

A Mini Project Report on

“Cancer Detection”

Submitted in partial fulfillment of the requirements for the award of the degree of

Bachelor of Engineering

in

Computer Engineering

by

Nilesh Parab (D221452)
Ananya Yadav (D221427)
Riya Singh (D221421)
Vivek Tiwari(D221422)

Under the Guidance of

Prof. Uma Ade



Department of Computer Engineering
Watumull Institute of Electronics Engineering and Computer
Technology
Ulhasnagar

UNIVERSITY OF MUMBAI

Academic Year 2022-2023

Approval Sheet

This Mini Project Report entitled “*Cancer Detection*” Submitted by “*Nilesh Parab*” (D221452), “*Ananya Yadav*” (D221427), “*Riya Singh*” (D221421), “*Vivek Tiwari*” (D221422) is approved for the partial fulfillment of the requirement for the award of the degree of *Bachelor of Engineering* in *Computer Engineering* from *University of Mumbai*.

Under The Guidance of

Prof. Uma Ade

Prof. Dhanajay Raut

Head Department of Computer Engineering

Place: Ulhasnagar

Date: 11.11.2022

CERTIFICATE

This is to certify that the mini project entitled “*Cancer Detection*” Submitted by “*Nilesh Parab*” (D221452), “*Ananya Yadav*” (D221427), “*Riya Singh* ” (D221421), “*Vivek Tiwari*” (D221422) or the partial fulfillment of the requirement for award of a degree *Bachelor of Engineering in Computer Engineering*, to the University of Mumbai, is a bonafide work carried out during academic year 2022-2023.

Guide Name & Signature

Examiners:

1.

2.

Prof. Dhanajay Raut

Head Department of Computer Engineering

Principal

Place: Ulhasnagar

Date: 11.11.2022

Declaration

We declare that this written submission represents our ideas in our own words and where others' ideas or words have been included, we have adequately cited and referenced the original sources. We also declare that we have adhered to all principles of academic honesty and integrity and have not misrepresented or fabricated or falsified any idea/data/fact/source in our submission. We understand that any violation of the above will be cause for disciplinary action by the Institute and can also evoke penal action from the sources which have thus not been properly cited or from whom proper permission has not been taken when needed.

(Signature)

Nilesh Parab (D221452)
Ananya Yadav (D221427)
Riya Singh (D221421)
Vivek Tiwari(D221422)

Date: 11.11.2022

Abstract

Cancer has been described as a diverse illness with a wide range of subgroups. Early cancer diagnosis and prognosis are essential for clinical patient treatment, which has become a requirement in cancer research. Numerous research teams from the biomedical and bioinformatics fields have studied the use of machine learning (ML) techniques due to the significance of categorizing cancer patients into high or low risk groups. These methods have been applied in an effort to simulate the development and management of malignant diseases. Furthermore, their significance is demonstrated by the fact that ML algorithms can recognize important features in complicated datasets. Even though it is evident that the use of ML methods can improve our understanding of cancer progression, an appropriate level of validation is needed in order for these methods to be considered in the everyday clinical practice.

The predictive models discussed here are based on various supervised ML techniques as well as on different input features and data samples. Given the growing trend on the application of ML methods in cancer research, we present here the most recent publications that employ these techniques as an aim to model cancer risk or patient outcomes.

CONTENTS

1 Introduction.....	7
2 Problem Statement	8
3 Implementation and Approach	9
4 Technology Used	10
5 Predicting Cancer Stage	11
5.1 American Cancer Society	11
5.1.1 CT scan for Cancer	11
5.1.2 What does a CT scan show?	11
5.1.3 How does a CT scan work?	11
5.2 Cancer.Net.....	12
5.2.1 Why does cancer stage matter?	12
5.2.2 Doctors can also use staging to	12
6 Description of Dataset	13
6.1 Name	13
6.2 Source.....	13
6.2.1 Creators	13
6.2.2 Donar	13
6.2.3 Date	13
6.3 Data Set Information	13
6.4 Attribute Information	14
6.5 Dataset	15
7 Machine Learning Algorithms	16
7.1 Random Forest model	16
7.2 SVC	16
7.3 Selection of model	17
Source Code	18
Testing of Model	27
Conclusion	28

Chapter 1: Introduction

Over the past few decades, there has been continuous development in cancer research. To identify certain cancer types before symptoms emerge, researchers have employed a number of strategies, such as early-stage screening. They have also developed novel techniques for the early prognosis of cancer therapy. Thanks to the advancement of modern medical technology, massive amounts of data on cancer have been acquired and are now accessible to the scientific community. Making an accurate forecast of a disease's trajectory, however, is one of the most exciting and challenging issues facing doctors. As a result, scientists working on medical research increasingly frequently use ML approaches. Even though these datasets are complex, these approaches can identify patterns and correlations among them.

According to statistics from around the world, breast cancer (BC) is one of the most prevalent malignancies in women and accounts for a considerable portion of new cancer cases and cancer-related fatalities in today's society. Early BC diagnosis might encourage patients to receive prompt clinical care, which can dramatically improve their prognosis and likelihood of survival. Patients may avoid receiving therapies they do not need if benign tumors are classified more accurately. Therefore, substantial research is being done on how to correctly diagnose BC and classify people into groups that are malignant or benign. Machine learning (ML) is widely acknowledged as the preferred technology in BC pattern categorization and forecast modeling due to its distinct benefits in identifying essential characteristics from complex BC datasets. Data may be effectively categorized using methods like data mining and classification. Particularly in the medical industry, where those techniques are frequently utilized to make judgments through diagnosis and analysis.

Currently, a technician operating the CT scan scanner at a clinic reports the cancer stage, but this information may be inaccurate because the technician may not be qualified to predict cancer stage. Our objective is to develop a model that can accurately predict cancer malignancy based on qualities without the requirement for qualifying.

Chapter 2: Problem Statement

OncoBREAST Dx - Final Report

Patient Information

Identification Data

HRN	BRE00000273	Concavity mean	0.3001	Symmetry se	0.03003
ID Card:	Unknown	Concave points mean	0.1471	Fractal dimension se	0.008193
Personal Data		Symmetry mean	0.2419	Radius worst	25.38
Age (years)	47	Fractal dimension mean	0.07871	Texture worst	17.23
Personal Data		Radius se	1.095	Perimeter worst	184.6
Radius mean	17.99	Texture se	0.9053	Area worst	2019
Texture mean	10.38	Perimeter se	8.589	Smoothness worst	0.1622
Perimeter mean	122.8	Area se	153.4	Compactness worst	0.6656
Area mean	1001	Smoothness se	0.006389	Concavity worst	0.7119
Smoothness mean	0.1184	Compactness se	0.04604	Concave points worst	0.2654
Compactness mean	0.2776	Concavity se	0.05372	Symmetry worst	0.4601
		Concave points se	0.01587	Fractal dimension worst	0.1189

Outcome

Results

Some Tumor Markers are outside the reference range and suggest malignancy.

Comments

In absence of comorbidities that are well-known source of False Positive (FP) in healthy patients by increasing some Tumor Markers levels, whole Tumor Markers levels suggest malignancy.

Conclusions

WE SUGGEST BREAST CANCER. THE 2016 ICD-10-CM DIAGNOSIS CODE IS C50 (MALIGNANT NEOPLASM OF BREAST).

This report has been generated from the data entered on 19-12-2017 10:55:02 UTC/GMT.

Disclaimer

This Multiple Biomarkers Disease Activity Algorithm (MBDAA) for Breast Cancer has been developed for the exclusive use by healthcare professionals, and solely as a Clinical Decision Support System (CDSS), not as an unique element for diagnosis. The algorithm bears a Sensitivity = 82.3%, Specificity = 85.7%. Please note, negativity of the Tumor Markers does not exclude at 100% the possibility of a malignant epithelial tumor.

Bioprognos SL
 Benet Mateu, 40 - 08034 Barcelona (Spain)
 Email: info@bioprognos.com
 Website: www.bioprognos.com

Technical Responsible
 Lorena Lendinez
 CSD Scientific Officer
 Bioprognos SL

- Early diagnosis of cancer focuses on detecting symptomatic patients as early as possible so they have the best chance for successful treatment.
- The technical member simply identified the data's patience stages.
- The technical person is unaware about cancer sickness, cancer malignancies, or any other issues associated with them.
- Therefore, patience sometimes addresses the subject of death as well as some other disease linked problems.
- Currently, the machine operator is forecasting the findings of the CT scan reports.

