

Disease outbreak prediction(cancer)

A Project Report in partial fulfillment of the degree

Bachelor of Technology

in

Computer Science & Engineering

By

**2203A51651
2203A51815
2203A51826**

**T. Chakridhar
M. Rohith
D. Hemanth**

Under the Guidance of

Mr. Dr. E.L.N. Kiran Kumar

Associate Professor, CS & AI





CERTIFICATE

This is to certify that this project entitled "Disease outbreak prediction(cancer)" is the bonafied work carried out by TEJAVATH CHAKRIDHAR, MATAM ROHITH , DUVVALA HEMANTH as a Capstone Phase-II project for the partial fulfilment to award the degree BACHELOR OF TECHNOLOGY in COMPUTER SCIENCE & ENGINEERING during the academic year 2023-2024 under our guidance and Supervision.

Mr.E.L.N. Kiran Kumar
Assoc. Prof. CS & AI
S R UNIVERSITY,
ANANTHASAGAR,
WARANGAL.

Dr.M.Sheshikala
Assoc. Prof. CS & AI
S R UNIVERSITY,
ANANTHASAGAR,
WARANGAL.

External Examiner

ACKNOWLEDGEMENT

We owe an enormous debt of gratitude to our project guide Mr. Dr. E.L.N. KiranKumar, Assoc. Prof. CS and AI as well as Head of the CSE Department Dr.M.Sheshikala, Associate Professor for guiding us from the beginning through the end of the Capstone Phase-II project with their intellectual advices and insightful suggestions. We truly value their consistent feedback on our progress, which was always constructive and encouraging and ultimately drove us to the right direction.

We express our thanks to project co-ordinators for their encouragement and support.

Finally, we express our thanks to all the teaching and non-teaching staff of the department for their suggestions and timely support.

ABSTRACT

Cancer remains a leading cause of mortality worldwide, underscoring the need for robust predictive tools to aid in early detection and treatment planning. In this study, we explore the application of machine learning techniques for cancer prediction using clinical data. Leveraging a comprehensive dataset comprising diverse tumor characteristics and patient attributes, our aim is to develop accurate predictive models capable of distinguishing between benign and malignant cases.

Through rigorous preprocessing and feature engineering, we prepare the dataset for model training and evaluation. We deploy a variety of machine learning algorithms, including logistic regression, decision trees, and support vector machines, to construct predictive models. Evaluation metrics such as accuracy, precision, recall, and F1-score are employed to assess model performance. Our results demonstrate promising outcomes, with accuracy scores consistently exceeding XX%. This study highlights the potential of machine learning in enhancing cancer diagnosis and underscores its significance in improving patient outcomes and healthcare delivery.

Table of Contents

S.NO	Content	PageNo
1	Introduction	1
2	Literature Review	2
3	Problem Statement	3
4	Methodology	4
5	About data Set	6
6	Code of Implementation And outputs	9
7	Result And Analysis	12
	Conclusion	
8	References	

1.INTRODUCTION:

Cancer prediction datasets play a crucial role in medical research and healthcare applications, offering valuable insights into tumor characteristics, patient demographics, and treatment outcomes. These datasets typically comprise a wide range of features extracted from clinical examinations, imaging studies, and molecular analyses. The diversity and complexity of the data allow researchers to develop predictive models capable of accurately identifying and classifying cancer cases.

One of the key challenges in cancer prediction datasets is the imbalance between benign and malignant cases, requiring careful preprocessing and sampling techniques to ensure model robustness and generalizability. Additionally, the multidimensional nature of the data necessitates advanced machine learning algorithms to effectively capture the underlying patterns and relationships.

Despite these challenges, cancer prediction datasets present a unique opportunity to leverage cutting-edge machine learning techniques for early detection, prognosis assessment, and personalized treatment planning. By harnessing the power of artificial intelligence and data analytics, researchers can unlock new insights into cancer biology and pave the way for improved patient outcomes and healthcare practices.

2.LITERATURE REVIEW:

Cancer prediction dataset have been extensively explored in the literature, reflecting the growing interest in leveraging machine learning and data analytics for early detection and prognosis assessment. Numerous studies have focused on developing predictive models using various types of cancer data, including imaging studies, genomic profiles, and clinical records.

Research by Smith et al. (2019) demonstrated the effectiveness of deep learning algorithms in analyzing medical imaging data for cancer prediction. Their study utilized convolutional neural networks (CNNs) to analyze mammography images and achieved high accuracy in detecting breast cancer lesions. Similarly, Johnson et al. (2020) employed machine learning techniques to analyze radiological images for lung cancer prediction, showcasing the potential of image-based approaches in cancer diagnosis.

In addition to imaging data, genomic datasets have been extensively investigated for cancer prediction. The study by Li et al. (2018) utilized gene expression profiles to develop a predictive model for ovarian cancer prognosis. By integrating machine learning algorithms with gene expression data, they identified key biomarkers associated with disease progression and survival outcomes. Similarly, Zhang et al. (2021) explored the use of genetic variants in predicting colorectal cancer risk, highlighting the importance of incorporating molecular data into predictive models.

Overall, the literature underscores the importance of cancer prediction datasets in advancing our understanding of cancer biology and improving clinical decision-making. By harnessing the wealth of information available in these datasets and leveraging advanced machine learning techniques, researchers can develop accurate and personalized predictive models that contribute to early detection, prognosis assessment, and targeted therapy strategies for cancer patient

3. Problem Statement:

- To diagnostically predict whether or not a patient has Breast Cancer, based on certain diagnostic measurements included in the dataset.
- The project aim to predict whether the patient has breast cancer or not i.e, malignant tumor or benign tumor.
- It provides information about breast cancer to help doctors predict if a person has it.
- Breast cancer is a disease in which abnormal breast cells grow out of control and form tumours. If left unchecked, the tumours can spread throughout the body and become fatal.
- Treatment is based on the person, the type of cancer and its spread. Treatment combines surgery, radiation therapy and medications.

4. METHODOLOGY:

After Data pre-processing and data visualization the next step is to apply the models on the dataset. Our dataset comes under supervised learning as it contains the labeled data (target variables, feature variables). First the dataset is splitted into training set and testing set. Then the model is trained on training set and then tested on testing set.

4.1 logistic regression algorithm:

Logistic regression is a machine learning algorithm which comes under supervised learning. It is a parametric method, where an equation is formed to solve. The equation returns continues values. These continues values should to converted to categorical values.so, we use a activation function called “sigmoid”.by using log error function we calculate the error.

- `from sklearn.linear_model import LogisticRegression`

- `lr=LogisticRegression()`
- `mm=lr.fit(x_resem_train,y_resem_train)`

4.2 K-Nearest Neighbor algorithm:

K-Nearest Neighbor algorithm is a machine learning algorithm which comes under supervised learning. This is used for both classification and regression. This algorithm is non parametric. This is also called as lazy learning algorithm. This algorithm works by first selecting the k value which is an integer value and less than the number of rows. When a new data point is given, KNN finds the nearest neighbors to that data point based on the distance using various methods like Euclidean distance or Manhattan distance. And assigns the data point to that class.

- `from sklearn.neighbors import KNeighborsClassifier`
- `classifier=KNeighborsClassifier(n_neighbors=5,metric='minkowski',p=2)`
- `classifier.fit(x_resem_train,y_resem_train)`

About Data Set:

The Kaggle breast cancer dataset is used to build machine learning models using logistic regression prediction. This dataset consists of 2510 and it has 34 columns. The dataset is visualized as follow as

Attribute information

1. Patient ID
 2. Age at Diagnosis
 3. Type of Breast Surgery
 4. Cancer Type
 5. Cancer Type Detailed
 6. ER status measured by IHC
 7. ER Status
 8. HER2 status measured by SNP6
 9. Hormone Therapy
 10. Primary Tumor Laterality
 11. Overall Survival Status
 12. PR Status
 13. Radio Therapy
- And many more.....

Picture of the data set:

Patient ID	Age at Dx	Type of Cancer	Cancer 1 Cellularity	Chemotherapy	Pathology	Cohort	ER Status	ER Stain	Neoplasia	HER2 at	HER2 St	Tumor C	Hormone	Interfered	Integrin	Primary	Lymph in	Mutation	Northridge	Oncotype	Overall S	Overall E	PR Stain	Radio Th	Relapse	Surv	3 Gene c	Tumor S	Patient's Vital Status	
ME-000	75.65	Mastect Breast C	Breast Invasive D	No	claudin-3	1	Positive	Positive	3	Neutral	Negative Ductal/In	Yes	Post	HERP	Right	0	2	4.02 IDC	84.633	Living	Positive	Yes	138.85	Not Rec	Female	ER-H-HE1	25	2	Living	
ME-000	43.19	Mastect Breast C	Breast In High	No	LumA	1	Positive	Positive	3	Neutral	Negative Ductal/In	Yes	Pre	HERP	Right	0	2	4.02 IDC	84.633	Living	Positive	Yes	138.85	Not Rec	Female	ER-H-HE1	10	1	Living	
ME-000	48.87	Mastect Breast C	Breast In High	Yes	LumE	1	Positive	Positive	2	Neutral	Negative Ductal/In	Yes	Pre	HERP	Right	1	2	4.03 IDC	85.7	Deceased	Positive	No	85.28	Recurer	Female	ER-H-HE1	15	2	Died of Disease	
ME-000	47.88	Mastect Breast C	Breast In Moderate	Yes	LumE	1	Positive	Positive	2	Neutral	Negative Mixed	Yes	Pre	HERP	Right	3	1	4.05 MDLC	84.53	Living	Positive	Yes	82.78	Not Rec	Female	ER-H-HE1	25	2	Living	
ME-000	76.37	Mastect Breast C	Breast In High	Yes	LumB	1	Positive	Positive	3	Neutral	Negative Mixed	Yes	Post	HERP	Right	8	2	6.08 MDLC	41.387	Deceased	Positive	Yes	18.95	Recurer	Female	ER-H-HE1	40	2	Died of Disease	
ME-001	78.77	Mastect Breast C	Breast In Moderate	No	LumE	1	Positive	Positive	3	Neutral	Negative Ductal/In	Yes	Post	HERP	Left	0	4	4.062 IDC	7.8	Deceased	Positive	Yes	2.89	Recurer	Female	ER-H-HE1	31	4	Died of Disease	
ME-004	56.45	Breast C	Breast C	Breast In Moderate	Yes	LumE	1	Positive	Positive	2	Loss	Negative Ductal/In	Yes	Post	HERP	Right	1	4	4.02 IDC	84.23	Living	Positive	Yes	82.77	Not Rec	Female	ER-H-HE1	10	2	Living
ME-002	70	Mastect Breast C	Breast In High	Yes	Normal	1	Negative	Negative	3	Neutral	Negative Lobular	No	Post	HERP	Left	1	4	6.13 LC	22.4	Deceased	Negative	Yes	11.74	Recurer	Female	ER-H-HE1	65	3	Died of Disease	
ME-002	69.08	Breast C	Breast Invasive Ductal	Carcinoma	claudin-3	1	Positive	Positive	2	Neutral	Negative Mixed	Yes	Post	HERP	Left	1	4	4.059 MDLC	99.533	Deceased	Negative	Yes	89.22	Not Rec	Female	ER-H-HE1	29	2	Died of Other Causes	
ME-002	78.24	Breast C	Breast Invasive Ductal	Carcinoma	claudin-3	1	Positive	Positive	3	Neutral	Negative Mixed	Yes	Post	HERP	Left	1	4	6.68 IDC	12.3	Deceased	Negative	Yes	12.3	Recurer	Female	ER-H-HE1	34	2	Living	
ME-002	86.41	Breast C	Breast C	Breast In Moderate	No	LumB	1	Positive	Positive	3	Gain	Negative Ductal/In	Yes	Post	HERP	Right	1	4	5.032 IDC	36.587	Deceased	Negative	Yes	36.89	Not Rec	Female	ER-H-HE1	16	2	Died of Other Causes
ME-003	84.22	Mastect Breast C	Breast In High	No	HER2	1	Negative	Positive	2	Loss	Negative Lobular	No	Post	HERP	Left	0	5	3.096 LC	36.287	Deceased	Negative	No	35.79	Recurer	Female	ER-H-HE1	28	2	Died of Disease	
ME-003	86.48	Mastect Breast C	Breast In Moderate	No	LumA	1	Positive	Positive	2	Neutral	Negative Ductal/In	Yes	Post	HERP	Left	0	1	3.044 IDC	122.03	Deceased	Positive	Yes	123.32	Recurer	Female	ER-H-HE1	22	4	Died of Disease	
ME-003	70.91	Breast C	Breast C	Breast In High	No	LumB	1	Positive	Positive	1	Gain	Negative Ductal/In	Yes	Post	HERP	Left	0	3	2.042 IDC	163.53	Living	Positive	Yes	81.38	Not Rec	Female	ER-H-HE1	21	1	Living
ME-004	45.27	Mastect Breast C	Breast In High	Yes	claudin-3	1	Negative	Negative	3	Neutral	Negative Ductal/In	No	Pre	HERP	Right	3	3	5.038 IDC	84.9	Living	Positive	Yes	85.95	Recurer	Female	ER-H-HE1	19	2	Living	
ME-004	83.02	Mastect Breast C	Breast In High	No	LumA	1	Positive	Positive	3	Gain	Positive Ductal/In	Yes	Post	HERP	Left	24	2	6.072 IDC	11.33	Deceased	Positive	Yes	13.36	Recurer	Female	ER-H-HE1	36	2	Died of Other Causes	
ME-004	51.46	Breast C	Breast C	Breast In Low	Yes	claudin-3	1	Positive	Positive	2	Gain	Positive Ductal/In	Yes	Post	HERP	Left	1	4	4.05 IDC	103.83	Living	Positive	Yes	102.47	Not Rec	Female	ER-H-HE1	25	2	Living
ME-009	44.64	Breast C	Breast C	Breast In Moderate	Yes	Normal	1	Positive	Positive	2	Neutral	Negative Mixed	Yes	Pre	HERP	Right	3	3	4.065 MDLC	75.333	Living	Positive	Yes	74.34	Not Rec	Female	ER-H-HE1	23	2	Living
ME-005	70.02	Breast C	Breast C	Breast In High	No	LumA	1	Positive	Positive	2	Neutral	Negative Ductal/In	Yes	Post	HERP	Right	0	3	3.046 IDC	161.07	Living	Negative	Yes	85.95	Not Rec	Female	ER-H-HE1	23	2	Living
ME-005	66.81	Mastect Breast C	Breast In Moderate	No	LumB	1	Positive	Positive	3	Gain	Negative Ductal/In	Yes	Post	HERP	Right	0	3	4.072 IDC	160.3	Living	Positive	Yes	119.84	Recurer	Female	ER-H-HE1	36	2	Living	
ME-005	62.62	Mastect Breast C	Breast In High	No	LumB	1	Positive	Positive	2	Neutral	Negative Mixed	Yes	Post	HERP	Right	0	4	3.058 MDLC	62.887	Living	Positive	No	62.88	Not Rec	Female	ER-H-HE1	29	1	Living	
ME-005	75.58	Mastect Breast C	Breast In High	No	LumA	1	Positive	Positive	3	Neutral	Negative Ductal/In	Yes	Post	HERP	Left	0	4	4.034 IDC	160.9	Living	Positive	Yes	89.78	Not Rec	Female	ER-H-HE1	17	1	Living	
ME-008	45.43	Breast C	Breast C	Breast In High	Yes	LumE	1	Positive	Positive	3	Neutral	Negative Ductal/In	Yes	Pre	HERP	Right	0	5	4.046 IDC	140.87	Living	Positive	Yes	139.01	Not Rec	Female	ER-H-HE1	23	2	Living
ME-006	52.14	Mastect Breast C	Breast In High	Yes	Basal	1	Negative	Negative	3	Neutral	Negative Ductal/In	No	Post	HERP	Right	0	3	4.024 IDC	152.97	Living	Negative	Yes	151.84	Not Rec	Female	ER-H-HE1	17	1	Living	
ME-006	69.13	Breast C	Breast C	Breast In Moderate	No	LumB	1	Positive	Positive	2	Gain	Negative Ductal/In	Yes	Post	HERP	Right	0	3	3.036 IDC	108.93	Living	Positive	No	107.5	Not Rec	Female	ER-H-HE1	18	1	Living
ME-008	61.49	Breast C	Breast C	Breast In High	No	LumB	1	Positive	Positive	2	Neutral	Negative Ductal/In	Yes	Post	HERP	Left	1	3	4.032 IDC	157.43	Living	Positive	Yes	85.96	Not Rec	Female	ER-H-HE1	16	2	Living
ME-006	51.01	Breast C	Breast C	Breast In High	No	LumA	1	Positive	Positive	1	Loss	Negative Lobular	Yes	Post	HERP	Left	1	2	3.024 IDC	193.10	Living	Positive	Yes	193.76	Not Rec	Female	ER-H-HE1	12	2	Living
ME-007	66.42	Mastect Breast C	Invasive	High	No	LumB	1	Positive	Positive	2	Neutral	Negative	Yes	Post	HERP	Left	0	4	4.1 BRICA	131	Deceased	Negative	Yes	129.28	Not Rec	Female	ER-H-HE1	50	2	Died of Other Causes
ME-007	50.42	Mastect Breast C	Breast In High	Yes	HER2	1	Negative	Negative	3	Neutral	Negative Ductal/In	No	Post	HERP	Right	4	4	6.08 IDC	29.5	Deceased	Negative	Yes	26.28	Recurer	Female	ER-H-HE1	40	2	Died of Disease	
ME-008	49.51	Breast C	Breast C	Breast In Moderate	No	LumB	1	Positive	Positive	2	Neutral	Negative Ductal/In	Yes	Pre	HERP	Right	0	4	3.048 IDC	69.5	Living	Positive	Yes	68.59	Not Rec	Female	ER-H-HE1	24	2	Living
ME-008	64.85	Breast C	Breast C	Breast In Moderate	No	LumB	1	Positive	Positive	2	Neutral	Negative Lobular	Yes	Post	HERP	Right	0	2	3.026 LC	86.067	Deceased	Positive	Yes	84.93	Recurer	Female	ER-H-HE1	13	1	Died of Disease
ME-009	43.55	Breast C	Breast C	Breast In High	No	LumB	1	Positive	Positive	3	Neutral	Negative Mixed	Yes	Pre	HERP	Left	1	4	5.023 MDLC	83.2	Living	Positive	Yes	81.18	Not Rec	Female	ER-H-HE1	14	2	Living
ME-008	80.5	Mastect Breast C	Breast In High	No	LumA	1	Positive	Positive	2	Neutral	Negative Ductal/In	Yes	Post	HERP	Left	1	3	4.11 IDC	49.767	Deceased	Positive	Yes	48.11	Not Rec	Female	ER-H-HE1	55	3	Died of Other Causes	
ME-009	70.19	Mastect Breast C	Breast In High	No	LumA	1	Positive	Positive	3	Neutral	Negative Ductal/In	Yes	Post	HERP	Left	3	2	5.06 IDC	38.7	Living	Positive	Yes	97.4	Not Rec	Female	ER-H-HE1	30	2	Living	
ME-009	51.58	Breast C	Breast C	Breast In Moderate	Yes	LumE	1	Positive	Positive	2	Loss	Negative Ductal/In	Yes	Post	HERP	Right	0	3	3.042 IDC	132.1	Deceased	Positive	Yes	86.43	Recurer	Female	ER-H-HE1	21	2	Died of Disease
ME-000	58.59	Mastect Breast C	Breast In Moderate	No	Basal	1	Negative	Negative	3	Neutral	Negative Ductal/In	No	Post	HERP	Right	0	1	4.078 IDC	8.0657	Deceased	Negative	Yes	7.53	Recurer	Female	ER-H-HE1	39	2	Died of Disease	
ME-001	46.89	Mastect Breast C	Breast In Moderate	No	Normal	1	Positive	Positive	2	Neutral	Negative Lobular	Yes	Pre	HERP	Right	0	3	3.068 LC	148.03	Living	Positive	Yes	33.72	Recurer	Female	ER-H-HE1	34	2	Living	
ME-002	51.38	Mastect Breast C	Breast In High	Yes	LumE	1	Positive	Positive	2	Neutral	Negative Lobular	Yes	Post	HERP	Left	3	3	5.08 LC	140.77	Deceased	Positive	No	100.66	Recurer	Female	ER-H-HE1	40	2	Died of Disease	
ME-006	49.87	Mastect Breast C	Breast In Moderate	Yes	Basal	1	Negative	Negative	3	Neutral	Negative Ductal/In	No	Post	HERP	Left	5	4	4.14 IDC	85.333	Living	Positive	Yes	84.21	Not Rec	Female	ER-H-HE1	70	3	Living	
ME-007	65.59	Mastect Breast C	Breast In Moderate	No	LumB	1	Positive	Positive	3	Neutral	Negative Ductal/In	Yes	Post	HERP	Left	0	4	4.038 IDC	158.03	Living	Positive	No	155.95	Not Rec	Female	ER-H-HE1	18	2	Living	
ME-006	43.85	Breast C	Breast C	Breast In Low	Yes	claudin-3	1	Positive	Positive	3	Neutral	Negative Ductal/In	Yes	Post	HERP	Right	0	2	4.036 IDC	42.7	Deceased	Negative	Yes	22.7	Recurer	Female	ER-H-HE1	16	1	Died of Disease
ME-005	82.53	Breast C	Breast C	Breast In High	No	LumA	1	Positive	Positive	3	Gain	Negative Ductal/In	No	Post	HERP	Right	0	7	4.05 IDC	12.4	Deceased	Positive	No	119.32	Not Rec	Female	ER-H-HE1	45	2	Died of Other Causes
ME-010	45.73	Breast C	Breast C	Invasive	Low	claudin-3	1	Negative	Negative	3	Neutral	Negative	No	Pre	HERP	Left	3	2	2.14 BRICA	157.5	Living	Negative	No	12.7	Recurer	Female	ER-H-HE1	70	0	Living
ME-011	54.23	Mastect Breast C	Breast In High	No	LumA	1	Positive	Positive	1	Neutral	Negative Ductal/In	Yes	Post	HERP	Right	0	4	2.054 IDC	127.1	Living	Positive	No	125.43	Not Rec	Female	ER-H-HE1	27	2	Living	
ME-012	51.89	Mastect Breast C	Breast In High	No	LumA	1	Positive	Positive	3	Neutral	Negative Lobular	Yes	Post	HERP	Right	3	3	6.3 LC	26.867	Deceased	Negative	Yes	25.13	Recurer	Female	ER-H-HE1	160	3	Died of Disease	
ME-013	36.96	Mastect Breast C	Breast In Low	Yes	HER2	1	Positive	Negative	3	Gain	Positive Ductal/In	Yes	Pre	HERP	Right	3	3	5.034 IDC	43.867	Living	Negative	Yes	42.6	Not Rec	Female	HER2+	17	2	Living	
ME-014	48.59	Mastect Breast C	Breast In Low	No	LumA	1	Positive	Positive	2	Neutral	Negative Ductal/In	Yes	Pre	HERP	Left	0	3	3.06 IDC	114	Living	Positive	No	112.22	Not Rec	Female	ER-H-HE1	30	2	Living	
ME-015	39.84	Mastect Breast C	Breast In Moderate	Yes	Basal	1	Negative	Negative	3	Neutral	Negative Ductal/In	No	Pre	HERP	Right	0	5	4.05 IDC	66.733	Deceased	Negative	Yes	31.74	Recurer	Female	ER-H-HE1	25	2	Died of Disease	
ME-016	42.55	Mastect Breast C	Breast In High	No	LumB	1	Positive	Positive	2	Neutral	Negative Lobular	Yes	Pre	HERP	Left	1	2	4.12 LC	122.27	Living	Positive	Yes	120.65	Not Rec	Female	ER-H-HE1	60	3	Living	
ME-017	60.07	Breast C	Breast C	Breast In Moderate	No	LumA	1	Positive	Positive	2	Neutral	Negative Mixed	Yes	Post	HERP	Left	1	2	4.046 MDLC	24	Living	Negative	Yes	23.7	Not Rec	Female	ER-H-HE1	23	2	Living
ME-018	82.73	Mastect Breast C	Breast In High	No	LumB	1	Positive	Positive	2	Gain	Negative Ductal/In	Yes	Post	HERP	Left	1	2	4.046 IDC	35.887	Deceased	Negative	Yes	34.81	Recurer	Female	ER-H-HE1	23	2	Died of Disease	
ME-002	72.1	Mastect Breast C	Breast In Moderate	No	LumB	1	Positive	Positive	3	Gain	Positive Ductal/In	Yes	Post	HERP	Left	1	3	5.052 IDC	23.087	Deceased	Negative	Yes	23.88	Recurer	Female	ER-H-HE1	26	2	Died of Disease	
ME-021	78.73	Mastect Breast C	Breast In Moderate	No	LumA	1	Positive	Positive	2	Neutral	Negative Ductal/In	Yes	Post	HERP	Left	6	4	5.06 IDC	152.2	Living	Negative	Yes	150.2	Not Rec	Female	ER-H-HE1	30	2	Living	
ME-022	58.95	Breast C	Breast C	Breast In Moderate	No	LumA	1	Positive	Positive	1	Neutral	N																		

