

*Hong Kong Baptist University
Department of Computer Science*

COMP 7810/4096 Business Intelligence (2019-20)

Data mining in SSAS

Introduction

Data Mining is a process to discover patterns for a large data set. It is an expert system that uses its historical experience (stored in relational databases or cubes) to predict the future.

Imagine that you own a company named Adventureworks. The company sells and manufactures bikes. You want to predict if a customer will buy a bike or not based on the customer information. Data Mining helps you to find the patterns and describe the characteristics of the customers with higher probability to buy the bikes or the lower probability.

Microsoft comes with a nice tool included in SQL Server Analysis Services (SSAS) for creating sophisticated data mining solutions. The tools in Analysis Services help you design, create, and manage data mining models that use either relational or cube data. You can manage client access to the data mining models and create prediction queries from multiple clients.

SSAS contains the features and tools you need to create complex data mining solutions.

- A set of industry-standard data mining algorithms.
- The Data Mining Designer, which you can use to create, manage, and explore data mining models, and then create predictions by using those models.
- The Data Mining Extensions (DMX) language, which you can use to manage mining models and to create complex prediction queries.

Some common data mining algorithms:

- *Clustering*: is a technique to create different groups of people according to their characteristics or patterns. It is a segmentation technique that divides the customers into different groups. We can identify natural groupings of customers.
- *Decision trees*: uses branches to classify the information.

Learning Outcomes

By finishing this lab session, you should be able to:

- Create a Data Source
- Create a Data View
- Create a Data Mining Project
- Predict information using the Mining Model

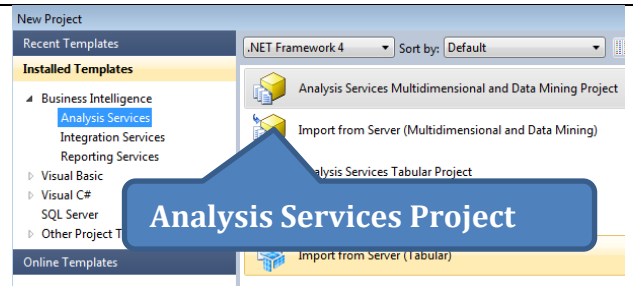
Tools

- Microsoft SQL Server Management Studio 2012
- Visual Studio 2010 with SQL Server Data Tools (SSDT)

Part A: Create data mining task

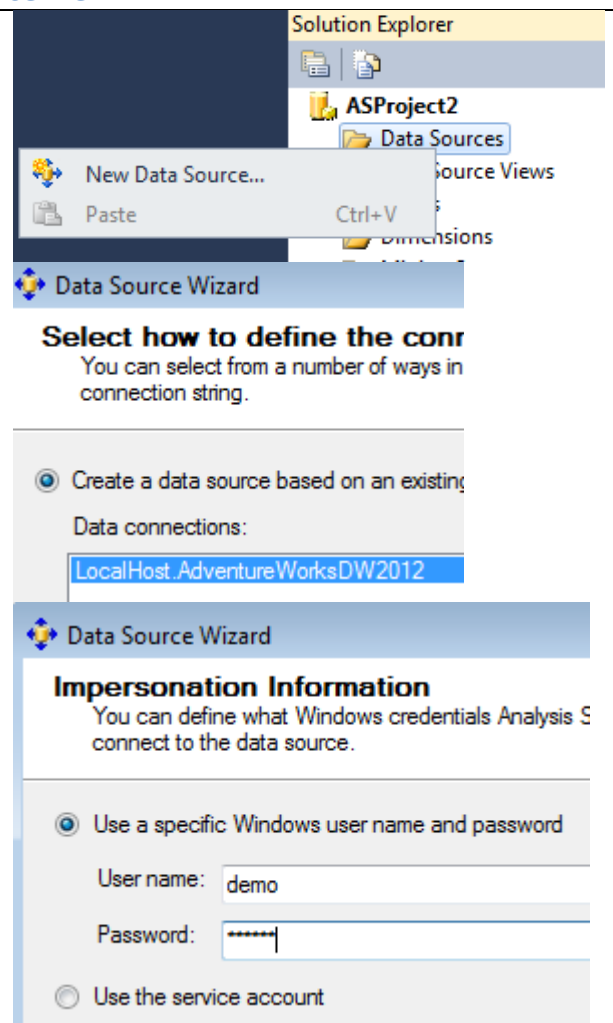
I. Create a new Analysis Services project

1. Open the **Visual Studio with SSDT**.
2. Select **File → New → Project**.
3. Expand **Business Intelligence → Analysis Services**, and then click **Analysis Services Multidimensional and Data Mining Project**.
4. Change the project name to **ASProject2**. Press **OK**.



