# R – Adding a new dimension to Analytics

Christiane Raffeiner
Advanced Analytics Consultant
Teradata Center of Expertise

**TERADATA®**
*Raising Intelligence*

# Agenda

- Introduction
- R IDEs and R GUIs
- Teradata R
- R use in DWH environment
  - > Preparation of data in DWH
  - > Modeling in R
  - > Applying of Models in DWH
- Advanced Analytics Use Case with R

# Introduction

# What is R?

- R is a system for statistical computation and graphics
  - Open Source Statistics Package
  - Core functions plus 1000s of different packages
  - Consist a programming language plus a run-time environment – R console

- R can be used
  - interactively via expressions from the command line
  - through support of related GUIs or Editors
  - by writing your own functions

- Use
  - Growing number of data analyst inside corporations & academia
  - Ideal starting point in Pilots & Proof of Concepts
  - R as add on to existing commercial products to extend functionality

# R IDEs and R GUIs
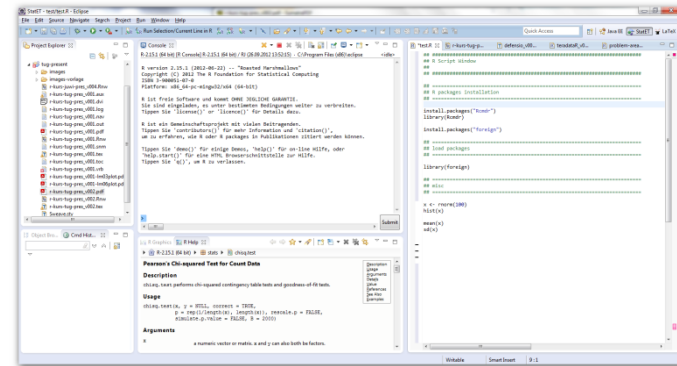
# Integrated Development Environments
(Free and open source)

- TinnR
  - > Probably one of the first IDEs for R
  - > Relatively simple (yet efficient) code editor
  - > Only available for Windows

- RStudio
  - > Comfortable to install
  - > Available for Linux, Mac, and Windows
  - > Integrated help, graphics, object browser, etc.
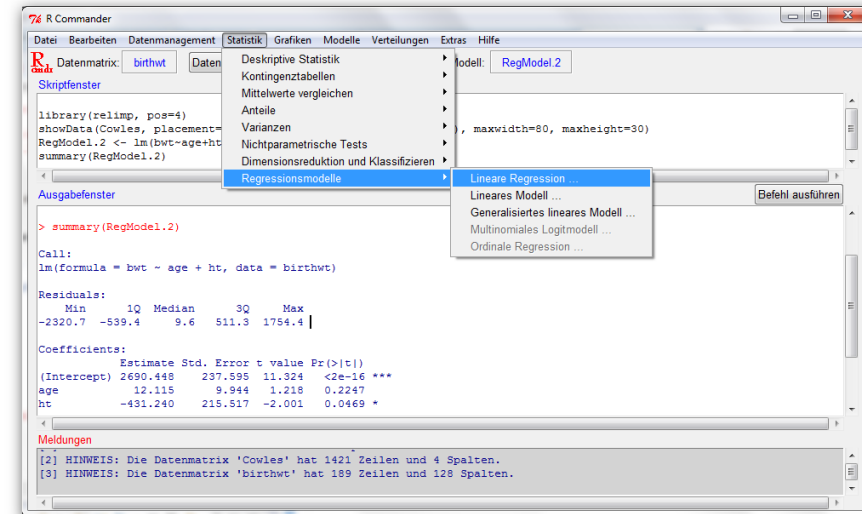


