

2023SP-COMP_SCI-5530-0002

Project Report: Breast Cancer Prediction

Using Machine Learning

Maruthi Phanindra Ayyagari 16323307

Mufrad Islam 16334246

Vamsi Kiran Murukutla 16326572

Ashok Karanam 16324075

Department of Computer Science

Submitted to: Yu Luo

Department of Computer Science

The University of Missouri, Kansas City, MO,64110.

Abstract

Breast cancer (BC) is one of the most common cancers among women worldwide, representing most new cancer cases and cancer-related deaths according to global statistics, making it a significant public health problem today.

The early diagnosis of BC can improve the prognosis and chance of survival significantly, as it can promote timely clinical treatment to patients. Further accurate classification of benign tumors can prevent patients from undergoing unnecessary treatments. Thus, the correct diagnosis of BC and classification of patients in malignant or benign groups are the subject of much research. The goal is to classify whether the breast cancer is benign or malignant.

We will use the UCI Machine Learning Repository for the breast cancer dataset

(<http://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+%28diagnostic%29>). The dataset we are going to use in this project is publicly available and was created by Dr. William H. Wolberg, a physician at the University Of Wisconsin Hospital at Madison, Wisconsin, USA

Keywords: Machine Learning, Pattern Recognition, Classification, Supervised Learning, Artificial Intelligence.

Introduction:

One of the most prevalent diseases in today's world is cancer among which breast cancer is caused by a variety of clinical, lifestyle, social, and economic variables. Breast cancer has become the most recurrent type of health issue among women especially for women in middle age. Early detection of breast cancer can help women cure this disease and the death rate can be reduced. The goal of this project is to predict breast cancer using machine learning algorithms.

Breast Cancer Classification:

Breast cancer classification is a supervised learning problem where the goal is to predict the diagnosis of breast cancer as either malignant or benign based on input features such as cell shape, size, and texture. It is a classification problem because the output variable is categorical, either malignant or benign. The dataset used in the Python program described earlier consists of several features related to breast cell characteristics, as well as the diagnosis of the cells as either malignant or benign.

Supervised learning is used because the dataset contains labeled data, which means that the output variable is known for each example in the dataset. In this case, the output variable is the diagnosis of breast cancer, and the input variables are the cell features.

The goal of the program is to train a machine learning model using the labeled dataset to predict the diagnosis of breast cancer for new, unlabeled examples. The program uses several machine learning algorithms to train and test the model and evaluate its performance using various metrics such as accuracy, precision, and recall.

Dataset: For this project, we have chosen the Breast Cancer Wisconsin Dataset of the following attributes and features.

The attributes are:

ID number

Diagnosis (M = Malignant, B = Benign)

Ten real-valued features are computed for each cell. They are:

Radius

Texture

Perimeter

Area

Smoothness

Compactness

Concavity

Concave points

Symmetry

Fractal Dimension

Dataset: "data.csv" is the dataset used in this program. It contains information about breast cancer diagnosis, including the characteristics of the cell nuclei present in the digitized image of a fine needle aspirate (FNA) of a breast mass. The dataset has 569 samples and 33 features. The targeted variables are the diagnosis column, which is either M (Malignant) or B (Benign).

Counting the number of Rows and Columns: Each row represents an individual patient, and the columns represent the features of each patient.

id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean
842302	M	17.99	10.38	122.88	1001	0.1184	0.2776	0.3001	0.1471

