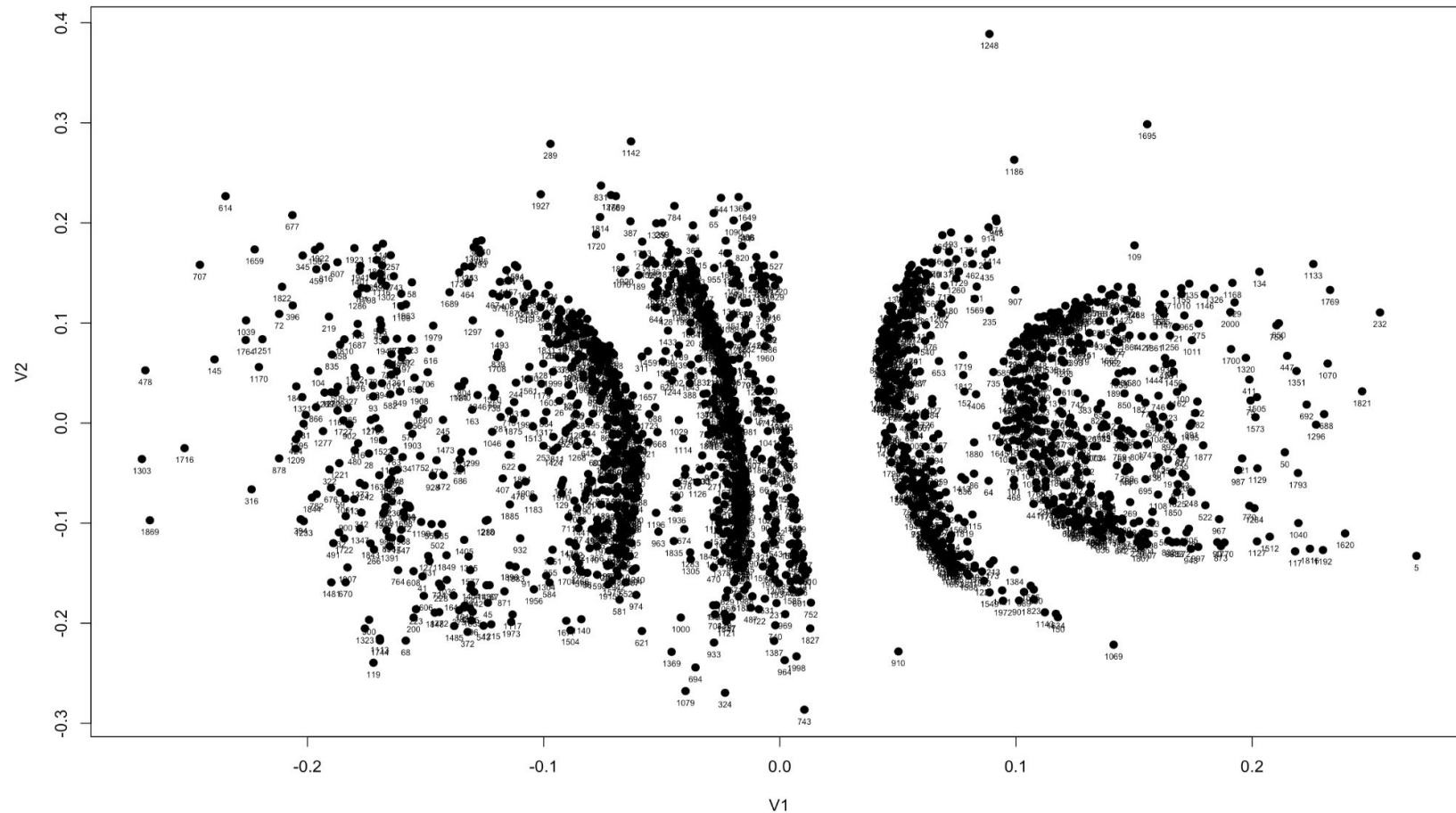
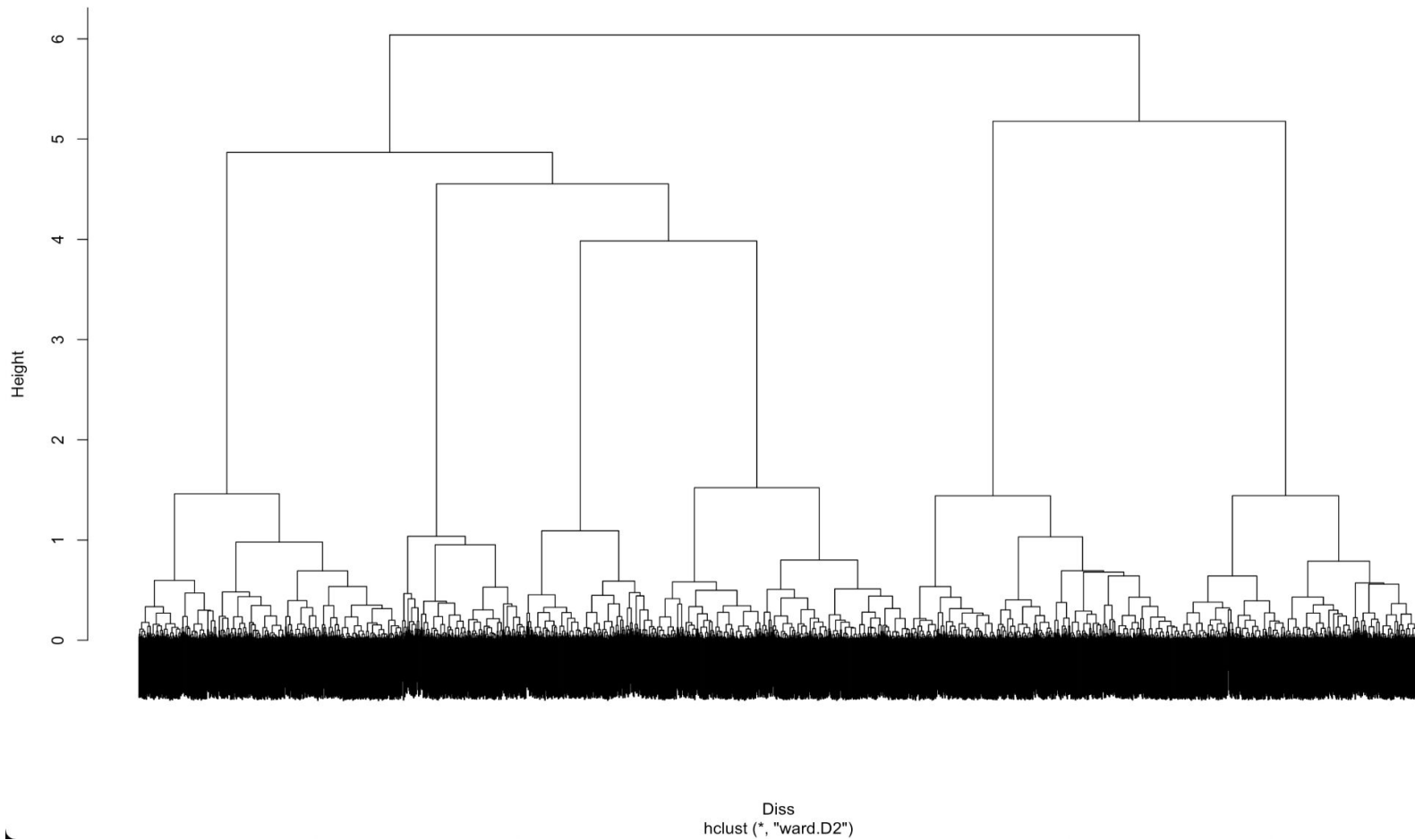


