Machine Learning with R and H2O

Spencer Aiello Eric Eckstrand Anqi Fu

Mark Landry Patrick Aboyoun

Edited by: Jessica Lanford

http://h2o.ai/resources/

April 2017: Sixth Edition

Machine Learning with R and H2O by Spencer Aiello, Eric Eckstrand, Anqi Fu, Mark Landry, & Patrick Aboyoun Edited by: Jessica Lanford

Published by H2O.ai, Inc. 2307 Leghorn St. Mountain View, CA 94043

©2017 H2O.ai, Inc. All Rights Reserved.

April 2017: Sixth Edition

Photos by ©H2O.ai, Inc.

All copyrights belong to their respective owners. While every precaution has been taken in the preparation of this book, the publisher and authors assume no responsibility for errors or omissions, or for damages resulting from the use of the information contained herein.

Printed in the United States of America.

Contents

1	Introduction	5
2	What is H2O?	6
3	Installation 3.1 Installing R	7 7 8 9
4	H2O Initialization4.1 Launching from R4.2 Launching from the Command Line4.3 Launching on Hadoop4.4 Checking Cluster Status	9 11 11 12
5	Data Preparation in R 5.1 Notes	12 12
6	Models 6.1 Supervised Learning	13 14 14 14
7	Demo: GLM	15
8	Data Manipulation in R 8.1 Importing Files	17 17 18 18 19 19 20 20 21 21
	8.11 Summarizing Data	22 23 24

	8.14	Splitting Frames	25
	8.15	Getting Frames	26
		Getting Models	26
		Listing H2O Objects	26
		Removing H2O Objects	27
		Adding Functions	27
9	Runi	ning Models	28
	9.1	Gradient Boosting Machine (GBM)	28
	9.2	Generalized Linear Models (GLM)	30
	9.3	K-means	32
	9.4	Principal Components Analysis (PCA)	33
	9.5	Predictions	33
10	App	endix: Commands	34
	• •	Dataset Operations	34
		General Data Operations	35
		Methods from Group Generics	36
		Other Aggregations	39
		Data Munging	39
		Data Modeling	40
			42
	10.7	H2O Cluster Operations	42
11	Auth	nors	44
12	Refe	erences	45

1 Introduction

This documentation describes how to use H2O in the R environment. More information on H2O's system and algorithms (as well as R user documentation) is available at the H2O website at http://docs.h2o.ai.

R uses a REST API to connect to H2O. To use H2O in R or launch H2O from R, specify the IP address and port number of the H2O instance in the R environment . Datasets are not directly transmitted through the REST API. Instead, commands (for example, importing a dataset at specified HDFS location) are sent either through the browser or the REST API to perform the specified task.

The dataset is then assigned an identifier (the .hex file type in H2O) used as a reference in commands to the web server. After preparing the dataset for modeling by defining significant data and removing insignificant data, H2O creates a model that represents the results of the data analysis. These models are assigned IDs used as references in commands. One of the most popular models for data analysis is GLM.

GLM estimates regression models for outcomes following exponential distributions in general. In addition to the Gaussian (i.e. normal) distribution, these include binomial, gamma, Poisson, and Tweedie distributions. Each serves a different purpose, and depending on distribution and link function, can be used for prediction or classification.

This booklet demonstrates H2O's implementation of GLM in an R environment. For more information on GLM, refer to **Generalized Linear Modeling with H2O** at http://h2o.ai/resources/.

H2O supports Spark, YARN, and most versions of Hadoop. Hadoop is a scalable open-source file system that uses clusters for distributed storage and dataset processing. Depending on the size of your data, H2O can run on your desktop or scale using multiple nodes with Hadoop, an EC2 cluster, or S3 storage.

H2O nodes run as JVM invocations on Hadoop nodes. For performance reasons, we recommend that you do not run an H2O node on the same hardware as the Hadoop NameNode. Because H2O nodes run as mapper tasks in Hadoop, administrators can view them in the normal JobTracker and TaskTracker frameworks, providing process-level (i.e. JVM instance-level) visibility.

H2O helps R users make the leap from laptop-based processing to large-scale environments. Hadoop lets H2O users scale their data processing capabilities

based on their current needs. Using H2O, R, and Hadoop, you can create a complete end-to-end data analysis solution.

This document describes the four steps of data analysis with H2O:

- 1. installing H2O
- 2. preparing your data for modeling (data munging)
- 3. creating a model using simple but powerful machine learning algorithms
- 4. scoring your models

2 What is H2O?

H2O is fast, scalable, open-source machine learning and deep learning for smarter applications. With H2O, enterprises like PayPal, Nielsen Catalina, Cisco, and others can use all their data without sampling to get accurate predictions faster. Advanced algorithms such as deep learning, boosting, and bagging ensembles are built-in to help application designers create smarter applications through elegant APIs. Some of our initial customers have built powerful domain-specific predictive engines for recommendations, customer churn, propensity to buy, dynamic pricing, and fraud detection for the insurance, healthcare, telecommunications, ad tech, retail, and payment systems industries.

Using in-memory compression, H2O handles billions of data rows in-memory, even with a small cluster. To make it easier for non-engineers to create complete analytic workflows, H2O's platform includes interfaces for R, Python, Scala, Java, JSON, and CoffeeScript/JavaScript, as well as a built-in web interface, Flow. H2O is designed to run in standalone mode, on Hadoop, or within a Spark Cluster, and typically deploys within minutes.

H2O includes many common machine learning algorithms, such as generalized linear modeling (linear regression, logistic regression, etc.), Naïve Bayes, principal components analysis, k-means clustering, and others. H2O also implements best-in-class algorithms at scale, such as distributed random forest, gradient boosting, and deep learning. Customers can build thousands of models and compare the results to get the best predictions.

H2O is nurturing a grassroots movement of physicists, mathematicians, and computer scientists to herald the new wave of discovery with data science by collaborating closely with academic researchers and industrial data scientists. Stanford university giants Stephen Boyd, Trevor Hastie, Rob Tibshirani advise the H2O team on building scalable machine learning algorithms. With

hundreds of meetups over the past three years, H2O has become a word-of-mouth phenomenon, growing amongst the data community by a hundred-fold, and is now used by 30,000+ users and is deployed using R, Python, Hadoop, and Spark in 2000+ corporations.

Try it out

- Download H2O directly at http://h2o.ai/download.
- Install H2O's R package from CRAN at https://cran.r-project.org/web/packages/h2o/.
- Install the Python package from PyPI at https://pypi.python.org/pypi/h2o/.

Join the community

- To learn about our meetups, training sessions, hackathons, and product updates, visit http://h2o.ai.
- Visit the open source community forum at https://groups.google.com/d/forum/h2ostream.
- Join the chat at https://gitter.im/h2oai/h2o-3.

3 Installation

H2O requires Java; if you do not already have Java installed, install it from https://java.com/en/download/ before installing H2O.

To use H2O with R, start H2O outside of R and connect to it, or launch H2O from R. However, if you launch H2O from R and close the R session, the H2O session closes as well. The H2O session directs R to the datasets and models located in H2O.

