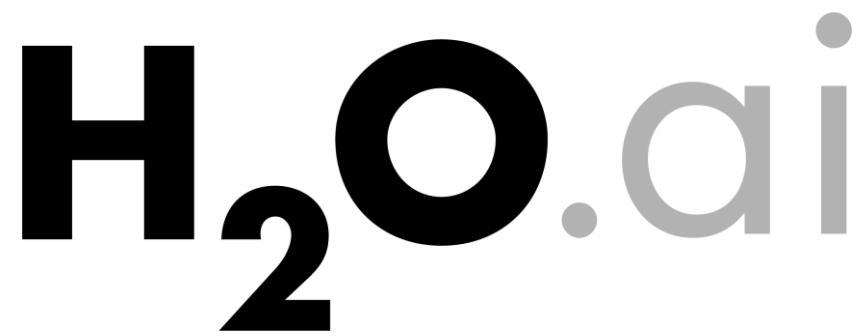


Latest Developments in H₂O



Jo-fai (Joe) Chow

Data Scientist

joe@h2o.ai

@matlabulous

H₂O at Booking.com – Amsterdam
6th April, 2017

About Me

- Civil (Water) Engineer
 - 2010 – 2015
 - Consultant (UK)
 - Utilities
 - Asset Management
 - Constrained Optimization
 - Industrial PhD (UK)
 - Infrastructure Design Optimization
 - Machine Learning + Water Engineering
 - Discovered H₂O in 2014
- Data Scientist
 - 2015
 - Virgin Media (UK)
 - Domino Data Lab (Silicon Valley, US)
 - 2016 – Present
 - H₂O.ai (Silicon Valley, US)

About Me

Domino Data Lab
At the intersection of data science and engineering.
Domino App Site | [Twitter](#) | [Email](#)

19 Sep 2014 • [Facebook Like](#) 0 | [Twitter Tweet](#) 21 | [Google+ 1](#) 4

How to use R, H2O, and Domino for a Kaggle competition

Guest post by Jo-Fai Chow

The sample project (code and data) described below is [available on Domino](#).

If you're in a hurry, feel free to skip to:

- Tutorial 1: [Using Domino](#)
- Tutorial 2: [Using H2O to Predict Soil Properties](#)
- Tutorial 3: [Scaling up your analysis](#)

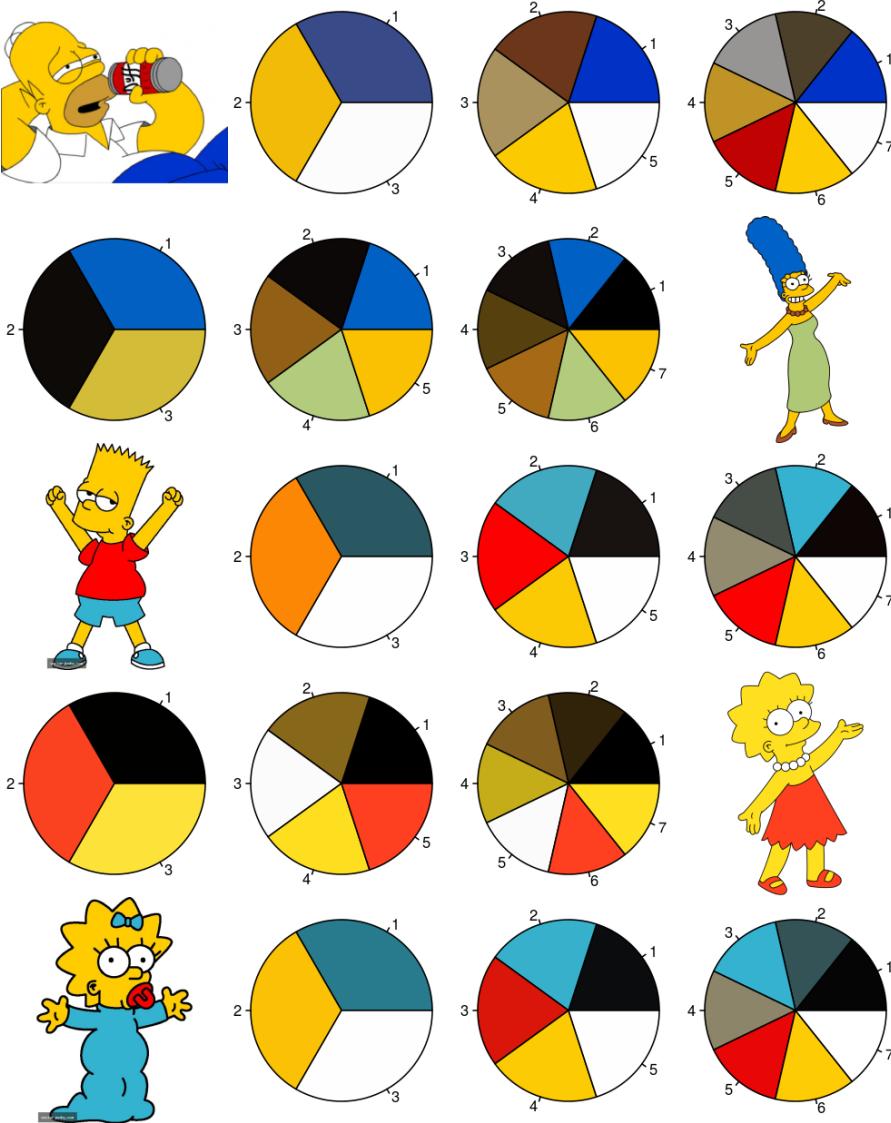
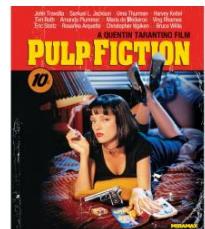
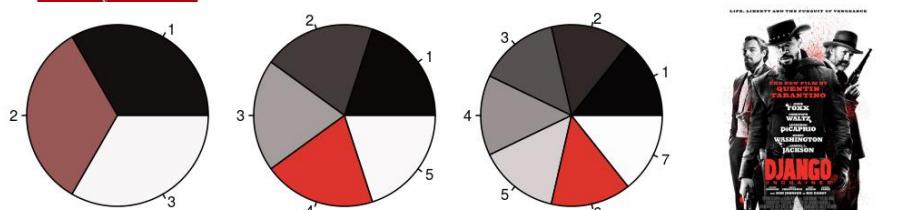
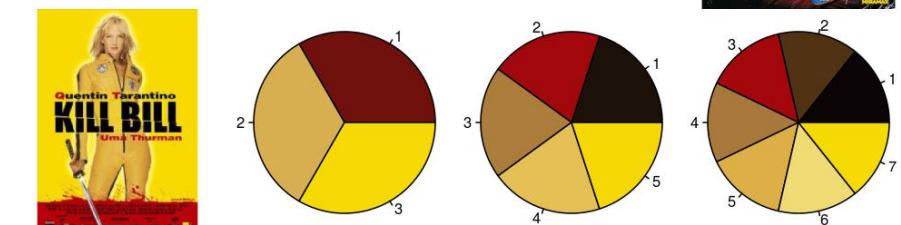
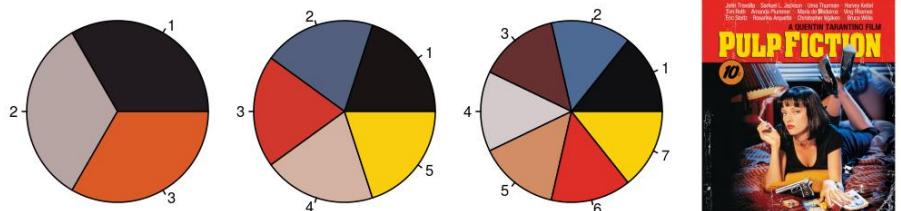
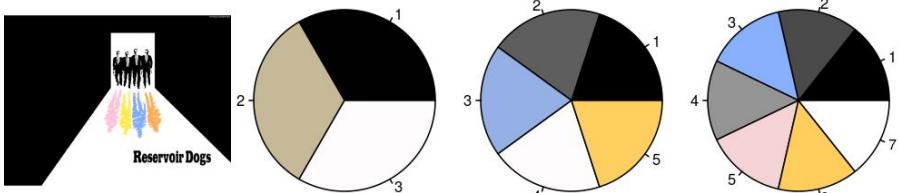
Introduction

