# 使用R机器学习打造智能零售行业的创新体验

董乃文 Nevin Dong 资深技术专家 微软开发工具及平台事业部(DX) nevind@microsoft.com

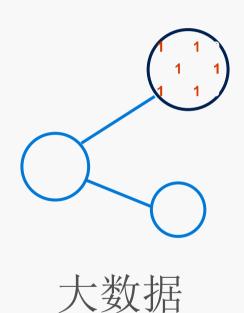


# 议题

- 零售业的挑战及数字化转型
- · 微软云高级分析服务及 R 平台
- Hadoop + Spark + R:三驾齐驱
- R + 机器学习:助力最佳零售数字化体验



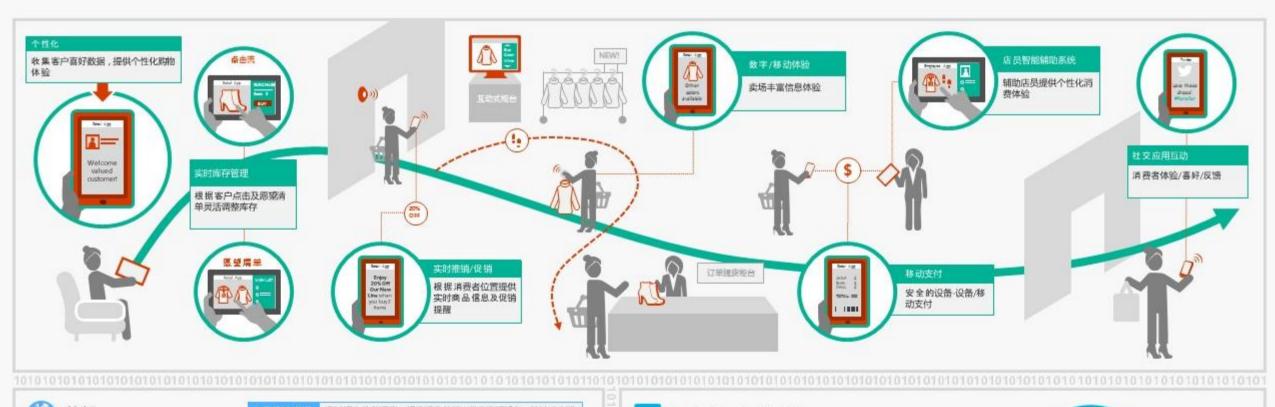
# 业务数字化转型趋势







# 消费者体验驱动的现代购物环境





# 微软云高级分析服务及R平台

# 微软云高级分析服务解决方案

### 解决方案

HDInsight/Spark

Microsoft R

Azure 机器学习

认知服务

解决方案











大数据平台

基于R的分析

云分析

分析API服务

预配置的应用及 解决方案

大规模并行计算 企业级规模,一次写入, 分布式部署

便于建模,分析 及对比

便捷的视觉、语音、语言、 知识等的分析 解决特定行业/业务/领域的应用/解决方案

数据科学家/ 数据工程师

数据科学家

数据科学家

开发者

BDM/TDM

# Microsoft R 兵器谱

- Microsoft R Open
- Microsoft R Client
- Microsoft R Server...

...for Hadoop, ...for Teradata, ...for Linux (SUSE, Red Hat/CentOS)

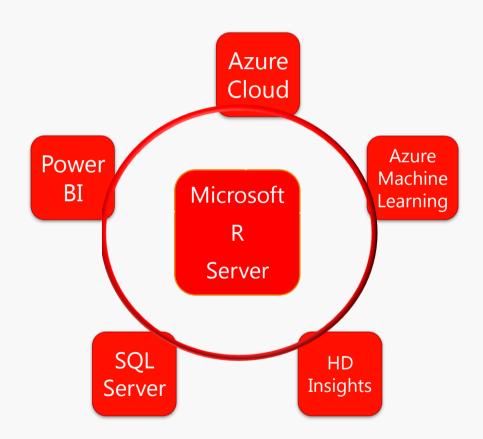
- Microsoft SQL Server 2016 R Services
- Microsoft R Open in Azure ML

# 至尊魔戒 Microsoft R: One Ring to Rule Them All

# Forbes / Tech JAN 15, 2016 © 06;14 AM 7,803 VIEWS Microsoft R: One Big Data Tool To Rule Them All? Adrian Bridgwater CONTRIBUTOR I Track enterprise software application development & data management. FOLLOW ON FORBES [72] FULL BID > Opinions expressed by Forbes.

Microsoft R Open'—a product name almost worth getting T-shirts printed for, were it not grammatically incorrect. Redmond's big data analytics dream builds one tool to rule them all, maybe... Image: Wikipedia

Microsoft MSFT +0.75% wants a slice of the big data analytics pie. Truth be told, it has already baked and served itself up a portion by acquiring the R-language and data crunching specialist Revolution Analytics, a purchase it completed in spring of 2015.



