

Paper 165-2011

## Eliminating Response Style Segments in Survey Data via Double Standardization Before Clustering

Murali Krishna Pagolu and Goutam Chakraborty, Oklahoma State University, Stillwater, OK

### ABSTRACT

Segmentation is the process of dividing a market into groups so that members within the groups are very similar with respect to their needs, preferences, and behaviors but members between groups are very dissimilar. Marketers often use clustering to find segments of respondents in data collected via surveys. However, such data often exhibits response styles of respondents. For example, if some respondents use only the extreme ends of scales for answering questions in a survey, the clustering method will identify that group as a unique segment, which cannot be used for segmentation.

In this paper, we first discuss the different data transformation methods that are commonly used before clustering. We then apply these different transformations to survey data collected from 959 customers of a business-to-business company. Both hierarchical and k-means clustering are then applied to the transformed data. Our results show that double-standardization performs better than other transformations in eliminating groups that identify response styles. We show how double-standardization can be achieved on any data using SAS<sup>®</sup> programs and SAS<sup>®</sup> macros.

### INTRODUCTION

All businesses rely on customers' continuous feedback in order to improve services or standards. To achieve this, customers' perceptions or attitudes towards the organization are measured and assessed. This is often achieved by means of questionnaire or survey methodology, which is usually a set of questions measuring customers' attitudes/perceptions. However, it is often difficult to record customers' true intentions and attitudes from surveys/questionnaires due to the bias effects introduced in the data (Bachman, 1984).

One of the main reasons for a bias in measuring customers' attitudes or perceptions is the response style behavior. Response styles in questionnaire/survey data is defined as the systematic inclination of responders to answer questions based on some unknown effect other than the content of the question (Paulhus, 1991, p. 17).

	Response Style	Completely disagree	Disagree	Somewhat disagree	Neither agree nor disagree	Somewhat agree	Agree	Completely agree
Truncated Scales	Optimal Responding	●	●	●	●	●	●	●
	Extreme Response Style (ERS)	●	○	○	○	○	○	●
	Response Range (RR)	●	●	○	○	○	●	●
	Mild Response Style	○	●	●	●	●	●	○
Social Styles	Midpoint Response (MPR)	○	○	○	●	○	○	○
	Acquiescence Response Style (ARS)	○	○	○	○	●	●	●
	Disacquiescence Response Style (DARS)	●	●	●	○	○	○	○
	Socially Desirable Responding (SDR)	●	●	●	○	●	●	●

Source: Vovici

Figure 1: Various types of response styles in attitudinal data

Various types of response styles are commonly found in attitudinal data. Figure 1 shows these on a sample Likert-type, seven-point scale measuring customers' attitudes towards a company (Henning, 2010). However, these types are not discussed in detail in this study.

Due to the nature of these response styles, the statistical properties of the data are easily distorted and the distributions of impacted attributes are no longer normal. The skewness, kurtosis, and bimodality coefficient values of attributes indicate the shape of their distribution. If the bimodality coefficient value is greater than 0.55, it indicates

that the distribution is no longer uniform and does not yield better results for clustering (Der & Everitt, p. 273, 2002). These values impact the performance of clustering procedures or any other unsupervised or supervised model in general. Various forms of standardization have been proposed and tested in last few years. However, very few of them stand out in terms of applicability and validity of data collected in diverse fields of study. Some of the key standardization forms are addressed in this paper.

Standardization refers to data transformation that involves correction of scores of either attributes or cases using either means or standard deviations or both. When the adjustment to the original scores is made using means, adjusted scores are produced as the mean across attributes or cases are subtracted from the original mean of either an attribute or a case (Fischer, 2004). The resulting adjusted score is often further adjusted using standard deviation. Standardization might be based on the adjustment of means of either cases, attributes, or both using either the mean across attributes for each case or across cases within an attribute, or both (Fischer, 2004). Thus the type of standardization to adjust raw scores relies on the type of the information (e.g., cases, attributes) and the kind of information used: means, standard deviations, or both, depending on the data and context of analysis (Fischer, 2004).

**RANGE STANDARDIZATION:** This form of standardization is most helpful when attributes are measured on different scales. Attributes with large values and a wide range of variation have significant effects on the final similarity measure. Hence, it is essential to make sure that each attribute is evenly constituted in the distance measurement by means of data standardization (Vickers, 2007). Range standardization is obtained by subtracting the minimum value of the attribute from each of its scores and dividing it by the range of the attribute (Maximum – Minimum). The resulting standardized data by this method produces attributes with values in the range of 0 and 1.

**CENTERING:** Centering refers to scores being adjusted using only the mean across the cases (Aiken & West, 1991). The attribute mean is subtracted from its original score. Standard deviation is not adjusted in this process.

**NORMAL STANDARDIZATION:** This form of standardization refers to correction of the scores using the attribute mean. The attribute mean is subtracted from its original score and then divided by the standard deviation (Howell, 1997). Hence, the resulting standardized score is the relative value of one specific case on one attribute relative to the value of other cases in that attribute. The mean across all the cases is zero and, assuming a normal distribution of responses, the resulting standard deviation will be equal to 1 (Fischer, 2004).

**ROW-CENTERING:** Row-centering refers to adjusting the scores using the mean of a case across all attributes (Fischer, 2004). The mean across all of the attributes for that particular case is subtracted from each individual's original measured value. Standard deviation is not adjusted in this process.

**ROW STANDARDIZATION:** This transformation refers to correction of scores for each case using the mean for that case across all attributes measured (Hofstede, 1980), which is then subtracted from each individual's original measured value. Thus the resulting standardized score is the relative value of the case on a variable corresponding to the other scores (Hicks, 1970). Hicks (1970) coined the term "ipsatization," which means the process of making the mean across all attributes for a case equal to zero. These newly assessed scores can be further adjusted for differences in the variation of the ratings around the mean by dividing the resulting score with the standard deviation across attributes for that case (Fischer, 2004).

**DOUBLE-CENTERING:** Double-centering refers to adjustments made using both the attribute mean and the mean of individual cases across all attributes (Fischer, 2004). Again, like centering and row-centering, this process does not involve further adjustment of scores using standard deviation.

**DOUBLE STANDARDIZATION:** In this process, the scores are first adjusted within the case and then the resulting scores are adjusted within the attribute (Fischer, 2004). Thus the mean for each case across attributes and the mean for each attribute across all cases will be zero. With the assumption that the raw data is normal, the correction using the standard deviation will produce standard deviations of 1 for both cases across attributes and attributes across cases. This combination of row and normal standardization was introduced by Leung and Bond (1989) and is named "Double Standardization." However, the order of using the row and normal standardization in combination is somewhat ambiguous. It is unclear whether the attributes should be standardized before the cases are standardized or vice versa. However, if data is standardized across attributes after standardizing across cases, then the properties of attributes are likely deformed and the true form of data is not available for study (SAS Institute Inc., p.989, 2004). The most commonly used standardization procedures are row standardization and double standardization (Fischer, 2004).