This blog post is the sequel to [TTTAR1](#) a.k.a. [An Introduction to H2O Deep Learning](#). If the previous blog post was a brief intro, this post is a proper machine learning case study based on a recent [Kaggle competition](#): I am leveraging [R](#), [H2O](#) and [Domino](#) to compete (and do pretty well) in a real-world data mining contest.

R + H₂O + Domino for Kaggle
[Guest Blog Post for Domino & H₂O \(2014\)](#)

- The Long Story
 - bit.ly/joe_kaggle_story

About Me



Developing R Packages for Fun
[rPlotter](#) (2014)

Agenda

- About H₂O.ai
 - Company
 - Machine Learning Platform
- What's New?
 - Deep Water
 - H₂O + xgboost
 - Stacked Ensembles
 - Auto Machine Learning
 - Model Interpretation
 - Community
- PyData Conference

Booking.com



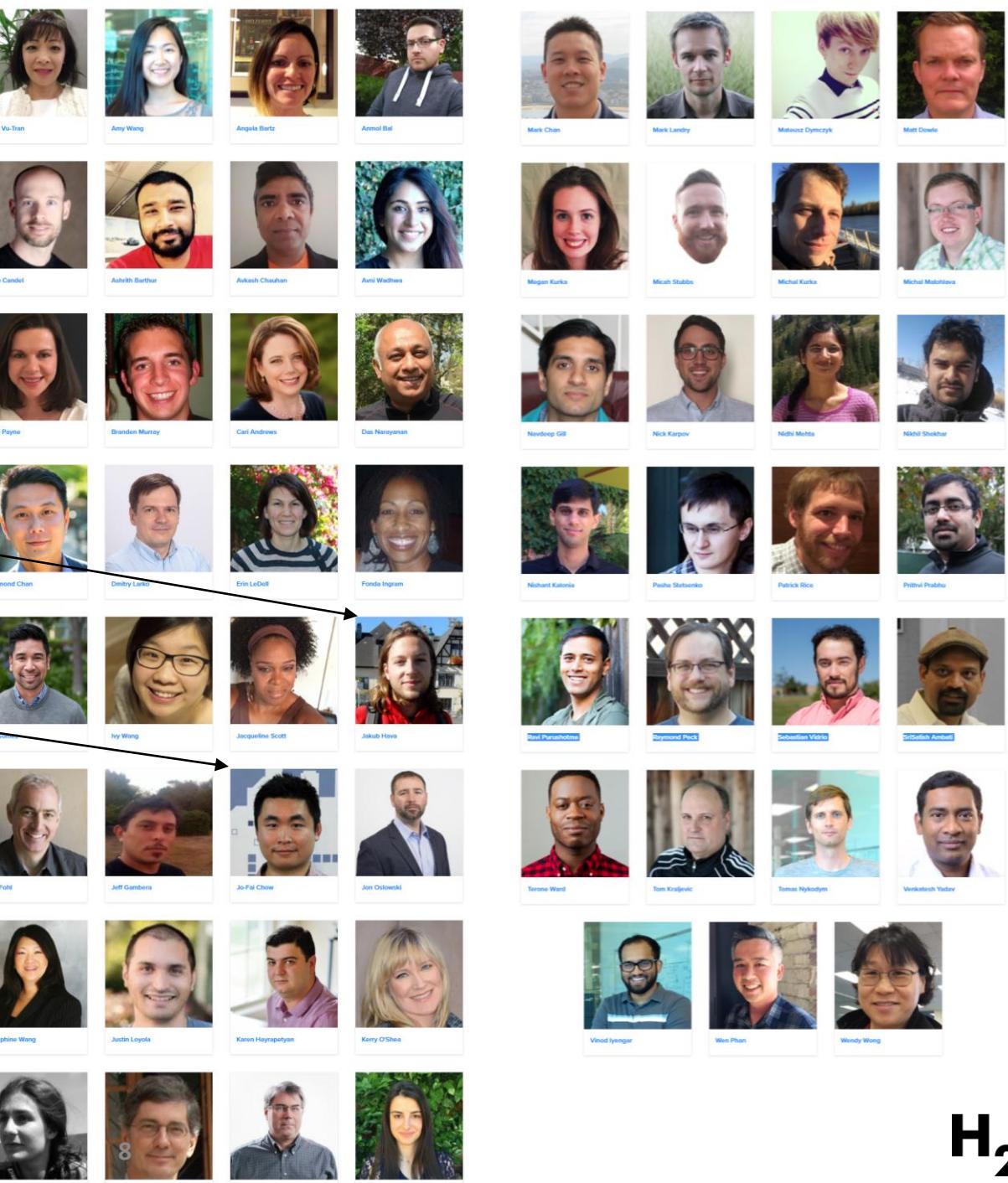
About H₂O.ai

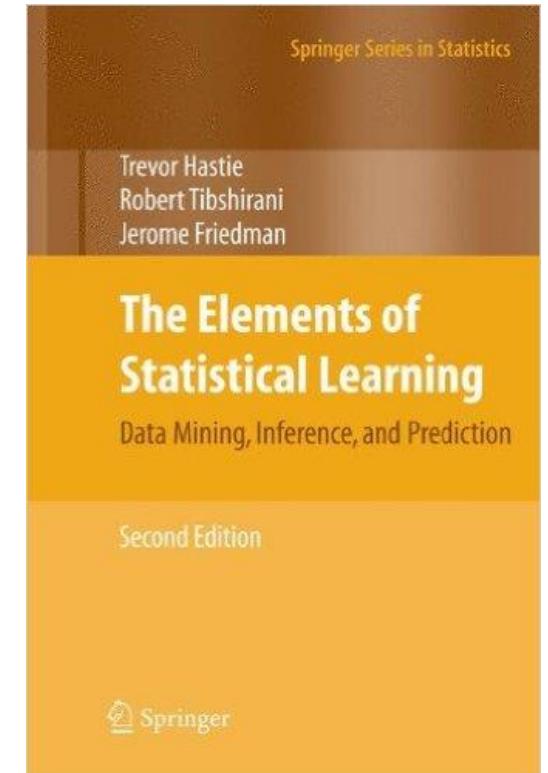
Company Overview

Founded	2011 Venture-backed, debuted in 2012
Products	<ul style="list-style-type: none">• H₂O Open Source In-Memory AI Prediction Engine• Sparkling Water• Steam
Mission	Operationalize Data Science, and provide a platform for users to build beautiful data products
Team	<p>70 employees</p> <ul style="list-style-type: none">• Distributed Systems Engineers doing Machine Learning• World-class visualization designers
Headquarters	Mountain View, CA



Our Team





Scientific Advisory Council



Dr. Trevor Hastie

- John A. Overdeck Professor of Mathematics, Stanford University
- PhD in Statistics, Stanford University
- Co-author, *The Elements of Statistical Learning: Prediction, Inference and Data Mining*
- Co-author with John Chambers, *Statistical Models in S*
- Co-author, *Generalized Additive Models*



Dr. Robert Tibshirani

- Professor of Statistics and Health Research and Policy, Stanford University
- PhD in Statistics, Stanford University
- Co-author, *The Elements of Statistical Learning: Prediction, Inference and Data Mining*
- Author, *Regression Shrinkage and Selection via the Lasso*
- Co-author, *An Introduction to the Bootstrap*