This following sections describe:

- installing R
- installing H2O from R

3.1 Installing R

To download R:

1. Go to http://cran.r-project.org/mirrors.html.

- 2. Select your closest local mirror.
- 3. Select your operating system (Linux, OS X, or Windows).
- 4. Depending on your OS, download the appropriate file, along with any required packages.
- 5. When the download is complete, unzip the file and install.

3.2 Installing H2O from R

To load a recent H2O package from CRAN, run:

Example in R

```
install.packages("h2o")
```

Note: The version of H2O in CRAN is often one release behind the current version.

For the latest recommended version, download the latest stable H2O-3 build from the H2O download page:

- 1. Go to http://h2o.ai/download.
- 2. Choose the latest stable H2O-3 build.
- 3. Click the "Install in R" tab.
- 4. Copy and paste the commands into your R session.

After H2O is installed on your system, verify the installation completed successfully by initializing H2O:

```
library (h2o)
1
2
   #Start H2O on your local machine using all available
3
      cores
   #(By default, CRAN policies limit use to only 2 cores)
   h2o.init(nthreads = -1)
5
   #Get help
7
   ?h2o.glm
8
   ?h2o.qbm
9
10
11
   #Show a demo
```

```
12 # demo(h2o.glm)
13 # demo(h2o.gbm)
```

3.3 Example Code

R code for the examples in this document is located here:

http://github.com/h2oai/h2o-3/tree/master/h2o-docs/src/booklets/v2_2015/source/R_Vignette_code_examples

3.4 Citation

To cite this booklet, use the following:

Aiello, S., Eckstrand, E., Fu, A., Landry, M., and Aboyoun, P. (Apr 2017). *Machine Learning with R and H2O*. http://h2o.ai/resources/.

4 H2O Initialization

This section describes how to launch H2O:

- from R
- from the command line
- on Hadoop

4.1 Launching from R

To specify the number of CPUs for the H2O session, use the nthreads = parameter in the h2o.init command. -2 uses the CRAN default of 2 CPUs. -1 uses all CPUs on the host, which is strongly recommended. To use a specific number of CPUs, enter a positive integer.

To specify the maximum amount of memory for the H2O session, use the max_mem_size parameter in the h2o.init command. The value must a multiple of 1024 greater than 2MB. Append the letter m or M to indicate megabytes, or q or G to indicate gigabytes.

If you do not specify a value for max_mem_size when you run h2o.init, the default heap size of the H2O instance running on 32-bit Java is 1g.

For best performance, the allocated memory should be 4x the size of your data, but never more than the total amount of memory on your computer. For larger datasets, we recommend running on a server or service with more memory available for computing.

H2O checks the Java version and suggests an upgrade if you are running 32-bit Java. On 64-bit Java, the heap size is 1/4 of the total memory available on the machine.

To launch H2O locally from R, run the following in R:

Example in R

```
library(h2o)
the starts H2O using localhost IP, port 54321, all CPUs,
and 4g of memory
h2o.init(ip = 'localhost', port = 54321, nthreads= -1,
max_mem_size = '4g')
```

After successfully launching, R displays output similar to the following example:

```
Successfully connected to http://localhost:54321
1
   R is connected to H2O cluster:
2
      H2O cluster uptime:
                                    11 minutes 35 seconds
3
      H2O cluster version:
                                    2.7.0.1497
4
      H2O cluster name:
                                    H2O_started_from_R
5
      H2O cluster total nodes:
6
      H2O cluster total memory:
                                    3.56 GB
7
      H2O cluster total cores:
8
      H2O cluster allowed cores:
9
                                    8
      H2O cluster healthy:
                                    TRUE
10
```

To launch H2O locally with default initialization arguments, use the following:

```
1 h2o.init()
```

To connect to an established H2O cluster (in a multi-node Hadoop environment, for example) specify the IP address and port number for the established cluster using the <code>ip</code> and <code>port</code> parameters in the <code>h2o.init()</code> command.

Example in R

```
1 h2o.init(ip = "123.45.67.89", port = 54321)
```

4.2 Launching from the Command Line

A simple way to launch H2O from the command line is to download the H2O zip file from the H2O download page. Unzip and launch H2O with the following:

```
unzip h2o-3.5.0.1-*.zip
cd h2o-3.5.0.1-*
java -jar h2o.jar
```

See the H2O Documentation for additional JVM and H2O command line options. After launching the H2O instance, connect to it from R with h2o.init() as described above.

4.3 Launching on Hadoop

To launch H2O nodes and form a cluster on the Hadoop cluster, run:

```
hadoop jar h2odriver.jar -nodes 1 -mapperXmx 6g -
output hdfsOutputDirName
```

 You must launch the Hadoop-specific H2O driver jar (h2odriver.jar) for your Hadoop distribution. Specific driver jar files are available for the following Hadoop versions:

```
- cdh5.2 - hdp2.1 - mapr4.0.1

- cdh5.3 - hdp2.2 - mapr5.0

- cdh5.4.2 - hdp2.3
```

 The above command launches exactly one 6g node of H2O; however, we recommend launching the cluster with 4 times the memory of your data file

- mapperXmx is the mapper size or the amount of memory allocated to each node.
- nodes is the number of nodes requested to form the cluster.
- output is the name of the directory created each time a H2O cloud is created so it is necessary for the name to be unique each time it is launched.

4.4 Checking Cluster Status

To check the status and health of the H2O cluster, use h2o.clusterInfo().

Example in R

```
library(h2o)
h2o.init()
h2o.clusterInfo()
```

An easy-to-read summary of information about the cluster displays.

```
1
   R is connected to H2O cluster:
2
     H2O cluster uptime:
                                   43 minutes 43 seconds
3
     H2O cluster version:
                                   2.7.0.1497
4
     H2O cluster name:
                                   H2O_started_from_R
5
     H2O cluster total nodes:
6
                                   3.56 GB
     H2O cluster total memory:
7
     H2O cluster total cores:
                                   8
8
     H2O cluster allowed cores:
9
     H2O cluster healthy:
                                   TRUE
10
```

5 Data Preparation in R

The following section contains information about data preparation (also known as data munging) and some of the tools and methods available in H2O, as well as a data training example.

5.1 Notes

 Although it may seem like you are manipulating the data in R, once the data has been passed to H2O, all data munging occurs in the H2O instance. The information is passed to R through JSON APIs, so some functions may not have another method.

- You are limited by the total amount of memory allocated to the H2O instance, not by R's ability to handle data. To process large datasets, make sure to allocate enough memory. For more information, refer to Launching from R.
- You can manipulate datasets with thousands of factor levels using H2O in R, so if you ask H2O to display a table in R with information from high cardinality factors, the results may overwhelm R's capacity.
- To manipulate data in R and not in H2O, use as.data.frame(), as.h2o(), and str().
 - as.data.frame() converts an H2O data frame into an R data frame. If your request exceeds the amount of data supported by R, the R session will crash. If possible, we recommend only taking subsets of the entire dataset (the necessary data columns or rows) instead of the whole dataset.
 - as.h2o() transfers data from R to the H2O instance. For successful data transfer, we recommend confirming enough memory is allocated to the H2O instance.
 - str.H2OFrame() returns the elements of the new object to confirm that the data transferred correctly. It's a good way to verify there were no data loss or conversion issues.

6 Models

The following section describes the features and functions of some common models available in H2O. For more information about running these models in R using H2O, refer to **Running Models**.

