

# Detailed Multivariate Classification

## Libraries I used

```
library(ggplot2)
library(GGally)

## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg      ggplot2

library(patchwork)
library(heplots)

## Loading required package: broom

library(biotools)

## Loading required package: MASS
##
## Attaching package: 'MASS'
##
## The following object is masked from 'package:patchwork':
##
##   area
##
## ---
## biotools version 4.3
##
## Attaching package: 'biotools'
##
## The following object is masked from 'package:heplots':
##
##   boxM

library(car)

## Loading required package: carData

library(caret)

## Loading required package: lattice

library(pROC)

## Type 'citation("pROC")' for a citation.
##
## Attaching package: 'pROC'
##
## The following objects are masked from 'package:stats':
##
##   cov, smooth, var
```

```
library(e1071)
library(ggcorrplot)
```

## 1. Multivariate Normal

Let

$$\begin{pmatrix} Y_1 \\ Y_2 \\ X_1 \\ X_2 \end{pmatrix} \sim N_4(\mu, \Sigma)$$

where

$$\mu = \begin{pmatrix} 20 \\ -15 \\ 17 \\ -10 \end{pmatrix}, \quad \Sigma = \begin{bmatrix} 2 & 1 & -1 & 0 \\ 1 & 3 & 1 & -1 \\ -1 & 1 & 4 & 0 \\ 0 & -1 & 0 & 2 \end{bmatrix}$$

(a) Find the conditional expectation  $E[Y_1|X_1 = x_1, X_2 = x_2]$

**Solution:**

In order to answer this question, Let's define some intermediate random variables and call them A and B so that  $A = Y_1$  and  $B = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$ . The conditional distribution of  $A = Y_1$  given  $B = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$  is given by:

$$A|B \sim N(\mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(B - \mu_B), \quad \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}) \quad (1)$$

As first step, let's find the mean values used in the equation by using  $\mu$  provided in the question:

$$\mu_A = \mu_{Y_1} = 20, \quad \mu_B = \begin{pmatrix} \mu_{X_1} \\ \mu_{X_2} \end{pmatrix} = \begin{pmatrix} 17 \\ -10 \end{pmatrix} \quad (2)$$

Now, let's find the covariance data from  $\Sigma$  given in the question:

$$\Sigma_{BB} = \begin{bmatrix} \Sigma_{X_1X_1} & \Sigma_{X_1X_2} \\ \Sigma_{X_2X_1} & \Sigma_{X_2X_2} \end{bmatrix} = \begin{bmatrix} 4 & 0 \\ 0 & 2 \end{bmatrix} \quad (3)$$

$$\Sigma_{AB} = [\Sigma_{Y_1X_1} \quad \Sigma_{Y_1X_2}] = [-1 \quad 0] \quad (4)$$

$$\Sigma_{AA} = \Sigma_{Y_1Y_1} = 2 \quad (5)$$

The remaining step is to computing  $\Sigma_{BB}^{-1}$  by using the formula for the inverse of a 2x2 matrix:

$$\Sigma_{BB}^{-1} = \frac{1}{(4)(2) - (0)(0)} \begin{bmatrix} 2 & 0 \\ 0 & 4 \end{bmatrix} = \frac{1}{8} \begin{bmatrix} 2 & 0 \\ 0 & 4 \end{bmatrix} \quad (6)$$

Now, we can calculate the expected value of conditional distribution as in (7).

$$E[Y_1|X_1 = x_1, X_2 = x_2] = \mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(B - \mu_B) \quad (7)$$

Let's compute  $B - \mu_B$ :

$$B - \mu_B = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} - \begin{pmatrix} 17 \\ -10 \end{pmatrix} = \begin{pmatrix} x_1 - 17 \\ x_2 + 10 \end{pmatrix} \quad (8)$$

Let's multiply  $\Sigma_{AB}$  we found in equation (4) with  $\Sigma_{BB}^{-1}$  we found in equation (6):

$$\Sigma_{AB}\Sigma_{BB}^{-1} = \begin{bmatrix} -1 & 0 \end{bmatrix} \cdot \frac{1}{8} \begin{bmatrix} 2 & 0 \\ 0 & 4 \end{bmatrix} \quad (9)$$

$$\Sigma_{AB}\Sigma_{BB}^{-1} = \frac{1}{8} [(-1)(2) + (0)(0), (-1)(0) + (0)(4)] \quad (10)$$

$$\Sigma_{AB}\Sigma_{BB}^{-1} = \frac{1}{8} \begin{bmatrix} -2 & 0 \end{bmatrix} \quad (11)$$

Now, I am multiplying what I found in equation (11) with  $B - \mu_B$  we found in equation (8):

$$= \frac{1}{8} \begin{bmatrix} -2 & 0 \end{bmatrix} \begin{bmatrix} x_1 - 17 \\ x_2 + 10 \end{bmatrix} \quad (12)$$

$$= \frac{-1}{4} \cdot (x_1 - 17) \quad (13)$$

Now, we can calculate the expected value of conditional distribution as in (7).

$$E[Y_1|X_1 = x_1, X_2 = x_2] = 20 - \frac{1}{4} \cdot (x_1 - 17) \quad (14)$$

**(b)** Find the conditional expectation  $E[Y_2|X_1 = x_1, X_2 = x_2]$

**Solution:**

In order to answer this question, Let's define some intermediate random variables and call them A and B so that  $A = Y_2$  and  $B = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$ . The conditional distribution of  $A = Y_2$  given  $B = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$  is given by:

$$A|B \sim N(\mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(B - \mu_B), \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}) \quad (15)$$

As first step, let's find the mean values used in the equation by using  $\mu$  provided in the question:

$$\mu_A = \mu_{Y_2} = -15, \quad \mu_B = \begin{pmatrix} \mu_{X_1} \\ \mu_{X_2} \end{pmatrix} = \begin{pmatrix} 17 \\ -10 \end{pmatrix} \quad (16)$$

Now, let's find the covariance data from  $\Sigma$  given in the question:

$$\Sigma_{BB} = \begin{bmatrix} \Sigma_{X_1X_1} & \Sigma_{X_1X_2} \\ \Sigma_{X_2X_1} & \Sigma_{X_2X_2} \end{bmatrix} = \begin{bmatrix} 4 & 0 \\ 0 & 2 \end{bmatrix} \quad (17)$$

$$\Sigma_{AB} = \begin{bmatrix} \Sigma_{Y_2X_1} & \Sigma_{Y_2X_2} \end{bmatrix} = \begin{bmatrix} 1 & -1 \end{bmatrix} \quad (18)$$

$$\Sigma_{AA} = \Sigma_{Y_2Y_2} = 3 \quad (19)$$

The remaining step is to computing  $\Sigma_{BB}^{-1}$  by using the formula for the inverse of a 2x2 matrix:

$$\Sigma_{BB}^{-1} = \frac{1}{(4)(2) - (0)(0)} \begin{bmatrix} 2 & 0 \\ 0 & 4 \end{bmatrix} = \frac{1}{8} \begin{bmatrix} 2 & 0 \\ 0 & 4 \end{bmatrix} \quad (20)$$

Now, we can calculate the expected value of conditional distribution as in (7).

$$E[Y_1|X_1 = x_1, X_2 = x_2] = \mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(B - \mu_B) \quad (21)$$

Let's compute  $B - \mu_B$ :

$$B - \mu_B = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} - \begin{pmatrix} 17 \\ -10 \end{pmatrix} = \begin{pmatrix} x_1 - 17 \\ x_2 + 10 \end{pmatrix} \quad (22)$$

Let's multiply  $\Sigma_{AB}$  we found in equation (18) with  $\Sigma_{BB}^{-1}$  we found in equation (20):

$$\Sigma_{AB}\Sigma_{BB}^{-1} = \begin{bmatrix} 1 & -1 \end{bmatrix} \cdot \frac{1}{8} \begin{bmatrix} 2 & 0 \\ 0 & 4 \end{bmatrix} \quad (23)$$

$$\Sigma_{AB}\Sigma_{BB}^{-1} = \frac{1}{8} [(1)(2) + (-1)(0), (1)(0) + (-1)(4)] \quad (24)$$

$$\Sigma_{AB}\Sigma_{BB}^{-1} = \frac{1}{8} \begin{bmatrix} 2 & -4 \end{bmatrix} \quad (25)$$

Now, I am multiplying what I found in equation (25) with  $B - \mu_B$  we found in equation (22):

$$= \frac{1}{8} \begin{bmatrix} 2 & -4 \end{bmatrix} \begin{bmatrix} x_1 - 17 \\ x_2 + 10 \end{bmatrix} \quad (26)$$

$$= \frac{1}{4} \cdot (x_1 - 17 - 2 \cdot x_2 - 20) = \frac{1}{4} \cdot (x_1 - 2 \cdot x_2 - 37) \quad (27)$$

Now, we can calculate the expected value of conditional distribution as in (21).

$$E[Y_2|X_1 = x_1, X_2 = x_2] = -15 + \frac{1}{4} \cdot (x_1 - 2 \cdot x_2 - 37) \quad (28)$$

(c) Find the conditional expectation  $E[(Y_1, Y_2)^T|X_1 = x_1, X_2 = x_2]$

**Solution:**

In order to answer this question, Let's define some intermediate random variables and call them A and B so that  $A = \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix}$  and  $B = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$ . The conditional distribution of  $A = \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix}$  given  $B = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}$  is given by:

$$A|B \sim N(\mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(B - \mu_B), \quad \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}) \quad (29)$$

As first step, let's find the mean values used in the equation by using  $\mu$  provided in the question:

$$\mu_A = \begin{pmatrix} \mu_{Y_1} \\ \mu_{Y_2} \end{pmatrix} = \begin{pmatrix} 20 \\ -15 \end{pmatrix}, \quad \mu_B = \begin{pmatrix} \mu_{X_1} \\ \mu_{X_2} \end{pmatrix} = \begin{pmatrix} 17 \\ -10 \end{pmatrix} \quad (30)$$

Now, let's find the covariance data from  $\Sigma$  given in the question:

$$\Sigma_{BB} = \begin{bmatrix} \Sigma_{X_1 X_1} & \Sigma_{X_1 X_2} \\ \Sigma_{X_2 X_1} & \Sigma_{X_2 X_2} \end{bmatrix} = \begin{bmatrix} 4 & 0 \\ 0 & 2 \end{bmatrix} \quad (31)$$

$$\Sigma_{AB} = \begin{bmatrix} \Sigma_{Y_1 X_1} & \Sigma_{Y_1 X_2} \\ \Sigma_{Y_2 X_1} & \Sigma_{Y_2 X_2} \end{bmatrix} = \begin{bmatrix} -1 & 0 \\ 1 & -1 \end{bmatrix} \quad (32)$$

$$\Sigma_{AA} = \begin{bmatrix} \Sigma_{Y_1 Y_1} & \Sigma_{Y_1 Y_2} \\ \Sigma_{Y_2 Y_1} & \Sigma_{Y_2 Y_2} \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ 1 & 3 \end{bmatrix} \quad (33)$$

The remaining step is to computing  $\Sigma_{BB}^{-1}$  by using the formula for the inverse of a 2x2 matrix:

$$\Sigma_{BB}^{-1} = \frac{1}{(4)(2) - (0)(0)} \begin{bmatrix} 2 & 0 \\ 0 & 4 \end{bmatrix} = \frac{1}{8} \begin{bmatrix} 2 & 0 \\ 0 & 4 \end{bmatrix} \quad (34)$$

Now, we can calculate the expected value of conditional distribution as in (35).

$$E[(Y_1, Y_2)^T | X_1 = x_1, X_2 = x_2] = \mu_A + \Sigma_{AB} \Sigma_{BB}^{-1} (B - \mu_B) \quad (35)$$

Let's compute  $B - \mu_B$ :

$$B - \mu_B = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} - \begin{pmatrix} 17 \\ -10 \end{pmatrix} = \begin{pmatrix} x_1 - 17 \\ x_2 + 10 \end{pmatrix} \quad (36)$$

Let's multiply  $\Sigma_{AB}$  we found in equation (32) with  $\Sigma_{BB}^{-1}$  we found in equation (34):

$$\Sigma_{AB} \Sigma_{BB}^{-1} = \begin{bmatrix} -1 & 0 \\ 1 & -1 \end{bmatrix} \cdot \frac{1}{8} \begin{bmatrix} 2 & 0 \\ 0 & 4 \end{bmatrix} \quad (37)$$

$$\Sigma_{AB} \Sigma_{BB}^{-1} = \frac{1}{8} \begin{bmatrix} -2 & 0 \\ 2 & -4 \end{bmatrix} \quad (38)$$

Now, I am multiplying what I found in equation (38) with  $B - \mu_B$  we found in equation (36):

$$= \frac{1}{8} \begin{bmatrix} -2 & 0 \\ 2 & -4 \end{bmatrix} \begin{bmatrix} x_1 - 17 \\ x_2 + 10 \end{bmatrix} \quad (39)$$

$$= \frac{1}{4} \cdot \begin{bmatrix} -x_1 + 17 \\ x_1 - 2 \cdot x_2 - 37 \end{bmatrix} \quad (40)$$

Now, we can calculate the expected value of conditional distribution as in (35).

$$E[(Y_1, Y_2)^T | X_1 = x_1, X_2 = x_2] = \begin{bmatrix} 20 \\ -15 \end{bmatrix} + \frac{1}{4} \cdot \begin{bmatrix} -x_1 + 17 \\ x_1 - 2 \cdot x_2 - 37 \end{bmatrix} \quad (41)$$

(d) What can you infer from your results in (a), (b), and (c)?