Dr. Steven Boyd

- Professor of Electrical Engineering and Computer Science, Stanford University
- PhD in Electrical Engineering and Computer Science, UC Berkeley
- Co-author, *Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers*
- Co-author, *Linear Matrix Inequalities in System and Control Theory*
- Co-author, *Convex Optimization*



wenphan
@wenphan

Following



So much brain power in one place:
[@ArnoCandel](#) and Stanford profs. Boyd,
Tibs, and Hastie. Hacking algos at [@h2oai](#)
HQ



Figure 1. Magic Quadrant for Data Science Platforms



H2O.ai recognized for completeness of vision and ability to execute

We are thrilled to be named a Visionary among the 16 vendors included in Gartner's 2017 Magic Quadrant for Data Science Platforms. As a Visionary we believe we are positioned highest in Ability to Execute for companies of our size and scale.

Since 2011, our mission has been to democratize data science through open source AI and [deep learning](#). Today, H2O.ai is focused on bringing AI to enterprises with a growing community of more than 8,500 organizations that depend on H2O for mission critical applications. H2O.ai was recently named [CB Insights AI 100](#) and is used by [107 of the Fortune 500 companies](#).

Disclaimer: This graphic was published by Gartner, Inc. as part of a larger research document and should be evaluated in the context of the entire document. The Gartner document is available upon request from H2O.ai.

Check out our website h2o.ai

World Record Performance for AI

H2O.ai is accelerating both machine learning and deep learning on GPUs, providing enterprises opportunities to build better models and enable new use cases.

[SEE DEEP WATER](#)

H2O in Action



▶ What data products mean and why H2O keeps this industry leader relevant.



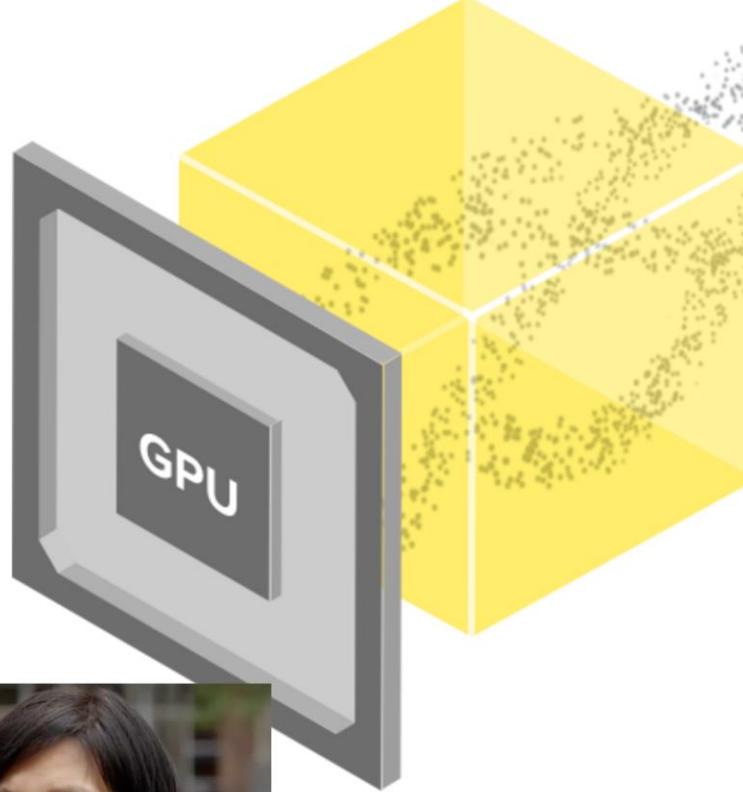
▶ Various data leaders discuss the transformative impact of H2O AI for ADP.



▶ Capital One team members find out how they can leverage H2O's leading technology.

