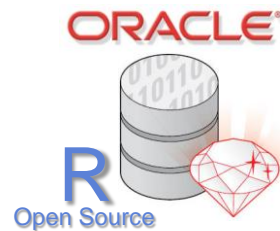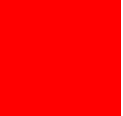**ORACLE**®

# Making Big Data Analytics accessible via the R environment

Vaishnavi Sashikanth (vaishnavi.sashikanth@oracle.com)

Vice President, Development, Database Technologies Division

**ORACLE**®

R
Open Source

The following is intended to outline our general product direction. It is intended for information purposes only, and may not be incorporated into any contract. It is not a commitment to deliver any material, code, or functionality, and should not be relied upon in making purchasing decisions.
The development, release, and timing of any features or functionality described for Oracle's products remain at the sole discretion of Oracle.

**ORACLE**

# Agenda

- What is Big Data Analytics?

- Oracle Big Data Analytics Architecture & Components

- Open Source R integration

ORACLE

# Big Data / Deep Analytics

- *Application of numerical, predictive and statistical techniques on big data*

### Financial Services

- Credit risk analysis
- Cross-LOB up-selling
- Fraud detection
- Retail banking personalization
- "Best customer" prediction & profiling

### Retail

- Real-time shopping cart recommendations
- Customer segmentation
- Customer profiling
- Market basket analysis
- Fraud detection

### Media & Entertainment

- Online ad placement
- Cable TV: option bundling
- Gaming: Targeting "right customer w/ "right product"
- Gambling: Fraud and anomaly detection

### Telecommunications

- Churn prevention
- Social network analysis
- Network monitoring
- Win-back analysis
- Fraud analysis

### Public Sector

- Healthcare Fraud prevention
- Infrastructure maintenance
- Constituent Sentiment
- Threat Identification
- Healthcare improvement

### Manufacturing

- Warranty analysis
- Quality improvement
- Product & process design and improvement

### Transportation and Logistics

- Anticipate bottlenecks
- Proactive resource planning
- Improved preventative maintenance strategies

### Utilities

- Customer loyalty management
- Fraud detection
- Product bundling
- Improved operations efficiencies

ORACLE

# Analytics driving the bottom line..

Oracle Copyright 2012, Oracle Proprietary & Confidential

ORACLE

# Analytics minimizing bad debt..

from American Express <AmericanExpress@welcome.aexp.com>

subject **Fraud Protection Alert**

reply-to American Express <alerts@service.americanexpress.com>

Internal records show a relationship between this IP address and Brazilian Organized Crime Groups associated with holding businesses hostage with malware / DDOS for ransom.

Transaction done on a French eCommerce site with the payment processed at a US gateway.

For your security, we regularly monitor accounts for possible fraud

Account 258 Bank 938

1 keywords selected    12 search results    threshold: 1.0

100% similarity    Account 7564 Bank 5144    #1

100% similarity    Account 6829 Bank 928    #2

100% similarity    Account 3 Bank 1593    #3

ORACLE

# Analytics permeating business operations..

ORACLE

# Analytics preserving high value customers

| originating_id | dialed_id | sou_sum | dialed_count | dialed_rank | originating_status | dialed_status |
|---|---|---|---|---|---|---|
| XXXXXXXXXXXXXXXXX | YYYYYYYYYYYYYYYYY | 793 | 35 | 1 | July | August |
| XXXXXXXXXXXXXXXXX | YYYYYYYYYYYYYYYYY | 360 | 30 | 2 | July | July |
| XXXXXXXXXXXXXXXXX | YYYYYYYYYYYYYYYYY | 407 | 25 | 3 | July | May |
| XXXXXXXXXXXXXXXXX | YYYYYYYYYYYYYYYYY | 167 | 14 | 4 | July | June |
| XXXXXXXXXXXXXXXXX | YYYYYYYYYYYYYYYYY | 142 | 8 | 5 | July | July |
| XXXXXXXXXXXXXXXXX | YYYYYYYYYYYYYYYYY | 124 | 6 | 6 | July | <active> |
| XXXXXXXXXXXXXXXXX | YYYYYYYYYYYYYYYYY | 268 | 4 | 7 | July | August |
| XXXXXXXXXXXXXXXXX | YYYYYYYYYYYYYYYYY | 124 | 4 | 8 | July | <active> |
| XXXXXXXXXXXXXXXXX | YYYYYYYYYYYYYYYYY | 92 | 3 | 9 | July | <active> |

