

## Reading in Data

```
In [1]: import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.preprocessing import LabelEncoder
import numpy as np
```

```
In [2]: #reading in data.csv file
tissue_pd = pd.read_csv('project/data.csv')
```

```
In [3]: #looking at first five rows of our data
tissue_pd.head()
```

		id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	symmetry_mean	fractal_dimension_mean	radius_se	texture_se	perimeter_se	area_se	smoothness_se	compactness_se	concavity_se	symmetry_se	fractal_dimension_se	radius_worst	texture_worst	perimeter_worst	area_worst	smoothness_worst	compactness_worst	concavity_worst	symmetry_worst	fractal_dimension_worst	Unnamed: 32	dtype
0	842302	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001																							
1	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869																							
2	8430903	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974																							
3	84348301	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414																							
4	84358402	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980																							
5 rows × 33 columns																																

```
In [4]: print(tissue_pd.shape)
print("569 rows and 33 columns")
```

(569, 33)  
569 rows and 33 columns

Our project dataset is from the UCI Machine Learning Repository. Each row represents a patient and information around their breast tissue imaging, especially with respect to cell nuclei.

## Missing Data?

```
In [5]: #loop through columns in dataframe
for col in any NaN values:
    #check for any NaN values
    print(col + " : " + str(tissue_pd[col].isnull().values.any()))
```

```
id: False
diagnosis: False
radius_mean: False
texture_mean: False
perimeter_mean: False
area_mean: False
smoothness_mean: False
compactness_mean: False
concavity_mean: False
concave points_mean: False
symmetry_mean: False
fractal_dimension_mean: False
radius_se: False
texture_se: False
perimeter_se: False
area_se: False
smoothness_se: False
compactness_se: False
concavity_se: False
concave points_se: False
symmetry_se: False
fractal_dimension_se: False
radius_worst: False
texture_worst: False
perimeter_worst: False
area_worst: False
smoothness_worst: False
compactness_worst: False
concavity_worst: False
concave points_worst: False
symmetry_worst: False
fractal_dimension_worst: False
Unnamed: 32: True
dtype: object
```

There are no missing values in our data!!

## Checking for Duplicates

```
In [6]: #checking for row-level duplicates in entire dataframe

print(tissue_pd.duplicated())

print("Sum of duplicate categories: ")
print(tissue_pd.duplicated().sum())
```

```
0      False
1      False
2      False
3      False
4      False
...
564     False
565     False
566     False
567     False
568     False
Length: 569, dtype: bool

Sum of duplicate categories:
0
```

There are no row-level duplicates in our dataset!

## Data Exploration

```
In [7]: #looking at column types
tissue_pd.dtypes
```

max 28.110000 39.280000 188.500000 2501.000000 0.163400 0.345400 0.426800 0.201200

8 rows × 31 columns

From looking at the description of the numerical variables, I am going to remove the column 'Unnamed:32', as it does not have any values besides NaN. It will not be helpful in our analysis or modeling.

```
#dropping Unnamed column
tissue_pd = tissue_pd.drop(['Unnamed: 32'],axis='columns')

tissue_pd.head()
```

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	point
0	842302	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	
1	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	
2	8430903	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	
3	84348301	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	
4	84358402	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	

5 rows × 32 columns

I also want to look at the class distribution of our target variable 'diagnosis', as it will help us understand if we need to account for class imbalance in our modeling.

There are three layers of our data: the mean (L\_mean), standard error (L\_se) and "worst" (L\_worst) which denotes the largest value i.e. the mean of the three largest values in the cell data.

Our target value in terms of logistic regression is the variable 'diagnosis'. It is currently denoted as an object type; however, for regression purposes, we will convert it to a numerical type with encoding towards the end of this analysis.

