# Divide & Recombine (D&R)
## with the DeltaRho D&R Software

### D&R + $\delta\rho$

http://deltarho.org

Meet the statistical and computational challenges
of deep analysis of big data and
high computational complexity of analytic methods.

# The D&R Framework: Choices of the Analyst

**Statistical Division Method: [D] Operations**
- divides data into subsets
- division persists, used for many analytic methods

**Analytic Method: [A] Operations**
- applied to each subset, resulting in subset outputs
- no communication among the subset computations
- embarrassingly parallel: simplest parallel computation

**Statistical Recombination Method for Each Analytic Method: [R] Operations**
- statistical recombination method applied to outputs
- this is the D&R result for the analytic method
- often has embarrassingly parallel component

**DeltaRho software implements D&R**

# MapReduce

Simple, powerful programming model for breaking a task into pieces and operating on those pieces in an embarrassingly parallel manner

MapReduce is the foundational computation for D&R operations

Map does parallel computations without communication; application of analytic methods to all subsets uses Map

Reduce enables communication; recombination applied to outputs use Reduce

DeltaRho software enables users to use MapReduce without having to program it directly

# Key-Value Pairs

The data structure used to store subsets or outputs of analytic methods are key-value pairs

The key is a label that uniquely identifies a subset or an output

The value is the subset or output corresponding to the key

# DeltaRho datadr R package

The user programs in R and uses the datadr R package

datadr is a domain specific language for D&R

First written by Ryan Hafen at PNNL (former grad student in Purdue Statistics)

1st implementation Jan 2013

Analyst R and datadr code specifies divisions, analytic methods, and recombinations

# DeltaRho datadr

Two major data types in datadr : the distributed data frame (ddf) and distributed data object (ddo)

A ddf can be thought of as a data frame that is split into subsets by rows

A distributed data object (ddo) is similar, but each subset can be an object with an arbitrary data structure

The data structure we use to store ddo/ddf objects are key-value pairs

Thus, a ddo/ddf is essentially a list, where each element of the list contains a key-value pair

datadr can be run on a multicore server and provide Map and Reduce

# Hadoop and RHIPE

Hadoop is a distributed computational environment running on a cluster with the MapReduce parallel compute engine and the Hadoop distributed file system (HDFS).

RHIPE is the R and Hadoop Integrated Programming Environment

Provides communication between R and Hadoop

Can provide programming of D&R but at a lower level than datadr

First written by Saptarshi Guha while a grad student in Purdue Statistics

1st implementation Jan 2009

# Hadoop and RHIPE

RHIPE enables datadr to have Hadoop at the back

datadr is back-end agnostic in that the code used to run on a single multicore server is the same as that used when Hadoop is the back end

# Hadoop

Runs the analyst's R code for divisions [D], analytic methods [A], and recombinations [R] in parallel on a cluster

Writes subsets and outputs to the HDFS
- R data structures
- specified by the analyst R + datadr code

Schedules computations: assigns core to a subset, e.g., trying to have both on the same cluster node

Operates sequentially on subsets and outputs until all operations are completed

Subsets or outputs do not need to the highly limiting requirement of being in memory all at the same time

# Subject-Matter Division

It is natural to divide data based on the subject matter

Divide by conditioning on the values of variables important to the analysis

Just as valid for small datasets
- widely practiced in the past
- a statistical best practice

D&R with DeltaRho takes advantage of this best practice for computational gain

There are other of division methods, but subject-matter division is the the most used in practice

# Subject-Matter Division

Wen-wen Tung in the next talk will give us examples

50,632 3-hr satellite rainfall measurements at each of 576,000 locations

Division (1) By Time Across Locations:
50,632 subsets, 576,000 measurements per subset

Division (2) By Location Across Time:
576,000 subsets, 50,632 measurements per subset

# Sampling Division: The Concept

Want one result, say a logistic regression, for all of the data

Each subset is seen as a sample of the data

Subsets are replicate samples, or replicates

For example, we can carry out random replicate division: choose subsets randomly

D&R research in statistical theory and methods seeks division and recombination methods that maximize statistical accuracy

# How Fast? The WSC Cluster

11 nodes, each a Hewlett Packard ProLiant DL165 G7

Collectively, the 11 nodes have
- $24 \times 11 = 264$ cores
- $22 \times 11 = 242$ <span style="color:blue">Hadoop</span> cores
- $48 \times 11 = 528$ GB total RAM
- $8 \times 11 = 88$ TB total disk

# How Fast? Analytic Method: Logistic Regression

Subset logistic regressions were carried out using the R function `glm.fit`

Outputs are estimates of regression coefficients

Recombination is taking the means of estimates across subsets

# How Fast? The Data

Number of observations $N = 2^{30}$

1 response and $p$ = 127 explanatory variables

All $128 = 2^7$ variables are numeric

1 TB of data, close to double the size of physical memory

Number of subsets, $R = 2^{20} \approx 1,000,000$

Number of observations per subset, $M = 2^{10} \approx 1,000$

# How Fast? Elapsed Times

18.1 min = total

12.1 min = read subsets into memory and form the subset data.frame objects

6.0 min = logistic regressions on subsets, compute recombination means, and write means to the HDFS

Solid State Disks would have greatly speeded things up, but could not afford them

# What Do We Get from D&R with DeltaRho

Deep analysis for data, big and small

Analyze data at their finest granularity, in detail, and not just summary statistics

The analyst can use any of the 1000s of methods of statistics, machine learning, and data visualization (R provides this)

High statistical accuracy, not substantially sacrificing accuracy to cope with big data and high computational complexity

# What Do We Get from D&R with DeltaRho

Visualization of the data at their finest granularity

Understand patterns in the data, critical for model building

In D&R, visualization method is applied to subsets, which contain the detailed data

Typically cannot look at plots for all subsets

Sample: sampling plans and cognostics

DeltaRho Trelliscope : provides a way to flexibly and interactively visualize large, complex data in great detail

# What do We Get from D&R with DeltaRho

High efficiency programming for the analyst

datadr makes it easy to program D&R

Protects the analyst from the details of parallel distributed computing

Interface is abstracted from the different back end choices, so that datadr code is the same whatever the back end

# What do We Get from D&R with DeltaRho

High performance computing for data analysis

From the parallel computing of Hadoop on a cluster

We are looking at Spark , a Hadoop competitor

Meets the challenge of big size and high computational complexity

The memory size of the data can exceed physical memory

For the TRMM data
- computational performance was not at all a problem
- the real challenge, as usual for big data and small, is discovering the right work flow of analytic methods

# What do We Get from D&R with DeltaRho

It's free, and easy to get

The software is all open source

Documentation will get you going quickly with datadr (try it on a multicore server)

Install DeltaRho on a Hadoop cluster

Use our appliance to spin up a cluster on AWS, the Amazon service

# Get DeltaRho Software and Documentation

Github is the development site: github.com/delta-rho

Open source with both a GPL and Apache license

Available for download from R CRAN for installation on a cluster

Appliance to spin up a cluster on the Amazon AWS service

Much documentation for datadr and RHIPE

Get code and documentation: deltarho.org

# Get Involved: See Invitation on deltarho.org

Contributors are welcome!

Visit our Github Organization page and feel free to fork any of our component repositories:

datadr , Trelliscope , RHIPE

And check out our contributing guide

Join and introduce yourself on our gitter dev chat room