

# Lab block 2

johmo870, nisra674

2024-12-07

## Statement of Contribution

- Nisal Amashan(nisra674) - Assignment 1, Assignment 3, Report
- John Möller (johmo870) - Assignment 2, Assignment 4, Report

## Assignment 1: Explicit Regularization

The goal is to investigate whether a near-infrared absorbance spectrum can predict fat content in meat samples.

The dataset contains:

- **Features:** A 100-channel spectrum of absorbance measurements.
- **Target:** Fat content in meat samples.

We perform:

1. Linear regression to model fat content.
2. LASSO regression to explore feature selection and sparsity.
3. Ridge regression for comparison.
4. Cross-validation to find the optimal penalty parameter.

## Results

### 1. Linear Regression

#### Training and Test Errors

- **Training MSE:** 0.0089  
This indicates a very good fit to the training data, suggesting that the model captures the relationships in the training set effectively.
- **Testing MSE:** 337.2898  
This high value for the test data suggests poor generalization to unseen data, indicating overfitting. The model performs well on the training set but fails to predict accurately on the test set, compromising its predictive capability.

## Cost Function for LASSO Regression

The LASSO regression cost function is defined as:

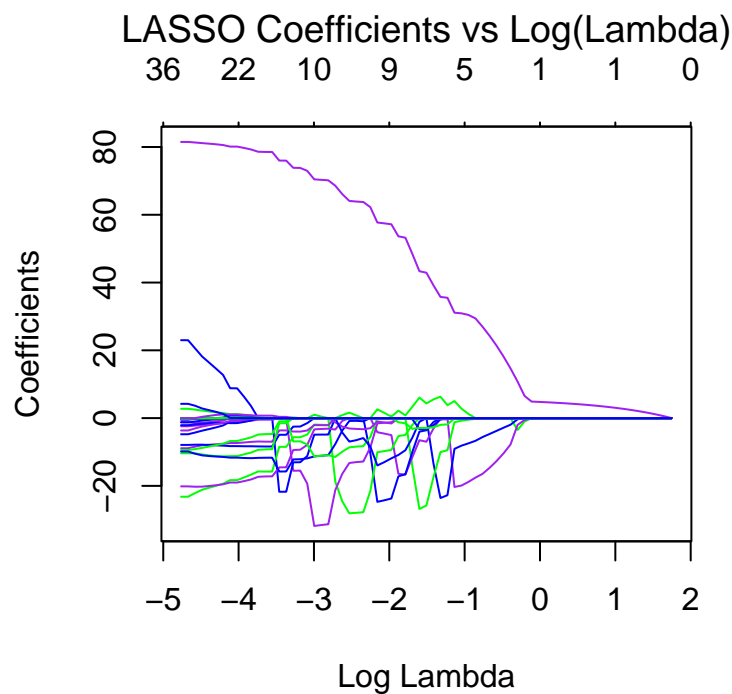
$$J(\beta) = \frac{1}{N} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

Where:

- $N$ : Number of samples.
- $p$ : Number of features (100 absorbance channels).
- $\beta_0$ : Intercept.
- $\beta_j$ : Coefficients of the features.
- $x_{ij}$ : Value of the  $j$ -th feature for the  $i$ -th sample.
- $\lambda$ : Regularization parameter (controls the strength of penalty on coefficients).

### 2. The parameter $\lambda$ determines the balance:

- **Larger  $\lambda$** : Stronger regularization, more coefficients shrink to zero.
- **Smaller  $\lambda$** : Weaker regularization,  $\beta$ s closer to ordinary least squares.



## Lambda for 3 features: 0.6764156

##	Feature	Coefficient
## 1	(Intercept)	18.362617
## 2	Channel17	-7.785416
## 3	Channel18	-1.940834
## 4	Channel141	14.654522

