# Questions 1
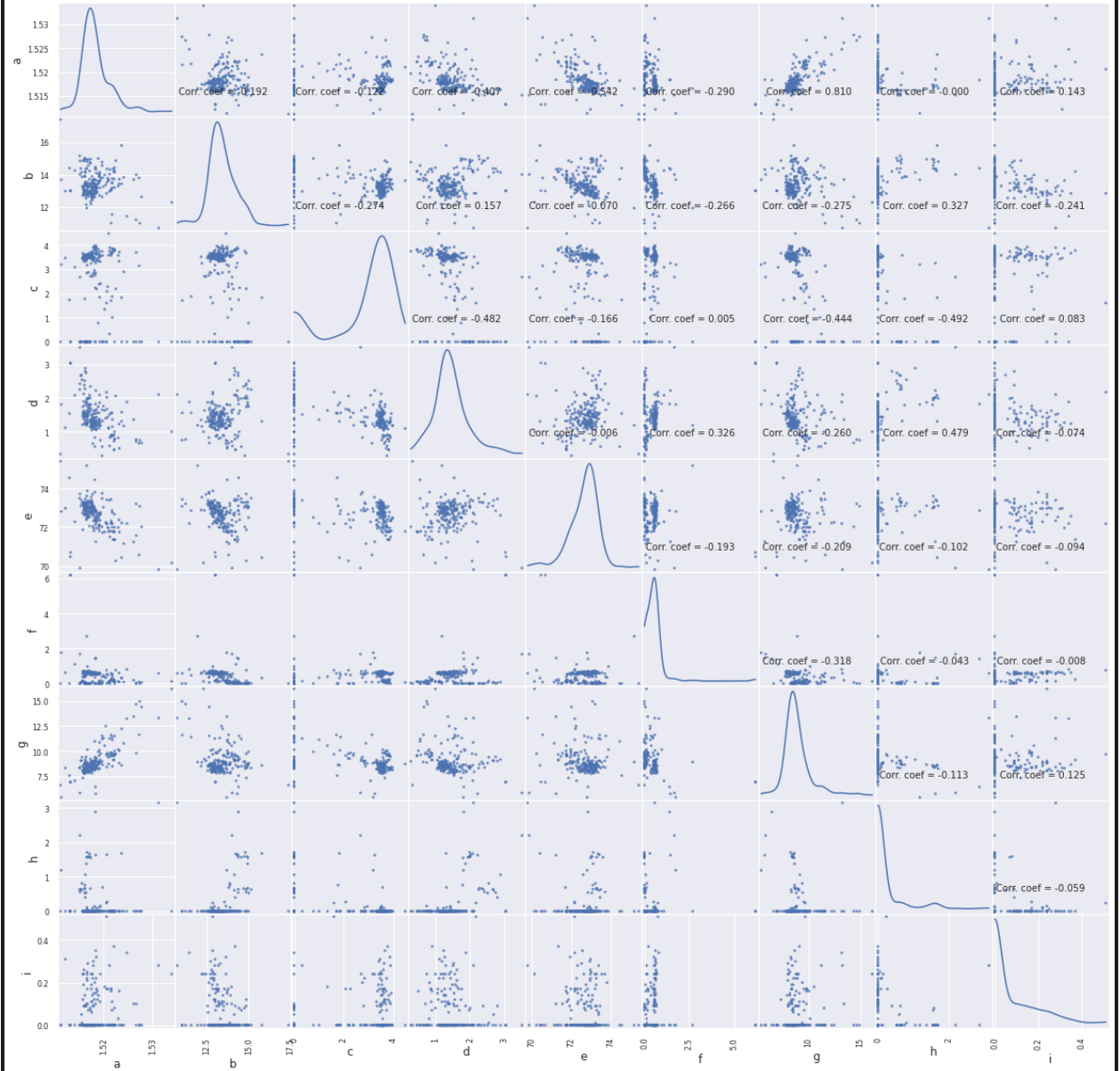
* All the features are in numerical format.

* There are no missing values.

* Mean and median are quite close to each other.

*There are no such outliners

* Not variance that much except for two columns c and g.