ct | Retain | Retain
ct | Retain | Retain

Likdihood

Sustain | Sustain | Sustain

Lifetime Value ----------------------->Highest

Social environment effects
– Peer commentary
– Social leader influence
– Promotions to a leader to influence group

# Analytics and the Art of Winning…

**ORACLE**

# Oracle Big Data Analytics

- Focus is on the Enterprise Data Scientist who engages in Quantitative Research

- Goals
  1. Improve user efficiency by enabling focus on analysis as opposed to data access
  2. Enable deep analytics with computations occurring closer to data
  3. Allow transparent access to Enterprise compute infrastructures
  4. Shorten the path to application of cutting edge ideas into practice
  5. Enable quick transition from analysis to mass consumption of results

# Oracle Big Data Analytics



SQL

R workspace console

Transparent access

Map-reduce R

Embedded R engine

Native implementation of several statistical & data mining techniques

Applications, BI, Web Services

External Tables

HDFS Files

Database Links

Other databases

External Tables

File systems

# Oracle Big Data Analytics



SQL

R workspace console

ORACLE R ENTERPRISE

Embedded R engine

Native implementation of several statistical & data mining techniques

Transparent access

Map-reduce R

ORACLE R CONNECTOR FOR HADOOP

Applications, BI, Web Services

External Tables

HDFS Files

Database Links

External Tables

Other databases

File systems

# What is R?

- **R is an Open Source language and environment for statistical computing and graphics**
**http://www.r-project.org/**

- **Started in 1994 as an Alternative to SAS, SPSS & Other proprietary Statistical Environments**

- **The R environment**
  - R is an integrated suite of software facilities for data manipulation, calculation and graphical display

- **Around 2 million R users worldwide**
  - Widely taught in Universities
  - Many Corporate Analysts know and use R

- **Hundreds of open sources packages to enhance productivity such as:**
  - Bioinformatics with R
  - Spatial Statistics with R
  - Financial Market Analysis with R
  - Linear and Non Linear Modeling



CRAN
Mirrors
What's new?
Task Views
Search

About R
R Homepage
The R Journal

Software
R Sources
R Binaries
Packages
Other

Documentation
Manuals
FAQs
Contributed

CRAN Task Views

| | |
|---|---|
| Bayesian | Bayesian Inference |
| ChemPhys | Chemometrics and Computational Physics |
| ClinicalTrials | Clinical Trial Design, Monitoring, and Analysis |
| Cluster | Cluster Analysis & Finite Mixture Models |
| Distributions | Probability Distributions |
| Econometrics | Computational Econometrics |
| Environmetrics | Analysis of Ecological and Environmental Data |
| ExperimentalDesign | Design of Experiments (DoE) & Analysis of Experimental Data |
| Finance | Empirical Finance |
| Genetics | Statistical Genetics |
| Graphics | Graphic Displays & Dynamic Graphics & Graphic Devices & Visualization |
| gR | gRaphical Models in R |
| HighPerformanceComputing | High-Performance and Parallel Computing with R |
| MachineLearning | Machine Learning & Statistical Learning |
| MedicalImaging | Medical Image Analysis |
| Multivariate | Multivariate Statistics |
| NaturalLanguageProcessing | Natural Language Processing |
| OfficialStatistics | Official Statistics & Survey Methodology |
| Optimization | Optimization and Mathematical Programming |
| Pharmacokinetics | Analysis of Pharmacokinetic Data |
| Phylogenetics | Phylogenetics, Especially Comparative Methods |
| Psychometrics | Psychometric Models and Methods |
| ReproducibleResearch | Reproducible Research |
| Robust | Robust Statistical Methods |
| SocialSciences | Statistics for the Social Sciences |
| Spatial | Analysis of Spatial Data |
| Survival | Survival Analysis |
| TimeSeries | Time Series Analysis |