## Selecting a Model with 3 Features

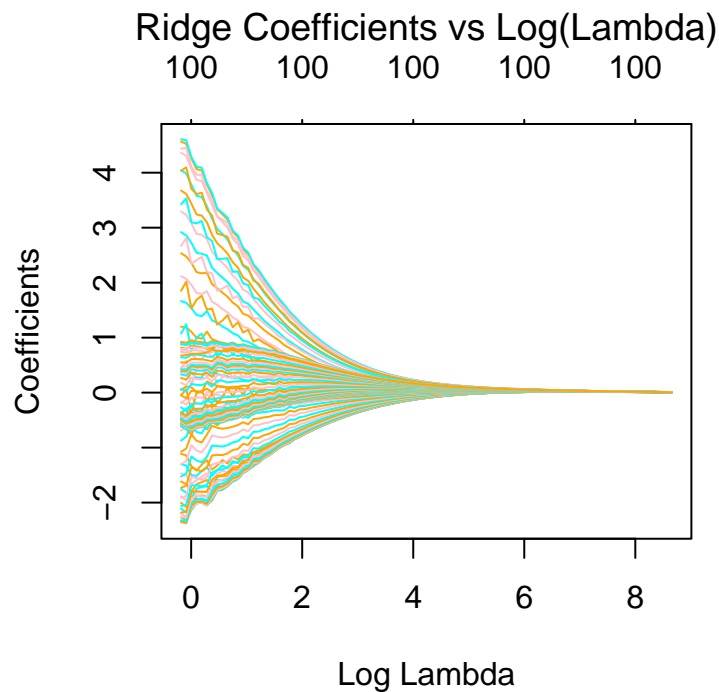
To identify a model with exactly 3 non-zero features, the coefficients were examined across different lambda values. The penalty factor Lambda for which only 3 features remain non-zero is:

Lambda = 0.6764

At this value of Lambda, the retained features and their coefficients are:

Feature	Coefficient
(Intercept)	18.36
Channel7	-7.79
Channel8	-1.94
Channel41	14.65

- The plot clearly illustrates the shrinking effect of Lambda on the coefficients and highlights the ability of LASSO to select a sparse set of features.
- At Lambda = 0.6764, the model selects **Channel7**, **Channel8**, and **Channel41** as the most significant predictors for Fat, along with the intercept term.
- This demonstrates LASSO's effectiveness in both feature selection and reducing model complexity.



## Key Observations from the Plots

### LASSO Plot

- Coefficients gradually shrink, with many becoming exactly zero as  $\lambda$  increases, indicating feature selection.
- Demonstrates sparsity, where only a small subset of features remains significant at higher  $\lambda$  values.

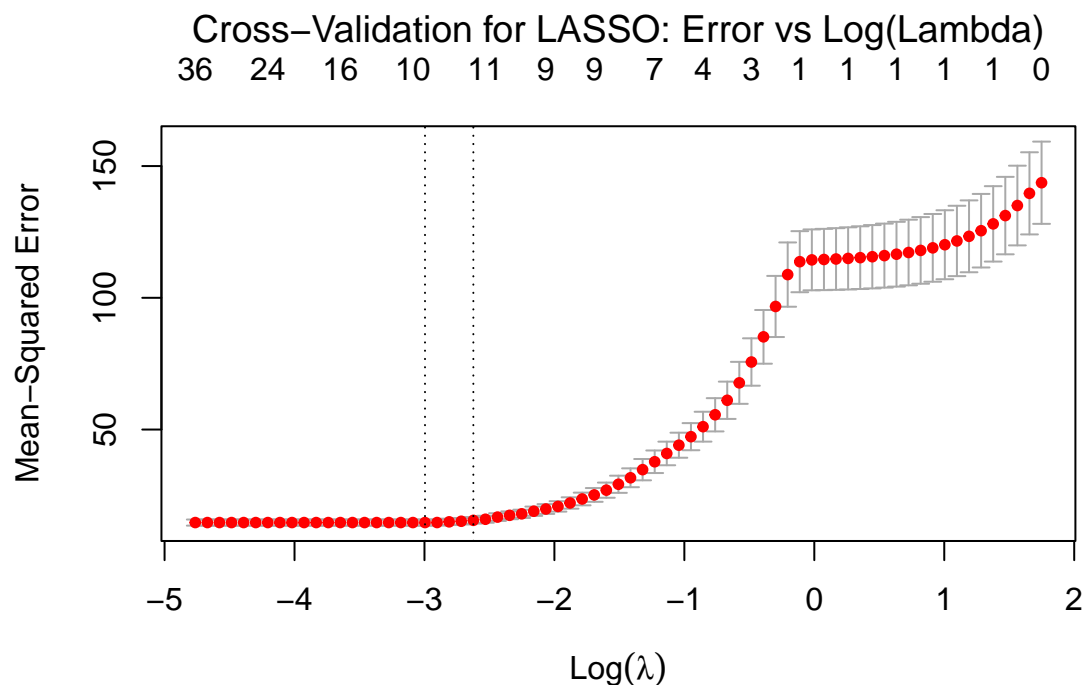
## Ridge Plot

- Coefficients decrease smoothly and uniformly, without any feature being completely eliminated.
- Highlights Ridge's ability to regularize the model while retaining all features, which is useful when all predictors are assumed to be relevant.

## Conclusions

- **LASSO** is more effective for feature selection in datasets with many irrelevant predictors. It provides a sparse model by retaining only a subset of features, which can improve interpretability.
- **Ridge** is more suitable for scenarios where all predictors are expected to contribute to the outcome. It minimizes overfitting by shrinking coefficients without eliminating any predictors.
- In this case, **LASSO** might be more appropriate if only a few absorbance channels are truly predictive of Fat. **Ridge**, on the other hand, would be beneficial if all features have some predictive value.

