

Evaluating Model Accuracy and Interpreting Variable Selection in Predicting Malignant Breast Tumors

Stephanie Daniella Hernandez Prado

December 2025

1 Introduction

Early detection of breast cancer is crucial to improve clinical outcomes and reduce mortality. In this study, we analyzed a dataset composed of morphological measurements obtained from digitized images of breast cells, with the aim of predicting whether a sample corresponds to a malignant or benign tumor. Due to the high dimensionality of the dataset and the strong correlations between many of the variables, it is necessary to carefully evaluate different modeling methods to determine which offer the best predictive performance and greatest stability.

The main objective of this work is to compare the performance of different classification models: Logistic Regression, Stepwise Selection, Backward Selection, Forward Selection, LASSO, and Ridge, with the intention of comparing them to a Random Forest machine learning model. To make the comparison as stable and comparable as possible, we will use a 5-fold Cross-Validation procedure, performing it 50 times with the same 5 folds for each model in each iteration.

Performance evaluation is performed by comparing the averages of the cross-validation Accuracy, which measures the percentage of correct classifications, and the Brier Score, which quantifies the confidence with which a model makes a prediction. We will also explore why models like LASSO choose certain variables to synthesize information and achieve simpler models in the presence of multicollinearity.

Finally, we will seek to identify the most accurate, stable and appropriate model for each situation to predict tumor diagnosis.

2 Dataset

The used dataset is the Breast Cancer Wisconsin (Diagnostic) Dataset [1], which is a widely used machine learning dataset containing 569 samples from fine needle aspirate (FNA) images of breast masses, collected at the University of Wisconsin Hospitals by Dr. William H. Wolberg. The dataset includes 30 numerical features derived from digitized cell nuclei images, where each of 10 cell characteristics (radius, texture, perimeter, area, smoothness, compactness, concavity, concave points, symmetry, and fractal dimension) is measured three ways: mean, standard error, and "worst" (mean of three largest values). Each sample is labeled as either malignant (212 cases, 37.3%) or benign (357 cases, 62.7%), making it a binary classification problem.

3 Methodology

3.1 Models

In this study, several algorithms were used to compare their predictive capacity and stability, measuring their performance under the same cross-validation conditions repeated multiple times. Each model was trained using the same folds during the same iteration, ensuring a fair comparison.

The models used are listed below along with a brief description:

3.1.1 Logistic Regression

Logistic regression is a type of model used to predict categorical variables based on a combination of predictor variables. This model uses the logit function as its link function, allowing us to estimate the probabilities of each possible outcome modeled as a function of the explanatory variables using the logistic function. This model is easy to understand, effective, and widely known, so we will use it as our base model.

Logistic regression looks like this:

$$\text{logit}(p_i) = \ln\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \sum_{j=1}^k \beta_j x_{j,i}$$

where x_i is the probability of success of the i -th observation, while $X_{j,i}$ is the j -th predictor variable of the i -th observation.

3.1.2 Stepwise

This model is a version of the previous model where we remove variables, add variables, or both actions based on their significance with the goal of identifying a set of predictors that improve the fit without overfitting.

The Forward method starts with a model containing only the intercept and adds variables with the smallest p-value, until no variables can be added due to having too large a p-value. The Backward method starts with the entire model and removes variables using the same criteria. The Both method combines both approaches.

3.1.3 LASSO

This is an extension of logistic regression, where an L1 constraint is incorporated into the size of the coefficients. This constraint favors simpler models and eliminates irrelevant predictors. The regularization parameter λ chosen was the one that performed best during cross-validation, giving us a balanced penalty parameter.

The parameters of the model applying LASSO penalty are found as follows:

$$\hat{\beta}_{LASSO} = \arg \min_{\beta} \{-\log \text{Likelihood}(\beta) + \lambda \sum_{j=1}^k |\beta_j|\}$$

where λ is the regularization parameter, which controls how much the coefficients are penalized.

3.1.4 Ridge

A model similar to Lasso, but using an L2 penalty that reduces the magnitude of the coefficients without making them exactly zero. It tends to work better when many predictors are correlated.

The parameters of the model applying Ridge penalty are found as follows:

$$\hat{\beta}_{Ridge} = \arg \min_{\beta} \{-\log \text{Likelihood}(\beta) + \lambda \sum_{j=1}^k \beta_j^2\}$$

3.1.5 Random Forest

Decision trees work by recursively splitting the dataset based on the values of the predictor variables. At each node, the variable and split point that best separate the classes are selected. The process continues until a stopping criterion is reached, either because the maximum depth of the tree has been reached or because the observations within a node are sufficiently homogeneous. At the end, each leaf of the tree assigns a predicted class based on the observations it contains.

A Random Forest involves training many decision trees, each constructed from a subset of observations obtained through bootstrap sampling and using only a random subset of variables at each split. Each tree produces its own prediction, and the final model response is obtained through majority voting or averaging. This approach reduces variance and improves the model's generalizability compared to a single tree.

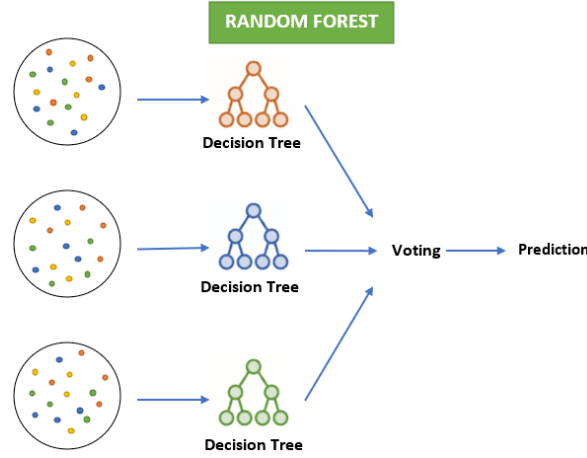


Figure 1: Diagram of the 5-Fold Cross-Validation process (Illustration of the Random Forest algorithm (Obtained from: [2])).

3.2 Metrics

The performance of the models was evaluated using the following metrics:

3.2.1 Accuracy

Accuracy is the proportion of correct predictions out of the total number of observations.

$$Accuracy = \frac{\text{Correct predictions}}{\text{Total predictions}}$$

3.2.2 Brier Score

Brier Score is the mean squared error of the forecast. It tells us how confidently the model predicts the response variable in each observation.

$$\text{Brier Score} = \frac{1}{N} \sum_{i=1}^N (p_i - y_i)^2$$

Where:

p_i : Predicted probability of the positive class.

y_i : Actual value encoded as 0 or 1.

3.3 K-Fold Repeated Cross Validation.

To obtain a robust evaluation, each model was trained and evaluated using a process that we will call K-Fold Repeated Cross Validation

In k-fold cross-validation, the original sample is randomly partitioned into k equal sized subsamples, often referred to as "folds". Of the k subsamples, a single subsample is retained as the validation data for testing the model, and the remaining k - 1 subsamples are used as training data. The cross-validation process is then repeated k times, with each of the k subsamples used exactly once as the validation data.

In our case, we used 5 folds and took the average performance (Accuracy and Brier Score) of these 5 folds for each model, then repeated this process 50 times. This allows for obtaining performance distributions, comparing variability, and avoiding conclusions based on a single partition of the data set.

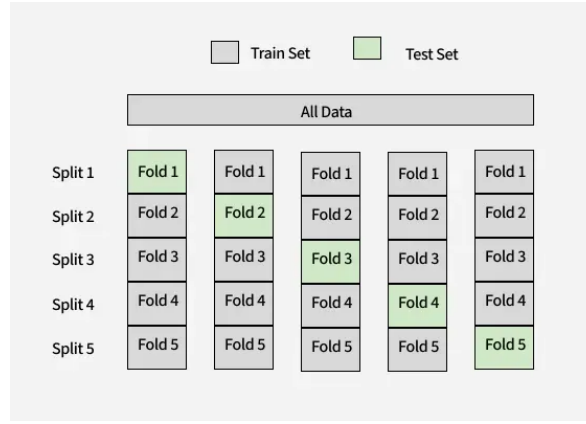


Figure 2: Illustration of the 5-Fold Cross-Validation process (obtained from [3]).

4 Results

4.1 Results of the first iteration.

4.1.1 Resulting Parameters by Model

Base model

The base model shows that most imaging features strongly influence the diagnosis, but the model is still unrefined and contains many significant and non-significant predictors that could later be simplified or improved.

```
##
## Call:
## NULL
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)    ## concavity2      1.529e+08  5.340e+06  28.624 < 2e-16 ***
## (Intercept) -2.881e+06  2.816e+05 -10.233 < 2e-16 ***    ## concave_points2 -1.260e+09  4.012e+07 -31.398 < 2e-16 ***
## radius1      2.427e+06  2.693e+05   9.014 < 2e-16 ***    ## symmetry2       2.890e+08  4.126e+06  70.055 < 2e-16 ***
## texture1     1.958e+05  1.471e+04  13.313 < 2e-16 ***    ## fractal_dimension2 1.512e+09  6.597e+07  22.921 < 2e-16 ***
## perimeter1   1.473e+06  2.464e+04  59.791 < 2e-16 ***    ## radius3        -6.130e+06  2.143e+05 -28.606 < 2e-16 ***
## area1       -1.301e+05  3.907e+03 -33.301 < 2e-16 ***    ## texture3       -5.832e+05  2.437e+04 -23.935 < 2e-16 ***
## smoothness1 -1.525e+08  8.361e+06 -18.234 < 2e-16 ***    ## perimeter3     -3.538e+05  1.219e+04 -29.023 < 2e-16 ***
## compactness1 -6.428e+06  3.213e+06  -2.001  0.04539 *    ## area3          8.950e+04  2.741e+03  32.658 < 2e-16 ***
## concavity1   1.042e+06  1.408e+06   0.740  0.45959    ## smoothness3    -2.161e+07  3.298e+06  -6.553 5.65e-11 ***
## concave_points1 -1.716e+07  5.382e+06  -3.188  0.00143 **    ## compactness3   8.986e+06  3.999e+05  22.470 < 2e-16 ***
## symmetry1    4.049e+07  7.772e+05  52.093 < 2e-16 ***    ## concavity3     -3.028e+07  1.523e+06 -19.876 < 2e-16 ***
## fractal_dimension1 -4.233e+07  2.169e+06 -19.519 < 2e-16 ***    ## concave_points3 1.431e+08  5.471e+06  26.162 < 2e-16 ***
## radius2      3.328e+07  1.169e+06  28.478 < 2e-16 ***    ## symmetry3      -2.474e+07  3.392e+05 -72.923 < 2e-16 ***
## texture2     6.368e+06  2.005e+05  31.763 < 2e-16 ***    ## fractal_dimension3 -3.698e+07  5.340e+06 -6.926 4.32e-12 ***
## perimeter2   1.701e+06  4.720e+04  36.032 < 2e-16 ***    ## ---
## area2       -6.393e+05  1.835e+04 -34.840 < 2e-16 ***    ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## smoothness2  7.492e+08  1.224e+07  61.213 < 2e-16 ***    ## (Dispersion parameter for binomial family taken to be 1)
## compactness2 -1.773e+08  5.732e+06 -30.931 < 2e-16 ***    ##
##                                         ## Null deviance:  751.44  on 568  degrees of freedom
##                                         ## Residual deviance: 32006.76  on 538  degrees of freedom
##                                         ## AIC: 32069
##                                         ## Number of Fisher Scoring iterations: 25
```

