# UNIVERSITÉ ANTONINE
## Faculté d'ingénieurs en Informatique, Multimédia, Réseaux & Télécommunications



## Data Mining – Implementation of Sequence Clustering

# Matière : System décisionnel

| Effectué par : | NOM Prénom | INF# | Option |
|---|---|---|---|
| | MATTA Elie | Privacy | OGL |
| | et al. | applied | |

# Introduction

This is a follow up report for our project: Sequence Clustering.

We would like to thank our teacher who gave us an extended week to complete our implementation. We based our essential work on the document of the "System decisionnel" which gave us a wider and a better idea about the right steps to complete our implementation.

Our project is about Sequence Clustering algorithm, to make it even easier to understand such an algorithm we decided to implement this algorithm in "Business Intelligence Development Studio" using the AdvanduteWorsDW database.

In the proceeding pages, we will guide you step-by-step into all the different parts and steps we ran into while implementing such algorithm using the upcoming screenshots.

Note that more valuable information is included in the video "**sequence clustering implementation.avi**" included in the DVD, so don't hesitate to check the video recorded to follow the easy step-by-step documentation.
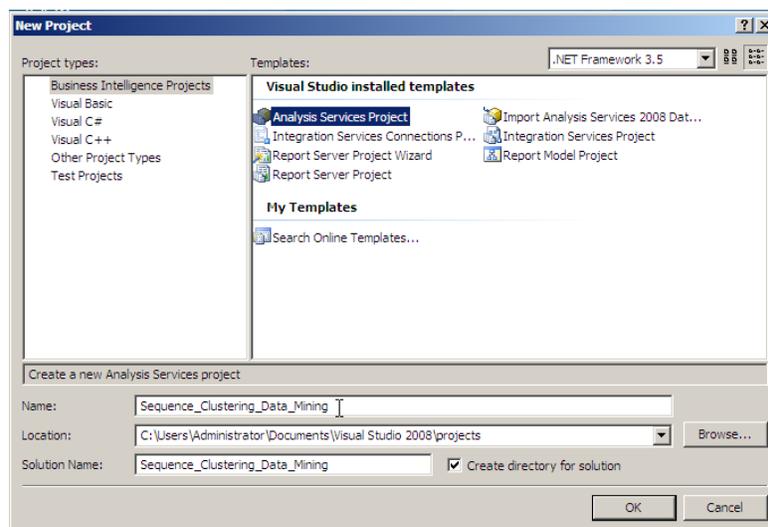
# Implementation

## I. Creating the Analysis Service Project

Step1:

- Open Business Intelligence Development Studio.
- Select New→ Project from the File menu.
- Select Analysis Services Project. We have named it "Sequence_Clustering_Data_Mining" as shown in Figure 1.
- Then click OK.



**Figure 1.** New Analysis Services Project

## II. Creating the Data Source

Step2:

- In the solution explorer, right click on "Data source" → new data source. "Data Source Wizard" dialog box opens.
- On the welcome page click Next.
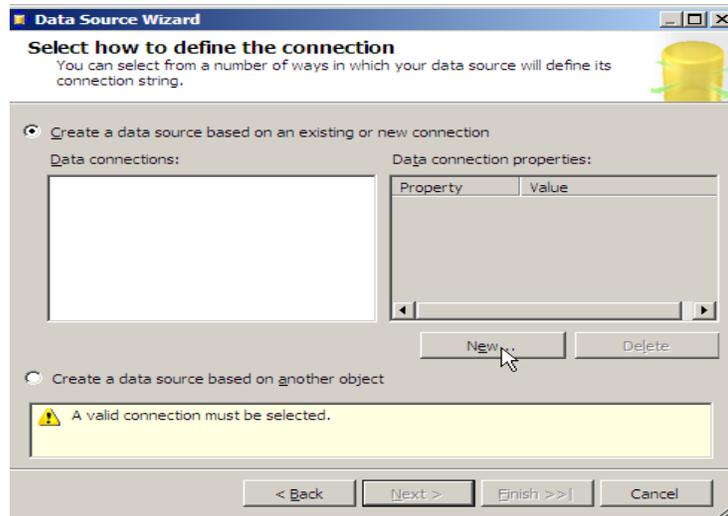- Click New to add a connection to the AdventureWorksDW database as shown in Figure 2.

---

**Figure 2.** New Data Source

Step3:

"Connection Manager" dialog box opens:

- In the Server name drop-down list, select the server where AdventureWorksDW is hosted (localhost="**.**")
- Use windows authentification.
- Select the database from the dropdown list → the **AdventureWorksDW** database
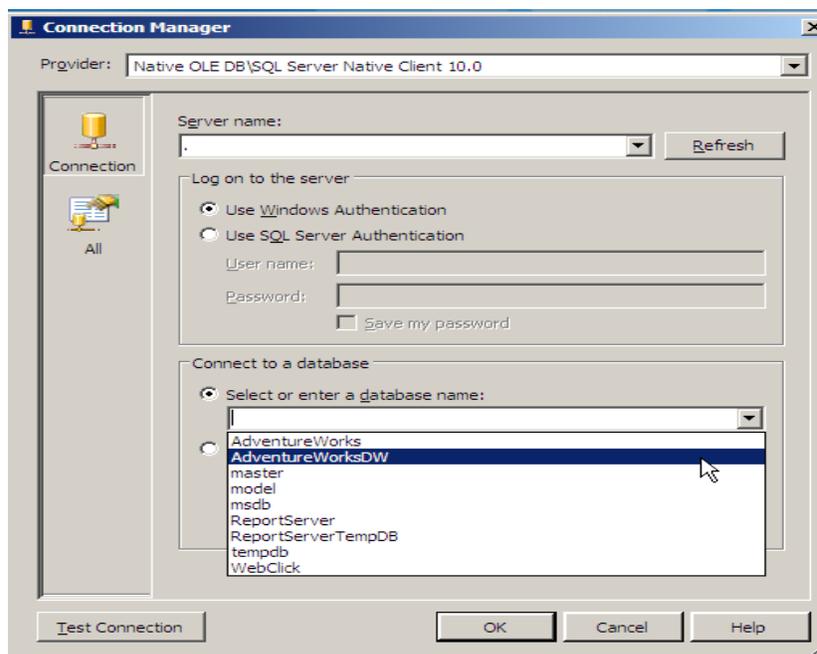- Click OK to close the "Connection manager" dialog box.