II. Define a new data source and data source view

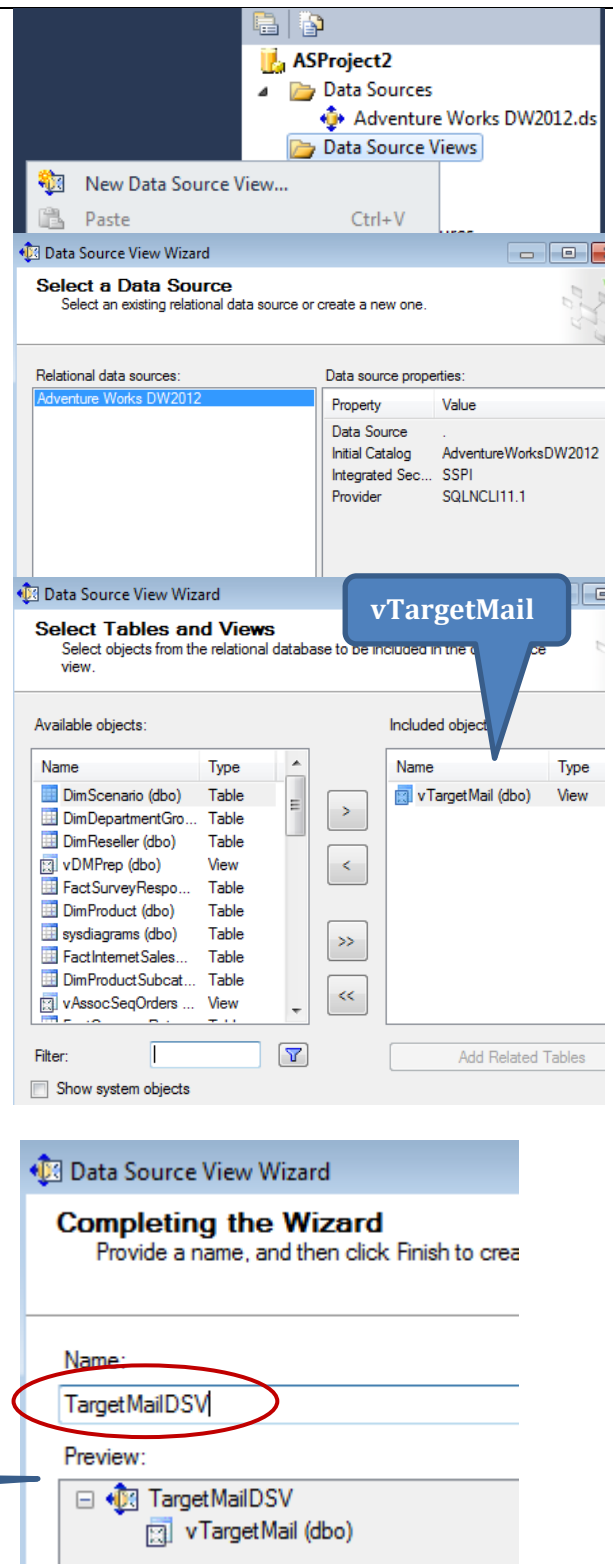
1. In **Solution Explorer**, right click **Data Sources**, select **New Data Source** then press **Next**.
2. Choose the data connections to connect with **AdventureWorksDW2012**, press **Next**.
3. In the *Impersonation Information* dialog box, select **the first option** and enter **your Windows account information**
4. Press **Next**. Then **Finish**.



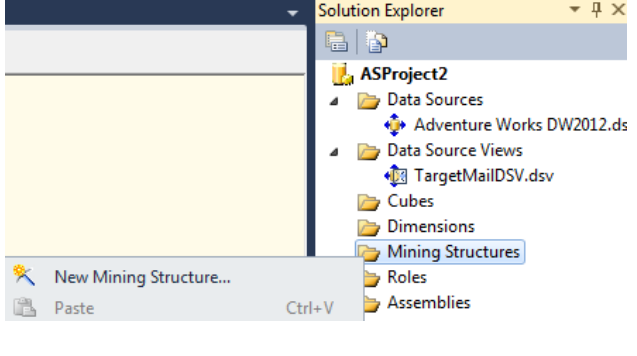
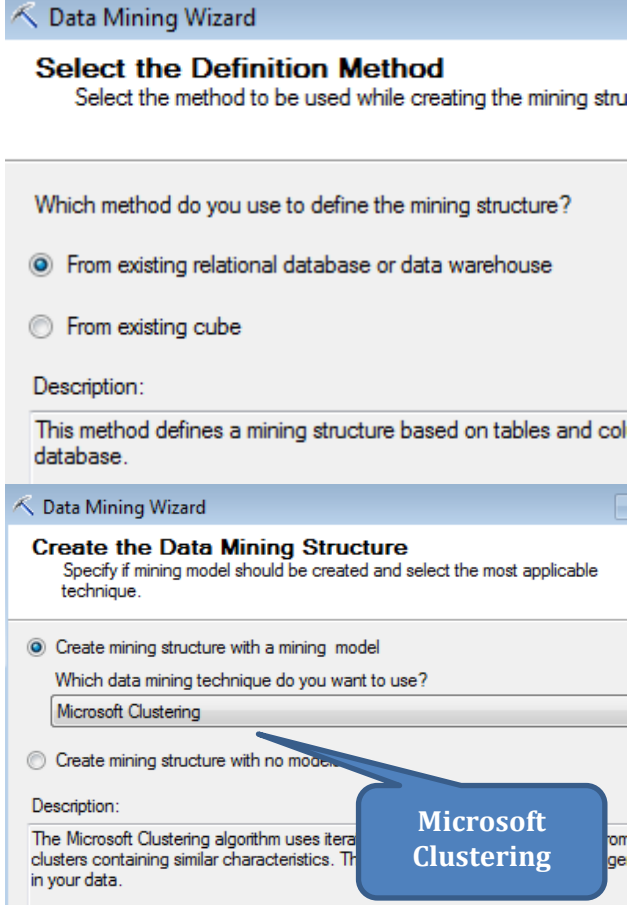
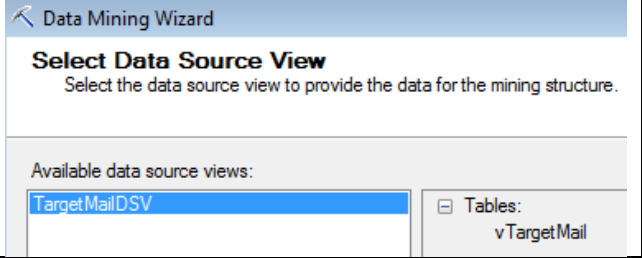
5. In the **Solution Explorer**, right click **Data Source Views** and select **New Data Source View**. Press **Next** in the Welcome page.
6. In the *Select a Data Source* window, select **Adventure Works DW 2012**. Press **Next**.
7. In the *Select Tables and Views*, choose **vTargetMail** as an included object.
8. **Rename** the DSV as **TargetMailDSV**, then press **Finish**.

