

Analysis F23

Abhi Thanvi, Jonathan Sneh

2023-12-04

Contents

Project Overview	2
Literature Review	3
Data Processing	4
Feature Engineering	4
Data Summary	5
Correlation Summary	6
Unsupervised Learning Algorithms	8
K-Means Algorithm	8
Hierarchial Clustering	24

Project Overview

- **Goals:**
- **Approach:**
- **Unique Methods:**
- **Conclusion:**

Literature Review

Data Processing

This section dives into the tasks performed for data processing. All the steps ensure the specifications of the projects were met, but some decisions were also made to ensure a more practical data to work with. To be considerate of the pages used for the Data Processing, we performed our Data Engineering steps in a `jupyter notebook` that you can view in our repo!

Feature Engineering

- **User ID** [`id`, `string`] — Unique IDs of each user.
 - We keep this to ensure tracking of user information for processing and analysis work.
- **Event ID** [`event_id`, `string`] — Incremental ID log of all events.
 - We keep this for processing steps, but remove it prior to analysis. The event IDs are useful as an ordinal feature of the log data.
- **Down Time / Up Time** [`down_time` / `up_time`, `integer`] — Time of event on down and up strokes of key or button, in seconds.
 - We summarize these features as an array of summary statistics; min, max, mean, median, and standard deviation. Measures of interest are max (i.e., how long a paper is written) and mean/median (i.e., when the center of most activity is).
- **Action Time** [`action_time`, `integer`] — Difference of time between down time and up time of event, i.e., duration of action in seconds.
 - Similarly, we summarize this feature as min, max, mean, median, and std. This gives insight into “major” consecutive actions, hesitancy, or other special behaviors.
- **Activity** [`activity`, `string`] — Actions to edit or modify the text (input, remove/cut, nonproduction, etc.)
 - We compute the proportions of each of these activities. All of the cursor “Move From” events are mapped to one category called “Move From”. We choose proportions over count to avoid undue influence of essays that take longer to write.
- **Down Event**
 - We compute the proportions of each of the activities. The events were pooled into four categories: alphanumeric, special_characters, control_keys, and unknown.
- **Up event**
 - Since these are the same events as down events, we ignore this feature.
- **Text Change**
 - We process and cluster these values into identified patterns of changes: many characters (at least 2 alphanumeric), at least one character (exactly one alphanumeric), non-zero characters (no alphanumeric). We also identified “transition” groups of “X to Y” for each of “many”, “single”, “none” (e.g., “many” to “many”, “many” to “single”, “many” to “none”, etc.). There was also a “no change” group. We created one additional group to represent the sum of all “transition” events because they coincided exclusively with “replacement” activities.
- **Cursor Position**
 - We computed an artificial array of cursor positions with the assumption that the text was streamed with no edits corresponding to what text changes there are (i.e., non-decreasing and doesn’t change if “no change” is observed in text change feature). Then we compute the MAE error metric between this stream version and the actual cursor positions to measure how much error exists between them. Greater errors imply more frequent and/or drastic changes.

- **Word Count**

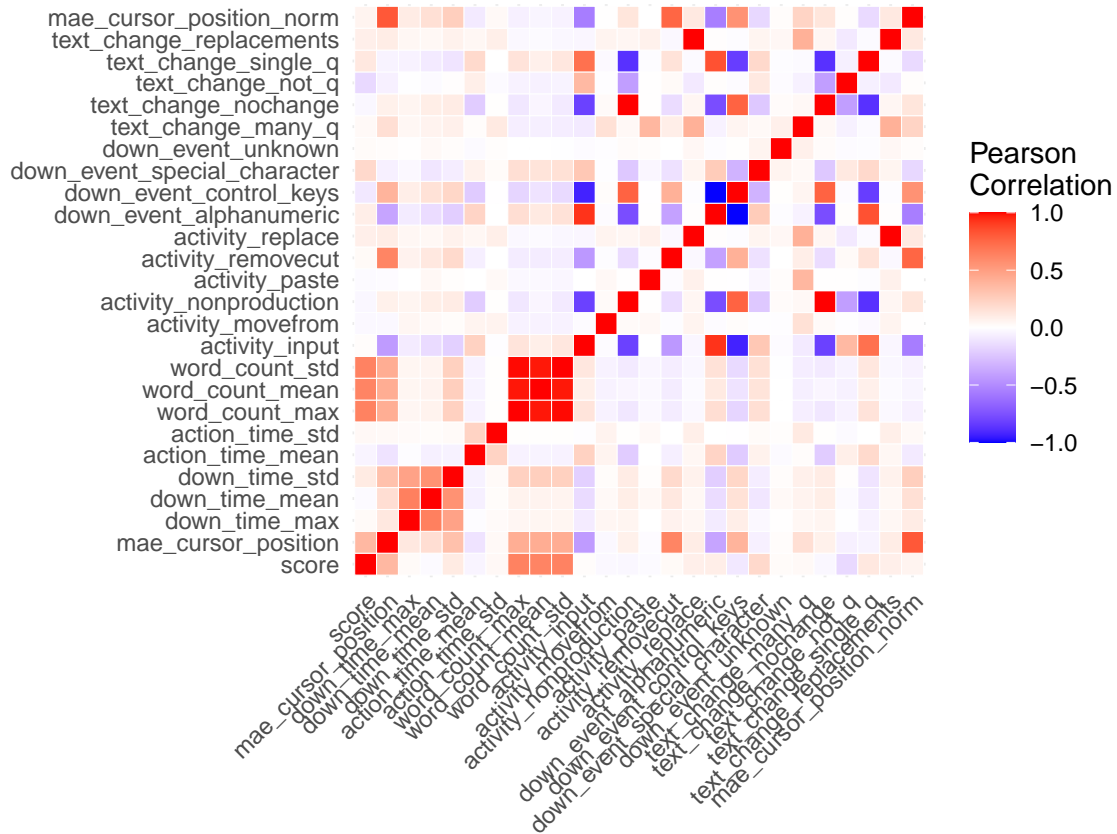
- We summarize these features as an array of summary statistics; min, max, mean, median, and standard deviation. We are primarily interested in the maximum measure as it indicates the length of the paper for each user.

Data Summary

Here is a sample of the processed data.

	001519c8	0022f953	0042269b	0059420b	0075873a	0081af50
id	001519c8	0022f953	0042269b	0059420b	0075873a	0081af50
score	3.5	3.5	6.0	2.0	4.0	2.0
mae_cursor_position	527.0469	380.7747	1238.7553	152.3933	640.1616	423.8706
down_time_max	1801877	1788842	1771219	1404394	1662390	1778845
down_time_mean	848180.8	518855.3	828491.8	785483.0	713354.2	544339.2
down_time_std	395112.7	384959.4	489500.8	385205.0	405576.4	484650.6
action_time_mean	116.24677	112.22127	101.83777	121.84833	123.94390	81.40434
action_time_std	91.79737	55.43119	82.38377	113.76823	62.08201	40.65305
word_count_max	256	323	404	206	252	275
word_count_mean	128.1162	182.7148	194.7727	103.6189	125.0830	132.9426
word_count_std	76.49837	97.76309	108.93507	61.88225	77.25505	81.20882
activity_input	0.7860774	0.7897311	0.8498549	0.8380463	0.7672857	0.8113976
activity_movefrom	0.00117325	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000
activity_nonproduction	0.04693000	0.10350448	0.04231141	0.06362468	0.02844725	0.03437359
activity_paste	0.0000000000	0.0004074980	0.0000000000	0.0006426735	0.0000000000	0.0000000000
activity_removecut	0.1630817	0.1059495	0.1061412	0.0970437	0.2042671	0.1528720
activity_replace	0.0027375831	0.0004074980	0.0016924565	0.0006426735	0.0000000000	0.0013568521
down_event_alphanumeric	0.6331639	0.6071720	0.7021277	0.6709512	0.6088503	0.6562641
down_event_control_keys	0.3523661	0.3712306	0.2860251	0.3155527	0.3646780	0.3364993
down_event_special_character	0.014470082	0.021597392	0.011847195	0.013496144	0.026471750	0.007236545
down_event_unknown	0	0	0	0	0	0
text_change_many_q	0.0007821666	0.0000000000	0.0007253385	0.0006426735	0.0000000000	0.0004522840
text_change_nochange	0.04693000	0.10350448	0.04231141	0.06362468	0.02844725	0.03437359
text_change_not_q	0.1908487	0.2041565	0.1677950	0.1985861	0.1955749	0.1804613
text_change_single_q	0.7587016	0.6919315	0.7874758	0.7365039	0.7759779	0.7833559
text_change_replacements	0.0027375831	0.0004074980	0.0016924565	0.0006426735	0.0000000000	0.0013568521

Correlation Summary



Discussion

We observe some pairs of features that show signs of multicollinearity.

- The word count metrics are highly correlated as expected, we could reasonably choose the maximum measure to use.
- Some features form parallel or perpendicular colinearity.
 - activity_input and activity_nonproduction (negative)
 - activity_input and down_event_alphanumeric (positive)
 - activity_input and down_event_control_keys (negative)
 - activity_input and text_change_nochange (negative)
- Importantly, we're interested in what's correlated with the user score feature
 - word count measures have a positive correlation with score, suggesting an association between longer essays and higher scores
 - Error rate of cursor positions against a “streamed” output also shows a positive correlation with score - i.e., essays written with less frequent or extreme edits is somewhat associated with higher scores.
 - Note: a positive correlation is also found between error rate of cursor positions with max word count, suggesting further that longer essays are associated with higher deviation from a “streamed” output. This suggests the possibility that interpretation of “streamed” deviation is influenced by the paper length (i.e., longer papers support possibility of edits being made “further away” from

the current “streamed” position, thus increasing the error rate). When we normalized the error rate by the paper size, we see that the correlation between the normalized error rate and the paper score is nearly zero. So, this feature is likely irrelevant for analysis.