ORACLE

# Why statisticians/data analysts use R

R is a statistics language similar to Base SAS or SPSS statistics

R environment is…

- Powerful - Extensive numerical techniques

- Extensible – 1000s of CRAN packages

- Exhaustive visualization

- Ease of installation and use
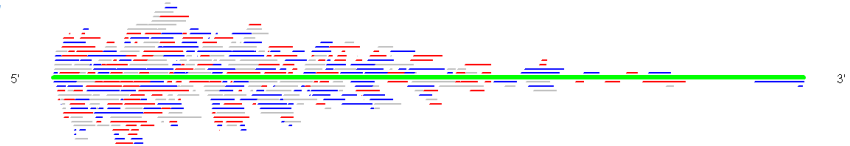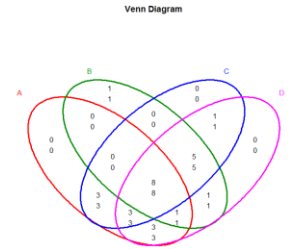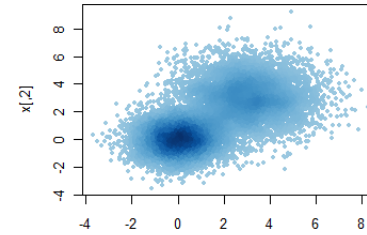
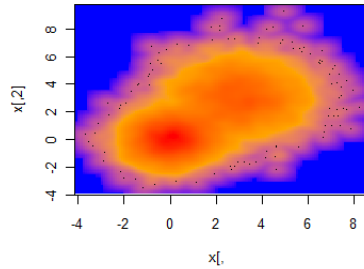- Is becoming the language of research

- ***Free***

Statisticians may not be

- SQL literate

- Familiar with DBA tasks

ORACLE

# Graph examples…

ORACLE

# R is the language of research

## Random forest

From Wikipedia, the free encyclopedia

*This article is about the machine learning technique. For other kinds of random tree, see Random tree (disambiguation).*

This article **is written like a personal reflection or essay** rather than an encyclopedic description of the subject. Please help improve it by rewriting it in an encyclopedic style. *(February 2012)*

**Random forest** (or **random forests**) is an ensemble classifier that consists of many decision trees and outputs the class that is the mode of the classes output by individual trees. The algorithm for inducing a random forest was developed by Leo Breiman[1] and Adele Cutler, and "Random Forests" is their trademark. The term came from **random decision forests** that was first proposed by Tin Kam Ho of Bell Labs in 1995. The method combines Breiman's "bagging" idea and the random selection of features, introduced independently by Ho[2][3] and Amit and Geman[4] in order to construct a collection of decision trees with controlled variation.

The selection of a random subset of features is an example of the random subspace method, which, in Ho's formulation, is a way to implement stochastic discrimination[5] proposed by Eugene Kleinberg.

**Contents** [hide]

The Comprehensive R Archive
cran.r-project.org

randomForest: Breiman and Cutler's random forests for classification and regression

Classification and regression based on a forest of trees using random inputs.

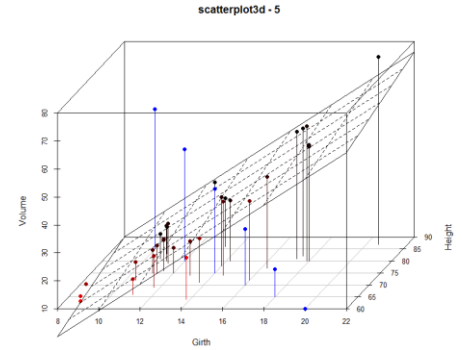| | |
|---|---|
| Version: | 4.6-6 |
| Depends: | R (≥ 2.5.0), stats |
| Suggests: | RColorBrewer, MASS |
| Published: | 2012-01-06 |
| Author: | Fortran original by Leo Breiman and Adele Cutler, R port by Andy Liaw and Matthew Wiener. |
| Maintainer: | Andy Liaw <andy_liaw at merck.com> |
| License: | GPL (≥ 2) |
| URL: | http://stat-www.berkeley.edu/users/breiman/RandomForests |
| Citation: | randomForest citation info |
| In views: | Environmetrics, MachineLearning |
| CRAN checks: | randomForest results |

