

The H2O Package

W. Evan Johnson, Ph.D.

Professor, Division of Infectious Disease

Director, Center for Data Science

Director, Center for Biomedical Informatics and Health AI

Rutgers University – New Jersey Medical School

2025-10-21

Describing H2O

There is a lot of buzz for machine learning algorithms as well as a requirement for its experts. We all know that there is a significant gap in the skill requirement. The motive of H2O is to provide a platform which made easy for the non-experts to do experiments with machine learning.

Describing H2O

H2O's core code is written in Java that enables the whole framework for multi-threading. Although it is written in Java, it provides interfaces for R, Python and others, thus enabling it to be used efficiently.

In short, we can say that H2O is an open source, in memory, distributed, fast and scalable machine learning and predictive analytics toolkit that facilitates the building and application of machine learning models.

Using H2O

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

```
##  
## -----  
##  
## Your next step is to start H2O:  
##   > h2o.init()  
##  
## For H2O package documentation, ask for help:  
##   > ??h2o  
##  
## After starting H2O, you can use the Web UI at http://localhost:54321  
## For more information visit https://docs.h2o.ai  
## -----  
  
##  
## Attaching package: 'h2o'  
  
## The following objects are masked from 'package:lubridate':  
##  
##   day, hour, month, week, year  
  
## The following objects are masked from 'package:stats':  
##
```

Using H2O

```
h2o.init()
```

```
## Connection successful!
##
## R is connected to the H2O cluster:
##   H2O cluster uptime:      2 minutes 23 seconds
##   H2O cluster timezone:    America/Denver
##   H2O data parsing timezone: UTC
##   H2O cluster version:     3.44.0.3
##   H2O cluster version age:  1 year and 10 months
##   H2O cluster name:        H2O_started_from_R_evan_owd174
##   H2O cluster total nodes:  1
##   H2O cluster total memory: 10.59 GB
##   H2O cluster total cores:  16
##   H2O cluster allowed cores: 16
##   H2O cluster healthy:      TRUE
##   H2O Connection ip:        localhost
##   H2O Connection port:      54321
##   H2O Connection proxy:     NA
##   H2O Internal Security:    FALSE
```

Using H2O

Note: Initializing H2O might throw an error in your system in the case where you don't have Jdk of 64 bit. If such issue arises, please install latest Jdk of 64 bits here, it should work without issue afterward.

Using H2O

The `h2o.init()` command is pretty smart and does a lot of work. At first, it looks for any active H2O instance before starting a new one and then starts a new one when instance are not present.

Using H2O

It does have arguments which helps to accommodate resources to the H2O instance frequently used are:

- ▶ **nthreads:** By default, the value of nthreads will be -1 which means the instance can use all the cores of the CPU, we can set the number of cores utilized by passing the value to the argument.
- ▶ **max_mem_size:** By passing a value to this argument you can restrict the maximum memory allocated to the instance. Its od string type can pass an argument as '2g' or '2G' for 2 GBs of memory, same when you want to allocate in MBs.

Using H2O

You can access the flow by typing `http://localhost:54321` in your browser.

Using H2O

Flow is the name of the web interface that is part of H2O which does not require any extra installations which is written in CoffeeScript (a JavaScript like language). You can use it for doing the following things:

- ▶ Upload data directly
- ▶ View data uploaded by the client
- ▶ Create models directly
- ▶ View models created by you or your client
- ▶ view predictions
- ▶ Run predictions directly

Untitled Flow









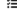






CS

assist

85ms

? Assistance

Routine	Description
 importFiles	Import file(s) into H ₂ O
 importSqlTable	Import SQL table into H ₂ O
 getFrames	Get a list of frames in H ₂ O
 splitFrame	Split a frame into two or more frames
 mergeFrames	Merge two frames into one
 getModels	Get a list of models in H ₂ O
 getGrids	Get a list of grid search results in H ₂ O
 getPredictions	Get a list of predictions in H ₂ O
 getJobs	Get a list of jobs running in H ₂ O
 runAutoML	Automatically train and tune many models
 buildModel	Build a model
 importModel	Import a saved model
 predict	Make a prediction

OUTLINE FLOWS CLIPS **HELP**

Help



Using Flow for the first time?

