

An R-based Machine Learning Protocol To Predict Diagnostic Biomarkers in Breast Cancer Using RNA-sequencing data

Paibali Shaw, Kirithana Viswanathani and Manjari Kiran*

2026-02-20

Loading Required Libraries

```
library(devtools)
library(DESeq2)
library(caret)
library(randomeForest)
library(PROC)
library(ggplot2)
library(ggrep2)
```

Raw Counts Data Loading and Filtering

Raw tumor (TP) and normal (NT) samples were loaded. They were filtered against genes having less than 10 counts across columns. To ensure it is prepared for next steps, it was transformed into a genes \times samples dataframe.

```
tumor_raw <- read.csv("raw_counts_tumor.csv", check.names = FALSE)
normal_raw <- read.csv("raw_counts_normal.csv", check.names = FALSE)

# For example: Visualizing the tumor counts file for reference
head(tumor_raw[, 1:5])

##          shortLetterCode ENSG00000000003 ENSG00000000005
## TC0A.A1.BB.01A.11B.A12D.07 TP 1036 21
## TC0A.E2.A14P.01A.31B.A12D.07 TP 1033 9
## TC0A.AH.A04A.01A.21B.A034.07 TP 1417 9
## TC0A.AH.A07A.01A.11B.A1M.07 TP 1117 29
## TC0A.AH.A071.01A.11B.A0U.07 TP 1117 15
## TC0A.E9.A1R4.01A.21B.A14D.07 TP 111 15
##          ENSG000000000419 ENSG000000000457
## TC0A.BB.A1B8.01A.11B.A12D.07 1439 1049
## TC0A.E2.A14P.01A.31B.A12D.07 4553 932
## TC0A.AH.A04A.01A.21B.A034.07 1982 1227
## TC0A.AH.A07A.01A.11B.A1M.07 2623 1727
## TC0A.AH.A071.01A.11B.A0U.07 1426 974
## TC0A.E9.A1R4.01A.21B.A14D.07 1943 1285
```

```
# Transpose to Genes  $\times$  Samples
tumor_raw <- as.data.frame(t(tumor_raw))
normal_raw <- as.data.frame(t(normal_raw))

# Combine Tumor + Normal
counts <- cbind(tumor_raw, normal_raw)
cat(paste("Counts DataFrame Dimensions:", paste(dim(counts), collapse = " x ")))

## Counts DataFrame Dimensions: 60619 x 1224
```

```
# Count Filtering
keep <- rowSums(counts) >= 10 >= 10
counts_filtered <- counts[keep, ]
cat(paste("Counts DataFrame Dimensions After Filtering:", paste(dim(counts_filtered), collapse = " x ")))

## Counts DataFrame Dimensions After Filtering: 39293 x 1224
```

Stratified Train-Test Split (70:30)

The filtered counts were split into training and testing datasets.

Note: The shortLetterCode column in `ml_data` contains the condition labels for tumor (TP) and normal (NT) samples.

```
ml_data <- as.data.frame(t(counts_filtered))

set.seed(123)
trainIndex <- createDataPartition(
  ml_data$shortLetterCode,
  p = 0.7,
  list = FALSE
)
trainData <- ml_data[trainIndex, ]
testData <- ml_data[-trainIndex, ]

# To visualize the 70:30 sample split
table(trainData$shortLetterCode)

##
## NT TP
## 80 778

table(testData$shortLetterCode)

##
## NT TP
## 33 333
```

