Mini Project Title

"BREAST CANCER DETECTION"

Submitted in partial fulfilment of the requirements of the degree

BACHELOR OF ENGINEERING IN COMPUTER ENGINEERING

By

Pooja Bhagat	104
Priyanka Korde	118
Raghuwardayal Maurya	123

Rohit Mishra 125

Supervisor

Prof. Manish Umale



Department of Computer Engineering

Lokmanya Tilak College of Engineering Koparkhairne, NaviMumbai - 400 709

University of Mumbai

(AY 2021-22)

CERTIFICATE

This is to certify that the Mini Project entitled "Breast Cancer Detection" is a Bonafide work of Pooja Bhagat(104), Priyanka Korde(118), Raghuwardayal Maurya(123), Rohit Mishra(125) submitted to the University of Mumbai in partial fulfilment of the requirement for the award of the degree of "Bachelor of Engineering" in "Computer Engineering".

Prof. Manish UmaleGuide

Prof. Rajendra Gawali

Head of Department

Dr. Vivek K. Sunnapawar

Principal

MIIII I I UJCCI APPIUVA	Mini Project A	pprova
-------------------------	----------------	--------

This Mini Project entitled "Breast Cancer Detection" by Pooja Bhagat(104), Priyanka Korde(118), Raghuwardayal Maurya(123), Rohit Mishra(125) is approved for the degree of Bachelor of Engineering in Computer Engineering.

Examiners

1	• • • • • •	• • • • •	••••	• • • • •	•••••		••••	••••
	(In	terna	al Ex	amin	er Na	ıme .	& Si	gn)

2.....

(External Examiner name & Sign)

Table of Contents

Chapter	Contents	Pg. No.
1.	Acknowledgements	05
2.	Introduction	06
	2.1) Introduction	06
	2.2) Problem Statement	07
3.	Literature Survey	08
	3.1) Literature Survey	08
	3.1) Survey of Existing System	09
4.	Proposed System	10
	4.1) System Architecture	10
	4.2) Datasets	11
	4.3) Results	12
	4.4) Conclusion	16
5.	References	17

Acknowledgement

I remain immensely obliged to Prof. Manish Umale for providing me with the idea

of this topic, and for his/her invaluable support in gathering resources for me either

by way of information or computer also his/her guidance and supervision which

made this project successful.

We would like to thank mini project Coordinators, Dr. S.K. Shinde Vice

Principal, Prof. Rajendra Gawli Head of Computer Engineering Department

and Dr. Vivek Sunnapwar Principal, LTCoE.

I am also thankful to faculty and staff of Computer Engineering Department and

Lokmanya Tilak College of Engineering, Navi Mumbai for their invaluable support.

I would like to say that it has indeed been a fulfilling experience for working out this

project topic

Pooja Bhagat

Priyanka Korde

Raghuwardayal Maurya

Rohit Mishra

5

INTRODUCTION

2.1 Introduction

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

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 undergoing unnecessary treatments. Thus, the correct diagnosis of BC and classification of patients into malignant or benign groups is the subject of much research. Because of its unique advantages in critical features detection from complex BC datasets, machine learning (ML), deep learning is widely recognized as the methodology of choice in BC pattern classification and forecast modeling.

2.3 Problem Statement

Breast Cancer is one of the leading cancer developed in many countries including India. Though the endurance rate is high — with early diagnosis 97% women can survive for more than 5 years. Statistically, the death toll due to this disease has increased drastically in last few decades. The main issue pertaining to its cure is early recognition. Hence, apart from medicinal solutions some Data Science solution needs to be integrated for resolving the death causing issue.

This analysis aims to observe which features are most helpful in predicting malignant or benign cancer and to see general trends that may aid us in model selection and hyper parameter selection. The goal is to classify whether the breast cancer is benign or malignant. To achieve this we have used Deep Learning classification methods to fit a function that can predict the discrete class of new input.

LITERATURE SURVEY

3.1 <u>Literature Survey</u>

Twenty-four recent research articles have been reviewed to explore the computational methods to predict breast cancer. The summaries of them are presented below. Chaurasia et al. developed prediction models of benign and malignant breast cancer. Wisconsin breast cancer data set was used. The dataset contained 699 instances, two classes (malignant and benign), and nine integer-valued clinical attributes such as uniformity of cell size. The researchers removed the 16 instances with missing values from the data set to become the data set of 683 instances. The benign were 458 (65.5%) and malignant were 241 (34.5%). The experiment was analyzed by the Waikato Environment for Knowledge Analysis (WEKA). Naive Bayes, RBF Network, and J48 are the three most popular data mining algorithms were used to develop the prediction models. The researchers used 10- fold cross-validation methods to measure the unbiased estimate of the three prediction models for performance comparison purposes. The models' performance evaluation was presented based on the methods" effectiveness and accuracy.

