

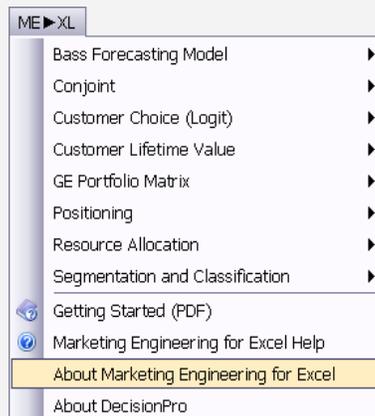
Tutorial

Segmentation and Classification



Marketing Engineering for Excel is a *Microsoft Excel* add-in. The software runs from within *Microsoft Excel* and only with data contained in an *Excel* spreadsheet.

After installing the software, simply open *Microsoft Excel*. A new menu appears, called "ME ▶ XL." This tutorial refers to the "ME ▶ XL/Segmentation and Classification" submenu.



Overview

Segmentation and classification is an analytic technique that helps firms compare and group customers who share common characteristics (i.e., segmentation variables) into homogeneous segments and identify those particular customers in a market on the basis of external variables (i.e., discriminant variables).

Segmentation refers to the process of classifying customers into homogenous groups (segments), such that each group of customers shares enough characteristics in common to make it viable for the firm to design specific offerings or products for it. This application identifies customer segments using needs-based variables called basis variables. Cluster analysis helps firms:

- ✓ Better understand their customers.
- ✓ Identify different segments in a market.
- ✓ Choose attractive customer segments for classification with its marketing programs.

Getting Started

To apply segmentation and classification analysis, you can use your own data directly or a template preformatted by the ME►XL software.

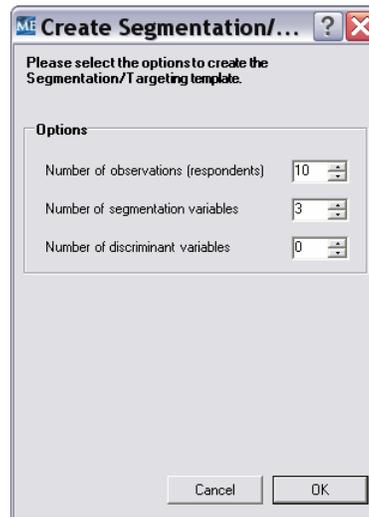


The next section explains how to create an easy-to-use template to enter your own data.

If you want to run a segmentation and classification analysis immediately, open the example file "*OfficeStar Data (Segmentation).xls*" and jump to "Step 3: Running analysis" (p. 4). By default, the example files install in "*My Documents/My Marketing Engineering/*."

Step 1 Creating a template

In Excel, if you click on ME►XL → SEGMENTATION AND CLASSIFICATION → CREATE TEMPLATE, a dialog box appears. This box represents the first step in creating a template to run the segmentation and classification analysis software.



The dialog box requests three pieces of information to design the template:

