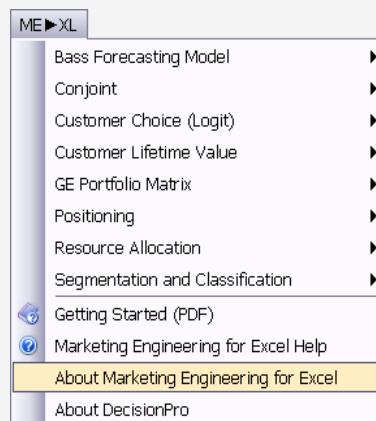**Tutorial**

# 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.

| ME▶XL | |
|---|---|
| Bass Forecasting Model | ▶ |
| Conjoint | ▶ |
| Customer Choice (Logit) | ▶ |
| Customer Lifetime Value | ▶ |
| GE Portfolio Matrix | ▶ |
| Positioning | ▶ |
| Resource Allocation | ▶ |
| Segmentation and Classification | ▶ |
| Getting Started (PDF) | |
| Marketing Engineering for Excel Help | |
| About Marketing Engineering for Excel | |
| About DecisionPro | |

## 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.
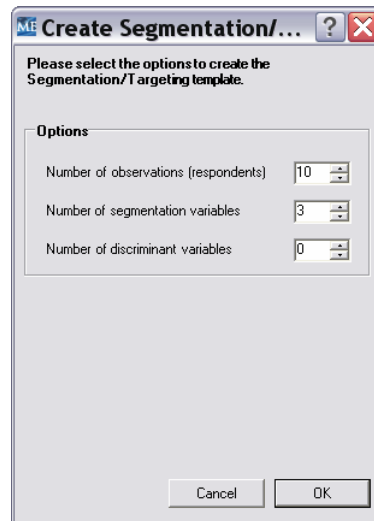
## Step 1   Creating a template

In Excel, if you click on ME►XL → S EGMENTATION AND C LASSIFICATION → C REATE T EMPLATE, 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:
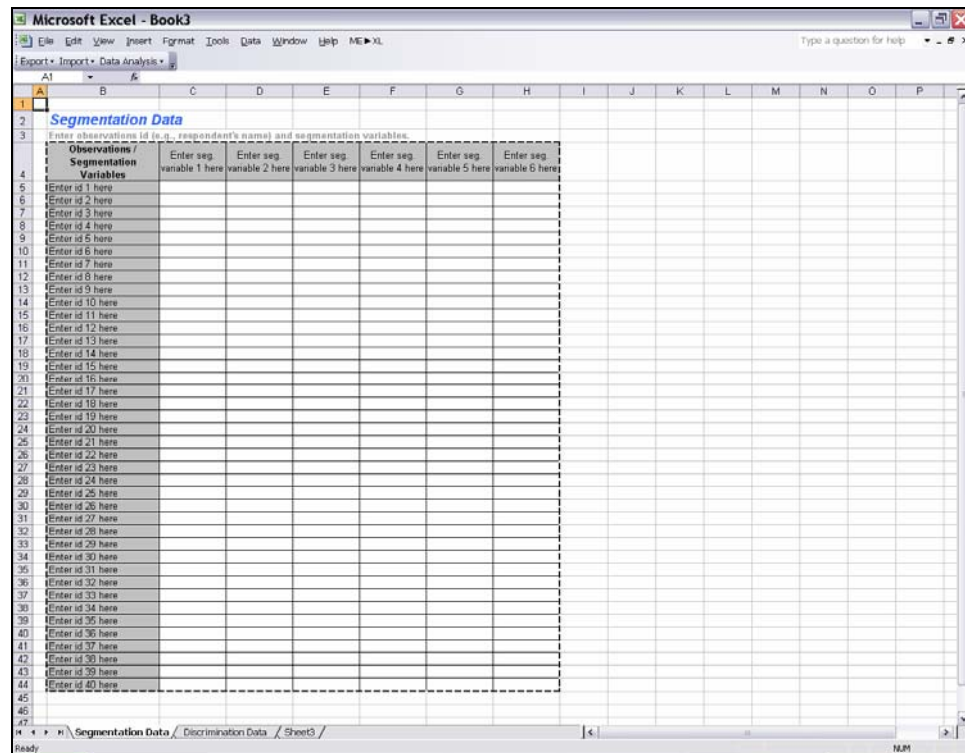
- **Observations** (respondents) indicate the number of customers or respondents in the data that need to be clustered.

