# Comparative Study of Regularization Techniques in Logistic Regression

Efua Ainooson Noonoo

March 2, 2025

# Introduction

The primary objective of this project is to analyze a leukemia gene expression dataset using various penalized regression methods. In high-dimensional settings—where the number of genes far exceeds the number of samples—traditional logistic regression often suffers from overfitting and multicollinearity. Penalized regression techniques, such as Lasso, Ridge, and Elastic Net, mitigate these challenges by incorporating regularization that shrinks coefficients, effectively selecting the most informative genes while improving model interpretability and predictive accuracy. As demonstrated by Greenwood et al. (2020) [Greenwood et al., 2020], penalized regression is a highly effective approach in such contexts. By applying these methods, this study aims to identify key genes associated with leukemia and to develop robust logistic regression models for classifying leukemia subtypes.

# Data Description

The leukemia dataset, originally published by [Golub et al., 1999], was used in this study. Referred to as `Golub_Merge`, it consists of 7129 genes (genes) measured across 72 samples. The dataset is structured as an `ExpressionSet` object in R, containing a gene expression matrix (`exprs`), where rows represent genes and columns correspond to patient samples. The accompanying metadata includes leukemia subtype classifications (`ALL` vs. `AML`), determined based on clinical diagnosis and molecular profiling.

While the dataset is well-suited for demonstrating the power of penalized regression techniques, it also has inherent limitations. With only 72 samples, the dataset poses challenges for generalizability and increases the risk of overfitting, even when using regularization methods. Additionally, the distribution of leukemia subtypes is uneven, with 47 samples of ALL and 25 of AML, which could introduce bias in classification results.

# Research Questions.

- How do LASSO, Ridge, and Elastic Net compare in terms of predictive performance and feature selection stability for classifying leukemia subtypes?

- How does the choice of regularization method influence the trade-off between model complexity and multicollinearity management in leukemia gene expression analysis?

# Modeling and Methods

## Modeling

### Logistic Regression Model

The logistic regression model was fitted to assess baseline performance. This model was selected as a benchmark due to its simplicity and interpretability, making it suitable for evaluating how well the genes in the dataset predict the binary response variable—leukemia subtype (ALL vs. AML). Logistic regression assumes a linear relationship between the genes and the log-odds of the response, providing an interpretable foundation for comparison against more complex models.

## Penalized Logistic Regression Models

Penalized logistic regression models, specifically LASSO, Ridge, and Elastic Net, were employed to address the challenges of high-dimensional data, where the number of genes far exceeds the number of samples. These models incorporate regularization techniques that impose penalties on the magnitude of the coefficients, thereby achieving three key objectives. LASSO, with its $L_1$-norm penalty, performs variable selection by shrinking some coefficients to exactly zero, thereby identifying the most predictive genes for leukemia subtypes. Ridge regression, using an $L_2$-norm penalty, minimizes the impact of collinearity among genes by distributing the influence across correlated variables. Elastic Net combines the strengths of LASSO and Ridge by blending $L_1$ and $L_2$ penalties. This approach is particularly useful when genes are highly correlated, as it retains feature selection capabilities while also managing multicollinearity [Friedman et al., 2010, Friedman et al., 2010].

## Cross-Validation and Hyperparameter Tuning

To ensure the reliability and generalizability of the models, 10-fold cross-validation was performed to optimize the regularization parameter ($\lambda$). Two key $\lambda$ values were identified: $\lambda_{\min}$, which minimizes the cross-validation error and favors predictive accuracy, and $\lambda_{1se}$, which is the largest value of $\lambda$ within one standard error of the minimum error and emphasizes model parsimony and interpretability. Additionally, the Elastic Net penalty was fine-tuned by optimizing the $\alpha$ parameter through cross-validation, balancing the sparsity of LASSO ($\alpha = 1$) with the grouping effect of Ridge ($\alpha = 0$).

## Feature Selection and Stability

Feature selection was conducted independently by each penalized regression method (LASSO, Ridge, and Elastic Net), with non-zero coefficients considered as selected genes. To ensure robustness, a secondary analysis using 10-fold cross-validation was performed to evaluate the frequency of feature selection across folds. Only those genes that were consistently selected (i.e., appearing in a majority of the folds) were considered stable genes. This two-step process ensures that the final set of genes is both reliable and informative.

## Prediction and Performance Evaluation

The final models were applied to classify leukemia subtypes, predicting whether samples belonged to ALL or AML. Model performance was evaluated using key metrics:

- **Sensitivity:** The proportion of actual ALL cases that were correctly identified.

- **Specificity:** The proportion of actual AML cases that were correctly identified.

- **Accuracy:** The overall correctness of the model's predictions.

- **Precision:** The proportion of predicted ALL cases that were truly ALL.

# Model Checking

Model fit and complexity were assessed using the Akaike Information Criterion (AIC) and Bayesian Information Criterion (BIC) on the test dataset. LASSO achieved the lowest AIC (54.3740) and BIC (62.0221), indicating its strong generalization performance and preference for sparse models. Elastic Net followed with AIC and BIC values of 99.5926 and 145.4812, respectively, demonstrating a balance between sparsity and feature retention. Ridge exhibited substantially higher AIC (14303.32) and BIC (27934.13) due to its lack of feature selection and high degrees of freedom, retaining all genes.

   Also, prediction metrics were also analysed. Overall, these diagnostics, together with key prediction metrics, support the selection of LASSO as the most parsimonious and predictive model.

# Results

Table 1 summarizes the optimal regularization parameters selected via 10-fold cross-validation for each penalized regression method. Figures 1a, 1c, and 1b display the corresponding cross-validation curves for Lasso, Ridge, and Elastic Net, respectively. Notably, Lasso regression achieved the lowest $\lambda_{\min}$ value of 0.004995635, reflecting its strong ability to shrink less informative coefficients to zero. In contrast, Ridge regression produced a much larger $\lambda_{\min}$ of 3.779015, consistent with its tendency to retain all genes while addressing multicollinearity. Additionally, the Elastic Net model, with an optimized $\alpha$ of 0.25, provided a balance between Lasso's sparsity and Ridge's grouping effects.

| Penalized Technique | $\lambda_{\min}$ | $\lambda_{1se}$ | $\alpha$ |
|---|---|---|---|
| Lasso | 0.004995635 | 0.089354407 | 1 |
| Ridge | 3.779015 | 13.900652 | 0 |
| Elastic Net | 0.01511606 | 0.20452751 | 0.25 |

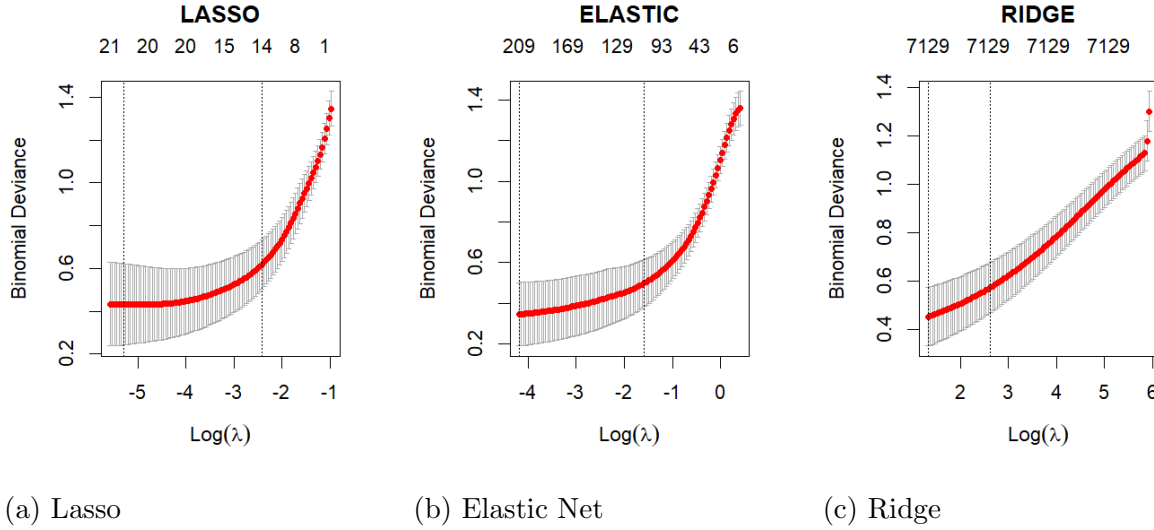Table 1: Optimal $\lambda$ values selected by Lasso, Ridge, and Elastic Net via 10-fold cross-validation.

(a) Lasso                     (b) Elastic Net                     (c) Ridge

Figure 1: Cross-validation curves for the penalized regression methods. The vertical lines indicate $\lambda_{\min}$ and $\lambda_{\mathrm{1se}}$ for each method.