## 4. Cross-Validation



The LASSO regression model was fitted using cross-validation to determine the optimal  $\lambda$  value. The plot below shows the mean-squared error (MSE) as a function of  $\log(\lambda)$  with error bars indicating variability across the folds.

- The CV score decreases as  $\lambda$  decreases ( $\log(\lambda)$  increases), reaching a minimum at the optimal  $\lambda$ .
- After this point, the CV score increases as the model becomes less regularized and overfitting occurs.
- The optimal  $\lambda$  is identified as:

$$\lambda_{\text{optimal}} = 0.0549$$

At this value, the model achieves the lowest cross-validation error.

## Selected Features and Test Performance

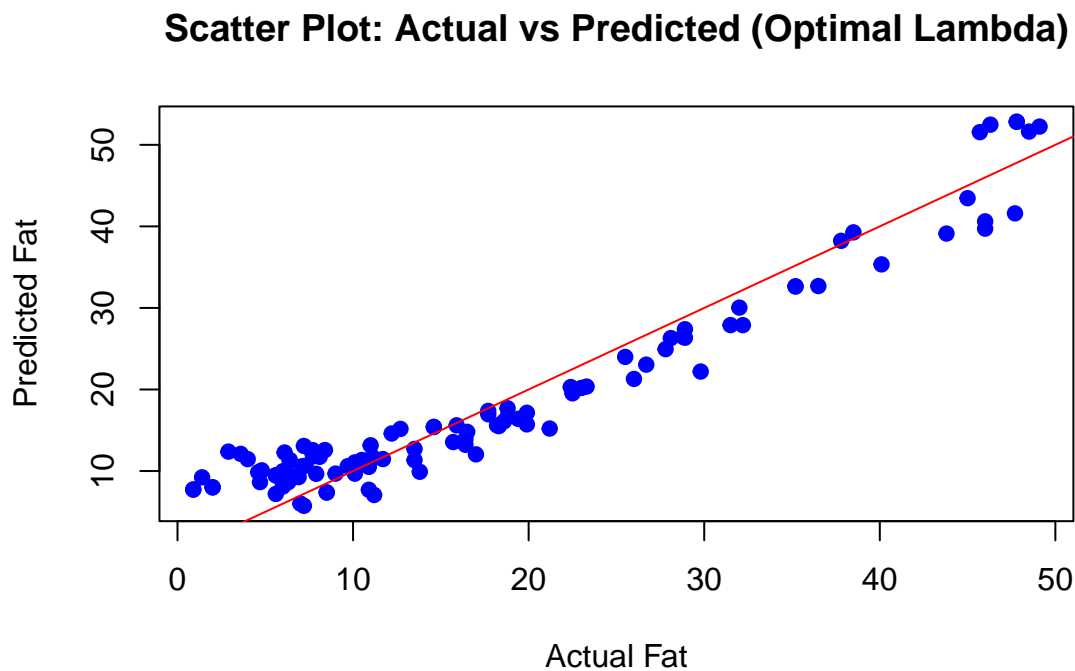
- **Number of Non-Zero Features:** At  $\lambda_{optimal}$ , the model retains 9 non-zero features, highlighting its ability to select important predictors while enforcing sparsity.
- **Test Set Performance:** The model's mean-squared error (MSE) on the test set at  $\lambda_{optimal}$  is 14.57.

### Comparison with $\log(\lambda) = -4$

- At  $\log(\lambda) = -4$ , the test MSE is 13.09, which is slightly lower than the test MSE for the optimal  $\lambda$ .
- While the test MSE for  $\lambda_{optimal}$  is slightly worse, the difference is minor, suggesting that both models perform comparably.
- However, the model at  $\lambda_{optimal}$  may be preferable due to its fewer retained features and reduced complexity.

## Scatter Plot: Actual vs Predicted (Optimal Lambda)

The scatter plot of actual versus predicted values for  $\lambda_{optimal}$  is shown below:



- The points lie close to the diagonal line, indicating good agreement between actual and predicted values.
- This suggests that the model provides reasonable predictions with the selected features.

## Final Comments

- The cross-validation process effectively balances model complexity and generalization, selecting  $\lambda_{optimal} = 0.0549$ .

- While the test MSE for  $\lambda_{optimal}$  is slightly higher than for  $\log(\lambda) = -4$ , the optimal  $\lambda$  model is simpler and less prone to overfitting.
- The scatter plot further confirms the predictive power of the model, as it closely aligns actual and predicted values.