H2O supports the following models:

- Deep Learning
- Naïve Bayes
- Principal Components Analysis (PCA)
- K-means

- Generalized Linear Models (GLM)
- Gradient Boosted Regression (GBM)
- Distributed Random Forest (DRF)

The list is growing quickly, so check www.h2o.ai to see the latest additions. The following list describes some common model types and features.

6.1 Supervised Learning

Generalized Linear Models (GLM): Provides flexible generalization of ordinary linear regression for response variables with error distribution models other than a Gaussian (normal) distribution. GLM unifies various other statistical models, including Poisson, linear, logistic, and others when using ℓ_1 and ℓ_2 regularization.

Distributed Random Forest (DRF): Averages multiple decision trees, each created on different random samples of rows and columns. It is easy to use, non-linear, and provides feedback on the importance of each predictor in the model, making it one of the most robust algorithms for noisy data.

Gradient Boosting (GBM): Produces a prediction model in the form of an ensemble of weak prediction models. It builds the model in a stage-wise fashion and is generalized by allowing an arbitrary differentiable loss function. It is one of the most powerful methods available today.

Deep Learning: Models high-level abstractions in data by using non-linear transformations in a layer-by-layer method. Deep learning is an example of supervised learning, which can use unlabeled data that other algorithms cannot.

Naïve Bayes: Generates a probabilistic classifier that assumes the value of a particular feature is unrelated to the presence or absence of any other feature, given the class variable. It is often used in text categorization.

6.2 Unsupervised Learning

K-means: Reveals groups or clusters of data points for segmentation. It clusters observations into k-number of points with the nearest mean.

Anomaly Detection: Identifies the outliers in your data by invoking the deep learning autoencoder, a powerful pattern recognition model.

6.3 Modeling Constructs

Grid Search: Performs standard hyper-parameter optimization to simplify model configuration.

After creating a model, use it to make predictions. For more information about predictions, refer to **Predictions**.

Demo: GLM

7 Demo: GLM

The following demo demonstrates how to:

- 1. Import a file
- 2. Define significant data
- 3. View data
- 4. Create testing and training sets using sampling
- 5. Define the model
- 6. Display the results

```
# Import dataset and display summary
   library(h2o)
  h2o.init()
3
   airlinesURL = "https://s3.amazonaws.com/h2o-airlines-
      unpacked/allyears2k.csv"
   airlines.hex = h2o.importFile(path = airlinesURL,
5
      destination_frame = "airlines.hex")
   summary(airlines.hex)
7
   # View quantiles and histograms
8
9
   #high_na_columns = h2o.ignoreColumns(data = airlines.
      hex)
   quantile(x = airlines.hex$ArrDelay, na.rm = TRUE)
10
   h2o.hist(airlines.hex$ArrDelay)
11
12
   # Find number of flights by airport
13
   originFlights = h2o.group_by(data = airlines.hex, by =
14
       "Origin", nrow("Origin"), qb.control=list(na.
      methods="rm"))
   originFlights.R = as.data.frame(originFlights)
15
16
   # Find number of flights per month
17
   flightsByMonth = h2o.group_by(data = airlines.hex, by
      = "Month", nrow("Month"), qb.control=list(na.
      methods="rm"))
   flightsByMonth.R = as.data.frame(flightsByMonth)
19
20
   # Find months with the highest cancellation ratio
21
```

```
which (colnames (airlines.hex) == "Cancelled")
   cancellations By Month = h2o.group by (data = airlines.
23
      hex, by = "Month", sum("Cancelled"), qb.control=
      list(na.methods="rm"))
   cancellation_rate = cancellationsByMonth$sum_Cancelled
24
      /flightsByMonth$nrow_Month
   rates table = h2o.cbind(flightsByMonth$Month,
25
      cancellation rate)
26
   rates table.R = as.data.frame(rates table)
27
   # Construct test and train sets using sampling
28
   airlines.split = h2o.splitFrame(data = airlines.hex,
29
      ratios = 0.85)
30
   airlines.train = airlines.split[[1]]
   airlines.test = airlines.split[[2]]
31
32
   # Display a summary using table-like functions
33
   h2o.table(airlines.train$Cancelled)
34
   h2o.table(airlines.test$Cancelled)
35
36
   # Set predictor and response variables
37
   Y = "IsDepDelayed"
38
   X = c("Origin", "Dest", "DayofMonth", "Year", "
39
      UniqueCarrier", "DayOfWeek", "Month", "DepTime", "
      ArrTime", "Distance")
   # Define the data for the model and display the
40
      results
   airlines.glm <- h2o.glm(training_frame=airlines.train,
41
       x=X, y=Y, family = "binomial", alpha = 0.5)
   # View model information: training statistics,
42
      performance, important variables
   summary(airlines.glm)
43
44
   # Predict using GLM model
45
   pred = h2o.predict(object = airlines.glm, newdata =
46
      airlines.test)
   # Look at summary of predictions: probability of TRUE
47
      class (p1)
   summary(pred$p1)
48
```

8 Data Manipulation in R

The following section describes some common R commands. For a complete command list, including parameters, refer to http://h2o-release.s3. amazonaws.com/h2o/latest_stable_Rdoc.html.

For additional help within R's Help tab, precede the command with a question mark (for example, ?h2o) for suggested commands containing the search terms. For more information on a command, precede the command with two question marks (??h2o).

8.1 Importing Files

The H2O package consolidates all of the various supported import functions using h2o.importFile(). There are a few ways to import files shown in the following examples:

```
#To import small iris data file from H2O's package:
   irisPath = system.file("extdata", "iris.csv", package=
      "h2o")
   iris.hex = h2o.importFile(path = irisPath, destination
3
      _frame = "iris.hex")
4
   #To import an entire folder of files as one data
5
   # pathToFolder = "/Users/data/airlines/"
6
   # airlines.hex = h2o.importFile(path = pathToFolder,
7
      destination frame = "airlines.hex")
8
   #To import from HDFS and connect to H2O in R using the
g
       IP and port of an H2O instance running on your
      Hadoop cluster:
10
   # h2o.init(ip= <IPAddress>, port =54321, nthreads =
      -1)
   # pathToData = "hdfs://mr-0xd6.h2oai.loc/datasets/
11
      airlines_all.csv"
12
   # airlines.hex = h2o.importFile(path = pathToData,
      destination_frame = "airlines.hex")
```

8.2 Uploading Files

To upload a file in a directory local or remote to your H2O cluster, use h2o.importFile(), which is a parallel/distributed read of the data (this is the recommended option). This requires the data to be visible to all H2O nodes in a multi-node cluster. h2o.uploadFile() uploads data local to your H2O cluster as well as uploading data local to your R session. The data only needs to be on the client machine (rather than all the nodes on the H2O cluster), and the upload is single threaded. In the parentheses, specify the H2O reference object in R and the complete URL or normalized file path for the file.

Example in R

```
irisPath = system.file("extdata", "iris.csv", package=
    "h2o")
iris.hex = h2o.uploadFile(path = irisPath, destination
    _frame = "iris.hex")
```

8.3 Finding Factors

To determine if any column contains categorical data (also known as a factor), use h2o.anyFactor(), with the R reference object in the parentheses.