Figure 3: Result of the base model

Stepwise, Backward and Forward

Firstly, the Stepwise model retained 25 of the 30 original variables, but importantly, all the remaining variables now have very high p-values, meaning they are not statistically significant. This is unusual and suggests potential overfitting or numerical issues despite the improved deviance.

```
##
## Call:
## NULL
##
## Coefficients:
##          Estimate Std. Error z value Pr(>|z|)
## (Intercept) -5.914e+03  2.619e+05 -0.023  0.982
## radius1     -6.630e+03  1.150e+05 -0.058  0.954
## texture1     1.913e+02  1.345e+03  0.142  0.887
## area1        6.077e+01  1.079e+03  0.056  0.955
## smoothness1  3.914e+04  2.517e+05  0.155  0.876
## compactness1 -8.621e+04  9.326e+05 -0.092  0.926
## concavity1   2.852e+04  2.402e+05  0.119  0.905
## concave_points1 5.886e+04  1.544e+06  0.038  0.970
## symmetry1    -1.964e+04  1.347e+05 -0.146  0.884
## fractal_dimension1 1.626e+05  1.120e+06  0.145  0.885
## perimeter2   -1.253e+03  1.822e+04 -0.069  0.945
## area2        1.562e+02  2.259e+03  0.069  0.945
## smoothness2  -9.793e+04  1.472e+06 -0.067  0.947
## compactness2  9.217e+04  7.142e+05  0.129  0.897

## concavity2      -8.131e+04  1.097e+06 -0.074  0.941
## concave_points2 4.398e+05  6.736e+06  0.065  0.948
## symmetry2      -1.038e+05  2.160e+06 -0.048  0.962
## fractal_dimension2 -1.092e+06  1.065e+07 -0.103  0.918
## radius3        2.226e+03  2.134e+04  0.104  0.917
## texture3       7.269e+01  3.150e+03  0.023  0.982
## perimeter3     1.267e+02  1.355e+03  0.093  0.926
## area3         -1.626e+01  1.165e+02 -0.140  0.889
## concavity3     6.737e+03  1.051e+05  0.064  0.949
## symmetry3      2.201e+04  3.283e+05  0.067  0.947
## fractal_dimension3 5.899e+04  1.032e+06  0.057  0.954

## (Dispersion parameter for binomial family taken to be 1)
## Null deviance: 7.5144e+02 on 568 degrees of freedom
## Residual deviance: 1.6713e+04 on 544 degrees of freedom
## AIC: 50
##
## Number of Fisher Scoring iterations: 25
```

Figure 4: Result of the stepwise model

Secondly, the Backward model, all retained variables have very high p-values, meaning none are statistically significant despite being kept in the model. The model required 25 iterations to converge, and the coefficient values remain in the thousands range, which is smaller than the base model but still relatively large. This suggests that backward elimination arrived at the same final model as stepwise selection, but the lack of significant p-values indicates potential overfitting or numerical instability issues.

```
##
## Call:
## NULL
##
## Coefficients:
##          Estimate Std. Error z value Pr(>|z|)
## (Intercept) -5.914e+03  2.619e+05 -0.023  0.982
## radius1     -6.630e+03  1.150e+05 -0.058  0.954
## texture1     1.913e+02  1.345e+03  0.142  0.887
## area1        6.077e+01  1.079e+03  0.056  0.955
## smoothness1  3.914e+04  2.517e+05  0.155  0.876
## compactness1 -8.621e+04  9.326e+05 -0.092  0.926
## concavity1   2.852e+04  2.402e+05  0.119  0.905
## concave_points1 5.886e+04  1.544e+06  0.038  0.970
## symmetry1    -1.964e+04  1.347e+05 -0.146  0.884
## fractal_dimension1 1.626e+05  1.120e+06  0.145  0.885
## perimeter2   -1.253e+03  1.822e+04 -0.069  0.945
## area2        1.562e+02  2.259e+03  0.069  0.945
## smoothness2  -9.793e+04  1.472e+06 -0.067  0.947
## compactness2  9.217e+04  7.142e+05  0.129  0.897

## concavity2      -8.131e+04  1.097e+06 -0.074  0.941
## concave_points2 4.398e+05  6.736e+06  0.065  0.948
## symmetry2      -1.038e+05  2.160e+06 -0.048  0.962
## fractal_dimension2 -1.092e+06  1.065e+07 -0.103  0.918
## radius3        2.226e+03  2.134e+04  0.104  0.917
## texture3       7.269e+01  3.150e+03  0.023  0.982
## perimeter3     1.267e+02  1.355e+03  0.093  0.926
## area3         -1.626e+01  1.165e+02 -0.140  0.889
## concavity3     6.737e+03  1.051e+05  0.064  0.949
## symmetry3      2.201e+04  3.283e+05  0.067  0.947
## fractal_dimension3 5.899e+04  1.032e+06  0.057  0.954

## (Dispersion parameter for binomial family taken to be 1)
## Null deviance: 7.5144e+02 on 568 degrees of freedom
## Residual deviance: 1.6713e+04 on 544 degrees of freedom
## AIC: 50
##
## Number of Fisher Scoring iterations: 25
```

Figure 5: Result of the backward model

Finally, the Forward model converged in only 11 iterations (compared to 25 for previous models) and has much smaller, more reasonable coefficient values, suggesting better numerical stability. Overall, forward selection produced a simpler, more interpretable model with genuinely significant predictors, though with a slightly lower fit than stepwise/backward.

```
##
## Call:
## NULL
##
## Coefficients:
##          Estimate Std. Error z value Pr(>|z|)
## (Intercept) -44.28703  11.47044  -3.861 0.000113 ***
## perimeter3   0.08387   0.07434   1.128 0.259268
## smoothness3 56.53029  30.08184   1.879 0.060215 .
## texture3     0.53014   0.13654   3.883 0.000103 ***
## radius2     -9.94253  16.81165  -0.591 0.554248
## symmetry3    16.97161   8.01255   2.118 0.034164 *
## compactness2 -130.42302  40.78433  -3.198 0.001384 **
## concavity1   39.87363  16.30611   2.445 0.014472 *

## texture2      -2.68316   1.51555  -1.770 0.076657 .
## area2         0.35067   0.20024   1.751 0.079892 .
## concave_points3 36.51134  19.42158   1.880 0.060117 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

## (Dispersion parameter for binomial family taken to be 1)
## Null deviance: 751.440 on 568 degrees of freedom
## Residual deviance: 54.713 on 558 degrees of freedom
## AIC: 76.713
##
## Number of Fisher Scoring iterations: 11
```

Figure 6: Result of the Forward model

4.2 LASSO

This is the simplest model so far, with LASSO effectively identifying the most important predictors while eliminating redundant or less informative variables, making it highly interpretable and efficient for prediction.

```
## 31 x 1 sparse Matrix of class "dgCMatrix"
##              s1 ## smoothness2      .
## (Intercept)  -14.96636168 ## compactness2      .
## radius1      .      ## concavity2      .
## texture1     .      ## concave_points2      .
## perimeter1   .      ## symmetry2      .
## area1        .      ## fractal_dimension2      .
## smoothness1  .      ## radius3      0.43324388
## compactness1 .      ## texture3      0.11076710
## concavity1   .      ## perimeter3      .
## concave_points1 13.57467042 ## area3      .
## symmetry1     .      ## smoothness3      6.03826173
## fractal_dimension1 .      ## compactness3      .
## radius2      0.78991905 ## concavity3      0.00699303
## texture2     .      ## concave_points3      16.78847079
## perimeter2   .      ## symmetry3      2.25946856
## area2        .      ## fractal_dimension3      .
```

Figure 7: Result of the LASSO model

Ridge

While Ridge is more complex than LASSO (using all 30 variables instead of just 9), the regularization penalty prevents overfitting by keeping coefficients small, resulting in a stable and well-behaved model that balances complexity with predictive performance.

```
## 31 x 1 sparse Matrix of class "dgCMatrix"
##              s1 ## smoothness2      3.463722e-01
## (Intercept)  -1.389314e+01 ## compactness2      -4.442906e+00
## radius1      8.081935e-02  ## concavity2      -1.176142e+00
## texture1     6.184018e-02  ## concave_points2      1.340849e+01
## perimeter1   1.152165e-02  ## symmetry2      -8.259528e+00
## area1        7.557894e-04  ## fractal_dimension2 -4.889345e+01
## smoothness1  8.107265e+00  ## radius3      6.990693e-02
## compactness1 1.603673e+00  ## texture3      5.460557e-02
## concavity1   2.801163e+00  ## perimeter3     9.566350e-03
## concave_points1 7.519044e+00 ## area3      5.233578e-04
## symmetry1     2.671773e+00  ## smoothness3     1.085056e+01
## fractal_dimension1 -2.010059e+01 ## compactness3     9.892990e-01
## radius2      8.894382e-01  ## concavity3      1.161959e+00
## texture2     -2.790517e-02  ## concave_points3     5.067763e+00
## perimeter2   9.921854e-02  ## symmetry3      3.876396e+00
## area2        4.526886e-03  ## fractal_dimension3 5.047323e+00
```

Figure 8: Result of the Ridge model

Random forest

The model automatically handles variable interactions and nonlinear relationships without requiring manual feature selection or dealing with multicollinearity issues. Random forest's ability to rank variable importance provides interpretable insights into which cell characteristics are most predictive of malignant versus benign tumors, making it a powerful complement to the regression-based approaches.

4.2.1 Parameters Comparison

We can note that the Lasso and Forward models are the simplest, with much smaller coefficient values. The Ridge model is slightly more complex, but its coefficient values are also very small. The Backward, Stepwise, and especially the Full model are more complex, and their coefficients are much larger.

The table with all the parameters is as follows:

##	variable	Full_Model	Stepwise	Backward	Forward	Lasso	Ridge
## 1	(Intercept)	-2881304.4	-5914.46	-5914.46	-44.29	-14.97	-13.89
## 2	radius1	2427013.0	-6629.52	-6629.52	-	-	0.08
## 3	texture1	195783.2	191.26	191.26	-	-	0.06
## 4	perimeter1	1473188.4	-	-	-	-	0.01
## 5	area1	-130118.2	60.77	60.77	-	-	-
## 6	smoothness1	-152452270.8	39142.9	39142.9	-	-	8.11
## 7	compactness1	-6428388.9	-86211.47	-86211.47	-	-	1.6
## 8	concavity1	1041553.7	28520.95	28520.95	39.87	-	2.8
## 9	concave_points1	-17156866.7	58858.72	58858.72	-	13.57	7.52
## 10	symmetry1	40485942.1	-19644.82	-19644.82	-	-	2.67
## 11	fractal_dimension1	-42329195.5	162556.55	162556.55	-	-	-20.1
## 12	radius2	33284830.2	-	-	-9.94	0.79	0.89
## 13	texture2	6368395.0	-	-	-2.68	-	-0.03

Figure 9: Parameters adjusted for each parametric model.

## 18	concavity2	152864835.6	-81310.72	-81310.72	-	-	-1.18
## 19	concave_points2	-1259854447.6	439795.13	439795.13	-	-	13.41
## 20	symmetry2	289010997.2	-103758.08	-103758.08	-	-	-8.26
## 21	fractal_dimension2	1512104102.9	-1092014.08	-1092014.08	-	-	-48.89
## 22	radius3	-6130234.0	2226.39	2226.39	-	0.43	0.07
## 23	texture3	-583246.0	72.69	72.69	0.53	0.11	0.05
## 24	perimeter3	-353820.5	126.66	126.66	0.08	-	0.01
## 25	area3	89504.5	-16.26	-16.26	-	-	-
## 26	smoothness3	-21611286.5	-	-	56.53	6.04	10.85
## 27	compactness3	8986340.5	-	-	-	-	0.99
## 28	concavity3	-30279288.9	6736.58	6736.58	-	0.01	1.16
## 29	concave_points3	143130487.1	-	-	36.51	16.79	5.07
## 30	symmetry3	-24735907.8	22008.73	22008.73	16.97	2.26	3.88
## 31	fractal_dimension3	-36983257.4	58988.94	58988.94	-	-	5.05

Figure 10: Parameters adjusted for each parametric model. (continued)

4.2.2 Accuracy results of the first repetitions

By obtaining the Accuracy of the 5 folds of each model, we can analyze its dispersion.