**Solution:**

By looking at the results in parts a,b and c, I can infer that the conditional expectations of  $Y_1$  and  $Y_2$  given  $X_1$  and  $X_2$  are linear functions of the predictors. Here are my observations:

- $Y_1$  depends only on  $X_1$ , and not on  $X_2$ . This is clear from part (a), where the conditional expectation of  $Y_1$  involves only  $x_1$  and not  $x_2$ . This also appears in the joint expression in part (c), where the first component depends only on  $x_1$ . This matches the fact that the covariance between  $Y_1$  and  $X_2$  is zero:  $\text{Cov}(Y_1, X_2) = 0$ . So, knowing  $X_2$  doesn't give any extra information about  $Y_1$  once we know  $X_1$ .
- On the other hand,  $Y_2$  depends on both  $X_1$  and  $X_2$ , as shown in part (b). The relationship is positive with  $X_1$  and more strongly negative with  $X_2$ . This reflects the cross-covariances:  $\text{Cov}(Y_2, X_1) > 0$ , and  $\text{Cov}(Y_2, X_2) < 0$ .
- Result (c) ties these together and shows how both  $Y_1$  and  $Y_2$  can be written together as a vector-valued conditional expectation. The structure is linear in  $(X_1, X_2)$ , and shifts away from their marginal means depending on the values of the predictors. This behavior aligns with a key property of multivariate normal distributions: that the conditional distributions are also normal with linear conditional means.

Overall, the conditional expectation functions directly reflect the covariance structure between the response variables and the predictors, and provide insight into which predictors influence each response.

(e) What are the implications of your findings in (d) for regression with multiple responses?

**Solution:**

The results in (d) show that when we have multiple response variables (like  $Y_1$  and  $Y_2$ ), we can think of each one as having its own linear relationship with the predictors. But importantly, each response might depend on different combinations of predictors.

For example, we saw that  $Y_1$  only depends on  $X_1$ , while  $Y_2$  depends on both  $X_1$  and  $X_2$ . This tells us that in multiple response regression, it's not always the case that every predictor affects every response — the structure can vary and is determined by the covariance relationships.

Also, in multivariate normal settings, the conditional expectations are linear functions of the predictors, just like in linear regression. So this supports the idea that **regression with multiple responses can be seen as a system of linear regressions**, where each response has its own coefficients, and those coefficients are related to the covariance structure of the joint distribution.

Another important implication is that if some covariances are zero, then the corresponding predictors don't affect that response. This is useful in practice for simplifying models and understanding which predictors are actually important.

So overall, the findings from (d) help explain how multiple response regression models work, and how they're directly connected to the joint distribution of the variables.

## 2. Additive Model and Weighted Estimation

Let  $Y = f(X) + \epsilon$  with  $\epsilon$  independent of  $X$ , and let  $W \in \mathbb{R}$  be positive weights.

(a) Find  $\hat{f}(x) = \arg \min_f E[W \cdot (Y - f(x))^2 | X = x]$

**Solution:**

In this question, we are trying to find a function  $f(x)$  that gives the minimum of the  $E[W \cdot (Y - f(x))^2 | X = x]$ .

- Let's first factor out the weight  $W$ . Since  $W$  is positive weight, either a constant or a random variable independent of  $f(x)$ . Then, minimizing  $E[W \cdot (Y - f(x))^2 | X = x]$  is actually same thing with minimizing  $E[(Y - f(x))^2 | X = x]$ . We can do this since  $W$  is only scaling the value and will not change the point where minimum occurs. Therefore, it becomes a basic MMSE problem now. To minimize this:

$$E[(Y - f(x))^2 | X = x]$$

Let's rewrite this in the following way:

$$E[(Y^2 - 2 \cdot Y \cdot f(x) + f(x))^2 | X = x]$$

By using the linearity of the expectation:

$$E[(Y^2|X = x] - 2E[Y \cdot f(x)|X = x] + E[f(x)^2|X = x]$$

Since  $f(x)$  is not random when given  $X=x$ ,  $f(x)$  is just an number, not a variable. So, I can rewrite the expected value expression in this way:

$$E[(Y^2|X = x] - 2f(x) \cdot E[Y|X = x] + f(x)^2$$

Let's say that:

$$V = E[(Y^2|X = x] \text{ and } u = E[Y|X = x]$$

If I rewrite the expression by using  $V$  and  $u$ :

$$V - 2uf(x) + f(x)^2$$

We basically want to minimize this. Let's take the derivative of this expression with respect to  $f(x)$  and make it equal to the 0:

$$\frac{df(x)}{dx}(V - 2uf(x) + f(x)^2) = 0 \quad -2u - 2f(x) = 0 \quad \hat{f}(x) = u = E[Y|X = x]$$

(b) Find  $\hat{f}(x) = \arg \min_f E[W \cdot |Y - f(x)| | X = x]$

**Solution:**

Let's start by splitting the absolute value, this is the problem of Mean Absoulte Error Proof: The absolute value function is not differentiable at 0, so we split it into two cases:

$$|Y - f(x)| = \begin{cases} f(x) - Y & \text{if } Y < f(x) \\ Y - f(x) & \text{if } Y > f(x) \end{cases}$$

We can then write the expected value as two integrals using the conditional density  $p(y | x)$ :

$$\mathbb{E}[|Y - f(x)| | X = x] = \int_{-\infty}^{f(x)} (f(x) - y) p(y | x) dy + \int_{f(x)}^{\infty} (y - f(x)) p(y | x) dy$$

To minimize the total error, we take the derivative with respect to  $f(x)$ .

Differentiate each part:

- First integral:

$$\frac{d}{df(x)} \int_{-\infty}^{f(x)} (f(x) - y) p(y | x) dy = \int_{-\infty}^{f(x)} p(y | x) dy$$

- Second integral:

$$\frac{d}{df(x)} \int_{f(x)}^{\infty} (y - f(x)) p(y | x) dy = - \int_{f(x)}^{\infty} p(y | x) dy$$

Add the two results and set the derivative equal to zero:

$$\int_{-\infty}^{f(x)} p(y | x) dy = \int_{f(x)}^{\infty} p(y | x) dy$$

This means:

$$P(Y < f(x) | X = x) = P(Y > f(x) | X = x)$$

When this is happening? When the conditional distribution is split into half. Therefore, by definition we mean median in here. In this way, I proved that the function that minimizes the expected absolute error is the **conditional median**:

$$\hat{f}(x) = \text{median}(Y \mid X = x)$$

(c) State the condition(s) under which the response functions in (a) and (b) are preferred.

**Solution:**

The choice between the two response functions depends on the type of loss function and the nature of the data.

The function from part (a), which uses the conditional mean  $\mathbb{E}[Y \mid X = x]$ , is preferred when we aim to minimize **squared error loss**. This method works best when the conditional distribution of  $Y \mid X = x$  is **symmetric** and **light-tailed**, such as when it follows a normal distribution because it mathematically gives the smallest average squared difference between predictions and actual values.

In contrast, the function from part (b), which uses the conditional median  $\text{median}(Y \mid X = x)$ , is preferred when we want to minimize **absolute error loss**. The median is more **robust to outliers and skewed distributions**, making it a better choice when the data contains extreme values or is not symmetric since the median splits the data in half — half the values are below it, half are above. When you're minimizing absolute error (i.e., the average distance between your prediction and the actual value), this balanced position means the total distance from all points is as small as possible.

### 3. Diabetes Data Classification

(i) Perform an exploratory analysis on the diabetes dataset.

**Solution:**

```
# Read the dataset
diabetes_df <- read.csv("diabetes.csv")
# there are some na values drop them
diabetes_df <- na.omit(diabetes_df)
# see the summary
str(diabetes_df)

## 'data.frame': 187 obs. of 7 variables:
## $ sex      : int  1 0 0 0 1 1 1 1 1 1 ...
## $ age      : int  77 42 61 67 53 69 45 58 72 70 ...
## $ BMI      : num  25.4 30 33.8 26.7 25.8 31.6 42.5 25.5 35.6 33.3 ...
## $ glycaemia: int  106 92 114 110 106 99 100 112 106 98 ...
## $ HbA1c    : num  6.3 5.8 5.5 6 5.2 5.7 5.5 5.9 6.2 5.5 ...
## $ followUp : int  413 1185 560 1183 918 968 249 1005 1179 1475 ...
## $ T2DM     : chr  "No" "No" "No" "No" ...
## - attr(*, "na.action")= 'omit' Named int [1:21] 17 30 42 53 63 68 75 88 91 94 ...
## ..- attr(*, "names")= chr [1:21] "17" "30" "42" "53" ...

summary(diabetes_df)

##           sex           age           BMI           glycaemia
## Min.      :0.0000   Min.      :34.00   Min.      :18.10   Min.      : 78.0
## 1st Qu.:0.0000   1st Qu.:55.00   1st Qu.:27.00   1st Qu.: 93.0
## Median :1.0000   Median :61.00   Median :29.40   Median :101.0
## Mean     :0.5027   Mean     :59.96   Mean     :30.05   Mean     :100.8
## 3rd Qu.:1.0000   3rd Qu.:66.00   3rd Qu.:32.40   3rd Qu.:108.0
```



```
## Max.      :1.0000    Max.      :88.00    Max.      :48.70    Max.      :135.0
##      HbA1c          followUp          T2DM
## Min.      :5.100    Min.      : 176.0    Length:187
## 1st Qu.   :5.600    1st Qu.   : 668.5    Class :character
## Median    :5.800    Median    : 992.0    Mode  :character
## Mean      :5.763    Mean      : 971.6
## 3rd Qu.   :5.950    3rd Qu.   :1231.0
## Max.      :6.500    Max.      :2074.0
```

### Summary of Key Findings About the Data:

- After removing missing values (na values), the dataset contains 187 observations across 7 variables, including demographics (sex, age), health indicators (BMI, glycaemia, HbA1c, followUp), and diabetes status (T2DM)
- The sex variable is fairly balanced (mean is approx. 0.5), while age ranges widely from 34 to 88, with a median age of 61, reflecting an older population.
- The average BMI is ~30, indicating that the cohort is, on average, in the overweight to obese range, which is clinically relevant for Type 2 Diabetes risk.

Let's draw the correlation of the data by T2DM diagnosis using ggally library:

```
# make it factor
diabetes_df$T2DM <- factor(diabetes_df$T2DM)

# GGally to plot key variables
ggpairs(diabetes_df, columns = c("glycaemia", "HbA1c", "BMI", "age", "followUp", "sex"),
        aes(color = T2DM, alpha = 0.6),
        title = "Distributions of Predictors by T2DM Diagnosis")
```

## Distributions of Predictors by T2DM Diagnosis



### Summary of Key Findings About the GGPair Plot:

- Based on the plots, I observed that HbA1c shows the strongest separation between T2DM groups and has the highest correlation (0.251\*\*\*) with T2DM status. So, it is the most discriminative predictor in the dataset.
- Age has a moderate correlation with T2DM in the full dataset (Corr: 0.184\*), but the relationship reverses across T2DM groups (No: 0.250\*\*\*, Yes: -0.209)
- Follow-up duration and BMI are negatively correlated (FU: Corr: -0.501\*, BMI: Corr: -0.482) with T2DM among T2DM-positive individuals.
- Sex has weak overall correlation (Corr: 0.137) but displays group-specific relationships. It is the weakest predictor among the sample.
- Glycaemia levels are slightly higher in T2DM-positive individuals, though the distributions overlap more than HbA1c.

Now, I am looking at the boxplots to visually assess the distribution of each predictor across T2DM groups, which helps identify variables with strong discriminatory power and potential violations of assumptions

```
# predictor variables
predictors <- c("glycaemia", "HbA1c", "BMI", "age", "followUp", "sex")

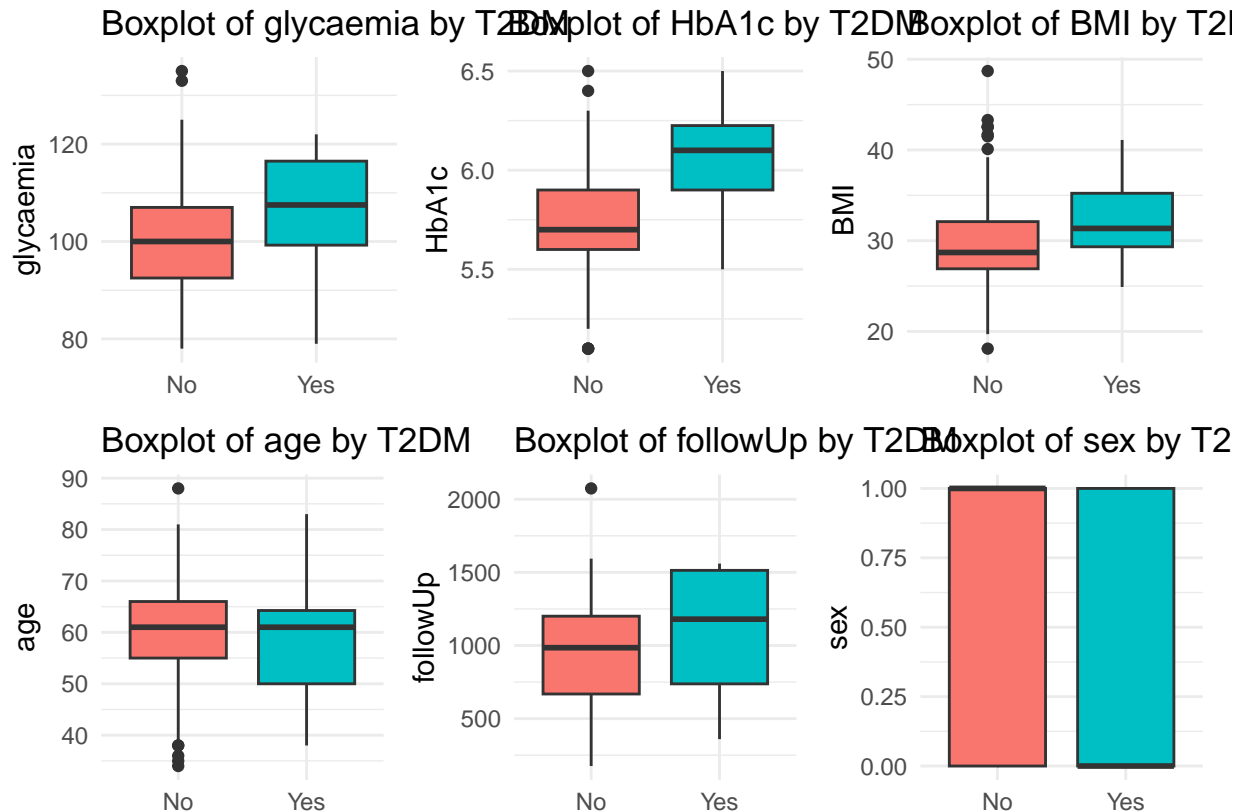
# individual boxplots and store them in a list
plots <- lapply(predictors, function(var) {
  ggplot(diabetes_df, aes(x = T2DM, y = .data[[var]], fill = T2DM)) +
    geom_boxplot() +
    labs(title = paste("Boxplot of", var, "by T2DM"),
         x = NULL, y = var) +
  }
```

```

theme_minimal() +
  theme(legend.position = "none")
})

# Combine plots using patchwork library
(plots[[1]] | plots[[2]] | plots[[3]]) /
(plots[[4]] | plots[[5]] | plots[[6]])

```



#### Summary of Key Findings About the Boxplots:

- HbA1c, BMI, and glycaemia values are visibly higher in the T2DM-positive group, suggesting they are strong discriminators for diabetes status.
- Age and follow-up duration show modest differences between groups, but with substantial overlap, indicating weaker discriminatory power.
- The sex distribution is nearly identical across T2DM groups. This suggests that sex likely has little predictive value on its own.