Unsupervised Learning Algorithms

This section dives into the tasks performed for the unsupervised learning algorithms. Currently, focusing on K-Means and Hierarchical Clustering. **THIS SECTION SUPER MESSY RN. Please feel free to edit or improve in any way**

```
# test train split
test_ids = sample(1:nrow(df), as.integer(0.2 * nrow(df)))
data = scale(df[, -1])
train = data[-test_ids, ]
test = data[test_ids, ]
```

K-Means Algorithm

using averaged inertia of a few clustering samples across a range of number of clusters (i.e., $k=1, \dots, 14$)

```
# Range of k values to try
cluster_num_list <- 1:14

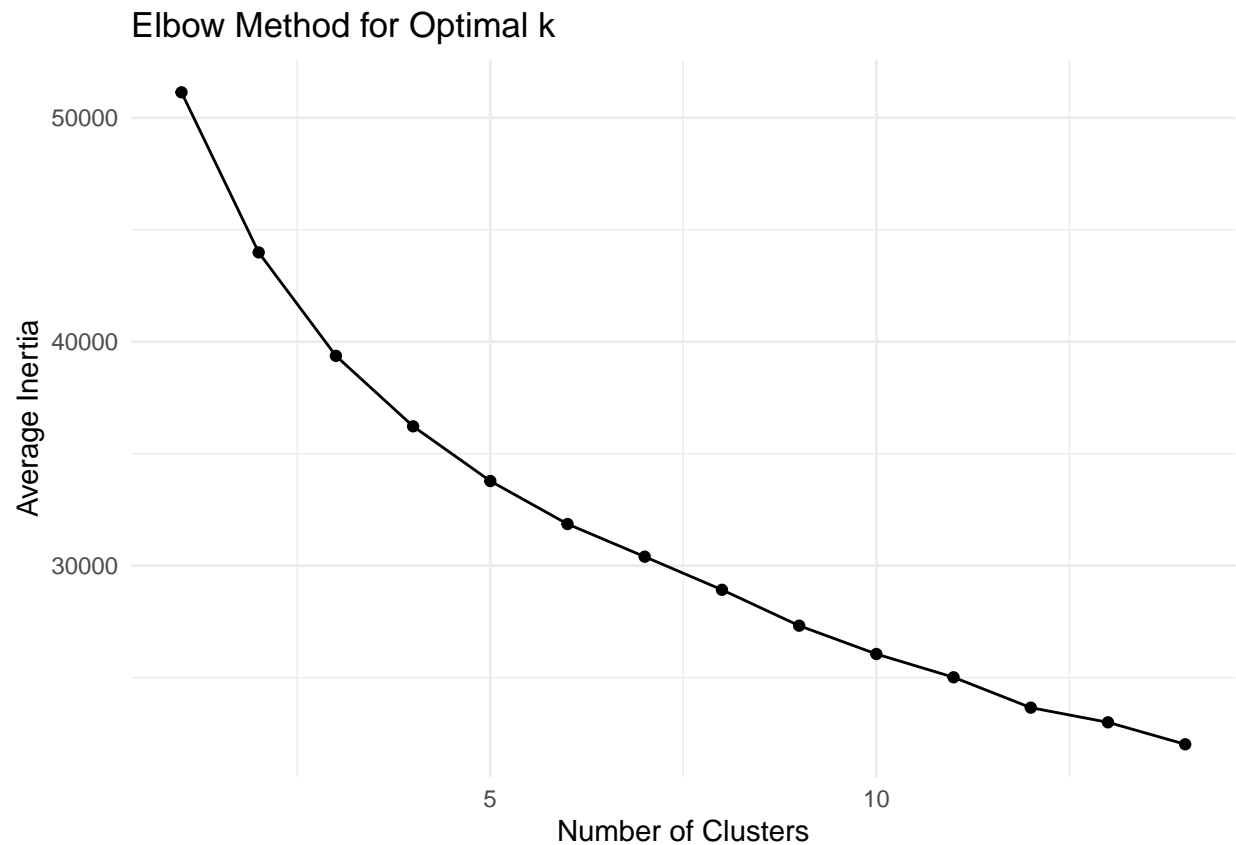
# Initialize a vector to store average inertias
avg_inertia_list <- numeric(length(cluster_num_list))

# Iterate over different values of k
for (k in cluster_num_list) {
  sub_inertia_list <- numeric(3) # For storing inertia of each trial

  for (i in 1:3) {
    set.seed(i) # Setting seed for reproducibility
    kmeans_result <- kmeans(train, centers=k, nstart=25, iter.max = 50)
    sub_inertia_list[i] <- kmeans_result$tot.withinss
  }

  avg_inertia_list[k] <- mean(sub_inertia_list)
}

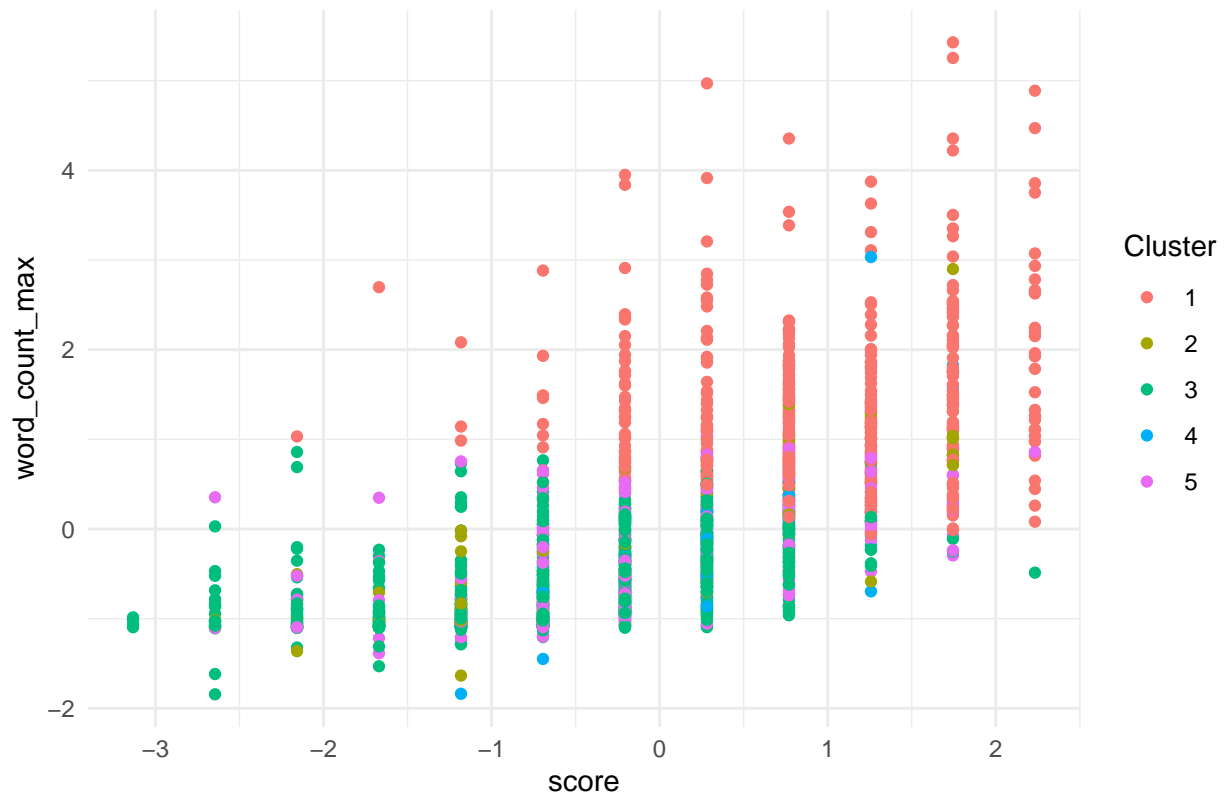
# Plotting the elbow plot
ggplot(data.frame(Clusters=cluster_num_list, Inertia=avg_inertia_list),
       aes(x=Clusters, y=Inertia)) +
  geom_line() +
  geom_point() +
  theme_minimal() +
  ggtitle("Elbow Method for Optimal k") +
  xlab("Number of Clusters") +
  ylab("Average Inertia")
```

selected $k=5$ clusters as closest to “elbow” of the inertia plot

```
# Perform K-means clustering  
# Here, we are specifying 3 clusters, but you can change this number  
result <- kmeans(train, centers=5)  
  
ggplot(data.frame(train), aes(x=score, y=word_count_max)) +  
  geom_point(aes(color=factor(result$cluster))) +  
  scale_color_discrete(name="Cluster") +  
  theme_minimal() +  
  ggtitle("K-means Clustering")
```

K-means Clustering



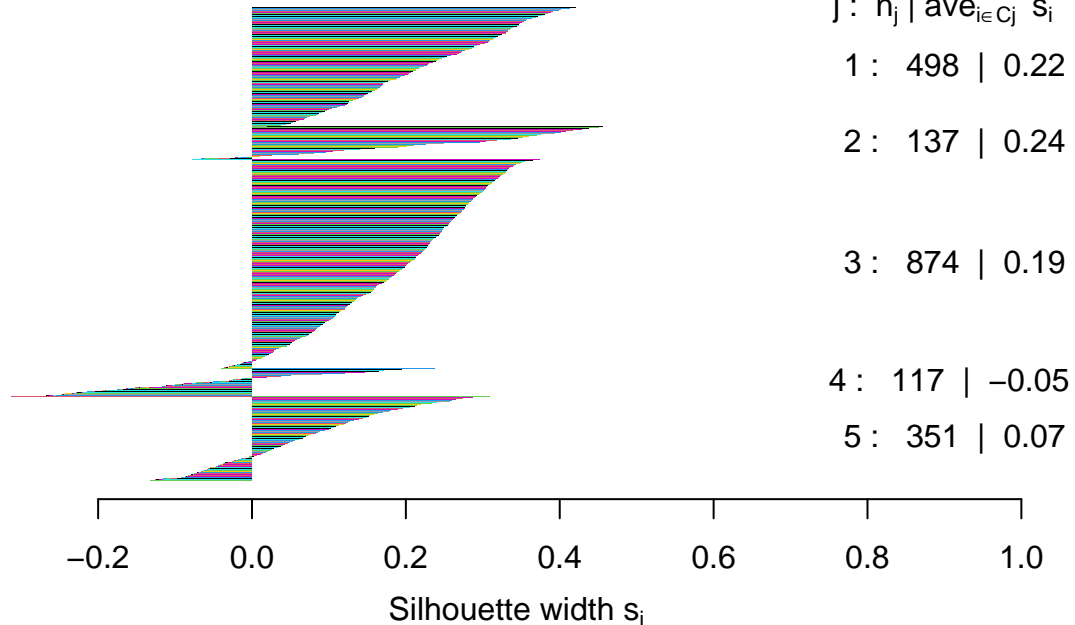
Using silhouette scores to evaluate how well the clustering structure fits in terms of similarity, i.e., more higher scores in a cluster imply greater similarity of points to its own cluster and poorer similarity to other clusters. The average silhouette score of 0.17 suggests that the clusters are somewhat favorably well separated and that the points within clusters aren't too dispersed. However, cluster 4 shows issues with cohesion and separation.

```
# Compute silhouette information
silhouette_info <- silhouette(result$cluster, dist(train))

# Plotting the silhouette plot
plot(silhouette_info, col=1:k, border=NA, main="Silhouette Plot")
```