Example in R

8.4 Converting to Factors

To convert an integer into a non-ordered factor (also called an enum or categorical), use as.factor() with the name of the R reference object in parentheses, followed by the number of the column to convert in brackets.

```
# Converts column 4 (RACE) to an enum
as.factor(prostate.hex[,4])

prostate.hex[,4] <- as.factor(prostate.hex[,4])
as.factor(prostate.hex[,4])

# Summary will return a count of the factors
summary(prostate.hex[,4])</pre>
```

8.5 Converting Data Frames

To convert an H2O parsed data object into an R data frame that can be manipulated using R commands, use as.data.frame() with the name of the R reference object in the parentheses. **Caution**: While this can be very useful, be careful when using this command to convert H2O parsed data objects. H2O can easily handle datasets that are often too large to be handled equivalently well in R.

Example in R

```
# Creates object that defines path
1
  prosPath <- system.file("extdata", "prostate.csv",</pre>
2
      package="h2o")
  # Imports data set
3
  prostate.hex = h2o.importFile(path = prosPath,
      destination_frame="prostate.hex")
5
  # Converts current data frame (prostate data set) to
      an R data frame
  prostate.R <- as.data.frame(prostate.hex)</pre>
7
  # Displays a summary of data frame where the summary
      was executed in R
  summary (prostate.R)
```

8.6 Transferring Data Frames

To transfer a data frame from the R environment to the H2O instance, use as.h2o(). In the parentheses, specify the object in the R environment to convert to an H2O object. Optionally, include the name of the destination frame in H2O. Precede the destination frame name with

destination_frame = and enclose the name in quotes as in the following example.

Example in R

```
# Import the iris data into H2O
data(iris)
iris

# Converts R object "iris" into H2O object "iris.hex"
iris.hex = as.h2o(iris, destination_frame= "iris.hex")
head(iris.hex)
```

8.7 Renaming Data Frames

To rename a dataframe on the server running H2O for a dataset manipulated in R, use h2o.assign(). In the following example, the prostate dataset was uploaded to the H2O instance and the data was manipulated to remove outliers. h2o.assign() saves the new dataset on the H2O server so it can be analyzed using H2O without overwriting the original dataset.

Example in R

```
prosPath <- system.file("extdata", "prostate.csv",</pre>
1
      package="h2o")
  prostate.hex<-h2o.importFile(path = prosPath)</pre>
2
3
  ## Assign a new name to prostate dataset in the KV
      store
  h2o.ls()
5
6
7
  prostate.hex <- h2o.assign(data = prostate.hex, key =
      "myNewName")
  h2o.ls()
8
```

8.8 Viewing Column Names

To view a list of the column names in the dataset, use colnames() or names() with the name of the R reference object in the parentheses.

Example in R

```
##Displays the titles of the columns
1
  > colnames(iris.hex)
2
  [1] "Sepal.Length" "Sepal.Width"
                                       "Petal.Length" "
3
      Petal.Width"
                     "Species"
  > names(iris.hex)
4
  [1] "Sepal.Length" "Sepal.Width"
                                       "Petal.Length" "
5
                     "Species"
      Petal.Width"
```

8.9 Getting Minimum and Maximum Values

To view the maximum values for the real-valued columns in a dataset, use \max () with the name of the R reference object in the parentheses. To obtain the minimum values for the real-valued columns in a dataset, use \min () with the name of the R reference object in the parentheses.

Example in R

```
1 > min(prostate.hex$AGE)
2 [1] 43
3 > max(prostate.hex$AGE)
4 [1] 79
```

8.10 Getting Quantiles

To request quantiles for an H2O parsed dataset, use quantile() with the name of the R reference object in the parentheses.

Use quantile (ReferenceObject\$ColumnName), where ReferenceObject represents the R reference object name and ColumnName represents the name of the specified column to request a quantile for a single numerical column.

When you request quantiles for a full parsed dataset consisting of a single column, quantile () displays a matrix with quantile information for the dataset.

Example in R

```
prosPath <- system.file("extdata", "prostate.csv",</pre>
1
      package="h2o")
2
  prostate.hex <- h2o.importFile(path = prosPath)</pre>
   # Returns the percentiles at 0, 10, 20, ..., 100%
3
  prostate.qs <- quantile(prostate.hex$PSA, probs =</pre>
       (1:10)/10)
5
  prostate.qs
   # Take the outliers or the bottom and top 10% of data
7
  PSA.outliers <- prostate.hex[prostate.hex$PSA <=
      prostate.qs["10%"] | prostate.hex$PSA >=
      prostate.gs["90%"],]
   # Check that the number of rows return is about 20% of
9
       the original data
   nrow(prostate.hex)
10
11
  nrow(PSA.outliers)
12
13
14
  nrow(PSA.outliers)/nrow(prostate.hex)
```

8.11 Summarizing Data

To generate a summary similar to the one in R for each of the columns in the dataset, use summary () with the name of the R reference object in the parentheses.

For continuous real functions, this produces a summary that includes information on quartiles, min, max, and mean. For factors, this produces information about counts of elements within each factor level.

```
> summary(prostate.hex)
1
                    CAPSULE
                                    AGE
                                                    RACE
2
   ID
                 DPROS
                   Min. :0.0000
                                    Min.
                                           :43.00
3
   Min. :
             1.00
                                                    Min.
       :0.000 Min. :1.000
   1st Qu.: 95.75 1st Qu.:0.0000
                                    1st Qu.:62.00
                                                    1st Ou
       .:1.000
                1st Qu.:1.000
   Median :190.50
                   Median :0.0000
                                    Median :67.00
                                                    Median
5
       :1.000 Median :2.000
```

6		Mean :0.4026 an :2.271	Mean :66.04	Mean
7	3rd Qu.:285.25	3rd Qu.:1.0000 rd Ou.:3.000	3rd Qu.:71.00	3rd Qu
8		Max. :1.0000	Max. :79.00	Max.
9	DCAPS		VOL	GLEASON
10		Min. : 0.300		
	:0.000			
11		1st Qu.: 4.900	1st Qu.: 0.00	1st Qu
12		Median : 8.664	Median :14.20	Median
13		Mean : 15.409	Mean :15.81	Mean
14		3rd Qu.: 17.063	3rd Qu.:26.40	3rd Qu
15		Max. :139.700	Max. :97.60	Max.