Downloads:

| | |
|---|---|
| Package source: | randomForest_4.6-6.tar.gz |
| MacOS X binary: | randomForest_4.6-6.tgz |
| Windows binary: | randomForest_4.6-6.zip |
| Reference manual: | randomForest.pdf |
| News/ChangeLog: | NEWS |
| Old sources: | randomForest archive |

CRAN
Mirrors
What's new?
Task Views
Search

About R
R Homepage
The R Journal

Software
R Sources
R Binaries
Packages
Other

Documentation
Manuals
FAQs
Contributed

ORACLE

# Limitations of R

- R is a client and server bundled together as 1 executable - like Excel
    - Single user tool
    - Not multi-threaded
    - Cannot leverage CPU capacity even on a user's laptop/desktop
- R requires data it operates on to be first loaded into memory
    - Loading data may not be a limitation given RAM available on laptops/desktops
    - R's *call by value* semantics means as data flows into functions, for each function invocation, a complete copy of the data is made
    - As a result you quickly run into memory limits

# R integration

R workspace console

**Density for Arrival Delay at SFO**

Transparent access via function push-down

ORACLE®

Oracle statistics engine,
Oracle data mining

Embedded R

**hadoop**

BI, Workflows, Web Services

R is integrated into SQL

No changes to the user experience

Lights out execution

Embed in operational systems

| Development | Production | Consumption |
| --- | --- | --- |

# Oracle's Approach – Comprehensive Enterprise-level Big Data Analytics based on R environment

1. **Oracle's Distribution and Support of Open Source R**

   - Only redistribution with comprehensive platform support – Linux, Solaris, AIX
     - Enhanced performance with Intel MKL, AMD ACML OR SUN perf libraries for x86 hardware
   - Certification of select CRAN packages
   - Distributed via public-yum.oracle.com, pkg.oracle.com

2. **Oracle R Enterprise**

   - Embedded component of the RDBMS
   - Eliminates R's memory constraint by enabling R to work transparently on database resident data
   - Brings R users closer to Oracle Database by transparently leveraging in-database analytics via R
   - Enables integration of R scripts into enterprise production applications and BI dashboards
   - Fully leverages the latest R algorithms and models contributed to R's CRAN

3. **Oracle R Connector For Hadoop**

   - Interactive R interface to HDFS data and Hadoop infrastructure
   - Only available solution to combine database, HDFS and local file system data into 1 hadoop R computation

**ORACLE**

# Licensing

1. **Oracle's Distribution and Support of Open Source R**

   - Free

2. **Oracle R Enterprise**

   - Available as part of Oracle Advanced Analytics Option to Oracle Database
   - Oracle Advanced Analytics Option = Oracle Data Mining + Oracle R Enterprise
   - Oracle Data Mining algorithms are available via Oracle R Enterprise interface, SQL and GUI

3. **Oracle R Connector For Hadoop**

   - Available as part of Oracle Big Data Connectors software suite

**ORACLE**

# 1. Collaborative Execution Model

1

2

3

R Engine

Other R packages

Oracle R Enterprise packages

SQL

**Oracle Database**

User tables

R

Results

Results

R Engine

Other R packages

Oracle R Enterprise packages

## User R Engine on desktop

- R-SQL Transparency Framework intercepts R functions for scalable in-database execution

- Interactive display of graphical results and flow control as in standard R

- Submit entire R scripts for execution by Oracle Database

Post processing of results

## Database Compute Engine

- Scale to large datasets

- Leverage database SQL parallelism

- Leverage new and existing in-database statistical and data mining capabilities

Collaborative execution with in-database R engine

## R Engine(s) spawned by Oracle DB

- Database can spawn multiple R engines for database-managed parallelism

- Efficient parallel data transfer to spawned R engines to emulate map-reduce style algorithms and applications

