



TIMi suite

integrated data mining solutions

***An automated predictive datamining tool
Data Preparation – File Format_{v1.04}***

December 2012



Data preparation - Introduction

When using **TIMi**, it is strongly suggested to accumulate the highest number of information (rows and columns) about the process to predict: don't reduce arbitrarily the number of rows or columns: always keep the "full dataset" (even if the target size is less than one percent). **TIMi** will use this extra information (that is not available to other "classical" datamining softwares that requires sampling in order to work) to produce great predictive models that, 99% of the time, outperform the predictive models constructed with any other datamining software. When using classical tools (such as SAS, SPSS, KXEN, etc.), we know that it's a very common miss-practice to perform a strong sampling on the dataset, this is why we insist here that you don't perform any kind of sampling.

Of course, to be completely rigorous, you should also not forget to "let aside" a TEST dataset that will be used to really assert the quality of the delivered predictive models (to be able to compare in an objective way the models constructed with different predictive datamining tools). See this web page that explains the importance of the TEST set:

http://www.business-insight.com/html/intelligence/bi_test_dataset.html

Creation/Learning dataset file format

TIMi can read dataset from many data sources: simple "txt" or ".csv" flat files, SAS files, ACCESS files, OLEDB link to any databases, ODBC link to any database. **But** for a first "quick benchmark", it's suggest to store the **creation dataset** inside a simple "txt" or ".csv" flat file (in order to prevent any inter-operability difficulties).

The "txt" or ".csv" flat file should follow the following format:

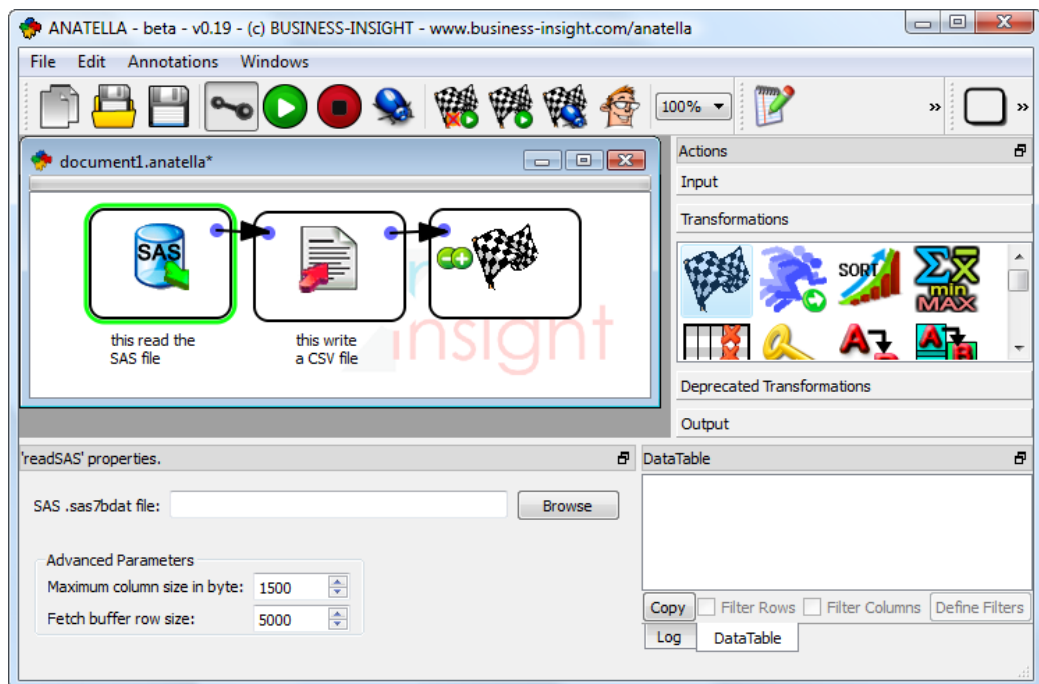
- The **creation dataset** is a "txt" or ".csv" file where the separator is a dot-comma ';'. The first line of the file must contain the column names.

WARNING: "txt" files exported from SAS have a size limitation: one line cannot exceed 65535 characters. If you encounter this bug in SAS, there are several solutions:

- Use the SAS plug-in for **TIMi** that allows you to read directly natively .sas7bdat SAS files (this is however very slow because the SAS driver is very slow).

To be able to read SAS files on any PC (even PC without any SAS installation), you should first install these Oledb drivers (that are originating from SAS):
<http://www.anatella.com/downloads/sasoledb.zip>

- Convert the .sas7bdat SAS file to a simple "txt" file using Anatella:
Use the following Anatella-data-transformation-script:



To be able to read SAS files on any PC (even PC without any SAS installation), you should first install these OleDB drivers (that are originating from SAS):
<http://www.anatella.com/downloads/sasoledb.zip>

- Column names must be unique.
WARNING: TIM is case IN-SENSITIVE (as is SQL)
- Column names are NOT within quotes.
- The data in the columns are NOT within quotes (never).
- The field separator character (here ';'') is not allowed (neither in the data, neither in the column names).
- The **creation dataset** contains one unique primary key.
- The decimal character is a dot and not a comma (Standard English notation or Scientific notation for numbers).
- If The **target** column (the column to predict) is:
 - Binary: then it must contains only '0' and '1' values (and the "one's" are the value to predict and must be the **minority case**).
 - Continuous: then it should not contain any "missing value".
- Missing values must always be encoded as empty values ("").

- OPTIONAL: The **creation dataset** should not contain any “consequence columns”. If the dataset nevertheless contains some “consequence columns”, it’s good to know their name in advance. However, you can always use **TIMi** to find all the “consequence columns” easily. To know more about “consequence columns”, see the pages 8,9,10,11 of the following document:
http://www.business-insight.com/downloads/DataPreparation_Churn.pdf
- OPTIONAL: the flat file can be compressed in RAR (.rar), GZip(.gz), Winzip(.zip)
- OPTIONAL: all the columns that represent a “True/False” information may contain only two different value: ‘0’ (for false) or ‘1’ (for true) or are empty (“”) if the value is missing.
- OPTIONAL: all the columns that represent either:
 - a number
 - an information that can be ordered
 ... should be encoded as pure number. For example:

number of cats		number of Cats
missing	→	
no cat		0
one cat		1
2 cats		2
3 or more		3

social class		social class
missing	→	
poor		0
middle		1
rich		2