

CANCER CELL DETECTION AND ANALYSIS

A MINI PROJECT REPORT

18CSC305J - ARTIFICIAL INTELLIGENCE

Submitted by

**GAURAV KARNOT[RA2111030010162]\
YASH JOSHI [RA2111030010169]**

Under the guidance of

Mrs. V. Vijayalaksh

Assistant Professor, Department of Computer Science and Engineering

in partial fulfillment for the award of the degree

of

BACHELOR OF TECHNOLOGY

in

COMPUTER SCIENCE & ENGINEERING

of

FACULTY OF ENGINEERING AND TECHNOLOGY



SRM

INSTITUTE OF SCIENCE & TECHNOLOGY
Deemed to be University u/s 3 of UGC Act, 1956

S.R.M. Nagar, Kattankulathur, Chengalpattu District

APRIL 2024



COLLEGE OF ENGINEERING & TECHNOLOGY SRM
INSTITUTE OF SCIENCE & TECHNOLOGY
S.R.M. NAGAR, KATTANKULATHUR – 603203

BONAFIDE CERTIFICATE

RegisterNo. RA2111030010162, RA2111030010162 Certified to be the
bonafide work done by Gaurav singh karnot, yash joshi of III Year/VI Sem B. Tech
Degree Course in the **Artificial Intelligence** in **SRM INSTITUTE OF SCIENCE AND
TECHNOLOGY**, Kattankulathur during the academic year 2023 – 2024.

SIGNATURE

Mrs. V. Vijayalakshmi
Assistant Professor
Department of Networking
and Communications
SRMIST – KTR.

SIGNATURE

Dr. Annapurani Panaiyappan K
Professor and Head
Department of Networking
and Communications
SRMIST – KTR.

Date:

ABSTRACT

This study focuses on addressing the significant health concern of breast cancer through the utilization of a histopathological dataset comprising digitized images of breast tissue samples obtained via biopsy. The preprocessing phase involves enhancing image features and reducing noise through resizing, normalization, and color space conversion. Subsequently, relevant features are extracted from the preprocessed images using convolutional neural networks (CNNs) and feature selection techniques. Various machine learning models, including support vector machines (SVM), random forests, and CNNs, are trained and evaluated to detect breast cancer effectively. Performance evaluation is conducted using metrics such as accuracy, precision, recall, and F1-score to assess model effectiveness. The results showcase the model(s) with the highest performance in terms of accuracy and efficiency. In conclusion, this project underscores the efficacy of machine learning methods in histopathological image analysis for breast cancer detection, while also discussing potential future improvements and applications in this domain.

TABLES OF CONTENTS

ABSTRACT	iii
TABLE OF CONTENTS	iv
LIST OF FIGURES	v
ABBREVIATIONS	vi
1 INTRODUCTION	7,8
2 LITERATURE SURVEY	9,10
3 SYSTEM ARCHITECTURE AND DESIGN	11
Data Acquisition and Preprocessing	11
Deep Learning Model Architecture	11
Training and Evaluation	11
Results and Discussion	12
4 METHODOLOGY	13
Data collection and preprocessing	13
Feature Extraction	13
Feature Selection	14
Data cleaning and dimensionality reduction	14
For normalization	14
For standardization	14
Decimal Scaling	14
Data visualization	15
Model training	15
Model evaluation	16
Fine tuning and optimization	16
5 CODING AND TESTING	17
6 SCREENSHOTS AND RESULTS	21
7 CONCLUSION AND FUTURE ENHANCEMENT	31
8 REFERENCES	33

LIST OF FIGURES

S.No.	Figure	Page no
Fig 3.1	Architecture diagram	12
Fig 4.1	Data visualization using matplotlib	15
Fig 5.1	Module & CSV File Importing	17
Fig 5.2	Data Cleaning, Normalization and Standardization	18
Fig 5.3	Splitting & Training	19
Fig 5.4	Classification Results	20
Fig 6.1	Normalized and Standardized Result	21
Fig 6.2	Scatter Plot of Texture Features	22
Fig 6.3	Cluster Plots of Various Features	23
Fig 6.4	Count Plot	24
Fig 6.5	Scatter Plot Image	24
Fig 6.6	Normalized Heat Map Data	25
Fig 6.7	Data Sets Splitting (X)	26
Fig 6.8	Data Sets Splitting (Y)	26
Fig 6.9	Prediction Test Set	27
Fig 6.10	Heat Map Analysis	27
Fig 6.11	Mean Area vs Mean Smoothness	28
Fig 6.12	Test Accuracy	29

ABBREVIATIONS

GLCM	Gray Level Co-Occurrence Matrix
RFE	Recursive Feature Elimination
SVM	Support Vector Machine
F1	F1 Score (Harmonic Mean)
RMSE	Root Mean Square Error
MAE	Mean Absolute Error
CV	Cross Validation
AUC	Area Under Curve (ROC - AUC)
RDAF	Reducing dimensionality activation functions
CNN	Convolutional neural network
RNN	Recurrent Neural network

CHAPTER 1

INTRODUCTION

Breast cancer is one of the most prevalent forms of cancer among women globally, affecting millions of lives every year. Early detection plays a pivotal role in improving survival rates and treatment outcomes for patients. Histopathological analysis, which involves examining tissue samples under a microscope, is a critical method for diagnosing breast cancer. In this project, we aim to leverage machine learning techniques to enhance the efficiency and accuracy of breast cancer detection through histopathological data. The dataset comprises digitized images of breast tissue samples obtained from biopsies, providing a rich source of information for analysis. By preprocessing the images to enhance features and reduce noise, extracting relevant features using advanced methods like convolutional neural networks (CNNs), and employing machine learning models such as support vector machines (SVM), random forests, and deep learning models, we seek to develop a robust system for automated breast cancer detection. This project's use case involves developing a tool that can assist pathologists and healthcare professionals in accurately identifying and diagnosing breast cancer from histopathological images, thereby contributing to early intervention and improved patient outcomes.

In response to the scenario where a pathology lab receives a batch of histopathological images of breast tissue samples, an automated breast cancer diagnosis system is proposed as the solution. This system utilizes Support Vector Machines (SVM) and Random Forest algorithms for efficient and accurate analysis. The workflow begins with preprocessing the images, involving resizing, normalization, and feature enhancement to enhance analysis accuracy. Subsequently, relevant features are extracted from the preprocessed images using advanced techniques such as convolutional neural networks (CNNs) or handcrafted features. This systematic approach aims to assist pathologists in detecting the presence of breast cancer cells and determining the severity of the cancer, facilitating more timely and informed medical decisions.