**Figure 3.** Connection Manager

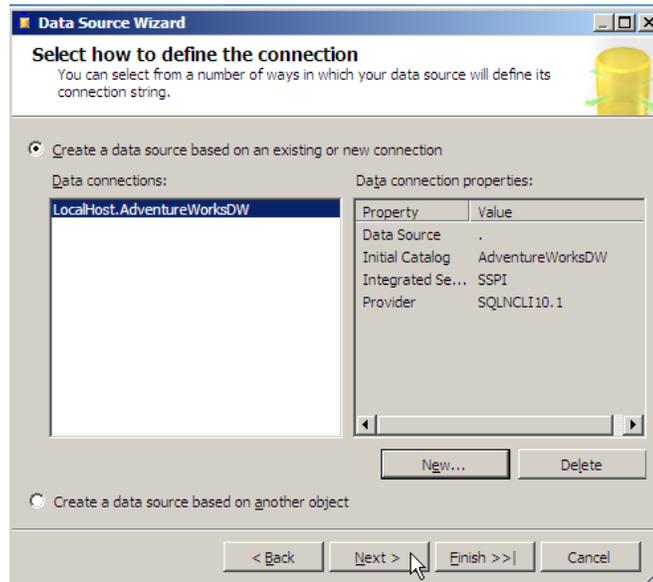Then you see the following figure then you click Next.



**Figure 4.** Choosing the Data Connection

Step 4:
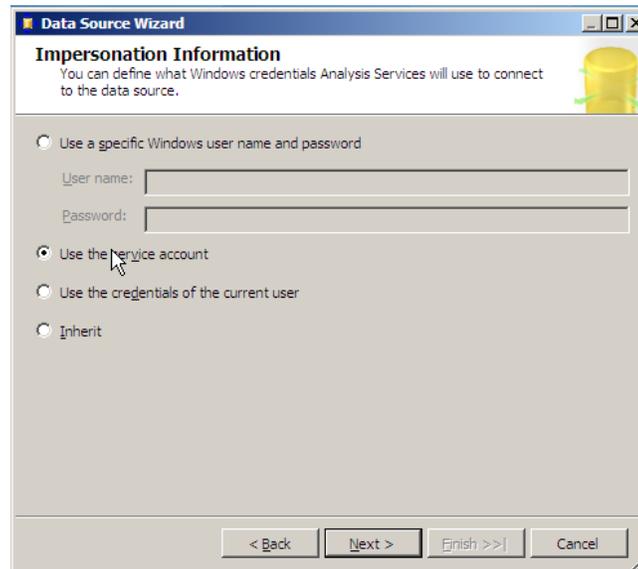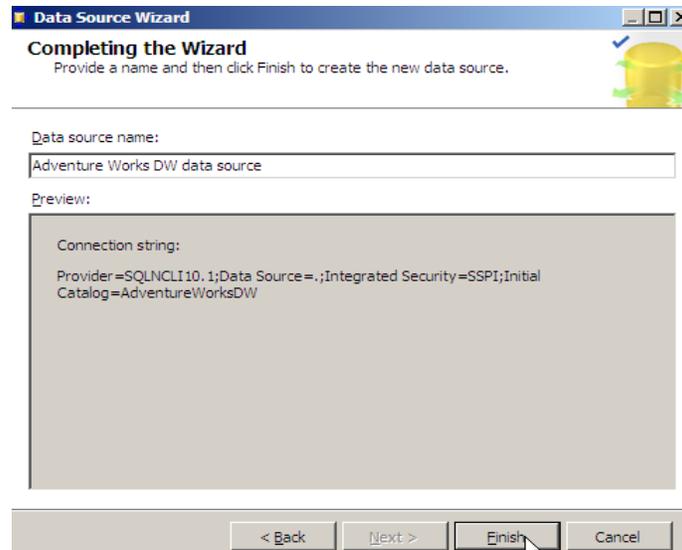- Choose "Use the service account" combo box as shown in Figure 5.
- Click Next.



**Figure 5.** using the service account.

Step 5:

Figure 6 shows this step.

- Give a name for the data source. We have named it "Adventure Works DW data source".
- Click Finish.



**Figure 6.** Naming the data source

## III. Creating Data Source View

Step 6:

→See figure 7.

In the solution explorer, right click on "Data Source View"→ new data source view.

"Data Source View Wizard" dialog box opens.

- On the welcome page click Next.
- The "Adventure Works DW data source" data source created before is selected by default in the Relational data sources window.
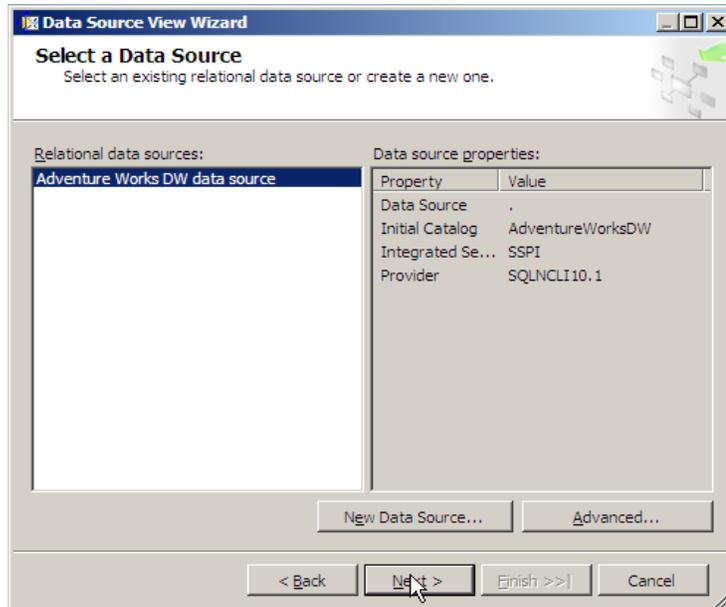- Click Next.

Préparé par Elie Matta et al.