8.12 Summarizing Data in a Table

To summarize data, use h2o.table(). Because H2O can handle larger datasets, it is possible to generate tables that are larger than R's capacity, so use caution when executing this command. To summarize multiple columns, use head(h2o.table (ObjectName[,

c(ColumnNumber, ColumnNumber)])), where ObjectName is the name of the object in R and ColumnNumber is the number of the column.

```
1
   # Counts of the ages of all patients
2
   > head(as.data.frame(h2o.table(prostate.hex[,"AGE"])))
3
      AGE Count
       43
               1
5
       47
               1
   3
       50
               2
7
       51
               3
   5
               2
       52
9
               4
       53
10
11
12
   # Two-way table of ages (rows) and race (cols) of all
      patients
```

```
# Example: For the first row there is one count of a
13
       43 year old that's labeled as RACE = 0
   > h2o.table(prostate.hex[,c("AGE","RACE")])
14
   H2OFrame with 53 rows and 3 columns
15
16
   First 10 rows:
17
      AGE RACE count
       53
               1
                      3
   2
       61
               1
                    12
20
   3
       70
               0
                     1
21
   4
       75
              1
                    11
22
   5
       74
23
               1
                    13
       76
               2
                      1
24
   7
               2
       53
                      1
25
   8
       52
              1
                      2
26
   9
               2
       61
                      1
27
   10
       60
               1
                      9
28
```

8.13 Generating Random Numbers

To append a column of random numbers to an H2O data frame for testing/training data splits that are used for analysis and validation in H2O, use h2o.runif() with the name of the R reference object in the parentheses. This method is best for customized frame splitting; otherwise, use h2o.splitFrame(). However, h2o.runif() is not as fast or stable as h2o.splitFrame().

```
> prosPath <- system.file("extdata", "prostate.csv",
1
      package="h2o")
   > prostate.hex <- h2o.importFile(path = prosPath)</pre>
2
3
   ## Creates object for uniform distribution on prostate
       data set
   > s <- h2o.runif(prostate.hex)</pre>
5
   > summary (s) ## Summarize the results of h2o.runif
   rnd
7
   Min.
          :0.000863
8
    1st Qu.: 0.239763
9
    Median : 0.507936
10
    Mean
          :0.506718
11
```

```
3rd Ou.: 0.765194
12
    Max.
           :0.993178
13
   ## Create training set with threshold of 0.8
14
   > prostate.train <- prostate.hex[s <= 0.8,]</pre>
15
   ##Assign name to training set
16
   > prostate.train <- h2o.assign(prostate.train, "</pre>
17
      prostate.train")
   ## Create test set with threshold to filter values
18
      greater than 0.8
19
   > prostate.test <- prostate.hex[s > 0.8,]
   ## Assign name to test set
20
   > prostate.test <- h2o.assign(prostate.test, "prostate</pre>
21
       .test")
   ## Combine results of test \& training sets, then
22
      display result
   > nrow(prostate.train) + nrow(prostate.test)
23
   [1] 380
24
   > nrow(prostate.hex) ## Matches the full set
25
   [1] 380
26
```

8.14 Splitting Frames

To generate two subsets (according to specified ratios) from an existing H2O dataset for testing/training, use h2o.splitFrame(), which returns contiguous sections of the data without random sampling.

Note: h2o.splitFrame() does not give an exact split. H2O is designed to be efficient on big data using a probabilistic splitting method rather than an exact split. For example specifying a split of 0.75/0.25, H2O will produce a test/train split with an expected value of 0.75/0.25 rather than exactly 0.75/0.25. On small datasets, the sizes of the resulting splits will deviate from the expected value more than on big data, where they will be very close to exact.

Example in R

8.15 Getting Frames

To create a reference object to the data frame in H2O, use h2o.getFrame(). This is helpful for switching between the web UI and the R API or for multiple users accessing the same H2O instance. The following example assumes prostate.hex is in the key-value (KV) store.

Example in R

```
prostate.hex <- h2o.getFrame(id = "prostate.hex")</pre>
```

8.16 Getting Models

To create a reference object for the model in H2O, use h2o.getModel(). This is helpful for users that alternate between the web UI and the R API or multiple users accessing the same H2O instance.

In the following example, it is assumed that a GBM with the ID GBM_8e4591a9b413407b983d73fbd9eb44cf is in the key-value (KV) store.

Example in R

8.17 Listing H2O Objects

To generate a list of all H2O objects generated during a session and each object's size in bytes, use h2o.ls().

8.18 Removing H2O Objects

To remove an H2O object on the server associated with the object in the R environment, use h2o.rm(). For optimal performance, we recommend removing the object from the R environment as well using remove(), with the name of the object in the parentheses. If you do not specify an R environment, then the current environment is used.

Example in R

```
h2o.rm(c("prostate.train", "prostate.test"))
h2o.ls()
```

8.19 Adding Functions

User-defined functions no longer need to be added explicitly to the H2O instance. An R function can be defined and executed against an H2OFrame.

```
# Create an R functional expression
1
   > simpleFun <- function(x) { 2*x + 5 }
2
   # Evaluate the expression across prostate's AGE column
3
   > calculated <- simpleFun(prostate.hex[,"AGE"])</pre>
   > h2o.cbind(prostate.hex[,"AGE"], calculated)
5
6
   H2OFrame with 380 rows and 2 columns
7
8
   First 10 rows:
9
      AGE AGEO
10
       65
            135
       72.
            149
   3
       70
          145
13
      76
           157
14
   5
       69
           143
15
16
       71
            147
   7
       68
            141
17
```

```
    18
    8
    61
    127

    19
    9
    69
    143

    20
    10
    68
    141
```

9 Running Models

This section describes how to run the following model types:

- Gradient Boosting Machine (GBM)
- Generalized Linear Models (GLM)
- K-means
- Principal Components Analysis (PCA)

as well as how to generate predictions.

9.1 Gradient Boosting Machine (GBM)

To generate gradient boosted models for developing forward-learning ensembles, use h2o.gbm(). In the parentheses, define x (the predictor variable vector), y (the integer or categorical response variable), the distribution type (multinomial is the default, gaussian is used for regression), and the name of the H2OParsedData object. For more information, use help(h2o.gbm).

```
> library(h2o)
1
  > h2o.init(nthreads = -1)
2
  > data(iris)
3
   > iris.hex <- as.h2o(iris,destination_frame = "iris.</pre>
4
      hex")
   > iris.gbm < h2o.gbm(y = 1, x = 2:5, training_frame =
5
       iris.hex, ntrees = 10,
       max_depth = 3, min_rows = 2, learn_rate = 0.2,
6
          distribution= "gaussian")
7
   # To obtain the Mean-squared Error by tree from the
8
      model object:
   > iris.gbm@model$scoring_history
9
   Scoring History:
10
                timestamp
                             duration number_of_trees
11
                    training MSE training deviance
```

12	1	2015-09-11 09:50:16		1
13	2	0.47256 2015-09-11 09:50:16 0.33494		2
14	3	2015-09-11 09:50:16		3
15	4	2015-09-11 09:50:16 0.18414	0.014 sec	4
16	5	2015-09-11 09:50:16 0.14363	0.017 sec	5
17	6	2015-09-11 09:50:16		6
18	7	2015-09-11 09:50:16 0.09916	0.023 sec	7
19	8	2015-09-11 09:50:16 0.08649	0.026 sec 0.08649	8
20	9	2015-09-11 09:50:16 0.07761	0.029 sec	9
21	10	2015-09-11 09:50:16 0.07071		10