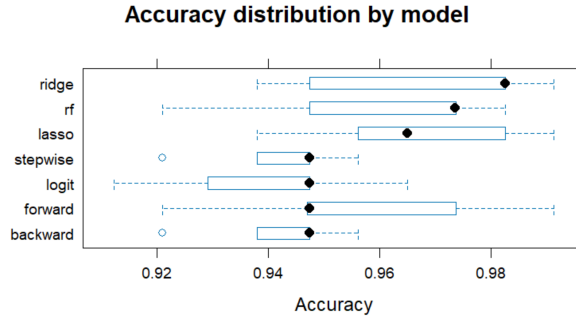


Figure 11: Parameters adjusted for each parametric model.

This chart shows a statistical summary of the accuracy of the 5 folds used.

The box shows the first and third quartiles; the black dot is the median. Points outside the whiskers are outliers.

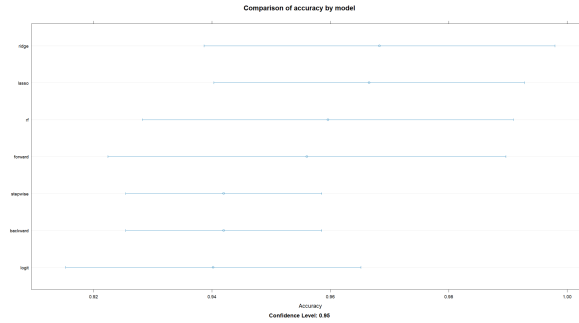


Figure 12: Parameters adjusted for each parametric model.

In this graph, the dot shows the average accuracy across the 5 folds for each model. The line around it shows the dispersion in accuracy. Shorter lines indicate a more consistent metric, while longer lines indicate less consistency and greater variability.

4.2.3 Training times for the first repetition

Note in 13 that the most accurate models are Ridge, followed by Lasso, and then Random Forest. Ridge's training time is slightly faster than Lasso, but Random Forest takes almost six times longer to be adjusted. Recall that in terms of complexity, the Lasso model was the simplest, but Ridge, despite being more complex, wasn't as complex as the other models.

The time required for stepwise, backward, and forward models is significantly longer due to the iterative process of removing or adding variables and checking their significance for the next step.

Similarly, the random forest model takes slightly longer due to its nature; the data subsetting process is time-consuming and is performed multiple times for different trees.

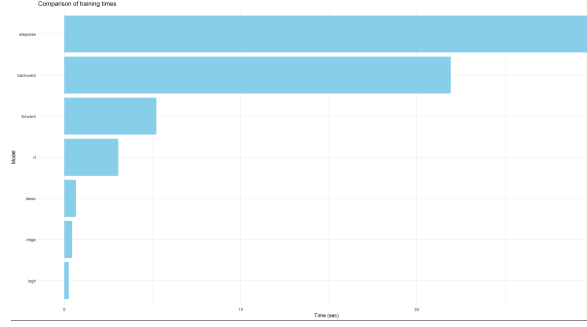


Figure 13: Time to adjust the different models in the first repetition.

4.2.4 Analysis of overall performance after a single repetition.

Below is a table summarizing the results of the first repetition.

##	Model	Accuracy_Mean	Time
## 1	logit	0.9402267	0.23
## 2	stepwise	0.9419966	29.74
## 3	backward	0.9419966	21.60
## 4	forward	0.9560472	5.00
## 5	lasso	0.9665580	0.61
## 6	ridge	0.9683124	0.35
## 7	rf	0.9596025	3.27

Figure 14: Summary of metrics for the first repetition.

Note that the most accurate models are Ridge, followed by Lasso, and then Random Forest. Ridge's training time is slightly faster than Lasso, but Random Forest takes almost six times longer to be adjusted. Recall that in terms of complexity, the Lasso model was the simplest, but Ridge, despite being more complex, was not as complex as the other models.

We can conclude that the Lasso model is the best option, as it offers simplicity, good accuracy, and a short adjustment time. If we are willing to sacrifice simplicity for a shorter training time and greater accuracy, Ridge is the best choice. Random Forest, being a machine learning model and therefore much less interpretable, has similar accuracy to the other models, which, while not bad, are inferior to Lasso and Ridge. However, it is faster and more accurate than the other models. The other models are not worth it in this situation; they have a significantly longer training time, their accuracy is worse, and they are much more complex.

4.3 Results after 50 repetitions

To obtain a more reliable assessment of model performance, we repeated the entire 5-fold cross-validation procedure 50 times. A single cross-validation split can introduce variability because the specific assignment of observations to folds may favor or disadvantage certain models. By repeating the process with different random fold partitions and averaging the results, we reduce the influence of any one particular split and obtain a more stable estimate of both accuracy and Brier score. This repetition also allows us to quantify the variability of each model's performance on different data partitions, providing a more robust comparison between methods.

Due to the time it takes for them to adjust, we will not be using the Backward and Stepwise models. These models did not demonstrate outstanding performance, so we decided to forget them at this stage.

4.3.1 Accuracy results after 50 repetitions

The box plot for the 50 repetitions shows that Ridge achieves the highest accuracy overall, followed by LASSO. Random Forest also performs well but is comparable to the Forward model and falls behind Ridge and LASSO.

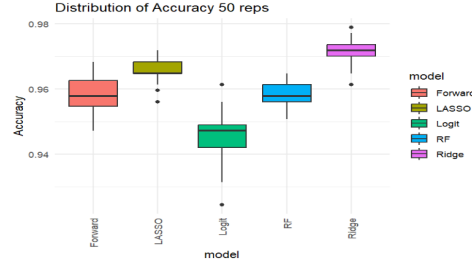


Figure 15: Accuracy summary after 50 repetitions.

4.3.2 Brier Score results after 50 repetitions

The Brier Score box plot indicates that all models achieve low error values, with Ridge and LASSO showing the lowest scores. Their box plots are also narrower, suggesting more consistent performance compared with the other models.

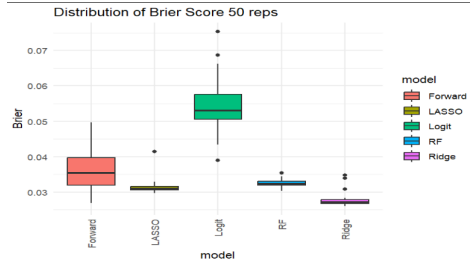


Figure 16: Brier Score summary after 50 repetitions.

5 Discussion

5.1 Interpretation and comparison of models

The overall results across the 50 repetitions show a consistent pattern in model performance. Ridge regression achieves the highest average accuracy, closely followed by LASSO, while Random Forest performs well but shows a larger spread in its predictive performance. This trend is also reflected in the Brier Score, where Ridge and LASSO obtain the smallest values and the tightest boxplots, indicating not only strong predictive accuracy but also high consistency across resampling. The stochastic nature of Random Forest, combined with the variability in bootstrap samples, likely explains why its performance exhibits greater dispersion.

5.2 LASSO Variable Selection Analysis

The drastic reduction in the choice of parameters by LASSO and the very different scale between the different estimated parameters of the different models lead us to analyze the correlation in the database.

A heat map showing the correlation of the explanatory variables will be displayed below.

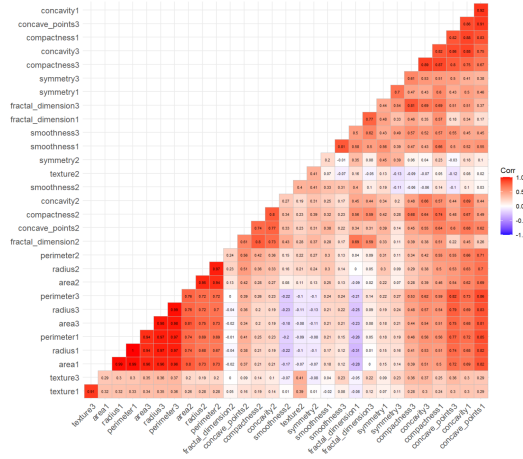


Figure 17: Heat map of the correlation in predictor variables.

There are too many variables, but we can easily notice that most of the squares are painted red, which indicates a strong correlation between the different variables in our data set.

This highly correlated data creates a multicollinearity problem in the logistic model that uses all the variables. This makes the model unstable and requires large coefficients to balance the repeated effects explained by other variables.

It is surprising that the base model has such high accuracy given its clear multicollinearity issues. But how does Lasso achieve slightly better precision without using so many variables and with much smaller coefficients?

By averaging the correlation of each variable with the others, we obtain the 10 variables with the highest average correlation and the 10 with the lowest average correlation. Comparing variables with those used by LASSO, we can note that 4 of the 10 most correlated variables are used by Lasso: concave_points1, radius3, concavity3 and concave_points3. Lasso, in turn, uses 3 of the variables with the lowest average absolute correlation: texture3, smoothness3 and symmetry3.

And finally there is the variable radius 2, which remains in the variables with intermediate average absolute correlation, specifically it is the 14th ordered from highest to lowest.

We can see that Lasso uses a combination of highly correlated variables, with poorly correlated variables and an intermediate variable.

The heat map of the variables used by Lasso is shown below.

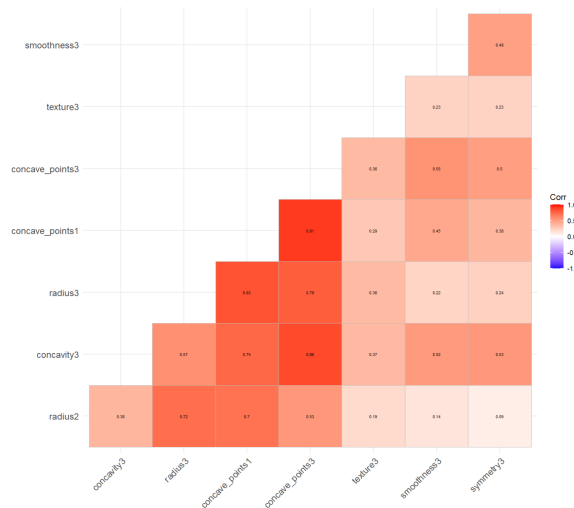


Figure 18: Heat map of the correlation in predictor variables used by LASSO.

We can see that variables with low correlation are in more orange or white tones, while variables with higher correlation tend to be more red.

This choice of variables is not accidental; in fact, it is quite logical. This behavior aligns with expectations for penalized regression: highly correlated variables tend to carry redundant information, and LASSO absorbs that shared signal by selecting only a few representatives. On the other hand, including some weakly correlated predictors may help the model capture additional, more specific patterns in the data that strong correlations do not explain. This suggests that the selected subset efficiently balances redundancy reduction with information specificity.

5.3 Variability and consistency of models

Repeating the full model fitting process 50 times offers insight into the stability of the algorithms. While a single cross-validation run can give a reasonable estimate of performance, it may be influenced by how the folds were partitioned. By varying the fold assignments across many repetitions, we gain a clearer view of the distribution of metrics and reduce the risk of over-interpreting results that may be tied to a particular random split. The narrow distributions of Ridge and LASSO demonstrate that their performance is robust to changes in resampling, whereas models like Logistic Regression and Forward selection exhibit more sensitivity.

6 Conclusion

The repeated evaluation of models across 50 iterations provides a robust assessment of predictive performance and variability. Ridge regression consistently achieves the highest accuracy and lowest Brier scores, indicating not only precise predictions but also stable performance across different data splits. LASSO performs similarly well, slightly below Ridge in accuracy, but with the added benefit of automatic variable selection, reducing the number of predictors while capturing the most informative variables. This behavior illustrates how penalization allows LASSO to summarize correlated information into fewer features while retaining predictive power. Forward selection and Random Forest also show strong predictive performance; however, they exhibit slightly higher variability across iterations, with Random Forest variability stemming from its stochastic nature due to bootstrap sampling and random feature selection at each split. Logistic regression without selection shows good average performance but lacks the refinement in variable selection and consistency offered by Ridge and LASSO. Overall, these results emphasize the importance of considering both predictive accuracy and model interpretability, suggesting that Ridge is ideal when stability and precision are the priority, while LASSO provides a practical trade-off between accuracy and feature reduction.