# R + CRAN

# (8,000+ packages)

### **CRAN Task Views**

CRAN Task Views are guides to the packages and functions useful for certain disciplines and methodologies. Many long-term R users I know have no idea they exist. As an effort to make them more widely known I thought Td jazz up the index page. Images are free to use, and got from SXC stock photo site. Visual puns are mine. Task View links go to the cran r-project org site and not a mirror



### Bayesian Inference

Applied researchers interested in Bayesian statistics are increasingly attracted to R. because of the ease of which one can code algorithms to sample...[more]



### Natural Language Processing

This CRAN task view contains a list of packages useful for natural language processing [more]



### Analysis of Spatial Data

Base R includes many functions that can b used for reading, vizualising, and analysing spatial data. The focus in this view is on "geographical" spatial [more]



### Chemometrics and Computational Physics

Chemometrics and computational physics are concerned with the analysis of data arising in chemistry and physics experiments, as well as the simulation



### Analysis of Pharmacokinetic Data

The primary goal of pharmacokinetic (PK) data analysis is to determine the relationship between the dosing regimen and the body's exposure to the drug as...[more]



### Clinical Trial Design, Monitoring, and Analysis

This task view gathers information on specific R packages for design, monitoring and analysis of data from clinical trials. It focuses on including [more]



### Official Statistics & Survey Methodology

This CRAN task view contains a list of packages that includes methods typically used in official statistics and survey methodology. Many packages provide.



### Survival Analysis

Survival analysis, also called event history analysis in social science, or reliability analysis in engineering, deals with time until occurrence of an ... [more]



### Cluster Analysis & Finite Mixture Models

This CRAN Task View contains a list of packages that can be used for finding groups in data and modelling unobserved ross-sectional heterogeneity. Many.



### Phylogenetics, Especially Comparative Methods

The history of life unfolds within a phylogenetic context. Comparative phylogenetic methods are statistical pproaches for analyzing historical...[me



### Time Series Analysis

Base R ships with a lot of functionality useful for time series, in particular in the stats package. This is complemented by many packages on CRAN, which are... [more]



### Probability Distributions

For most of the classical distributions, base R provides probability distribution functions (p), density functions (d), quantile functions (q), and ... [more]



### Multivariate Statistics

Base R contains most of the functionality for classical multivariate analysis. somewhere. There are a large number of packages on CRAN which extend this.



### Robust Statistical Methods

Robust (or "resistant") methods for statistics modelling have been available in S from the start, in R in package stats (e.g., median(), mean(\*, trim = .),...[more]



### Computational Econometrics

Base R ships with a lot of functionality useful for computational econometrics, in particular in the stats package. This functionality is complemented by many [more]



### Optimization and

optimization problems. Although every

regression model in statistics...[more]

Mathematical Programming This CRAN task view contains a list of packages which offer facilities for solving



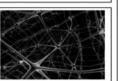
### Statistics for the Social

Social scientists use a wide range of statistical methods. To make the burden carried by this task view lighter, I have suppressed detail in some areas that...



### Analysis of Ecological and **Environmental Data**

This Task View contains information about using R to analyse ecological and environmental data [more]



### Machine Learning & Statistical Learning

Several add-on packages implement ideas and methods developed at the borderline between computer science and statistics this field of research is usually...[more]



### Sciences



### (DoE) & Analysis of **Experimental Data**

This task many collects information on P. packages for experimental design and analysis of data from experiments. Please feel free to suggest enhancements,...[more]

Graphic Displays &

Graphic Devices &

Visualization

**Dynamic Graphics &** 

R is rich with facilities for creating and

developing interesting graphics. Base R.

including coplots, mosaic...[more]

gRaphical Models in R

graph that represents independencies

among random variables by a graph in

which each node is a random variable,

and [more]

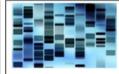
Wikipedia defines a graphical model as a

contains functionality for many plot types



### **Empirical Finance**

This CRAN Task View contains a list of packages useful for empirical work in Finance, grouped by topic....[more]



### Statistical Genetics

Great advances have been made in the field of genetic analysis over the last years. The availability of millions of single nucleotide polymorphisms (SNPs)...[more]



### Medical Image Analysis

This task view is for input, output, and analysis of medical imaging files....[more] packages, grouped by topic, that are useful



Parallel Computing with R

This CRAN task view contains a list of

for high-performance computing (HPC)

with R. In this context, we are [more]

### Reproducible Research

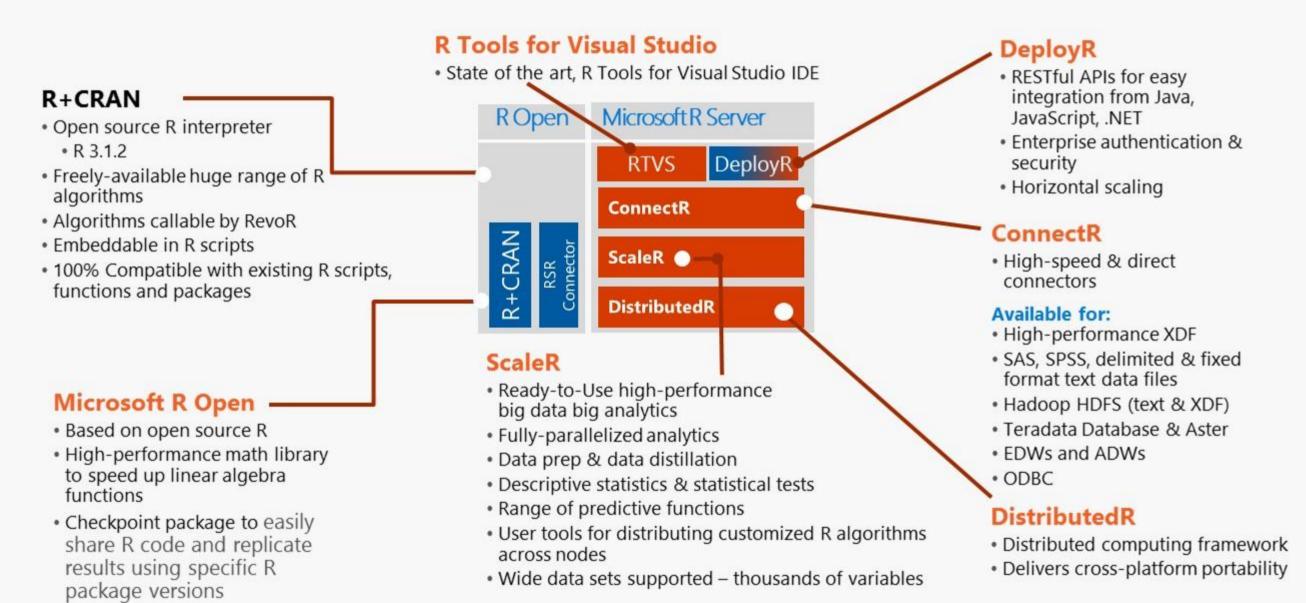
The goal of reproducible research is to tie specific instructions to data analysis and experimental data so that scholarship can e recreated, better [more]