All of the other variables are floats, and they are rounded to 4 significant figures.

```
In [8]: #converting id variable to type string
tissue_pd['id'] = tissue_pd['id'].astype(str)
```

```
In [9]: #describing the data
tissue_pd.describe()
```

	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	fractal_dimension_mean
count	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000
mean	14.127292	19.289649	91.969033	654.889104	0.096360	0.104341	0.089799	0.048919	0.038803	
std	3.524049	4.301056	24.289881	351.914129	0.014064	0.052813	0.079720	0.038803		
min	6.981000	9.710000	43.790000	143.500000	0.052630	0.019380	0.000000	0.000000		
25%	11.700000	16.170000	75.170000	420.300000	0.086370	0.064920	0.029560	0.020310		
50%	13.370000	18.840000	86.240000	551.100000	0.092630	0.092630	0.061540	0.033500		
75%	15.780000	21.800000	104.100000	782.700000	0.105300	0.130400	0.130700	0.074200		
max	28.110000	39.280000	188.500000	2501.000000	0.163400	0.345400	0.426800	0.201200		
8 rows × 31 columns										

From looking at the description of the numerical variables, I am going to remove the column 'Unnamed:32', as it does not have any values besides NaN. It will not be helpful in our analysis or modeling.

```
In [10]: #dropping Unnamed column
tissue_pd = tissue_pd.drop(['Unnamed: 32'],axis='columns')
```

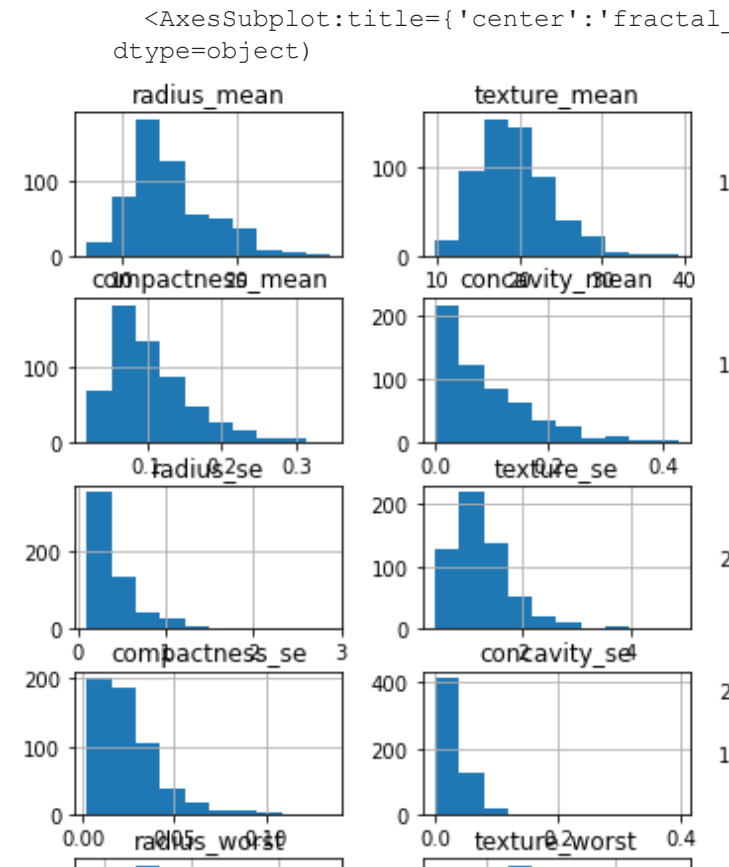
```
In [11]: tissue_pd.head()
```

[<AxesSubplot:title='(center:'compactness_mean')',
<AxesSubplot:title='(center:'concavity_mean')',
<AxesSubplot:title='(center:'concave points_mean')',
<AxesSubplot:title='(center:'symmetry_mean')',
<AxesSubplot:title='(center:'fractal_dimension_mean')'],
[<AxesSubplot:title='(center:'radius_se')',
<AxesSubplot:title='(center:'texture_se')',
<AxesSubplot:title='(center:'perimeter_se')',
<AxesSubplot:title='(center:'area_se')',
<AxesSubplot:title='(center:'smoothness_se')'],
<AxesSubplot:title='(center:'compactness_se')',

I also want to look at the class distribution of our target variable 'diagnosis', as it will help us understand if we need to account for class imbalance in our modeling.

```
In [12]: #looking at value counts of diagnosis variable
print(tissue_pd['diagnosis'].value_counts())
tissue_pd['diagnosis'].value_counts().sort_values().plot(kind = 'barh')
```

```
0      357
M      212
Name: diagnosis, dtype: int64
Out[12]:
```



Out of the 569 patients that were screened for breast cancer, 357 were diagnosed with benign masses, meaning they do not have breast cancer. The other 212 patients were found to have malignant or cancerous tumors.

