

Homework on Newton's methods

Zhuodiao Kuang(zk2275)

Problem 1:

Design a Golden-Search algorithm to find the MLE in the Example 1.2 (see the lecture notes) and implement it into **R**.

Answer: your answer starts here...

Example 1.2 in the lecture notes is about finding the Maximum Likelihood Estimation (MLE) for a trinomial distribution where the probabilities are determined by a single parameter, θ . The optimization problem involves maximizing a likelihood function expressed in terms of θ . The likelihood function is given by the sum of $(x_1 \log(\theta + 2) + x_2 \log(1 - \theta) + x_3 \log(\theta))$, with θ constrained between 0 and 1.

```
set.seed(229)
# Define the likelihood function for the trinomial distribution
likelihood_function <- function(theta, x1, x2, x3) {
  term1 <- x1 * log(theta + 2)
  term2 <- x2 * log(1 - theta)
  term3 <- x3 * log(theta)
  L<-sum(term1 + term2 + term3)
  return(L)
}

# Golden-Search algorithm
golden_search <- function(f, lower, upper, tol = 0.0001) {
  gconstant <- (sqrt(5) - 1)/2
  c <- upper - (upper - lower) * gconstant
  d <- lower + (upper - lower) * gconstant
  while (abs(c - d) > tol) {
    if (f(c) < f(d)) {
      lower <- c
    } else {
      upper <- d
    }
    c <- upper - (upper - lower) * gconstant
    d <- lower + (upper - lower) * gconstant
  }
  return((lower + upper) / 2)
}

# Example usage
# Generate Xdata
```

```

n = 100
theta <- 0.2
Xdata <- rmultinom(n, 40, c((2 + theta)/4, (1 - theta)/2, theta/4))

x1 <- Xdata[1,]
x2 <- Xdata[2,]
x3 <- Xdata[3,]
# Find the MLE of theta
theta_mle <- golden_search(function(theta) likelihood_function(theta, x1, x2, x3), 0, 1)
print(theta_mle)

## [1] 0.1936355

```

Problem 2:

Design an optimization algorithm to find the minimum of the continuously differentiable function

$$f(x) = -e^{-x} \sin(x)$$

on the closed interval $[0, 1.5]$. Write out your algorithm and implement it into **R**.

Answer: your answer starts here...

Algorithm Steps:

Using Newton's method: - For each iteration, compute $f'(x)$ and $f''(x)$ at the current point x . - Update the guess using a modified Newton step: $x_{new} = x - \frac{f'(x)}{f''(x)}$. - Check for convergence: if $|x_{new} - x| < \epsilon$, then stop and return x_{new} as the minimum. - If the maximum number of iterations is reached without convergence, consider adjusting the initial guess or the damping factor.

```

f <- function(x) {
  -exp(-x) * sin(x)
}

f_deriv <- function(x) {
  exp(-x) * (sin(x) - cos(x))
}

f_double_deriv <- function(x) {
  2 * cos(x) * exp(-x)
}

newton_method <- function(f_deriv, f_double_deriv,
                          start_point, tol = 1e-5, max_iter = 100) {
  x_current <- start_point
  for (i in 1:max_iter) {
    x_next <- x_current - f_deriv(x_current) / f_double_deriv(x_current)
    print(x_next)
    if (abs(x_next - x_current) < tol) {

```

```

    if (x_next >= 0 && x_next <= 1.5) {
      return(x_next)
    } else {
      break # If outside the interval, break and return NA
    }
  }
  x_current <- x_next
}
return(NA) # Return NA if no convergence within the interval or max iterations reached
}

# Choose an initial guess within the interval [0, 1.5]
start_point <- 0.5

# Apply the Newton Method
solution <- newton_method(f_deriv, f_double_deriv, start_point)

```

```

## [1] 0.7268488
## [1] 0.7822195
## [1] 0.7853881
## [1] 0.7853982
## [1] 0.7853982

```

```

# Display the solution
if (!is.na(solution)) {
  cat("The minimum is located at x =", solution, "\n")
  cat("The minimum value of the function is f(x) =", f(solution), "\n")
} else {
  cat("The algorithm did not converge to a solution within the interval.\n")
}

```

```

## The minimum is located at x = 0.7853982
## The minimum value of the function is f(x) = -0.3223969