[Quickstart Videos](#)

Or, [view example Flows](#) to explore and learn H₂O.

STAR H₂O ON GITHUB!

GENERAL

- [Flow Web UI ...](#)
- ... [Importing Data](#)
- ... [Building Models](#)
- ... [Making Predictions](#)
- ... [Using Flows](#)
- ... [Troubleshooting Flow](#)

EXAMPLES

Flow packs are a great way to explore and learn H₂O. Try out these Flows and run them in your browser.

[Browse installed packs...](#)H₂O REST API

- [Routes](#)
- [Schemas](#)

H2O AutoML

AutoML helps in automatic training and tuning of many models within a user-specified time limit.

The current version of AutoML function can train and cross-validate a Random Forest, an Extremely-Randomized Forest, a random grid of Gradient Boosting Machines (GBMs), a random grid of Deep Neural Nets, and then trains a Stacked Ensemble using all of the models.

H2O AutoML

When we say AutoML, it should cater to the aspects of data preparation, Model generation, and Ensembles and also provide few parameters as possible so that users can perform tasks with little confusion.

AutoML inputs required arguments **y** and the **training_frame**, with the **x** and **validation frame** as optional arguments. The user can also configure values for **max_runtime_sec** and **max_models**.

H2O AutoML

Additional optional parameters include:

- ▶ `leaderboard_frame`
- ▶ `nfolds`
- ▶ `fold_columns`
- ▶ `weights_column`
- ▶ `ignored_columns`
- ▶ `stopping_metric`
- ▶ `sort_metric`

H2O Kmeans on the iris data

```
iris_h2o <- as.h2o(iris)
```

```
## |
```

```
iris_h2o['Species'] <- as.factor(iris_h2o['Species'])
```

```
predictors <- colnames(iris_h2o)[-length(iris_h2o)]
```

```
iris_splits <- h2o.splitFrame(data = iris_h2o,  
                             ratios = 0.7, seed = 1234)
```

```
train <- iris_splits[[1]]
```

```
valid <- iris_splits[[2]]
```

H2O Kmeans on the iris data

```
kmeans_model <- h2o.kmeans(training_frame = train,  
                             x = predictors, k = 3,  
                             seed = 1)
```

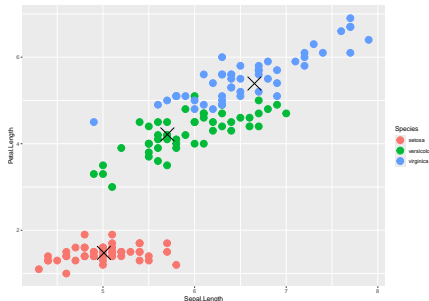
```
##      |
```

```
centers <- h2o.centers(kmeans_model)  
centers
```

```
##      sepallength  sepalwidth  petallength  petalwidth  
## 1      5.702857    2.637143    4.214286    1.340000  
## 2      6.653333    3.051111    5.391111    1.937778  
## 3      5.010000    3.436667    1.476667    0.250000
```


H2O Kmeans on the iris data

```
iris %>% ggplot(aes(Sepal.Length, Petal.Length)) +  
  geom_point(aes(col=Species), size=5) +  
  geom_point(aes(sepallength, petallength),  
             col=1, size=10, pch = 4, data=centers)
```



H2O AutoML on the iris data

```
iris_automl <- h2o.automl(x = predictors, y = "Species",
  training_frame = train,
  max_runtime_secs = 20, seed = 1,
  validation_frame = valid)
```

```
## |
## 09:53:15.514: User specified a validation frame with cross-validation still enabled. Please note that the models will still be vali
## 09:53:15.515: AutoML: XGBoost is not available; skipping it.
## 09:53:15.806: _min_rows param, The dataset size is too small to split for min_rows=100.0: must have at least 200.0 (weighted) rows,
```

```
iris_automl
```

```
## AutoML Details
## =====
## Project Name: AutoML_3_20251021_95315
## Leader Model ID: DeepLearning_grid_1_AutoML_3_20251021_95315_model_1
## Algorithm: deeplearning
##
## Total Number of Models Trained: 59
## Start Time: 2025-10-21 09:53:16 UTC
## End Time: 2025-10-21 09:53:36 UTC
## Duration: 20 s
##
## Leaderboard
## =====
##
```