- Enables "lights-out" execution of R scripts

Analytic techniques not available in-database

**ORACLE**

# 2. Deferred execution

```
#Filter rows that correspond to American Airlines
#Flights
ONTIME <- ONTIME[ONTIME$UNIQUECARRIER=='AA']
```

```
select *
from ONTIME
where uniquecarrier = 'AA'
```

```
#Calculate median arrival delay for all flights grouped
#by destination
aggdata <- aggregate(ONTIME$ARRDELAY,
                     by = list(ONTIME$DEST),
                     FUN = median)
```

```
select dest, median(arrdelay)
from ONTIME
where uniquecarrier = 'AA'
group by dest
```

EXECUTE!
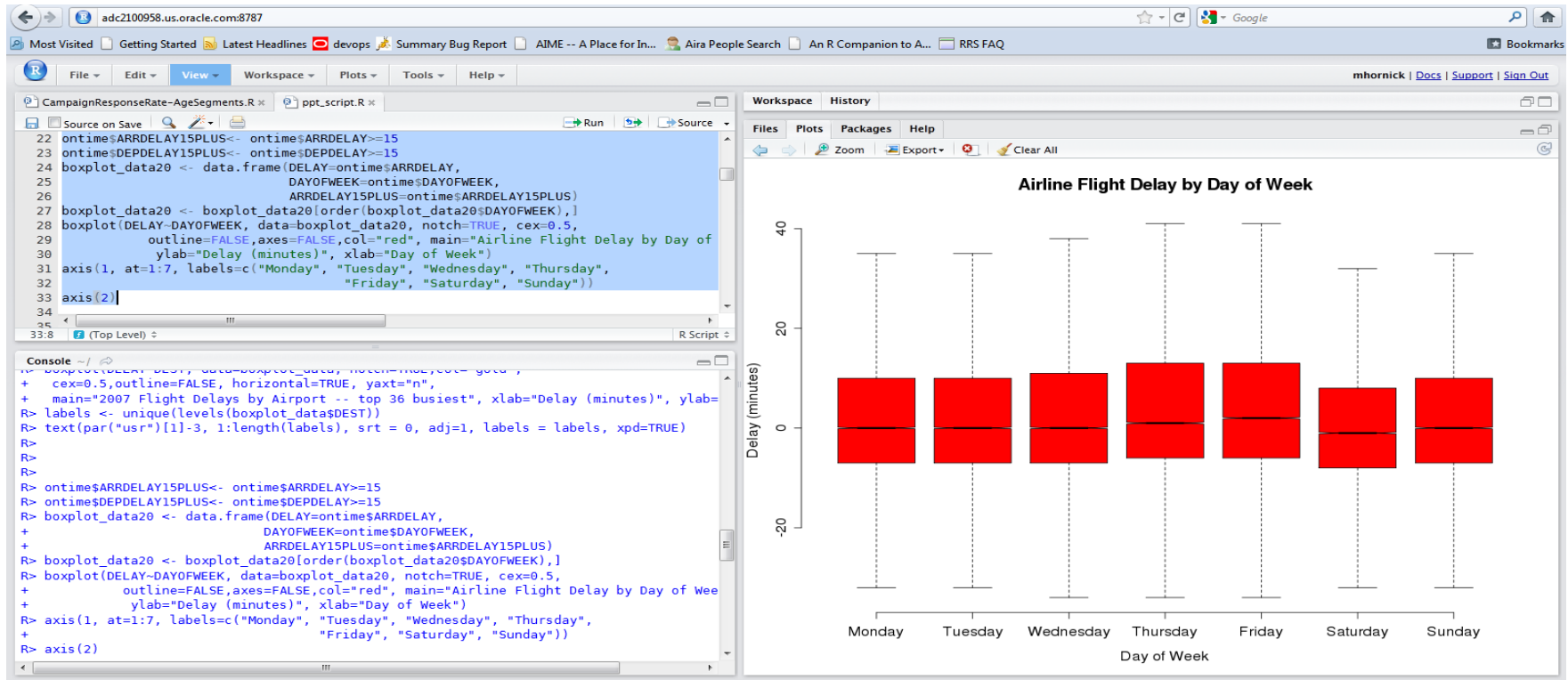
```
plot(aggdata)
```

**ORACLE**

# 3. Collaborative Visualization

Push computations into SQL and render using R

ORACLE

# 4. R is integrated into SQL

```
select * from table(rqTableEval(
    cursor(select * from fish),
