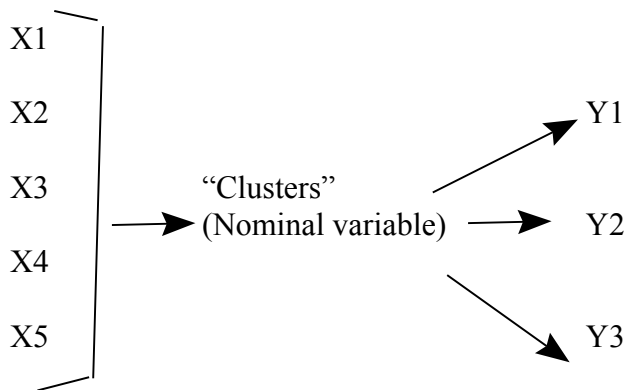


*Neuendorf*  
Cluster Analysis

Model:



Assumes:

1. Actually, any level of measurement (nominal, ordinal, interval/ratio) is acceptable for certain types of clustering. The typical methods, though, require metric (I/R) data. The most basic assumption is that there are two or more (clustering or “internal” or “independent”) variables, and that cases (often people, sometimes products or other identifiable units; Hair et al. call them “objects”) will be placed into groups on the basis of how similar they are on the selected variables. [Then, usually, these clusters are used to predict status on one or more “dependent” or “external” variables.]
2. Sample representativeness. That means some type of probability/random sampling.
3. Not very much else is assumed! As Hair et al. note, even multicollinearity is not technically prohibited (although it’s hard to imagine *why* you would want to cluster on a set of variables that could better be represented by a scale). As Aldenderfer and Blashfield point out in their Sage QASS book:
  - “Most cluster analysis methods are relatively simple procedures that. . . are not supported by an extensive body of statistical reasoning.”
  - “Different clustering methods can and do generate different solutions to the same data set.”

Decisions to Make:

1. Whether you need to run PermuCLUSTER or not. This add-on to SPSS will help determine the “optimal” ordering for your data set to avoid problems that can occur when there are distance “ties” in your data. This step must be conducted *before* running Cluster in SPSS. See the references below for more info.
2. Standardized variables vs. original (“raw”) data. This choice is very similar to the choice of whether to standardize variables before putting them together in a scale. If the variables are measured on difference response scale, we typically standardize. If the variables are measured on the same response scale, we typically do not standardize. But, even when data are unstandardized, the variable(s) with the largest variance will dominate.

### 3. Similarity measure

- A. Correlation coefficients (Pearson or Spearman)--the unit of analysis is now the variable, and the "variable" is now a person or object. For example, there will be a correlation between person A and person B, for their values across, say, ten selected variables. This is quite similar to a factor analysis of persons or objects! This method of analysis is sensitive to patterns of relationships, rather than absolute values on variables (i.e., the magnitude of the difference).
- B. Association coefficients--for categorical/nominal data; there are more than 30 of these stats. available; e.g., the simple matching coefficient and Jaccard's coefficient (see Aldenderfer and Blashfield, p. 29).
- C. Distance coefficients--assumes I/R data; except for #4, can select either original variables or standardized variables. See the attached two-variable example of these distance measures:

1. Euclidean distance

$$d_{ij} = \sqrt{\sum_{k=1}^p (x_{ik} - x_{jk})^2}$$

2. Squared Euclidean distance

$$d_{ij} = \sum_{k=1}^p (x_{ik} - x_{jk})^2$$

3. Manhattan or city-block distance

$$d_{ij} = \sum_{k=1}^p |x_{ik} - x_{jk}|$$

4. Mahalanobis  $D^2$  or generalized distance

$$d_{ij} = (X_i - X_j)' \Sigma^{-1} (X_i - X_j)$$

where  $\Sigma$  is the pooled within-groups variance-covariance matrix, and  $X_i$  and  $X_j$  are vectors of the values for the variables for cases  $i$  and  $j$ . When the correlations between variables are all zero, Mah.  $D^2$  = squared Euclidean distance. The noteworthiness of Mah.  $D^2$  is that it takes intercorrelations among the variables into account. Also, it “automatically” standardizes the variables.

