



Evolving your Business with Data, Advanced Analytics, and AI

Joe Sack, Principal Program Manager,
Microsoft

Predicting Operating Room Capacity

Scenario

- **Healthcare company** with separate hospitals, uses a confederation model in which affiliate medical centers jointly borrow and purchase common services and information technology support, such as **operating room booking systems**
- **Last-minute operating room cancellations** result in rooms going unused, over-staffing and lengthy patient wait-lists.
- Customer uses a combination of SSAS, T-SQL calculations and a iterative Data Mining Model to predict operating room capacity based on past booking history
- This solution has been effective in reducing operating room “spoilage”, but the solution is higher in complexity than desired by the customer.

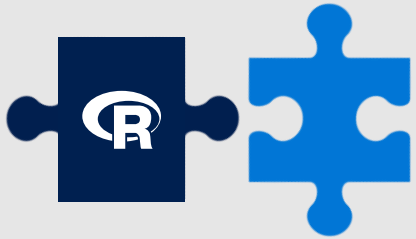
Solution

- Reduce complexity and reduce data movement by implementing the full modelling and prediction process within the SQL Server engine via SQL Server R Services
- Used rxLinMod function to locally train the model based on past booking history and predict future bookings behavior for future dates

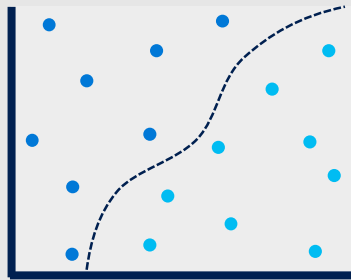


Deploying predictive analytics

Develop, explore and experiment in your favorite IDE or language



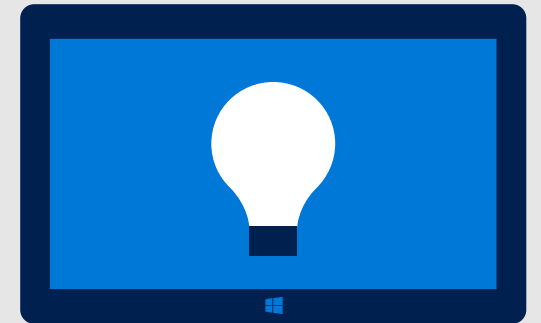
Train a model



Deploy and predict with the model



Make your apps **intelligent** by **consuming** predictions



Develop

Train

Deploy

Consume



Language
platform

Statistics programming language
Data visualization tool
Open source

Community

2.5+M users
Taught in most universities
Popular with new and recent grads
Thriving user groups worldwide

Ecosystem

10,000+ packages in CRAN
Scalable to big data
Rich application and platform integration

Challenges of using R

Data movement



- **Moving data from the DB to R**
- **Runtime becomes painful as data volumes grow**
- **Movement carries security risks**

Deployment



- **How do I call the R script from my production application?**

Scale and performance



- **Most R functions are single threaded and only accommodate datasets that fit into available memory**

Pain Points

Performance of data exploration activities may be slow (data locality | single-threaded | memory constraints)

Data may leave "trust boundary"

Solutions may not be easily operationalized, requiring re-coding in other solutions

Data Scientist may not be aware of already-cleansed and collected data

Data Scientist may be using data cleansing techniques which would be more efficiently performed by SQL Server

Database Engineer may be facilitating analytic and statistical operations which would be better performed using R

Microsoft Advanced Analytics Landscape

Azure Machine Learning

- Fully managed cloud service that enables you to easily build, deploy, and share predictive analytics solutions

SQL Server Machine Learning Services (In-Database)

- Supports both R and Python pushed SQL Server compute-context

Microsoft R Server

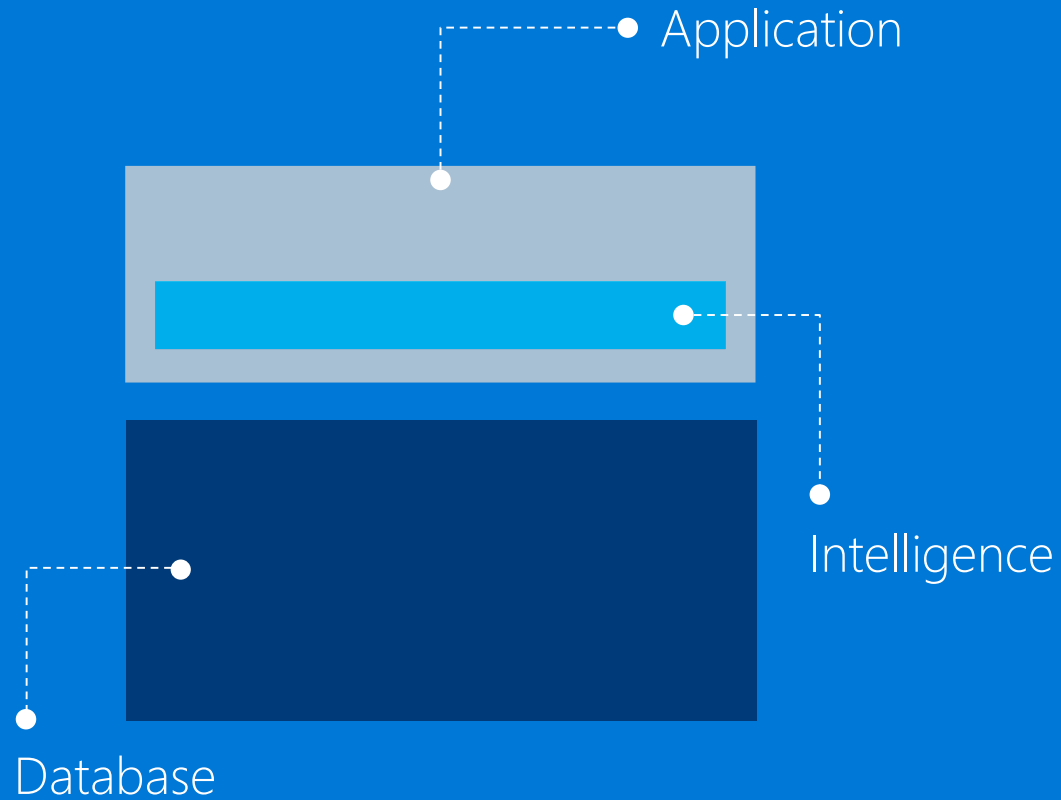
- For enterprise-level R deployments on Windows and Linux servers