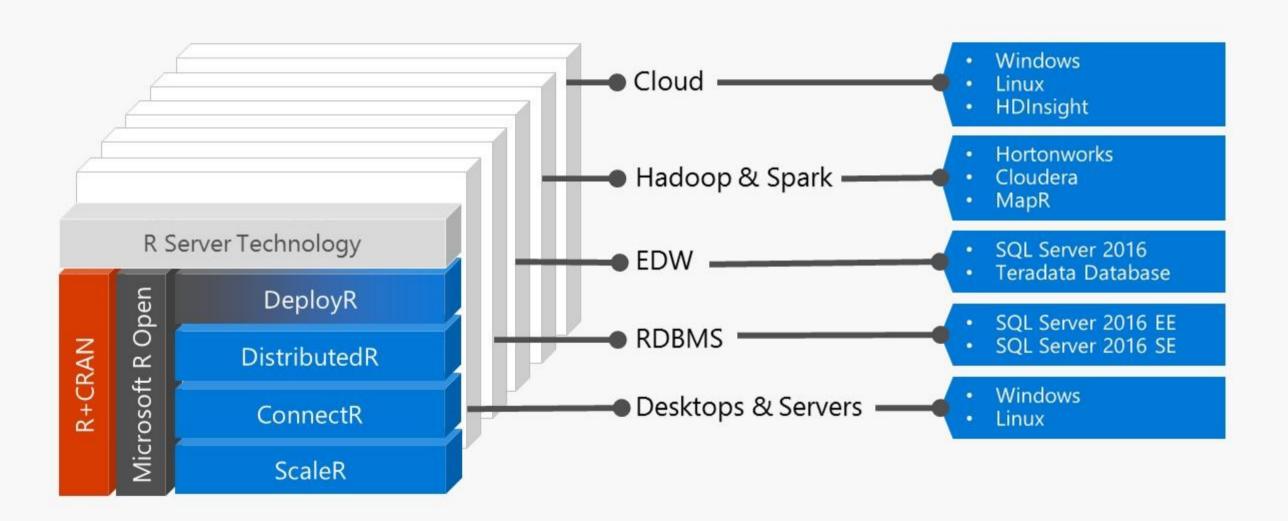
### Psychometric Models and Methods

Psychometrics is concerned with the design and analysis of research and the measurement of human characteristics Psychometricians have also worked...