## 4. Methods of clustering

## A. Hierarchical

1. Agglomerative (the most common; *see below*)
2. Divisive

B. Nonhierarchical--includes several options; see Hair p. 590. Such procedures have grown in popularity with increasing computational power, but are not a good choice without theory. That's because nonhierarchical procedures require the choice of “cluster seeds” (with values on all variables) to begin the clustering process. They also assume that the researcher can intelligently select the number of clusters to be derived. . . prior to the clustering.

- A. 1. Agglomerative, cont'd.: Assuming that Hierarchical Agglomerative is chosen, there are at least 5 options from which to choose, in terms of how cases are chosen for inclusion in clusters:
- a. Single linkage ("nearest neighbor")--can be a problem in that it may produce snake-like clusters
  - b. Complete linkage ("furthest neighbor")
  - c. Average linkage--criterion is average distance from individuals in one cluster to individuals in another; probably the most common method
  - d. Ward's method--uses squared Euclidean distances from each case to the mean of its cluster across all variables, summed across all cases; at each step, the two clusters that merge are those that result in the smallest increase in the overall sum of the squared within-cluster distances; useful in creating clusters of relatively equal size, one of the two most common methods
  - e. Centroid method--uses distances between cluster centroids (multivariate means); for every new grouping, centroids are recomputed, thus the disadvantage is that "reversals" may occur (where centroid distance between two unlinked clusters can be smaller than that of two clusters linked earlier)

5. Number of clusters
  - A. A priori--A decision based on theory, past work, etc. (With the "K-Means" procedure, this is the only option.) With the "Hierarchical Cluster" procedure, there are two additional options:
  - B. Intuitive, data-based procedures--(1) look at the dendrogram or icycle plot, and/or (2) look at n's for clusters (need to save cluster groupings and then run frequencies!).
  - C. Formal tests--SPSS provides an agglomeration coefficient. Small values indicate that fairly homogeneous clusters are being merged at that step. We may look for large increases as a "cut-off," similar to the scree test in factor analysis. See Hair et al. p. 604 for an example.
6. Treatment of outliers--Single-case clusters may be candidates to be dropped from the analysis, particularly in a solution with a small number of clusters.
7. Validation and interpretation techniques (Profiling)--Hair et al. indicate that one may validate by splitting the sample and comparing cluster analyses. Another procedure commonly used to interpret or "make sense" of the clusters is to look at cluster means on (1) the "internal" (clustering) variables used to produce the clusters, and (2) "external" variables of interest, not used to produce the clusters. Hair et al. call this "profiling." We can simply profile each cluster as an heuristic in order to get a clearer picture of that group, or we can use the cluster groupings as an independent variable for an ANOVA or MANOVA procedure, or as a dependent variable in a discriminant analysis procedure. To use the cluster groupings for further analyses, use the "save" function in cluster analysis, and cluster membership variables will be added to the data set.

#### Statistics:

Not much! Just what's already been mentioned:

1. Agglomerative coefficient--the within-cluster sum of squares (Ward's method) or the squared Euclidean distance between the two cases of clusters being combined (other methods), this coefficient indicates whether the clusters being joined are very homogenous (small coefficient) or different (large coefficient). Again, a scree-type inspection may be used.
2. ANOVA/MANOVA analyses--to "profile" the cluster groupings with both "internal" and "external" variables. Must be run *after* clustering, with cluster solution(s) saved. I tend to use "MEANS" with the ANOVA stats clicked under Options. The cluster variable is the IV.