Experimental results showed that the Naive Bayes had gained the best performance with a classification accuracy of 97.36%; followed by RBF Network with a classification accuracy of 96.77% and the J48 was the third with a classification accuracy of 93.41%. In addition, the researchers conducted sensitivity analysis and specificity analysis of the three algorithms to gain insight into the relative contribution of the independent variables to predict survival. The sensitivity results indicated that the prognosis factor

3.2 Survey of Existing System

Author & Ref.	Method	Findings	Dataset
Shubham Sharma et al.	Random Forest, KNN	KNN was a good	Wisconsin Breast
[16]	and Naïve Bayes.	classifier in terms of	Cancer dataset from
		accuracy.	UCI Repository.
R. Preetha et al. [17]	Data Mining	Detect the hidden	Wisconsin breast
	techniques	cancer associated for	cancer dataset.
		classification.	
Majid Nawazet al. [18]	Deep Learning	It got 95.4% accuracy	BreakHis Dataset is
	Convolution neural	when compared with	used
	network	state-of-art models	
		and DenseCNN model	
	_ , .	u sed for this.	
NareshKhuriwal et al.	Deep learning	It achieved 98%	Mammogram MIAS
[19]		accuracy by using	database.
Ajay kumar et al. [20]	Classification	CNN. By using BCDW11, it	BCDW11 and
Ajay kumaret al. [20]	techniques like SVM,	gave 97.13% accuracy	WBCD32 dataset from
	KNN, Naïve Bayes	and using WBCD32,	UCI Repository.
	and Decision Tree.	SVM gave 97.89%	OCI Repository.
	and Decision free.	accuracy.	
Sri Hari Nallamala et al.	Machine learning	It achieved the	Wisconsin Breast
[21]	techniques	98.50% precision.	Cancer dataset.
R.Chtihrakkannan,	Machine learning	It achieved 96%	Mammogram images.
P.Kavitha et al. [22]	techniques	accuracy by using	
	•	DNN.	
Weal E.Fathy et al. [23]	Deep learning	It achieved 96% area	Digital Database for
		under ROC and 99.8%	Screening
		sensitivity and 82.1%	Mammography
		specificity.	dataset.
Nikita Rane et al. [24]	Machine learning	According to this,	Wisconsin Breast
	techniques	enhancement in	Cancer Dataset.
		machine learning gave	
D	D1	better results.	Donat Com
PamuwatMekha et al. [25]	Deep learning	The anthor compared	Breast Cancer
		the machine learning techniques and deep	Wisconsin dataset.
		learning. It achieved	
		the 96.99% accuracy	
		with deep learning.	
Mahmoud Khademiet al.	Probabilistic	It used the graphical	Netherlands Cancer
[26]	Graphical models and	model and deep belief	Institute dataset,
	deep belief network	network with manifold	METABRIC breast
		learning to find out the	cancer dataset,
		better accuracy.	Lju bljana breast
			cancer dataset and
			WDBC.

PROPOSED SYSTEM

4.1 System Architecture

In this project we will use Data Mining and Machine Learning Algorithms to detect breast cancer, based off of data. Breast Cancer (BC) is a common cancer for women around the world. Early detection of BC can greatly improve prognosis and survival chances by promoting clinical treatment to patients. We will use the UCI Machine Learning Repository for breast cancer dataset. Url:http://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+%28diagnosti

Url:http://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+%28diagnosti c%29 The dataset used in this story is publicly available and was created by Dr. William H. Wolberg, physician at the University Of Wisconsin Hospital at Madison, Wisconsin, USA. To create the dataset Dr. Wolberg used fluid samples, taken from patients with solid breast masses and an easy-to-use graphical computer program called Xcyt, which is capable of perform the analysis of cytological features based on a digital scan. The program uses a curve-fitting algorithm, to compute ten features from each one of the cells in the sample, than it calculates the mean value, extreme value and standard error of each feature for the image, returning a 30 real-valuated vector.

Attribute Information:

1. ID number 2) Diagnosis (M = malignant, B = benign) 3–32)

Ten real-valued features are computed for each cell nucleus:

- 2. radius (mean of distances from center to points on the perimeter)
- 3. texture (standard deviation of gray-scale values)
- 4. perimeter
- 5. area
- 6. smoothness (local variation in radius lengths)
- 7. compactness (perimeter² / area 1.0)
- 8. concavity (severity of concave portions of the contour)
- 9. concave points (number of concave portions of the contour)
- 10. symmetry
- 11. fractal dimension ("coastline approximation" 1)

The mean, standard error and "worst" or largest (mean of the three largest values) of these features were computed for each image, resulting in 30 features. For instance, field 3 is Mean Radius, field 13 is Radius SE, field 23 is Worst Radius.