**Figure 7.** Selecting the data source

Step 7:
Here we choose what tables do we need for the model. So we have chosen:
- vAssocSeqOrders table and
- vAssocSeqLineItems table.
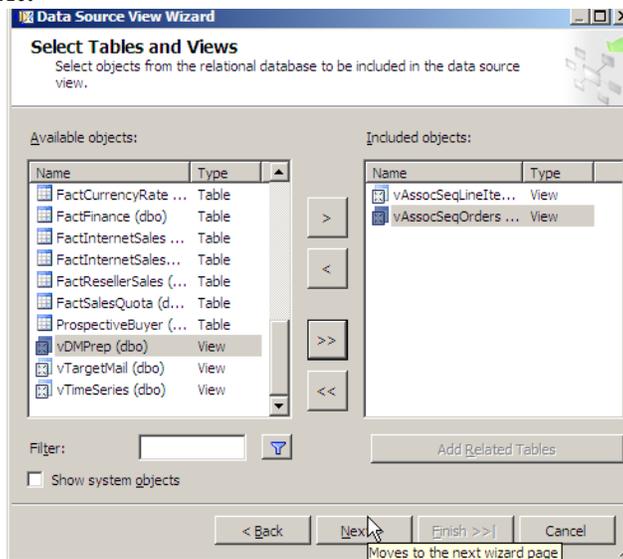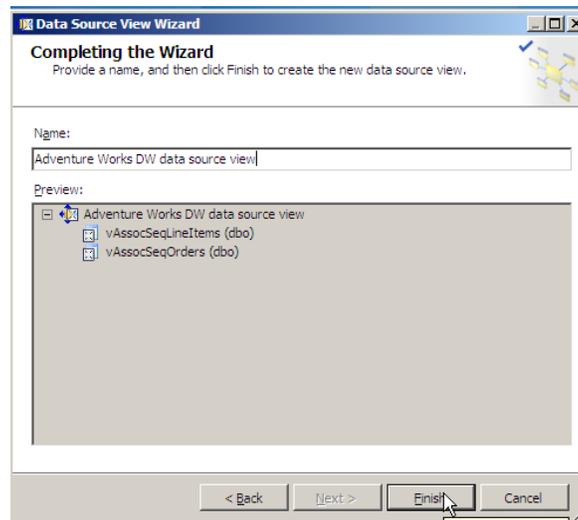
As shown in Figure 8.
- Then click next.


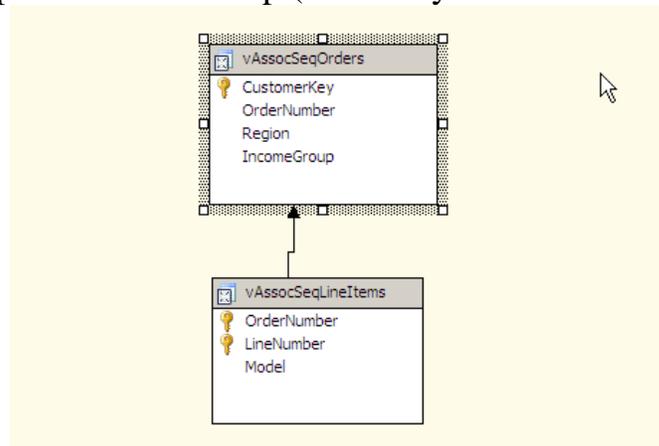**Figure 8.** Select tables and views

Step 8:
This step is shown in Figure 9.
- Give a name for the data source view. We have named it "Adventure Works DW data source view".
- Click Finish.

**Figure 9.** Naming the data source view

So we have the following view shown in Figure 10. Do not forget to relate these two tables with a specific relationship (related by OrderNumber).


**Figure 10.** Adventure Works DW Data Source View

## IV. Creating the mining structure

Step 9:
- In the solution explorer, right click "Mining Stucture" → new mining structure. "Data Mining Wizard" dialog box opens.
- On the welcome page click Next.
- Click "From existing relational database or data warehouse".
- Click Next.

- Under "What data mining technique do you want to use?" click Microsoft Sequence Clustering as shown in Figure 11. Then click Next.
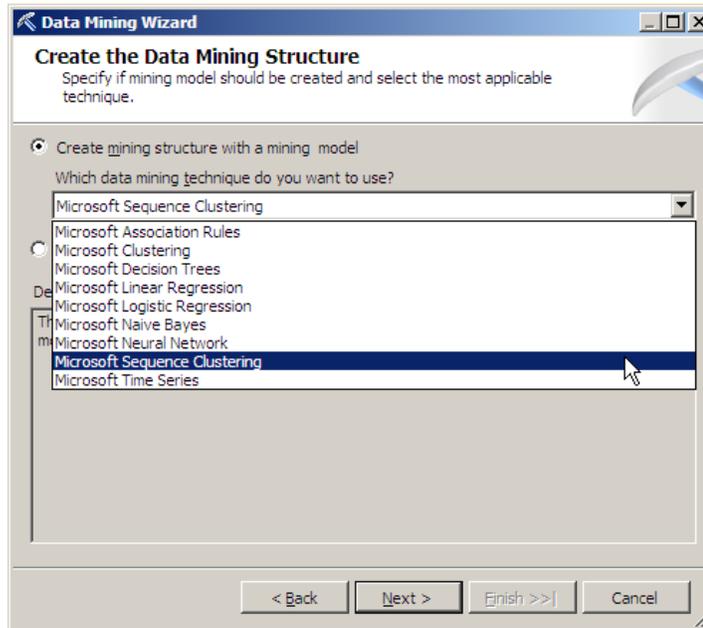


**Figure 11.** Choosing the data mining algorithm

Step 10: (see Figure 12)