## METHODOLOGY

XYZ is a leading supplier of hydraulic and pneumatic products serving more than 50,000 customers in U.S. (The name of the company is kept anonymous for privacy reasons.) We used survey data collected from 959 customers based on their perception of important factors in selecting a supplier for the hydraulic and pneumatic products (Table 1). XYZ would like to segment their customers based on that perception. This study is focused on going through various forms of transformations for the given data and then applying hierarchical clustering and k-means procedures to find clean and stable cluster solutions feasible in the business world to segment and profile the customers. Throughout the clustering procedures performed in this study, the average linkage method is used for fair comparison across the transformed forms of data. This method is based on average similarity of all observations or cases within a cluster and is likely to produce clusters with small within-cluster variation, which is less impacted by the presence of outliers in the data.

How important are the following issues to customers in choosing a supplier for hydraulic, pneumatic, and related products?	Attribute	Scale	Not at all Important	Extremely Important
1. The reliability of the supplier	reliab	9 point	1	9
2. The timeliness of the deliveries by the supplier	time	9 point	1	9
3. The availability of a large breadth of products to choose from	av_br	9 point	1	9
4. The availability of well documented technical specification	av_spec	9 point	1	9
5. The price of products	price	9 point	1	9
6. The credit policy of the supplier	credit	9 point	1	9
7. The availability of electronic payment/debit option	av_pay	9 point	1	9
8. The return policy of the supplier	return	9 point	1	9
9. The warranty coverage provided by the supplier	warranty	9 point	1	9
10. The ability to talk directly to a salesperson about your needs	talk_dir	9 point	1	9

Table 1: Important factors for customers in selecting a supplier for the hydraulic and pneumatic products

## HIERARCHICAL CLUSTERING

Performing simple hierarchical clustering using the average linkage method with no data transformations on the original data set produced the results below. The summary statistics, including skewness, kurtosis, and bimodality coefficients, of all rating attribute data in the data set are displayed in the result (Figure 1.a). The mean values vary between attributes, as do the standard deviations. Few of the attributes have bimodality coefficients greater than 0.55, which indicates non-uniform distributions. Skewness and kurtosis values for a few variables are much higher, indicating nonsymmetrical distributions. Figure 1.b shows the last 10 generations of hierarchical clustering on the original data set; it is evident that few observations join the clusters very late. This shows that the original data set should not be used straight away for clustering procedures and that it definitely needs some sort of data transformation.

Variable	Mean	Std Dev	Skewness	Kurtosis	Bimodality	Cluster History										
reliab	8.4911	0.9045	-2.9934	14.8482	0.5578	NCL	Clusters Joined	FREQ	SPRSQ	RSQ	ERSQ	CCC	PSF	PST2	Norm RMS Dist	T i e
time	8.4828	0.8384	-2.4000	9.7831	0.5284	10	CL11 CL38	919	0.0048	.121	578	-76	14.5	5.2	1.1082	
av_br	6.8436	1.6944	-0.5396	-0.1078	0.4450	9	CL10 CL14	926	0.0099	.111	565	-76	14.8	10.7	1.2148	
av_spec	7.6569	1.4238	-1.1007	0.9209	0.5627	8	CL22 CL24	19	0.0048	.106	552	-75	16.1	6.0	1.2248	
price	7.6830	1.4331	-1.1849	1.5503	0.5272	7	CL9	790	0.0026	.103	535	-61	18.3	2.7	1.2946	
credit	6.0480	2.1489	-0.5075	-0.2898	0.4624	6	CL8 CL13	24	0.0075	.096	515	-60	20.2	7.4	1.3186	
av_pay	3.4025	2.2938	0.5915	-0.6496	0.5720	5	CL7 CL20	930	0.0069	.089	490	-60	23.3	7.4	1.3527	
return	6.8780	1.7979	-0.7735	0.2252	0.4941	4	CL5 CL6	954	0.0471	.042	455	-55	13.9	49.4	1.4177	
warranty	7.7414	1.4749	-1.3266	1.6398	0.5936	3	CL19	132	0.0031	.039	403	-44	19.2	3.3	1.5016	
talk_dir	8.2711	1.2017	-2.6377	9.7391	0.6242	2	CL3	316	0.0049	.034	306	-30	33.5	3.0	1.8725	
						1	CL4 CL2	959	0.0338	.000	.000	0.00	.	33.5	2.1625	

Figure 1.a

Figure 1.b

**RANGE STANDARDIZATION:** Results of hierarchical clustering show no improvement in skewness, kurtosis, or bimodality coefficient values of the rating attribute data when compared to the results from hierarchical clustering on the original data set (Figures 1.a, 2.a). Figure 2.b shows the last 10 generations of hierarchical clustering, and we can see that the observations still join the clusters very late. Hence, we see that centering has proven useless with regard to dealing with distributions of attributes or the outliers in the data.

Variable	Mean	Std Dev	Skewness	Kurtosis	Bimodality
reliab	0.9364	0.1131	-2.9934	14.8482	0.5578
time	0.9261	0.1198	-2.4000	9.7831	0.5284
av_br	0.7304	0.2118	-0.5396	-0.1078	0.4450
av_spec	0.8081	0.2034	-1.1007	0.9209	0.5627
price	0.8354	0.1791	-1.1849	1.5503	0.5272
credit	0.6310	0.2686	-0.5075	-0.2898	0.4624
av_pay	0.3003	0.2867	0.5915	-0.6496	0.5720
return	0.7347	0.2247	-0.7735	0.2252	0.4941
warranty	0.8427	0.1844	-1.3266	1.6398	0.5936
talk_dir	0.9089	0.1502	-2.6377	9.7391	0.6242

Figure 2.a

Cluster History										
NCL	Clusters Joined	FREQ	SPRSQ	RSQ	ERSQ	CCC	PSF	PST2	Norm RMS Dist	T i e
10	CL76 CL15	6	0.0030	.109	.572	-.77	12.9	4.2	1.1848	
9	CL11 790	929	0.0025	.106	.560	-.76	14.1	2.7	1.284	
8	CL9 CL10	935	0.0101	.096	.547	-.76	14.5	10.7	1.3127	
7	CL12 CL119	14	0.0044	.092	.530	-.69	16.1	4.7	1.3194	
6	CL8 CL22	938	0.0068	.085	.510	-.61	17.7	7.2	1.3431	
5	CL20 132	4	0.0030	.082	.485	-.60	21.3	3.2	1.4912	
4	CL6 CL7	952	0.0359	.046	.451	-.54	15.4	37.4	1.4945	
3	CL4 CL19	954	0.0064	.040	.399	-.44	19.8	6.5	1.5148	
2	CL5 316	5	0.0050	.035	.303	-.30	34.5	3.1	1.8906	
1	CL3 CL2	959	0.0348	.000	.000	0.00	.	34.5	2.1854	

Figure 2.b

**CENTERING:** When compared to the results from hierarchical clustering on the original data set (Figures 1.a, 3.a), the results of hierarchical clustering using centered data show no improvement in skewness, kurtosis, or bimodality coefficient values of the rating attributes. Figure 3.b shows the last 10 generations of hierarchical clustering; we can see that the observations still join the clusters very late. Hence, we find that centering has proven useless with regard to dealing with distributions of attributes or the outliers in the data.