```

```

# Define the function f(x) = -exp(-x) * sin(x)
f <- function(x) {
  -exp(-x) * sin(x)
}

# Create a sequence of x values from 0 to 1.5
x_values <- seq(-1.5, 1.5, by = 0.01)

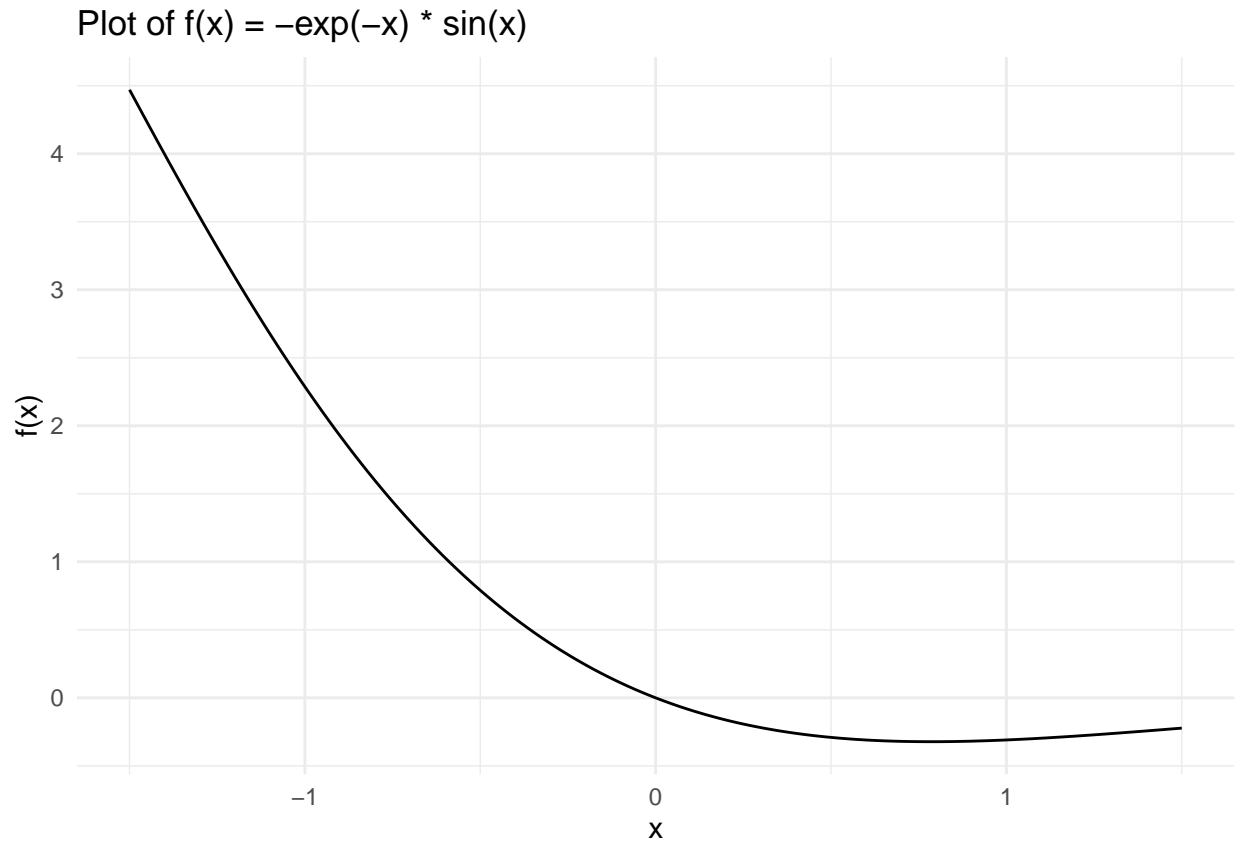
# Calculate y values using the defined function
y_values <- f(x_values)

# Create a data frame for plotting
data <- data.frame(x = x_values, y = y_values)

# Plot the function using ggplot2
ggplot(data, aes(x = x, y = y)) +
  geom_line() +

```

```
theme_minimal() +
ggtitle("Plot of f(x) = -exp(-x) * sin(x)") +
xlab("x") +
ylab("f(x)")
```



Problem 3

The Poisson distribution is often used to model “count” data — e.g., the number of events in a given time period.

