

732A99/732A68/TDDE01 Machine Learning Computer Lab 1

Erik Lundqvist (erilu777)
Axel Strid (axest556)
Oliver Solvang Stoltz (oliso325)

December 13, 2024

Statement of Contribution

Erik Lundqvist, Axel Strid and Oliver Solvang Stoltz started working on the beginning of task 1.1 and 1.2, Oliver Solvang Stoltz then completed assignment 1 and wrote the report about assignment 1.

Erik Lundqvist coded and wrote the part of the report for assignment 2.

Axel Strid coded and wrote the report regarding Assignment 3.

Oliver Solvang Stoltz did theory part for assignment 4 and wrote the report for this assignment .

1 Assignment 1: Handwritten Digit Recognition with K-nearest Neighbors

This assignment was about classifying handwritten numbers, 0-9 using the k-nearest neighbors method (KNN). The dataset contained handwritten numbers from 43 different people and a validation number of the actual number.

1.1 Data Preparation and Initial Model

The data was imported and divided into train-, validation- and test by 50%, 25%, 25% respectively. The seed 12345 was used to ensure reproducibility in the random splitting process when splitting the data. After the data was split, the KNN model was now fitted with a K value = 30 and a rectangular kernel. The model was fitted twice: first using the training data for both training and testing and second, using the training data for training and test data for testing. The resulting confusion matrices and misclassification errors were compared to analyze the differences.

1.2 Classification Results

The result of the confusion matrix for the training data can be seen in Table: 1 and the confusion matrix for the test data can be seen in Table: 2

Table 1: Confusion Matrix for Training Data

Truth \ Predicted	0	1	2	3	4	5	6	7	8	9
0	202	0	0	0	0	0	0	0	0	0
1	0	179	11	0	0	0	0	1	1	3
2	0	1	190	0	0	0	0	1	0	0
3	0	0	0	185	0	1	0	1	0	1
4	1	3	0	0	159	0	0	7	1	4
5	0	0	0	1	0	171	0	1	0	8
6	0	2	0	0	0	0	190	0	0	0
7	0	3	0	0	0	0	0	178	1	0
8	0	10	0	2	0	0	2	0	188	2
9	1	3	0	5	2	0	0	3	3	183

Table 2: Confusion Matrix for Test Data

Truth \ Predicted	0	1	2	3	4	5	6	7	8	9
0	77	0	0	0	1	0	0	0	0	0
1	0	81	2	0	0	0	0	0	0	3
2	0	0	98	0	0	0	0	0	3	0
3	0	0	0	107	0	2	0	0	1	1
4	0	0	0	0	94	0	2	6	2	5
5	0	1	1	0	0	93	2	1	0	5
6	0	0	0	0	0	0	90	0	0	0
7	0	0	0	1	0	0	0	111	0	0
8	0	7	0	1	0	0	0	0	70	0
9	0	1	1	1	0	0	0	1	0	85

The misclassification error for the training data was calculated to be $0.04500262 \approx 4.50\%$, and the misclassification error for the test data was calculated to be $0.05329154 \approx 5.33\%$. This implies that the overall accuracy of the model is about 95%. This level of accuracy might be acceptable for use cases where some error is tolerable, such as scanning documents, where occasional misclassification may not significantly impact the outcome. However, if the model is used for more critical tasks, such as grading exams, a higher accuracy would likely be required to minimize errors and ensure fairness. The model sometimes misclassifies 4 as 7, 5 as 9, and 8 as 1, this is likely since these numbers have some things in common. However, a high majority is still accurate which makes us believe that this model is accurate.

1.3 Analysis of Digit "8" Classification

The next step was to identify two cases that the model found easiest to classify as eight and three cases it struggled the most to classify as eight. After arranging the probabilities for the prediction of the number eight, heatmaps were created to visualize different cases. The easiest cases to classify as an eight were at indices 129 and 195, with index 129 being the easiest to classify. For the hardest cases to classify as an eight, the indices were 520, 431, and 1294, with index 520 being the hardest to classify. In figure: 1, 2, 3, 4 and 5 the different heat maps can be seen.

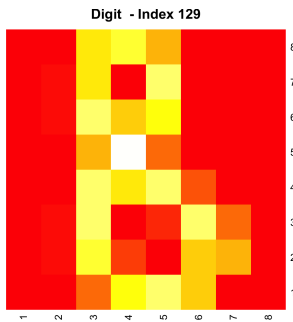


Figure 1: The easiest case for the model to predict as an eight.

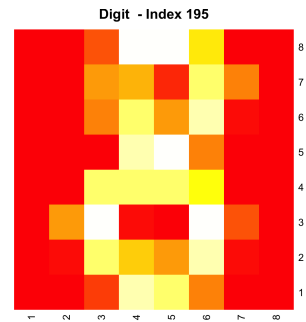


Figure 2: The second easiest case for the model to predict as an eight.

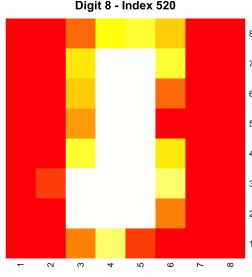


Figure 3: The hardest case for the model to predict as an eight.



Figure 4: The second hardest case for the model to predict as an eight.

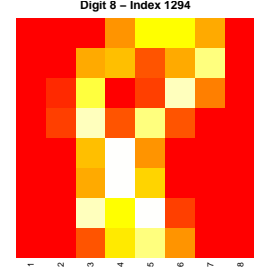


Figure 5: The third hardest case for the model to predict as an eight.

Visually, the easy cases can be identified as eights. The hard cases are indeed challenging to predict, and these numbers could be 8, 9, or even 1.

1.4 Different K-values and their effect on misclassification error

The next step was to calculate the misclassification error for different K-values for the training data and the validation data to see what value corresponded with the lowest misclassification error. This resulted in a graph with training error, validation error and an optimal k-value, this graph can be found in figure: 6.

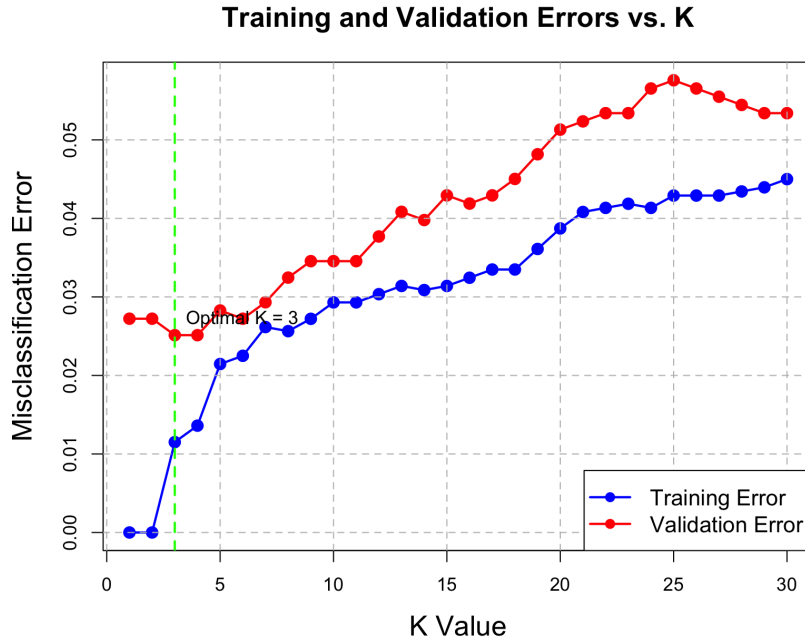


Figure 6: An graph for different K values and their corresponding misclassification error, for training- and validation data.