Variable	Mean	Std Dev	Skewness	Kurtosis	Bimodality
reliab	0	0.9045	-2.9934	14.8482	0.5578
time	0	0.8384	-2.4000	9.7831	0.5284
av_br	0	1.6944	-0.5396	-0.1078	0.4450
av_spec	0	1.4238	-1.1007	0.9209	0.5627
price	0	1.4331	-1.1849	1.5503	0.5272
credit	0	2.1489	-0.5075	-0.2898	0.4624
av_pay	0	2.2938	0.5915	-0.6496	0.5720
return	0	1.7979	-0.7735	0.2252	0.4941
warranty	0	1.4749	-1.3266	1.6398	0.5936
talk_dir	0	1.2017	-2.6377	9.7391	0.6242

Figure 3.a

Cluster History										
NCL	Clusters Joined	FREQ	SPRSQ	RSQ	ERSQ	CCC	PSF	PST2	Norm RMS Dist	T i e
10	CL11 CL38	919	0.0048	.121	.578	-.76	14.5	5.2	1.1082	
9	CL10 CL14	926	0.0099	.111	.565	-.76	14.8	10.7	1.2148	
8	CL22 CL24	19	0.0048	.106	.552	-.75	16.1	6.0	1.2248	
7	CL9 790	927	0.0026	.103	.535	-.61	18.3	2.7	1.2946	
6	CL8 CL13	24	0.0075	.096	.515	-.60	20.2	7.4	1.3186	
5	CL7 CL20	930	0.0069	.089	.490	-.60	23.3	7.4	1.3527	
4	CL5 CL6	954	0.0471	.042	.455	-.55	13.9	49.4	1.4177	
3	CL19 132	4	0.0031	.039	.403	-.44	19.2	3.3	1.5016	
2	CL3 316	5	0.0049	.034	.306	-.30	33.5	3.0	1.8725	
1	CL4 CL2	959	0.0338	.000	.000	0.00	.	33.5	2.1625	

Figure 3.b

**NORMAL STANDARDIZATION:** From the results of hierarchical clustering, we see no change in skewness, kurtosis, or bimodality coefficient values of the rating attribute data when compared to the results from hierarchical clustering on the original data set (Figures 1.a, 4.a). Figure 4.b shows the last 10 generations of hierarchical clustering; we can see the observations join clusters very late. Thus, we find that this method proves less useful with regard to dealing with distributions of attributes or the outliers in the data.

Variable	Mean	Std Dev	Skewness	Kurtosis	Bimodality
reliab	0	1.0000	-2.9934	14.8482	0.5578
time	0	1.0000	-2.4000	9.7831	0.5284
av_br	0	1.0000	-0.5396	-0.1078	0.4450
av_spec	0	1.0000	-1.1007	0.9209	0.5627
price	0	1.0000	-1.1849	1.5503	0.5272
credit	0	1.0000	-0.5075	-0.2898	0.4624
av_pay	0	1.0000	0.5915	-0.6496	0.5720
return	0	1.0000	-0.7735	0.2252	0.4941
warranty	0	1.0000	-1.3266	1.6398	0.5936
talk_dir	0	1.0000	-2.6377	9.7391	0.6242

Figure 4.a

Cluster History										
NCL	Clusters Joined	FREQ	SPRSQ	RSQ	ERSQ	CCC	PSF	PST2	Norm RMS Dist	T i e
10	CL15 CL17	922	0.0158	.167	.531	-.71	21.1	18.5	1.3064	
9	CL10 480	923	0.0028	.164	.519	-.65	23.3	3.2	1.3271	
8	CL12 CL13	21	0.0085	.156	.505	-.65	25.0	7.3	1.4299	
7	132 CL35	3	0.0027	.153	.490	-.64	28.6	3.2	1.451	
6	CL9 CL11	931	0.0187	.134	.471	-.59	29.5	21.4	1.4617	
5	CL6 126	932	0.0038	.130	.447	-.52	35.8	4.2	1.4975	
4	CL5 CL8	953	0.0501	.080	.414	-.49	27.8	55.0	1.515	
3	CL39 CL7	5	0.0047	.076	.371	-.39	39.1	3.3	1.6261	
2	CL3 316	6	0.0058	.070	.287	-.26	71.8	2.6	2.0485	
1	CL4 CL2	959	0.0698	.000	.000	0.00	.	71.8	2.6927	

Figure 4.b

**ROW-CENTERING:** Based on the results of hierarchical clustering, we see significant improvement in the skewness, kurtosis, and bimodality coefficient values of the rating attribute data when compared to the results from hierarchical clustering on the original data set and also in comparison with the normal standardized data set (Figures 1.a, 4.a, and 5.a). Figure 5.b shows the last 10 generations of hierarchical clustering; we find that the observations still join the clusters very late. Hence, we find that though row-centering has better skewness, kurtosis, and bimodality coefficient values than that of the original data and within-case standardized data, it is still not very effective in handling the outliers in the data.

Variable	Mean	Std Dev	Skewness	Kurtosis	Bimodality
reliab	1.3413	0.9311	-0.2101	1.1663	0.2500
time	1.3330	0.9099	0.0657	0.6397	0.2752
av_br	-0.3063	1.3354	-0.4791	0.8129	0.3217
av_spec	0.5071	1.0858	-0.3410	0.4134	0.3261
price	0.5332	1.1665	-0.4323	0.5706	0.3315
credit	-1.1019	1.6591	-0.6492	0.6618	0.3872
av_pay	-3.7473	1.9658	0.4053	-0.1287	0.4042
return	-0.2718	1.2847	-0.8627	1.4708	0.3893
warranty	0.5916	1.0310	-0.6974	1.0815	0.3633
talk_dir	1.1213	1.0637	-1.3245	6.5327	0.2887

Figure 5.a

Cluster History											
NCL	Clusters Joined	FREQ	SPRSQ	RSQ	ERSQ	CCC	PSF	PST2	Norm RMS Dist	T	i
10	CL11 CL15	934	0.0236	.086	.528	-.76	9.9	25.1	1.2965		
9	CL10	480	0.0026	.083	.515	-.76	10.8	2.7	1.3099		
8	CL9 CL20	941	0.0110	.072	.499	-.69	10.5	11.4	1.3106		
7	CL8 CL33	947	0.0132	.059	.480	-.70	9.9	13.5	1.349		
6	126 CL16	4	0.0023	.057	.457	-.69	11.4	1.8	1.3664		
5	CL7 181	948	0.0036	.053	.428	-.61	13.4	3.6	1.4811		
4	CL5 CL13	951	0.0106	.042	.385	-.53	14.1	10.7	1.6048		
3	CL4 CL6	955	0.0133	.029	.327	-.43	14.3	13.3	1.6305		
2	CL3 CL14	958	0.0206	.009	.217	-.32	8.2	20.3	2.0531		
1	CL2	316	0.0085	.000	.000	0.00	.	8.2	2.1396		

Figure 5.b

**ROW STANDARDIZATION:** Based on the results of hierarchical clustering, we see significant improvement in the skewness, kurtosis, and bimodality coefficient values of the rating attribute data when compared to the results from hierarchical clustering on the original data set (Figures 1.a, 6.a).