## Feature Selection and Stability Analysis

The pure Ridge regression rarely drives any coefficient to exactly zero, it does not explicitly 'select' variables the way Lasso does. To ensure a fair comparison, we impose a small numeric threshold—here $\beta_j > 1 \times 10^{-3}$ and treat any coefficient with an absolute value below that as effectively zero. This way, all three methods (Lasso, Ridge, Elastic Net) are evaluated using the same cutoff criterion, making their selected feature sets more directly comparable.

The Lasso, Ridge, and Elastic Net methods were employed to identify genes relevant for classifying leukemia subtypes (ALL vs. AML) without incorporating cross-validation. Lasso identified 19 genes, optimizing for sparsity by shrinking coefficients of less important genes to zero. Ridge retained all 7129 genes, effectively addressing multicollinearity using its $L_2$-norm penalty. Elastic Net, with an optimized alpha value of 0.25, achieved a balance between Lasso's sparsity and Ridge's inclusiveness by selecting 210 genes. Table 2 summarizes the number of genes selected and the deviance values for each method.

| Method | genes Selected | Deviance |
|---|---|---|
| Lasso | 19 | 98.69 |
| Ridge | 1179 | 92.80 |
| Elastic Net | 193 | 98.61 |

Table 2: genes selected and deviance values for Lasso, Ridge, and Elastic Net.

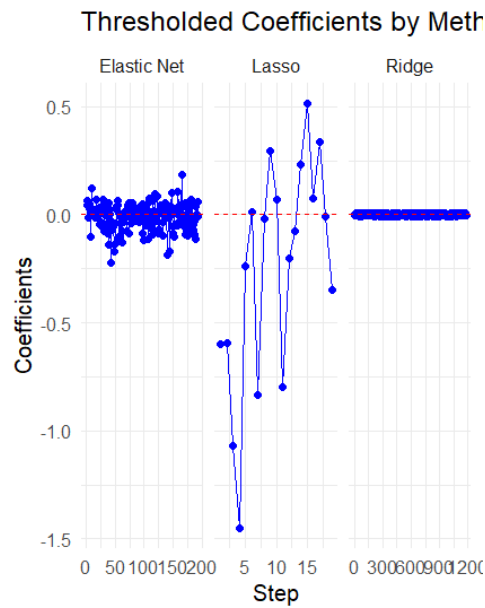Number of genes selected are displayed in Figures 2.

Figure 2: Plot showing genes selected by the penalised regression methods

Feature overlaps among the three methods were examined using a threshold of $1 \times 10^{-3}$ for Ridge. Lasso and Ridge shared 19 common genes.
Lasso and Elastic Net shared 19 genes, while Elastic Net and Ridge shared 193 genes.
All three methods identified the same 19 core genes; L20861_at, M15990_at, M19483_at, M23161_at, M27878_at, M31606_at, M31951_at, M84349_at, X51466_at, X59405_at, X95715_at, Y07596_at, Y12556_at, M60828_at, X06182_s_at, U40002_s_at, M92843_s_at, Z31690_s_at, and X83490_s_at underscoring their importance in classifying ALL and AML. These overlaps are visually represented in Figure 3
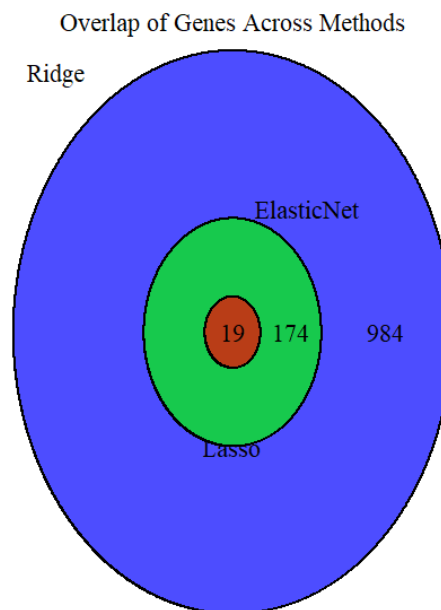


Figure 3: Venn diagram showing overlap and unique genes selected by Lasso, Ridge, and Elastic Net across folds.

To evaluate the stability and robustness of feature selection, 10-fold cross-validation was performed. The Lasso, Ridge, and Elastic Net methods were employed to identify stable and robust genes relevant for classifying leukemia subtypes (ALL vs. AML). genes that appeared in at least 80% of the cross-validation folds (8 out of 10 folds) were considered stable and selected for the final models. Lasso identified 10 genes. Ridge identified all 806 genes. Elastic Net identified 133 genes. Table 3 summarizes the number of genes selected and the deviance values for each method. Feature overlaps among the methods were examined using Venn diagrams. All three methods consistently identified a core set of 10 genes; `HG2562-HT2658_s_at M23197_at X95735_at M19507_at Y07604_at M16038_at M31994_at M63838_s_at X51521_at`. Figure 4 shows the overlap among the stable genes (selected in $\geq 8$ folds).

| Method | Stable genes Selected |
|---|---|
| Lasso | 10 |
| Ridge | 806 |
| Elastic Net | 133 |

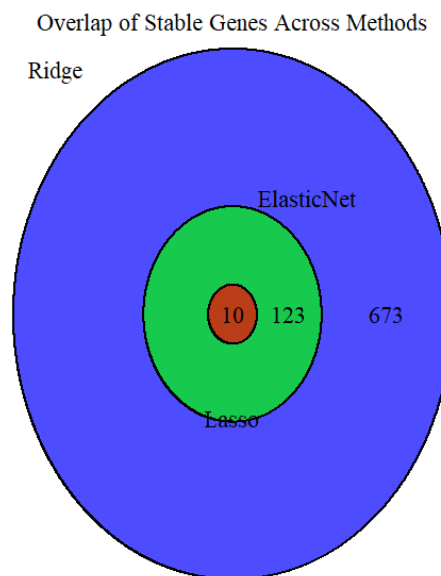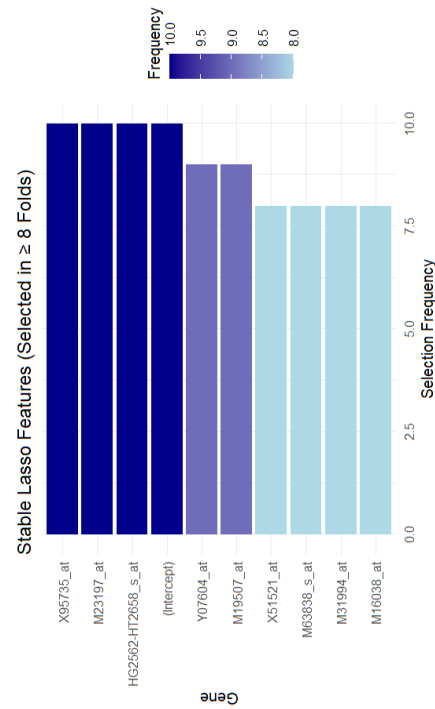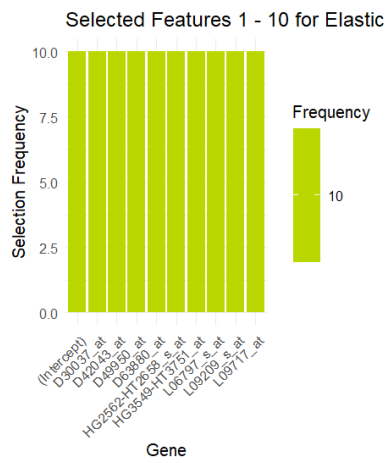Table 3: Stable genes selected and deviance values for Lasso, Ridge, and Elastic Net.



Figure 4: Venn diagram showing stable genes selected by Lasso, Ridge, and Elastic Net across folds.

Figure 5 shows the stable genes selected by Lasso. The stable Elastic Net genes are displayed in a grid format in Figure **??** shows the first 80 stable genes selected by Elastic. For clarity, the genes have been grouped into subsets. Figure 11b presents the first 80 stable Ridge genes in a grid, grouped for clarity.

Figure 5: Stable Lasso genes (selected in ≥8 folds).



(a) Group 1



(b) Group 2



(c) Group 3

(a) Group 4                    (b) Group 5                    (c) Group 6



(a) Group 7                                        (b) Group 8

Figure 8: Stable Elastic Net genes (groups of 10).