Microsoft Machine Learning Server

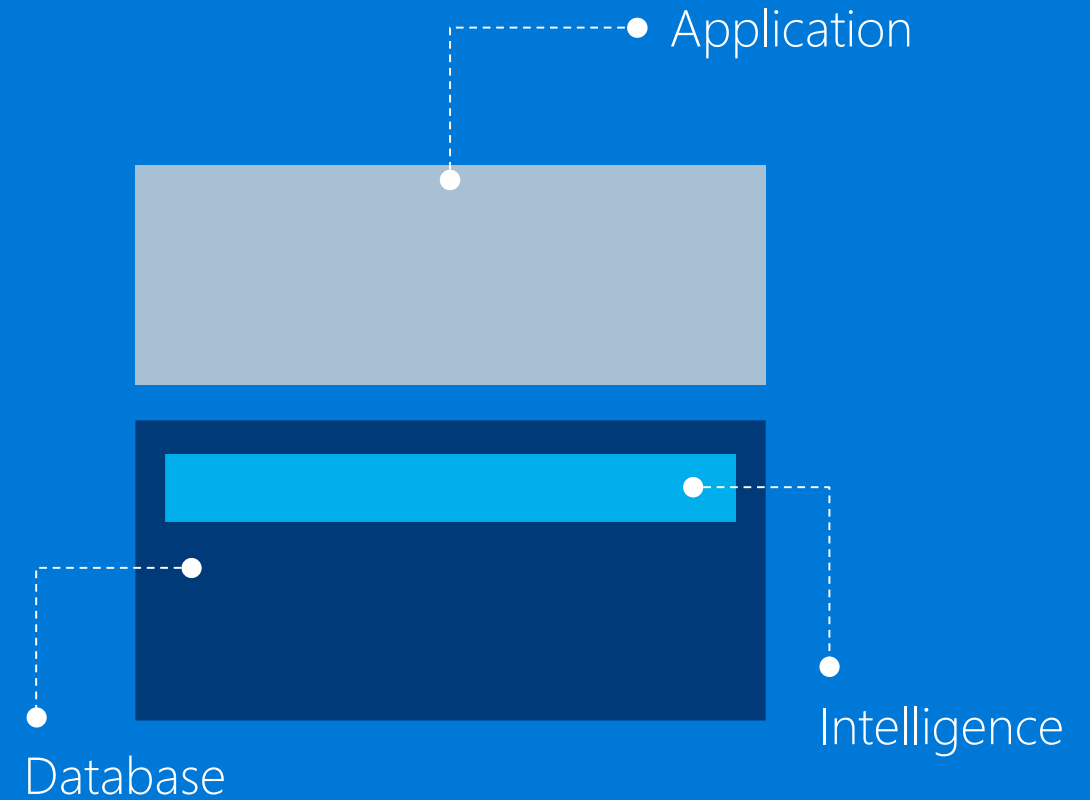
- Supports R and Python deployments on Windows servers, with expansion to other supported platforms planned for late 2017

In-database advanced analytics

Pushing intelligence to where data lives



Before



Intelligence built in to the DB

SQL Server ML Services solves problems

Reduce or eliminate data movement with in-database analytics



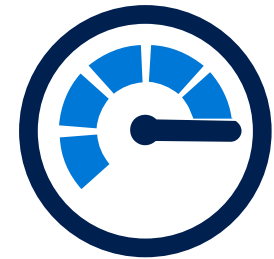
- **Leverage built-in extensibility mechanisms to allow secure execution of scripts**

Deploy R or Python scripts and models



- **Use familiar T-SQL stored procedures to invoke scripts from your app**
- **Embed the returned predictions and plots**

Achieve enterprise scale and performance



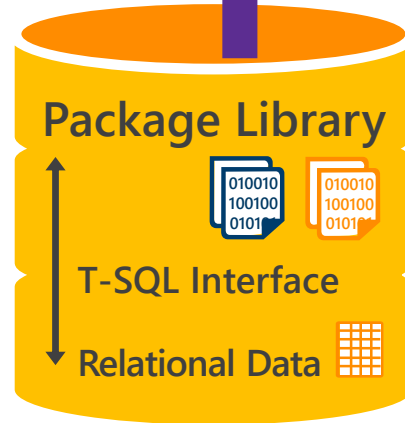
- **Use parallelism query capabilities of in-memory and ColumnStore indexes**
- **Leverage RevoScaleR and RevoScalePy support for large datasets and parallel algorithms**

SQL Server Machine Learning Services

Example Solutions

- Fraud detection
- Sales forecasting
- Warehouse efficiency
- Predictive maintenance

