

Assignment #2

K-means Clustering & Principal Component Analysis

[illegible]

Question 1

To ensure my test set and my training set used equivalent classes of “Cultivar” I used stratified random sampling with the strata classified by the “Cultivar” types (1, 2, and 3)

```
for (i in 1:3) {  
  print(paste("Percentage of Cultivar type",as.integer(i),"in Training  
Set",(length(which(d.wine.trainset$Cultivar==i))/nrow(d.wine.trainset))*100), quote=F)  
}  
for (i in 1:3) {  
  print(paste("Percentage of Cultivar type",as.integer(i),"in Test  
Set:",(length(which(d.wine.testset$Cultivar==i))/nrow(d.wine.testset))*100), quote=F)  
}
```

Out:

[1] Percentage of Cultivar type 1 in Training Set 33.0508474576271

[1] Percentage of Cultivar type 2 in Training Set 39.8305084745763

[1] Percentage of Cultivar type 3 in Training Set 27.1186440677966

[1] Percentage of Cultivar type 1 in Test Set: 33.3333333333333

[1] Percentage of Cultivar type 2 in Test Set: 40

[1] Percentage of Cultivar type 3 in Test Set: 26.6666666666667

As you can see the representation of cultivar is fairly constant across both sets.

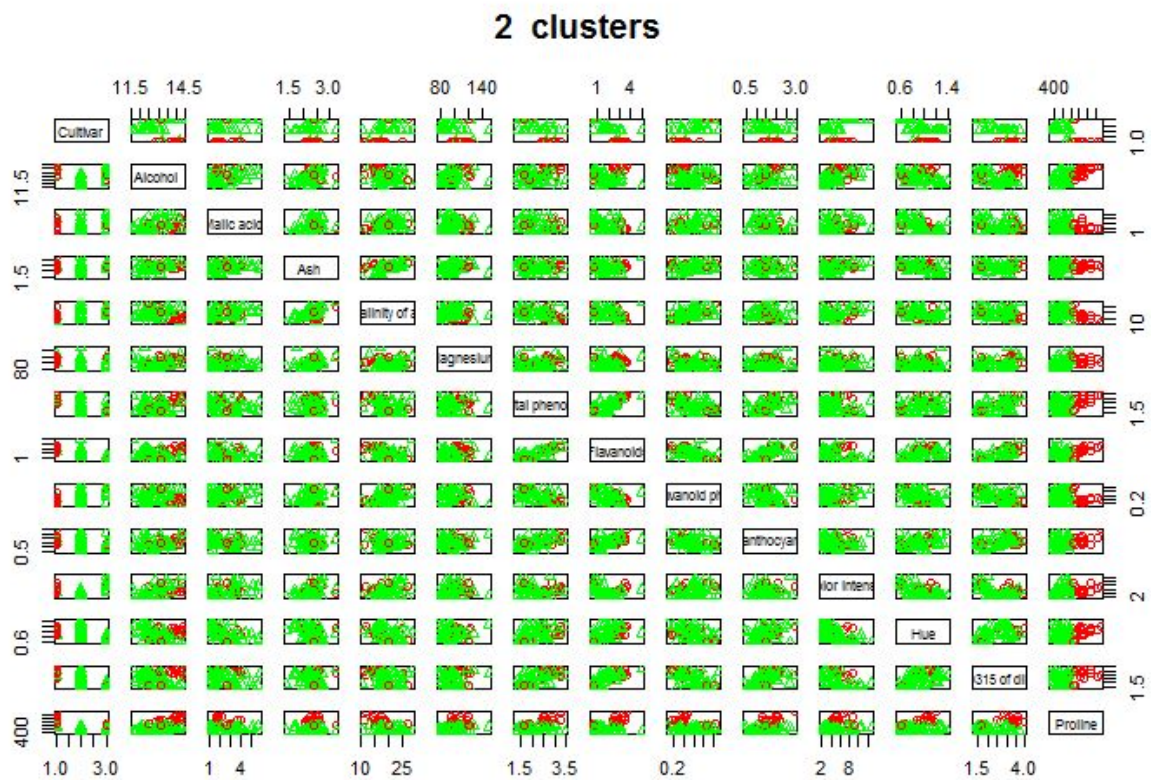
Question 2

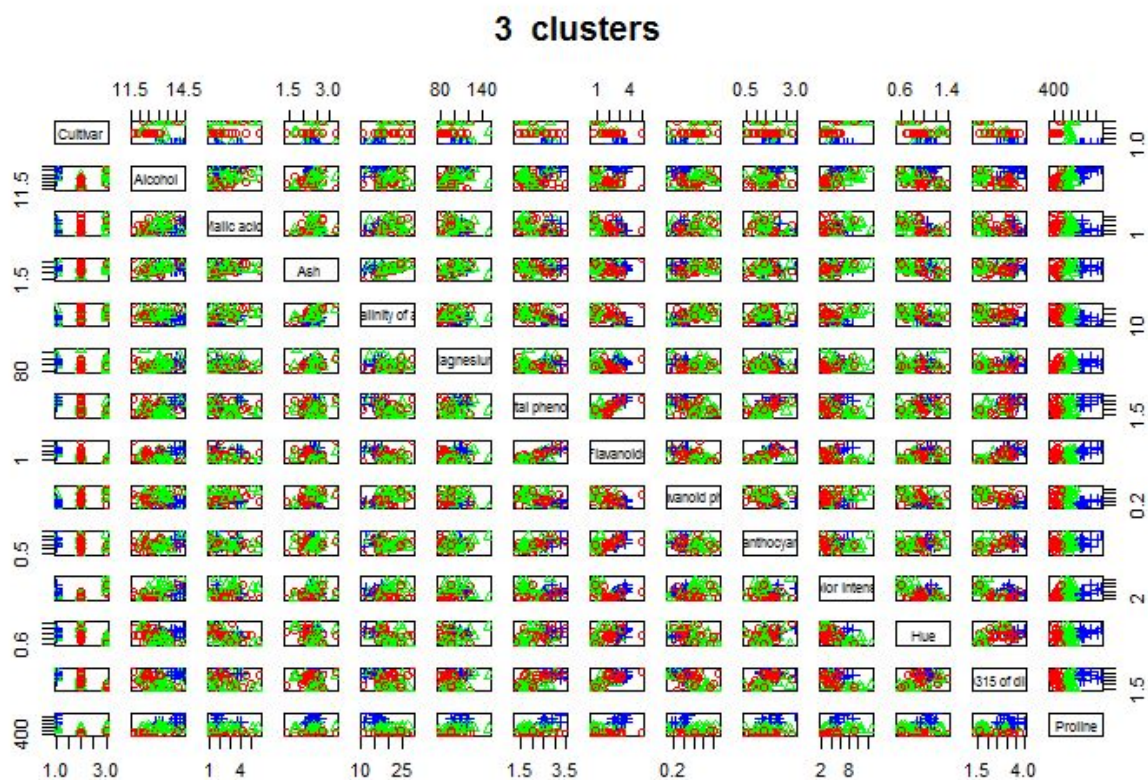
Here is output from a k-means clustering algorithm run 13 times with SSE and Davies-Bouldin averages applied.

Each runthrough was done 10 times with 10 different seeds applied.