Variable	Mean	Std Dev	Skewness	Kurtosis	Bimodality
reliab	0.6748	0.4154	-0.6260	2.5282	0.2513
time	0.6700	0.4114	-0.5490	1.6631	0.2785
av_br	-0.1751	0.6854	-0.4756	-0.0252	0.4109
av_spec	0.2410	0.5785	-0.8443	1.6421	0.3682
price	0.2597	0.6001	-0.6223	0.5903	0.3854
credit	-0.5327	0.7613	-0.0845	-0.7156	0.4391
av_pay	-1.8744	0.7407	1.3681	2.8193	0.4927
return	-0.1421	0.6234	-0.5128	-0.1453	0.4410
warranty	0.3070	0.4998	-0.8598	0.7726	0.4599
talk_dir	0.5719	0.4844	-1.4510	4.6343	0.4063

Figure 6.a

Cluster History											
NCL	Clusters Joined	FREQ	SPRSQ	RSQ	ERSQ	CCC	PSF	PST2	Norm RMS Dist	T	i
10	654 674	2	0.0017	.113	.488	-.75	13.4	.	1.2664		
9	CL19 CL35	9	0.0055	.107	.473	-.75	14.3	6.3	1.3496		
8	CL33 CL21	7	0.0038	.103	.457	-.69	15.7	4.4	1.3571		
7	CL20 CL10	4	0.0024	.101	.438	-.68	17.8	1.6	1.3677		
6	CL18 CL11	935	0.0282	.073	.413	-.65	15.0	30.2	1.4394		
5	CL6 CL8	942	0.0156	.057	.380	-.58	14.5	16.1	1.4441		
4	CL5 CL9	951	0.0205	.037	.338	-.52	12.1	20.8	1.4811		
3	CL4 402	952	0.0038	.033	.272	-.42	16.3	3.8	1.5168		
2	CL27 CL7	7	0.0054	.028	.172	-.29	27.2	3.7	1.5585		
1	CL3 CL2	959	0.0276	.000	.000	0.00	.	27.2	1.8011		

Figure 6.b

Figure 6.b shows the last 10 generations of hierarchical clustering; we see that the observations still join the clusters very late. Hence, we find that though within-case standardization shows a significant impact in improving the skewness, kurtosis, and bimodality coefficient values of the original data, it is not so effective in handling the outliers in the data.

**DOUBLE-CENTERING:** Comparing the results of hierarchical clustering with that of other forms of standardized data, we see that the skewness and kurtosis are better than most forms of transformed data. However, the bimodality coefficient values of the rating attribute data are the best among all forms of transformed data or the original data itself (Figures 1.a, 4.a, 5.a, 6.a, and 7.a). Observing the last 10 generations of the results from hierarchical clustering in Figure 7.b, we see that there are some cases that join the clusters very late in the final stages of clustering. Thus, we can infer that double-centered data helps reduce the skewness, kurtosis, and bimodality coefficient values of data compared to various other forms of transformed data.

Variable	Mean	Std Dev	Skewness	Kurtosis	Bimodality
reliab	0	0.9311	-0.2101	1.1663	0.2500
time	0	0.9099	0.0657	0.6397	0.2752
av_br	0	1.3354	-0.4791	0.8129	0.3217
av_spec	0	1.0858	-0.3410	0.4134	0.3261
price	0	1.1665	-0.4323	0.5706	0.3315
credit	0	1.6591	-0.6492	0.6618	0.3872
av_pay	0	1.9658	0.4053	-0.1287	0.4042
return	0	1.2847	-0.8627	1.4708	0.3893
warranty	0	1.0310	-0.6974	1.0815	0.3633
talk_dir	0	1.0637	-1.3245	6.5327	0.2887

Figure 7.a

Cluster History											
NCL	Clusters Joined	FREQ	SPRSQ	RSQ	ERSQ	CCC	PSF	PST2	Norm RMS Dist	T	i
10	CL12 CL15	936	0.0214	.083	.528	-.77	9.5	22.6	1.2901		
9	CL10	480	0.0026	.080	.515	-.76	10.4	2.7	1.3107		
8	CL9 CL19	943	0.0110	.069	.499	-.70	10.1	11.4	1.3115		
7	CL8 CL34	949	0.0132	.056	.480	-.70	9.4	13.5	1.3498		
6	CL11 132	3	0.0020	.054	.457	-.69	10.9	1.2	1.3632		
5	172 CL57	3	0.0026	.051	.428	-.61	12.9	3.7	1.4327		
4	CL7 CL5	952	0.0104	.041	.385	-.53	13.6	10.5	1.6351		
3	CL4 CL6	955	0.0118	.029	.327	-.43	14.3	11.8	1.724		
2	CL3 CL14	958	0.0206	.009	.217	-.32	8.2	20.3	2.0531		
1	CL2	316	0.0085	.000	.000	0.00	.	8.2	2.1396		

Figure 7.b

**DOUBLE STANDARDIZATION:** From the results of hierarchical clustering, we see that the skewness, kurtosis, and bimodality coefficient values of the rating attribute data are best when compared to the results from hierarchical clustering on various other forms of transformed data or the original data set itself (Figures 1.a, 4.a, 5.a, 6.a, 7.a, and 8.a).

Variable	Mean	Std Dev	Skewness	Kurtosis	Bimodality
reliab	-0.0464	0.8761	-0.4115	-0.4945	0.4650
time	-0.0437	0.9123	-0.4868	-0.3733	0.4693
av_br	0.0235	1.0024	-0.3666	-0.6322	0.4772
av_spec	0.00521	0.8953	-0.5137	-0.2853	0.4640
price	-0.00599	0.9547	-0.5666	-0.4428	0.5147
credit	0.0424	0.9921	-0.4564	-0.3944	0.4621
av_pay	0.00114	1.2564	-0.00249	-0.6765	0.4286
return	0.0336	0.8769	-0.5438	-0.2434	0.4685
warranty	0.0174	0.7989	-0.6052	-0.1205	0.4729
talk_dir	-0.0271	0.8443	-0.7106	0.1654	0.4740

Figure 8.a

Cluster History										
NCL	Clusters Joined	FREQ	SPRSQ	RSQ	ERSQ	CCC	PSF	PST2	Norm RMS Dist	T i e
10	CL18 CL30	62	0.0119	.399	.468	-17	70.1	21.9	0.9392	
9	CL14 CL38	102	0.0084	.391	.453	-16	76.2	13.3	0.9403	
8	CL22 CL28	153	0.0268	.364	.437	-19	77.8	52.9	0.9443	
7	CL12 CL9	314	0.0362	.328	.418	-22	77.4	54.4	0.947	
6	CL15 CL16	239	0.0396	.288	.394	-24	77.2	63.1	0.9558	
5	CL7 CL10	376	0.0260	.262	.363	-22	84.8	33.7	0.978	
4	CL13 CL11	191	0.0240	.238	.318	-17	99.6	34.4	0.9895	
3	CL4 CL6	430	0.0497	.189	.258	-14	111	61.8	0.9962	
2	CL5 CL8	529	0.0747	.114	.164	-11	123	94.3	1.026	
1	CL2 CL3	959	0.1139	.000	.000	0.00	.	123	1.0567	

Figure 8.b

Observing the last 10 generations of the results from hierarchical clustering in Figure 8.b, we see that there are no instances where single observations join the clusters late. Hence, we can conclude that the double-standardized data proves to be better than the original data or other forms of transformed data in effectively handling the skewness, kurtosis, or bimodality coefficient of the rating attribute data and outliers in the data without actually trimming down the data. Based on the local peaking of Pseudo F and Pseudo T-squared plots, the number of clusters suggested by the results of hierarchical clustering is 6 or 7 (Figure 8.c). The code in the Appendix labeled “code for creating the double standardized data and applying hierarchical clustering” shows the method for creating the required double standardized data and applying hierarchical clustering.