## Assignment3: Principal components and implicit regularization

### Part1

## Number of components needed for 95% variance: 35

## Proportion of variance explained by PC1: 0.2501699

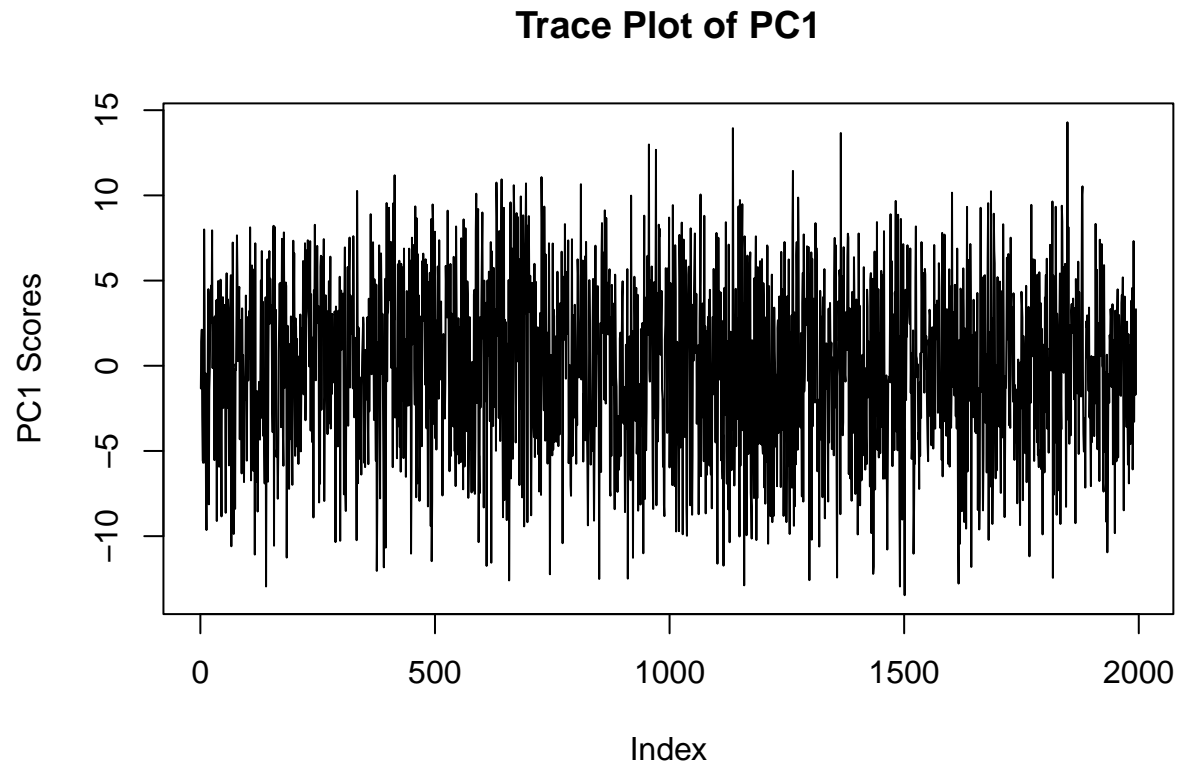
## Proportion of variance explained by PC2: 0.1693597

### Interpretation

- **Number of Components for 95% Variance:** 35.
  - This indicates that 35 principal components are required to capture at least 95% of the variance in the data.
- **Proportion of Variance Explained:**
  - **PC1:** 0.2502 (25.02%)
  - **PC2:** 0.1694 (16.94%)
  - Together, these two components explain approximately 42% of the total variance in the dataset.

## Part2

PCA using princomp()