By the correlation method  "Pearson",  All the features are normally distributed except for f, h, and i.

strongly positive correlation with

*  a & g
*  b & h
*  d & h
*  d & f

```
Below graphs are Distribution of Scatter and density
plots
```
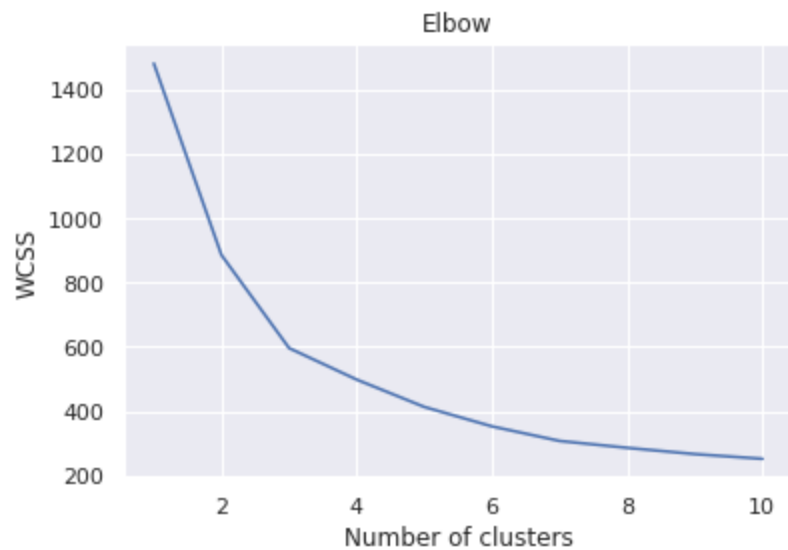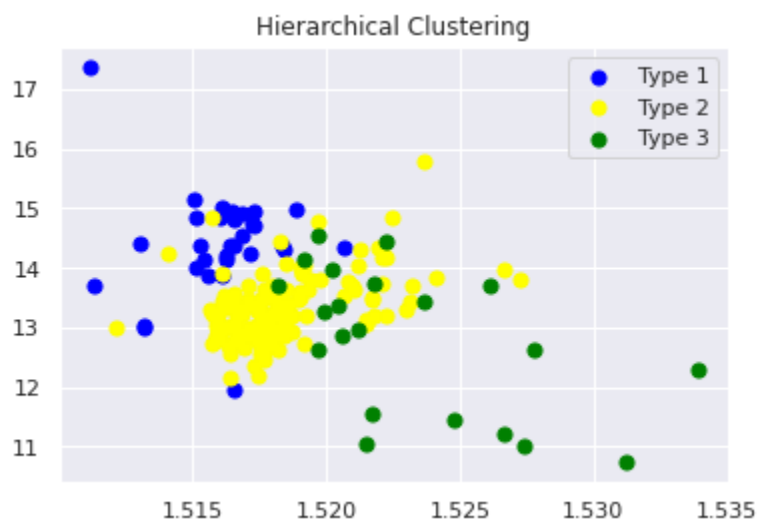
Scatter and Density Plot

# K-means clustering

In below graph, "Elbow" method mentions we should choose 3 clusters

Elbow

- Below image is hierarchical clustering clusters



Hierarchical Clustering

# Questions 2 (Oil palm data)

As a pairwise correlation, against FFB_yield, one Vs all correlation has been used to study all external factors correlation. So, We can reject the hypothesis that the two variables are not correlated if the p-value is below 0.05, generally. So we could mention, that there is a significant correlation between all the variables against FFB_yield.
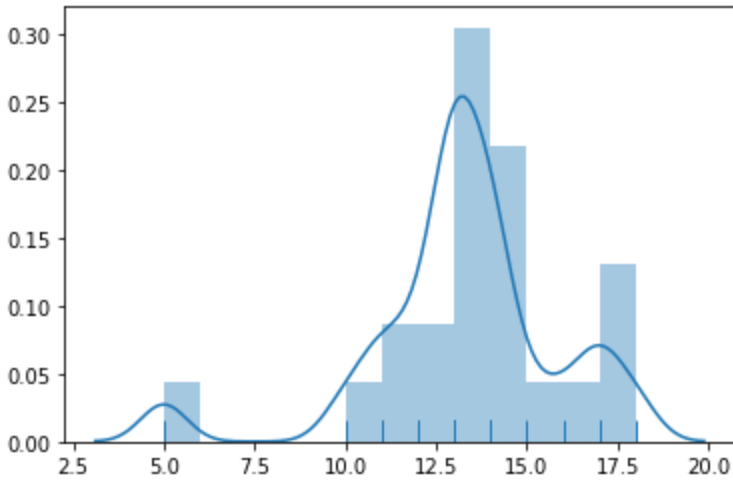
**FFB_Yield has**

•weak positive correlation with Max_Temp and Min_Temp

•strong negative correlation with HA_Harvested

•strong positive correlation with Precipitation

•strong moderate correlation with Working days.

 - Soil Moisture is negatively associated with FFB_Yield.
 - Temperature variables show no significant association with FFB_Yield.

# Questions 3(NLP- paragraph)

- Probability of "data" appearing in every line id 0.782608695652174

- The number of times "data analytics" appear together is 6

- The number of times only "analytics" appear in the complete text is 10

**Below the graph is the distribution of the distinct word counts in every line -**

The last question of Question3 is the Probability of "analytics" appearing after "data".

- Its Probability is **0.6**