In Figure 6, we can see that the optimal K value is 3, with a misclassification error of approximately 2.5%, as the validation error is the lowest. With a low K -value, the model becomes more complex and tries to fit every inconsistently, which will lead to overfitting. As K increases, the model becomes less complex because it generalizes more, leading to

underfitting. With a K -value = 3 the test error was calculated to $0.03239289 \approx 3.24\%$. The test error shows that the model would be useful on unseen data with a high accuracy.

1.5 Different K-values and their effect on Cross-entropy error

The final task was to use cross-entropy error instead of misclassification error and then see which K -value gave the optimal solution. In figure 7 the result of this is shown.

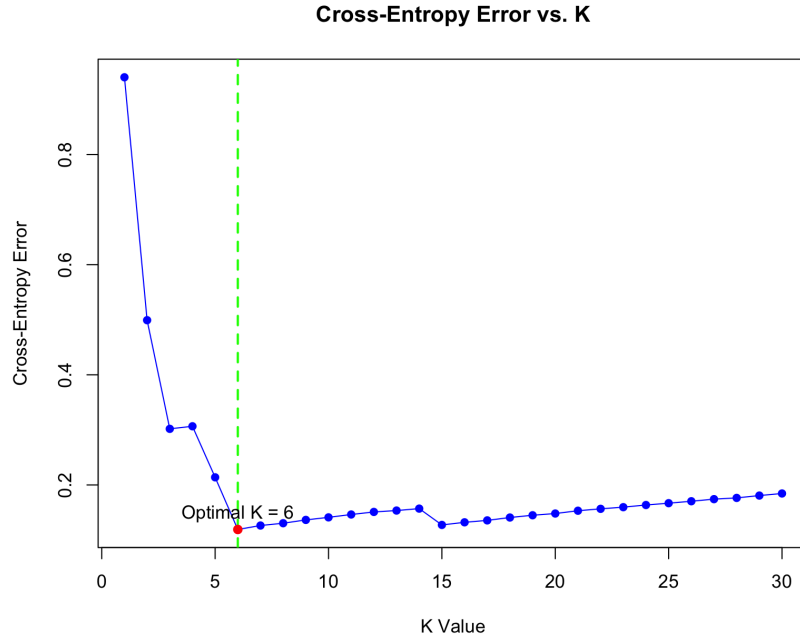


Figure 7: An graph for different K values and their corresponding cross-entropy error, for validation data.

We can see that the optimal solution was $k=6$ with a cross entropy error to be approximately 0.07. If you assume that the response has a multinomial distribution, cross-entropy error could be more useful than misclassification error. This is because cross-entropy error takes into account the prediction probability, not just the predicted class like the misclassification error does, this makes it better to calibrate the prediction and penalize wrong predictions [1, pp. 109–119].

2 Assignment 2: Linear Regression and Ridge Regression

2.1 Introduction

This assignment focuses on analyzing Parkinson's disease symptom data using linear regression and ridge regression techniques. The main objective is to predict the motor UPDRS score (a measure of Parkinson's disease symptom severity) using various biomedical voice measurements. The dataset contains voice recordings from 42 people with early-stage Parkinson's disease, collected during a six-month trial of a telemonitoring device. The given data set consists of 5875 observations and 22 variables.

The assignment consists of four main tasks:

1. Data preparation and scaling
2. Initial linear regression analysis
3. Implementation of ridge regression functions
4. Analysis of ridge regression with different penalty parameters

2.2 Data Preparation

The data was split into 60% training data and 40% testing data. Then, the training data was scaled using caret (z-score standardization) so that every column has the same mean ($=0$) and standard deviation ($=1$) and using the same parameters, the test data was also scaled.

2.3 Linear Regression

A linear regression model was fitted using the scaled training data. Table 3 shows the coefficient estimates and their statistical significance for all predictor variables.

Table 3: Summary of linear regression coefficients and their statistical significance, ordered by absolute t-value

Variable	Estimate	Std. Error	t value	p-value
DFA	-0.280318	0.020136	-13.921	$< 2 \times 10^{-16}$ ***
PPE	0.226467	0.032881	6.887	6.70×10^{-12} ***
HNR	-0.238543	0.036395	-6.554	6.41×10^{-11} ***
Shimmer.APQ11	0.305546	0.061236	4.990	6.34×10^{-7} ***
Jitter.Abs.	-0.169609	0.040805	-4.157	3.31×10^{-5} ***
NHR	-0.185387	0.045567	-4.068	4.84×10^{-5} ***
Shimmer.APQ5	-0.387507	0.113789	-3.405	6.68×10^{-4} ***
Shimmer	0.592436	0.205981	2.876	4.05×10^{-3} **
Jitter...	0.186931	0.149561	1.250	2.11×10^{-1}
Shimmer.dB.	-0.172655	0.139316	-1.239	2.15×10^{-1}
Jitter.PPQ5	-0.074568	0.087766	-0.850	3.96×10^{-1}
Shimmer.DDA	-32.387241	77.158814	-0.420	6.75×10^{-1}
Shimmer.APQ3	32.070932	77.159242	0.416	6.78×10^{-1}
Jitter.DDP	5.249558	18.837525	0.279	7.81×10^{-1}
Jitter.RAP	-5.269544	18.834160	-0.280	7.80×10^{-1}
RPDE	0.004068	0.022664	0.179	8.58×10^{-1}
Significance codes: *** $p < 0.001$, ** $p < 0.01$, * $p < 0.05$				
$R^2 = 0.1212$, Adjusted $R^2 = 0.1172$, F-statistic = 30.25 on 16 and 3509 DF				
Residual standard error: 0.9394 on 3509 degrees of freedom				

The linear regression model shows an adjusted R-squared of 0.1172, indicating that about 12% of the variance in motor_UPDRS is explained by the voice characteristics. The most significant predictors are DFA ($t = -13.921$), PPE ($t = 6.887$), and HNR ($t = -6.554$).

MSE values:

- Training data: 0.8785
- Testing data: 0.9354

2.4 Ridge Regression Implementation

To perform ridge regression analysis, four core functions that work together to find optimal parameters while incorporating a ridge penalty were implemented. Each function serves a specific purpose in the implementation:

2.4.1 Log-likelihood Function

A function to compute the log-likelihood for normally distributed data was implemented:

- **Input:** Parameter vector θ (value of the coefficients for how much each feature contributes to the target variable), dispersion σ , and training data
- **Implementation:**
 - Constructs design matrix X with removed target value column
 - Makes predictions using $X\theta$

- Computes log-likelihood using the normal distribution formula:

$$-\frac{n}{2}\log(2\pi) - \frac{n}{2}\log(\sigma^2) - \frac{\sum(y - X\theta)^2}{2\sigma^2}$$

- **Purpose:** Measures how well our parameters explain the observed data, want to maximize this

2.4.2 Ridge Function

Built upon the log-likelihood function to create the ridge regression, it adds a ridge penalty to parameters:

- **Input:** Parameters θ , σ (the same θ and σ are used as input to the log-likelihood function), penalty parameter λ (bigger lambda -> coefficients will be smaller, and training data)
- **Implementation:**
 - Computes negative log-likelihood
 - Adds ridge penalty $\lambda \sum \theta_{(i)}^2$
- **Purpose:** Creates the complete objective function to be minimized

2.4.3 Ridge Optimization Function

This function finds the optimal parameters for the ridge regression model:

- **Input:** Penalty parameter λ and training data
- **Implementation:**
 - Creates a helper function that uses the ridge function (which already combines log-likelihood and ridge penalty)
 - Uses R's optimization function `optim()` with BFGS method
 - Starts with initial values (zeros for coefficients, 1 for σ)
 - Iteratively finds parameters that minimize the ridge function value
- **Purpose:** Automatically finds the best parameter values (θ and σ) that:
 - Minimize the ridge function value for a given λ
 - Return optimal coefficients and sigma for the model
- **Output:** Returns the optimal values for all parameters ($\theta_0, \theta_1, \dots, \theta_p, \sigma$)

2.4.4 Degrees of Freedom Function

This function calculates how flexible the ridge regression model is, tells us "how many effective parameters" we're using with the current values of λ :

- **Input:** Penalty parameter λ and training data
- **Implementation:**
 - Uses the formula $\text{trace}(X(X'X + \lambda I)^{-1}X')$ to calculate effective degrees of freedom
 - Larger λ leads to fewer effective degrees of freedom
 - Represents how much freedom the model has to fit the data
- **Purpose:** Measures how complex or flexible our model is after ridge penalty is applied

These functions work together in sequence to perform ridge regression:

1. Log-likelihood measures fit of parameters
2. Ridge function adds penalty for complexity
3. Optimization finds best parameters
4. Degrees of freedom quantifies model complexity

2.5 Ridge Regression Analysis

To evaluate the effect of different ridge penalties, we tested the model with three different values of λ (1, 100, and 1000) and compared them with the original linear regression model ($\lambda = 0$). For each model, we calculated:

- Training and test Mean Squared Error (MSE) for both the training and test data
- Effective degrees of freedom to measure model complexity

Table 4 shows the comparison of MSE values and degrees of freedom for different λ values:

Table 4: Comparison of model performance with different ridge penalties (λ)			
λ	Training MSE	Test MSE	Degrees of Freedom
0 (Initial Linear Regression)	0.8785	0.9354	16.0000
1	0.8786	0.9350	13.8607
100	0.8844	0.9323	9.9249
1000	0.9211	0.9539	5.6439

Among the tested values, $\lambda=100$ appears most appropriate as it achieves the lowest test MSE (0.9323). As λ increases, we can see a trade-off: the degrees of freedom decrease significantly (from 16 to 5.6), indicating reduced model complexity, while the training MSE gradually increases. The results suggest that moderate regularization ($\lambda=100$) provides slightly better generalization performance compared to both the unregularized model

($\lambda=0$) and strong regularization ($\lambda=1000$). The test MSE improves slightly from 0.9354 ($\lambda=0$) to 0.9323 ($\lambda=100$), while training MSE gets slightly worse (0.8785 to 0.8844). This suggests there was a small amount of overfitting that regularization helped with, but not much.

3 Assignment 3: Logistic Regression and Basis Function Expansion

3.1 Initial Data Visualization

The first task of assignment 3 was to create a scatterplot showing Plasma glucose concentration versus Age, with points colored based on Diabetes status. Based on the figure, diabetes is not easy to classify using a standard logistic regression model that uses these two variables because the data shows significant overlap between the classes.

3.2 Basic Logistic Regression

The code trains a logistic regression model using Plasma Glucose Concentration and Age as features to predict Diabetes. It calculates predicted probabilities, and then classifies observations using a threshold of 0.5. It also extracts the probabilistic equation of the estimated model, indicating how the target depends on the features and the estimated model parameters probabilistically. The training misclassification error has also been calculated. Finally, it visualizes the predicted classifications on a scatterplot, similar as the one in assignment 3.1, but now showing the predicted classifications of diabetes.

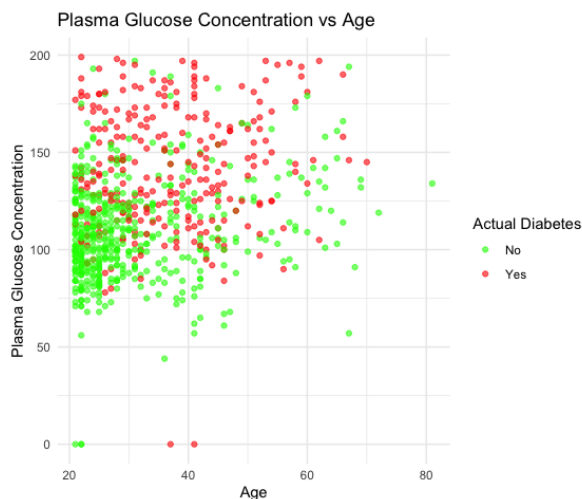


Figure 8: Scatterplot with points colored based on actual Diabetes status.



Figure 9: Scatterplot of trained logistic regression model with points based on predicted Diabetes status.

Logistic Regression - Probabilistic Equation:

$$\text{logit}(P(y = 1)) = -5.8979 + 0.0356 \cdot \text{Glucose} + 0.0245 \cdot \text{Age}$$

Training Misclassification Error: 26.60%

The model indicates that higher glucose and age increase the probability of diabetes, but the 26.60% misclassification error suggests moderate accuracy, implying room for improvement in feature selection or model complexity.

3.3 Decision Boundary Analysis

The decision boundary between the two classes has been calculated for the model used in assignment 3.2. This line has been implemented in the scatterplot.

Decision Boundary Equation: The decision boundary is determined where the probability of $y = 1$ equals 0.5, which corresponds to the logit function being equal to 0:

$$\text{logit}(P(y = 1)) = w_0 + w_1 \cdot \text{Glucose} + w_2 \cdot \text{Age} = 0$$

Rearranging for Age:

$$\text{Age} = -\frac{w_0}{w_2} - \frac{w_1}{w_2} \cdot \text{Glucose}$$

Substituting the coefficients from the model ($w_0 = -5.8979$, $w_1 = 0.0356$, $w_2 = 0.0245$):

$$\text{Age} = -\frac{-5.8979}{0.0245} - \frac{0.0356}{0.0245} \cdot \text{Glucose}$$