# Microsoft R Server



# Microsoft R Server

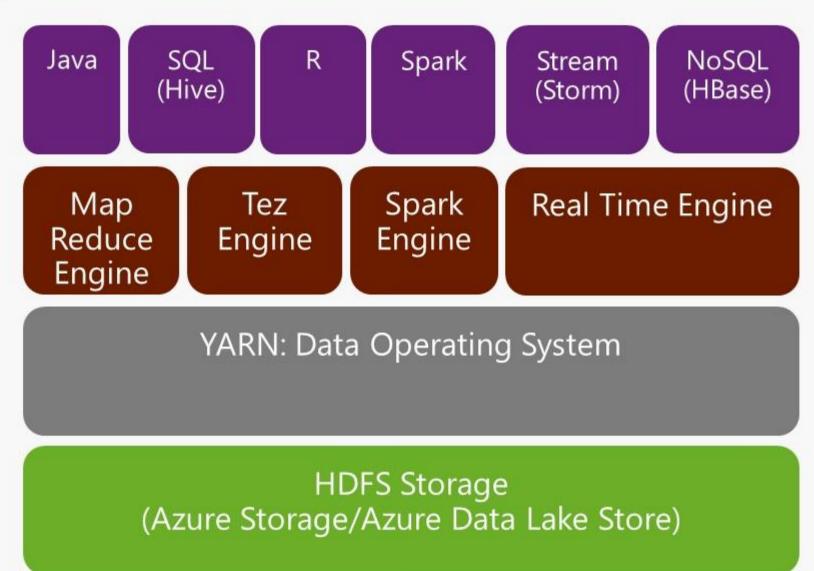


# Hadoop + Spark + R: 三驾齐驱

# Microsoft Hadoop 技术栈

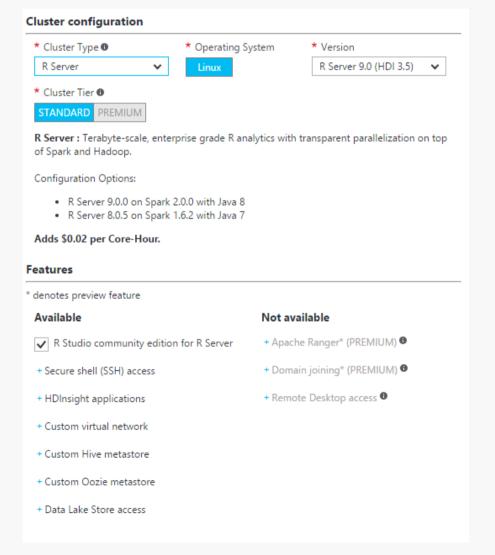
Hadoop Distributions running in Azure VMs **Azure HDInsight** Spark R Server Interactive HBase Storm cloudera MAPR. Streaming Map reduce, Machine Real Time Interactive Learning Batch Local (HDFS) or Cloud (Azure Blob/Azure Data Lake Store)