Data Preparation and Running DESeq2 analysis

DESeq2 analysis was run only on the training dataset to obtain differentially expressed genes (DEGs). This is done to prevent data leakage to the testing dataset. DEGs were visualized using a volcano plot.

```
# Preparing Training data for DESeq2
trainData <- as.data.frame(t(trainData)) # Converting a genes  $\times$  samples
train_group <- as.vector(unlist(trainData[, 3])) # Storing the counts
train_counts <- trainData[, 3, ] # Storing the labels separately (TP or NT)
train_counts <- as.data.frame(
  apply(train_counts, 2, FUN=function(x) as.numeric(trimws(x))),
  row.names = rownames(train_counts) # Converting to numeric without changing rownames
train_counts[is.na(train_counts)] <- 0 # Converting NA counts to 0
train_counts <- as.matrix(train_counts) # DESeq2 requires a matrix as input for counts

col_data <- data.frame(
  row.names = colnames(train_counts),
  condition = factor(train_group)
)

# Creating a DESeq2 object
dds_train <- DESeqDataSetFromMatrix(
  countData = train_counts,
  colData = col_data,
  design = ~ condition
)

# Preparing the Testing dataset similarly
testData <- as.data.frame(t(testData))
test_group <- as.vector(unlist(testData[, 3]))
test_counts <- testData[, 3, ]
test_counts <- as.data.frame(
  apply(test_counts, 2, FUN=function(x) as.numeric(trimws(x))),
  row.names = rownames(test_counts)
test_counts[is.na(test_counts)] <- 0
test_counts <- as.matrix(test_counts)

col_data <- data.frame(
  row.names = colnames(test_counts),
  condition = factor(test_group)
)
dds_test <- DESeqDataSetFromMatrix(
  countData = test_counts,
  colData = col_data,
  design = ~ condition
)

# Running DESeq2 on the Training dataset only
dds_train <- DESeq(dds_train)
res_df_train <- results(dds_train)
res_df_train <- as.data.frame(res_df_train)
```

```
# Volcano Plot
# Reconvert clean dataframe from DESeq2 results
res_df_train <- as.data.frame(res_df_train)

# Remove NA padj
res_df_train <- res_df_train[is.na(res_df_train$padj), ]

# Convert columns explicitly
res_df_train$padj <- as.numeric(as.res_df_train$padj)
res_df_train$log2FoldChange <- as.numeric(as.res_df_train$log2FoldChange)

# Add gene ID column
res_df_train$gene_id <- rownames(res_df_train)

res_df_train$threshold <- ifelse(
  res_df_train$padj < 0.05 & abs(res_df_train$log2FoldChange) > 1,
  "Significant",
  "Non-Significant"
)

res_df_train$threshold <- factor(
  res_df_train$threshold,
  levels = c("Non-Significant", "Significant")
)

# Order by padj
res_df_train <- res_df_train[order(res_df_train$padj), ]

# Keep only significant
res_sig <- res_df_train[res_df_train$threshold == "Significant", ]

# Taking the top 10 genes
top_genes <- head(res_sig, 10)

ggplot(res_df_train,
  aes(x = log2FoldChange,
      y = -log10(padj),
      color = threshold)) +
  geom_point(alpha = 0.5) +
  theme_minimal()
# Gene names
res_sig$gene_id <- top_genes$gene_id,
aes(label = gene_id,
  size = 2,
  theme_minimal() +
  scale_color_manual(values = c("grey", "red")) +
  labs(title = "Volcano Plot (Training Data)",
    x = "Log2 Fold Change",
    y = "-Log10 Adjusted p-value",
    color = "Differential Expression")

# Volcano Plot (Training Data)
```

The top 200 differentially expressed genes (DEGs), which are our features, were selected based on a threshold of p -value < 0.05 and \log_2 fold change ≥ 1 . Variance stabilizing transformation (VST) normalization was applied to both the training and testing datasets. To evaluate the ability of these selected DEGs to distinguish between tumor and normal samples, principal component analysis (PCA) was performed, and the resulting clusters were visualized.

```
# Selecting Top 200 Genes
res_sig <- res_df_train %>%
  filter(padj < 0.05 & abs(log2FoldChange) > 1) %>%
  arrange(padj)
top_genes <- rownames(res_sig[1:200])
cat("Number of significant DEGs:", nrow(res_sig), "\n")

## Number of significant DEGs: 9087
```

```
# Normalizing the Training dataset
# Using Variance Stabilizing Transformation on Raw counts (training data)
vst_train <- vst(dds_train, blind=FALSE)
train_expr <- assay(vst_train)
train_ml <- as.data.frame(t(train_expr))
train_ml$shortLetterCode <- train_group
train_ml <- train_ml %>% select(shortLetterCode, everything())

# Normalizing the Testing dataset
test_expr <- assay(vst_test)(top_genes, )
test_ml <- as.data.frame(t(test_expr))
test_ml$shortLetterCode <- test_group
test_ml <- test_ml %>% select(shortLetterCode, everything())

# PCA Plot (Training Data)
pca_data <- plotPCA(vst_train,
  intgroup="condition",
  label="none")
percentVar <- round(100 * attr(pca_data, "percentVar"))

ggplot(pca_data, aes(PC1, PC2, color=condition)) +
  geom_point(size=3) +
  labs(title = "PCA Plot (Training Data)",
    x = "PC1: 20%",
    y = "PC2: 12%",
    color = "condition")
```

The PCA plot showed clear separation between tumor and normal samples, supporting effective DEG-based feature selection, though some overlap of normal samples within the tumor cluster may reflect data imbalance.