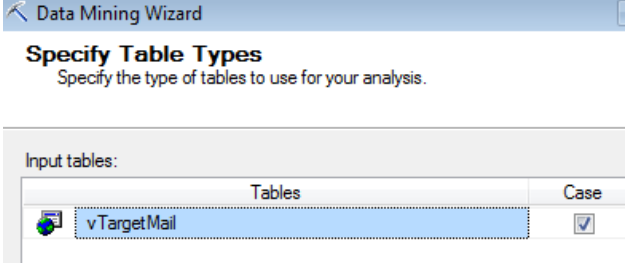
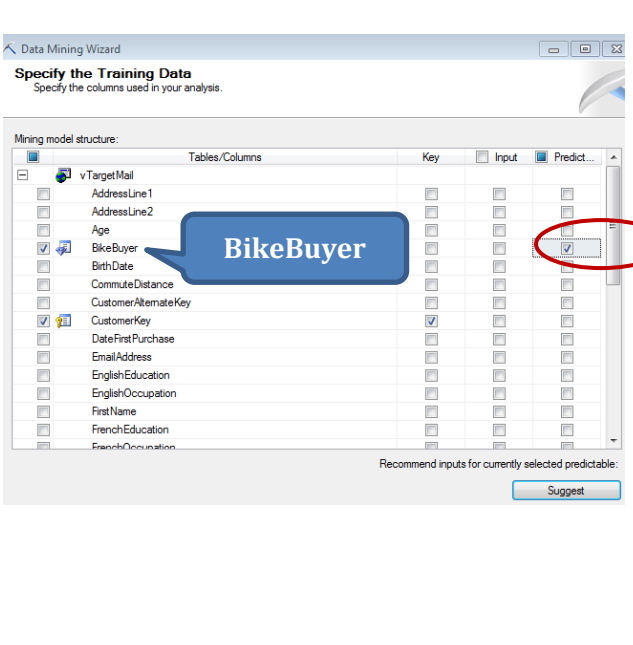
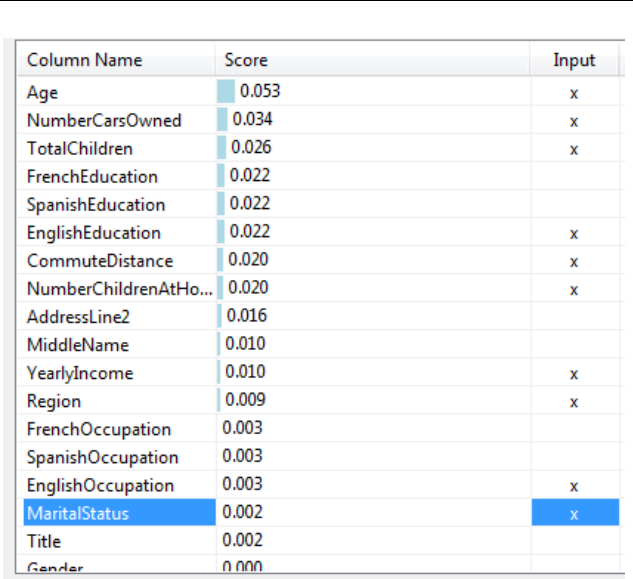
(You may explore the data inside)

We just created a Data View with the view to give experience to our Data Mining Model. The vTargetMail is a view that contains historical data about the customers. Using that experience, our mining model, will predict the future.



III. Create New Mining Structures and New Mining Model

<ol style="list-style-type: none"> 1. Right click the Mining Structures folder and select New Mining Structure in the Solution Explorer. 2. It opens up a <i>Data Mining Wizard</i>, you can use this wizard to create new mining model for making prediction. Press Next to continue. 	
<ol style="list-style-type: none"> 3. In the <i>select the Definition Method</i> window, select the option From existing relational database or data warehouse to define mining structure. Press Next to continue. 4. In the <i>Create the Data Mining Structure</i> window, select Create mining structure with a mining model → choose Microsoft Clustering as the data mining technique. Press Next. Microsoft Clustering algorithm uses iterative techniques to <u>group records from a dataset into clusters containing similar characteristics</u>. It is used to find general groupings in your data. 	
<ol style="list-style-type: none"> 5. In the <i>Select Data Source View</i> window, choose TargetMailDSV. Press Next to continue. 	