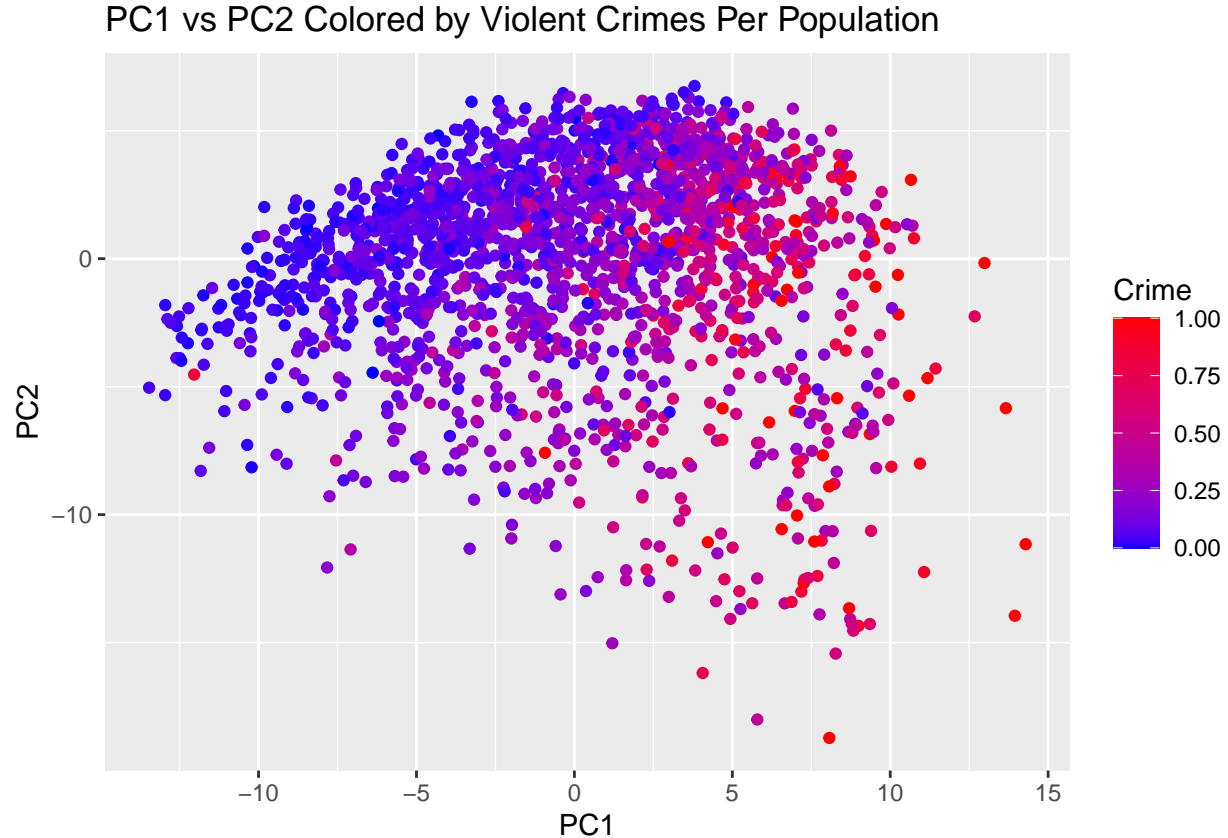


## Top Contributing Features

## Top 5 features contributing to PC1 (absolute values):

```
## [1] "medFamInc"      "medIncome"      "PctKids2Par"    "pctWInvInc"
## [5] "PctPopUnderPov"
```

## PC1 and PC2 Plot



## Results and Analysis

### 1. Top 5 Contributing Features:

- The top 5 features contributing to PC1 (in absolute value) are:
  - medFamInc (Median Family Income)
  - medIncome (Median Income)
  - PctKids2Par (Percentage of Kids in Two-Parent Households)
  - pctWInvInc (Percentage of Households with Investment Income)
  - PctPopUnderPov (Percentage of Population Under Poverty Line)

### 2. Interpretation:

- These features relate to socio-economic factors such as income, family structure, and poverty level.
- There seems to be a logical connection between these factors and the crime rate (**ViolentCrimesPerPop**).
- For instance, higher poverty rates or lower family income might be associated with higher crime levels, while stable family structures might reduce crime.

### 3. PC1 vs PC2 Plot:

- The plot of PC1 vs PC2 reveals clusters and trends where areas with high **ViolentCrimesPerPop** (in red) may be distinguishable.
- This analysis suggests that principal components capture socio-economic patterns related to crime.



