# Microsoft R Server and SQL Server 2016

Tomaž Kaštrun

e: tomaz.kastrun@gmail.com
t: @tomaz_tsql
b: https://tomaztsql.wordpress.com

satRday #1

September 3 2016
MTA TTK, Budapest, Hungary

Microsoft

# About

- BI Developer and data analyst
- 15years experience with MSSQL Server
- 15years experience data analysis and DM
- Working:  Spar ICS Austria, Spar Slovenija
- MCT, MCPT, MCSE SQL Server
-  tomaz.kastrun@gmail.com
-  @tomaz_tsql

https://tomaztsql.wordpress.com

- Frequent community speaker at SQL and Microsoft events
- Blogger, Avid Coffee Lover, Bicycle junkie

# Analytical Barriers

Common Challenges

Addressing Challanges with R from Microsoft

| | | | |
|---|---|---|---|
| Uncertain total cost of ownership | Inadequate access to important business data | Limited business agility | Limited business value |
| Peace of mind | Efficiency | Speed and scalability | Flexibility and agility |

# What is R?

- ## A Language Platform
  - A Procedural Language optimized for Statistics and data science (and much more)
  - A Data Visualization framework
  - Provided as Free Software

- ## A Community and a system
  - Taught on universitieis and many active user groups across the world
  - Estimated 3Mio Users
  - Repositories (CRAN, BioConductor, Github,…)

  In fact, <u>R is a movement</u>!

# Limitations of R as a free software

- Memory Based Data access model
- Interpreted vs. Compiled Performance
- Lack of parallel computation
- Data movement & Duiplication Costs
- Governance and providence oversight
- Community support vs. Enterprise utilization

# Revolution Analytics Product Integration

# Microsoft R SQL Server platform



-> Free and open Source R distribution
-> Enhanced and distributed by
Revolution analytics

-> Built in Advanced Analytics and Standalone
Server Capability
-> Leverages the benefits of SQL Server 2016EE

# Microsoft R Platform

Microsoft R Open

Microsoft R Client

Microsoft SQL  R Services

Microsoft R Server

Different flavors:
Microsoft R server for Linux,
Microsoft R Server for  Teradata,
Microsoft R Server for Hadoop,
Microsoft R HDInsight

# Microsoft R Server

- Evolved from Revolution R Enterprise
- Based on open Source R
- Adapted for Enterprise Scale
- For multiple platforms
  - Hadoop
  - Teradata
  - LinuX
  - Azure
  - Windows

- Interoperable
- On-premises + Cloud + Hybrid
- Operationzalize analytics for Big scale datasets and big data

# Based on Open Source R

- Open source based

- Runs your normal R Script

- MetaCran / CRAN / Github / Bioconductor

# Microsoft R Server

## DeployR

- Web service - API integration
- Compatible with array of tools
- Abstract usage of R without knowing it

## DevelopR

- R IDE based on Visual studio
- Rstudio for linux Users
- Client Based

# Microsoft R Server

## ConnectR

- Serier of connectors for consistent access to scaleR algorithms

## DistributedR

- Normalization layer for ScaleR algoritms (SQLServer,Win, Lin, TeraData, Hadoop, HDI)

## ScaleR

- Typical statistical approaches refactored for parallel computation
- Block-wise computation; No In-Memory constraints

# ScaleR algorithms

## Data Preparation

- Data import – Delimited, Fixed, SAS, SPSS, OBDC
- Variable creation & transformation
- Recode variables
- Factor variables
- Missing value handling
- Sort, Merge, Split
- Aggregate by category (means, sums)

## Descriptive Statistics

- Min / Max, Mean, Median (approx.)
- Quantiles (approx.)
- Standard Deviation
- Variance
- Correlation
- Covariance
- Sum of Squares (cross product matrix for set variables)
- Pairwise Cross tabs
- Risk Ratio & Odds Ratio
- Cross-Tabulation of Data (standard tables & long form)
- Marginal Summaries of Cross Tabulations

## Statistical Tests

- Chi Square Test
- Kendall Rank Correlation
- Fisher's Exact Test
- Student's t-Test

