# Hands-On Tutorial SAP Predictive Analytics, Automated Mode: Data Manager

Sharpen your data to produce stronger predictive models. This guide shows how to use the Data Manager in SAP Predictive Analytics to derive additional variables for your analysis, leading to a more detailed picture and better predictions.

The data used in this guide is publicly available so that the readers can follow hands-on and carry out the same analysis themselves.

June 2016

**Andreas Forster**
Predictive Presales Expert
**SAP Switzerland**
andreas.forster@sap.com

# TABLE OF CONTENTS

## INTRODUCTION

You may know, or guess, that the "Automated Analytics" in SAP Predictive Analytics is all about automating the process of creating predictive models.

Part of that process is to help the user create additional variables that help the model predict more accurately. This is what the Data Manager is primarily for. Without having to create new columns in the database, which might or might not be used in a model, the Data Manager creates a semantic view on top of the physical layer of a database. This is done in a graphical interface, which can produce hundreds or thousands of variables.

Think of creating such variables like a digital camera that zooms in on a situation. Creating these variables is very similar for the predictive model. They help to have a closer look and to better understand the situation. The model can take this information into account to get a clearer picture. Also similar to a digital camera, that picture is taken in "Automatic" mode. There is no need to fiddle with the detailed settings. Just keep producing models with little effort.

A very comprehensive framework underpins this creation, ensuring that powerful models are created without putting any constraints on the data. If you like to know more about the overall concept, you may want to read this article.[1]

This document explains how to create these detailed variables with the Data Manager, laying the foundations for strong predictive models. You can simply read this document or even implement your first models based on logic from Data Manager by following the steps hands-on. As example, we will predict which customers of a bank are likely to be interested in signup up for a credit card. Whilst this tutorial uses Data Manager in combination with a classification, Data Manager can also enhance other models, such as Regressions, Clustering or Time Series.

You may benefit most from this guide if you have already got some first experience with SAP Predictive Analytics. In case you have not used SAP Predictive Analytics before, I suggest to have a look at another tutorial first, in which you create predictions based on existing information/columns.[2] You can also get an overview with the official tutorials.[3]

To learn more about Data Manager please see the "Getting Started with Data Manipulation" help file.[4]

I would like to thank Abdel Dadouche, Gaëtan Saulnier and Orla Cullen, whose content of a hands-on session at a TechEd && d-code event inspired the creation of this guide.

---

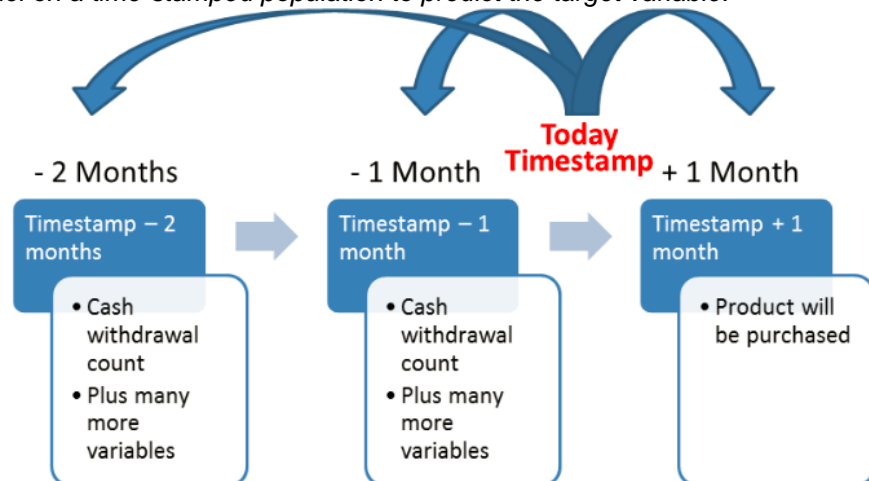### COMPONENTS

The Data Manager creates a semantic view called **"Analytical Data"**, or "Analytical Dataset". It is a graphical interface to join existing database tables and to create additional columns that can help the predictions. Note that the columns to not have to be created one-by-one. You will see that many detailed variables can be created at once.

The view has an understanding of time built in. It is able to relate to a certain point in time (called time-stamp). Information from before the time-stamp describes past behaviour and is used to train a predictive model. The view also contains a target variable describing a behaviour after the time-stamp, on which a model will be trained. The trained model can then be applied with a more recent timestamp (or even today's datetime) to predict future behaviour.

*Example for training a model on past data with a time-stamped population.*



*Example for applying a model on a time-stamped population to predict the target variable.*



The Analytical Dataset created by the Data Manager consists of the following three components. Only once all these components have been created, the dataset can be used for predictions.

| Component | Description |
|---|---|
| **Entity** | Whose behaviour to predict |
| **Time-stamped Population** | What behaviour to predict |
| **Analytical Record** | Historic information the predictions are based on |

**Entity**
Everything starts with the entity. The entity is the person or item whose behaviour we want to forecast. Based on the business requirement, it is immediately clear what the entity is going to be. In our example we want to forecast the behaviour of customers. So in that case the entity is a customer.

Note that Data Manager requires that you have a table with unique entries for all relevant customers. Duplicates are not allowed in this table.

**Time-stamped Population**
The time-stamped population enriches the entity with the business context of what behaviour is to be forecasted. This contains two pieces of information
- A filter on the entities to select only the ones relevant for the prediction
- A target variable, which describes during training what behaviour the entity showed. When applying the model, it describes the probability of showing the behaviour in future. For a classification, as shown in this tutorial, this indicator has to be binary, ie yes/no, 1/0 or any other combination of two values. To carry out a regression or time-series forecast the variable must be numeric. The target variable is the only variable that is based on information after the time-stamp.

**Analytical Record**
The analytical record describes the entity's behaviour (so the customer's in our case) before the time stamp. Such descriptions can be rather static (ie "Gender", "Natonality") or dynamic (ie "Age", "Income last month", Cash withdrawal count last month").

Typically you will spend most of the time trying to improve a model by providing additional variables in the Analytical Record. If the relevant data is spread over multiple columns, you can join these and dynamically derive new columns, ie through aggregation or pivoting.

As a quick recap: Columns from the Analytical Dataset describe the past. The target column from the time-stamped population is the only variable from after the time-stamp describing the future behaviour to predict.