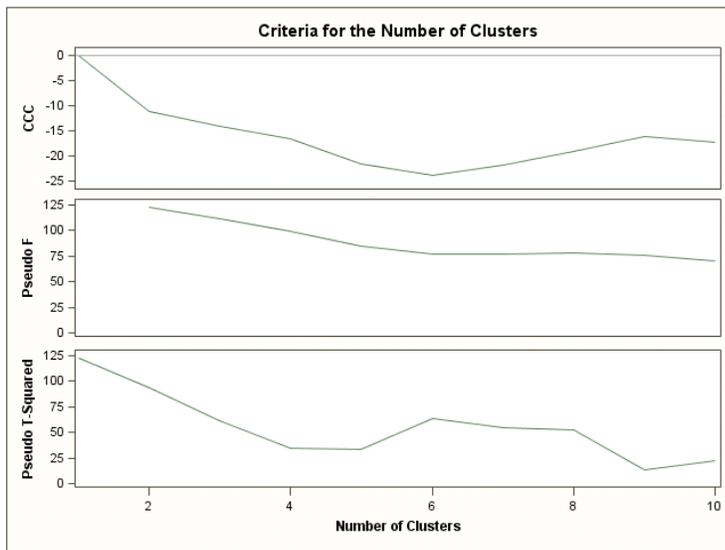


Figure 8.c

### K-MEANS CLUSTERING

Below are the various forms of transformed data used in k-means algorithms for the clustering procedure. The resulting output after running the k-means algorithm for each form of the transformed data using an average linkage method (until otherwise stated) is presented below.

**UNTRANSFORMED DATA:** When the original data is used, 5 clusters are formed but one cluster shows one observation (See Figure 9.a).

**RANGE STANDARDIZED DATA:** The range standardized form of survey data produced an 11-cluster solution, which is not feasible in real-world applications and with the size of smaller clusters with respect to the overall size of the data (See Figure 9.b).

**CENTERED DATA:** The centered form of survey data produced a 5-cluster solution and it appears worse compared to the clusters generated by original data. Also, we still have one cluster with only one observation (See Figure 9.c).

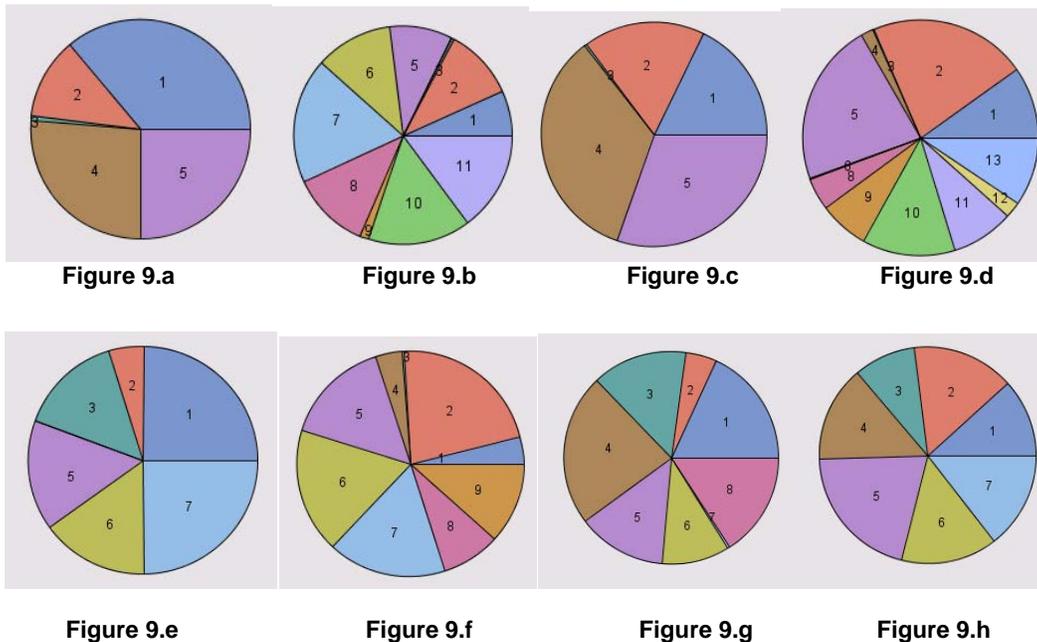
**NORMAL STANDARDIZED DATA:** The normal standardized form of survey data produced a 13-cluster solution, which is not feasible in real-world applications and with the size of smaller clusters with respect to the overall size of the data (See Figure 9.d).

**ROW-CENTERED DATA:** The row centered form of survey data produced a 7-cluster solution, which is better than clusters generated by the earlier forms of data (See Figure 9.e).

**ROW STANDARDIZED DATA:** Using row standardized data produces a 9-cluster solution, but again with one observation each in two of the clusters. This is not a feasible solution given the size of the smaller clusters relative to the size of the larger clusters or the overall data (See Figure 9.f).

**DOUBLE-CENTERED DATA:** Using double-centered data produces an 8-cluster solution, but again with one observation each in two of the clusters. This is not a feasible solution given the size of the smaller clusters relative to the size of the larger clusters or the overall data (See Figure 9.g).

**DOUBLE-STANDARDIZED DATA:** Using double-standardization produces a 7-cluster solution where the cluster sizes are comparable and significant in relation to the overall data. This is similar to the clusters generated from row-centered data. However, these cluster solutions appear more feasible (see Figure 9.h) as the frequency of observations in each of the clusters indicates a more stable and cleaner cluster solution than row-centered data (see Figure 9.e).