Extensibility



Data Scientist

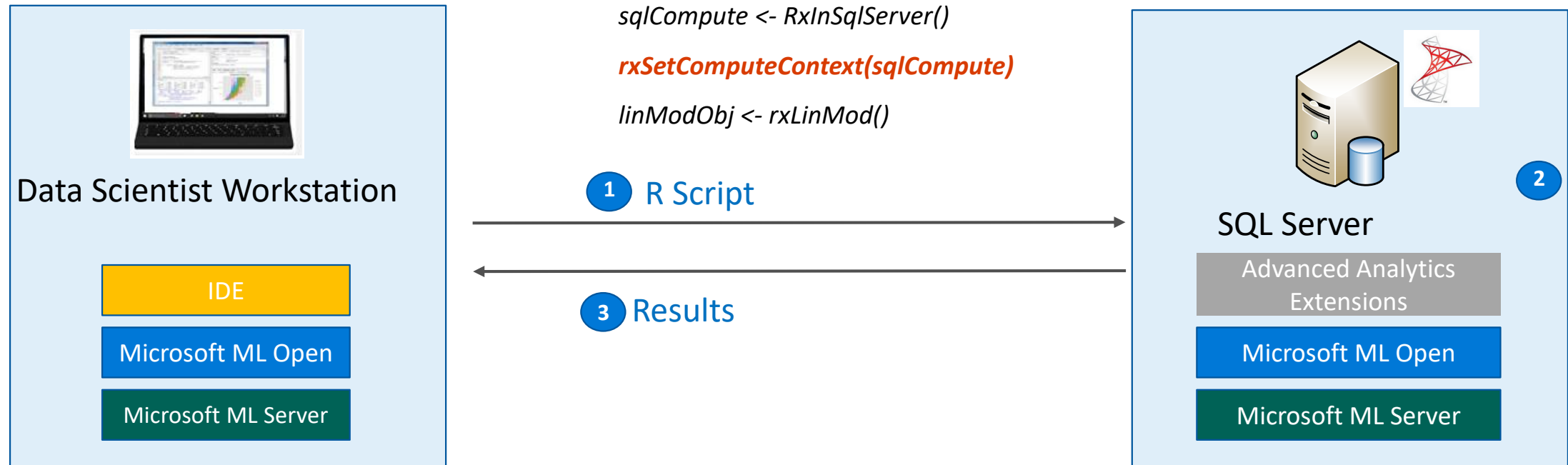
Interact directly with data



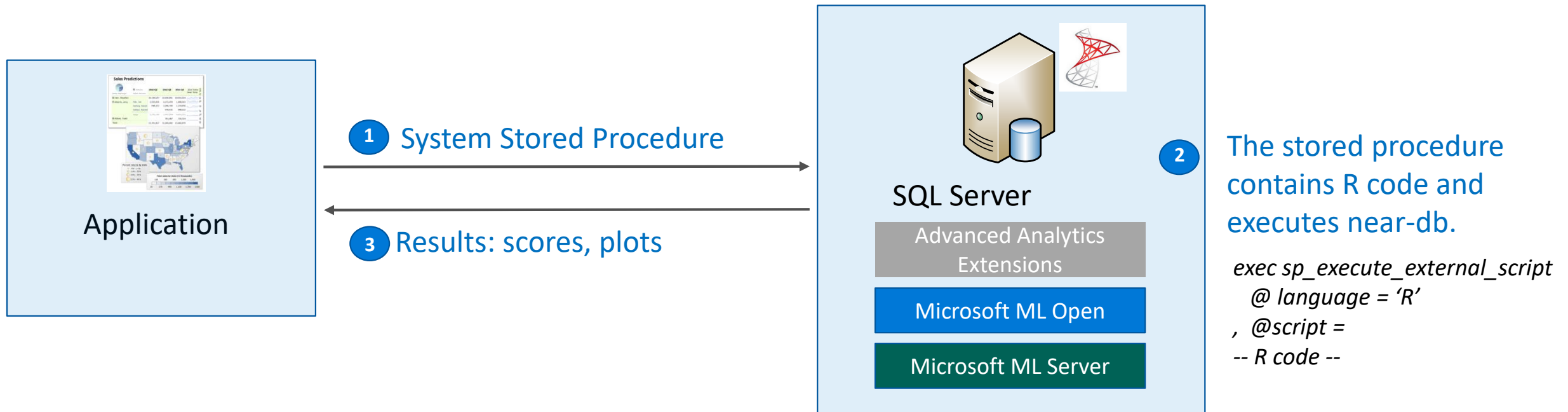
Data Developer/DBA

Manage data and analytics together

Data Scientist Scenario



SQL Server Developer Scenario



ML Server functions & algorithms

Data step

Data import – Delimited, fixed, SAS, SPSS, ODBC
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

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

Custom parallelization

PEMA-R API
rxDataStep
rxExec

sp_execute_external_script

```
sp_execute_external_script
  @language = N'language' ,
  @script = N'script',
  @input_data_1 = ] 'input_data_1'
[ , @input_data_1_name = ] N'input_data_1_name' ]
[ , @output_data_1_name = 'output_data_1_name' ]
[ , @parallel = 0 | 1 ]
[ , @params = ] N'@parameter_name data_type [ OUT | OUTPUT ] [ ,...n ]'
[ , @parameter1 = ] 'value1'[OUT | OUTPUT] [ ,...n ]
[ WITH <execute_option> [ ,...n ] ]
[;]
```



Demo

Customer Usage Drivers

Reduce data movement

Reduce complexity by
writing once in R

Use RevoScale optimized
function equivalents

Use compute context for
exploration activities,
leveraging server's
resources

Leverage SQL Server
design optimizations

Operationalize with
`sp_execute_external_script`

PROS, Inc.



Revenue & Profit Realization

PROS provides a real-time software solution platform to help companies drive pricing & sales effectiveness



Reduce Quote
Turnaround Time



Simplify Complex
Product Catalogs



Raise Win Rates



Grow Deal Sizes



Increase Margins



Reduce Sales Cycle
Duration

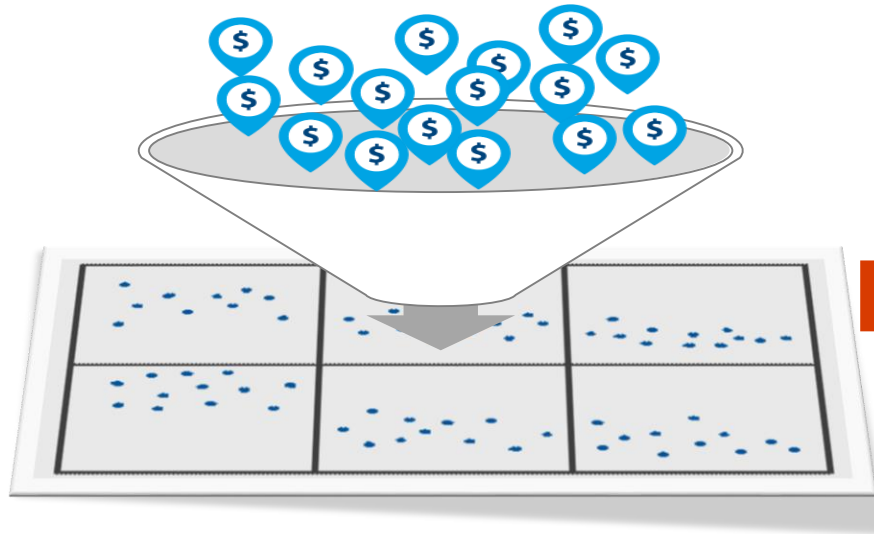


Increase Quota
Attainment

Segmentation

Foundation for PROS smart pricing guidance

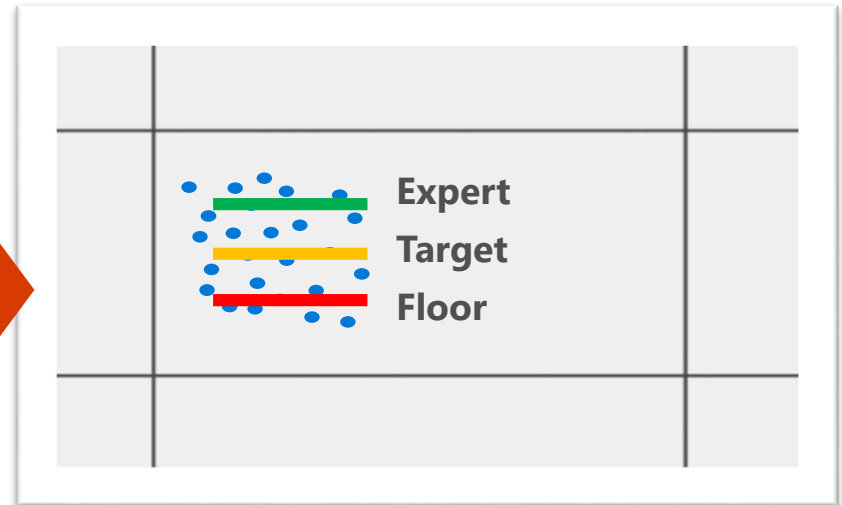
Segmentation



Group customers, products & transactions into micro-segments of similar willingness-to-pay