# Predict Earnings

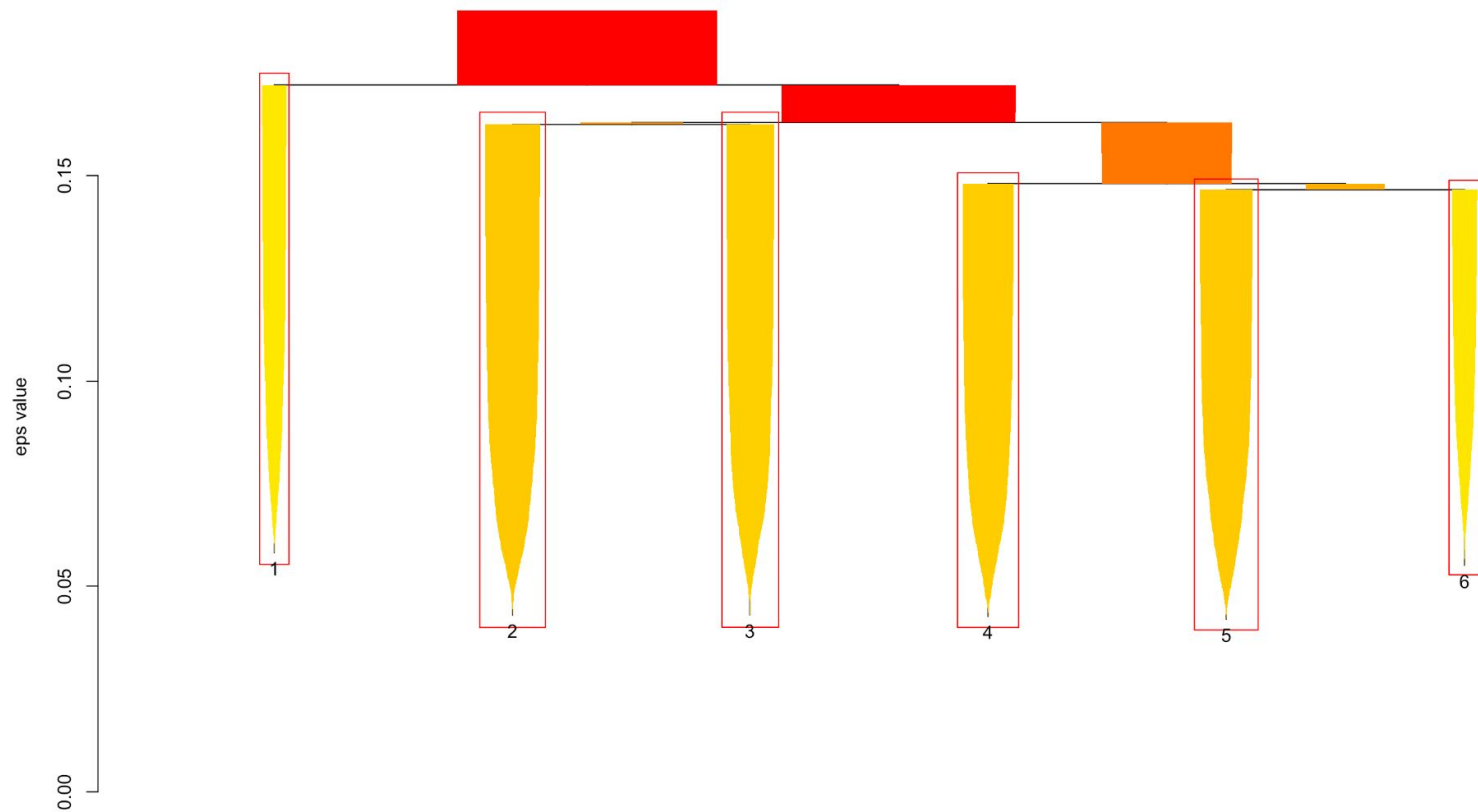
## V1



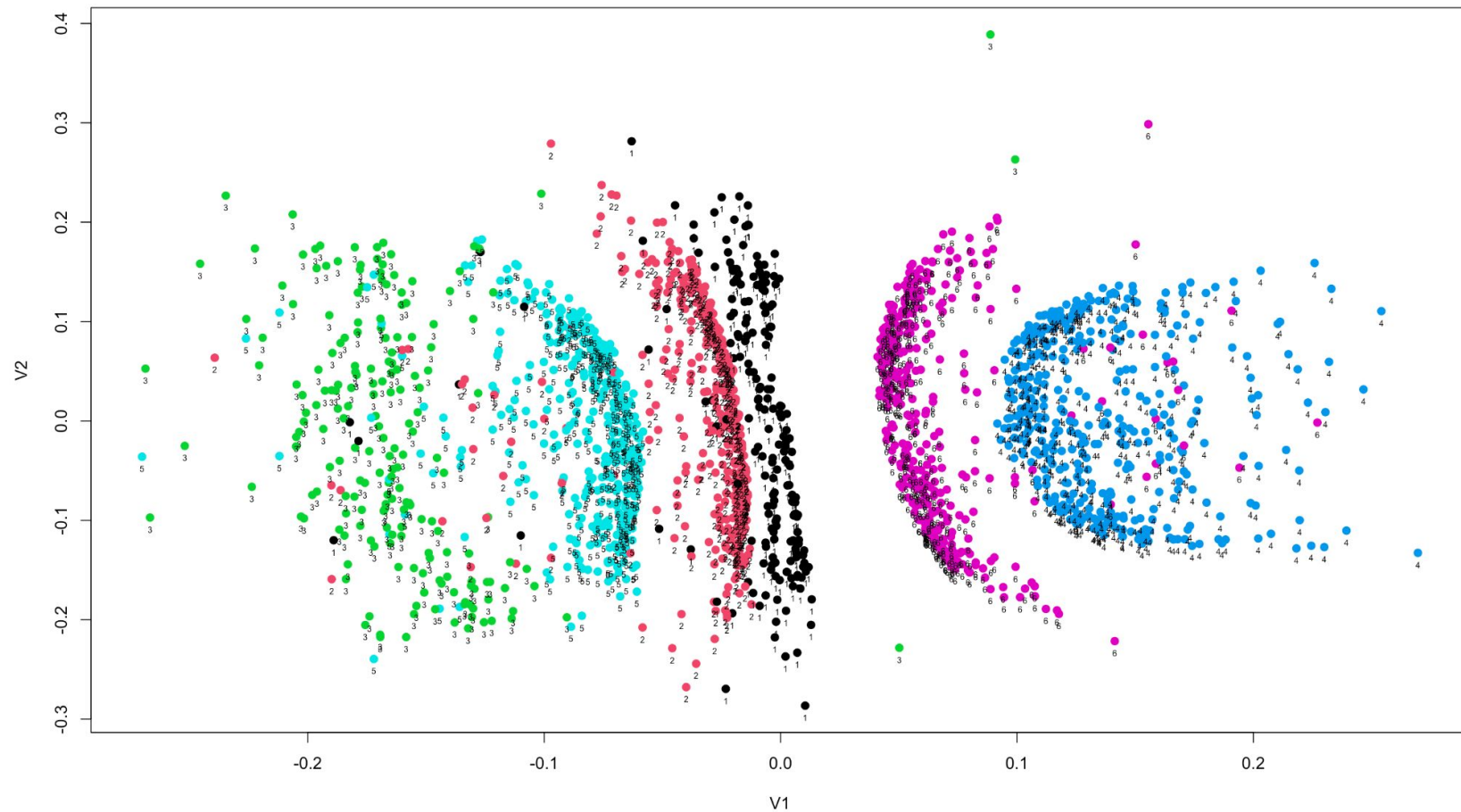
Cluster Dendrogram



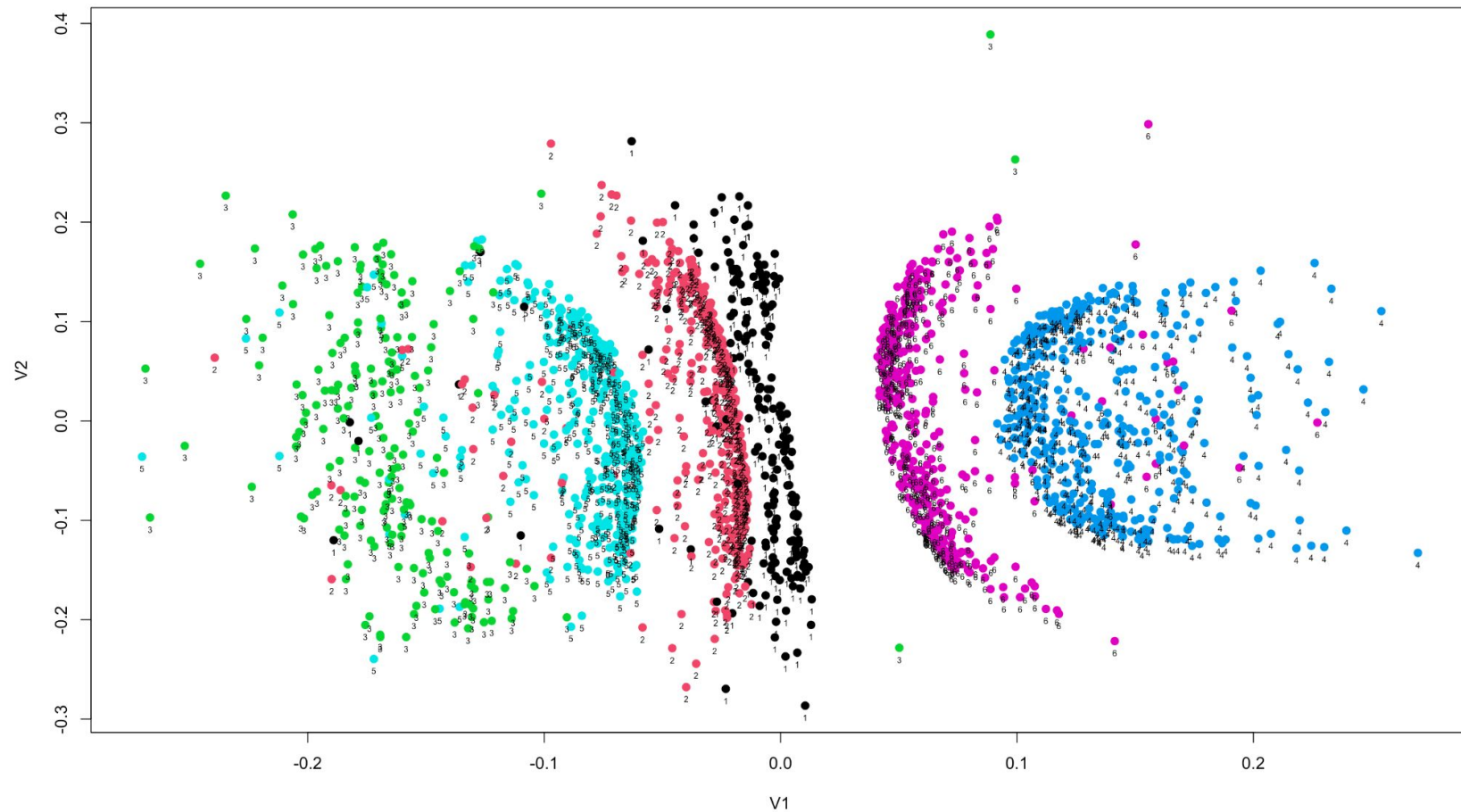
# HDBSCAN\*



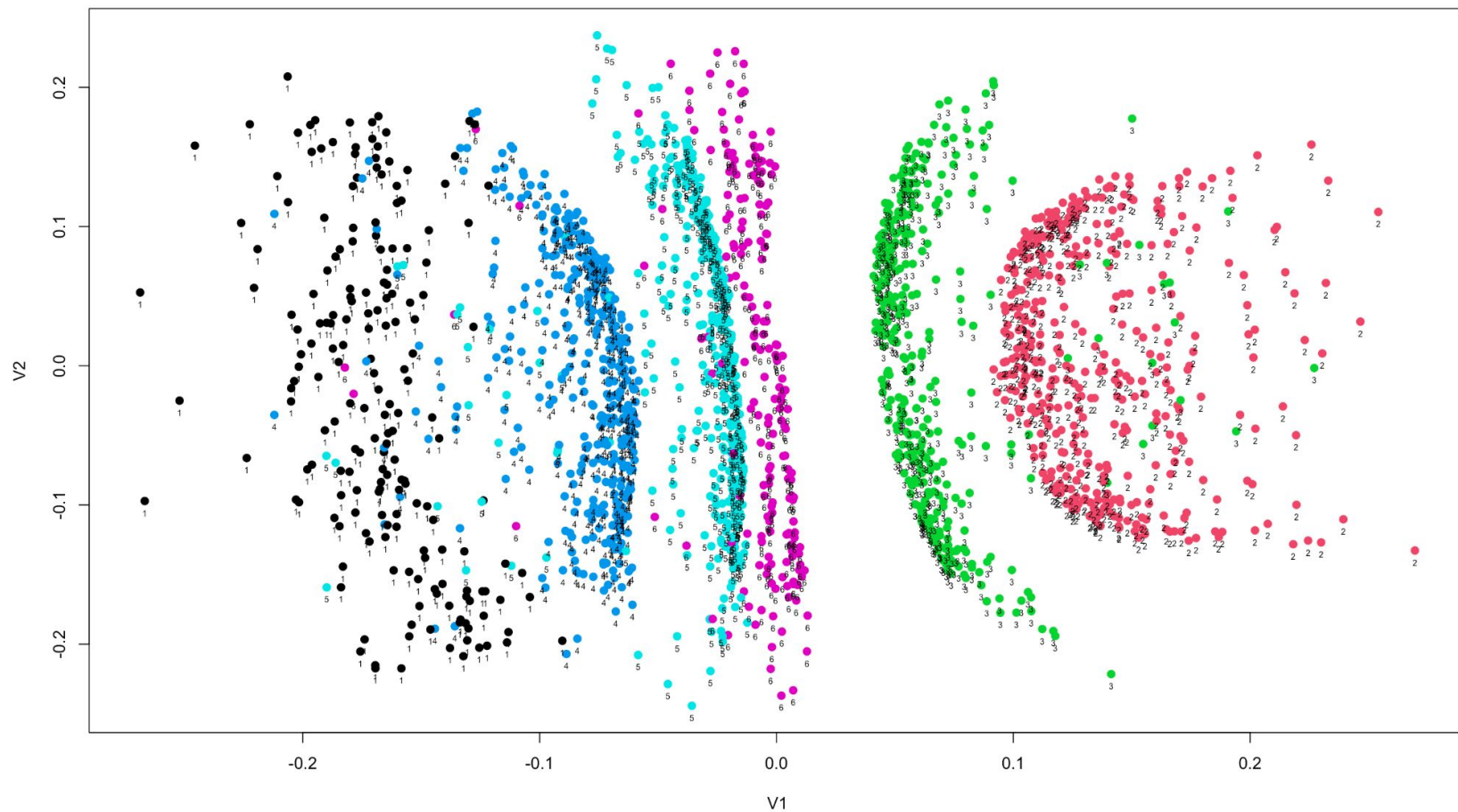
# PAM



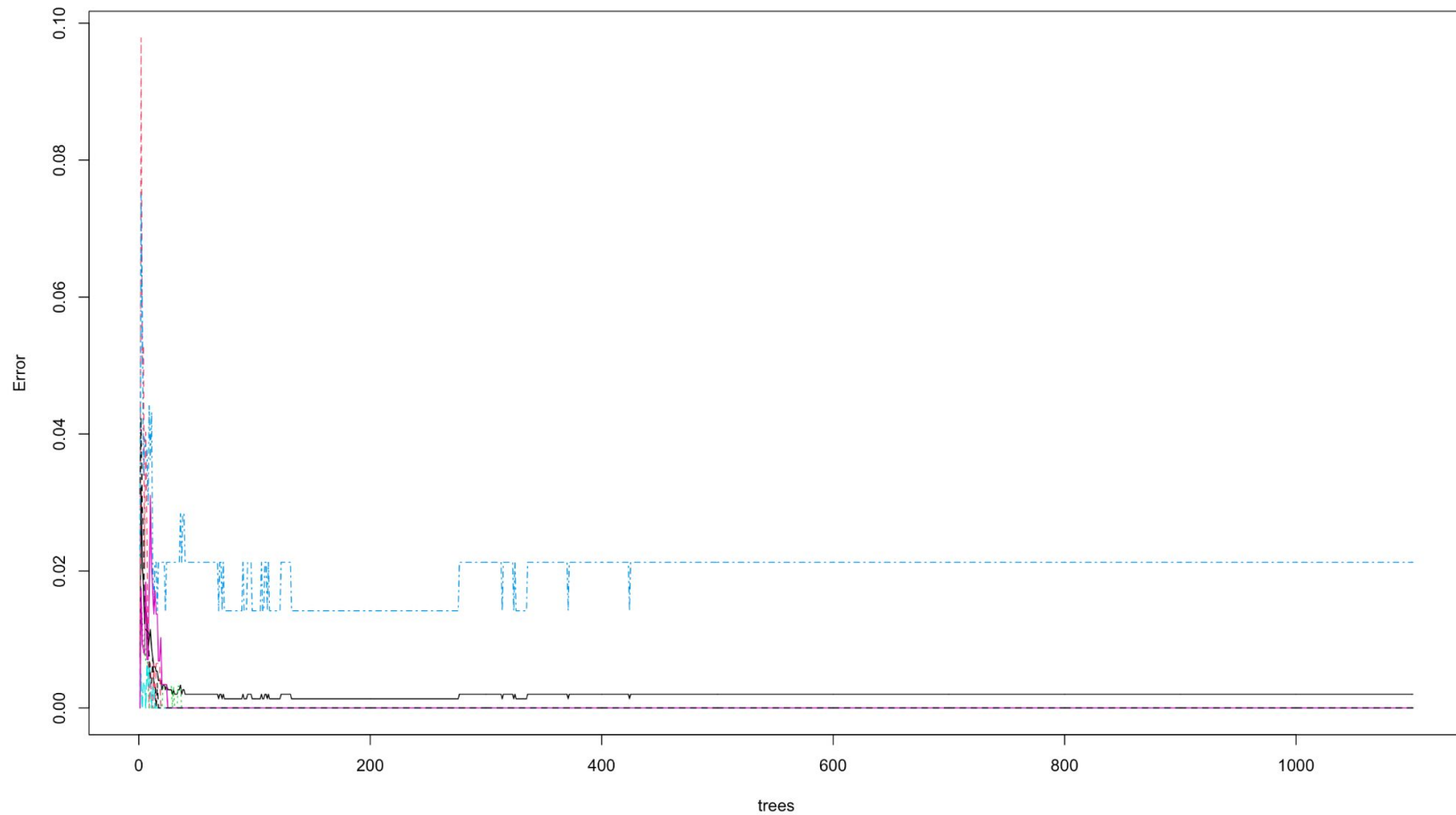
# HCLUST



# HDBSCAN

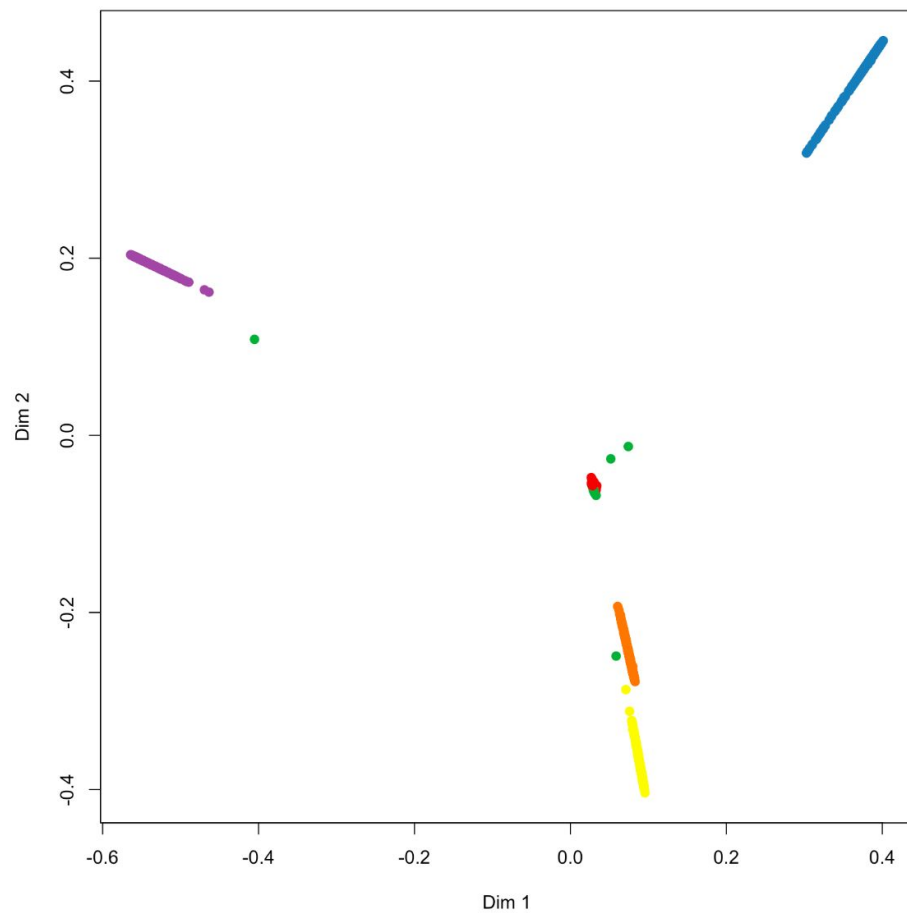