Silhouette Plot

n = 1977



Average silhouette width : 0.17

Additionally, with projecting the data using UMAP, we see that the clusters might not be very well separated and aren't globular, so it is reasonable to conclude that KMeans algorithm may struggle with this data.

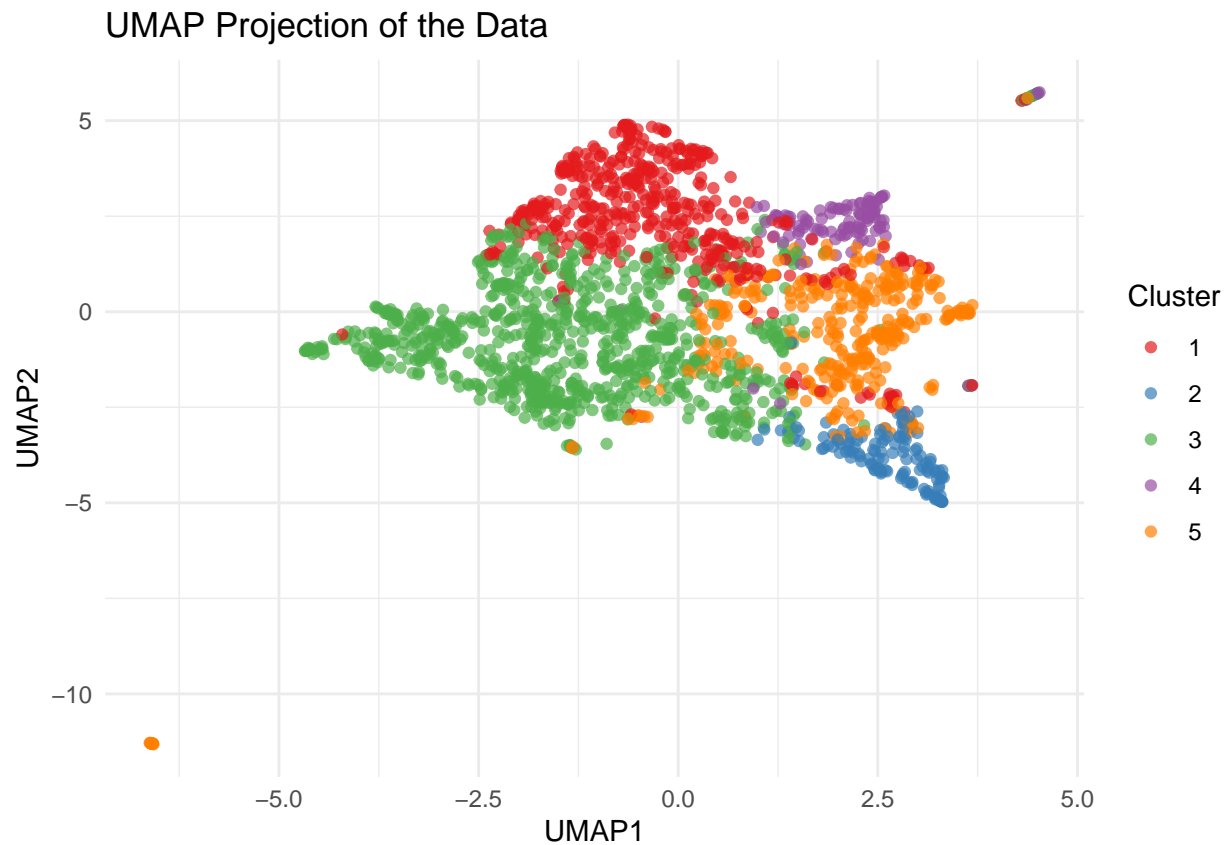
```
umap_result <- umap(train, n_components = 2)

umap_data <- as.data.frame(umap_result$layout)
colnames(umap_data) <- c("UMAP1", "UMAP2")

clusters <- as.factor(result$cluster)
umap_data$Cluster <- clusters

# Choose a palette
palette <- brewer.pal(n = 5, name = "Set1") # Adjust 'name' as needed

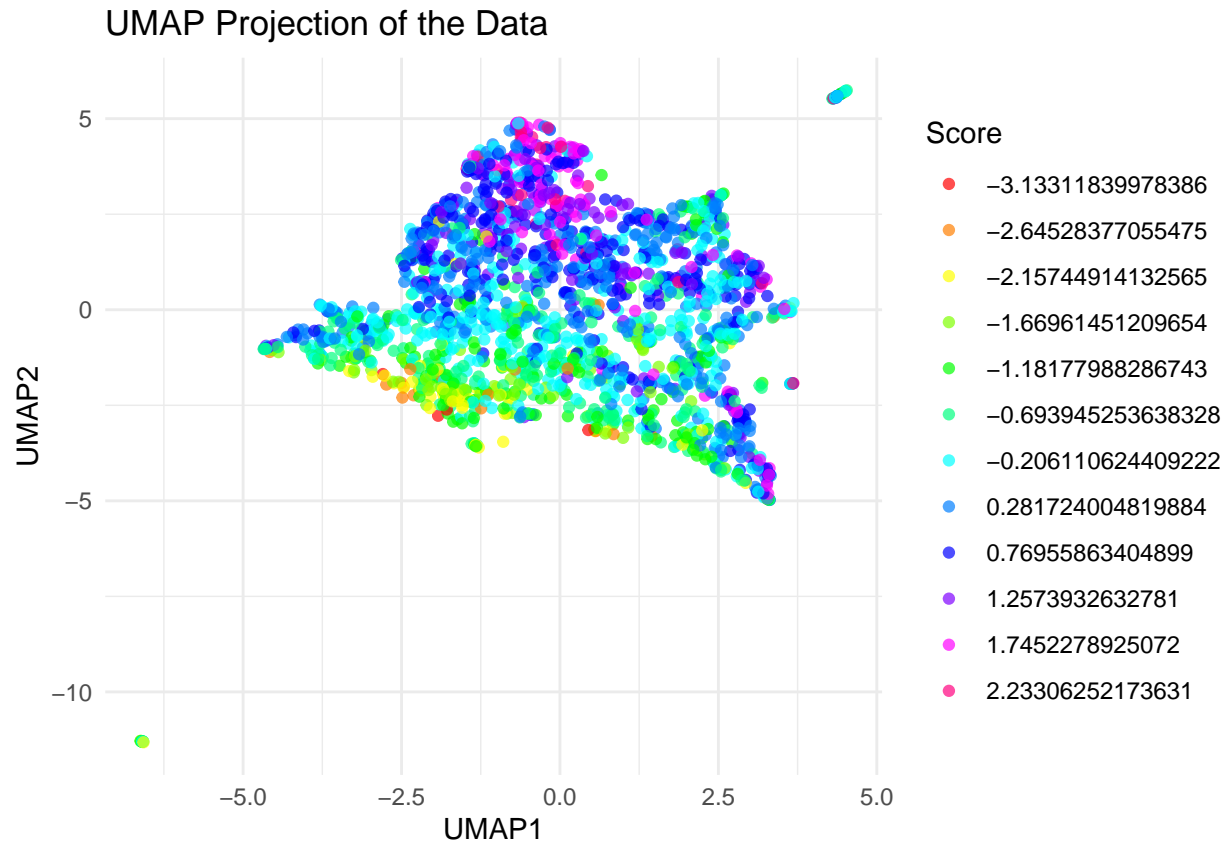
ggplot(umap_data, aes(x = UMAP1, y = UMAP2, color = Cluster)) +
  geom_point(alpha = 0.7) +
  scale_color_manual(values = palette) +
  theme_minimal() +
  ggtitle("UMAP Projection of the Data")
```



Here's distribution of scores in UMAP

```
scores <- as.factor(as.data.frame(train)$score)
umap_data$Score <- scores

ggplot(umap_data, aes(x = UMAP1, y = UMAP2, color = Score)) +
  geom_point(alpha = 0.7) +
  scale_color_manual(values = rainbow(12)) +
  theme_minimal() +
  ggtitle("UMAP Projection of the Data")
```



investigation of the score distribution suggests that the underlying clustering structure of the data does not closely align with the structure of scores.

spread of data between kmeans clusters and scores. characterizations are as follows:

- cluster 1 tends to capture those who score 3 to 6.
- cluster 2 tends to capture those who score between 3.5 to 4.5
- cluster 3 tends to capture those who score 1.5 to 4.5
- cluster 4 tends to capture those who score 3.5 to 4
- cluster 5 tends to capture those who score 2.5 to 4.5

e.g., a user who scores around 5 is likely to be in cluster 1

while there are some relationships, it would appear that the dispersions of scores between groups tend to overlap heavily and are not well separate to segment the groups in a very meaningful way. However clusters 1 and 2 vs. clusters 3, 4, and 5 seem to show some disparity.

```
t(table(result$cluster, df[-test_ids, ]$score)) / nrow(df[-test_ids, ])
```

```
##
##           1           2           3           4           5
## 0.5 0.0000000000 0.0000000000 0.0020232676 0.0000000000 0.0000000000
## 1   0.0000000000 0.0005058169 0.0096105210 0.0000000000 0.0015174507
## 1.5 0.0005058169 0.0020232676 0.0227617602 0.0000000000 0.0030349014
## 2   0.0005058169 0.0030349014 0.0288315630 0.0000000000 0.0060698027
## 2.5 0.0015174507 0.0096105210 0.0485584219 0.0060698027 0.0161861406
```

```
## 3 0.0040465352 0.0080930703 0.0870005058 0.0065756196 0.0298431968
## 3.5 0.0212443096 0.0111279717 0.1087506323 0.0141628730 0.0419828022
## 4 0.0445118867 0.0171977744 0.0864946889 0.0192210420 0.0359129995
## 4.5 0.0773899848 0.0101163379 0.0394537178 0.0070814365 0.0298431968
## 5 0.0460293374 0.0040465352 0.0070814365 0.0045523520 0.0080930703
## 5.5 0.0409711684 0.0035407183 0.0010116338 0.0015174507 0.0045523520
## 6 0.0151745068 0.0000000000 0.0005058169 0.0000000000 0.0005058169
```

to predict on new data in test...