Simplifying:

$$\text{Age} = 240.73 - 1.453 \cdot \text{Glucose}$$

Thus, the **decision boundary equation is:**

$$\text{Age} = 240.73 - 1.453 \cdot \text{Glucose}$$

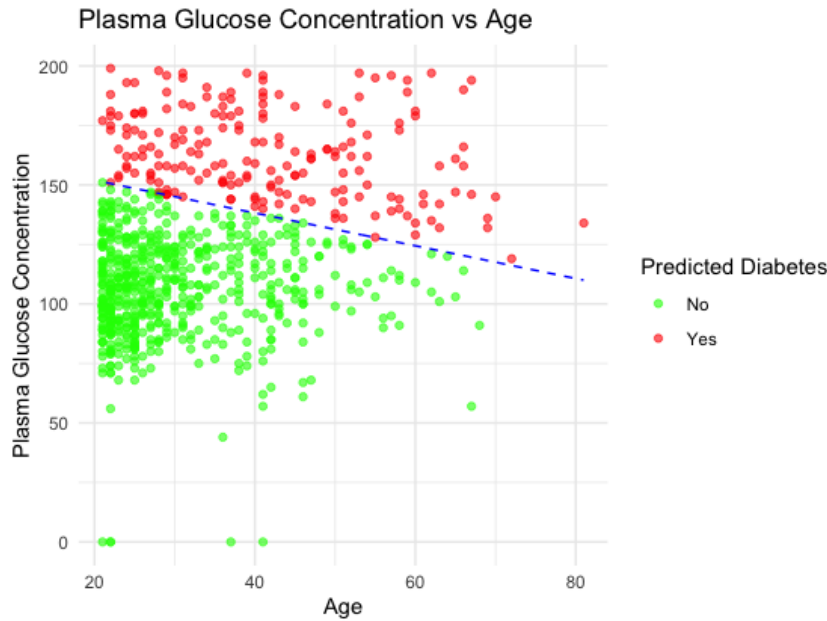


Figure 10: Scatterplot from assignment 3.2 with decision boundary line

The decision boundary in the scatterplot generally separates the two classes (diabetes vs. no diabetes) based on glucose and age. It captures the data distribution reasonably well, with most of the red points (diabetes) above the line and green points (no diabetes) below. However, there is some overlap, particularly near the boundary, indicating some misclassified points.

3.4 Threshold Analysis

The code script adjusts the decision threshold (r) for classifying diabetes and generates scatterplots for $r = 0.2$ and $r = 0.8$.

When $r = 0.2$, most predictions are classified as diabetes (red) since any probability greater than 0.2 is classified as positive ($=1$). Conversely, when $r = 0.8$, most predictions are classified as no diabetes (green), as only probabilities above 0.8 are classified as positive. This demonstrates the sensitivity of predictions to the chosen threshold.

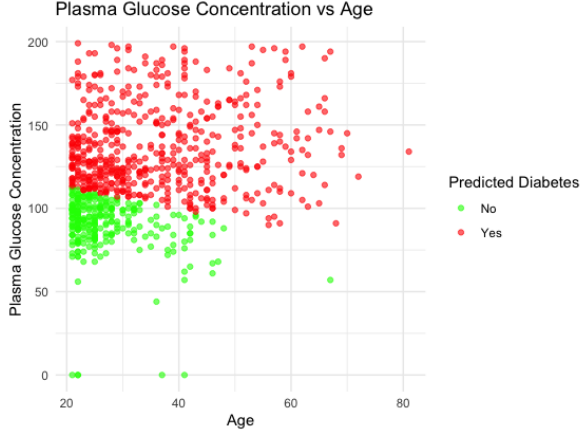


Figure 11: Scatterplot with $r = 0.2$: Majority classified as diabetes.

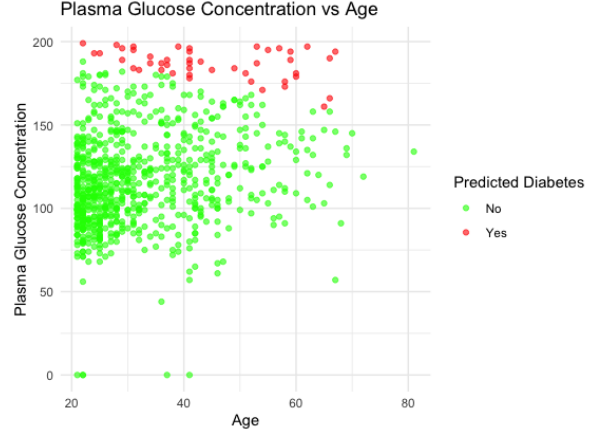


Figure 12: Scatterplot with $r = 0.8$: Majority classified as no diabetes.

3.5 Basis Function Expansion

The task involves applying a basis function expansion to the logistic regression model by creating new polynomial interaction features: $z_1 = x_1^4$, $z_2 = x_1^3 \cdot x_2$, $z_3 = x_1^2 \cdot x_2^2$, $z_4 = x_1 \cdot x_2^3$, and $z_5 = x_2^4$. Here, $x_1 = \text{PlasmaGlucoseConcentration}$ and $x_2 = \text{Age}$. These features were added to the dataset and used to train an extended logistic regression model.

Training Misclassification Error (Expanded Model): 24.64%

The basis function expansion slightly reduces the training misclassification error from 26.60% in the previous model to 24.64% in the extended model. While the new model fits the data better and captures more complexity, it risks overfitting due to its complexity. This confirms that predicting diabetes based on Plasma Glucose Concentration and Age alone is challenging, and the initial model may not be the most reliable.

The decision boundary for the basis expansion model is determined by solving the equation:

$$w_0 + w_1x_1 + w_2x_2 + w_3x_1^4 + w_4x_1^3x_2 + w_5x_1^2x_2^2 + w_6x_1x_2^3 + w_7x_2^4 = 0$$

Given the complexity of the equation, numerical methods are necessary to plot the decision boundary. The basis expansion trick has transformed the decision boundary from linear to a curve, reflecting the added polynomial interaction terms.

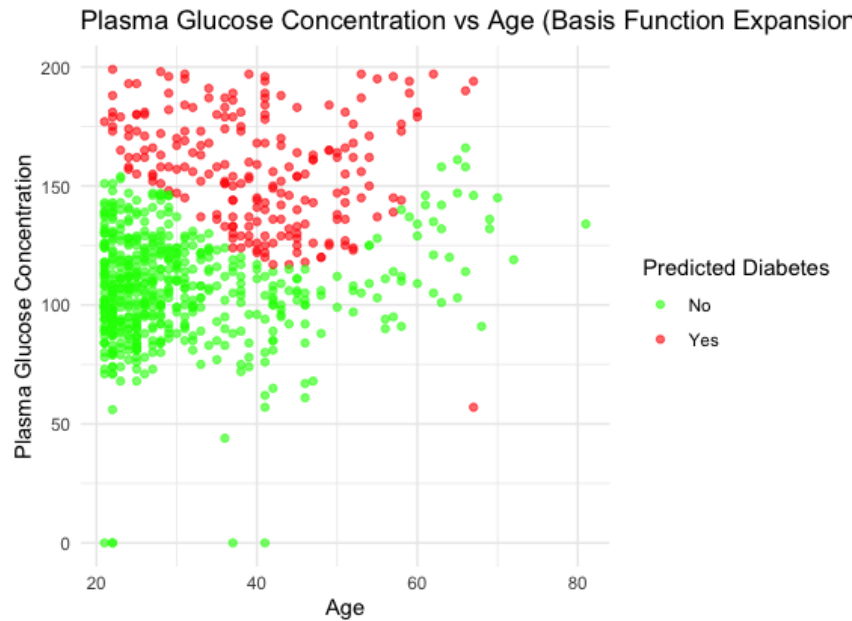


Figure 13: Scatterplot of predicted Diabetes status with Basis Expansion Model.

4 Assignment 4: Theory

In this assignment, three questions will be answered using the course book, *Machine Learning: A First Course for Engineers and Scientists* [1].

4.1 Probability Thresholds in Classification

Why can it be important to consider various probability thresholds in the classification problems, according to the book? In the book, it is stated that it could be a good idea to consider different probability thresholds since this could let the classifier achieve different things. The standard probability threshold is 0.5 which will make the misclassification error the smallest but real-world problems might not be the most important factor, by changing the threshold we can better predict asymmetric events. For example, in the medical field, it is more important to not miss a diagnosis of a sick person rather than false classifying a healthy person with a disease [1, pp. 49–50].