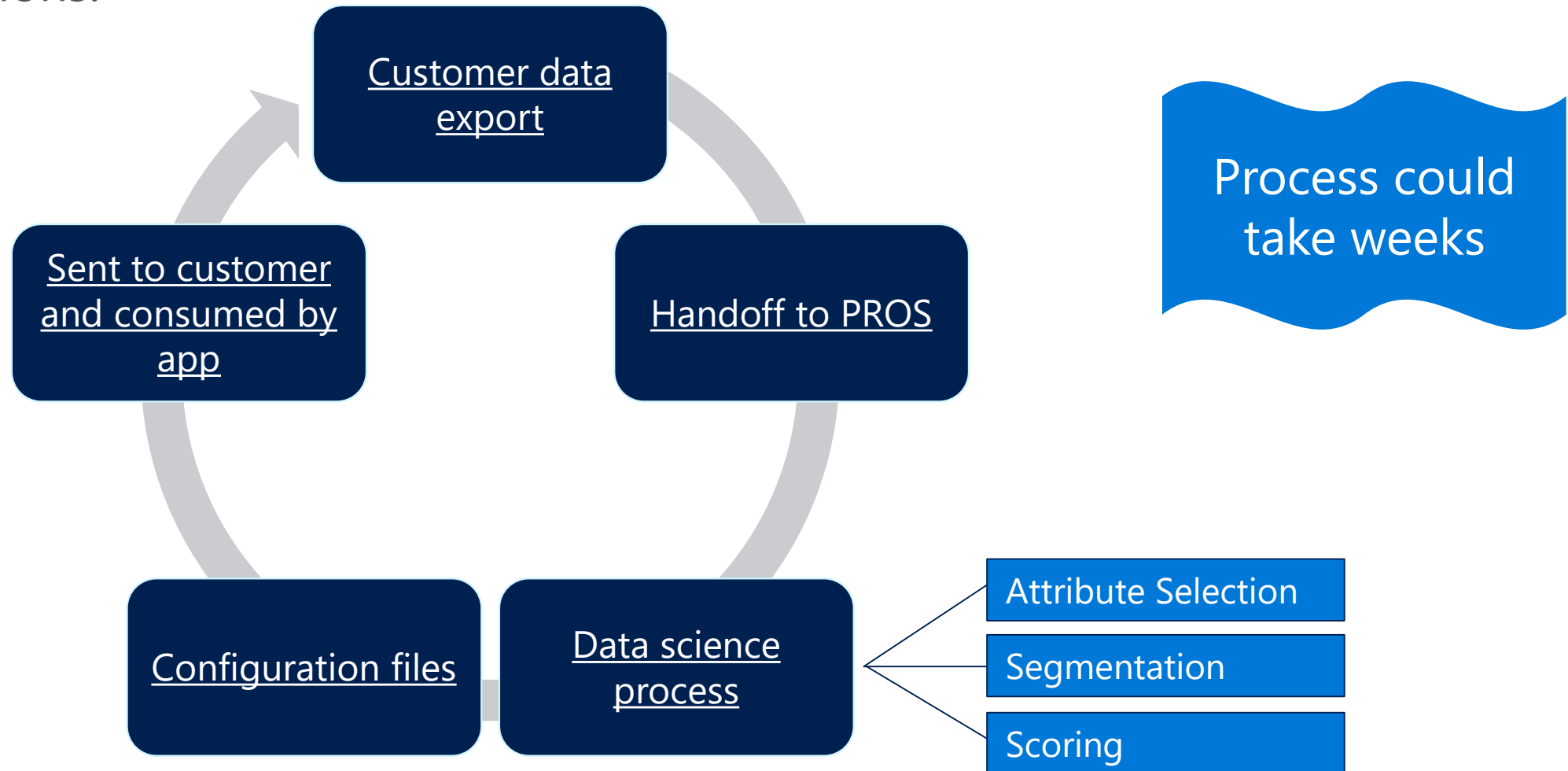
Pricing Optimization



Apply optimization algorithms to target the pricing envelope 'sweet spot' in every segment

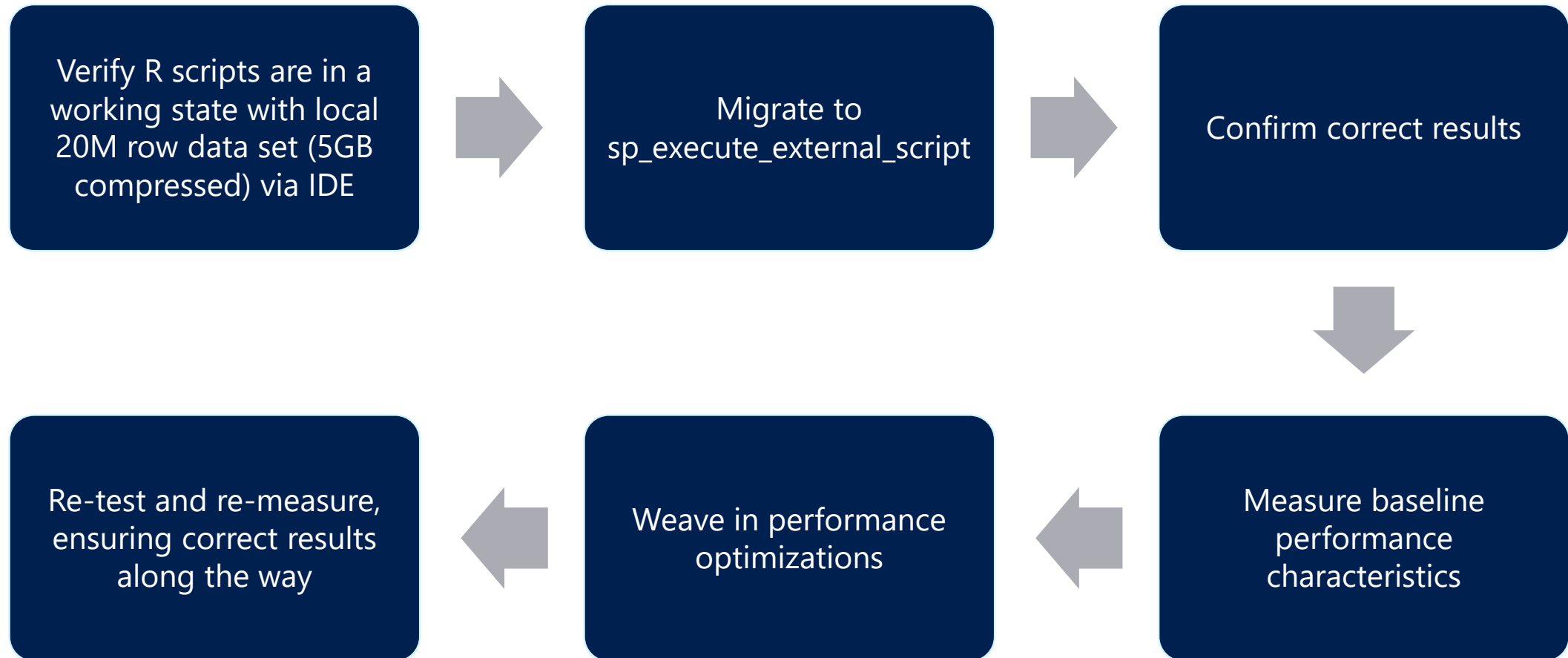
Pricing segmentation: Process

Identify appropriate set of business attributes and model how they can be used with statistical and practical soundness to provide intelligent **pricing benchmarks** for future transactions.