- In this step we have to choose on which data source view we want to apply the data mining structure. So we choose the data source view created above the "Adventure Works DW data source view".
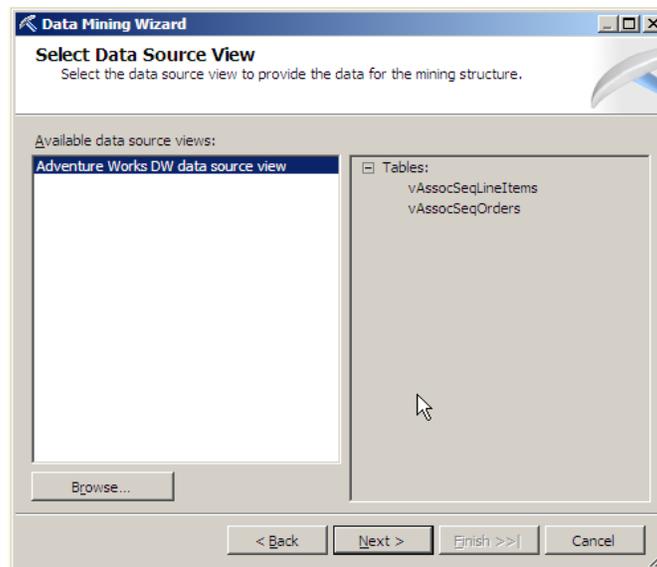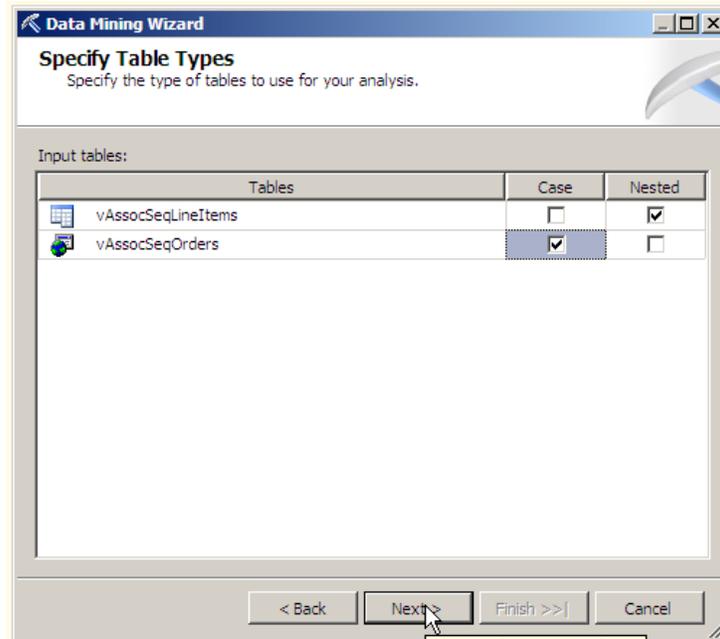- Then click Next.



|**Figure 12.** Selecting Data Source View

Step 11:
See Figure 13.

- In this step we have to specify the tables type. Which is the nested and which are the cases. In our example vAssocSeqOrders is a case table and vAssocSeqLineItems is the nested table.
- Then click Next



**Figure 13.** Specifying table types

Step 12:
In this step we can specify the columns used in our analysis.

- Clear the **key** check box next to **CustomerKey.**
- By default, OrderNumber and LineNumber are listed as Key types, which is correct.
- Select the **Input** and **Predictable** check boxes next to the **Model** columns. Make sure that the selection is the same as shown in Figure 14.
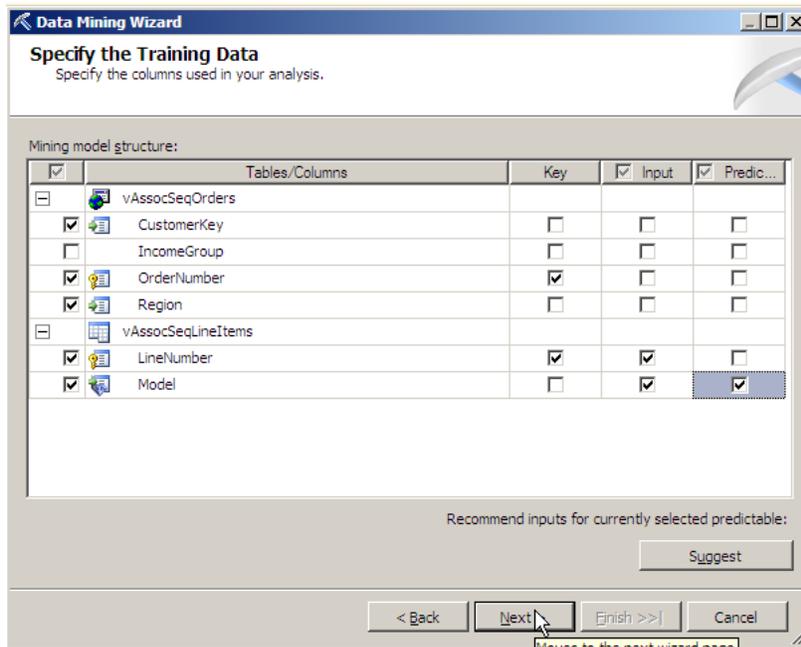- Then click Next.

**Figure 14.** Attribute specification for the sequence clustering mining structure.

Step 13:
See Figure 15.
- In this step we specify mining structure column's content and data type.
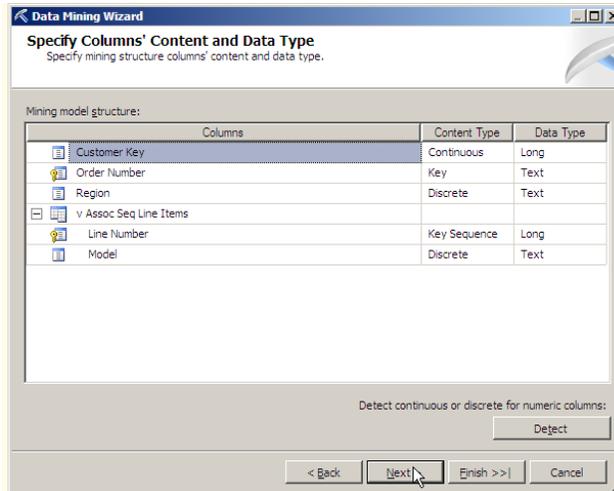- Then click Next.



**Figure 15.** Specifying column's content and data type

Step 14:
In this step we specify the number of cases to be reserved for model testing.
We determine the percentage of data for testing. By default it's 30%.
And also we determine the maximum number of cases in testing data set.
Then we click Next.

Step 15:
In this step we provide a name for the mining structure and the mining model.
We have named the both "Sequence Clustering"
Then we click Finish

Finally the new Sequence Clustering mining structure is displayed as shown in Figure 16.