<p>6. Select vTargetMail as the input table by default. Press Next to continue.</p> <p>(This becomes the only one case table)</p>	 <p>Data Mining Wizard Specify Table Types Specify the type of tables to use for your analysis.</p> <p>Input tables:</p> <table border="1"> <thead> <tr> <th>Tables</th> <th>Case</th> </tr> </thead> <tbody> <tr> <td>vTargetMail</td> <td><input checked="" type="checkbox"/></td> </tr> </tbody> </table>	Tables	Case	vTargetMail	<input checked="" type="checkbox"/>																																																																
Tables	Case																																																																				
vTargetMail	<input checked="" type="checkbox"/>																																																																				
<p>7. Specify the column(s) used in the analysis (which information you want to predict?). BikeBuyer is the column in the mining model used for prediction (predictable column). It represents whether the customer will buy the bike or not.</p> <p>8. Press the button Suggest to specify the input columns for the mining model. After sampling, it shows the importance of the attributes within the model with scores greater than or equal to 0.</p>	 <p>Data Mining Wizard Specify the Training Data Specify the columns used in your analysis.</p> <p>Mining model structure:</p> <table border="1"> <thead> <tr> <th>Tables/Columns</th> <th>Key</th> <th>Input</th> <th>Predict...</th> </tr> </thead> <tbody> <tr><td>vTargetMail</td><td></td><td></td><td></td></tr> <tr><td>AddressLine1</td><td></td><td></td><td></td></tr> <tr><td>AddressLine2</td><td></td><td></td><td></td></tr> <tr><td>Age</td><td></td><td></td><td></td></tr> <tr><td>BikeBuyer</td><td></td><td></td><td><input checked="" type="checkbox"/></td></tr> <tr><td>BirthDate</td><td></td><td></td><td></td></tr> <tr><td>CommuteDistance</td><td></td><td></td><td></td></tr> <tr><td>CustomerAlternateKey</td><td></td><td></td><td></td></tr> <tr><td>CustomerKey</td><td><input checked="" type="checkbox"/></td><td></td><td></td></tr> <tr><td>DateFirstPurchase</td><td></td><td></td><td></td></tr> <tr><td>EmailAddress</td><td></td><td></td><td></td></tr> <tr><td>EnglishEducation</td><td></td><td></td><td></td></tr> <tr><td>EnglishOccupation</td><td></td><td></td><td></td></tr> <tr><td>FirstName</td><td></td><td></td><td></td></tr> <tr><td>FrenchEducation</td><td></td><td></td><td></td></tr> <tr><td>FrenchOccupation</td><td></td><td></td><td></td></tr> </tbody> </table> <p>Recommend inputs for currently selected predictable:</p> <p>Suggest</p>	Tables/Columns	Key	Input	Predict...	vTargetMail				AddressLine1				AddressLine2				Age				BikeBuyer			<input checked="" type="checkbox"/>	BirthDate				CommuteDistance				CustomerAlternateKey				CustomerKey	<input checked="" type="checkbox"/>			DateFirstPurchase				EmailAddress				EnglishEducation				EnglishOccupation				FirstName				FrenchEducation				FrenchOccupation			
Tables/Columns	Key	Input	Predict...																																																																		
vTargetMail																																																																					
AddressLine1																																																																					
AddressLine2																																																																					
Age																																																																					
BikeBuyer			<input checked="" type="checkbox"/>																																																																		
BirthDate																																																																					
CommuteDistance																																																																					
CustomerAlternateKey																																																																					
CustomerKey	<input checked="" type="checkbox"/>																																																																				
DateFirstPurchase																																																																					
EmailAddress																																																																					
EnglishEducation																																																																					
EnglishOccupation																																																																					
FirstName																																																																					
FrenchEducation																																																																					
FrenchOccupation																																																																					
<p>9. Check the following ten attributes as input columns, the press OK:</p> <ul style="list-style-type: none"> • Age • NumberCarsOwned • TotalChildren • EnglishEducation • CommuteDistance • NumberChildrenAtHome • YearlyIncome • Region • EnglishOccupation • MaritalStatus <p>Press Next to continue.</p>	 <table border="1"> <thead> <tr> <th>Column Name</th> <th>Score</th> <th>Input</th> </tr> </thead> <tbody> <tr><td>Age</td><td>0.053</td><td>x</td></tr> <tr><td>NumberCarsOwned</td><td>0.034</td><td>x</td></tr> <tr><td>TotalChildren</td><td>0.026</td><td>x</td></tr> <tr><td>FrenchEducation</td><td>0.022</td><td></td></tr> <tr><td>SpanishEducation</td><td>0.022</td><td></td></tr> <tr><td>EnglishEducation</td><td>0.022</td><td>x</td></tr> <tr><td>CommuteDistance</td><td>0.020</td><td>x</td></tr> <tr><td>NumberChildrenAtHome</td><td>0.020</td><td>x</td></tr> <tr><td>AddressLine2</td><td>0.016</td><td></td></tr> <tr><td>MiddleName</td><td>0.010</td><td></td></tr> <tr><td>YearlyIncome</td><td>0.010</td><td>x</td></tr> <tr><td>Region</td><td>0.009</td><td>x</td></tr> <tr><td>FrenchOccupation</td><td>0.003</td><td></td></tr> <tr><td>SpanishOccupation</td><td>0.003</td><td></td></tr> <tr><td>EnglishOccupation</td><td>0.003</td><td>x</td></tr> <tr><td>MaritalStatus</td><td>0.002</td><td>x</td></tr> <tr><td>Title</td><td>0.002</td><td></td></tr> <tr><td>Gender</td><td>0.000</td><td></td></tr> </tbody> </table>	Column Name	Score	Input	Age	0.053	x	NumberCarsOwned	0.034	x	TotalChildren	0.026	x	FrenchEducation	0.022		SpanishEducation	0.022		EnglishEducation	0.022	x	CommuteDistance	0.020	x	NumberChildrenAtHome	0.020	x	AddressLine2	0.016		MiddleName	0.010		YearlyIncome	0.010	x	Region	0.009	x	FrenchOccupation	0.003		SpanishOccupation	0.003		EnglishOccupation	0.003	x	MaritalStatus	0.002	x	Title	0.002		Gender	0.000												
Column Name	Score	Input																																																																			
Age	0.053	x																																																																			
NumberCarsOwned	0.034	x																																																																			
TotalChildren	0.026	x																																																																			
FrenchEducation	0.022																																																																				
SpanishEducation	0.022																																																																				
EnglishEducation	0.022	x																																																																			
CommuteDistance	0.020	x																																																																			
NumberChildrenAtHome	0.020	x																																																																			
AddressLine2	0.016																																																																				
MiddleName	0.010																																																																				
YearlyIncome	0.010	x																																																																			
Region	0.009	x																																																																			
FrenchOccupation	0.003																																																																				
SpanishOccupation	0.003																																																																				
EnglishOccupation	0.003	x																																																																			
MaritalStatus	0.002	x																																																																			
Title	0.002																																																																				
Gender	0.000																																																																				

10. Press **Detect** button to change the data types of the columns. Press **Next** to continue.

Specify Columns' Content and Data Type

Specify mining structure columns' content and data type.

Mining model structure:

Columns	Content Type	Data Type
Age	Continuous	Long
Bike Buyer	Discrete	Long
Commute Distance	Discrete	Text
Customer Key	Key	Long
English Education	Discrete	Text
English Occupation	Discrete	Text
Marital Status	Discrete	Text
Number Cars Owned	Discrete	Long
Number Children At Home	Discrete	Long
Region	Discrete	Text
Total Children	Discrete	Long
Yearly Income	Continuous	Double

Detect continuous or discrete for numeric columns:

Detect

11. Here you can create testing data for testing the accuracy of the DM model, putting **30%** of data for testing means using 70% of data for training the model.

Note: Input data will be randomly split into two sets, a training set and a testing set, based on the percentage of data for testing

- *Training set* is used to create the mining model.
- *Testing set* is used to check model accuracy.

Press **Next** to continue.

Data Mining Wizard

Create Testing Set

Specify the number of cases to be reserved for model testing.

Percentage of data for testing:

30 %

Maximum number of cases in testing data set:

Description:

Input data will be randomly split into two sets, a training set and a testing set, based on the percentage of data for testing and maximum number of cases in testing data set you provide. The training set is used to create the mining model. The testing set is used to check model accuracy.