# HDInsight 开放的技术架构

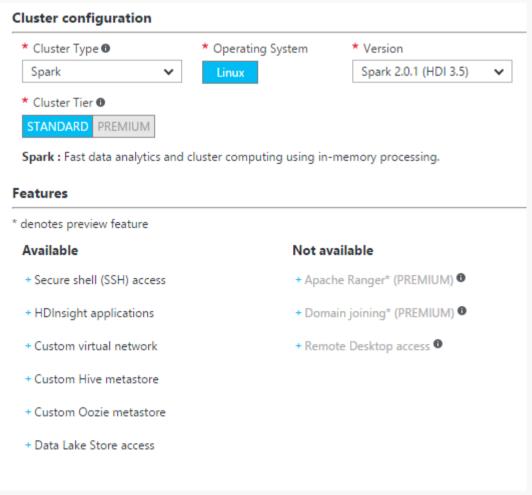


# 基于HDInsight提供高伸缩性的机器学习平台

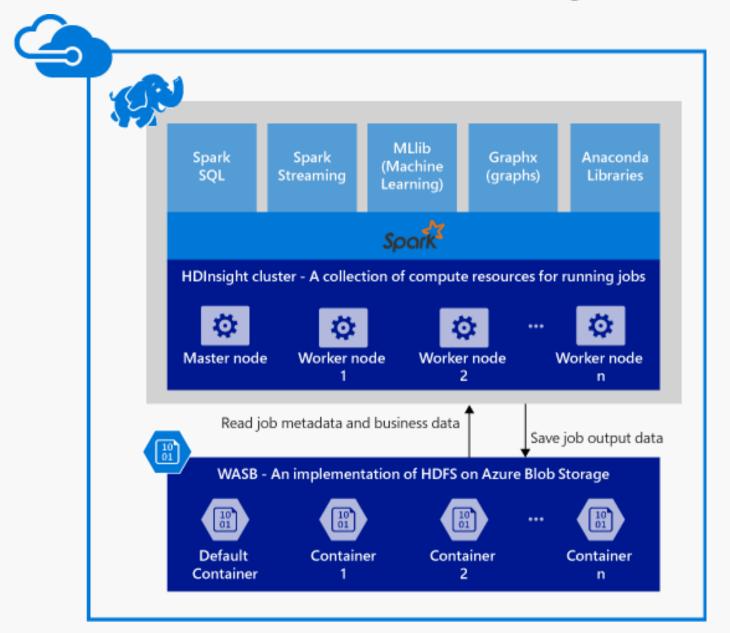




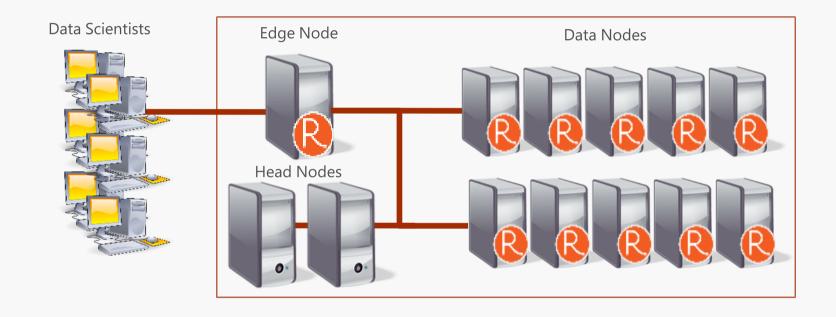




# Spark clusters in Azure HDInsight

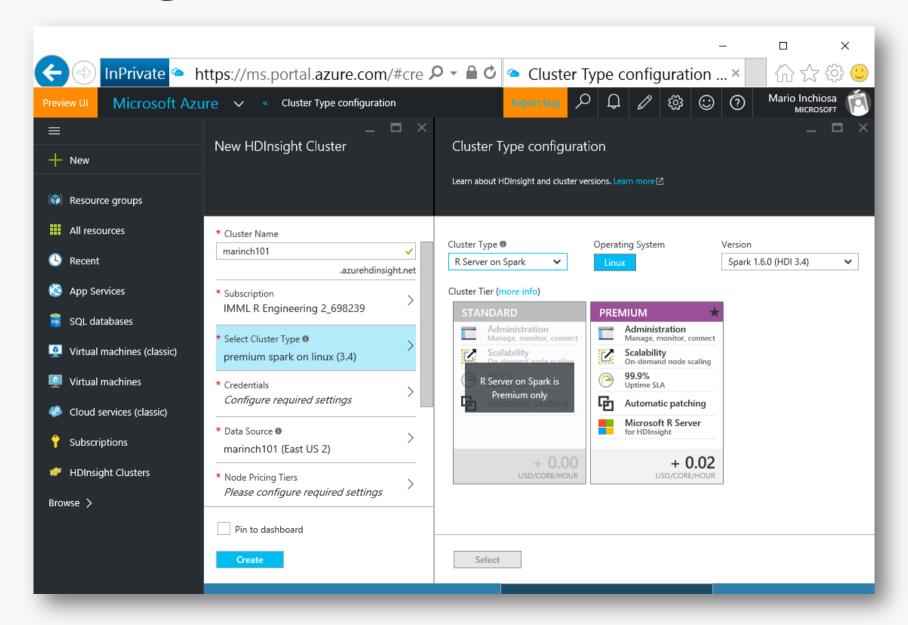


# R Server on HDInsight

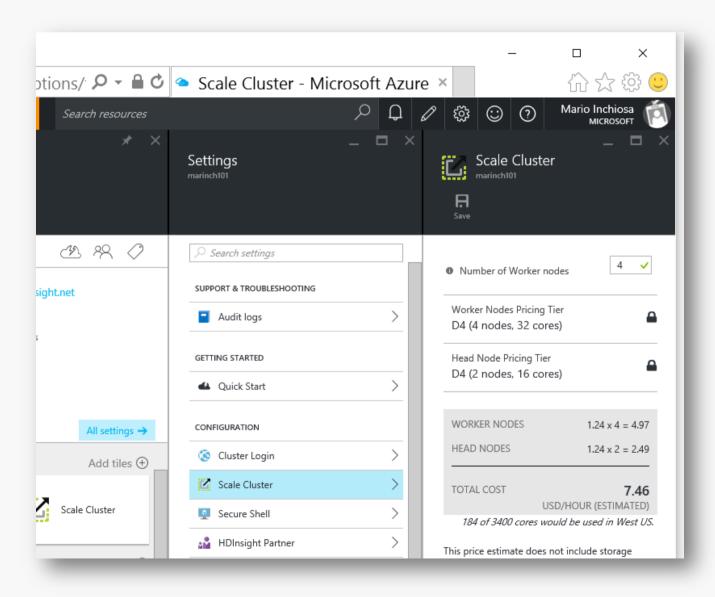




# 创建 HDInsight cluster - R Server 类型



# 集群扩展



# 开源R

### R Server

### Switch functions

```
mydata <- RxTextData(<u>"/data/binary.csv"</u>, fileSystem = hdfsFS)
mylogit <- rxLogit(admit ~ gre + gpa + rank, data = mydata)</pre>
```

# R Server 通过 Spark 实现并行计算

```
切换计算上下文
```

```
rxSetComputeContext( RxSpark(...) )
mydata <- RxTextData(<u>"/data/binary.csv"</u>, fileSystem = hdfsFS)