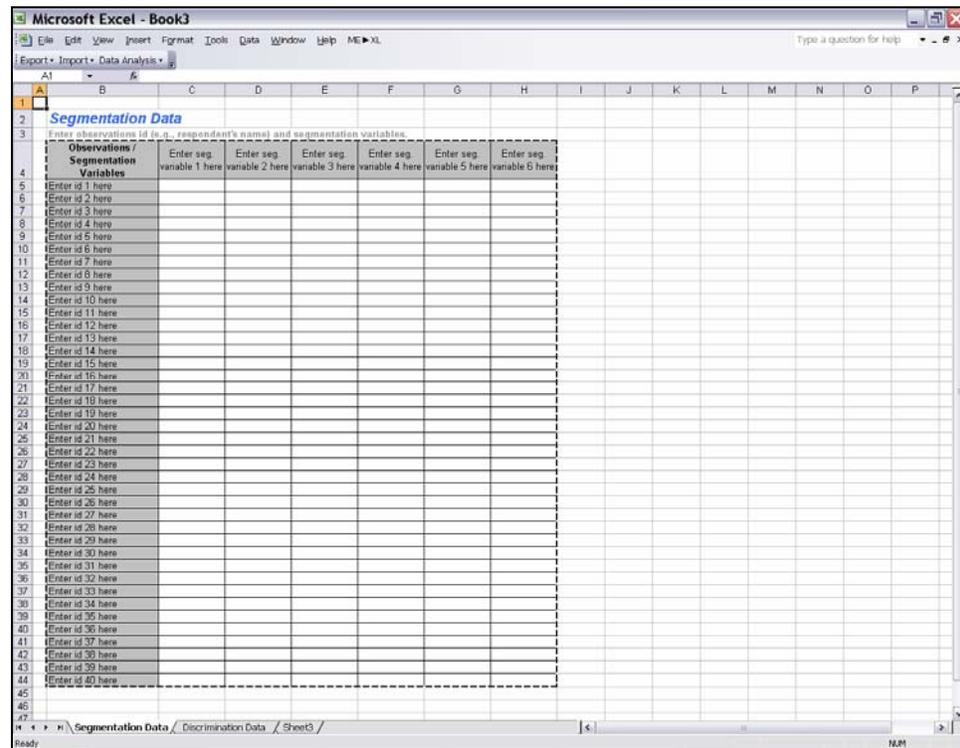
- **Observations** (respondents) indicate the number of customers or respondents in the data that need to be clustered.
- **Segmentation variables** help us assess the similarity between two respondents. These variables serve as the basis for segmentation and are often called **basis variables**. They might include customer's needs, wants, expectations, or preferences.
- **Discriminant variables**, also called **descriptors**, are optional variables that can describe the segments formed on the basis of the segmentation variables. These include demographic variables, such as educational level, gender, income, media consumption, and the like.

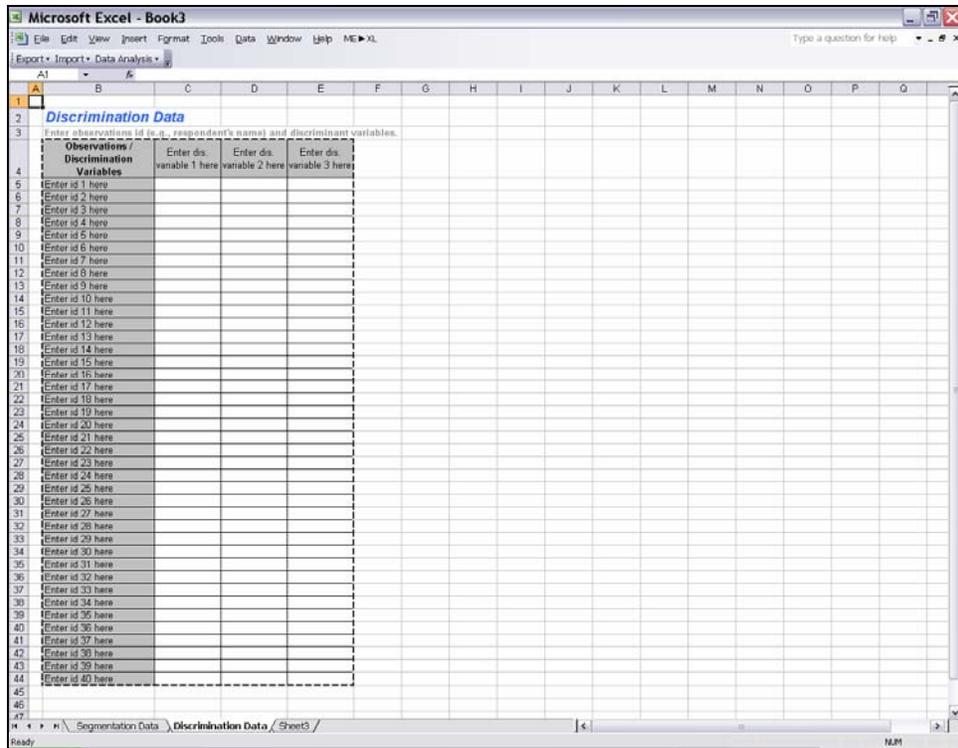


It is not always clear whether a specific variable should be treated as segmentation variable or discriminant variable. This choice might depend on the context, the managerial question, or the product category.

When in doubt, ask yourself the following questions: (1) Would this piece of information tell me what that customer wants, in which case it should be treated as segmentation variable, or (2) does this piece of information tell me who that customer is and therefore should be treated as discriminant variable? For example, "gender" would fall in the second category most of the time, whereas "need for timely information" usually falls in the former category.