The Poisson regression model states that

$$Y_i \sim \text{Poisson}(\lambda_i),$$

where

$$\log \lambda_i = \alpha + \beta x_i$$

for some explanatory variable x_i . The question is how to estimate α and β given a set of independent data $(x_1, Y_1), (x_2, Y_2), \dots, (x_n, Y_n)$.

1. Modify the Newton-Raphson function from the class notes to include a step-halving step.
2. Further modify this function to ensure that the direction of the step is an ascent direction. (If it is not, the program should take appropriate action.)
3. Write code to apply the resulting modified Newton-Raphson function to compute maximum likelihood estimates for α and β in the Poisson regression setting.

The Poisson distribution is given by

$$P(Y = y) = \frac{\lambda^y e^{-\lambda}}{y!}$$

for $\lambda > 0$.

Answer: your answer starts here...

To address the given tasks, we'll modify the Newton-Raphson method for Poisson regression to include a step-halving step and ensure the step is in an ascent direction for estimating α and β .

The likelihood function for a set of n independent observations given this model is:

$$L(\alpha, \beta | y, x) = \prod_{i=1}^n \frac{e^{-\lambda_i} \lambda_i^{y_i}}{y_i!}$$

Taking the logarithm of the likelihood function gives the log-likelihood:

$$\ell(\alpha, \beta | y, x) = \sum_{i=1}^n (-\lambda_i + y_i \log(\lambda_i) - \log(y_i!))$$

Substituting $\lambda_i = e^{\alpha + \beta x_i}$ into the log-likelihood gives:

$$\ell(\alpha, \beta | y, x) = \sum_{i=1}^n (-e^{\alpha + \beta x_i} + y_i(\alpha + \beta x_i) - \log(y_i!))$$

1. **Include Step-Halving:** To ensure convergence, we'll incorporate a step-halving procedure. If an iteration does not decrease the log-likelihood, we reduce the step size by half and reevaluate.
2. **Ensure Ascent Direction:** Before updating the parameters, we'll check if the step increases the log-likelihood. If not, we adjust the direction or reduce the step size.

```
set.seed(229) # For reproducibility

# Generate data
n <- 1000
x <- runif(n, 0, 10)
alpha <- 0.5 # True alpha
beta <- 0.1 # True beta

# Calculate lambda for the Poisson distribution
lambda <- exp(alpha + beta * x)
Y <- rpois(n, lambda)

# Check the first few values
head(x)
```

```
## [1] 7.887599 5.782350 4.363164 7.602662 4.431796 2.349943
```

```
head(Y)
```

```
## [1] 5 0 4 7 3 1
```

```
# Log-likelihood function for Poisson regression
log_likelihood <- function(alpha, beta, x, Y) {
  lambda <- exp(alpha + beta * x)
  sum(Y * log(lambda) - lambda)
}

# Gradient of the log-likelihood
gradient <- function(alpha, beta, x, Y) {
  lambda <- exp(alpha + beta * x)
  grad_alpha <- sum(Y - lambda)
  grad_beta <- sum((Y - lambda) * x)
  return(c(grad_alpha, grad_beta))
}

# Hessian matrix of the log-likelihood
hessian <- function(alpha, beta, x, Y) {
  lambda <- exp(alpha + beta * x)
  hess <- matrix(c(-sum(lambda), -sum(lambda * x),
                  -sum(lambda * x), -sum(lambda * x^2)), nrow=2)
  return(hess)
}

# Newton-Raphson method with step-halving and ascent direction check
newton_raphson <- function(x, Y, alpha_init, beta_init, tol = 1e-6, max_iter = 100) {
  alpha <- alpha_init
  beta <- beta_init
  for (i in 1:max_iter) {
    grad <- gradient(alpha, beta, x, Y)
    hess <- hessian(alpha, beta, x, Y)
    update <- solve(hess) %*% grad

    # Double check if the direction is an ascent direction
    if (t(grad) %*% update > 0) {
      cat("Iteration", i, ": Direction is not an ascent direction. Adjusting step...\n")
      update <- -update
    }

    # Step-halving
    step_size <- 1
    new_alpha <- alpha - step_size * update[1]
    new_beta <- beta - step_size * update[2]

    # Ensure Ascent Direction
    while (log_likelihood(new_alpha, new_beta, x, Y) < log_likelihood(alpha, beta, x, Y)) {
      # research for directions
      step_size <- step_size / 2
      new_alpha <- alpha - step_size * update[1]
      new_beta <- beta - step_size * update[2]
      if (step_size < tol) break
    }
  }
}
```

```

    }

    # Check for convergence
    if (sqrt(sum(update^2)) < tol) {
      cat("Converged in", i, "iterations\n")
      break
    }

    alpha <- new_alpha
    beta <- new_beta
  }
  return(c(alpha, beta))
}

alpha_init <- beta_init <- 0
estimates <- newton_raphson(x, Y, alpha_init, beta_init)