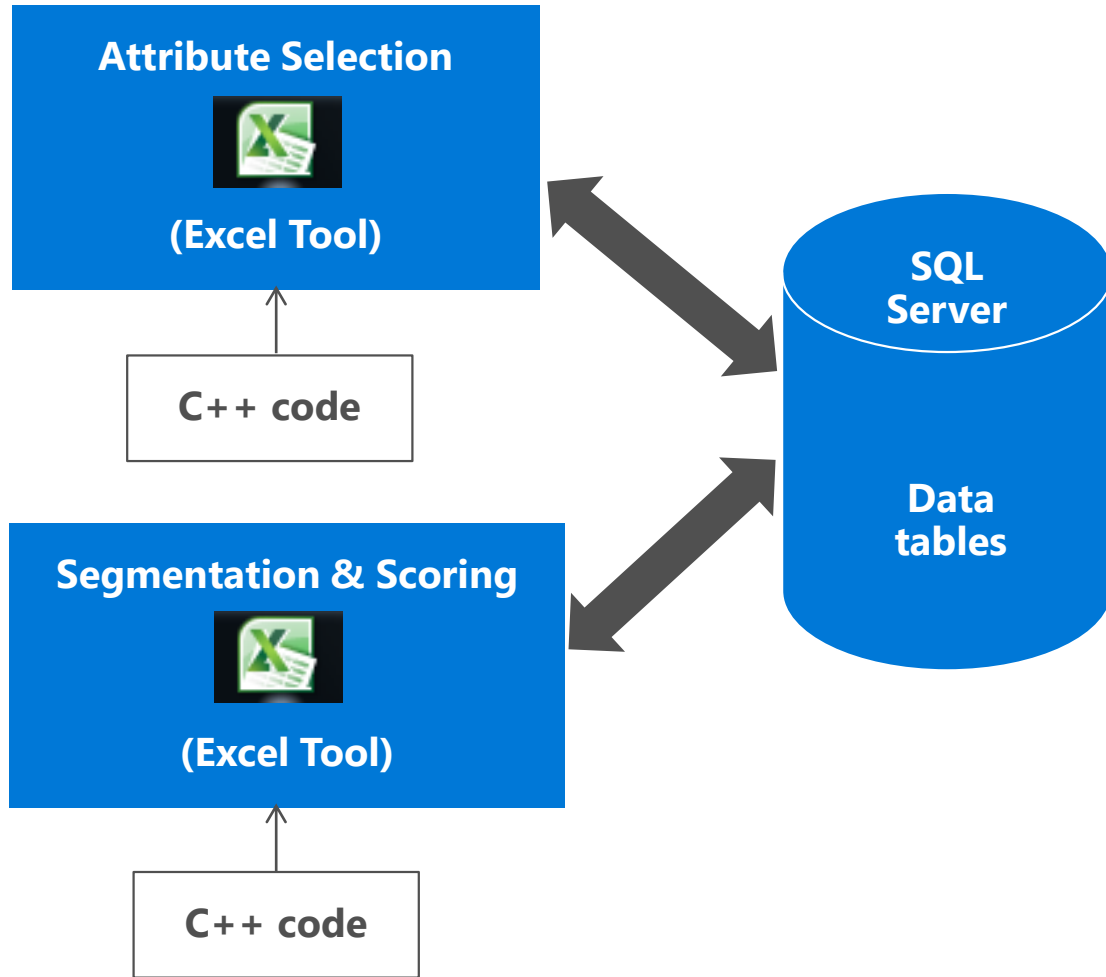


Solution

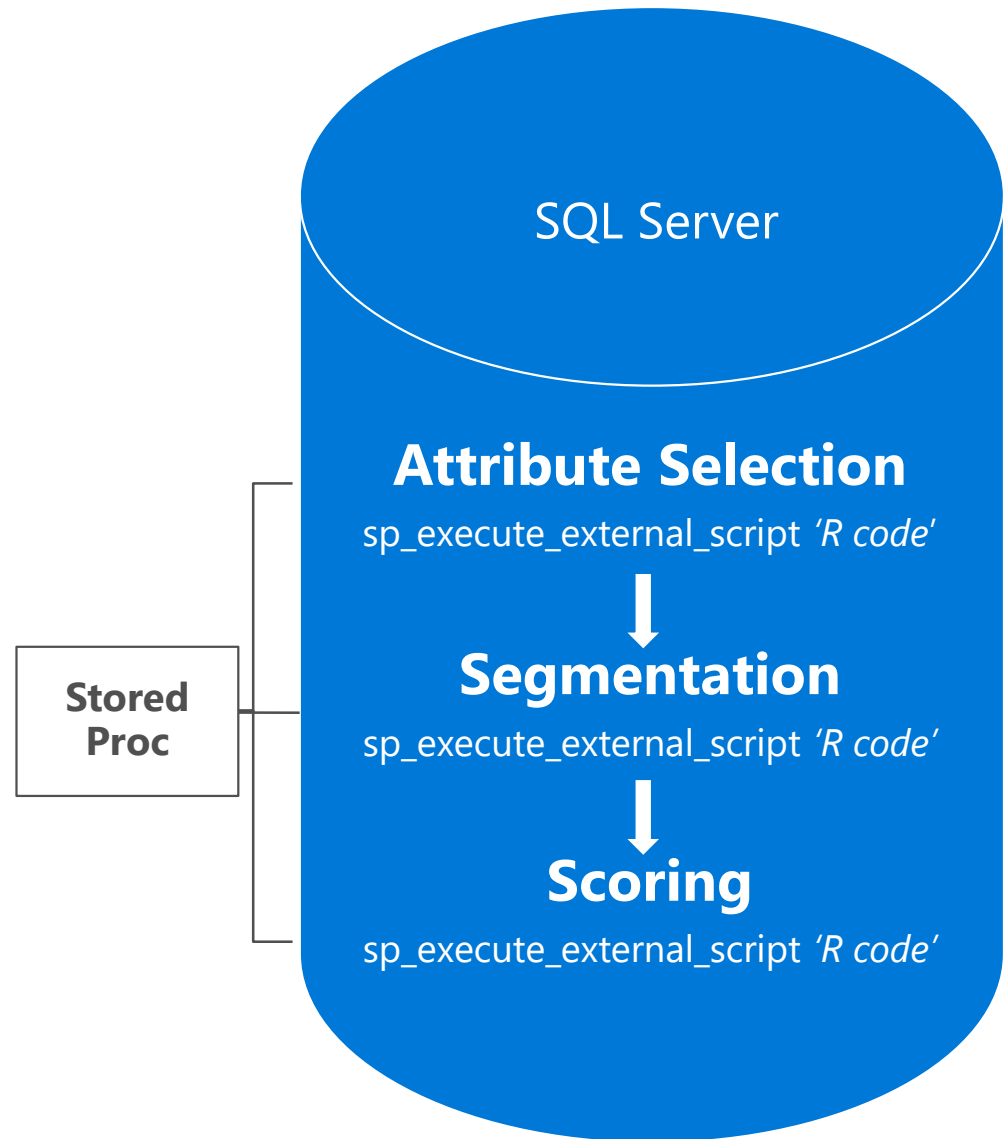
Using SQL Server R Services to improve Segmentation process performance
(Excel/C++/Python tools as baseline)