In terms of class distribution, they are pretty equal in sample size, and I would say that the classes are not imbalanced.

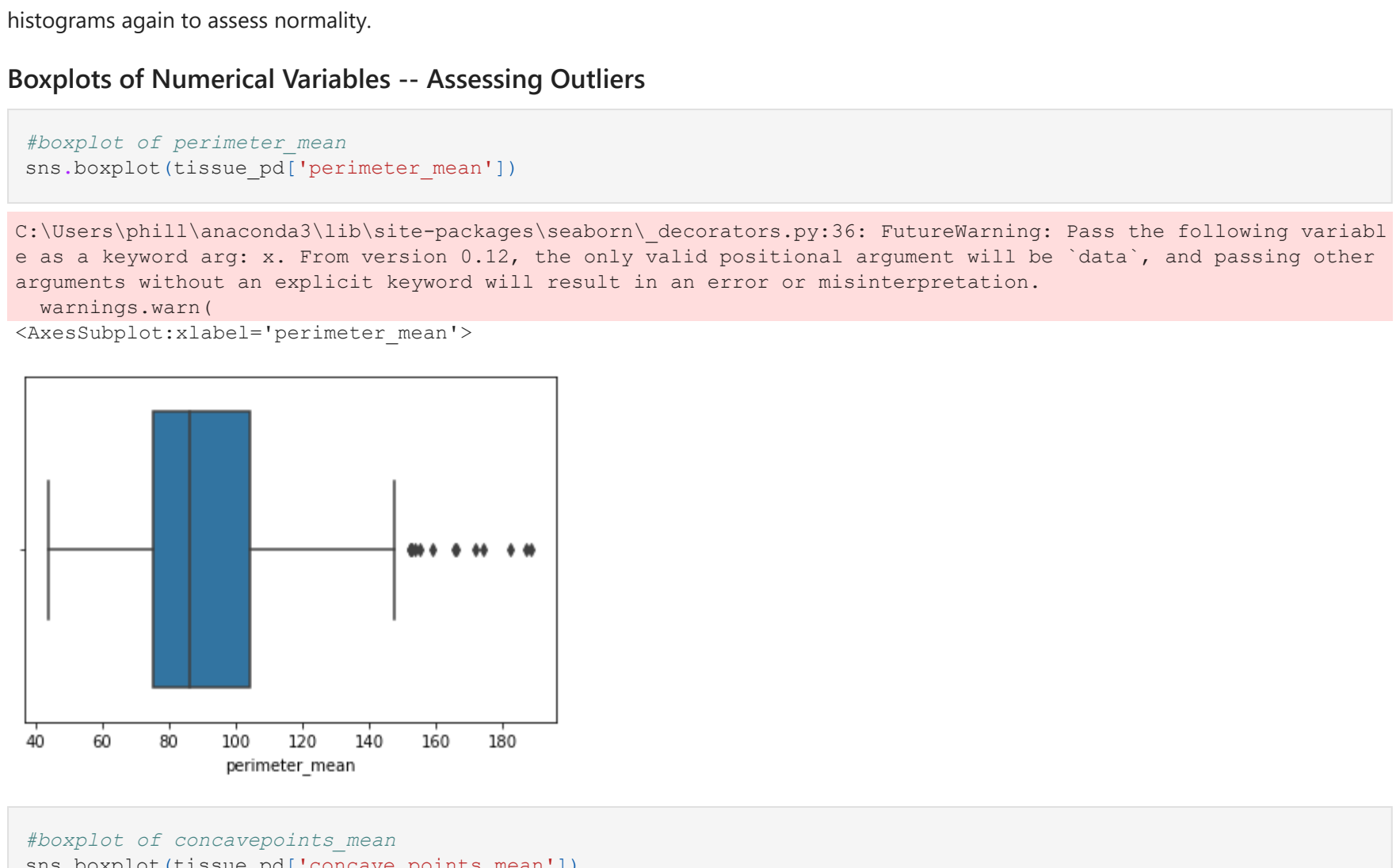
$$1/2 * 569 = 284.5$$

- Benign class is slightly above half
- Malignant class is slightly below half

## Data Visualizations

### Histograms of Numeric Variables

```
In [13]: tissue_pd.hist(bins=10, figsize=(15, 10))
```



Most of the distributions for the numerical variables in the dataset lean towards a positive skew (majority of the points are at the lower value range).