### Part3

- **Training Error:** 0.01505
- **Test Error:** 0.02008

### Interpretation

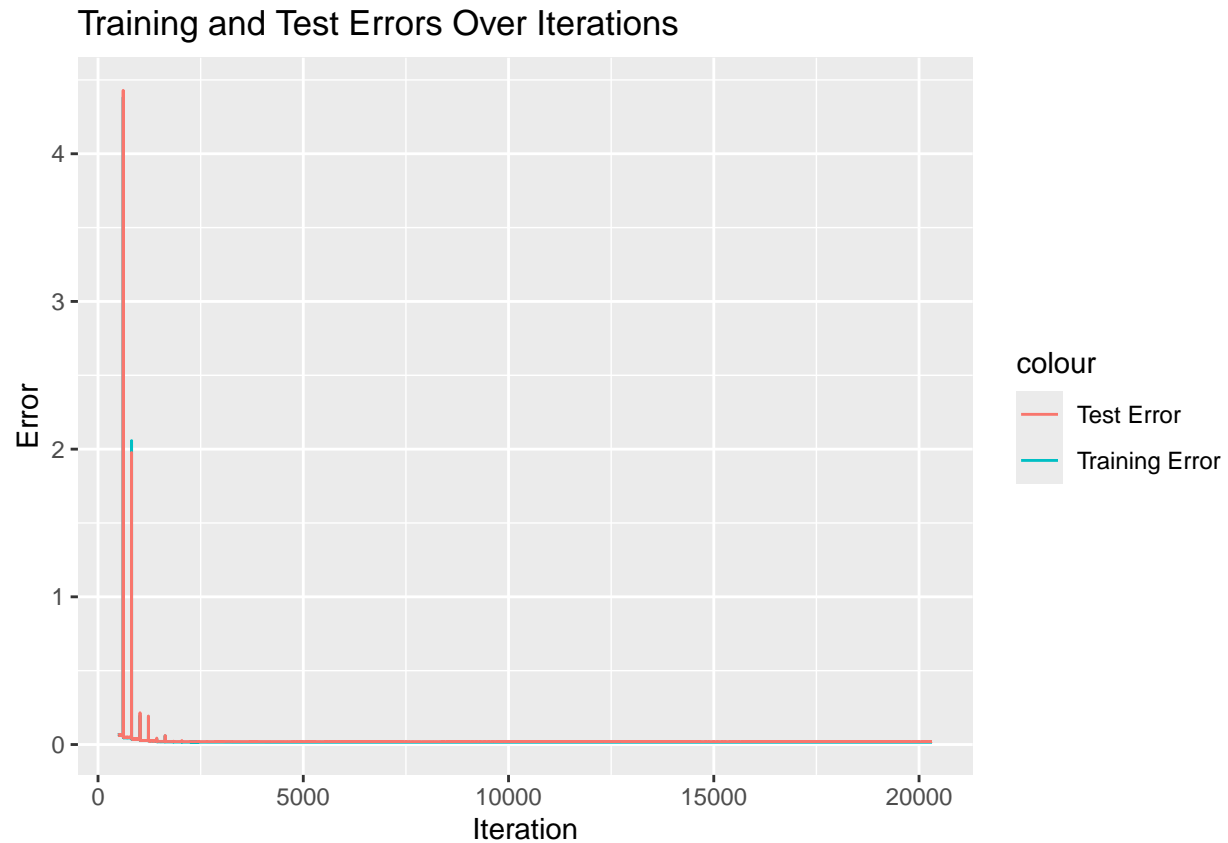
1. **Training Error:** The training error is relatively low, indicating that the model fits the training data well.
2. **Test Error:** The test error is slightly higher than the training error, but it is still low. This suggests that the model generalizes well to unseen data.
3. **Model Quality:**
  - The low test error indicates that the linear regression model is able to capture the relationship between the features and the target variable (**ViolentCrimesPerPop**) effectively.
  - However, linear regression assumes a linear relationship between predictors and the response. If non-linear relationships exist, other methods like regularized regression or non-linear models may perform better.

### Part4

#### Plot: Training and Test Errors Over Iterations

The plot below shows the training and test errors over the iterations of the BFGS optimization process. Initial iterations (first 500) are discarded for better visibility.

```
## initial  value 0.108368
## iter   10 value 0.016762
## iter   20 value 0.015241
## iter   30 value 0.015109
## iter   40 value 0.015073
## iter   50 value 0.015058
## iter   60 value 0.015051
## iter   70 value 0.015048
## iter   80 value 0.015047
## iter   90 value 0.015046
## iter  100 value 0.015046
## final   value 0.015046
## stopped after 100 iterations
```



#### Optimal Iteration and Errors

- **Optimal Iteration:** 1895
- **Optimal Training Error:** 0.01595
- **Optimal Test Error:** 0.01975

#### Comparison with Step 3

- **Step 3 Errors:**
  - Training Error: 0.01505
  - Test Error: 0.02008
- **Step 4 Errors:**
  - Training Error (Optimal): 0.01595
  - Test Error (Optimal): 0.01975

#### Interpretation and Conclusions

##### 1. Early Stopping:

- The optimal iteration (1895) minimizes the test error, ensuring the model generalizes well to unseen data. Early stopping prevents overfitting, as evidenced by the slightly higher training error compared to Step 3.

## 2. Improved Test Error:

- The test error in Step 4 is slightly lower than in Step 3, demonstrating improved generalization due to the early stopping criterion.

## 3. Model Quality:

- The model performs well, balancing training and test errors effectively.
- The small difference between training and test errors indicates good generalization capability.