After specifying the number of observations and variables, click OK to proceed. The software generates a template that contains either one or two sheets, depending on whether you have included discriminant data.





Step 2 Entering your data



In this tutorial, we use the example file "OfficeStar Data (Segmentation).xls," which appears by default in "My Documents/My Marketing Engineering/."

To view a proper data format, open that spreadsheet in Excel. A snapshot is shown below.

Microsoft Excel - OfficeStar (Segmentation).xls

File Edit View Insert Format Tools Data Window Help ME XL

Type a question for help

1 A B C D E F G H I J K L M N O P

2 **Segmentation Data**

3 Enter observations id (e.g., respondent's name) and segmentation variables.

Observations / Segmentation Variables	Variety of choice	Electronics	Furniture	Quality of service	Low prices	Return policy
Respondent 1	8	6	6	3	2	2
Respondent 2	6	3	1	4	7	8
Respondent 3	6	1	2	4	9	6
Respondent 4	8	3	3	4	8	7
Respondent 5	4	6	3	9	2	5
Respondent 6	8	4	3	5	10	6
Respondent 7	7	2	2	2	8	7
Respondent 8	7	5	7	2	2	3
Respondent 9	7	7	5	1	5	4
Respondent 10	8	4	0	4	9	8
Respondent 11	9	8	5	1	5	2
Respondent 12	4	4	2	8	2	3
Respondent 13	10	6	6	1	3	3
Respondent 14	6	5	2	9	3	6
Respondent 15	7	3	0	2	7	6
Respondent 16	9	7	4	5	2	3
Respondent 17	10	6	7	4	4	3
Respondent 18	5	2	1	3	8	7
Respondent 19	10	5	4	4	3	3
Respondent 20	5	5	2	9	2	6
Respondent 21	3	7	1	9	2	3
Respondent 22	8	6	6	2	5	4
Respondent 23	8	4	1	4	7	8
Respondent 24	4	3	0	7	1	3
Respondent 25	10	5	7	1	4	4
Respondent 26	10	6	6	2	2	2
Respondent 27	10	5	7	2	5	2
Respondent 28	4	5	2	8	4	5
Respondent 29	7	1	1	5	9	5
Respondent 30	10	8	4	4	5	5
Respondent 31	4	5	2	5	10	5
Respondent 32	10	5	4	1	2	2
Respondent 33	7	6	5	3	5	3
Respondent 34	10	5	7	1	2	5
Respondent 35	7	3	2	2	10	5
Respondent 36	8	2	4	2	7	2
Respondent 37	7	1	0	2	7	5
Respondent 38	6	4	2	9	4	4
Respondent 39	9	6	6	4	3	3
Respondent 40	10	8	5	3	4	5

Segmentation Data / Discrimination Data / Sheet3

File Edit View Insert Format Tools Data Window Help ME XL Adobe PDF

J32

1 A B C D E F G H I J K L M N O P Q R

2 **Discrimination Data**

3 Enter observations id (e.g., respondent's name) and discriminant variables.

Observations / Discrimination Variables	Professional	Income (000's)	Age
Respondent 1	1	40	49
Respondent 2	0	20	41
Respondent 3	0	20	34
Respondent 4	1	38	34
Respondent 5	1	45	58
Respondent 6	1	35	28
Respondent 7	1	45	30
Respondent 8	0	65	59
Respondent 9	0	45	59
Respondent 10	0	45	23
Respondent 11	1	50	34
Respondent 12	0	25	63
Respondent 13	1	65	32
Respondent 14	1	60	58
Respondent 15	1	30	24
Respondent 16	0	45	38
Respondent 17	0	55	43
Respondent 18	0	25	30
Respondent 19	0	40	32
Respondent 20	1	70	55
Respondent 21	1	55	58
Respondent 22	0	25	38
Respondent 23	1	15	28
Respondent 24	1	50	30
Respondent 25	1	70	31
Respondent 26	1	70	58
Respondent 27	0	55	50
Respondent 28	0	65	54
Respondent 29	0	50	60
Respondent 30	0	30	32
Respondent 31	0	50	23
Respondent 32	0	20	43
Respondent 33	1	55	38
Respondent 34	0	65	59
Respondent 35	1	25	33
Respondent 36	0	20	24
Respondent 37	1	40	20
Respondent 38	1	20	30
Respondent 39	0	45	54
Respondent 40	0	70	44

Segmentation Data / Discrimination Data / Classification Data

A typical segmentation spreadsheet contains one or two spreadsheets that contain segmentation and/or discrimination data.