```

```
## Converged in 5 iterations
```

```
cat("True alpha:", alpha, "\nTrue beta:", beta, "\n")
```

```
## True alpha: 0.5
## True beta: 0.1
```

```
cat("Estimated alpha:", estimates[1], "\nEstimated beta:", estimates[2], "\n")
```

```
## Estimated alpha: 0.5143557
## Estimated beta: 0.09938042
```

Problem 4: Breast Cancer Diagnosis

Background

The dataset “breast-cancer.csv” consists of 569 rows and 33 columns. The first column, labeled “ID”, contains individual breast tissue images. The second column, labeled “Diagnosis”, identifies whether the image is from a cancerous or benign tissue (M=malignant, B=benign). There are 357 benign and 212 malignant cases. The other 30 columns represent the mean, standard deviation, and largest values of the distributions of 10 features computed for the cell nuclei.

- radius (mean of distances from center to points on the perimeter)
- texture (standard deviation of gray-scale values)
- perimeter
- area
- smoothness (local variation in radius lengths)
- compactness ($\text{perimeter}^2 / \text{area} - 1.0$)
- concavity (severity of concave portions of the contour)
- concave points (number of concave portions of the contour)
- symmetry
- fractal dimension (“coastline approximation” - 1)

The goal of the exercise is to build a predictive model based on logistic regression to facilitate cancer diagnosis;

Tasks:

The dataset “breast-cancer.csv” consists of 569 rows and 33 columns. The first column, labeled “ID”, contains individual breast tissue images. The second column, labeled “Diagnosis”, identifies whether the image is from a cancerous or benign tissue (M=malignant, B=benign). There are 357 benign and 212 malignant cases. The other 30 columns represent the mean, standard deviation, and largest values of the distributions of 10 features computed for the cell nuclei.

Build a logistic-LASSO model to classify images as malignant or benign. Use Lasso to automatically select features and implement a path-wise coordinate-wise optimization algorithm to obtain a path of solutions with a sequence of descending λ values.

Write a report to summarize your findings.

Report

Path-wise coordinate descent

Path-wise coordinate descent is a popular optimization technique used for fitting Lasso models. Lasso (Least Absolute Shrinkage and Selection Operator) is a regression analysis method that performs both variable selection and regularization in order to enhance the prediction accuracy and interpretability of the statistical model it produces.

How Path-wise Coordinate Descent Works

Path-wise coordinate descent optimizes the Lasso problem by iteratively updating one parameter at a time while keeping the others fixed. This approach simplifies the optimization problem into a series of one-dimensional problems that are much easier to solve.

Algorithm Steps

1. **Initialization:** Start with an initial estimate for β , often $\beta = 0$.
2. **Cyclically Update Each Coordinate:**
 - For each coefficient β_j in the vector β , update the value of β_j by solving the one-dimensional optimization problem, which involves minimizing the objective function with respect to β_j , keeping all other coefficients fixed.
 - This step often uses the soft-thresholding function due to the L1 penalty in the Lasso objective. The updated value for β_j is given by:

$$\beta_j = S\left(\frac{1}{n}X_j^T(y - X_{-j}\beta_{-j}), \lambda\right)$$

where X_j is the j^{th} column of X , X_{-j} is the matrix X excluding X_j , β_{-j} is the coefficient vector excluding β_j , and S is the soft-thresholding operator.

3. **Convergence Check:** Repeat the cyclic update until the changes in β are below a certain threshold, indicating convergence.

Key Features

- **Efficiency:** The path-wise coordinate descent method is computationally efficient, especially for high-dimensional datasets, because it breaks down the complex optimization problem into simpler one-dimensional updates.
- **Sparsity:** By leveraging the L1 penalty, it effectively produces sparse models, where irrelevant variables are given coefficients of zero, thus performing variable selection.
- **Path Solution:** It can efficiently compute the entire path of solutions for a sequence of λ values, which is useful for selecting the best model through cross-validation.