- Eclipse (with StatET plugin)
  - > Very powerful IDE
  - > Highly configurable, lots of functions and shortcuts
  - > Eclipse can be used for various purposes, including Teradata Access
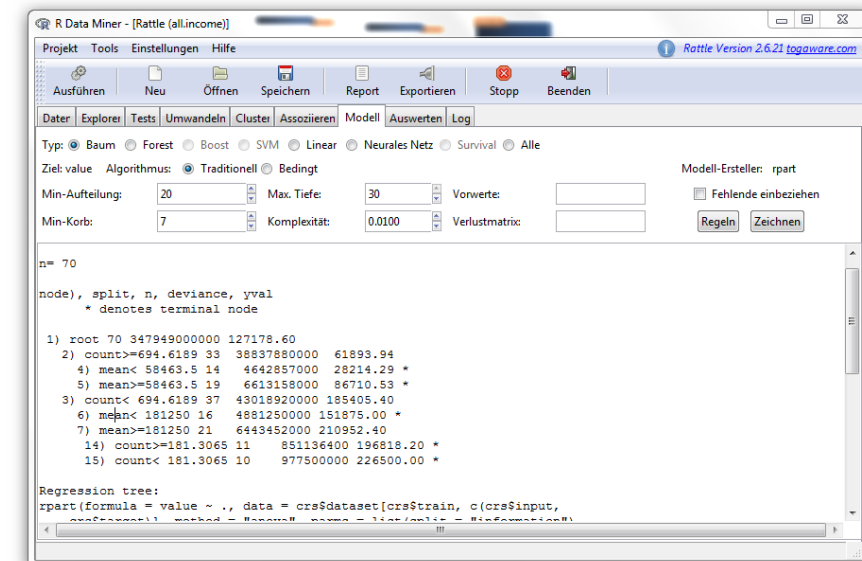
# R GUIs
(Free and open source)

- **RCOMMANDER**
  - > Offers point-and-click graphics surface
  - > Contains basic analyses
  - > Good for starters



- **RATTLE**
  - > Graphics User Interface
  - > For data mining
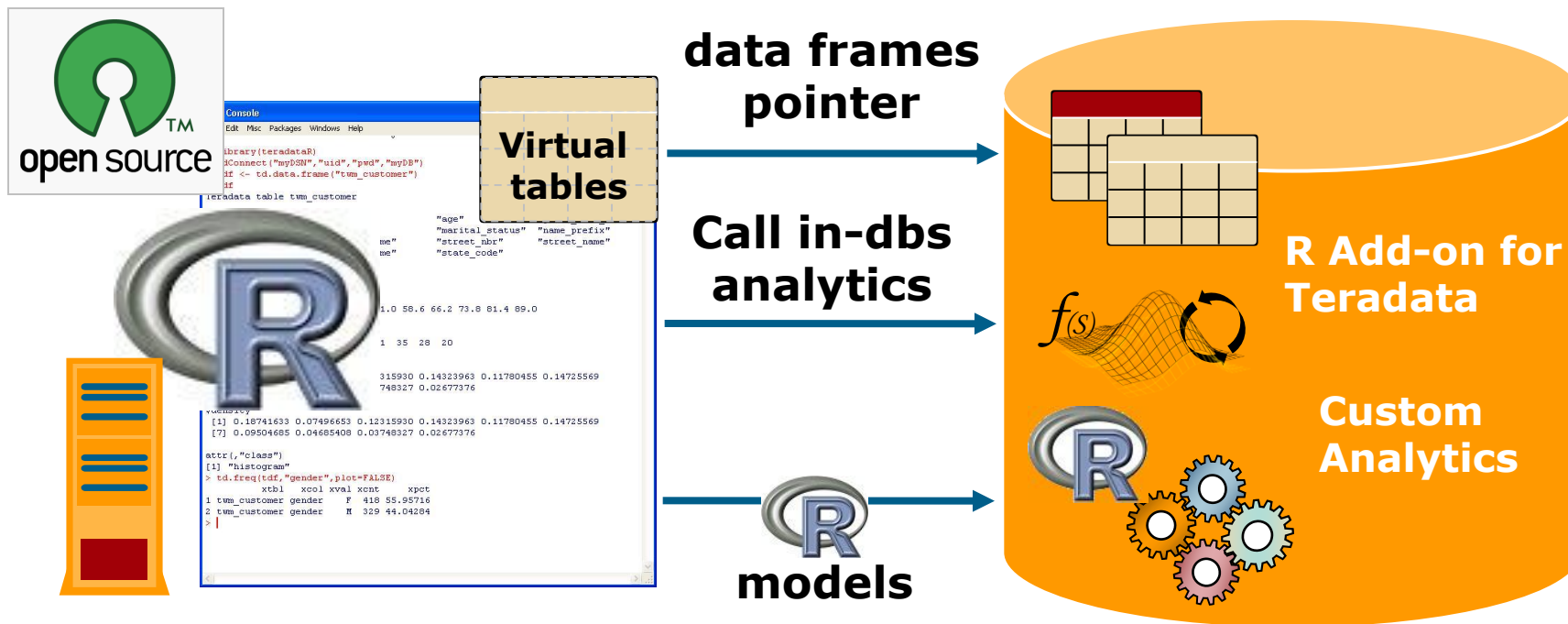  - > Can access Teradata by ODBC