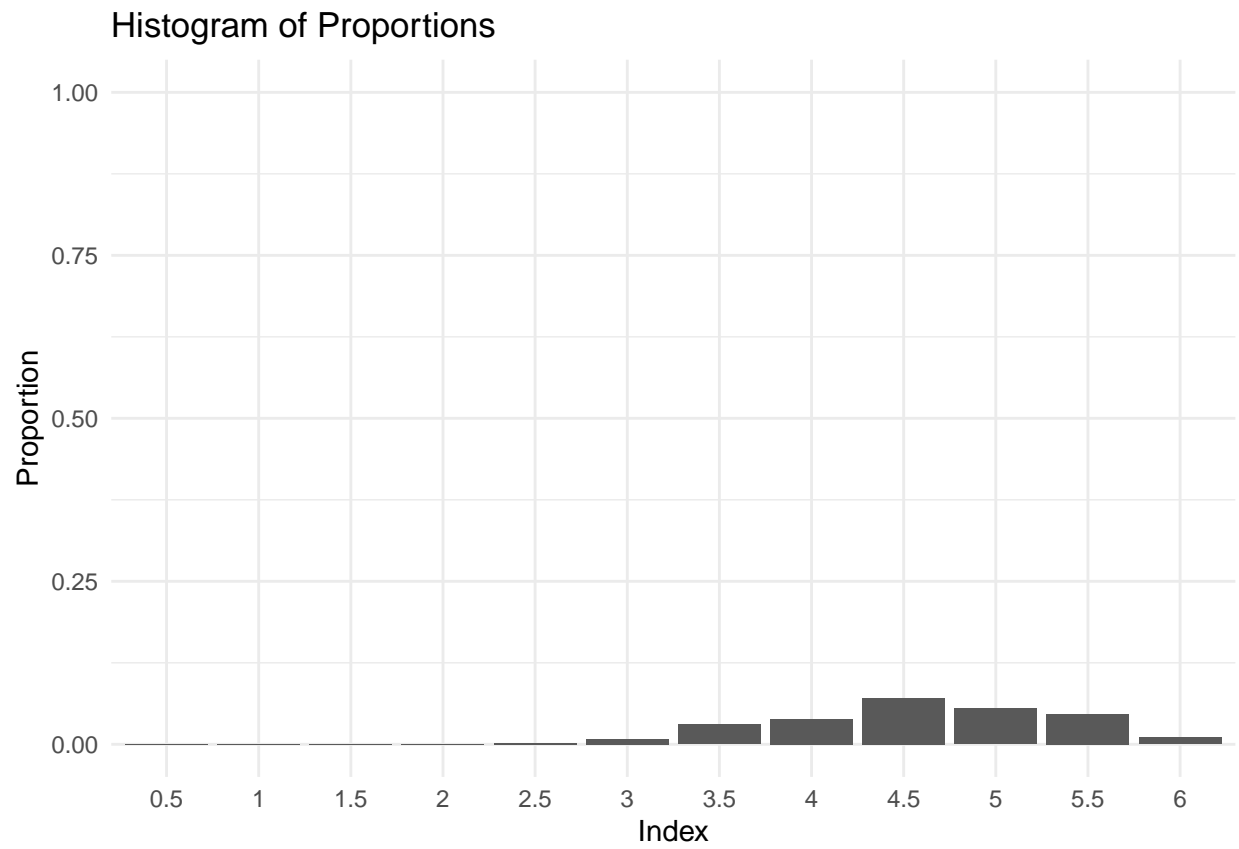
```
assign_cluster <- function(new_data, centers) {
  # Calculate Euclidean distances from each new data point to each cluster center
  distances <- as.matrix(dist(rbind(centers, new_data), method = "euclidean"))
  distances <- distances[(nrow(centers)+1):nrow(distances), 1:nrow(centers)]

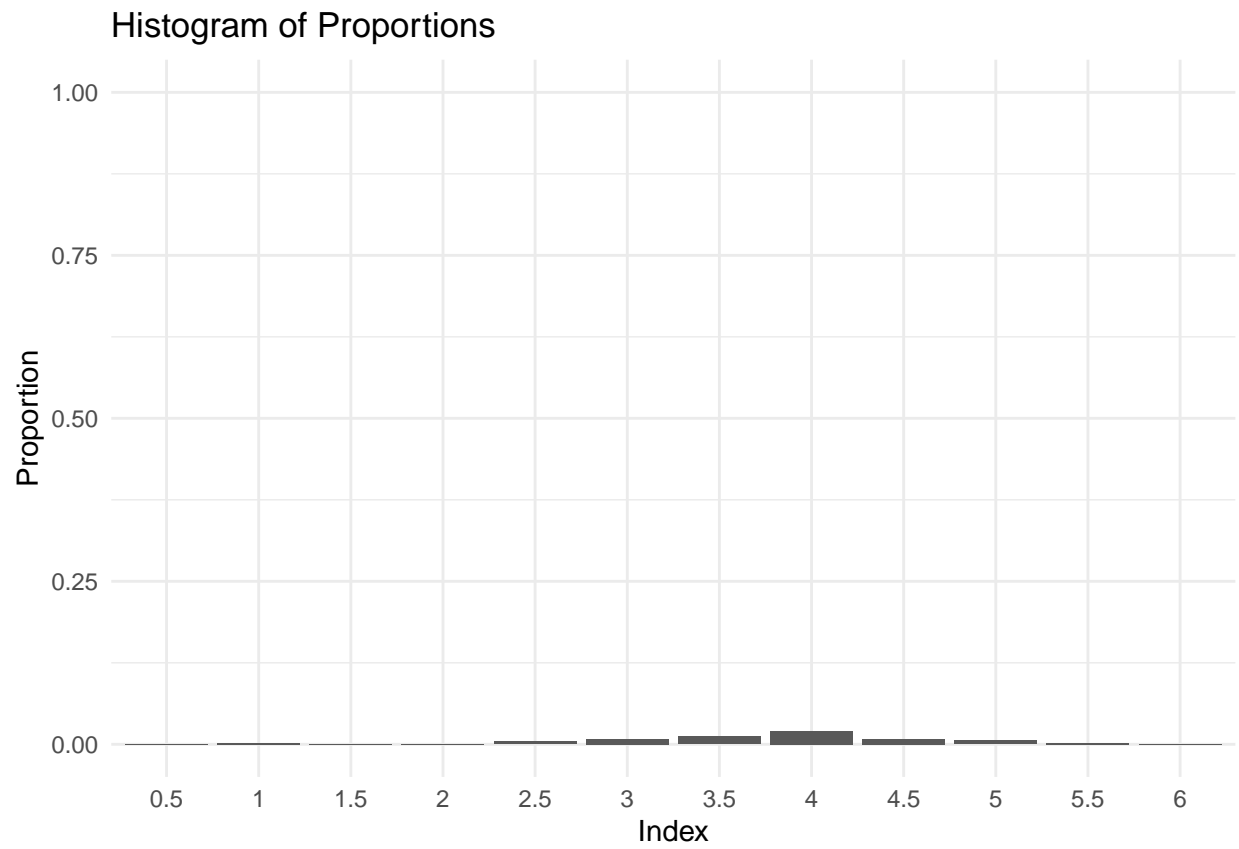
  # Assign each new data point to the nearest cluster
  max.col(-distances)
}

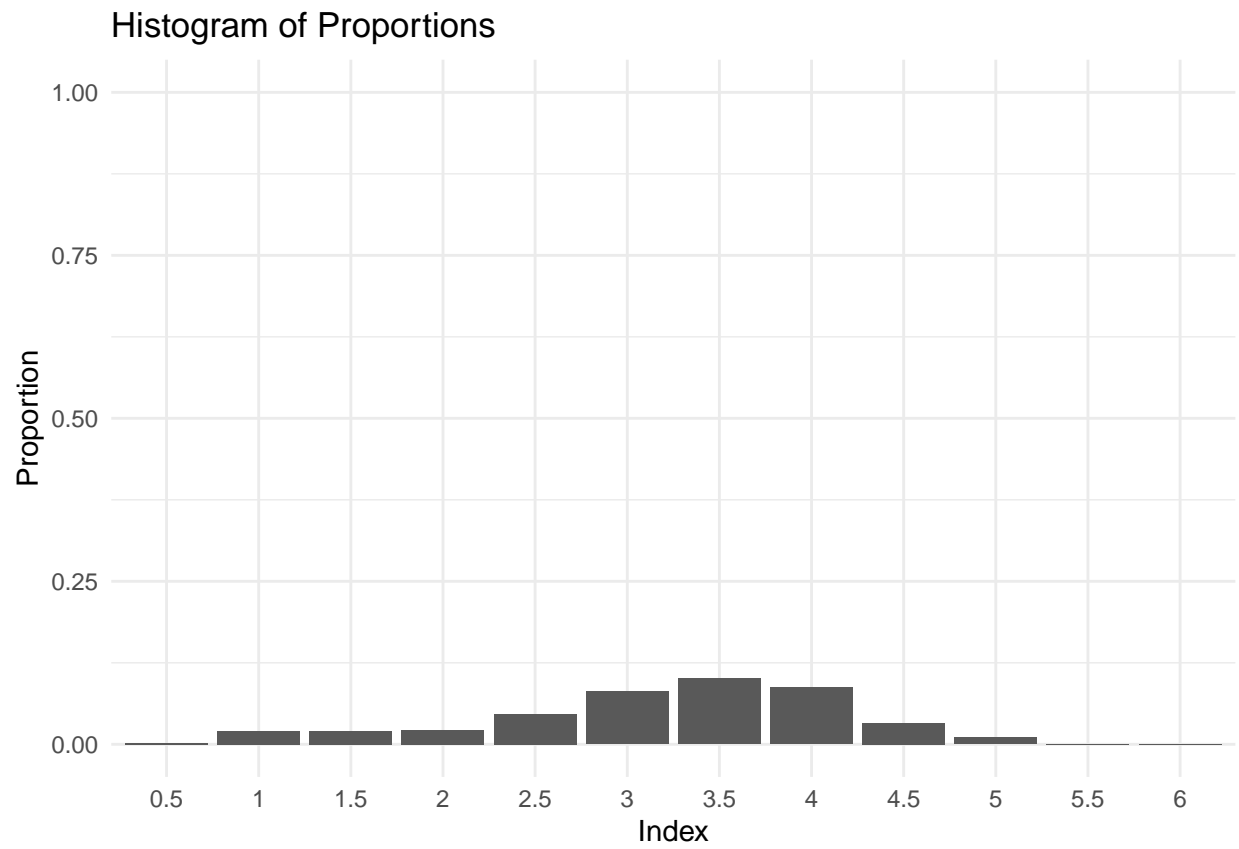
new_clusters <- assign_cluster(test, result$centers)
```