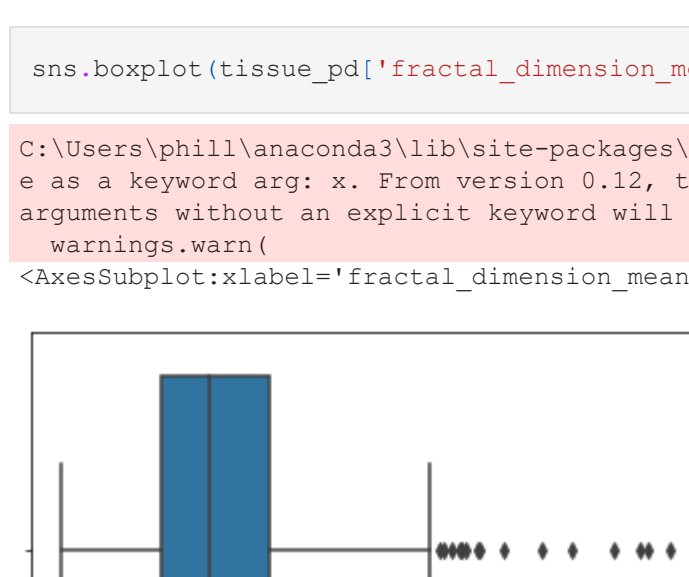
There are a few distributions that are relatively normal though: texture\_mean, smoothness\_mean, symmetry\_mean, texture\_worst and smoothness\_worst. They are mostly centered around the mean of the column values, and follow a relatively similar Gaussian distribution.

There are also some variables which may contain outliers, and when those are assessed and possibly removed, we can look at their boxplots again to assess normality.

### Boxplots of Numerical Variables -- Assessing Outliers

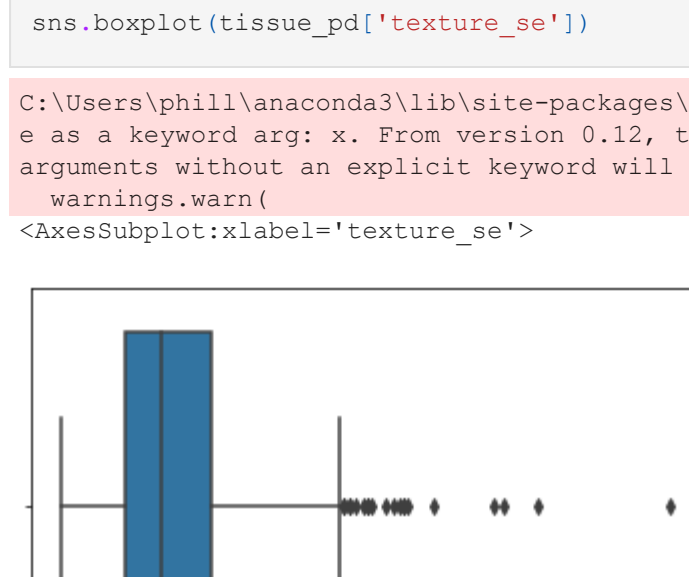
```
In [14]: #boxplot of perimeter_mean
sns.boxplot(tissue_pd['perimeter_mean'])

C:\Users\phill\anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable
as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other
arguments without an explicit keyword will result in an error or misinterpretation.
warnings.warn(
<AxesSubplot:xlabel='perimeter_mean'>
```