84 25 17	M	20. 57	17.7 7	132. 9	13 26	0.084 74	0.078 64	0.08 69	0.070 17
84 30 09 03	M	19. 69	21.2 5	130	12 03	0.109 6	0.159 9	0.19 74	0.127 9
84 34 83 01	M	11. 42	20.3 8	77.5 8	38 6.1	0.142 5	0.283 9	0.24 14	0.105 2
84 35 84 02	M	20. 29	14.3 4	135. 1	12 97	0.100 3	0.132 8	0.19 8	0.104 3
84 37 86	M	12. 45	15.7	82.5 7	47 7.1	0.127 8	0.17	0.15 78	0.080 89
84 43 59	M	18. 25	19.9 8	119. 6	10 40	0.094 63	0.109	0.11 27	0.074
84 45 82 02	M	13. 71	20.8 3	90.2	57 7.9	0.118 9	0.164 5	0.09 366	0.059 85
84 49 81	M	13	21.8 2	87.5	51 9.8	0.127 3	0.193 2	0.18 59	0.093 53
84 50 10 01	M	12. 46	24.0 4	83.9 7	47 5.9	0.118 6	0.239 6	0.22 73	0.085 43
84 56 36	M	16. 02	23.2 4	102. 7	79 7.8	0.082 06	0.066 69	0.03 299	0.033 23
84 61 00 02	M	15. 78	17.8 9	103. 6	78 1	0.097 1	0.129 2	0.09 954	0.066 06
84 62 26	M	19. 17	24.8	132. 4	11 23	0.097 4	0.245 8	0.20 65	0.111 8
84 63 81	M	15. 85	23.9 5	103. 7	78 2.7	0.084 01	0.100 2	0.09 938	0.053 64
84 66	M	13. 73	22.6 1	93.6	57 8.3	0.113 1	0.229 3	0.21 28	0.080 25

74 01									
84 79 90 02	M	14. 54	27.5 4	96.7 3	65 8.8	0.113 9	0.159 5	0.16 39	0.073 64
84 84 06	M	14. 68	20.1 3	94.7 4	68 4.5	0.098 67	0.072	0.07 395	0.052 59
84 86 20 01	M	16. 13	20.6 8	108. 1	79 8.8	0.117	0.202 2	0.17 22	0.102 8
84 90 14	M	19. 81	22.1 5	130	12 60	0.098 31	0.102 7	0.14 79	0.094 98
85 10 42 6	B	13. 54	14.3 6	87.4 6	56 6.3	0.097 79	0.081 29	0.06 664	0.047 81
85 10 65 3	B	13. 08	15.7 1	85.6 3	52 0	0.107 5	0.127	0.04 568	0.031 1
85 10 82 4	B	9.5 04	12.4 4	60.3 4	27 3.9	0.102 4	0.064 92	0.02 956	0.020 76
85 11 13 3	M	15. 34	14.2 6	102. 5	70 4.4	0.107 3	0.213 5	0.20 77	0.097 56
85 15 09	M	21. 16	23.0 4	137. 2	14 04	0.094 28	0.102 2	0.10 97	0.086 32
85 25 52	M	16. 65	21.3 8	110	90 4.6	0.112 1	0.145 7	0.15 25	0.091 7
85 26 31	M	17. 14	16.4	116	91 2.7	0.118 6	0.227 6	0.22 29	0.140 1
85 27 63	M	14. 58	21.5 3	97.4 1	64 4.8	0.105 4	0.186 8	0.14 25	0.087 83
85 27 81	M	18. 61	20.2 5	122. 1	10 94	0.094 4	0.106 6	0.14 9	0.077 31

85									
29		15.	25.2	102.	73	0.108	0.169	0.16	0.087
73	M	3	7	4	2.4	2	7	83	51

Data Cleaning: Before properly using the data, we used the function called `isna()` to count the number of missing values in each column. We found a few columns with missing values and removed them using the `dropna()` function. We again checked the number of missing values using `isna()`, after removing the column with missing values. To make sure that no missing values remained.

Data Processing: To encode the categorical data in the diagnosis column, We then used the `LabelEncoder()` function of Scikit-Learn. To visualize the distribution of features and their correlation with each other, we also created a pair plot of the dataset using the `pairplot()` function of Seaborn. By using the `corr()` function of Pandas DataFrame the correlation of the columns in the dataset was obtained, and the correlation matrix was visualized using the `heatmap()` function of Seaborn [5].