Path-wise coordinate descent is favored for its simplicity, efficiency, and effectiveness in handling large, sparse datasets common in machine learning and statistical applications.

```
# Define the sigmoid function
sigmoid <- function(z) {
  1 / (1 + exp(-z))
}

Pathwise = function(x, y, lambda = 0.001, beta = rep(0, ncol(x)), tol = 0.01){
  step = 1
  # n = nrow(x)
  # Calculate the initial loglike loss
  iniLL = -sum((y * (x %%% beta) - log(1 + exp(x %%% beta))))
  + lambda * sum(beta[2:length(beta)])

  # Optimize log likelihood loss
  while (1) {
    # Notations for difference in beta now and before
    Changed = rep(TRUE, ncol(x))
    # restore the initial value
    beta_temp = beta

    # coordinate descent
    for (j in 1:ncol(x)) {
      yhat = 1 / (1 + exp(-(x %%% beta)))
      e = y - yhat
      w = yhat * (1 - yhat)

      diff = x %%% beta - x[, -j] %%% beta[-j]
      beta_new = sum(x[, j] * (w * diff + e))

      m <- nrow(x)
      predictions <- sigmoid(x %%% beta)
      gradient <- t(x) %%% (predictions - y) / m

      if (j == 1) { # Skip regularization for intercept
        beta_new <- beta_new - gradient[j]
      } else {
        ##use soft threshold for covariates##
        beta_new = sign(beta_new) * pmax(abs(beta_new) - lambda, 0)
      }
    }
  }
}
```

```

    beta_new = beta_new / (t(w) %*% (x[, j]^2))

    # Label if not changed too much
    if (abs(beta_new - beta[j]) < tol) {
        Changed[j] = FALSE
    }

    # save the new beta
    beta[j] = beta_new
}

### coordinate descendant ends at this turn ###

loglike_loss = -sum(w * (y * (x %*% beta) - log(1 + exp(x %*% beta))))
+ lambda * sum(beta[2:length(beta)])
if (iniLL < loglike_loss || sum(Changed) == 0){
    return(list(beta_temp = beta_temp, loglike_loss = loglike_loss))
}

# default setting
iniLL = loglike_loss
step = step + 1
}
}

```

The LASSO function provided appears to implement logistic regression with LASSO regularization, allowing for variable selection and regularization of coefficients to prevent overfitting. Let's break down the code into its components and describe the mathematical notations and operations corresponding to each significant part.

1. Including Zero Lambda:

- If `include_zero_lambda` is TRUE and zero is not already in the lambda sequence, zero is added to the sequence.
 - $\lambda = \lambda \cup \{0\}$ if $\lambda \not\ni 0$ and `include_zero_lambda` is TRUE.

2. Standardizing Features:

- Calculate column means and scale of \mathbf{x} (features matrix), then standardize \mathbf{x} by subtracting the mean and dividing by its scale.
 - $\text{colmean}_i = \frac{1}{n} \sum_{j=1}^n x_{ji}$ for each column i of \mathbf{x} .
 - Standardization: $x_{ji} = \frac{x_{ji} - \text{colmean}_i}{\sqrt{\sum_{j=1}^n (x_{ji} - \text{colmean}_i)^2}}$.

3. Augmenting Feature Matrix:

- Augment \mathbf{x} with a column of ones for the intercept.
 - $\mathbf{x} = [1, \mathbf{x}]$, where 1 represents a column vector of ones.

4. Initialization of Coefficients:

- Initialize a coefficients vector `beta` with zeros.
 - $\beta = \mathbf{0}$, where $\mathbf{0}$ is a zero vector of length equal to the number of features plus one (for the intercept).

5. Handling Single vs. Multiple Lambda Values:

- If only one lambda value is provided, compute **beta** using the Pathwise Coordinate Descent method for that lambda.
- For multiple lambda values, iterate over them in decreasing order, updating **beta** for each, optionally using warm starts (reusing the **beta** from the previous lambda as the starting point for the next).

