

Machine Learning With H2O — Hands-On Guide for Data Scientists

by Pavel Pscheidl · Jun. 27, 18 · AI Zone · Analysis

H2O is the world's number one machine learning platform. It is an open-source software, and the H2O-3 GitHub repository is available for anyone to start hacking. This hands-on guide aims to explain the basic principles behind H2O and get you as a data scientist started as quickly as possible in the most simple way. The rest is just machine learning.

After reading this guide, you'll be able to:

- understand which basic problems H2O solves and why,
- play with H2O — explore data and create and tune models,
- see beyond the horizon. Understand where H2O can take you.

As a data scientist, you're most likely to use R and/or Python. H2O integrates with both. Interestingly, H2O makes it easy to seamlessly switch Python, R, and other data science tools while still working on the same project. This allows data scientists to interact more easily, as well as use the best tool for the job, but the possibilities do not stop there. H2O also offers its own web-based interface named Flow. By means of Flow, data scientists are able to import, explore, and modify datasets, play with models, verify models performances, and much more. Flow is beautiful and a quick way to do machine learning. Flows can be saved and given to other data scientists, making cooperation easy.

H2O respects habits of data scientists and does not get into their way. Using Python, data scientists are familiar with Pandas, Scikit, NumPy and others. H2O's syntax is very similar to those. H2O is able to work directly with Pandas' data structures, as well as being compatible with NumPy's arrays and primitive Python lists and collections. H2O follows the same pattern with R, respecting the naming and syntax R developers are used to.

Getting Started

The preparations before take-off are short. H2O is extremely easy to start with. All it takes is a common laptop to get started. Once H2O is installed, it is very easy and convenient to import a dataset and create a model out of it.

In the examples below, the famous Airlines Delay dataset is used. There is no need to download it, as H2O is going to take care of downloading the dataset for you. The dataset is very intuitive to work with, however if

We've Updated Our Site Policies.

you're unfamiliar with it or simply want to know more, visit Kaggle website for a brief description of each dataset. **We have recently updated our terms of service and privacy policy.**

For additional information, visit:

Once fluent in R, working with H2O in Python's environment is easy. And it always works the other way around, as the APIs are very similar, while respecting the target platforms.

CLOSE

Getting Started in R

Open up R CLI (by typing `R` in terminal on most systems) or start R-Studio. It takes just a few lines to install H₂O in R.

Before installing H₂O itself, H₂O requires two packages: `RCurl` and `jsonlite`. Install those by entering the following command into R console.

Installation

```
1 install.packages("RCurl","jsonlite")
```

After `RCurl` and `jsonlite` are installed, one last step is to install H₂O itself. Installation of latest stable release is done as demonstrated in the following snippet. During the installation, a one-time download of the H₂O backend containing all the algorithms and computing know-how will occur.

```
1 install.packages("h2o", type="source", repos=c("http://h2o-release.s3.amazonaws.com"))
```

That's it. H₂O is now installed and ready to be used. As a first step, it is required to tell R to import the H₂O library with `library(h2o)` command. Once the library is imported, instruct H₂O to start itself by calling `h2o.init()`. Both commands are placed in the following code snippet for clarity.

```
1 library(h2o)
2 h2o.init()
```

The `h2o.init()` command is pretty smart and does a lot of things. First, an attempt is made to search for an existing H₂O instance being started already before starting a new one. When none is found automatically or specified manually with argument available, a new instance of H₂O is started. As this is a fresh installation and it is highly unlikely there is an instance of H₂O already running in your environment, a new instance is started right away. During startup, H₂O is going to print some useful information. The R version it is running on, H₂O's version, and how to connect to H₂O's Flow interface or where error logs reside, just to name a few. As usual, an example of H₂O's output during startup is to be found below this text in the snippet.

```
1 > h2o.init()
2
3 H2O is not running yet, starting it now...
4
5 Note: In case of errors look at the following log files:
6     /tmp/RtmpYs7uDC/h2o_pavel_started_from_r.out
7     /tmp/RtmpYs7uDC/h2o_pavel_started_from_r.err
8
9 java version "1.8.0_171"
10 Java(TM) SE Runtime Environment (build 1.8.0_171-b11)
11 Java HotSpot(TM) 64-Bit Server VM (build 25.171-b11, mixed mode)
12
```

We've Updated Our Site Policies. Connection successful!

We have recently updated our terms of service and privacy policy.

For additional information, visit:
H2O cluster uptime: 3 seconds 44 milliseconds

<https://dzone.com/pages/tos> <https://dzone.com/pages/privacy>

H2O data parsing timezone: UTC

CLOSE

```

18  h2o.data.parsing.timezone:  UTC
19  H2O cluster version:      3.20.0.2
20  H2O cluster version age:   8 days
21  H2O cluster name:         H2O_started_from_R_pavel_yuw261
22  H2O cluster total nodes:   1
23  H2O cluster total memory:  5.21 GB
24  H2O cluster total cores:   8
25  H2O cluster allowed cores: 8
26  H2O cluster healthy:      TRUE
27  H2O Connection ip:        localhost
28  H2O Connection port:      54321
29  H2O Connection proxy:      NA
30  H2O Internal Security:     FALSE
31  H2O API Extensions:        XGBoost, Algos, AutoML, Core V3, Core V4
32  R Version:                 R version 3.4.4 (2018-03-15)

```

Data Import

Let's import a dataset and train a model on it very quickly!

```
1  airlinesTrainData <- h2o.importFile("https://s3.amazonaws.com/h2o-airlines-unpacked/airlines_train.csv")
```

H₂O will automatically download the dataset and parse it. It will also try to guess the datatype of each column automatically. H₂O does a great job at datatype recognition, however, each decision can be overridden manually by the user, if required. The imported dataset can also be given a name using `destination_frame` argument. For example, `h2o.importFile("https://s3.amazonaws.com/h2o-airlines-unpacked/allyears2k.csv", destination_frame='airlines_train')` imports the very same dataset with airplane delays that can be further addressed by name `airlines_train`, even from other interfaces like Python, another R console, Flow, Java, or direct API calls. If no name is provided, H₂O will generate an artificial name. Simply put, an imported dataset is called `Frame` in H₂O. List of frames can be shown by using the `h2o.ls()` function. An example output of calling `h2o.ls()` function can be found in the following code snippet. The very first record in the example is a named frame, and the second one is a frame name generated automatically by H₂O.

```

1  > h2o.ls()
2
3  key
4  1      airlines_train
5  2 allyears2k.hex_sid_ae6f_1

```

A preview of the data imported can be displayed with by typing the variable pointing to the `H2OFrame`, in this case `airlinesTrainData`.

```

1  > airlinesTrainData
2  Year Month DayOfMonth DayOfWeek DepTime CRSDepTime ArrTime CRSArrTime UniqueCarrier
3  1 1987    10         14         3      741          730      912          849
4  2 1987    10         15         4      729          730      903          849
5  3 1987    10         16         4      741          730      918          849
6  4 1987    10         18         7      729          730      847
7  5 1987    10         19         1      749          730      922
8  6 1987    10         21         3      728          730      848          849
9  Origin Dest Distance TaxiIn TaxiOut Cancelled CancellationCode Diverted CarrierDe

```

We've Updated Our Site Policies.