**Examples**

Here are some examples for such Analytical Datasets:

**Purchasing Affinity**: Bank's customers' interest in credit card

| Component | Description |
|---|---|
| Entity | Customer |
| Time-stamped Population | As seen from the time stamp:<br>Filter: Existing customers who have not yet signed up for a credit card<br>Target Variable: Indicator whether the customer did or will sign up in the future time-window |
| Analytical Record | Past behaviour of the customer, as seen from the timestamp, ie<br>- Age<br>- Number of years the customer has been a client of the bank<br>- Average amount on current account previous quarter<br>- Average amount on current account the quarter before that<br>- Change in amount between the quarters as absolute values and in percent<br>- … |

**Customer Churn**: Risk of losing an insurance customer

| Component | Description |
|---|---|
| Entity | Customer |
| Time-stamped Population | As seen from the time stamp:<br>Filter: Existing insurance customer<br>Target Variable: Indicator whether the customer did or will cancel the existing contract in the future time-window |
| Analytical Record | Past behaviour of the customer, as seen from the timestamp, ie<br>- Age<br>- Number of years the customer has been a client of the insurer<br>- Number of claims last year<br>- Number of claims last year that were not reimbursed in full<br>- … |

**Preventive Maintenance**: Risk of a manufacturing machine/device showing unexpected behaviour

| Component | Description |
|---|---|
| Entity | Manufacturing machine |
| Time-stamped Population | As seen from the time stamp:<br>Filter: Machines that have been running without anomalies after the last maintenance interval<br>Target Variable: Indicator whether the machine did or will show an unusual pattern in the future time-window (ie unexpectedly high vibrations) |
| Analytical Record | Past behaviour of the machine, as seen from the timestamp, ie<br>- Duration since last maintenance<br>- Count of items produced since last maintenance<br>- Maximum temperature previous day<br>- Maximum temperature two days ago<br>- Change in temperature between the two days in degrees Celsius or Fahrenheit<br>- Change in temperature between the two days in percent<br>- … |

**MODELING STEPS**

The modeling of the Analytical Dataset is closely linked with creating / testing the predictive model. You have to create the entity, time-stamped population and a first Analytical record before testing the logic with a predictive model. If this step is successful, you enrich the Analytical Record to better describe the time-stamped population. You keep testing the prediction and enriching the Analytical Record until you are happy to deploy the model.



**Create Entity**
The entity (ie Customer) is defined with a single table with unique entries for each entity. The table can have multiple columns, you just need to specify an identifier column for the entity (ie an ID column). Automated Analytics will rename the ID to "KxId", clearly indicating the importance of this column.

**Configure time-stamped population**
Now link the entity table with further tables to filter the dataset according to the use case and to calculate the target variable. One of the tables must contain a date column so that the time-stamp concept can be applied. The time-stamp is used to
- Create the target variable to indicate the future behaviour. Note that the variable must return only two values, ie "Yes" and "No or 1 and 0 or any other value pair. So to predict whether a customer will sign up for a credit card, this variable will indicate whether the customer did sign up in the time-frame defined after the time-stamp.
- Apply a filter so that the dataset is restricted according to the use case, ie exclude any customers that had already signed up for a credit card before the time-stamp. As they have already purchased the product, their profile in the dataset would distort the model.

**Create Analytical Record**
Next, you need to create the Analytical Record to describe the entity's past behaviour. In the Analytical Record you link further tables to the entity. Columns from these tables are immediately added. Further columns can be calculated. If you need to use tables that have already been used in the time-stamped population, these tables also need to be added again to the Analytical Record.

At this point you just need to create a very basic Analytical Record. Once the Analytical Record is saved, you can test the model by creating a classification in the Automated Modeler.

**Test Analytical Record**
Open the Automated Modeler to test the Analytical Record. Have a look how the data is described. If the model does not give any warnings or error you know you are on a good track and you can further enhance the Analytical Record.

**Enrich Analytical Record & Test**
Back to the Analytical Record, keep adding additional columns to describe the entity and its behaviour in more detail. Keep testing the model every so often. Models can be saved under different versions, so it is easy to go back to an earlier model if needed.

Some examples of the types of variables that can be created are
- **Aggregates**: Your data might hold detailed transactional data about activities in your customers' accounts. Out of this history, Data Manager can create aggregates such as the count of the transactions or the average amount per transaction.
- **Pivots**: Pivoting builds on the above aggregation and makes it easy to graphically create a large number of more detailed aggregates. Sticking to the same example of activities in a customer account, pivoting can create individual counts of transaction types. So you can have transaction counts by cash withdrawals, standing orders, and so on.
- **Understanding of time**: Typically it is crucial to put the historic data into a context of time. Therefore the Data Manager has an in-built concept to relate the various measures to moments or ranges in time. Again, without the need for any coding, Data Manager can create detailed variables such as
  - Count of cash withdrawals in the previous quarter.
  - Count of cash withdrawals in the same quarter the year before.
  - Change in cash withdrawal counts in absolute values.
  - Change in cash withdrawal counts in percent.

The more variables you have available the better. Automated Analytics will find out which variables are needed for the model and which ones can be eliminated.

Aim to have the "Prediction Confidence" equal to or greater than 0.95 and ty to get the "Predictive Power" as high as possible (close to 1). In order to increase "Predictive Power" you can try adding additional variables. To increase "Predictive Confidence" you can try adding additional rows of data.

**Deploy**

Once you are happy with the model you put it into action. This could be scoring your current customer base and to write the probabilities in the database. It could also mean to add a new scoring column to an existing database view for real-time scoring or you could embed the model in various programming languages such as C++ or JavaScript into your application.

In productive environment you will use the "Model Manager" to monitor the model's accuracy over time and to retrain if needed. This step is not covered in this guide.

**Reuse Analytical Record**

Often a large part of the first Analytical Record you have created can be reused for further predictions. The input columns used for additional predictions, ie probability to upgrade from a single product to a bundle, will be very similar. Just make sure that the target variable is the only column using data from after the time-stamp.

**HANDS-ON IMPLEMENTATION**

**Background**

We want to help a bank optimize its marketing for credit cards and we will predict which customers are most likely interested in signing up for a credit card next quarter. We will use a dataset that was shared as part of the "3rd European Conference on Principles and Practice of Knowledge Discovery in Databases" held in Prague in the year 1999, or abbreviated PKDD '99.[5]

This dataset contains information about individual (anonymous) customers, some general data about the person and some activity on the accounts. For a more detailed description you can see the Appendix.

