## STATISTICAL MODELING

```
1 exp(coef(log_model))
```

## TERMINAL OUTPUT

(Intercept)	Glucose	BMI
↪ Age	Pregnancies DiabetesPedigreeFunction	
0.0001009765	1.0381215630	1.0830855110
↪ 1.0096399484	1.1393803946	1.9340669910

## 4.3 Model Performance

## 5 Linear Regression (Predict Insulin)

## STATISTICAL MODELING

```
1 linear_simple <- lm(Insulin ~ Glucose + BMI + Age, data = final_data)
2 summary(linear_simple)
```

## TERMINAL OUTPUT

```
Call:
lm(formula = Insulin ~ Glucose + BMI + Age, data = final_data)

Residuals:
    Min       1Q   Median       3Q      Max
-289.02  -45.91  -14.70   25.01  567.44

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -167.6456    19.1120  -8.772  < 2e-16 ***
Glucose       1.8543     0.1125  16.481  < 2e-16 ***
BMI           2.0651     0.4791   4.310 1.84e-05 ***
Age           0.5693     0.2841   2.004  0.0455 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 89.01 on 764 degrees of freedom
Multiple R-squared:  0.3401,    Adjusted R-squared:  0.3375
F-statistic: 131.2 on 3 and 764 DF,  p-value: < 2.2e-16
```

## STATISTICAL MODELING

```
1 linear_spline <- lm (Insulin ~ Glucose + BMI + bs(Age, degree = 3), data =
  ↪ final_data)
2 summary(linear_spline)
```

## TERMINAL OUTPUT

```
Call:
lm(formula = Insulin ~ Glucose + BMI + bs(Age, degree = 3), data = final_data)

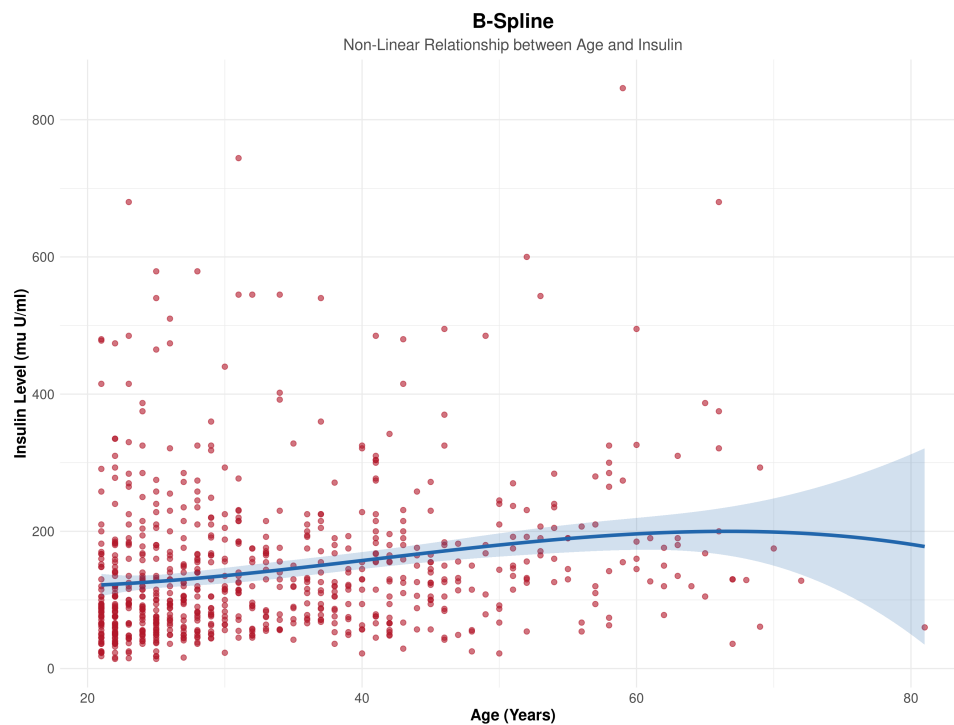
Residuals:
    Min       1Q   Median       3Q      Max
-299.47  -47.78  -14.42   26.00  554.90

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  -151.3071    18.2230  -8.303 4.61e-16 ***
Glucose        1.8612     0.1123  16.571  < 2e-16 ***
BMI            2.2633     0.4871   4.646 3.98e-06 ***
bs(Age, degree = 3)1  -57.2504    32.8240  -1.744  0.0815 .
bs(Age, degree = 3)2   80.9994    50.5390   1.603  0.1094
bs(Age, degree = 3)3   5.4944    62.7432   0.088  0.9302
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 88.82 on 762 degrees of freedom
Multiple R-squared:  0.3447,    Adjusted R-squared:  0.3404
F-statistic: 80.16 on 5 and 762 DF,  p-value: < 2.2e-16
```

## 5.1 B-Spline Justification

Refer to Appendix Section for the full R script.