The segments formed by running the k-means procedure on various forms of the original data are then assessed for response styles. The code in the Appendix labeled **“SAS Code nodes to compare and assess the response style impact of data on customer segments”** is used to calculate the overall mean of the perception attributes in addition to the individual segment level mean for these attributes for each form of the transformed data as shown in Figure 10.

	reliab	time	av_br	av_spec	price	credit	av_pay	return	warranty	talk_dir
Mean	8.49	8.48	6.84	7.66	7.68	6.05	3.40	6.88	7.74	8.27
Segment Id										
1	Mean 8.23	8.30	6.11	7.40	6.46	3.37	2.23	4.70	6.44	8.00
2	Mean 8.14	8.13	5.23	5.86	7.42	5.51	2.48	5.99	6.56	7.46
3	Mean 2.75	3.75	6.00	4.00	2.25	5.25	7.25	3.00	3.00	2.50
4	Mean 8.69	8.68	7.76	8.27	8.26	7.37	5.90	7.83	8.41	8.65
5	Mean 8.72	8.66	7.20	8.19	8.01	6.44	1.76	7.67	8.49	8.55
6	Mean 2.00	2.00	1.00	9.00	1.00	9.00	4.00	3.00	8.00	9.00

**Figure10: Overall mean and segment-wise mean of perception attributes in the original survey data**

Figure 11 shows the individual segment level mean and overall mean of perception attributes resulting from running the code in Appendix labeled **“SAS Code nodes to compare and assess the response style impact of data on customer segments”** using the double standardized data. If eight out of the ten perception attributes in a segment show segment level means much higher or lower when compared to the overall mean, then that segment is regarded as a problematic segment still exhibiting the response styles.

		reliab	time	av_br	av_spec	price	credit	av_pay	return	warranty	talk_dir
Mean		-0.05	-0.04	0.02	0.01	-0.01	0.04	0.00	0.03	0.02	-0.03
Segment Id											
1	Mean	0.08	0.22	0.55	0.49	0.53	-1.55	-0.72	-0.02	0.28	0.14
2	Mean	0.21	0.18	-1.35	-0.37	0.38	0.24	-0.66	0.50	0.54	0.33
3	Mean	0.31	0.54	-0.05	-0.38	0.44	0.36	0.50	-0.00	-0.10	-1.62
4	Mean	0.67	0.66	0.11	-0.16	0.24	-0.02	-0.03	-1.00	-0.99	0.52
5	Mean	-0.83	-1.00	0.29	-0.15	-0.03	0.38	1.52	0.18	-0.12	-0.25
6	Mean	-0.62	-0.59	0.76	0.24	0.22	0.95	-1.31	0.39	0.17	-0.21
7	Mean	0.32	0.32	-0.12	0.39	-1.57	-0.40	0.16	0.08	0.38	0.44

Figure 11: Overall mean and segment-wise mean of perception attributes in the double standardized data

We can see that the values of individual segment-level means are so small and move on either side of 0, which makes them little confusing to compare and assess. In order to make a better comparison, we used the original values of these attributes rather than the double-standardized values. Please see the code available in the Appendix labeled “Merging k-means output data with original attribute values.”

### TRAFFIC LIGHTING THE MULTI-SHEET MICROSOFT EXCEL WORKBOOK

We incorporated this method to produce the desired mean statistics in an html output format, and we used the code from the paper, which explains how to perform multi-sheet traffic lighting using Microsoft Excel workbooks (DelGobbo, 2010). This gives better a visual presentation of how the individual segment level means for each perception attribute falls either above or below the overall means for that attribute. Please see the code available in the Appendix labeled “Preparing data for multi-sheet traffic lighting on Microsoft Excel workbooks.” Figure 12 shows the traffic lighted worksheet that indicates the range flags and the colors that represent them. Figure 13 shows the traffic lighted worksheet that shows the individual segment means of each attribute in its original values and whether they are above or below the overall mean of that attribute.

**Note:** Though the code is provided in the Appendix for creating the multi-sheet traffic lighting, you will still need to follow the original paper (DelGobbo, 2010) to create the modified styles and set up the ExcelXP tagset, ODS environment before using this code.

Range Flags	Color
Low	
High	

Figure 12: Range flags worksheet of multi-sheet Excel workbook generated

If eight or more cells within a segment have the same color, then it is probably exhibiting a response style indicating a problematic segment. Please see the code available in the Appendix labeled “code for generating the segment means in excel workbook and traffic lighting.”

Segment Id	reliab_Mean	time_Mean	av_br_Mean	av_spec_Mean	price_Mean	credit_Mean	av_pay_Mean	return_Mean	warranty_Mean	talk_dir_Mean
1	8.6814159	8.7168142	7.6371681	8.3539823	8.3893805	3.619469	2.2743363	7.0088496	8.159292	8.4955752
2	8.6482759	8.6	5.0551724	7.2	8.1241379	6.3793103	2.2965517	7.4758621	8.3241379	8.6068966
3	8.5113636	8.6704545	6.3977273	6.8409091	7.9318182	6.1704545	3.5795455	6.5227273	7.2159091	6.1818182
4	8.7482014	8.7194245	6.5755396	7.0647482	7.5539568	5.3309353	2.647482	4.7697842	6.0503597	8.4676259
5	7.9744898	7.9234694	7.4081633	7.6836735	7.7806122	6.9438776	6.1020408	7.3214286	7.7959184	8.1989796
6	8.4532374	8.4532374	8.1438849	8.3309353	8.381295	7.9280576	2.294964	7.9208633	8.4460432	8.6043165
7	8.6690647	8.6330935	6.5179856	7.9640288	5.7841727	5.1726619	3.4172662	6.8129496	8.0359712	8.6330935

Figure 13: Segment means worksheet of multi-sheet Excel workbook generated

## RESULTS

A summary of the results from hierarchical clustering and k-means procedures are shown in a single table to allow comparison of the relative performance of these segmentation procedures on various forms of transformed data (see Figure 14).

Method	Skewness*	Kurtosis*	Bimodality*	Hierarchical Clustering		K-means	
				Outliers**	Clusters***	Segments Formed	Problematic Segments #
No transformation	-2.9934	14.8482	0.5578	3	3 or 5	6	5
Range Standardized	-2.9934	14.8482	0.5578	3	5 or 6	11	6
Centering	-2.9934	14.8482	0.5578	3	3 or 5	6	5
Normal Standardization	-2.9934	14.8482	0.5578	4	5	13	5
Row Centering	-0.2101	1.1663	0.25	4	3	7	0
Row Standardization	-0.626	2.5282	0.2513	3	7	9	0
Double Centering	-0.2101	1.1663	0.25	4	3	8	0
Double Standardization	-0.4115	-0.4945	0.465	0	6 or 7	7	0

**Figure 14: Comparison table assessing relative performance of hierarchical clustering and k-means procedure on various forms of transformed data.**

\*For the perception attribute "Reliability," which is impacted the most because of response styles.

\*\*Number of observations joining very late in the cluster generations (last 10 cluster generations from cluster history).

\*\*\*Number of possible clusters using local peaks in Pseudo F or Pseudo T-Squared and/or surge in Normalized RMS Distance values.

#Number of segments where more than 80% of attributes continue to exhibit response styles.

### KEY OBSERVATIONS:

- The skewness, kurtosis, and bimodality coefficient values have significantly improved for double-standardized data.
- The number of possible outliers is zero for double-standardized data, observing the last 10 generations of cluster history from hierarchical clustering.
- The possible number of clusters is either 6 or 7 for double-standardized data, considering results from both hierarchical clustering and k-means procedures, indicating the reliability of double-standardized data in forming a stable number of clusters.
- No response style segments are found among the segments formed using the k-means procedure for double-standardized data.

### CONCLUSION

Based on the descriptive statistics and results from hierarchical and k-means clustering, we see that double-standardized data performs better than any other form of transformed data (standardized or centered). The results from hierarchical clustering suggests a six or seven cluster solution (Figure 8.c), whereas k-means results suggest a seven-cluster solution (Figure 9.h) on the double-standardized data. For all other forms of transformed data, either the number of suggested clusters is unclear or the number of cases in one or more clusters formed are not reasonable. Hence, transformation of survey data using the double-standardization method helps improve the chances of getting a better (stable and clean) cluster solution that can be used for profiling customers or raters based on their perceptions where response styles and outliers are inevitable.