mylogit <- rxLogit(admit ~ gre + gpa + rank, data = mydata)</pre>
```

# R Server 通过MapReduce 实现并行计算

切换计算上下文

```
rxSetComputeContext( RxHadoopMR(...) )
mydata <- RxTextData(<u>"/data/binary.csv"</u>, fileSystem = hdfsFS)
mylogit <- rxLogit(admit ~ gre + gpa + rank, data = mydata)</pre>
```

# 在R Server 中使用 SparkR

# 在 R Server 中训练及评估模型

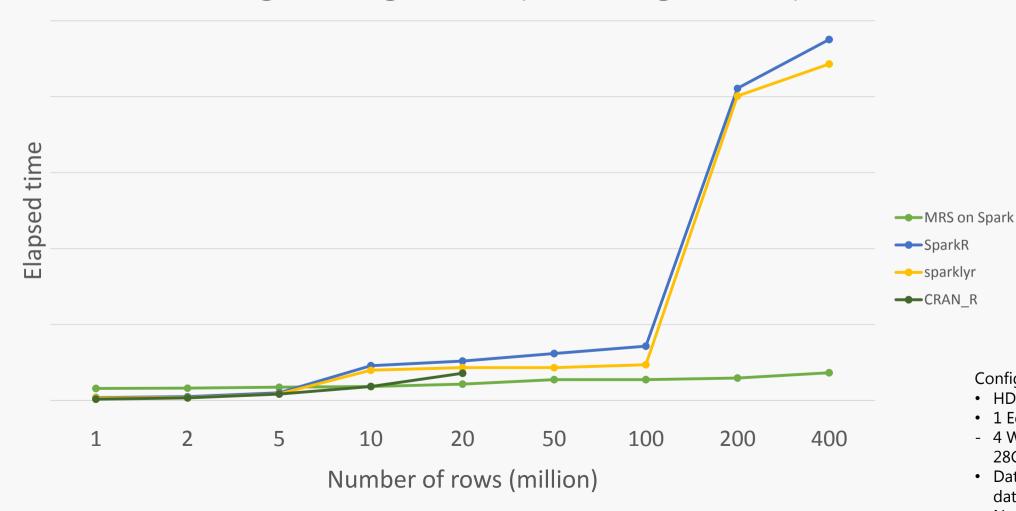
```
# Train and Test a Decision Tree model
# Train using the scalable rxDTree function
dTreeModel <- rxDTree(formula, data = trainDS,</pre>
                 maxDepth = 6, pruneCp = "auto")
# Test using the scalable rxPredict function
rxPredict(dTreeModel, data = testDS, outData = treePredict,
        extraVarsToWrite = c("ArrDel15"), overwrite = TRUE)
```