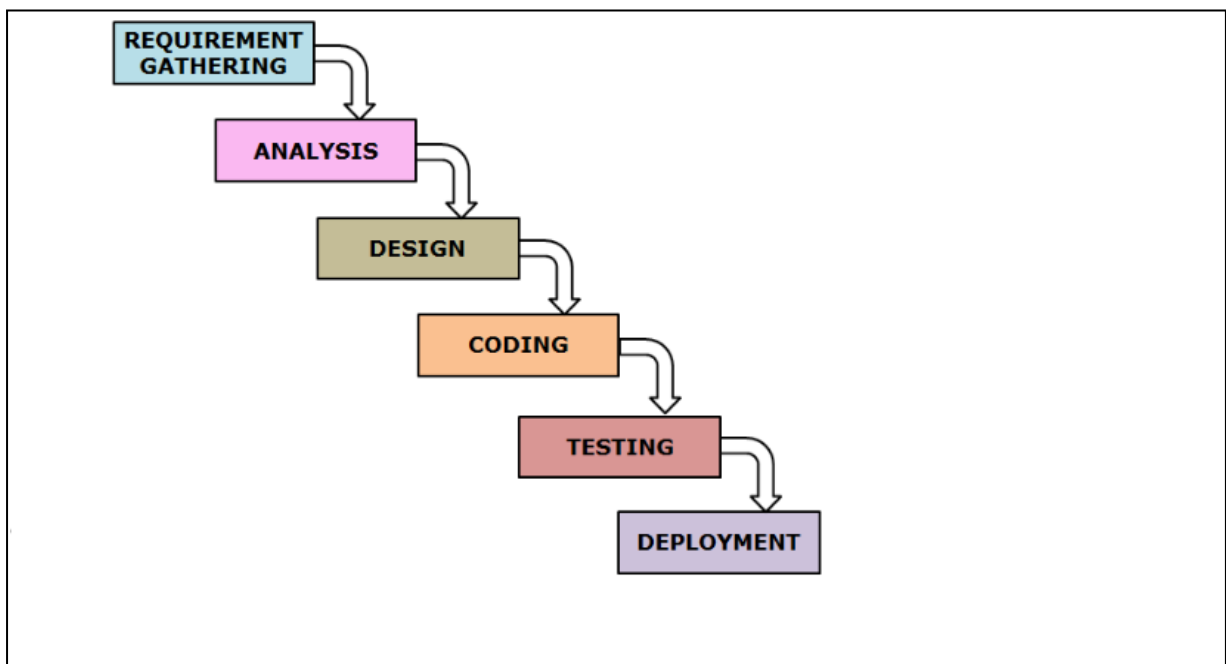
- Machine operators are not very skilled in identifying the stage of malignancy.
- In this situation, there is a strong probability that the outcome will be incorrect.

Chapter 3: Implementation and Approach

We started with choosing the tools first, in terms of how we'll build the web applications I started doing feasibility study about that we found python will work well for us in terms of strong predictions and analyzing. The data used in this project was downloaded from UCI We started finding some of the frontend development tools which can help us to give good look and feel for project, We discovered some of the tools like Google Collab, Jupyter Notebook etc. We decided to go with Python. Python is a cross-functional, maximally interpreted language that has lots of advantages to offer. The object-oriented programming language is commonly used to streamline large complex data sets. Over and above, having a dynamic semantics plus unmeasured capacities of RAD(*rapid application development*), Python is heavily utilized to script as well. We started deciding modules like analyzing the highest rate of crime taking place in which area. Now we had to choose my software development approach which will be better for my idea, I decided to choose waterfall model.

Waterfall Approach:

Development activities are performed in order, with possibly minor overlap, but with little or no iteration between activities. User needs are determined, requirements are defined, and the full system is designed, built, and tested for ultimate delivery at one point in time. A document driven approach best suited for highly precedence systems with stable requirements. The waterfall model is often also referred to as the linear and sequential model, for the flow of activities in this model are rather linear and sequential as the name suggests. In this model, the software development activities move to the next phase only after the activities in the current phase are over. However, like is the case with a waterfall, one cannot return to the previous stage



Chapter 4: Technology Used



Pandas:

Pandas is an open-source Python package that is most widely used for data science/data analysis and machine learning tasks. It is built on top of another package named NumPy, which provides support for multi-dimensional arrays.



NumPy

NumPy is a very popular python library for large multi-dimensional array and matrix processing, with the help of a large collection of high-level mathematical functions. It is very useful for fundamental scientific computations in Machine Learning



Matplotlib

In Machine learning, it helps to understand the huge amount of data through different visualizations.



scikit learn

It provides a selection of efficient tools for machine learning and statistical modeling including classification, regression, clustering and dimensionality reduction via a consistence interface in Python.

Chapter 5: Predicting Cancer Stage

5.1 American Cancer Society

5.1.1 CT scan for Cancer

A CT scan (also known as a computed tomography scan, CAT scan, and spiral or helical CT) can help doctors find cancer and show things like a tumor's shape and size. CT scans are most often an outpatient procedure. The scan is painless and takes about 10 to 30 minutes.

5.1.2 What does a CT scan show?

CT scans show a slice, or cross-section, of the body. The image shows your bones, organs, and soft tissues more clearly than standard x-rays.

CT scans can show a tumor's shape, size, and location. They can even show the blood vessels that feed the tumor – all without having to cut into the patient.

Doctors often use CT scans to help them guide a needle to remove a small piece of tissue. This is called a CT-guided biopsy. CT scans can also be used to guide needles into tumors for some types of cancer treatments, such as radiofrequency ablation (RFA), which uses heat to destroy a tumor.

By comparing CT scans done over time, doctors can see how a tumor is responding to treatment or find out if the cancer has come back after treatment.

5.1.3 How does a CT scan work?

In a way, CT scans are like standard x-ray tests. But an x-ray test aims a broad beam of radiation from only one angle. A CT scan uses a pencil-thin beam to create a series of pictures taken from different angles. The information from each angle is fed into a computer, which then creates a black and white picture that shows a slice of a certain area of the body – much like looking at a single slice from a loaf of bread.

Special contrast materials can be used to get a clearer picture. These can be swallowed as a liquid, put into a vein, or put into the intestines through the rectum as an enema.

By layering CT image slices on top of each other, the machine can create a 3-dimensional (3-D) view. The 3-D image can be rotated on a computer screen to look at different angles.

Doctors are now taking CT technology one step further in a technique called virtual endoscopy. They can look at the inside surfaces of organs such as the lungs (virtual bronchoscopy) or colon (virtual colonoscopy or CT colonography) without actually having to put scopes into the body. The 3-D CT images are arranged to create a black and white view on the computer screen. This looks a lot like it would if they were doing an actual endoscopy.

5.2 Cancer.Net

Staging is a way to describe a cancer. The cancer's stage tells you where a cancer is located and its size, how far it has grown into nearby tissues, and if it has spread to nearby lymph nodes or other parts of the body. Before starting any cancer treatment, doctors may use physical exams, imaging scans, and other tests to determine a cancer's stage. Staging may not be completed until all the tests are finished.

5.2.1 Why does cancer stage matter?

Staging helps your doctor plan the best treatment. This may include choosing a type of surgery and whether or not to use chemotherapy or radiation therapy. Knowing the cancer stage lets your entire health care team talk about your diagnosis in the same way.

5.2.2 Doctors can also use staging to:

- Understand the chance that the cancer will come back or spread after the original treatment.
- Help forecast the prognosis, which is the chance of recovery
- Help determine which cancer clinical trials may be open to you.
- See how well a treatment worked
- Compare how well new treatments work among large groups of people with the same diagnosis

Chapter 6: Description of Dataset

6.1 Name: Diagnostic Wisconsin Breast Cancer Database

6.2 Source:

6.2.1 Creators:

1. Dr. William H. Wolberg, General Surgery Dept.

University of Wisconsin, Clinical
Sciences Center

Madison, WI 53792

2. W. Nick Street, Computer
Sciences Dept.

University of Wisconsin, 1210 West
Dayton St., Madison, WI 53706

3. Olvi L. Mangasarian, Computer
Sciences Dept.

University of Wisconsin, 1210 West
Dayton St., Madison, WI 53706

6.2.2 Donor: Nick Street

6.2.3 Date: November 1995

6.3 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. A few of the images can be found at [Web Link]

Separating plane described above was obtained using Multi surface 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/

6.4 Attribute Information:

1) ID number

2) Diagnosis (M = malignant, B = benign)

3-32)

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

a) radius (mean of distances from centre to points on the perimeter)

b) texture (standard deviation of grey-scale values)

c) perimeter

d) area

e) smoothness (local variation in radius lengths)

f) compactness ($\text{perimeter}^2 / \text{area} - 1.0$)

g) concavity (severity of concave portions of the contour)

h) concave points (number of concave portions of the contour)

i) symmetry

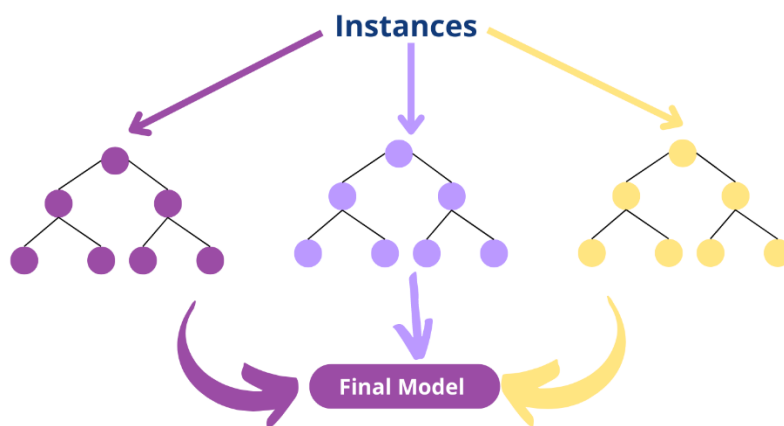
j) fractal dimension ("coastline approximation" - 1)