We have recently updated our terms of service and privacy policy.

For additional information, visit:

<https://dzone.com/pages/tos> | <https://dzone.com/pages/privacy>

CLOSE

10	1	SAN	SFO	447	NaN	NaN	0	NA	0
11	2	SAN	SFO	447	NaN	NaN	0	NA	0
12	3	SAN	SFO	447	NaN	NaN	0	NA	0
13	4	SAN	SFO	447	NaN	NaN	0	NA	0
14	5	SAN	SFO	447	NaN	NaN	0	NA	0
15	6	SAN	SFO	447	NaN	NaN	0	NA	0
16		IsDepDelayed							
17	1							YES	
18	2							NO	
19	3							YES	
20	4							NO	
21	5							YES	
22	6							NO	

Model Training

On top of the data imported, a model can be built quickly. There are many algorithms available in H₂O. For the purpose of this tutorial, a widely known Gradient Boosting Machines method will be used. Let's train a model that is able to predict if the plane arrives late based on month, day of week, and distance the plane has to travel before reaching its destination. By invoking `h2o.gbm(...)`, H₂O will run a gradient boosting algorithm on the data. There are many variables to play with that each and every data scientist can explore on his/her own. Overriding the default hyperparameters would only make this tutorial more complicated. H₂O only needs to know three things:

- predictor columns,
- response variable column,
- training frame — a dataset to train the model on.

Nothing more. It is even able to guess the distribution of the response variable, even though as stated before, everything can be overridden manually by the data scientist, if required. After a model is trained, basic information about the model can be shown just by typing the name of the variable pointing to the trained model, in this case `gbmModel`.

```

1 > gbmModel <- h2o.gbm(x=c("Month", "DayOfWeek", "Distance"), y="IsArrDelayed", tra:
2 |=====
3 > gbmModel
4 Model Details:
5 =====
6
7 H2OBinomialModel: gbm
8 Model ID: GBM_model_R_1529834562303_470
9 Model Summary:

```

We've Updated Our Site Policies.

We have recently updated our terms of service and privacy policy.

For additional information, visit:

<https://dzone.com/pages/tos> | <https://dzone.com/pages/privacy>

** Reported on training data. **

CLOSE

```

16
17 MSE: 0.2349663
18 RMSE: 0.4847332
19 LogLoss: 0.6609351
20 Mean Per-Class Error: 0.4892883
21 AUC: 0.6237765
22 Gini: 0.2475531
23
24 Confusion Matrix (vertical: actual; across: predicted) for F1-optimal threshold:
25      NO    YES    Error    Rate
26 NO      604 18933 0.969084 =18933/19537
27 YES     232 24209 0.009492  =232/24441
28 Totals 836 43142 0.435786 =19165/43978
29
30 Maximum Metrics: Maximum metrics at their respective thresholds
31      metric threshold    value idx
32 1      max f1    0.431588 0.716423 368
33 2      max f2    0.355217 0.862241 395
34 3      max f0point5 0.513454 0.633681 278
35 4      max accuracy 0.511993 0.594661 279
36 5      max precision 0.972534 1.000000 0
37 6      max recall    0.347469 1.000000 397
38 7      max specificity 0.972534 1.000000 0
39 8      max absolute_mcc 0.605839 0.175948 151
40 9      max min_per_class_accuracy 0.539888 0.582464 234
41 10 max mean_per_class_accuracy 0.547177 0.584595 225
42
43 Gains/Lift Table: Extract with `h2o.gainsLift(<model>, <data>)` or `h2o.gainsLift(

```

Overall, this model is not expected to perform very well, given the huge error rate to be observed in the confusion matrix. By playing with different GBM hyperparameters and including different predictors into the model, much better results can be achieved. As a data scientist, the task of making the model perform better is easy for you, and that's certain. In the `h2o.` package, there are many additional functions to work with the model and help a data scientist understand what happened during the training phase and much more. As an example, the `h2o.varimp` function shows importances of variables (relative, percentage) taken into account in the model. As you begin exploring H₂O, the reference guide will guide you through all the H₂O's functionality.

```

1 > h2o.varimp(gbmModel)
2 Variable Importances:
3      variable relative_importance scaled_importance percentage
4 1 Distance          1379.994019          1.000000    0.501313
5 2 Month              970.739441          0.703437    0.352643
6 3 ...                ...              ...          ...
7 4 ...                ...              ...          ...
8 5 ...                ...              ...          ...

```

We've Updated Our Site Policies.

We have recently updated our terms of service and privacy policy. When it comes to data science, it's important to have a clear understanding of the data you are working with. Of course, according to this very basic model given the default parameters. A delayed response is not acceptable. For additional information, visit <https://dzone.com/pages/tos> or <https://dzone.com/pages/privacy>. We support XGBoost. It is trivial to swap GBM for XGBoost in this phase and see how the model changes with default hyperparameters:

CLOSE

```
1 xgBoostModel <- h2o.xgboost(x=c("Month", "DayOfWeek", "Distance"), y="IsArrDelayed")
```

Prediction

Prediction is very simple as well. Calling `h2o.predict(model, data)`, where `model` is the variable pointing to the model trained and `data` is the `H2OFrame` with data to do the prediction on. To test the prediction is functional in a very simple way, let's use the `gbmModel` and let it predict the original training dataset.

```
1 > h2o.predict(gbmModel, airlinesTrainData)
2 |=====|
3 predict      NO      YES
4 1      YES 0.1419937 0.8580063
5 2      YES 0.1015739 0.8984261
6 3      YES 0.2036055 0.7963945
7 4      YES 0.1239904 0.8760096
8 5      YES 0.1384360 0.8615640
9 6      YES 0.1419937 0.8580063
10
11 [43978 rows x 3 columns]
```

The prediction is not very accurate in case there was no delay. This is expected, as the model is very basic. Of course, the confusion matrix seen earlier in this tutorial gave out the information about such “bad” performance beforehand.

Getting Started in Python

In order to get started in Python, only few lines of code are required. A common way of installing dependencies in python is `pip` or `anaconda`. In this tutorial, `pip` is preferred due to most users being familiar with it. If you'd like to use Conda, please follow the tutorial in H2O documentation. Python 2.7 up to Python 3.x are supported. The differences are minimal and this guide should work on both versions.

Installation

Before installing H2O itself, a few dependencies are required. Please install them using `pip install`. On some systems, super-user privileges may be required. If so, adding `sudo` before `pip install` will solve the problem.

```
1 pip install requests
2 pip install tabulate
3 pip install scikit-learn
4 pip install colorama
5 pip install future
```

Once the dependencies required are installed, one last step is to install H2O itself. Installation of the latest stable release is done as demonstrated in the following snippet. During the installation, a one-time download of the

We've Updated Our Site Policies

H2O becomes a part of the algorithmic computing know-how will occur.