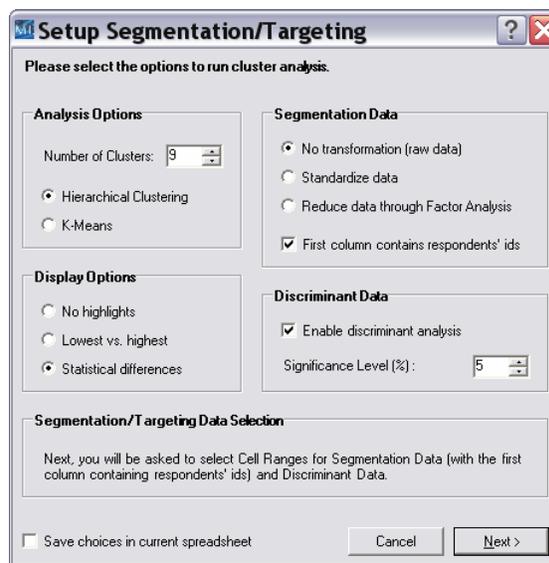
- **Segmentation data** are required for the segmentation model. This data set contains the respondent identifier and a column for each segmentation variable collected in the study. The data within each column must be

scaled using the same scale (e.g., 1–10), but each column can have a different scale (e.g., 1–10 for satisfaction, 1–5 for convenience). Typically, segmentation variables are numerical values (interval or ratio scale). The data set contains one row per respondent in your study.

- **Discriminant data** constitute an optional data set, depending on whether your study has collected discrimination data. Recall that discrimination data enables you to differentiate one customer from another (e.g., age, income, gender). Again, data within a column must be scaled using the same scale, but different columns may use different scales. Typically, discriminant variables are numerical (interval or ratio scale) or nominal (“male”, “female”). Each respondent in your study appears in a separate row.

Step 3 Running segmentation analyses

After you enter your data in an Excel spreadsheet with the appropriate format, click on ME ► XL → SEGMENTATION AND CLASSIFICATION → RUN SEGMENTATION. The dialog box that appears indicates the next steps required to perform a segmentation analysis of your data.



Analysis options

You may specify the number of segments (clusters) to develop during the analysis. For the segmentation method, you can choose either K-means or hierarchical clustering.

- **Hierarchical clustering** builds up or breaks down the data, customer by customer (row by row).
- **K-means** partitioning breaks the data into a prespecified number of segments and then reallocates or swaps customers to improve some measure of effectiveness.



Usually, a segmentation analysis consists of two steps. First, you run the analysis with a large number of segments (up to 9). Second, on the basis of a dendrogram analysis (discussed subsequently), you determine the optimal number of segments to retain.

Segmentation data

This section enables you to specify how to treat the data and whether a first column of respondent identifiers exists.

- **No transformation.** This button indicates you want to use the original data.
- **Standardize data.** This option scales all variables to 0 mean and unit variance before the analysis. Choosing this option is a good idea if you have measured the variables on different scales.
- **Reduce data through Factor Analysis.** This button combines related variables into unique factors.

Display options

In this section, you specify how you want the cluster data presented.

- **No highlights.** The data are unformatted.
- **Lowest vs. highest.** For each variable, colors highlight the value of the cluster with the highest (green) and lowest (red) values.
- **Statistical differences.** For each variable, colors highlight clusters whose values are statistically different from the overall mean at a 95% confidence level. Those that are different from the mean at a 99% confidence level appear in italics.

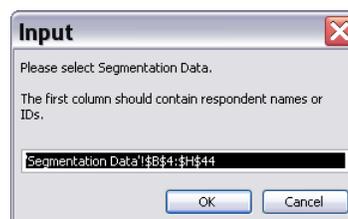
Discriminant data

Decide whether you want the analysis to include a discriminant analysis. Check this button if you wish to perform discriminant analysis, and indicate the level of statistical significance you wish to use.