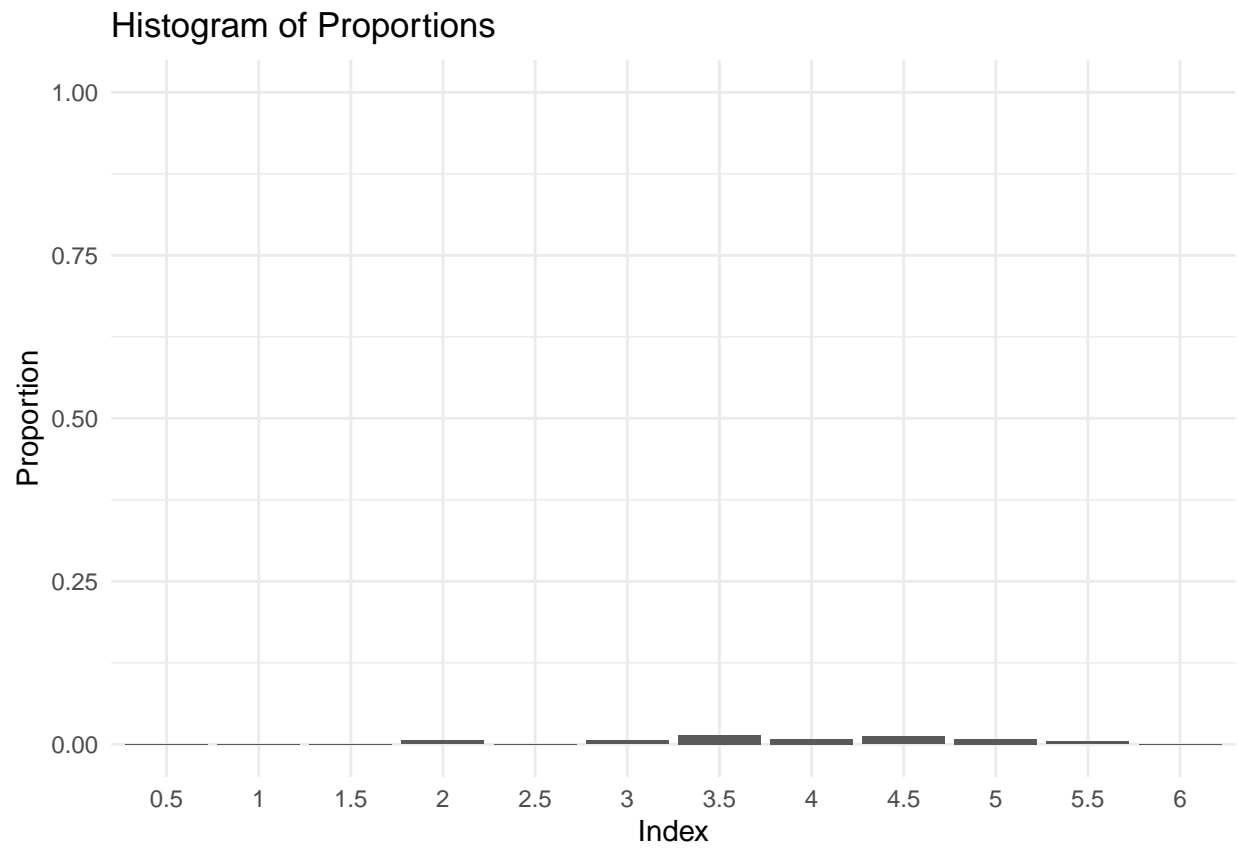
```
acc_data = as.matrix(
  t(table(new_clusters, df[test_ids, ]$score)) / nrow(df[test_ids, ])
)

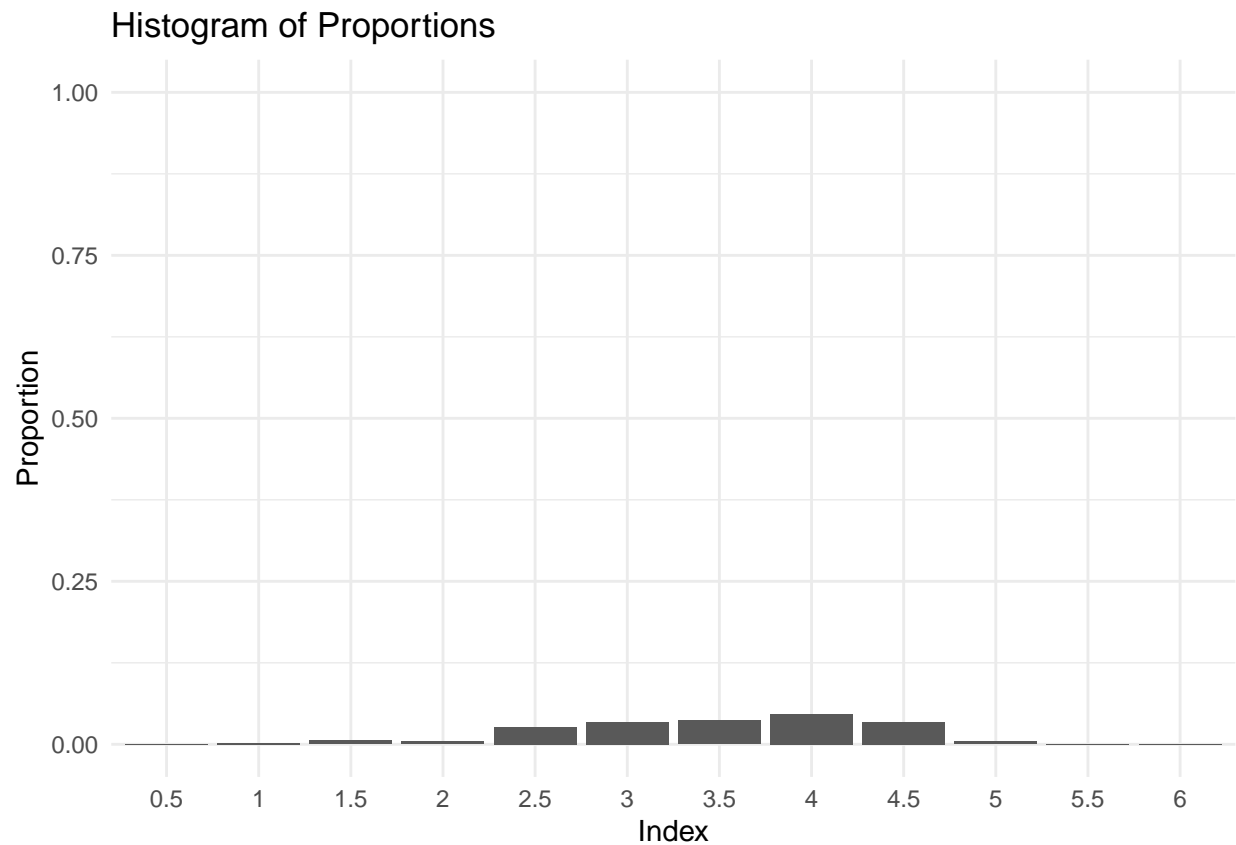
for (i in 1:5) {
  prop_plot = ggplot(
    data.frame(
      Index = rownames(acc_data), Proportion = acc_data[, i]
    ),
    aes(x = Index, y = Proportion)) +
  geom_bar(stat = "identity") +
  theme_minimal() +
  xlab("Index") +
  ylab("Proportion") +
  ggtitle("Histogram of Proportions") +
  ylim(c(0, 1))
  print(prop_plot)
}
```