4.2 Target Variable Collection Methods

What ways of collecting correct values of the target variable for the supervised learning problems are mentioned in the book? In the book two different methods to get the target variable are mentioned: Joint recording and manually labeling the output. Joint recording is when one input gives an output that is easier to obtain than manually labeling the output. To manually label the output a branch expert needs to examine every input and then classify the most likely output [1, pp. 13–14].

4.3 Linear Regression Cost Function Matrix Form

How can one express the cost function of the linear regression in the matrix form, according to the book? In the book [1, pp. 41], one can find the formula 3.11 Which expresses the

cost function in linear regression in matrix form. In 3.11 we have

$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (\hat{y}(x_i; \boldsymbol{\theta}) - y_i)^2 = \frac{1}{n} \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2 = \frac{1}{n} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 = \frac{1}{n} \|\boldsymbol{\epsilon}\|_2^2, \quad (3.11)$$

where $\|\cdot\|_2$ denotes the usual Euclidean vector norm and $\|\cdot\|_2^2$ its square. Note that this is directly taken from the book.

References

- [1] A. Lindholm, N. Wahlström, F. Lindsten, and T. B. Schön, *Machine Learning: A First Course for Engineers and Scientists*. Cambridge University Press, Jul. 2022. Pre-publication version. Accessed online 24 Nov 2024.

A Code

```
1 ##### INSTALL NECESSARY PACKAGES #####
2 install.packages("kknn")
3 library(kknn)
4
5 ##### DIVIDE THE DATA #####
6
7 # Load the data into a variable
8 data <- read.csv("optdigits.csv", header=F)
9
10 # Get the number of rows in the dataset
11 n <- dim(data)[1]
12
13 # Set a random seed for reproducibility
14 set.seed(12345)
15
16 # Partition 50% of the data for the training set
17 id <- sample(1:n, floor(n * 0.5))
18 train <- data[id, ]
19
20 # Partition 25% of the data for the validation set
21 id1 <- setdiff(1:n, id)
22 set.seed(12345)
23 id2 <- sample(id1, floor(n * 0.25))
24 valid <- data[id2, ]
25
26 # Use the rest for the test set
27 id3 <- setdiff(id1, id2)
28 test <- data[id3, ]
29
30
31 ##### FIT A K-NEAREST NEIGHBOR MODEL TO TRAIN DATA #####
32
33 # Fitting = Learning = Training...?
34
35 # Fit the model on training data and test on training data
36 model_train <- kknn(as.factor(V65) ~ ., train = train, test =
37   ↪ train, k = 30, kernel = "rectangular")
38 train_pred <- fitted(model_train)
39 confusion_matrix_train <- table(Truth = train$V65, Predicted =
40   ↪ train_pred)
41 confusion_matrix_train
42
43 # Define the misclassification error function
44 missclass <- function(truth, predicted) {
45   n <- length(truth)
46   return(1 - sum(diag(table(truth, predicted))) / n)
47 }
48
49 # Misclassification errors for the training data
```

```

48 train_misclass_error <- missclass(train$V65, train_pred)
49 print(train_misclass_error) #0.04500262
50
51 # Fit the model on training data and test on test data
52 model_test <- kknn(as.factor(V65) ~ ., train = train, test = test
53   ↪ , k = 30, kernel = "rectangular")
54 test_pred <- fitted(model_test)
55 confusion_matrix_test <- table(Truth = test$V65, Predicted = test
56   ↪ _pred)
57 confusion_matrix_test
58
59 # Misclassification errors for the training data
60 test_misclass_error <- missclass(test$V65, test_pred)
61 print(test_misclass_error) #0.05329154
62
63 ##### TASK 1.3 #####
64 # Get probabilities of each class for each case in the training
65   ↪ data
66 train_probs <- model_train$prob
67 # Filter cases where the true label is "8"
68 eight_indices <- which(train$V65 == 8)
69
70 # Get the probabilities for the correct class (8) for each case
71 eight_probs <- train_probs[eight_indices, "8"]
72
73 # Sort by probability: highest for easiest, lowest for hardest
74 sorted_indices <- order(eight_probs, decreasing = TRUE)
75 easiest_indices <- eight_indices[sorted_indices[1:2]] # Two
76   ↪ highest probabilities
77 hardest_indices <- eight_indices[sorted_indices[(length(sorted_
78   ↪ indices)-2):length(sorted_indices)]] # Three lowest
79   ↪ probabilities
80
81 # Function to visualize an 8x8 matrix of a digit with index
82 visualize_digit <- function(data_row, label, index) {
83   # Reshape to 8x8 matrix
84   digit_matrix <- matrix(as.numeric(data_row), nrow = 8, ncol =
85     ↪ 8, byrow = TRUE)
86   # Plot heatmap
87   heatmap(digit_matrix, Rowv = NA, Colv = NA, scale = "none", col
88     ↪ = heat.colors(256),
89     main = paste("Digit", label, "- Index", index))
90 }
91
92 # Visualize easiest cases for digit "8" with indices
93 for (i in easiest_indices) {
94   visualize_digit(train[i, -ncol(train)], train$XV65[i], i) #
95     ↪ Include index
96 }
97
98 # Visualize hardest cases for digit "8" with indices

```

```

90 for (i in hardest_indices) {
91   visualize_digit(train[i, -ncol(train)], train$V65[i], i) #
92   ↪ Include index
93 }
94
95 ##### TASK 1.4 #####
96
97 Kvalues <- 1:30
98 train_errors <- numeric(length(Kvalues)) #creates a empty set of
99   ↪ vectors
100 valid_errors <- numeric(length(Kvalues)) #creates a empty set of
101   ↪ vectors
102
103 for (k in Kvalues) {
104   # Fit KNN model on training data and predict on training data
105   model_train <- knn(as.factor(V65) ~ ., train = train, test =
106     ↪ train, k = k, kernel = "rectangular")
107   train_pred <- fitted(model_train)
108   train_errors[k] <- missclass(train$V65, train_pred) # Store
109     ↪ training error
110
111   # Fit KNN model on training data and predict on validation data
112   model_valid <- knn(as.factor(V65) ~ ., train = train, test =
113     ↪ valid, k = k, kernel = "rectangular")
114   valid_pred <- fitted(model_valid)
115   valid_errors[k] <- missclass(valid$V65, valid_pred) # Store
116     ↪ validation error
117 }
118
119 # Set up plot for training errors with improved readability
120 plot(Kvalues, train_errors, type = "o", col = "blue", pch = 16,
121   ↪ lwd = 2, cex = 1.5,
122   xlab = "K Value", ylab = "Misclassification Error",
123   main = "Training and Validation Errors vs. K",
124   ylim = range(c(train_errors, valid_errors)), cex.lab = 1.5,
125   ↪ cex.main = 1.5)
126
127 # Add validation errors with a thicker line and larger points
128 lines(Kvalues, valid_errors, type = "o", col = "red", pch = 16,
129   ↪ lwd = 2, cex = 1.5)
130
131 # Add grid lines for readability
132 grid(nx = NULL, ny = NULL, lty = 2, col = "gray")
133
134 # Add a legend with larger font size
135 legend("bottomright", legend = c("Training Error", "Validation
136   ↪ Error"),
137   col = c("blue", "red"), pch = 16, lty = 1, lwd = 2, cex =
138     ↪ 1.2)
139
140