To generate a classification model that uses labels, use distribution= "multinomial":

```
> iris.gbm2 <- h2o.gbm(y = 5, x = 1:4, training_frame</pre>
1
      = iris.hex, ntrees = 15, max_depth = 5, min_rows =
       2, learn rate = 0.01, distribution= "multinomial"
2
3
  > iris.gbm2@model$training_metrics
4
5
  H2OMultinomialMetrics: gbm
  ** Reported on training data. **
6
  Training Set Metrics:
  ______
10
  Extract training frame with 'h2o.getFrame("iris.hex") '
11
  MSE: (Extract with 'h2o.mse') 0.3293958
12
  R^2: (Extract with 'h2o.r2') 0.5059063
13
  Logloss: (Extract with 'h2o.logloss') 0.8533637
14
```

```
Confusion Matrix: Extract with 'h2o.confusionMatrix(<
     model>,train=TRUE) ')
  _____
16
            setosa versicolor virginica
                                           Error
17
               Rate
                50
  setosa
                           0
                                    0 0.00000000
                                                 0 /
18
      50
                                    1 0.02000000
  versicolor
                 0
                          49
                                                 1 /
19
      50
  virginica
                 0
                           1
                                   49 0.02000000
                                                 1
20
      50
  Totals
                          50
                                   50 0.01333333 2 /
21
                50
     150
22
  Hit Ratio Table: Extract with 'h2o.hit_ratio_table(<
23
     model>,train=TRUE) '
  _____
24
  Top-3 Hit Ratios:
25
    k hit ratio
26
      0.986667
27
  2. 2.
       1.000000
  3 3 1.000000
29
```

9.2 Generalized Linear Models (GLM)

Generalized linear models (GLM) are some of the most commonly-used models for many types of data analysis use cases. While some data can be analyzed using general linear models, general linear models may not be as accurate if the variables are more complex. For example, if the dependent variable has a non-continuous distribution or if the effect of the predictors is not linear, generalized linear models will produce more accurate results than general linear models.

Generalized Linear Models (GLM) estimate regression models for outcomes following exponential distributions in general. In addition to the Gaussian (i.e. normal) distribution, these include Poisson, binomial, gamma and Tweedie distributions. Each serves a different purpose, and depending on distribution and link function choice, it can be used either for prediction or classification.

H2O's GLM algorithm fits the generalized linear model with elastic net penalties. The model fitting computation is distributed, extremely fast, and scales extremely well for models with a limited number (\sim low thousands) of

predictors with non-zero coefficients. The algorithm can compute models for a single value of a penalty argument or the full regularization path, similar to glmnet. It can compute Gaussian (linear), logistic, Poisson, and gamma regression models.

To generate a generalized linear model for developing linear models for exponential distributions, use h2o.glm(). You can apply regularization to the model by adjusting the lambda and alpha parameters. For more information, use help(h2o.glm).

```
> prostate.hex <- h2o.importFile(path = "https://raw.
1
      github.com/h2oai/h2o/master/smalldata/logreg/
      prostate.csv" , destination_frame = "prostate.hex"
2
   > prostate.glm<-h2o.glm(y = "CAPSULE", x = c("AGE","
3
      RACE", "PSA", "DCAPS"), training_frame = prostate.
      hex, family = "binomial", nfolds = 10, alpha =
      0.5)
  > prostate.glm@model$cross_validation_metrics
4
  H2OBinomialMetrics: glm
5
  ** Reported on cross-validation data. **
  Description: 10-fold cross-validation on training data
7
  MSE: 0.2093902
8
  R^2: 0.1294247
9
  LogLoss: 0.6095525
10
  AUC: 0.6909965
11
  Gini: 0.381993
  Null Deviance: 513.8229
13
  Residual Deviance: 463.2599
14
  AIC: 473.2599
15
  Confusion Matrix for F1-optimal threshold:
16
17
                     Error
                                 Rat.e
          122 105 0.462555
                            =105/227
18
           41 112 0.267974
                            =41/153
19
  Totals 163 217 0.384211 =146/380
20
  Maximum Metrics:
21
                         metric threshold
                                              value idx
22
23
  1
                         max f1 0.312978 0.605405 216
  2
                         max f2 0.138305 0.772727 377
24
  3
                   max f0point5 0.400689 0.628141 110
25
   4
                   max accuracy 0.400689 0.700000 110
26
```

```
      27
      5
      max precision
      0.998848
      1.000000
      0

      28
      6
      max absolute_MCC
      0.400689
      0.357638
      110

      29
      7 max min_per_class_accuracy
      0.330976
      0.621145
      181
```

9.3 K-means

To generate a K-means model for data characterization, use h2o.kmeans(). This algorithm does not rely on a dependent variable. For more information, use help(h2o.kmeans).

```
> h2o.kmeans(training_frame = iris.hex, k = 3, x =
      1:4)
   Model Details:
2
   =========
3
   H2OClusteringModel: kmeans
  Model ID: K-means model R 1441989204383 30
  Model Summary:
     number of rows number of clusters number of
        categorical columns number of iterations within
        cluster_sum_of_squares
   1
                150
                                       3
8
                                    0
                                                          8
                            139.09920
     total_sum_of_squares between_cluster_sum_of_squares
9
                                                 456.90080
                596,00000
10
   H2OClusteringMetrics: kmeans
11
   ** Reported on training data. **
12
   Total Within SS: 139.0992
13
   Between SS:
                456.9008
14
   Total SS: 596
   Centroid Statistics:
                  size within_cluster_sum_of_squares
     centroid
17
   1
            1 44.00000
                                              43.34674
18
            2 50.00000
                                              47.35062
   2
19
20
   3
            3 56.00000
                                              48.40184
```

9.4 Principal Components Analysis (PCA)

To map a set of variables onto a subspace using linear transformations, use h2o.prcomp(). This is the first step in Principal Components Regression. For more information, use help(h2o.prcomp).

Example in R

```
> ausPath = system.file("extdata", "australia.csv",
1
      package="h2o")
  > australia.hex = h2o.importFile(path = ausPath)
2
  > australia.pca <- h2o.prcomp(training_frame =</pre>
3
      australia.hex, transform = "STANDARDIZE", k = 3)
  > australia.pca
  Model Details:
5
  _____
  H2ODimReductionModel: pca
7
  Model Key: PCA_model_R_1441989204383_36
  Importance of components:
9
                                        pc2
10
                               pc1
                                                 рс3
  Standard deviation 1.750703 1.512142 1.031181
11
  Proportion of Variance 0.383120 0.285822 0.132917
12
  Cumulative Proportion 0.383120 0.668942 0.801859
```

9.5 Predictions

The following section describes some of the prediction methods available in H2O.

Predict: Generate outcomes of a dataset with any model. Predict with GLM, GBM, Decision Trees or Deep Learning models.

Confusion Matrix: Visualize the performance of an algorithm in a table to understand how a model performs.

Area Under Curve (AUC): A graphical plot to visualize the performance of a model by its sensitivity, true positives and false positives to select the best model.

Hit Ratio: A classification matrix to visualize the ratio of the number of correctly classified and incorrectly classified cases.

PCA Score: Determine how well your feature selection fits a particular model.

Multi-Model Scoring: Compare and contrast multiple models on a dataset to find the best performer to deploy into production.

To apply an H2O model to a holdout set for predictions based on model results, use h2o.predict(). In the following example, H2O generates a model and then displays the predictions for that model. For classification, the predict column is the model's discrete prediction, based on maximum F1 by default; the individual class probabilities are the remaining columns in the data frame. It is common to utilize the p1 column for binary classification, if a raw probability is desired.