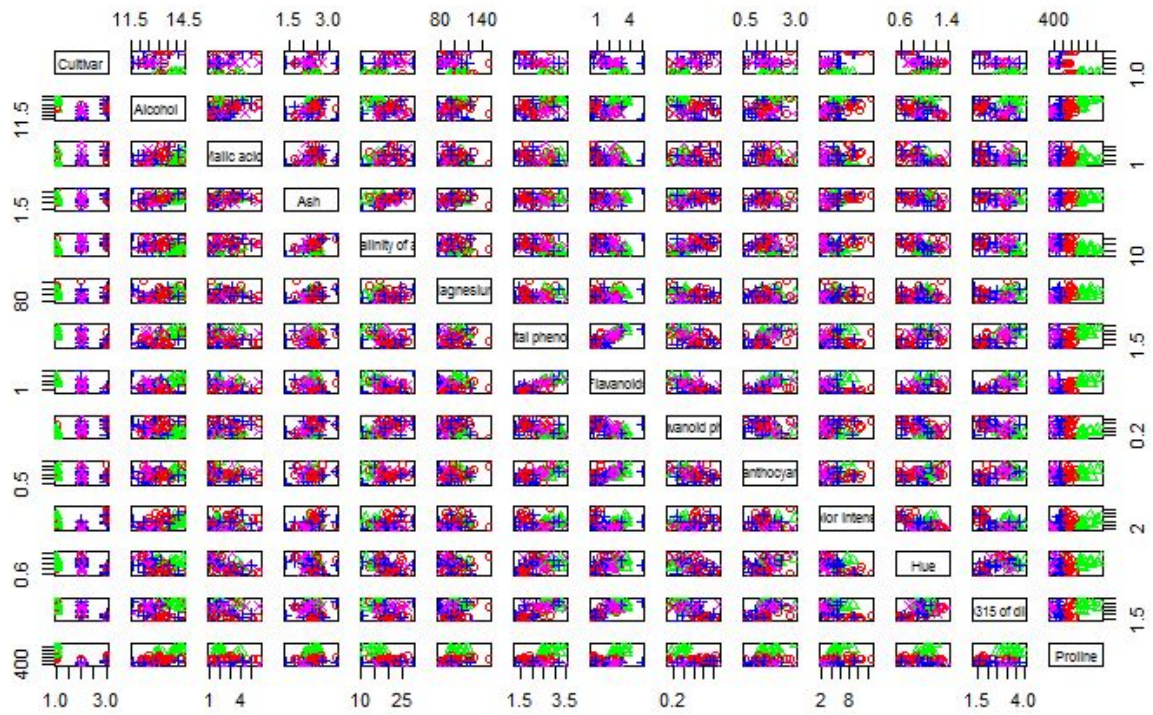
Test Set:

Seeds: 1 - 8

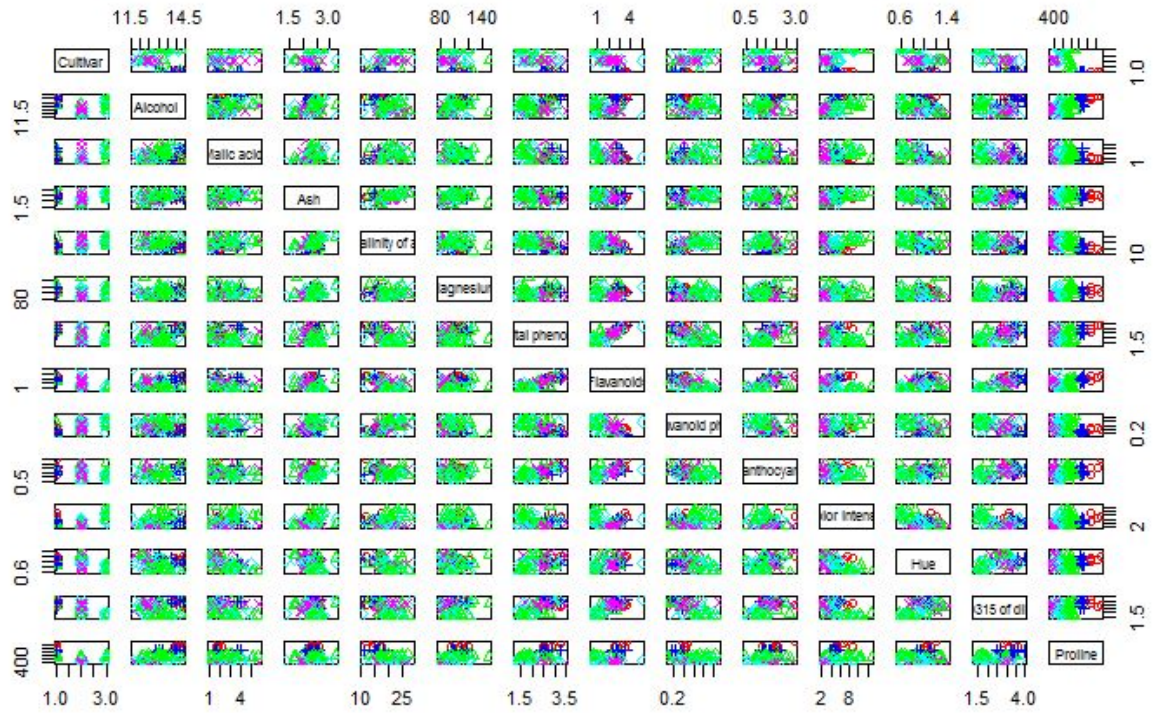




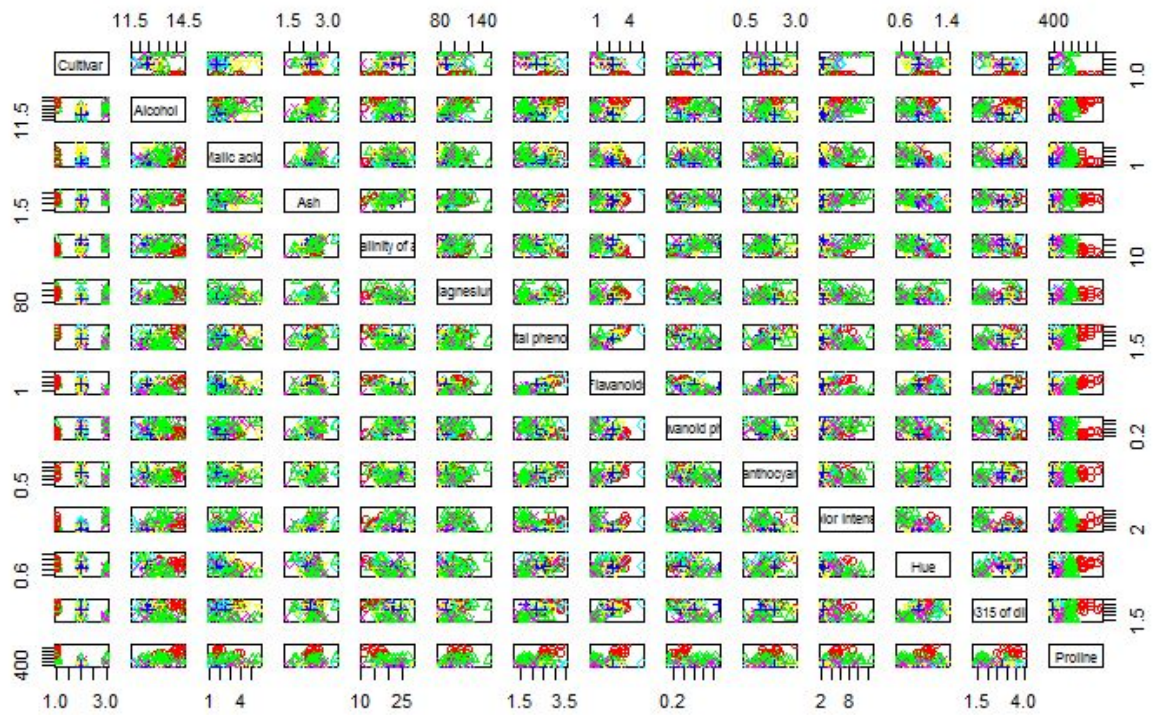