**FIGURE 8**

A standard centered figure.

Source: Generated by LaTeX

## 5.2 Comparison Linear vs. B-Spline

### STATISTICAL MODELING

```
1 linear_simple <- lm(Insulin ~ Glucose + BMI + Age, data = final_data)
2 linear_spline <- lm (Insulin ~ Glucose + BMI + bs(Age, degree = 3), data =
  ↪ final_data)
3 anova(linear_simple, linear_spline)
```

### TERMINAL OUTPUT

Analysis of Variance Table

Model 1: Insulin ~ Glucose + BMI + Age

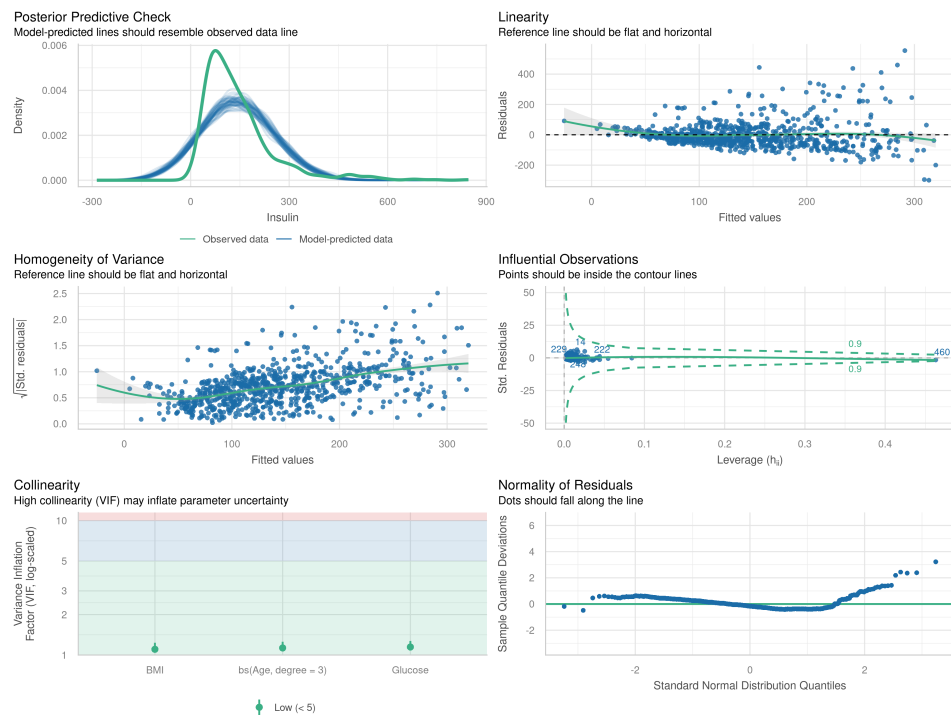
Model 2: Insulin ~ Glucose + BMI + bs(Age, degree = 3)

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	764	6052997				
2	762	6010880	2	42116	2.6696	0.06993 .

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

### 5.3 Residual Diagnostics



**FIGURE 9**

A standard centered figure.

Source: Generated by LaTeX

#### STATISTICAL MODELING

```
1 image_rend <- check_model(linear_spline)
2 png("splien_model.png", width = 12, height = 9, units = "in", res = 300)
3 print(image_rend)
4 dev.off()
```

## 6 Discussions

### 6.1 Validation of Literature Claims

### 6.2 Key Findings

### 6.3 Limitations (if any)

## 7 Appendix

### HISTOGRAM COMPARISON CODE EXAMPLE

```

1 p1 <- gg_histogram(raw_data, x = "Pregnancies",
2   title = "Original Pregnancies (Raw)",
3   xlab = "Pregnancies", ylab = "Count",
4   fill = "#b2182b",
5   color = "#b2182b",
6   add = "mean", rug = TRUE, add_density = TRUE,
7   ggtheme = theme_pubr()
8 )
9
10 p2 <- gg_histogram(final_data, x = "Pregnancies",
11   title = "Imputed Pregnancies (Complete)",
12   xlab = "Pregnancies", ylab = "Count",
13   fill = "#2166ac",
14   color = "#2166ac",
15   add = "mean", rug = TRUE, add_density = TRUE,
16   ggtheme = theme_pubr()
17 )
18
19 ggarrange(p1, p2, ncol = 2, nrow = 1)