```

```

129 # Add a vertical line at the optimal K value
130 optimal_k <- which.min(valid_errors)
131 abline(v = Kvalues[optimal_k], col = "green", lty = 2, lwd = 2)
132 text(Kvalues[optimal_k], valid_errors[optimal_k] + 0.002, labels
      ↪ = paste("Optimal K =", Kvalues[optimal_k]), col = "black",
      ↪ pos = 4)
133
134 #Estimate the test error for k=optimal_k
135 model_test <- kkn(as.factor(V65) ~ ., train = train, test = test
      ↪ , k = optimal_k, kernel = "rectangular")
136 test_pred <- fitted(model_test)
137 test_misclass_error <- missclass(test$V65, test_pred)
138 print(test_misclass_error) #0.03239289
139
140
141 ##### TASK 1.5 #####
142
143 # Define the cross-entropy error function
144 CE <- function(truth, probability_matrix) {
145   epsilon <- 1e-15 # Small constant to avoid log(0)
146   cross_entropy <- 0
147
148   # Loop over each sample
149   for (i in 1:length(truth)) {
150     # Get the true class for the i-th sample
151     true_class <- as.integer(truth[i])+1
152
153     # Compute cross-entropy for this sample
154     cross_entropy <- cross_entropy - log(probability_matrix[i,
      ↪ true_class] + epsilon)
155   }
156
157   # Return the average cross-entropy error
158   return(cross_entropy / length(truth))
159 }
160
161 # Define the range of K values
162 Kvalues <- 1:30
163 cross_entropy_errors <- numeric(length(Kvalues)) # Create an empty
      ↪ vector for cross-entropy errors
164 epsilon <- 1e-15 # Small constant to avoid log(0)
165
166 # Loop over each K value
167 for (k in Kvalues) {
168   # Fit KNN model on training data and predict on validation data
169   model_valid <- kkn(as.factor(V65) ~ ., train = train, test =
      ↪ valid, k = k, kernel = "rectangular")
170
171   # Extract predicted probabilities for each class
172   probs <- model_valid$prob
173   valid_labels <- valid$V65

```

```

174
175
176
177 # Calculate cross-entropy error for the validation set and
178   ↪ store it
179 cross_entropy_errors[k] <- CE(valid_labels, probs)
180 }
181
182 # Plot Cross-Entropy Error vs. K
183 plot(Kvalues, cross_entropy_errors, type = "o", col = "blue", pch
184   ↪ = 16,
185       xlab = "K Value", ylab = "Cross-Entropy Error", main = "
186   ↪ Cross-Entropy Error vs. K")
187
188 # Identify and mark the optimal K (minimum cross-entropy)
189 optimal_k <- which.min(cross_entropy_errors)
190 abline(v = Kvalues[optimal_k], col = "green", lty = 2, lwd = 2)
191 points(Kvalues[optimal_k], cross_entropy_errors[optimal_k], col =
192   ↪ "red", pch = 19)
193 text(Kvalues[optimal_k], cross_entropy_errors[optimal_k], labels
194   ↪ = paste("Optimal K =", Kvalues[optimal_k]), pos = 3)

```

```

1 # ASSIGNMENT 2
2
3 # Package imports
4 library(caret)
5
6 ### TASK 1 ###
7
8 # read the data
9 data <- read.csv("parkinsons.csv")
10
11 voice_features <- c("Jitter...", "Jitter.Abs.", "Jitter.RAP", "
12   ↪ Jitter.PPQ5", "Jitter.DDP",
13   "Shimmer", "Shimmer.dB.", "Shimmer.APQ3", "
14   ↪ Shimmer.APQ5", "Shimmer.APQ11", "Shimmer
15   ↪ .DDA",
16   "NHR", "HNR",
17   "RPDE",
18   "DFA",
19   "PPE")
20
21 model_data <- data[, c(voice_features, "motor_UPDRS")]
22
23 # Get the number of rows in the data set
24 n <- dim(model_data)[1]
25
26 # set random seed for reproducibility
27 set.seed(12345)
28
29 # Partition 60% of data for the training set

```

```

27 id <- sample(1:n, floor(n * 0.6))
28 train <- model_data[id, ]
29
30 # Partition 40% of the data for the testing set
31 id1 <- setdiff(1:n, id)
32 test <- model_data[id1, ]
33
34 # Scale training data with caret
35 scaler <- preProcess(train) # Learns scaling from training data
36 trainScaled <- predict(scaler, train) # Applies scaling to
    ↪ training data
37
38 # Scale test data using same parameters from training data
39 testScaled <- predict(scaler, test)
40
41 ### TASK 2, Linear Regression Model ###
42
43 # Fit linear regression model using scaled training data
44 # motor_UPDRS is the target variable, use all other columns as
    ↪ predictors
45 # motor_UPDRS means we're trying to predict motor_UPDRS
46 fit <- lm(formula = motor_UPDRS ~ . - 1, data=trainScaled) # -1
    ↪ removes intercept
47
48 # Get summary to see which variables are significant
49 summary(fit) # DFA has a t-value of -13.921, the predictor with
    ↪ the highest absolute t-value
50
51 # Calculate MSE for training data
52 pred_train <- predict(fit, trainScaled) # get predictions for
    ↪ training data
53 mse_train <- mean((trainScaled$motor_UPDRS - pred_train)^2) #
    ↪ calculate training MSE
54
55 # Calculate MSE for test data
56 pred_test <- predict(fit, testScaled) # get predictions for test
    ↪ data
57 mse_test <- mean((testScaled$motor_UPDRS - pred_test)^2) #
    ↪ calculate test MSE
58
59 # Print MSE results
60 print(paste("Training MSE:", round(mse_train, 4)))
61 print(paste("Test MSE:", round(mse_test, 4)))
62
63 ### TASK 3 ###
64
65 # 3a #
66 loglikelihood <- function(theta, sigma, traindata){
67   # Get X matrix (predictors) from training data
68   X <- as.matrix(traindata[, -which(names(traindata) == "motor_
    ↪ UPDRS")])

```

```

69
70 # Get y (actual values) from training data
71 y <- traindata$motor_UPDRS
72
73 # Calculate predicted values (X * theta)
74 y_pred <- X %*% theta
75
76 # Calculate log-likelihood using normal distribution formula
77 n <- length(y) # Number of observations
78 loglik <- -n/2 * log(2*pi) - n/2 * log(sigma^2) - sum((y - y_
    ↪ pred)^2)/(2*sigma^2) # Minus for minimization
79
80 return(loglik)
81 }
82
83 # 3b #
84 ridge <- function(theta, sigma, lambda, traindata){
85   # Get negative log-likelihood (minus because we want to
    ↪ minimize)
86   neg_loglik <- -loglikelihood(theta, sigma, traindata)
87
88   # Calculate ridge penalty using: lambda * sum(theta^2)
89   # We exclude first theta (intercept) from penalty
90   ridge_penalty <- lambda * sum(theta^2)
91
92   # Return total negative log-likelihood + ridge penalty
93   return(neg_loglik + ridge_penalty)
94 }
95
96 # 3c #
97 # Function that takes lambda and finds optimal theta and sigma
98 ridgeOpt <- function(lambda, traindata) {
99   # Number of parameters we need (number of columns - 1 for motor
    ↪ _UPDRS)
100   n_params <- ncol(traindata) - 1
101
102   # Function that optim will minimize
103   # Takes parameters and returns ridge value
104   ridge_value <- function(params) {
105     theta <- params[1:n_params] # first n_params values are
    ↪ theta
106     sigma <- params[n_params + 1] # last value is sigma
107     return(ridge(theta, sigma, lambda, traindata))
108   }
109
110   # Starting values for optim: zeros for theta, 1 for sigma
111   start_values <- c(rep(0, n_params), 1)
112
113   # Run optimization
114   result <- optim(par = start_values,          # starting values
115                  fn = ridge_value,            # function to