4 clusters



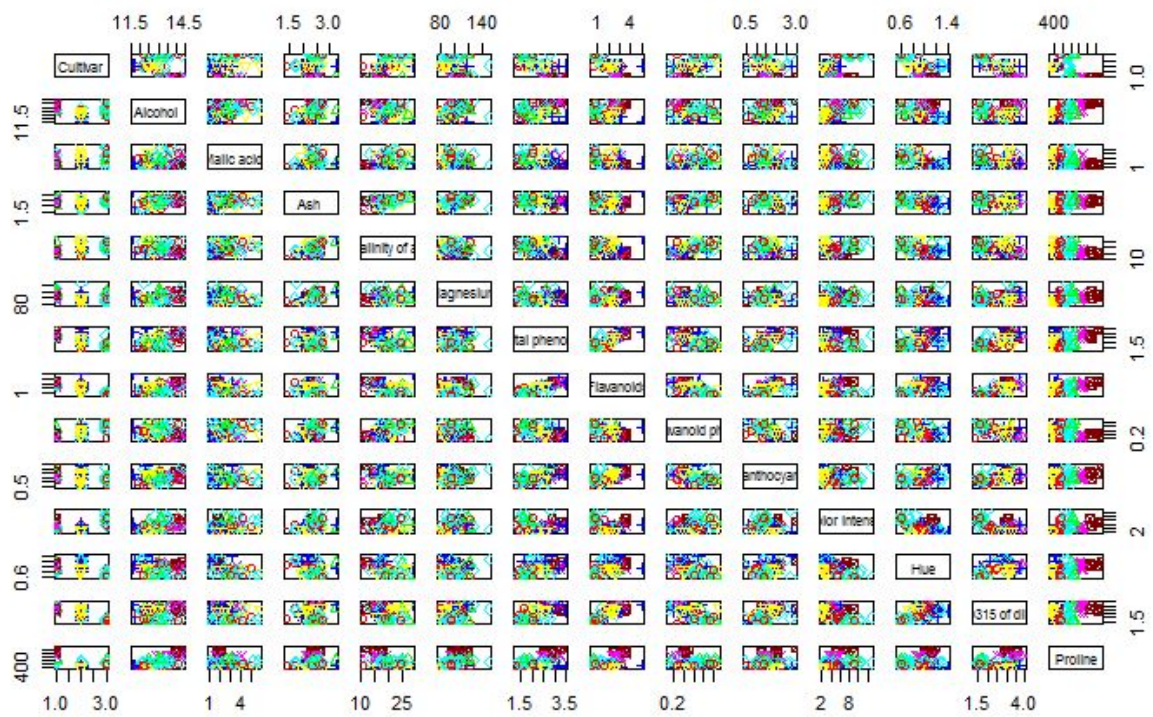
5 clusters



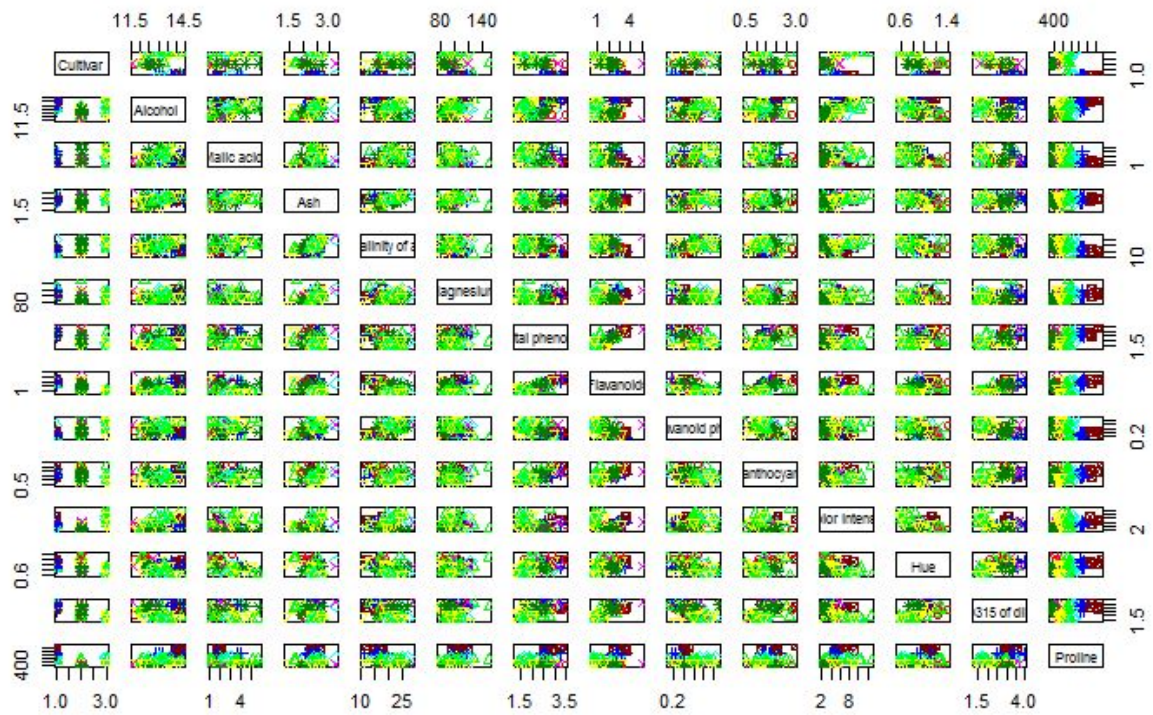