(a) Group 1


(b) Group 2


(c) Group 3


(a) Group 4


(b) Group 5


(c) Group 6


(a) Group 7
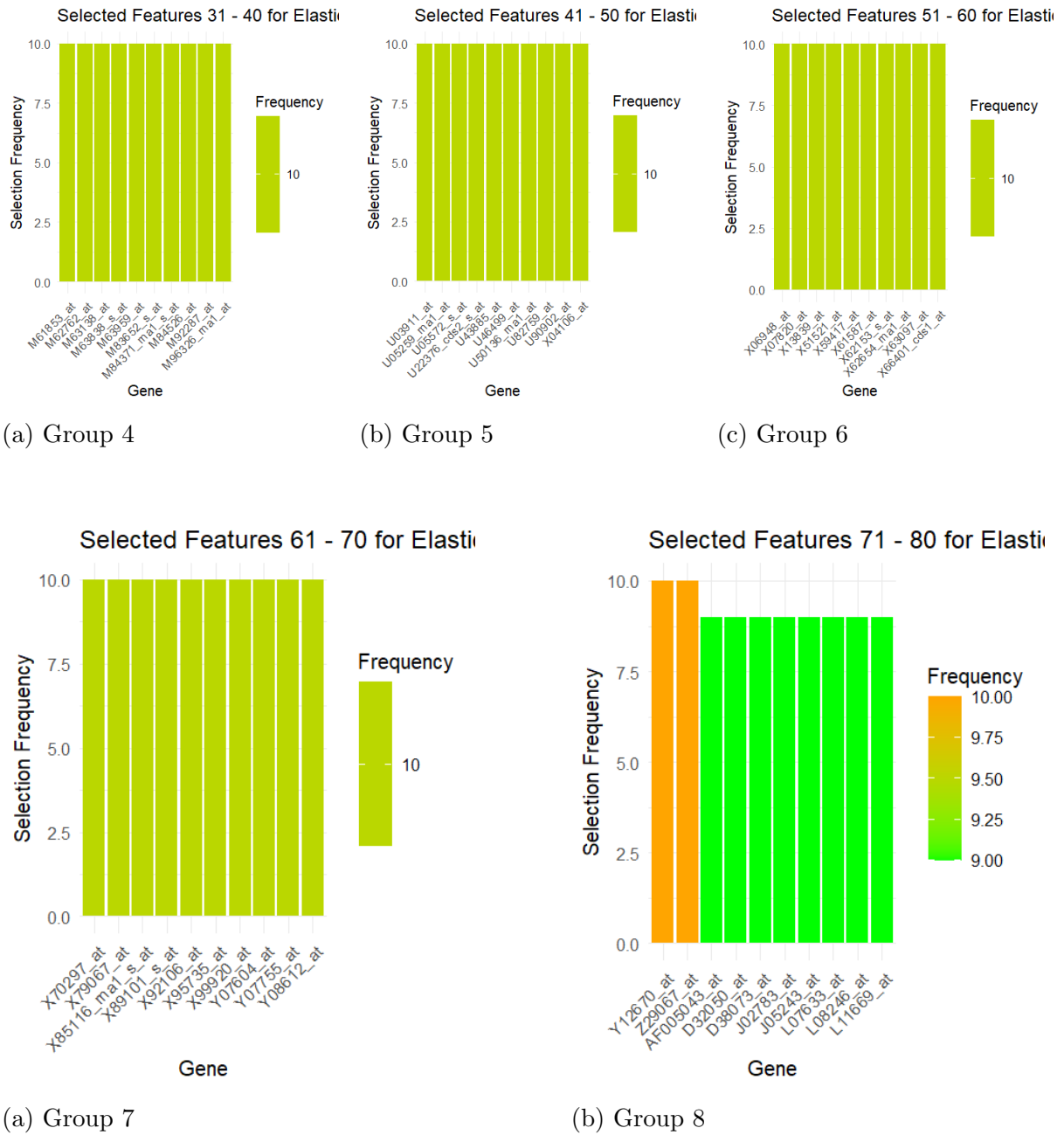

(b) Group 8

Figure 11: Stable Ridge genes (groups of 10).

## Prediction and Performance Metrics

The models were first evaluated without cross-validation to assess their ability to classify leukemia subtypes (ALL vs. AML) using key performance metrics, including accuracy, precision, sensitivity, and specificity. Table 4 summarizes the results. The Lasso and Elastic Net models achieved a perfect classification performance on the test data, with an accuracy, precision, sensitivity, and specificity of 100%. These models correctly classified all samples with no false positives (AML cases incorrectly predicted as ALL) or false negatives (ALL cases incorrectly predicted as AML). This performance highlights their capability to balance precision and sensitivity effectively, ensuring robust predictions for both ALL and AML subtypes.

The Ridge model, while highly effective in identifying ALL cases with a sensitivity of 100%, showed some limitations in specificity. It achieved an overall accuracy of 91%, with a precision of 88% and specificity of 71%. This reduced specificity is attributed to two false positives, where AML samples were incorrectly classified as ALL. Despite these false positives, Ridge maintained a strong ability to correctly classify ALL cases, achieving perfect sensitivity. Figure 13 show the methods and their performance.

| Method | Counts | | | | Metrics | | | |
|---|---|---|---|---|---|---|---|---|
| | TN | FP | FN | TP | Accuracy (%) | Precision (%) | Sensitivity (%) | Specificity (%) |
| **Elastic Net** | 7 | 0 | 0 | 15 | 100 | 100 | 100 | 100 |
| **Ridge** | 5 | 2 | 0 | 15 | 90.9 | 88.2 | 100 | 71.4 |
| **Lasso** | 7 | 0 | 0 | 15 | 100 | 100 | 100 | 100 |

Table 4: Comparison of Elastic Net, Ridge, and Lasso Prediction Results.



Figure 12: Performance metrics for Lasso, Ridge, and Elastic Net models for stable genes.

To ensure robustness, the final models were evaluated using the stable genes obtained through cross-validation. The results confirmed the initial findings, both Lasso and Elastic Net maintaining perfect performance (100% accuracy, precision, sensitivity, and specificity). Ridge regression, also maintained 100% performance. These findings underscore the ability of Lasso and Elastic Net to provide robust and reliable predictions using the most informative and stable genes, as shown in Table 5 and Figure 12.

| Method | Counts | | | | Metrics | | | |
|---|---|---|---|---|---|---|---|---|
| | TN | FP | FN | TP | Accuracy % | Precision % | Sensitivity % | Specificity % |
| **LASSO & Elastic & Ridge** | 7 | 0 | 0 | 15 | 100 | 100 | 100 | 100 |

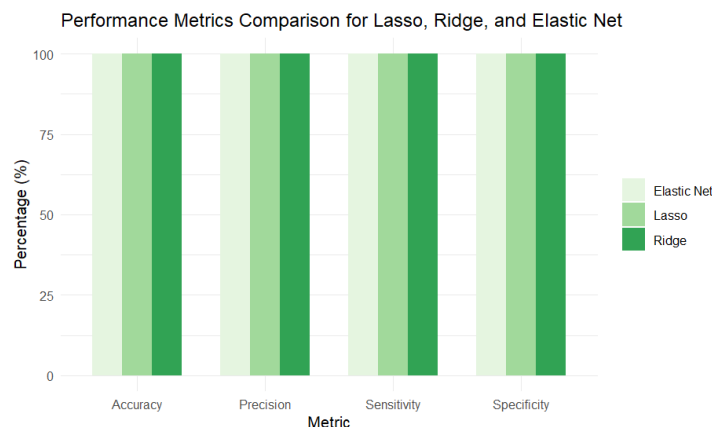Table 5: Comparison of LASSO & Elastic Net and Ridge Prediction Results.



Figure 13: Performance metrics for Lasso, Ridge, and Elastic Net models.

Lasso demonstrated exceptional performance in classifying leukemia subtypes, achieving perfect accuracy while selecting the sparsest set of genes, making it highly interpretable. Elastic Net also achieved perfect classification but retained more genes, balancing sparsity and handling correlated genes effectively. Ridge, while managing multicollinearity, struggled with specificity and retained all genes, limiting its interpretability and practical applicability. The stability analysis confirmed Lasso as the most robust model, consistently selecting a smaller, stable feature set. Lasso is the recommended approach, with Elastic Net as a viable alternative when additional feature retention is desired.

## Discussion of Results

Lasso and Elastic Net models achieved perfect classification performance on the test data, demonstrating both high predictive accuracy and effective feature selection. Ridge, while achieving 91% accuracy, struggled with specificity due to a few false positives. The stability analysis confirmed that Lasso consistently selected a small, robust set of genes, whereas Ridge retained a larger number of genes. The overlaps in selected genes (illustrated by the Venn diagrams) reinforce the importance of a core set of genes in distinguishing between ALL and AML.

Overall, the results provide a comprehensive comparative analysis of penalized regression methods, with robust feature selection and performance evaluation. Further enhancements in model diagnostics and biological interpretation could elevate the project even more.

## Potential Pitfalls

The perfect classification performance reported for Lasso and Elastic Net (100% accuracy, sensitivity, and specificity) is a critical concern. Achieving such flawless results in biomed-

ical data is highly unusual and often indicative of potential issues such as data leakage, overfitting, or improper train-test split methodology. Further scrutiny is necessary to ensure the validity of these findings.

