NERDERY®

# R Studio Server on Amazon EMR

Chad Dvoracek

Brandon Veber

NERDERY.

- **Leading IDE for R**
- **Integrated console**
- **Code completion**
- **Syntax highlighting editor**
- **Direct code execution**
- **Plotting, history, debugging**
- **Open source**

NERDERY.

3

# EMR (Elastic MapReduce)
**Managed Hadoop Framework**

- Easy to use

- Quick set up

- Low cost (Spot Instance)

- Elastic

- Flexible

# GUI Interactive Environments

# Case Study

Speeding data deliver to data scientists.

NERDERY.

# Data Science Team

- Temporary projects
- Long term data persistence not needed
- Data sets received in .zip files
- Multiple complex joins
- Aggregation
- Process predictive algorithms

NERDERY.

# Initial Challenges

- Joins in RDBMS taking to long
- Moving data between systems
- Time to delivery on changes
- Data scientists only use R

NERDERY.

VOLUME
DATA SIZE

VELOCITY
SPEED OF CHANGE

VARIETY
DIFFERENT FORMS
OF DATA SOURCES

VERACITY
UNCERTAINTY OF
DATA

NERDERY.

# The Hadoop ecosystem can run in Amazon EMR

MapReduce · Avro · hadoop HDFS · Hive

Flume · Spark · Apache Drill

R · Apache Phoenix · Tez · presto

cascading · MAPR · gradle · accumulo · Apache HBASE · mahout

NERDERY.

# Case Study: Results

✓ **Reduced time for data delivery**

Moving to EMR allowed for faster processing and preparing of the data. By incorporating automating data load scripts and utilizing Hive for batch processing and Spark for in memory processing it allowed the data scientist to focus on solutions rather than time constraints.

✓ **Confidence utilizing Big Data systems**

By providing assistance in set up and training for data processing it allowed the data engineering team the ability to gain confidence in preparing data on distributed systems.

✓ **Process change**

Having the option to transform data in distributed systems provided an opportunity to re-think the process and time needed for data and solution delivery.

✓ **Future considerations**

With the data scientist working primarily in R, exploring R Studio Server and SparkR may be a logical next step in improving the teams workflow.

NERDERY.

# Bootstrap Problems

# What about Spark R?

# EMR Set Up

Quick & Simple

NERDERY.

# Simple Start Up Script

```bash
1   #/bin/bash
2
3   USER="rstudio"
4   USERPW="rstudio"
5
6
7   # create rstudio user on all machines
8   # we need a unix user with home directory
9   # and password and hadoop permission
10  sudo adduser $USER
11  sudo sh -c "echo '$USERPW' | passwd $USER --stdin"
12
13  # fix hadoop tmp permission on all machines
14  sudo chmod 777 -R /mnt/var/lib/hadoop/tmp
```