The *Save choices in current worksheet* option allows you to save cell range selections when you perform Run Analysis. If you are using your own data or have modified a Marketing Engineering for Excel template, you should choose this check box to save your selections.

After selecting all the options, you must select the cells containing the data. When you click Next, the following dialog box appears:



The software requests a range for the segmentation data. If you are using a *Marketing Engineering for Excel* template, the software preselects the cell ranges.

If you have specified the inclusion of discriminant data, the following dialog box appears, which allows to select your discrimination data. The cell ranges might be pre-selected.



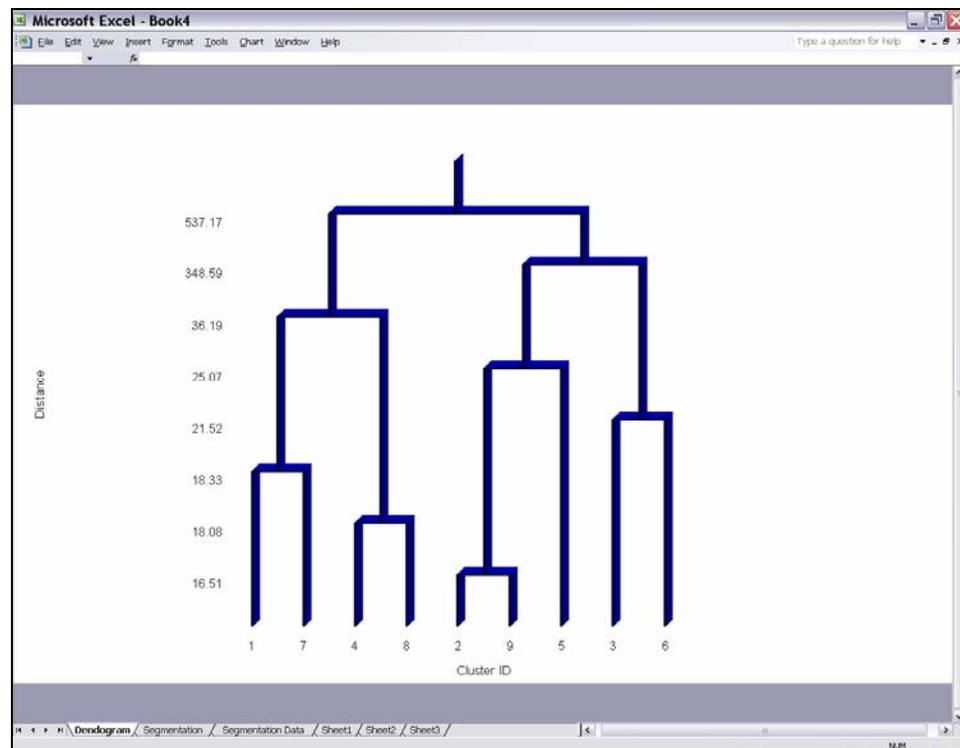
The newly generated workbook contains the results of your segmentation analysis.

Step 4 Interpreting the segmentation results

The workbook generated by segmentation analysis may contain several worksheets, depending on whether your study has included discriminant data.

Dendrogram

Dendograms provide graphical representations of the loss of information generated by grouping different clusters (or customers) together.



At one extreme (upper part of the dendrogram), all customers group into one cluster, and the loss of information is maximum, because they all receive undifferentiated treatment, regardless of their characteristics.

At the other extreme (lower part of the dendrogram), customers appear in separate, small clusters, and only those customers very similar to one another group together ("similar" or "close" in this context refers to the distance between two customers in terms of the segmentation variables).

When reviewing a dendrogram, look for significant distances or "jumps" in the distances. For example, the *OfficeMax* example contains a very large jump when moving from three to two clusters. Grouping these three clusters into two generates a significant loss of information; in other words, it results in grouping within the same cluster customers who are very dissimilar. In the preceding example, a three-cluster solution seems to be the best approach.

A dendrogram is simply a graphical representation of the clustering output. For a more detailed understanding of cluster members and attributes, you must analyze the other tabs in the segmentation output as well.

Segmentation

The tab contains the statistical output of the cluster process and shows cluster sizes (number of members), cluster means, and the placement of each member in clusters (highlighted in yellow). This tab also provides columns that represent individual members and where they would be clustered in a 2–9 cluster solution.