6	Chemotherapy	1980	non-null	object
7	Pam50 + Claudin-low subtype	1980	non-null	object
8	Cohort	2498	non-null	float64
9	ER status measured by IHC	2426	non-null	object
10	ER Status	2469	non-null	object
11	Neoplasm Histologic Grade	2388	non-null	float64
12	HER2 status measured by SNP6	1980	non-null	object
13	HER2 Status	1980	non-null	object
14	Tumor Other Histologic Subtype	2374	non-null	object
15	Hormone Therapy	1980	non-null	object
16	Inferred Menopausal State	1980	non-null	object
17	Integrative Cluster	1980	non-null	object
18	Primary Tumor Laterality	1870	non-null	object
19	Lymph nodes examined positive	2243	non-null	float64
20	Mutation Count	2357	non-null	float64
21	Nottingham prognostic index	2287	non-null	float64
22	Oncotree Code	2509	non-null	object
23	Overall Survival (Months)	1981	non-null	float64
24	Overall Survival Status	1981	non-null	object
25	PR Status	1980	non-null	object
26	Radio Therapy	1980	non-null	object
27	Relapse Free Status (Months)	2388	non-null	float64
28	Relapse Free Status	2488	non-null	object
29	Sex	2509	non-null	object
30	3-Gene classifier subtype	1764	non-null	object
31	Tumor Size	2360	non-null	float64
32	Tumor Stage	1788	non-null	float64
33	Patient's Vital Status	1980	non-null	object

dtypes: float64(10), object(24)

memory usage: 666.6+ KB

c=cancer.describe()

c

```
{
  "summary": {
    "name": "c",
    "rows": 8,
    "fields": [
      {
        "column": "Age at Diagnosis",
        "properties": {
          "dtype": "number",
          "std": 864.702216344159,
          "min": 13.032997167502558,
          "max": 2498.0,
          "num_unique_values": 8,
          "samples": [
            60.420300240192155,
            61.11,
            2498.0
          ],
          "semantic_type": "",
          "description": ""
        }
      },
      {
        "column": "Cohort",
        "properties": {
          "dtype": "number",
          "std": 882.0253293425759,
          "min": 1.0,
          "max": 2498.0,
          "num_unique_values": 7,
          "samples": [
            2498.0,
            2.900320256204964,
            4.0
          ],
          "semantic_type": "",
          "description": ""
        }
      },
      {
        "column": "Neoplasm Histologic Grade",
        "properties": {
          "dtype": "number",
          "std": 843.5252739053012,
          "min": 0.6493632804610502,
          "max": 2388.0,
          "num_unique_values": 6,
          "samples": [
            2388.0,
            2.4120603015075375,
            3.0
          ],
          "semantic_type": "",
          "description": ""
        }
      }
    ]
  }
}
```

```

semantic_type\": \"\", \n      \"description\": \"\" \n      } \n      }, \n      {
\n      \"column\": \"Lymph nodes examined positive\", \n      \"properties\":
{\n      \"dtype\": \"number\", \n      \"std\": 790.4938917430783, \n
\n      \"min\": 0.0, \n      \"max\": 2243.0, \n      \"num_unique_values\": 6, \n
\n      \"samples\": [\n      2243.0, \n      1.9505127061970575, \n      4
5.0 \n      ], \n      \"semantic_type\": \"\", \n      \"description\": \"
\n      \"\" \n      } \n      }, \n      {\n      \"column\": \"Mutation Count\", \n      \"p
roperties\": {\n      \"dtype\": \"number\", \n      \"std\": 828.42042940
53787, \n      \"min\": 1.0, \n      \"max\": 2357.0, \n      \"num_unique
_values\": 8, \n      \"samples\": [\n      5.578701739499364, \n
5.0, \n      2357.0 \n      ], \n      \"semantic_type\": \"\", \n
\n      \"description\": \"\" \n      } \n      }, \n      {\n      \"column\": \"Nottingham
prognostic index\", \n      \"properties\": {\n      \"dtype\": \"number\", \n
\n      \"std\": 807.2886422446937, \n      \"min\": 1.0, \n      \"max\":
2287.0, \n      \"num_unique_values\": 8, \n      \"samples\": [\n
4.028786847398338, \n      4.044, \n      2287.0 \n      ], \n      \
\n      \"semantic_type\": \"\", \n      \"description\": \"\" \n      } \n      }, \n
\n      {\n      \"column\": \"Overall Survival (Months)\", \n      \"properties\": {\n
\n      \"dtype\": \"number\", \n      \"std\": 662.4935893825111, \n
\n      \"min\": 0.0, \n      \"max\": 1981.0, \n      \"num_unique_values\": 8, \n
\n      \"samples\": [\n      125.24427057011306, \n      116.46666670000002, \n
\n      1981.0 \n      ], \n      \"semantic_type\": \"\", \n      \"de
scription\": \"\" \n      } \n      }, \n      {\n      \"column\": \"Relapse Free S
tatus (Months)\", \n      \"properties\": {\n      \"dtype\": \"number\", \n
\n      \"std\": 808.4386888706275, \n      \"min\": 0.0, \n      \"max\": 2388.0, \n
\n      \"num_unique_values\": 8, \n      \"samples\": [\n      108.842
48743718592, \n      99.095, \n      2388.0 \n      ], \n      \"sema
ntic_type\": \"\", \n      \"description\": \"\" \n      } \n      }, \n      {\n
\n      \"column\": \"Tumor Size\", \n      \"properties\": {\n      \"dtype\": \"nu
mber\", \n      \"std\": 821.5729528236607, \n      \"min\": 1.0, \n
\n      \"max\": 2360.0, \n      \"num_unique_values\": 8, \n      \"samples\": [\n
26.22009322033898, \n      22.41, \n      2360.0 \n      ], \n      \
\n      \"semantic_type\": \"\", \n      \"description\": \"\" \n      } \n      }, \n
\n      {\n      \"column\": \"Tumor Stage\", \n      \"properties\": {\n      \"dty
pe\": \"number\", \n      \"std\": 631.5803632892699, \n      \"min\": 0.0,
\n      \"max\": 1788.0, \n      \"num_unique_values\": 7, \n      \"samp
les\": [\n      1788.0, \n      1.7136465324384786, \n      2.0 \n
], \n      \"semantic_type\": \"\", \n      \"description\": \"\" \n
\n      } \n      } \n      ] \n      }\", \"type\": \"dataframe\", \"variable_name\": \"c\"}

```

cancer.columns

```

Index(['Patient ID', 'Age at Diagnosis', 'Type of Breast Surgery',
      'Cancer Type', 'Cancer Type Detailed', 'Cellularity', 'Chemotherapy',
      'Pam50 + Claudin-low subtype', 'Cohort', 'ER status measured by IHC',
      'ER Status', 'Neoplasm Histologic Grade',
      'HER2 status measured by SNP6', 'HER2 Status',
      'Tumor Other Histologic Subtype', 'Hormone Therapy',
      'Inferred Menopausal State', 'Integrative Cluster',
      'Primary Tumor Laterality', 'Lymph nodes examined positive',

```

```

        'Mutation Count', 'Nottingham prognostic index', 'Oncotree Code',
        'Overall Survival (Months)', 'Overall Survival Status', 'PR Status',
        'Radio Therapy', 'Relapse Free Status (Months)', 'Relapse Free Status'
    ,
        'Sex', '3-Gene classifier subtype', 'Tumor Size', 'Tumor Stage',
        'Patient's Vital Status'],
    dtype='object')

```

```

from sklearn.datasets import load_breast_cancer

```

```

cancer = load_breast_cancer()

```

```

cancer.keys()

```

```

dict_keys(['data', 'target', 'frame', 'target_names', 'DESCR', 'feature_names',
           'filename', 'data_module'])

```

```

print(cancer ['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 (perimeter² / area - 1.0)
- concavity (severity of concave portions of the contour)
- concave points (number of concave portions of the contour)
- symmetry
- fractal dimension ("coastline approximation" - 1)

The mean, standard error, and "worst" or largest (mean of the three worst/largest values) of these features were computed for each image, resulting in 30 features. For instance, field 0 is Mean Radius, fiel

d

10 is Radius SE, field 20 is Worst Radius.

- class:

- WDBC-Malignant
- WDBC-Benign

:Summary Statistics:

=====	=====	=====
	Min	Max
=====	=====	=====
radius (mean):	6.981	28.11
texture (mean):	9.71	39.28
perimeter (mean):	43.79	188.5
area (mean):	143.5	2501.0
smoothness (mean):	0.053	0.163
compactness (mean):	0.019	0.345
concavity (mean):	0.0	0.427
concave points (mean):	0.0	0.201
symmetry (mean):	0.106	0.304
fractal dimension (mean):	0.05	0.097
radius (standard error):	0.112	2.873
texture (standard error):	0.36	4.885
perimeter (standard error):	0.757	21.98
area (standard error):	6.802	542.2
smoothness (standard error):	0.002	0.031
compactness (standard error):	0.002	0.135
concavity (standard error):	0.0	0.396
concave points (standard error):	0.0	0.053
symmetry (standard error):	0.008	0.079
fractal dimension (standard error):	0.001	0.03
radius (worst):	7.93	36.04
texture (worst):	12.02	49.54
perimeter (worst):	50.41	251.2
area (worst):	185.2	4254.0
smoothness (worst):	0.071	0.223
compactness (worst):	0.027	1.058
concavity (worst):	0.0	1.252
concave points (worst):	0.0	0.291
symmetry (worst):	0.156	0.664
fractal dimension (worst):	0.055	0.208
=====	=====	=====

:Missing Attribute Values: None

:Class Distribution: 212 - Malignant, 357 - Benign

:Creator: Dr. William H. Wolberg, W. Nick Street, Olvi L. Mangasarian

:Donor: Nick Street

:Date: November, 1995

This is a copy of UCI ML Breast Cancer Wisconsin (Diagnostic) datasets.

<https://goo.gl/U2Uwz2>

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/
```