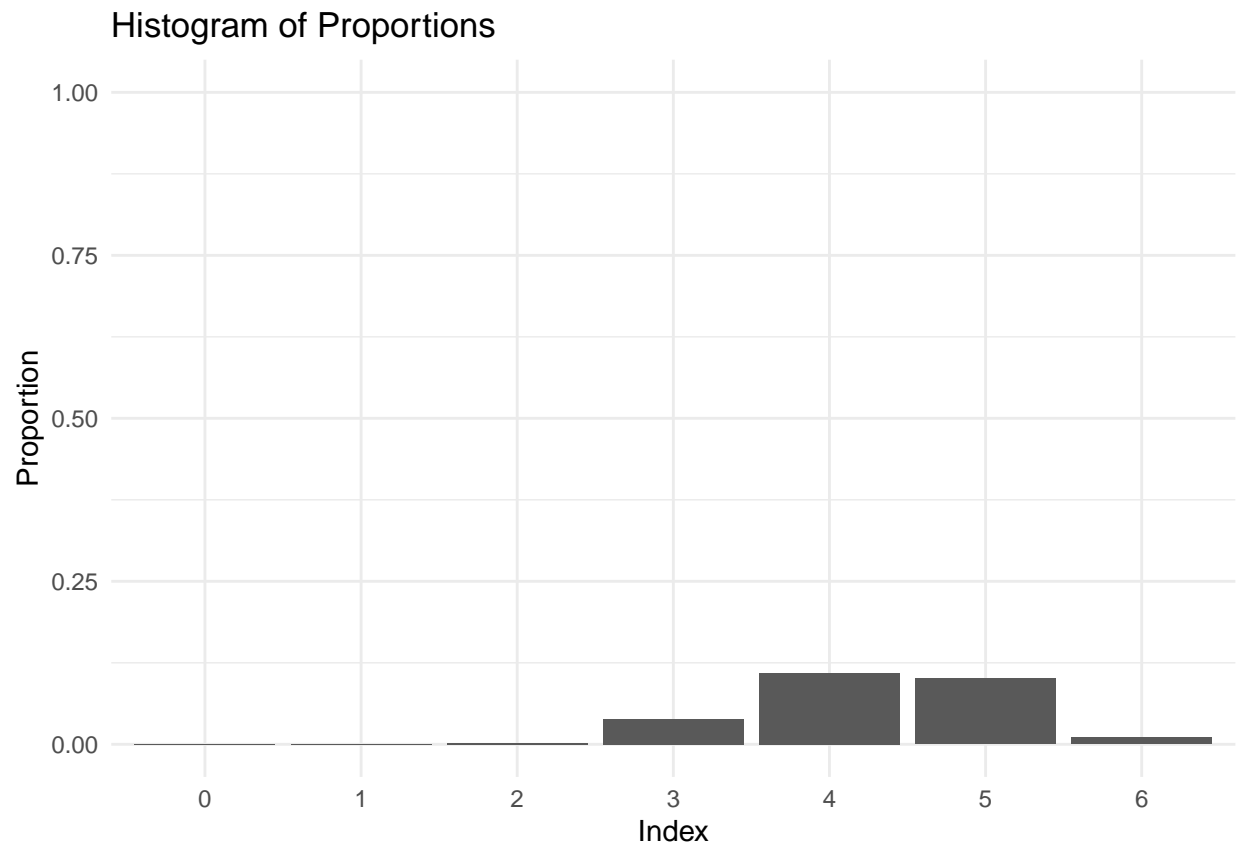


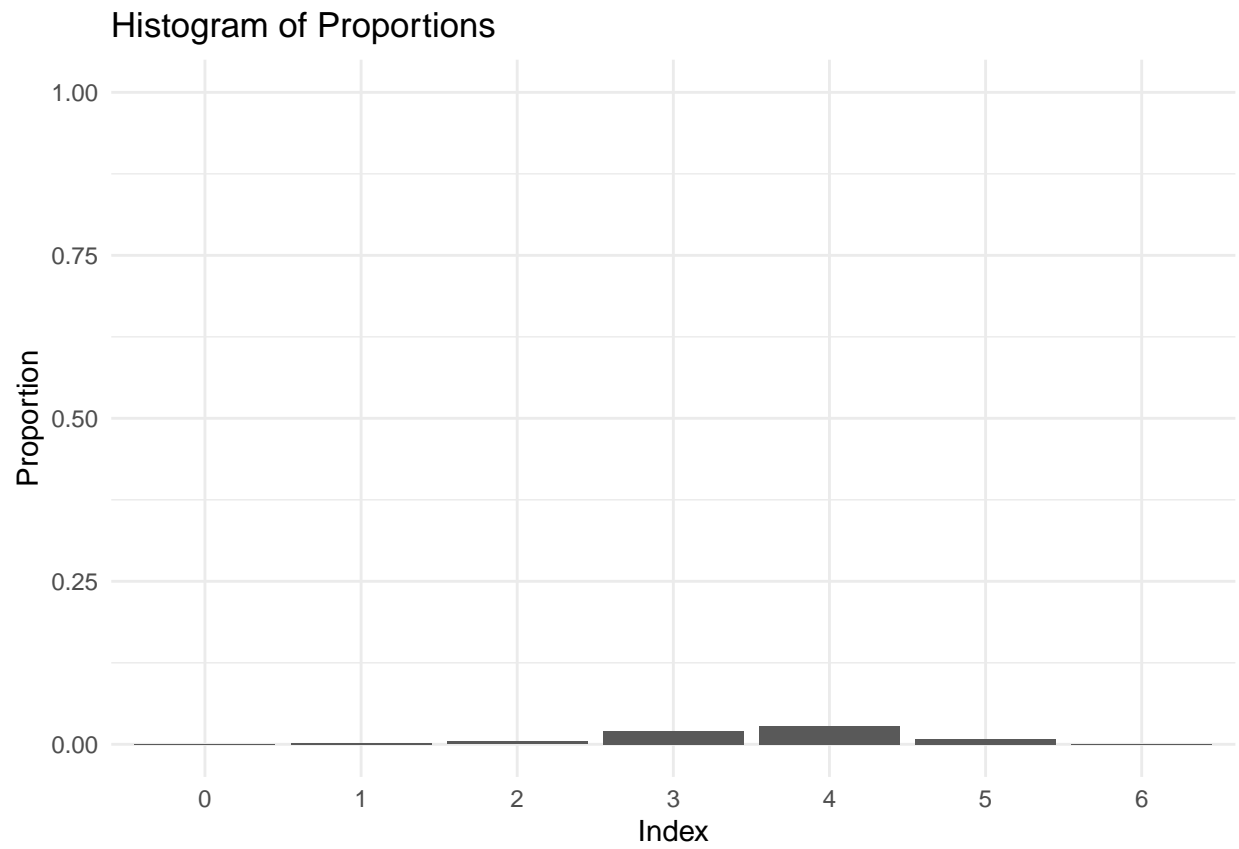


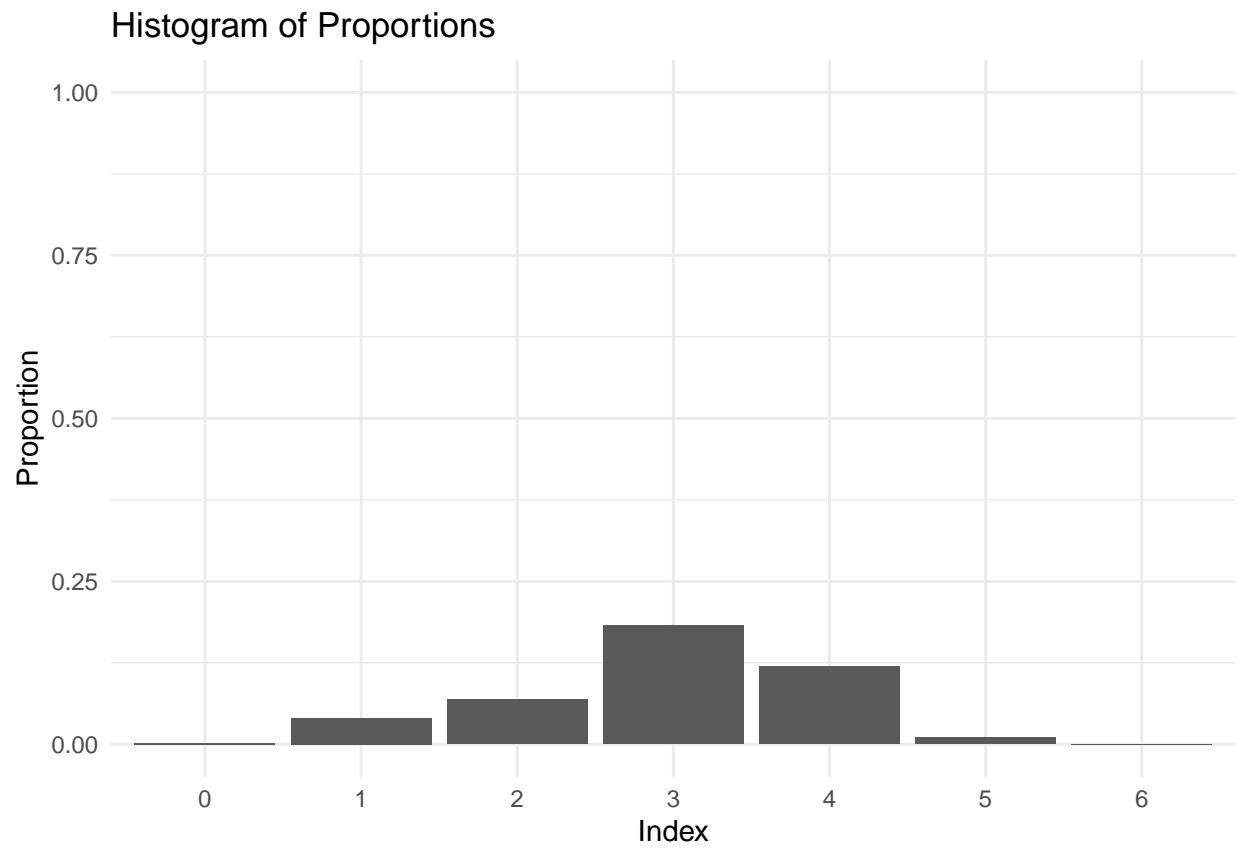


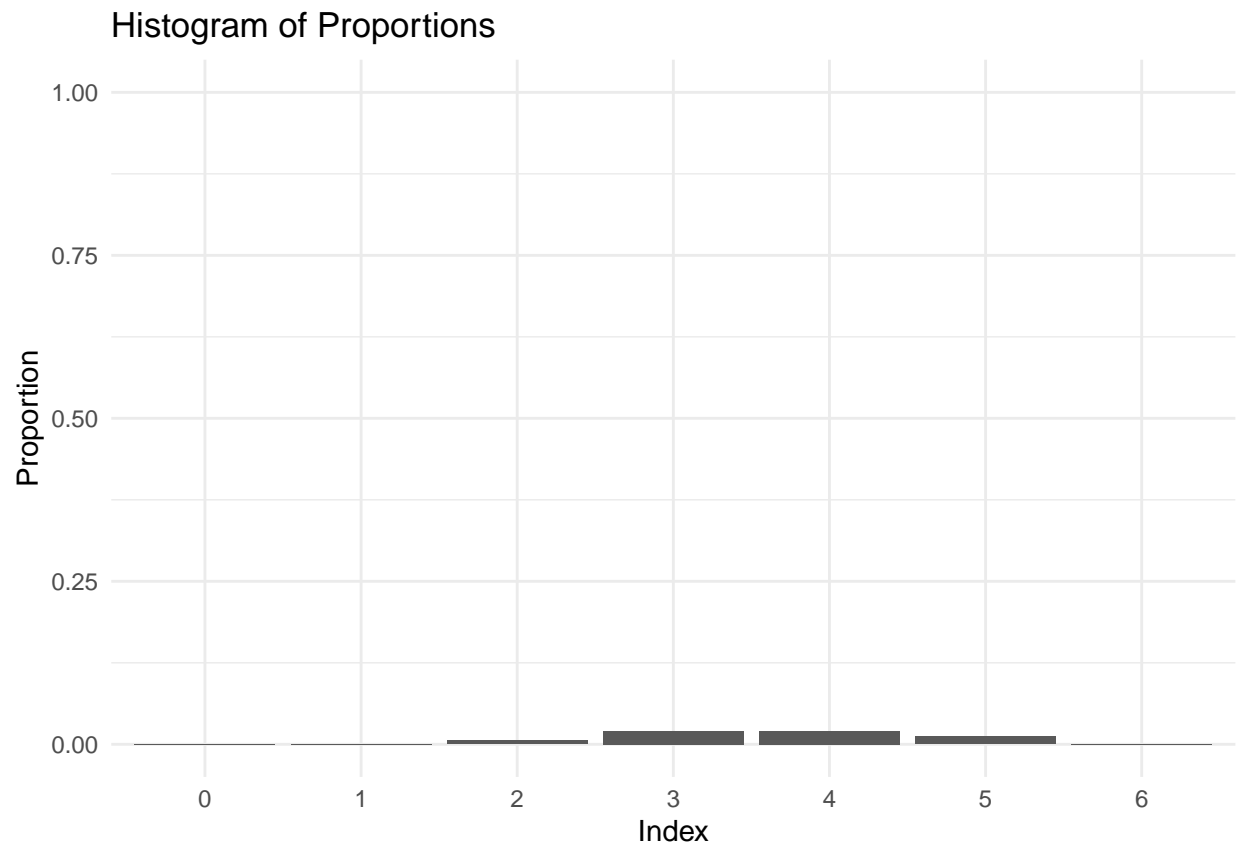
```
acc_data = as.matrix(  
  t(table(new_clusters, as.integer(df[test_ids, ]$score))) / nrow(df[test_ids, ])  
)
```

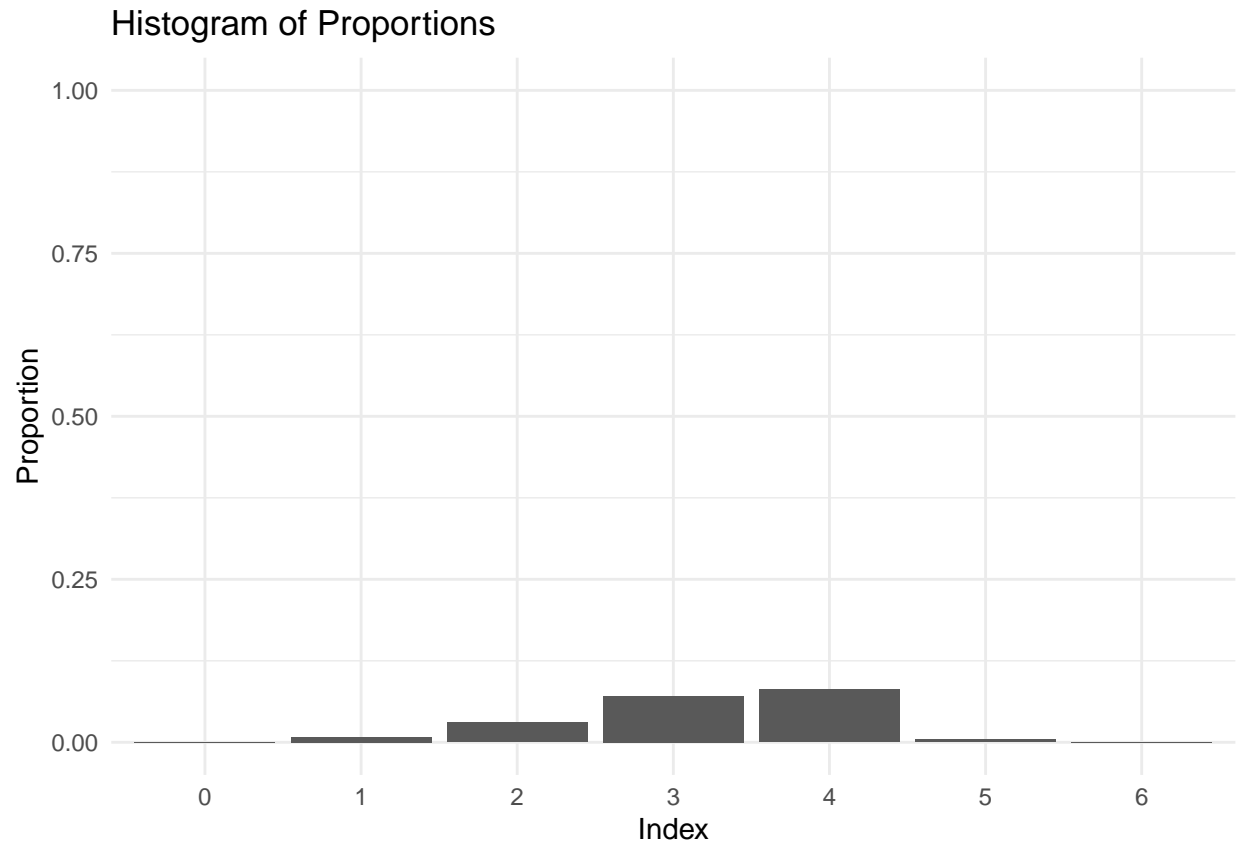
```
for (i in 1:5) {  
  prop_plot = ggplot(  
    data.frame(  
      Index = rownames(acc_data), Proportion = acc_data[, i]  
    ),  
    aes(x = Index, y = Proportion)) +  