The following table lists the size of the population and of each segment, in both absolute and relative terms.

Size / Cluster	Overall	Cluster 1	Cluster 2	Cluster 3	Cluster 4	Cluster 5	Cluster 6	Cluster 7	Cluster 8	Cluster 9
Number of observations	40	8	3	5	4	7	3	4	2	4
Proportion	1	0.2	0.075	0.125	0.1	0.175	0.075	0.1	0.05	0.1

Means of each segmentation variable for each segment.

Segmentation variable / Cluster	Overall	Cluster 1	Cluster 2	Cluster 3	Cluster 4	Cluster 5	Cluster 6	Cluster 7	Cluster 8	Cluster 9
Variety of choice	7.53	9.38	7.33	5	7.75	7	5	9.5	10	6.5
Electronics	4.57	5.38	3.67	5	6.75	2.86	3.33	5.75	8	2
Furniture	3.45	5.88	0.667	2.2	5.25	2.14	7	6.75	4.5	0.75
Quality of service	4	1.88	4	8.8	1.75	4	8	3.5	3.5	2.25
Low prices	5.05	2.5	7.67	3	5	9	1.67	4.25	4.5	7.5
Return policy	4.5	3	8	5.2	3.25	5.57	3	2.5	5	6.25

The following table lists the cluster number to which each observation belongs for varying cluster solutions. For example, the column "for 2 clusters" gives the cluster number of each observation in a 2-cluster solution. The cluster solution you have selected is in bold with a yellow background.

Observation / Cluster solution	With 2 clusters	With 3 clusters	With 4 clusters	With 5 clusters	With 6 clusters	With 7 clusters	With 8 clusters	With 9 clusters
Respondent 1	1	1	1	1	1	1	1	1
Respondent 2	2	2	2	2	2	2	2	2
Respondent 3	2	2	2	5	5	5	5	5
Respondent 4	2	2	2	5	5	5	5	5
Respondent 5	2	3	3	3	3	3	3	3
Respondent 6	2	2	2	5	5	5	5	5
Respondent 7	2	2	2	2	2	2	2	9
Respondent 8	1	1	1	1	1	1	1	1
Respondent 9	1	1	4	4	4	4	4	4
Respondent 10	2	2	2	2	2	2	2	2

Discrimination

This optional spreadsheet reflects the output of the discrimination analysis. The matrices included on this sheet are as follows:

- **Cluster sizes** depicts the number of respondents who appear in each cluster, along with the proportion of the whole population that each cluster represents.
- **Discriminant variables** depict the means of each discriminant variable for each cluster.
- **Discriminant function** reflects the correlation of the variables with each significant discriminant function and thus indicates the predictive ability of each discriminant function.
- **Confusion matrix** depicts how well the discriminant data predict correct clusters. Two matrices are available, one showing the actual data counts and the other showing percentages for these same data.
- **Classification weights** and **classification coefficients** are intermediary results required to run further classification analyses on external data. These matrices are of no particular interest as is, and cannot be easily interpreted, but are necessary to carry over further classification analyses.

Microsoft Excel - Book4

File Edit View Insert Format Tools Data Window Help ME XL Adobe PDF

H20

Size / Cluster	Overall	Cluster 1	Cluster 2	Cluster 3
Number of observations	40	18	14	8
Proportion	1	0.45	0.35	0.2

Discriminant Variables

Means of each discriminant variable for each segment.

Discriminant variable / Cluster	Overall	Cluster 1	Cluster 2	Cluster 3
Age	40.525	44.222	36.525	45
Income (000's)	42.5	48.333	32.143	47.5
Professional	0.475	0.333	0.5	0.75

Discriminant Function

Correlation of variables with each significant discriminant function

(Significance level < 0.05)

Discriminant variable / Function	Function 1	Function 2
Age	0.91	0.013
Income (000's)	0.696	0.336
Professional	0.068	-0.771
Variance explained	71.36	28.64
Cumulative variance explained	71.36	100
Significance level	0	0.042

Confusion Matrix

Comparison of cluster membership predictions based on discriminant data, and actual cluster memberships. High values in the diagonal of the confusion matrix (in bold) indicates that discriminant data is good at predicting cluster membership.