Additionally, the test set comprises only 22 samples (15 ALL and 7 AML), which is both small and imbalanced. This combination raises concerns about the robustness and representativeness of the reported performance metrics. A small test set limits the ability to generalize these results to larger populations reliably. Leukemia datasets are typically small, exacerbating the risk of overfitting and undermining the stability of feature selection. Additional evaluation with larger or independent datasets is recommended to confirm the robustness of these results.

Another challenge lies in the disparity in the number of stable genes selected by the models (Lasso: 10, Elastic Net: 133, Ridge: 806). While Lasso's aggressive sparsity enhances interpretability, it risks excluding potentially important genes. Conversely, Ridge retains an excessively large number of genes, reducing its practicality and interpretability. The threshold for defining stable genes (genes appearing in at least 80% of cross-validation folds) and $1 \times 10^{-3}$ may also require further validation to confirm its suitability and relevance for this analysis.

# Appendix

```
1  library(glmnet)
2  library(tidyr)
3  library(dplyr)
4  library(ggplot2)
5  library(caret)
6  library(kableExtra)
7  library(RColorBrewer)
```

Listing 1: R libraries

```
1  data(Golub_Merge)
2
3  expression_data = exprs(Golub_Merge)
4
5  sample_labels = pData(Golub_Merge)$ALL.AML
6
7  summary(sample_labels)
8
9  dataset_summary = as.data.frame(table(sample_labels))
10  colnames(dataset_summary) = c("Leukemia_Type", "Count")
11
12
13  table_visual <- dataset_summary %>%
14    kable("html", col.names = c("Leukemia Type", "Sample Count"),
          align = "c") %>%
15    kable_styling(bootstrap_options = c("striped", "hover", "
          condensed", "responsive")) %>%
16    row_spec(0, bold = TRUE, background = "#f2f2f2") %>%
17    add_header_above(c("Dataset Breakdown" = 2))
```

```r
18
19  table_visual
20
21
22  # A bar chart for dataset breakdown
23
24  ggplot(dataset_summary, aes(x = Leukemia_Type, y = Count, fill =
       Leukemia_Type)) +
25    geom_bar(stat = "identity", width = 0.5) +
26    theme_minimal() +
27    scale_fill_manual(values = c("ALL" = "purple", "AML" = "
         turquoise")) +
28    labs(
29      title = "Dataset Breakdown by Leukemia Type",
30      x = "Leukemia Type",
31      y = "Sample Count"
32    ) +
33    theme(
34      plot.title = element_text(hjust = 0.5, face = "bold"),
35      axis.text.x = element_text(face = "bold"),
36      axis.text.y = element_text(face = "bold")
37    )
38
39
40    expression_df= as.data.frame(t(expression_data[1:5, ]))
41  colnames(expression_df) = paste0("Gene", 1:5)
42  expression_df$LeukemiaType = sample_labels
43
44
45  long_data = pivot_longer(expression_df, cols = starts_with("Gene"
       ),
46                           names_to = "Gene", values_to = "
                              Expression")
47
48  # Boxplots with colors for leukemia types
49  ggplot(long_data, aes(x = Gene, y = Expression, fill =
       LeukemiaType)) +
50    geom_boxplot() +
51    scale_fill_manual(values = c("ALL" = "purple", "AML" = "
         turquoise")) +
52    theme_minimal() +
53    labs(title = "Gene Expression by Leukemia Type",
54         x = "Gene",
55         y = "Expression Level",
56         fill = "Leukemia Type") +
57    theme(axis.text.x = element_text(angle = 45, hjust = 1))
58
59
60  scaled_expression_data = scale(expression_data)
61
62
```

```r
63  ## Encoding the response variable as 0 for AML and 1 for ALL
64
65  response = numeric(length(sample_labels))
66
67  for (i in seq_along(sample_labels)){
68    if (sample_labels[i] == "ALL"){
69
70    response[i] = 1
71    }
72    else {
73    response[i] = 0
74    }
75  }
76
77  ## : Split the dataset(sample) into training (70%) and testing
        (30%) sets for model.
78
79  n_samples = ncol(Golub_Merge)
80
81  set.seed(4321)
82
83  train_indices = sample(1:n_samples, size = 0.7*n_samples)
84
85
86  training_data = Golub_Merge[,train_indices]
87
88  testing_data = Golub_Merge[,-train_indices]
89
90  response_train = response[train_indices]
91
92  response_test = response[-train_indices]
93
94  expression_train = as.data.frame(t(exprs(training_data)))
95
96  expression_train$response = response_train
97
98  X = subset(expression_train, select = -response)
99
100 X_mat = as.matrix(X)
101
102 X_mat = scale(X_mat)
103
104 expression_test = as.data.frame(t(exprs(testing_data)))
```

Listing 2: Data Processing and Exploration

```r
1
2  ## Linear Regression
3
4  model_lm =
5  lm(response~., data = expression_train)
6
```

```r
7   summary(model_lm)
8
9   ## Generalised Linear Regression
10
11  model_glm =
12  glm(response~., data = expression_train, family=binomial(link=
        logit))
13
14  summary(model_glm)
15
16
17  ## Deviance goodness-of-fit test
18
19  residual_deviance = model_glm$deviance
20
21  df = model_glm$df.residual
22
23  p_value = 1 - pchisq(residual_deviance, df)
24
25  residual_deviance
26
27  p_value
28
29
30  deviance_residuals = residuals(model_glm, type = "deviance")
31  fitted_values = fitted(model_glm)
32
33  plot(fitted_values,deviance_residuals,
34      xlab = "Fitted Values",
35      ylab = "Deviance Residuals",
36      main = "Deviance Residuals vs Fitted Values",
37      pch = 19,
38      col = "turquoise")
39
40  # Test on a simulated binary classification dataset
41  set.seed(42)
42  X_test = X_mat
43  y_test = expression_train$response
44
45  # Fit Lasso, Ridge, and Elastic Net
46  fit_lasso <- glmnet(X_test, y_test, family = "binomial", alpha =
        1)
47  fit_ridge <- glmnet(X_test, y_test, family = "binomial", alpha =
        0)
48  fit_elastic <- glmnet(X_test, y_test, family = "binomial", alpha
        = 0.25)
49
50  # Calculate AIC/BIC
51  aic_bic_lasso = calculate_aic_bic(fit_lasso, fit_lasso$lambda
        [10], X_test, y_test)
52  aic_bic_ridge = calculate_aic_bic(fit_ridge, fit_ridge$lambda
```

```
53      [10], X_test, y_test)
   aic_bic_elastic = calculate_aic_bic(fit_elastic, fit_elastic$
       lambda[10], X_test, y_test)
54
55 cat("Test Lasso AIC:", aic_bic_lasso$AIC, "BIC:", aic_bic_lasso$
       BIC, "\n")
56 cat("Test Ridge AIC:", aic_bic_ridge$AIC, "BIC:", aic_bic_ridge$
       BIC, "\n")
57 cat("Test Elastic Net AIC:", aic_bic_elastic$AIC, "BIC:", aic_bic
       _elastic$BIC, "\n")
```

Listing 3: Model Checking

```
1
2  ## mean deviance for lasso regression
3
4  set.seed(4321)
5
6  fit_cv_lasso = cv.glmnet(
7    x = X_mat,
8    y = expression_train$response,
9   family = "binomial",
10   alpha = 1,
11   nfolds = 10,
12   #type.measure = "class"
13   )
14
15 plot(fit_cv_lasso)
16 title("LASSO", line = 2.5)
17
18 list(c( "lambda.min.lasso" = fit_cv_lasso$lambda.min,
19
20 "lambda.1se.lasso" = fit_cv_lasso$lambda.1se))
21
22 ## mean deviance for ridge regression
23
24 set.seed(4321)
25
26 fit_cv_ridge = cv.glmnet(
27   x = X_mat,
28   y = expression_train$response,
29  family = "binomial",   # logistic regression
30   alpha = 0,             # lasso penalty
31   nfolds = 10,           # 10-fold CV
32   #type.measure = "class" # classification error for binomial
33   )
34
35 plot(fit_cv_ridge)
36 title("RIDGE", line = 2.5)
37
38 list(c( "lambda.min.ridge" = fit_cv_ridge$lambda.min,
39
```