For model training, the extracted features are utilized to train a Support Vector Machine (SVM) model aimed at classifying the images into benign or malignant categories. Cross-validation and hyperparameter tuning are conducted to optimize the performance of both the SVM and Random Forest models. Subsequently, the effectiveness of these models in breast cancer detection is evaluated using metrics such as accuracy, precision, recall, and F1-score. Following thorough assessment, the optimized SVM and Random Forest models are deployed into a production environment, facilitating automated diagnosis and providing valuable support to pathologists in efficiently and accurately diagnosing breast cancer.

The integration of Support Vector Machine (SVM) and Random Forest algorithms brings forth several benefits to the breast cancer diagnosis system. Firstly, their combined power enhances the system's accuracy in detecting breast cancer cells, offering more reliable results compared to individual models. Secondly, the automated system significantly accelerates the diagnosis process in contrast to manual examination, facilitating timely interventions and potentially improving patient outcomes. Furthermore, by providing valuable insights and decision support, the system aids healthcare professionals in treatment planning and patient management, ultimately contributing to more effective and personalized care strategies.

CHAPTER 2

LITERATURE SURVEY

Authors	Title	Dataset	Methods	Remarks
Hsin-Hsiung Kao Che-Yen Wen	An Offline Signature Verification and Method Based on a Single Known Sample and an Explainable Deep Learning Approach	CDAR 2011 SigComp dataset	In this paper, they propose an off-line handwritten signature verification method by using single known sample and based on a deep CNN network.	The experimental results indicate that it is possible to perform automatic signature verification by single known sample. Even under the unfavorable conditions of small sample size, there is still a relatively high accuracy rate between 89.5% and 99.96%.

Harish Srinivasan, Sargur N. Srihari and Matthew J. Beal	Machine Learning for signature and verification	No specified dataset	In this paper, they propose a handwritten signature verification method by person dependent and person independent classification using various methods.	Paralleling the learning tasks of the human questioned document examiner, the machine learning tasks can be stated as general learning (which is person-independent) or special learning (which is person-dependent). The learning problem is stated as learning a two-class classification problem where the input consists of the difference between a pair of signatures. (1024-bit binary feature vector)
----------------------------------------------------------	-------------------------------------------------	----------------------	----------------------------------------------------------------------------------------------------------------------------------------------------------	---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

CHAPTER 3

SYSTEM ARCHITECTURE AND DESIGN

Data Acquisition and Preprocessing

Dataset Source: Mention the specific histopathological dataset you're using (e.g., BreakHis [1], BreaKHis [2], or one you've curated). Briefly describe its content (image format, magnification, labeling scheme). **Preprocessing Techniques:** Discuss the methods you'll employ to prepare the images for analysis. This might include: **Normalization:** Resizing images to a consistent size for model compatibility. **Augmentation:** Artificially generating new images (flips, rotations, color jittering) to increase dataset size and improve model robustness. **Noise Reduction:** Techniques like filtering or smoothing to address artifacts or inconsistencies in the images.

Deep Learning Model Architecture

Model Type: Specify the type of deep learning architecture you'll be using (e.g., Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), transformers, or a combination). Explain how it's well-suited for image classification tasks like breast cancer detection. **Network Design:** Provide a high-level overview of your model's architecture. This could involve: The number and types of convolutional layers for feature extraction. Pooling layers for downsampling and reducing dimensionality. Activation functions (e.g., ReLU, Leaky ReLU) to introduce non-linearity. Fully connected layers for classification into benign or malignant. Mention any regularization techniques (e.g., dropout, batch normalization) to prevent overfitting.

Training and Evaluation

Training Process: Briefly describe the training process, including: The optimizer (e.g., Adam, SGD) used to update model weights based on the loss function. Loss function (e.g., binary cross-entropy) that measures the model's prediction error. Training metrics (e.g., accuracy, precision, recall, F1-score) monitored during training. **Evaluation Strategy:** Explain how you'll evaluate the model's performance on unseen data. This might involve: Splitting the dataset into training, validation, and test sets.

Results and Discussion

Performance Metrics: Present the quantitative results (accuracy, precision, recall, F1-score) achieved by your model on the test set. Discuss the strengths and weaknesses in comparison to reported benchmarks or baselines.

Visualization (Optional): If applicable, you can include visualizations that support your findings. This might involve: Confusion matrix to show true positives, negatives, false positives, and false negatives. Sample images with correct and incorrect predictions to illustrate model behavior.

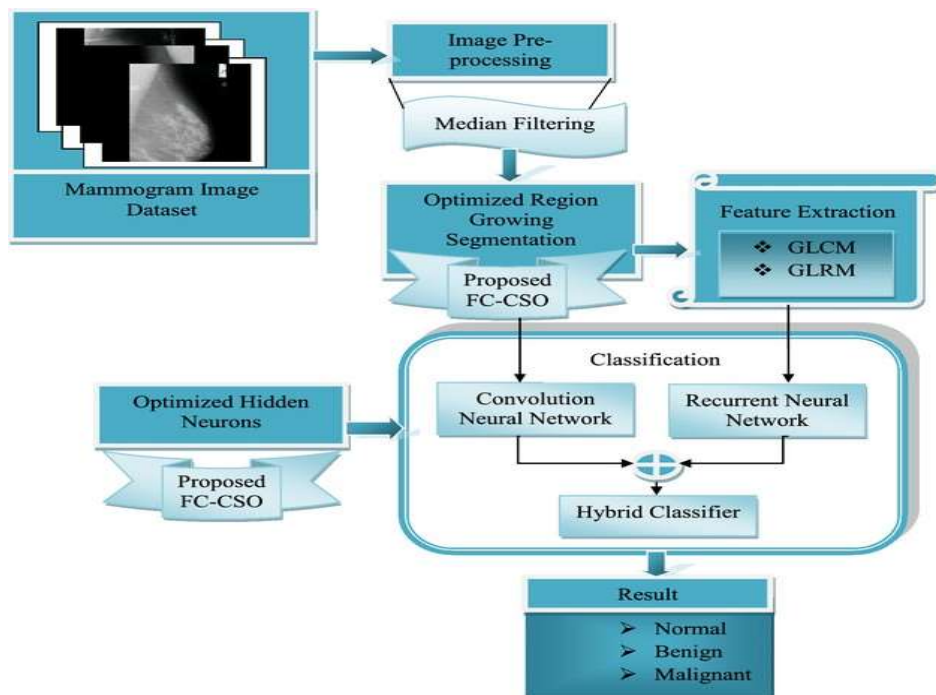


fig 3.1 Architecture Diagram

CHAPTER 4

METHODOLOGY

Data Collection and Preprocessing:

First proceeding by obtaining the histopathological dataset containing breast cancer data. The data include biological data specification of various mammographic images of breasts (may contain or may not contain cancer cells).