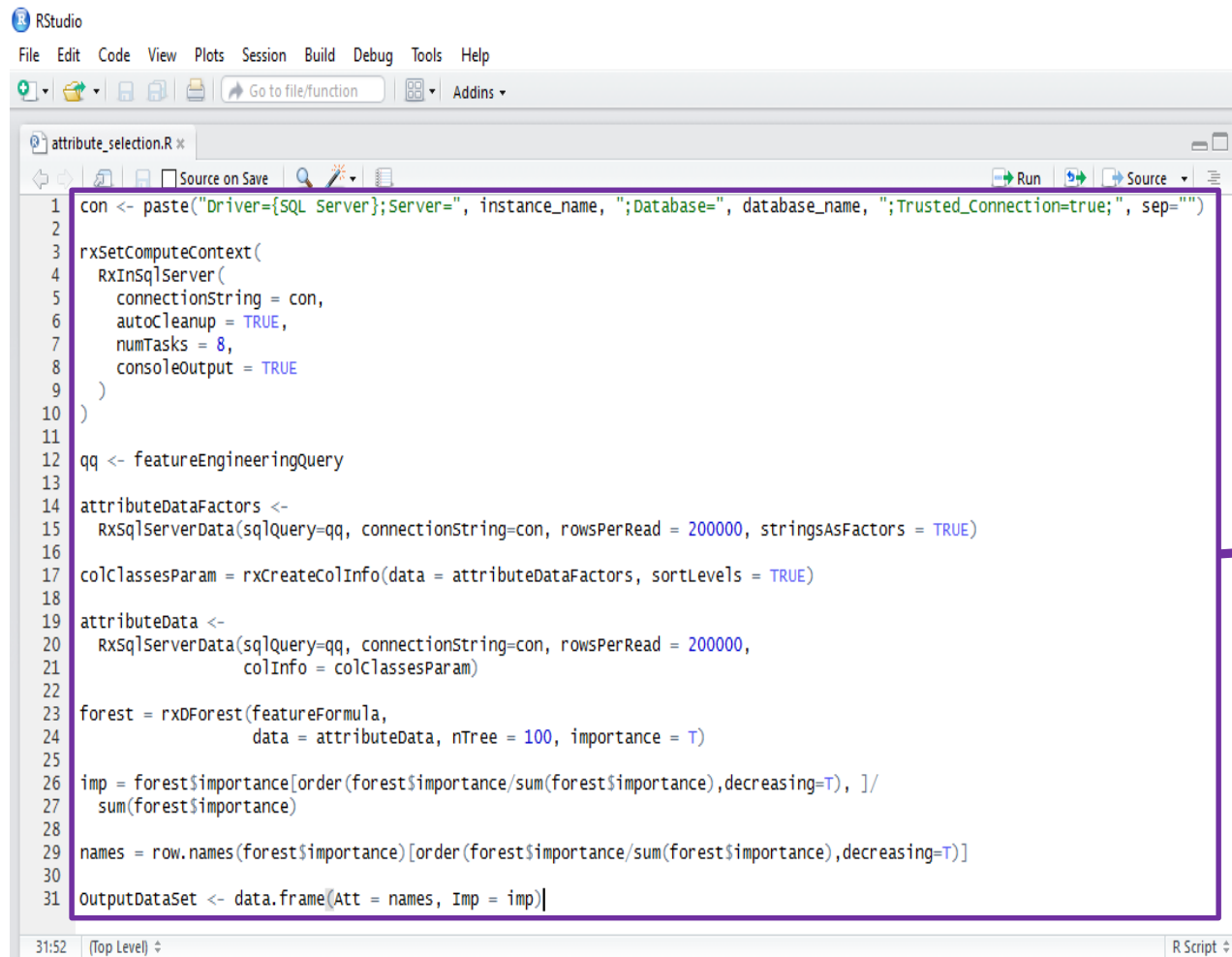
Before



After



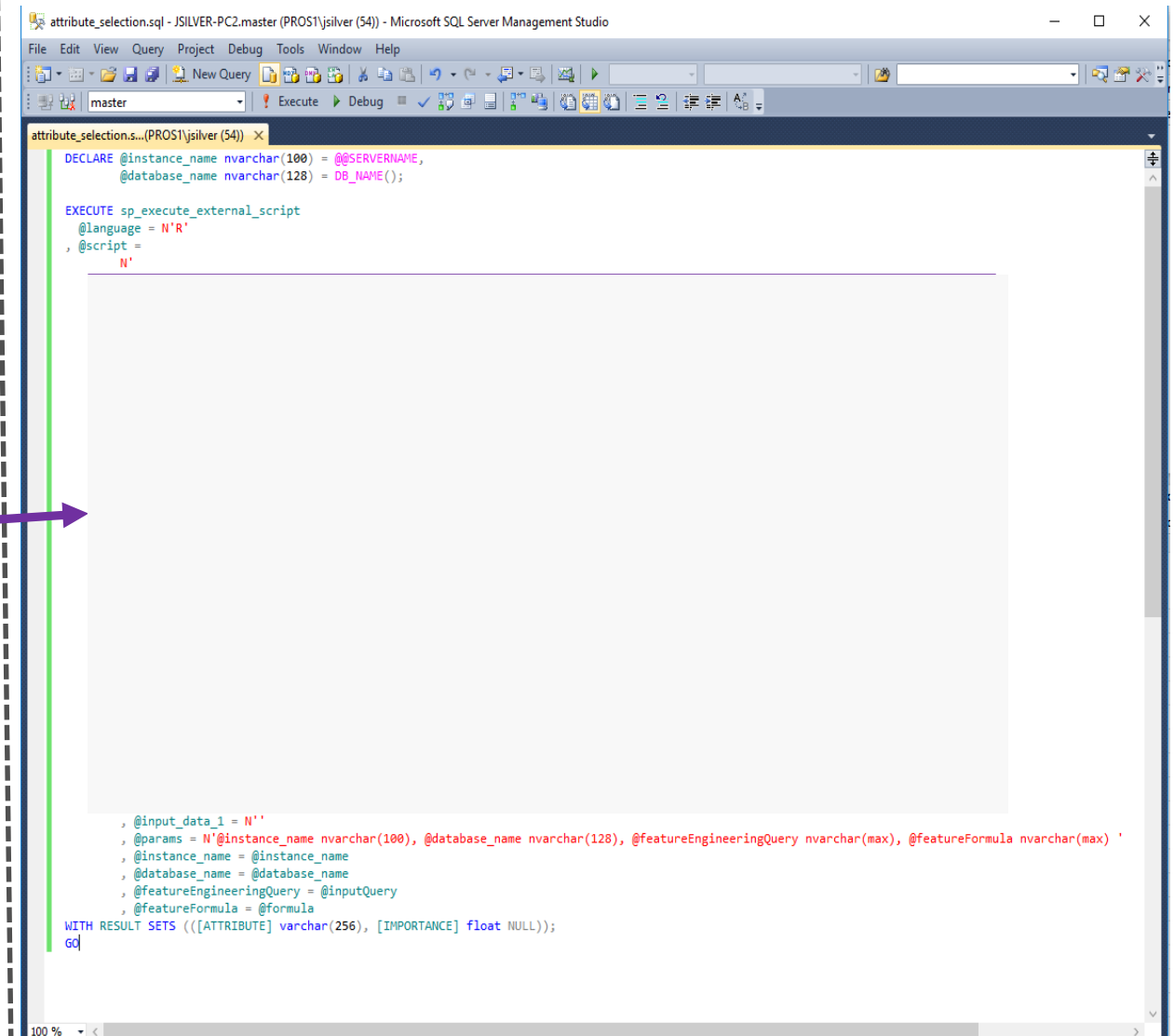
Before



The screenshot shows the RStudio interface with a script editor containing R code. A purple box highlights the first 31 lines of code, which define a connection string, set the compute context, read data from a SQL Server, and train a random forest model. The code is as follows:

```
1 con <- paste("Driver={SQL Server};Server=", instance_name, ";Database=", database_name, ";Trusted_Connection=true;", sep="")
2
3 rxSetComputeContext(
4   RxSqlServer(
5     connectionString = con,
6     autoCleanup = TRUE,
7     numTasks = 8,
8     consoleOutput = TRUE
9   )
10 )
11
12 qq <- featureEngineeringQuery
13
14 attributeDataFactors <-
15   RxSqlServerData(sqlQuery=qq, connectionString=con, rowsPerRead = 200000, stringsAsFactors = TRUE)
16
17 colClassesParam = rxCreateColInfo(data = attributeDataFactors, sortLevels = TRUE)
18
19 attributeData <-
20   RxSqlServerData(sqlQuery=qq, connectionString=con, rowsPerRead = 200000,
21     colInfo = colClassesParam)
22
23 forest = rxDForest(featureFormula,
24   data = attributeData, nTree = 100, importance = T)
25
26 imp = forest$importance[order(forest$importance/sum(forest$importance),decreasing=T), ]/
