



2903: Supervisor Meeting

Date/Time 29 March 2018, 1:00PM – 1:30PM

Attendees : Shubhangi, Tanushree

Sr. No.	Notes	Actors	Follow up Action
1	<p>Running clustering variables after 2803 pointers:</p> <ul style="list-style-type: none"> - High Skew → Recently, Monetary, Frequency, <p>Don't Transform → Resto, cuisine variety.</p> <p>Freq too small to transform Frequency → Multiply by 1,000 to be able to conduct a transformation</p>	All	Work on clustering
2	<p>Notes on the Paper:</p> <ul style="list-style-type: none"> - Remember to define how you've calculated each - And talk about the cut-off q 	All	
3	<p>Notes on Standardization:</p> <ul style="list-style-type: none"> - Only use standardization if the data range is very big <p>Our data is not that wide. Therefore, we need to fix skewness not standardization. Therefore we should do transformation.</p> <p>Log Transformation: When you log, range becomes wider, then you don't need standardization</p>	All	

<p>4</p>	<p>Notes on using different data types:</p> <ul style="list-style-type: none"> - We have lots of variables of different types of data, some proportions, some continuous and numeric. Need to know how to deal with it. <p>Prof Comments</p> <p>Option #1:</p> <ul style="list-style-type: none"> - Use RFM Analysis <ul style="list-style-type: none"> - Then Profile. Identify common booking behaviour <p>Option #2:</p> <ul style="list-style-type: none"> - Standardization doesn't always have to be Z-Score - Look at the variable that has maximum values - And can be scaled up 	<p>All</p>	
<p>5</p>	<p>Using Clustering:</p> <p>K-Means:</p> <ul style="list-style-type: none"> - Use - Step #2: <p>Based on CCE (statistical method to choose the best number of clusters), decide on the number of clusters. CCE should not be negative. Two possibilities of negative:</p> <ul style="list-style-type: none"> - Lots of outliers: 	<p>All</p>	

	<p>If there's a cluster with only 1 variable, then it's wrong. It should ideally have equal number in each cluster</p> <ul style="list-style-type: none"> - Data Skewed <p>Deciding on important variables:</p> <ul style="list-style-type: none"> - Look at cluster mean and cluster SD. - But since it's standardized, it's not reliable to use it to interpret. - Save the cluster. - Go back to the main table, map out distribution of cluster inputs and cluster itself. - Remove the cluster. - Then at the menu on top of graphs, find data filter, then click cluster. - Use that for interpretation. <p>Interpretation #2:</p> <ul style="list-style-type: none"> - Graph→ Parallel Plot, put all cluster inputs. - Click Data Filter, Cluster. Then look at the parallel plot, colored based on clusters - Uncheck Include and show. <p>Question 2. Outlier is important</p>		
6	Clustering Technique #2:	All	

	<p>If there are outliers, use Normal Mixtures.</p> <p>If the clusters are equal sizes, then use this.</p> <p>Look at BIC and AIC and find the smallest.</p> <p>Rerun the iteration just in case and expand the range.</p> <p>Then check the AIC and BIC.</p> <p>If the small BIC is in one cluster and small AIC is another cluster then compare both the clusters.</p>		
--	--	--	--