**Pre-Requisites**

You need the following elements to be able to follow the steps hands-on

- It is assumed you have a fair understanding of SAP HANA with some hands-on experience of loading data.
- It is also assumed you have had some first experience creating a classification with SAP Predictive Analytics, Automated Mode In case you are new to SAP Predictive Analytics, see the introduction of this guide for some initial tutorials to get you started.
- You need to have access to a SAP HANA system in which you can load the PKDD '99 data. See the chapter "Data Load" in the appendix. You are given the database content, which can be loaded easily into SAP HANA. Just note, that some of the data has been transformed from its original structure. For instance the client's birthdate and gender are stored in a single column in the original data. This has been separated in two columns.
  In case you don't have a SAP HANA system available, you might be able to use subscription-based system at relatively low cost, for instance on Amazon Web Services.[6]
- Most obviously, you need an installation of SAP Predictive Analytics, which includes the Data Modeler and Automated Mode. This guide has been written with SAP Predictive Analytics 2.1. Evaluation copies are currently available on the SAP Community Network.[7]
- Finally, you need to set up an ODBC Source from the computer with SAP Predictive Analytics to the SAP HANA system. In this document the ODBC source is called "My HANA" and the schema name is "I056450". In your environment those details will be different.

---

[5] PKDD '99, http://lisp.vse.cz/pkdd99/
[6] Amazon Web Services, http://aws.amazon.com/sap/
[7] SCN, http://scn.sap.com/community/predictive-analytics

**Step 1 – Create Entity**

Start by opening up SAP Predictive Analytics. The version used in this document is SAP Predictive Analytics 2.1.

Click on "Data Manager" → "Create or Edit Analytical Data". This Analytical Data will contain out Entity, the Time-stamped population and the Analytical Record.

Ensure the "Data Type" is set to "Data Base". Click "Browse" and select the ODBC source to the SAP HANA system. Remember to enter the SAP HANA user name and password in this window. Hit "OK".

Now specify where to save the metadata of the "Analytical Data" model. By default it is saved in the same database, that we use for the prediction. In that case SAP Predictive Analytics creates additional tables in the database which will contain the metadata. However, we save the model locally as files on our local computer.

Click "Metadata". Select "Store the metadata in a single place". Set "Data Type" to "Text Files" and select a folder.



Click "OK" → "Next".

Now we are ready to start creating our data model. The Entity, the Time-stamped Population and the Analytical Record are all created from this window. We will see this screen many times.



You have to start with the entity whose behaviour we want to predict. Click "New" → "Entity". Name it "ENT_CUSTOMER". Select the "CLIENT" table and set the "Id Field" to the column that uniquely identifies each entitiy. Here this is "CLIENT_ID".

No further changes are needed. Hit "Next". You will be asked for the Date type. Select "Date and Time" and "OK". Most databases support joins on DateTime-stamps. If your database does not support such joins, you have to select "Date Only".



And we are back in the main screen, where you can see the entity we have just created.

You can also see that a time-stamped population has been created. Before configuring this time-stamp population, look into the content of the entity. Select the entity "ENT_CUSTOMER" and click "View Data".



You see a preview of the data. Notice that the entity's ID got renamed from "CLIENT_ID" to "KxId". In the database the column is still called "CLIENT_ID",but within Data Manager this very important column is now referred to as "KxId".

Click "Close".

**Step 2 – Configure Time-stamped Population**
You will remember that the time-stamp population filters the dataset and that it provides the target variable. It does not matter in which sequence you create these elements, we start with the filtering.

In our case, we need to reduce the dataset to customers that have not taken out a credit card before the timestamp.

Select the Time-stamped population called "ENT_CUSTOMER\Population" and click "Edit".



Notice how the population is linked to the entity we created earlier. Change the object's name to "TSP_BuysCreditCardNextQuarter" and give it a meaningful description so we can easily find out content once we have created more objects.

Now we create the filter. This filter ensures that based on the date-time parameter that is used, the dataset is filtered dynamically to include only customers without credit cards at that point in time.

Click on "Edit Filters" and you see the current structure of the data. So far it is very empty, because we have not created much. However, see the option "Display only visible fields". Unticking this would display further fields from our entity. Keep it ticked.



The filter for the time-stamped population must be applied on an aggregate we still need to define. This aggregate will indicate whether the client does or does not have a credit card at the moment of the time-stamp. The information we are looking for is in the "CARD" table. Our entity is based on the "CLIENT" table. To get from the "CLIENT" table to the "CARD" table we have to link in the "DISPOSITION" table.

In order to join in the "DISPOSITION" table, click "Merge".



Here click "New" and select the "DISPOSITION" table as target.

Now click the plus-sign on the right hand side to define the join. On the left-hand side select the "KxId" column. On the right select the "CLIENT_ID".



Hit "OK" → "OK" → "Close" and you should see the columns of the "DISPOSITION" table in the time-stamped population editor.



Now we can create the aggregate that indicates whether a customer has a credit card at the given time stamp. The aggregate we create will include the link to the "CARD" table.

Click "New" → "New Aggregate". Use the following settings as shown in the screenshot below:
- Link to the "CARD" table, which contains a date column named "ISSUED". This date column will be used in relation with the time-stamp parameters.
- Link from the "DISP_ID" column of the "DISPOSITION" table to the "DISP_ID" of the "CARD" table.
- As aggregate use "NotExists" on the "CARD_ID" column.

If a client has not purchased a credit card, he will not have an entry in the "CARD" table and the aggregate will return TRUE.



Click on the "Period Settings" tab, tick "Define Periods" and select "Successive Periods". Here we define the time-range the aggregate relates to. Use the settings from the following screenshots.

We create one aggregate for a period of 100 years before the time-stamp. Be careful to get these settings correct. We use a period of 100 years starting 100 years before the time-stamp. The values can be clicked with the mouse. See below on how to specify the "TimeStampPrompt".





Please double-check your period settings are exactly as shown in the two screenshots above. Then hit "OK" and enter the new aggregate's name "HasNoCreditCardAtTimeStamp".

Click "OK" and you see the new column. Should you not see the column, then untick "Display only visible fields" and it should appear. Now you can change the column's visibility with a tick in the corresponding column.



Now we can create the actual filter. So far we have created a column that indicates whether we want to include the entity (meaning person in this case) in our dataset. So we apply a filter on the column as we are interested only in clients that have no credit cards at the time of the time-stamp.

Click the "Filters" icon, then "New Condition". You just need to double-click on the aggregate we have just created. You will find it on the right-hand side under "Variables".



Hit "OK" → "Close" and the filter is active.

Now we need to define the target-variable. This is another aggregate. You should still be in the "Time-stamped Population Editor".



