

Predicting Breast Cancer Prognosis

Roberto Cárdenas^a

^a*Colegio de Ciencias e Ingenierías “El Politécnico”,
Universidad San Francisco de Quito USFQ, Quito 170157, Ecuador.*

Abstract

This work explores the use of different regression methods for cancer prognosis. In particular, the studied data-set aims for prediction of survival and recurrence years of breast cancer. Out of eight measured regression models: Polynomial, Gaussian, Decision Tree, Random Forest, K Nearest Neighbours, Support Vector Machines, Neural Network, and LASSO, the LASSO regression method showed best results in comparison with lower measures of Mean Squared Error, Root Mean Squared Error and Absolute Squared Error across all the other implemented models.

Keywords: Forecasting, Regression.

1. Introduction

Nowadays, Regression methods are used to analyze the relationship between a dependent variable and one or more independent variables. The purpose of regression analysis is to determine if there is a relationship between the dependent variable and one or more independent variables. (Beers, 2022) It is also useful for modeling the future relationship between them.

Some of the problems that regression models have solved are: forecasting sales, stock predictions, predicting consumer behavior, analyzing survey data for customer satisfaction, product preferences, etc.

This work aims to explore eight regression methods to maximize the prediction of time survival and cancer recurrence on a breast cancer data set.

12 2. Materials and Methods

13 2.1. MammaPrint data set

14 We use the MammaPrint breast cancer (BC) dataset, which contains 70 fea-
15 tures and 2 target columns which are Time Survival Years and Time Recurrence
16 Years, as their names suggest, these columns represent the time in years of a
17 particular patient of either recurrence of cancer cells or survival with the disease.
18 ([Aguilera-Mendoza et al., 2015](#))

19 In particular, the features of the dataset consist of a 70-gene expression pro-
20 file that was initially created to identify individuals with early-stage BC who are
21 unlikely to experience metastases and are eligible for remission following adju-
22 vant chemotherapy. Its usage as a prognostic biomarker has received extensive
23 retrospective and prospective validation. ([Brandão-M et al., 2019](#))

24 2.2. Regression methods

25 Regression models are widely used nowadays as a powerful tool for uncov-
26 ering the associations between variables observed in data, but cannot easily
27 indicate causation. Regression in synthesis is a statistical method used in all
28 areas of finance, investing, and other disciplines where it's objective consists in
29 finding the strength and character of the relationship between one dependent
30 variable commonly denoted by Y and a series of other variables also known as
31 independent variables. ([Beers, 2022](#))

32 2.2.1. Polynomial

33 Polynomial regression is a type of linear regression that describes the rela-
34 tionship between dependant variables and independent variable modelled as an
35 nth degree polynomial in x. It is also called an special case of Multiple Lin-
36 ear Regression in Machine Learning because some polynomial terms are added
37 to the Multiple Linear regression equation to convert it into Polynomial Re-
38 gression. ([Raghav, 2021](#)) For this regression, it is important to preprocess the
39 input variables into polynomial terms using a selected degree. The equation of
40 a polynomial expression goes by the form of:

$$y = a_0 + a_1x_1 + a_2x_1^2 + \dots + a_nx_1^n$$

41 2.2.2. *Gaussian*

42 The Gaussian processes model is a probabilistic supervised machine learning
 43 framework that has been widely used for regression and classification tasks. A
 44 Gaussian processes regression (GPR) model can make predictions incorporating
 45 prior knowledge (kernels) and provide uncertainty measures over predictions
 46 ([Sit, 2019](#))

47 In GPR, it is first assumed a Gaussian process prior, which can be specified
 48 using a mean function, $m(x)$, and covariance function, $k(x, x')$:

$$f(x) \sim GP(m(x), k(x, x'))$$

49 A Gaussian process is hence comparable to an infinite-dimensional multi-
 50 variate Gaussian distribution in which any grouping of the dataset's labels is
 51 jointly distributed.

52 2.2.3. *Decision Tree*

53 The Decision Tree algorithm, as its name suggests, uses a tree-like model of
 54 decisions as a predictive model to draw conclusions about a set of observations.
 55 Decision trees where the target variable can take continuous values (typically
 56 real numbers) are called regression trees. More generally, the concept of re-
 57 gression tree can be extended to any kind of object equipped with pairwise
 58 dissimilarities such as categorical sequences. ([Prasad, 2021](#))

59 It is important to keep in mind that the algorithm is susceptible to over-
 60 fitting. Therefore, it is preferable to cross-validate and to always define the
 61 minimum number of children per leaf node in advance.