RF PAM Predictors Error

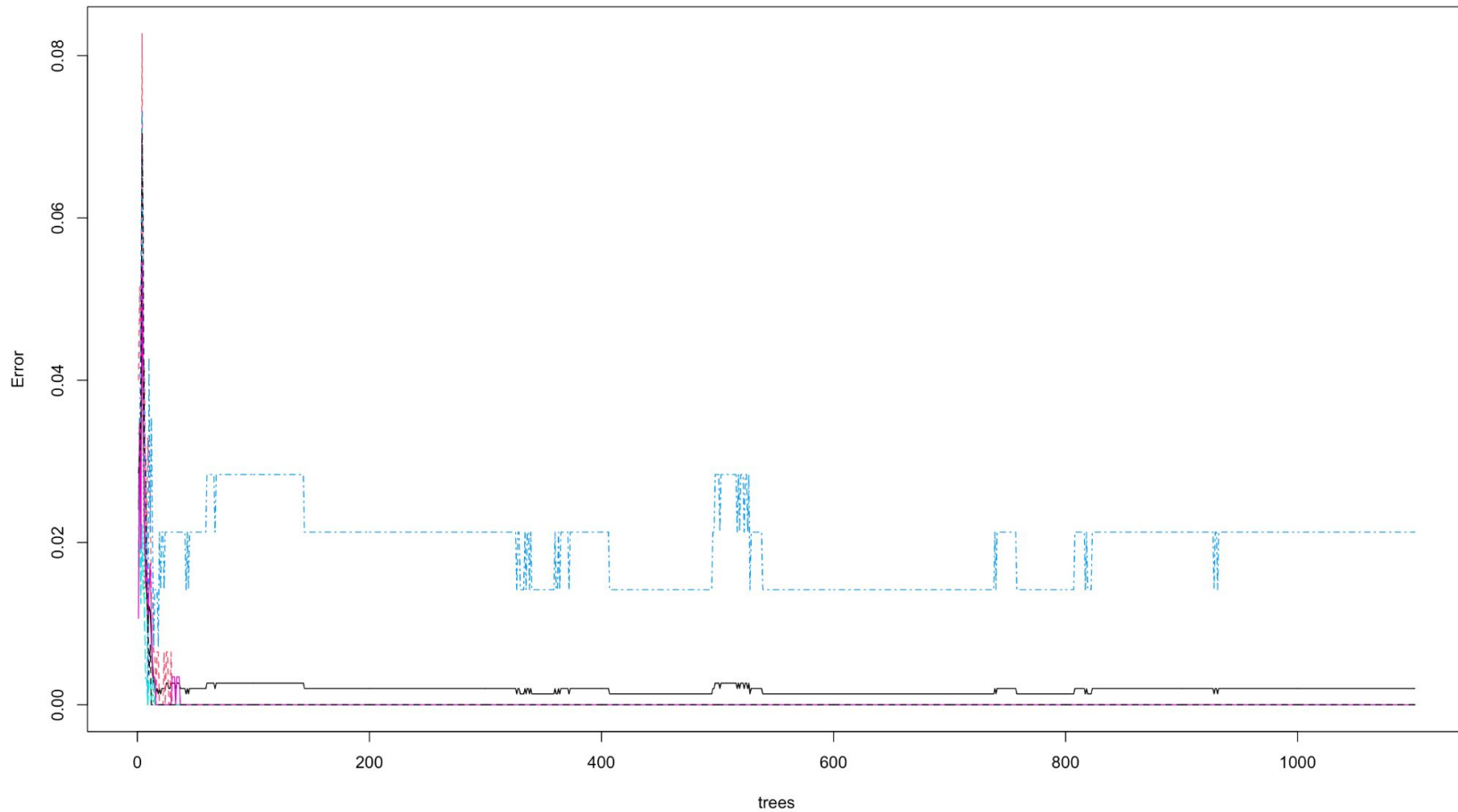




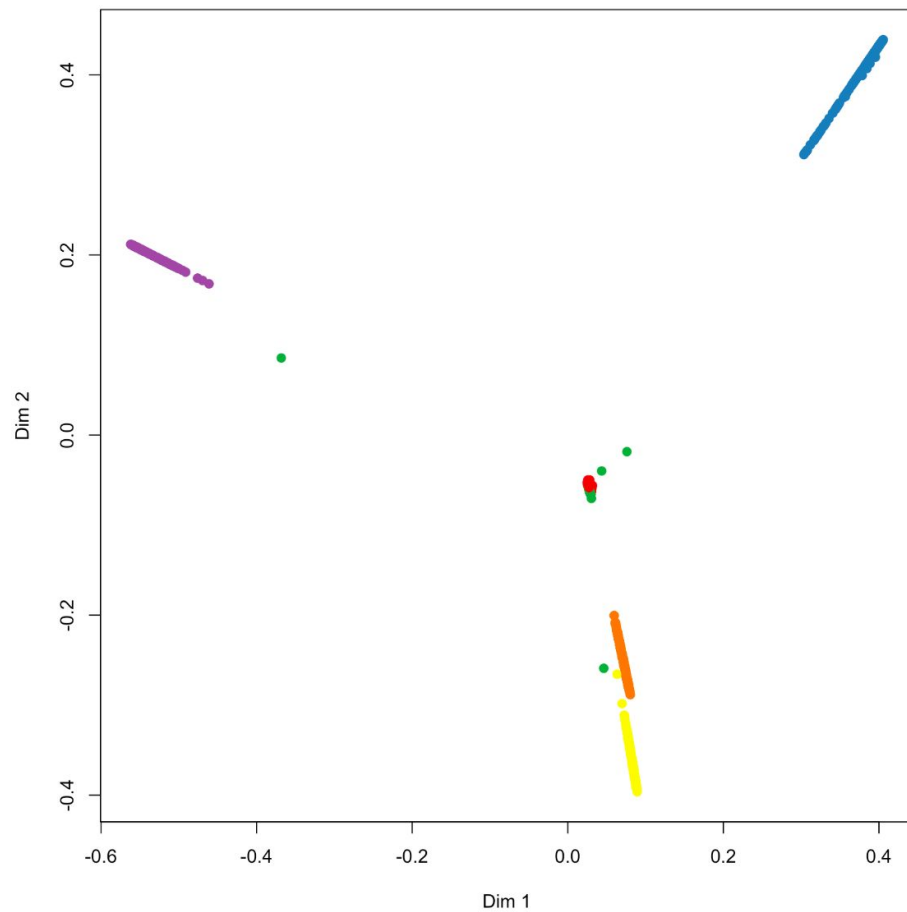
RF PAM



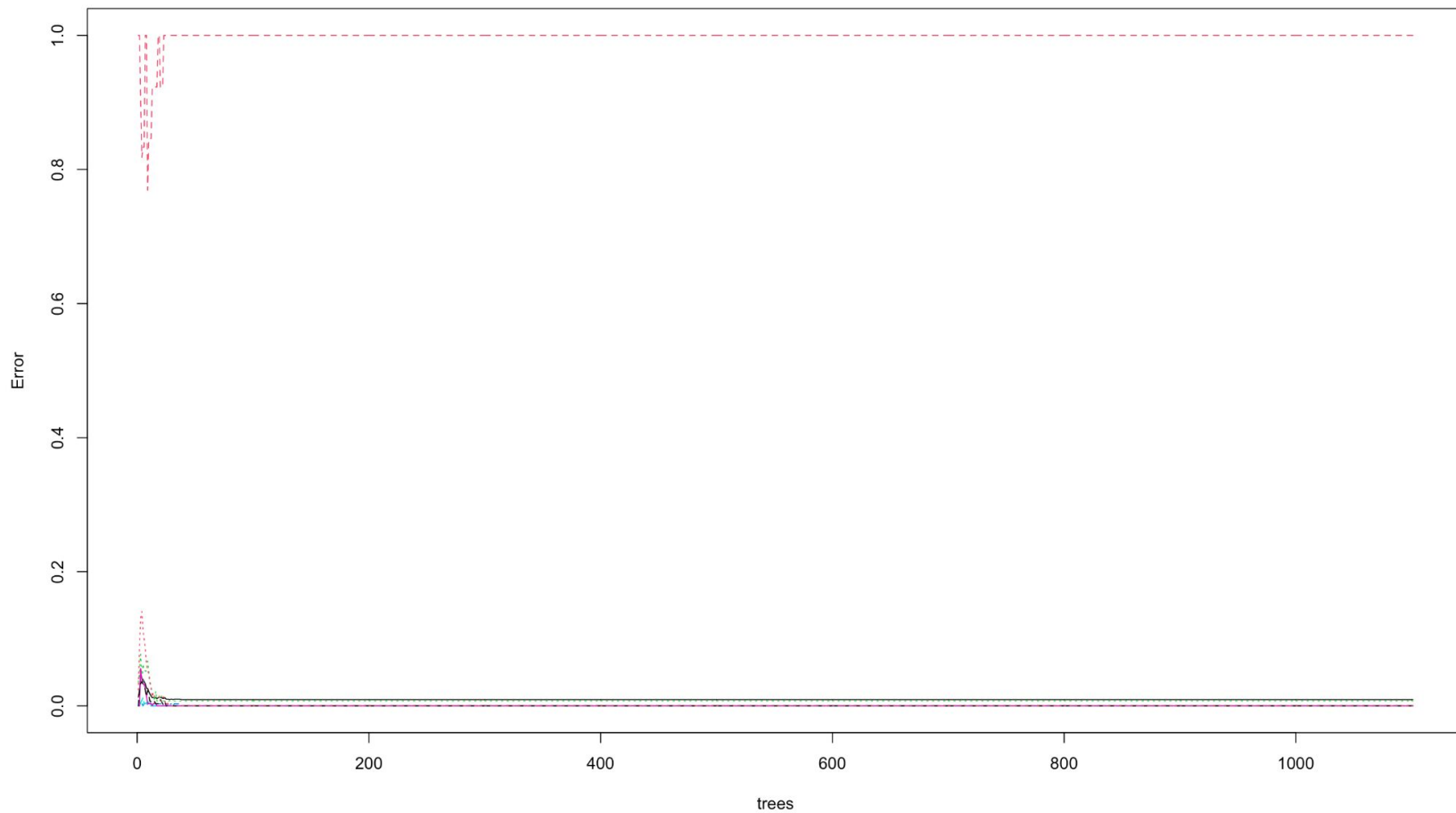
RF HCLUST Predictors Error



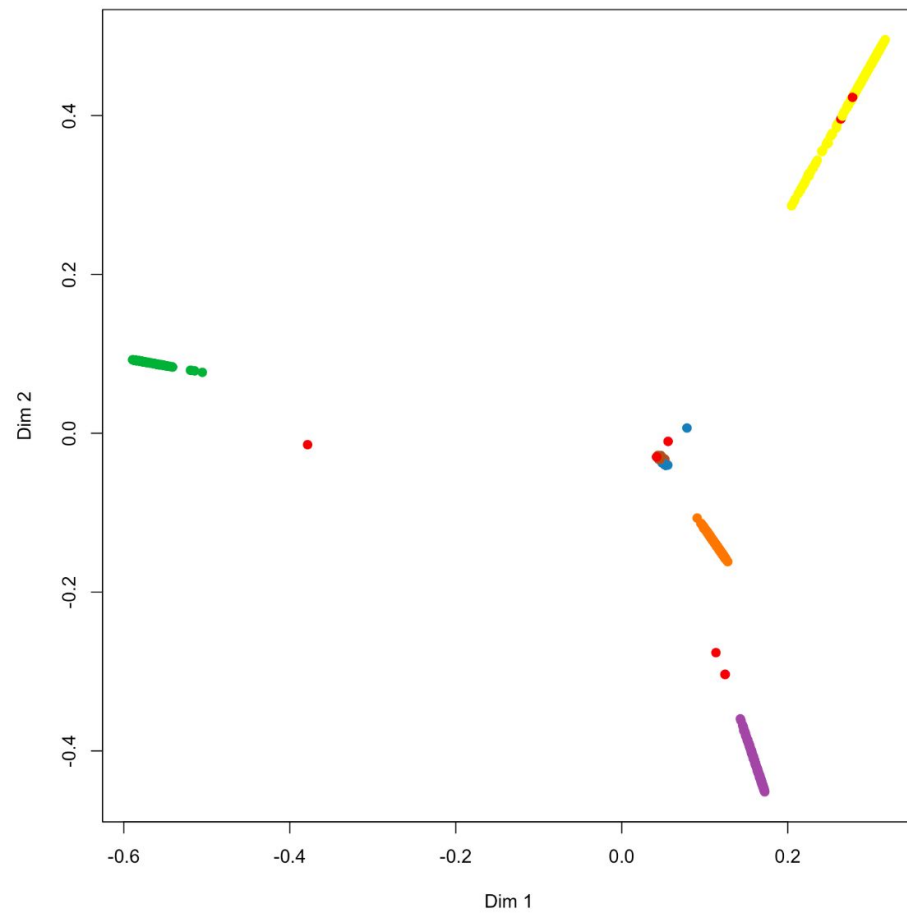
RF HCLUST



RF HDBSCAN Predictors Error



RF HDBSCAN