Click "New" → "New Aggregate" and use the following settings. Notice that we are now using "Exists" as aggregate. Now we look for people who have signed up for a credit card after the time stamp.

As we want to identify those customers that bought the credit card in the quarter after the time stamp, set the following "Period Settings".



Click "OK" and name the aggregate "BuysCreditCardInQuarterAfterTimeStamp".

Click "OK" and you should see the following:

Click "Next", go to the "Target" tab and set the target to our new aggregate "BuysCreditCardInQuarterAftertTimeStamp".



Click "Next" and we are back in the Data Manager.



Since we saved the time-stamped population under a new name, there are two different time-stamped populations. The default version cannot be deleted, so just ignore it.

Before proceeding we quickly test our time-stamped population. Select "TSP_BuysCreditCardNextQuarter" and hit "View Data". The time prompt comes up, enter the 5th of August 2014 and hit "OK".

You see a first preview of the data.

We want to have a look at the target variable. How many positive cases do we have in the historic quarter after the timestamp?

Select the "Statistics" tab. Click "Compute statistics over the whole dataset". This analyses the whole dataset.

Click the "Category Frequency" tab and select our target variable from the drop-down. Click the pie-chart icon and you should see the following distribution.



The dataset contains 4700 rows. Only 2.91% of these are positive cases, so 137 customer purchased their first credit card in the following quarter. Ideally we would like to have 1000 positive cases to train on. However, our data set if fairly small with under 5000 customers and we continue trying to create a strong and robust model. Should you encounter a similar situation with your own data, you could try to increase the size of the dataset or the length of the time frame the target variable is based on.

Click "Close".

**Step 3 – Add Analytical Record**

Now we add the analytical record, which will describe our entity, the customer. Typically this is where you will spend most of the time, creating variables that describe your entity.

In "Data Manager" click "New" → "Analytical Record". Give it the name "ANA_CUSTOMER". Notice that the analytical record is already linked to our entity "ENT_CUSTOMER".

To add columns to the analytical record, hit "Edit Attributes". So far the analytical record contains only the entity ID and the time-stamp. The columns from the time-stamped population are not listed.

As the entity does not yet have the additional columns from the "CLIENT" table, on which the entity is based, we start by adding these. Click "Merge" → "New". Select the "CLIENT" table as done earlier.

Click the plus-sign and join from the entity's "KxId" column to the "CLIENT_ID" column. In effect, this is a self-join, as both columns are the same.



Hit "OK" → "OK" → "Close". You now see the additional columns. They are shown in the same color to indicate they come from the same table. Should you not like the color, you can change it in the "Edition" tab on top.

However, we should not use the date of birth. Instead we will later calculate the client's age at the given time-stamp. Make sure the option "Display only visible fields" is unticked and deselect the option "Visibility" for the "DATE_BIRTH" column.



Now we need to instruct Data Manager to sort the data consistently. This is important to obtain reproducible models, which we will try to improve in multiple iterations. This is possible by instructing Data Manager to sort the data on a specific column. If not instructed, the database can return the data in different sequence every time we query it. This would change how data is split between the estimation and validation subsets, which as a result can change the model that is produced.

To obtain consistent results, enter the value 1 in the Order for the "Client_ID" column. Make sure the column's "Visibility" is ticked.

Now we are ready to create a first predictive model. At this stage we are just testing that everything works well, we are not yet expecting a strong result. We will improve the model later by adding more content to the analytical record.

Click "Next" twice and you are back on the main screen of the Data Manager.

**Step 4 – Create First Prediction to Test Data Model**

Click "Cancel" to get to the main screen of SAP Predictive Analytics. This document assumes that you have created some classification model beforehand. Therefore these steps will be described a little shorter.

On the left hand side, select "Modeler" in the "Automated Analytics" section and click "Create Classification/Regression" in the center area.



Now select "Use Data Manager". Choose the ODBC source, select the analytical record we have created followed by the time-stamped population.

Eventually we want to test the model on previously unseen data. It is best to get in the habit of changing the cutting strategy, so that it sets some data aside to test the model's performance at the end on new data.

Change the cutting strategy to "Random with test at the end".



"OK" this and hit "Next". You will be prompted for the time-stamp. Select the same date as before, 5th of August 2014. Click "OK".

In the "Data Description" screen hit "Analyze" to retrieve the data structure.



Click "Next" and chose "GENDER" and "DISTRICT_ID" as explanatory variables. Exclude "KxID", "KxTimeStamp" and "CLIENT_ID" as they are not adding any value. The target variable should already be set to "BuysCreditCardInQuarterAfterTimeStamp".

Hit "Next" followed by "Generate" to create the predictive model. You should get a summary screen, confirming that a very first model has been created.



The test has been successful in that a very basic model was created without any errors. Note the Predictive Power of 0.2706 and the Predictive Confidence of 0.8865.

If you would like to see the model's ability graphically, click on "Next" → "Model Graphs". To see the model's performance on all data sets, click the "Data Sets" icon and tick "All Data Sets". Now you also see the model performing on data that has not been used at all in producing the model.



You can now "Cancel" this and go to the "Home Panel" to continue enhancing the analytical record.

**Step 5 – Add static columns to the Analytical Record**
Now we continue adding columns to the analytical record to better describe our entity, the customer, and its behaviour. At any time, feel free to test your data model with a classification model.

Back in the "Data Manager", edit the analytical record "ANA_CUSTOMER".



Click "Edit Attributes".

Now we add the data from the "DISTRICT" table. Select "Merge" "→ "New". Join in the "DISTRICT" table using the "DISTRICT_ID" columns as keys.



Click "OK" → "OK" → "Close" and the additional columns are listed.



Feel free to test the data model with a new Classification model if you like. When done, continue on the next page.

**Step 6 – Add calculated column to the Analytical Record**
Earlier on we excluded the client's birthdate. Now we use this information to calculate the person's age at the moment of the time-stamp.

Click "New" → "Expression Editor". Construct the following formula:
dateDiffNbDays(KxTimeStamp, DATE_BIRTH) / 365.25

Click "OK" and name the column "Age". Hit "OK".

Our view of the customer is now more detailed.



Click "Next" twice to end editing the analytical record. When asked whether to save the new logic as a new version, you could select "Yes" of course. For simplicity, in this document we overwrite the existing version. So select "No".



You are now back in the "Data Manager".

Let's create a new predictive model and see if it has improved. Select "Cancel" and create a classification as before.