17

```r
40  "lambda.1se.ridge" = fit_cv_ridge$lambda.1se))
41
42  set.seed(4321)
43
44  alpha_values = c(0, 0.25, 0.5, 0.75, 1)
45
46  results = data.frame(alpha = numeric(),
47                       lambda.min = numeric(),
48                       lambda.1se = numeric(),
49                       cv.error.min = numeric(),
50                       cv.error.1se = numeric(),
51                       stringsAsFactors = FALSE)
52
53  for (a in alpha_values) {
54
55    # Perform cross-validation for given alpha
56
57    cv_fit = cv.glmnet(X_mat, expression_train$response, family = "
          binomial", alpha = a, type.measure = "class")
58
59   results = rbind(results,
60                     data.frame(
61                       alpha = a,
62                       lambda.min = cv_fit$lambda.min,
63                       lambda.1se = cv_fit$lambda.1se,
64                       cv.error.min = min(cv_fit$cvm),
65                       cv.error.1se = cv_fit$cvm[which(cv_fit$
                         lambda == cv_fit$lambda.1se)]
66                     ))
67  }
68
69  print(results)
70
71  best <- results[which.min(results$cv.error.min), ]
72  cat("Best alpha:", best$alpha,
73      "with lambda.min:", best$lambda.min,
74      "achieved a cross-validation error of:", best$cv.error.min, "
          \n")
75
76  ## mean deviance for elastic regression
77
78  set.seed(4321)
79
80  fit_cv_elastic = cv.glmnet(
81    x = X_mat,
82    y = expression_train$response,
83   family = "binomial",
84    alpha = 0.25,
85    nfolds = 10,
86    #type.measure = "class" # classification error for binomial
87    )
```

```r
88
89  plot(fit_cv_elastic)
90  title("ELASTIC", line = 2.5)
91
92  list(c( "lambda.min.elastic" = fit_cv_elastic$lambda.min,
93
94  "lambda.1se.elastic" = fit_cv_elastic$lambda.1se))
```
Listing 4: Mean Deviance

```r
1
2
3  rlas = glmnet(X_mat, expression_train$response,  family = "
       binomial", alpha = 1,  lambda = fit_cv_lasso$lambda.min)
4
5  rrid = glmnet(X_mat, expression_train$response,  family = "
       binomial", alpha = 0,  lambda = fit_cv_ridge$lambda.min)
6
7  renet = glmnet(X_mat, expression_train$response, family = "
       binomial", alpha = .25, lambda = fit_cv_elastic$lambda.min)
```
Listing 5: Fitting the methods

```r
1
2
3  ## non-zero coefficients for lasso
4
5  coef_matrix_lasso = coef(rlas)
6
7  coef_values_lasso = coef_matrix_lasso[-1,,drop = FALSE]
8
9  nonzero_idx_lasso = which(coef_values_lasso != 0)
10
11 nonzero_coefs_lasso = coef_values_lasso[nonzero_idx_lasso, , drop
       = FALSE]
12
13
14
15 ## non-zero coefficients for ridge
16
17
18 coef_matrix_ridge = coef(rrid)
19
20 coef_values_ridge = coef_matrix_ridge[-1,,drop = FALSE]
21
22 nonzero_idx_ridge = which(coef_values_ridge != 0)
23
24 nonzero_coefs_ridge = coef_values_ridge[nonzero_idx_ridge, , drop
       = FALSE]
25
26
27 ## non-zero coefficients for elastic
```

```r
28
29  coef_matrix_elastic = coef(renet)
30
31  coef_values_elastic = coef_matrix_elastic[-1,,drop = FALSE]
32
33  nonzero_idx_elastic = which(coef_values_elastic != 0)
34
35  nonzero_coefs_elastic = coef_values_elastic[nonzero_idx_elastic,
       , drop = FALSE]
36
37  # Remove intercept row for each method
38  coef_values_lasso_noIntercept = coef_values_lasso[-1,, drop=FALSE
       ]
39  coef_values_ridge_noIntercept = coef_values_ridge[-1,, drop=FALSE
       ]
40  coef_values_elastic_noIntercept = coef_values_elastic[-1,, drop=
       FALSE]
41
42  # Threshold
43
44
45  # Lasso
46  threshold_idx_lasso = which(abs(coef_values_lasso_noIntercept) >
       threshold)
47  threshold_coefs_lasso = coef_values_lasso_noIntercept[threshold_
       idx_lasso, , drop=FALSE]
48  selected_genes_lasso = rownames(coef_values_lasso_noIntercept)[
       threshold_idx_lasso]
49
50  # Ridge
51  threshold_idx_ridge = which(abs(coef_values_ridge_noIntercept) >
       threshold)
52  threshold_coefs_ridge = coef_values_ridge_noIntercept[threshold_
       idx_ridge, , drop=FALSE]
53  selected_genes_ridge = rownames(coef_values_ridge_noIntercept)[
       threshold_idx_ridge]
54
55  # Elastic Net
56  threshold_idx_elastic = which(abs(coef_values_elastic_noIntercept
       ) > threshold)
57  threshold_coefs_elastic = coef_values_elastic_noIntercept[
       threshold_idx_elastic, , drop=FALSE]
58  selected_genes_elastic = rownames(coef_values_elastic_noIntercept
       )[threshold_idx_elastic]
59
60  # Compare coeff selected
61  list(
62    lasso   = length(threshold_coefs_lasso),
63    ridge   = length(threshold_coefs_ridge),
64    elastic = length(threshold_coefs_elastic)
65  )
```

```r
66
67
68  df_lasso <- data.frame(
69    Step        = seq_len(nrow(threshold_coefs_lasso)),
70    Coefficient = threshold_coefs_lasso[, 1],
71    Method      = "Lasso"
72  )
73
74  df_elastic <- data.frame(
75    Step        = seq_len(nrow(threshold_coefs_elastic)),
76    Coefficient = threshold_coefs_elastic[, 1],
77    Method      = "Elastic Net"
78  )
79
80  df_ridge <- data.frame(
81    Step        = seq_len(nrow(threshold_coefs_ridge)),
82    Coefficient = threshold_coefs_ridge[, 1],
83    Method      = "Ridge"
84  )
85
86
87  df_all <- rbind(df_lasso, df_elastic, df_ridge)
```

Listing 6: Selection of non-zero coefficients

```r
1
2
3   y_min = min(df_all$Coefficient, na.rm = TRUE)
4   y_max = max(df_all$Coefficient, na.rm = TRUE)
5
6   ggplot(df_all, aes(x = Step, y = Coefficient)) +
7
8     geom_line(color = "blue") +
9     geom_point(color = "blue") +
10
11    geom_hline(yintercept = 0, linetype = "dashed", color = "red")
        +
12
13    facet_wrap(~ Method, scales = "free_x") +
14
15    coord_cartesian(ylim = c(y_min, y_max)) +
16    labs(
17      title = "Thresholded Coefficients by Method",
18      x = "Step",
19      y = "Coefficients"
20    ) +
21    theme_minimal()
```

Listing 7: Plot Threshold Coefficients by Method

```r
1
2
```

```r
3  variable_names = colnames(X_mat)
4
5  # Common coefficients between LASSO and Ridge
6  common_coefficients1 = Reduce(intersect, list(threshold_idx_lasso
       , threshold_idx_ridge))
7  common_vars_lasso_ridge = variable_names[common_coefficients1]
8  cat("Common variables between LASSO and Ridge:\n", common_vars_
       lasso_ridge, "\n")
9
10 # Common coefficients between LASSO and Elastic Net
11 common_coefficients2 = Reduce(intersect, list(threshold_idx_lasso
       , threshold_idx_elastic))
12 common_vars_lasso_elastic = variable_names[common_coefficients2]
13 cat("Common variables between LASSO and Elastic Net:\n", common_
       vars_lasso_elastic, "\n")
14
15 # Common coefficients between Elastic Net and Ridge
16 common_coefficients3 = Reduce(intersect, list(threshold_idx_
       elastic, threshold_idx_ridge))
17 common_vars_elastic_ridge = variable_names[common_coefficients3]
18 cat("Common variables between Elastic Net and Ridge:\n", common_
       vars_elastic_ridge, "\n")
19
20 # Common coefficients across all three methods
21 common_coefficients4 = Reduce(intersect, list(threshold_idx_lasso
       , threshold_idx_ridge, threshold_idx_elastic))
22 common_vars_all = variable_names[common_coefficients4]
23 cat("Common variables across all three methods:\n", common_vars_
       all, "\n")
24
25 ### venn diagram of associated genes
26
27 venn.plot =  venn.diagram(
28   x = list(
29     Lasso = threshold_idx_lasso,
30     Ridge = threshold_idx_ridge,
31     ElasticNet = threshold_idx_elastic
32   ),
33   filename = NULL,
34   fill = c("red", "blue", "green"),
35   alpha = 0.7,
36   main = "Overlap of Genes Across Methods"
37 )
38 grid.newpage()
39 grid.draw(venn.plot)
```