# Teradata R

# Advanced Analytics

- Teradata in-database advanced analytics includes the following:
  - > Partner optimizations with SAS, IBM SPSS Modeler, KXEN
  - > Teradata Warehouse Miner
  - > Emerging technology: R



- Benefits
  - > Eliminate data movement to accelerate the process
  - > Lifts all "big data" limitations with Teradata's scalability
  - > Leverages the parallel processing of the database

TERADATA
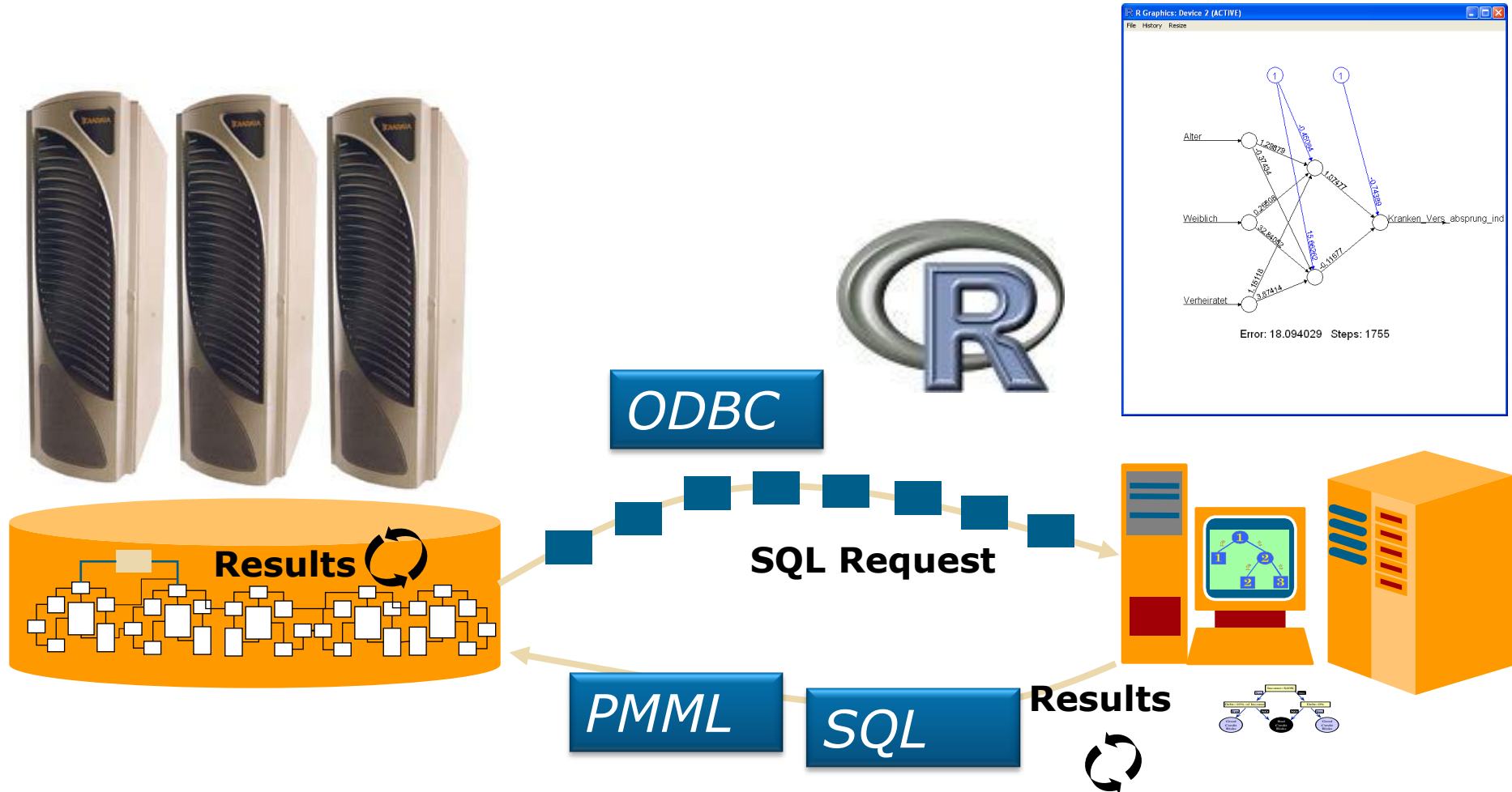
Raising Intelligence

# R Add-on for Teradata



**R Package** *teradataR*

- Simplifies connection to Teradata
- Establishes a data frame pointer (virtual table) to Teradata tables
- Provides over 40 in-database analytical functions
- Custom Analytics (udfs) created are callable by R
- R models can run in-database via PMML

# R Use in DWH Environment

# R and Teradata Architecture

# SQL – RODBC package

- RODBC package implements ODBC database connectivity
- RODBC allows direct access to tables stored in the DWH

- Functions provided
  - > Internal odbc* commands
  - > Sql* functions to read, save, copy and manipulate data

- SQL* arguments can be used within R and are passed through to the DWH
  - > sqlSave
  - > sqlCopy
  - > sqlQuery
  - > sqlFetch

R Package *RODBC*

TERADATA
Raising Intelligence

# PMML – an open data format



- Predictive Model Markup Language (PMML)
  - > is an XML based standard to describe statistical and data mining models
  - > since 1998 developed by the Data Mining Group (www.dmg.org)
  - > describes the data inputs to data mining models, some algorithm specific transformations used to prepare data for data mining, and the parameters which define the models themselves

- Application-independent method of defining models
  - > Users can develop models within one vendor's application, and use the model in any other application supporting PMML
  - > Requires a PMML consumer to read PMML output from a PMML producer and create an executable form of the model for scoring purposes

R Package *pmml*

TERADATA
*Raising Intelligence*

# Advanced Analytics Use Case

# RETAIL Price Optimization

**TERADATA**

*Raising Intelligence*

# Pilot Objectives

- Main Goals:

  > Understand how to define promotional sales prices