Encoding the categorical data: Here, the values consisting in the “diagnosis” column are the categorical data. The values are Malignant (M) and Benign (B). We have encoded this value to the integer-type data

Creating a pair plot: We have created a pair plot for each row based on the data points of the “diagnosis” column. The following screenshot is provided for the first 5 indexes

Determining the correlation between the features of the dataset

The statistical relationship between two variables is referred to as their correlation. A correlation could be positive, meaning both variables move in the same direction. Negative, meaning that when one variable’s value increases, the other variable’s value decrease.

Visualizing the Correlation

A correlation heatmap is a graphical representation of a correlation matrix representing the correlation between different variables. Method **corr()** is invoked on the **Pandas DataFrame** to determine the correlation between different variables including predictor and response variables

Data normalization Techniques:

For the data normalization technique we used the **StandardScaler()** function of Scikit-Learn to scale the data. Scaling data ensures that each feature contributes equally to the analysis. Feature engineering or dimensionality reduction in this program is not performed.

train_test_split() function: The **train_test_split()** method is used to split our data into train and test sets. First, we need to divide our data into features (X) and labels (y). The dataframe gets divided into **X_train**, **X_test**, **y_train**, and **y_test**. **X_train** and **y_train** sets are used for training and fitting the model

Models Analysis: We have tested and trained our datasets using 3 different models.

Logistic Regression model: Logistic regression aims to solve classification problems. It does this by predicting categorical outcomes, unlike linear regression which predicts a continuous outcome.

Decision Tree Classifier: The decision tree classifier creates the classification model by building a decision tree. Each node in the tree specifies a test on an attribute, each branch descending from that node corresponds to one of the possible values for that attribute.

Random Forest Classifier: One of the most important features of the Random Forest Algorithm is that it can handle the data set containing continuous variables as in the case of regression and categorical variables as in the case of classification. It performs better results for classification problems.

We have tested the accuracy of each model. The results are shown in the following screenshot:

A table that is used in classification problems to assess where errors in the model were made. The rows represent the actual classes the outcomes should have been. While the columns represent the predictions we have made

Visualization Analysis: We used the `countplot()` function of Seaborn to visualize the count of the two types of cells as shown above. We also created a pair plot of the dataset using the `pairplot()` function of Seaborn to visualize the distribution of features and their correlation with each other as shown above. The correlation matrix was visualized using the `heatmap()` function of Seaborn. These visualizations helped us to understand the distribution of data and the correlation between features.

Conclusion: The proposed machine-learning models could predict breast cancer as the early detection of this disease could help slow down the progress of the disease and reduce the mortality rate through appropriate therapeutic interventions at the right time. Applying different machine learning approaches, accessibility to bigger datasets from different institutions (multi-center study) and considering key features from a variety of relevant data sources could improve the performance of modeling.

References

- [1] "Intro to Machine Learning | Udacity." Intro to Machine Learning | Udacity. Accessed April 27, 2016. <https://www.udacity.com/course/intro-to-machine-learning--ud120>.
- [2] "Elements of Statistical Learning: Data Mining, Inference, and Prediction. 2nd Edition. Datasets:Coronary Heart Disease Dataset." Elements of Statistical Learning: Data Mining, Inference, and Prediction. 2nd Edition. Accessed April 27, 2016. <http://statweb.stanford.edu/~tibs/ElemStatLearn/>.
- [3] Schölkopf, Bernhard, Christopher J. C. Burges, and Alexander J. Smola. Advances in Kernel Methods: Support Vector Learning. Cambridge, MA: MIT Press, 1999.
- [4] Norving, Peter, and Stuart Russel. Artificial Intelligence: A Modern Approach. S.l.:Pearson Education Limited, 2013.
- [5] Witten, I. H., and Eibe Frank. Data Mining: Practical Machine Learning Tools and Techniques. Amsterdam: Morgan Kaufman, 2005.