Actual / Predicted cluster	Cluster 1	Cluster 2	Cluster 3
Cluster 1	10	3	5
Cluster 2	0	13	1
Cluster 3	2	2	4

Actual / Predicted cluster	Cluster 1	Cluster 2	Cluster 3
Cluster 1	55.60%	16.70%	27.80%
Cluster 2	00.00%	92.90%	07.10%
Cluster 3	25.00%	25.00%	50.00%
Hit Rate (percent of total cases correctly classified)	67.50%		

Classification Weights

H \ Dendrogram \ Discrimination \ Segmentation \ Discrimination Data \ Segmentation Data \ Sheet1 \ Sheet2 \ Sheet3

Classification Weights		
Sum of each segment's projection on each function.		
This matrix was used internally, and will be required to run further discriminant analysis (i.e., classification) on external data.		
Clusters / Discriminant Functions	Function 1	Function 2
Segment 1	2.549721	-0.0937281
Segment 2	1.818049	-0.3086759
Segment 3	2.823745	-0.5134836

Classification Coefficients		
Coefficient for each variable in the discrimination function.		
Coefficient for each variable in the discrimination function.		
Discriminant Variables /	Function 1	Function 2
Professional	0.2166553	-0.8049017
Income (000's)	0.0138589	0.0170437
Age	0.0408766	-0.0146871

Segmentation and discriminant data

These tabs contain the original segmentation and discriminant data used for the segmentation analysis, included in the output for your convenience. The original spreadsheet used for the analysis remains intact, so you can modify it for subsequent analysis runs. The data preserved with this tab always reflect the data represented in the dendrogram and segmentation tabs.

Step 5 Running classification analyses

Introduction

If you ran segmentation analysis with discriminant data, the software estimated the best way to predict to which cluster an individual is most likely to belong based solely on discriminant data. This is very useful to predict whether young people (age as a discriminant factor) are more likely to be more price sensitive (price sensitivity as a segmentation variable); or if businesses in certain industries require more support than others.

The ability of recouping segment membership based on discriminant variables is best summarized by the confusion matrix and hit rate (see above).

Once this discriminant analysis has been applied to the original dataset, it can be applied again to external customers for whom discriminant data –but no segmentation data- is available. The process of classifying customers among segments, based on a preceding segmentation analysis, but using discriminant data only, is called **classification analysis**.



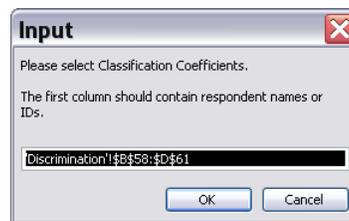
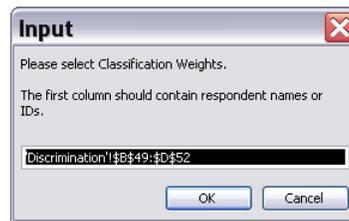
Classification analysis is usually applied to new customers, for whom segmentation data is not available. For learning purpose, you can also apply it to discriminant data of customers for whom segmentation data is available, and see how well segment memberships are recouped. This analysis is automatically done when you run a segmentation analysis, and its results are summarized by the confusion matrix.

Selecting data

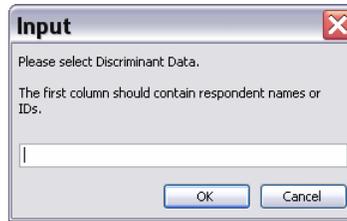
Click on ME ▶ XL → SEGMENTATION AND CLASSIFICATION → RUN CLASSIFICATION. The dialog box that appears indicates the next steps required to perform a classification analysis of your data.