The source:



The time-stamp:

The selected variables:



Create the model and we can see that it has already improved by a big margin.

| | |
|---|---|
| Predictive Power (KI) | 0.5303 |
| Prediction Confidence (KR) | 0.9357 |

**Step 7 – Add time-related behavioural data**

So far we have described the customer itself, ie the gender, age or where the person lives. Now we focus on the person's behavior. We will look at the sum of the transaction amounts on the account, broken down by the four quarters before the time-stamp.

Get back into the "Data Manager" and edit the analytical record.



Click "Edit Attributes".

To include the behavior in the model we have to go through the "DISPOSITON" and "ACCOUNT" table to get to the "TRANSACTIONS" table. Therefore merge in first the "DISPOSITION" table followed by the "ACCOUNT" table.

The "DISPOSITION" table:

| Source Field | | | | | Target Field | | |
|---|---|---|---|---|---|---|---|
| Alias | Type | Source | Order | Key | Allowed Type: Integer | | |
| Age | continuous | Function | 0 | 0 | 🔑 ACCOUNT_ID | | |
| 🔑 KxId | nominal | ENT_CUSTOM... | 0 | 1 | 🔑 CLIENT_ID | | |
| 🔑 KxTimeSta... | nominal | ENT_CUSTOM... | 0 | 1 | 🔑 DISP_ID | | |
| 🔑 CLIENT_ID | continuous | I056450CLIENT | 1 | 1 | | | |
| 🔑 DISTRICT_ID | continuous | I056450CLIENT | 0 | 2 | | | |
| GENDER | nominal | I056450CLIENT | 0 | 0 | | | |
| DATE_BIRTH | continuous | I056450CLIENT | 0 | 0 | | | |
| DISTRICT_ID_1 | continuous | I056450DISTRI... | 0 | 0 | | | |
| A2 | nominal | I056450DISTRI... | 0 | 0 | | | |
| A3 | nominal | I056450DISTRI... | 0 | 0 | | | |
| A4 | continuous | I056450DISTRI... | 0 | 0 | | | |
| A5 | continuous | I056450DISTRI... | 0 | 0 | | | |
| A6 | continuous | I056450DISTRI... | 0 | 0 | | | |
| A7 | continuous | I056450DISTRI... | 0 | 0 | | | |
| A8 | ordinal | I056450DISTRI... | 0 | 0 | | | |
| A9 | ordinal | I056450DISTRI... | 0 | 0 | | | |
| A10 | continuous | I056450DISTRI... | 0 | 0 | | | |
| A11 | continuous | I056450DISTRI... | 0 | 0 | | | |
| A12 | continuous | I056450DISTRI... | 0 | 0 | | | |
| A13 | continuous | I056450DISTRI... | 0 | 0 | | | |
| A14 | continuous | I056450DISTRI... | 0 | 0 | | | |
| A15 | continuous | I056450DISTRI... | 0 | 0 | | | |
| A16 | continuous | I056450DISTRI... | 0 | 0 | | | |

OK    Cancel

The "ACCOUNT" table:

This leads to the following collection of merges:



Close this window and the analytical record contains the new columns. You may have to scroll down to see them.

With these tables available, we can now use the "TRANSACTION" table. Select "New" → "New Aggregate".
Link to the "TRANSACTION" table as follows, calculating the sum of the transaction amounts:



Continue in the "Period Settings". Have the sum calculated individually for the 4 quarters before the time-stamp.

These quarterly aggregates could already help produce a better model. However, let's make it even more detailed. Go into "Filters and Pivots Settings".



Here we break down the quarterly value in more detailed aggregates. The "TRANSACTION" table specifies the transaction's operation, ie a cash withdrawal or credit card withdrawal. We will create individual variables by quarter by operation for the last 4 quarters.

At the bottom left in the "Pivot" section, set "Variable" to "OPERATION". Click the binoculars next to it and chose "from the whole dataset". This will load all possible values from the "OPERATION" column. Just be patient in case this takes a few seconds. With larger datasets you will want to pick one of the other options to prevent that all data gets downloaded on your client. Each of these values will become a quarterly aggregation variable, for each of the 4 specified quarters.

Once complete, the "Variable" drop down might change. Set it back to "OPERATION" and you should see the following:

Select "Also create aggregates without pivot". This will create another quarterly variable with the total sum. So in addition to the 5 values in "OPERATION" a 6th column with the total value is calculated. 6 columns by 4 quarters makes 24 new fields, which is also indicated at the bottom left of the screen.



Click "OK" and name the variable "call it QuarterlyTransaction". Hit "OK" and you will see the new variables. Each variable has the prefix that we had chosen followed by a description of the variable's logic.

Click "Next" once. If you are interested to have a look at the new colums, click the "View Data" icon.



Confirm the date prompt and see the first 100 rows.



Close this window with "OK". Hit "Next" and save the model when prompted.

Now create a new classification model. You should only have to manually exclude the "DATE" column.



The model has been greatly improved further! The Predictive Power incrased from 0.0.53 to 0.64 and the Prediction Confidence is now above the desired 0.95, indicating a robust model.

Now we want to apply the trained model on new data. Which customer is likely to buy a credit card in the next quarter? Click "Next" to continue.



A short look into the ""Contribution Variables" shows the selected variables and their importance.

Click "Previous", then go into "Run" section.



Click "Apply Model" and set the following options. "Data" is the time-stamped population "TSP_BuysCreditCardNextQuarter". Generate "Probability & Error Bars". Write the results in a table called "MyScores".

Before proceeding, go into "Advanced Apply Settings…" and copy the "KxId" column into the output so you can identify the scores with the related clients.



Click "OK" and "Apply". Now enter a time-stamp that is 3 months after the previously used time-stamps. So pick 5[th] November 2014. The model learned from the past and now tries to predict a future it has not yet seen.

Click "OK". You should get a confirmation that the data manipulations are materialized temporarily in a SAP HANA table.



Close this window and you should see a confirmation, that the model was applied successfully.

Click "View Output" and you see the first rows of the results.



The colum "KXID" identifies the client. The column starting with "PROBA_RR" holds the estimated probability that the client will sign up for a credit card in the three months after 5$^{th}$ November 2011. You can use this information, to optimize your Marketing campaign.

Since the scores are written into a SAP HANA table, you can also see the results directly in the SAP HANA Studio.



Well done! You have used the Data Manager to improve a classification model created with "Automated Analytics". The additional insight from the newly created variables based around the understanding of time has provided a much clearer picture, resulting in a stronger and more robust model.