## 4. Practical Insights:

- Early stopping based on test error is a practical and effective strategy to improve the model's performance without overfitting.
- 

# Appendix

## Code for Assignment 1

```
# Load the data
meat_data <- read.csv("data/tecator.csv")

# Set seed for reproducibility
set.seed(123)

# Extract features and target variable
absorbance_features <- meat_data[, 2:101]
fat_content <- meat_data$Fat

# Split data into training (50%) and testing (50%)
train_indices <- sample(1:nrow(meat_data), size = nrow(meat_data) * 0.5)
train_features <- absorbance_features[train_indices, ]
train_target <- fat_content[train_indices]
test_features <- absorbance_features[-train_indices, ]
test_target <- fat_content[-train_indices]

# Standardize features
train_features_scaled <- scale(train_features)
test_features_scaled <- scale(test_features)

# Linear Regression
linear_model <- lm(train_target ~ ., data = as.data.frame(cbind(train_target, train_features_scaled)))
train_predictions <- predict(linear_model, newdata = as.data.frame(train_features_scaled))
test_predictions <- predict(linear_model, newdata = as.data.frame(test_features_scaled))
mse_train <- mean((train_target - train_predictions)^2)
mse_test <- mean((test_target - test_predictions)^2)
cat("Training MSE:", mse_train, "\n")
cat("Testing MSE:", mse_test, "\n")

# Load glmnet library
library(glmnet)
```

```

train_features_matrix <- as.matrix(train_features_scaled)
train_target_vector <- as.vector(train_target)

# LASSO Regression
lasso_model <- glmnet(train_features_matrix, train_target_vector, alpha = 1)
plot(lasso_model, xvar = "lambda", label = TRUE, col = c("blue", "green", "purple"))
title("LASSO Coefficients vs Log(Lambda)")

# Lambda with 3 non-zero features
lambda_values <- lasso_model$lambda
for (lambda in lambda_values) {
  coef_nonzero <- coef(lasso_model, s = lambda)[-1]
  num_features <- sum(coef_nonzero != 0)
  if (num_features == 3) {
    cat("Lambda for 3 features:", lambda, "\n")
    break
  }
}

# Extract and display coefficients for 3-feature lambda
selected_coefficients <- coef(lasso_model, s = lambda)
selected_coefficients_matrix <- as.matrix(selected_coefficients)
nonzero_indices <- which(selected_coefficients_matrix != 0)
selected_features <- data.frame(
  Feature = rownames(selected_coefficients_matrix)[nonzero_indices],
  Coefficient = selected_coefficients_matrix[nonzero_indices]
)
print(selected_features)

# Ridge Regression
ridge_model <- glmnet(train_features_matrix, train_target_vector, alpha = 0)
plot(ridge_model, xvar = "lambda", label = TRUE, col = c("orange", "cyan", "pink"))
title("Ridge Coefficients vs Log(Lambda)")

# Cross-Validation for LASSO
cv_lasso <- cv.glmnet(train_features_matrix, train_target_vector, alpha = 1)
plot(cv_lasso, col = "darkgreen")
title("Cross-Validation for LASSO: Error vs Log(Lambda)")
optimal_lambda <- cv_lasso$lambda.min
cat("Optimal Lambda:", optimal_lambda, "\n")

# Non-zero features for optimal lambda
optimal_coefficients <- coef(cv_lasso, s = "lambda.min")
nonzero_features_count <- sum(optimal_coefficients != 0) - 1
cat("Number of non-zero features at optimal lambda:", nonzero_features_count, "\n")

# Test performance for optimal lambda
lasso_test_predictions <- predict(cv_lasso, s = optimal_lambda, newx = as.matrix(test_features_scaled))
mse_test_lasso <- mean((test_target - lasso_test_predictions)^2)
cat("Test MSE for LASSO with Optimal Lambda:", mse_test_lasso, "\n")

# Comparison with log(lambda) = -4
lambda_at_log_minus4 <- exp(-4)

```

```

test_predictions_log_minus4 <- predict(lasso_model, s = lambda_at_log_minus4, newx = as.matrix(test_fea
mse_test_log_minus4 <- mean((test_target - test_predictions_log_minus4)^2)
cat("Test MSE for LASSO with log(lambda) = -4:", mse_test_log_minus4, "\n")

# Scatter plot for optimal lambda
plot(test_target, lasso_test_predictions,
      main = "Scatter Plot: Actual vs Predicted (Optimal Lambda)",
      xlab = "Actual Fat", ylab = "Predicted Fat", pch = 19, col = "blue")
abline(0, 1, col = "red")