We have recently updated our terms of service and privacy policy.

```
1 pip install --r http://h2o-release.s3.amazonaws.com/h2o/latest_stable
```

For additional information, visit:

That's it. H2O is now installed and ready to be used. As a first step, it is required to tell Python to import the H2O module with `import h2o` command. Once the module is imported, instruct H2O to start itself by calling

CLOSE

`h2o.init()` . Both commands are placed in the following code snippet for clarity. The process of setup is very similar to R.

```
1 import h2o
2 h2o.init()
```

The `h2o.init()` command is pretty smart and does a lot of things. First, an attempt is made to search for an existing H₂O instance being started already, before starting a new one. When none is found automatically or specified manually with argument available, a new instance of H₂O is started. As this is a fresh installation and it is highly unlikely there is an instance of H₂O already running in your environment, a new instance is started right away. During startup, H₂O is going to print some useful information. Version of the Python it is running on, H₂O's version, how to connect to H₂O's Flow interface or where error logs reside, just to name a few. As usual, an example of H₂O's output during startup is to be found below this text in the snippet.

```
1 >>> h2o.init()
2 Checking whether there is an H2O instance running at http://localhost:54321..... no
3 Attempting to start a local H2O server...
4   Java Version: java version "1.8.0_171"; Java(TM) SE Runtime Environment (build 1
5   Starting server from /usr/local/lib/python2.7/dist-packages/h2o/backend/bin/h2o.
6   Ice root: /tmp/tmp7R80vB
7   JVM stdout: /tmp/tmp7R80vB/h2o_pavel_started_from_python.out
8   JVM stderr: /tmp/tmp7R80vB/h2o_pavel_started_from_python.err
9   Server is running at http://127.0.0.1:54321
10 Connecting to H2O server at http://127.0.0.1:54321... successful.
11 versionFromGradle='3.19.0',projectVersion='3.19.0.99999',branch='pavel_pubdev-5336
12 -----
13 H2O cluster uptime:          01 secs
14 H2O cluster timezone:       Europe/Prague
15 H2O data parsing timezone:   UTC
16 H2O cluster version:        3.19.0.99999
17 H2O cluster version age:     3 months and 5 days
18 H2O cluster name:           H2O_from_python_pavel_oy9g50
19 H2O cluster total nodes:     1
20 H2O cluster free memory:     5.207 Gb
21 H2O cluster total cores:     8
22 H2O cluster allowed cores:   8
23 H2O cluster status:         accepting new members, healthy
24 H2O connection url:         http://127.0.0.1:54321
25 H2O connection proxy:
26 H2O internal security:      False
27 H2O API Extensions:         XGBoost, Algos, AutoML, Core V3, Core V4
28 Python version:             2.7.15 candidate
29 -----
```

We've Updated Our Site Policies.

Data Import We have recently updated our terms of service and privacy policy.

Let's import a dataset and train a model on it very quickly.

For additional information, visit: <https://dzone.com/pages/terms> | <https://dzone.com/pages/privacy>

CLOSE

```
1 airlines_train_data = h2o.import_file("https://s3.amazonaws.com/h2o-airlines-unpack
```


H₂O will automatically download the dataset and parse it. It will also try to guess the datatype of each column automatically. H₂O does a great job at datatype recognition, however, each decision can be overridden manually by the user, if required. The imported dataset can also be given a name using `destination_frame` argument. For example, `h2o.import_file("https://s3.amazonaws.com/h2o-airlines-unpacked/allyears2k.csv", destination_frame='airlines_train')` imports the very same dataset with airplane delays that can be further addressed by name `airlines_train`, even from other interfaces like Python, another R console, Flow, Java, or direct API calls. If no name is provided, H₂O will generate an artificial name. Simply put, an imported dataset is called `Frame` in H₂O. List of frames can be shown by using the `h2o.ls()` function. An example output of calling `h2o.ls()` function can be found in the following code snippet. The very first record in the example is a named frame, the second one is a frame name generated automatically by H₂O.

```

1  >>> h2o.ls()
2
3      key
4  0  airlines_train
5  1  allyears2k.hex

```

A preview of the data imported can be displayed with by typing the variable pointing to the H₂O Frame, in this case `airlines_train_data`.

```

1  >>> airlines_train_data
2
3      Year Month DayOfMonth DayOfWeek DepTime CRSDepTime ArrTime CRSArrTime UniqueCarr:
4  1 1987      10           14          3      741          730      912          849
5  2 1987      10           15          4      729          730      903          849
6  3 1987      10           17          6      741          730      918          849
7  4 1987      10           18          7      729          730      847          849
8  5 1987      10           19          1      749          730      922          849
9  6 1987      10           21          3      728          730      848          849
10
11     Origin Dest Distance TaxiIn TaxiOut Cancelled CancellationCode Diverted CarrierD:
12 1      SAN  SFO      447    NaN    NaN          0              NA          0
13 2      SAN  SFO      447    NaN    NaN          0              NA          0
14 3      SAN  SFO      447    NaN    NaN          0              NA          0
15 4      SAN  SFO      447    NaN    NaN          0              NA          0
16 5      SAN  SFO      447    NaN    NaN          0              NA          0
17 6      SAN  SFO      447    NaN    NaN          0              NA          0
18
19     IsDepDelayed
20 1              YES
21 2              NO
22 3              YES
23 4              NO
24 5              YES
25 6              NO

```

Model Training

We've Updated Our Site Policies.

On top of the data imported, a model can be built quickly. There are many algorithms available in H₂O. For the purpose of this tutorial, a widely known Gradient Boosting Machines method will be used. For additional information, visit <https://h2oai.com/pages/privacy> before reaching its destination. GBM resides on the `h2o.ensemble.gbm` package. The first step is to import `H2OGradientBoostingEstimator` to avoid the need for typing a fully qualified name in future.