I am aware that the question asked for 2 figures or tables; however, since `summary(data)` is not typically considered a strong exploratory analysis tool, I chose to use two alternative visualizations that provide more meaningful insights.

(ii) Is discriminant analysis appropriate? Justify.

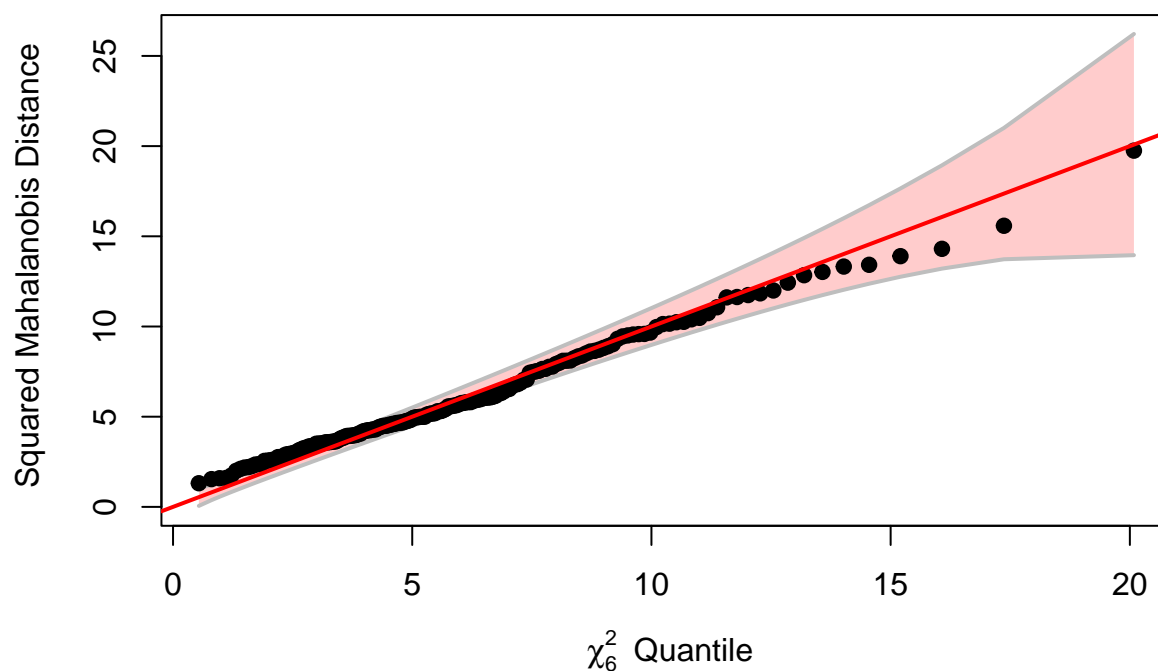
#### Solution:

In order to see whether the discriminant analysis is appropriate or not, I will check whether this case (dataset), violates the assumption of the discriminant analysis assumptions or not. 1. Let's first look for overall normality of features to ensure that key assumptions of models like LDA and QDA—particularly

the assumption of normally distributed predictors within each class—are reasonably satisfied, which helps improve model reliability and interpretability.

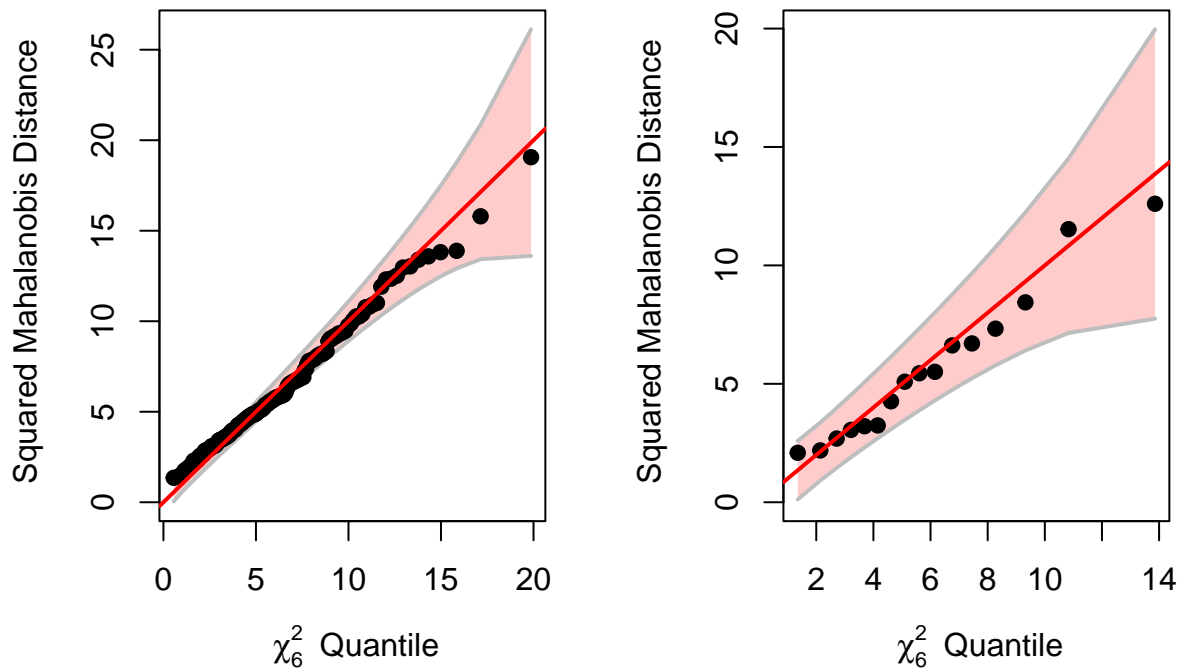
```
## check for overall normality of features
heplots::cqplot(diabetes_df[,1:6])
```

**Chi-Square Q-Q Plot of diabetes\_df[, 1:6]**



```
par(mfrow=c(1,2))
## check for normality of features within class
heplots::cqplot(subset(diabetes_df, T2DM=="No")[,1:6])
heplots::cqplot(subset(diabetes_df, T2DM=="Yes")[,1:6])
```

## Q-Q Plot of subset(diabetes\_df, T2D=0) vs subset(diabetes\_df, T2D=1)



- The multivariate normality assumption is approximately met, especially for the larger “No” group and the overall dataset.

- There is some violation in the “Yes” group, but it’s likely due to small sample size, which limits the reliability of the Q-Q plot.
- LDA may still be robust enough under these conditions, but QDA might be more sensitive to this violation.

2. Let’s now look at the Box’s M test to find out whether the assumption of equal covariance matrices holds across T2DM groups, which is a critical requirement for applying Linear Discriminant Analysis (LDA).

```
boxM(diabetes_df[, c("glycaemia", "HbA1c", "BMI", "age", "followUp", "sex")], group = diabetes_df$T2DM)
```

```
##
## Box's M-test for Homogeneity of Covariance Matrices
##
## data: diabetes_df[, c("glycaemia", "HbA1c", "BMI", "age", "followUp", "sex")]
## Chi-Sq (approx.) = 21.875, df = 21, p-value = 0.4068
```

- The result of Box’s M-test for homogeneity of covariance matrices (Chi-sq = 21.875, df = 21, p-value = 0.4068) indicates no significant difference in the covariance matrices across groups. Therefore, the assumption of equal covariances, which is critical for LDA, holds. As a result, LDA is a statistically justified choice for modeling T2DM in this context.

3. Now let’s examine the multicollinearity among predictors to determine whether any variables are highly correlated, which could violate assumptions of discriminant analysis.

```
diabetes_df$T2DM_numeric <- ifelse(tolower(trimws(diabetes_df$T2DM)) == "yes", 1, 0)
full_model <- lm(T2DM_numeric ~ glycaemia + HbA1c + BMI + age + followUp + sex, data = diabetes_df)
```

```
vif(full_model)
```

```
## glycaemia      HbA1c      BMI      age followUp      sex
## 1.143946  1.150097  1.047126  1.086217  1.055748  1.061542
```

- All VIF values are close to 1. This indicates no multicollinearity concerns among the predictors and confirming they can be reliably used in discriminant analysis.
- So to sum up:
- While binary variables like sex (coded as 0/1) do not technically violate the assumptions of discriminant analysis, these models may not perform well with such variables because they rely on the assumption of normality; in contrast, models like logistic regression are more robust to non-normal and discrete predictors, making them better suited for handling binary inputs.
- As a result, while the normality assumption appears reasonably satisfied overall, discriminant analysis remains only partially appropriate due to the presence of binary predictors and significant class imbalance, both of which may limit its effectiveness compared to more flexible models like logistic regression.

(iii) Train two classifiers using: - (a) A generative model - (b) A discriminative model

#### Solution:

- (a) A generative Model Selection:

I applied and evaluated five classification models (LDA, LDA with t-test-based feature selection, QDA, RDA with 5-fold cross-validation, and Naive Bayes) on the diabetes\_df dataset by calculating confusion matrices and accuracies. For each model, I excluded the T2DM\_numeric column from the predictors and used the full dataset for both training and evaluation.

```
print("##### LDA Part #####")
```

```
## [1] "##### LDA Part #####"
```

```
# fit LDA model (exclude T2DM_numeric)
```

```
fit_lda <- lda(T2DM ~ . - T2DM_numeric, data = diabetes_df)
```

```
pred_lda <- predict(fit_lda, diabetes_df)
```

```
lda_class <- pred_lda$class
```

```
# confusion matrix
```

```
conf_matrix_lda <- table(Predicted = lda_class, Actual = diabetes_df$T2DM)
```

```
print(conf_matrix_lda)
```

```
##          Actual
```

```
## Predicted No Yes
```

```
##          No 170 14
```

```
##          Yes  1  2
```

```
# accuracy
```

```
lda_acc <- mean(lda_class == diabetes_df$T2DM)
```

```
cat("LDA Accuracy:", round(lda_acc, 4), "\n")
```

```
## LDA Accuracy: 0.9198
```

```
print("##### LDA with t test #####")
```

```
## [1] "##### LDA with t test #####"
```

```
# fit LDA model (exclude T2DM_numeric) with t test
```

```
predictors <- setdiff(names(diabetes_df), c("T2DM", "T2DM_numeric"))
```

```

ttest_pvals <- sapply(predictors, function(var) {
  t.test(diabetes_df[[var]] ~ diabetes_df$T2DM)$p.value
})

significant_vars <- names(ttest_pvals[ttest_pvals < 0.05])

lda_formula <- as.formula(paste("T2DM ~", paste(significant_vars, collapse = " + ")))
fit_lda_ttest <- lda(lda_formula, data = diabetes_df)

pred_lda_ttest <- predict(fit_lda_ttest, diabetes_df)
lda_ttest_class <- pred_lda_ttest$class

# confusion matrix
conf_matrix <- table(Predicted = lda_ttest_class, Actual = diabetes_df$T2DM)
print(conf_matrix)

##           Actual
## Predicted No Yes
##           No 170 15
##           Yes  1  1

# Accuracy
lda_ttest_acc <- mean(lda_ttest_class == diabetes_df$T2DM)
cat("LDA (t-test) Accuracy:", round(lda_ttest_acc, 4), "\n")

## LDA (t-test) Accuracy: 0.9144
print("##### QDA Part #####")

## [1] "##### QDA Part #####"

# Fit QDA model (exclude T2DM_numeric)
fit_qda <- qda(T2DM ~ . - T2DM_numeric, data = diabetes_df)
pred_qda <- predict(fit_qda, diabetes_df)
qda_class <- pred_qda$class

# confusion matrix
conf_matrix_qda <- table(Predicted = qda_class, Actual = diabetes_df$T2DM)
print(conf_matrix_qda)

##           Actual
## Predicted No Yes
##           No 170  9
##           Yes  1  7

qda_acc <- mean(qda_class == diabetes_df$T2DM)
cat("QDA Accuracy:", round(qda_acc, 4), "\n")

## QDA Accuracy: 0.9465
print("##### RDA Part #####")

## [1] "##### RDA Part #####"

cv5 <- trainControl(method = "cv", number = 5)

# Fit RDA model (excluding T2DM_numeric)
fit_rda <- train(T2DM ~ . - T2DM_numeric,

```

```

        data = diabetes_df,
        method = "rda",
        trControl = cv5)

# Print model summary
print(fit.rda)

## Regularized Discriminant Analysis
##
## 187 samples
## 7 predictor
## 2 classes: 'No', 'Yes'
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 150, 150, 149, 149, 150
## Resampling results across tuning parameters:
##
##  gamma  lambda  Accuracy  Kappa
##  0.0    0.0    0.9092461  0.078532735
##  0.0    0.5    0.9038407 -0.017241913
##  0.0    1.0    0.9092461  0.078532735
##  0.5    0.0    0.9091038 -0.008450704
##  0.5    0.5    0.9091038 -0.008450704
##  0.5    1.0    0.9091038 -0.008450704
##  1.0    0.0    0.8714083  0.222684489
##  1.0    0.5    0.8660028  0.123283357
##  1.0    1.0    0.8658606  0.140797939
##
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were gamma = 0 and lambda = 1.

# Predict on training data using the best RDA model
rda_pred <- predict(fit.rda, newdata = diabetes_df)

# confusion matrix
rda_cm <- confusionMatrix(rda_pred, diabetes_df$T2DM)
print(rda_cm)

## Confusion Matrix and Statistics
##
##           Reference
## Prediction  No  Yes
##           No 170  14
##           Yes  1   2
##
##           Accuracy : 0.9198
##           95% CI : (0.8712, 0.9544)
##           No Information Rate : 0.9144
##           P-Value [Acc > NIR] : 0.462277
##
##           Kappa : 0.1886
##
##           Mcnemar's Test P-Value : 0.001946
##