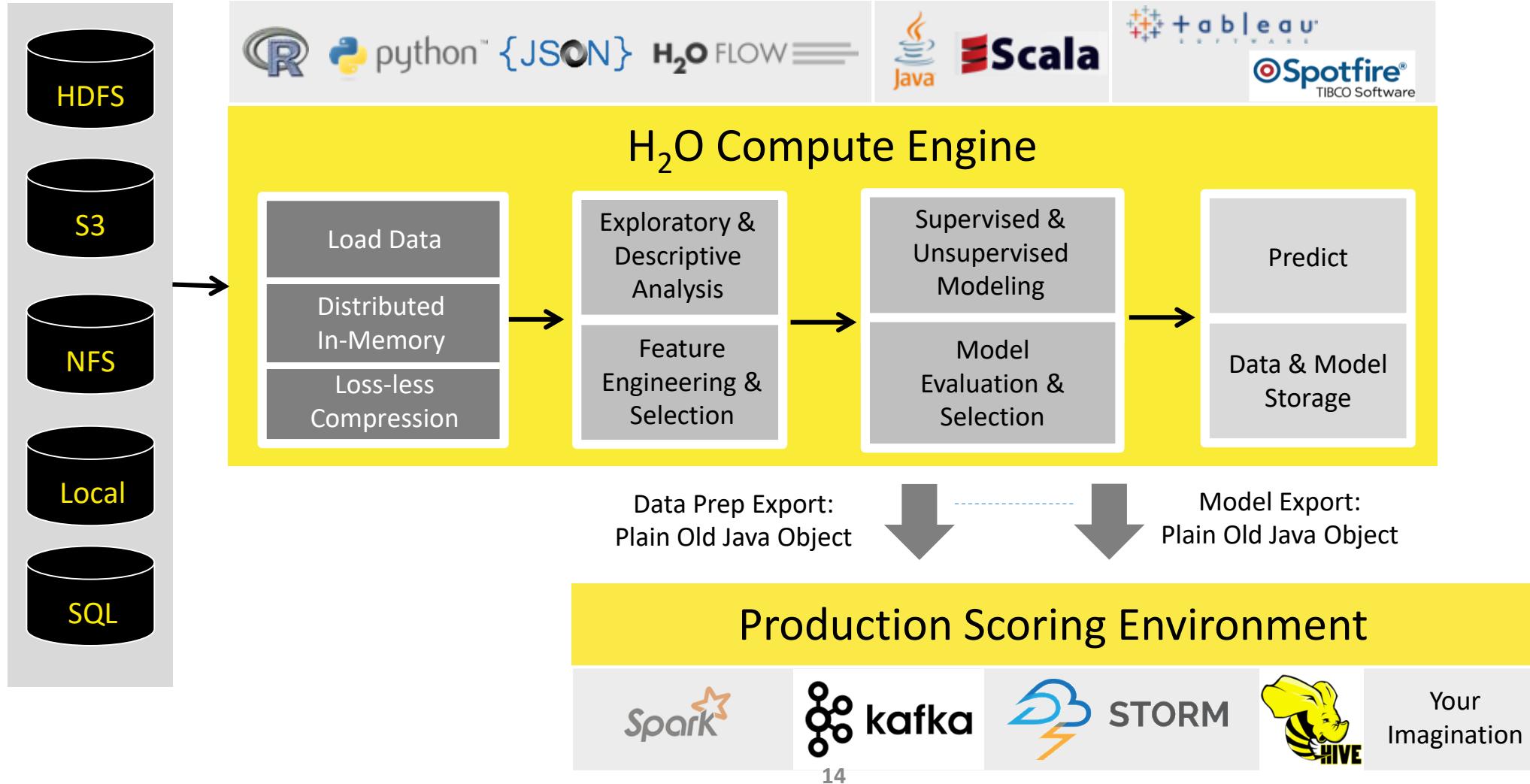


▶ Kaiser uses H2O machine learning to save lives.