Future work could include validating these models on external datasets to assess generalizability, exploring Elastic Net to combine the strengths of LASSO and Ridge, and performing more extensive hyperparameter tuning for tree-based models like Random Forest. Such approaches could enhance model reliability, improve prediction stability, and provide deeper insights into the variable relationships driving tumor classification.

Appendix A

The Rmd document containing the code ready to be executed in R, as well as its compiled PDF version and the database used, can be found in the following GitHub repository: <https://github.com/Stephanie-Daniella/Breast-Tumors>

Appendix B: Code used

Stephanie Daniella Hernandez Prado

2025-12-13

```
library(readxl)
```

```
## Warning: package 'readxl' was built under R version 4.4.3
```

```
library(caret)
```

```
## Warning: package 'caret' was built under R version 4.4.3
```

```
## Cargando paquete requerido: ggplot2
```

```
## Warning: package 'ggplot2' was built under R version 4.4.3
```

```
## Cargando paquete requerido: lattice
```

```
library(randomForest)
```

```
## Warning: package 'randomForest' was built under R version 4.4.3
```

```
## randomForest 4.7-1.2
```

```
## Type rfNews() to see new features/changes/bug fixes.
```

```
##
```

```
## Adjuntando el paquete: 'randomForest'
```

```
## The following object is masked from 'package:ggplot2':
```

```
##
```

```
##     margin
```

```
library(ggplot2)
```

```
library(ggcorrplot)
```

```
## Warning: package 'ggcorrplot' was built under R version 4.4.3
```

```
library(Metrics)
```

```
## Warning: package 'Metrics' was built under R version 4.4.3
```

```
##
```

```
## Adjuntando el paquete: 'Metrics'
```

```
## The following objects are masked from 'package:caret':
```

```
##
```

```
##      precision, recall
```

Data Set

```
# Reading the File
```

```
Breast_Cancer <- read_excel("~/Stephanie/Maestria/Fall 2025/STST 578/Breas_Cancer/Breast_Cancer.xlsx")
```

```
head(Breast_Cancer)
```

```
## # A tibble: 6 x 31
```

```
##   radius1 texture1 perimeter1 area1 smoothness1 compactness1 concavity1
```

```
##   <dbl>    <dbl>      <dbl> <dbl>      <dbl>      <dbl>      <dbl>
```

```
## 1   18.0     10.4      123. 1001      0.118      0.278      0.300
```

```
## 2   20.6     17.8      133. 1326      0.0847     0.0786     0.0869
```

```
## 3   19.7     21.2      130 1203      0.110      0.160      0.197
```

```
## 4   11.4     20.4       77.6 386.      0.142      0.284      0.241
```

```
## 5   20.3     14.3      135. 1297      0.100      0.133      0.198
```

```
## 6   12.4     15.7      82.6 477.      0.128      0.17       0.158
```

```
## # i 24 more variables: concave_points1 <dbl>, symmetry1 <dbl>,
```

```
## #   fractal_dimension1 <dbl>, radius2 <dbl>, texture2 <dbl>, perimeter2 <dbl>,
```

```
## #   area2 <dbl>, smoothness2 <dbl>, compactness2 <dbl>, concavity2 <dbl>,
```

```
## #   concave_points2 <dbl>, symmetry2 <dbl>, fractal_dimension2 <dbl>,
```

```
## #   radius3 <dbl>, texture3 <dbl>, perimeter3 <dbl>, area3 <dbl>,
```

```
## #   smoothness3 <dbl>, compactness3 <dbl>, concavity3 <dbl>,
```

```
## #   concave_points3 <dbl>, symmetry3 <dbl>, fractal_dimension3 <dbl>, ...
```

```
data <- Breast_Cancer
```

```
data$Diagnosis <- as.factor(data$Diagnosis) #Changeing response variable to a factor variable
```

K-folds CV set up

```
set.seed(996) #Save the seed for replicability
```

```
folds <- createFolds(data$Diagnosis, k = 5, returnTrain = TRUE) #Set the 5 folds that will  
#be used for k-fold CV
```

```
Method <- trainControl(method = "cv", index = folds) #We set the K-folds CV method to fit  
#the models
```

Base Model

```
suppressWarnings({ #Supress Warning Messages

  Base_Time <- system.time({ # Count the execution time
    Base_model <- train(
      Diagnosis ~ ., data = data, #model and data used
      method = "glm", #Logistic model setting
      family = binomial, #Binomial response
      trControl = Method #CV setting
    ) #Logistic Model Fitting

  })["elapsed"]
})

summary(Base_model$finalModel)
```



```
##
## Call:
## NULL
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -2.881e+06  2.816e+05 -10.233 < 2e-16 ***
## radius1      2.427e+06  2.693e+05   9.014 < 2e-16 ***
## texture1     1.958e+05  1.471e+04  13.313 < 2e-16 ***
## perimeter1   1.473e+06  2.464e+04  59.791 < 2e-16 ***
## area1        -1.301e+05  3.907e+03 -33.301 < 2e-16 ***
## smoothness1  -1.525e+08  8.361e+06 -18.234 < 2e-16 ***
## compactness1 -6.428e+06  3.213e+06  -2.001 0.04539 *
## concavity1    1.042e+06  1.408e+06   0.740 0.45959
## concave_points1 -1.716e+07  5.382e+06  -3.188 0.00143 **
## symmetry1     4.049e+07  7.772e+05  52.093 < 2e-16 ***
## fractal_dimension1 -4.233e+07  2.169e+06 -19.519 < 2e-16 ***
## radius2       3.328e+07  1.169e+06  28.478 < 2e-16 ***
## texture2      6.368e+06  2.005e+05  31.763 < 2e-16 ***
## perimeter2    1.701e+06  4.720e+04  36.032 < 2e-16 ***
## area2        -6.393e+05  1.835e+04 -34.840 < 2e-16 ***
## smoothness2   7.492e+08  1.224e+07  61.213 < 2e-16 ***
## compactness2  -1.773e+08  5.732e+06 -30.931 < 2e-16 ***
## concavity2    1.529e+08  5.340e+06  28.624 < 2e-16 ***
## concave_points2 -1.260e+09  4.012e+07 -31.398 < 2e-16 ***
## symmetry2     2.890e+08  4.126e+06  70.055 < 2e-16 ***
## fractal_dimension2 1.512e+09  6.597e+07  22.921 < 2e-16 ***
## radius3       -6.130e+06  2.143e+05 -28.606 < 2e-16 ***
## texture3      -5.832e+05  2.437e+04 -23.935 < 2e-16 ***
## perimeter3    -3.538e+05  1.219e+04 -29.023 < 2e-16 ***
## area3         8.950e+04  2.741e+03  32.658 < 2e-16 ***
## smoothness3   -2.161e+07  3.298e+06  -6.553 5.65e-11 ***
## compactness3   8.986e+06  3.999e+05  22.470 < 2e-16 ***
## concavity3    -3.028e+07  1.523e+06 -19.876 < 2e-16 ***
## concave_points3 1.431e+08  5.471e+06  26.162 < 2e-16 ***
```

```
## symmetry3          -2.474e+07  3.392e+05 -72.923 < 2e-16 ***
## fractal_dimension3 -3.698e+07  5.340e+06  -6.926 4.32e-12 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 751.44 on 568 degrees of freedom
## Residual deviance: 32006.76 on 538 degrees of freedom
## AIC: 32069
##
## Number of Fisher Scoring iterations: 25
```

We can see that all variables are being used, even if some are not significant, as they have a very low p-value.

Stepwise

```
suppressWarnings({
  Stepwise_Time <- system.time({
    Stepwise_Model <- train(
      Diagnosis ~ ., data = data,
      method = "glmStepAIC",
      trControl = Method,
      direction = "both", #Set stepwise method
      trace = FALSE # Suppress training trace
    )
  })["elapsed"]
})
summary(Stepwise_Model$finalModel)
```

```
##
## Call:
## NULL
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -5.914e+03  2.619e+05  -0.023  0.982
## radius1      -6.630e+03  1.150e+05  -0.058  0.954
## texture1      1.913e+02  1.345e+03   0.142  0.887
## area1         6.077e+01  1.079e+03   0.056  0.955
## smoothness1   3.914e+04  2.517e+05   0.155  0.876
## compactness1  -8.621e+04  9.326e+05  -0.092  0.926
## concavity1     2.852e+04  2.402e+05   0.119  0.905
## concave_points1 5.886e+04  1.544e+06   0.038  0.970
## symmetry1     -1.964e+04  1.347e+05  -0.146  0.884
## fractal_dimension1 1.626e+05  1.120e+06   0.145  0.885
## perimeter2    -1.253e+03  1.822e+04  -0.069  0.945
## area2         1.562e+02  2.259e+03   0.069  0.945
## smoothness2   -9.793e+04  1.472e+06  -0.067  0.947
## compactness2    9.217e+04  7.142e+05   0.129  0.897
## concavity2    -8.131e+04  1.097e+06  -0.074  0.941
```

```
## concave_points2      4.398e+05  6.736e+06   0.065   0.948
## symmetry2           -1.038e+05  2.160e+06  -0.048   0.962
## fractal_dimension2 -1.092e+06  1.065e+07  -0.103   0.918
## radius3             2.226e+03  2.134e+04   0.104   0.917
## texture3            7.269e+01  3.150e+03   0.023   0.982
## perimeter3          1.267e+02  1.355e+03   0.093   0.926
## area3              -1.626e+01  1.165e+02  -0.140   0.889
## concavity3          6.737e+03  1.051e+05   0.064   0.949
## symmetry3           2.201e+04  3.283e+05   0.067   0.947
## fractal_dimension3  5.899e+04  1.032e+06   0.057   0.954
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 7.5144e+02  on 568  degrees of freedom
## Residual deviance: 1.6713e-04  on 544  degrees of freedom
## AIC: 50
##
## Number of Fisher Scoring iterations: 25
```

With a quick inspection, we can see that the p-values of this model are high for all the variables used, so they are highly significant.

Backward

```
suppressWarnings({
  Backward_Time <- system.time({
    Backward_Model <- train(
      Diagnosis ~ ., data = data,
      method = "glmStepAIC",
      trControl = Method,
      direction = "backward", # Set backward method
      trace = FALSE
    )
  })["elapsed"]
})
summary(Backward_Model$finalModel)
```

```
##
## Call:
## NULL
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -5.914e+03  2.619e+05  -0.023   0.982
## radius1      -6.630e+03  1.150e+05  -0.058   0.954
## texture1      1.913e+02  1.345e+03   0.142   0.887
## area1         6.077e+01  1.079e+03   0.056   0.955
## smoothness1   3.914e+04  2.517e+05   0.155   0.876
## compactness1 -8.621e+04  9.326e+05  -0.092   0.926
## concavity1    2.852e+04  2.402e+05   0.119   0.905
## concave_points1 5.886e+04  1.544e+06   0.038   0.970
```



```
## symmetry1      -1.964e+04  1.347e+05  -0.146   0.884
## fractal_dimension1  1.626e+05  1.120e+06   0.145   0.885
## perimeter2      -1.253e+03  1.822e+04  -0.069   0.945
## area2           1.562e+02  2.259e+03   0.069   0.945
## smoothness2     -9.793e+04  1.472e+06  -0.067   0.947
## compactness2     9.217e+04  7.142e+05   0.129   0.897
## concavity2      -8.131e+04  1.097e+06  -0.074   0.941
## concave_points2  4.398e+05  6.736e+06   0.065   0.948
## symmetry2       -1.038e+05  2.160e+06  -0.048   0.962
## fractal_dimension2 -1.092e+06  1.065e+07  -0.103   0.918
## radius3          2.226e+03  2.134e+04   0.104   0.917
## texture3         7.269e+01  3.150e+03   0.023   0.982
## perimeter3       1.267e+02  1.355e+03   0.093   0.926
## area3            -1.626e+01  1.165e+02  -0.140   0.889
## concavity3        6.737e+03  1.051e+05   0.064   0.949
## symmetry3         2.201e+04  3.283e+05   0.067   0.947
## fractal_dimension3 5.899e+04  1.032e+06   0.057   0.954
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 7.5144e+02  on 568  degrees of freedom
## Residual deviance: 1.6713e-04  on 544  degrees of freedom
## AIC: 50
##
## Number of Fisher Scoring iterations: 25
```

With a quick inspection, we can see that the p-values of this model are high for all the variables used, so they are highly significant.