CLOSE


```
1 from h2o.estimators.gbm import H2OGradientBoostingEstimator
```

First, it is required to construct a new GBM estimator instance by calling `gbm_model = H2OGradientBoostingEstimator()` constructor. By invoking `gbm_model.train(...)`, H2O will run a gradient boosting algorithm on the data. There are many variables to play with, and each and every data scientist can explore on his/her own. Overriding the default hyperparameters would only make this tutorial more complicated. H2O only needs to know three things:

- predictor columns,
- response variable column,
- training frame — a dataset to train the model on.

```
1 gbm_model = H2OGradientBoostingEstimator()
2 gbm_model.train(x = ["Month", "DayOfWeek", "Distance"], y = "IsArrDelayed", traini
```

Nothing more. It is even able to guess the distribution of the response variable, even though as stated before, everything can be overridden manually by the data scientist, if required. After a model is trained, basic information about the model can be shown just by typing the name of the variable pointing to the trained model, in this case `gbm_model`. In the next figure, there is a demonstration of previous steps regarding model training, with H₂O's output included.

```

1 >>> from h2o.estimators.gbm import H2OGradientBoostingEstimator
2 >>> gbm_model = H2OGradientBoostingEstimator();
3 >>> gbm_model.train(x = ["Month", "DayOfWeek", "Distance"], y = "IsArrDelayed", tra
4 gbm Model Build progress: |███████████████████████████████████████████████████
5 >>> gbm_model
6 Model Details
7 =====
8 H2OGradientBoostingEstimator : Gradient Boosting Machine
9 Model Key: GBM_model_python_1529844691141_1
10
11
12 ModelMetricsBinomial: gbm
13 ** Reported on train data. **
14
15 MSE: 0.234966307798
16 RMSE: 0.484733233643
17 LogLoss: 0.660935092712
18 Mean Per-Class Error: 0.41540494848
19 AUC: 0.623776525749
20 Gini: 0.247553051497
21 Confusion Matrix (Act/Pred) for max f1 @ threshold = 0.43158830279:

```

We've Updated Our Site Policies.

~~We have recently updated our terms of service and privacy policy.~~

For additional information, visit: <https://www.noaa.gov/data/monitoring-assessments/air-quality/air-quality-models>

	YES	NO	OR	95% CI
1.25	232	24209	0.0095	(232.0/24441.0)

<https://dzone.com/pages/tos> | <https://dzone.com/pages/privacy>

26	Total	836	43142	0.4358	(19165.0/43978.0)
----	-------	-----	-------	--------	-------------------

CLOSE

Maximum Metrics: Maximum metrics at their respective thresholds

metric	threshold	value	idx
max f1	0.431588	0.716423	368
max f2	0.355217	0.862241	395
max f0point5	0.513454	0.633681	278
max accuracy	0.511993	0.594661	279
max precision	0.972534	1	0
max recall	0.347469	1	397
max specificity	0.972534	1	0
max absolute_mcc	0.605839	0.175948	151
max min_per_class_accuracy	0.539888	0.582464	234
max mean_per_class_accuracy	0.547177	0.584595	225

Gains/Lift Table: Avg response rate: 55,58 %

group	cumulative_data_fraction	lower_threshold	lift	cumulative_lift
1	0.0100732	0.897096	1.70593	1.70593
2	0.0221247	0.864273	1.64658	1.6736
3	0.0302651	0.83021	1.54302	1.63848
4	0.040111	0.753167	1.40873	1.58208
5	0.0514803	0.696079	1.34952	1.53072
6	0.103688	0.649372	1.30093	1.41502
7	0.150348	0.618087	1.225	1.35605
8	0.20101	0.600111	1.17103	1.30942
9	0.3004	0.578542	1.07854	1.23303
10	0.400632	0.557964	1.03234	1.18282
11	0.500568	0.540445	1.00633	1.14758
12	0.606894	0.525153	0.973175	1.11703
13	0.701874	0.508853	0.921	1.0905
14	0.800673	0.492089	0.869653	1.06325
15	0.90261	0.465819	0.80997	1.03465
16	1	0.279442	0.678906	1

Scoring History:

timestamp	duration	number_of_trees	training_rmse	training_loss
2018-06-24 15:03:17	0.087 sec	0.0	0.49688163904	0.6869169
2018-06-24 15:03:17	0.368 sec	1.0	0.495487418342	0.6841010
2018-06-24 15:03:18	0.487 sec	2.0	0.494360221437	0.6818030
2018-06-24 15:03:18	0.659 sec	3.0	0.493435134409	0.6798920
2018-06-24 15:03:18	0.659 sec	4.0	0.492679667171	0.6783050
2018-06-24 15:03:19	1.667 sec	46.0	0.485057007111	0.6610100
2018-06-24 15:03:19	1.710 sec	48.0	0.484945381131	0.6613770
2018-06-24 15:03:19	1.710 sec	48.0	0.484866365041	0.6612160

We've Updated Our Site Policies.

We have recently updated our terms of service and privacy policy.

For additional information, visit:

<https://dzone.com/pages/tos> <https://dzone.com/pages/privacy>

CLOSE

```

73      2018-06-24 15:03:17 1.710 sec 48.0      0.484800303041 0.661216
74      2018-06-24 15:03:19 1.734 sec 49.0      0.484810585256 0.661098
75      2018-06-24 15:03:19 1.760 sec 50.0      0.484733233643 0.660935
76
77 See the whole table with table.as_data_frame()
78 Variable Importances:
79 variable      relative_importance      scaled_importance      percentage
80 -----
81 Distance      1379.99      1      0.501313
82 Month          970.739      0.703437      0.352643
83 DayOfWeek      402.024      0.291323      0.146044

```

Overall, this model is not expected to perform very well, given the huge error rate to be observed in the confusion matrix. By playing with different GBM hyperparameters and including different predictors into the model, much better results can be achieved. As a data scientist, the task of making the model perform better is easy for you, that's certain. To get detailed information about model and its scoring history, invoke the `print(gbm_model)` command. The output contains table with importances of variables (relative, percentage, scaled) taken into account in the model. Also, a detailed scoring history is available, as well as basic measures like mean squared error (MSE). As you begin exploring H₂O, the reference guide will guide you through all the H₂O's functionality. A shortened example of a detailed view on GBM model is to be found in the next figure. Some of the text was omitted.

```

1  print(gbm_model)
2  Model Details
3  =====
4  H2OGradientBoostingEstimator : Gradient Boosting Machine
5  Model Key:  GBM_model_python_1529844691141_1
6
7
8  ModelMetricsBinomial: gbm
9  ** Reported on train data. **
10
11 Scoring History:
12      timestamp      duration      number_of_trees      training_rmse      training_
13  ---
14      2018-06-24 15:03:17 0.087 sec 0.0      0.49688163904 0.686916
15      2018-06-24 15:03:17 0.368 sec 1.0      0.495487418342 0.684101
16      2018-06-24 15:03:18 0.487 sec 2.0      0.494360221437 0.681803
17      2018-06-24 15:03:18 0.571 sec 3.0      0.493435134409 0.679892
18      2018-06-24 15:03:18 0.659 sec 4.0      0.492679667171 0.678305
19  ---
20      2018-06-24 15:03:19 1.667 sec 46.0      0.485057037114 0.661619
21      2018-06-24 15:03:19 1.688 sec 47.0      0.484945381131 0.661377
22      2018-06-24 15:03:19 1.710 sec 48.0      0.484866365041 0.661216
23      2018-06-24 15:03:19 1.734 sec 49.0      0.484810585256 0.661098
24      2018-06-24 15:03:19 1.760 sec 50.0      0.484733233643 0.660935
25 See the whole table with table.as_data_frame()
26

```

We've Updated Our Site Policies.