#### Other notes:

1. SPSS provides "K-Means" Cluster Analysis, which is faster and simpler; but it limits the user to simple Euclidean distances only. This procedure asks for you to indicate how many clusters you want, and to provide "centers" (i.e., values on the clustering variables to serve as the comparison points). If you do not specify the centers, SPSS selects some for you. (SPSS says, "a number of well-spaced cases equal to the number of clusters is selected from the data.")
2. When using regular "Hierarchical Cluster" in SPSS, the user might need to wait a few minutes--it's much slower to run than other procedures we've used. It is found under "Analyze->Classify->Hierarchical Cluster."

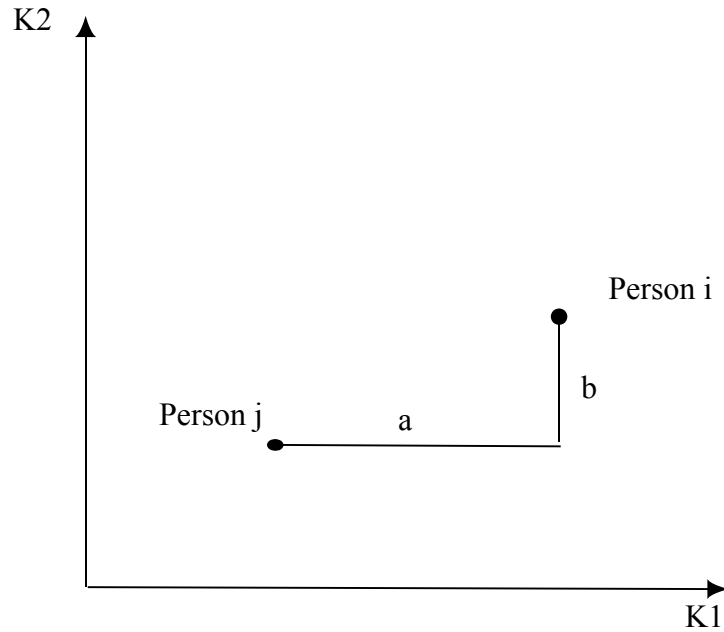
References:

Aldenderfer, M. S., & Blashfield, R. K. (1984). *Cluster analysis*. Beverly Hills, CA: Sage.

Spaans, A., & van der Kloot, W. (2004). *PermuCLUSTER 1.0 user's guide*. The Netherlands: Department of Psychology, University of Leiden. (Available at: <http://www.leidenuniv.nl/fsw/mmlab/software/PermuCluster%201.0%20User%20Guide.pdf>)

van der Kloot, W. A., Spaans, A. M. J., & Heiser, W. J. (2005). Instability of hierarchical cluster analysis due to input order of the data: The PermuCLUSTER solution. *Psychological Methods, 10*, 468-476.

Attachment A: Examples of distances with two variables (K1 and K2):



$$a = x_{iK1} - x_{jK1}$$

$$b = x_{iK2} - x_{jK2}$$

1. Euclidean distance  $d = \sqrt{a^2 + b^2}$
2. Squared Euclidean distance  $d = a^2 + b^2$
3. Manhattan or city-block distance  $d = a + b$
4. Mahalanobis  $D^2$  or generalized distance (IFF K1 and K2 are orthogonal)  $d = a^2 + b^2$

