

Name : Aman Bashar

Contact: Amanbashar96@gmail.com

Mobile: 9085429899

Title: Python Exploratory Data analysis

---

**Answer the following questions below and upload them in the google form.**

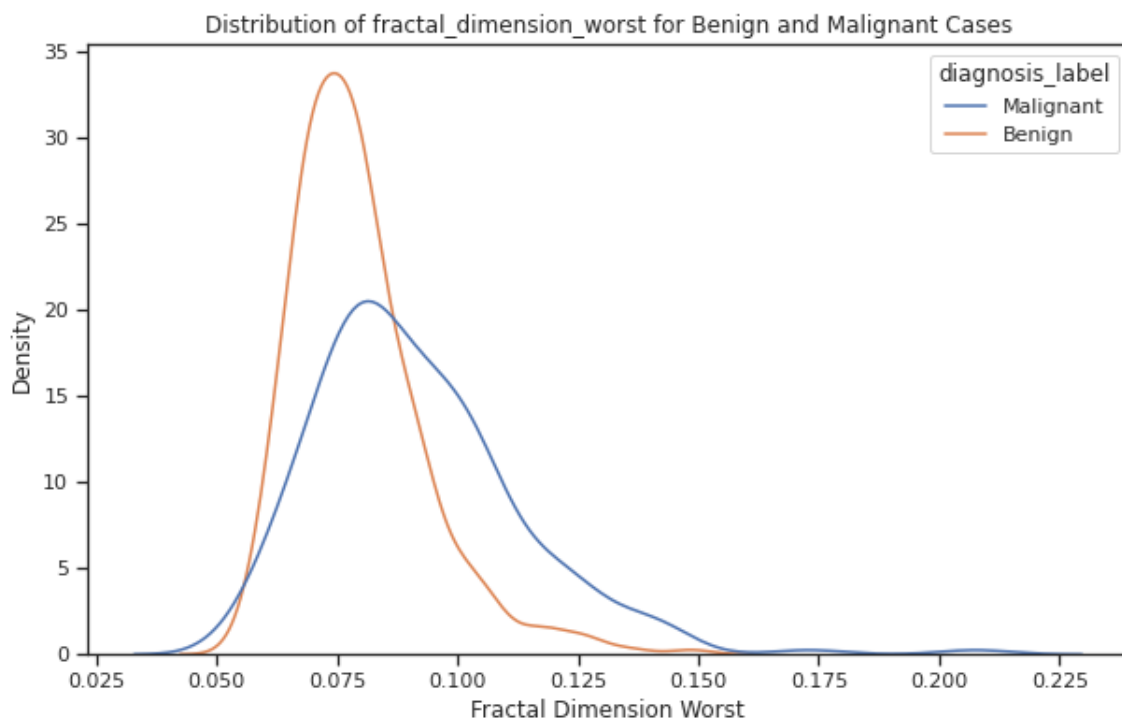
1. How does the distribution of feature “fractal\_dimension\_worst” differ between benign and malignant cases?
2. What is the range of values for the feature “radius\_mean” and how skewed is its distribution?
3. Are there any outliers in feature “area\_mean” and how might they affect analysis?
4. Based on the EDA, what factors seem to be most relevant to predicting breast cancer diagnosis?
5. What limitations are there in the data, and how might they affect our conclusions?

Answer 1.

By Investigating the distribution of the feature 'fractal\_dimension\_worst' between benign and malignant cases reveals these insights:-

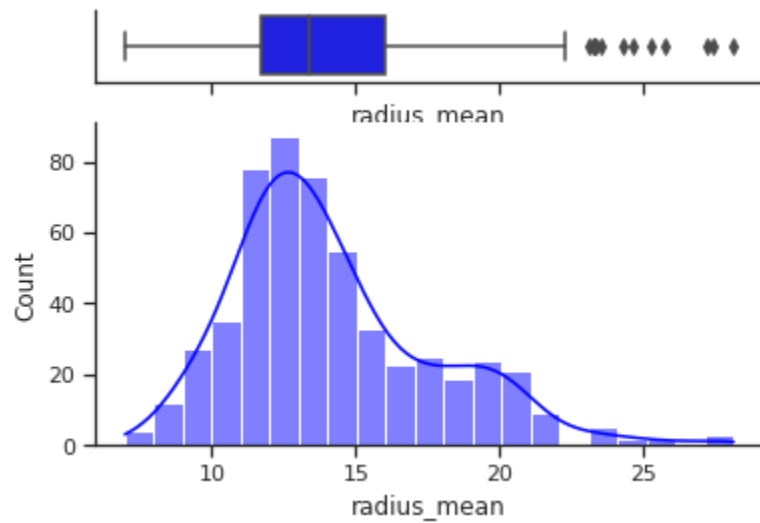
Benign cases tend to concentrate around a certain range of fractal dimension values, while malignant cases exhibit a broader distribution with higher density in some regions.

This suggests that 'fractal\_dimension\_worst' may be a discriminative feature for distinguishing between benign and malignant tumors, that is it can be a potential indicator for identifying type of tumor.