Forward

```
suppressWarnings({
  Forward_Time <- system.time({
    Forward_Model <- train(
      Diagnosis ~ ., data = data,
      method = "glmStepAIC",
      trControl = Method,
      direction = "forward", # Set forward method
      trace = FALSE
    )
  })["elapsed"]
})
summary(Forward_Model$finalModel)

##
## Call:
## NULL
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -44.28703    11.47044  -3.861 0.000113 ***
```

```
## perimeter3      0.08387      0.07434      1.128 0.259268
## smoothness3    56.53029    30.08184      1.879 0.060215 .
## texture3       0.53014      0.13654      3.883 0.000103 ***
## radius2       -9.94253    16.81165     -0.591 0.554248
## symmetry3      16.97161      8.01255      2.118 0.034164 *
## compactness2  -130.42302    40.78433     -3.198 0.001384 **
## concavity1     39.87363    16.30611      2.445 0.014472 *
## texture2      -2.68316      1.51555     -1.770 0.076657 .
## area2          0.35067      0.20024      1.751 0.079892 .
## concave_points3 36.51134    19.42158      1.880 0.060117 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 751.440  on 568  degrees of freedom
## Residual deviance:  54.713  on 558  degrees of freedom
## AIC: 76.713
##
## Number of Fisher Scoring iterations: 11
```

With a quick inspection, we can see that the p-values of this model are high for all the variables used, so they are highly significant.

Lasso

```
suppressWarnings({
  Lasso_Time <- system.time({
    Lasso_Model <- train(
      Diagnosis ~ ., data = data,
      method = "glmnet",
      trControl = Method,
      tuneGrid = expand.grid(alpha = 1,
                             lambda = seq(0.0001, 1, length = 50)) # alpha=1 means lasso
                             #model and we fit for different lambdas
    )
  })["elapsed"]
})

best_lambda_lasso <- Lasso_Model$bestTune$lambda #Get the lambda of the best model

coef(Lasso_Model$finalModel, s = best_lambda_lasso) #Coef of the model with the best lambda

## 31 x 1 sparse Matrix of class "dgCMatrix"
##              s1
## (Intercept)  -14.96636168
## radius1      .
## texture1     .
## perimeter1   .
## area1        .
## smoothness1  .
```

```
## compactness1      .
## concavity1        .
## concave_points1   13.57467042
## symmetry1         .
## fractal_dimension1 .
## radius2           0.78991905
## texture2          .
## perimeter2        .
## area2             .
## smoothness2       .
## compactness2      .
## concavity2        .
## concave_points2   .
## symmetry2         .
## fractal_dimension2 .
## radius3           0.43324388
## texture3          0.11076710
## perimeter3        .
## area3             .
## smoothness3       6.03826173
## compactness3      .
## concavity3        0.00699303
## concave_points3   16.78847079
## symmetry3         2.25946856
## fractal_dimension3 .
```

With the Lasso penalty, using the best lambda, we can see how too many variables were discarded; now we need to compare its performance and see if this much simpler model is at least as good as the others.

Ridge

```
suppressWarnings({
  Ridge_Time <- system.time({
    Ridge_Model <- train(
      Diagnosis ~ ., data = data,
      method = "glmnet",
      trControl = Method,
      tuneGrid = expand.grid(alpha = 0, lambda = seq(0.0001, 1, length = 50)) # alpha=0
      #means Ridge regression and we fit for different lambdas
    )
  })["elapsed"]
})

best_lambda_ridge <- Ridge_Model$bestTune$lambda #Best Ridge lambda

coef(Ridge_Model$finalModel, s = best_lambda_ridge) #Coef of the model with the best lambda

## 31 x 1 sparse Matrix of class "dgCMatrix"
##                               s1
## (Intercept)      -1.389314e+01
## radius1          8.081935e-02
```

```
## texture1          6.184018e-02
## perimeter1        1.152165e-02
## area1             7.557894e-04
## smoothness1       8.107265e+00
## compactness1      1.603673e+00
## concavity1        2.801163e+00
## concave_points1   7.519044e+00
## symmetry1         2.671773e+00
## fractal_dimension1 -2.010059e+01
## radius2           8.894382e-01
## texture2          -2.790517e-02
## perimeter2        9.921854e-02
## area2             4.526886e-03
## smoothness2       3.463722e-01
## compactness2      -4.442906e+00
## concavity2        -1.176142e+00
## concave_points2   1.340849e+01
## symmetry2         -8.259528e+00
## fractal_dimension2 -4.889345e+01
## radius3           6.990693e-02
## texture3          5.460557e-02
## perimeter3        9.566350e-03
## area3             5.233578e-04
## smoothness3       1.085056e+01
## compactness3      9.892990e-01
## concavity3        1.161959e+00
## concave_points3   5.067763e+00
## symmetry3         3.876396e+00
## fractal_dimension3 5.047323e+00
```

Regression Rigge uses many variables, resulting in a more complex model; however, many of the coefficients show very small values. We will have to see if this complexity and the minimal contribution of some variables help its performance.

Random Forest

```
suppressWarnings({
  Random_Forest_Time <- system.time({
    RF_model <- train(
      Diagnosis ~ ., data = data,
      method = "rf", # Train a random forest model
      trControl = Method
    )
  })["elapsed"]
})

varImp(RF_model) # Importance of the 20 most important variables

## rf variable importance
##
## only 20 most important variables shown (out of 30)
```

```
##
## Overall
## area3 100.000
## radius3 95.167
## concave_points3 94.321
## perimeter3 80.770
## concavity1 70.969
## concave_points1 70.938
## radius1 68.643
## area1 53.316
## concavity3 51.989
## perimeter1 49.280
## area2 43.858
## perimeter2 27.614
## compactness3 27.437
## compactness1 25.748
## texture3 20.682
## radius2 20.054
## texture1 14.380
## smoothness3 11.988
## symmetry3 9.938
## concavity2 9.432
```

Table of coefficient comparison

```
#Get the names of all the variables
vars <- names(coef(Base_model$finalModel))

# Data frame with the first column the
#name of all the variables
df_coefs <- data.frame(
  variable = vars,
  stringsAsFactors = FALSE
)

#Function to extract the value of the coef,
#make sure that all the columns have the
# same variables in the same order
extract_coefs <- function(m, vars) {

  co <- m

  # Check if it is a sparse matrix
   #(Lasso and Ridge)
  if ("dgCMatrix" %in% class(co)) {
    #If so, extract the first column
    #that contains the names.
    values <- as.numeric(co[,1])
    names(values) <- rownames(co)
  } else {
    #If not, extract the names of the
```

```

#coefficients.
values <- as.numeric(co)
names(values) <- names(co)
}

return(values[vars])
}

#Extract the coefficient values for each model
df_coefs$Full_Model <- extract_coefs(coef(Base_model$finalModel), vars)

df_coefs$Stepwise <- extract_coefs(coef(Stepwise_Model$finalModel), vars)

df_coefs$Backward <- extract_coefs(coef(Backward_Model$finalModel), vars)

df_coefs$Forward <- extract_coefs(coef(Forward_Model$finalModel), vars)

df_coefs$Lasso <- extract_coefs(
  coef(Lasso_Model$finalModel, s = Lasso_Model$bestTune$lambda),
  vars
)

df_coefs$Ridge <- extract_coefs(
  coef(Ridge_Model$finalModel, s = Ridge_Model$bestTune$lambda),
  vars
)

#Round coefficients to 2 decimal places
df_coefs[, -1] <- round(df_coefs[, -1], 2)

#Replace NA and 0 with -
df_coefs[is.na(df_coefs)|df_coefs==0] <- "-"

df_coefs

```

##	variable	Full_Model	Stepwise	Backward	Forward	Lasso
## 1	(Intercept)	-2881304.4	-5914.46	-5914.46	-44.29	-14.97
## 2	radius1	2427013.0	-6629.52	-6629.52	-	-
## 3	texture1	195783.2	191.26	191.26	-	-
## 4	perimeter1	1473188.4	-	-	-	-
## 5	area1	-130118.2	60.77	60.77	-	-
## 6	smoothness1	-152452270.8	39142.9	39142.9	-	-
## 7	compactness1	-6428388.9	-86211.47	-86211.47	-	-
## 8	concavity1	1041553.7	28520.95	28520.95	39.87	-
## 9	concave_points1	-17156866.7	58858.72	58858.72	-	13.57
## 10	symmetry1	40485942.1	-19644.82	-19644.82	-	-
## 11	fractal_dimension1	-42329195.5	162556.55	162556.55	-	-
## 12	radius2	33284830.2	-	-	-9.94	0.79
## 13	texture2	6368395.0	-	-	-2.68	-
## 14	perimeter2	1700716.8	-1253.11	-1253.11	-	-
## 15	area2	-639346.7	156.21	156.21	0.35	-
## 16	smoothness2	749172456.3	-97931.96	-97931.96	-	-
## 17	compactness2	-177307723.7	92173.51	92173.51	-130.42	-

```

## 18      concavity2    152864835.6   -81310.72   -81310.72      -      -
## 19   concave_points2 -1259854447.6   439795.13   439795.13      -      -
## 20      symmetry2    289010997.2  -103758.08  -103758.08      -      -
## 21 fractal_dimension2 1512104102.9 -1092014.08 -1092014.08      -      -
## 22      radius3     -6130234.0     2226.39     2226.39      -    0.43
## 23      texture3     -583246.0       72.69       72.69     0.53    0.11
## 24      perimeter3   -353820.5      126.66      126.66     0.08      -
## 25      area3        89504.5      -16.26      -16.26      -      -
## 26      smoothness3  -21611286.5          -          -    56.53    6.04
## 27      compactness3   8986340.5          -          -      -      -
## 28      concavity3   -30279288.9     6736.58     6736.58      -    0.01
## 29   concave_points3  143130487.1          -          -    36.51   16.79
## 30      symmetry3   -24735907.8     22008.73     22008.73   16.97    2.26
## 31 fractal_dimension3 -36983257.4     58988.94     58988.94      -      -
##      Ridge
## 1   -13.89
## 2    0.08
## 3    0.06
## 4    0.01
## 5      -
## 6    8.11
## 7    1.6
## 8    2.8
## 9    7.52
## 10   2.67
## 11  -20.1
## 12   0.89
## 13  -0.03
## 14   0.1
## 15      -
## 16   0.35
## 17  -4.44
## 18  -1.18
## 19   13.41
## 20  -8.26
## 21 -48.89
## 22   0.07
## 23   0.05
## 24   0.01
## 25      -
## 26   10.85
## 27   0.99
## 28   1.16
## 29   5.07
## 30   3.88
## 31   5.05

```

We can see that the Lasso and Forward models are the simplest, with much smaller coefficient values. The Ridge model is slightly more complex, but its coefficient values are also very small. The Backward, Stepwise, and especially the Full model are more complex, and their coefficients are much larger.