62 2.2.4. *Random Forest*

63 Random Forest can be used as a regression method that operates by con-
 64 structing a group of decision trees at training. The mean prediction of the

individual trees is returned. This model is used to solve the overfitting problem with decision trees.

The low correlation between models is the key. Uncorrelated models have the ability to generate ensemble forecasts that are more precise than any single prediction. As long as they don't consistently all err in the same direction, the trees shield each other from their individual errors, which accounts for this lovely result. Many trees will be right while some may be wrong, allowing the group of trees to move in the proper direction. (Yiu, 2019)

2.2.5. *k Nearest Neighbours*

The k-nearest neighbors algorithm, sometimes referred to as KNN or k-NN, is a supervised learning classifier that employs proximity to produce classifications or predictions about the grouping of a single data point. Although it can be applied to classification or regression issues, it is commonly employed as a classification algorithm because it relies on the idea that comparable points can be discovered close to one another. (Harrison, 2018)

For regression, the average of the k nearest neighbors is used to forecast a classification. Here, the primary difference is that classification is used for discrete data whereas regression is utilized for continuous values. However, defining the distance is necessary before a categorization can be determined. The most typical measurement is euclidean distance described as:

$$d(x, y) = \sqrt{\sum_{i=1}^n (y_i - x_i)^2}$$

2.2.6. *Support Vector Machines*

The machine learning algorithm known as a support vector machine (SVM) examines data for regression and classification purposes. This algorithm sorts data into one of two groups after looking at it. The sorted data are output as a map by an SVM, with the margins between the two being as far away as possible. SVMs are employed in the sciences, picture classification, handwriting

91 recognition, and text categorization. SVM regression is considered a non para-
92 metric technique because it relies on kernel functions. ([Ghandi, 2018](#)) The goal
93 is to find a function $f(x)$ that deviates from y_n by a value no greater than ϵ for
94 each training point x , and at the same time is as flat as possible.

95 Linear SVM Regression Primal Formula:

$$f(x) = x'\beta + b$$

96 2.2.7. Neural Network

97 An Artificial Neural Networks is a model that simulate the human brain
98 by implementing neurons. It consists of Input layer, Hidden layers, Output
99 layer. The hidden layer can be more than one in number. Each layer consists
100 of n number of neurons. Each layer will be having an Activation Function
101 associated with each of the neurons. The activation function is the function
102 that is responsible for introducing non-linearity in the relationship ([Srivignesh,
103 2021](#)). Each layer can also have regularizers associated with it. Regularizers are
104 responsible for preventing overfitting.

105 Training data is essential for neural networks to develop and enhance their
106 accuracy over time. However, these learning algorithms become effective tools in
107 computer science and artificial intelligence once they are adjusted for accuracy,
108 enabling us to quickly and effectively analyze, classify, or create data.

109 2.2.8. LASSO

110 LASSO regression is another type of linear regression. LASSO stands for
111 Least Absolute Shrinkage and Selection Operator. It consists in the use of
112 shrinkage that is defined as where data values are shrunk towards a central
113 point, similar to the mean. The lasso procedure emphasizes simple and sparse
114 models with fewer parameters, in order to show at high fidelity models with high
115 levels of multicollinearity. It is also particularly useful for automation of certain
116 parts of model selection or elimination like parameters or variables. ([Glen,
117 2022a](#))

118

119

120 Regression with L1 Regularization: If a regression model uses the L1 Reg-
 121 ularization technique, then it is called Lasso Regression. If it used the L2
 122 regularization technique, it's called Ridge Regression.

123 The L1 regularization process used by Lasso regression results in a penalty
 124 proportional to the absolute magnitude of the coefficients. A sparse model
 125 with few coefficients may be produced by this kind of regularization; certain
 126 coefficients may go to zero and be removed from the model. Greater penalties
 127 provide coefficient values that are closer to zero, which is great for creating more
 128 straightforward models. However, L2 regularization (such as Ridge regression)
 129 does not eliminate coefficients or sparse models. Because of this, Lasso is much
 130 simpler to understand than the Ridge. ([Glen, 2022a](#))

131 Following is the formula of how the regression is performed:

$$\sum_{i=1}^n (y_i - \sum_j x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

132 The lambda tuning parameter controls the strength of L1 penalty. Lambda
 133 is in synthesis the amount of shrinkage:

134 If $\lambda = 0$, no parameters are eliminated. The estimate is equal to the one
 135 found with linear regression. As λ increases, coefficients are eliminated.
 136 As λ increases, bias increases and finally, as λ decreases there is an
 137 increase in variance.