## REFERENCES

- Aiken, L.S., & West, S.G. (1991). *Multiple regression: Testing and interpreting interactions*. Newbury Park: Sage Publications.
- Bachman, J.G., & O'Malley, P.M. (1984). Yea-saying, nay-saying, and going to extremes: Blackwhite differences in response styles. *Public Opinion Quarterly*, 48, 491-509.
- DelGobbo, V. (2009). Traffic lighting your multi-sheet Microsoft Excel workbooks the easy way with SAS®. Available at <http://support.sas.com/resources/papers/proceedings10/153-2010.pdf>.
- Der, G., & Everitt, B.S. (2002). *Handbook of statistical analyses using SAS, 2nd edn*. Chapman & Hall/CRC.
- Fischer, R. (2004). Standardization to account for cross-cultural response bias - A classification of score adjustment procedures and review of research in JCCP. *Journal of Cross-Cultural Psychology*, 35(3), 263-282.
- Henning, J. (2010). Response styles confound cross-cultural comparisons. Voice of Vovici. Vovici, 23 JUL 2010. Web. 24 Aug 2010. <<http://blog.vovici.com/blog/bid/38361/Response-Styles-Confound-Cross-Cultural-Comparisons>>.
- Hicks, L.E. (1970). Some properties of ipsative, normative and forced-choice normative measures. *Psychological Bulletin*, 74, 167-184.
- Hofstede, G. (1980). *Culture's consequences: International differences in work-related values*. Beverly Hills, CA: Sage.
- Howell, D.C. (1997). *Statistical methods for psychology*. Belmont, CA: Duxbury Press.
- Leung, K., & Bond, M.H. (1989). On the empirical identification of dimensions for cross-cultural comparisons. *Journal of Cross-Cultural Psychology*, 20, 133-151.
- Paulhus, D. L. (1991). Measurement and control of response bias. In J.P. Robinson, P.R. Shaver, & L.S. Wrightman (Eds.), *Measures of Personality and Social Psychological Attitudes* (Vol. 1). San Diego, CA: Academic Press.
- SAS Institute Inc. (2004). *SAS/STAT® 9.1 User's Guide*. Cary, NC: SAS Institute Inc.
- Vickers, D., & Rees, P. (2007). Creating the UK national statistics 2001 output area classification. *Journal of the Royal Statistical Society Series a-Statistics in Society*, 170, 379-403.

## TRADEMARKS

SAS and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration. Other brand and product names are trademarks of their respective companies.

## ACKNOWLEDGEMENTS

Our thanks to Mr. Will Neafsey, Brand DNA and Consumer Segmentation Manager for Ford Motor Company, for suggesting the idea of this comparison in Dr. Chakraborty's BKS course in M2009.

## CONTACT INFORMATION

The authors welcome and encourage any questions, feedback, remarks, both on- and off-topic via email.

Murali Krishna Pagolu, Oklahoma State University, Stillwater, OK, [murali.k.pagolu@okstate.edu](mailto:murali.k.pagolu@okstate.edu)

Murali Krishna Pagolu is a graduate student in Management Information Systems at Spears School of Business, Oklahoma State University. Before joining the graduate program he worked as an Information Technology consultant. He is a BASE SAS® 9 and Advanced SAS® 9 certified professional and a certified SAS® predictive modeler using

Enterprise Miner 6. He won honorable mention award in the annual Data Mining Shootout competition in a team event and student poster competition at the individual level held by M2010 annual data mining conference.

Goutam Chakraborty, Oklahoma State University, Stillwater, OK, [goutam.chakraborty@okstate.edu](mailto:goutam.chakraborty@okstate.edu)

Goutam Chakraborty is a professor of marketing and founder of SAS® and OSU data mining certificate program at Oklahoma State University. He has published in many journals such as *Journal of Interactive Marketing*, *Journal of Advertising Research*, *Journal of Advertising*, *Journal of Business Research*, etc. He has chaired the national conference for direct marketing educators for 2004 and 2005 and co-chaired M2007 data mining conference. He is also a Business Knowledge Series instructor for SAS®.

## APPENDIX

**Note:** In order to ensure that identification values are retrieved correctly on to the double-standardized data set, a macro is written to store the identification attribute values to macro variables initially. COLn (where n is the observation position number) macro variables are stored in the global macro table based on the observation position. Once the transposed data set is standardized on mean and standard deviation and again transposed to original form, the COLn values stored in the \_NAME\_ attribute are used for mapping the original identification attribute values using the stored macro variables.

### CODE FOR CREATING THE DOUBLE STANDARDIZED DATA AND APPLYING HIERARCHICAL CLUSTERING

```
* create library for the data set *;
libname project '<Type in your path for the data stored here>';

* Sort the data set by identification number *;
proc sort data=project.surveydata out=project.surveydata_sorted;
    by id;
run;

* create macro variables for each identification number *;
proc sql noprint;
    select count(*) into :n
        from project.surveydata_sorted;
    %let n=&n;
    %let pref=COL;
    %put Total number of observations in the data set = &n;
    select id label='id'
    into :&pref.1-:&pref&n
    from project.surveydata_sorted
    order by id;
quit;

* Standardize columns with 0 mean and 1 standard deviation *;
proc standard data=project.surveydata_sorted mean=0 std=1
    out=project.surveydata_s;
run;

* Transpose the normal standardized data *;
proc transpose data=project.surveydata_s
    out=project.surveydata_st;
run;

* Standardize columns (now columns in the transposed data set) with 0 mean and 1
standard deviation *;
proc standard data=project.surveydata_st mean=0 std=1
    out=project.surveydata_sts;
run;

* Transpose the double standardized data set to original form *;
proc transpose data=project.surveydata_sts
    out=project.surveydata_stst;
run;
```

```

* Retrieve the stored identification numbers from macro variables *;
data project.smallexample_dc_final (drop=_NAME_);
  length id $ 8;
  set project.surveydata_stst;
  id = symget(_NAME_);
  label id = 'id';
run;

*Run hierarchical clustering procedure *;
ods graphics on;
proc cluster data=project.smallexample_dc_final
  method=average
  simple
  ccc
  pseudo
  outtree=project.clustreesmallexampdc(label="cluster tree data for
project.smallexample_dc")
  print=30
  plots=psf
  plots=pst2
  plots=ccc;
var reliab time av_br av_spec price credit av_pay return warranty talk_dir;
run;

proc tree data=project.clustreesmallexampdc
  out=project.treedatasmallexampdc(label="disjoint cluster data (from
proc tree) for project.smallexample_dc")
  nclusters=10;
run;
ods graphics off;

```

## SAS CODE NODES TO COMPARE AND ASSESS THE RESPONSE STYLE IMPACT OF DATA ON CUSTOMER SEGMENTS

```

* Create a cross-tabulation of segments vs. average rating in each segment and overall
data *;
data &em_export_train;
  set &em_import_data;
proc tabulate;
var reliab time av_br av_spec price credit av_pay return warrant talk_dir;
  class _segment_;
  table mean _segment_*mean,
  reliab time av_br av_spec price credit av_pay return warranty talk_dir;
run;