4.2 Datasets

Data Set Information:

Features are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image.

Separating plane described above was obtained using Multisurface Method-Tree (MSM-T) [K. P. Bennett, "Decision Tree Construction Via Linear Programming." Proceedings of the 4th Midwest Artificial Intelligence and Cognitive Science Society, pp. 97-101, 1992], a classification method which uses linear programming to construct a decision tree. Relevant features were selected using an exhaustive search in the space of 1-4 features and 1-3 separating planes.

The actual linear program used to obtain the separating plane in the 3-dimensional space is that described in: [K. P. Bennett and O. L. Mangasarian: "Robust Linear Programming Discrimination of Two Linearly Inseparable Sets", Optimization Methods and Software 1, 1992, 23-34].

This database is also available through the UW CS ftp server: ftp ftp.cs.wisc.edu cd math-prog/cpo-dataset/machine-learn/WDBC/

Creators:

- 1. Dr. William H. Wolberg, General Surgery Dept. University of Wisconsin, Clinical Sciences Center Madison, WI 53792 wolberg '@' eagle.surgery.wisc.edu
- 2. W. Nick Street, Computer Sciences Dept. University of Wisconsin, 1210 West Dayton St., Madison, WI 53706 street '@' cs.wisc.edu 608-262-6619
- 3. Olvi L. Mangasarian, Computer Sciences Dept. University of Wisconsin, 1210 West Dayton St., Madison, WI 53706 olvi '@' cs.wisc.edu

Donor:

Nick Stree

4.3 Results

```
+ Code + Text
     ▼ Importing the packages
\{x\}
             import pandas as pd
import numpy as np
               from sklearn.datasets import load_breast_cancer
               from sklearn.preprocessing import StandardScaler
              import matplotlib.pyplot as plt
              import seaborn as sns
               from keras.models import Sequential
               from keras.layers import Dense
              import tensorflow as tf

▼ Loading the dataset

     [2] data = load_breast_cancer()

→ Having a look on dataset

      [3] data.keys()
               dict_keys(['data', 'target', 'frame', 'target_names', 'DESCR', 'feature_names', 'filename', 'data_module'])
     / [4] print(data['DESCR'])
              .. _breast_cancer_dataset:
              Breast cancer wisconsin (diagnostic) dataset
              **Data Set Characteristics:**
                   :Number of Instances: 569
                   :Number of Attributes: 30 numeric, predictive attributes and the class
                   :Attribute Information:

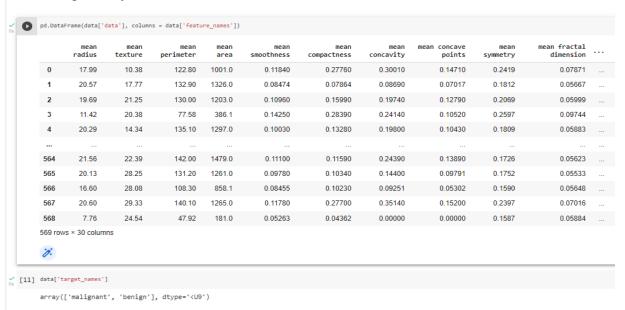
    radius (mean of distances from center to points on the perimeter)
    texture (standard deviation of gray-scale values)

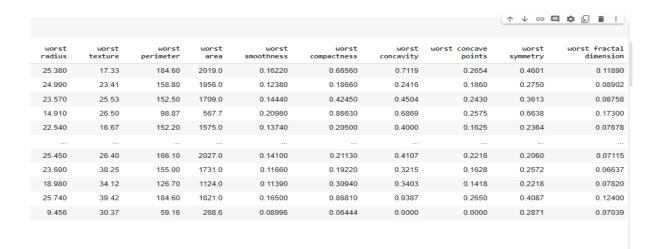
<>
                        - perimeter
                        - area
- smoothness (local variation in radius lengths)
- compactness (perimeter^2 / area - 1.0)
- concavity (severity of concave portions of the contour)
```

```
( [8] data['data'][-1]
     array([7.760e+00, 2.454e+01, 4.792e+01, 1.810e+02, 5.263e-02, 4.362e-02,
          7.000+00, 2.1000+00, 1.507e-01, 5.854e-02, 3.857e-01, 1.428e-00, 2.548e+00, 1.915e+01, 7.189e-01, 4.660e-01, 0.000e+00, 0.000e+00, 2.676e-02, 2.783e-03, 9.456e+00, 3.837e+01, 5.916e+01, 2.686e+02,
          8.995e-02, 6.444e-02, 0.000e+00, 0.000e+00, 2.871e-01, 7.039e-02])
[9] data['target'][::-1]
     1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1,
          0, 1, 1, 0, 1, 1, 1, 0, 1,
                                 1, 1, 0, 1, 1, 1,
          1, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
          0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0,
          0, 0, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 0, 0, 0, 1, 1, 0, 1, 1, 1, 1,
          1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1,
          1, 0, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 1, 0, 0, 0,
          0, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0,

    Converting the arrays to dataframe
```