```

```

116         ↪ minimize
           method = "BFGS")           # optimization method
117
118     # Return optimized parameters
119     return(result$par) # returns optimal theta and sigma
120 }
121
122 # 3d #
123 DF <- function(lambda, traindata){
124     # Remove motor_UPDRS
125     X <- as.matrix(traindata[, -which(names(traindata) == "motor_
           ↪ UPDRS")])
126
127     # Matrix Magic
128     p <- ncol(X)
129     XtX <- t(X) %*% X
130     ridge_matrix <- XtX + lambda * diag(p)
131     df <- sum(diag(X %*% solve(ridge_matrix) %*% t(X)))
132     return(df)
133 }
134
135 ### 4 ###
136
137 # ridgeOpt returns the optimal parameters using the ridge
           ↪ function
138
139 result1 <- ridgeOpt(lambda=1, trainScaled)           # lambda=1
140 result100 <- ridgeOpt(lambda=100, trainScaled)       # lambda=100
141 result1000 <- ridgeOpt(lambda=1000, trainScaled)     # lambda=1000
142
143 # Function to calculate MSE
144 calculate_mse <- function(theta, sigma, data) {
145     # Make X matrix without intercept
146     X <- as.matrix(data[, -which(names(data) == "motor_UPDRS")])
147     pred <- X %*% theta # Calculate predictions
148     mse <- mean((data$motor_UPDRS - pred)^2)
149     return(mse)
150 }
151
152 ##### Calculate MSE for each lambda #####
153
154 # Calculate how many predictors we have
155 n_params <- ncol(model_data) - 1 # -1 for motor_UPDRS
156
157 # Calculate MSE for each lambda
158 # Lambda = 1
159 mse_train1 <- calculate_mse(result1[1:n_params], result1[n_params
           ↪ + 1], trainScaled)
160 mse_test1 <- calculate_mse(result1[1:n_params], result1[n_params
           ↪ + 1], testScaled)
161

```



```

162 # Lambda = 100
163 mse_train100 <- calculate_mse(result100[1:n_params], result100[n_
    ↪ params + 1], trainScaled)
164 mse_test100 <- calculate_mse(result100[1:n_params], result100[n_
    ↪ params + 1], testScaled)
165
166 # Lambda = 1000
167 mse_train1000 <- calculate_mse(result1000[1:n_params], result1000
    ↪ [n_params + 1], trainScaled)
168 mse_test1000 <- calculate_mse(result1000[1:n_params], result1000[
    ↪ n_params + 1], testScaled)
169
170 # Calculate degrees of freedom (DF) for each lambda
171 df1 <- DF(1, trainScaled)
172 df100 <- DF(100, trainScaled)
173 df1000 <- DF(1000, trainScaled)
174
175 # Print the MSE results and degrees of freedom
176 print(paste("Lambda = 0 (without ridge function):    Train MSE:",
    ↪ round(mse_train, 4),
177          "Test MSE:", round(mse_test, 4),
178          "DF:", n_params)) # DF is number of predictors (no
    ↪ intercept)
179 print(paste("Lambda = 1:    Train MSE:", round(mse_train1, 4),
180          "Test MSE:", round(mse_test1, 4),
181          "DF:", round(df1, 4)))
182 print(paste("Lambda = 100:  Train MSE:", round(mse_train100, 4),
183          "Test MSE:", round(mse_test100, 4),
184          "DF:", round(df100, 4)))
185 print(paste("Lambda = 1000: Train MSE:", round(mse_train1000, 4),
186          "Test MSE:", round(mse_test1000, 4),
187          "DF:", round(df1000, 4)))
188
189
190 # Higher lambda -> Higher MSE -> Predictions get worse with
    ↪ increasing lambda
191 # This suggests that original model wasn't overfitting, no need
    ↪ for regularization
192
193 # Data frame of coefficients for comparison
194 coef_comparison <- data.frame(
195   Variable = colnames(trainScaled)[-which(names(trainScaled) == "
    ↪ motor_UPDRS)],
196   Lambda0 = fit$coefficients, # All coefficients (no intercept
    ↪ to exclude)
197   Lambda1 = result1[1:n_params],
198   Lambda100 = result100[1:n_params],
199   Lambda1000 = result1000[1:n_params]
200 )
201
202 # Print to look at the coefficients

```

```

203 print("Coefficient comparison:")
204 print(coef_comparison)
205
206
207
208 # ASSIGNMENT 3
209 # Logistic regression and basis function expansion
210
211 # Load necessary library
212 library(ggplot2)
213
214 # Load the data
215 data <- read.csv("pima-indians-diabetes.csv")
216
217 # Rename the columns
218 colnames(data) <- c("Pregnancies", "PlasmaGlucoseConcentration",
219   ↪ "BloodPressure", "SkinfoldThickness",
220   ↪ "Insulin", "BMI", "DiabetesPedigreeFunction",
221   ↪ "Age", "Diabetes")
222
223 ##### Assignment 1.1 #####
224 # Scatterplot of diabetes observations in a Glucose vs Age graph
225
226 # Create the scatterplot
227 ggplot(data, aes(x = Age, y = PlasmaGlucoseConcentration, color =
228   ↪ factor(Diabetes))) + # Define x/y-axes
229   geom_point(alpha = 0.6) + # Add points with 40% transparency
230   labs( # Labels for the plot
231     title = "Plasma Glucose Concentration vs Age",
232     x = "Age",
233     y = "Plasma Glucose Concentration",
234     color = "Actual Diabetes"
235   ) +
236   theme_minimal() + # Clean theme
237   scale_color_manual(values = c("green", "red"), labels = c("No",
238     ↪ "Yes")) # Customize color/label for Diabetes
239
240 ##### Assignment 1.2 #####
241 # Predictions of Diabetes based on a logistic regression model
242   ↪ trained with Glucose & Age as features
243
244 # Train the logistic regression model
245 logistic_model <- glm(Diabetes ~ PlasmaGlucoseConcentration + Age
246   ↪ , data = data, family = binomial)
247
248 # Make predictions for all observations!
249
250 # Predicted probabilities for target variable Diabetes
251 data$Predicted_Probabilities <- predict(logistic_model, type = "