[Percentage of data for testing] specifies percentages of cases reserved for testing set. [Maximum number of cases in testing data set] limits total number of cases in the testing set. If both values are specified, both limits are enforced.

12. Finally provide a name for the structure. Then press **Finish**.

Mining structure name:

- **TargetMailMS**

Mining model name:

- **TargetMailMM**

13. Check the box **Allow drill through** and press **Finish**

Data Mining Wizard

Completing the Wizard

Completing the Data Mining Wizard by providing a name for the mining structure.

Mining structure name:

TargetMailMS

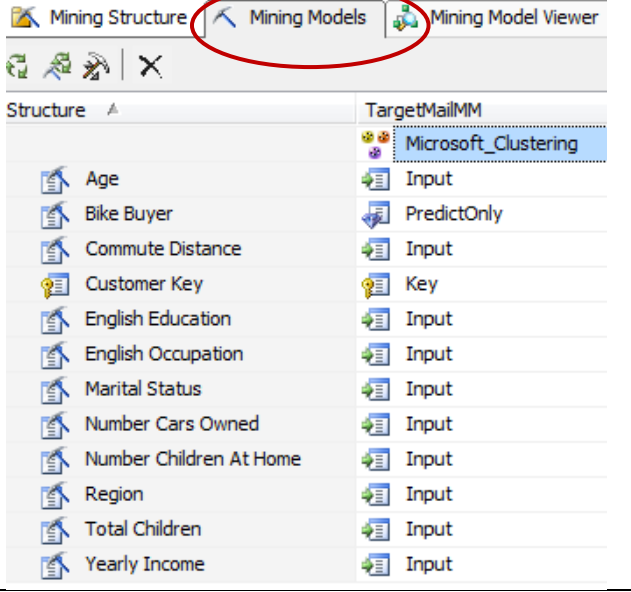
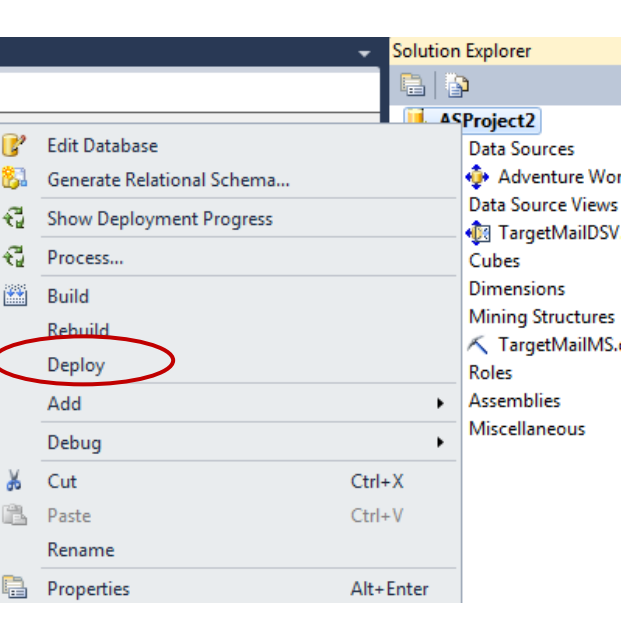
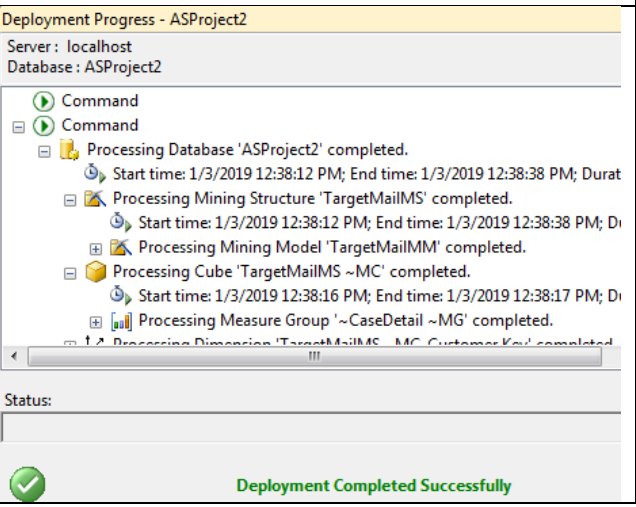
Mining model name:

TargetMailMM

☒ Allow drill through

Preview:

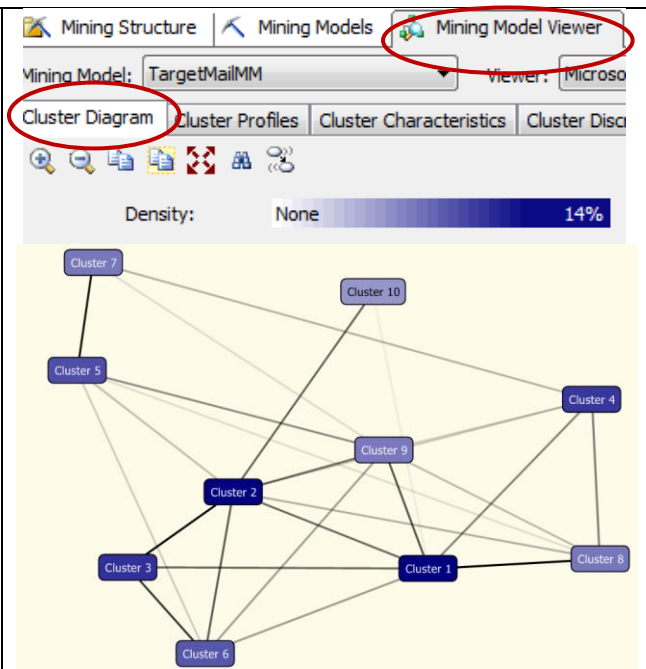
TargetMailMS
Columns
Age
Bike Buyer
Commute Distance
Customer Key

<p>14. Go to Mining Models tab, here you can create another data mining model by choosing another data mining algorithm like <i>Decision Tree</i>.</p>	
<p>15. Go to the Solution Explorer, deploy the project ASProject2 to the server. The process will take 70% of data as training data.</p>	
<p>16. You will see the deployment is successful.</p>	