**HINTS AND TIPPS**

Some basic points, which you might find helpful when working with Data Manager

- Data Manager needs to connect to a database with SQL. It cannot be used directly on flat files.
- It is preferred to use capital letters for all table names and column names, ie use CLIENT_ID instead of ClientID.
- The entity table must not contain duplicates. So you cannot base the entity on a transaction table, which will have duplicates.
- When working with date columns, typically the date itself should not be included in the model. Instead derive further columns from it, ie a duration from the time-stamp to the date variable. This is the case when calculating the age or the duration of how long an account has been active for instance.
- Only the target variable must use information from after the timestamp
- It is best practice to have 1000 positive cases to train a model on. Less cases can be ok, just try to get the Predictive Confidence (KR) above 0.95. This value can often be improved through additional variables, which impact the model positively.

**APPENDIX**

**Data Description**
This chapter describes the data we are working with. Most information in this chapter is taken from the PKDD '99 website[8].

The following tables are available:

| Table | Numer of Rows | |
|---|---:|---|
| **ACCOUNT** | 4500 | static characteristics of an account |
| **CLIENT** | 5369 | characteristics of a client |
| **DISPOSITION** | 5369 | relates together a client with an account i.e. this relation describes the rights of clients to operate accounts |
| **PERMANENT_ORDER** | 6471 | describes characteristics of a payment order |
| **TRANSACTION** | 1056320 | describes one transaction on an account |
| **LOAN** | 682 | describes a loan granted for a given account |
| **CARD** | 892 | describes a credit card issued to an account |
| **DISTRICT** | 77 | describes demographic characteristics of a district |

The database model is as follows:



---

The tables contain the following columns.

ACCOUNT

| Column | Description | Content |
|---|---|---|
| ACCOUNT_ID | Account identifier | |
| DISTRICT_ID | Branch's district identifier | |
| FREQUENCY | frequency of issuance of statements | 'AFTER TRANSACTION', 'MONTHLY', 'WEEKLY' |
| DATE | Date of account opening | |

CLIENT

| Column | Description | Content |
|---|---|---|
| CLIENT_ID | Client identifier | |
| DISTRICT_ID | Branch's district identifier | |
| GENDER | Gender | 'FEMALE' , 'MALE' |
| DATE_BIRTH | Date of birth | |

DISPOSITION

| Column | Description | Content |
|---|---|---|
| DISP_ID | Disposition identifier | |
| CLIENT_ID | Client identifier | |
| ACCOUNT_ID | Account identifier | |
| TYPE | Type of disposition | 'DISPONENT', 'OWNER' |

PERMANENT_ORDER

| Column | Description | Content |
|--------|-------------|---------|
| ORDER_ID | Order identifier | |
| ACCOUNT_ID | Account identifier | |
| BANK_TO | Bank of the recipient | 'AB', 'CD', 'EF', 'GH', 'IJ', 'KL', 'MN', 'OP', 'QR', 'ST', 'UV', 'WX', 'YZ' |
| ACCOUNT_TO | Account of the recipient | |
| AMOUNT | Debited amount | |
| K_SYMBOL | Characterization of the payment | 'HOUSEHOLD PAYMENT', 'INSURANCE PAYMENT', 'LEASING', 'LOAN PAYMENT', <Null> |

TRANSACTION

| Column | Description | Content |
|--------|-------------|---------|
| TRANS_ID | Transaction identifier | |
| ACCOUNT_ID | Account identifier | |
| TYPE | Transaction type | 'CREDIT', 'VYBER', 'WITHDRAWAL' |
| OPERATION | Mode of transaction | 'COLLECTION FROM ANOTHER BANK', 'CREDIT CARD WITHDRAWAL', 'CREDIT IN CASH', 'REMITTANCE TO ANOTHER BANK', 'WITHDRAWAL IN CASH', <Null> |
| AMOUNT | Amount of money | |
| BALANCE | Balance after transaction | |
| K_SYMBOL | Characterization of the transaction | 'INTEREST CREDITED', 'PAYMENT FOR STATEMENT', 'HOUSEHOLD', 'OLD-AGE PENSION', 'INSURANCE PAYMENT', 'LOAN PAYMENT', 'SANCTION INTEREST IF NEGATIVE BALANCE', <Null> |
| BANK | Bank of the partner | 'AB', 'CD', 'EF', 'GH', 'IJ', 'KL', 'MN', 'OP', 'QR', 'ST', 'UV', 'WX', 'YZ', <Null> |
| ACCOUNT | Account of the partner | |
| DATE | Date of the transaction | |

LOAN

| Column | Description | Content |
|---|---|---|
| LOAN_ID | Loan identifier | |
| ACCOUNT_ID | Account identifier | |
| AMOUNT | Amount of money | |
| DURATION | Duration of the loan in months | |
| PAYMENTS | Monthly payments | |
| STATUS | Status of paying off the loan | 'A': contract finished, no problems<br>'B': contract finished, loan not paid<br>'C': running contract, OK so far<br>'D': running contract, client in debt |
| DATE | Date when the loan was granted | |

CARD

| Column | Description | Content |
|---|---|---|
| CARD_ID | Credit card identifier | |
| DISP_ID | Disposition identifier | |
| TYPE | Type of card | 'classic', 'gold', 'junior' |
| ISSUED | Date when the card was issued | |

DISTRICT

| Column | Description | Content |
|---|---|---|
| DISTRICT_ID | Branch's district identifier | |
| A2 | District name | 'Benesov', 'Beroun', 'Blansko', … |
| A3 | District region | 'Prague', 'central Bohemia', 'east Bohemia', … |
| A4 | Number of inhabitants | |
| A5 | Number of municipalities with inhabitants < 499 | |
| A6 | Number of municipalities with inhabitants 500-1999 | |
| A7 | Number of municipalities with inhabitants 2000-9999 | |
| A8 | Number of municipalities with inhabitants >10000 | |
| A9 | Number of cities | |
| A10 | Ratio of urban inhabitants | |
| A11 | Average salary | |
| A12 | Unemployment rate previous year | |
| A13 | Unemployment rate most recent year | |
| A14 | Number of entrepreneurs per 1000 inhabitants | |
| A15 | Number of committed crimes previous year | |
| A16 | Number of committed crimes most recent year | |

DISTRICT

**Data Load**
You have two options to load the data into SAP HANA. The easiest would be to import the whole content with a few clicks and you have the content exactly as described in this document as option 1 (Import SAP HANA Content). This upload includes a number of transformations made on the original data, for instance a translation to English values and an update of the dates to recent values.