    NULL,
    'select t.*, 1 rowsum from fish t',
    'function(x, param) {
      dat <- data.frame(x, stringsAsFactors=F)
      cbind(dat, ROWSUM = apply(dat,1,sum))
     }'));

select * from table(rqRowEval(
    cursor(select * from fish),
    NULL,
    'select t.*, 1 rowsum from fish t',
    1,
    'function(x, param) {
      dat <- data.frame(x, stringsAsFactors=F)
      cbind(dat, ROWSUM = apply(dat,1,sum)+10)
     }'));
```
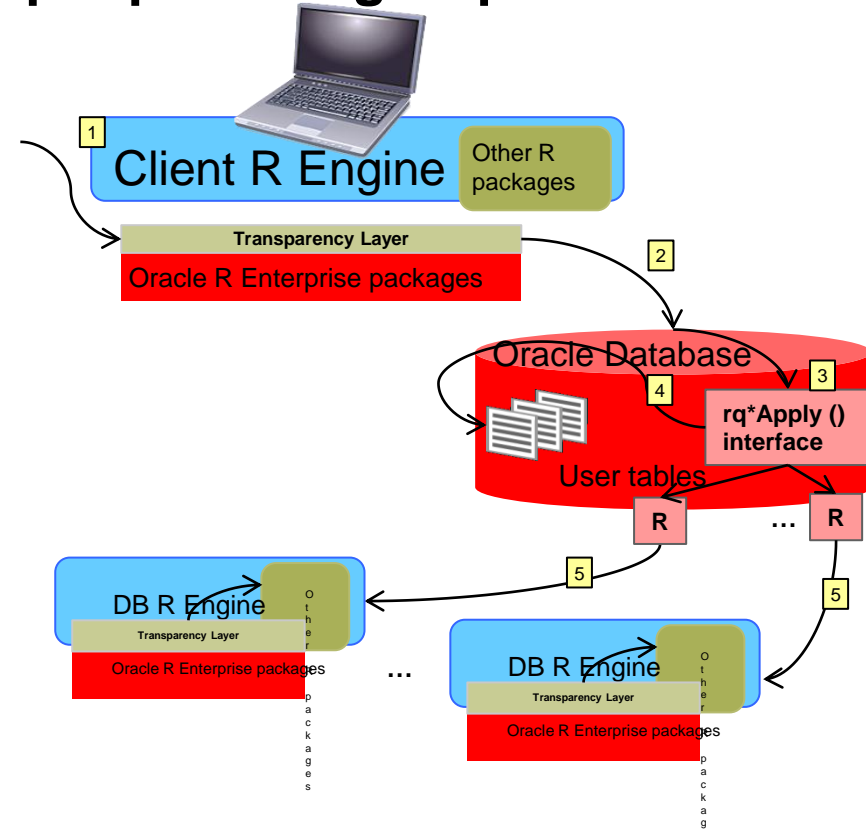
- R closure (script) is the integration point
- Different types of inputs
  - **Parallel row streams**
  - **Parallel groups of rows**
  - **Parallel iterations**
- Run-time parameters
  - **e.g. Date Filters, R objects**
- Flexible outputs
  - **Vertical or Horizontal addition to an existing table**
    - **Data or models**
  - **Frames, Vectors, Graphics**