Example in R

```
> prostate.fit = h2o.predict(object = prostate.glm,
1
      newdata = prostate.hex)
   > prostate.fit
2
   H2OFrame with 380 rows and 3 columns. First 10 rows:
3
      predict
                       0g
             0 0.74476265 0.2552373
   1
5
   2
             1 0.39763451 0.6023655
6
   3
             1 0.41268532 0.5873147
7
             1 0.37270563 0.6272944
   4
8
   5
             1 0.64649990 0.3535001
g
   6
             1 0.43367145 0.5663285
10
   7
             1 0.26542251 0.7345775
11
             1 0.06143281 0.9385672
12
13
             0 0.73057373 0.2694263
             1 0.46709293 0.5329071
   10
```

10 Appendix: Commands

The following section lists some common R commands by function and a brief description of each command.

10.1 Dataset Operations

Data Import/Export

h2o.downloadCSV: Download a H2O dataset to a CSV file on local disk.

h2o.exportFile: Export H2O Data Frame to a file.

h2o.importFile: Import a file from the local path and parse it.

Appendix: Commands

h2o.parseRaw: Parse a raw data file.

h2o.uploadFile: Upload a file from the local drive and parse it.

Native R to H2O Coercion

as.h2o: Convert an R object to an H2O object.

H2O to Native R Coercion

as.data.frame: Check if an object is a data frame, or coerce it if possible.

Data Generation

h2o.createFrame: Create an H2O data frame, with optional randomization.

h2o.runif: Produce a vector of random uniform numbers.

h2o.interaction: Create interaction terms between categorical features of an H2O Frame.

Data Sampling/Splitting

h2o.splitFrame: Split an existing H2O dataset according to user-specified ratios.

Missing Data Handling

h20.impute: Impute a column of data using the mean, median, or mode.

h2o.insertMissingValues: Replaces a user-specified fraction of entries in a H2O dataset with missing values.

10.2 General Data Operations

Subscripting example to pull pieces from data object.

```
## note: chooses column J, not row J
     x[i]
1
     x[i, j]
2
     x[[i]]
3
     x$name
4
     x[i] \leftarrow value
     x[i, j, \ldots] \leftarrow value
6
     x[[i]] \leftarrow value
7
     x$i <- value
8
```

Subsetting

head, tail: Return the First or Last Part of an Object

Concatenation

c: Combine Values into a Vector or List

h2o.cbind: Take a sequence of H2O datasets and combine them by column.

Data Attributes

colnames: Return column names for a parsed H2O data object.

colnames<-: Retrieve or set the row or column names of a matrix-like object.

names: Get the name of an object.

names<-: Set the name of an object.

dim: Retrieve the dimension of an object.

length: Get the length of vectors (including lists) and factors.

nrow: Return a count of the number of rows in an H2OParsedData object.

ncol: Return a count of the number of columns in an H2OParsedData object.

h20.anyFactor: Check if an H2O parsed data object has any categorical data columns.

is.factor: Check if a given column contains categorical data.

Data Type Coercion

as.factor: Convert a column from numeric to factor.

as.Date: Converts a column from factor to date.

10.3 Methods from Group Generics

Math (H2O)

abs: Compute the absolute value of \boldsymbol{x} .

sign: Return a vector with the signs of the corresponding elements of x (the sign of a real number is 1, 0, or -1 if the number is positive, zero, or negative, respectively).

sqrt: Computes the principal square root of x, \sqrt{x} .

ceiling: Take a single numeric argument x and return a numeric vector containing the smallest integers not less than the corresponding elements of x.

floor: Take a single numeric argument x and return a numeric vector containing the largest integers not greater than the corresponding elements of x.

trunc: Take a single numeric argument x and return a numeric vector containing the integers formed by truncating the values in x toward 0.

log: Compute logarithms (by default, natural logarithms).

exp: Compute the exponential function.

Math (generic)

cummax: Display a vector of the cumulative maxima of the elements of the argument.

cummin: Display a vector of the cumulative minima of the elements of the argument.

cumprod: Display a vector of the cumulative products of the elements of the argument.

cumsum: Display a vector of the cumulative sums of the elements of the argument.

log10: Compute common (i.e., base 10) logarithms

log2: Compute binary (i.e., base 2) logarithms.

log1p: Compute log(1+x) accurately also for |x| << 1.

acos: Compute the trigonometric arc-cosine.

acosh: Compute the hyperbolic arc-cosine.

asin: Compute the trigonometric arc-sine.

asinh: Compute the hyperbolic arc-sine.

atan: Compute the trigonometric arc-tangent.

atanh: Compute the hyperbolic arc-tangent.

expm1: Compute $\exp(x) - 1$ accurately also for |x| << 1.

 ${\tt cos:}$ Compute the trigonometric cosine.

cosh: Compute the hyperbolic cosine.

cospi: Compute the trigonometric two-argument arc-cosine.

sin: Compute the trigonometric sine.

sinh: Compute the hyperbolic sine.

sinpi: Compute the trigonometric two-argument arc-sine.

tan: Compute the trigonometric tangent.

tanh: Compute the hyperbolic tangent.

tanpi: Compute the trigonometric two-argument arc-tangent.

gamma: Display the gamma function γx

lgamma: Display the natural logarithm of the absolute value of the gamma function.

digamma: Display the first derivative of the logarithm of the gamma function.

trigamma: Display the second derivative of the logarithm of the gamma function.

Math2 (H2O)

round: Round the values to the specified number of decimal places. The default is 0.

 ${\tt signif:}$ Round the values to the specified number of significant digits.

Summary (H2O)

max: Display the maximum of all the input arguments.

min: Display the minimum of all the input arguments.

range: Display a vector containing the minimum and maximum of all the given arguments.

sum: Calculate the sum of all the values present in its arguments.

Summary (generic)

prod: Display the product of all values present in its arguments.

any: Given a set of logical vectors, determine if at least one of the values is true.

all: Given a set of logical vectors, determine if all of the values are true.

Appendix: Commands

39

10.4 Other Aggregations

Non-Group Generic Summaries

mean: Generic function for the (trimmed) arithmetic mean.

sd: Calculate the standard deviation of a column of continuous real valued data

var: Compute the variance of x.

summary: Produce result summaries of the results of various model fitting functions.

quantile: Obtain and display quantiles for H2O parsed data.

Row / Column Aggregation

apply: Apply a function over an H2O parsed data object (an array).

Group By Aggregation

h2o.group_by: Apply an aggregate function to each group of an H2O dataset.

Tabulation

h2o.table: Use the cross-classifying factors to build a table of counts at each combination of factor levels.

10.5 Data Munging

General Column Manipulations

is.na: Display missing elements.

Element Index Selection

h2o.which: Display the row numbers for which the condition is true.

Conditional Element Value Selection

h2o.ifelse: Apply conditional statements to numeric vectors in H2O parsed data objects.

Numeric Column Manipulations

h2o.cut: Convert H2O Numeric Data to Factor.

Character Column Manipulations

h2o.strsplit: Splits the given factor column on the input split.

h2o.tolower: Change the elements of a character vector to lower case.

h2o.toupper: Change the elements of a character vector to lower case.

h2o.trim: Remove leading and trailing white space.

h20.gsub: Match a pattern & replace all instances of the matched pattern with the replacement string globally.

h2o.sub: Match a pattern & replace the first instance of the matched pattern with the replacement string.

Factor Level Manipulations

h2o.levels: Display a list of the unique values found in a column of categorical data.

Date Manipulations

h2o.month: Convert the entries of a H2OParsedData object from milliseconds to months (on a 0 to 11 scale).

h2o.year: Convert the entries of a H2OParsedData object from milliseconds to years, indexed starting from 1900.