Listing 8: Associated Genes

```r
1
2
3  # Extract genes from test data and convert to matrix
4
```

```r
5   X_test = as.matrix(expression_test)
6
7   ## Predict probabilities on test data; (threshold = 0.5)
8
9   pred_probs_lasso = predict(lasso_model, newx = X_test, type = "
        response")
10
11  pred_class_lasso = ifelse(pred_probs_lasso > 0.5, 1, 0)
12
13  ## Model performance on test data
14
15  confusion_matrix_lasso = table(Predicted = pred_class_lasso,
        Actual = response_test)
16  print(confusion_matrix_lasso)
17
18
19  X_test = as.matrix(expression_test)
20
21  ## Predict probabilities on test data; (threshold = 0.5)
22  pred_probs_ridge = predict(ridge_model, newx = X_test, type = "
        response")
23
24
25  pred_class_ridge = ifelse(pred_probs_ridge > 0.5, 1, 0)
26
27  ## Model performance on test data
28
29  confusion_matrix_ridge = table(Predicted = pred_class_ridge,
        Actual = response_test)
30  print(confusion_matrix_ridge)
31
32  X_test = as.matrix(expression_test)
33
34  ## Predict probabilities on test data; (threshold = 0.5)
35  pred_probs_elastic = predict(elastic_model, newx = X_test, type =
         "response")
36
37  pred_class_elastic = ifelse(pred_probs_elastic > 0.5, 1, 0)
38
39  # Model performance on test data
40
41  confusion_matrix_elastic = table(Predicted = pred_class_elastic,
        Actual = response_test)
42  print(confusion_matrix_elastic)
```

Listing 9: Prediction and Classification

```r
1
2   ## Model Performance for Lasso
3
4   TN_lasso = confusion_matrix_lasso[1,1]
5
```

```r
6   FN_lasso = confusion_matrix_lasso[1,2]
7
8   FP_lasso = confusion_matrix_lasso[2,1]
9
10  TP_lasso = confusion_matrix_lasso[2,2]
11
12  all_samples_lasso = TN_lasso + FN_lasso + FP_lasso + TP_lasso
13
14
15  ### Accuracy
16
17  Accuracy_lasso = ((TN_lasso + TP_lasso)/ all_samples_lasso)*100
18
19
20
21
22  ### Sensitivity
23
24  Sensitivity_lasso = ((TP_lasso)/ (TP_lasso + FN_lasso))*100
25
26
27
28
29  ### Specificity
30
31  Specificity_lasso = ((TN_lasso)/ (TN_lasso + FP_lasso))*100
32
33
34
35
36  ### Precision
37
38  Precision_lasso = ((TP_lasso)/ (TP_lasso + FP_lasso))*100
39
40
41
42
43  list(c(Accuracy_lasso = Accuracy_lasso, Sensitivity_lasso =
        Sensitivity_lasso, Specificity_lasso = Specificity_lasso,
        Precision_lasso = Precision_lasso ))
44
45
46  ## Model performance for Ridge
47
48  TN_ridge = confusion_matrix_ridge[1,1]
49
50  FN_ridge = confusion_matrix_ridge[1,2]
51
52  FP_ridge = confusion_matrix_ridge[2,1]
53
54  TP_ridge = confusion_matrix_ridge[2,2]
```

```
55
56
57
58
59  ### Accuracy
60
61  Accuracy_ridge = ((TN_ridge + TP_ridge)/ (TN_ridge + TP_ridge +
        FN_ridge + FP_ridge) )*100
62
63
64
65
66  ### Sensitivity
67
68  Sensitivity_ridge = ((TP_ridge)/ (TP_ridge + FN_ridge))*100
69
70
71
72
73  ### Specificity
74
75  Specificity_ridge = ((TN_ridge)/ (TN_ridge + FP_ridge))*100
76
77
78
79
80  ### Precision
81
82  Precision_ridge = ((TP_ridge)/ (TP_ridge + FP_ridge))*100
83
84
85
86
87  list(c(Accuracy_ridge = Accuracy_ridge, Sensitivity_ridge =
        Sensitivity_ridge, Specificity_ridge = Specificity_ridge,
        Precision_ridge = Precision_ridge ))
88
89
90  ## Model Performance for Elastic Methods
91
92  TN_elastic = confusion_matrix_elastic[1,1]
93
94  FN_elastic = confusion_matrix_elastic[1,2]
95
96  FP_elastic = confusion_matrix_elastic[2,1]
97
98  TP_elastic = confusion_matrix_elastic[2,2]
99
100
101
102
```

```r
103  ### Accuracy
104
105  Accuracy_elastic = ((TN_elastic + TP_elastic)/ (TN_elastic + TP_
         elastic + FN_elastic + FP_elastic) )*100
106
107
108
109
110  ### Sensitivity
111
112  Sensitivity_elastic = ((TP_elastic)/ (TP_elastic + FN_elastic))*
         100
113
114
115
116
117  ### Specificity
118
119  Specificity_elastic = ((TN_elastic)/ (TN_elastic + FP_elastic))*
         100
120
121
122
123
124  ### Precision
125
126  Precision_elastic = ((TP_elastic)/ (TP_elastic + FP_elastic))*100
127
128
129  list(c(Accuracy_elastic = Accuracy_elastic, Sensitivity_elastic =
          Sensitivity_elastic, Specificity_elastic = Specificity_
         elastic, Precision_elastic = Precision_elastic ))
130
131  metrics_df = data.frame(
132    Method = rep(c("Lasso", "Ridge", "Elastic Net"), each = 4),
133    Metric = rep(c("Accuracy", "Sensitivity", "Specificity", "
            Precision"), times = 3),
134    Value  = c(Accuracy_lasso, Sensitivity_lasso, Specificity_lasso
            , Precision_lasso,
135               Accuracy_ridge, Sensitivity_ridge, Specificity_ridge
                    , Precision_ridge,
136               Accuracy_elastic, Sensitivity_elastic, Specificity_
                    elastic, Precision_elastic)
137  )
138
139  ggplot(metrics_df, aes(x = Metric, y = Value, fill = Method)) +
140    geom_bar(stat = "identity", position = "dodge", width = 0.7) +
141    ylim(0, 100) +
142    labs(title = "Performance Metrics Comparison for Lasso, Ridge,
            and Elastic Net",
143          x = "Metric",
```

```
144        y = "Percentage (%)") +
145    scale_fill_brewer(palette = "Greens") +
146    theme_minimal() +
147    theme(legend.title = element_blank())
```

Listing 10: Model Performance

```
1
2
3  y_min = min(df_all$Coefficient, na.rm = TRUE)
4  y_max = max(df_all$Coefficient, na.rm = TRUE)
5
6  ggplot(df_all, aes(x = Step, y = Coefficient)) +
7
8    geom_line(color = "blue") +
9    geom_point(color = "blue") +
10
11    geom_hline(yintercept = 0, linetype = "dashed", color = "red")
         +
12
13    facet_wrap(~ Method, scales = "free_x") +
14
15    coord_cartesian(ylim = c(y_min, y_max)) +
16    labs(
17      title = "Thresholded Coefficients by Method",
18      x = "Step",
19      y = "Coefficients"
20    ) +
21    theme_minimal()
```