```

```

    ↪ response") # Values between 0 and 1
248 # Convert predicted probabilities into binary classifications
    ↪ based on threshold = 0.5
249 data$Predicted_Classifications <- ifelse(data$Predicted_
    ↪ Probabilities > 0.5, 1, 0) # Values either 0 or 1
250
251 # Probabilistic equation of the estimated model
252 coefficients <- coef(logistic_model) # Extract the coefficients
253 intercept <- coefficients[1]
254 coef_glucose <- coefficients[2]
255 coef_age <- coefficients[3]
256 cat("Logistic Regression - Probabilistic Equation:\n")
257 cat(sprintf("logit(P(y = 1)) = %.4f + %.4f * Glucose + %.4f * Age
    ↪ \n",
258             intercept, coef_glucose, coef_age)) # Print the
    ↪ result
259
260 # Training misclassification error
261 error <- mean(data$Predicted_Classifications != data$Diabetes)
262 cat(sprintf("\nTraining Misclassification Error: %.2f%\n", error
    ↪ * 100)) # Print the result
263
264 # Plot the predicted Diabetes values based on Glucose & Age
265 ggplot(data, aes(x = Age, y = PlasmaGlucoseConcentration, color =
    ↪ factor(Predicted_Classifications))) + # Define x/y-axes
266 geom_point(alpha = 0.6) + # Add points with 40% transparency
267 labs( # Labels for the plot
268       title = "Plasma Glucose Concentration vs Age",
269       x = "Age",
270       y = "Plasma Glucose Concentration",
271       color = "Predicted Diabetes"
272     ) +
273     theme_minimal() + # Clean theme
274     scale_color_manual(values = c("green", "red"), labels = c("No",
    ↪ "Yes")) # Customize color/label for Diabetes
275
276
277 ##### Assignment 1.3 #####
278 # Decision boundary of the estimated logistic regression model
279
280 # The line where probability of y=1 is equal to 0.5 <==> logit
    ↪ function = 0
281 # logit(P(y=1)) = w0 + w1*Glucose + w2*Age = 0 <==> Age = -(w0/w2
    ↪ )-(w1*Glucose)/w2
282
283 # Function to calculate decision boundary between classes Glucose
    ↪ & Age
284 decision_boundary <- function(glucose) {
285   -(intercept/coef_age)-(coef_glucose/coef_age)*glucose
286 }
287

```

```

288 # Define Glucose range for boundary line
289 glucose_range <- data$PlasmaGlucoseConcentration # Use all
    ↳ Glucose values
290 age_boundary <- decision_boundary(glucose_range)
291
292 # Add the decision boundary in predicted Diabetes plot
293 ggplot(data, aes(x = Age, y = PlasmaGlucoseConcentration, color =
    ↳ factor(Predicted_Classifications))) + # Define x/y-axes
294   geom_point(alpha = 0.6) + # Add points with 40% transparency
    # Decision boundary line
295   geom_line(aes(x = age_boundary, y = glucose_range), color="blue
    ↳ ", linetype="dashed", size=0.5) +
296   labs( # Labels for the plot
297     title = "Plasma Glucose Concentration vs Age",
298     x = "Age",
299     y = "Plasma Glucose Concentration",
300     color = "Predicted Diabetes"
301   ) +
302   xlim(min(data$Age), max(data$Age)) + # Restrict the x-axis to
    ↳ the observed range of Age
303   theme_minimal() + # Clean theme
304   scale_color_manual(values = c("green", "red"), labels = c("No",
    ↳ "Yes")) # Customize color/label for Diabetes
305
306
307
308 ##### Assignment 1.4 #####
309 # Plotted graphs based on different r-values
310
311 # Convert predicted probabilities into binary classifications
    ↳ based on threshold = 0.2
312 data$Predicted_Classifications_02 <- ifelse(data$Predicted_
    ↳ Probabilities > 0.2, 1, 0) # Values either 0 or 1
313 # Convert predicted probabilities into binary classifications
    ↳ based on threshold = 0.8
314 data$Predicted_Classifications_08 <- ifelse(data$Predicted_
    ↳ Probabilities > 0.8, 1, 0) # Values either 0 or 1
315
316 # Plot r = 0.2
317 ggplot(data, aes(x = Age, y = PlasmaGlucoseConcentration, color =
    ↳ factor(Predicted_Classifications_02))) + # Define x/y-axes
318   geom_point(alpha = 0.6) + # Add points with 40% transparency
319   labs( # Labels for the plot
320     title = "Plasma Glucose Concentration vs Age",
321     x = "Age",
322     y = "Plasma Glucose Concentration",
323     color = "Predicted Diabetes"
324   ) +
325   theme_minimal() + # Clean theme
326   scale_color_manual(values = c("green", "red"), labels = c("No",
    ↳ "Yes")) # Customize color/label for Diabetes
327

```

```

328 # Plot r = 0.8
329 ggplot(data, aes(x = Age, y = PlasmaGlucoseConcentration, color =
    ↪ factor(Predicted_Classifications_08))) + # Define x/y-axes
330 geom_point(alpha = 0.6) + # Add points with 40% transparency
331 labs( # Labels for the plot
332     title = "Plasma Glucose Concentration vs Age",
333     x = "Age",
334     y = "Plasma Glucose Concentration",
335     color = "Predicted Diabetes"
336 ) +
337 theme_minimal() + # Clean theme
338 scale_color_manual(values = c("green", "red"), labels = c("No",
    ↪ "Yes")) # Customize color/label for Diabetes
339
340
341 ##### Assignment 1.5 #####
342 # Basis function expansion to introduce polynomial interaction
    ↪ terms into the logistic regression model
343 # x1 = PlasmaGlucoseConcentration, x2 = Age
344
345 # Compute new features for basis expansion and add these to the
    ↪ dataset as new columns
346 data$z1 <- data$PlasmaGlucoseConcentration^4
347 data$z2 <- data$PlasmaGlucoseConcentration^3 * data$Age
348 data$z3 <- data$PlasmaGlucoseConcentration^2 * data$Age^2
349 data$z4 <- data$PlasmaGlucoseConcentration * data$Age^3
350 data$z5 <- data$Age^4
351
352 # Train the new logistic regression model using expanded new
    ↪ features z1...z5
353 logistic_model_expanded <- glm(Diabetes ~
    ↪ PlasmaGlucoseConcentration + Age + z1 + z2 + z3 + z4 + z5,
354     data = data,
355     family = binomial)
356
357 # Extract predicted values and classify them to 0 or 1 with r =
    ↪ 0.5
358 data$Predicted_Probabilities_Expanded <- predict(logistic_model_
    ↪ expanded, type = "response") # Values between 0 and 1
359 data$Predicted_Classifications_Expanded <- ifelse(data$Predicted_
    ↪ Probabilities_Expanded > 0.5, 1, 0) # Values either 0 or 1
360
361 # Training misclassification error
362 error_expanded <- mean(data$Predicted_Classifications_Expanded !=
    ↪ data$Diabetes)
363 cat(sprintf("\nTraining Misclassification Error (Expanded Model):
    ↪ %.2f%\n", error_expanded * 100)) # Print the result
364
365 # Scatterplot of predicted Diabetes with new model
366 ggplot(data, aes(x = Age, y = PlasmaGlucoseConcentration, color =
    ↪ factor(Predicted_Classifications_Expanded))) + # Define x/y

```

```

↪ -axes
367 geom_point(alpha = 0.6) + # Add points with 40% transparency
368 labs( # Labels for the plot
369     title = "Plasma Glucose Concentration vs Age (Basis Function
↪ Expansion)",
370     x = "Age",
371     y = "Plasma Glucose Concentration",
372     color = "Predicted Diabetes"
373 ) +
374 theme_minimal() + # Clean theme
375 scale_color_manual(values = c("green", "red"), labels = c("No",
↪ "Yes")) # Customize color/label for Diabetes

```