We have recently updated our terms of service and privacy policy.

For additional information, visit.

<https://dzone.com/pages/tos> | <https://dzone.com/pages/privacy>

CLOSE

27	Variable Importances:			
28	variable	relative_importance	scaled_importance	percentage
29	-----	-----	-----	-----
30	Distance	1379.99	1	0.501313
31	Month	970.739	0.703437	0.352643
32	DayOfWeek	402.024	0.291323	0.146044

Looks like distance is much more important than month or day of week when it comes to the plane being delayed. Of course, according to this very basic model given the default parameters. A pro-tip at the end: H₂O supports XGBoost. It is trivial to swap GBM for XGBoost in this phase and see how the model changes with default hyperparameters:

```
1 from h2o.estimators.xgboost import H2OXGBoostEstimator
2 xgb_model.train(x = ["Month", "DayOfWeek", "Distance"], y = "IsArrDelayed", train_data=train_data)
```

Prediction

Once a model is created, predictions are simply done by calling `predict(data)` method on a model, where the `data` argument is the variable pointing to an `H2OFrame` with data to do the prediction on. To test the prediction is functional in a very simple way, let's use the `gbm_model` and let it predict the original training dataset by issuing `gbm_model.predict(airlines_train_data)` command.

```

1 >>> gbm_model.predict(airlines_train_data)
2 gbm prediction progress: |██████████████████████████████████████████████████████████████████████████████|
3 predict          NO          YES
4 -----
5 YES             0.141994    0.858006
6 YES             0.101574    0.898426
7 YES             0.203606    0.796394
8 YES             0.12399     0.87601
9 YES             0.138436    0.861564
10 YES            0.141994    0.858006
11 YES            0.101574    0.898426
12 YES            0.103421    0.896579
13 YES            0.203606    0.796394
14 YES            0.12399     0.87601
15
16 [43978 rows x 3 columns]

```

A result of the prediction is a `H2OFrame`. Pointer to it can be saved into a variable as well, e.g. `prediction = gbm_model.predict(airlines_train_data)`. The table printed is only a preview of the first few predictions made. As the above example demonstrates, the prediction is not very accurate in case there was no delay. This is expected, as the model is very basic. Of course, the confusion matrix seen earlier in this tutorial gave out the information about such “bad” performance beforehand.

We've Updated Our Site Policies.

Getting Started With Flow

Getting Started With Flow

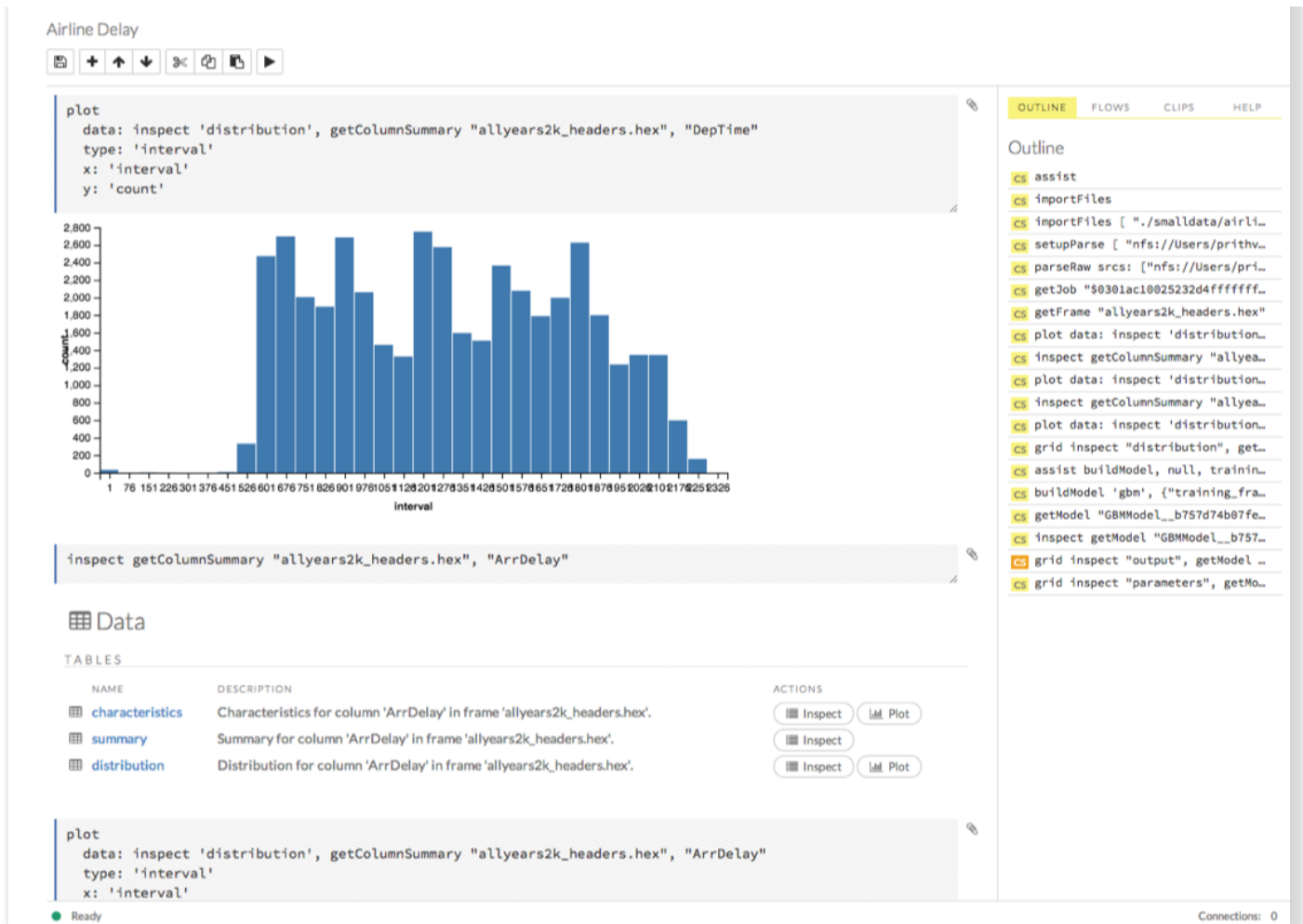
We have recently updated our terms of service and privacy policy.

Flow is H₂O's web interface. It is very powerful and oriented on visuals. It has it all. For additional information, visit:

new models, visualizing, and refining existing achievements to scoring.