```

### STATISTICAL MODELING

```

1 p1 <- gg_histogram(raw_data, x = "Glucose",
2   title = "Original Glucose (Raw)",
3   xlab = "Glucose", ylab = "Count",
4   fill = "#b2182b",
5   color = "#b2182b",
6   add = "mean", rug = TRUE, add_density = TRUE,
7   ggtheme = theme_pubr()
8 )
9
10 p2 <- gg_histogram(final_data, x = "Glucose",
11   title = "Imputed Glucose (Complete)",
12   xlab = "Glucose", ylab = "Count",
13   fill = "#2166ac",
14   color = "#2166ac",
15   add = "mean", rug = TRUE, add_density = TRUE,
16   ggtheme = theme_pubr()
17 )
18
19 ggarrange(p1, p2, ncol = 2, nrow = 1)

```

## STATISTICAL MODELING

```

1 p1 <- gg_histogram(raw_data, x = "BloodPressure",
2   title = "Original BloodPressure (Raw)",
3   xlab = "BloodPressure", ylab = "Count",
4   fill = "#b2182b",
5   color = "#b2182b",
6   add = "mean", rug = TRUE, add_density = TRUE,
7   ggtheme = theme_pubr()
8 )
9
10 p2 <- gg_histogram(final_data, x = "BloodPressure",
11   title = "Imputed BloodPressure (Complete)",
12   xlab = "BloodPressure", ylab = "Count",
13   fill = "#2166ac",
14   color = "#2166ac",
15   add = "mean", rug = TRUE, add_density = TRUE,
16   ggtheme = theme_pubr()
17 )
18
19 ggarrange(p1, p2, ncol = 2, nrow = 1)

```

## STATISTICAL MODELING

```

1 p1 <- gg_histogram(raw_data, x = "SkinThickness",
2   title = "Original SkinThickness (Raw)",
3   xlab = "SkinThickness", ylab = "Count",
4   fill = "#b2182b",
5   color = "#b2182b",
6   add = "mean", rug = TRUE, add_density = TRUE,
7   ggtheme = theme_pubr()
8 )
9
10 p2 <- gg_histogram(final_data, x = "SkinThickness",
11   title = "Imputed SkinThickness (Complete)",
12   xlab = "SkinThickness", ylab = "Count",
13   fill = "#2166ac",
14   color = "#2166ac",
15   add = "mean", rug = TRUE, add_density = TRUE,
16   ggtheme = theme_pubr()
17 )
18
19 ggarrange(p1, p2, ncol = 2, nrow = 1)

```

## STATISTICAL MODELING

```

1 p1 <- gghistogram(raw_data, x = "Insulin",
2   title = "Original Insulin (Raw)",
3   xlab = "Insulin", ylab = "Count",
4   fill = "#b2182b",
5   color = "#b2182b",
6   add = "mean", rug = TRUE, add_density = TRUE,
7   ggtheme = theme_pubr()
8 )
9
10 p2 <- gghistogram(final_data, x = "Insulin",
11   title = "Imputed Insulin (Complete)",
12   xlab = "Insulin", ylab = "Count",
13   fill = "#2166ac",
14   color = "#2166ac",
15   add = "mean", rug = TRUE, add_density = TRUE,
16   ggtheme = theme_pubr()
17 )
18
19 ggarrange(p1, p2, ncol = 2, nrow = 1)

```

## STATISTICAL MODELING

```

1 p1 <- gghistogram(raw_data, x = "BMI",
2   title = "Original BMI (Raw)",
3   xlab = "BMI", ylab = "Count",
4   fill = "#b2182b",
5   color = "#b2182b",
6   add = "mean", rug = TRUE, add_density = TRUE,
7   ggtheme = theme_pubr()
8 )
9
10 p2 <- gghistogram(final_data, x = "BMI",
11   title = "Imputed BMI (Complete)",
12   xlab = "BMI", ylab = "Count",
13   fill = "#2166ac",
14   color = "#2166ac",
15   add = "mean", rug = TRUE, add_density = TRUE,
16   ggtheme = theme_pubr()
17 )
18
19 ggarrange(p1, p2, ncol = 2, nrow = 1)

```



## STATISTICAL MODELING

```
1 ggplot(final_data, aes(x = Age, y = Insulin)) +  
2   geom_point(color = "#b2182b", alpha = 0.6, size = 2) +  
3   geom_smooth(method = "lm",  
4               formula = y ~ bs(x, degree = 3),  
5               color = "#2166ac",  
6               fill = "#2166ac",  
7               alpha = 0.2,  
8               size = 1.5) +  
9  
10  labs(title = "B-Spline",  
11       subtitle = "Non-Linear Relationship between Age and Insulin",  
12       x = "Age (Years)",  
13       y = "Insulin Level (mu U/ml)") +  
14  
15  theme_minimal() +  
16  theme(  
17    plot.title = element_text(face = "bold", hjust = 0.5, size = 18),  
18    plot.subtitle = element_text(hjust = 0.5, color = "gray30", size = 14),  
19    axis.title = element_text(face = "bold", size = 14),  
20    axis.text = element_text(size = 12))
```