  > Simulation of price changes and their impact on sales quantity using a predefined set of stores and chosen articles of predefined categories

  > Estimation of the products price elasticity and cross price elasticity between products of the chosen category:
    > Price Elasticity, Cross Selling & Cannibalization

# Step 1: Data Preparation

- Create Analytical Data Set for Modeling in R
  - > Combination of very close price points
- Data Filtering
  - > Articles/stores with minimum number of observations and minimum price spread
  - > Elimination of Outliers using descriptive statistics

# Step 2: ODBC connection to Teradata database

- R libraries used:
  - > „RODBC" for ODBC connection to TD
  - > „stats" for regression
  - > „gplot" for plotting results

# Step 3: Price Elasticity Models with R

- Four different models created per article/region:
  - Linear model
  - Multiplicative model
  - Two Gutenberg-models

```
i = 1
      while (i <= article_cnt)
      {
      article <- grp_by_art[[i]]
      store_id <- paste(article$Region[1]) #_CAN
      art_id <- max(article$ID_ART)
      art_des <- paste(article$DES_ART[1])
      max_QTA <- max(article$QTA)
      cod_cat <- paste(article$COD_CAT_NEW_DES[1])
      min_PRZ <- min(article$PRZ)
      max_PRZ <- max(article$PRZ)

# linear model
      try_default(nls_lin <- nls(QTA ~ a + b*PRZ, data = article),default=NULL)
      std_error_lin=sqrt(deviance(nls_lin)/df.residual(nls_lin))

# multiplicative model
      try_default(nls_mult <- nls(QTA ~ a*PRZ^b, data = article),default=NULL)
      std_error_mult=sqrt(deviance(nls_mult)/df.residual(nls_mult))

# Gutenberg
      avg_price=max(article$PRZ_avg)
      try_default(nls_Gut <- nls(QTA ~ a + b*PRZ + sinh(d*(avg_price-PRZ)),
      data = article, start = list(a=coef(nls_lin)[1], b=coef(nls_lin)[1], d=1),alg = "port"),default=NULL)
      try_default(std_error_Gut <- sqrt(deviance(nls_Gut)/df.residual(nls_Gut)),default=NULL)

# Gutenberg2
      try_default(nls_Gut2 <- nls(QTA ~ a + b*PRZ + d*sinh(avg_price-PRZ),
      data = article, start = list(a=coef(nls_lin)[1], b=coef(nls_lin)[1], d=1),alg = "port"),default=NULL)
      std_error_Gut2=sqrt(deviance(nls_Gut2)/df.residual(nls_Gut2))

i = i +1
}
```