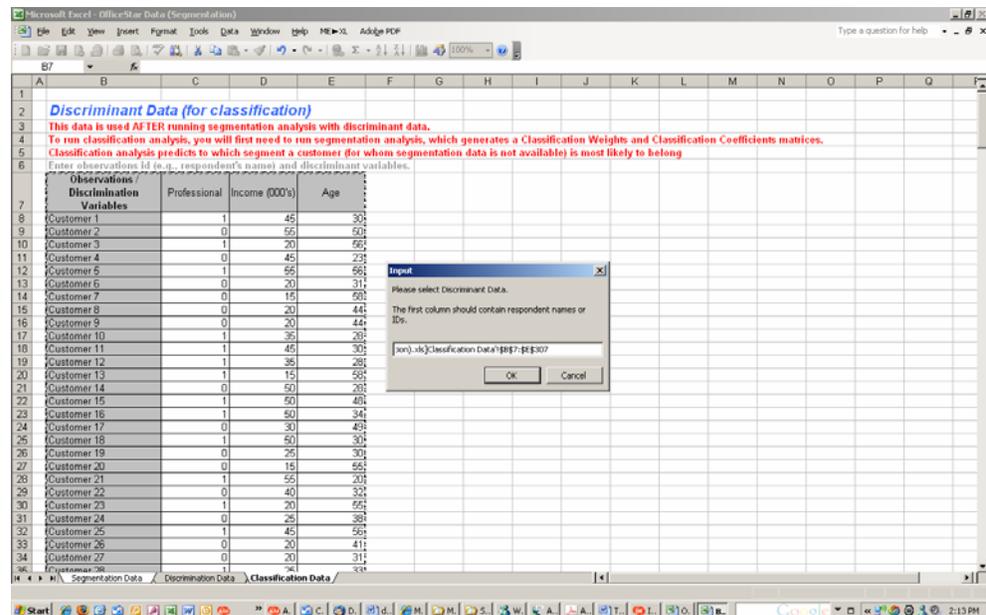
The first steps of the classification analysis consist of selecting two cell ranges: classification weights and classification coefficients. You can find that data at the bottom of the “discriminant” sheet in the analysis workbook generated by segmentation analysis. The cell ranges might be pre-selected.



The last step is to select discriminant data. In most cases, that consists of data about new customers for whom no segmentation data is available. It is important that formatting of the discriminant data matches exactly the format of the discriminant data (both variables, orders and ranges) used for the original segmentation analysis.



Discriminant data of "new" customers is available on the original OfficeStars workbook, in the last sheet. Go back to the OfficeStar workbook, and manually select the discriminant data available for the 300 additional customers (the last sheet of the workbook, named classification data).



F.Y.I.

Once you are in selecting mode, Excel might not allow you to easily switch between two workbooks. If you require selecting data in different workbooks (as it is usually the case with classification analysis), simply use the Window menu of Excel to select and open another workbook.



Interpreting the results

When you click Ok, a new workbook is generated. This workbook contains the discriminant data used to run classification analysis, and the segment to which each customer is most likely to belong.

Microsoft Excel - Book6

File Edit View Insert Format Tools Data Window Help ME XL Adobe PDF

75%

Respondents / Discriminant variables and predicted cluster	Professional	Income (000's)	Age	Predicted Cluster
Customer 1	1	45	30	2
Customer 2	0	55	50	1
Customer 3	1	20	56	3
Customer 4	0	45	23	2
Customer 5	1	55	56	3
Customer 6	0	20	31	2
Customer 7	0	15	58	3
Customer 8	0	20	44	2
Customer 9	0	20	44	2
Customer 10	1	35	28	2
Customer 11	1	45	30	2
Customer 12	1	35	28	2
Customer 13	1	15	58	3
Customer 14	0	50	28	2
Customer 15	1	50	49	3
Customer 16	1	50	34	1
Customer 17	0	30	49	1
Customer 18	1	50	30	2
Customer 19	0	25	30	2
Customer 20	0	15	55	3
Customer 21	1	55	20	2
Customer 22	0	40	32	2
Customer 23	1	20	55	3
Customer 24	0	25	30	2
Customer 25	1	45	56	3
Customer 26	0	20	41	2
Customer 27	0	20	31	2
Customer 28	1	25	33	2
Customer 29	1	50	60	3
Customer 30	1	30	23	2
Customer 31	1	30	22	2
Customer 32	0	45	32	2
Customer 33	1	55	34	1
Customer 34	1	45	30	2
Customer 35	1	15	26	2
Customer 36	1	35	34	2
Customer 37	1	70	31	1
Customer 38	0	20	41	2
Customer 39	0	25	30	2
Customer 40	0	55	57	1
Customer 41	1	30	20	2
Customer 42	0	45	55	1

Note that this classification of customers across segments is our best guess based on discriminant analysis. It is not perfect, and some customers might be misclassified, that is, they are the closest to segment A in terms of needs, but their discriminant variables send us astray and predict they are more likely to belong to segment B.