6. Pathwise Coordinate Descent Call:

- For each lambda value:
 - If **warm_start** is TRUE, update **beta** using the Pathwise Coordinate Descent method, starting from the current **beta**.
 - If **warm_start** is FALSE, reinitialize **beta** to zeros and then update **beta** using the Pathwise Coordinate Descent method.
 - $\beta^{(\lambda)} = \text{Pathwise}(x, y, \lambda, \beta, \text{tol})$, where $\beta^{(\lambda)}$ denotes the coefficients vector obtained for a given λ .

7. Returning Results:

- Return a list containing the sequence of lambda values, the list of **beta** coefficients for each lambda, and the column means and scales used for standardization.

In mathematical terms, this function configures and executes a series of optimization problems, adjusting the model's complexity by varying λ , the regularization strength. The function standardizes the input features to ensure consistent regularization effects across them and utilizes the Pathwise Coordinate Descent method to efficiently find sparse solutions in the context of LASSO regularization.

```
LASSO = function(x, y, lambda, tol = 0.01, warm_start = FALSE,
                include_zero_lambda = TRUE){
  # If include_zero_lambda is TRUE and lambda does
  # not already include 0, add 0 to the lambda sequence
  if (include_zero_lambda && sum(lambda == 0) == 0) {
    lambda = c(lambda, 0)
  }

  # Calculate column means of X (feature matrix) and standardize features
  colmean = colMeans(x) # _j for each feature j
  colscale = c() # Initialize vector to store scaling factors (_j for each feature j)
  for (i in 1:ncol(x)) {
    x[, i] = x[, i] - colmean[i] # Standardize:  $X_{ij} = X_{ij} - \_j$ 
    colscale = c(colscale, sqrt(sum(x[, i] * x[, i]))) # Calculate _j for feature j
    x[, i] = x[, i] / sqrt(sum(x[, i] * x[, i])) # Scale:  $X_{ij} = X_{ij} / \_j$ 
  }

  # Add intercept term to X
  x = cbind(rep(1, nrow(x)), x) # Augment X with a column of 1s for intercept _0
  beta = matrix(rep(0, ncol(x))) # Initialize vector (model coefficients) with zeros

  beta_list = list() # Initialize list to store vectors for different values

  # Single lambda case: perform pathwise coordinate descent and store the result
  if (length(lambda) == 1) {
    beta <- Pathwise(x, y, lambda[1], beta, tol)$beta_temp
    beta_list[[paste("beta <- lambda:", lambda[k])]] = beta
    return(list(lambda = lambda, beta = beta_list, colmean = colmean, colscale = colscale))
  }
```

```

}
# Multiple lambda case: sort lambda in decreasing order and iterate over it
else{
  lambda = sort(lambda, decreasing = TRUE)
  for (k in 1:length(lambda)) {
    if (warm_start)
      beta = Pathwise(x, y, lambda[k], beta, tol)$beta_temp
    # Use previous as starting point if warm_start is TRUE
    else{
      beta = matrix(rep(0, ncol(x)))
      # Reset to zeros if not using warm_start
      beta = Pathwise(x, y, lambda[k], beta, tol)$beta_temp
    }
    beta_list[[paste("beta <- lambda:", lambda[k])]] = beta
    # Store vector for current
  }
  return(list(lambda = lambda, beta = beta_list,
             colmean = colmean, colscale = colscale))
  # Return sequence, vectors, column means and scales
}
}

```

Criterion to select lambda

This code defines two primary functions and an example score function related to evaluating the performance of predictive models. To describe this code using notations, let's break down each function:

1. RMSE (Root Mean Square Error) Score Function

This function calculates the root mean square error, a common measure for the differences between values predicted by a model and the values actually observed. The function is defined as follows:

$$\text{rmse_score}(\mathbf{y}, \hat{\mathbf{y}}) = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

- Where:
 - \mathbf{y} represents the vector of true values (*true_y*),
 - $\hat{\mathbf{y}}$ represents the vector of predicted values (*predict_y*),
 - n is the number of observations (*length(true_y)*).