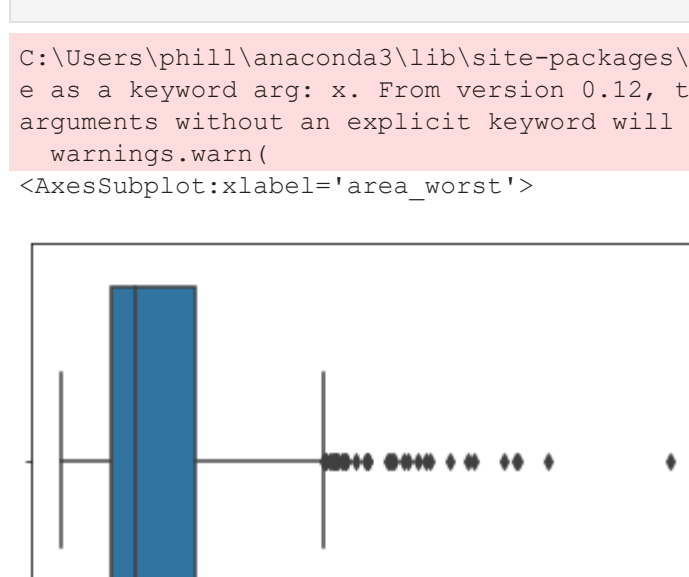
```
In [15]: #boxplot of concavepoints_mean
sns.boxplot(tissue_pd['concave points_mean'])

C:\Users\phill\anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable
as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other
arguments without an explicit keyword will result in an error or misinterpretation.
warnings.warn(
<AxesSubplot:xlabel='concave points_mean'>
```



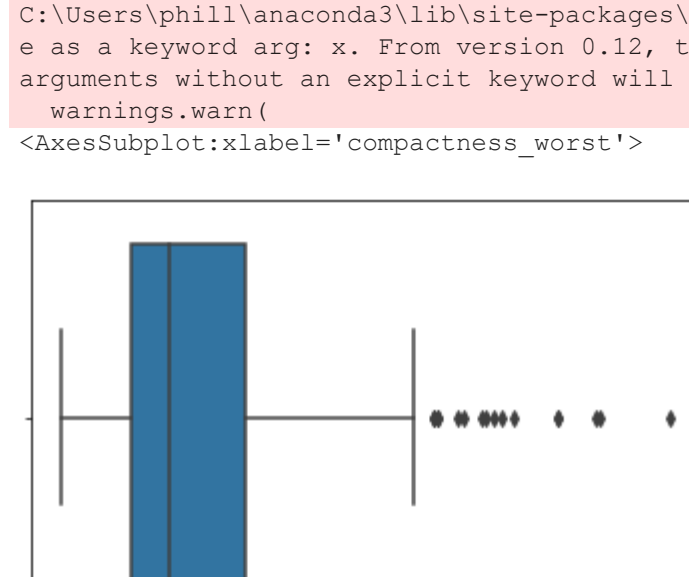
```
In [16]: #boxplot of fractal_dimension_mean
sns.boxplot(tissue_pd['fractal_dimension_mean'])

C:\Users\phill\anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable
as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other
arguments without an explicit keyword will result in an error or misinterpretation.
warnings.warn(
<AxesSubplot:xlabel='fractal_dimension_mean'>
```



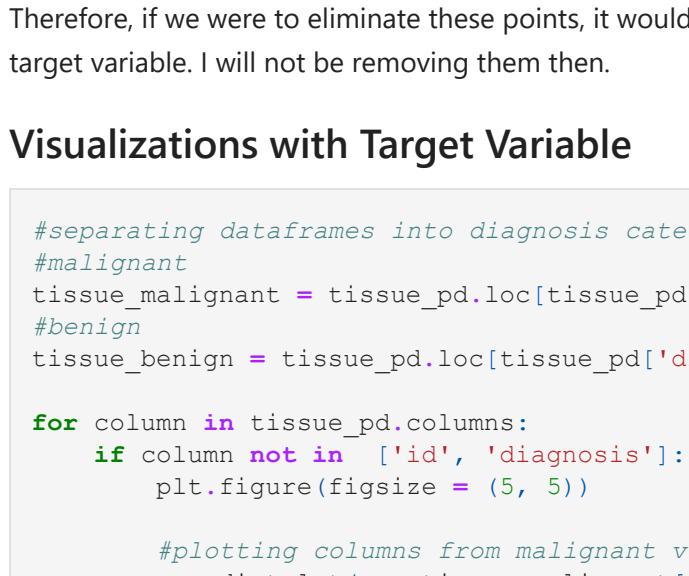
```
In [17]: sns.boxplot(tissue_pd['texture_se'])

C:\Users\phill\anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable
as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other
arguments without an explicit keyword will result in an error or misinterpretation.
warnings.warn(
<AxesSubplot:xlabel='texture_se'>
```



```
In [18]: sns.boxplot(tissue_pd['area_worst'])

C:\Users\phill\anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable
as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other
arguments without an explicit keyword will result in an error or misinterpretation.
warnings.warn(
<AxesSubplot:xlabel='area_worst'>
```



```
In [19]: sns.boxplot(tissue_pd['compactness_worst'])

C:\Users\phill\anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable
as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing other
arguments without an explicit keyword will result in an error or misinterpretation.
warnings.warn(
<AxesSubplot:xlabel='compactness_worst'>
```