6.5 Dataset:

842002	1799	1038	1228	1001	011184	02776	30001	01471	02419	078781	1095	90053	8589	1534	0064	00844	05373	01587	030003	000619	2538	1733	1846	209	01822	06856	07119	02654	04601	01189		
842117	2057	1777	1322	1326	08474	07864	08860	07017	01812	05687	5045	97339	3398	7408	00023	01308	01866	0134	01389	000933	2499	2541	1888	1956	1238	1286	01416	0186	0735	08902		
842179	150	14607	179	1599	1458	0799	1458	0799	1458	0799	1458	0799	1458	0799	1458	0799	1458	0799	1458	0799	1458	0799	1458	0799	1458	0799	1458	0799	1458	0799		
842700	1142	208	7758	5861	01425	2839	02414	01052	2957	094474	40956	156	1445	2723	020911	02458	00661	01867	05980	0021	1491	265	9887	577	0298	08663	06869	02753	6658	0173		
842706	1249	154	1351	1297	01030	01328	0198	0403	01809	08083	07532	18803	2438	9444	01149	02461	05686	01885	07518	005012	2254	1607	1522	1755	01374	0205	04	01623	2394	01778		
842786	2025	147	837	1171	01278	017	01578	08089	02087	07613	0535	08802	2217	2739	01591	02453	05682	01017	02165	00508	1477	2575	1035	7416	01791	05249	05555	01741	3085	01424		
842804	1184	108	108	108	108	108	108	108	108	108	108	108	108	108	108	108	108	108	108	108	108	108	108	108	108	108	108	108	108	108		
842807	1371	2083	802	5777	0189	01645	08596	05985	0196	07451	05835	1377	8756	5036	00081	03029	02488	01448	0486	00041	1706	2106	1106	897	01654	03682	02878	01556	0396	01151		
844801	13	2182	875	5198	01275	01625	08596	05953	0235	07389	03060	1002	2406	2432	00072	03052	03553	01226	01241	00375	149	3073	1062	7393	0703	05401	03959	026	04378	01072		
845000	84	897	897	897	897	897	897	897	897	897	897	897	897	897	897	897	897	897	897	897	897	897	897	897	897	897	897	897	897	897		
845366	1602	2324	1027	7978	07806	06069	02399	03203	01528	05967	03795	187	2466	4051	00040	00927	01011	02759	0104	000304	1919	388	1238	150	0181	01551	01499	009975	0248	08452		
845378	1578	1789	1036	781	02901	01292	09954	06606	01842	06082	0558	09849	5564	5416	00577	00461	02791	01282	02008	00414	202	7728	365	1299	0396	05609	03965	0181	3792	0104		
846226	197	248	1324	1123	010974	02486	02065	01138	0297	0788	02255	3568	1117	1162	01341	02637	00889	0409	04484	01284	2096	294	1517	1332	0707	09003	03639	0167	3176	01023		
846381	1598	293	1037	782	0778	01002	09948	05364	01467	05838	01839	1078	9031	3658	01077	01814	05053	01992	00928	00019	2768	21	8795	1315	0194	01322	012	0289	02897	00687		
846397	1973	261	98	5783	0131	0249	0239	0218	0205	0208	07682	01251	169	1061	1921	00043	09938	05501	01828	0161	0809	105	3201	1087	6977	01651	07725	04603	0208	0396	01431	
846404	1534	274	9673	658	01139	0599	01809	0784	0203	07077	37	1033	1875	8255	00016	0424	0741	01009	01807	00547	147	3713	1241	9432	01678	06577	07026	01712	04218	01341		
846405	94	74	1072	5127	01437	01535	025417	01681	00841	00841	00841	00841	00841	94	74	1072	5127	01437	01535	025417	01681	00841	00841	94	74	1072	5127	01437	01535	025417	01681	
846406	1613	2068	1081	7988	0117	02022	01722	01028	0164	05736	05692	1075	8354	5418	00703	02051	03188	0297	01689	00414	2096	148	1368	153	11789	04238	04784	02073	0706	01142		
846407	181	251	150	1260	00883	01027	01479	09498	01582	05935	07582	107	5865	1124	00649	01893	03091	01521	03556	002	2732	3088	1886	2998	01512	0315	05372	0388	02768	07615		
8510426	134	1636	87	466	563	009779	02089	06664	04781	0185	085796	2699	07886	0508	00846	0146	02087	01515	0198	00203	151	1626	997	7114	0416	03773	0239	01888	02707	07259		
8510427	134	1636	87	466	563	009779	02089	06664	04781	0185	085796	2699	07886	0508	00846	0146	02087	01515	0198	00203	151	1626	997	7114	0416	03773	0239	01888	02707	07259		
8510428	134	1636	87	466	563	009779	02089	06664	04781	0185	085796	2699	07886	0508	00846	0146	02087	01515	0198	00203	151	1626	997	7114	0416	03773	0239	01888	02707	07259		
8511234	9504	124	6034	7739	01024	04962	02956	02076	0181	06905	0273	07968	1909	157	00961	01432	01985	0421	02027	000397	1023	1656	6513	149	01324	0118	00867	06227	0445	07996		
8511235	1534	1624	1023	704	01073	01335	0207	009736	0251	02073	04388	07906	3384	4491	00979	05328	04646	02252	03627	000439	1027	1908	1511	9809	0139	05954	06039	03923	0465	07996		
8511236	2116	2034	1372	1404	04828	01022	01097	06832	01769	05278	0807	1127	8035	8999	04473	01259	01575	01038	01088	00099	247	3559	188	3615	1401	014	0355	02009	02822	05765		
8511237	1598	293	1037	782	0778	01002	09948	05364	01467	05838	01839	1078	9031	3658	01077	01814	05053	01992	00928	00019	2768	21	8795	1315	0194	01322	012	0289	02897	00687		
852831	1714	164	116	1027	01186	0276	0229	0101	0404	07413	1046	976	2736	114	00803	07999	03732	0297	01689	00414	2096	148	1368	153	11789	04238	04784	02073	0706	01142		
852832	1458	213	9741	646	01004	01868	01425	08783	0252	09824	2345	9832	211	2105	04045	00955	02088	01892	01544	00371	1762	3321	1224	8899	01528	06843	05509	0701	04264	01075		
852833	1458	213	9741	646	01004	01868	01425	08783	0252	09824	2345	9832	211	2105	04045	00955	02088	01892	01544	00371	1762	3321	1224	8899	01528	06843	05509	0701	04264	01075		
852834	1534	1624	1023	704	01073	01335	0207	009736	0251	02073	04388	07906	3384	4491	00979	05328	04646	02252	03627	000439	1027	1908	1511	9809	0139	05954	06039	03923	0465	07996		
853001	1863	2511	1248	1038	01044	01887	0239	01144	01783	06919	83007	1446	5574	105	00025	03034	01936	0158	02007	00456	2315	1401	1605	1870	0491	02497	01633	01848	0446	07919		
853002	1158	77	58	118	01037	01335	0207	009736	0251	02073	04388	07906	3384	4491	00979	05328	04646	02252	03627	000439	1027	1908	1511	9809	0139	05954	06039	03923	0465	07996		
853003	1507	258	1298	1128	00983	01197	0466	0217	01203	0248	06382	0609	1399	999	6778	02027	0308	05042	0112	02102	00385	2088	3209	161	1344	0164	0359	05588	01847	0535	08482	
853004	1927	247	167	1279	1162	00401	01719	01657	03979	01853	06821	0558	0662	528	617	00502	03318	0497	00664	0554	0039	245	309	1614	0183	0509	0599	06091	01785	03672	01123	
853005	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	107	
854523	1674	2510	1011	8695	0106	0356	0138	0164	06018	0186	06566	0451	0937	3008	4519	00578	02499	03695	0195	02789	00627	2001	2902	1355	1229	1565	0385	05409	01813	0483	08633	
854524	1425	212	9363	63	009633	01098	0139	0519	05598	01885	02615	0286	019	2657	2491	00588	02995	04815	01161	02028	00402	1589	3036	1562	7996	0446	0428	0186	0447	3591	01014	
854941	6	1303	82	61	5238	00893	02766	02562	02933	0147	05883	01839	1078	9031	3658	01077	01814	05053	01992	00928	00019	2768	21	8795	1315	0194	01322	012	0289	02897	00687	
855133	1499	252	9554	698	01809	0097	05131	02398	02899	01565	05054	1214	2188	8077	106	00808	01094	01818	0197	00788	00715	1499	252	9554	698	01809	0097	05131	02398	02899	01565	05054
855134	134	158	884	563	01082	01335	0207	009736	0251	02073	04388	07906	3384	4491	00979	05328	04646	02252	03627	000439	1027	1908	1511	9809	0139	05954	06039	03923	0465	07996		
855167	1534	1624	1023	704	01073	01335	0207	009736	0251	02073	04388	07906	3384	4491	00979	05328	04646	02252	03627	000439	1027	1908	1511	9809	0139	05954	06039	03923	0465	07996		
855205	1095	2135	719	1017	01217	01128	0104	05669	01895	0487	02366	1428	1822	1822	1822	00082	01764	02039	0107	01357	00039	1284	554	8722	514	0909	02698</					

Chapter 7: Machine Learning Algorithms

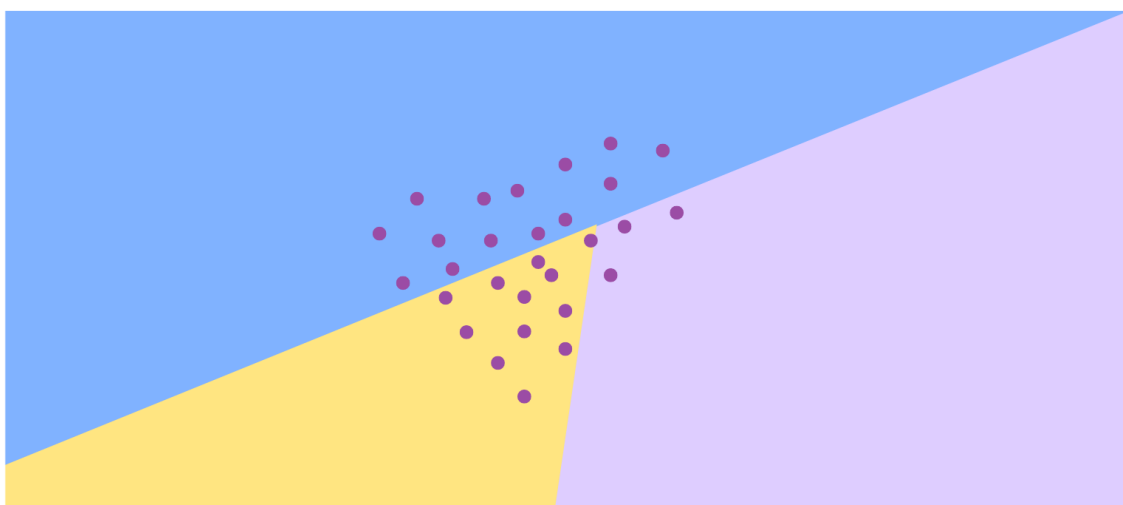
7.1 Random Forest Classifier



Random Forest is a popular machine learning algorithm that belongs to the supervised learning technique. It can be used for both Classification and Regression problems in ML. It is based on the concept of ensemble learning, which is a process of combining multiple classifiers to solve a complex problem and to improve the performance of the model.