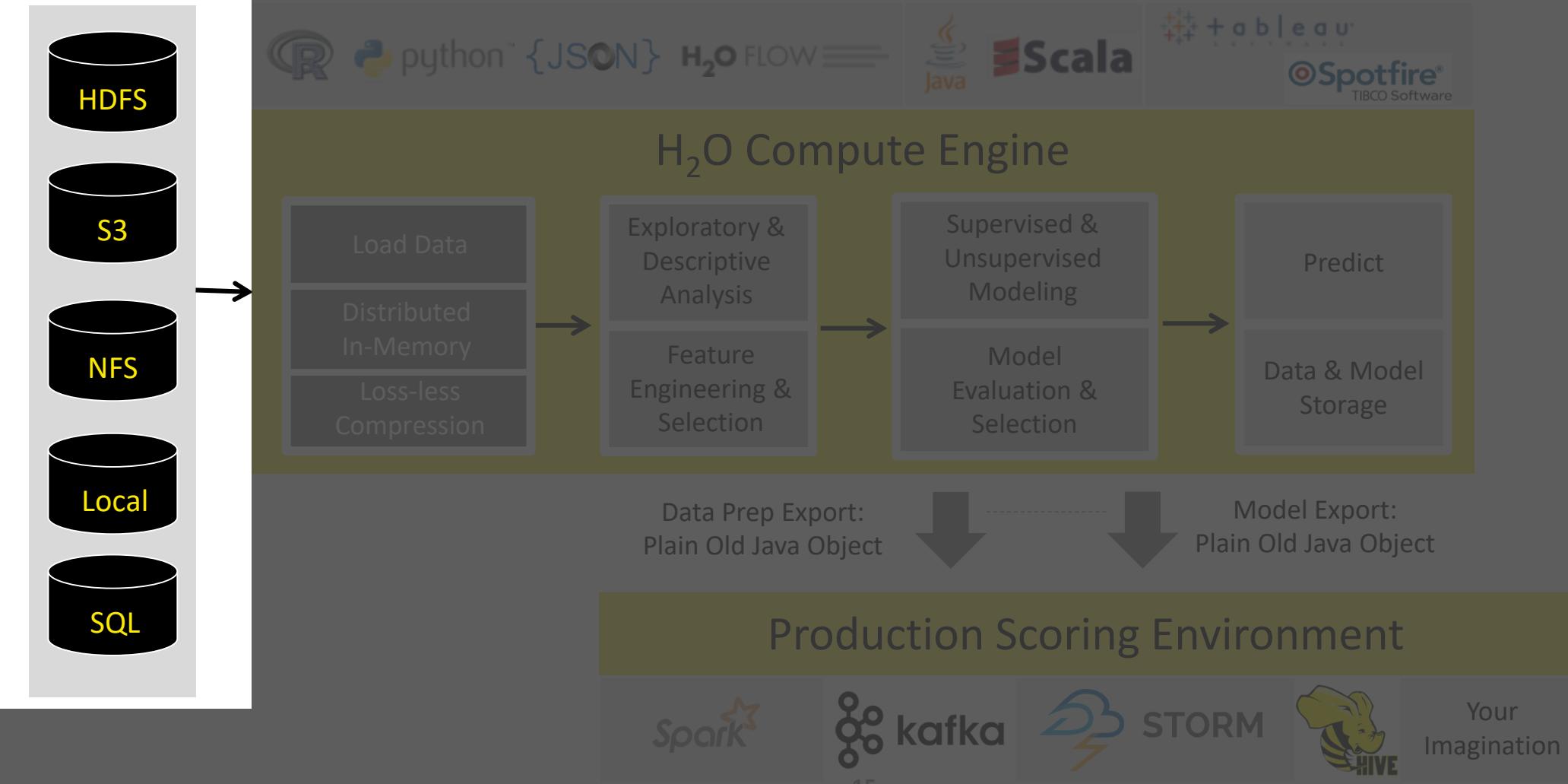
H₂O Machine Learning Platform

High Level Architecture



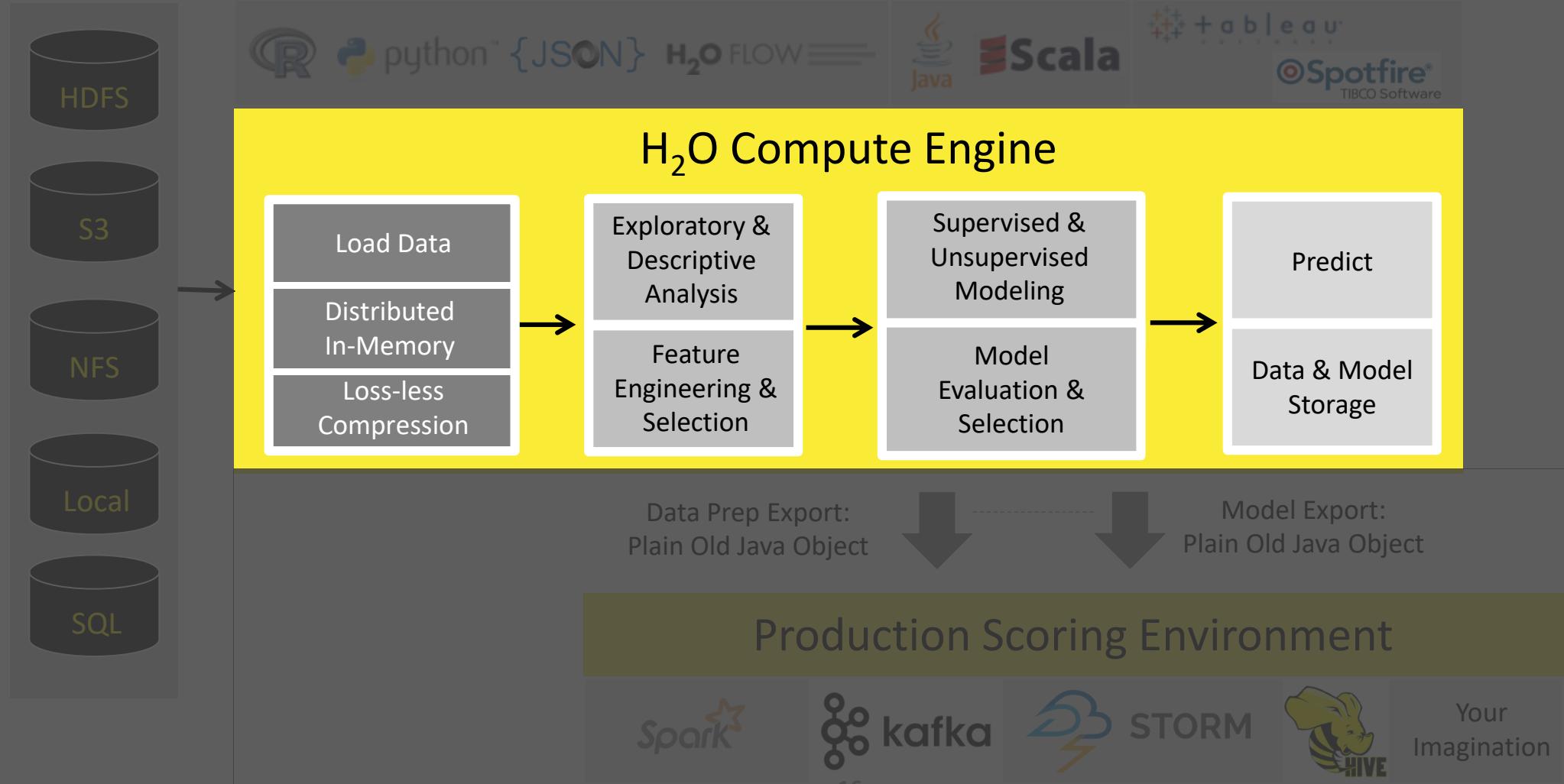
High Level Architecture

Import Data from
Multiple Sources



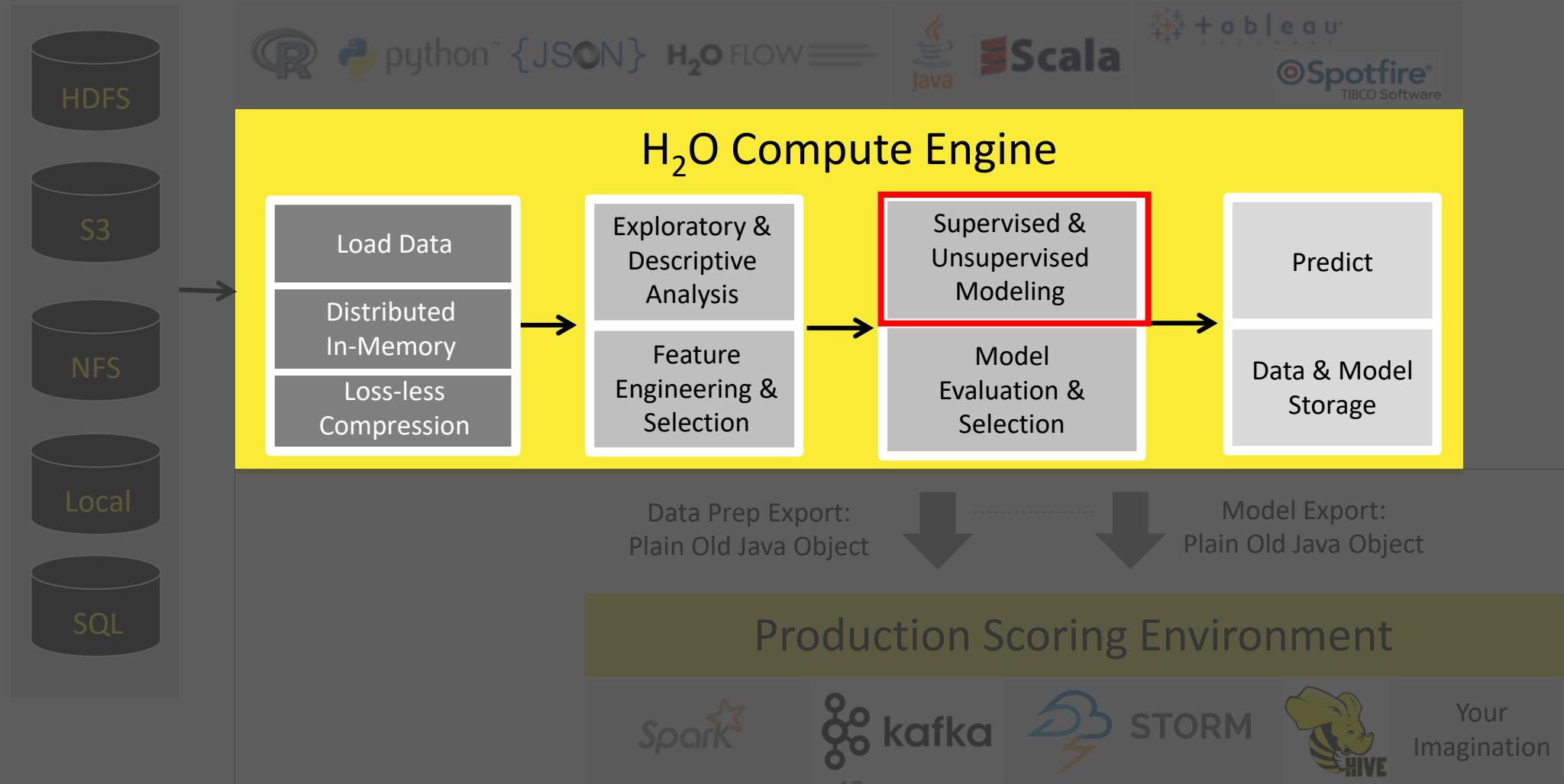
High Level Architecture

Fast, Scalable & Distributed
Compute Engine Written in
Java



High Level Architecture

Fast, Scalable & Distributed
Compute Engine Written in
Java



Algorithms Overview

Supervised Learning

Statistical Analysis

- **Generalized Linear Models:** Binomial, Gaussian, Gamma, Poisson and Tweedie
- **Naïve Bayes**

Ensembles

- **Distributed Random Forest:** Classification or regression models
- **Gradient Boosting Machine:** Produces an ensemble of decision trees with increasing refined approximations

Deep Neural Networks

- **Deep learning:** Create multi-layer feed forward neural networks starting with an input layer followed by multiple layers of nonlinear transformations

Unsupervised Learning

Clustering

- **K-means:** Partitions observations into k clusters/groups of the same spatial size. Automatically detect optimal k

Dimensionality Reduction

- **Principal Component Analysis:** Linearly transforms correlated variables to independent components
- **Generalized Low Rank Models:** extend the idea of PCA to handle arbitrary data consisting of numerical, Boolean, categorical, and missing data

Anomaly Detection

- **Autoencoders:** Find outliers using a nonlinear dimensionality reduction using deep learning

H₂O Deep Learning in Action

116M rows, 6GB CSV file
800+ predictors (numeric + categorical)

airlines_all_selected_cols.hex

Actions: View Data, Split..., Build Model..., Predict, Download, Export

Rows	Columns	Compressed Size
116695259	12	2GB



Job

Run Time 00:00:36.712

Remaining Time 00:00:17.188

Type Model

Key Q deeplearning-dd2f42f7-81f7-42e8-9d98-e34437309828

Description DeepLearning

Status RUNNING

Progress 69%

Iterations: 12. Epochs: 0.628821. Speed: 2,243,735 samples/sec. Estimated time left: 21.849 sec