ORACLE

# 5. Data Flow parallelism at work
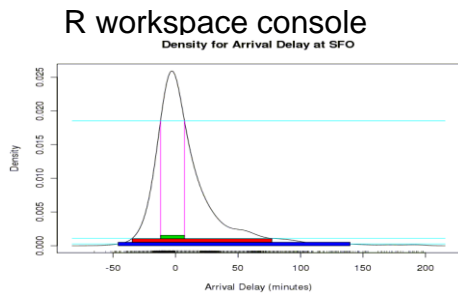## Partitioned model builds: 1 model per product group

```
modList <- ore.groupApply(
   ONTIME_S,
   INDEX=ONTIME_S$DEST,
   function(dat) {
     library(randomforest)
     reg(ARRDELAY ~ DISTANCE + DEPDELAY, dat)
   });
modList_local <- ore.pull(modList)
summary(modList_local$BOS) ## return model for BOS
```

Goal: Build models in parallel on partitions of dataset
Function loaded to DB R Engine
Parallelism enabled through INDEX column

Data group subset for 1 INDEX value passed to DB R Engine

Result "modList" returned as a list of model objects, one per group



Client R Engine — Other R packages

Transparency Layer — Oracle R Enterprise packages

Oracle Database — User tables

rq*Apply () interface

DB R Engine — Transparency Layer — Oracle R Enterprise packages — Other packages

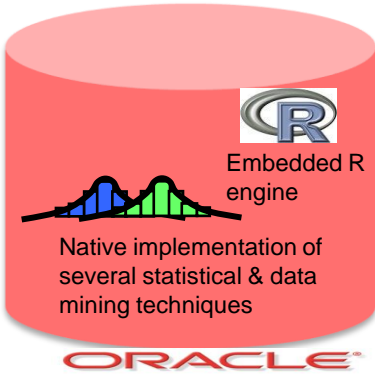# Oracle R Connector for Hadoop



R workspace console

Oracle transparency layer: R to Oracle SQL

HIVE transparency layer: R to HIVE QL

R-Hadoop map-reduce framework

R-HDFS interface for file exploration, sampling

Embedded R engine

Native implementation of several statistical & data mining techniques

Data transfer

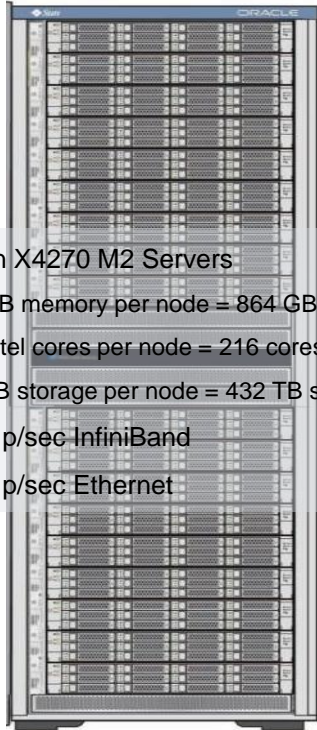Native & Mahout based algorithms

HDFS Files

# Oracle R Connector for Hadoop (ORCH) Concepts

1. Access to HDFS files
   1. Auto discovery of metadata
   2. Sampling
2. HIVE SQL connectivity
   1. R to SQL
3. Hadoop Analytics
   1. Open source Mahout
   2. Home grown techniques