# 发布 Web Service

```
# Publish the scoring function as a web service
library(AzureML)
workspace <- workspace(config = "azureml-settings.json")</pre>
endpoint <- publishWebService(workspace, scoringFn,</pre>
                   name="Delay Prediction Service",
                   inputSchema = exampleDF)
# Score new data via the web service
scores <- consume(endpoint, dataToBeScored)</pre>
```

# R Server on Spark – 高性能、高伸缩性

Logistic Regression (executing models)



### Configuration:

- HDI cluster size: 7 nodes
- 1 Edge Node: 8 cores, 28GB
- 4 Worker Nodes: 8 cores, 28GB
- Dataset: Duplicated Airlines data (.csv)
- Number of columns: 26

R + 机器学习:助力最佳零售数字化体验

# 千人干面:打造完美的个性化体验



### 推荐、促销、导购

满足个性化购买喜好/愿望,随时、随地提供丰富的商品信息/推荐

### 跨区间销售,实时响应

深入洞察消费者,基于社交/移动的销售协作,实时响应消费者的购买需求/意愿,线上-线下互动

### 提供个性化消费体验

消费历史追踪,跨渠道整合消费偏好,推动交叉销售/追加销售

## 解决方案:交叉销售/追加销售

### 销售商品

- 店面商品
- 库存商品
- 网店商品
- 加盟店商品

• 首发新品



- 专卖商品
- 折扣商品



### 顾客画像

- 销售记录
- 客户记录
- 社交信息
- 面部识别

VIP客户

会员客户

新客户

潜在客户

流失客户

销售 渠道







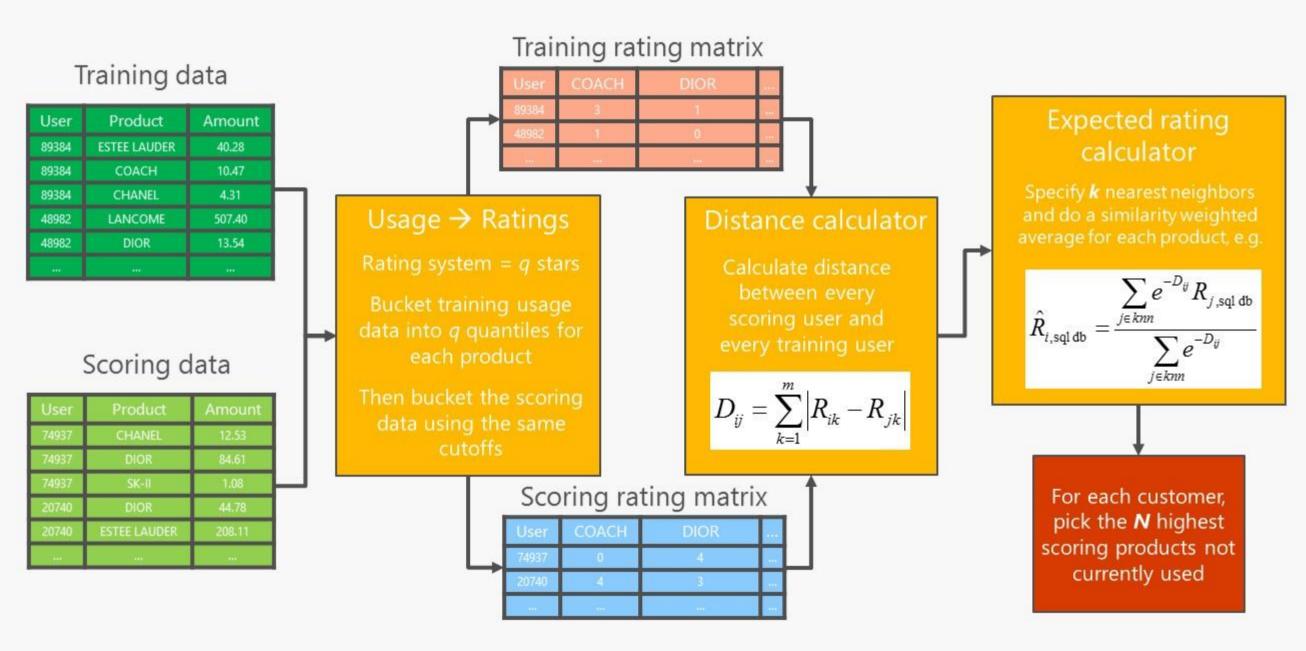




目标

提升消费者体验,通过个性化推荐,实现交叉营销、追加营销、精准营销

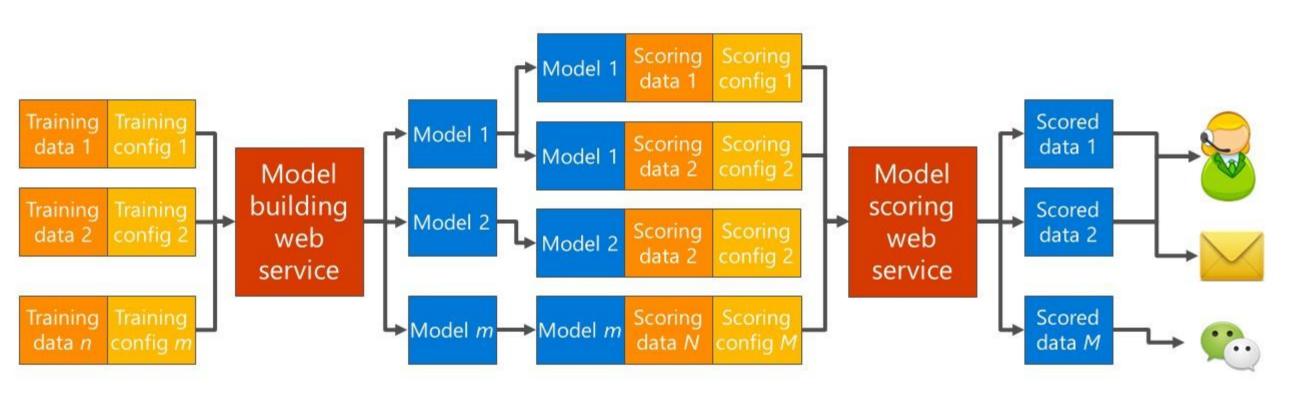
# 搭建模型 – User-User Collaborative Filtering