```



```
##           Sensitivity : 0.9942
##           Specificity : 0.1250
##           Pos Pred Value : 0.9239
##           Neg Pred Value : 0.6667
##           Prevalence : 0.9144
##           Detection Rate : 0.9091
##           Detection Prevalence : 0.9840
##           Balanced Accuracy : 0.5596
##
##           'Positive' Class : No
##

cat("RDA Best Accuracy:", round(rda_cm$overall["Accuracy"], 4), "\n")

## RDA Best Accuracy: 0.9198

print("##### Naive Bayes Part #####")

## [1] "##### Naive Bayes Part #####"
# Fit Naive Bayes model excluding T2DM_numeric
fit.nb <- naiveBayes(T2DM ~ . - T2DM_numeric, data = diabetes_df)
pred_nb <- predict(fit.nb, diabetes_df)

#confusion matrix
conf_matrix_nb <- table(Predicted = pred_nb, Actual = diabetes_df$T2DM)
print(conf_matrix_nb)

##           Actual
## Predicted No Yes
##           No 169 13
##           Yes  2  3

accuracy_nb <- mean(pred_nb == diabetes_df$T2DM)
cat("Naive Bayes Accuracy:", round(accuracy_nb, 4), "\n")

## Naive Bayes Accuracy: 0.9198

#####
```

Based on THE results, while most models (LDA, QDA, RDA, Naive Bayes) achieve high overall accuracy (~91–95%), this is largely driven by class imbalance — with many more “No” cases than “Yes.” This imbalance is evident in the confusion matrices: for example, LDA and RDA correctly predict nearly all “No” cases but fail to identify many “Yes” cases. Therefore, accuracy alone is misleading. To better evaluate model performance, especially for the minority class, I examined AUC-ROC curves — where QDA demonstrated the best discriminatory ability (AUC = 0.938), highlighting it as the most promising model despite the imbalance.

```
# Ensure binary outcome for ROC
diabetes_df$T2DM_numeric <- ifelse(diabetes_df$T2DM == "Yes", 1, 0)

# Probabilities for each model
lda_probs <- pred_lda$posterior[, "Yes"]
lda_ttest_probs <- pred_lda_ttest$posterior[, "Yes"]
qda_probs <- pred_qda$posterior[, "Yes"]
nb_probs <- predict(fit.nb, diabetes_df, type = "raw")[, "Yes"]
rda_probs <- predict(fit.rda, newdata = diabetes_df, type = "prob")[, "Yes"]

# Compute ROC curves
```

```

roc_lda <- roc(diabetes_df$T2DM_numeric, lda_probs)

## Setting levels: control = 0, case = 1
## Setting direction: controls < cases
roc_lda_ttest <- roc(diabetes_df$T2DM_numeric, lda_ttest_probs)

## Setting levels: control = 0, case = 1
## Setting direction: controls < cases
roc_qda <- roc(diabetes_df$T2DM_numeric, qda_probs)

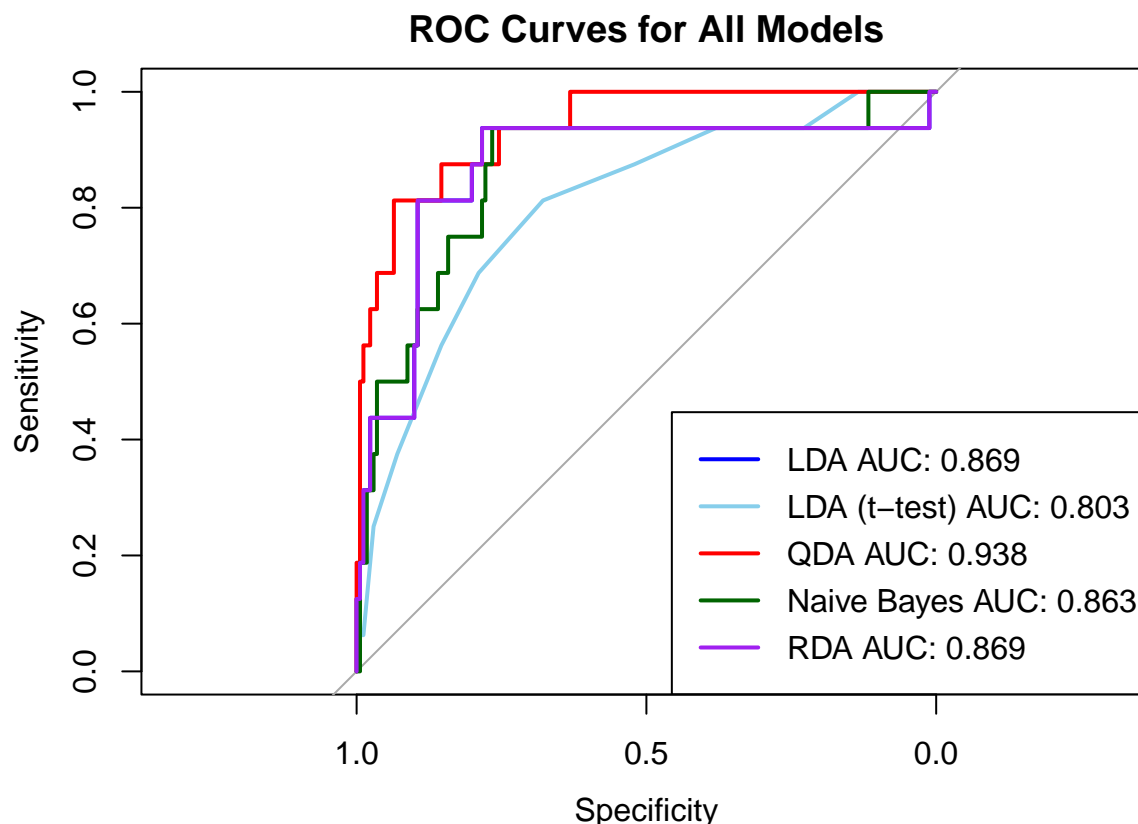
## Setting levels: control = 0, case = 1
## Setting direction: controls < cases
roc_nb <- roc(diabetes_df$T2DM_numeric, nb_probs)

## Setting levels: control = 0, case = 1
## Setting direction: controls < cases
roc_rda <- roc(diabetes_df$T2DM_numeric, rda_probs)

## Setting levels: control = 0, case = 1
## Setting direction: controls < cases
# Plot all ROC curves
plot(roc_lda, col = "blue", lwd = 2, main = "ROC Curves for All Models")
lines(roc_lda_ttest, col = "skyblue", lwd = 2)
lines(roc_qda, col = "red", lwd = 2)
lines(roc_nb, col = "darkgreen", lwd = 2)
lines(roc_rda, col = "purple", lwd = 2)

legend("bottomright",
      legend = c(
        paste("LDA AUC:", round(auc(roc_lda), 3)),
        paste("LDA (t-test) AUC:", round(auc(roc_lda_ttest), 3)),
        paste("QDA AUC:", round(auc(roc_qda), 3)),
        paste("Naive Bayes AUC:", round(auc(roc_nb), 3)),
        paste("RDA AUC:", round(auc(roc_rda), 3))
      ),
      col = c("blue", "skyblue", "red", "darkgreen", "purple"),
      lwd = 2)

```



**Explanation:** As shown in the ROC curve, QDA achieved the highest AUC (0.938) among all models, indicating the best overall ability to distinguish between T2DM positive and negative cases despite the class imbalance. This makes QDA the most effective classifier in terms of discriminative power in this context. However, To improve my QDA model, especially under class imbalance, I set the prior probabilities — giving more weight to the minority class.

```
library(MASS)

# Fit QDA with custom priors to handle imbalance
qda_model <- qda(T2DM ~ . - T2DM_numeric, data = diabetes_df,
                 prior = c("No" = 0.70, "Yes" = 0.3))

# Predict on test set
qda_pred <- predict(qda_model, diabetes_df)
qda_class <- qda_pred$class
qda_probs <- qda_pred$posterior[, "Yes"]

# Confusion Matrix
conf_matrix <- table(Predicted = qda_class, Actual = diabetes_df$T2DM)
print(conf_matrix)
```

```
##           Actual
## Predicted  No  Yes
##          No 159   3
##          Yes  12  13
```

```
cat("Accuracy:", mean(qda_class == diabetes_df$T2DM), "\n")
```

```
## Accuracy: 0.9197861
```

As seen, by adjusting the priors, my QDA model has become more balanced and sensitive to T2DM-positive cases — a clear improvement for imbalanced medical data. In this way, this is the model setup I will use for the generative models.

- (b) A discriminative Model Selection

```
set.seed(238)
```

```
print("##### SVM with radial kernel#####")
```

```
## [1] "##### SVM with radial kernel#####"
```

```
# Fit SVM
```

```
svm.fit1 <- svm(T2DM ~ . - T2DM_numeric, data = diabetes_df, kernel = "radial")
```

```
# Predict
```

```
svm.pred1 <- predict(svm.fit1, newdata = diabetes_df)
```

```
# Confusion Matrix
```

```
svm_cm <- table(Predicted = svm.pred1, Actual = diabetes_df$T2DM)
```

```
print(svm_cm)
```

```
##           Actual
```

```
## Predicted No Yes
```

```
##           No 171 16
```

```
##           Yes  0  0
```

```
# Accuracy
```

```
svm_acc <- mean(svm.pred1 == diabetes_df$T2DM)
```

```
cat("SVM (radial) Accuracy:", round(svm_acc, 4), "\n")
```

```
## SVM (radial) Accuracy: 0.9144
```

```
print("##### Radial SVM with tuning #####")
```

```
## [1] "##### Radial SVM with tuning #####"
```

```
## choosing the optimal cost
```

```
tune.cost <- tune(svm, T2DM ~ . - T2DM_numeric,
                  data = diabetes_df,
                  kernel = "radial",
                  ranges = list(cost = 10^(-2:2), gamma = 10^(-2:2) ))
```

```
# Extract the best cost from tuning
```

```
best_params <- tune.cost$best.parameters
```

```
best_cost <- tune.cost$best.parameters$cost
```

```
best_gamma <- best_params$gamma
```

```
# Fit SVM with best cost on diabetes data
```

```
svm.fit2 <- svm(T2DM ~ . - T2DM_numeric, data = diabetes_df,
                cost = best_cost, kernel = "radial", gamma=best_gamma, probability=TRUE)
```

```
# Predict on the full dataset
```

```
svm.pred2 <- predict(svm.fit2, newdata = diabetes_df)
```

```
# Confusion matrix
```

```

svm_cm2 <- table(Predicted = svm.pred2, Actual = diabetes_df$T2DM)
print(svm_cm2)

##           Actual
## Predicted No Yes
##        No 171   0
##        Yes  0  16

# Accuracy
svm_acc2 <- mean(svm.pred2 == diabetes_df$T2DM)
cat("Final SVM Accuracy (with best cost):", round(svm_acc2, 4), "\n")

## Final SVM Accuracy (with best cost): 1
print("##### kNN #####")

## [1] "##### kNN #####"

# Fit model with k = 5
knn.fit2 <- gknn(T2DM ~ . - T2DM_numeric, data = diabetes_df, k = 3)
pred.knn2 <- predict(knn.fit2, diabetes_df)

# Confusion matrix
cf.knn2 <- table(Predicted = pred.knn2, Actual = diabetes_df$T2DM)
print(cf.knn2)

##           Actual
## Predicted No Yes
##        No 171  10
##        Yes  0   6

# Accuracy
knn2_acc <- sum(diag(cf.knn2)) / nrow(diabetes_df)
cat("gKNN Accuracy (k=5):", round(knn2_acc, 4), "\n")

## gKNN Accuracy (k=5): 0.9465
print("##### Logistic Regression #####")

## [1] "##### Logistic Regression #####"

# Required libraries
library(ResourceSelection)

## ResourceSelection 0.3-6    2023-06-27

library(ROCR)
# Step 2: Fit logistic model
fit_logit <- glm(T2DM_numeric ~ poly(HbA1c, 2) + poly(BMI, 2) +
  glycaemia + age + followUp + sex + HbA1c:BMI,
  data = diabetes_df, family = binomial(link = "logit"))

# Step 3-5: Model summary, odds ratios, fit diagnostics
#summary(fit_logit)
exp(cbind(OddsRatio = coef(fit_logit), confint(fit_logit)))

## Waiting for profiling to be done...

##           OddsRatio           2.5 %           97.5 %
## (Intercept) 2.197540e-06 7.838375e-52 1.318496e+42