2. Criterion Function

This function iterates over a list of predictions, calculates a score for each using the `rmse_score` function, and returns a list of scores. The function can be expressed in a more general form as:

$$\text{criterion}(\mathbf{y}, \{\hat{\mathbf{y}}_1, \hat{\mathbf{y}}_2, \dots, \hat{\mathbf{y}}_k\}, \text{score}) = \{\text{score}(\mathbf{y}, \hat{\mathbf{y}}_1), \text{score}(\mathbf{y}, \hat{\mathbf{y}}_2), \dots, \text{score}(\mathbf{y}, \hat{\mathbf{y}}_k)\}$$

- Where:
 - \mathbf{y} is the vector of true values,
 - $\{\hat{\mathbf{y}}_1, \hat{\mathbf{y}}_2, \dots, \hat{\mathbf{y}}_k\}$ is the list of prediction vectors,
 - `score` is the scoring function used to evaluate predictions, which in this example is `rmse_score`.

3. Accuracy Score Function

This function calculates the accuracy, the proportion of true results (both true positives and true negatives) among the total number of cases examined. It is defined as follows:

$$\text{accuracy}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(y_i = \hat{y}_i)$$

- Where:
 - \mathbf{y} represents the vector of true values,
 - $\hat{\mathbf{y}}$ represents the vector of predicted values,
 - $\mathbf{1}(\cdot)$ is the indicator function, equal to 1 when $y_i = \hat{y}_i$ and 0 otherwise,
 - n is the number of observations.

These notations provide a mathematical representation of the R functions described in the code, translating programming constructs into mathematical formulas that describe the operations performed by the code.

```
rmse_score = function(true_y, predict_y){
  return(sqrt(sum((true_y - predict_y) ^ 2) / length(true_y)))
}

criterion = function(true_y, predict_y_list, score){
  result = list()
  for (i in names(predict_y_list)) {
    predict_y = predict_y_list[[i]]
    result[[i]] = rmse_score(true_y, predict_y)
  }
  return(result)
}

#### a score function example --- Accuracy:

accuracy = function(true_y, predict_y){
  return(sum(true_y == predict_y) / length(true_y))
}

predict = function(model, x){

  # use the scale factors for training data to standarize predicting data
  beta_list = model$beta
  colmean = model$colmean
  colscale = model$colscale
  predict_y_list = list()
  for (i in 1:ncol(x)) {
    x[, i] = x[, i] - colmean[i]
    x[, i] = x[, i] / colscale[i]
  }

  # add intercept
  x = cbind(rep(1, nrow(x)), x)

  # calculate the predicting reponse variable for each beta (or lambda)
  for (i in names(beta_list)) {
```

```

# "i" is the value of lambda
predict_y = 1 / (1 + exp(-x %*% beta_list[[i]]))
predict_y[predict_y < 0.5] = 0
predict_y[predict_y >= 0.5] = 1
predict_y_list[[i]] = predict_y
}

return(predict_y_list)
}

```

Cross Validation

1. Cross-Validation Loop

For each fold i from 1 to $k - 1$:

- Split D into training set $D_{\text{train}}^{(i)}$ and test set $D_{\text{test}}^{(i)}$.
- Prepare feature matrices $X_{\text{train}}^{(i)}$ and $X_{\text{test}}^{(i)}$, and target vectors $y_{\text{train}}^{(i)}$ and $y_{\text{test}}^{(i)}$.

2. Model Training and Prediction

Train a logistic regression model with LASSO regularization on $D_{\text{train}}^{(i)}$, using a specified sequence of λ values for regularization strength. Predict outcomes for $D_{\text{test}}^{(i)}$.

3. Evaluation Metrics Calculation

- **RMSE** for fold i :

$$\text{RMSE}^{(i)} = \sqrt{\frac{1}{n_{\text{test}}} \sum_{j=1}^{n_{\text{test}}} (\hat{y}_j - y_j)^2}$$

- **Accuracy** for fold i :

$$\text{Accuracy}^{(i)} = \frac{TP + TN}{n_{\text{test}}}$$

where TP and TN are the counts of true positives and true negatives, respectively, and n_{test} is the size of the test set.

- **AUC** for fold i is calculated using the `ROCR` package, comparing the predicted probabilities against the true binary outcomes to evaluate the model's discriminative ability.

4. Return Results

Return a list containing vectors of RMSE, accuracy, and AUC values for the $k - 1$ folds.