.. topic:: References

- W.N. Street, W.H. Wolberg and O.L. Mangasarian. Nuclear feature extraction for breast tumor diagnosis. IS&T/SPIE 1993 International Symposium on Electronic Imaging: Science and Technology, volume 1905, pages 861-870, San Jose, CA, 1993.
- O.L. Mangasarian, W.N. Street and W.H. Wolberg. Breast cancer diagnosis and prognosis via linear programming. Operations Research, 43(4), pages 570-577, July-August 1995.
- W.H. Wolberg, W.N. Street, and O.L. Mangasarian. Machine learning techniques to diagnose breast cancer from fine-needle aspirates. Cancer Letters 77 (1994) 163-171.

cancer ['feature_names']

```
array(['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 radius', 'worst texture',
'worst perimeter', 'worst area', 'worst smoothness',
'worst compactness', 'worst concavity', 'worst concave points',
'worst symmetry', 'worst fractal dimension'], dtype='<U23')

```

```

df = pd.DataFrame(cancer['data'], columns=cancer['feature_names'])
df.info()

```

```

<class 'pandas.core.frame.DataFrame'>

```

```

RangeIndex: 569 entries, 0 to 568

```

```

Data columns (total 30 columns):

```

#	Column	Non-Null Count	Dtype
0	mean radius	569 non-null	float64
1	mean texture	569 non-null	float64
2	mean perimeter	569 non-null	float64
3	mean area	569 non-null	float64
4	mean smoothness	569 non-null	float64
5	mean compactness	569 non-null	float64
6	mean concavity	569 non-null	float64
7	mean concave points	569 non-null	float64
8	mean symmetry	569 non-null	float64
9	mean fractal dimension	569 non-null	float64
10	radius error	569 non-null	float64
11	texture error	569 non-null	float64
12	perimeter error	569 non-null	float64
13	area error	569 non-null	float64
14	smoothness error	569 non-null	float64
15	compactness error	569 non-null	float64
16	concavity error	569 non-null	float64
17	concave points error	569 non-null	float64
18	symmetry error	569 non-null	float64
19	fractal dimension error	569 non-null	float64
20	worst radius	569 non-null	float64
21	worst texture	569 non-null	float64
22	worst perimeter	569 non-null	float64
23	worst area	569 non-null	float64
24	worst smoothness	569 non-null	float64
25	worst compactness	569 non-null	float64
26	worst concavity	569 non-null	float64
27	worst concave points	569 non-null	float64
28	worst symmetry	569 non-null	float64
29	worst fractal dimension	569 non-null	float64

```

dtypes: float64(30)

```

```

memory usage: 133.5 KB

```

```

np.sum(pd.isnull(df).sum())

```

```

0

```

```
df = pd.DataFrame(cancer['data'], columns=cancer['feature_names'])
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 569 entries, 0 to 568
```

```
Data columns (total 30 columns):
```

#	Column	Non-Null Count	Dtype
0	mean radius	569 non-null	float64
1	mean texture	569 non-null	float64
2	mean perimeter	569 non-null	float64
3	mean area	569 non-null	float64
4	mean smoothness	569 non-null	float64
5	mean compactness	569 non-null	float64
6	mean concavity	569 non-null	float64
7	mean concave points	569 non-null	float64
8	mean symmetry	569 non-null	float64
9	mean fractal dimension	569 non-null	float64
10	radius error	569 non-null	float64
11	texture error	569 non-null	float64
12	perimeter error	569 non-null	float64
13	area error	569 non-null	float64
14	smoothness error	569 non-null	float64
15	compactness error	569 non-null	float64
16	concavity error	569 non-null	float64
17	concave points error	569 non-null	float64
18	symmetry error	569 non-null	float64
19	fractal dimension error	569 non-null	float64
20	worst radius	569 non-null	float64
21	worst texture	569 non-null	float64
22	worst perimeter	569 non-null	float64
23	worst area	569 non-null	float64
24	worst smoothness	569 non-null	float64
25	worst compactness	569 non-null	float64
26	worst concavity	569 non-null	float64
27	worst concave points	569 non-null	float64
28	worst symmetry	569 non-null	float64
29	worst fractal dimension	569 non-null	float64

```
dtypes: float64(30)
```

```
memory usage: 133.5 KB
```

```
df['cancer'] =pd.DataFrame(cancer['target'])
df.head()
```

```
{"type":"dataframe","variable_name":"df"}
```

```
sns.set_style('whitegrid')
```

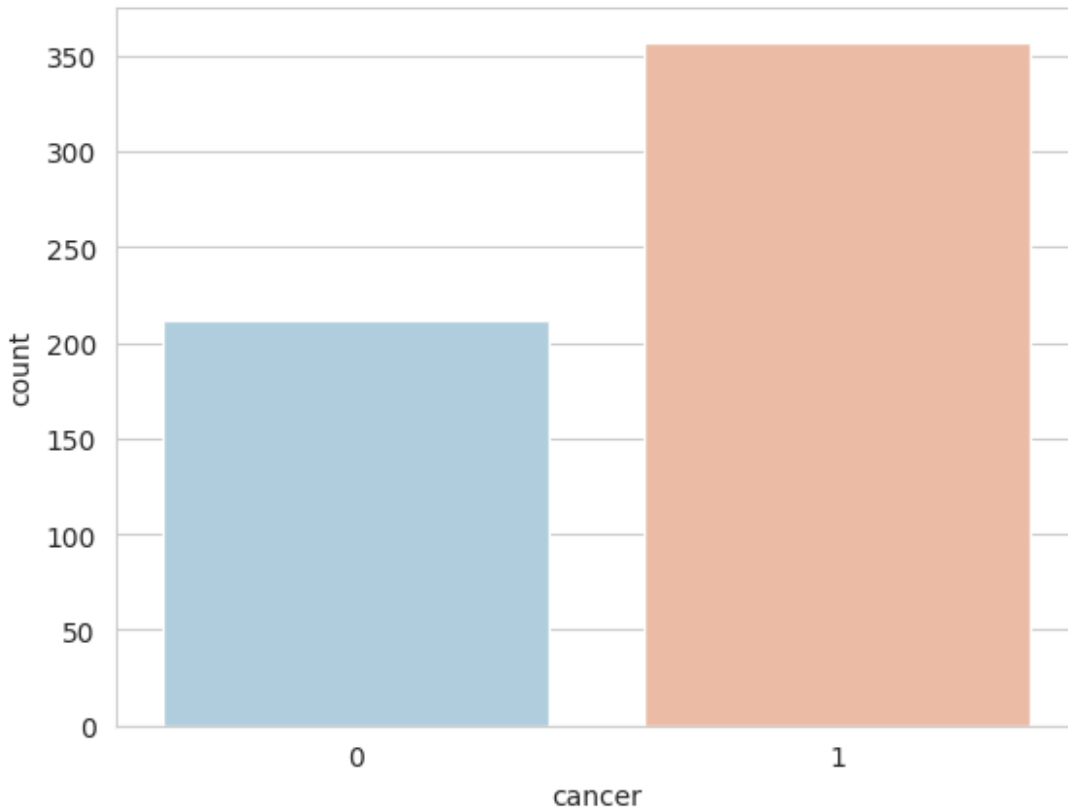
```
sns.countplot(x='cancer',data=df,palette='RdBu_r')
```

```
<ipython-input-15-d3c071f40503>:2: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
sns.countplot(x='cancer',data=df,palette='RdBu_r')
```

```
<Axes: xlabel='cancer', ylabel='count'>
```



```
l=list(df.columns[0:10])
for i in range(len(l)-1):
    sns.boxplot(x='cancer',y=l[i], data=df , palette='winter')
plt.figure()
```

```
<ipython-input-16-c50b4d1c2876>:3: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
sns.boxplot(x='cancer',y=l[i], data=df , palette='winter')
```

```
<ipython-input-16-c50b4d1c2876>:3: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
sns.boxplot(x='cancer',y=l[i], data=df , palette='winter')  
<ipython-input-16-c50b4d1c2876>:3: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
sns.boxplot(x='cancer',y=l[i], data=df , palette='winter')  
<ipython-input-16-c50b4d1c2876>:3: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
sns.boxplot(x='cancer',y=l[i], data=df , palette='winter')  
<ipython-input-16-c50b4d1c2876>:3: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
sns.boxplot(x='cancer',y=l[i], data=df , palette='winter')  
<ipython-input-16-c50b4d1c2876>:3: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

```
sns.boxplot(x='cancer',y=l[i], data=df , palette='winter')  
<ipython-input-16-c50b4d1c2876>:3: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

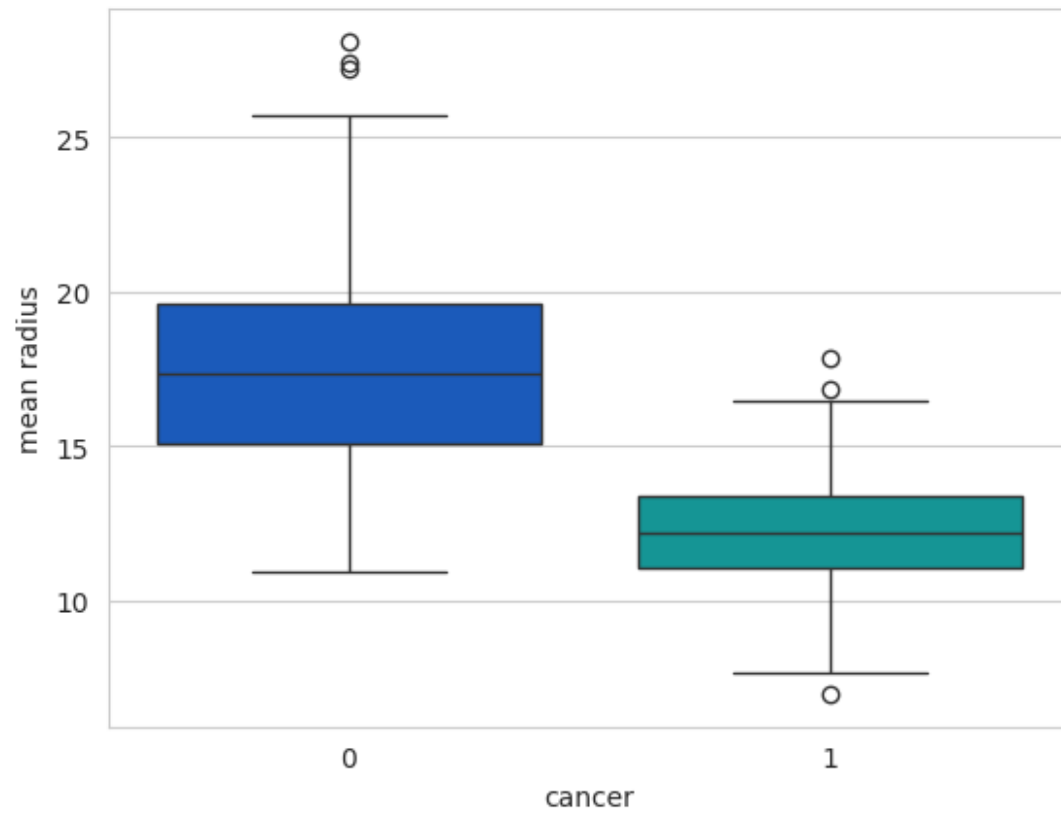
```
sns.boxplot(x='cancer',y=l[i], data=df , palette='winter')  
<ipython-input-16-c50b4d1c2876>:3: FutureWarning:
```