Listing 11: Plot Threshold Coefficients by Method

```
1
2
3  selected_genes_lasso = list()
4  selected_genes_ridge = list()
5  selected_genes_elastic = list()
6
7  # Cross-validation
8  set.seed(4321)
9  k = 10
10 folds = createFolds(expression_train$response, k = k, list = TRUE
     )
11
12
13 fit_cv_lasso = cv.glmnet(X_mat, expression_train$response,
14                        family = "binomial",
15                        alpha = 1)
16
17 fit_cv_ridge = cv.glmnet(X_mat, expression_train$response,
18                        family = "binomial",
19                        alpha = 0)
```

```r
20
21  fit_cv_elastic = cv.glmnet(X_mat, expression_train$response,
22                             family = "binomial",
23                             alpha = 0.5)
24
25
26  for (i in seq_along(folds)) {
27      train_indices = unlist(folds[-i])
28      valid_indices = folds[[i]]
29
30      X_train = X_mat[train_indices, ]
31      y_train = expression_train$response[train_indices]
32
33      y_train = as.vector(y_train)
34
35      # Train Lasso model
36
37      fit_lasso = glmnet(X_train, y_train,
38                         family = "binomial",
39                         alpha = 1,
40                         lambda = fit_cv_lasso$lambda.min)
41
42      coef_lasso = coef(fit_lasso, s = fit_cv_lasso$lambda.min)
43
44      coef_lasso_matrix = as.matrix(coef_lasso)
45
46      selected_genes_lasso = rownames(coef_lasso_matrix)[abs(coef_
          lasso_matrix[, 1]) > threshold]
47
48      selected_genes_lasso[[i]] = selected_genes_lasso
49
50      # Train Ridge model
51
52      fit_ridge = glmnet(X_train, y_train,
53                         family = "binomial",
54                         alpha = 0,
55                         lambda = fit_cv_ridge$lambda.min)
56
57      coef_ridge = coef(fit_ridge, s = fit_cv_ridge$lambda.min)
58
59      coef_ridge_matrix = as.matrix(coef_ridge)
60
61      selected_genes_ridge = rownames(coef_ridge_matrix)[abs(coef_
          ridge_matrix[, 1]) > threshold]
62
63      selected_genes_ridge[[i]] = selected_genes_ridge
64
65      # Train Elastic Net model
66
67      fit_elastic = glmnet(X_train, y_train,
68                           family = "binomial",
```

```r
69                             alpha = 0.25,
70                             lambda = fit_cv_elastic$lambda.min)
71
72     coef_elastic = coef(fit_elastic, s = fit_cv_elastic$lambda.
           min)
73
74     coef_elastic_matrix = as.matrix(coef_elastic)
75
76     selected_genes_elastic = rownames(coef_elastic_matrix)[abs(
           coef_elastic_matrix[, 1]) > threshold]
77
78     selected_genes_elastic[[i]] = selected_genes_elastic
79  }
80
81  list(c("lasso genes" = fit_lasso$df, "lasso Deviance" = fit_
        elastic$dev.ratio,
82          "ridge genes" = fit_ridge$df,
83          "ridge Deviance" = fit_ridge$dev.ratio,
84          "elastic genes" = fit_elastic$df, "elastic Deviance" = fit
               _elastic$dev.ratio))
85
86
87  #  All selected genes from all folds
88
89  selected_genes_all_lasso = unlist(selected_genes_lasso)
90  selected_genes_all_elastic = unlist(selected_genes_elastic)
91  selected_genes_all_ridge = unlist(selected_genes_ridge)
92
93
94  selected_genes_summary_lasso = sort(table(selected_genes_all_
        lasso), decreasing = TRUE)
95  selected_genes_summary_elastic = sort(table(selected_genes_all_
        elastic), decreasing = TRUE)
96  selected_genes_summary_ridge = sort(table(selected_genes_all_
        ridge), decreasing = TRUE)
97
98
99  lasso_stability = as.data.frame(selected_genes_summary_lasso)
100 colnames(lasso_stability) = c("Gene", "Frequency")
101 lasso_stability$Method = "Lasso"
102
103 elastic_stability = as.data.frame(selected_genes_summary_elastic)
104 colnames(elastic_stability) = c("Gene", "Frequency")
105 elastic_stability$Method = "Elastic Net"
106
107 ridge_stability = as.data.frame(selected_genes_summary_ridge)
108 colnames(ridge_stability) = c("Gene", "Frequency")
109 ridge_stability$Method = "Ridge"
110
111
112 combined_stability = rbind(lasso_stability, elastic_stability,
```

```r
      ridge_stability)

head(combined_stability)

## genes that appear in >= 8 out of 10 folds.

#  Lasso:
stable_lasso = names(selected_genes_summary_lasso[
                      selected_genes_summary_lasso >= 8
                      ])

#  Ridge:
stable_ridge = names(selected_genes_summary_ridge[
                      selected_genes_summary_ridge >= 8
                      ])

#  Elastic Net:
stable_elastic = names(selected_genes_summary_elastic[
                      selected_genes_summary_elastic >= 8
                        ])


list(c( stable_lasso = length(stable_lasso), stable_ridge =
    length(stable_ridge), stable_elastic = length(stable_elastic))
    )

# Venn diagram with stable genes

venn.plot = venn.diagram(
    x = list(
        Lasso = stable_lasso,
        Ridge = stable_ridge,
        ElasticNet = stable_elastic
    ),
    filename = NULL,
    fill = c("red", "blue", "green"),
    alpha = 0.7,
    main = "Overlap of Stable Genes Across Methods"
)
grid.newpage()
grid.draw(venn.plot)

## Lasso Stable genes selected

lasso_stability_df = as.data.frame(selected_genes_summary_lasso)
colnames(lasso_stability_df) <- c("Gene", "Frequency")


stable_lasso_df = subset(lasso_stability_df, Frequency >= 8)
```

```r
161  stable_lasso_df = stable_lasso_df[order(stable_lasso_df$Frequency
        , decreasing = TRUE), ]
162
163
164  stable_lasso_df = stable_lasso_df[order(stable_lasso_df$Frequency
        , decreasing = TRUE), ]
165
166
167  ggplot(stable_lasso_df, aes(x = reorder(Gene, Frequency), y =
        Frequency, fill = Frequency)) +
168    geom_bar(stat = "identity") +
169    coord_flip() +
170    scale_fill_gradient(low = "lightblue", high = "darkblue") +
171    labs(title = "Stable Lasso genes (Selected in    8 Folds)",
172         x = "Gene",
173         y = "Selection Frequency") +
174    theme_minimal()
175
176  elastic_stability_df = as.data.frame(selected_genes_summary_
        elastic)
177
178  colnames(elastic_stability_df) = c("Gene", "Frequency")
179
180  stable_elastic_df = subset(elastic_stability_df, Frequency >= 8)
181
182  stable_elastic_df = stable_elastic_df[order(stable_elastic_df$
        Frequency, decreasing = TRUE), ]
183
184
185  stable_elastic_df = stable_elastic_df[order(stable_elastic_df$
        Frequency, decreasing = TRUE), ]
186
187
188  create_plot = function(data, title) {
189    ggplot(data, aes(x = reorder(Gene, -Frequency), y = Frequency,
        fill = Frequency)) +
190      geom_bar(stat = "identity") +
191      theme_minimal() +
192      labs(
193        title = title,
194        x = "Gene",
195        y = "Selection Frequency"
196      ) +
197      theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
198      scale_fill_gradient(low = "green", high = "orange")
199  }
200
201
202  plot_list = list()
203  for (i in 1:8) {
204    subset_data = stable_elastic_df[((i - 1) * 10 + 1):(i * 10), ]
```

```
205     plot_title = paste("Selected genes", (i - 1) * 10 + 1, "-", i *
            10, "for Elastic Net")
206     plot_list[[i]] = create_plot(subset_data, plot_title)
207 }
208
209
210 for (p in plot_list) {
211   print(p)
212 }
213
214 ridge_stability_df = as.data.frame(selected_genes_summary_ridge)
215
216 colnames(ridge_stability_df) = c("Gene", "Frequency")
217
218
219 stable_ridge_df = subset(ridge_stability_df, Frequency >= 8)
220
221 stable_ridge_df = stable_ridge_df[order(stable_ridge_df$Frequency
        , decreasing = TRUE), ]
222
223
224 stable_ridge_df = stable_ridge_df[order(stable_ridge_df$Frequency
        , decreasing = TRUE), ]
225
226
227 plot_list = list()
228 for (i in 1:8) {
229   subset_data <- stable_ridge_df[((i - 1) * 10 + 1):(i * 10), ]
230   plot_title <- paste("Selected genes", (i - 1) * 10 + 1, "-", i
          * 10, "for Ridge Net")
231   plot_list[[i]] <- create_plot(subset_data, plot_title)
232 }
233
234 for (p in plot_list) {
235   print(p)
236 }
```