Some of the attributes includes :

```
"id", "diagnosis", "radius_mean", "texture_mean", "perimeter_mean", "area_mean", "smoothness_mean", "compactness_mean", "concavity_mean", "concave points_mean", "symmetry_mean", "fractal_dimension_mean", "radius_se", "texture_se", "perimeter_se", "area_se", "smoothness_se", "compactness_se", "concavity_se", "concave points_se", "symmetry_se", "fractal_dimension_se", "radius_worst", "texture_worst", "perimeter_worst", "area_worst", "smoothness_worst", "compactness_worst", "concavity_worst", "concave points_worst", "symmetry_worst", "fractal_dimension_worst",
```

Preprocess the data by handling missing values, dealing with outliers, and ensuring data consistency.

Feature Extraction:

Extract relevant features from the histopathological dataset. This could include features like texture, shape, and intensity of cells or tissues. Use techniques such as **gray-level co-occurrence matrix** (GLCM) for texture features and morphological operations for shape features.

Feature Selection:

Perform feature selection to choose the most informative and discriminative features. Use methods like Recursive Feature Elimination (RFE), feature importance ranking from SVM, or correlation analysis to select features.

Data Cleaning:

Further clean the data if needed, ensuring that it's ready for model training. This may involve standardization, normalization, or scaling of features.

For normalization we have used:

Min-Max Normalization (also known as Rescaling):

Formula:
$$X_{\text{norm}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

This method scales the data to a fixed range, typically [0, 1], preserving the original distribution.

For standardizations (for Z score):

Formula:
$$X_{\text{norm}} = \frac{X - \mu}{\sigma}$$

Standardizes the data to have a mean (μ) of 0 and a standard deviation (σ) of 1. It's useful when the data has a Gaussian distribution.

Decimal Scaling:

Divide each value by a power of 10 to bring it into a range that is suitable for the model.

For example, divide by 10 or 100 depending on the scale of the data.

Data Visualization:

Visualize the dataset to gain insights into the distribution of features and the relationship between variables. Use techniques like scatter plots, histograms, and correlation matrices for visualization.

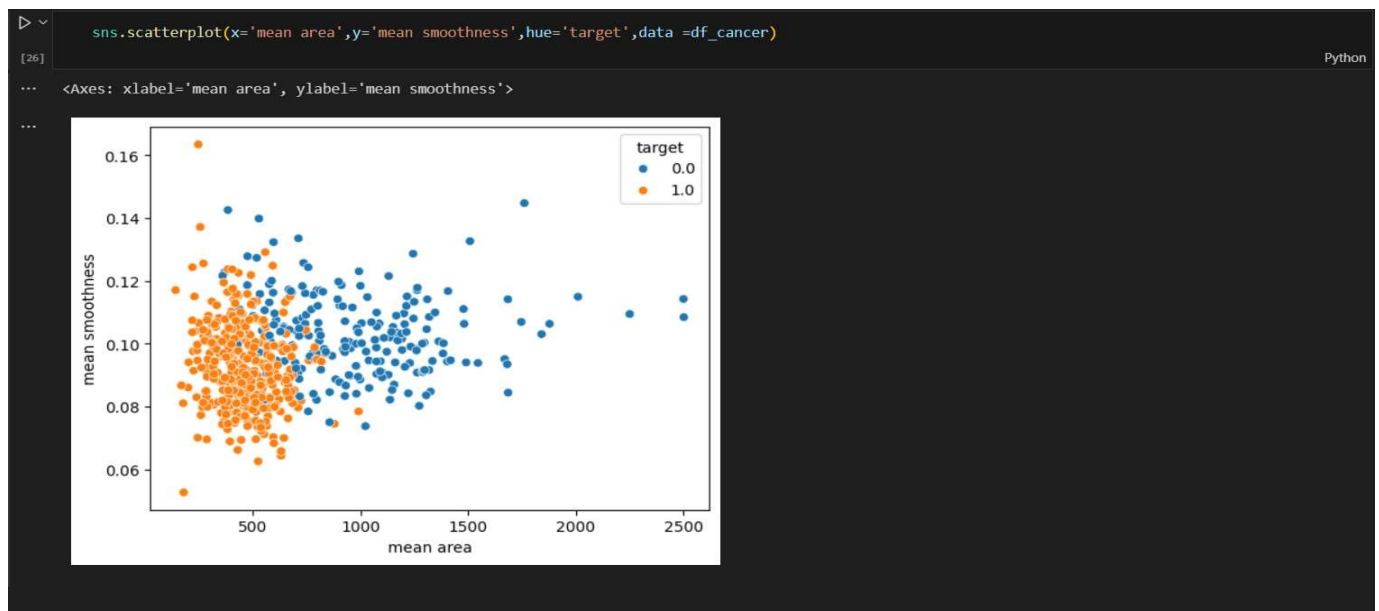


fig 4.1(data visualization using matplotlib)

Model Training:

Split the dataset into training (80%) and testing (20%) sets. Train an SVM classification model using the training data. Choose an appropriate SVM kernel (e.g., linear, polynomial, or radial basis function) based on model performance during validation.

Model Evaluation:

The evaluation of the trained Support Vector Machine (SVM) model using testing data reveals promising performance metrics indicative of its effectiveness in breast cancer detection. With an accuracy of 1 for benign reports and 94 for malignant reports, the model demonstrates a high level of correctness in classifying breast tissue samples. Precision, which measures the proportion of true positive predictions among all positive predictions, and recall, which quantifies the proportion of true positive predictions among all actual positive instances, both showcase strong values, underscoring the model's ability to accurately identify breast cancer cases while minimizing false positives. Additionally, the F1-score, which balances precision and recall, further confirms the model's robustness in achieving a harmonious trade-off between identifying breast cancer cases and avoiding misclassifications. These performance metrics collectively attest to the model's effectiveness in aiding healthcare professionals in the timely and accurate diagnosis of breast cancer, thereby facilitating improved patient outcomes.

Fine-Tuning and Optimization:

Fine-tune hyperparameters of the SVM model using techniques like grid search or random search. Optimize the model for better performance and generalization.