As the name suggests, "Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset." Instead of relying on one decision tree, the random forest takes the prediction from each tree and based on the majority votes of predictions, and it predicts the final output.

7.2 SVC



The most applicable machine learning algorithm for our problem is Linear SVC. Before hopping into Linear SVC with our data, we're going to show a very simple example that should help solidify your understanding of working with Linear SVC.

The objective of a Linear SVC (Support Vector Classifier) is to fit to the data you provide, returning a "best fit" hyperplane that divides, or categorizes, your data. From there, after getting the hyperplane, you can then feed some features to your classifier to see what the "predicted" class is. This makes this specific algorithm rather suitable for our uses, though you can use this for many situations. Let's get started.

7.3 Selection of model

Random Forest

98.24%

Accuracy

SVC

97.36%

Accuracy

For random forest, we obtained an accuracy of 98.24%, while for SVC, we obtained an accuracy of 97.36%. As a result, we decided to use random forest as our primary model because it provided us with better accuracy.

Chapter 8: Source Code

8.1 Importing Libraries:

```
In [1]: %matplotlib inline
import numpy as np
import pandas as pd
import seaborn as s
from sklearn import model_selection
from sklearn.model_selection import train_test_split
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
from sklearn.ensemble import RandomForestClassifier
from sklearn import metrics
from sklearn.metrics import confusion_matrix, classification_report, accuracy_score
from sklearn import svm
import matplotlib.pyplot as plt
from sklearn.preprocessing import StandardScaler, LabelEncoder
#optimum parameter choosing
from sklearn.model_selection import GridSearchCV
from sklearn.svm import SVC
from xgboost import XGBClassifier
import pickle
import os
import warnings
from pandas import MultiIndex, Int64Index
warnings.filterwarnings('ignore')
```

8.2 Importing Dataset:

```
In [2]: os.chdir('C:\\Users\\Nilesh\\Documents\\Cancer Detection Project')
data = pd.read_csv('data.csv')
data
```

```
Out[2]:
```

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	...
0	842302	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.30010	0.14710	...
1	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.08690	0.07017	...
2	84300903	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.19740	0.12790	...
3	84348301	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.24140	0.10520	...
4	84358402	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.19800	0.10430	...
...
564	926424	M	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.24390	0.13890	...
565	926682	M	20.13	28.25	131.20	1261.0	0.09780	0.10340	0.14400	0.09791	...
566	926954	M	16.60	28.08	108.30	858.1	0.08455	0.10230	0.09251	0.05302	...
567	927241	M	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0.15200	...
568	92751	B	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0.00000	...

569 rows × 32 columns

8.3 Creating x variable:

```
In [9]: x= df.drop (labels='diagnosis' ,axis =1 )
x
```

Out[9]:

	id	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean
0	842302	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.30010	0.14710	0.2415
1	842517	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.08690	0.07017	0.1812
2	84300903	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.19740	0.12790	0.2065
3	84348301	11.42	20.38	77.58	386.1	0.14250	0.28390	0.24140	0.10520	0.2597
4	84358402	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.19800	0.10430	0.1805
...
564	926424	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.24390	0.13890	0.1725
565	926682	20.13	28.25	131.20	1261.0	0.09780	0.10340	0.14400	0.09791	0.1752
566	926954	16.60	28.08	108.30	858.1	0.08455	0.10230	0.09251	0.05302	0.1590
567	927241	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0.15200	0.2397
568	92751	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0.00000	0.1587

569 rows × 31 columns

8.4 Normalizing x variable:

```
In [13]: df_norm = (x- x.mean()) / (x.max()- x.min())
df_norm= pd.concat ([df_norm,y], axis =1 )
df_norm
```

Out[13]:

	id	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean
0	0.032403	0.182815	-0.301307	0.213053	0.146813	0.198968	0.531437	0.495081	0.487976	0.306758
1	-0.032403	0.304923	-0.051392	0.282848	0.284671	-0.104905	-0.078833	-0.004450	0.105621	0.000190
2	0.059177	0.263274	0.066295	0.262808	0.232497	0.119524	0.170416	0.254453	0.392549	0.129997
3	0.059229	-0.128132	0.036874	-0.099434	-0.114014	0.416536	0.550761	0.357546	0.279726	0.396657
4	0.059241	0.291671	-0.167388	0.298051	0.272369	0.035567	0.087292	0.255859	0.275253	-0.001320
...
564	-0.032311	0.351778	0.104848	0.345733	0.349570	0.132163	0.035455	0.363404	0.447221	-0.043240
565	-0.032311	0.284098	0.303022	0.271101	0.257099	0.012997	-0.002886	0.129336	0.243493	-0.030110
566	-0.032310	0.117029	0.297273	0.112853	0.086198	-0.106620	-0.006260	0.008694	0.020382	-0.111925
567	-0.032310	0.306342	0.339545	0.332603	0.258796	0.193552	0.529596	0.615278	0.512330	0.295647
568	-0.032326	-0.301353	0.177557	-0.304395	-0.201013	-0.394785	-0.186249	-0.208058	-0.243137	-0.113440

569 rows × 32 columns

8.5 Creating y variable:

```
In [10]: y = df['diagnosis']
y
```

Out[10]:

```
0    M
1    M
2    M
3    M
4    M
..
564  M
565  M
566  M
567  M
568  B
Name: diagnosis, Length: 569, dtype: category
Categories (2, object): ['B', 'M']
```

8.6 Normalizing y variable:

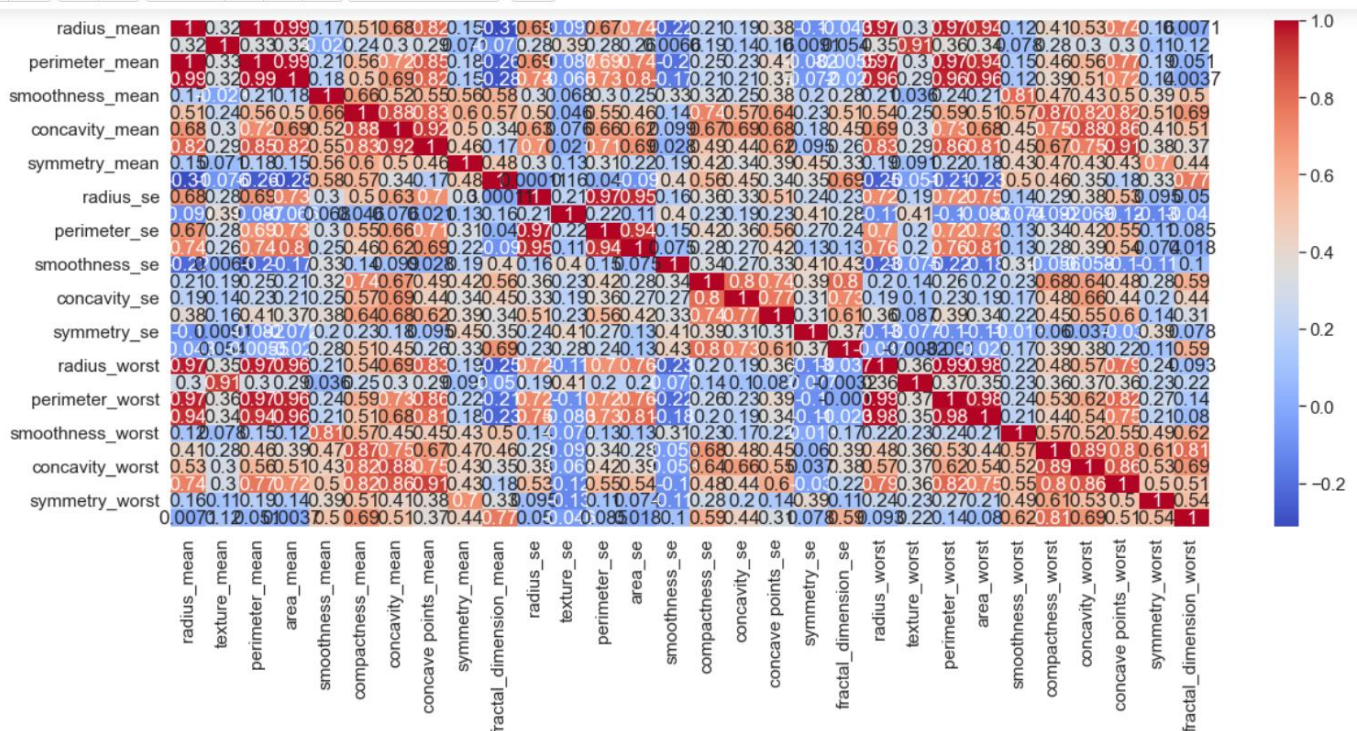
```
In [21]: le = LabelEncoder()
le.fit(y_norm)
y_norm = le.transform(y_norm)
y_norm = pd.DataFrame(y_norm)
print(y_norm)
```

```
0
0 1
1 1
2 1
3 1
4 1
.. ..
564 1
565 1
566 1
567 1
568 0

[569 rows x 1 columns]
```