```

## Code for Assignment 2

## Code for Assignment 3

```

# PART 1: PCA Using eigen()

# Load necessary libraries
library(dplyr)

# Data Preprocessing
data <- read.csv("data/communities.csv")
scaled_data <- data %>% select(-ViolentCrimesPerPop) %>% scale()

# PCA using eigen()
cov_matrix <- cov(scaled_data)
eigen_results <- eigen(cov_matrix)

# Calculate variance explained
eigenvalues <- eigen_results$values
explained_variance <- eigenvalues / sum(eigenvalues)
cumulative_variance <- cumsum(explained_variance)

# Find number of components needed for 95% variance
num_components_95 <- which(cumulative_variance >= 0.95)[1]
pc1_variance <- explained_variance[1]
pc2_variance <- explained_variance[2]

# Output results
cat("Number of components for 95% variance:", num_components_95, "\n")
cat("Proportion of variance explained by PC1:", pc1_variance, "\n")
cat("Proportion of variance explained by PC2:", pc2_variance, "\n")

# PART 2: PCA Using princomp()

# Perform PCA
pca_results <- princomp(scaled_data, cor = TRUE)

# Plot PC1 trace
plot(pca_results$scores[, 1], type = "l", main = "Trace Plot of PC1", ylab = "PC1 Scores")

```

```

# Identify top 5 contributing features to PC1
loadings <- pca_results$loadings[, 1]
top_5_features <- sort(abs(loadings), decreasing = TRUE)[1:5]

# Output results
cat("Top 5 features contributing to PC1 (absolute values):\n")
print(names(top_5_features))

# Plot PC1 vs PC2 with crime levels
library(ggplot2)
plot_data <- data.frame(PC1 = pca_results$scores[, 1],
                        PC2 = pca_results$scores[, 2],
                        Crime = data$ViolentCrimesPerPop)
ggplot(plot_data, aes(x = PC1, y = PC2, color = Crime)) +
  geom_point() +
  scale_color_gradient(low = "blue", high = "red") +
  ggtitle("PC1 vs PC2 Colored by Violent Crimes Per Population")

# PART 3: Linear Regression

# Split data into training and test sets
set.seed(123)
train_indices <- sample(1:nrow(data), nrow(data) / 2)
train_data <- data[train_indices, ]
test_data <- data[-train_indices, ]

# Scale features and prepare responses
scaled_train <- scale(train_data %>% select(-ViolentCrimesPerPop))
scaled_test <- scale(test_data %>% select(-ViolentCrimesPerPop),
                    center = attr(scaled_train, "scaled:center"),
                    scale = attr(scaled_train, "scaled:scale"))
train_response <- train_data$ViolentCrimesPerPop
test_response <- test_data$ViolentCrimesPerPop

# Fit linear regression
model <- lm(train_response ~ ., data = as.data.frame(scaled_train))

# Compute training and test errors
train_predictions <- predict(model, as.data.frame(scaled_train))
test_predictions <- predict(model, as.data.frame(scaled_test))
train_error <- mean((train_response - train_predictions)^2)
test_error <- mean((test_response - test_predictions)^2)

# Output results
cat("Training error:", train_error, "\n")
cat("Test error:", test_error, "\n")

# PART 4: BFGS Optimization with Early Stopping

# Create cost function for tracking errors
iteration_errors <- list()

```

```

iteration_test_errors <- list()

cost_function_with_tracking <- function(theta, X, y, X_test, y_test) {
  train_predictions <- X %*% theta
  train_error <- y - train_predictions
  train_mse <- mean(train_error^2)

  test_predictions <- X_test %*% theta
  test_error <- y_test - test_predictions
  test_mse <- mean(test_error^2)

  iteration_errors <- append(iteration_errors, list(train_mse))
  iteration_test_errors <- append(iteration_test_errors, list(test_mse))

  return(train_mse)
}

# Prepare matrices for optimization
X_train <- as.matrix(cbind(1, scaled_train))
X_test <- as.matrix(cbind(1, scaled_test))

# Run optimization
optim_results <- optim(rep(0, ncol(X_train)),
  fn = function(theta) cost_function_with_tracking(theta, X_train, train_response,
    method = "BFGS",
    control = list(trace = 1, maxit = 100))

# Prepare error data for plotting
errors <- data.frame(Iteration = 1:length(iteration_errors),
  TrainingError = unlist(iteration_errors),
  TestError = unlist(iteration_test_errors))
filtered_errors <- errors[-(1:500), ]

# Plot training and test errors
ggplot(filtered_errors, aes(x = Iteration)) +
  geom_line(aes(y = TrainingError, color = "Training Error")) +
  geom_line(aes(y = TestError, color = "Test Error")) +
  ggtitle("Training and Test Errors Over Iterations") +
  ylab("Error") + xlab("Iteration")

# Identify optimal iteration
optimal_iteration <- which.min(filtered_errors$TestError)
cat("Optimal iteration:", optimal_iteration, "\n")
cat("Optimal Training Error:", filtered_errors$TrainingError[optimal_iteration], "\n")
cat("Optimal Test Error:", filtered_errors$TestError[optimal_iteration], "\n")

```