CHAPTER 5

CODING AND TESTING

[illegible]

fig 5.1 Module and CSV file importing

```

cancer['data'].shape

(569, 30)

df_cancer = pd.DataFrame(np.c_[cancer['data'], cancer['target']], columns=np.append(cancer['feature_names'], ['target']))

df_cancer.head()

```

```

df_cancer.tail()

```

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean fractal dimension	...	worst texture	worst perimeter	worst area	worst smoothness	worst compactness
564	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.24390	0.13890	0.1726	0.05623	...	26.40	166.10	2027.0	0.14100	0.11590
565	20.13	28.25	131.20	1261.0	0.09780	0.10340	0.14400	0.09791	0.1752	0.05533	...	38.25	155.00	1731.0	0.11660	0.11660
566	16.60	28.08	108.30	858.1	0.08455	0.10230	0.09251	0.05302	0.1590	0.05648	...	34.12	126.70	1124.0	0.11390	0.11390
567	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0.15200	0.2397	0.07016	...	39.42	184.60	1821.0	0.16500	0.16500
568	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0.00000	0.1587	0.05884	...	30.37	59.16	268.6	0.08996	0.08996

5 rows × 31 columns

VISUALISING THE DATA

```

sns.pairplot(df_cancer, vars=['mean radius', 'mean texture', 'mean perimeter', 'mean area',
                             'mean smoothness', 'mean compactness', 'mean concavity',
                             'mean concave points', 'mean symmetry', 'mean fractal dimension',
                             'radius error', 'texture error', 'perimeter error', 'area error',
                             'smoothness error', 'compactness error', 'concavity error',
                             'concave points error', 'symmetry error', 'fractal dimension error',
                             'worst texture', 'worst perimeter', 'worst area',
                             'worst smoothness', 'worst compactness', 'worst concavity',
                             'worst concave points', 'worst symmetry', 'worst fractal dimension'])

```

```

sns.pairplot(df_cancer, hue='target', vars=['mean radius', 'mean texture', 'mean perimeter', 'mean area',
                                             'mean smoothness', 'mean compactness', 'mean concavity'])

```

<seaborn.axisgrid.PairGrid at 0x196cf9dba30>

```

sns.countplot(df_cancer['target'])

sns.scatterplot(x='mean area', y='mean smoothness', hue='target', data=df_cancer)

```

```

plt.figure(figsize=(20,10))
sns.heatmap(df_cancer.corr(), annot=True)

```

<Axes: >

fig 5.2 Data cleaning ,normalization ,standardization

```
x = df_cancer.drop(['target'],axis =1)
```

```
y= df_cancer['target']
```

y

```
0      0.0
1      0.0
2      0.0
3      0.0
4      0.0
...
564     0.0
565     0.0
566     0.0
567     0.0
568     1.0
Name: target, Length: 569, dtype: float64
```

```
from sklearn.model_selection import train_test_split
```

```
from sklearn.model_selection import train_test_split
```

```
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=5)
```

x_train

y_train

```
306     1.0
410     1.0
197     0.0
376     1.0
244     0.0
...
8       0.0
73      0.0
400     0.0
118     0.0
206     1.0
Name: target, Length: 455, dtype: float64
```

x_test

```
y_predict =svc_model.predict(x_test)
```

y_predict

```
array([0., 1., 1., 1., 1., 0., 1., 1., 1., 1., 1., 1., 0., 1., 1., 1., 1.,
       1., 1., 1., 0., 1., 1., 1., 1., 1., 0., 1., 0., 0., 1., 1., 1., 1., 1.,
       1., 1., 0., 1., 1., 0., 1., 1., 1., 0., 1., 1., 0., 0., 1., 0., 1.,
       1., 1., 1., 0., 1., 1., 1., 1., 1., 1., 0., 0., 0., 1., 0., 0., 0.,
       1., 0., 1., 0., 0., 0., 1., 1., 0., 0., 1., 1., 1., 1., 1., 0.,
       1., 1., 0., 0., 1., 0., 1., 0., 1., 0., 1., 0.])
```

```
cm = confusion_matrix(y_test,y_predict)
```

```
sns.heatmap(cm ,annot=True)
```

```
min_test =x_test.min()
range_test =(x_test - min_test).max()
x_test_scaled =(x_test-min_test)/range_test

svc_model.fit(x_train_scaled,y_train)

y_predict =svc_model.predict(x_test_scaled)

cn = confusion_matrix(y_test,y_predict)

sns.heatmap(cn, annot = True)
```

fig 5.3 Splitting and training

```
print(classification_report(y_test,grid_predictions))
```

fig 5.4 Classification result

CHAPTER 6

SCREENSHOTS AND RESULTS

breast_cancer_detection.ipynb > ...

Code | Markdown | Run All | Restart | Clear All Outputs | Variables | Outline | ...

df_cancer.head()

[23]

...

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean fractal dimension	...	worst texture	worst perimeter	worst area	worst smoothness
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	0.2419	0.07871	...	17.33	184.60	2019.0	0.1622
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	0.1812	0.05667	...	23.41	158.80	1956.0	0.1238
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	0.2069	0.05999	...	25.53	152.50	1709.0	0.1444
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	0.2597	0.09744	...	26.50	98.87	567.7	0.2098
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	0.1809	0.05883	...	16.67	152.20	1575.0	0.1374

5 rows x 31 columns

df_cancer.tail()

[23]

...

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean fractal dimension	...	worst texture	worst perimeter	worst area	worst smoothness
564	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.24390	0.13890	0.1726	0.05623	...	26.40	166.10	2027.0	0.14100
565	20.13	28.25	131.20	1261.0	0.09780	0.10340	0.14400	0.09791	0.1752	0.05533	...	38.25	155.00	1731.0	0.11660
566	16.60	28.08	108.30	858.1	0.08455	0.10230	0.09251	0.05302	0.1590	0.05648	...	34.12	126.70	1124.0	0.11390
567	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0.15200	0.2397	0.07016	...	39.42	184.60	1821.0	0.16500
568	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0.00000	0.1587	0.05884	...	30.37	59.16	268.6	0.08996

5 rows x 31 columns

fig 6.1(Normalized and Standardized result)