8.7 Heatmap:

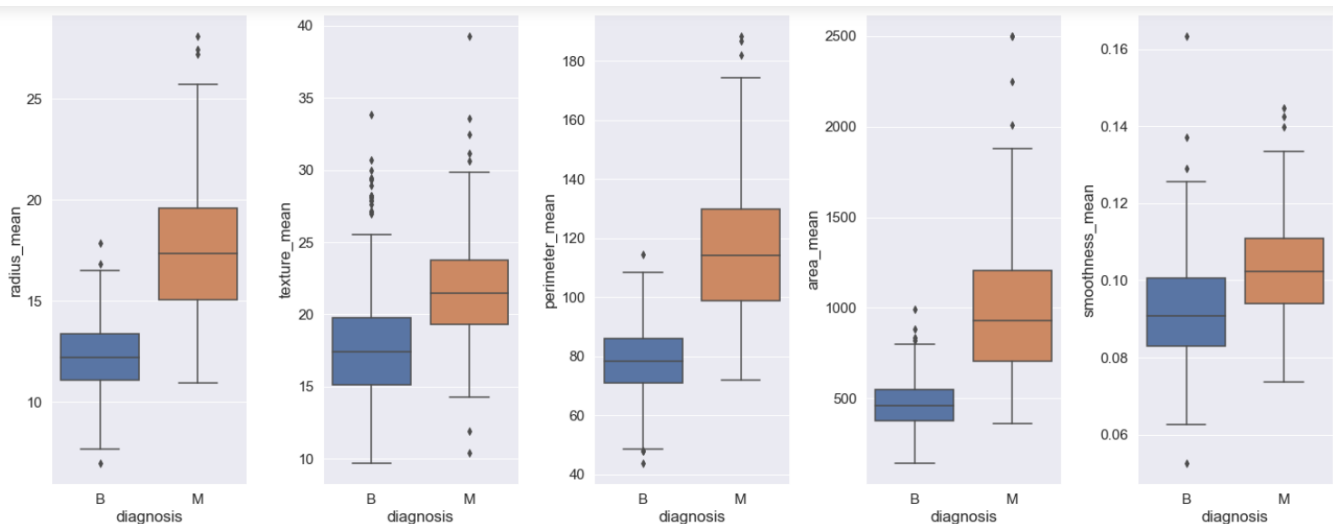
```
In [15]: plt.rcParams['figure.figsize']=(20,8)
s.set(font_scale=1.4)
# In co relation 1 is the highest and -1 is lowest
s.heatmap(df.drop('diagnosis',axis =1).drop('id',axis =1).corr(),cmap = 'coolwarm',annot = True)
```



8.8 Boxplot:

```
In [16]: plt.rcParams['figure.figsize']=(20,8)
f, (ax1,ax2,ax3,ax4,ax5) = plt.subplots (1,5)
s.boxplot ('diagnosis', y = 'radius_mean',data = df , ax = ax1)
s.boxplot ('diagnosis', y = 'texture_mean',data = df , ax = ax2)
s.boxplot ('diagnosis', y = 'perimeter_mean',data = df , ax = ax3)
s.boxplot ('diagnosis', y = 'area_mean',data = df , ax = ax4)
s.boxplot ('diagnosis', y = 'smoothness_mean',data = df , ax = ax5)
f.tight_layout()

f, (ax1,ax2,ax3,ax4,ax5) = plt.subplots (1,5)
s.boxplot ('diagnosis', y = 'compactness_mean',data = df , ax = ax1)
s.boxplot ('diagnosis', y = 'concavity_mean',data = df , ax = ax2)
s.boxplot ('diagnosis', y = 'concave points_mean',data = df , ax = ax3)
s.boxplot ('diagnosis', y = 'symmetry_mean',data = df , ax = ax4)
s.boxplot ('diagnosis', y = 'fractal_dimension_mean',data = df , ax = ax5)
f.tight_layout()
```