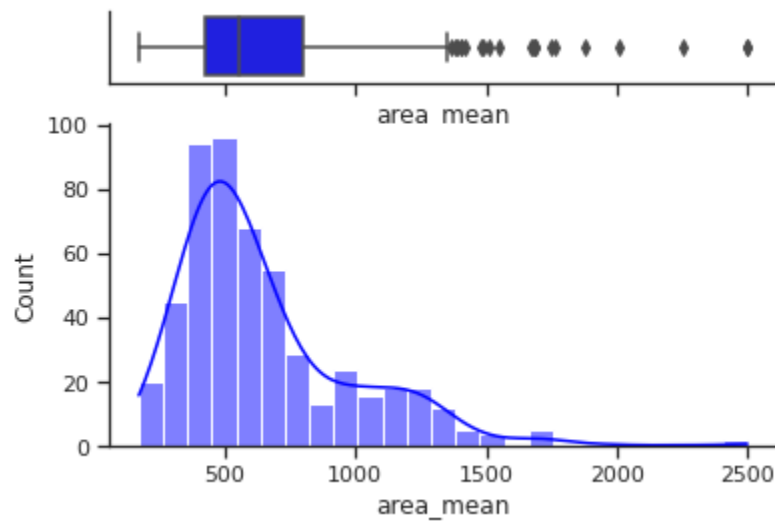


## Answer 2:

The range of values of the feature 'radius\_mean' is 6.981 to 28.11 and Skewness of radius\_mean is 0.998492, which is positively skewed.



Answer 3:



Yes, there are outliers present in the feature 'area\_mean' as we can see it in the box plot. We also calculated that there are 39 outliers present in 'area\_mean'.

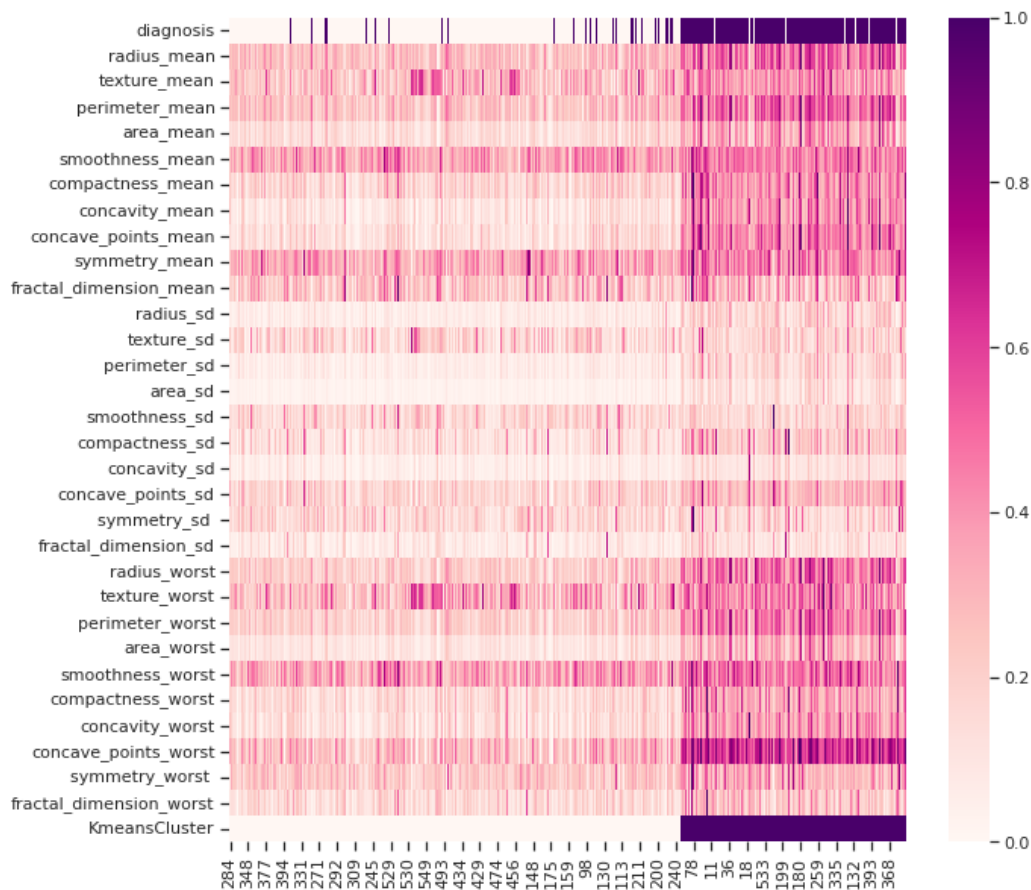
Presence of outliers in dataset can have several impacts on the analysis, they are :-

- Outliers can impact summary statistics such as mean and Standard deviation.
- Outliers may violate the assumptions of certain statistical tests, leading to inaccurate results.
- Outliers can distort the shape and spread of data distributions, affecting the interpretation of patterns and relationships.

## Answer 4:

Based on the correlation analysis of the data-set through heat-map, we can conclude that following factors seem to be most relevant to predicting breast cancer diagnosis :-

- radius\_mean
- perimeter\_mean
- smoothness\_mean
- symmetry\_mean
- fractal\_dimension\_mean
- radius\_worst
- texture\_worst
- Smoothness\_worst
- Concave\_points\_worst



**Answer 5 :**

**limitation in the data and their impact on Conclusions:-**

**1. Imbalanced Data:-**

we need to check if the dataset has an imbalance in the distribution of benign and malignant cases.

imbalance data can bias the model towards the majority class, potentially affecting the accuracy of predictions.

**2. Missing Values:**

We need to check if there's any missing values in the dataset.

This missing values can lead to biased or incomplete analysis.

**3. Data Source and Bias-ness:-**

Biases in data collection may introduce inaccuracies or limit the generalization of findings.