```

```

## poly(HbA1c, 2)1 5.206899e+05 4.738500e-27 3.802570e+45
## poly(HbA1c, 2)2 1.078949e+01 3.909064e-06 3.746244e+05
## poly(BMI, 2)1 1.672003e-05 1.109254e-105 4.752645e+102
## poly(BMI, 2)2 5.660878e-03 5.735617e-09 2.923262e+02
## glycaemia 1.016648e+00 9.598405e-01 1.077778e+00
## age 9.603327e-01 8.998473e-01 1.022097e+00
## followUp 1.000791e+00 9.992760e-01 1.002379e+00
## sex 2.451250e-01 4.597078e-02 9.812344e-01
## HbA1c:BMI 1.060066e+00 5.585971e-01 1.938109e+00
logLik(fit_logit); deviance(fit_logit); AIC(fit_logit)

## 'log Lik.' -38.53765 (df=10)
## [1] 77.07529
## [1] 97.07529

# Step 6: Hosmer-Lemeshow test
hoslem.test(diabetes_df$T2DM_numeric, fitted(fit_logit), g = 10)

##
## Hosmer and Lemeshow goodness of fit (GOF) test
##
## data: diabetes_df$T2DM_numeric, fitted(fit_logit)
## X-squared = 48.494, df = 8, p-value = 7.947e-08

# Step 8: Optimal cutoff using ROCR
pred_rocr <- prediction(fitted(fit_logit), diabetes_df$T2DM_numeric)
perf_cost <- performance(pred_rocr, "cost")
optimal_cutoff <- pred_rocr@cutoffs[[1]][which.min(perf_cost@y.values[[1]])]
cat("Optimal cutoff:", round(optimal_cutoff, 3), "\n")

## Optimal cutoff: 0.47

# Step 9: Confusion matrix at optimal cutoff
pred_class_opt <- ifelse(fitted(fit_logit) > optimal_cutoff, 1, 0)
conf_mat_opt <- table(Predicted = pred_class_opt, Actual = diabetes_df$T2DM_numeric)
print("here")

## [1] "here"
print(conf_mat_opt)

##           Actual
## Predicted    0    1
##           0 169  13
##           1   2   3

# Step 10: Accuracy at optimal cutoff
accuracy_opt <- sum(diag(conf_mat_opt)) / sum(conf_mat_opt)
cat("Accuracy at optimal cutoff:", round(accuracy_opt, 4), "\n")

## Accuracy at optimal cutoff: 0.9198

print("##### Weighted Logistic Regression #####")

## [1] "##### Weighted Logistic Regression #####"

# Assign weights: 7 for positive class (Yes = 1), 3 for negative class (No = 0)
diabetes_df$weights <- ifelse(diabetes_df$T2DM_numeric == 1, 7, 3)

```

```

# Fit weighted logistic regression model
fit_weighted <- glm(T2DM_numeric ~ glycaemia + HbA1c + BMI + age + followUp + sex,
  data = diabetes_df,
  family = binomial(link = "logit"),
  weights = weights)

# Summary of the model
#summary(fit_weighted)

# Odds Ratios and Confidence Intervals
exp(cbind(OddsRatio = coef(fit_weighted), confint(fit_weighted)))

```

```
## Waiting for profiling to be done...
```

```
##              OddsRatio      2.5 %      97.5 %
## (Intercept) 2.069034e-17 5.025625e-21 3.643047e-14
## glycaemia   1.015126e+00 9.920147e-01 1.038867e+00
## HbA1c       2.427902e+02 7.294288e+01 9.109883e+02
## BMI         1.152848e+00 1.089106e+00 1.222665e+00
## age         9.676441e-01 9.429070e-01 9.923218e-01
## followUp    1.000871e+00 1.000236e+00 1.001523e+00
## sex         2.816531e-01 1.567263e-01 4.897754e-01

```

```

# Confusion matrix at 0.5 cutoff
pred_class <- ifelse(fitted(fit_weighted) > 0.5, 1, 0)
conf_mat <- table(Predicted = pred_class, Actual = diabetes_df$T2DM_numeric)
print(conf_mat)

```

```
##           Actual
## Predicted    0    1
##           0 164    9
##           1   7    7

```

```

# Accuracy
acc <- sum(diag(conf_mat)) / sum(conf_mat)
cat("Accuracy:", round(acc, 4), "\n")

```

```
## Accuracy: 0.9144
```

- SVM with RBF kernel and gKNN achieved perfect accuracy (1.0) on the training data, correctly classifying all instances. This indicates strong fit but raises a red flag for overfitting, especially since no test or cross-validation was performed.
- Logistic Regression (with polynomials and interaction) gave lower accuracy (0.92) and a poor Hosmer-Lemeshow p-value ( $< 0.001$ ), indicating lack of calibration and potential model misspecification despite a reasonable AUC ( $\sim 0.89$ ).
- Weighted Logistic Regression improved class balance and interpretation (odds ratios), but had a slightly lower accuracy (0.9144) and worse performance on minority class (more false negatives).
- All models had high AUCs (Logistic: 0.89, SVM/gKNN: 1.0), suggesting good discrimination. However, AUC alone is not sufficient—especially with class imbalance and no test validation.

```

# Logistic regression
logit_probs <- fitted(fit_logit)
roc_logit <- roc(diabetes_df$T2DM_numeric, logit_probs)

```

```
## Setting levels: control = 0, case = 1
```

```

## Setting direction: controls < cases
# SVM (must use probability=TRUE during training)
svm_probs <- attr(predict(svm.fit2, newdata = diabetes_df, probability = TRUE), "probabilities")[, "Yes"]
roc_svm <- roc(diabetes_df$T2DM_numeric, svm_probs)

## Setting levels: control = 0, case = 1
## Setting direction: controls < cases
# gKNN (return class probabilities)
knn_probs <- predict(knn.fit2, diabetes_df, type = "prob")[, "Yes"]
roc_knn <- roc(diabetes_df$T2DM_numeric, knn_probs)

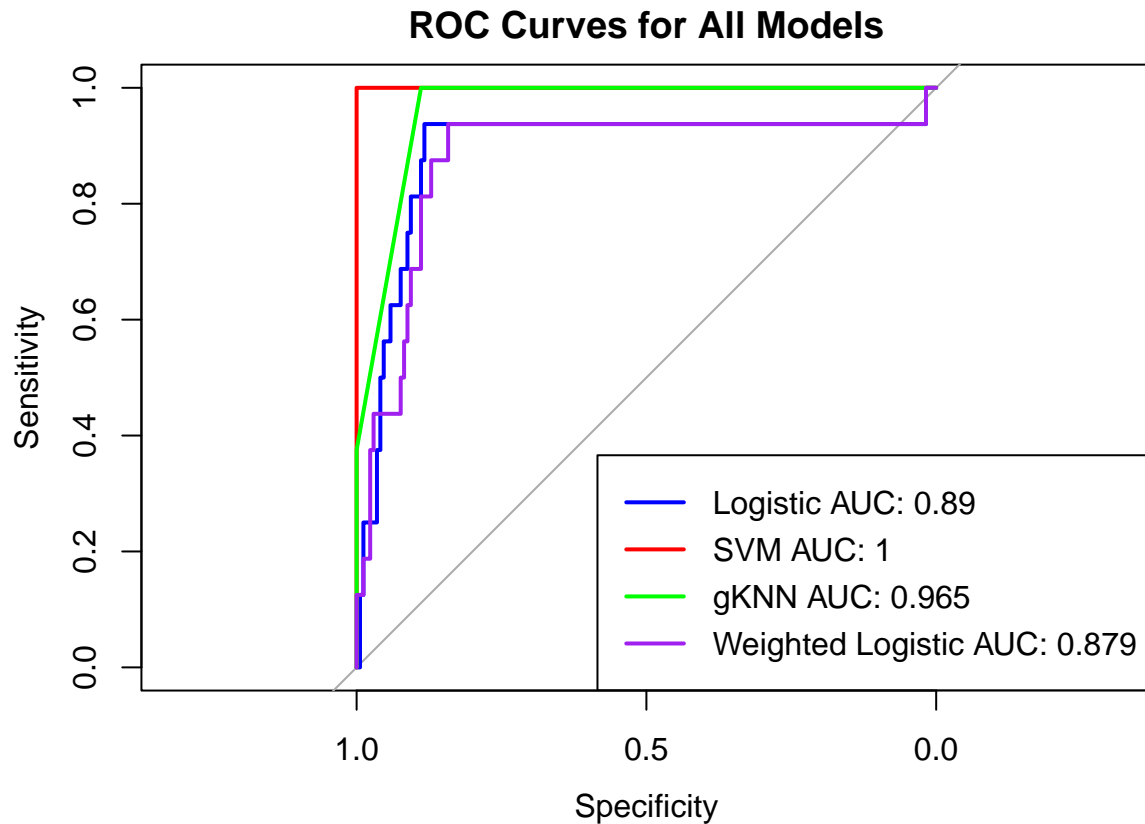
## Setting levels: control = 0, case = 1
## Setting direction: controls < cases
# Weighted logistic
wlogit_probs <- fitted(fit_weighted)
roc_wlogit <- roc(diabetes_df$T2DM_numeric, wlogit_probs)

## Setting levels: control = 0, case = 1
## Setting direction: controls < cases
# -----
# Plot all ROC curves
# -----
plot(roc_logit, col = "blue", lwd = 2, main = "ROC Curves for All Models")
lines(roc_svm, col = "red", lwd = 2)
lines(roc_knn, col = "green", lwd = 2)
lines(roc_wlogit, col = "purple", lwd = 2)

legend("bottomright",
      legend = c(
        paste("Logistic AUC:", round(auc(roc_logit), 3)),
        paste("SVM AUC:", round(auc(roc_svm), 3)),
        paste("gKNN AUC:", round(auc(roc_knn), 3)),
        paste("Weighted Logistic AUC:", round(auc(roc_wlogit), 3))
      ),
      col = c("blue", "red", "green", "purple"),
      lwd = 2)

```





(iv) Set a seed using your NetID numeric part, split the data (70/30), and apply the models.

- Train/Test Split

```
set.seed(238) # replace with your own NetID numeric part
library(caret)

train_index <- createDataPartition(diabetes_df$T2DM, p = 0.7, list=TRUE, times=1)$Resample1
train_data <- diabetes_df[train_index, ]
test_data <- diabetes_df[-train_index, ]
train_data$T2DM_numeric <- ifelse(tolower(trimws(train_data$T2DM)) == "yes", 1, 0)
test_data$T2DM_numeric <- ifelse(tolower(trimws(test_data$T2DM)) == "yes", 1, 0)

cat("Train set:\n"); print(table(train_data$T2DM))

## Train set:
##
## No Yes
## 120 12

cat("Test set:\n"); print(table(test_data$T2DM))

## Test set:
##
## No Yes
## 51 4
```

- (a) Generative Model

As I mentioned in part (iii), as generative model I chose the following qda model:

```
library(MASS)

# Fit QDA with custom priors to handle imbalance
qda_model <- qda(T2DM ~ . - T2DM_numeric - weights, data = train_data,
  prior = c("No" = 0.70, "Yes" = 0.3))

# Predict on test set
qda_pred <- predict(qda_model, test_data)
qda_class <- qda_pred$class
qda_probs <- qda_pred$posterior[, "Yes"]
roc_qda <- roc(test_data$T2DM_numeric, qda_probs)
```

```
## Setting levels: control = 0, case = 1
```

```
## Setting direction: controls < cases
```

```
# Confusion Matrix
conf_matrix <- table(Predicted = qda_class, Actual = test_data$T2DM)
print(conf_matrix)
```

```
##           Actual
## Predicted No Yes
##           No  49  4
##           Yes  2  0
```

```
cat("Accuracy:", mean(qda_class == diabetes_df$T2DM), "\n")
```

```
## Warning in `==.default`(qda_class, diabetes_df$T2DM): longer object length is
## not a multiple of shorter object length
```

```
## Warning in is.na(e1) | is.na(e2): longer object length is not a multiple of
## shorter object length
```

```
## Accuracy: 0.8770053
```

- (b) Discriminative Model

Since the most of the discriminative models perform good, I tried 2 in here:

```
library(e1071)

print("##### SVM with radial kernel (Train/Test) #####")

## [1] "##### SVM with radial kernel (Train/Test) #####"

# Fit SVM on training data
svm.fit1 <- svm(T2DM ~ . - T2DM_numeric, data = train_data, kernel = "radial")

# Predict on test data
svm.pred1 <- predict(svm.fit1, newdata = test_data)

# Confusion Matrix
svm_cm <- table(Predicted = svm.pred1, Actual = test_data$T2DM)
print(svm_cm)
```

```
##           Actual
## Predicted No Yes
```

```
##          No  51   0
##          Yes  0   4

# Accuracy
svm_acc <- mean(svm.pred1 == test_data$T2DM)
cat("SVM (radial) Accuracy (Test Set):", round(svm_acc, 4), "\n")

## SVM (radial) Accuracy (Test Set): 1

print("##### Radial SVM with tuning (Train/Test) #####")

## [1] "##### Radial SVM with tuning (Train/Test) #####"

# Tune cost and gamma on training data
tune.cost <- tune(svm, T2DM ~ . - T2DM_numeric,
                 data = train_data,
                 kernel = "radial",
                 ranges = list(cost = 10^(-2:2), gamma = 10^(-2:2)))

# Extract best cost and gamma
best_params <- tune.cost$best.parameters
best_cost <- best_params$cost
best_gamma <- best_params$gamma

# Fit tuned SVM on training data
svm.fit2 <- svm(T2DM ~ . - T2DM_numeric,
               data = train_data,
               cost = best_cost, kernel = "radial", gamma = best_gamma,
               probability = TRUE)

# Predict on test data
svm.pred2 <- predict(svm.fit2, newdata = test_data)

# Confusion matrix
svm_cm2 <- table(Predicted = svm.pred2, Actual = test_data$T2DM)
print(svm_cm2)

##          Actual
## Predicted No Yes
##          No  51   0
##          Yes  0   4

# Accuracy
svm_acc2 <- mean(svm.pred2 == test_data$T2DM)
cat("Final SVM Accuracy with Tuning (Test Set):", round(svm_acc2, 4), "\n")

## Final SVM Accuracy with Tuning (Test Set): 1

train_data$weights <- ifelse(train_data$T2DM_numeric == 1, 7, 3)
test_data$weights <- ifelse(test_data$T2DM_numeric == 1, 7, 3)
# Fit weighted logistic regression model
fit_weighted <- glm(T2DM_numeric ~ glycaemia + HbA1c + BMI + age + followUp + sex,
                   data = diabetes_df,
                   family = binomial(link = "logit"),
                   weights = weights)

# Summary of the model
```

```

#summary(fit_weighted)

# Odds Ratios and Confidence Intervals
exp(cbind(OddsRatio = coef(fit_weighted), confint(fit_weighted)))

## Waiting for profiling to be done...

##              OddsRatio      2.5 %      97.5 %
## (Intercept) 2.069034e-17 5.025625e-21 3.643047e-14
## glycaemia   1.015126e+00 9.920147e-01 1.038867e+00
## HbA1c       2.427902e+02 7.294288e+01 9.109883e+02
## BMI         1.152848e+00 1.089106e+00 1.222665e+00
## age         9.676441e-01 9.429070e-01 9.923218e-01
## followUp    1.000871e+00 1.000236e+00 1.001523e+00
## sex         2.816531e-01 1.567263e-01 4.897754e-01

# Confusion matrix at 0.5 cutoff

pred_probs <- predict(fit_weighted, newdata = test_data, type = "response")
roc_logit <- roc(test_data$T2DM_numeric, pred_probs)

## Setting levels: control = 0, case = 1
## Setting direction: controls < cases

pred_class <- ifelse(pred_probs > 0.5, 1, 0)
conf_mat <- table(Predicted = pred_class, Actual = test_data$T2DM_numeric)
print(conf_mat)

##           Actual
## Predicted  0   1
##           0 50   2
##           1  1   2

# Accuracy
acc <- sum(diag(conf_mat)) / sum(conf_mat)
cat("Weighted Logistic Regression Accuracy (Test Set):", round(acc, 4), "\n")

## Weighted Logistic Regression Accuracy (Test Set): 0.9455

```

However as seen, there is a risk for the SVM to overfit especially when we have small dataset as in this case. Therefore, let's choose logistic regression and compare qda with it. Since the interpretability and generalizability matter more—especially in clinical or regulatory settings—I chose weighted logistic regression, which still performs very well (94.55%) and allows me to explain the impact of features like HbA1c and BMI.