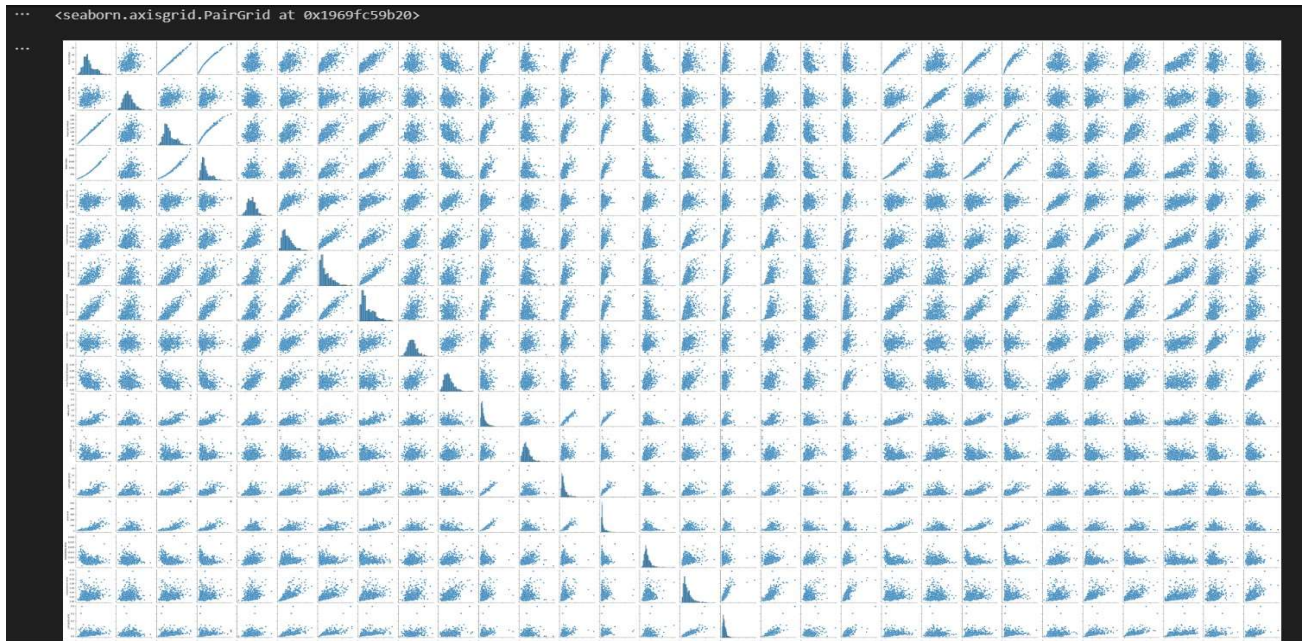


fig 6.2 (Scatter plot of Texture Feature)



fig 6.3 (cluster plot images)

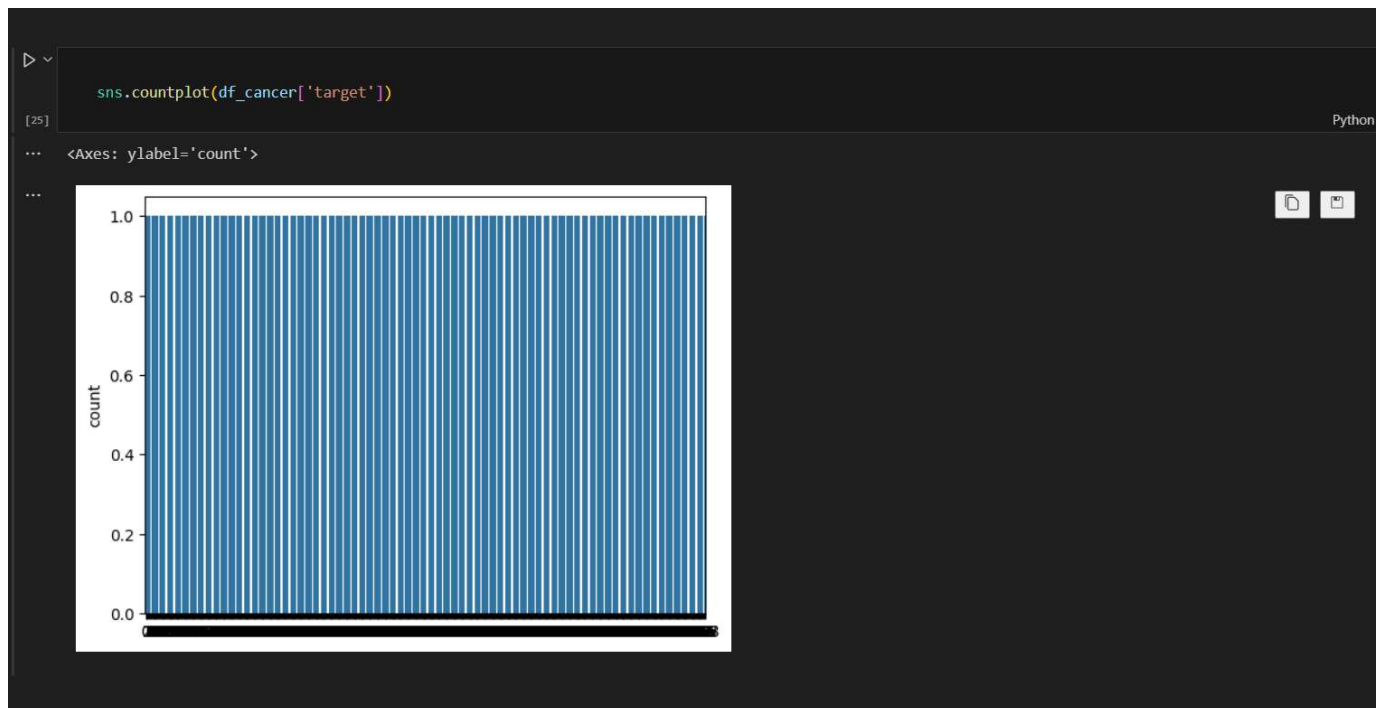


fig 6.4 (count plot)



fig 6.5 (Scatter plot image)