## Attachment B: General Population Cluster Profiling, Humor Study.

Cluster Name:	_____	_____	_____	_____	_____	_____	_____	_____	Total		
Variable:	1 (20)	2 (44)	3 (71)	4 (66)	5 (20)	6 (52)	7 (16)	8 (32)	(312)	F	Sig.
<b>Internal Variables:</b>											
Mean-spirited humor	-14 <sup>c</sup>	-57 <sup>c</sup>	1.10 <sup>a</sup>	.61 <sup>b</sup>	-.72 <sup>c</sup>	-.94 <sup>d</sup>	-.54 <sup>c,d</sup>	-.38 <sup>c</sup>	0.00	55.61 <sub>(7,313)</sub>	.000
Visual/Verbal humor	-1.69 <sup>a</sup>	-.04 <sup>d</sup>	-.10 <sup>d</sup>	.24 <sup>d</sup>	1.01 <sup>b</sup>	-.65 <sup>c</sup>	.44 <sup>b,c,d</sup>	1.04 <sup>b</sup>	0.00	36.45 <sub>(7,313)</sub>	.000
Absurd/Stupid humor	-.63 <sup>a,b</sup>	-1.11 <sup>a</sup>	-.30 <sup>b</sup>	.41 <sup>c</sup>	.86 <sup>c</sup>	.75 <sup>c</sup>	.94 <sup>c</sup>	-.38 <sup>b</sup>	0.00	33.46 <sub>(7,313)</sub>	.000
Social humor	-.78 <sup>a,c</sup>	.29 <sup>b</sup>	-.21 <sup>c,d</sup>	.75 <sup>b</sup>	.35 <sup>b,d</sup>	.31 <sup>b</sup>	-1.20 <sup>a</sup>	-1.12 <sup>a</sup>	0.00	29.46 <sub>(7,313)</sub>	.000
Satire/Death humor	-.71 <sup>a,b</sup>	-.37 <sup>b,d</sup>	.69 <sup>c,e</sup>	-.79 <sup>a,d</sup>	1.05 <sup>c</sup>	.45 <sup>e</sup>	-1.26 <sup>a</sup>	.30 <sup>e</sup>	0.00	40.92 <sub>(7,313)</sub>	.000
<b>External Variables:</b>											
Age (in years)	45.0 <sup>a,b</sup>	46.0 <sup>a</sup>	34.5 <sup>b,c,d</sup>	32.0 <sup>d</sup>	48.5 <sup>a</sup>	48.4 <sup>a</sup>	48.4 <sup>a,c</sup>	49.4 <sup>a</sup>	41.6	8.50 <sub>(7,297)</sub>	.000
Gender (% female)	75% <sup>a</sup>	68% <sup>a</sup>	32% <sup>b</sup>	62% <sup>a</sup>	45% <sup>a,b</sup>	77% <sup>a</sup>	88% <sup>a</sup>	66% <sup>a</sup>	60%	6.25 <sub>(7,313)</sub>	.000
Education (1-6 scale; 5=college grad)	3.7 <sup>a,b</sup>	3.9 <sup>a,b</sup>	4.1 <sup>a</sup>	3.5 <sup>b</sup>	4.2 <sup>a,b</sup>	4.2 <sup>a</sup>	3.9 <sup>a,b</sup>	4.1 <sup>a,b</sup>	4.0	2.47 <sub>(7,302)</sub>	.018
Race (% non-white)	42% <sup>a</sup>	9% <sup>b</sup>	14% <sup>a,b</sup>	23% <sup>a,b</sup>	15% <sup>a,b</sup>	19% <sup>a,b</sup>	38% <sup>a,b</sup>	16% <sup>a,b</sup>	19%	2.20 <sub>(7,312)</sub>	.034
Marital status (% married)	45% <sup>a,b</sup>	48% <sup>a,b</sup>	41% <sup>a,b</sup>	21% <sup>a</sup>	45% <sup>a,b</sup>	48% <sup>a,b</sup>	38% <sup>a,b</sup>	53% <sup>b</sup>	41%	2.18 <sub>(7,313)</sub>	.036
Days/week read newspaper	3.6 <sup>a,b</sup>	4.6 <sup>a,b</sup>	3.9 <sup>a,b</sup>	3.2 <sup>a</sup>	5.5 <sup>b</sup>	4.1 <sup>a,b</sup>	3.6 <sup>a,b</sup>	4.1 <sup>a,b</sup>	4.0	2.11 <sub>(7,307)</sub>	.042