Step 16:
- In the solution explorer, Right click on the project "Sequence_Clustering_Data_Mining" → Deploy
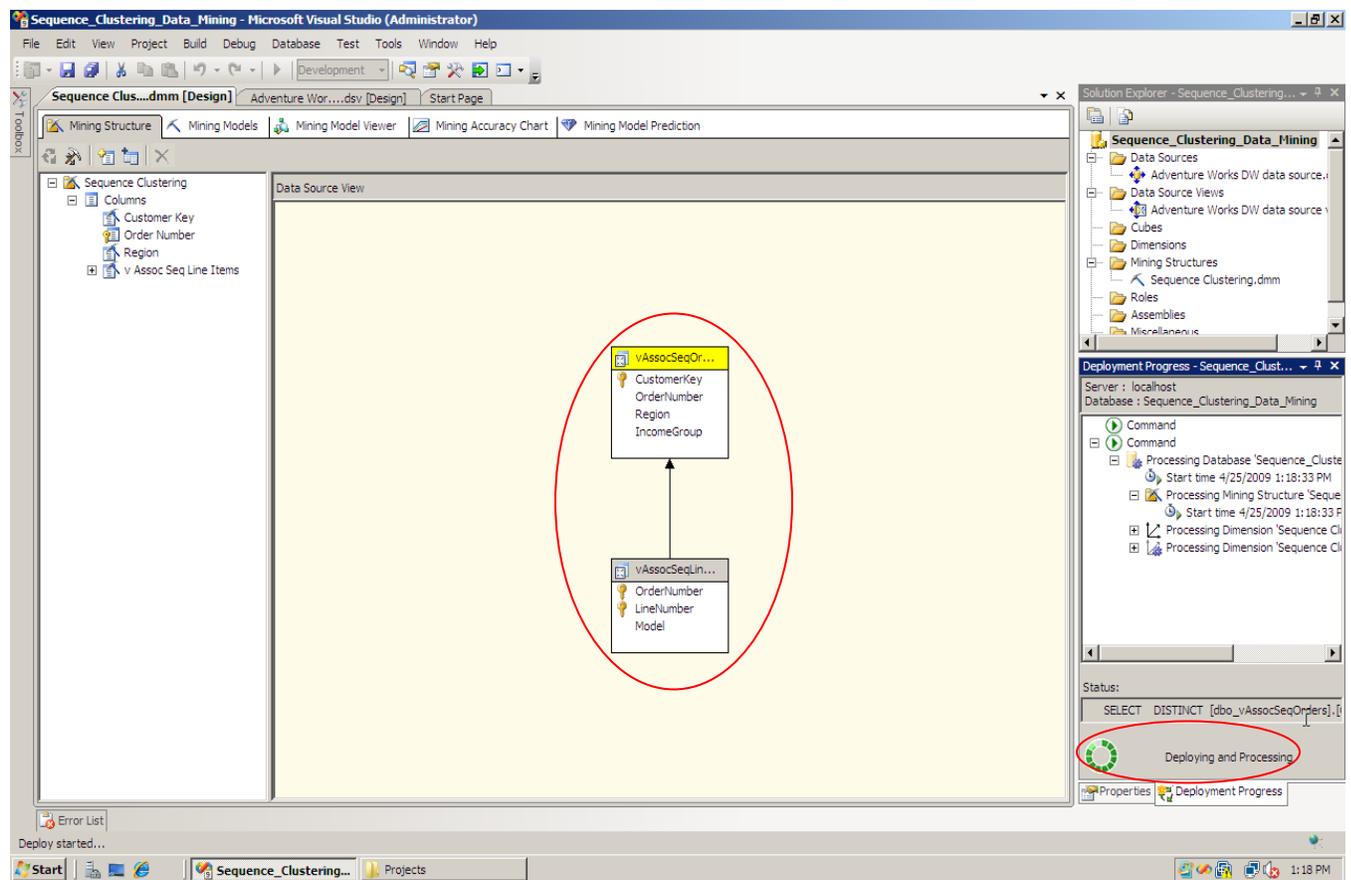


**Figure 16.** Sequence Clustering mining structure.

## IV. Exploring mining model

We use the Sequence Clustering viewer to explore the mining model we created. To open the Sequence Clustering viewer, click Mining Model Viewer.

The Sequence Clustering viewer contains five tabs: Cluster Diagram, Cluster Profiles, Cluster Characteristics, Cluster Discrimination and State Transitions.