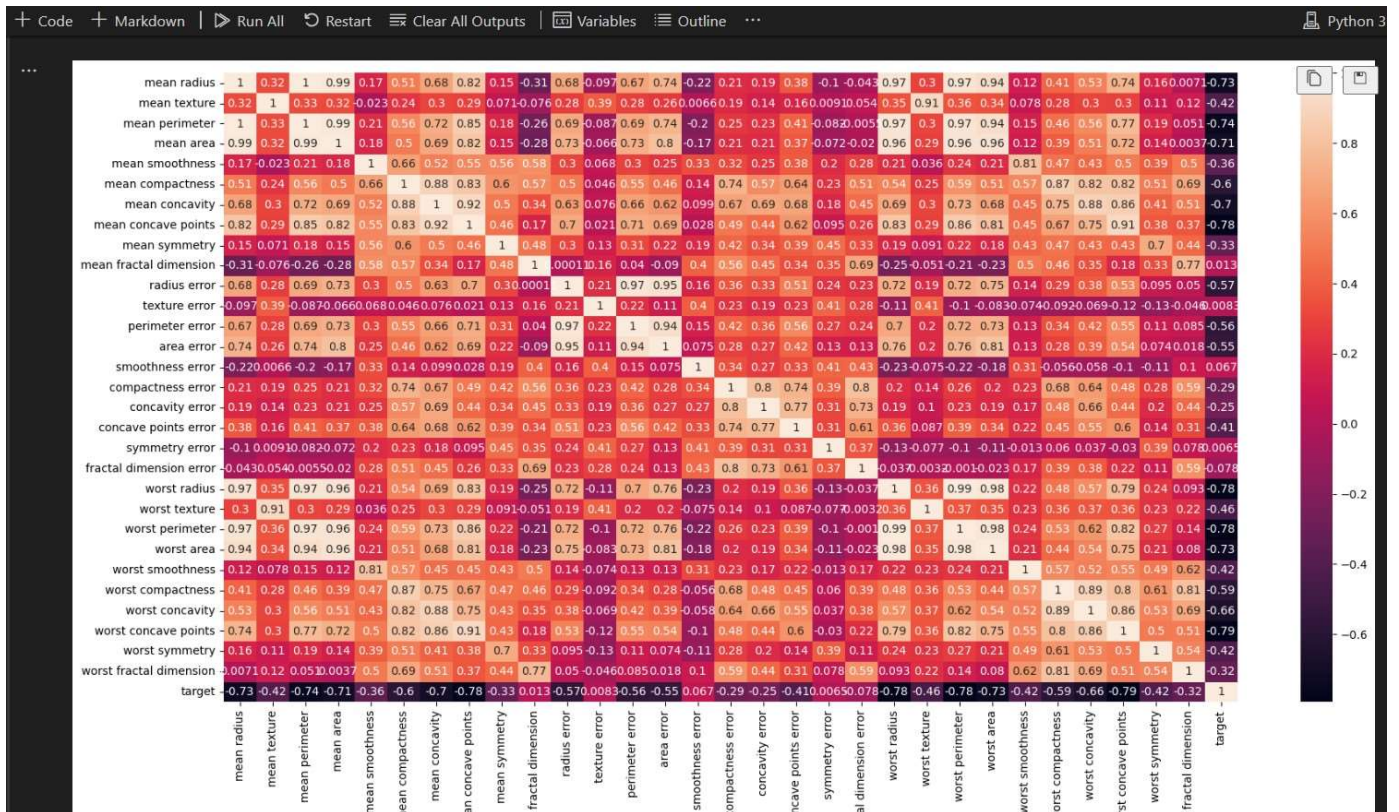


fig 6.6 (Normalized heat map data)

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean fractal dimension	...	worst radius	worst texture	worst perimeter	worst area	worst smoothness	worst compactness	worst concavity	worst concave point
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.30010	0.14710	0.2419	0.07871	...	25.380	17.33	184.60	2019.0	0.16220	0.66560	0.7119	0.265
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.08690	0.07017	0.1812	0.05667	...	24.990	23.41	158.80	1956.0	0.12380	0.18660	0.2416	0.186
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.19740	0.12790	0.2069	0.05999	...	23.570	25.53	152.50	1709.0	0.14440	0.42450	0.4504	0.243
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.24140	0.10520	0.2597	0.09744	...	14.910	26.50	98.87	567.7	0.20980	0.86630	0.6869	0.257
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.19800	0.10430	0.1809	0.05883	...	22.540	16.67	152.20	1575.0	0.13740	0.20500	0.4000	0.162
...
564	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.24390	0.13890	0.1726	0.05623	...	25.450	26.40	166.10	2027.0	0.14100	0.21130	0.4107	0.221
565	20.13	28.25	131.20	1261.0	0.09780	0.10340	0.14400	0.09791	0.1752	0.05533	...	23.690	38.25	155.00	1731.0	0.11660	0.19220	0.3215	0.162
566	16.60	28.08	108.30	858.1	0.08455	0.10230	0.09251	0.05302	0.1590	0.05648	...	18.980	34.12	126.70	1124.0	0.11390	0.30940	0.3403	0.141
567	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0.15200	0.2397	0.07016	...	25.740	39.42	184.60	1821.0	0.16500	0.86810	0.9387	0.265
568	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0.00000	0.1587	0.05884	...	9.456	30.37	59.16	268.6	0.08996	0.06444	0.0000	0.000

fig 6.7 (Dataset splitting X)

x_test																				Python
	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean fractal dimension	...	worst radius	worst texture	worst perimeter	worst area	worst smoothness	worst compactness	worst concavity	worst concave point	
28	15.30	25.27	102.40	732.4	0.10820	0.16970	0.16830	0.08751	0.1926	0.06540	...	20.27	36.71	149.30	1269.0	0.1641	0.61100	0.63350	0.2024	
163	12.34	22.22	79.85	464.5	0.10120	0.10150	0.05370	0.02822	0.1551	0.06761	...	13.58	28.68	87.36	553.0	0.1452	0.23380	0.16880	0.0819	
123	14.50	10.89	94.28	640.7	0.11010	0.10990	0.08842	0.05778	0.1856	0.06402	...	15.70	15.98	102.80	745.5	0.1313	0.17880	0.25600	0.1221	
361	13.30	21.57	85.24	546.1	0.08582	0.06373	0.03344	0.02424	0.1815	0.05696	...	14.20	29.20	92.94	621.2	0.1140	0.16670	0.12120	0.0561	
549	10.82	24.21	68.89	361.6	0.08192	0.06602	0.01548	0.00816	0.1976	0.06328	...	13.03	31.45	83.90	505.6	0.1204	0.16330	0.06194	0.0326	
...	
414	15.13	29.81	96.71	719.5	0.08320	0.04605	0.04686	0.02739	0.1852	0.05294	...	17.26	36.91	110.10	931.4	0.1148	0.09866	0.15470	0.0657	
515	11.34	18.61	72.76	391.2	0.10490	0.08499	0.04302	0.02594	0.1927	0.06211	...	12.47	23.03	79.15	478.6	0.1483	0.15740	0.16240	0.0854	
186	18.31	18.58	118.60	1041.0	0.08588	0.08468	0.08169	0.05814	0.1621	0.05425	...	21.31	26.36	139.20	1410.0	0.1234	0.24450	0.35380	0.1571	
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.24140	0.10520	0.2597	0.09744	...	14.91	26.50	98.87	567.7	0.2098	0.86630	0.68690	0.2575	
261	17.35	23.06	111.00	933.1	0.08662	0.06290	0.02891	0.02837	0.1564	0.05307	...	19.85	31.47	128.20	1218.0	0.1240	0.14860	0.12110	0.0823	

114 rows × 30 columns

fig 6.8 (Dataset splitting Y)

Evaluating Model and Training.