8.9 Distplot:

```
In [17]: g = s.FacetGrid (df,col = 'diagnosis', hue = 'diagnosis')
g.map (s.distplot, "radius_mean", hist = False, rug = True)

g = s.FacetGrid (df,col = 'diagnosis', hue = 'diagnosis')
g.map (s.distplot, "texture_mean", hist = False, rug = True)

g = s.FacetGrid (df,col = 'diagnosis', hue = 'diagnosis')
g.map (s.distplot, "perimeter_mean", hist = False, rug = True)

g = s.FacetGrid (df,col = 'diagnosis', hue = 'diagnosis')
g.map (s.distplot, "area_mean", hist = False, rug = True)

g = s.FacetGrid (df,col = 'diagnosis', hue = 'diagnosis')
g.map (s.distplot, "smoothness_mean", hist = False, rug = True)

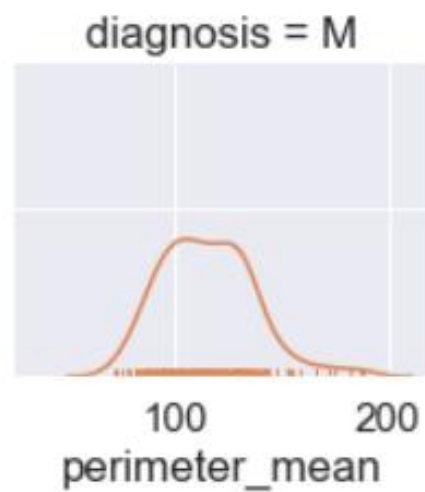
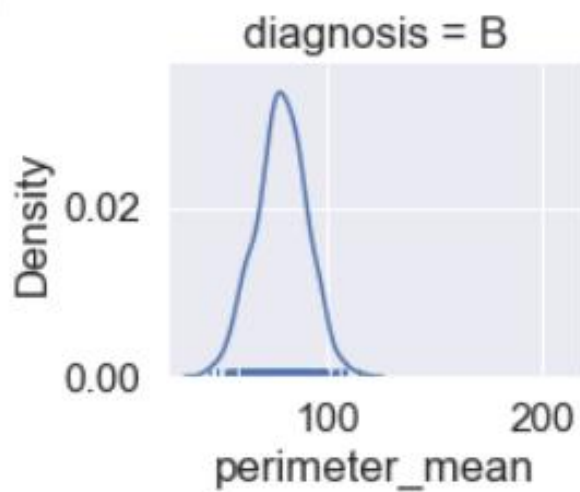
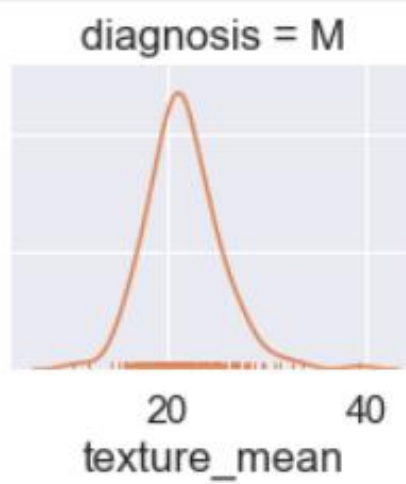
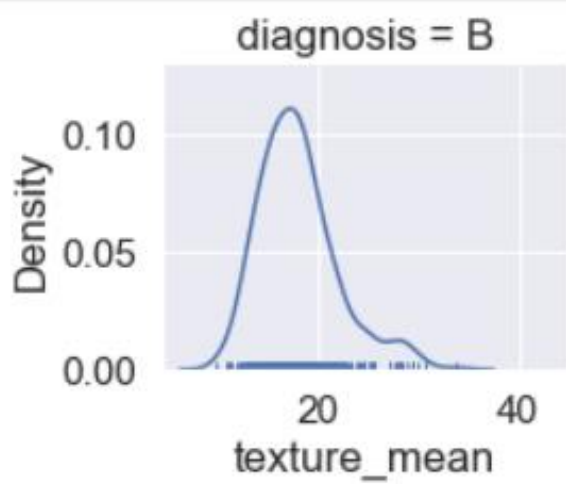
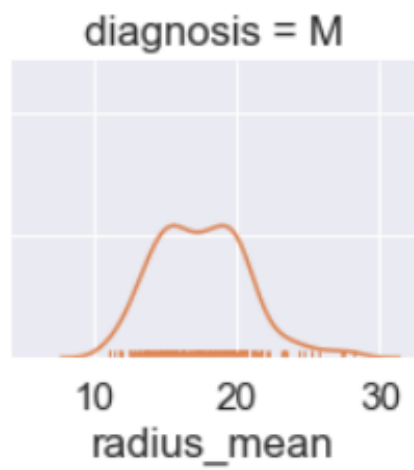
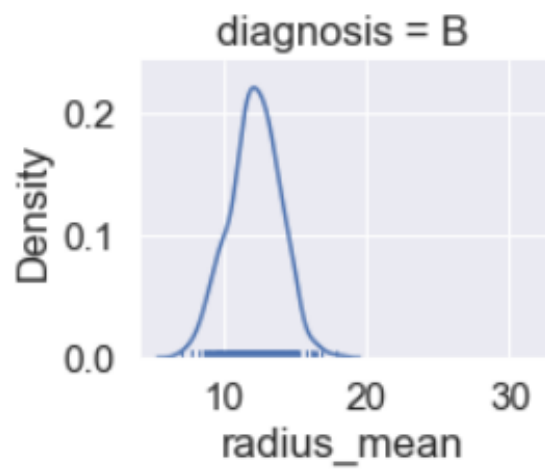
g = s.FacetGrid (df,col = 'diagnosis', hue = 'diagnosis')
g.map (s.distplot, "compactness_mean", hist = False, rug = True)

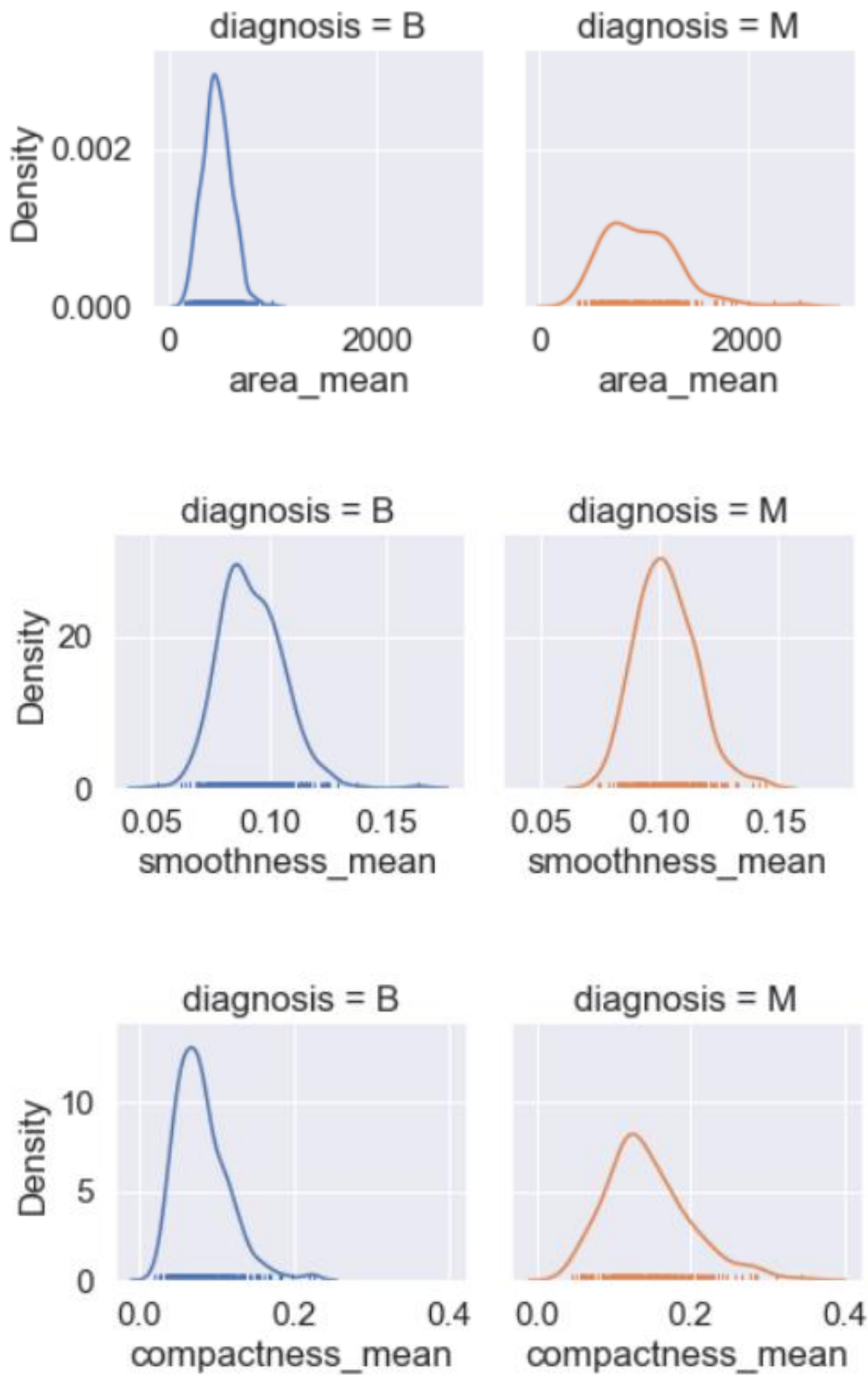
g = s.FacetGrid (df,col = 'diagnosis', hue = 'diagnosis')
g.map (s.distplot, "concavity_mean", hist = False, rug = True)

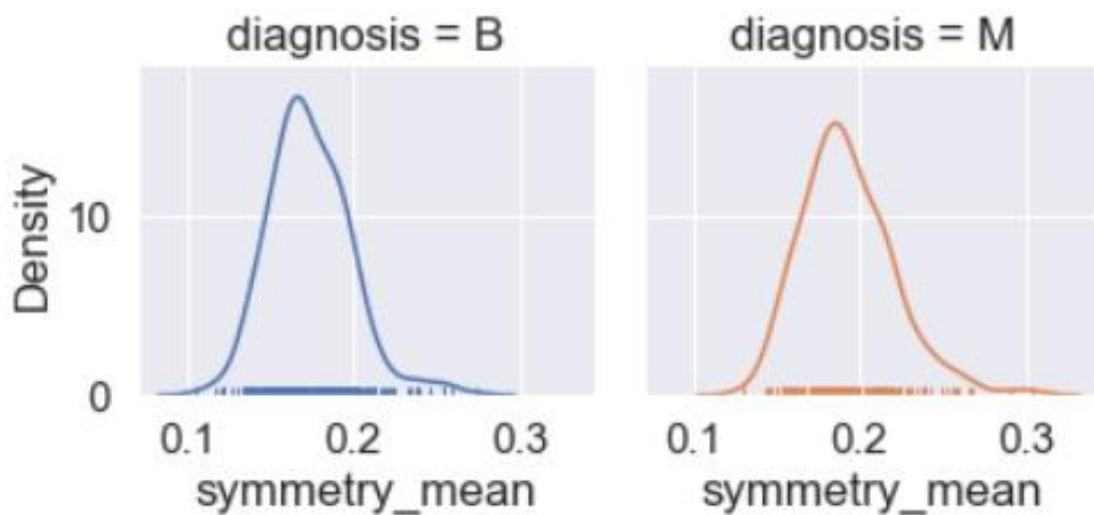
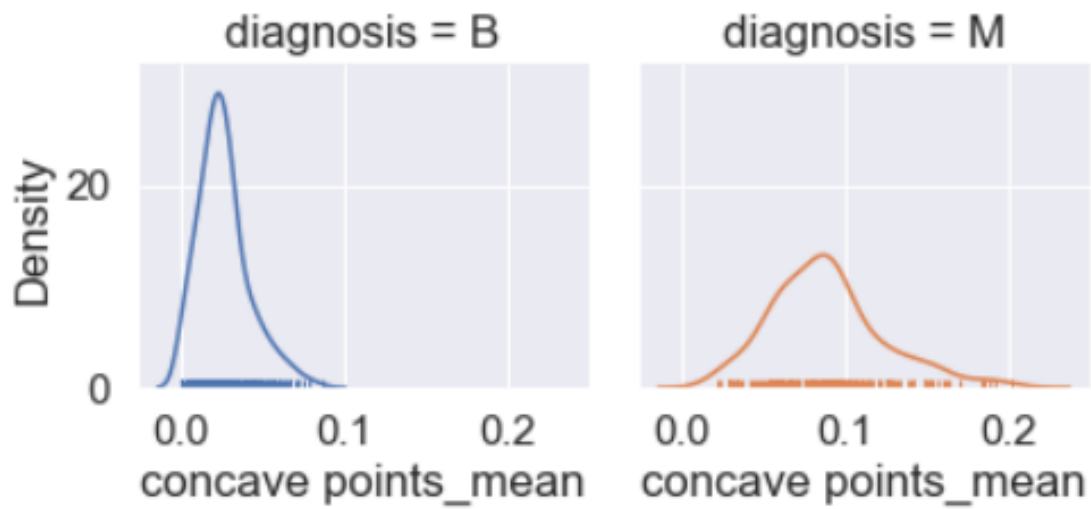
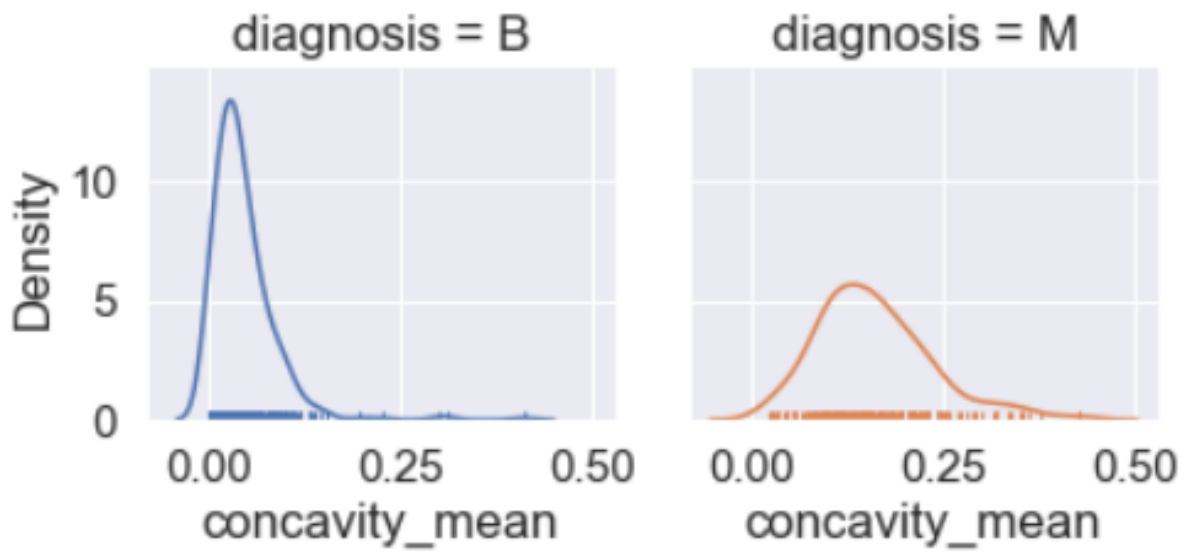
g = s.FacetGrid (df,col = 'diagnosis', hue = 'diagnosis')
g.map (s.distplot, "concave points_mean", hist = False, rug = True)

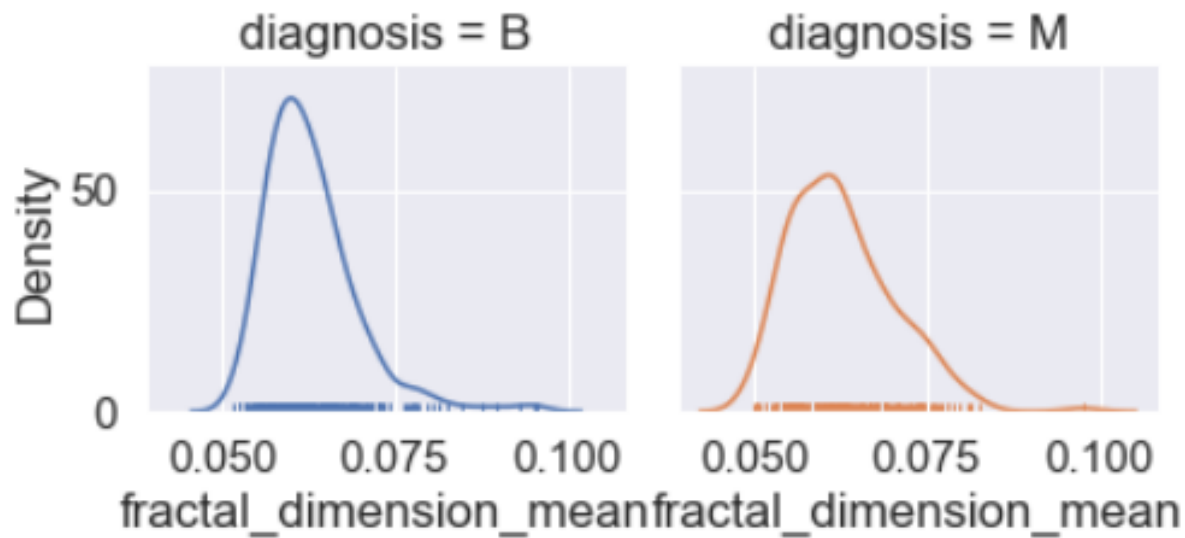
g = s.FacetGrid (df,col = 'diagnosis', hue = 'diagnosis')
g.map (s.distplot, "symmetry_mean", hist = False, rug = True)

g = s.FacetGrid (df,col = 'diagnosis', hue = 'diagnosis')
g.map (s.distplot, "fractal_dimension_mean", hist = False, rug = True)
```