Videos watched in past month	3.5 <sup>a,b</sup>	3.9 <sup>a</sup>	6.1 <sup>a,b</sup>	9.5 <sup>b</sup>	7.7 <sup>a,b</sup>	5.0 <sup>a,b</sup>	3.5 <sup>a,b</sup>	3.6 <sup>a,b</sup>	5.8	2.39 <sub>(7,307)</sub>	.021
Have DVD player	15% <sup>a,b</sup>	7% <sup>a,b</sup>	13% <sup>a,b</sup>	25% <sup>a</sup>	11% <sup>a,b</sup>	0% <sup>b</sup>	25% <sup>a,b</sup>	6% <sup>a,b</sup>	12%	3.20 <sub>(7,306)</sub>	.003
Have camcorder	15% <sup>a</sup>	36% <sup>a,b</sup>	56% <sup>b</sup>	58% <sup>b</sup>	37% <sup>a,b</sup>	40% <sup>a,b</sup>	63% <sup>a,b</sup>	34% <sup>a,b</sup>	46%	3.12 <sub>(7,306)</sub>	.003
Have satellite dish	10% <sup>a,b</sup>	0% <sup>a</sup>	4% <sup>a,b</sup>	17% <sup>b</sup>	5% <sup>a,b</sup>	6% <sup>a,b</sup>	13% <sup>a,b</sup>	0% <sup>a</sup>	7%	2.51 <sub>(7,306)</sub>	.016
How eager for DTV (0-10 scale)	0.8 <sup>a</sup>	2.1 <sup>a,b</sup>	3.0 <sup>a,b</sup>	3.3 <sup>b</sup>	1.5 <sup>a,b</sup>	1.7 <sup>a,b</sup>	2.3 <sup>a,b</sup>	2.1 <sup>a,b</sup>	2.4	2.43 <sub>(7,289)</sub>	.020
Depression scale (0-140 scale)	34.6 <sup>a,b</sup>	27.7 <sup>a,b</sup>	31.5 <sup>a,b</sup>	39.0 <sup>a</sup>	22.4 <sup>a,b</sup>	27.2 <sup>a,b</sup>	26.2 <sup>a,b</sup>	19.4 <sup>b</sup>	29.9	3.28 <sub>(7,304)</sub>	.001
Wallowing scale (0-100 scale)	29.7 <sup>a</sup>	27.3 <sup>a</sup>	33.1 <sup>a</sup>	33.4 <sup>a</sup>	29.3 <sup>a</sup>	34.9 <sup>a</sup>	26.9 <sup>a</sup>	26.9 <sup>a</sup>	31.3	2.25 <sub>(7,312)</sub>	.030
Favorite TV show is comedy	27% <sup>a</sup>	44% <sup>a,b</sup>	63% <sup>a,b</sup>	63% <sup>a,b</sup>	83% <sup>b</sup>	46% <sup>a,b</sup>	64% <sup>a,b</sup>	57% <sup>a,b</sup>	57%	2.60 <sub>(7,265)</sub>	.013
Favorite TV show is weepy/melodrama	47% <sup>a</sup>	31% <sup>a</sup>	15% <sup>a</sup>	11% <sup>a</sup>	11% <sup>a</sup>	11% <sup>a</sup>	7% <sup>a</sup>	19% <sup>a</sup>	17%	2.13 <sub>(7,265)</sub>	.041
Favorite movie is weepy (2), melodrama (1), or not (0)	1.2 <sup>a,b</sup>	0.8 <sup>a,b</sup>	0.4 <sup>a</sup>	0.7 <sup>a,b</sup>	0.9 <sup>a,b</sup>	0.9 <sup>a,b</sup>	1.5 <sup>b</sup>	0.9 <sup>a,b</sup>	0.8	3.41 <sub>(7,272)</sub>	.002
Violence in favorite movie (graphic=2, lite=1; no=0)	.46 <sup>a,b</sup>	.63 <sup>a,b</sup>	1.07 <sup>a</sup>	1.00 <sup>a</sup>	.68 <sup>a,b</sup>	.36 <sup>b</sup>	.58 <sup>a,b</sup>	.74 <sup>a,b</sup>	.78	3.69 <sub>(7,272)</sub>	.001

Means that do not share a superscript are significantly different at  $p < .05$  using Tukey's HSD post hoc test.