# Step 4: Model Evaluation - R plots

- Developed R-script to build price elasticity models automatically for all articles and plot results
- Results are automatically plotted and saved for evaluation

# Overall Analytical Process Project Steps



1. Derive, Aggregate, Transform to create Analytical Data Set (ADS)

2. Access ADS from R via RODBC

3. Create Model in R

4. Evaluate Model

5. Matrix with regression results and coefficients is written back to TD database

6. Results can now be accessed by Front End via ODBC connection

# Using Results – What if Analysis

## WHAT IF (price --> quantity)

| Article ID | Price Input | PriceElasticity | EstimatedQuantity |
|---|---|---|---|
| 986 | 3.7 | -3.01 | 63.17 |

| AVG Price A | Price article_B | Quantity article_A | Turnover (Price*Q TA) | Cross Price Elasticity | Elasticity |
|---|---|---|---|---|---|
| 4.05 | 4.12 | 52.97 | 214.75 ⬆ | 4.86 | -7.88 |

| AVG Price A | Price article_B | Quantity article_A | Turnover (Price*Q TA) | Cross Price Elasticity | Elasticity |
|---|---|---|---|---|---|
| 4.05 | 3.17 | 14.84 | 60.14 ⬆ | 4.86 | -7.88 |

| AVG Price A | Price article_B | Quantity article_A | Turnover (Price*Q TA) | Cross Price Elasticity | Elasticity |
|---|---|---|---|---|---|
| 4.05 | 2.22 | 2.63 | 10.67 ⬆ | 4.86 | -7.88 |

**Sales vs Price**

- price B=2.22
- price B=3.17
- price B=4.12

TERADATA. THE BEST DECISION POSSIBLE

# Summary

- Advantages
  - > Open source - free
  - > Choice of thousands of analytical functions
  - > Flexibility and extendability
  - > Possibility to go beyond commercial products
  - > Use R as add-on or single data mining tool
  - > Easy fit into existing working processes
  - > R can interact with other programs easily
  - > Expert Interaction and Support through Forums
  - > Fast implementation of new algorithms due to active community

- Challenges
  - > Unsupported tools need strong solution partnerships
  - > Code based work maintenace requires well defined processes
  - > Buy in traditional GUI users and grow R skills
  - > Sound statistical knowledge important for exploiting full functionality
  - > Limitations in-memory processing
  - > Rather slow in computation time for some calculations

TERADATA
Raising Intelligence

**Christiane Raffeiner**
Analytical Consultant
Vienna Advanced Analytics
Center of Expertise

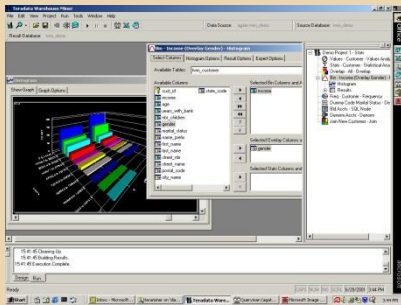Teradata GmbH
christiane.raffeiner@teradata.com