    geom_bar(stat = "identity") +  
    theme_minimal() +  
    xlab("Index") +  
    ylab("Proportion") +  
    ggtitle("Histogram of Proportions") +  
    ylim(c(0, 1))  
    print(prop_plot)  
}
```









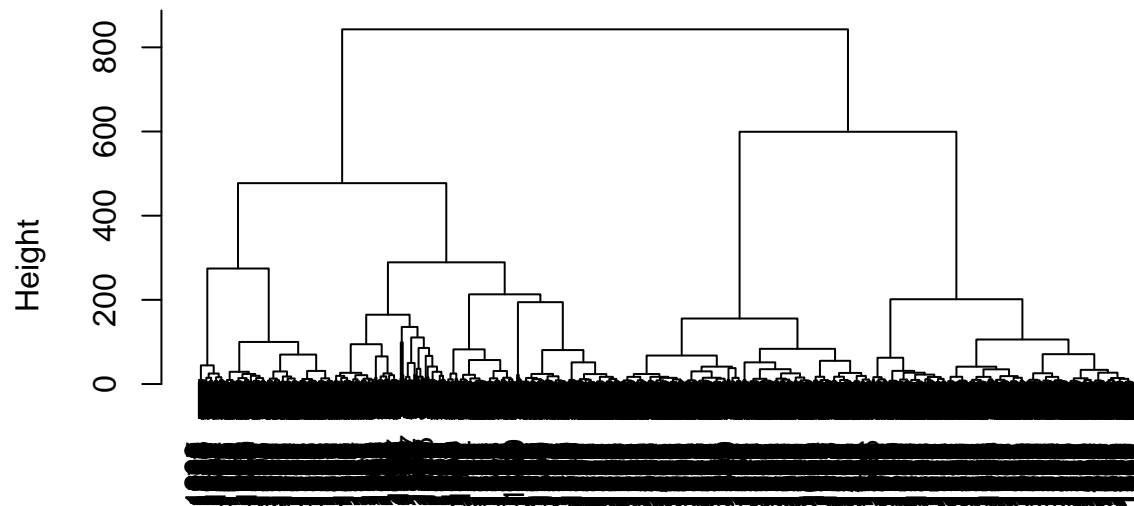


selected Ward's linkage due to cohesion and separability issues in the data. we select 7 clusters because there are 7 integer scores (after rounding half scores). The clustering structure in the data seems to support this number of clusters.