### 1) Cluster Diagram

The Cluster Diagram tab displays the clusters discovered by the algorithm in the database. The layout represents the cluster relationships. Similar clusters are grouped close together. By default, the node color represents the density of all cases in the cluster (the darker one contains the highest number of cases). We can also change the meaning of node color-coding so that it represents an attribute and state. For example, to generate the diagram shown in Figure 17, in the Shading Variable list, click Model, and in the State list, choose for example "Bike Wash". We can see that Cluster 5 contains the highest density of Bike Wash.
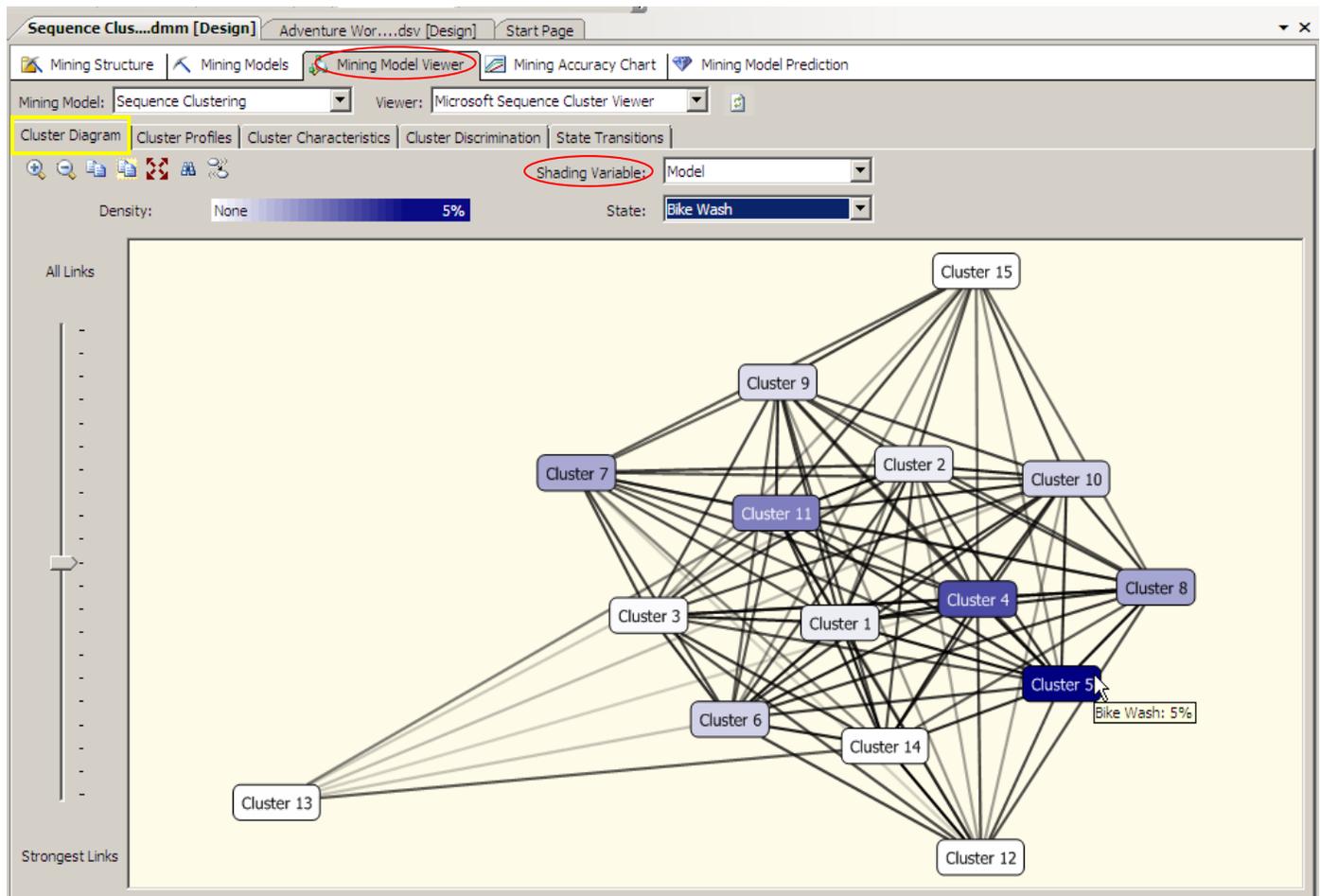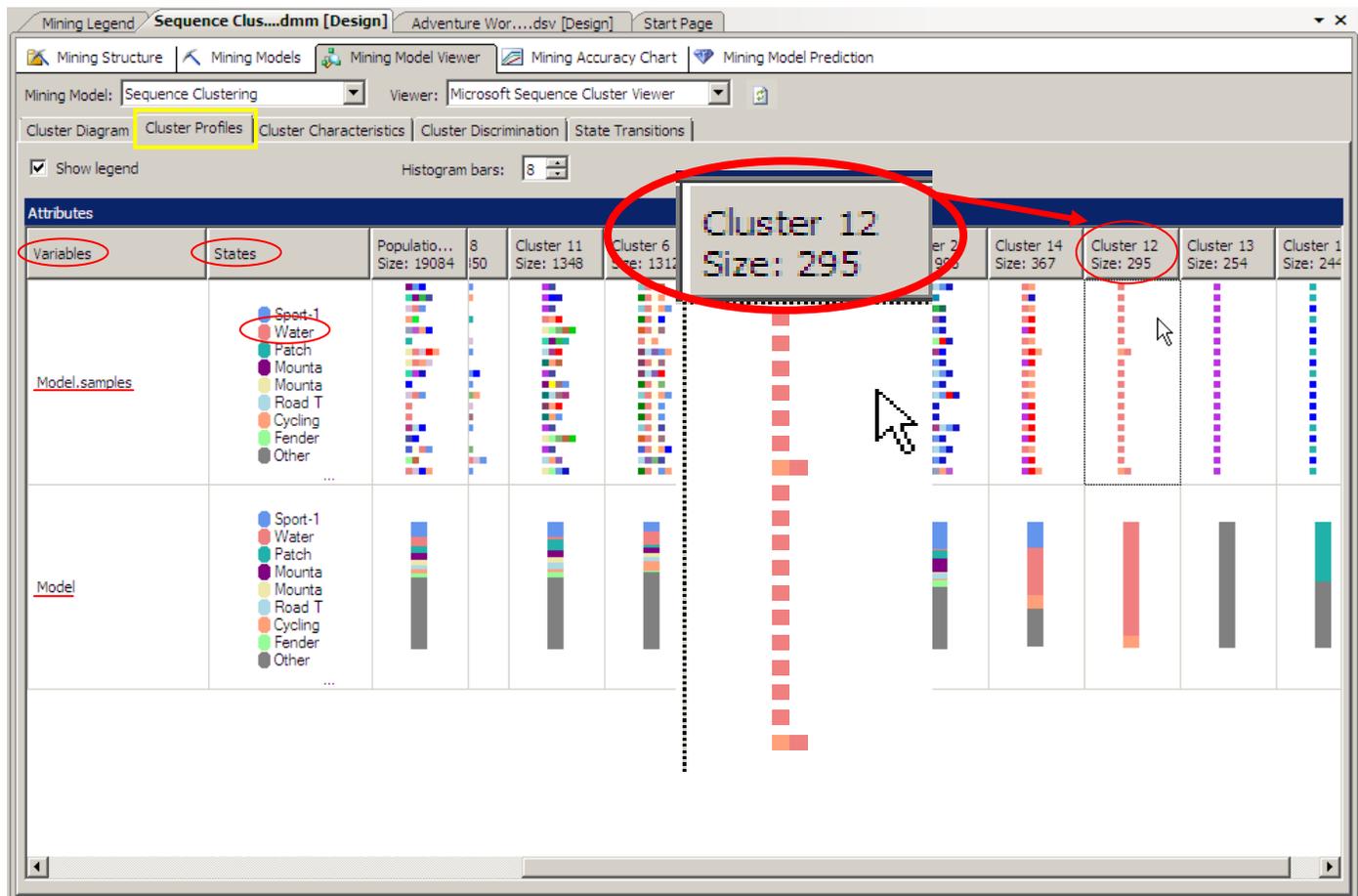


**Figure 17.** Cluster Diagram tab of the Microsoft Sequence Clustering model

## 2) Cluster Profiles

The Cluster Profiles tab displays the sequences that exist in each cluster. The rows listed in the Variables column show the variable distributions for a cluster. In Figure 18, the Model.samples row represents sequence data, and the Model row describes the overall distribution of items in a cluster. Each line of the color sequences displayed in each cell of the Model.samples row represents the behavior of a randomly selected user in the cluster. Each color in the sequence histogram represents a product model.



**Figure 18.** Cluster Profiles tab of the Microsoft Sequence Clustering model

For example, the pink color in cluster 12 represents the water bottle. The first color in most of the sequences is pink means that a customer is very likely to place water bottle in the shopping basket first.