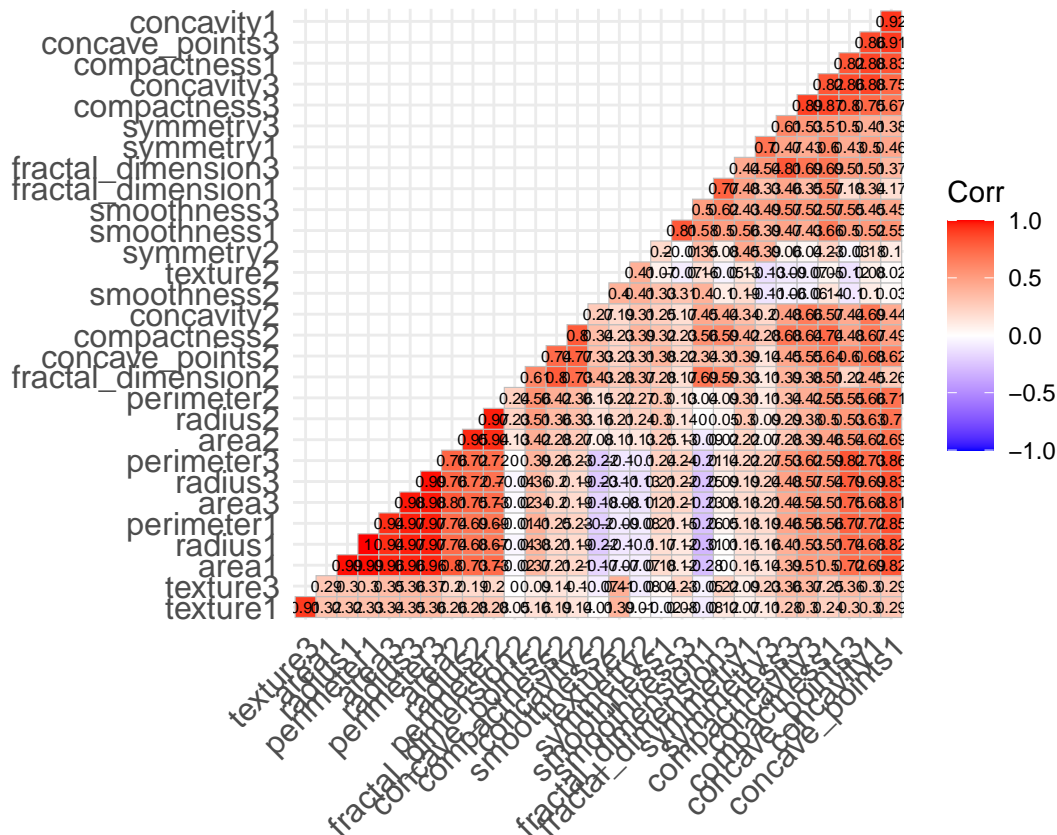
Correlation analysis of variables

A heat map showing the correlation of the explanatory variables will be displayed below.

```
M <- cor(data[sapply(data, is.numeric)])

ggcorrplot(M, hc.order = TRUE, type = "lower",
            lab = TRUE, lab_size = 2)

## Warning: 'aes_string()' was deprecated in ggplot2 3.0.0.
## i Please use tidy evaluation idioms with 'aes()'.
## i See also 'vignette("ggplot2-in-packages")' for more information.
## i The deprecated feature was likely used in the ggcorrplot package.
## Please report the issue at <https://github.com/kassambara/ggcorrplot/issues>.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.
```



There are too many variables, but we can easily notice that most of the squares are painted red, which indicates a strong correlation between the different variables in our data set.

This highly correlated data creates a multicollinearity problem in the logistic model that uses all the variables. This makes the model unstable and requires large coefficients to balance the repeated effects explained by other variables.

It's surprising that the base model has such high accuracy given its clear collinearity issues. But how does Lasso achieve slightly better precision without using so many variables and with much smaller coefficients?

First, let's make a list of the variables with the absolute value lowest average correlation.

```
numeric_data <- data[, sapply(data, is.numeric)]

#Correlation matrix
M <- cor(numeric_data, use = "pairwise.complete.obs")

diag(M) <- NA

#We take the average of the correlations
avg_corr <- apply(abs(M), 1, mean, na.rm = TRUE)

# We ordered from highest to lowest correlation
avg_corr_sorted <- sort(avg_corr, decreasing = TRUE)

most_correlated <- names(avg_corr_sorted)[tail(order(avg_corr_sorted), 10)]

most_correlated
```

```
## [1] "area3"          "compactness3"    "perimeter1"      "radius3"
## [5] "perimeter3"      "concavity3"      "concave_points3" "compactness1"
## [9] "concave_points1" "concavity1"
```

We can see that 4 of the 10 most correlated variables are used by Lasso: concave_points1, radius3, concavity3 and concave_points3. Let's see if we can extract those with the lowest absolute value average correlation.

```
most_correlated <- names(avg_corr_sorted)[20:30]

most_correlated
```

```
## [1] "symmetry1"          "fractal_dimension1" "fractal_dimension3"
## [4] "smoothness3"        "symmetry3"          "fractal_dimension2"
## [7] "texture3"           "texture1"           "smoothness2"
## [10] "symmetry2"          "texture2"
```

We can see that Lasso uses 3 of the variables with the least average correlation in absolute value: texture3, smoothness3 and symmetry3.

And finally there is the variable radius 2, which remains in the variables with intermediate average absolute correlation, specifically it is the 14th ordered from highest to lowest.

We can see that Lasso uses a combination of highly correlated variables, with poorly correlated variables and an intermediate variable.

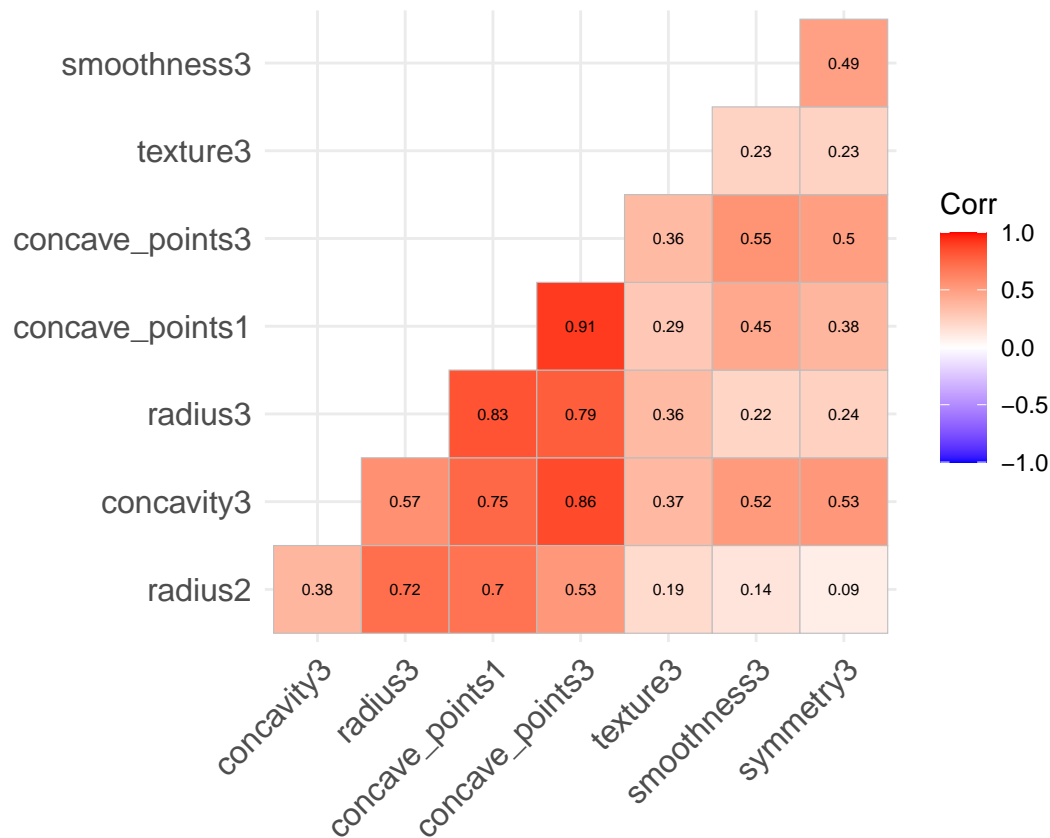
The heat map of the variables used by Lasso is shown below.

```
vars <- c(
  "concave_points1", "radius2", "radius3", "texture3",
  "smoothness3", "concavity3", "concave_points3", "symmetry3"
)

data_sub <- data[, vars]

M <- cor(data_sub[sapply(data_sub, is.numeric)])
```

```
ggcorrplot(M, hc.order = TRUE, type = "lower",
           lab = TRUE, lab_size = 2)
```



We can see that variables with low correlation are in more orange or white tones, while variables with higher correlation tend to be more red.

This choice of variables is not accidental; in fact, it's quite logical. The explanation I can offer is that Lasso, being penalized, uses highly correlated variables to explain the greater amount of information contained in many variables, synthesizing it into a few. Meanwhile, he uses some less correlated variables to explain more specific information not found in other variables.

Accuracy comparison

```
#Statistical summary of the 5-fold CV
#for each model
```

```
results <- resamples(list(
  logit = Base_model,
  stepwise = Stepwise_Model,
  backward = Backward_Model,
  forward = Forward_Model,
  lasso = Lasso_Model,
  ridge = Ridge_Model,
  rf = RF_model
```

```
))
```