<https://dzone.com/pages/tos> | <https://dzone.com/pages/privacy>

CLOSE



Flow is active whenever H₂O is started. If you've followed previous Python or R tutorials, during the active Python or R session, Flow can be reached from your browser. Whenever `h2o.init()` is called, Flow is started with H₂O. In the output after `h2o.init()`, the URL and are printed. When ran locally, the Flow is usually bound to `localhost`, with the default port being `54321`. Address or port may be changed by `h2o.init()` arguments, e.g. `h2o.init(ip='127.0.0.1', port='10001')`.

```

1 >>> h2o.init()
2 Checking whether there is an H2O instance running at http://localhost:54321..... n
3 Attempting to start a local H2O server...
4 //Some output omitted
5 -----
6 H2O cluster uptime:      01 secs
7 H2O cluster timezone:   Europe/Prague
8 H2O data parsing timezone: UTC
9 H2O cluster version:    3.19.0.99999
10 H2O cluster version age: 3 months and 5 days
11 H2O cluster name:       H2O_from_python_pavel_oy9g50
12 H2O cluster total nodes: 1
13 H2O cluster free memory: 5.207 Gb
14 H2O cluster total cores: 8
15 H2O connection url:     http://127.0.0.1:54321
16 H2O connection proxy:
17
18

```

We've Updated Our Site Policies.
We have recently updated our terms of service and privacy policy.
For additional information, visit: <https://dzone.com/pages/tos> accepting new members, healthy
<https://dzone.com/pages/privacy>
H2O connection url: <http://127.0.0.1:54321> <<<<--- URL to FLOW

CLOSE

```

19 H2O internal security:      False
20 H2O API Extensions:        XGBoost, Algos, AutoML, Core V3, Core V4
21 Python version:            2.7.15 candidate
22 -----

```

From the shortened output, the actual URL can be observed: `H2O connection url:`
`http://127.0.0.1:54321` . As `127.0.0.1` resolves to localhost, copying the given URL or simply typing
`localhost:54321` into a web browser opens H2O Flow.

Starting Flow Without Python/R

There is no need to install the Python module or R library in order to run H₂O. Just download the latest stable H2O release and run it. Java is required to run H₂O.

1. Download package with latest stable release
2. Unpack it
3. Run `java -jar h2o.jar`

The `h2o.jar` file is to be found in the extracted directory. If you're on Windows, double-clicking the `h2o.jar` file should be sufficient to run it if there is Java installed. At the time this tutorial is written, Java version 7 is the minimum required version. Java 8 is recommended.

Importing Data With Flow

On top of the page, click on `Data` . As the menu appears, choose `Import files` . Do not interchange with “Upload file” option, which is only useful to upload files from your very own machine. After clicking on the `Import files` option, a form appears at the bottom of the page. Please fill the `search` input box with `https://s3.amazonaws.com/h2o-airlines-unpacked/allyears2k.csv` to download the Airlines dataset described at the beginning of this chapter. By clicking on the `magnifying glass` next to the input box, H2O is going to verify that it can reach the file. This way, even local files or files from remote filesystems (including Hadoop filesystem) can be imported. By inserting whole folders, H2O will pop-up all the files available, thus enabling multi-file import with ease. In this simple case, after clicking on the `magnifying glass` , a single file named exactly as the URL inserted into the input box will appear a little bit lower. By clicking on the `plus` icon , the file is marked for import. By clicking `Import` button , H2O will automatically download the file.

After the file is imported, another form will appear. This time, H₂O confirms that the file has been imported and asks to user to starting parsing it. By clicking on `Parse these files` , H2O will begin parsing.

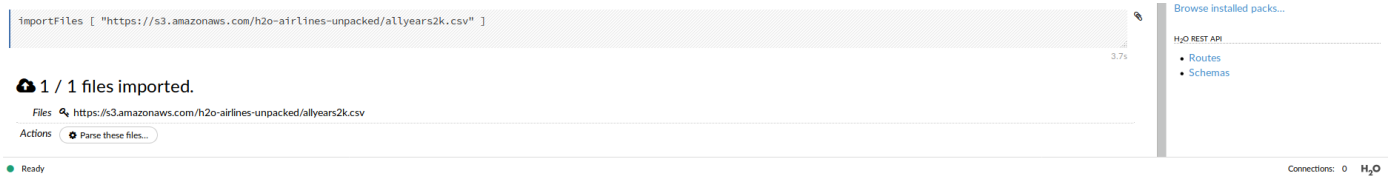
We've Updated Our Site Policies.

We have recently updated our terms of service and privacy policy.

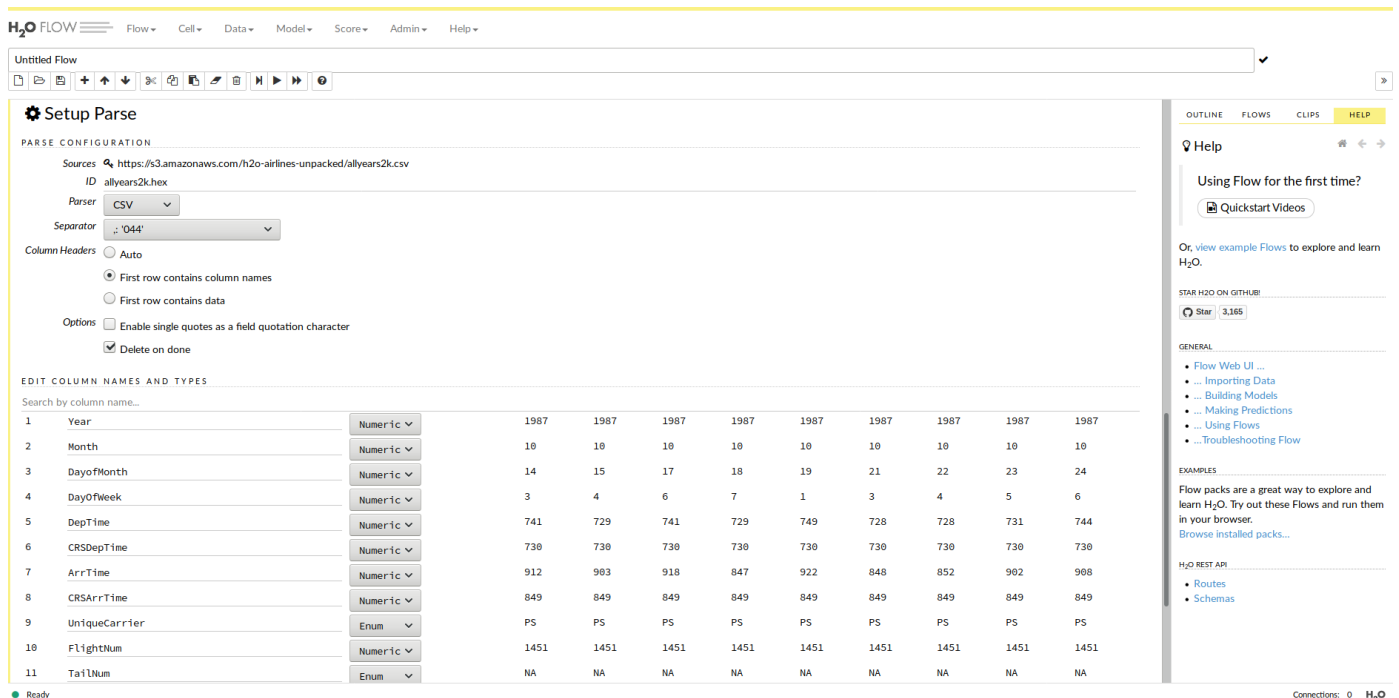
For additional information, visit:

<https://dzone.com/pages/tos> | <https://dzone.com/pages/privacy>

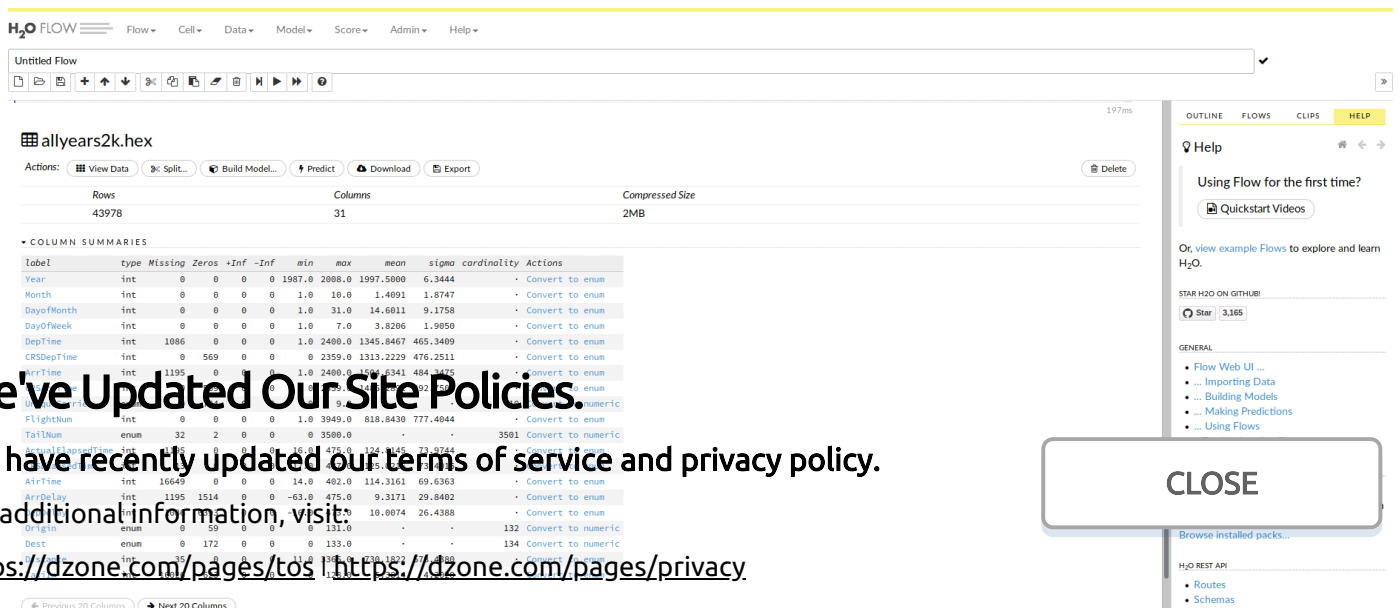
CLOSE



However, H2O is not going to parse the whole file right away. The could be potentially very big (or there may be multiple files). Instead, it will automatically detect the format and suggests best parsing settings based on internal heuristics. In case of CSV imported in this tutorial, it will correctly detect the format, set separator to be a comma, detect column headers on first row. However, there is more — columns datatypes are detected. Yet the data scientist has the power to override every decision, e.g. change column type. By clicking on **Parse** button at the very bottom, the file is finally parsed by H2O.



The sample dataset is very small (4,4 MB) and H2O's CSV parser is very fast, so the parsing should be quick. After the parsing phase, H2O will inform you about the dataset being ready for an inspection. By clicking on the **view** button, a table view preview of the parsed data appears. Here, the user has several choices, including viewing all the data, exporting it, splitting it into several smaller chunks (e.g. for testing), or most importantly to **Build Model**.



By clicking on the `Build model` button above the data overview, a dialogue asking the user to select an algorithm appears. In this tutorial, GBM is used as a simple example, therefore, select `Gradient Boosting Machine`. As with Python and R, let's train a model that is able to predict if the plane arrives late based on month, day of week, and distance the plane has to travel before reaching its destination. From top to bottom, only few of the choices need your attention right now. Most of these are hyperparameters, which are of course important in real world, yet for a getting started guide, the default values will suffice.

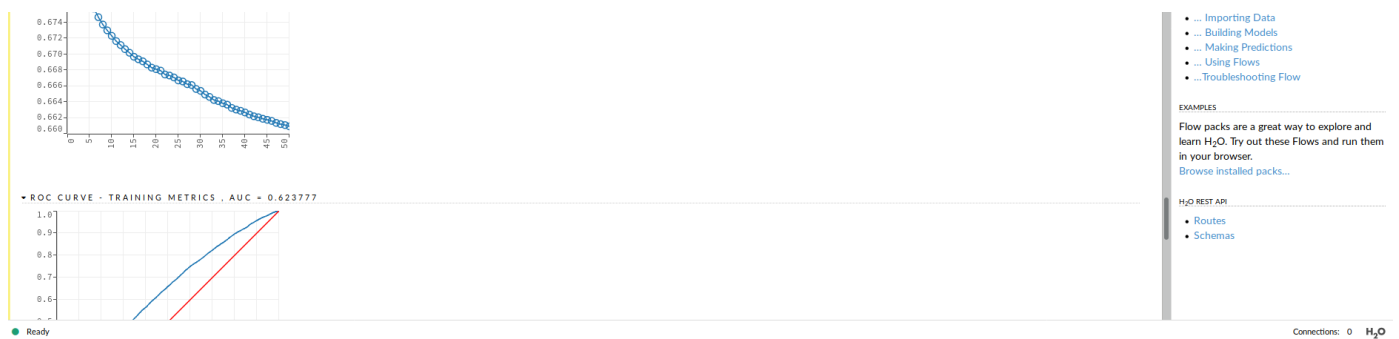
The `training_frame` option is already pre-set by flow to point to the data imported. Creating validation frames from training is easy using the following option, which we'll leave intact for now. Since this model is trying to predict if the plane will be delayed on arrival at its destination, the response column should be set to `IsArrDelayed`. The prediction is based on month, day of week, and traveling distance (chosen arbitrarily). Therefore, in the ignored column section, first check all the columns as ignored and leave only `Month`, `DayOfWeek`, `Distance`, and `IsArrDelayed` (response variable) unchecked.

After the parameters are set, click on `Build model` button at the bottom to start the model training.