27   sum(forest$importance)
28
29 names = row.names(forest$importance)[order(forest$importance/sum(forest$importance),decreasing=T)]
30
31 OutputDataSet <- data.frame(Att = names, Imp = imp)
```

After



The screenshot shows the Microsoft SQL Server Management Studio interface. A SQL query is entered in the query editor, which declares variables for instance name and database name, and then executes an external script (the R script) using the `sp_execute_external_script` function. The query is as follows:

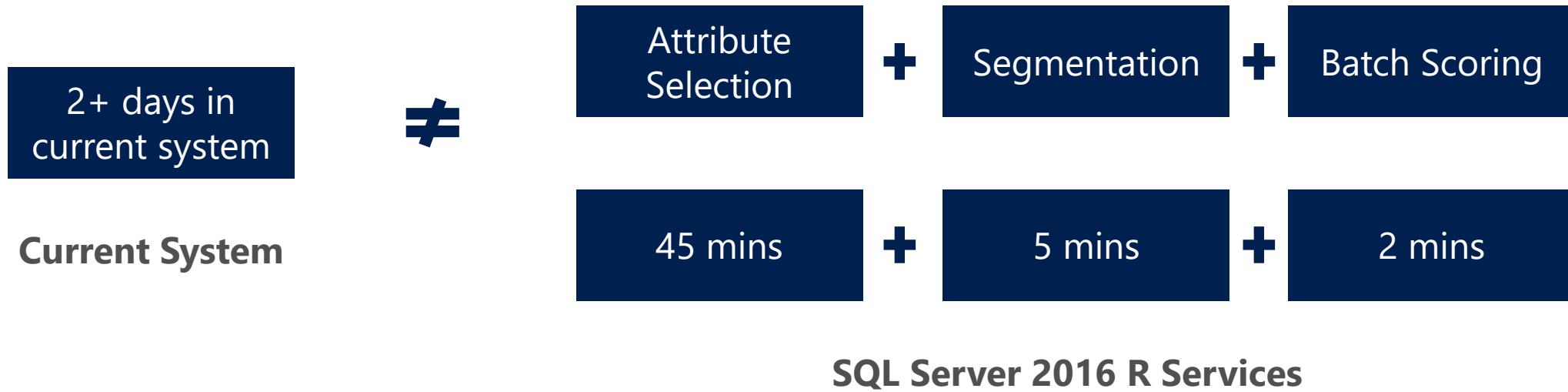
```
DECLARE @instance_name nvarchar(100) = @@SERVERNAME,
        @database_name nvarchar(128) = DB_NAME();

EXECUTE sp_execute_external_script
  @language = N'R'
  , @script =
    N'
    '

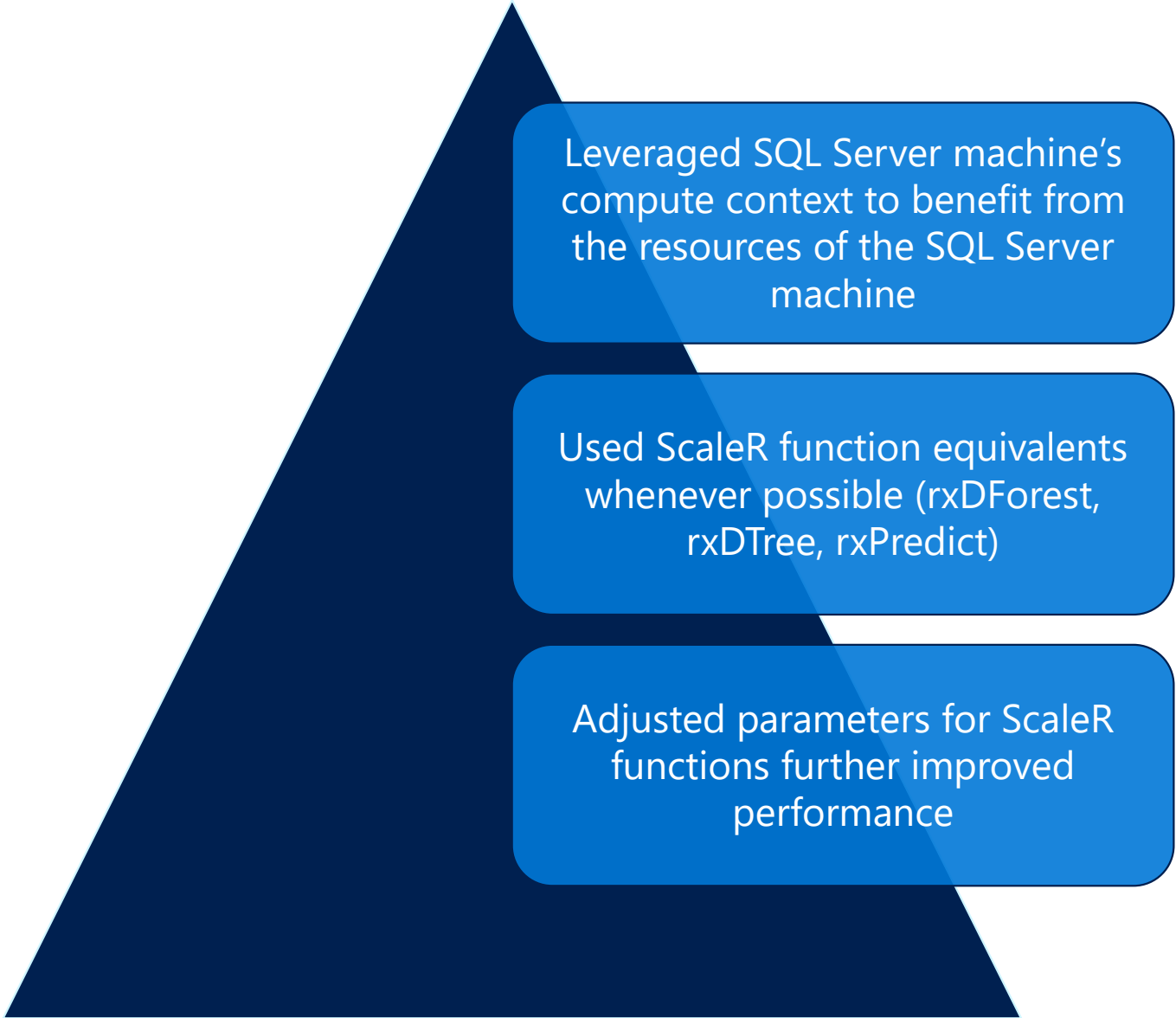
    , @input_data_1 = N''
    , @params = N'@instance_name nvarchar(100), @database_name nvarchar(128), @featureEngineeringQuery nvarchar(max), @featureFormula nvarchar(max) '
    , @instance_name = @instance_name
    , @database_name = @database_name
    , @featureEngineeringQuery = @inputQuery
    , @featureFormula = @formula

WITH RESULT SETS (([ATTRIBUTE] varchar(256), [IMPORTANCE] float NULL));
GO
```

The results



How was this achieved?