▼ Converting the arrays to dataframe





```
+ Code + Text
Q
       Unsupported Cell Type. Double-Click to inspect/edit the content.
\{x\}

    Evaluating the model on testing data

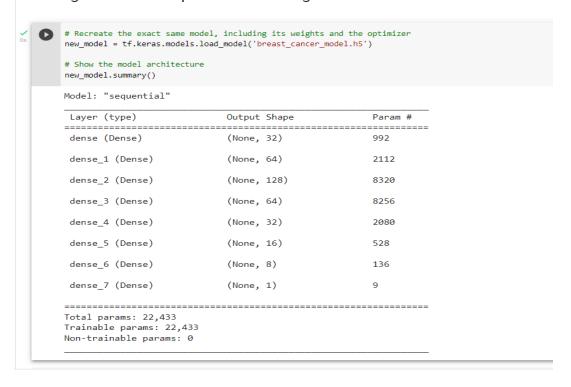
 [21] model.evaluate(x_test , y_test)
           2/2 [===========] - 0s 5ms/step - loss: 0.0096 - accuracy: 1.0000 [0.009568220935761929, 1.0]
    [22] model.evaluate(x_val, y_val)
           Got an Accuracy of 0.9706

    Display the predicted labels and actual label

       for i in range(10):
sample = x_test[i]
<>
             sample = np.reshape(sample, (1,30))
\blacksquare
             if (model.predict(sample)[0][0] > 0.5):
    print("-Benign")

✓ 0s completed at 1:57 PM
```

Using the model to perdict label for given features



```
+ Code + Text
:=
             ______
    v [25] Total params: 22,433
Q
            Trainable params: 22,433
            Non-trainable params: 0
{x}
           # 0 -> maglignat 1->beginine
features - pp. array([[1.288e+01, 2.892e+01, 8.250e+01, 5.143e+02, 8.123e-02, 5.824e-02, 6.195e-02, 2.343e-02, 1.566e-01, 5.788e-02, 2.116e-01, 1.360e+00, 1.502e+00, 1.683e+01, 8.412e-03, 2.153e-02, 3.898e-02, 7.620e-03, 1.695e-02, 2.801e-03, 1.389e+01, 3.574e+01, 8.884e+01, 5.957e+02,
                  1.227e-01, 1.620e-01, 2.439e-01, 6.493e-02, 2.372e-01, 7.242e-02]])
            prediction = new_model.predict(features)
print(prediction)
            if ( prediction > 0.5):
               print("-Benign")
print('NO fear of Heart Disease But for Better Understanding can Consult your Doctor!')
               print("-Malignant")
                print("There are Chances of Heart Disease! Consult your Doctor Soon!")
             Malignant
            There are Chances of Heart Disease! Consult your Doctor Soon!
       data['target'][::]
        <>
                   1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0,
                   >_
                                                                                           completed at 1:57 PM
```

4.4 Conclusion

Breast cancer detection is a challenging problem because it is most popular and harmful disease. Breast cancer is growing every year and there is less chance to recover from this disease. For detection of breast cancer, machine learning and deep learning techniques are used. It is concluded from the previous research, the machine learning techniques give better results in their own field.

The previous research is conducted through many machine learning techniques with some enhancement and augmentation in dataset for the better performance. But it is concluded that machine learning gives better results on linear data. It is also concluded from the previous research, when the data is in the form of images where the machine is failed. To solve the problem of machine learning techniques, an innovative technique is used. Deep learning is recently developed technique that frequently used in data science. For the classification of the breast cancer images data, a deep learning based technique CNN is used. CNN mostly works on the images dataset. In the previous research, it is also concluded that CNN gives better results as compared to machine learning techniques.

REFERENCES

https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic)

 $\frac{\text{https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0262349\#:} \sim :te}{xt=Deep\%20learning\%20is\%20one\%20of,the\%20diagnosis\%20of\%20breast\%} 20cancer.$

https://www.researchgate.net/publication/348604972_A_Review_Paper_on_Breast_Cancer_Detection_Using_Deep_Learning

https://www.researchgate.net/publication/342303246_Breast_Cancer_Detection_and_Prediction_using_Machine_Learning