# EMR Configuration: Step 1

# EMR Configuration: Step 2

# EMR Configuration: Step 2

# EMR Configuration: Step 3

# EMR Configuration: Step 4

# EMR Web Connection

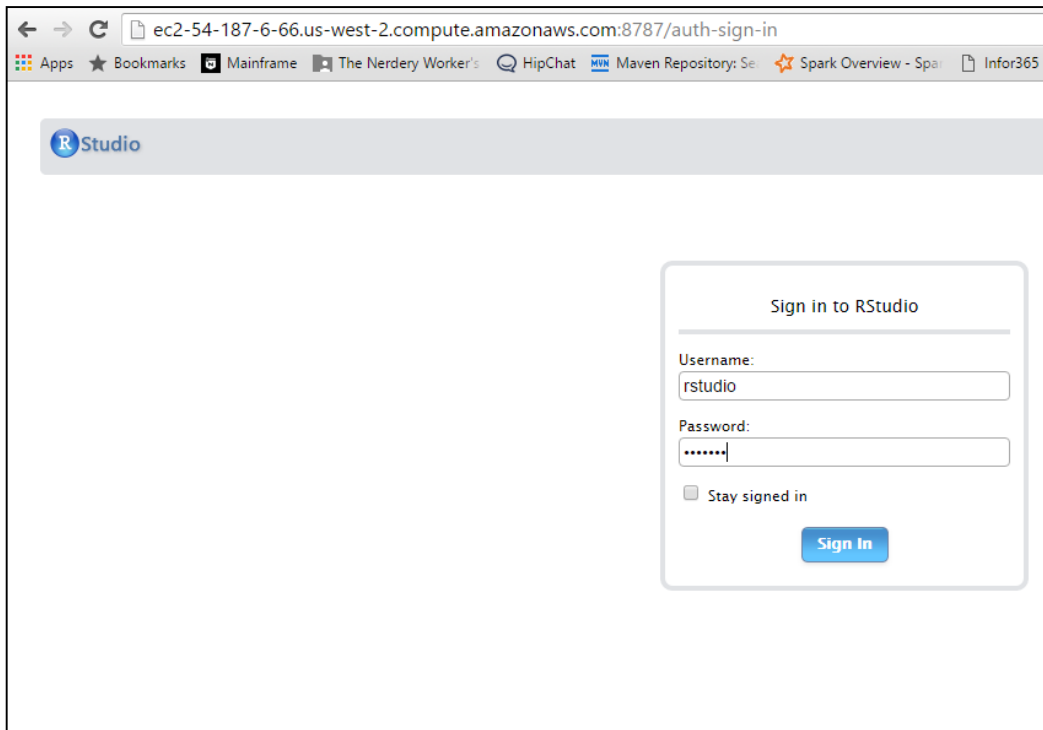# EMR Web Connection

# Final Step

```
8    #SSH Into the EMR Cluster
9
10   #Add rstudio user to hadoop group
11   sudo usermod -a -G hadoop rstudio
12
13
14   #Install R Studio Server
15   wget https://download2.rstudio.org/rstudio-server-rhel-0.99.902-x86_64.rpm
16   sudo yum install --nogpgcheck rstudio-server-rhel-0.99.902-x86_64.rpm
17   sudo rstudio-server verify-installation
18
19   ## Open Web Browser
20   ## Accessing the Server
21   http://<server-ip>:8787
22   example:  http://ec2-54-187-11-218.us-west-2.compute.amazonaws.com:8787/
23   #Username: rstudio
24   #Password: rstudio
25
```

# R Studio Server

R and SparkR

# R Studio Server: Sign In

# R Studio Server

# Environment Set Up



```
bigdata2016.R ×

Source on Save            Run    Source

1
2
3
4
5   # Set this to where Spark is installed
6   Sys.setenv(SPARK_HOME="/usr/lib/spark")
7   Sys.getenv()
8
9   #Load Library and initialize spark context
10  library(SparkR, lib.loc = c(file.path(Sys.getenv("SPARK_HOME"), "R", "lib")))
11
12  #Get local environment info
13  Sys.info()
14  library("parallel", lib.loc="/usr/lib64/R/library")
15  detectCores(all.tests = FALSE, logical = TRUE)
16
17  #Load Library
18  library(SparkR, lib.loc = c(file.path(Sys.getenv("SPARK_HOME"), "R", "lib")))
19
20  #Initialize Spark Contect
21  sc <- sparkR.init(master = "yarn-client", sparkPackages="com.databricks:spark-csv_2.11:
22
23  #Create SQLContext
24  sqlContext <- sparkRSQL.init(sc)
25
26

1:1   (Top Level)                                    R Script
```

# Data Analysis

NERDERY.