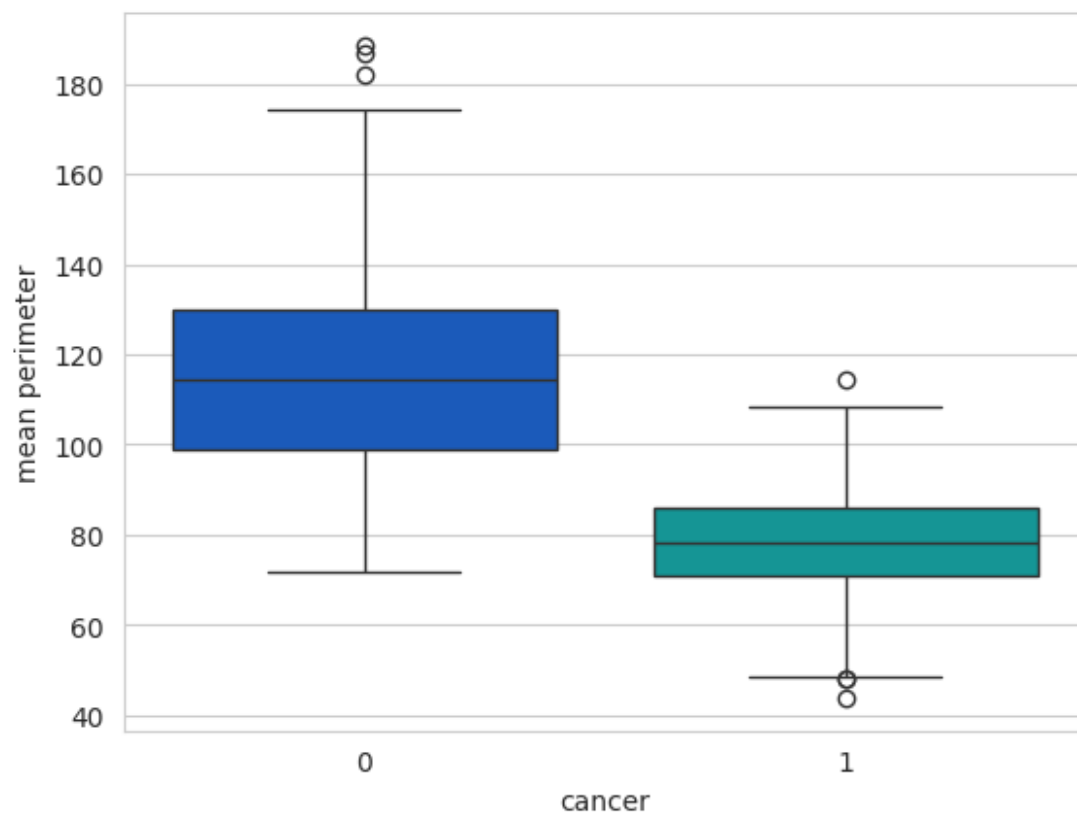
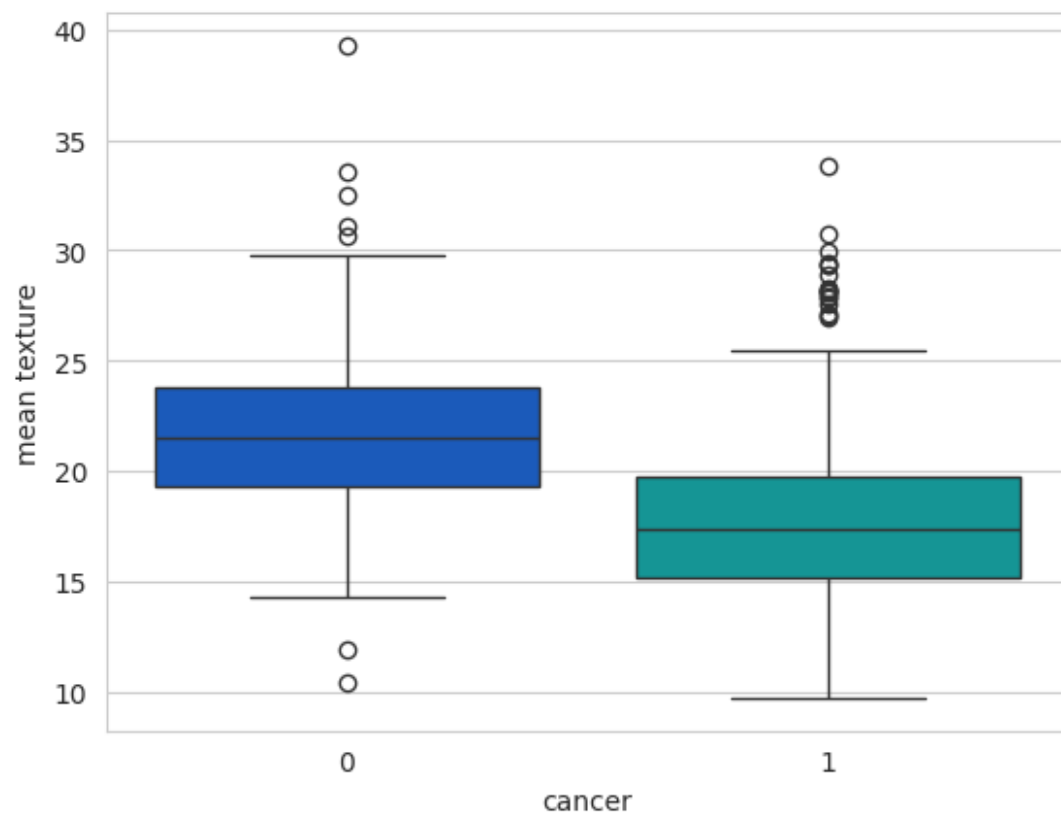
Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

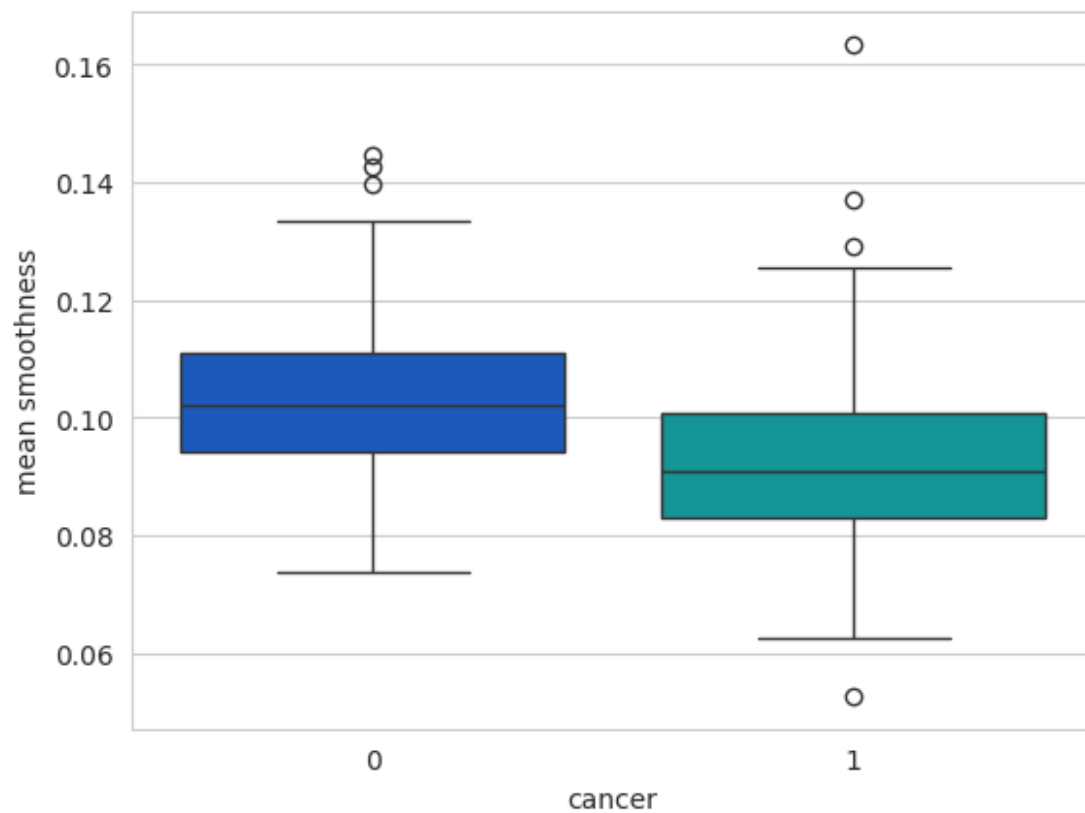
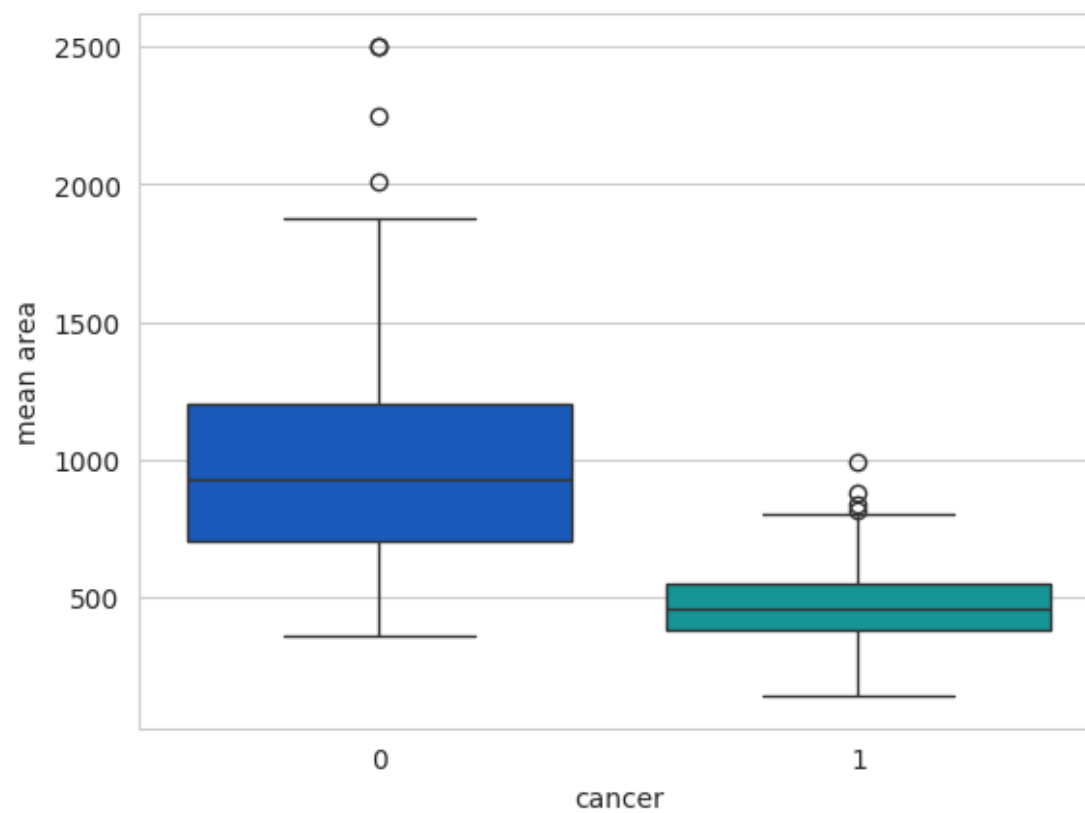
```
sns.boxplot(x='cancer',y=l[i], data=df , palette='winter')  
<ipython-input-16-c50b4d1c2876>:3: FutureWarning:
```