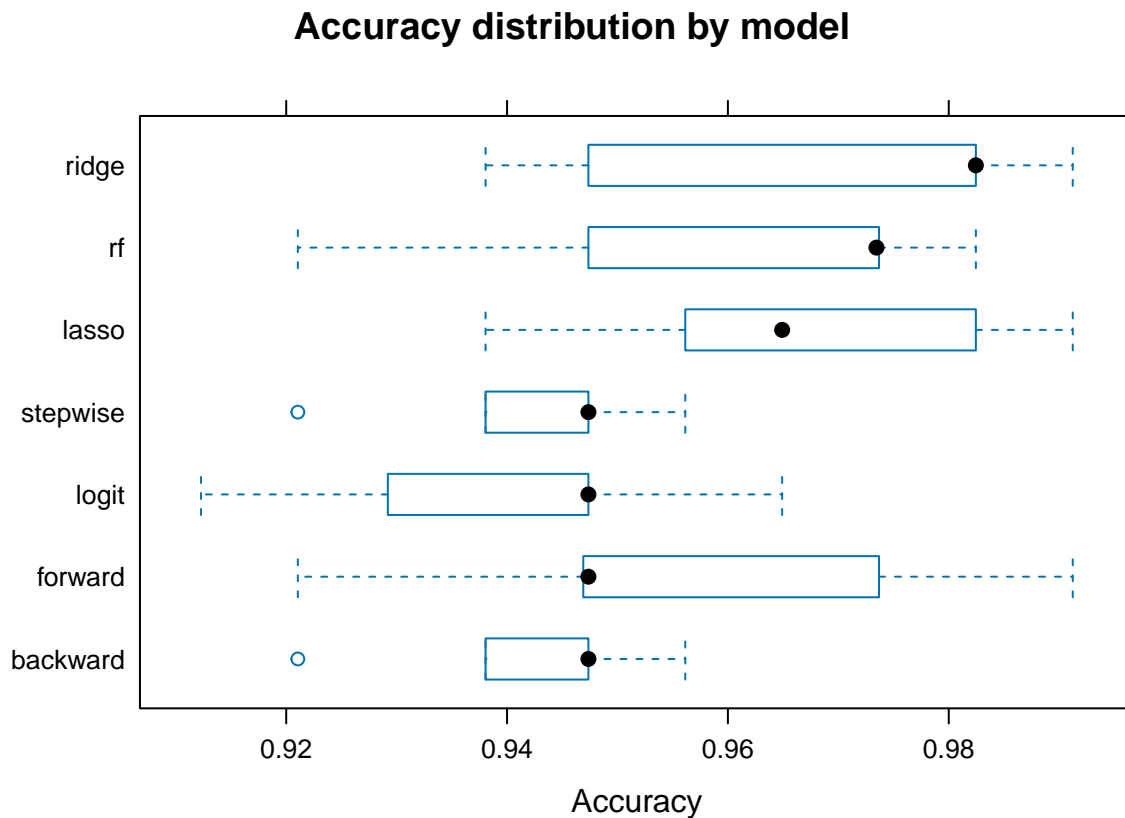
```
summary(results)$statistics$Accuracy
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
## logit	0.9122807	0.9292035	0.9473684	0.9402267	0.9473684	0.9649123	0
## stepwise	0.9210526	0.9380531	0.9473684	0.9419966	0.9473684	0.9561404	0
## backward	0.9210526	0.9380531	0.9473684	0.9419966	0.9473684	0.9561404	0
## forward	0.9210526	0.9469027	0.9473684	0.9560472	0.9736842	0.9912281	0
## lasso	0.9380531	0.9561404	0.9649123	0.9665580	0.9824561	0.9912281	0
## ridge	0.9380531	0.9473684	0.9824561	0.9683124	0.9824561	0.9912281	0
## rf	0.9210526	0.9473684	0.9734513	0.9596025	0.9736842	0.9824561	0

Accuracy is the proportion of correct predictions out of the total number of observations.

$$Accuracy = \frac{\text{Correct predictions}}{\text{Total predictions}}$$

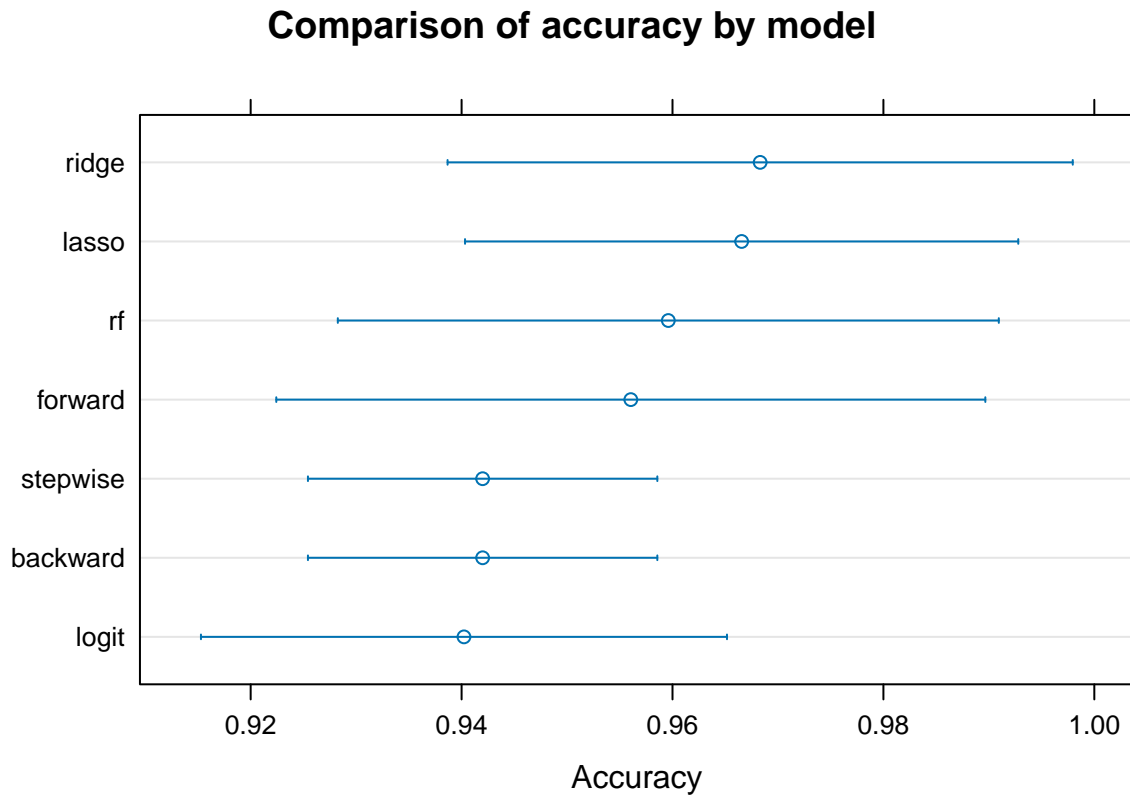
```
bwplot(results, metric = "Accuracy", main = "Accuracy distribution by model")
```



This chart shows a statistical summary of the accuracy of the 5 folds used.

The box shows the first and third quartiles; the black dot is the median. Points outside the whiskers are outliers.

```
dotplot(results, metric = "Accuracy", main = "Comparison of accuracy by model")
```



Confidence Level: 0.95

In this graph, the dot shows the average accuracy across the 5 folds for each model. The line around it shows the dispersion in accuracy. Shorter lines indicate a more consistent metric, while longer lines indicate less consistency and greater variability.

Fitting Time

```
#Combine the training times into one data frame
Times_DF <- data.frame(
  #Models column
  Model = c("Logistic", "Stepwise", "Backward", "Forward", "Lasso", "Ridge", "RandomForest"),
  #Column of Times
  Fitting_Time = c(Base_Time, Stepwise_Time, Backward_Time, Forward_Time,
                    Lasso_Time, Ridge_Time, Random_Forest_Time)
)

#Order from smallest to largest
Times_DF <- Times_DF[order(Times_DF$Fitting_Time), ]

Times_DF
```

```
##           Model Fitting_Time
```

```
## 1      Logistic      0.28
## 6      Ridge        0.37
## 5      Lasso         0.61
## 7 RandomForest      3.22
## 4      Forward      4.82
## 3      Backward     21.24
## 2      Stepwise     31.78
```

This table shows the time it takes to train each model.

Comparison

```
#We extract the average of Accuracy
Metrics<- summary(results)
metrics_summary <- Metrics$statistics

accuracy_mean <- metrics_summary$Accuracy[, 'Mean']

#We create a vector with the times
tiempos <- c(Base_Time, Stepwise_Time, Backward_Time, Forward_Time, Lasso_Time, Ridge_Time,
              Random_Forest_Time)

#Date Frame with the values to compare the models
Comparison <- data.frame(
  Model = names(accuracy_mean),
  Accuracy_Mean = unname(accuracy_mean),
  Time = tiempos
)

Comparison
```

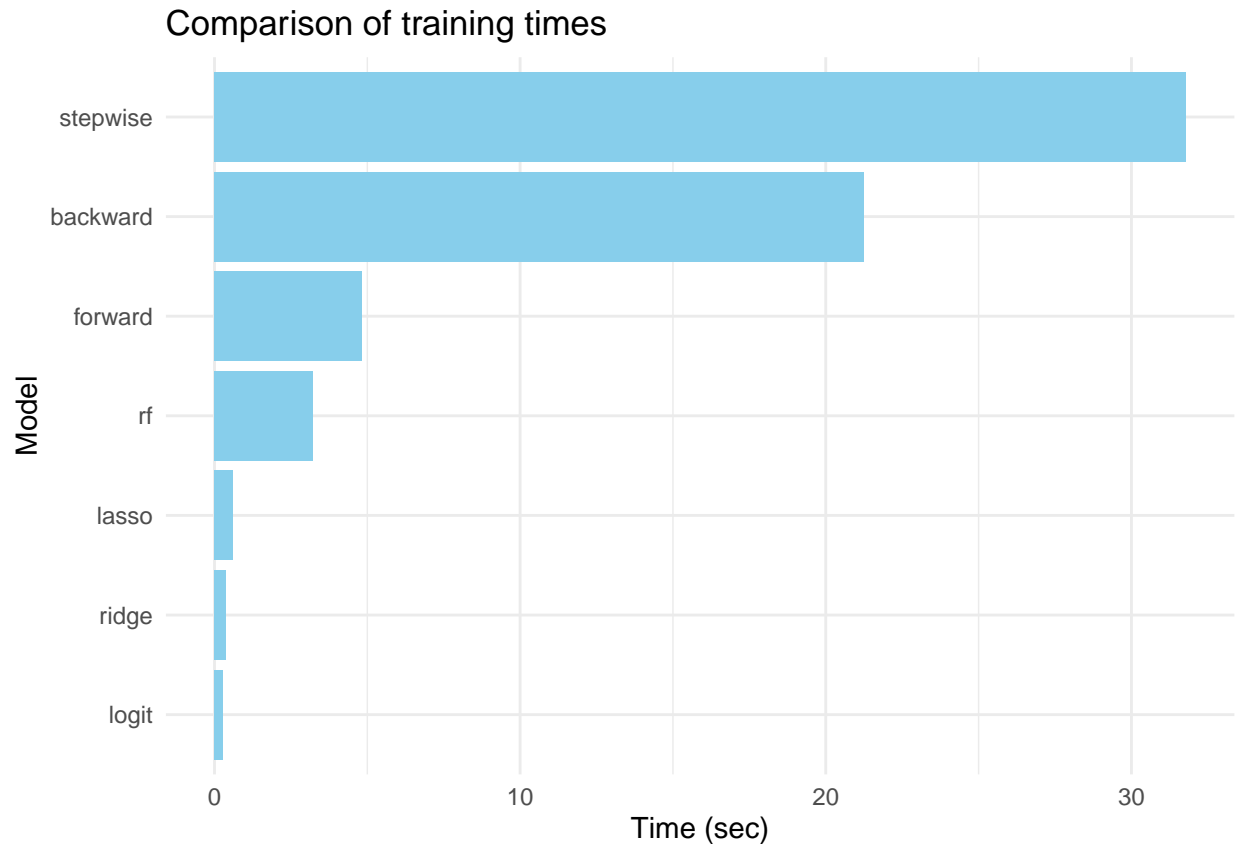
```
##      Model Accuracy_Mean  Time
## 1   logit      0.9402267  0.28
## 2 stepwise      0.9419966 31.78
## 3 backward      0.9419966 21.24
## 4 forward       0.9560472  4.82
## 5  lasso        0.9665580  0.61
## 6  ridge        0.9683124  0.37
## 7    rf         0.9596025  3.22
```

Time Graph

```
ggplot(Comparison,
       #Order the bars from shortest to longest time
       aes(x = reorder(Model, Time), y = Time)) +

  #draw bars according to the time values
  geom_bar(stat = "identity", fill = "skyblue") +
```

```
coord_flip() + #makes the bars horizontal
#Labels
labs(title = "Comparison of training times",
      x = "Model",
      y = "Time (sec)") +
theme_minimal()
```



Note that the most accurate models are Ridge, followed by Lasso, and then Random Forest. Ridge's training time is slightly faster than Lasso, but Random Forest takes almost six times longer to be adjusted. Recall that in terms of complexity, the Lasso model was the simplest, but Ridge, despite being more complex, wasn't as complex as the other models.

We can conclude that the Lasso model is the best option, as it offers simplicity, good accuracy, and a short adjustment time. If we are willing to sacrifice simplicity for a shorter training time and greater accuracy, Ridge is the best choice. Random Forest, being a machine learning model and therefore much less interpretable, has similar accuracy to the other models, which, while not bad, are inferior to Lasso and Ridge. Even so, it is faster and more accurate than the other models. The other models are not worthwhile in this situation; they have a significantly longer training time, their accuracy is worse, and they are much more complex.

Let's repeat this analysis 50 times

Due to the time it takes for them to adjust, we will not be using the Backward and Stepwise models. These models did not demonstrate outstanding performance, so I decided to forget them at this stage.

To obtain more conclusive data, we will repeat this same experiment 50 times, with different folds, to obtain a more precise sample of accuracy and not attribute it to chance.

Creation of the folds for the different iterations

```
n_repeats <- 50
k <- 5

set.seed(996)

#Vector that stores the folds for each iteration.
folds_list <- vector("list", n_repeats)

for (i in 1:n_repeats) {
  folds_list[[i]] <- createFolds(y = data$Diagnosis, k = k, returnTrain = TRUE)
}
```

Training and storage of results