Teradata.at

# Backup

# Advanced Analytics Best Practice

## Data Exploration

**Explore all the data directly in the database with Teradata Profiler**

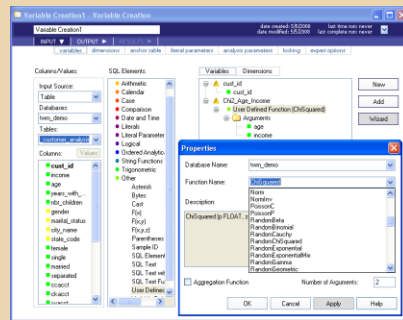

**Profiler**

## Data Preparation

**Transform and aggregate data in the database with Teradata ADS Generator**



**ADS Generator**

## Model Development
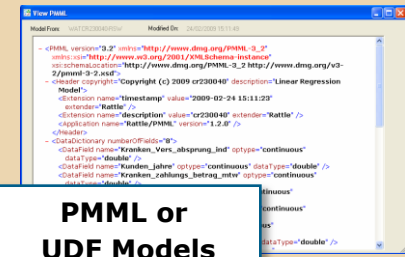
**Sample your ADS data and build your model on an R client**
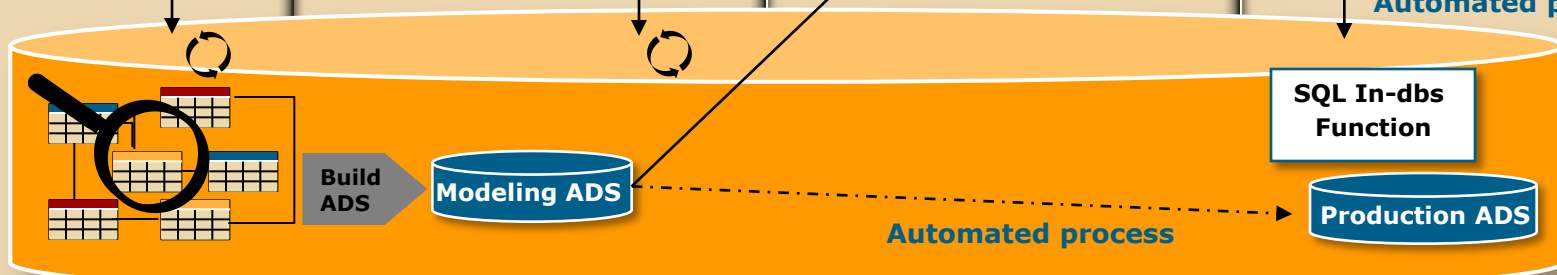


**Sample Data**

## Model Deployment

**Converts your R PMML model to SQL automatically generates the production ADS**



**PMML or UDF Models**

**ADS Generator**

**Automated process**

**SQL In-dbs Function**

**Build ADS** → **Modeling ADS** ⋯ **Automated process** → **Production ADS**

TERADATA
*Raising Intelligence*

# Sample R Code



```
R Console

File  Edit  Misc  Packages  Windows  Help

> library(teradataR)
> tdConnect("myDSN","uid","pwd","myDB")
> tdf <- td.data.frame("twm_customer")
> tdf
Teradata table twm_customer

 [1] "cust_id"        "income"       "age"            "years_with_bank"
 [5] "nbr_children"   "gender"       "marital_status" "name_prefix"
 [9] "first_name"     "last_name"    "street_nbr"     "street_name"
[13] "postal_code"    "city_name"    "state_code"

747 rows
> td.hist(tdf,"age",plot=FALSE)
$breaks
 [1] 13.0 20.6 28.2 35.8 43.4 51.0 58.6 66.2 73.8 81.4 89.0

$counts
 [1] 140  56  92 107  88 110  71  35  28  20

$intensities
 [1] 0.18741633 0.07496653 0.12315930 0.14323963 0.11780455 0.14725569
 [7] 0.09504685 0.04685408 0.03748327 0.02677376

$density
 [1] 0.18741633 0.07496653 0.12315930 0.14323963 0.11780455 0.14725569
 [7] 0.09504685 0.04685408 0.03748327 0.02677376

attr(,"class")
[1] "histogram"
> td.freq(tdf,"gender",plot=FALSE)
         xtbl   xcol xval xcnt      xpct
1 twm_customer gender    F  418 55.95716
2 twm_customer gender    M  329 44.04284
> |
```

1) Teradata library

2) Connect to Teradata

3) Establish a pointer to a Teradata table

4) View Teradata table variables

5) Call in-database analytics

6) View results

TERADATA
Raising Intelligence