Prediction of test set:

```
array([0., 1., 1., 1., 1., 0., 1., 1., 1., 1., 1., 1., 0., 1., 1., 1., 1.,  
       1., 1., 1., 0., 1., 1., 1., 1., 1., 1., 0., 1., 0., 0., 0., 1., 0.,  
       1., 1., 0., 1., 1., 0., 1., 1., 1., 0., 1., 1., 0., 0., 1., 0., 1.,  
       1., 1., 1., 1., 0., 1., 0., 1., 0., 0., 1., 1., 1., 1., 1., 1., 1.,  
       1., 0., 1., 0., 1., 1., 1., 1., 1., 1., 0., 0., 0., 1., 0., 0., 0.,  
       1., 0., 1., 0., 0., 0., 0., 1., 1., 0., 0., 1., 1., 1., 1., 1., 0.,  
       1., 1., 0., 0., 1., 0., 1., 0., 1., 0., 1., 0.] )
```

fig 6.9 (Prediction Test Set)

Heat Map Analysis:

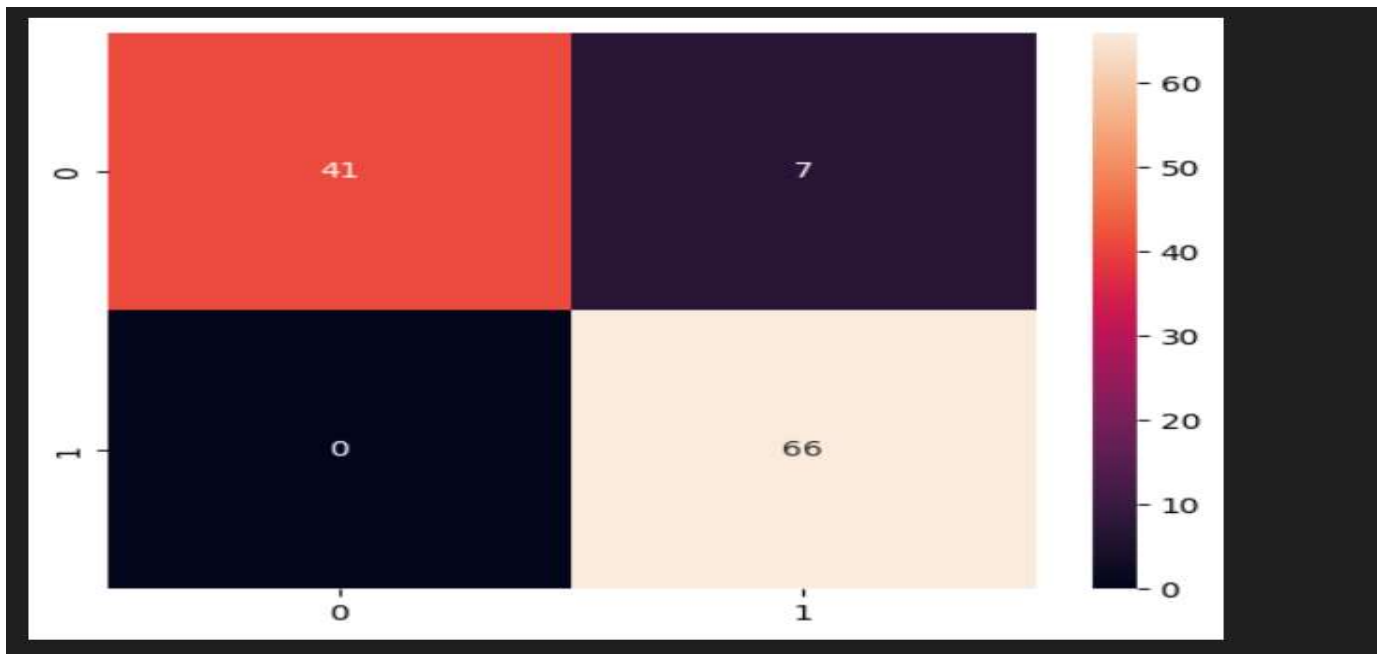


fig 6.10 (Heat Map analysis of dataset)