17. Select **Mining Model Viewer** tab, the mining model TargetMailMM is created. A **Cluster Diagram** is shown.

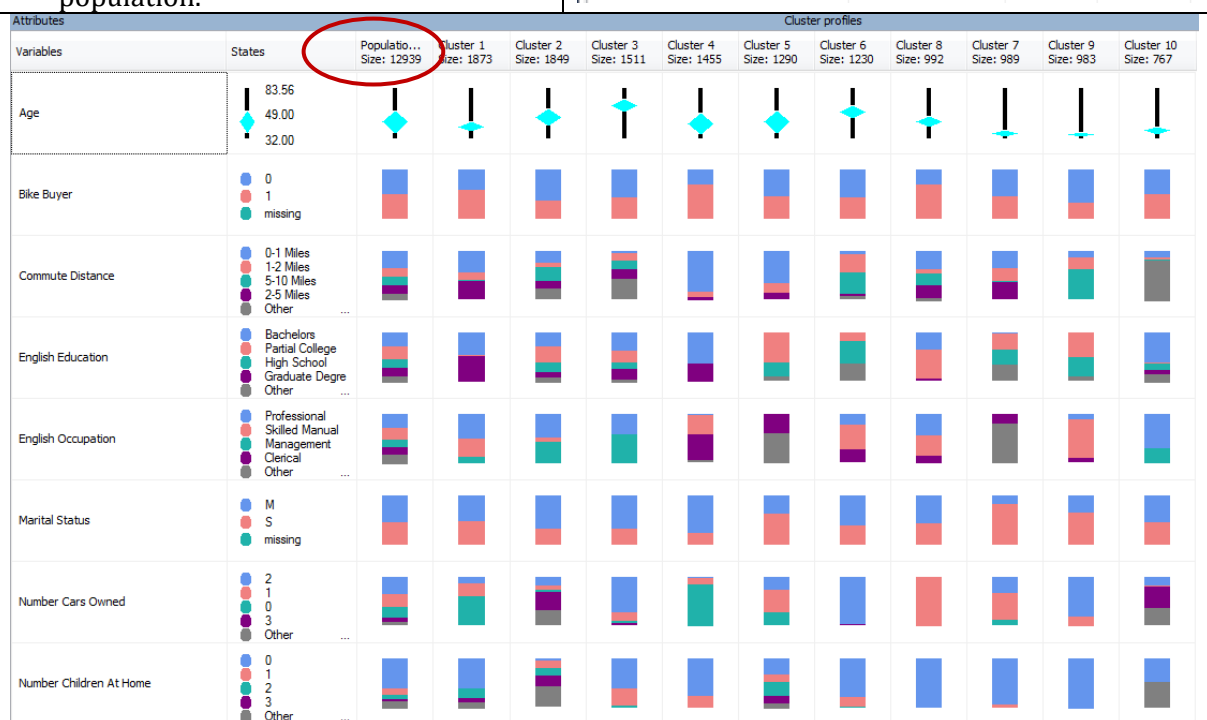
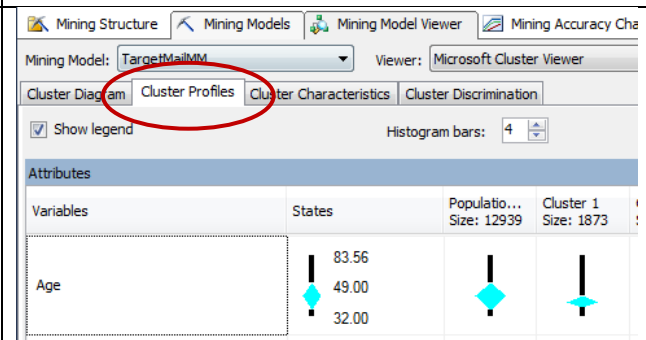
18. Clustering creates different groups for all the customers. The groups are named cluster 1, cluster 2 and so on. The clusters are created based on customers' characteristics. Darker colors means **higher density clusters**.

(You can change the shading variable to Bike Buyer and state to 1)

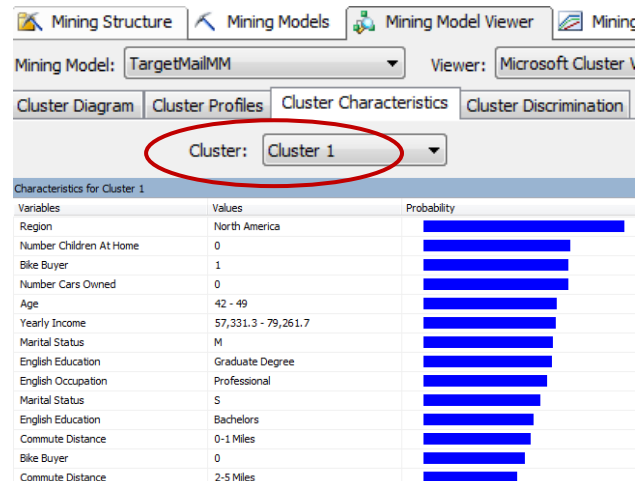


19. Go to **Cluster Profiles** tab, it shows different **variables** and the population for each cluster. The **variables** show the customer's characteristics like the age, yearly income etc