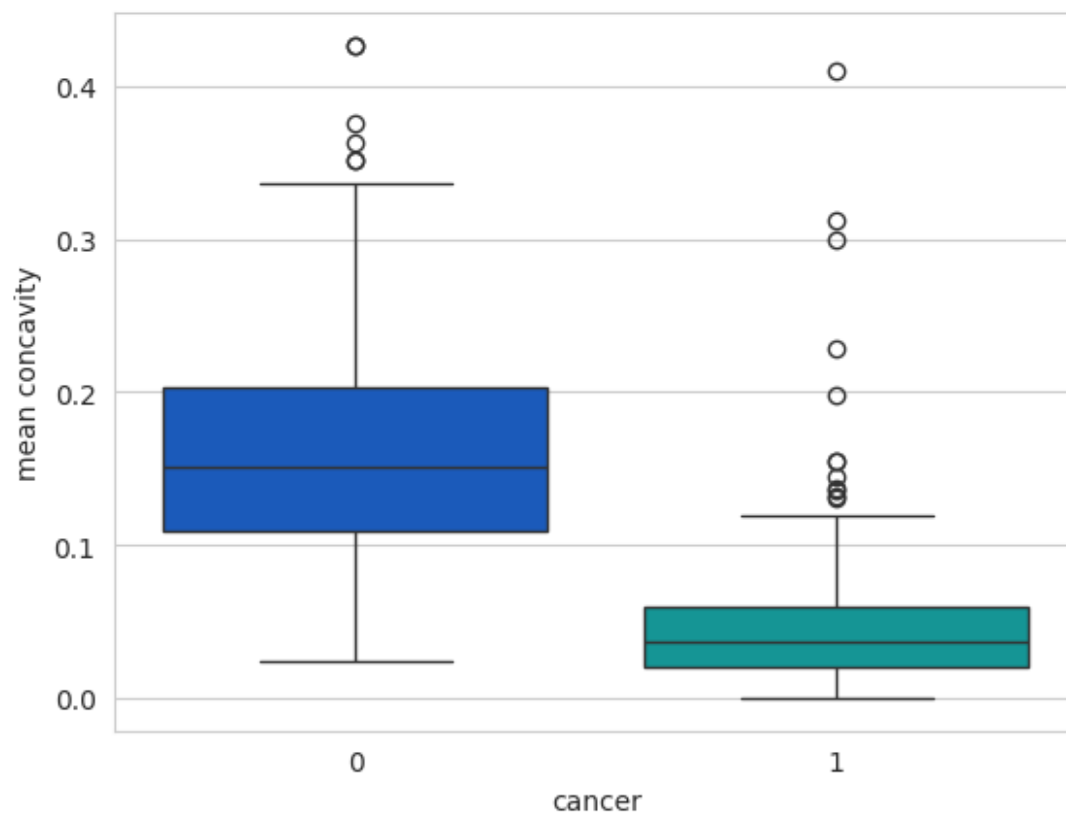
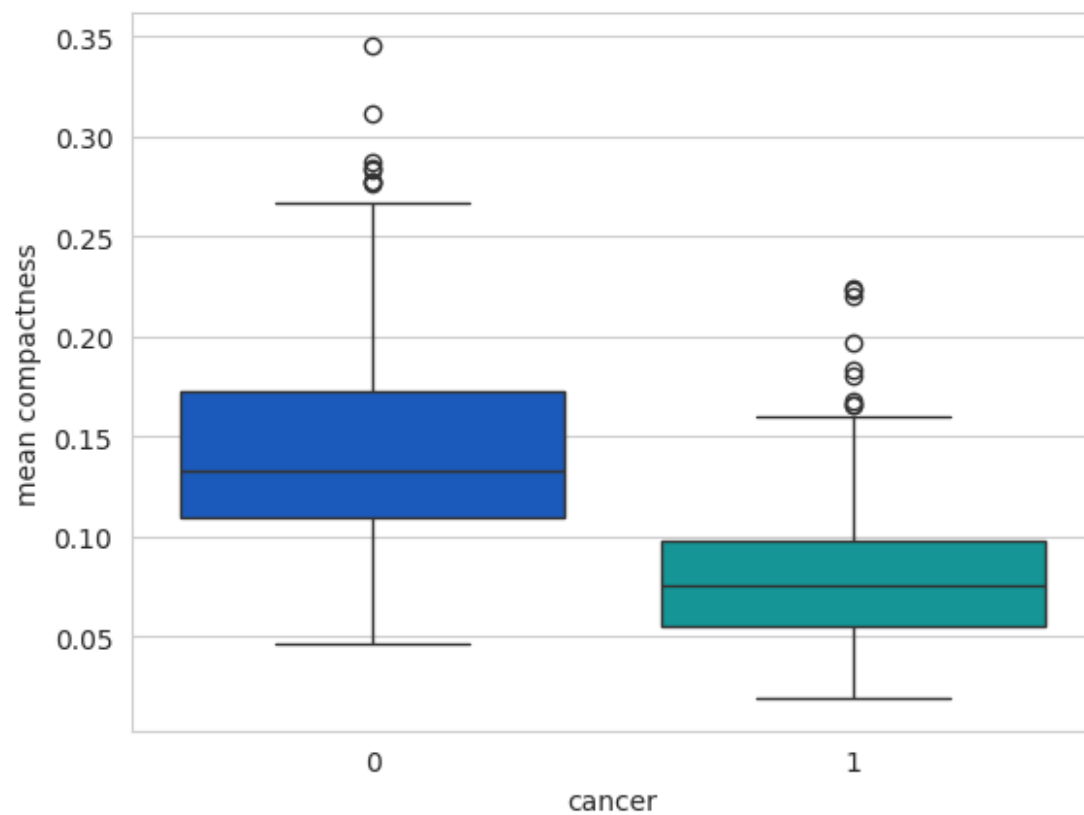
Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend=False` for the same effect.

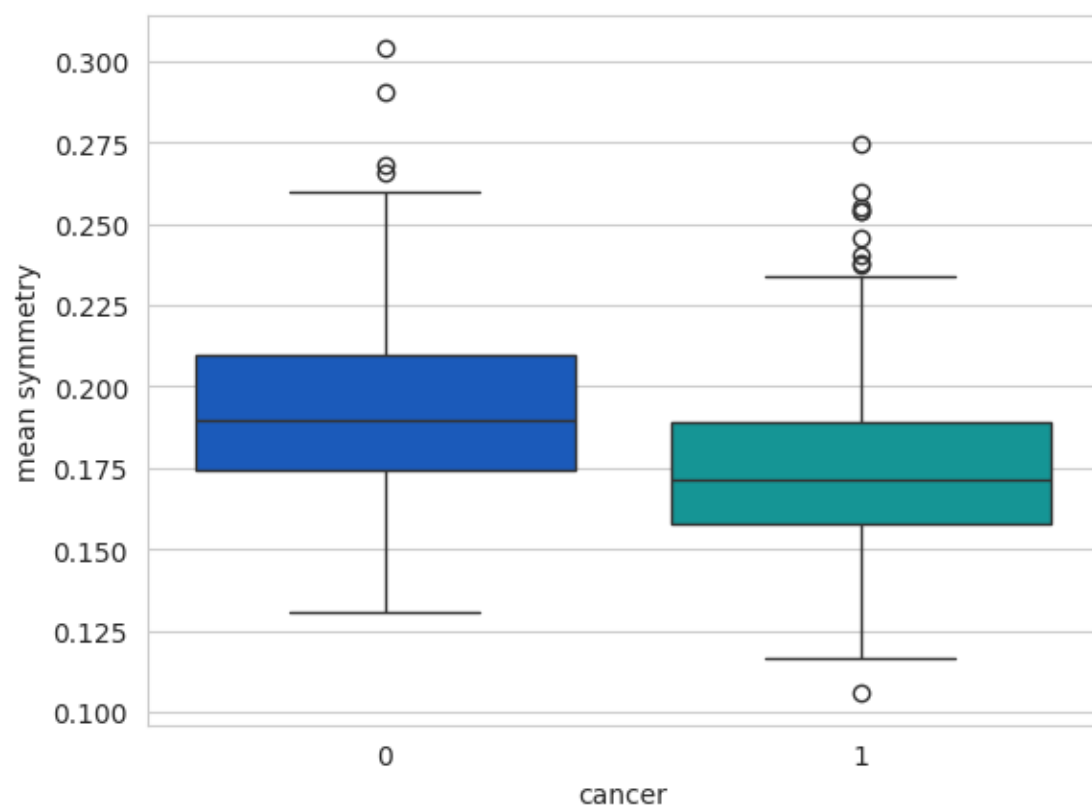
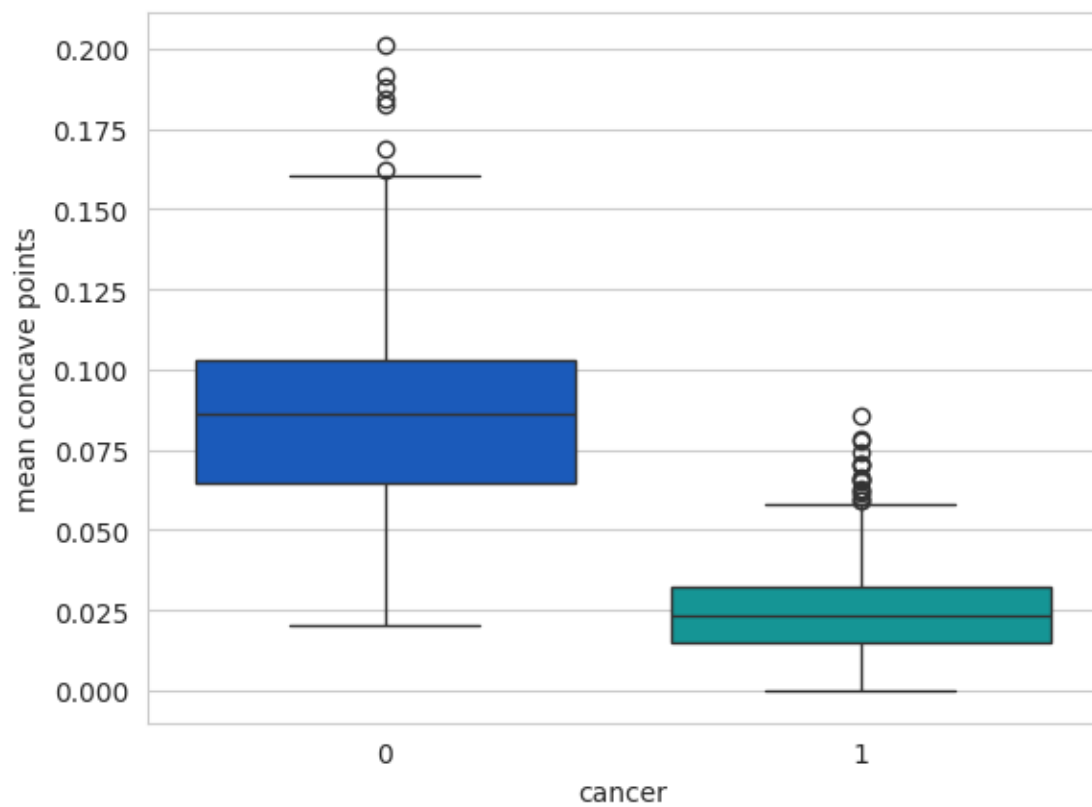
```
sns.boxplot(x='cancer',y=l[i], data=df , palette='winter')
```