6 clusters



7 clusters



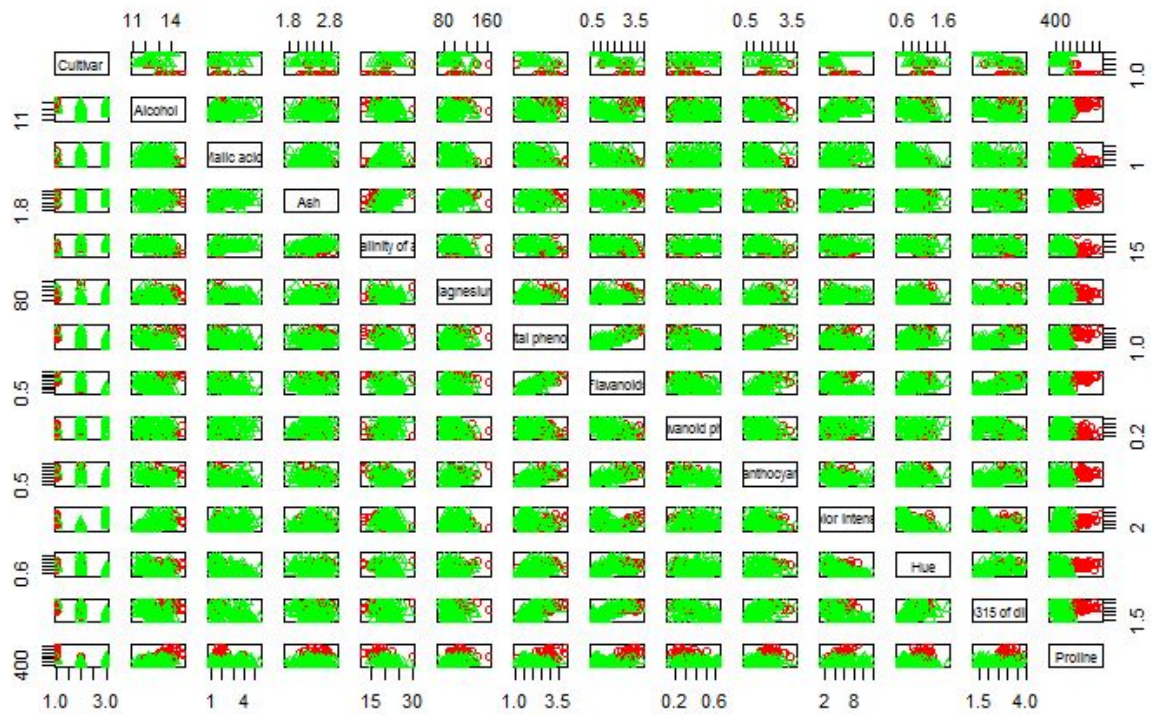
8 clusters



Training Set:

Seeds: 1 - 8

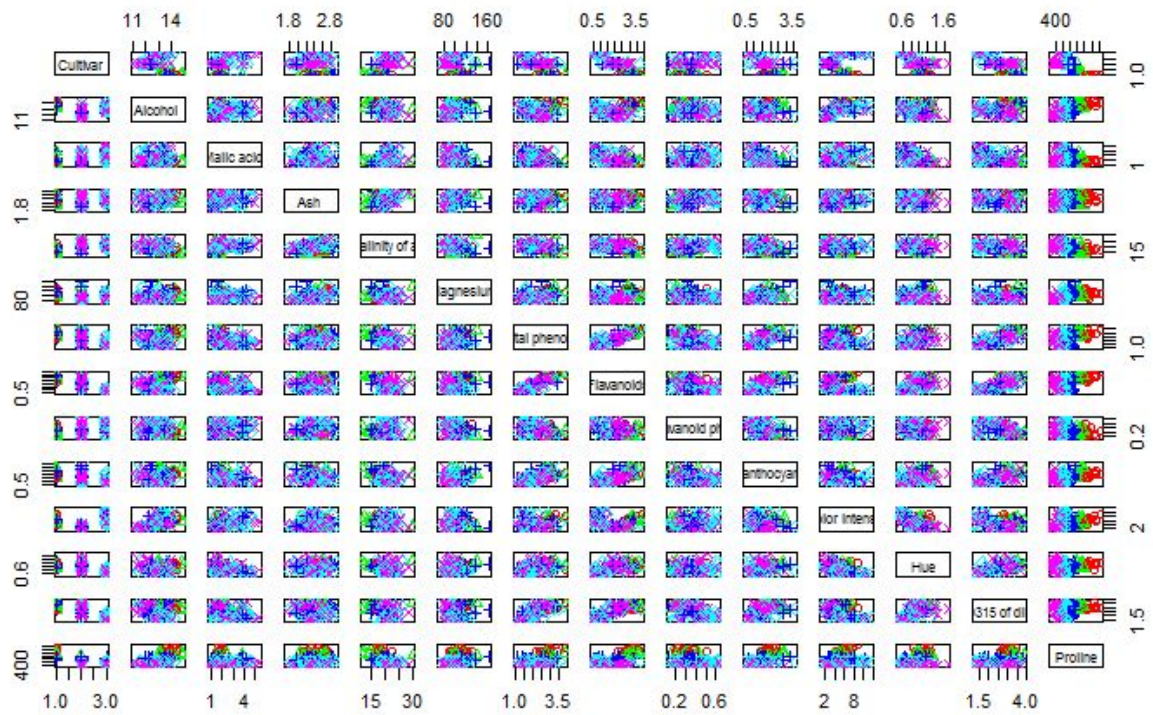
2 clusters





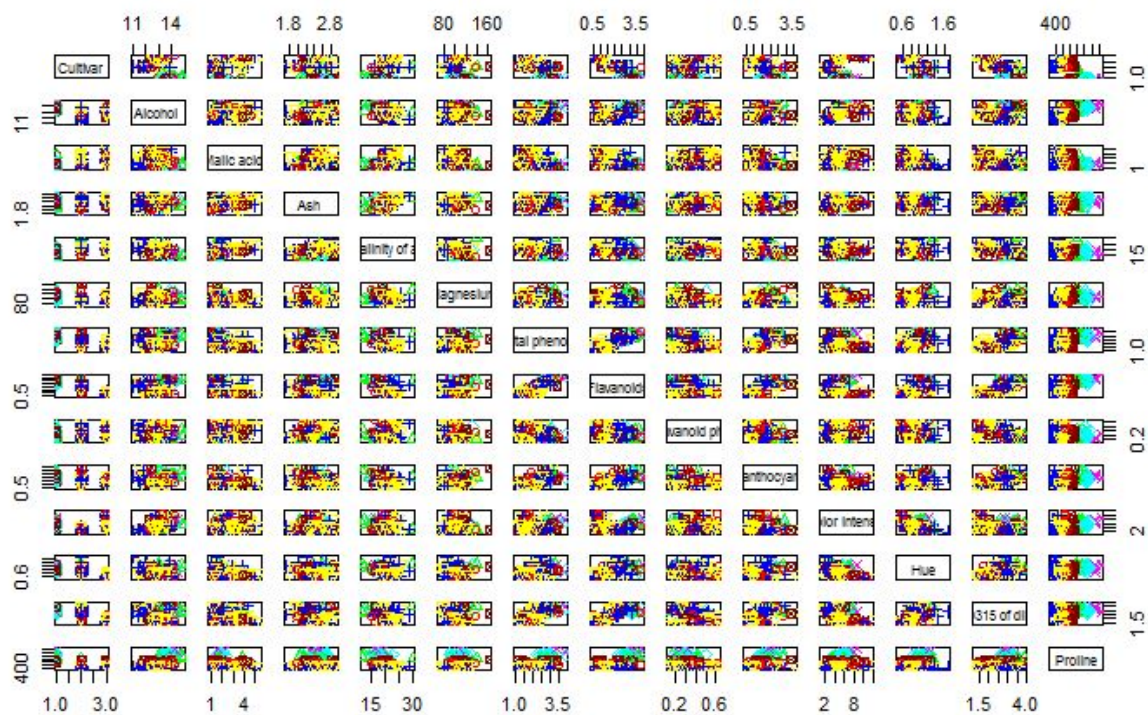


5 clusters

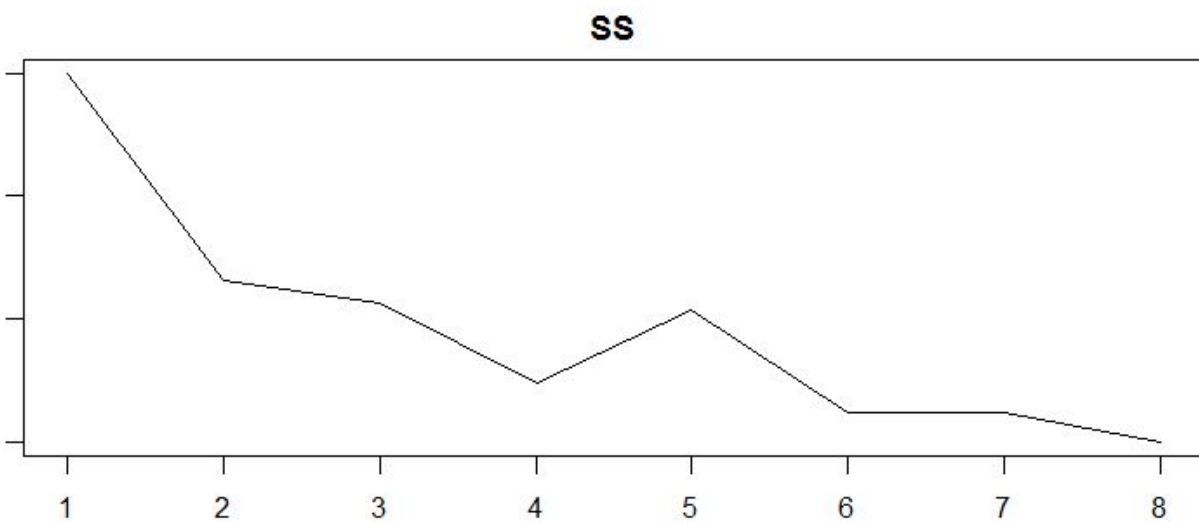




7 clusters



Here is the standard error plot, the elbow shows 4 clusters for the test data.



Here is the standard error plot for the training data.