### 3) Cluster Characteristics

The Cluster Characteristics tab summarizes the transitions between states in a cluster, with bars describing the importance of the attribute value for the selected cluster. For example, in Cluster 2, on of the most important profiles is that customers tend to place a Mountain tire tube in the shopping cart first.
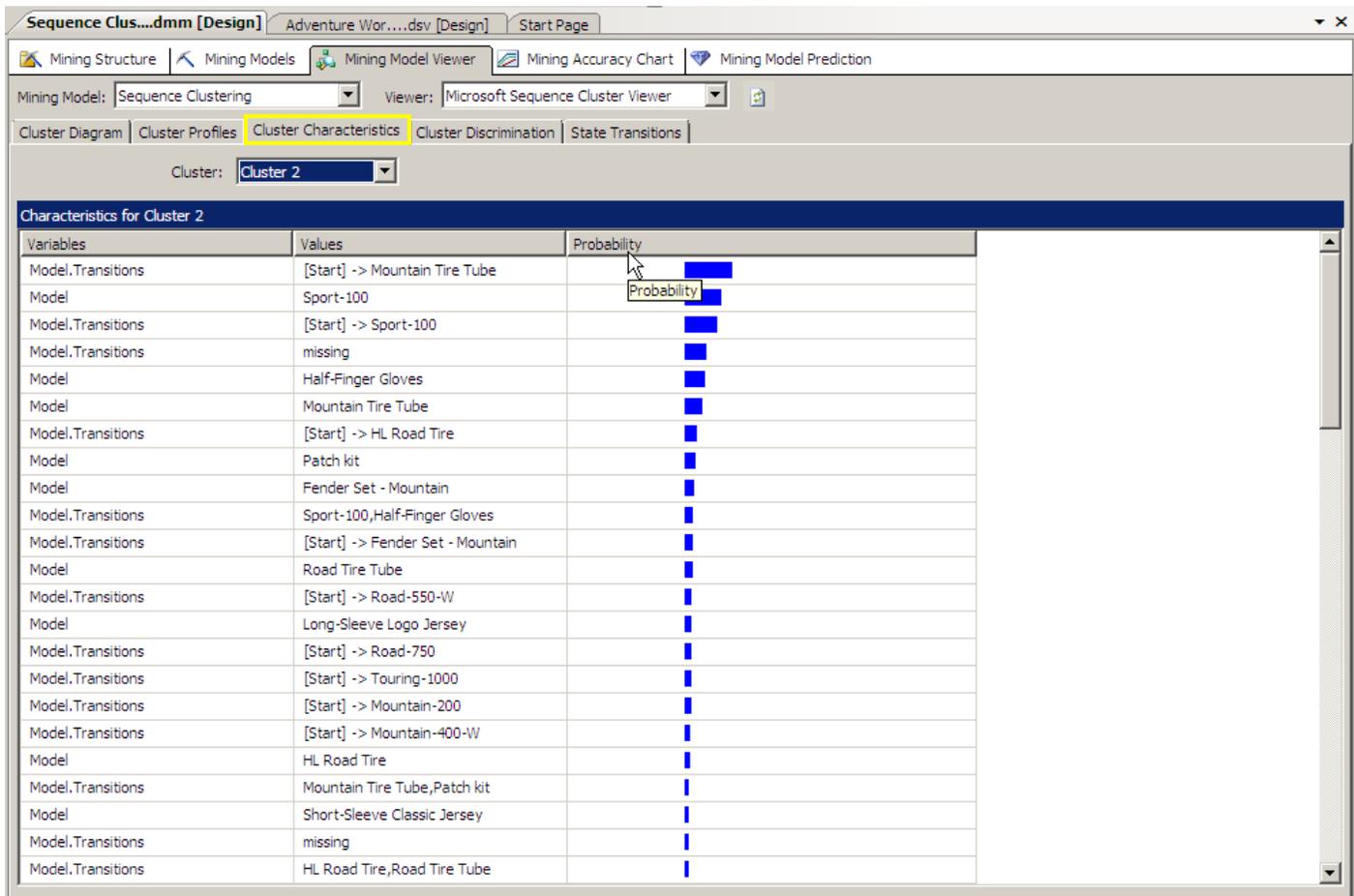


**Figure 19.** Cluster Characteristics tab of the Microsoft Sequence Clustering model

## 4) Cluster discrimination

In the Cluster Discrimination tab, we can compare two clusters, determining which models favor which clusters. The tab contains four columns: Variables, Values, Favors Cluster (i), Favors Cluster (i). If the cluster favors a specific model, a blue bar appears in one of the Favors Cluster(i) columns, in the row of the model listed in the Variables column. The longer the blue bar, the more the model favors the cluster.