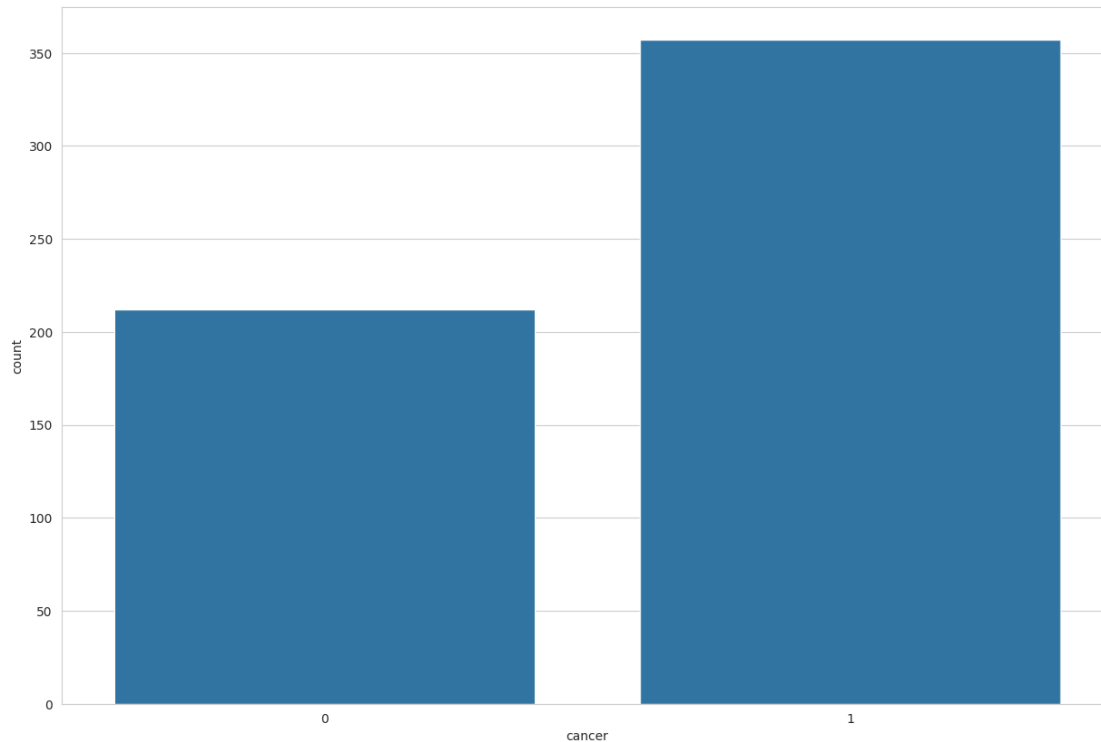




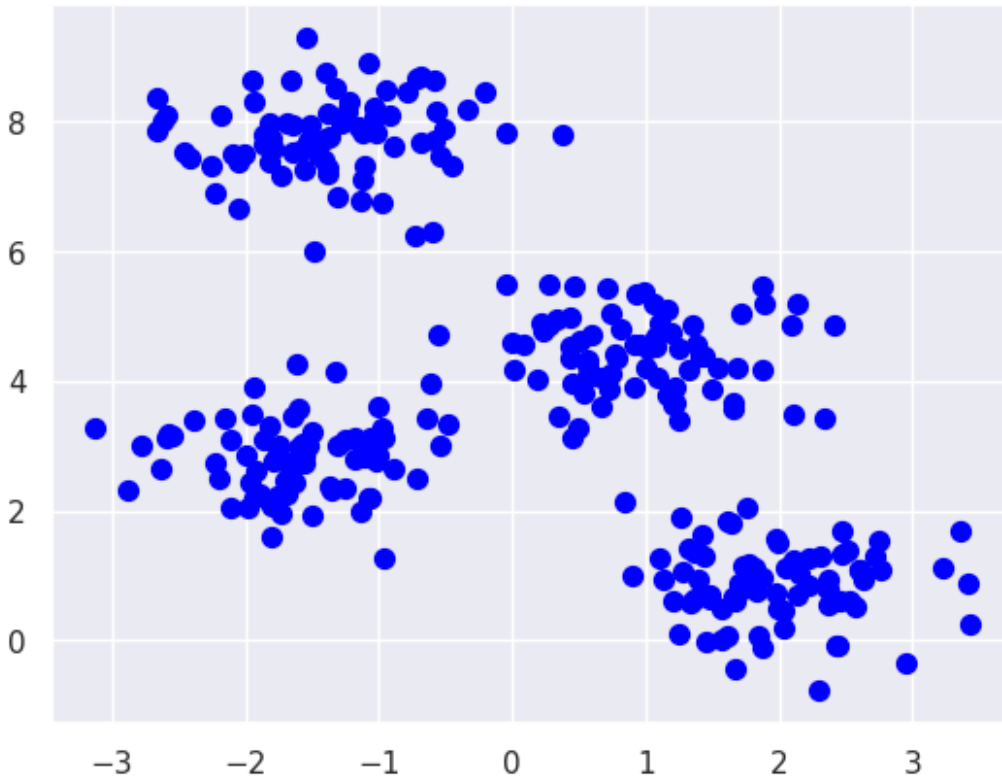
<Figure size 640x480 with 0 Axes>

```
plt.figure(figsize=(15,10))  
sns.countplot(data=df, x='cancer')
```

<Axes: xlabel='cancer', ylabel='count'>



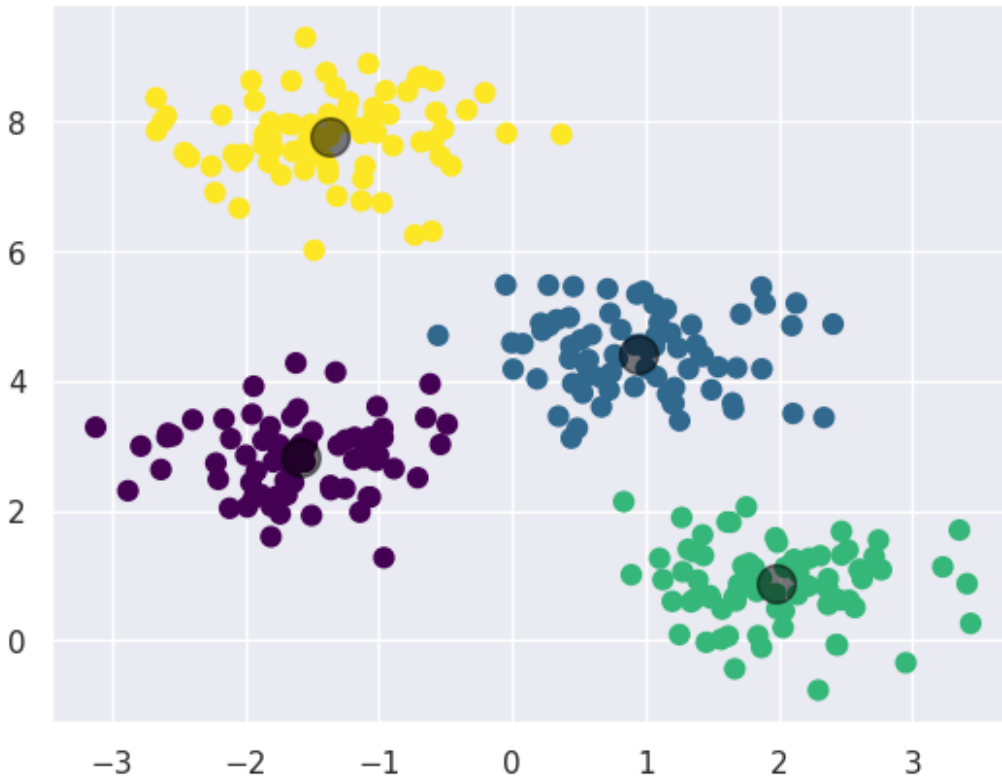
```
import matplotlib.pyplot as plt  
import seaborn as sns  
sns.set() #Plot styling  
import numpy as np  
  
from sklearn.datasets import make_blobs  
X, y_true = make_blobs(n_samples=300, centers=4,  
                        cluster_std=0.60, random_state=0)  
plt.scatter(X[:, 0], X[:, 1], s=50, color='blue');
```



```
from sklearn.cluster import KMeans
kmeans = KMeans(n_clusters=4, n_init=10)
kmeans.fit(X)
y_kmeans = kmeans.predict(X)

plt.scatter(X[:, 0], X[:, 1], c=y_kmeans, s=50, cmap='viridis')

centers = kmeans.cluster_centers_
plt.scatter(centers[:, 0], centers[:, 1], c='black', s=200, alpha=0.5);
```



```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 569 entries, 0 to 568
```

```
Data columns (total 31 columns):
```

#	Column	Non-Null Count	Dtype
0	mean radius	569 non-null	float64
1	mean texture	569 non-null	float64
2	mean perimeter	569 non-null	float64
3	mean area	569 non-null	float64
4	mean smoothness	569 non-null	float64
5	mean compactness	569 non-null	float64
6	mean concavity	569 non-null	float64
7	mean concave points	569 non-null	float64
8	mean symmetry	569 non-null	float64
9	mean fractal dimension	569 non-null	float64
10	radius error	569 non-null	float64
11	texture error	569 non-null	float64
12	perimeter error	569 non-null	float64
13	area error	569 non-null	float64
14	smoothness error	569 non-null	float64
15	compactness error	569 non-null	float64

```

16  concavity error          569 non-null    float64
17  concave points error    569 non-null    float64
18  symmetry error          569 non-null    float64
19  fractal dimension error  569 non-null    float64
20  worst radius            569 non-null    float64
21  worst texture           569 non-null    float64
22  worst perimeter         569 non-null    float64
23  worst area              569 non-null    float64
24  worst smoothness        569 non-null    float64
25  worst compactness       569 non-null    float64
26  worst concavity         569 non-null    float64
27  worst concave points    569 non-null    float64
28  worst symmetry          569 non-null    float64
29  worst fractal dimension  569 non-null    float64
30  cancer                  569 non-null    int64
dtypes: float64(30), int64(1)
memory usage: 137.9 KB

scaler.fit(df.drop('cancer',axis=1))
scaled_features = scaler.transform(df.drop('cancer',axis=1))

df_feat = pd.DataFrame(scaled_features,columns=df.columns[:-1])
df_feat.head()

{"type":"dataframe","variable_name":"df_feat"}

from sklearn.model_selection import train_test_split
X = df_feat
y = df['cancer']
X_train, X_test, y_train, y_test = train_test_split(scaled_features,df['cancer'],
                                                    test_size=0.30, random_state=101)

from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors=1)
knn.fit(X_train,y_train)

KNeighborsClassifier(n_neighbors=1)

pred = knn.predict(X_test)

from sklearn.metrics import classification_report,confusion_matrix
conf_mat=confusion_matrix(y_test,pred)
print(conf_mat)

[[ 61   5]
 [  3 102]]

print(classification_report(y_test,pred))