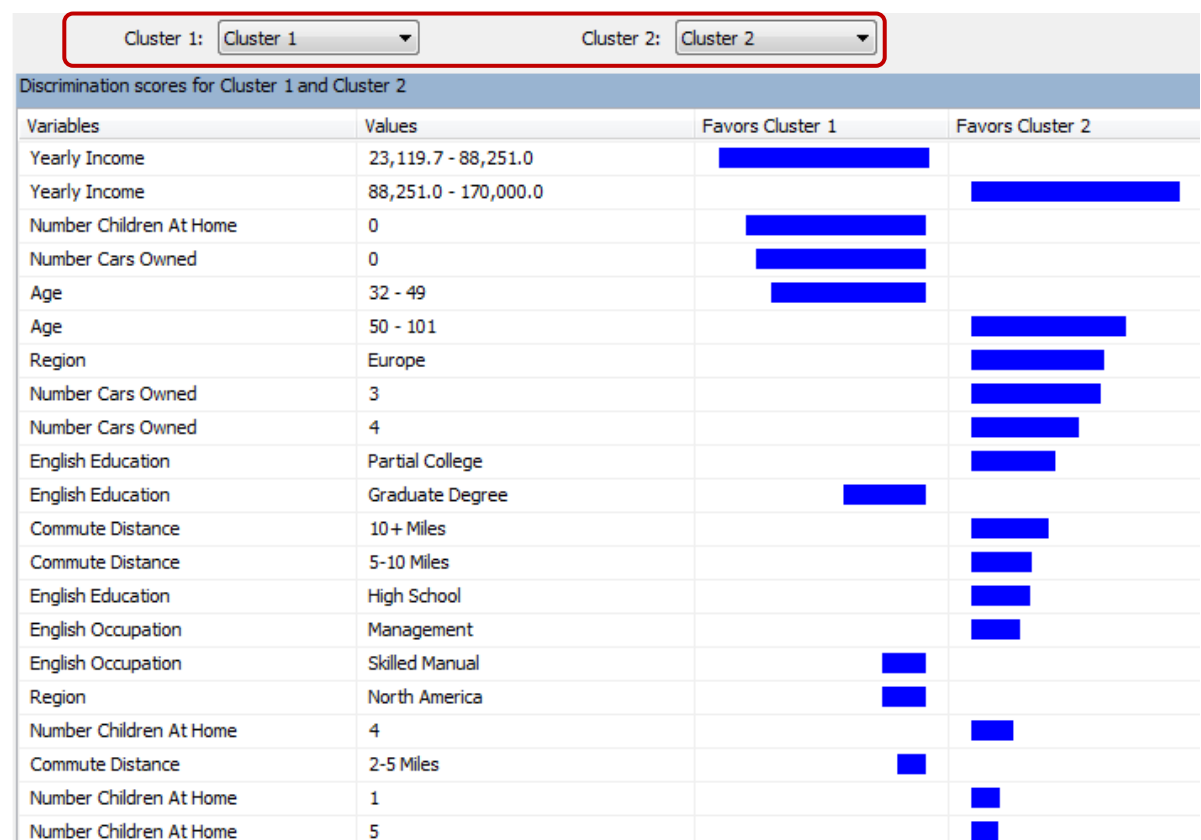
20. The total population is 12939. The *cluster 1* is the *most populated* cluster. 10 clusters are formed according to the population.



21. **Cluster characteristics** tab shows the characteristics per cluster. Select cluster 1, to find out the main characteristics of this cluster.

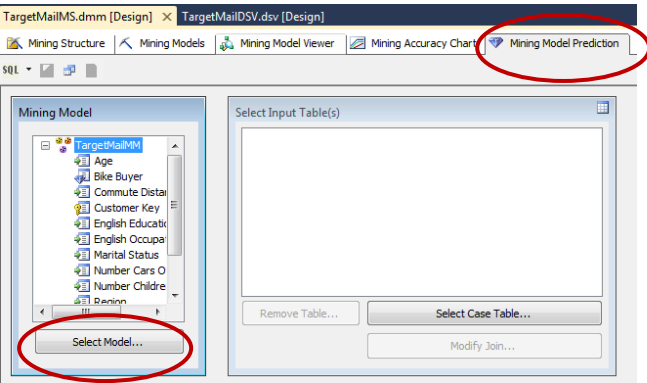
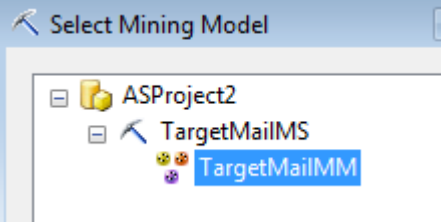
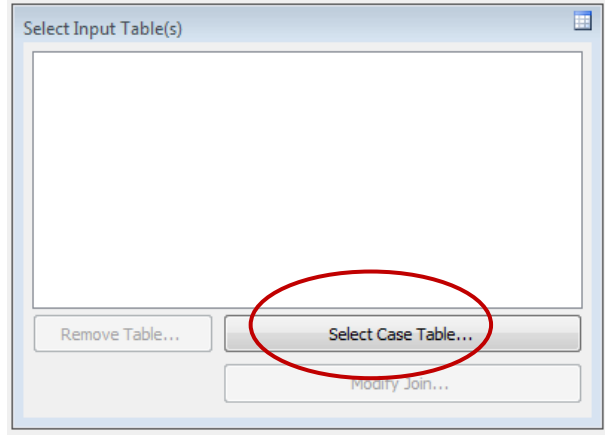
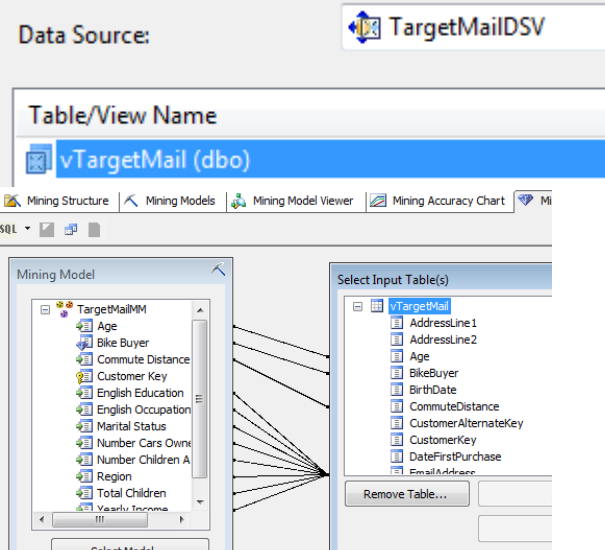


22. **Cluster Discrimination** tab allows you to compare the differences between two clusters.



IV. Make a Prediction

Here are the steps to predict the probability of the customer to buy a bike. 1 represents buy and 0 represents not buy.