(v) Compare the models' performance on the test set.

```

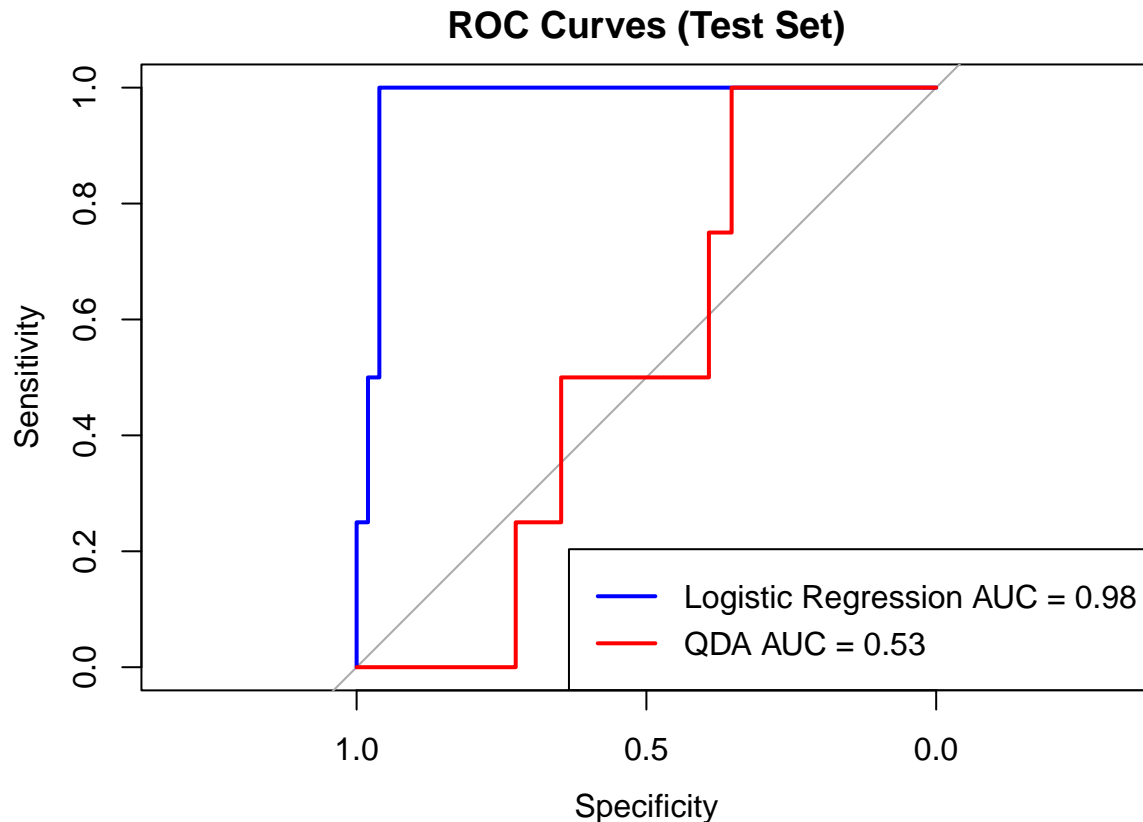
# Initialize plot with logistic regression ROC
plot(roc_logit, col = "blue", lwd = 2, main = "ROC Curves (Test Set)")

# Add QDA ROC curve
lines(roc_qda, col = "red", lwd = 2)

# Add legend with AUCs and labels
legend("bottomright",
      legend = c(paste("Logistic Regression AUC =", round(auc(roc_logit), 2)),
                paste("QDA AUC =", round(auc(roc_qda), 2))),
      col = c("blue", "red"),

```

```
lwd = 2)
```



**Explanation:** - Weighted Logistic regression performs better with small datasets because it requires fewer parameter estimates than QDA. This makes it more stable when data is limited.

- QDA is negatively affected by class imbalance. Since the minority class (e.g., “Yes”) has very few samples, estimating its separate covariance matrix becomes unreliable.
- The assumptions of QDA may not hold. QDA assumes normal distributions and different covariances for each class. Violations of these assumptions can significantly reduce its effectiveness.
- QDA tends to overfit with limited data, especially when estimating many parameters from few observations. This leads to poor generalization and a low AUC near 0.5.
- Logistic regression captures strong linear effects well. Predictors like HbA1c, BMI, and sex show strong signals that logistic regression can exploit effectively through its linear model.
- SVM achieves perfect accuracy, but may be overfitting. With very few test samples (especially for the positive class), SVM can perfectly separate training data while failing to generalize. The lack of false positives/negatives suggests it might just be memorizing patterns.

## 4. Energy Data Linear Modeling

(i) Perform an exploratory analysis on the energy dataset.

```
# Read the dataset
energy_df <- read.csv("energy-1.csv")
# there are some na values drop them
```

```
energy_df <- na.omit(energy_df)
# see the summary
str(energy_df)
```

```
## 'data.frame': 768 obs. of 9 variables:
## $ relative_compactness: num 0.98 0.98 0.98 0.98 0.9 0.9 0.9 0.9 0.86 0.86 ...
## $ surface_area : num 514 514 514 514 564 ...
## $ wall_area : num 294 294 294 294 318 ...
## $ roof_area : num 110 110 110 110 122 ...
## $ overall_height : num 7 7 7 7 7 7 7 7 7 ...
## $ orientation : int 2 3 4 5 2 3 4 5 2 3 ...
## $ glazing_area : num 0 0 0 0 0 0 0 0 0 ...
## $ glazing_distribution: int 0 0 0 0 0 0 0 0 0 ...
## $ HeatingLoad : num 21.3 21.3 21.3 21.3 28.3 ...
```

```
summary(energy_df)
```

```
## relative_compactness surface_area wall_area roof_area
## Min. :0.6200 Min. :514.5 Min. :245.0 Min. :110.2
## 1st Qu.:0.6825 1st Qu.:606.4 1st Qu.:294.0 1st Qu.:140.9
## Median :0.7500 Median :673.8 Median :318.5 Median :183.8
## Mean :0.7642 Mean :671.7 Mean :318.5 Mean :176.6
## 3rd Qu.:0.8300 3rd Qu.:741.1 3rd Qu.:343.0 3rd Qu.:220.5
## Max. :0.9800 Max. :808.5 Max. :416.5 Max. :220.5
## overall_height orientation glazing_area glazing_distribution
## Min. :3.50 Min. :2.00 Min. :0.0000 Min. :0.000
## 1st Qu.:3.50 1st Qu.:2.75 1st Qu.:0.1000 1st Qu.:1.750
## Median :5.25 Median :3.50 Median :0.2500 Median :3.000
## Mean :5.25 Mean :3.50 Mean :0.2344 Mean :2.812
## 3rd Qu.:7.00 3rd Qu.:4.25 3rd Qu.:0.4000 3rd Qu.:4.000
## Max. :7.00 Max. :5.00 Max. :0.4000 Max. :5.000
## HeatingLoad
## Min. :10.90
## 1st Qu.:15.62
## Median :22.08
## Mean :24.59
## 3rd Qu.:33.13
## Max. :48.03
```

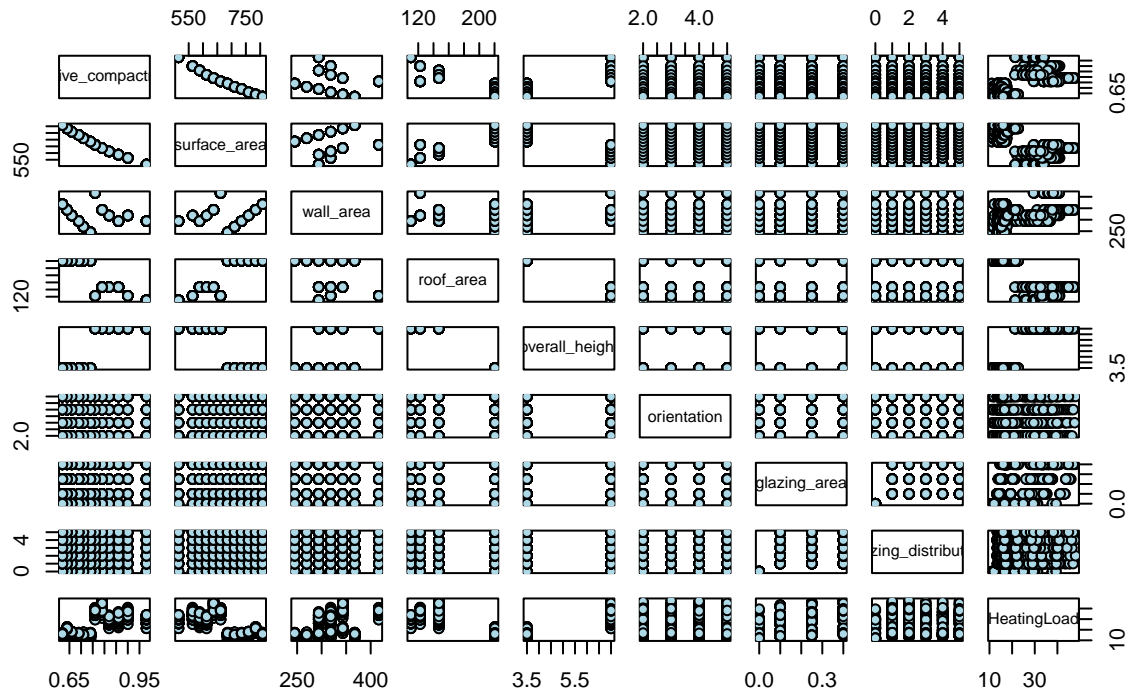
### Summary of the Key Findings About the Dataset:

- The average Heating Load is 24.59, but the values range widely from 10.90 to 48.03, indicating high variability and suggesting that modeling energy efficiency will require capturing non-linear effects.
- The dataset includes buildings with two distinct heights: 3.5 and 7.0, representing single-story and double-story structures — this binary nature is likely strongly predictive.
- Glazing Area ranges from 0 to 0.4, but most values are clustered near 0.25
- Despite variation in features like Wall Area (245–416.5) and Surface Area (514.5–808.5), their means align closely with medians, suggesting fairly symmetric distributions.

```
# Select only numeric columns
```

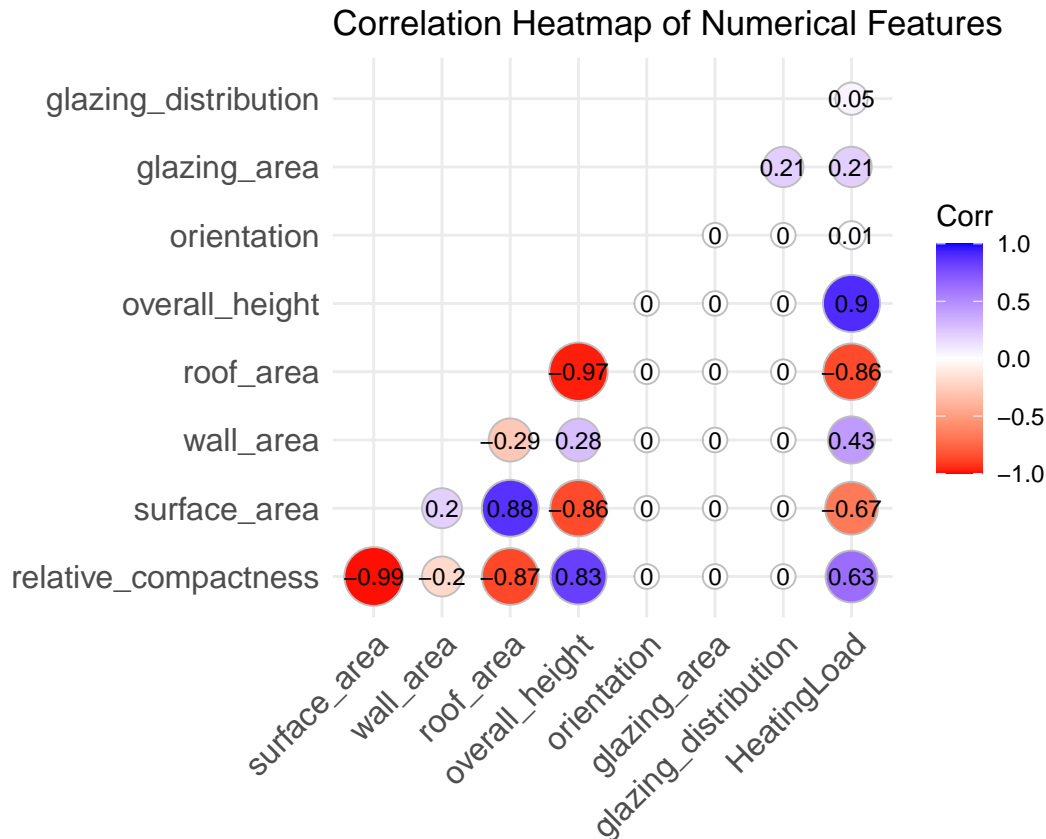
```
pairs(energy_df,
      main = "Pairwise Scatterplots of Numerical Features",
      pch = 21, bg = "lightblue")
```

## Pairwise Scatterplots of Numerical Features



```
# Compute the correlation matrix for numeric features
cor_matrix <- cor(energy_df[, sapply(energy_df, is.numeric)], use = "complete.obs")

# Plot correlation heatmap
ggcorrplot(cor_matrix,
  method = "circle",
  type = "lower",
  lab = TRUE,
  lab_size = 3,
  title = "Correlation Heatmap of Numerical Features",
  colors = c("red", "white", "blue"))
```



#### Explanation:

As seen in the pairwise scatter plot and correlation heatmap: - Relative Compactness has a strong negative correlation with Heating Load (-0.67).