# 实现及运作



# 多通道的个性化推荐



# 智胜之道:深入业务洞察,智能商业模式创新



### 销量分析预测

订单销量、多层级销量、新品销售消费者满意度、客户流失情况

### 市场营销投入

渠道投入、渠道贡献 渠道拆解与重构、渠道协同

### 物流、库存分析预测

物流流转、区域配送 拆单发货、跨区发货 仓库吞吐、分仓配货补货

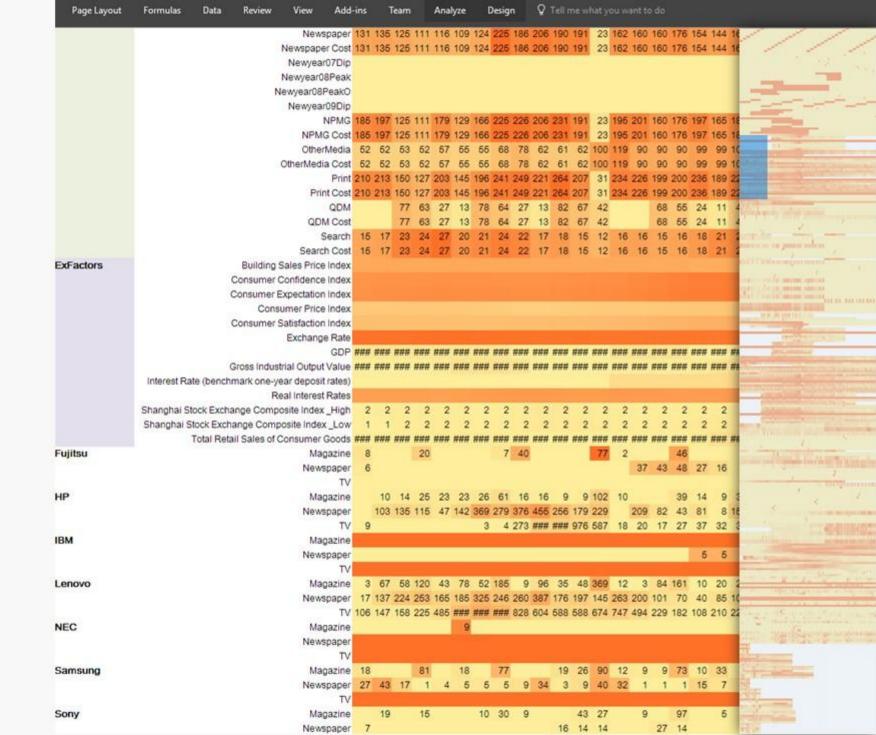
# 案例分析: 市场渠道及回报

超过不同来源与不同格式的四百多个变量才能描述客户所处的真实商业环境。

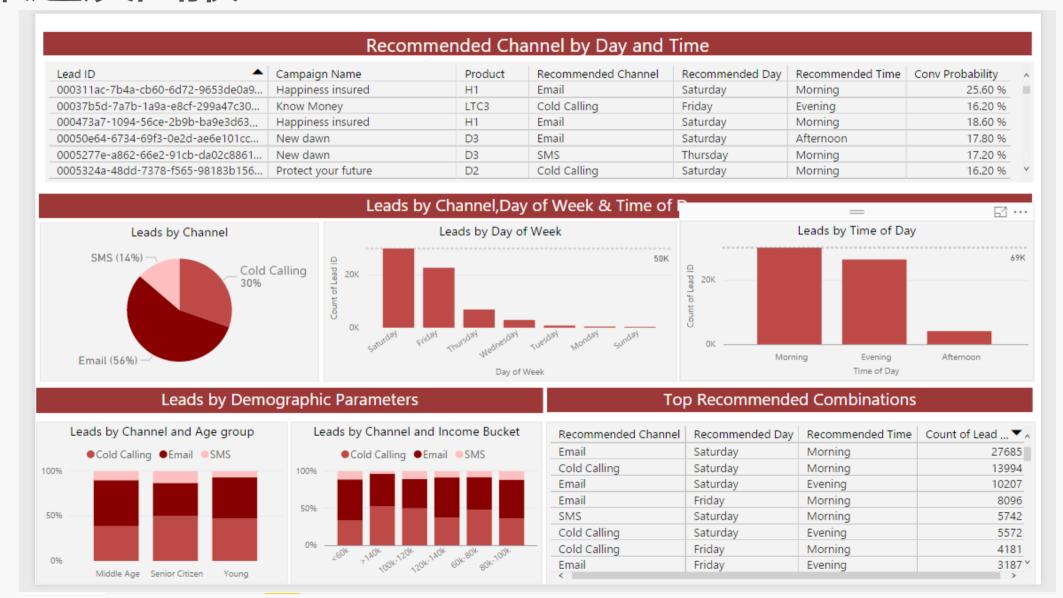
其中包含所有广告与促销投资及其所产生的影响。

以及不在客户控制范围的外部因 素,比如竞品市场活动、宏观经济 环境甚至天气。

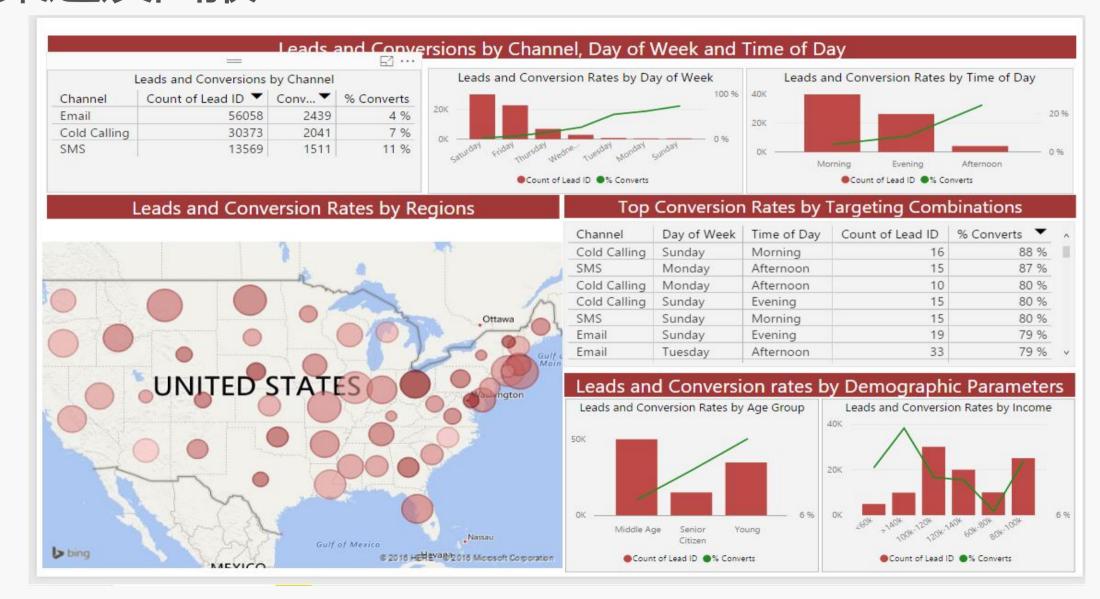
Source: HighCloud MIA Solution on Azure ML, Microsoft Ignite China 2016



# 案例分析: 市场渠道及回报



# 案例分析: 市场渠道及回报

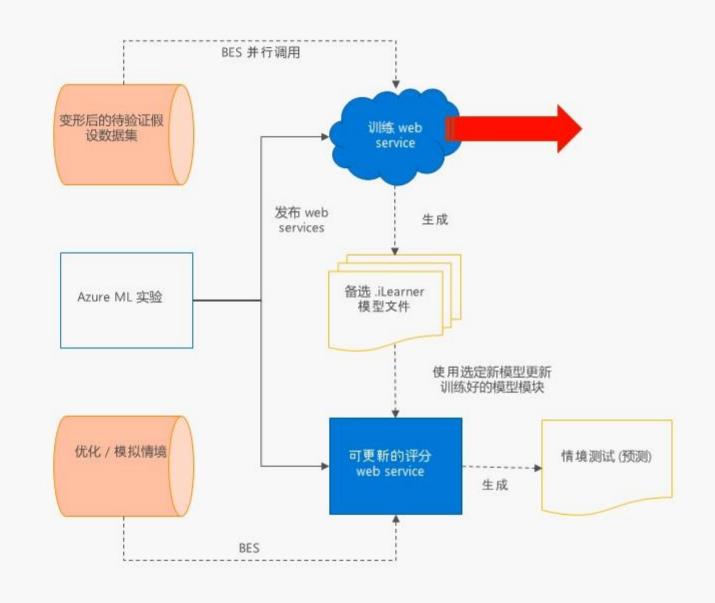


# 案例分析: 市场渠道及回报

每endpoint最高200个并发调用 每秒4000次回归运算(50ms任务) 可通过增加endpoints扩展计算能力

VS.

传统统计软件每秒10-30次回归运算



### 总结

最佳体验最关键

Hadoop + Spark + R

学R以致用,广纳Open之源

消费者:全面的线上线下个性化体验

零售商:深入业务洞察,商业模式创新

微软云高级分析服务及 R 平台

Hadoop, R and Spark are better together

Performance, Scalability are critical

使用R打造端到端的大数据应用

使用 R , 为 R 开源社区奉献

开发推动零售行业数字化转型的ML算法

# 谢谢!