Actions View, Cancel Job

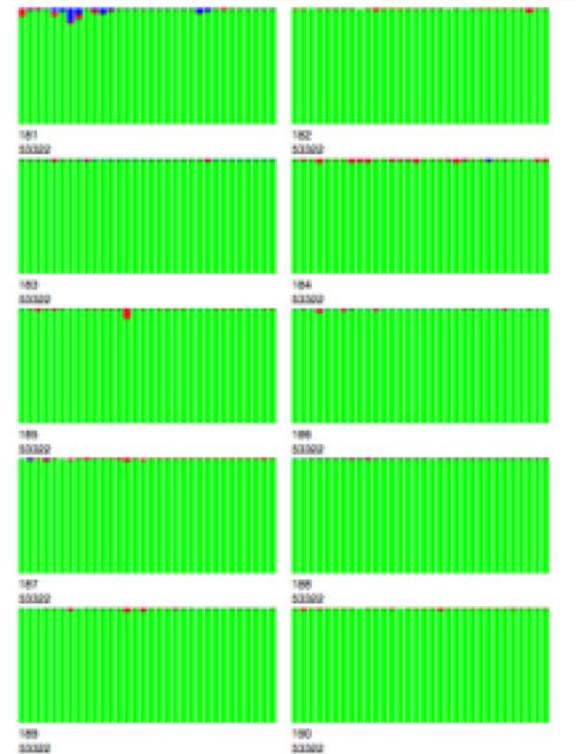
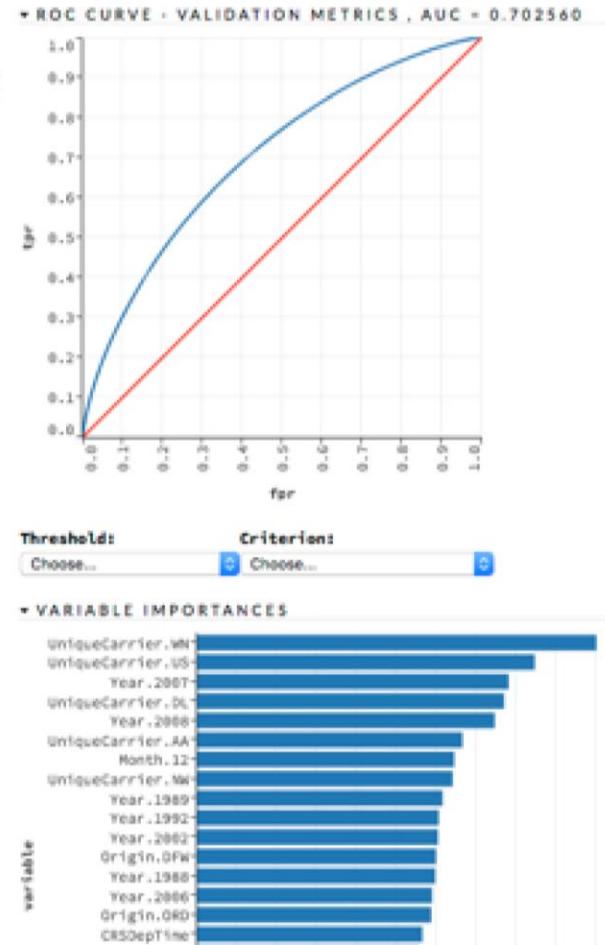
* OUTPUT - STATUS OF NEURON LAYERS (PREDICTING ISDELAYED, 2-CLASS CLASSIFICATION, BERNoulli DISTRIBUTION, CROSSENTROPY LOSS, 17,462 WEIGHTS/BIASES, 221.3 KB, 106,585,385 TRAINING SAMPLES, MINI-BATCH SIZE 1)

layer	units	type	dropout	l1	l2	mean_rate	rate_rms	momentum	mean_weight	weight_rms	mean_bias	bias_rms
1	887	Input	0									
2	20	Rectifier	0	0	0	0.0493	0.2020	0	-0.0021	0.2111	-0.9139	1.0036
3	20	Rectifier	0	0	0	0.0157	0.0227	0	-0.1833	0.5362	-1.3988	1.5259
4	20	Rectifier	0	0	0	0.0517	0.0446	0	-0.1575	0.3068	-0.8846	0.6046
5	20	Rectifier	0	0	0	0.0761	0.0844	0	-0.0374	0.2275	-0.2647	0.2481
6	2	Softmax	0	0	0	0.0161	0.0083	0	0.0741	0.7268	0.4269	0.2056

H₂O.ai

Deep Learning Model

real-time, interactive
model inspection in Flow



Legend

Each bar represents one CPU.

Blue: idle time

Green: user time

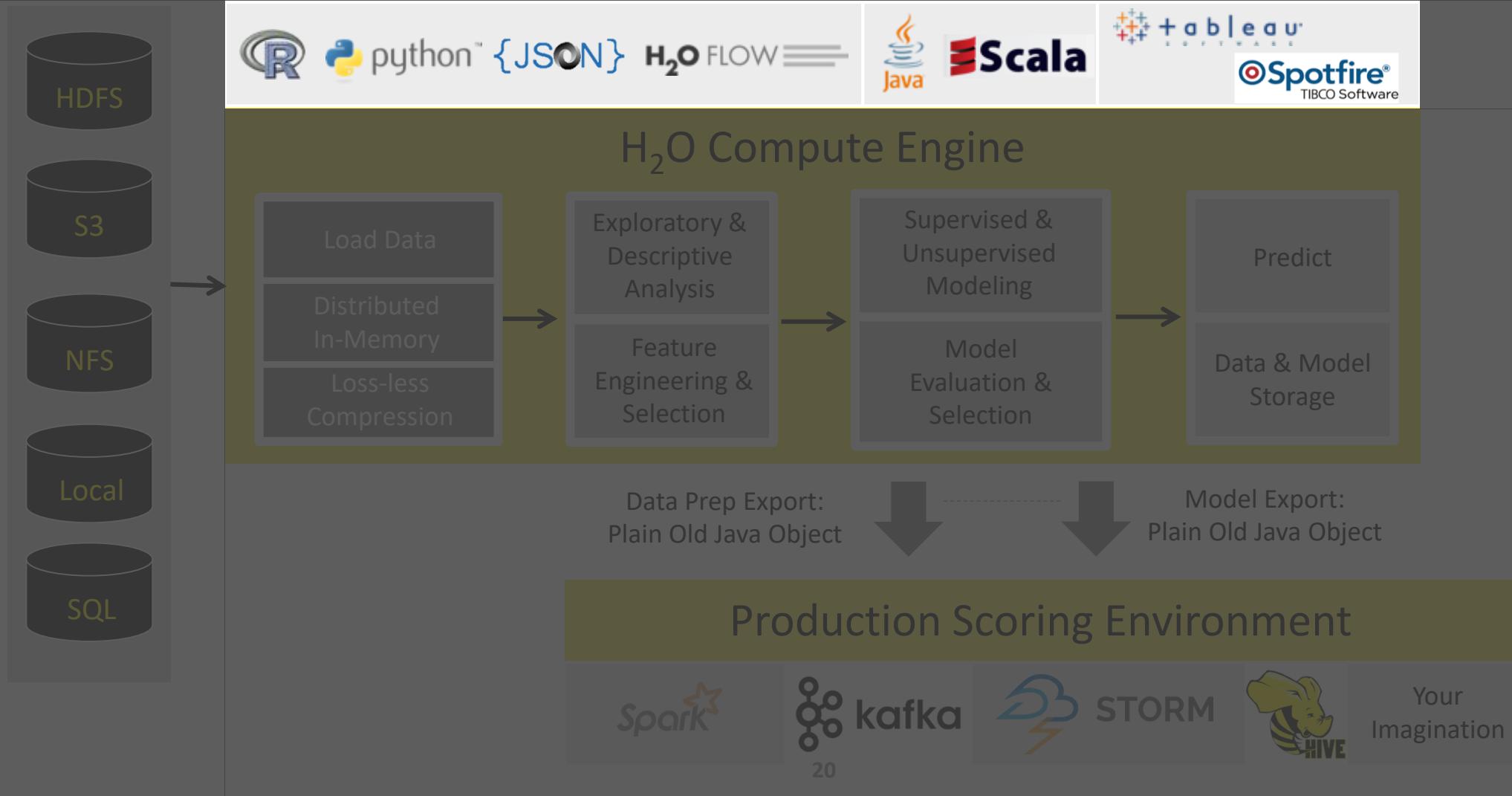
Red: system time

White: other time (e.g. Io)