ORACLE®

# Key Highlights

1. Supports interactive access to HDFS data and Hadoop infrastructure

2. Allows database resident data to be used within a Hadoop calculation

3. Supports local execution and debugging of R code – disconnected from Hadoop

4. Treats metadata and data separately when possible

    - Samples HDFS files to create metadata description

5. Provides flexible output options

    - As R object to user session

    - Load to database

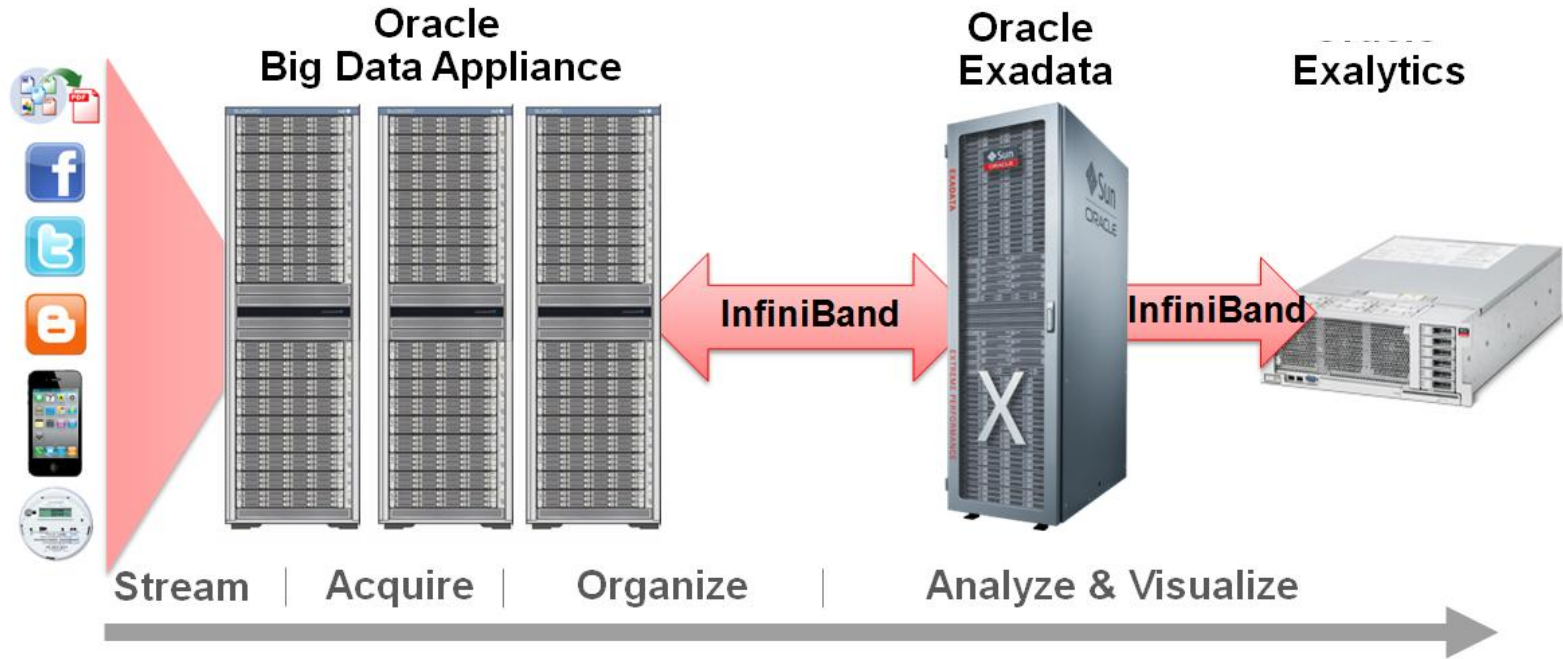    - Continue to post-process

ORACLE

# Big Data Appliance

18 Sun X4270 M2 Servers

   48 GB memory per node = 864 GB memory

   12 Intel cores per node = 216 cores

   24 TB storage per node = 432 TB storage

40 Gb p/sec InfiniBand

10 Gb p/sec Ethernet

- An engineered system optimized for capturing and integrating "low density" data into Exadata
  - High-performance Hardware
    - Optimized for Hadoop and NoSQL workloads
    - InfiniBand Networking for integration with **Exadata**
  - Software:
    - Oracle Hadoop
    - Oracle R Hadoop Connector
    - Oracle R Enterprise client (optional)
    - Oracle NoSQL DB
    - Oracle Data Integrator (Hadoop capabilities)
    - Oracle Loader for Hadoop

ORACLE

# Big Data Appliance Usage Model

ORACLE®

# Key take-aways

1. Improve user efficiency by enabling focus on analysis as opposed to data access
   - Transparent support for R language on database and HDFS objects

2. Enable deep analytics with computations occurring closer to data
   - Native implementation of statistics and data mining algorithms
   - R engine as an embedded component of database

3. Allow transparent access to Compute Infrastructures
   - Database & Hadoop platforms

4. Shorten path to application of cutting edge ideas into practice
   - Oracle's R Distribution & Embedded R engine

5. Enable quick transition from analysis to mass consumption of results
   - R integrated into SQL

**ORACLE**