## Sampling

- Subsample (observations & variables)
- Random Sampling

## Predictive Models

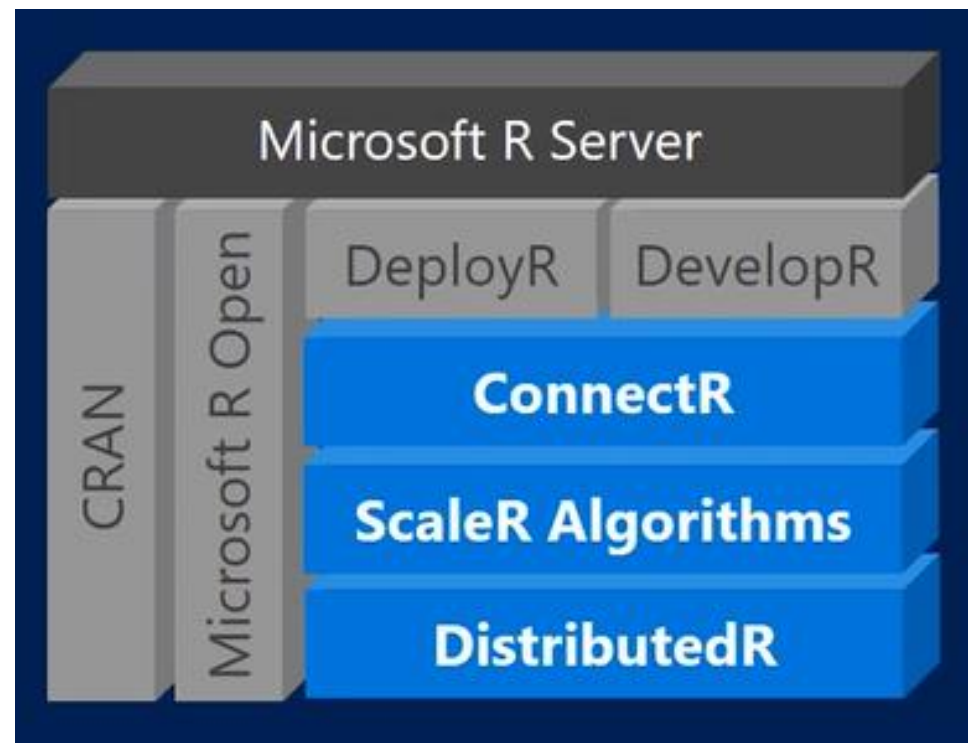- Sum of Squares (cross product matrix for set variables)
- Multiple Linear Regression
- Generalized Linear Models (GLM) exponential family distributions: binomial, Gaussian, inverse Gaussian, Poisson, Tweedie. Standard link functions: cauchit, identity, log, logit, probit. User defined distributions & link functions.
- Covariance & Correlation Matrices
- Logistic Regression
- Classification & Regression Trees
- Predictions/scoring for models
- Residuals for all models

## Variable Selection

- Stepwise Regression

## Simulation

- Simulation (e.g. Monte Carlo)
- Parallel Random Number Generation

## Cluster Analysis

- K-Means

## Classification

- Decision Trees
- Decision Forests
- Gradient Boosted Decision Trees
- Naïve Bayes

## Combination

- rxDataStep
- rxExec
- PEMA-R API Custom Algorithms

# Microsoft R Server - architecture

Open Source R Provides:
- R Scripts
- CRAN Packages

Microsoft R Server Adds:
- Scalable Parallel Algorithms
- No Memory Limits
- High Quality R Development Tools
- Web Services Deployment
- Support and Services

Data Connectors

Data

Databases

hadoop Map Reduce

EDWs

File Systems

# Parallelizing data process

# R code in SQL Server as T-SQL

```sql
EXECUTE sp_execute_external_script
  @language = N'R'
  ,@script = N'
          library(e1071);
          irismodel <-naiveBayes(iris_data[,1:4], iris_data[,5]);
          trained_model <- data.frame(payload = as.raw(serialize(irismodel, connection=NULL)));'
  ,@input_data_1 = N'select "Sepal.Length", "Sepal.Width","Petal.Length","Petal.Width","Species" from iris_data'
  ,@input_data_1_name = N'iris_data'
  ,@output_data_1_name = N'trained_model'

WITH RESULT SETS ((model VARBINARY(MAX)));
```

# R code in SQL Server using Scale R algorithms

```sql
EXECUTE sp_execute_external_script
    @language = N'R'
    ,@script = N'require("RevoScaleR");
                irisLinMod <- rxLinMod(Sepal.Length ~ Sepal.Width + Petal.Length + Petal.Width + Species,
                    data = iris_rx_data);
                trained_model <- data.frame(payload = as.raw(serialize(irisLinMod, connection=NULL)));'
    ,@input_data_1 = N'select "Sepal.Length", "Sepal.Width", "Petal.Length", "Petal.Width", "Species" from
                    iris_rx_data'
    ,@input_data_1_name = N'iris_rx_data'
    ,@output_data_1_name = N'trained_model'

WITH result SETS ((model VARBINARY(MAX)));
```