Listing 12: Stability Analysis

```
1
2
3 min_frequency = 8
4
5
6 stable_genes_lasso = names(selected_genes_summary_lasso[selected_
      genes_summary_lasso >= min_frequency])
7
8 stable_genes_lasso = setdiff(stable_genes_lasso, "(Intercept)")
9
10
11 X_train_stable_lasso = expression_train[, stable_genes_lasso,
      drop = FALSE]
```

```
12
13  X_test_stable_lasso = expression_test[, stable_genes_lasso, drop
       = FALSE]
14
15
16  if (!identical(colnames(X_train_stable_lasso), colnames(X_test_
       stable_lasso))) {
17    stop("Mismatch between training and testing genes.")
18  }
19
20
21  cv_fit_lasso_final = cv.glmnet(
22    x = as.matrix(X_train_stable_lasso),
23    y = expression_train$response,
24    family = "binomial",
25    alpha = 1
26  )
27
28
29  lasso_model = glmnet(
30    x = as.matrix(X_train_stable_lasso),
31    y = expression_train$response,
32    family = "binomial",
33    alpha = 1,
34    lambda = cv_fit_lasso_final$lambda.min
35  )
36
37
38  pred_probs = predict(lasso_model,
39                       newx = as.matrix(X_test_stable_lasso),
40                       type = "response")
41
42
43  if (length(pred_probs) != length(response_test)) {
44    stop("Mismatch between predicted probabilities and actual
         responses.")
45  }
46
47
48  pred_class = ifelse(pred_probs > 0.5, 1, 0)
49
50
51  confusion_matrix_lasso = table(Predicted = pred_class, Actual =
       response_test)
52  print(confusion_matrix)
53
54
55  stable_genes_ridge = names(selected_genes_summary_ridge[selected_
       genes_summary_ridge >= min_frequency])
56
57  stable_genes_ridge = setdiff(stable_genes_ridge, "(Intercept)")
```

```r
58
59
60  X_train_stable_ridge = expression_train[, stable_genes_ridge,
        drop = FALSE]
61
62  X_test_stable_ridge = expression_test[, stable_genes_ridge, drop
        = FALSE]
63
64
65  if (!identical(colnames(X_train_stable_ridge), colnames(X_test_
        stable_ridge))) {
66    stop("Mismatch between training and testing genes.")
67  }
68
69
70  cv_fit_ridge_final = cv.glmnet(
71    x = as.matrix(X_train_stable_ridge),
72    y = expression_train$response,
73    family = "binomial",
74    alpha = 0
75  )
76
77
78  ridge_model = glmnet(
79    x = as.matrix(X_train_stable_ridge),
80    y = expression_train$response,
81    family = "binomial",
82    alpha = 0,
83    lambda = cv_fit_ridge_final$lambda.min
84  )
85
86
87  pred_probs = predict(ridge_model,
88                       newx = as.matrix(X_test_stable_ridge),
89                       type = "response")
90
91
92  if (length(pred_probs) != length(response_test)) {
93    stop("Mismatch between predicted probabilities and actual
          responses.")
94  }
95
96
97  pred_class = ifelse(pred_probs > 0.5, 1, 0)
98
99
100 confusion_matrix_ridge= table(Predicted = pred_class, Actual =
        response_test)
101 print(confusion_matrix)
102
103
```

```r
104  stable_genes_elastic = names(selected_genes_summary_elastic[
         selected_genes_summary_elastic >= min_frequency])
105
106  stable_genes_elastic = setdiff(stable_genes_elastic, "(Intercept)
         ")
107
108
109  X_train_stable_elastic = expression_train[, stable_genes_elastic,
          drop = FALSE]
110
111  X_test_stable_elastic = expression_test[, stable_genes_elastic,
         drop = FALSE]
112
113
114  if (!identical(colnames(X_train_stable_elastic), colnames(X_test_
         stable_elastic))) {
115    stop("Mismatch between training and testing genes.")
116  }
117
118
119  cv_fit_elastic_final = cv.glmnet(
120    x = as.matrix(X_train_stable_elastic),
121    y = expression_train$response,
122    family = "binomial",
123    alpha = 0.25
124  )
125
126
127  elastic_model = glmnet(
128    x = as.matrix(X_train_stable_elastic),
129    y = expression_train$response,
130    family = "binomial",
131    alpha = 0.25,
132    lambda = cv_fit_elastic_final$lambda.min
133  )
134
135
136  pred_probs = predict(elastic_model,
137                       newx = as.matrix(X_test_stable_elastic),
138                       type = "response")
139
140
141  if (length(pred_probs) != length(response_test)) {
142    stop("Mismatch between predicted probabilities and actual
         responses.")
143  }
144
145
146  pred_class = ifelse(pred_probs > 0.5, 1, 0)
147
148
```

```
149  confusion_matrix_elastic = table(Predicted = pred_class, Actual =
         response_test)
150  print(confusion_matrix)
```

Listing 13: Classification and Prediction using Stable genes

```
1
2
3    TN_lasso = confusion_matrix_lasso[1,1]
4
5    FN_lasso = confusion_matrix_lasso[1,2]
6
7    FP_lasso = confusion_matrix_lasso[2,1]
8
9    TP_lasso = confusion_matrix_lasso[2,2]
10
11   all_samples_lasso = TN_lasso + FN_lasso + FP_lasso + TP_lasso
12
13
14   ### Accuracy
15
16   Accuracy_lasso = ((TN_lasso + TP_lasso)/ all_samples_lasso)*100
17
18
19
20
21   ### Sensitivity
22
23   Sensitivity_lasso = ((TP_lasso)/ (TP_lasso + FN_lasso))*100
24
25
26
27
28   ### Specificity
29
30   Specificity_lasso = ((TN_lasso)/ (TN_lasso + FP_lasso))*100
31
32
33
34
35   ### Precision
36
37   Precision_lasso = ((TP_lasso)/ (TP_lasso + FP_lasso))*100
38
39
40
41
42   list(c(Accuracy_lasso = Accuracy_lasso, Sensitivity_lasso =
         Sensitivity_lasso, Specificity_lasso = Specificity_lasso,
         Precision_lasso = Precision_lasso ))
43
44
```

```r
45  TN_ridge = confusion_matrix_ridge[1,1]
46
47  FN_ridge = confusion_matrix_ridge[1,2]
48
49  FP_ridge = confusion_matrix_ridge[2,1]
50
51  TP_ridge = confusion_matrix_ridge[2,2]
52
53
54
55
56  ### Accuracy
57
58  Accuracy_ridge = ((TN_ridge + TP_ridge)/ (TN_ridge + TP_ridge +
        FN_ridge + FP_ridge) )*100
59
60
61
62
63  ### Sensitivity
64
65  Sensitivity_ridge = ((TP_ridge)/ (TP_ridge + FN_ridge))*100
66
67
68
69
70  ### Specificity
71
72  Specificity_ridge = ((TN_ridge)/ (TN_ridge + FP_ridge))*100
73
74
75
76
77  ### Precision
78
79  Precision_ridge = ((TP_ridge)/ (TP_ridge + FP_ridge))*100
80
81
82
83
84  list(c(Accuracy_ridge = Accuracy_ridge, Sensitivity_ridge =
        Sensitivity_ridge, Specificity_ridge = Specificity_ridge,
        Precision_ridge = Precision_ridge ))
85
86
87  TN_elastic = confusion_matrix_ridge[1,1]
88
89  FN_elastic = confusion_matrix_ridge[1,2]
90
91  FP_elastic = confusion_matrix_ridge[2,1]
92
```

```
93  TP_elastic = confusion_matrix_ridge[2,2]
94
95
96
97
98  ### Accuracy
99
100 Accuracy_elastic = ((TN_elastic + TP_elastic)/ (TN_elastic + TP_
        elastic + FN_elastic + FP_elastic) )*100
101
102
103
104
105 ### Sensitivity
106
107 Sensitivity_elastic = ((TP_elastic)/ (TP_elastic + FN_elastic))*
        100
108
109
110
111
112 ### Specificity
113
114 Specificity_elastic = ((TN_elastic)/ (TN_elastic + FP_elastic))*
        100
115
116
117
118
119 ### Precision
120
121 Precision_elastic = ((TP_elastic)/ (TP_elastic + FP_elastic))*100
122
123
124
125
126 list(c(Accuracy_elastic = Accuracy_elastic, Sensitivity_elastic =
        Sensitivity_elastic, Specificity_elastic = Specificity_
        elastic, Precision_elastic = Precision_elastic ))
127
128
129 metrics_df = data.frame(
130   Method = rep(c("Lasso", "Ridge", "Elastic Net"), each = 4),
131   Metric = rep(c("Accuracy", "Sensitivity", "Specificity", "
        Precision"), times = 3),
132   Value  = c(Accuracy_lasso, Sensitivity_lasso, Specificity_lasso
        , Precision_lasso,
133             Accuracy_ridge, Sensitivity_ridge, Specificity_ridge
                  , Precision_ridge,
134             Accuracy_elastic, Sensitivity_elastic, Specificity_
                  elastic, Precision_elastic)
```

```
135  )
136
137  ggplot ( metrics_df , aes (x = Metric , y = Value , fill = Method )) +
138    geom_bar ( stat = "identity", position = "dodge", width = 0.7) +
139    ylim (0 , 100) +
140    labs ( title = "Performance Metrics Comparison for Lasso , Ridge ,
         and Elastic Net",
141        x = "Metric",
142        y = "Percentage (%)") +
143    scale_fill_brewer ( palette = "Greens") +
144    theme_minimal () +
145    theme ( legend.title = element_blank ())
```

Listing 14: Model Performance Evaluation

# References

[Greenwood et al., 2020] Christopher J. Greenwood, George J. Youssef, Primrose Letcher, Jacqui A. Macdonald, Lauryn J. Hagg, Ann Sanson, Jenn Mcintosh, Delyse M. Hutchinson, John W. Toumbourou, Matthew Fuller-Tyszkiewicz, and Craig A. Olsson. *A comparison of penalised regression methods for informing the selection of predictive markers.* PLoS ONE, 15(10):e0241570, 2020. doi: 10.1371/journal.pone.0241570.

[Friedman et al., 2010] Jerome Friedman, Trevor Hastie, and Rob Tibshirani. *Regularization Paths for Generalized Linear Models via Coordinate Descent.* Journal of Statistical Software, 33(1):1-22, 2010. doi: 10.18637/jss.v033.i01.

[Golub et al., 1999] Todd R. Golub, Donna K. Slonim, Pablo Tamayo, Christine Huard, Michael Gaasenbeek, Jill P. Mesirov, Hilary Coller, Michelle L. Loh, James R. Downing, Michael A. Caligiuri, Clara D. Bloomfield, and Eric S. Lander. *Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring.* Science, 286(5439):531–537, 1999. doi: 10.1126/science.286.5439.531.