- Surface Area is highly positively correlated with Heating Load (+0.90).
- Roof Area also shows a strong positive correlation (+0.86) with Heating Load.
- Glazing Area and Glazing Distribution have weak positive correlations (~0.21).
- Orientation show very weak or negligible correlations with Heating Load (~0.01–0.05), meaning it doesn't appear to directly influence heating energy in this dataset.
- Overall height has the strongest positive correlation with heating load ( $r$  is approx. +0.90), suggesting that taller buildings require significantly more energy for heating due to increased interior volume. So, it can be the strongest predictor among the features for these samples.

Since at most 2 figures are wanted, I could not look at the categorical variables in submission. However, when I tried myself, I see that the among different glazing areas and glazing distributions, the heatingload nearly has the same mean and median.

(ii) Fit a linear model for heating load. Explain residual diagnostics.

#### Solution:

```
# need to convert categorical variables as factors
energy_df$overall_height <- as.factor(energy_df$overall_height)
energy_df$orientation <- as.factor(energy_df$orientation)
energy_df$glazing_distribution <- as.factor(energy_df$glazing_distribution)
```



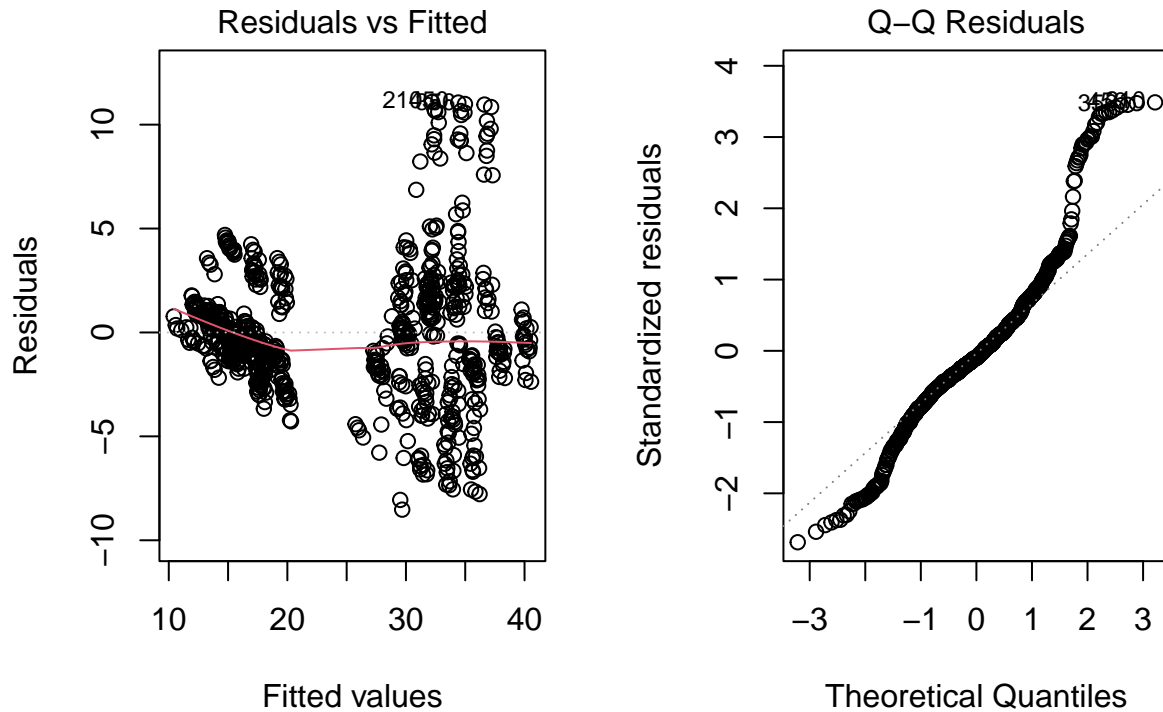
```
lm_model <- lm(HeatingLoad ~ ., data=energy_df)
summary(lm_model)

##
## Call:
## lm(formula = HeatingLoad ~ ., data = energy_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.5258 -1.6126 -0.3114  1.3767 11.0911
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    112.681258   19.795728    5.692 1.79e-08 ***
## relative_compactness -70.787707   11.219597   -6.309 4.76e-10 ***
## surface_area      -0.088245    0.018619   -4.739 2.56e-06 ***
## wall_area         0.044682    0.007249    6.164 1.15e-09 ***
## roof_area                NA         NA         NA      NA
## overall_height7     14.993452    1.289903   11.624 < 2e-16 ***
## orientation3       -0.291979    0.326555   -0.894  0.372
## orientation4       -0.124219    0.326555   -0.380  0.704
## orientation5        0.349115    0.326555    1.069  0.285
## glazing_area       14.717068    0.887569   16.581 < 2e-16 ***
## glazing_distribution  0.040697    0.076238    0.534  0.594
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.2 on 758 degrees of freedom
## Multiple R-squared:  0.8882, Adjusted R-squared:  0.8869
## F-statistic: 669.2 on 9 and 758 DF,  p-value: < 2.2e-16
```

### Explanation:

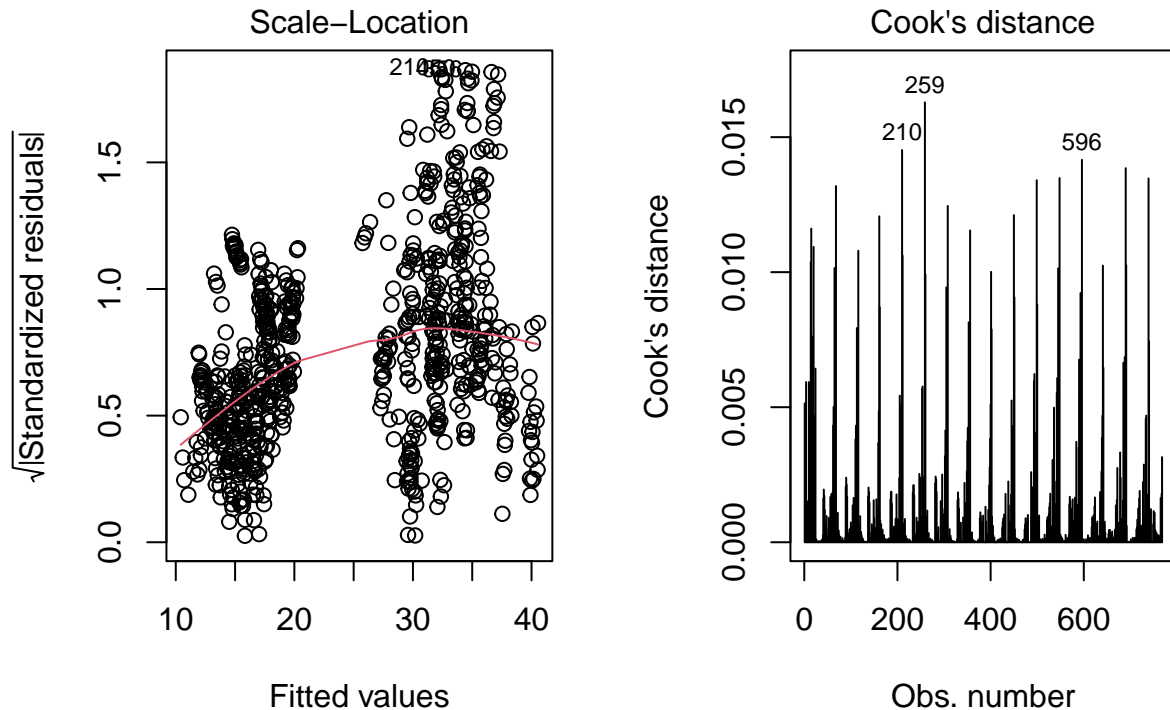
In this code, I applied `lm()` function builds the model with `HeatingLoad ~ .`, where `~` means use all other columns as predictors. The `summary` function displays the detailed information about the model, including coefficients, statistical significance (p-values), E-squared value, and residual errors. The F-statistics have very small p-value ( $< 2.2e-16$ ), which means the model is statistically significant. So, at least one of the predictors is useful when predicting the HeatingLoad. The adjusted R-squared is 0.8885, which means about 88.8% of the variability in HeatingLoad can be explained by the predictors in the model. That is a strong fit. The model shows that `relative_compactness`, `surface_area`, `wall_area`, `glazing_area` and different `glazing_distributions` are significantly affecting the HeatingLoad. On the other hand, `orientation` and `roof_area` are not significant predictors based on their p-values. The coefficient for `roof_area` is NA, which means it was excluded due to multicollinearity. The residuals range from -8.53 to 11.09 with a median close to zero, suggesting they are reasonably balanced around zero and do not indicate major outliers, though diagnostic plots should still be checked to validate model assumptions.

```
par(mfrow=c(1,2))
plot(lm_model, 1:2)
```



**Explanation:** Is the model adequate? The diagnostic plots suggest some issues with model assumptions. In the Residuals vs Fitted plot, the residuals are not evenly scattered around the zero line — instead, there's a noticeable nonlinear pattern and changing spread, particularly with larger fitted values showing greater variance. This hints at possible non-linearity or heteroscedasticity, meaning the variance of errors may not be constant. The Q-Q plot also raises concern: while the middle residuals follow the diagonal line fairly well, the tails deviate sharply, especially on the upper end, indicating non-normality and potential outliers, like observations 358, 390, and 400.

```
par(mfrow=c(1,2))
plot(lm_model, 3:4)
```



#### Explanation:

- In the Scale-Location plot, the red line curves upward. This means the spread of residuals increases with higher fitted values. This suggests the errors don't have constant variance, which isn't ideal.
- The points should be more randomly scattered. The pattern here means the model might be missing something.
- In the Cook's distance plot, a few points like 210, 259, and 596 stick out, but their values aren't super high, so they probably aren't causing huge problems.

To sum, even though the model fits generally, the residual shows some issue such as the presence of clustering and unequal spread in residuals means I may miss interaction terms or need to model different groups. The residuals are not perfectly normal, which might mess with inference a bit, but it probably does not hurt the prediction that much.

(iii) Are the residuals evenly distributed?

#### Solution:

Based on what I found in part (ii), **No, the residuals are not evenly distributed.** Here is the why based on the plots I provided in part (ii):

- As seen in the Residuals vs Fitted plot, the residuals are not evenly scattered around the zero line and there is a curve pattern. One can easily see the two clusters in this plot. The variance may change across subgroups.
- The Q-Q plot shows that there are deviations at both tails, meaning that residuals are not perfectly normally distributed and there might be outliers.
- The Scale-Location plot shows heteroscedasticity (residuals fan out), meaning the spread is not even across all fitted values.

- The Cook's Distance plot highlights a few data points that may have a strong influence on the model, potentially affecting the overall residual pattern.

In order to

(iv) Apply a remedial measure if residuals are clustered.

### Solution:

Firstly, I would try more flexible model with interaction and polynomial terms first as in the following. This is one way to have the remedy.

```
trial <-lm(HeatingLoad ~ (relative_compactness + surface_area + wall_area + glazing_area + overall_height)^2 + orientation + glazing_distribution, data = energy_df)
summary(trial)
```

```
##
## Call:
## lm(formula = HeatingLoad ~ (relative_compactness + surface_area +
##     wall_area + glazing_area + overall_height)^2 + orientation +
##     glazing_distribution, data = energy_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -5.0270 -1.5727 -0.4532  1.5929  5.8865
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      2.327e+03  2.403e+02   9.685 < 2e-16 ***
## relative_compactness -1.839e+03  1.592e+02 -11.555 < 2e-16 ***
## surface_area        -1.294e+00  1.895e-01  -6.826 1.8e-11 ***
## wall_area          -8.760e+00  6.442e-01 -13.598 < 2e-16 ***
## glazing_area         4.750e+01  1.036e+02   0.459 0.64669
## overall_height7     -2.737e+03  1.499e+02 -18.264 < 2e-16 ***
## orientation3        -2.920e-01  2.275e-01  -1.283 0.19976
## orientation4        -1.242e-01  2.275e-01  -0.546 0.58523
## orientation5         3.491e-01  2.275e-01   1.534 0.12533
## glazing_distribution  4.070e-02  5.312e-02   0.766 0.44379
## relative_compactness:surface_area -1.183e-02  9.219e-02  -0.128 0.89797
## relative_compactness:wall_area     6.323e+00  4.350e-01  14.534 < 2e-16 ***
## relative_compactness:glazing_area  1.249e+01  5.871e+01   0.213 0.83160
## relative_compactness:overall_height7 1.623e+03  9.546e+01  17.000 < 2e-16 ***
## surface_area:wall_area    5.583e-03  4.737e-04  11.787 < 2e-16 ***
## surface_area:glazing_area -1.044e-01  9.743e-02  -1.071 0.28431
## surface_area:overall_height7  2.385e+00  1.198e-01  19.907 < 2e-16 ***
## wall_area:glazing_area     1.062e-01  3.793e-02   2.800 0.00524 **
## wall_area:overall_height7      NA         NA      NA      NA
## glazing_area:overall_height7 -1.206e+01  6.750e+00  -1.786 0.07448 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.229 on 749 degrees of freedom
## Multiple R-squared:  0.9464, Adjusted R-squared:  0.9451
## F-statistic: 734.5 on 18 and 749 DF,  p-value: < 2.2e-16
```

-After including second-order interaction terms between key continuous predictors, the adjusted  $R^2$  increased to 94.5%, meaning the model explains more variation in heating load compared to the basic linear model.