Note: The provided code has a logical discrepancy as it omits the last fold from evaluation and initializes metrics arrays with length $k - 1$ instead of k . This may not fully utilize the provided data for cross-validation and could be an area for correction.

```

CrossValidation = function(data, k = 10) {
  # Create k folds from the data, stratified by the 'diagnosis' variable
  folds = createFolds(data$diagnosis, k = k, list = TRUE)

  # Initialize vectors to store the Root Mean Square Error (RMSE)
  #, accuracy, and Area Under the Curve (AUC) for each fold
  rmse = rep(NA, k - 1)
  accuracy = rep(NA, k - 1)
}

```

```

auc = rep(NA, k - 1)

# Iterate over each fold, leaving one out as test set in each iteration
for (fold in 1:(length(folds) - 1)) {
  # Define the test set for the current fold and the training set as the rest
  test = data[folds[[fold]],]
  train = data[-folds[[fold]],]

  # Prepare the training data: select features
  # (excluding 'id' and 'diagnosis') and convert to matrix
  x_train = train |> select(-id, -diagnosis) |> as.matrix()
  # Prepare the training labels: binary encoding of 'diagnosis' (M = 1, B = 0)
  y_train <- matrix(rep(0, nrow(train)))
  y_train[train$diagnosis == "M"] = 1

  # Prepare the test data similarly to training data
  x_test = test |> select(-id, -diagnosis) |> as.matrix()
  y_test <- matrix(rep(0, nrow(test)))
  y_test[test$diagnosis == "M"] = 1

  # Train the LASSO model on the training data
  train_model = LASSO(x_train, y_train,
    lambda = seq(exp(-8), exp(-7), exp(-9)),
    tol = 0.01, warm_start = TRUE,
    include_zero_lambda = FALSE)

  # Make predictions on the test set
  predict_result = predict(train_model, x_test)

  # Calculate and store the RMSE for the current fold
  error = c(predict_result[[1]]) - y_test
  rmse[fold] = sqrt(mean(error^2))

  # Calculate and store the accuracy for the current fold
  tp = sum(ifelse(c(predict_result[[1]]) == 1 & y_test == 1, 1, 0))
  # True Positives
  tn = sum(ifelse(c(predict_result[[1]]) == 0 & y_test == 0, 1, 0))
  # True Negatives
  fp = sum(ifelse(c(predict_result[[1]]) == 1 & y_test == 0, 1, 0))
  # False Positives
  fn = sum(ifelse(c(predict_result[[1]]) == 0 & y_test == 1, 1, 0))
  # False Negatives

  accuracy[fold] = (tp + tn) / length(y_test)

  # Calculate and store the AUC for the current fold
  pred = ROCR::prediction(c(predict_result[[1]]), y_test)
  auc_perf = ROCR::performance(pred, "measure" = "auc")
  auc[fold] = auc_perf@y.values[[1]]
}

# Return a list containing the computed RMSE, accuracy, and AUC for all folds
return(list("RMSE" = rmse, "Accuracy" = accuracy, "AUC" = auc))

```

```

}

# load data
cancer = read.csv("breast-cancer.csv")

x = cancer |> select(-id, -diagnosis) |> as.matrix()

y <-matrix(rep(0,nrow(cancer)))
y[cancer$diagnosis=="M"] = 1

# calculate beta_hat
beta = Pathwise(x, y, tol = 0.01)$beta_temp
beta

## [1] -2.54365344 0.36908757 -0.07913149 0.04479890 71.27163517
## [6] 4.23884754 7.84833165 70.08539219 14.31575292 -105.64397826

# check the data in glm
model = glm( y~x, family = "binomial")
model$coefficients[-1]

```

```

##          xradius_mean      xtexture_mean      xperimeter_mean
##          -2.04930490          0.38473434          -0.07151042
##          xarea_mean      xsmoothness_mean      xcompactness_mean
##          0.03979620          76.43227376          -1.46242225
##          xconcavity_mean      xconcave.points_mean      xsymmetry_mean
##          8.46869976          66.82175685          16.27824232
## xfractal_dimension_mean
##          -68.33702689

```

The beta's are fairly close!

```

set.seed(314)
fit_stats = CrossValidation(cancer, k = 10)
rmse = mean(fit_stats$RMSE)
rmse

```

```
## [1] 0.253325
```

```

acc = mean(fit_stats$Accuracy)
acc

```

```
## [1] 0.9334075
```

```

auc = mean(fit_stats$AUC)
auc

```

```
## [1] 0.9256373
```

High AUC(92.56%) models are particularly valuable in medical diagnostics for prioritizing follow-up tests or interventions for higher-risk individuals.