```

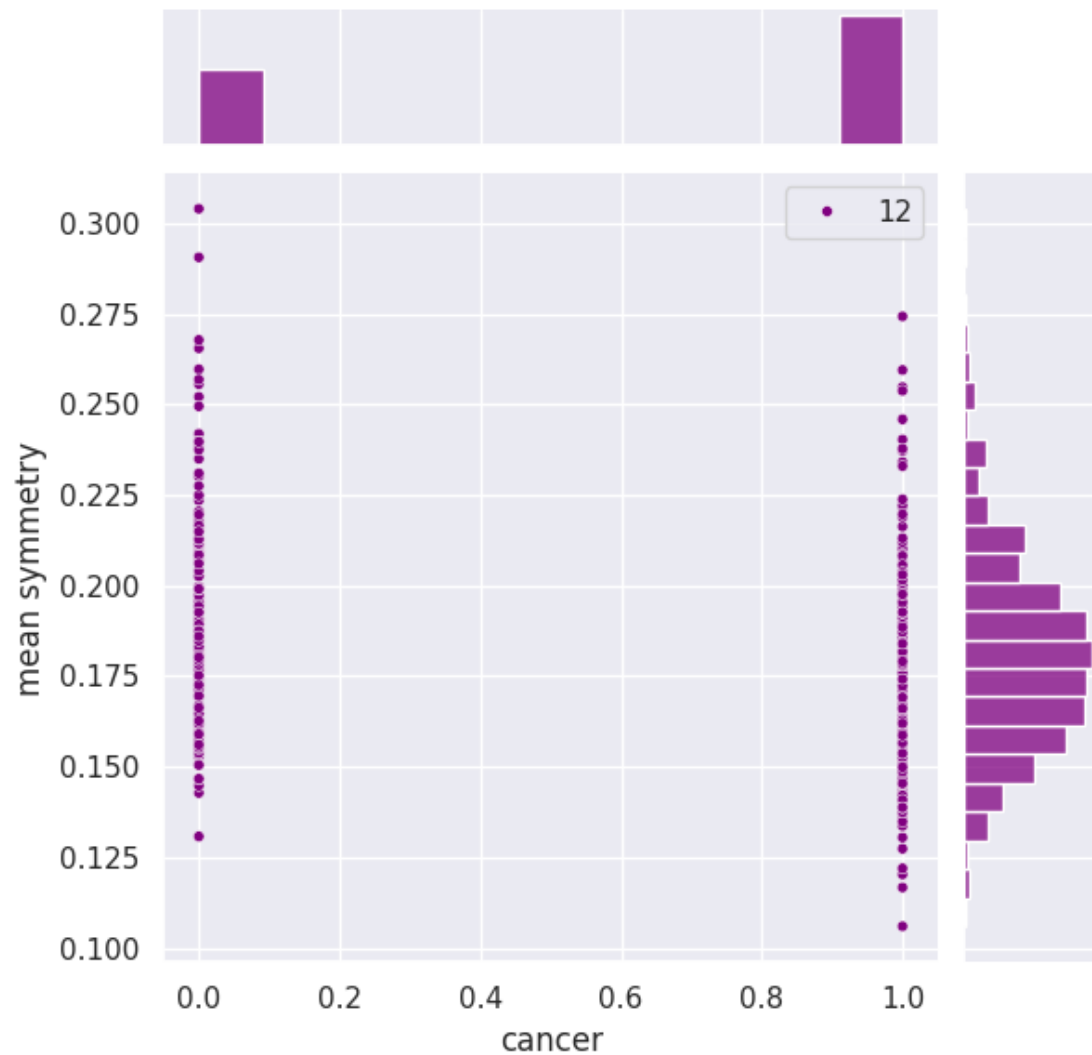
	precision	recall	f1-score	support
0	0.95	0.92	0.94	66
1	0.95	0.97	0.96	105
accuracy			0.95	171
macro avg	0.95	0.95	0.95	171
weighted avg	0.95	0.95	0.95	171

cancer['target']

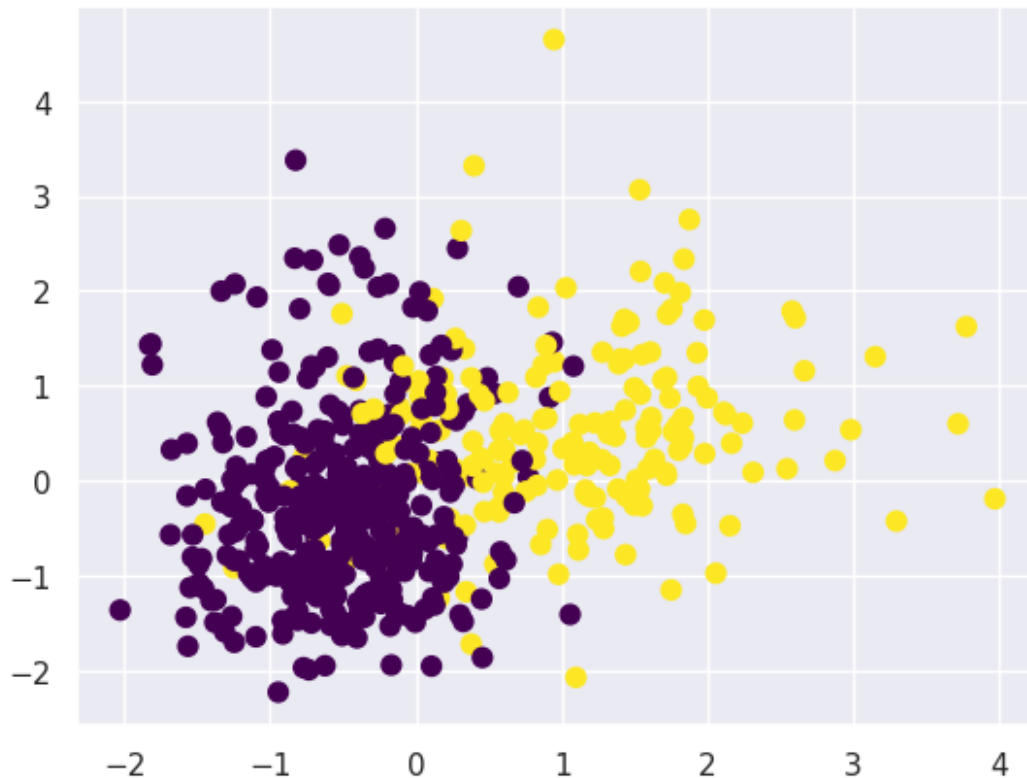
```
array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1,
       0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
       0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0,
       1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 0,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

```
sns.jointplot(x='cancer',y=l[i],data=df, color='purple', size=12)
```

```
<seaborn.axisgrid.JointGrid at 0x7b56d5b46ad0>
```



```
import matplotlib.pyplot as plt
from sklearn.cluster import SpectralClustering
import numpy as np
model = SpectralClustering(n_clusters=2, affinity='nearest_neighbors', assign_labels='kmeans')
labels = model.fit_predict(X)
plt.scatter(X[:,0], X[:,1], c=labels, s=50, cmap='viridis')
plt.show()
```



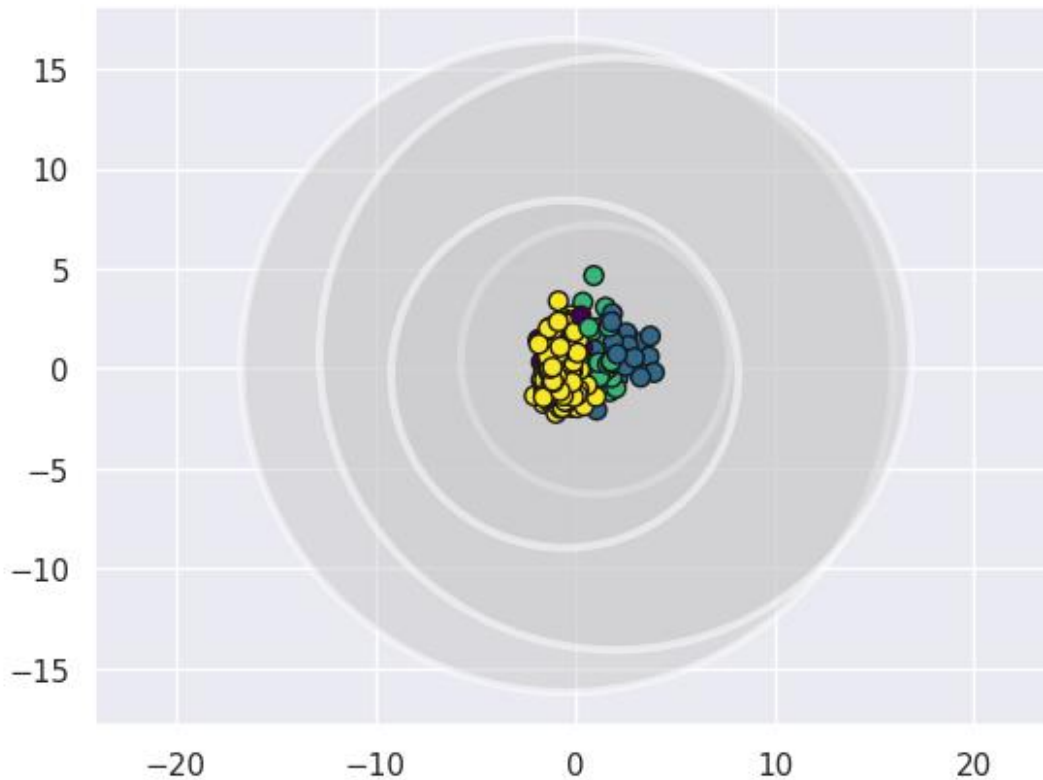
```

from sklearn.cluster import KMeans
from scipy.spatial.distance import cdist

def plot_kmeans(kmeans, X, n_clusters=4, rseed=0, ax=None):
    labels = kmeans.fit_predict(X)
    ax = ax or plt.gca()
    ax.axis('equal')
    ax.scatter(X[:, 0], X[:, 1], c=labels, s=50, cmap='viridis', edgecolor='k', zorder=2)
    centers = kmeans.cluster_centers_
    radii = [cdist(X[labels == i], [center]).max()
              for i, center in enumerate(centers)]
    for c, r in zip(centers, radii):
        ax.add_patch(plt.Circle(c, r, fc='#CCCCCC', lw=3, alpha=0.5, zorder=1))
    ))

kmeans = KMeans(n_clusters=4, random_state=0, n_init=10)
plot_kmeans(kmeans, X)

```

Results And Analysis :

Breast cancer is the second leading cause of cancer death in women, second only to lung cancer.

The leading risk factor for breast cancer is simply being a woman. Though breast cancer does occur in men, the disease is 100 times more common in women.

Men can also get breast cancer. In 2017, the American Cancer Society estimates 2,470 new cases of invasive breast cancer will be diagnosed in men in the U.S. A woman has about a one in eight chance of being diagnosed with breast cancer in her lifetime, according to the National Cancer Institute.

Most women (about eight out of 10) who get breast cancer do not have a family history of the disease.

Conclusion

By performing suitable morphological operations, system computes the suitable region properties such as Area, Euler number etc., and displays the boundary detected image along with the tumor area. These techniques improve accuracy in tracking the breast cancer cells. To assess the correctness in classifying data with respect to efficiency and effectiveness of each algorithm in terms of accuracy, precision, sensitivity, and specificity. Hence the design is to provide high accuracy and maximum efficiency in prediction and tracking of breast cancer. The combination of Multi-Level Wavelet Conversion strategy associated to PCA with 13 features extracted and then classified gives an average accuracy of nearly 92%.

As a future improvement, the system can add more features such as recommendation of medicines/treatments based on the severity of the patient. This prediction and recommendation system can help doctors to diagnose and cure the disease more efficiently.

References:

- <https://www.kaggle.com/code/niteshyadav3103/breast-cancer-classification/input>
- Some of the data from the github accounts
- Some of codes by previous lab activities