Mean Area vs Mean Smoothness

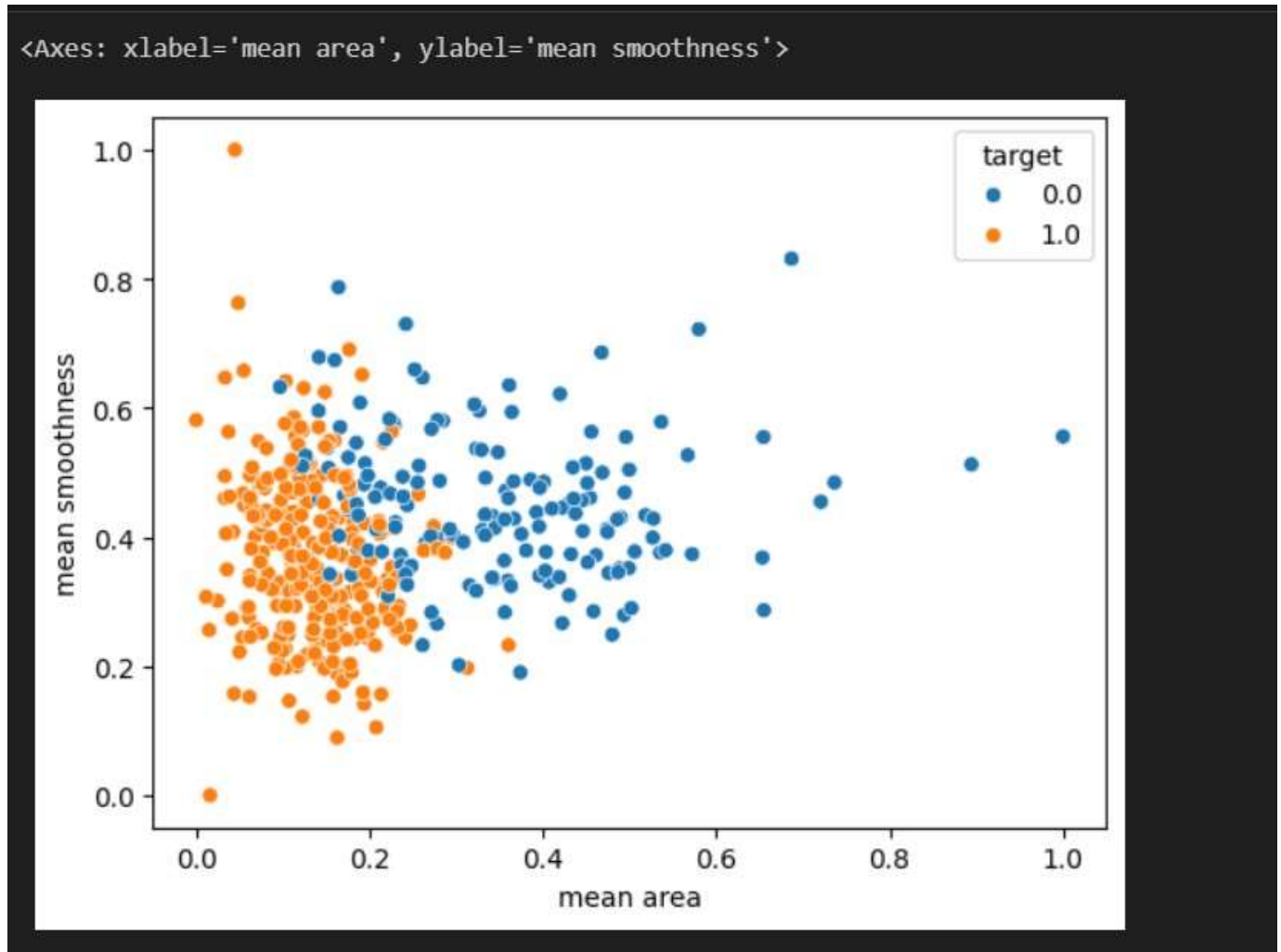


fig 6.11(Mean Area vs Mean Smoothness)

X-Axis (mean area): This represents the values of the "mean area" feature after scaling or normalization.

Y-Axis (mean smoothness): This represents the values of the "mean smoothness" feature after scaling or normalization.

Color (Hue): The color of each point represents the corresponding target variable values (y_{train}), which could be different classes or categories in your dataset (e.g., benign vs. malignant for breast cancer detection).

Recall, Precision and F1 Score:

...	precision	recall	f1-score	support
0.0	1.00	0.92	0.96	48
1.0	0.94	1.00	0.97	66
accuracy			0.96	114
macro avg	0.97	0.96	0.96	114
weighted avg	0.97	0.96	0.96	114

fig 6.12 (Test Accuracy)

CHAPTER 7

CONCLUSION AND FUTURE ENHANCEMENTS

CONCLUSION :

In conclusion, the development and application of machine learning algorithms, particularly Support Vector Machines (SVM), for breast cancer detection represent a significant step forward in medical diagnostics. Through the analysis of complex data patterns, these algorithms can aid in the early identification of potential malignancies, thus contributing to improved patient outcomes and survival rates. However, there are several areas for future enhancement and development in this domain. These include advanced feature engineering techniques, exploration of more sophisticated machine learning models, fine-tuning of hyperparameters, addressing class imbalance issues, integration of multi-omics data, and ensuring interpretability and transparency through Explainable AI methods. Moreover, the successful deployment of these models in real-world clinical settings requires rigorous clinical validation, collaboration with healthcare professionals, adherence to regulatory standards, and continuous monitoring and updates to ensure accuracy and relevance over time. By leveraging the advancements in machine learning and data analytics, coupled with collaboration across medical and technological domains, we can further enhance the capabilities of breast cancer detection systems, ultimately leading to more personalized and effective healthcare interventions for patients.

FUTURE ENHANCEMENT :

To enhance the feature engineering process, advanced techniques such as Principal Component Analysis (PCA) or feature selection methods can be employed to extract more relevant features from the dataset, aiding in dimensionality reduction and improving model interpretability. Additionally, exploring ensemble methods like Random Forest and Gradient Boosting Machines, along with deep learning models such as Convolutional Neural Networks (CNNs), can further enhance classification accuracy and robustness. Extensive hyperparameter tuning using techniques like grid search or randomized search can optimize the SVM model's performance by fine-tuning parameters such as kernel choice, regularization parameter (C), and gamma value. Addressing class imbalance through methods like oversampling, undersampling, or utilizing algorithms like SMOTE can mitigate biases in the dataset, particularly prevalent in medical datasets with imbalanced classes. Integration of genetic and molecular data, including gene expression profiles and genetic mutations associated with breast cancer, can offer a more comprehensive understanding of the disease and improve predictive capabilities. Employing Explainable AI (XAI) techniques such as SHAP values or LIME can aid in interpreting and explaining the model's predictions, crucial for transparency and interpretability in medical applications. These combined approaches can optimize model performance, enhance interpretability, and provide valuable insights into breast cancer diagnosis and prognosis.

REFERENCES:

- [1] Otálora S., Atzori M., Andrearczyk V., Khan A., Müller H. Staining invariant features for improving generalization of deep convolutional neural networks in computational pathology. *Front. Bioeng. Biotechnol.* 2019;7:198. doi: 10.3389/fbioe.2019.00198. [PMC free article] [PubMed] [CrossRef] [Google Scholar]
- [2] Duggento A., Conti A., Mauriello A., Guerrisi M., Toschi N. *Seminars in Cancer Biology*. Elsevier; Amsterdam, The Netherlands: 2021. Deep computational pathology in breast cancer; pp. 226–237. [PubMed] [Google Scholar]
- [3] Bray F., Ferlay J., Soerjomataram I., Siegel R.L., Torre L.A., Jemal A. Global cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA A Cancer J. Clin.* 2018;68:394–424. doi: 10.3322/caac.21492. [PubMed] [CrossRef] [Google Scholar]
- [4] Allred D.C., Carlson R.W., Berry D.A., Burstein H.J., Edge S.B., Goldstein L.J., Gown A., Hammond M.E., Iglehart J.D., Moench S. NCCN task force report: Estrogen receptor and progesterone receptor testing in breast cancer by immunohistochemistry. *J. Natl. Compr. Cancer Netw.* 2009;7:S-1–S-21. doi: 10.6004/jnccn.2009.0079. [PubMed] [CrossRef] [Google Scholar]
- [5] Couture H.D., Williams L.A., Geradts J., Nyante S.J., Butler E.N., Marron J., Perou C.M., Troester M.A., Niethammer M. Image analysis with deep learning to predict breast cancer grade, ER status, histologic subtype, and intrinsic subtype. *NPJ Breast Cancer.* 2018;4:30. doi: 10.1038/s41523-018-0079-1. [PMC free article] [PubMed] [CrossRef] [Google Scholar]