- **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.
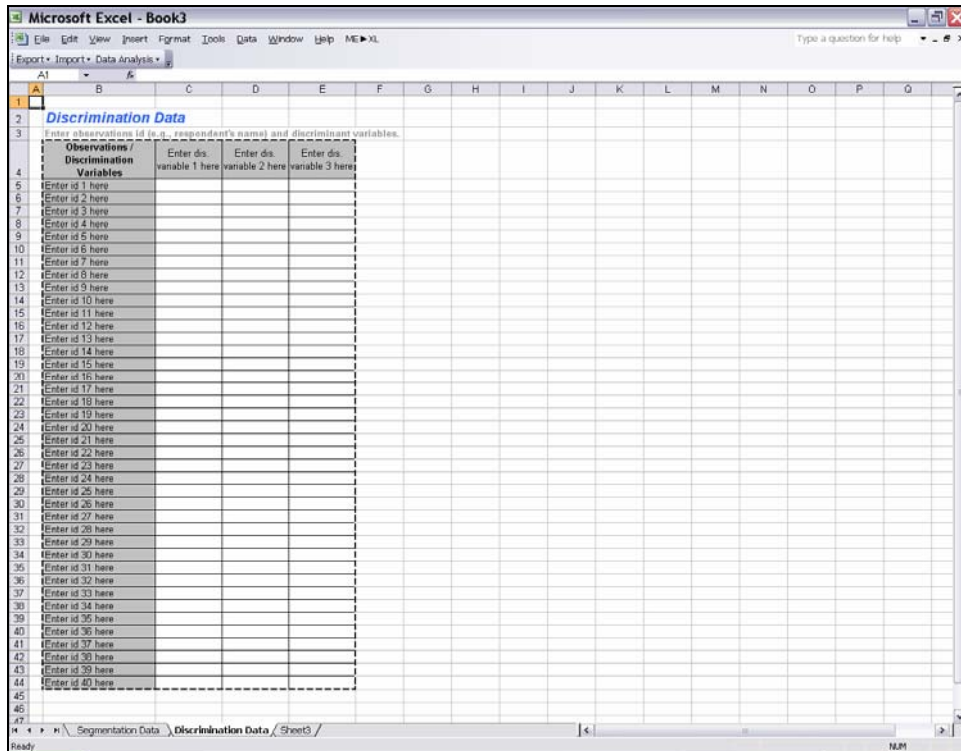
**F.Y.I.**

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 (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**

## Segmentation Data

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 | 6 | 7 | 4 | 5 | 2 |
| Respondent 17 | 10 | 6 | 7 | 4 | 4 | 3 |
| Respondent 18 | 5 | 2 | 1 | 3 | 0 | 7 |
| Respondent 19 | 10 | 5 | 4 | 4 | 3 | 3 |
| Respondent 20 | 5 | 5 | 2 | 9 | 2 | 6 |
| Respondent 21 | 7 | 3 | 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 | 6 |
| Respondent 29 | 7 | 1 | 1 | 5 | 9 | 5 |
| Respondent 30 | 10 | 8 | 4 | 4 | 5 | 6 |
| Respondent 31 | 5 | 4 | 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 | 4 | 2 | 3 | 7 | 5 |
| 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 | 6 |