Building the Random Forest Model

The data were split into training and testing sets with x (features) and y (labels) variables. A 5-fold cross-validation was applied to the training set to prevent overfitting, and the model was trained using 200 trees.

```
# x_train and y_test contain the features or counts data (input variable)
# y_train and y_test contain the sample labels (output variable)
y_train <- as.factor(train_ml$shortLetterCode)
y_test <- as.factor(test_ml$shortLetterCode)
x_train <- train_ml %>% dplyr::select(-shortLetterCode)
x_test <- test_ml %>% dplyr::select(-shortLetterCode)
cat("x_train =", paste(dim(x_train), collapse = " x "),
  "\n" | y_train = ", length(y_train), "\n")

## x_train = 858 x 200 | y_train = 858

# Test split
# Only selecting the data from the x_test set
y_test <- as.factor(y_test)
y_test <- factor(y_test, levels = c("NT", "TP"))
x_test <- test_ml %>% dplyr::select(-shortLetterCode)
x_test <- as.matrix(x_test)
cat("x_test =", paste(dim(x_test), collapse = " x "),
  "\n" | y_test = ", length(y_test), "\n")

## x_test = 366 x 200 | y_test = 366

# Ensuring identical feature order between x_train and x_test
r <- identical(colnames(x_train), colnames(x_test))
# Output should be true

# Building the Random Forest Model
set.seed(123)

# Train Control for 5-Fold CV
ctrl <- trainControl(
  method = "cv",
  number = 5,
  classProbs = TRUE,
  summaryFunction = twoClassSummary,
  savePredictions = "final"
)

# Train Random Forest with 5-Fold CV
x = x_train,
y = y_train,
method = "rf",
metric = "ROC",
ctrl = ctrl,
cross = 20,
importance = TRUE
)

print(rf_model)
```

```
## Random Forest
## 858 samples
## 200 predictors
## 2 classes: 'NT', 'TP'
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 687, 686, 686, 686, 687
## Resampling results across tuning parameters:
```

```
## mtry ROC Sens Spec
## 2 0.9987565 0.9000 0.9948553
## 101 0.9995983 0.9375 0.9961456
## 200 0.9996384 0.9375 0.9935732

## ROC was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 200.
```

Evaluating Model Performance: Testing Data

The trained model was evaluated on the testing data, and the Confusion Matrix, Receiver Operating Characteristic (ROC) curve and Area Under the Curve (AUC) were computed.

```
# Prediction on Test Data
pred_class <- predict(rf_model, x_test)

# Evaluation of Model Performance
# Confusion matrix
confusionMatrix(pred_class, y_test)

## Confusion Matrix and Statistics
##
## Reference
## Prediction NT TP
## NT 32 0
## TP 1 333
##
## Accuracy : 0.9973
## 95% CI : (0.9849, 0.9999)
## No Information Rate : 0.9098
## P-Value (Acc > NIR) : 3.563e-14
##
## Kappa : 0.9931
##
## Mcnemar's Test P-Value : 1
##
## Sensitivity : 0.96970
## Specificity : 1.00000
## Pos Pred Value : 1.00000
## Neg Pred Value : 0.99701
## Prevalence : 0.00016
## Detection Rate : 0.08743
## Detection Rate : 0.08743
## Balanced Accuracy : 0.98485
##
## 'Positive' Class : NT

## ROC Curve + AUC
pred_prob <- predict(rf_model, x_test, type = "prob")[,2]
roc_obj <- roc(y_test, pred_prob)
auc(roc_obj)

## Area under the curve: 0.9999

plot(roc_obj,
  col = "darkred",
  lwd = 3,
  main = "ROC Curve - Random Forest")
```

The ROC curve shows a high area under the curve (AUC) of approximately 0.9999, indicating excellent model performance in distinguishing between tumor and normal samples.