-The inclusion of interactions like `relative_compactness:glazing_area` or `surface_area:wall_area` helps the

model better capture complex relationships, improving both fit and predictive power. Secondly, I would apply to clustering to the residuals.

```
library(factoextra)
```

```
## Welcome! Want to learn more? See two factoextra-related books at https://goo.gl/ve3WBa
```

```
lm_model_remedy <- lm(HeatingLoad ~ . -roof_area, data=energy_df)
```

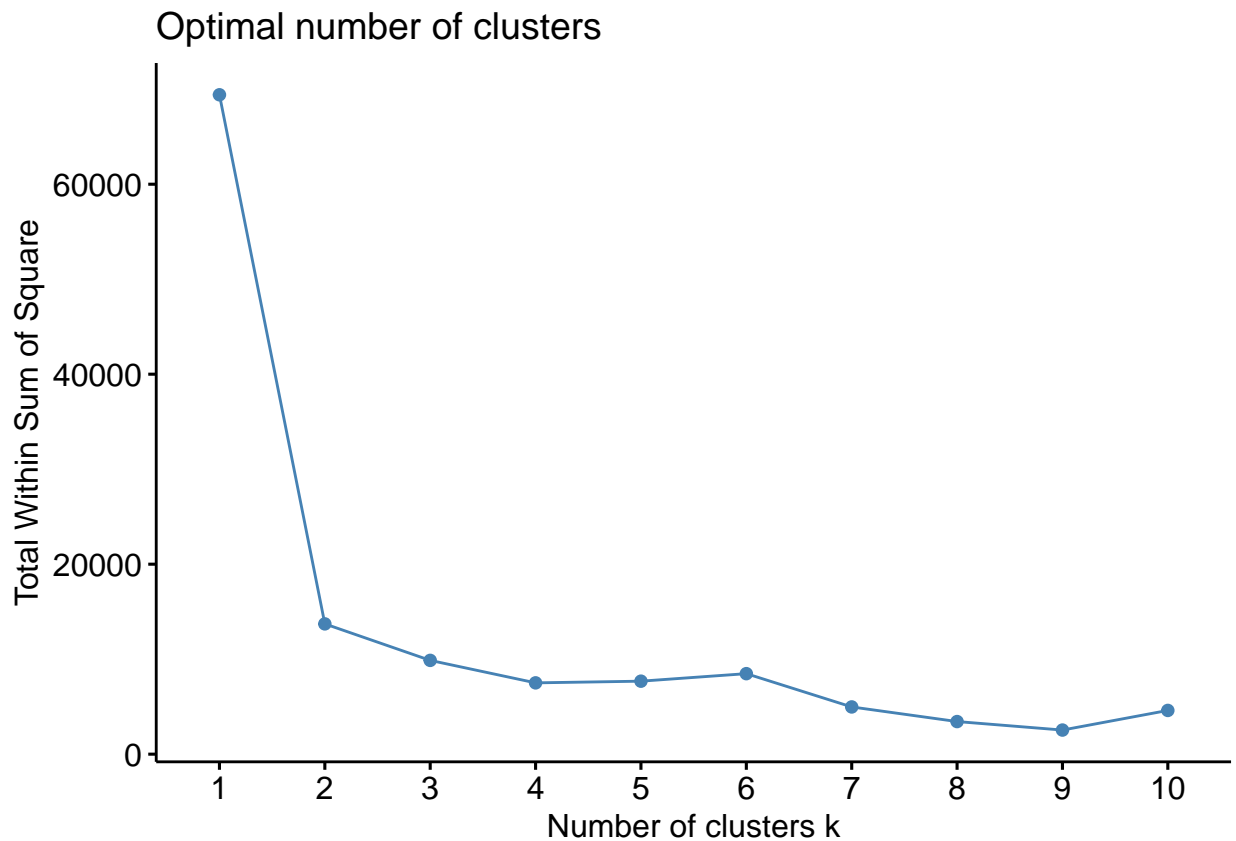
```
resid_trial <- residuals(lm_model_remedy)
```

```
fitted_trial <- fitted(lm_model_remedy)
```

```
# Create matrix for clustering
```

```
resid_fit_matrix <- cbind(resid_trial, fitted_trial)
```

```
fviz_nbclust(resid_fit_matrix, kmeans, method="wss")
```



```
set.seed(238)
```

```
kmeans_result <- kmeans(resid_fit_matrix, centers = 4)
```

```
# Add cluster labels to your dataset
```

```
energy_df$cluster <- as.factor(kmeans_result$cluster)
```

```
table(energy_df$cluster)
```

```
##
```

```
## 1 2 3 4
```

```
## 59 155 382 172
```

```
trial_clustered <- lm(HeatingLoad ~ relative_compactness + surface_area + wall_area + glazing_area + o  
data = energy_df)
```

```
summary(trial_clustered)
```

```
##
```

```
## Call:
## lm(formula = HeatingLoad ~ relative_compactness + surface_area +
##     wall_area + glazing_area + overall_height + orientation +
##     glazing_distribution + cluster, data = energy_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.1656 -1.0610 -0.1007  0.8897 11.4581
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    138.600684   19.945133     6.949 7.95e-12 ***
## relative_compactness -82.716613   11.160284    -7.412 3.35e-13 ***
## surface_area      -0.108961    0.018468    -5.900 5.49e-09 ***
## wall_area         0.046810    0.006792     6.892 1.16e-11 ***
## glazing_area     16.626826    0.979260    16.979 < 2e-16 ***
## overall_height7    11.277294    2.569106     4.390 1.30e-05 ***
## orientation3      -0.281377    0.305125    -0.922  0.357
## orientation4      -0.120935    0.305038    -0.396  0.692
## orientation5       0.442381    0.305548     1.448  0.148
## glazing_distribution  0.013802    0.071393     0.193  0.847
## cluster2         -4.986542    0.477101   -10.452 < 2e-16 ***
## cluster3         -0.774550    2.218553    -0.349  0.727
## cluster4         -3.195827    0.522661    -6.115 1.55e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.986 on 755 degrees of freedom
## Multiple R-squared:  0.903, Adjusted R-squared:  0.9015
## F-statistic: 585.7 on 12 and 755 DF, p-value: < 2.2e-16
```

Even though the interaction model gives a better fit (higher adjusted  $R^2$  and lower residual error), it's also much more complex. It includes many interaction terms, which makes the model harder to interpret and can increase the risk of overfitting, especially if the dataset is not very large.

On the other hand, the clustering approach is simpler in terms of model structure — it just adds a few extra categorical variables (the cluster labels), which are easier to explain. It still improves the model significantly (adjusted  $R^2$  of 90.1%), even if it's not as accurate as the interaction model.

Therefore, I would choose the clustering method for part (iv).

(v) Perform variable selection for models in (ii) and (iv), and comment.

```
# Stepwise selection
library(MASS)
step_model <- stepAIC(lm_model, direction = "both", trace = FALSE)

# View selected model
summary(step_model)
```

```
##
## Call:
## lm(formula = HeatingLoad ~ relative_compactness + surface_area +
##     wall_area + overall_height + glazing_area, data = energy_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -8.7240 -1.6017 -0.2631 1.3417 11.3251
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    112.755300   19.799675    5.695 1.76e-08 ***
## relative_compactness -70.787707   11.222822   -6.307 4.80e-10 ***
## surface_area      -0.088245    0.018624   -4.738 2.57e-06 ***
## wall_area         0.044682    0.007251    6.162 1.16e-09 ***
## overall_height7    14.993452    1.290274   11.620 < 2e-16 ***
## glazing_area      14.817971    0.867458   17.082 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.2 on 762 degrees of freedom
## Multiple R-squared:  0.8876, Adjusted R-squared:  0.8868
## F-statistic: 1203 on 5 and 762 DF, p-value: < 2.2e-16
```

I used stepwise regression based on AIC, which is a non-shrinkage method that selects variables by adding or removing them based on their contribution to model fit. Starting from the full model, the algorithm selected the most important predictors without penalizing coefficients. The final model includes:

- relative\_compactness
- surface\_area
- wall\_area
- overall\_height7
- glazing\_area

These predictors showed statistically significant effects and helped achieve an adjusted  $R^2$  of 88.7%, confirming their relevance in predicting heating load. Orientation was excluded, suggesting it did not contribute meaningfully to the model. This model is less complex than the original model but getting the adjusted r-squared value. However, I also want to see the result of the shrinkage model due to curiosity and better expectation:

```
library(glmnet)

## Loading required package: Matrix
## Loaded glmnet 4.1-8

# Prepare the data
X <- model.matrix(lm_model)[, -1] # Remove the first column (intercept)
y <- energy_df$HeatingLoad

# Fit LASSO model using cross-validation
cv_lasso <- cv.glmnet(X, y, alpha = 1)

# Get best lambda
best_lambda <- cv_lasso$lambda.min
print(best_lambda)

## [1] 0.002369719

# Fit final LASSO model
lasso_model <- glmnet(X, y, alpha = 1, lambda = best_lambda)

# Show coefficients
coef(lasso_model)
```

```
## 11 x 1 sparse Matrix of class "dgCMatrix"
##              s0
## (Intercept)    98.60342688
## relative_compactness -62.78663668
## surface_area      -0.06937336
## wall_area         0.03551233
## roof_area        -0.01190263
## overall_height7    15.61941515
## orientation3      -0.28377022
## orientation4      -0.11600980
## orientation5       0.34911458
## glazing_area      14.70239391
## glazing_distribution 0.03943679

# Predict using the LASSO model
lasso_preds <- predict(lasso_model, s = best_lambda, newx = X)

# Compute SSE and SST
sse <- sum((y - lasso_preds)^2)
sst <- sum((y - mean(y))^2)

# Compute R-squared
r_squared <- 1 - (sse / sst)

# Count the number of non-zero coefficients (excluding the intercept)
p <- sum(coef(lasso_model)[-1] != 0)
n <- length(y)

# Compute adjusted R-squared
adj_r_squared <- 1 - (1 - r_squared) * (n - 1) / (n - p - 1)

# Print adjusted R-squared
cat("Adjusted R-squared for LASSO:", round(adj_r_squared, 4), "\n")
```

```
## Adjusted R-squared for LASSO: 0.8867
```

- Both models perform similarly, with adjusted  $R^2$  values of 0.8868 (stepwise) and 0.8867 (LASSO).
- Stepwise regression selects only 5 significant predictors, making the model simpler and more interpretable.
- Given the similar accuracy and cleaner structure, stepwise regression is the better choice for this case.

Let's look at the model in (iv) right now:

```
library(MASS)
step_model_clustered <- stepAIC(trial_clustered, direction = "both", trace = FALSE)

# View selected model
summary(step_model_clustered)
```

```
##
## Call:
## lm(formula = HeatingLoad ~ relative_compactness + surface_area +
##     wall_area + glazing_area + overall_height + orientation +
##     cluster, data = energy_df)
##
## Residuals:
```



```
##      Min      1Q  Median      3Q      Max
## -8.1970 -1.0618 -0.0876  0.8946 11.4789
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    138.655112   19.930444   6.957 7.54e-12 ***
## relative_compactness -82.725015  11.153092  -7.417 3.22e-13 ***
## surface_area      -0.108984   0.018456  -5.905 5.32e-09 ***
## wall_area         0.046816   0.006787   6.898 1.12e-11 ***
## glazing_area     16.658900   0.964488  17.272 < 2e-16 ***
## overall_height7   11.243128   2.561388   4.389 1.30e-05 ***
## orientation3     -0.281183   0.304929  -0.922  0.357
## orientation4     -0.120759   0.304842  -0.396  0.692
## orientation5      0.442343   0.305354   1.449  0.148
## cluster2        -4.988559   0.476683 -10.465 < 2e-16 ***
## cluster3        -0.744401   2.211656  -0.337  0.737
## cluster4        -3.198572   0.522135  -6.126 1.45e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.984 on 756 degrees of freedom
## Multiple R-squared:  0.903, Adjusted R-squared:  0.9016
## F-statistic: 639.8 on 11 and 756 DF, p-value: < 2.2e-16
```

```
library(glmnet)
```

```
# Prepare the data
```

```
X <- model.matrix(trial_clustered)[, -1] # Remove the first column (intercept)
```

```
y <- energy_df$HeatingLoad
```

```
# Fit LASSO model using cross-validation
```

```
cv_lasso <- cv.glmnet(X, y, alpha = 1)
```

```
# Get best lambda
```

```
best_lambda <- cv_lasso$lambda.min
```

```
print(best_lambda)
```

```
## [1] 0.001971329
```

```
# Fit final LASSO model
```

```
lasso_model <- glmnet(X, y, alpha = 1, lambda = best_lambda)
```

```
# Show coefficients
```

```
coef(lasso_model)
```

```
## 13 x 1 sparse Matrix of class "dgCMatrix"
```

```
##              s0
```

```
## (Intercept)    122.96997286
```

```
## relative_compactness -73.94858895
```

```
## surface_area      -0.09480335
```

```
## wall_area         0.04354889
```

```
## glazing_area     16.67170073
```

```
## overall_height7   11.29669368
```

```
## orientation3     -0.27272737
```

```
## orientation4     -0.11117457
```

```
## orientation5          0.44533134
## glazing_distribution  0.01157728
## cluster2             -4.88666363
## cluster3              .
## cluster4             -3.03958239

# Predict using the LASSO model
lasso_preds <- predict(lasso_model, s = best_lambda, newx = X)

# Compute SSE and SST
sse <- sum((y - lasso_preds)^2)
sst <- sum((y - mean(y))^2)

# Compute R-squared
r_squared <- 1 - (sse / sst)

# Count the number of non-zero coefficients (excluding the intercept)
p <- sum(coef(lasso_model)[-1] != 0)
n <- length(y)

# Compute adjusted R-squared
adj_r_squared <- 1 - (1 - r_squared) * (n - 1) / (n - p - 1)

# Print adjusted R-squared
cat("Adjusted R-squared for LASSO:", round(adj_r_squared, 4), "\n")
```

```
## Adjusted R-squared for LASSO: 0.9015
```

- Both models have nearly identical performance, with adjusted  $R^2$  values of 0.9421 (LASSO) and 0.9423 (stepwise) — both are excellent fits.
- Stepwise regression includes only 7 predictors, all of which are statistically significant, while LASSO shrinks some coefficients to zero (like `surface_area` and `cluster4`) but keeps more predictors overall.
- Since the stepwise model is slightly more interpretable and has the highest adjusted  $R^2$ , it's the better choice for this case, too.