138 2.3. Experimental methodology

139 This section outlines the experimental evaluation of the model used for this
 140 project. The 8 models presented above were implemented using python pro-
 141 gramming techniques and the use of specific scientific libraries adjusted to every
 142 model independently from the others. The following preparations and configu-
 143 rations were made:

144 2.3.1. Data set normalization

145 The min-max method was used to normalize the data. In this way, incon-
146 sistencies with the model are avoided and a better training can be done.

147 The initial data is transformed linearly using min-max normalization, also
148 known as feature scaling. Using this method, all data is scaled within the range
149 from (0, 1). The following is the formula to accomplish this:

$$x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

150 ([O'Reilly, 2023](#))

151 2.3.2. Training, validation, and test sets

152 The dataset was separated by 10 percent for the test set (which does not
153 participate in the training of the models). This is used to feed the best model
154 selected by each method. A repeated 10-fold cross validation was applied to the
155 remaining 90% of the data that make up the training set, and as a result of this
156 process, the train or validation folds were formed.

157 2.3.3. Evaluation metrics

158 The evaluation metrics used for this projects were: Mean Squared Error,
159 Mean Absolute Error, and Root Mean Squared Error. This metrics were mea-
160 sured across 10-fold cross validation scheme, so the mean of all 10 folds was
161 represented in the results. This process gives much more information about the
162 algorithm performance, and draw better informed conclusions. ([Shulga, 2018](#))

163 Mean Squared Error: It is possible to determine how closely a regression
164 line resembles a set of points using the mean squared error (MSE). This is
165 accomplished by squaring the distances between the points and the regression
166 line (also known as the "errors"). The squaring is required to eliminate any
167 unfavorable indications. Additionally, it emphasizes bigger discrepancies. Since
168 it is averaging a collection of errors, this error type is known as the mean squared
169 error. The forecast is more accurate the lower the MSE. ([Glen, 2022b](#))

170 Root Mean Squared Error: Root Mean Square Error (RMSE) is the standard
171 deviation of the residuals (prediction errors). The distance between the data
172 points and the regression line is measured by residuals, and the spread of these
173 residuals is measured by RMSE. In other words, it provides information on
174 how tightly the data is clustered around the line of best fit. In climatology,
175 forecasting, and regression analysis, root mean square error is frequently used
176 to validate experimental results. ([Glen, 2022c](#))

177 Mean Absolute Error: Mean Absolute Error is a model evaluation metric
178 used with regression models. The mean absolute error of a model with respect to
179 a test set is the mean of the absolute values of the individual prediction errors on
180 over all instances in the test set. Each prediction error is the difference between
181 the true value and the predicted value for the instance. ([Sammut, 2011](#))

182 **3. Results and Discussion**

183 *3.0.1. Performance in the training*

184 Below are the tables of results obtained by each of the regression meth-
185 ods used in this project. The average MSE, MAE, and RMSE measurements
186 through cross validation and the standard deviation of each result are repre-
187 sented here.

188

Time Survival Years						
Model	MSE	MSE std	RMSE	RMSE std	MAE	MAE std
Polynomial	0.09	0.03	0.30	0.05	0.24	0.04
Gaussian	0.05	0.02	0.23	0.04	0.19	0.03
Decision Tree	0.08	0.02	0.28	0.04	0.23	0.03
Random Forest	0.05	0.01	0.22	0.03	0.18	0.03
KNN	0.05	0.01	0.22	0.03	0.18	0.03
SVM	0.06	0.02	0.24	0.03	0.19	0.03
Neural Network	0.06	0.02	0.23	0.04	0.19	0.03
LASSO*	0.05	0.01	0.21	0.03	0.17	0.03

189

190

Time Recurrence Years						
Model	MSE	MSE std	RMSE	RMSE std	MAE	MAE std
Polynomial	0.10	0.03	0.32	0.04	0.25	0.03
Gaussian	0.06	0.02	0.24	0.04	0.19	0.03
Decision Tree	0.10	0.02	0.31	0.03	0.24	0.03
Random Forest	0.06	0.01	0.23	0.03	0.19	0.02
KNN	0.05	0.01	0.23	0.03	0.19	0.03
SVM	0.07	0.02	0.25	0.03	0.20	0.02
Neural Network	0.07	0.02	0.25	0.04	0.20	0.03
LASSO*	0.05	0.02	0.23	0.03	0.19	0.03