model id	mean per class error
----------	----------------------

H2O AutoML on the iris data

```
lb <- h2o.get_leaderboard(iris_automl)
head(lb)
```

```
##                                     model_id mean_per_class_error
## 1   DeepLearning_grid_1_AutoML_3_20251021_95315_model_1          0.01587302
## 2                                     GLM_1_AutoML_3_20251021_95315          0.02595453
## 3   GBM_grid_1_AutoML_3_20251021_95315_model_34          0.02595453
## 4   GBM_grid_1_AutoML_3_20251021_95315_model_32          0.02595453
## 5   GBM_grid_1_AutoML_3_20251021_95315_model_12          0.02595453
## 6 StackedEnsemble_BestOfFamily_2_AutoML_3_20251021_95315          0.02595453
##      logloss      rmse      mse
## 1 0.07784368 0.1371276 0.01880398
## 2 0.05974883 0.1381892 0.01909625
## 3 0.11595687 0.1664132 0.02769336
## 4 0.12050457 0.1681316 0.02826823
## 5 0.11449570 0.1667693 0.02781201
## 6 0.08217182 0.1521377 0.02314588
```

H2O AutoML on the iris data

```
iris_automl@leader
```

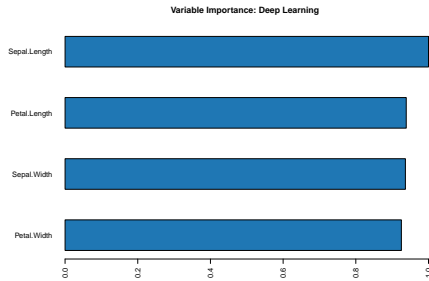
```
## Model Details:
## =====
##
## H2OMultinomialModel: deeplearning
## Model ID: DeepLearning_grid_1_AutoML_3_20251021_95315_model_1
## Status of Neuron Layers: predicting Species, 3-class classification, multinomial distribution, CrossEntropy loss, 803 weights/biases
## layer units      type dropout    l1      l2 mean_rate rate_rms
## 1      1      4      Input 15.00 %      NA      NA      NA      NA
## 2      2      100 RectifierDropout 10.00 % 0.000000 0.000000 0.002010 0.004091
## 3      3      3      Softmax      NA 0.000000 0.000000 0.001685 0.000812
## momentum mean_weight weight_rms mean_bias bias_rms
## 1      NA      NA      NA      NA      NA
## 2 0.000000 -0.001885 0.141170 0.497611 0.032311
## 3 0.000000 -0.001110 0.499156 -0.006928 0.035866
##
##
## H2OMultinomialMetrics: deeplearning
## ** Reported on training data. **
## ** Metrics reported on full training frame **
##
## Training Set Metrics:
## =====
##
## Extract training frame with 'h2o.getFrame("AutoML_3_20251021_95315_training_RTMP_sid_b510_4")'
```

H2O AutoML on the iris data

```
h2o.varimp(iris_automl@leader)
```

```
## Variable Importances:
##      variable relative_importance scaled_importance percentage
## 1 Sepal.Length      1.000000      1.000000      0.263120
## 2 Petal.Length      0.938738      0.938738      0.247000
## 3 Sepal.Width       0.936442      0.936442      0.246396
## 4 Petal.Width       0.925374      0.925374      0.243484
```

```
h2o.varimp_plot(iris_automl@leader)
```



H2O AutoML on the iris data

```
pred <- h2o.predict(iris_automl@leader, valid)
```

```
##      |
```

```
pred
```

```
##      predict      setosa versicolor      virginica
## 1  setosa 0.9952966 0.004228243 0.0004751776
## 2  setosa 0.9890041 0.010706118 0.0002898031
## 3  setosa 0.9850540 0.014291734 0.0006542809
## 4  setosa 0.9967064 0.001697642 0.0015959823
## 5  setosa 0.9947530 0.004704620 0.0005424038
## 6  setosa 0.9934720 0.003539749 0.0029882760
##
## [40 rows x 4 columns]
```