```
results_list <- list()

for (i in 1:n_repeats) {
  # I establish the folds that correspond
  #to each iteration.
  Method <- trainControl(
    method = "cv",
    number = k,
    index = folds_list[[i]],
    summaryFunction = twoClassSummary,
    classProbs = TRUE,
    savePredictions = "final"
  )

  suppressWarnings({

    Base_Time <- system.time({
      Base_model <- train(
        Diagnosis ~ ., data = data,
        method = "glm",
        family = binomial,
        trControl = Method
      )
    })["elapsed"]

    Forward_Time <- system.time({
      Forward_Model <- train(
        Diagnosis ~ ., data = data,
        method = "glmStepAIC",
        trControl = Method,

```

```

        direction = "forward",
        trace = FALSE
    )
  })["elapsed"]

Lasso_Time <- system.time({
  Lasso_Model <- train(
    Diagnosis ~ ., data = data,
    method = "glmnet",
    trControl = Method,
    tuneGrid = expand.grid(alpha = 1,
                           lambda = seq(0.0001, 1, length = 50))
  )
  })["elapsed"]

Lasso_Model$pred <- Lasso_Model$pred[
  Lasso_Model$pred$lambda == Lasso_Model$bestTune$lambda, ]

Ridge_Time <- system.time({
  Ridge_Model <- train(
    Diagnosis ~ ., data = data,
    method = "glmnet",
    trControl = Method,
    tuneGrid = expand.grid(alpha = 0, lambda = seq(0.0001, 1, length = 50))
  )
  })["elapsed"]

Ridge_Model$pred <- Ridge_Model$pred[
  Ridge_Model$pred$lambda == Ridge_Model$bestTune$lambda, ]

Random_Forest_Time <- system.time({
  RF_model <- train(
    Diagnosis ~ ., data = data,
    method = "rf",
    trControl = Method
  )
  })["elapsed"]
})

results_list[[i]] <- list(Logit = Base_model, Forward=Forward_Model,
                          LASSO=Lasso_Model, Ridge=Ridge_Model, RF = RF_model)
}

```

Summary the results

```
all_results <- data.frame()
```



```

for (i in 1:n_repeats) {
  for (model_name in names(results_list[[i]])) {

    pred <- results_list[[i]][[model_name]]$pred

    # Accuracy average
    acc <- mean(pred$pred == pred$obs)

    # Brier Score
    probs <- pred$M
    obs <- ifelse(pred$obs == "M", 1, 0)
    brier <- mean((probs - obs)^2)

    all_results <- rbind(all_results,
                        data.frame(
                          iteration = i,
                          model = model_name,
                          Accuracy = acc,
                          Brier = brier
                        ))
  }
}

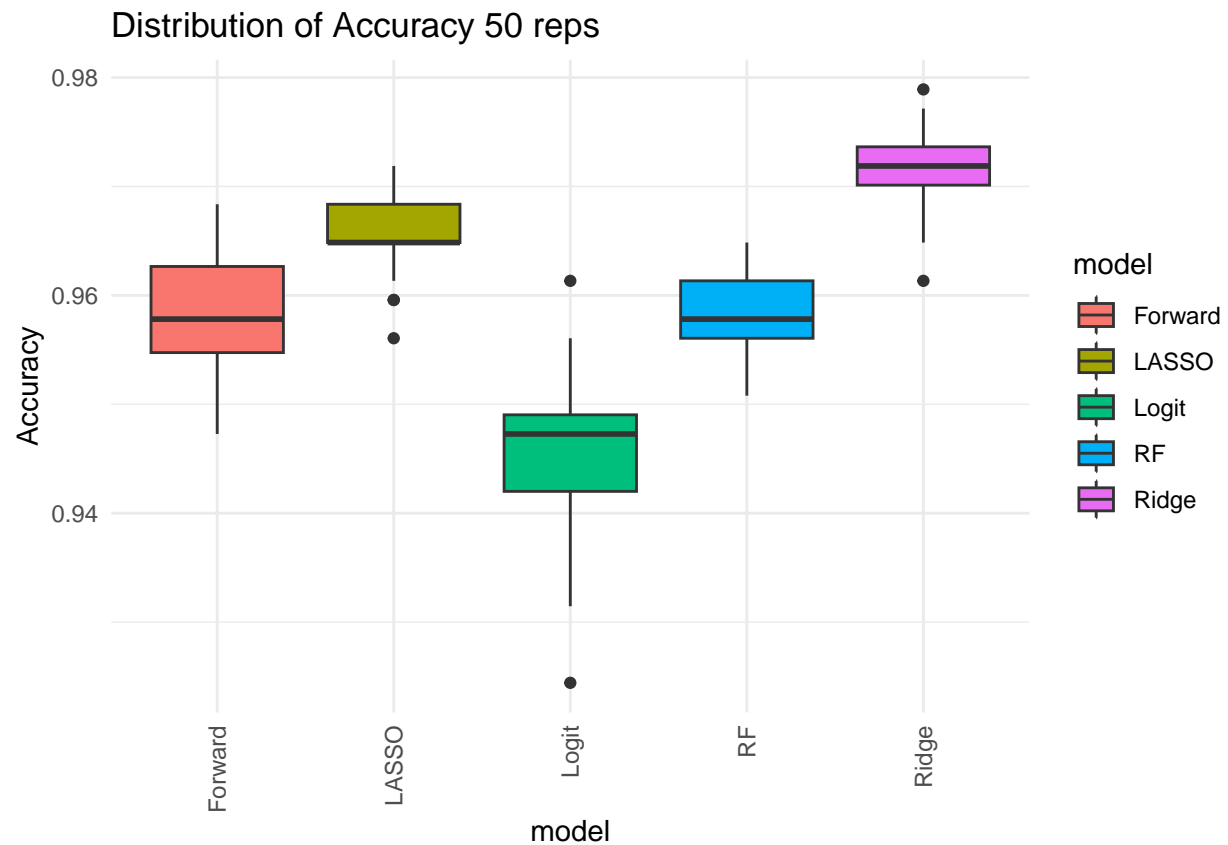
```

Graphs

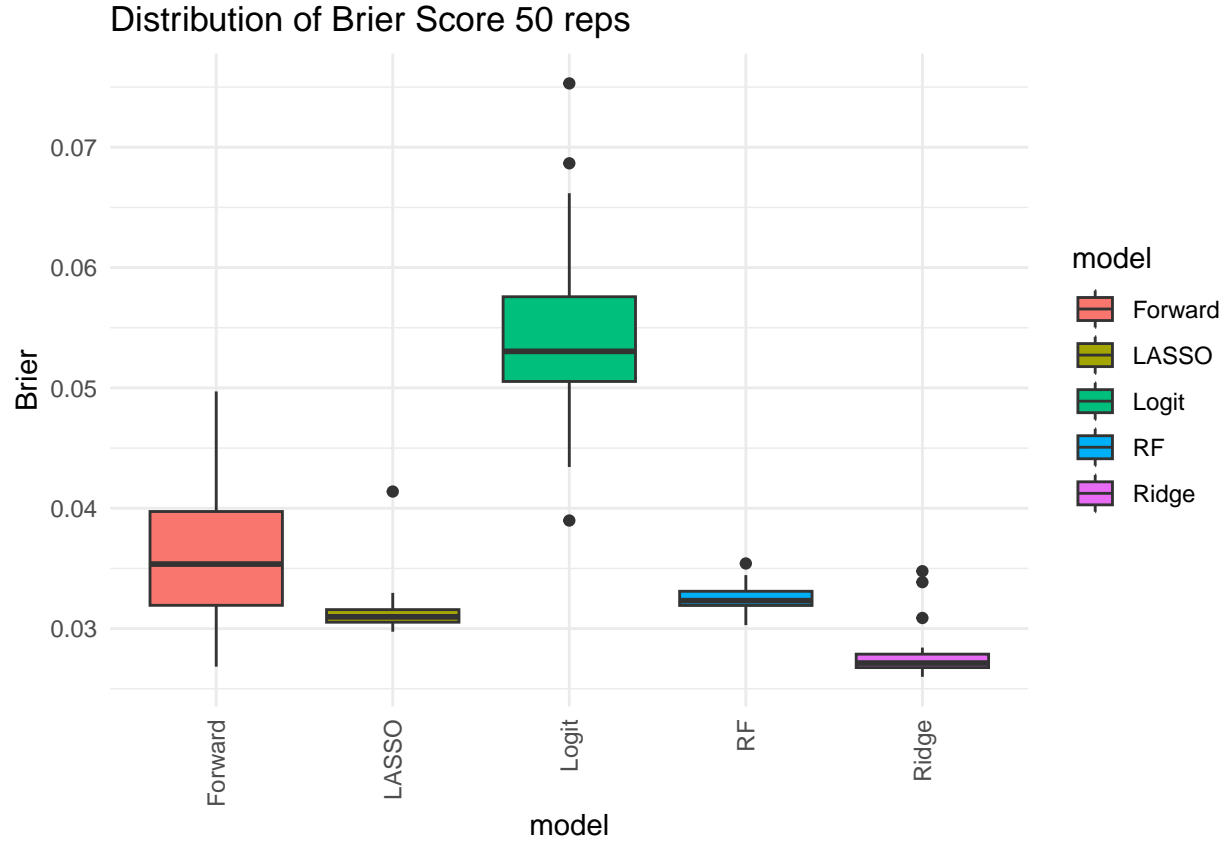
```

ggplot(all_results, aes(x = model, y = Accuracy, fill = model)) +
  geom_boxplot() +
  theme_minimal() +
  ggtitle("Distribution of Accuracy 50 reps")+
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust = 1))

```



```
# Boxplot Brier
ggplot(all_results, aes(x = model, y = Brier, fill = model)) +
  geom_boxplot() +
  theme_minimal() +
  ggtitle("Distribution of Brier Score 50 reps")+
  theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust = 1))
```



Accuracy is the proportion of correct predictions out of the total number of observations.

$$Accuracy = \frac{\text{Correct predictions}}{\text{Total predictions}}$$

Brier Score is the mean squared error of the forecast. It tells us how confidently the model predicts the response variable in each observation.

$$\text{Brier Score} = \frac{1}{N} \sum_{i=1}^N (p_i - y_i)^2$$

Where:

p_i : Predicted probability of the positive class.

y_i : Actual value encoded as 0 or 1.

Final Reflection

Looking at the box plot for the 50 repetitions, we can conclude that Ridge remains the best-performing model, followed by LASSO. Random Forest is still a good candidate but is on par with Forward, falling behind these two.

Using the Brier Score, we can conclude that the models predict with high accuracy, particularly Ridge and Lasso, in the same order of performance. Ridge and Lasso have smaller boxes, suggesting that their prediction accuracy is more consistent than the others.

Our final conclusion is the same as with the single-repetition model. The most accurate and reliable predictor is Ridge, but if we don't mind sacrificing a little precision for a much simpler model, LASSO is undoubtedly the clear winner. Although RF is also a good candidate, it still lags behind the two mentioned above and completely lacks interpretability.

In my opinion, the best model is LASSO; it is simple, accurate, reliable, and interpretable, not to mention that its penalty is much easier to understand.

References

- [1] Wolberg W, Mangasarian OL, Street N, Street W. Breast Cancer Wisconsin (Diagnostic) [Database]. UCI Machine Learning Repository; 1993. DOI: 10.24432/C5DW2B. Available from: <https://www.kaggle.com/datasets/uciml/breast-cancer-wisconsin-data>
- [2] An illustration of Random Forest [Internet]. ResearchGate; [cited 2025 Dec 10]. Available from: <https://www.researchgate.net/publication/372809468/figure/fig3/>
- [3] GeeksforGeeks. K-Fold Cross-Validation in Machine Learning [Internet]. GeeksforGeeks; [cited 2025 Dec 10]. Available from: <https://www.geeksforgeeks.org/machine-learning/k-fold-cross-validation-in-machine-learning/>