Hierarchial Clustering

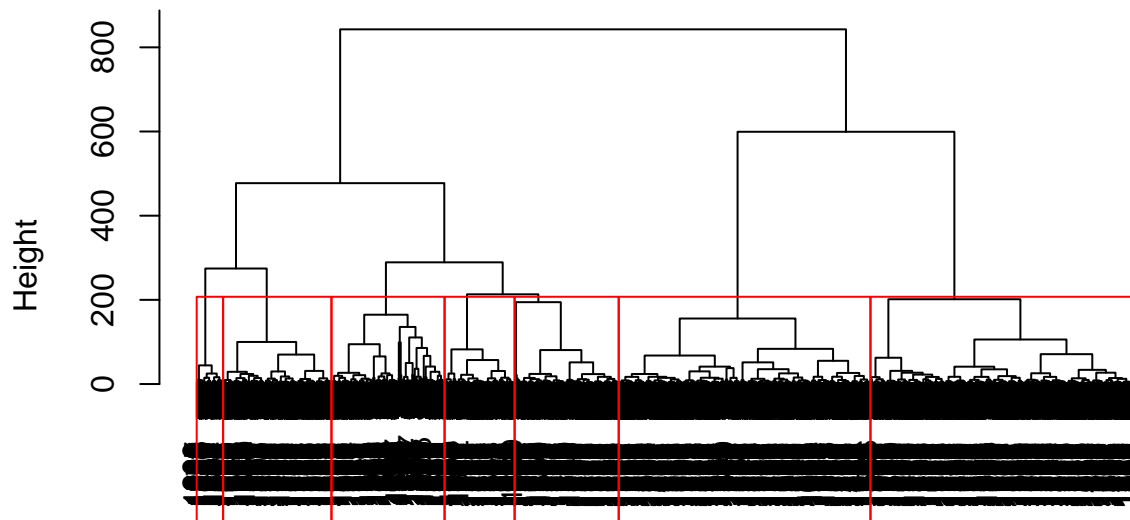
```
distance_matrix <- dist(train, method = "euclidean")  
hc <- hclust(distance_matrix, method = "ward.D")  
plot(hc, main = "Hierarchical Clustering with Complete Linkage",  
      xlab = "", sub = "", cex = 0.9)
```


Hierarchical Clustering with Complete Linkage



```
k <- 7 # Number of clusters
plot(hc, main = "Hierarchical Clustering with Complete Linkage",
     xlab = "", sub = "", cex = 0.9)
rect.hclust(hc, k = k, border = "red") # You can change the border color if you like
```

Hierarchical Clustering with Complete Linkage

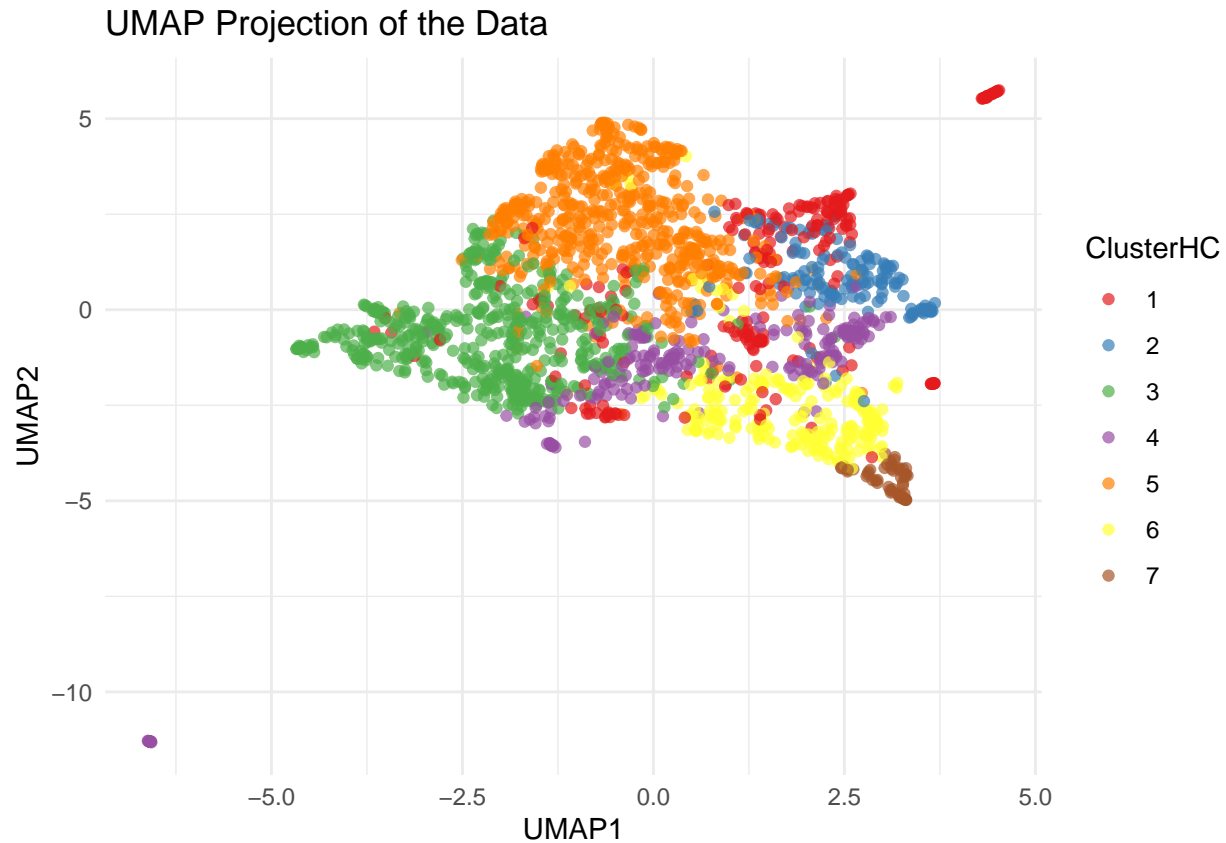


clustering structure in UMAP

```
clusters <- as.factor(cutree(hc, k = k))
umap_data$ClusterHC <- clusters

# Choose a palette
palette <- brewer.pal(n = k, name = "Set1") # Adjust 'name' as needed

ggplot(umap_data, aes(x = UMAP1, y = UMAP2, color = ClusterHC)) +
  geom_point(alpha = 0.7) +
  scale_color_manual(values = palette) +
  theme_minimal() +
  ggtitle("UMAP Projection of the Data")
```



investigation of the score distribution suggests that the underlying clustering structure of the data ... spread of data between kmeans clusters and scores. characterizations are as follows:

- cluster 1 tends to capture those who score 2.5 to 4.5
- cluster 2 tends to capture those who score 3.5 to 4.5
- cluster 3 tends to capture those who score 1.5 to 4.5
- cluster 4 tends to capture those who score 2.5 to 4
- cluster 5 tends to capture those who score 3.5 to 6
- cluster 6 tends to capture those who score 2.5 to 4.5
- cluster 7 tends to capture those who score 3 to 4.5

e.g., a user who scores around 5 is likely to be in cluster 5

```
t(table(umap_data$ClusterHC, df[-test_ids, ]$score)) / nrow(df[-test_ids, ])
```

```
##
##           1           2           3           4           5
## 0.5 0.0000000000 0.0000000000 0.0010116338 0.0005058169 0.0000000000
## 1   0.0005058169 0.0005058169 0.0075872534 0.0010116338 0.0005058169
## 1.5 0.0010116338 0.0005058169 0.0156803237 0.0060698027 0.0015174507
## 2   0.0060698027 0.0000000000 0.0171977744 0.0075872534 0.0010116338
## 2.5 0.0151745068 0.0005058169 0.0298431968 0.0171977744 0.0035407183
## 3   0.0232675771 0.0065756196 0.0480526050 0.0283257461 0.0085988872
## 3.5 0.0298431968 0.0126454224 0.0667678300 0.0323722812 0.0313606474
## 4   0.0278199292 0.0161861406 0.0536165908 0.0141628730 0.0596863935
```

```
## 4.5 0.0116337886 0.0202326758 0.0222559433 0.0040465352 0.0839656045
## 5 0.0035407183 0.0096105210 0.0060698027 0.0000000000 0.0419828022
## 5.5 0.0010116338 0.0070814365 0.0010116338 0.0000000000 0.0338897319
## 6 0.0010116338 0.0010116338 0.0000000000 0.0000000000 0.0136570561
##
##           6           7
## 0.5 0.0005058169 0.0000000000
## 1 0.0015174507 0.0000000000
## 1.5 0.0025290845 0.0010116338
## 2 0.0065756196 0.0000000000
## 2.5 0.0146686899 0.0010116338
## 3 0.0161861406 0.0045523520
## 3.5 0.0187152251 0.0055639858
## 4 0.0263024785 0.0055639858
## 4.5 0.0161861406 0.0055639858
## 5 0.0060698027 0.0025290845
## 5.5 0.0060698027 0.0025290845
## 6 0.0005058169 0.0000000000
```

to predict on new data in test...

```
# Function to calculate centroids of clusters
calculate_centroids <- function(data, clusters) {
  aggregate(data, by=list(cluster=clusters), FUN=mean)
}

# Function to predict the cluster of new data
predict_cluster <- function(new_data, train_data, clusters) {
  centroids <- calculate_centroids(train_data, clusters)
  # Remove the cluster column
  centroids <- centroids[, -1]

  # Function to find nearest centroid
  find_nearest_centroid <- function(point, centroids) {
    dists <- apply(centroids, 1, function(centroid) dist(rbind(centroid, point)))
    which.min(dists)
  }

  apply(new_data, 1, find_nearest_centroid, centroids = centroids)
}

new_clusters <- predict_cluster(test, train, clusters)
```

```
acc_data = as.matrix(
  t(table(new_clusters, as.integer(df[test_ids, ]$score))) / nrow(df[test_ids, ]))
)
```

```
for (i in 1:7) {
  prop_plot = ggplot(
    data.frame(
      Index = rownames(acc_data), Proportion = acc_data[, i]
    ),
    aes(x = Index, y = Proportion)) +
```

```

geom_bar(stat = "identity") +
theme_minimal() +
xlab("Index") +
ylab("Proportion") +
ggtitle("Histogram of Proportions") +
ylim(c(0, 1))
print(prop_plot)
}

```