For example, Figure 20 compares Cluster 1 with Cluster 2. A customer who purchases a **Touring Tire Tube** is more likely to be in Cluster 2, and a customer who purchases a **Classic Vest** is more likely to be grouped into Cluster 1.
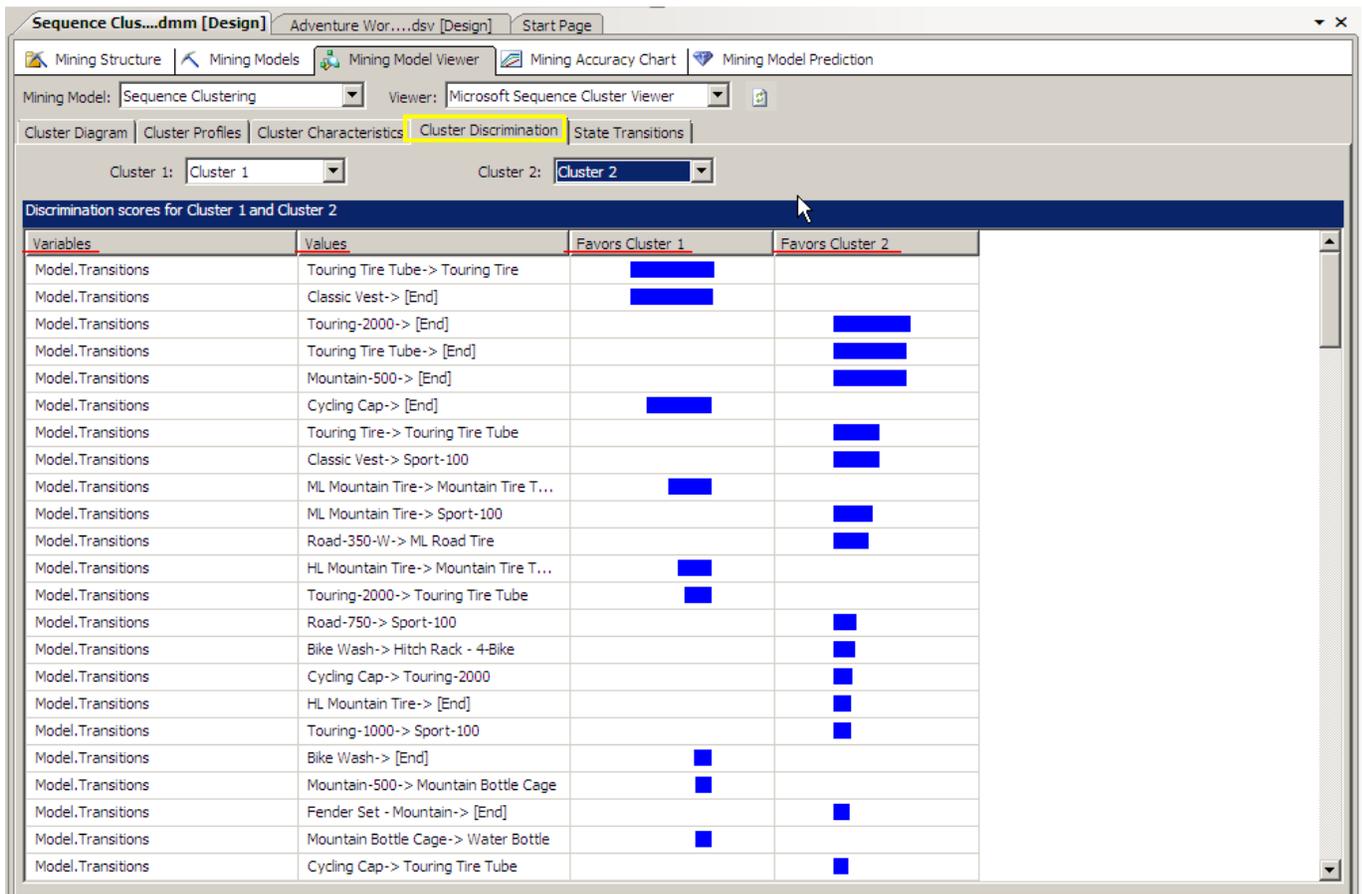


**Figure 20.** Cluster Discrimination tab of the Microsoft Sequence Clustering model

## 5) State transactions

On the State Transitions tab, we can select a cluster and browse through its state transitions. Each node represents a state of the model (such as Touring-2000). A line represents the transition between states, and each node is based on the probability of a transition. The background color represents the frequency of the node in the cluster.

For example, select **Cluster 2** from **Cluster.** As we can see in Figure 21, if users put a **Touring-2000** into his shopping cart, there is a probability of 0.60=60% (indicated by the blue arrow) that he will next put a **Touring Tire Tube** into the cart, and a probability of 1.00=100% that he will end his shopping by placing a **Touring Tire** into his shopping cart.
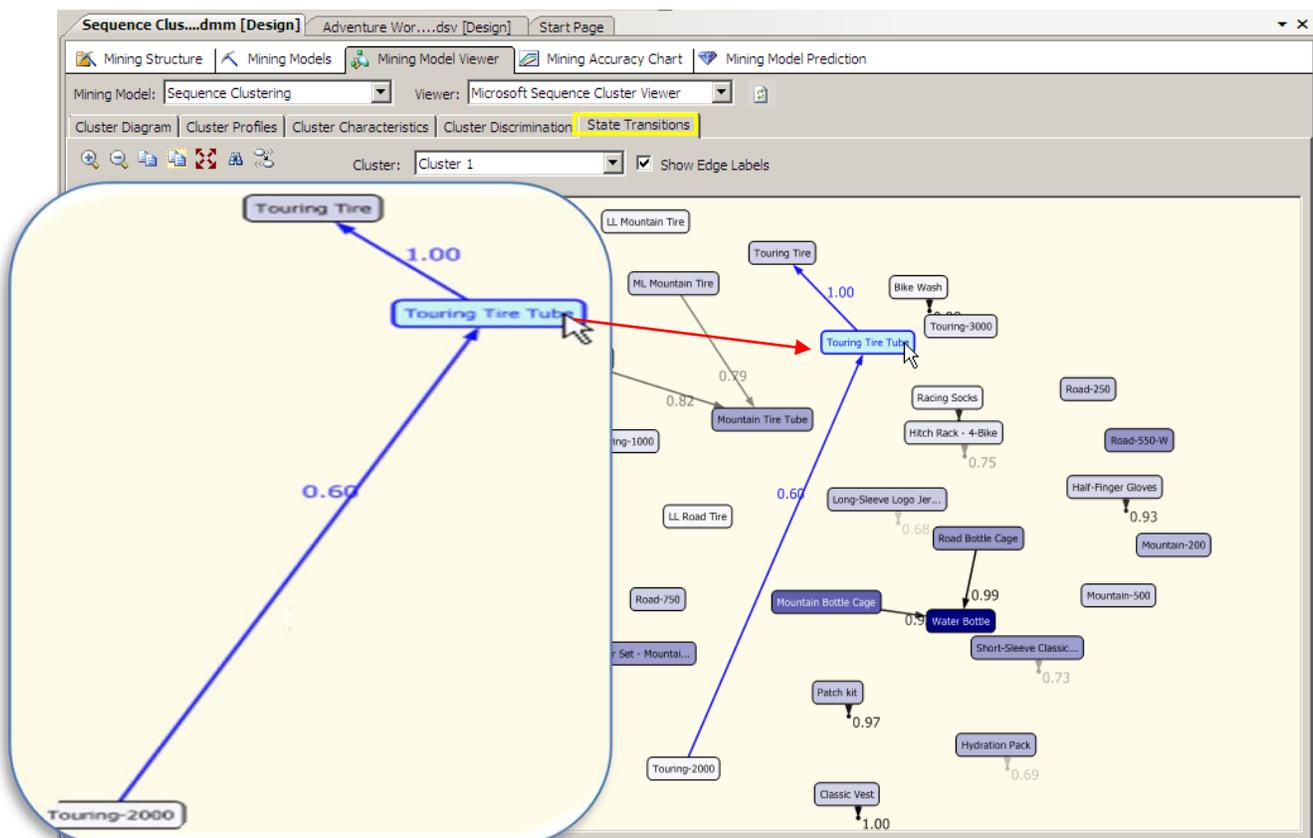


**Figure 21.** Cluster Transitions tab of the Microsoft Sequence Clustering model