# 2803: Supervisor Briefing: Cluster Analysis Review

Date/Time    28 March 2018, 2PM – 3PM
Attendees    Arushi, Shubhangi, Tanushree

| Sr. No. | Notes | Actors | Follow up Action |
|---|---|---|---|
| 1 | Ran RFM Variables and Clusters with Prof. : Frequency Comments: <br> - Problem with current Frequency measure→ Average time b/w bookings. <br> - We have frequency values like 0. But then that doesn't make sense. <br> - Make definition clear first <br> - One Way: 3 months/6months/year→ Within that time period, how many times does the customer come? How many times the customer booked/the days since the customer joined the eatigo application. | Shubhangi | Remember to update the frequency variables |

| | | | |
|---|---|---|---|
| | - For everyone who were there from earlier, total number of bookings/365.<br>- If the base is not meaningful, you can consider dividing by one month or 3 months or the period you think meaningful | | |
| 2 | Recency:<br>Calculated as last date of year- Max booking date.<br>Variable passed | | |
| 3 | Monetary:<br>Currently calculating the average number of diners across the bookings.<br>Look at Total Number averaged number of diners | Shubhangi | Revise the monetary values |
| 4 | Notes on the Paper:<br>- Remember to define how you've calculated each<br>- And talk about the cut-off q | Shubhangi | Keep in mind |
| 5 | Notes on Standardization:<br>- Only use standardization if the data range is very big<br>Our data is not that wide. Therefore, we need to fix skewness not standardization. Therefore we should do transformation.<br>Log Transformation: | Arushi | Figure out the standardization method and complete by next meeting |

| | When you log, range becomes wider, then you standardize | | |
|---|---|---|---|
| 6 | Notes on using different data types: <br> - We have lots of variables of different types of data, some proportions, some continuous and numeric. Need to know how to deal with it. <br> Prof Comments <br> Option #1: <br> - Use RFM Analysis <br>   - Then Profile. Identify common booking behaviour <br> Option #2: <br> - Standardization doesn't always have to be Z-Score <br> - Look at the variable that has maximum values <br> - And can be scaled up | All | Decide on which option we would like to go ahead with |
| 7 | Using Clustering: <br> K-Means: <br> - Use <br> - Step #2: <br> Based on CCE (statistical method to choose the best number of clusters), decide on the number of clusters. | All | By next meeting have the distributions of variables standardized and transformed ready. Have the initial clustering done |

| | | | |
|---|---|---|---|
| | CCE should not be negative.<br>Two possibilities of negative:<br>- Lots of outliers:<br>If there's a cluster with only 1 variable, then it's wrong.<br>It should ideally have equal number in each cluster<br>- Data Skewed<br>Deciding on important variables:<br>- Look at cluster mean and cluster SD.<br>- But since it's standardized, it's not reliable to use it to interpret.<br>- Save the cluster.<br>- Go back to the main table, map out distribution of cluster inputs and cluster itself.<br>- Remove the cluster.<br>- Then at the menu on top of graphs, find data filter, then click cluster.<br>- Use that for interpretation.<br>Interpretation #2:<br>- Graph→ Parallel Plot, put all cluster inputs.<br>- Click Data Filter, Cluster. Then look at the parallel plot, colored based on clusters<br>- Uncheck Include and show. | | with parallel plots ready to show prof. In case any questions email prof with the doubts. |

| | Question 2. Outlier is important | | |
| --- | --- | --- | --- |
| | Clustering Technique #2: If there are outliers, use Normal Mixtures. If the clusters are equal sizes, then use this. Look at BIC and AIC and find the smallest. Rerun the iteration just in case and expand the range. Then check the AIC and BIC. If the small BIC is in one cluster and small AIC is another cluster then compare both the clusters. | | |