8.10 Creating function for fitting model:

```
In [22]: def FitModel (X,Y, algo_name , algorithm, gridSearchParams, cv):
    np.random.seed(10)
    x_train, x_test, y_train, y_test = train_test_split (X,Y,test_size = 0.2)

    # Find the Parameters , then choose best parameters

    grid = GridSearchCV(estimator = algorithm, param_grid = gridSearchParams,
                        cv = cv, scoring = 'accuracy', verbose = 1 , n_jobs = -1 )

    grid_result = grid.fit(x_train, y_train)
    best_params = grid_result.best_params_
    pred = grid_result.predict (x_test)
    cm = confusion_matrix (y_test,pred)

    print (pred)
    pickle.dump(grid_result,open(algo_name,'wb'))

    print ('Best Params :', best_params)
    print ('Classification Report:',classification_report(y_test,pred))
    print ('Accuracy Score', (accuracy_score(y_test,pred)))
    print ('Confusion Matrix :\n',cm)
```

8.11 SVC Model:

```
In [30]: param = {
    'C': [0.1,1,100,1000],
    'gamma':[0.0001,0.001, 0.005, 0.1,1, 3,5,10, 100]
}

FitModel (x_norm,y_norm,'SVC',SVC(), param, cv =5)

Fitting 5 folds for each of 36 candidates, totalling 180 fits
[1 0 0 1 0 1 0 0 1 1 0 0 1 0 1 1 0 0 0 0 1 1 0 0 1 1 0 1 0 0 0 0 1 0 0 0 0
 0 0 0 0 0 0 1 0 1 1 1 1 1 0 0 1 0 0 1 1 0 0 0 1 0 1 1 0 0 0 0 1 1 0 0 0 0
 0 1 0 0 1 1 0 0 1 0 0 0 0 1 0 0 0 0 1 0 0 0 1 0 1 0 1 1 0 0 0 0 1 1 0 0 0
 0 1 1]
Best Params : {'C': 1, 'gamma': 1}
Classification Report:

```

		precision	recall	f1-score	support
	0	1.00	0.96	0.98	75
	1	0.93	1.00	0.96	39
accuracy				0.97	114
macro avg		0.96	0.98	0.97	114
weighted avg		0.98	0.97	0.97	114

```

Accuracy Score 0.9736842105263158
Confusion Matrix :
[[72  3]
 [ 0 39]]

```

8.12 Random Forest Model:

```
In [31]: param = { 'n_estimators': [100,500,1000,2000] }
FitModel (x_norm,y_norm,'Random Forest',RandomForestClassifier(), param, cv =10)

Fitting 10 folds for each of 4 candidates, totalling 40 fits
[1 0 0 1 0 0 0 0 1 1 0 0 1 0 1 1 0 0 0 1 1 1 0 0 1 1 0 1 1 0 0 0 0 1 0 0 0 0
 0 0 0 0 0 0 1 0 1 1 1 1 0 0 0 1 0 0 1 1 0 0 0 1 0 1 1 0 0 0 0 1 1 0 0 0 0
 0 1 0 0 1 1 0 0 1 0 0 0 0 1 0 0 0 0 1 0 0 0 1 0 1 0 1 1 0 0 0 0 1 1 0 0 0
 0 1 1]
Best Params : {'n_estimators': 100}
Classification Report:

```

			precision	recall	f1-score	support
	0	1.00	0.97	0.99	75	
	1	0.95	1.00	0.97	39	
accuracy				0.98	114	
macro avg		0.98	0.99	0.98	114	
weighted avg		0.98	0.98	0.98	114	

```

Accuracy Score 0.9824561403508771
Confusion Matrix :
[[73  2]
 [ 0 39]]
```

Chapter 9: Testing of Model

9.1 Loading pickle file for execution

```
In [2]: import pickle
import os
import pandas as pd

load_model = pickle.load(open("SVC", "rb"))

os.chdir('C:\\Users\\Wilesh\\Documents\\Cancer Detection Project')
data = pd.read_csv('Reports.csv')

x = data.drop(labels='Actual diagnosis', axis=1)
x_norm = (x - x.mean()) / (x.max() - x.min())

load_model = pickle.load(open("SVC", "rb"))

pred = load_model.predict(x_norm)
pred = pred.astype(str)

pred[pred == "1"] = "Malignant"
pred[pred == "0"] = "Benign"

data['Predictions'] = pred
data.to_csv('Result.csv')
```

9.2 Input

		id	radius_m	texture_m	perimeter_area_m	smoothness	compactness	concavity	concave	symmetry	fractal_d	radius_s	texture_s	perimeter_area_s	smoothness	compactness	concavity	concave	symmetry	fractal_d	radius_w	texture_w	perimeter_area_w	smoothness	compactness	concavity	concave	symmetry	fractal_d	Actual diagnosis		
1	859575	18.94	21.31	123.6	1130	0.09009	0.1029	0.108	0.07951	0.1582	0.05461	0.7888	0.7975	5.486	96.05	0.00444	0.01652	0.02269	0.0137	0.01386	0.0017	24.86	26.58	165.9	1866	0.1193	0.2336	0.2687	0.1789	0.2551	0.06589	M
2	8611161	13.34	15.86	86.49	520	0.1078	0.1535	0.1169	0.06987	0.1942	0.06902	0.286	1.016	1.535	12.96	0.00679	0.03575	0.0398	0.01383	0.02134	0.0046	15.53	23.19	96.66	614.9	0.1536	0.4791	0.4858	0.1708	0.3527	0.1016	B
3	859717	17.2	24.52	114.2	929.4	0.1071	0.1183	0.1692	0.07944	0.1927	0.06487	0.5907	1.041	3.705	69.47	0.00582	0.0616	0.04252	0.01127	0.01527	0.0063	23.32	33.82	151.6	1681	0.1585	0.7394	0.5566	0.1899	0.3313	0.1339	M
4	8611163	17.93	24.48	115.2	998.9	0.08855	0.07027	0.05699	0.04744	0.1538	0.0551	0.4212	1.433	2.765	45.81	0.00544	0.01169	0.01622	0.00852	0.01419	0.00275	20.92	34.69	135.1	1320	0.1315	0.1806	0.208	0.1136	0.2504	0.07948	M
5	913102	14.64	16.85	94.21	666	0.08641	0.06698	0.05192	0.02791	0.1409	0.05355	0.2204	1.006	1.471	19.98	0.00354	0.01393	0.018	0.00614	0.01254	0.00122	16.46	25.44	106	831	0.1142	0.207	0.2437	0.07828	0.2455	0.06596	B
6	8610908	12.86	18	83.19	506.3	0.09934	0.09546	0.03889	0.0315	0.1718	0.05997	0.2655	1.095	1.778	20.35	0.00529	0.01661	0.02071	0.00818	0.01748	0.00285	14.24	24.82	91.88	622.1	0.1289	0.2141	0.1731	0.07918	0.2779	0.07918	B
7	912193	12.16	18.03	78.29	453.3	0.09087	0.07838	0.02916	0.01527	0.1464	0.06284	0.2194	1.19	1.678	16.26	0.00491	0.01666	0.01397	0.00516	0.01454	0.00186	13.34	27.87	88.83	547.4	0.1208	0.2379	0.162	0.0569	0.2406	0.07729	B
8	8610637	18.05	16.15	120.2	1006	0.1065	0.2146	0.1684	0.108	0.2152	0.06673	0.9806	0.5505	6.311	134.8	0.00794	0.05839	0.04658	0.0207	0.02591	0.00705	22.39	18.91	150.1	1610	0.1478	0.5634	0.3786	0.2102	0.3751	0.1108	M
9	8610862	20.18	23.97	143.7	1245	0.1286	0.3454	0.3754	0.1604	0.2906	0.08142	0.9317	1.885	8.649	116.4	0.01038	0.06835	0.1091	0.02593	0.07895	0.00599	23.37	31.72	170.3	1623	0.1639	0.6164	0.7681	0.2508	0.544	0.09964	M
10	91227	13.9	19.24	88.73	602.9	0.07991	0.05326	0.02995	0.0207	0.1579	0.05594	0.3316	0.9264	2.056	28.41	0.0037	0.01082	0.0153	0.00628	0.01062	0.00222	16.41	26.42	104.4	830.5	0.1064	0.1415	0.1673	0.0815	0.2356	0.07603	B
11	912519	13.47	14.06	87.32	546.3	0.1071	0.1155	0.05786	0.05266	0.1779	0.06639	0.1588	0.5733	1.102	12.84	0.00445	0.01452	0.01334	0.00879	0.01698	0.00279	14.83	18.32	94.94	660.2	0.1393	0.2499	0.1848	0.1335	0.3227	0.09326	B
12	927241	20.6	29.33	140.1	1265	0.1178	0.277	0.3514	0.152	0.2397	0.07016	0.726	1.595	5.772	86.22	0.00652	0.06158	0.07117	0.01664	0.02324	0.00619	25.74	39.42	184.6	1821	0.165	0.8681	0.9387	0.265	0.4087	0.124	M
13	912558	13.7	17.64	87.76	571.1	0.0995	0.07957	0.04548	0.0316	0.1732	0.06088	0.2431	0.9462	1.564	20.64	0.00325	0.00819	0.01698	0.00923	0.01285	0.00152	14.96	23.53	95.78	686.5	0.1199	0.1346	0.1742	0.09077	0.2518	0.0696	B
14	8610404	16.07	19.65	104.1	817.7	0.09168	0.08424	0.09769	0.06638	0.1798	0.05391	0.7474	1.016	5.029	79.25	0.01082	0.02203	0.035	0.01809	0.0155	0.00195	19.77	24.56	128.8	1223	0.15	0.2045	0.2829	0.152	0.265	0.06387	M
15	884180	19.4	23.5	129.1	1155	0.1027	0.1558	0.2049	0.08886	0.1978	0.06	0.5243	1.802	4.037	60.41	0.01061	0.03252	0.03915	0.01559	0.02186	0.00395	21.65	30.53	144.9	1417	0.1463	0.2968	0.3458	0.1564	0.292	0.07614	M