H2O AutoML on the prostate cancer dataset

```
prostate_path <- system.file("extdata", "prostate.csv",
                             package = "h2o")
prostate <- h2o.importFile(path = prostate_path, header = TRUE)
```

```
## |
```

```
y <- "CAPSULE"
prostate[,y] <- as.factor(prostate[,y])
knitr::kable(prostate)
```

ID	CAPSULE	AGE	RACE	DPROS	DCAPS	PSA	VOL	GLEASON
1	0	65	1	2	1	1.40	0.00	6
2	0	72	1	3	2	6.70	0.00	7
3	0	70	1	1	2	4.90	0.00	6
4	0	76	2	2	1	51.20	20.00	7
5	0	69	1	1	1	12.30	55.90	6
6	1	71	1	3	2	3.30	0.00	8
7	0	68	2	4	2	31.90	0.00	7
8	0	61	2	4	2	66.70	27.20	7
9	0	69	1	1	1	3.90	24.00	7
10	0	68	2	1	2	13.00	0.00	6
11	1	68	2	4	2	4.00	0.00	7
12	1	72	1	2	2	21.20	0.00	7
13	1	72	1	4	2	22.70	0.00	9

H2O AutoML on the prostate cancer dataset

```
aml <- h2o.automl(y = y, training_frame = prostate,
  max_runtime_secs = 20, seed = 1)
```

```
## |
## 09:53:38.39: AutoML: XGBoost is not available; skipping it. |
```

```
lb <- h2o.get_leaderboard(aml)
head(lb)
```

```
##                               model_id      auc  logloss
## 1 StackedEnsemble_BestOfFamily_6_AutoML_4_20251021_95338 0.8160433 0.5159649
## 2               GBM_grid_1_AutoML_4_20251021_95338_model_11 0.8144309 0.5289923
## 3 StackedEnsemble_BestOfFamily_4_AutoML_4_20251021_95338 0.8140854 0.5206670
## 4   StackedEnsemble_AllModels_2_AutoML_4_20251021_95338 0.8135095 0.5173276
## 5               GBM_grid_1_AutoML_4_20251021_95338_model_4 0.8123290 0.5222005
## 6   StackedEnsemble_AllModels_1_AutoML_4_20251021_95338 0.8091331 0.5216481
##      aucpr mean_per_class_error      rmse      mse
## 1 0.7392318          0.2324004 0.4137657 0.1712020
## 2 0.7171983          0.2298379 0.4190206 0.1755782
## 3 0.7291787          0.2391523 0.4156836 0.1727929
## 4 0.7348310          0.2424203 0.4152257 0.1724124
## 5 0.7152434          0.2413550 0.4169814 0.1738735
## 6 0.7373510          0.2436296 0.4172881 0.1741293
```


Session Info

```
sessionInfo()
```

```
## R version 4.5.1 (2025-06-13)
## Platform: aarch64-apple-darwin20
## Running under: macOS Sequoia 15.6.1
##
## Matrix products: default
## BLAS:   /Library/Frameworks/R.framework/Versions/4.5-arm64/Resources/lib/libRblas.0.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/4.5-arm64/Resources/lib/libRlapack.dylib; LAPACK version 3.12.1
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## time zone: America/Denver
## tzcode source: internal
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods   base
##
## other attached packages:
## [1] h2o_3.44.0.3  lubridate_1.9.4 forcats_1.0.0  stringr_1.5.1
## [5] dplyr_1.1.4   purrr_1.1.0    readr_2.1.5    tidyr_1.3.1
## [9] tibble_3.3.0  tidyverse_2.0.0 caret_7.0-1     lattice_0.22-7
## [13] ggplot2_3.5.2
##
## loaded via a namespace (and not attached):
## [1] rstan_2.21.0  rstanarm_2.21.0  rstanarm_2.21.0  rstanarm_2.21.0
```