Alternatively, you can also manually import and transform the data as described in option 2 (Manually Import and Transform).

***Option 1: Import SAP HANA Content***
You can download the SAP HANA tables from GitHub[9]. Extract the zip file.

Now import the content into SAP HANA. In HANA Studio, select "File" → "Import" → "SAP HANA" → "Catalog Objects".



Hit "Next".

_____

[9] Data for Tutorial: https://github.com/AndreasForster/Predictive/raw/master/DataManagerTutorialData.zip

Select your HANA system and click "Next".

On the next screen select "Import Catalog Objects from current client". Select the folder you have extracted DataManagerTutorialData.zip into.

Click "Next".

Select all tables



Click "Next" and "Finish". The tables should now be loaded.

### *Option 2: Manually Import and Transform*
*Download original data and Upload into SAP HANA*
In order to download and transform the original data, follow these steps. However, I strongly recommend to simply upload the prepared package into SAP HANA whenever possible (see above for Option 1).

**Step 1**:
Download the raw data from the "Past ECML/PKDD Discovery Challenges" 10. From the PKD '99 challenge select "Data" → "Original data download" → "Download the Financial Data"

**Step 2:**
Extract the data:berka.zip file and rename all file extensions from .asc to .csv.

**Step 3:**
One row in the district.csv file contains a '?' in two columns to indicate missing data. Remove these question-marks so that the column is automatically identified as numerical.

**Step 4:**
Load each file as new table into SAP HANA. This can be done for instance with the SAP HANA Studio. Carry out the following steps for each table:

In the SAP HANA Studio, select "File" → "Import" → "SAP HANA Content" → "Data from Local File".



---

10 Past Challenges, http://lisp.vse.cz/challenge/PAST/index2.htm

Click "Next" → Select the Target System  (most likely you will only see one) → Now specify the table to load.
- Select the .csv file
- Set "Field Delimiter" to "Semi Colon"
- Tick "Header row exists"
- Tick "Import all data"
- Select the database Schema from the dropdown
- As "Table Name" enter the file name (without file extension) in upper case, ie CLIENT

Example:



Click "Next" and change the column names from lower case to upper case. If indicated in the table below, you may have to set additional settings, ie a Key or different data types.

For the "CLIENT" table just set the "CLIENT_ID" and "DISTRICT_ID" as Keys.



Click "Finish" and you should see a success status in the job log.

The table has been created and the data has been uploaded. You can see the table in the Systems browser. You can now view the data, for instance by right-clicking on the table to open the data preview.



These are the specific table settings you have to apply, in addition to changing the table and column names from lower case to upper case.

| File Name | Table Name | Key Columns | Additional |
|---|---|---|---|
| **account.csv** | ACCOUNT | ACCOUNT_ID, DISTRICT_ID | |
| **card.csv** | CARD | CARD_ID, DISP_ID | |
| **client.csv** | CLIENT | CLIENT_ID, DISTRICT_ID | |
| **disp.csv** | DISPOSITION | ACCOUNT_ID, CLIENT_ID, DISP_ID | |
| **dstrict.csv** | DISTRICT | DISTRICT_ID | Set name of A1 to DISTRICT_ID |
| **loan.csv** | LOAN | ACCOUNT_ID, LOAN_ID | |
| **order.csv** | PERMANENT_ORDER | ACCOUNT_ID, ORDER_ID | Set type of ORDER_ID to INTEGER |
| **trans.csv** | TRANSACTION | ACCOUNT_ID, TRANS_ID | Set length of K_SYMBOL to 17 |

You should now see all tables in SAP HANA.

```
▲ 📂 Tables
       ⊞ ACCOUNT
       ⊞ CARD
       ⊞ CLIENT
       ⊞ DISPOSITION
       ⊞ DISTRICT
       ⊞ LOAN
       ⊞ PERMANENT_ORDER
       ⊞ TRANSACTION
```

*Data Transformation*
Some of the data needs some minor transformation. The columns with dates for instance need to be transformed from numeric type to Date type.

Open the "SQL Console" in the SAP HANA Studio to execute the following SQL code, using your schema name instead of "I056450". Execution can be started with the green arrow on the top right hand side.



```
/*
Table: ACCOUNT
Turn "DATE" column into DATE format
*/
ALTER TABLE "I056450"."ACCOUNT" ADD ("DATE_TEMP" NVARCHAR(6));
UPDATE "I056450"."ACCOUNT"  SET "DATE_TEMP" = TO_NVARCHAR("DATE");
ALTER TABLE "I056450"."ACCOUNT" DROP ("DATE");
ALTER TABLE "I056450"."ACCOUNT" ADD ("DATE" DATE);
UPDATE "I056450"."ACCOUNT"  SET "DATE" = TO_DATE("DATE_TEMP", 'YYMMDD');
ALTER TABLE "I056450"."ACCOUNT" DROP ("DATE_TEMP");
```

```sql
/*
Table: CARD
Turn "ISSUED" column into DATE format
*/
ALTER TABLE "I056450"."CARD" ADD ("ISSUED_TEMP" NVARCHAR(15));
UPDATE "I056450"."CARD"  SET "ISSUED_TEMP" = "ISSUED";
ALTER TABLE "I056450"."CARD" DROP ("ISSUED");
ALTER TABLE "I056450"."CARD" ADD ("ISSUED" DATE);
UPDATE "I056450"."CARD"  SET "ISSUED" =
      TO_DATE(LEFT(TO_NVARCHAR("ISSUED_TEMP"), 6), 'YYMMDD');
ALTER TABLE "I056450"."CARD" DROP ("ISSUED_TEMP");


/*
Table: CLIENT
Add "GENDER" column
*/
ALTER TABLE "I056450"."CLIENT" ADD ("GENDER" NVARCHAR(6));
UPDATE "I056450"."CLIENT"  SET "GENDER" = 'MALE'
      WHERE MOD("BIRTH_NUMBER", 10000) < 2000;
UPDATE "I056450"."CLIENT"  SET "GENDER" = 'FEMALE'
      WHERE MOD("BIRTH_NUMBER", 10000) > 2000;


/*
Table: CLIENT
Add "DATE_BIRTH" column
Drop "BIRTH_DATE" column
*/
ALTER TABLE "I056450"."CLIENT" ADD ("DATE_BIRTH" DATE);
UPDATE "I056450"."CLIENT"  SET "DATE_BIRTH" =
       TO_DATE(TO_NVARCHAR("BIRTH_NUMBER" - mod("BIRTH_NUMBER" , 10000)
       + mod("BIRTH_NUMBER" , 5000)+19000000), 'YYYYMMDD');
ALTER TABLE "I056450"."CLIENT" DROP ("BIRTH_NUMBER");


/*
Table: LOAN
Turn "DATE" column into DATE format
*/
ALTER TABLE "I056450"."LOAN" ADD ("DATE_TEMP" NVARCHAR(6));
UPDATE "I056450"."LOAN"  SET "DATE_TEMP" = "DATE";
ALTER TABLE "I056450"."LOAN" DROP ("DATE");
ALTER TABLE "I056450"."LOAN" ADD ("DATE" DATE);
UPDATE "I056450"."LOAN"  SET "DATE" =
      TO_DATE(LEFT(TO_NVARCHAR("DATE_TEMP"), 6), 'YYMMDD');
ALTER TABLE "I056450"."LOAN" DROP ("DATE_TEMP");


/*
Table: TRANSACTION
Turn "DATE" column into DATE format
*/
ALTER TABLE "I056450"."TRANSACTION" ADD ("DATE_TEMP" NVARCHAR(6));
UPDATE "I056450"."TRANSACTION"  SET "DATE_TEMP" = TO_NVARCHAR("DATE");
ALTER TABLE "I056450"."TRANSACTION" DROP ("DATE");
ALTER TABLE "I056450"."TRANSACTION" ADD ("DATE" DATE);
UPDATE "I056450"."TRANSACTION"  SET "DATE" = TO_DATE("DATE_TEMP", 'YYMMDD');
ALTER TABLE "I056450"."TRANSACTION" DROP ("DATE_TEMP");
```