Segmentation Data / Discrimination Data / Sheet3 /

**Microsoft Excel - OfficeStar (Segmentation).xls**

## Discrimination Data

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

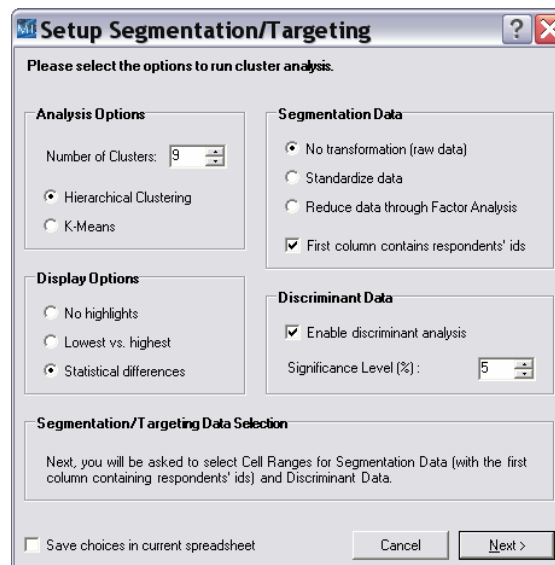| Observations / Discrimination Variables | Professional | Income | Age |
|---|---|---|---|
| Respondent 1 | 1 | 40000 | 45 |
| Respondent 2 | 0 | 20000 | 41 |
| Respondent 3 | 0 | 20000 | 31 |
| Respondent 4 | 1 | 30000 | 37 |
| Respondent 5 | 1 | 45000 | 56 |
| Respondent 6 | 1 | 35000 | 29 |
| Respondent 7 | 1 | 45000 | 30 |
| Respondent 8 | 0 | 15000 | 58 |
| Respondent 9 | 0 | 45000 | 66 |
| Respondent 10 | 0 | 45000 | 23 |
| Respondent 11 | 1 | 50000 | 34 |
| Respondent 12 | 0 | 25000 | 52 |
| Respondent 13 | 1 | 65000 | 32 |
| Respondent 14 | 1 | 60000 | 66 |
| Respondent 15 | 1 | 30000 | 22 |
| Respondent 16 | 0 | 45000 | 36 |
| Respondent 17 | 0 | 55000 | 48 |
| Respondent 18 | 0 | 25000 | 30 |
| Respondent 19 | 0 | 40000 | 32 |
| Respondent 20 | 1 | 70000 | 55 |
| Respondent 21 | 1 | 55000 | 56 |
| Respondent 22 | 0 | 25000 | 36 |
| Respondent 23 | 1 | 15000 | 26 |
| Respondent 24 | 1 | 50000 | 30 |
| Respondent 25 | 1 | 70000 | 31 |
| Respondent 26 | 1 | 70000 | 68 |
| Respondent 27 | 0 | 55000 | 60 |
| Respondent 28 | 0 | 55000 | 57 |
| Respondent 29 | 0 | 50000 | 60 |
| Respondent 30 | 0 | 30000 | 32 |
| Respondent 31 | 0 | 50000 | 28 |
| Respondent 32 | 0 | 30000 | 49 |
| Respondent 33 | 1 | 55000 | 38 |
| Respondent 34 | 0 | 65000 | 59 |
| Respondent 35 | 1 | 25000 | 30 |
| Respondent 36 | 0 | 20000 | 24 |
| Respondent 37 | 1 | 40000 | 20 |
| Respondent 38 | 1 | 20000 | 30 |
| Respondent 39 | 0 | 45000 | 59 |
| Respondent 40 | 0 | 70000 | 44 |

Segmentation Data / Discrimination Data / Sheet3 /

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 <u>within</u> 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 <u>within</u> 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 dendogram 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:

**Input**

Please select Segmentation Data.

The first column should contain respondent names or IDs.

Segmentation Data'!$B$4:$H$44

OK    Cancel

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 preselected.



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.

**Dendogram**

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

At one extreme (upper part of the dendogram), 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 dendogram), 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 dendogram, 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 dendogram 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.

**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.

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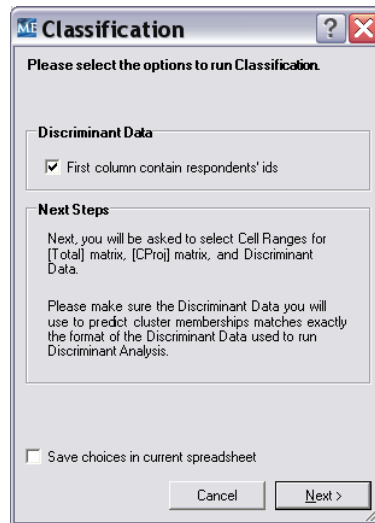
Once this discriminant analysis has been applied to the original dataset, it can be applied again to external customers for whom discriminant data –but no segmentation data- is available. The process of classifying customers among segments, based on a preceding segmentation analysis, but using discriminant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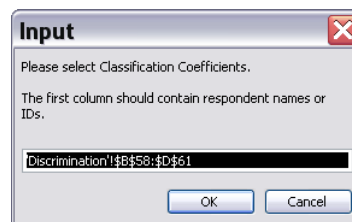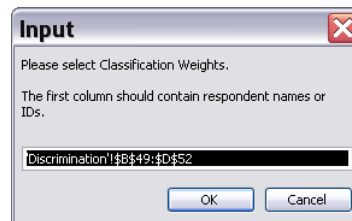
F.Y.I.