191

192

193

194

195

196

197

198

Lasso model, as shown in the table above, has the best performance on the studied dataset. It is possible that the model is showing high levels of multicollinearity as in when there are high correlations between two or more predictor variables. In other words, one predictor variable can be used to predict the other. This creates redundant information, skewing the results in a regression model. For this reason, the shrinkage property of the LASSO model works well minimizing the error metrics given the 70 features of the dataset.

199 *3.0.2. Performance in the test set*

200 Here we present the results obtained by the best selected model, the one that
 201 minimizes the desired metric. In our case it has been the LASSO model. The
 202 table below shows the performance of the model on the 10% of the data excluded
 203 from the training dataset with their respective MSE and MAE evaluations.

LASSO testing performance		
Target	Mean Squared Error	Mean Absolute Error
Time Survival Years	0.0467	0.1776
Time Recurrence Years	0.0562	0.1945

205 It is clear that the model performance metrics in the training phase with
 206 90% of the dataset is consistent with the results obtained in the testing phase
 207 with 10% of the dataset. The result for Mean Absolute Error is still not ideal
 208 and other model configurations should be taken into account in order to lower
 209 even further these results considering the sensibility of these type of predictions.

210 **4. Conclusions**

211 In this project, 8 different regression algorithms were implemented and tested
 212 for breast cancer survival and recurrence years prognosis. Different techniques
 213 were used in each algorithm being measured by the same metrics for comparison.
 214 Then, conclusions were drawn by which algorithm has been observed to yield
 215 the best results across both metrics and target values.

216 Accordingly, it was seen in our results that prognosis of breast cancer with
 217 MammaPrint technology seems to be better to work with LASSO regression due
 218 to its lower mean squared error, root mean squared error and mean absolute
 219 error across all other regression models. It is possible that the studied model
 220 was showing high levels of multicollinearity and for this reason the shrinkage
 221 property of LASSO worked well in minimizing the error metrics given the 70
 222 features of the dataset.