From looking over the boxplots, it is clear that there are some outlier points which would need to be further analyzed. However, when I was reading the original research paper from where this data came from, it mentioned that higher and lower value points for the various measures were usually indicative of malignancy in patients' breast masses.

Therefore, if we were to eliminate these points, it could disrupt our data and possibly minimize the number of malignant diagnoses in our target variable. I will not be removing them then.

## Visualizations with Target Variable

```
In [20]: #separating dataframes into diagnosis categories
#malignant
tissue_malignant = tissue_pd.loc[tissue_pd['diagnosis'] == 'M']
#benign
tissue_benign = tissue_pd.loc[tissue_pd['diagnosis'] == 'B']

for column in tissue_pd.columns:
    if column not in ['id', 'diagnosis']:
        plt.figure(figsize = (5, 5))

        #plotting columns from malignant vs. benign as distplot (distribution)
        sns.distplot(a = tissue_malignant[column], hist = False, color = 'red')
        sns.distplot(a = tissue_benign[column], hist = False, color = 'blue')

        plt.title(column + " Benign (blue) vs. Malignant (red)")
        plt.show()
```



For the most part, the distributions for the numerical variables between the diagnosis categories are relatively equal, i.e. they overlap a great deal. There are some cases where there is a more drastic difference in the skew and distribution means, but there is nothing concerning the need to feature-engineer.

Encoding for Target Variable

```

diagnosis = ['diagnosis']
cat_cols = ['diagnosis']

#turning categories into their numerical counterparts using LabelEncoder
for var in cat_cols:
    number = LabelEncoder()
    tissue_pd[var+ '_cat'] = number.fit_transform(tissue_pd[var].astype('str'))

tissue_pd[['diagnosis','diagnosis_cat']]

```

```
number = LabelEncoder()
tissue_pd[vars*cat"] = number.fit_transform(tissue_pd[vars].astype('str'))

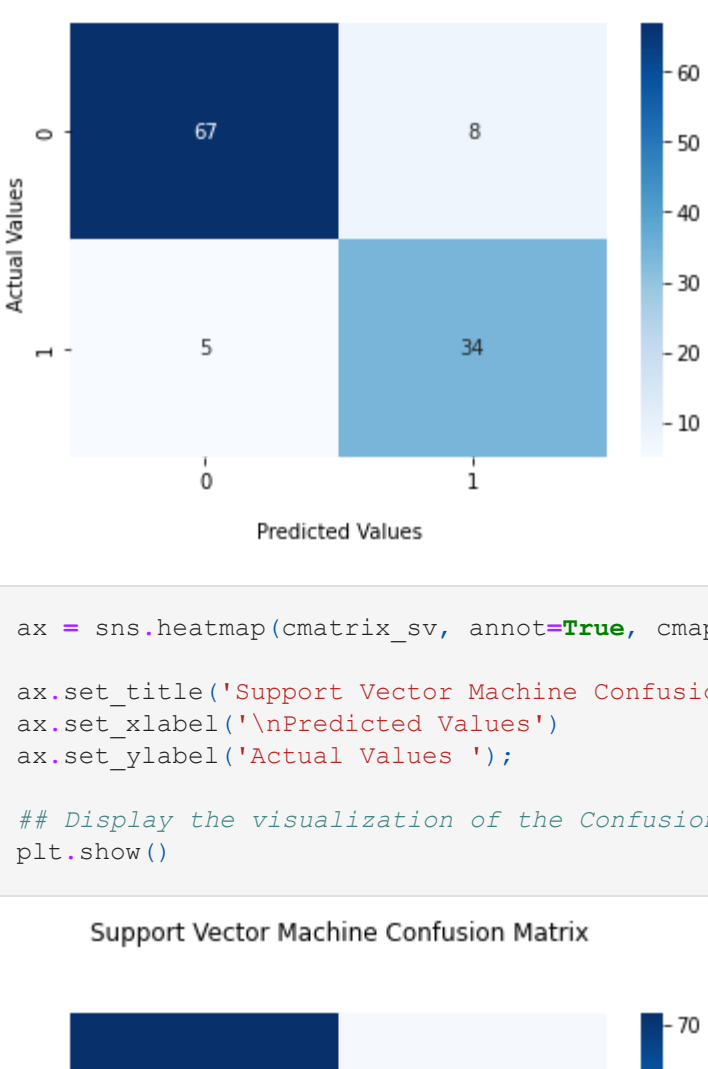
tissue_pd[['diagnosis', 'diagnosiscat']]
```