# RevoScaleR Code

```r
12
13  #####~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~
14  #
15  #
16  #      LOADING DATA (small sample)
17  #      178 MB
18  #      8.4Mio Rows
19  #####~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~
20
21
22
23  ptm <- proc.time()
24  #inFile <- file.path(rxGetOption("sampleDataDir"), "AirlineDemoSmall.csv")
25  inFile <- file.path(rxGetOption("sampleDataDir"), "airsample.csv")|
26  rxTextToXdf(inFile = inFile, outFile = "airline.xdf",  stringsAsFactors = T, rowsPerRead = 200000, overwrite=TRUE)
27  proc.time() - ptm
28  # ~ 22 seconds!
29  # - 42 Chunks per 200.000 Rows; Total: 8.400.000 Rows
30
31  ###############################
32  #   EXPLORING DATA (small sample)
33  ###############################
34
35
36  rxGetInfo(data="airline.xdf", getVarInfo = TRUE, numRows = 5)
37
38  #Histograms by day of week
39  ptm <- proc.time()
40  rxHistogram( ~ ArrDelay|DayOfWeek, data = "airline.xdf")
41  proc.time() - ptm
42
43  #summary
44  rxSummary( ~ ArrDelay, data = "airline.xdf")
45
46
47  rxSort(inData="airline.xdf", outFile = "sortFlights.xdf", sortByVars="ArrDelay",  decreasing = TRUE,overwrite=TRUE)
48  # ~ 4 Seconds!
49  mostflights5 <- rxGetInfo(data = "sortFlights")
50  mostflights5
51  top5f <- as.data.frame(mostflights5[[5]])
52  topOA <- unique(as.vector(top5f$ArrDelay))
53  topOA
54
55
56
57  ###############################
58  #   Linear Model with ReportProgress!
59  ###############################
60
61  # Linear Model using rxLinMod
62  sampleDataDir <- rxGetOption("sampleDataDir")
63  airlineDemoSmall <- file.path(sampleDataDir, "AirlineDemoSmall.xdf")
64
        ineLinMod <- rxLinMod(ArrDelay ~ CRSDepTime, data = airlineDemoSmall,
.evel)
```

# R code in SQL Server as T-SQL to generate graphs

```sql
DECLARE @RScript nvarchar(max)
DECLARE @SQLScript nvarchar(max)

SET @RScript = N'library(plotly)
                 library(ggplot2)
                 library(htmlwidgets)
                 #setwd("C:/DataTK/HTML")
                  image_file <- tempfile()
                  jpeg(filename = image_file, width = 500, height = 500)
                  df <- InputDataSet
                  d <- df[sample(nrow(df), 10), ]
                  p <- plot_ly(d, x = OrderQty, y = DiscountPct, text = paste("OrderQty: ", OrderQty),
                        mode = "markers", color = OrderQty, size = OrderQty)
                 saveWidget(as.widget(p), "index.html")
                 OutputDataSet <- data.frame(data=readBin(file(image_file, "rb"), what=raw(), n=1e6))'


SET @SQLScript = N'SELECT
                    ps.[Name]
                   ,AVG(sod.[OrderQty]) AS OrderQty
                   ,so.[DiscountPct]
                   ,pc.name AS Category
               FROM  Adventureworks.[Sales].[SalesOrderDetail] sod
               INNER JOIN Adventureworks.[Sales].[SpecialOffer] so
               ON so.[SpecialOfferID] = sod.[SpecialOfferID]
               INNER JOIN Adventureworks.[Production].[Product] p
               ON p.[ProductID] = sod.[ProductID]
               INNER JOIN Adventureworks.[Production].[ProductSubcategory] ps
               ON ps.[ProductSubcategoryID] = p.ProductSubcategoryID
               INNER JOIN Adventureworks.[Production].[ProductCategory] pc
               ON pc.ProductCategoryID = ps.ProductCategoryID
               GROUP BY ps.[Name],so.[DiscountPct],pc.name'

EXECUTE sp_execute_external_script
@language = N'R',
@script = @RScript,
@input_data_1 = @SQLScript
WITH RESULT SETS ((Plot varbinary(max)))
```

# Benefits of R integration

- Based on Open source R
- Different versions available (Open, Client and Server)
- Distributed workloads, multi-threading and parallelization
- Interoperable (Windows, Linux, MacOS) with different flavors (Hadoop, Teradata, HDInsight)
- Faster model prediction and model deployment
- No „in-memory" constraints, less data movement, less bottlenecks in performance, no data size limitations
- Hybrid topologies, agile development, stable platform for data operationalization, investment protection (SLA, Terms and agreements)
- R Code is available in SSMS environment
- Community and commercial support
- R Language is growing in popularity

# Questions?

## Contacts:
Email: **tomaz.kastrun@gmail.com**
Twitter: **@tomaz_tsql**
Blog: https://tomaztsql.wordpress.com