*Data Translation*

Finally let's translate all Czech data content into English. Execute the following SQL statements, having replaced "I056450" with your own schema name.

```sql
/*
Table: ACCOUNT
Translate "FREQUENCY" content into english
*/
UPDATE "I056450"."ACCOUNT" SET "FREQUENCY" =
    (
    CASE
      WHEN ("FREQUENCY" = 'POPLATEK MESICNE') THEN 'MONTHLY'
        WHEN ("FREQUENCY" = 'POPLATEK TYDNE') THEN 'WEEKLY'
        WHEN ("FREQUENCY" = 'POPLATEK PO OBRATU') THEN 'AFTER TRANSACTION'
        ELSE "FREQUENCY"
      END
    );


/*
Table: PERMANENT_ORDER
Translate "K_SYMBOL" content into english
*/
ALTER TABLE "I056450"."PERMANENT_ORDER" ALTER ("K_SYMBOL" NVARCHAR(20));
UPDATE "I056450"."PERMANENT_ORDER" SET "K_SYMBOL" =
    (
    CASE
      WHEN ("K_SYMBOL" = 'POJISTNE') THEN 'INSURANCE PAYMENT'
        WHEN ("K_SYMBOL" = 'SIPO') THEN 'HOUSEHOLD PAYMENT'
        WHEN ("K_SYMBOL" = 'LEASING') THEN 'LEASING'
        WHEN ("K_SYMBOL" = 'UVER') THEN 'LOAN PAYMENT'
        ELSE "K_SYMBOL"
    END
    );


/*
Table: TRANSACTION
Translate "K_SYMBOL" content into english
*/
ALTER TABLE "I056450"."TRANSACTION" ALTER ("K_SYMBOL" NVARCHAR(40));
UPDATE "I056450"."TRANSACTION" SET "K_SYMBOL" =
    (
    CASE
      WHEN ("K_SYMBOL" = 'POJISTNE') THEN 'INSURANCE PAYMENT'
        WHEN ("K_SYMBOL" = 'SLUZBY') THEN 'PAYMENT FOR STATEMENT'
        WHEN ("K_SYMBOL" = 'UROK') THEN 'INTEREST CREDITED'
        WHEN ("K_SYMBOL" = 'SANKC. UROK') THEN 'SANCTION INTEREST IF NEGATIVE
BALANCE'
        WHEN ("K_SYMBOL" = 'SIPO') THEN 'HOUSEHOLD'
        WHEN ("K_SYMBOL" = 'DUCHOD') THEN 'OLD-AGE PENSION'
        WHEN ("K_SYMBOL" = 'UVER') THEN 'LOAN PAYMENT'
        ELSE "K_SYMBOL"
    END
    );
```

```
/*
Table: TRANSACTION
Translate "OPERATION" content into english
*/
ALTER TABLE "I056450"."TRANSACTION" ALTER ("OPERATION" NVARCHAR(30));
UPDATE "I056450"."TRANSACTION" SET "OPERATION" =
      (
    CASE
      WHEN ("OPERATION" = 'VYBER KARTOU') THEN 'CREDIT CARD WITHDRAWAL'
        WHEN ("OPERATION" = 'VKLAD') THEN 'CREDIT IN CASH'
        WHEN ("OPERATION" = 'PREVOD Z UCTU') THEN 'COLLECTION FROM ANOTHER BANK'
        WHEN ("OPERATION" = 'VYBER') THEN 'WITHDRAWAL IN CASH'
        WHEN ("OPERATION" = 'PREVOD NA UCET') THEN 'REMITTANCE TO ANOTHER BANK'
        ELSE "OPERATION"
    END
    );


/*
Table: TRANSACTION
Translate "TYPE" content into english
*/
ALTER TABLE "I056450"."TRANSACTION" ALTER ("TYPE" NVARCHAR(10));
UPDATE "I056450"."TRANSACTION" SET "TYPE" =
      (
    CASE
      WHEN ("TYPE" = 'PRIJEM') THEN 'CREDIT'
        WHEN ("TYPE" = 'VYDAJ') THEN 'WITHDRAWAL'
        WHEN ("TYPE" = 'VYBER') THEN 'WITHDRAWAL'
        ELSE "TYPE"
    END
    );
```

*Update Dates*

The dates used in the original data go up to the end of 1998. Add 16 years to these dates so that we work with more recent timeframes.

```
UPDATE "I056450"."ACCOUNT" SET "DATE" = ADD_YEARS("DATE", 16);
UPDATE "I056450"."CARD" SET "ISSUED" = ADD_YEARS("ISSUED", 16);
UPDATE "I056450"."CLIENT" SET "DATE_BIRTH" = ADD_YEARS("DATE_BIRTH", 16);
UPDATE "I056450"."LOAN" SET "DATE" = ADD_YEARS("DATE", 16);
UPDATE "I056450"."TRANSACTION" SET "DATE" = ADD_YEARS("DATE", 16);
```