The screenshot shows the H2O FLOW interface with the 'Build a Model' dialog open. The dialog is for a Gradient Boosting Machine model. The 'training_frame' is set to 'allyears2k.hex'. The 'response_column' is set to 'IsArrDelayed'. The 'ignored_columns' section shows a list of columns with checkboxes to select or deselect them. The 'Month' and 'DayOfWeek' columns are selected. The 'Distance' column is also selected. The 'IsArrDelayed' column is the response variable. The 'Build a Model' button is at the bottom right of the dialog.

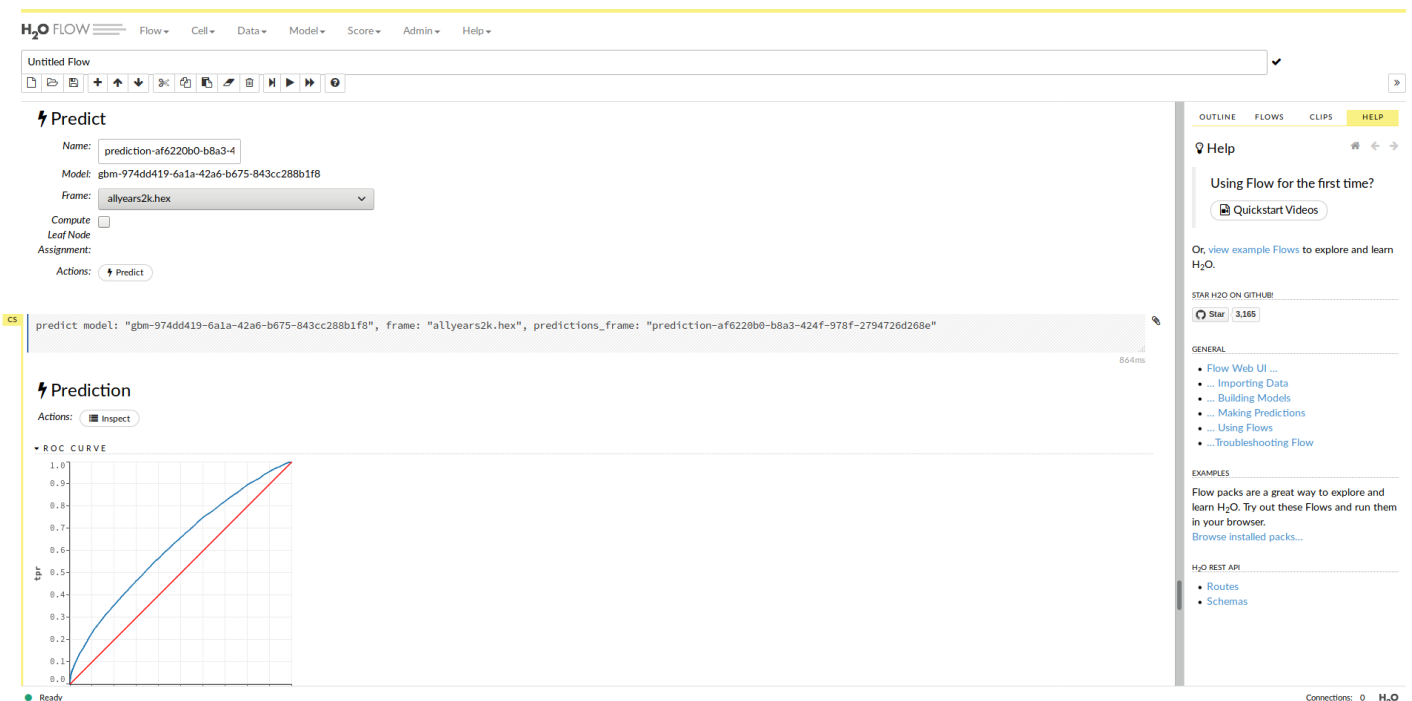
Model training in this small case should be done in an instant. H₂O will again display a message about training progress and after the training is finished, the model can be displayed by clicking on the `view` button. A complete model description appears. Variable importance, logloss, complete history per each round of training (including durations), and much more can be observed.

The screenshot shows the H2O FLOW interface with the 'Model Details' dialog open. The dialog displays the model ID, algorithm, and training history. A large overlay message is present: "We've Updated Our Site Policies. We have recently updated our terms of service and privacy policy. For additional information, visit: <https://dzone.com/pages/tos> | <https://dzone.com/pages/privacy>." A "CLOSE" button is visible in the bottom right corner of the overlay.



Predicting With Flow

After the model is built, click the `Predict` button (with a lightning icon) on top of model summary reached in previous steps. To test the prediction is functional in a very simple way, let's use the current model and let it predict the original training dataset. After clicking on the `Predict` button, select the training frame in the dropdown menu's `Frame` option and click on `Predict`. The resulting predictions appear, including various metrics. The model is not very accurate, especially the error rate regarding prediction of a plane not arriving late are very inaccurate. Creating a better model by tuning the hyperparameters or playing with various predictors is up to each and every data scientist to try now.



Where to Go Next?

In this hands-on guide, very basic ways to interact with H₂O were shown. However, there is so much more to H₂O. You can take H₂O from your laptop to large clusters, creating models and predicting on extremely large datasets. Cross-validation, stacked ensemble, tuning parameters with grid search and many, many other expected functionality. After data scientists tune the model, deploying it into production is easy with H₂O's POJO/MOJO functionality. This way, the model is simply exported into a self-containing package, making it easy for the engineers to plug it into any Java-based production environment. Model's performance can be

We've Updated Our Site Policies

We have recently updated our terms of service and privacy policy.

Data scientists and developers are best to visit H₂O documentation. It provides step-by-step guides for various tasks. For additional information for visitors interested specifically in R language or Python language, please visit our website.

For a comprehensive overview of what H₂O is capable of, besides a list of available algorithms, methods of cross-validation, or the ability to build models on top of existing ones, ways of productionizing created models quickly are described.

CLOSE

In the documentation, everything from data import, exploration, and filtering to deployment into production is described. Running on Hadoop? No problem. Using Apache Spark? Sparkling water is at your service. H2O also integrates well with various cloud services. H2O also offers countless tutorials on GitHub.

Video tutorials, speeches, and expert advices are available on H2O's YouTube channel. There is so much to explore.

H₂O also holds conferences named H2O World. Come and visit us, we'll be glad to talk to you.

Like This Article? Read More From DZone



Inspecting Decision Trees in H2O



What's New in H2O Machine Learning: Christmas 2018



How We Built Tagger News: Machine Learning on a Tight Schedule



Free DZone Refcard Introduction to TensorFlow

Topics: ARTIFICIAL INTELLIGENCE , MACHINE LEARNING , H2O.AI , PYTHON , R , XGBOOST , GBM , STATISTICS

Published at DZone with permission of Pavel Pscheidl . [See the original article here.](#) 
Opinions expressed by DZone contributors are their own.

IN PROGRESS

We've Updated Our Site Policies.

We have recently updated our terms of service and privacy policy.

For additional information, visit:

<https://dzone.com/pages/tos> | <https://dzone.com/pages/privacy>

CLOSE