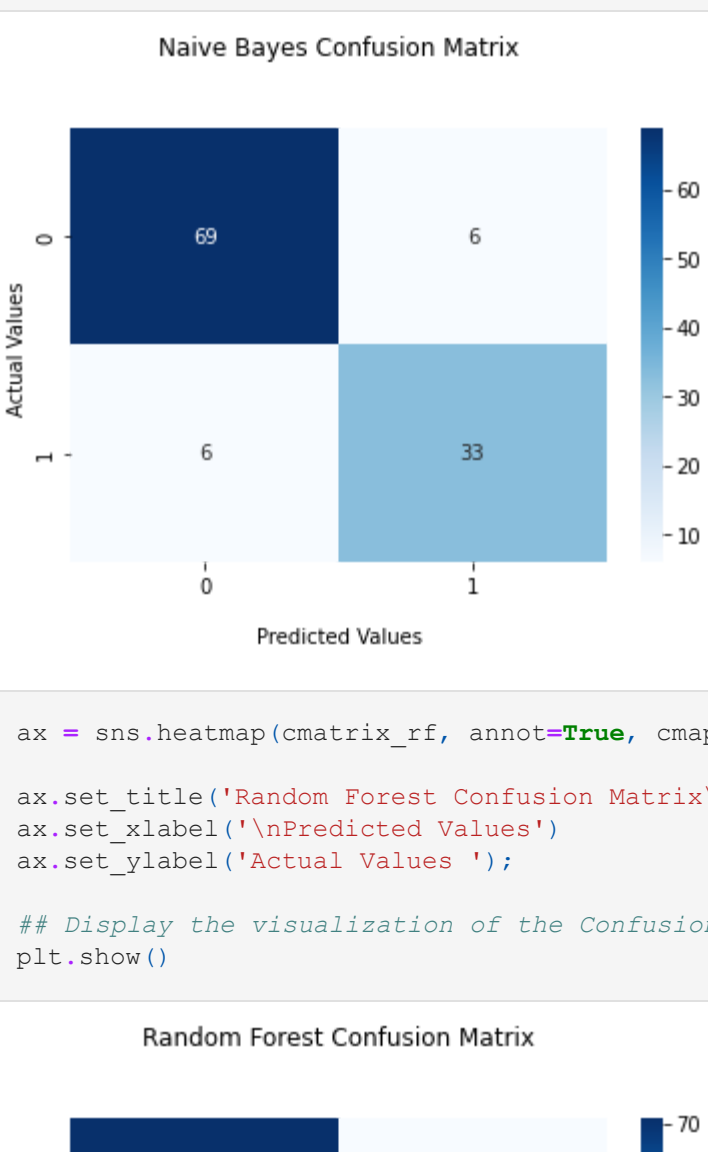


Decision Tree Classifier Confusion Matrix



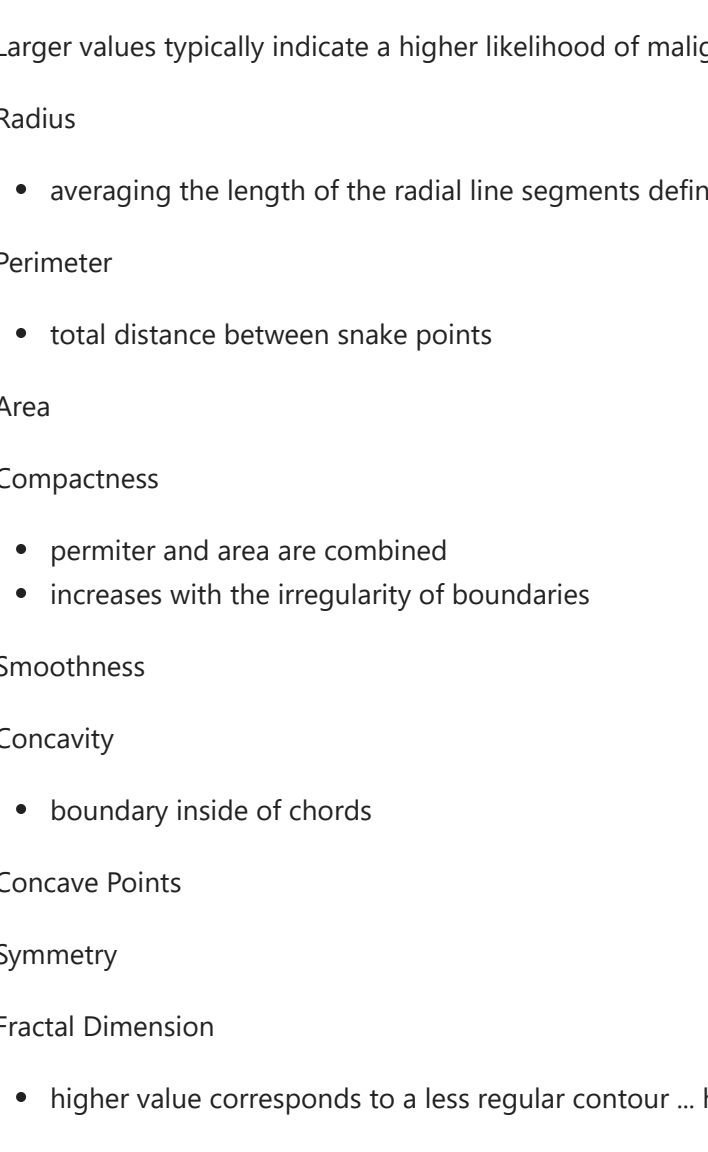
```
In [69]: ax = sns.heatmap(matrix_sv, annot=True, cmap='Blues')
ax.set_title('Support Vector Machine Confusion Matrix\n\n');
ax.set_xlabel('\nPredicted Values');
ax.set_ylabel('Actual Values ');
# Display the visualization of the Confusion Matrix.
plt.show()
```

Support Vector Machine Confusion Matrix



```
In [68]: ax = sns.heatmap(matrix_gnb, annot=True, cmap='Blues')
ax.set_title('Naive Bayes Confusion Matrix\n\n');
ax.set_xlabel('\nPredicted Values');
ax.set_ylabel('Actual Values ');
# Display the visualization of the Confusion Matrix.
plt.show()
```

Naive Bayes Confusion Matrix



```
In [70]: ax = sns.heatmap(matrix_rf, annot=True, cmap='Blues')
ax.set_title('Random Forest Confusion Matrix\n\n');
ax.set_xlabel('\nPredicted Values');
ax.set_ylabel('Actual Values ');
# Display the visualization of the Confusion Matrix.
plt.show()
```

Random Forest Confusion Matrix



Paper Information

Curvature

- deviations from a circle are problematic for malignancy
- can lead to quickly-dividing cells which usually points to cancer
- high or low pints are concerning

Larger values typically indicate a higher likelihood of malignancy

Radius

- averaging the length of the radial line segments defined by the centroid of the snake and the individual snake points

Perimeter

- total distance between snake points

Area

Compactness

- perimeter and area are combined
- increases with the irregularity of boundaries

Smoothness

Concavity

- boundary inside of chords

Concave Points

Symmetry

Fractal Dimension

- higher value corresponds to a less regular contour ... higher probability of malignancy

Texture

- variance of the gray scale intensities

Extreme values are the most intuitively useful for the problem at hand, since only a few malignant cells may occur in a given sample