Evaluating Model Performance: Independent dataset

An independent dataset comprising 18 normal and 35 breast tumor samples was used to evaluate the performance of the model trained on the top 200 DEGs.

```
# Reading and preparing independent data
indep_data <- read.csv("test11.csv", header = T)
indep_data <- as.data.frame(indep_data)
indep_data$class <- as.character(indep_data$class)
head(indep_data[, 1:5])

##          ID class ENSG000000000003 ENSG000000000005 ENSG000000000419
## 1 SRR2148235 Tumor 8.175975 2.639261 8.916178
## 2 SRR2148236 Tumor 9.749349 6.410359 9.210448
## 3 SRR2148237 Tumor 6.704148 3.852458 0.876623
## 4 SRR2148238 Normal 9.212099 6.027375 6.585745
## 5 SRR2148239 Normal 10.297149 1.873562 10.229853
## 6 SRR2148240 Normal 10.273146 5.919515 9.324600

indep_data <- indep_data %>%
  mutate(class = case_when(
    class == "Tumor" ~ "TP",
    class == "Normal" ~ "NT",
    TRUE ~ class
  ))
indep_data %>%
  group_by(class) %>%
  summarise(n = n())

## # A tibble: 2 x 2
## class n
## <chr> <int>
## 1 NT 18
## 2 TP 15

# Preparing an independent dataset as the test data
indep_data <- as.data.frame(indep_data)
rownames(indep_data) <- indep_data$ID
indep_data <- indep_data %>% select(-c(ID))
cat(paste("Independent data counts dimensions:", paste(dim(indep_data), collapse = " x ")))

## Independent data counts dimensions: 33 x 58736

y_indep <- as.factor(indep_data$class)
y_indep <- factor(y_indep, levels = c("NT", "TP"))
x_indep <- indep_data %>% dplyr::select(-class)
x_indep <- as.matrix(x_indep)

cat("x_indep =", paste(dim(x_indep), collapse = " x "),
  "\n" | y_indep = ", length(y_indep), "\n")

## x_indep = 33 x 200 | y_indep = 33

x_indep <- as.matrix(x_indep)
r <- identical(colnames(x_train), colnames(x_indep))

# Running the model on independent dataset
pred_class <- predict(rf_model, x_indep)

# Evaluation of Model Performance
# Confusion matrix
confusionMatrix(pred_class, y_indep, positive = "TP")

## Confusion Matrix and Statistics
##
## Reference
## Prediction NT TP
## NT 12 6
## TP 6 9
##
## Accuracy : 0.6364
## 95% CI : (0.4512, 0.736)
## No Information Rate : 0.5453
## P-Value (Acc > NIR) : 0.1916
##
## Kappa : 0.2667
##
## Mcnemar's Test P-Value : 1.0000
##
## Sensitivity : 0.6000
## Specificity : 0.6667
## Pos Pred Value : 0.6000
## Neg Pred Value : 0.6667
## Prevalence : 0.4545
## Detection Rate : 0.2727
## Detection Rate : 0.4545
## Balanced Accuracy : 0.6333
##
## 'Positive' Class : TP

## ROC Curve + AUC
pred_prob <- predict(rf_model, x_indep, type = "prob")[,2]
roc_obj <- roc(y_indep, pred_prob)
auc_value <- auc(roc_obj)
cat("AUC value =", auc_value, "\n")

## AUC value = 0.651819

plot(roc_obj,
  col = "darkred",
  lwd = 3,
  legacy.axes=TRUE,
  print.auc=TRUE,
  main="ROC Curve - Independent Dataset")
```

The ROC curve for the independent dataset shows an AUC of 0.652, indicating moderate model performance on new, unseen data.