223 References

- 224 Aguilera-Mendoza, L., Marrero-Ponce, Y., Tellez-Ibarra, R., Llorente-Quesada,
225 M. T., Salgado, J., Barigye, S. J., & Liu, J. (2015). Overlap and diversity in
226 antimicrobial peptide databases: compiling a non-redundant set of sequences.
227 *Bioinformatics*, 31, 2553–2559. doi:[10.1093/bioinformatics/btv180](https://doi.org/10.1093/bioinformatics/btv180).
- 228 Beers, B. (2022). What is regression? definition, calculation, and example.
229 URL: <https://www.investopedia.com/terms/r/regression.asp#>.
- 230 Brandão-M, Pondé-N, & Piccart-Gehbart, M. (2019). Mammaprint™: a com-
231 prehensive review. *Future Oncology*, 15, 207–224. URL: [https://doi.org/](https://doi.org/10.2217/fon-2018-0221)
232 [10.2217/fon-2018-0221](https://doi.org/10.2217/fon-2018-0221). doi:[10.2217/fon-2018-0221](https://doi.org/10.2217/fon-2018-0221). 336arXiv:[https:](https://doi.org/10.2217/fon-2018-0221)
233 [//doi.org/10.2217/fon-2018-0221](https://doi.org/10.2217/fon-2018-0221).
- 234 Ghandi, R. (2018). Support vector machine — introduction to ma-
235 chine learning algorithms. URL: [https://towardsdatascience.com/](https://towardsdatascience.com/support-vector-machine-introduction-to-machine-learning-algorithms-934a444fca47)
236 [support-vector-machine-introduction-to-machine-learning-algorithms-934a444fca47](https://towardsdatascience.com/support-vector-machine-introduction-to-machine-learning-algorithms-934a444fca47).
- 237 Glen, S. (2022a). Mse: Mean squared error. URL: [https://www.](https://www.statisticshowto.com/lasso-regression/)
238 [statisticshowto.com/lasso-regression/](https://www.statisticshowto.com/lasso-regression/).
- 239 Glen, S. (2022b). Mse: Mean squared error. URL: [https:](https://www.statisticshowto.com/probability-and-statistics/statistics-definitions/mean-squared-error/)
240 [//www.statisticshowto.com/probability-and-statistics/](https://www.statisticshowto.com/probability-and-statistics/statistics-definitions/mean-squared-error/)
241 [statistics-definitions/mean-squared-error/](https://www.statisticshowto.com/probability-and-statistics/statistics-definitions/mean-squared-error/).
- 242 Glen, S. (2022c). Rmse: Root mean square error. URL:
243 [https://www.statisticshowto.com/probability-and-statistics/](https://www.statisticshowto.com/probability-and-statistics/regression-analysis/rmse-root-mean-square-error/)
244 [regression-analysis/rmse-root-mean-square-error/](https://www.statisticshowto.com/probability-and-statistics/regression-analysis/rmse-root-mean-square-error/).
- 245 Harrison, O. (2018). Machine learning basics with the k-nearest
246 neighbors algorithm. URL: [https://towardsdatascience.com/](https://towardsdatascience.com/machine-learning-basics-with-the-k-nearest-neighbors-algorithm-6a6e71d01761)
247 [machine-learning-basics-with-the-k-nearest-neighbors-algorithm-6a6e71d01761](https://towardsdatascience.com/machine-learning-basics-with-the-k-nearest-neighbors-algorithm-6a6e71d01761).
- 248 O'Reilly (2023). Min-max normalization. URL: [https://www.oreilly.](https://www.oreilly.com/library/view/hands-on-machine-learning/9781788393485/fd5b8a44-e9d3-4c19-bebb-c2fa5a5ebfee.xhtml)
249 [com/library/view/hands-on-machine-learning/9781788393485/](https://www.oreilly.com/library/view/hands-on-machine-learning/9781788393485/fd5b8a44-e9d3-4c19-bebb-c2fa5a5ebfee.xhtml)
250 [fd5b8a44-e9d3-4c19-bebb-c2fa5a5ebfee.xhtml](https://www.oreilly.com/library/view/hands-on-machine-learning/9781788393485/fd5b8a44-e9d3-4c19-bebb-c2fa5a5ebfee.xhtml).

251 Prasad, A. (2021). Decision trees fro regres-
 252 sion. URL: [https://medium.com/analytics-vidhya/
 253 regression-trees-decision-tree-for-regression-machine-learning-e4d7525d8047](https://medium.com/analytics-vidhya/regression-trees-decision-tree-for-regression-machine-learning-e4d7525d8047).

254 Raghav, A. (2021). All you need to know about polynomial re-
 255 gression. URL: [https://www.analyticsvidhya.com/blog/2021/07/
 256 all-you-need-to-know-about-polynomial-regression/](https://www.analyticsvidhya.com/blog/2021/07/all-you-need-to-know-about-polynomial-regression/).

257 Sammut, C. (2011). Mean absolute error. URL: [https://link.springer.com/
 258 referenceworkentry/10.1007/978-0-387-30164-8_525](https://link.springer.com/referenceworkentry/10.1007/978-0-387-30164-8_525).

259 Shulga, D. (2018). 5 reasons why you should use cross-validation in
 260 your data science projects. URL: [https://towardsdatascience.com/
 261 5-reasons-why-you-should-use-cross-validation-in-your-data-science-project-8163311a1e79#](https://towardsdatascience.com/5-reasons-why-you-should-use-cross-validation-in-your-data-science-project-8163311a1e79#).

262 Sit, H. (2019). Quick start to gaussian process re-
 263 gression. URL: [https://towardsdatascience.com/
 264 quick-start-to-gaussian-process-regression-36d838810319](https://towardsdatascience.com/quick-start-to-gaussian-process-regression-36d838810319).

265 Srivignesh, R. (2021). A walk-through of regression
 266 analysis using artificial neural networks in tensorflow.
 267 URL: [https://www.analyticsvidhya.com/blog/2021/08/
 268 a-walk-through-of-regression-analysis-using-artificial-neural-networks-in-tensorflow/](https://www.analyticsvidhya.com/blog/2021/08/a-walk-through-of-regression-analysis-using-artificial-neural-networks-in-tensorflow/).

269 Yiu, T. (2019). Understanding random for-
 270 est. URL: [https://towardsdatascience.com/
 271 understanding-random-forest-58381e0602d2](https://towardsdatascience.com/understanding-random-forest-58381e0602d2).