Matrix Operations

**: Multiply two matrices, if they are conformable. t: Given a matrix or data.frame x, t returns the transpose of x.

10.6 Data Modeling

Model Training: Supervised Learning

h2o.deeplearning: Perform Deep Learning neural networks on an H2OParsedData object.

h20.gbm: Build gradient boosted classification trees and gradient boosted regression trees on a parsed dataset.

h2o.glm: Fit a generalized linear model, specified by a response variable, a set of predictors, and a description of the error distribution.

Appendix: Commands

h2o.naiveBayes: Build gradient boosted classification trees and gradient boosted regression trees on a parsed dataset.

h2o.prcomp: Perform principal components analysis on the given dataset.

h2o.randomForest: Perform random forest classification on a dataset.

Model Training: Unsupervised Learning

h20.anomaly: Detect anomalies in a H2O dataset using a H2O deep learning model with auto-encoding.

h20.deepfeatures: Extract the non-linear features from a H2O dataset using a H2O deep learning model.

h2o.kmeans: Perform k-means clustering on a dataset.

Grid Search

h2o.grid: Efficient method to build multiple models with different hyperparameters.

Model Scoring

h2o.predict: Obtain predictions from various fitted H2O model objects.

Model Metrics

ho2.model_metrics: Given predicted values (target for regression, class-1 probabilities, or binomial or per-class probabilities for multinomial), compute a model metrics object.

Classification Model Helpers

h2o.accuracy: Get the between cluster sum of squares.

h2o.auc: Retrieve the AUC (area under ROC curve).

h2o.confusionMatrix: Display prediction errors for classification data from a column of predicted responses and a column of actual (reference) responses in H2O.

h2o.hit_ratio_table: Retrieve the Hit Ratios. If train, valid, and xval parameters are FALSE (default), then the training Hit Ratios value is returned. If more than one parameter is set to TRUE, then a named list

of Hit Ratio tables are returned, where the names are train, valid, or xval.

h2o.performance: Evaluate the predictive performance of a model via various measures.

Regression Model Helper

h20.mse: Display the mean squared error calculated from a column of predicted responses and a column of actual (reference) responses in H20.

Clustering Model Helper

h2o.betweenss: Get the between cluster sum of squares.

h2o.centers: Retrieve the Model Centers.

10.7 H2O Cluster Operations

H2O Key Value Store Access

h2o.assign: Assign H2O hex.keys to objects in their R environment.

h2o.getFrame: Get a reference to an existing H2O dataset.

h2o.getModel: Get a reference to an existing H2O model.

h20.1s: Display a list of object keys in the running instance of H20.

h20.rm: Remove H2O objects from the server where the instance of H2O is running, but does not remove it from the R environment.

H2O Object Serialization

h2o.loadModel: Load an H2OModel object from disk.

h2o.saveModel: Save an H2OModel object to disk to be loaded back into H2O using h2o.loadModel.

H2O Cluster Connection

h2o.init (nthreads = -1): Connect to a running H2O instance using all CPUs on the host and check the local H2O R package is the correct version.

h2o.shutdown: Shut down the specified H2O instance. All data on the server will be lost!

H2O Load Balancing

h2o.rebalance: Rebalance (repartition) an existing H2O dataset into given number of chunks (per Vec), for load-balancing across multiple threads or nodes.

Appendix: Commands

H2O Cluster Information

h2o.clusterInfo: Display the name, version, uptime, total nodes, total memory, total cores and health of a cluster running H2O.

 ${\tt h2o.clusterStatus}$: Retrieve information on the status of the cluster running H2O.

H2O Logging

h2o.clearLog: Clear all H2O R command and error response logs from the local disk.

h2o.downloadAllLogs: Download all H2O log files to the local disk.

h20.logAndEcho: Write a message to the H2O Java log file and echo it back.

h2o.openLog: Open existing logs of H2O R POST commands and error responses on the local disk.

h2o.getLogPath: Get the file path for the H2O R command and error response logs.

 ${\tt h2o.startLogging: Begin \ logging \ H2O \ R \ POST \ commands \ and \ error \ responses.}$

h2o.stopLogging: Stop logging H2O R POST commands and error responses.

H2O String Manipulation

h20.gsub: String global substitution (all occurrences).

h2o.strsplit: String Split.

h2o.sub: String substitution (first occurrence).

h2o.tolower: Convert characters to lower case.

h2o.toupper: Convert characters to upper case.

h2o.trim: Trim spaces.

11 Authors

Spencer Aiello

Spencer comes from an unconventional background. After studying Physics and Math as an undergraduate at UCSC, he came to San Francisco to continue his education, earning his MS in Analytics from USF. Spencer has worked on a number of projects related to analytics in R, Python, Java, and SQL. At H2O, he works primarily on the R front-end and the backend Java R-interpreter.

Eric Eckstrand

Eric is a Quality and Performance Hacker at H2O.ai. Eric has formal education in computer science and systems engineering. Prior to joining H2O, Eric was a submariner in the US Navy. His roles included Reactor Controls Assistant and Communications Officer (USS Pittsburgh SSN-720) and Submarine Operations Officer (Destroyer Squadron One).

Angi Fu

Anqi is a data hacker at H2O.ai. After completing her undergraduate at University of Maryland: College Park, she attended Stanford University. While there she earned her master's degree in Economics, and a second masters degree in Statistics. Her interests include machine learning and optimization.

Mark Landry

Mark Landry is a competition data scientist and product manager at H2O. He enjoys testing ideas in Kaggle competitions, where he is ranked in the top 100 in the world (top 0.03%) and well-trained in getting quick solutions to iterate over. Interests are multi-model architectures and helping the world make fewer models that perform worse than the mean.

Patrick Aboyoun

Patrick is a former math and data hacker at H2O.ai who has made a career out of creating and delivering software and training for data scientists, particularly those who love R. Patrick received an M.S. in Statistics from the University of Washington and a B.S. in Statistics from Carnegie Mellon University.

Jessica Lanford

Jessica is a word hacker and seasoned technical communicator at H2O.ai. She brings our product to life by documenting the many features and functionality of H2O. Having worked for some of the top companies in technology including Dell, AT&T, and Lam Research, she is an expert at translating complex ideas to digestible articles.

12 References

H2O.ai Team. H2O website, 2017. URL http://h2o.ai

H2O.ai Team. H2O documentation, 2017. URL http://docs.h2o.ai

H2O.ai Team. **H2O GitHub Repository**, 2017. URL https://github.com/h2oai

H2O.ai Team. H2O Datasets, 2017. URL http://data.h2o.ai

H2O.ai Team. H2O JIRA, 2017. URL https://jira.h2o.ai

H2O.ai Team. H2O R Package Documentation, 2017. URL

http://h2o-release.s3.amazonaws.com/h2o/latest_stable_ Rdoc.html

H2O.ai Team. **H2Ostream**, 2017. URL https://groups.google.com/d/forum/h2ostream

Erin LeDell. R Ensemble Documentation, 2014. URL

http://www.stat.berkeley.edu/~ledell/R/h2oEnsemble.pdf

H2O.ai Team. **h2o:** R Interface for H2O, 2017. URL http://www.h2o.ai. R package version 3.1.0.99999