9.3 Output

		radius_m	texture_m	perimeter_area_m	smoothness_m	compactness_m	concavity_m	concave_m	symmetry_m	fractal_d_m	radius_s	texture_s	perimeter_area_s	smoothness_s	compactness_s	concavity_s	concave_s	symmetry_s	fractal_d_s	radius_w	texture_w	perimeter_area_w	smoothness_w	compactness_w	concavity_w	concave_w	symmetry_w	fractal_d_w	Actual diagnosis	Predictions		
1	859575	18.94	21.31	123.6	1130	0.09009	0.1029	0.108	0.07951	0.1582	0.05461	0.7888	0.7975	5.486	96.05	0.00444	0.01652	0.02269	0.0137	0.01386	0.001698	24.86	26.58	165.9	1866	0.1193	0.2336	0.2687	0.1789	0.2551	0.06589	Malignant
2	8611161	13.34	15.86	86.49	520	0.1078	0.1535	0.1169	0.06987	0.1942	0.06902	0.286	1.016	1.535	12.96	0.00679	0.03575	0.0398	0.01383	0.02134	0.004603	15.53	23.19	96.66	614.9	0.1536	0.4791	0.4858	0.1708	0.3527	0.1016	Benign
3	859717	17.2	24.52	114.2	929.4	0.1071	0.1183	0.1692	0.07944	0.1927	0.06487	0.5907	1.041	3.705	69.47	0.00582	0.0616	0.04252	0.01127	0.01527	0.006299	23.32	33.82	151.6	1681	0.1585	0.7394	0.5566	0.1899	0.3313	0.1339	Malignant
4	8611163	17.93	24.48	115.2	998.9	0.08855	0.07027	0.05699	0.04744	0.1538	0.0551	0.4212	1.433	2.765	45.81	0.00544	0.01169	0.01622	0.00852	0.01419	0.002751	20.92	34.69	135.1	1320	0.1315	0.1806	0.208	0.1136	0.2504	0.07948	Malignant
5	913102	14.64	16.85	94.21	666	0.08641	0.06698	0.05192	0.02791	0.1409	0.05355	0.2204	1.006	1.471	19.98	0.00353	0.01393	0.018	0.00614	0.01254	0.001219	16.46	25.44	106	831	0.1142	0.207	0.2437	0.07828	0.2455	0.06596	Benign
6	8610908	12.86	18	83.19	506.3	0.09934	0.09546	0.03889	0.0315	0.1718	0.05997	0.2655	1.095	1.778	20.35	0.00529	0.01661	0.02071	0.008179	0.01748	0.002848	14.24	24.82	91.88	622.1	0.1289	0.2141	0.1731	0.07926	0.2779	0.07918	Benign
7	912193	12.16	18.03	78.29	453.3	0.09087	0.07838	0.02916	0.01527	0.1464	0.06284	0.2194	1.19	1.678	16.26	0.00491	0.01666	0.01397	0.005161	0.01454	0.001858	13.34	27.87	88.83	547.4	0.1208	0.2379	0.162	0.0569	0.2406	0.07729	Benign
8	8610637	18.05	16.15	120.2	1006	0.1065	0.2146	0.1684	0.108	0.2152	0.06673	0.9806	0.5505	6.311	134.8	0.00794	0.05839	0.04658	0.0207	0.02591	0.007054	22.39	18.91	150.1	1610	0.1478	0.5634	0.3786	0.2102	0.3751	0.1108	Malignant
9	8610862	20.18	23.97	143.7	1245	0.1286	0.3454	0.3754	0.1604	0.2906	0.08142	0.9317	1.885	8.649	116.4	0.01038	0.06835	0.1091	0.02593	0.07895	0.005987	23.37	31.72	170.3	1623	0.1639	0.6164	0.7681	0.2508	0.544	0.09964	Malignant
10	91227	13.9	19.24	88.73	602.9	0.07991	0.05326	0.02995	0.0207	0.1579	0.05594	0.3316	0.9264	2.056	28.41	0.003704	0.01082	0.0153	0.006275	0.01062	0.002217	16.41	26.42	104.4	830.5	0.1064	0.1415	0.1673	0.0815	0.2356	0.07603	Benign
11	912519	13.47	14.06	87.32	546.3	0.1071	0.1155	0.05786	0.05266	0.1779	0.06639	0.1588	0.5733	1.102	12.84	0.00445	0.01452	0.01334	0.008791	0.01698	0.002787	14.83	18.32	94.94	660.2	0.1393	0.2499	0.1848	0.1335	0.3227	0.09326	Benign
12	927241	20.6	29.33	140.1	1265	0.1178	0.277	0.3514	0.152	0.2397	0.07016	0.726	1.595	5.772	86.22	0.00652	0.06158	0.07117	0.01664	0.02324	0.006185	25.74	39.42	184.6	1821	0.165	0.8681	0.9387	0.265	0.4087	0.124	Malignant
13	912558	13.7	17.64	87.76	571.1	0.0995	0.07957	0.04548	0.0316	0.1732	0.06088	0.2431	0.9462	1.564	20.64	0.003245	0.008186	0.01698	0.009233	0.01285	0.001524	14.96	23.53	95.78	686.5	0.1199	0.1346	0.1742	0.09077	0.2518	0.0696	Benign
14	8610404	16.07	19.65	104.1	817.7	0.09168	0.08424	0.09769	0.06638	0.1798	0.05391	0.7474	1.016	5.029	79.25	0.01082	0.02203	0.035	0.01809	0.0155	0.001948	19.77	24.56	128.8	1223	0.15	0.2045	0.2829	0.152	0.265	0.06387	Malignant
15	884180	19.4	23.5	129.1	1155	0.1027	0.1558	0.2049	0.08886	0.1978	0.06	0.5243	1.802	4.037	60.41	0.01061	0.03252	0.03915	0.01559	0.02186	0.003949	21.65	30.53	144.9	1417	0.1463	0.2968	0.3458	0.1564	0.292	0.07614	Malignant

Chapter 10: Conclusion

In this project, we talked about machine learning (ML) ideas and described how they may be used to predict cancer. It is clear from a study of their findings that combining multidimensional heterogeneous data with the use of various feature selection and classification algorithms might result in useful inference tools for the cancer domain. The proposed machine-learning approaches 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.

Acknowledgement

We have great pleasure in presenting the mini project report on “**Cancer Detection**”. We take this opportunity to express our sincere thanks towards our guide **Prof.Uma Ade** Department of Computer Engineering, Watumull thane for providing the technical guidelines and suggestions regarding line of work. We would like to express our gratitude towards his constant encouragement, support and guidance through the development of project.

We thank **Prof. Dhanjay Raut** Head of Department, Computer Engineering, Watumull for his encouragement during progress meeting and providing guidelines to write this report.

We also thank the entire staff of Watumull for their invaluable help rendered during the course of this work. We wish to express our deep gratitude towards all our colleagues of Watumull for their encouragement.

Nilesh Parab (D221452)
Ananya Yadav (D221427)
Riya Singh (D221421)
Vivek Tiwari (D221422)

Bibliography

- [https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+\(diagnostic\)](https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+(diagnostic))
- <https://www.cancer.net/navigating-cancer-care/diagnosing-cancer/stages-cancer>
- <https://www.javatpoint.com/machine-learning-random-forest-algorithm>
- <https://www.cancer.org/treatment/understanding-your-diagnosis/tests/ct-scan-for-cancer.html>