> 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.
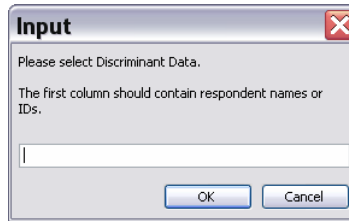
### Selecting data

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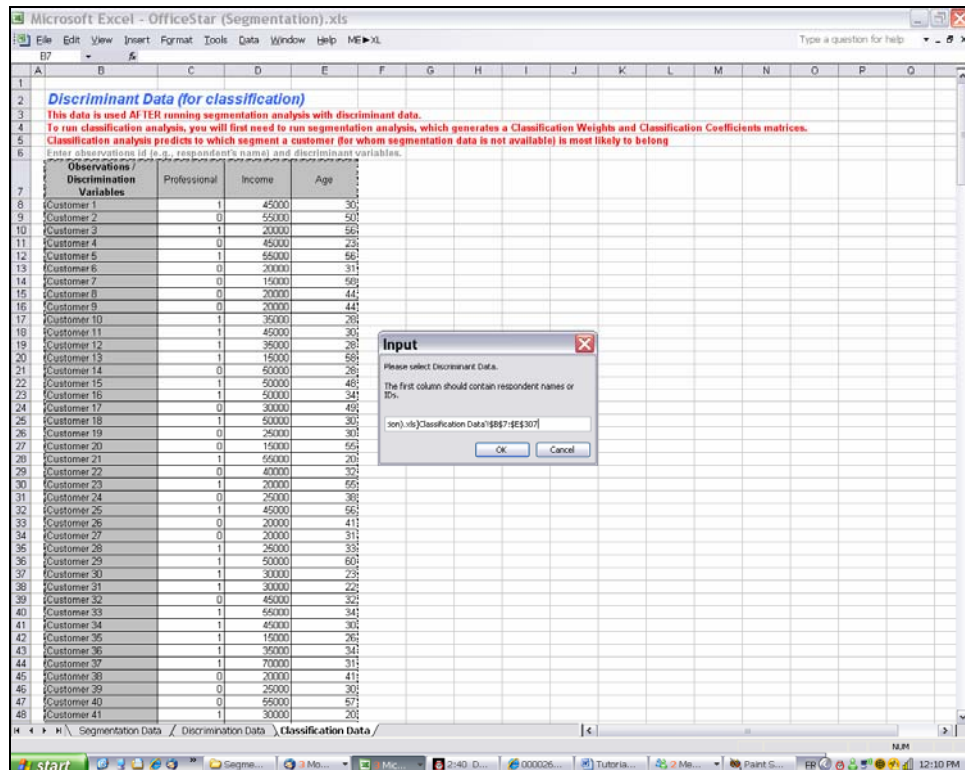
The first steps of the classification analysis consist of selecting two cell ranges: classification weights and classification coefficients. You can find that data at the bottom of the "discriminant" sheet in the analysis workbook generated by segmentation analysis. The cell ranges might be preselected.





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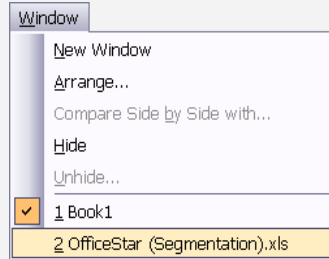
Discriminant data of "new" customers is available on the original OfficeStars workbook, in the last sheet. Go back to the OfficeStar workbook, and manually select the discriminant data available for the 300 additional customers (the last sheet of the workbook, named classification data).

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.



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*.