Leveraged SQL Server machine's compute context to benefit from the resources of the SQL Server machine

Used ScaleR function equivalents whenever possible (rxDForest, rxDTree, rxPredict)

Adjusted parameters for ScaleR functions further improved performance

ATTOM Data Solutions & Greenfield Advisors



Company overview



- ATTOM Data Solutions is a leading provider of property data - including tax, deed, mortgage, foreclosure, environmental risk, natural hazard, health hazard, neighborhood characteristics and property characteristics – for more than 150 million U.S. properties.



- Greenfield Advisors is a real estate and business consulting firm headquartered in Seattle, Washington. They are internationally recognized in the real estate appraisal profession as the leading authorities on the analysis and valuation of property impacted by environmental factors.

Property valuation

The property valuation process needs to provide timely and accurate calculations

- Over 100 million subject residential properties

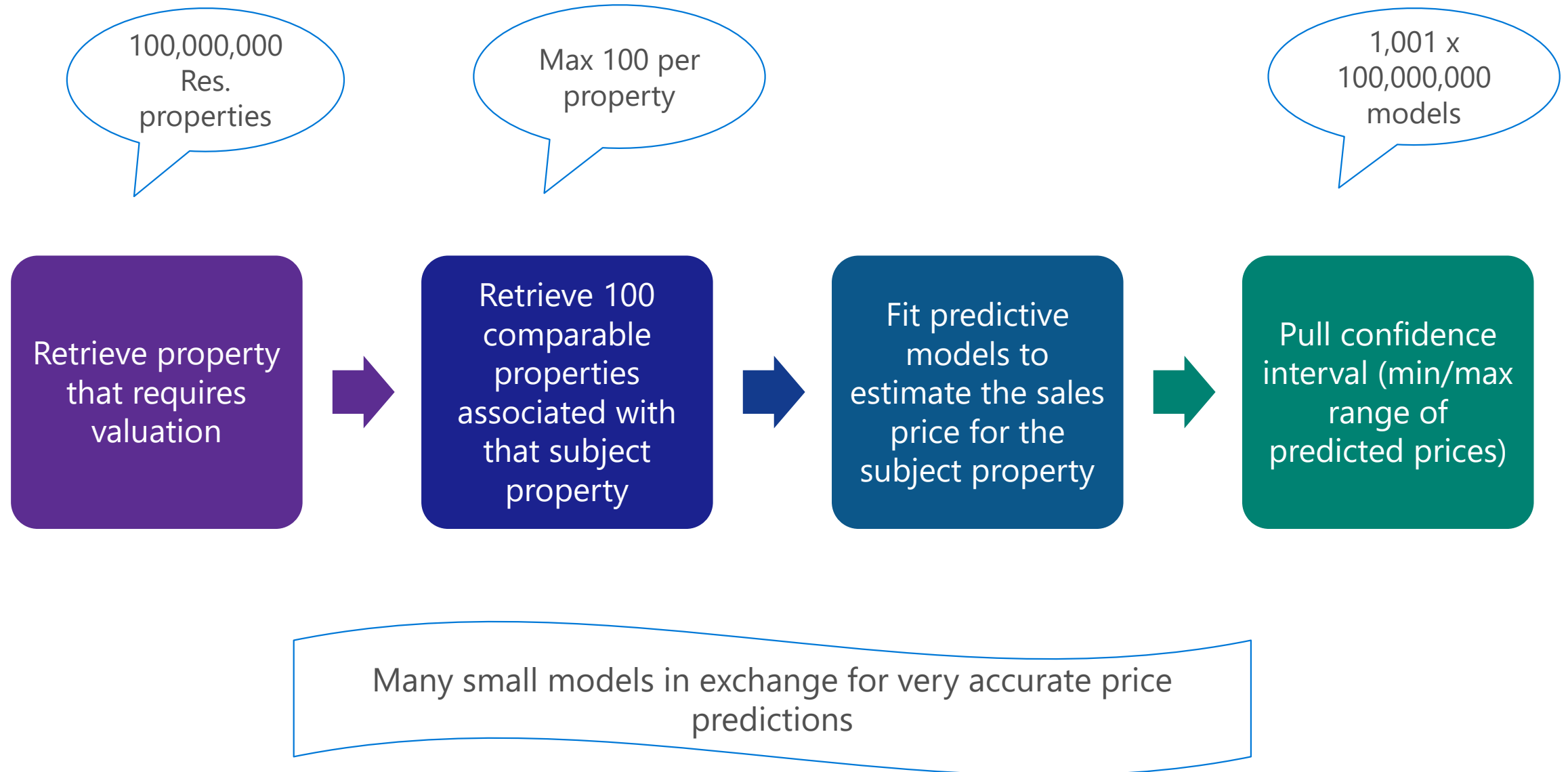
Initial implementation: limited to ~ 200 valuations per hour

- 7 R script files, all Open R, ~ 600 lines of code

Objectives

- Port the existing legacy code to SQL Server R Services
- Gain deep understanding of existing code
- Look for refactoring opportunities (R and SQL)
- Test maximum throughput capabilities

Data flow



The results

- How was this achieved?

- Putting data and computation in same place was automatic win
- Pre-calculation of the (n-1) 100 comparable companies using T-SQL and a cache
 - At the time of the exercise, the cache was in a jump to 660k
- Concurrence to execute external data sets
- DBA and expedited a
- Doing the concentrated exercise (which was not possible before), allowed for progress in a short time window

Before: limited to ~
200 valuations per
hour



After: Using SQL
Server R Services, able
to scale to ~ 720,000
valuations per hour

What about DBA use cases?

Predicting
capacity over the
holiday

Identifying
outliers in error
logs

Characterizing
workloads

Analysis of
upgrade-
regression testing

Query Store data
mining

Rich visualizations

Microsoft Advanced Analytics Landscape

Azure Machine Learning

- Fully managed cloud service that enables you to easily build, deploy, and share predictive analytics solutions

SQL Server Machine Learning Services (In-Database)

- Supports both R and Python pushed SQL Server compute-context

Microsoft R Server

- For enterprise-level R deployments on Windows and Linux servers

Microsoft Machine Learning Server

- Supports R and Python deployments on Windows servers, with expansion to other supported platforms planned for late 2017

Thank you