```
> RF.EARN[["importance"]]
```

	%IncMSE	IncNodePurity
GPA	4880.2691	141039235
Number_Of_Professional_Connections	584015.2996	2057388875
Major	8809639.8968	42349063726
Graduation_Year	-204.6793	81038460
Number_Of_Credits	-2150.7916	52026449
Number_Of_Parking_Tickets	218.2527	29991946

```
> |
```

So, the importance of the variables in predicting the earnings are in this order based on IncNodePurity: Major => Connections => GPA => Graduation year => Credits => Parking Tickets

	ERROR	PAM	HCLUST	HDBSCAN
1	OOB	0.2	0.2	0.9333333
2	TEST	0.2	0.2	0.6000000

PAM and HCLUST had the lowest test set error rates of 0.2%. We will examine the medoids of PAM next.

	▲ GPA	Number_Of_Professional_Connections	Earnings	Graduation_Year	Number_Of_Credits	Number_Of_Parking_Tickets	Major	PAM	HCLUST	HDBSCAN
1	2.65	10	10265.00	1992	122	1	Buisness	1	1	6
2	2.60	10	10254.77	1987	122	1	Humanities	2	2	5
3	2.61	11	5135.74	1986	121	1	Other	3	3	1
4	2.57	14	13258.34	1989	122	1	Vocational	4	4	2
5	2.52	7	9748.50	1987	122	1	STEM	5	5	4
6	2.37	7	11761.69	1990	121	1	Professional	6	6	3

These are the medoids of PAM and their respective cluster assignments



```
> RF.EARN
```

```
Call:
```

```
randomForest(x = Data.Original[, .SD, .SDcols = !c("Earnings")], y = Data.Original[, Earnings],  
ntree = 2777, mtry = 4, importance = TRUE, proximity = TRUE, keep.forest = TRUE)  
Type of random forest: regression  
Number of trees: 2777  
No. of variables tried at each split: 4  
  
Mean of squared residuals: 16869.37  
% Var explained: 99.62
```

**MSE = 16869.37**

**% of the Variance explained = 99.62%**