# Data Set

# Data Set

# Create RDD Data Frame



```r
25
26  #create rdd data frame
27  loan <- read.df(sqlContext, "/data/clean/loan.txt",source = "com.databricks.spark.csv", header="true", inferSch
28
29  #see head
30  head(loan)
31  take(loan, 10)
32
33  #information on a column
34  typeof(take(loan, 2) [["loan_amnt"]])
35  # [1] "double"
36
37  #rows in the data set
38  count(loan)
39  nrow(loan)
40
41
42  printSchema(loan)
43
44  #Register dataframe as Table
45  registerTempTable(loan, "loanTemp")
46
47
48  #Test TempTable
49  print(head(sql(sqlContext, "Select * from loanTemp limit 5")))
50
```

# SparkR API

# sqlContext

# Access Hive

# Full Power of R

# Full Power of R

# SparkR: In Memory Processing

# Machine Learning

NERDERY

# SparkR

## Limitations

- Machine Learning limited to glm

- Distributed processing constrained by existing API.

## Future

- ✓ Databricks

- ✓ Alteryx

NERDERY.

# Machine Learning in R*

1. **e1071** Functions for latent class analysis, short time Fourier transform, fuzzy clustering, support vector machines, shortest path computation, bagged clustering, naive Bayes classifier etc *(142479 downloads)*
2. **rpart** Recursive Partitioning and Regression Trees. (135390)
3. **igraph** A collection of network analysis tools. (122930)
4. **nnet** Feed-forward Neural Networks and Multinomial Log-Linear Models. (108298)
5. **randomForest** Breiman and Cutler's random forests for classification and regression. (105375)
6. **caret** package (short for Classification And Regression Training) is a set of functions that attempt to streamline the process for creating predictive models. (87151)
7. **kernlab** Kernel-based Machine Learning Lab. (62064)
8. **glmnet** Lasso and elastic-net regularized generalized linear models. (56948)
9. **ROCR** Visualizing the performance of scoring classifiers. (51323)
10. **gbm** Generalized Boosted Regression Models. (44760)
11. **party** A Laboratory for Recursive Partitioning. (43290)
12. **arules** Mining Association Rules and Frequent Itemsets. (39654)
13. **tree** Classification and regression trees. (27882)
14. **klaR** Classification and visualization. (27828)
15. **RWeka** R/Weka interface. (26973)
16. **ipred** Improved Predictors. (22358)
17. **lars** Least Angle Regression, Lasso and Forward Stagewise. (19691)
18. **earth** Multivariate Adaptive Regression Spline Models. (15901)
19. **CORElearn** Classification, regression, feature evaluation and ordinal evaluation. (13856)
20. **mboost** Model-Based Boosting. (13078)

* KDnuggets Top 20 R Machine Learning Packages and Data Science Packages.
Source: http://www.kdnuggets.com/2015/06/top-20-r-machine-learning-packages.html

NERDERY.

# Machine Learning in Spark

- **Classification and regression**
  - **linear models (SVMs, logistic regression, linear regression)**
  - **naive Bayes**
  - **decision trees**
  - **ensembles of trees (Random Forests and Gradient-Boosted Trees)**
  - **isotonic regression**
- **Collaborative filtering**
  - **alternating least squares (ALS)**
- **Clustering**
  - **k-means**
  - **Gaussian mixture**
  - **power iteration clustering (PIC)**
  - **latent Dirichlet allocation (LDA)**
  - **bisecting k-means**
  - **streaming k-means**

- **Dimensionality reduction**
  - **singular value decomposition (SVD)**
  - **principal component analysis (PCA)**
- **Feature extraction and transformation**
- **Frequent pattern mining**
  - **FP-growth**
  - **association rules**
  - **PrefixSpan**
- **Evaluation metrics**
- **PMML model export**
- **Optimization (developer)**
  - **stochastic gradient descent**
  - **limited-memory BFGS (L-BFGS)**

**NERDERY.**

# Resources

Slides:    @the_nerdery

Code:  https://github.com/thenerdery/SparkRTalk

NERDERY.

# Questions?

# Contact

The Nerdery

info@nerdery.com

(877) 664.6373