Finding the Top 10 Biomarkers from our dataset

```
library(Biomart)

# Connecting to Biomart to parse for the human genes dataframe
mart <- useEnsembl(
  biomart = "ENSEMBL",
  dataset = "hsapiens_gene_ensembl"
)

# Taking the previously obtained significant genes dataframe
# Getting gene names for the obtained IDs
gene_annotation <- getBM(
  attributes = c("ensembl_gene_id", "hgnc_symbol"),
  filters = c("ensembl_gene_id"),
  values = res_sig$gene_id,
  mart = mart
)

# Merging with the DESeq2 table (res_sig)
res_sig_annotated <- res_sig %>%
  left_join(gene_annotation, by = c("gene_id" = "ensembl_gene_id"))

# Removing Unannotated genes
res_sig_annotated <- as.data.frame(res_sig_annotated)
filter(hgnc_symbol != "")

# Ordering based on padj
top10_biomarkers <- res_sig_annotated %>%
  arrange(padj) %>%
  dplyr::select(1:10) %>%
  pull(hgnc_symbol)

# Printing the top 10 biomarkers
cat(top10_biomarkers, sep = "\n")

## HNF1L
## HNF2
## HNF3
## TNSI
## COL11A1
## LINC01614
## LYVE1
## RIF1A
## SIRT
## PRMT11

## sessionInfo()

## R version 4.5.2 (2025-10-31)
## Platform: x86_64-pc-linux-gnu
## Running under: Ubuntu 22.04.4 LTS
##
## Matrix products: default
## BLAS: /usr/lib/x86_64-linux-gnu/blas/libblas.so.3.10.0
## LAPACK: /usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.10.0 LAPACK version 3.10.0
##
## locale:
## [1] LC_CTYPE=en_US.UTF-8 LC_NUMERIC=C
## [3] LC_TIME=en_US.UTF-8 LC_COLLATE=en_US.UTF-8
## [5] LC_MONETARY=en_US.UTF-8 LC_MESSAGES=en_US.UTF-8
## [7] LC_PAPER=en_US.UTF-8 LC_NAME=C
## [9] LC_ADDRESS=C LC_TELEPHONE=C
## [11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
##
## time zone: Asia/Kolkata
## RStudio: source:github (lib)
##
## attached base packages:
## [1] stats4 stats graphics grDevices utils datasets methods
## [8] base
##
## other attached packages:
## [1] biomart 2.46.0 ggrep2 0.9.6
## [3] proc 1.19.0.1 randomForest 4.7-1.2
## [5] caret 7.0-1 lattice 0.22-5
## [7] ggplot2 0.0.2 DESeq2 1.59.2
## [9] SumatranExperiment 1.40.0 Biobase 2.70.0
## [11] MatrixGenerics 1.22.0 matrixData 1.5.0
## [13] GenomicRanges 1.62.1 pillar 1.21.1
## [15] Ranges 2.44.0 S4Vectors 0.48.0
## [17] Bioconductor 0.56.0 generics 0.1-4
## [19] dplyr 1.2.0
##
## loaded via a namespace (and not attached):
## [1] DBI 1.2.3 http2 1.2.2 rlang 1.1.7
## [4] Rcpp 1.2.0-4 openssl 2.0.4 okei 0.2-0 e1071 7.1-7
## [7] compiler 4.5.2 RSQLite 2.4.6 png 0.1-8
## [10] vctrs 0.7.1 reshape2 1.4.5 stringi 1.6.0
## [13] crayon 1.5.3 plotly 4.10.0 fontawesome 0.5.6
## [16] dplyr 2.5.2 XVector 0.50.0 labeling 0.4-3
## [19] utf8 1.2-6 rmarkdown 2.30 stringr 1.4.0
## [22] bit 4.0-4 pillar 1.21.1 forcats 0.5.0
## [25] cachem 1.1.0 jsonlite 2.0.0 progress 1.2-3
## [28] blob 1.3.0 rscapex 1.3.1 DelayedArray 0.36.0
## [31] Bioconductor 1.44.0 RColorbrewer 1.2.2 RColorbrewer 1.2.2
## [34] R6 2.6.1 palisade 1.0.0 stringr 1.4.7
## [37] RColorbrewer 1.1-3 parallel 1.46.1 spat 4.1-24
## [40] lubridate 1.9.0 RColorbrewer 1.2.2 RColorbrewer 1.2.2
## [43] iterators 1.0.14 knitr 1.51 future 1.32.0
## [46] Matrix 1.7-4 splines 4.5-2 mnet 7-3-20
## [49] timechange 0.4.0 tidyselect 1.2.1 rsudopkg 1.0.18
## [52] GenomicRanges 1.62.1 yaml 2.3.12 forecast 1.20.0
## [55] code2tools 0.2-19 curl 7.0.0 listenr 0.10.0
## [58] tibble 3.3.1 plyr 1.8.9 RSGDR 1.50.0
## [61] withr 3.0.2 sass 0.4.10 evaluate 1.0.5
## [64] future 1.69.0 survival 3.8-3 proxy 0.4-29
## [67] BiocFileCache 3.0.0 xml2 1.5.2 Biobase 2.78.0
## [70] RColorbrewer 1.0.3 RColorbrewer 1.2.2 RColorbrewer 1.2.2
## [73] hms 1.1.4 scales 1.4.0 global 0.19.0
## [76] class 7.3-23 glue 1.8.0 tools 4.5.2
## [79] RColorbrewer 1.0.3 RColorbrewer 1.2.2 RColorbrewer 1.2.2
## [82] RColorbrewer 1.0.3 RColorbrewer 1.2.2 RColorbrewer 1.2.2
## [85] RColorbrewer 1.0.3 RColorbrewer 1.2.2 RColorbrewer 1.2.2
## [88] RColorbrewer 1.0.3 RColorbrewer 1.2.2 RColorbrewer 1.2.2
## [91] RColorbrewer 1.0.3 RColorbrewer 1.2.2 RColorbrewer 1.2.2
## [94] RColorbrewer 1.0.3 RColorbrewer 1.2.2 RColorbrewer 1.2.2
## [97] RColorbrewer 1.0.3 RColorbrewer 1.2.2 RColorbrewer 1.2.2
## [100] RColorbrewer 1.0.3 RColorbrewer 1.2.2 RColorbrewer 1.2.2
```