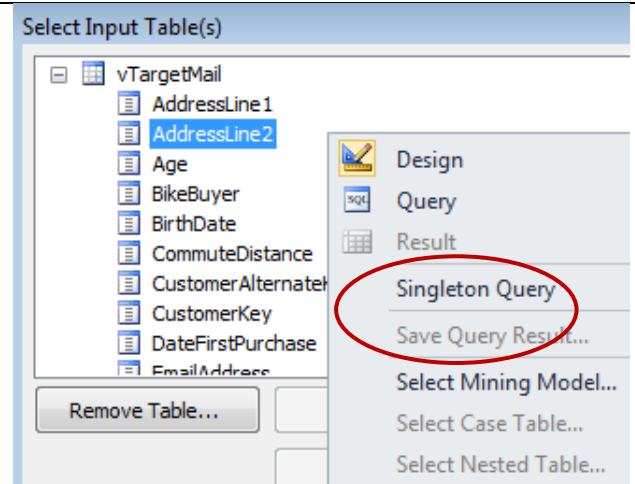
<p>1. Click the Mining Model Prediction tab, in the Mining Model box, press the button Select Model.</p>	
<p>2. In the Select Mining Model box, expand TargetMailMS and select the model TargetMailMM. Then press OK.</p>	
<p>3. In the <i>Select Input Table(s)</i> box, click the Select Case Table</p>	
<p>4. Select the vTargetMail and press OK</p>	

5. **Right click** in the **Select Input Table(s)** and select **Singleton Query**

6. In the **Singleton Query**, specify the following information:

- Age: 45
- Commute Distance: 0-1 Miles
- English Education: Bachelors
- English Occupation: Professional
- Marital Status: M
- Number Cars Owned: 1
- Number Children At Home: 1
- Region: Europe
- Total Children: 1
- Yearly Income: Missing

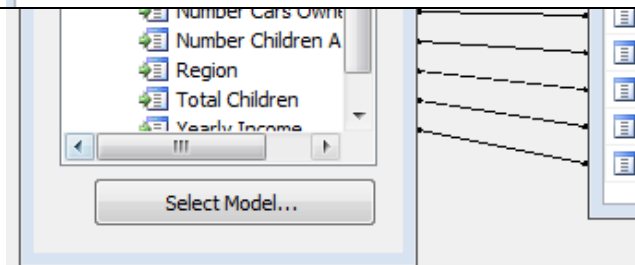
In this step we are specifying the customer characteristics.



Singleton Query Input

Mining Model Column	Value
Age	45
Bike Buyer	
Commute Distance	0-1 Miles
English Education	Bachelors
English Occupation	Professional
Marital Status	M
Number Cars Owned	1
Number Children At Home	1
Region	Europe
Total Children	1
Yearly Income	

7. In the query grid, for **Source**, select the **TargetMailMM** mining model. (Field is BikeBuyer)



Source	Field	Alias
TargetMailMM	Bike Buyer	

Source	Field	Alias	Show	Group	And/Or	Criteria/Argument
TargetMailMM	Bike Buyer		<input checked="" type="checkbox"/>			
Prediction Function	Cluster		<input checked="" type="checkbox"/>			
Prediction Function	PredictHistogram	Results	<input checked="" type="checkbox"/>			[TargetMailMM].[Bike Buyer]

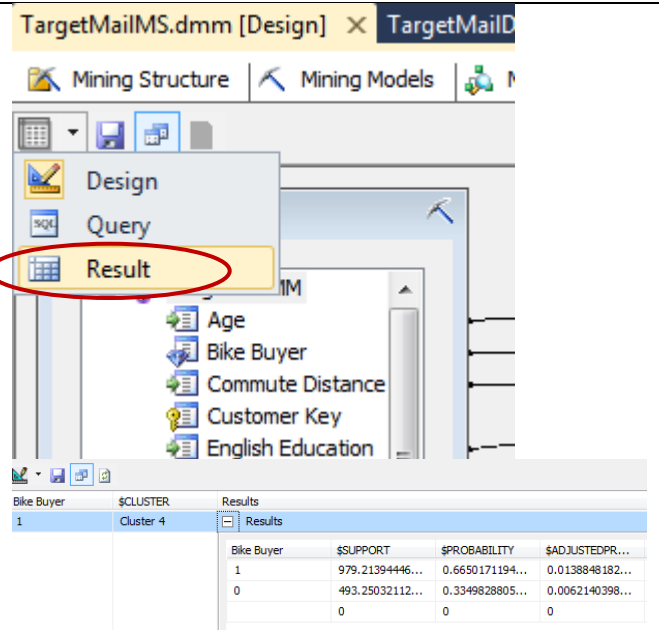
8. In the second row of the **Source** column, select the **Prediction Function**, and add the function **Cluster**.

9. In the second row of the **Source** column, select **Prediction function**, add the function **PredictHistogram**, drag the model column **[Bike Buyer]** into the **Criteria/Argument** box

10. Type **Results** as alias



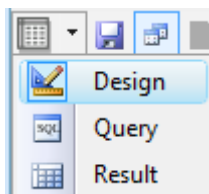
11. Click on the **Switch to query result view** icon and select **Result** to verify the results of the query.
12. The results show the probability to buy a bike is ~66%
13. Select **File** → **Save All** to save the project **ASProject2**



Bike Buyer	\$SUPPORT	\$PROBABILITY	\$ADJUSTEDPR...
1	979.21394446...	0.6650171194...	0.0138848182...
0	493.25032112...	0.3349828805...	0.0062140398...
0	0	0	0

V. Exercise

1. Go back to the **Design** View, try to predict another user with the following characteristics will buy a bike or not, submit the result as **online texts** in **buelearning** website.



- Age: 35
- Commute Distance: 1-2 Miles
- English Education: Bachelors
- English Occupation: Clerical
- Marital Status: S
- Number Cars Owned: 1
- Number Children At Home: 0
- Region: North America
- Total Children: 0

The probability that the customer would buy a bike is:

55.6%

The probability that the customer would NOT buy a bike is:

44.3%

VI. Answer Submission

1. **Zip** your Analysis Services projects (**ASProject2 folder with ASProject2.sln**) that you created in C:\users\demo\documents\visual studio 2010\Projects. The default file name is **ASProject2.zip**
2. Submit the following files to the site <http://buelearning.hkbu.edu.hk/>
 - **ASProject2.zip**
 - **Online texts with probabilities**