```

## MERGING K-MEANS OUTPUT DATA WITH ORIGINAL ATTRIBUTE VALUES

```

* Both the original and the transformed variables have same names in common. Hence, we
need to rename the variables with a suffix or a prefix to ensure both the transformed
and the original variables are available after the merge To avoid any issues while
merging, converting the id variable to numeric *;
data project.kmeansdsdataset_new
  ( rename= (reliab=reliab_ds time=time_ds av_br=av_br_ds av_spec=av_spec_ds
  price=price_ds credit=credit_ds av_pay=av_pay_ds
  return=return_ds warranty=warranty_ds talk_dir=talk_dir_ds) );
  length num_id 8;
  set project.kmeansdsdataset;
  num_id = input(id, 8.);
  drop id ;
  rename num_id=id ;
run;

```

```

* Sorting the data set with untransformed variables before match merging *;
proc sort data=project.xyz_filtered out=project.xyz_filtered_sorted;
    by id;
run;

* Sorting the data set which contains the Double standardized variables and their k-
means cluster memberships *;
proc sort data=project.kmeansdsdataset_new out=project.kmeansdsdataset_sorted;
    by id;
run;

* Match-merging the original filtered data set with the double standardized data set
*;
data project.kmeansdsmerged;
    merge project.xyz_filtered_sorted project.kmeansdsdataset_sorted;
    by id;
run;

```

## PREPARING DATA FOR MULTI-SHEET TRAFFIC LIGHTING ON MICROSOFT EXCEL WORKBOOKS

```

* producing the cross-tabulation of segments vs. average rating in each segment and
overall data in a data set *;
proc tabulate data=project.kmeansdsmerged out=project.dsmeanstats;
var reliab time av_br av_spec price credit av_pay return warranty talk_dir;
class _segment_;
table mean _segment_ * mean,
    reliab time av_br av_spec price credit av_pay return warranty talk_dir;
run;

* drop unnecessary attributes from the data set *;
data project.dsmeanstats_new;
    set project.dsmeanstats (drop= _type_ _page_ _table_);
run;

* create the range flags in a new data set required for traffic lighting by comparing
the individual mean of each attribute within a segment to its overall mean.
* flag = 0 if segment level mean is less than the overall mean
* flag = 1 if segment level mean is greater than or equal to overall mean *;
data project.dsmeanstats_final
    (drop=x i rename=(flag1=reliab_flag flag2=time_flag
        flag3=av_br_flag flag4=av_spec_flag
        flag5=price_flag flag6=credit_flag
        flag7=av_pay_flag flag8=return_flag
        flag9=warranty_flag flag10=talk_dir_flag));
set project.dsmeanstats_new;
array segment_means(10) reliab_mean time_mean
    av_br_mean av_spec_mean
    price_mean credit_mean
    av_pay_mean return_mean
    warranty_mean talk_dir_mean;
array compare (10) _temporary_ ;
array flag(*) flag1-flag10;
if _segment_ eq . then
do x = 1 to dim(segment_means);
    compare(x) = segment_means(x);
end;
else
do i = 1 to dim(segment_means);
    if segment_means{i} < compare{i} then flag{i} = 0;
    else if segment_means{i} >= compare{i} then flag{i} = 1;
    else flag{i}=.;
end;
if _segment_ ne .;
run;

```

## CODE FOR GENERATING THE SEGMENT MEANS IN EXCEL WORKBOOK AND TRAFFIC LIGHTING

```

* Create a SAS table that is used for the range flags worksheet *;
data project.Legend;
length SegmentFlag 8 Range $30;
SegmentFlag = 0; Range = 'Low'; output;
SegmentFlag = 1; Range = 'High'; output;
label Range = 'Range Flags'
SegmentFlag = 'Color';
run;

* Create a format for traffic lighting, based on the lab flag values
* 0: #CCFFFF - Low
* 1: #9999FF - High *;
proc format;
value FlagFmt
0 = '#CCFFFF'
1 = '#9999FF';
run; quit;

ods listing close;
ods tagsets.ExcelXP path=<Type the path where you want to output the file here>
file='SegmentsReport.xml'
style=XLSansPrinter;
title; footnote;

* Create the range flags worksheet *;
ods tagsets.ExcelXP options(sheet_name='Range Flags');

proc print data=project.Legend noobs label;
var Range;
var SegmentFlag / style(column)=[foreground=FlagFmt. background=FlagFmt.];
run; quit;

* Need to reset the option value *;
ods tagsets.ExcelXP options(sheet_name='Segment Means');

* Create the lab results worksheets *;
ods tagsets.ExcelXP options(sheet_label='Big Picture'
suppress_bylines='yes'
absolute_column_width='8,10,10,10,10,10,10,10,10,12,10,12,12,12,12'
autofit_height='yes'
frozen_headers='yes' frozen_rowheaders='11');

proc report data=project.dsmeansstats_final split='*' nowindows;

* 'ID' columns *;
column _SEGMENT_;

* Data columns with spanned headers *;
column reliab_Mean time_Mean av_br_Mean av_spec_Mean price_Mean
credit_Mean av_pay_Mean return_Mean warranty_Mean talk_dir_Mean;

* Hidden columns containing the range flags *;
column reliab_flag time_flag av_br_flag av_spec_flag price_flag credit_flag
av_pay_flag return_flag warranty_flag talk_dir_flag;

* Dummy column to perform traffic lighting *;
column dummy;

* 'ID' columns *;
define _SEGMENT_ / display order style(Column)=data_center;

```

```

* Data columns *;
define reliab_Mean / display;
define time_Mean / display;
define av_br_Mean / display;
define av_spec_Mean / display;
define price_Mean / display;
define credit_Mean / display;
define av_pay_Mean / display;
define return_Mean / display;
define warranty_Mean / display;
define talk_dir_Mean / display;

* Hidden columns containing the range flags *;
define reliab_flag / display noprint;
define time_flag / display noprint;
define av_br_flag / display noprint;
define av_spec_flag / display noprint;
define price_flag / display noprint;
define credit_flag / display noprint;
define av_pay_flag / display noprint;
define return_flag / display noprint;
define warranty_flag / display noprint;
define talk_dir_flag / display noprint;

* Dummy column to perform traffic lighting *;
define dummy / computed noprint;

* Traffic light the data columns based on the hidden columns *;
compute dummy;

array name(10) $31 ('reliab_Mean' 'time_Mean' 'av_br_Mean' 'av_spec_Mean' 'price_Mean'
'credit_Mean' 'av_pay_Mean' 'return_Mean' 'warranty_Mean'
'talk_dir_Mean');

array flag(10) reliab_flag time_flag av_br_flag av_spec_flag price_flag credit_flag
av_pay_flag return_flag warranty_flag talk_dir_flag;

* Loop over all the _Result columns ('name' array), and set the BACKGROUND style
attribute based on the value of the corresponding _Flag column ('flag' array) *;
do i = 1 to dim(name);
if (flag(i) ge 0) then call define(name(i), 'style',
'style=[background=' || put(flag(i), FlagFmt.) || ']');
end;
endcomp;
run; quit;
ods tagsets.ExcelXP close;

```