10 nodes: all
320 cores busy



High Level Architecture



H₂O + R

```
# -----  
# Train a H2O Model  
# -----  
  
# Train three basic H2O models  
model_drf <- h2o.randomForest(x = features,  
.....y = target,  
.....model_id = "iris_random_forest",  
.....training_frame = d_iris)  
  
model_gbm <- h2o.gbm(x = features,  
.....y = target,  
.....model_id = "iris_gbm",  
.....training_frame = d_iris)  
  
model_dnn <- h2o.deeplearning(x = features,  
.....y = target,  
.....model_id = "iris_deep_learning",  
.....training_frame = d_iris)
```

H₂O + Python

Gradient Boosting Machines

```
# Build a Gradient Boosting Machines (GBM) model with default settings

# Import the function for GBM
from h2o.estimators.gbm import H2OGradientBoostingEstimator

# Set up GBM for regression
# Add a seed for reproducibility
gbm_default = H2OGradientBoostingEstimator(model_id = 'gbm_default', seed = 1234)

# Use .train() to build the model
gbm_default.train(x = features,
                   y = 'quality',
                   training_frame = wine_train)

gbm Model Build progress: |██████████| 100%
```



Flow ▾ Cell ▾ Data ▾

Model ▾ Score ▾ Admin ▾ Help ▾

Iris Demo



CS

Expression...

- Aggregator...
- Deep Learning...
- Distributed Random Forest...
- Gradient Boosting Machine... 🕒
- Generalized Linear Modeling...
- Generalized Low Rank Modeling...
- K-means...
- Naive Bayes...
- Principal Components Analysis...

- List All Models
- List Grid Search Results
- Import Model...
- Export Model...

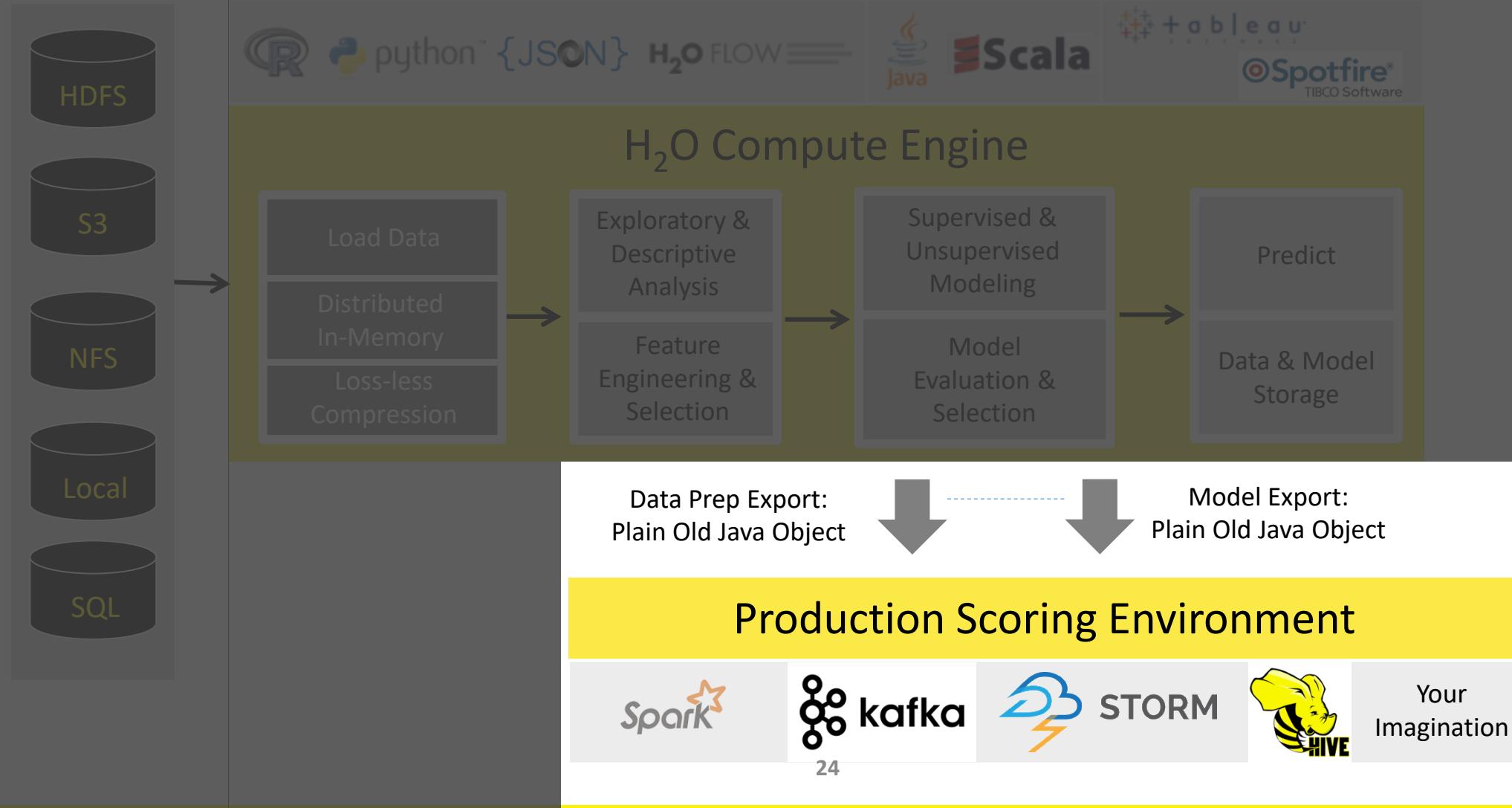
H₂O Flow (Web) Interface



Connections: 0 H₂O

High Level Architecture

Export Standalone Models
for Production



Languages

R

[Quick Start Video - R](#)
[R Package Docs](#)
[R Booklet](#)
[Examples and Demos](#)
[R FAQ](#)
[Ensemble R Package Readme](#)
[RSparkling Readme](#)
[Migrating from H2O-2](#)

Python

[Quick Start Video - Python](#)
[Python Module Docs](#)
[Python Booklet](#)
[Examples and Demos](#)
[Python FAQ](#)
[PySparkling Readme](#) [2.0](#) | [1.6](#)
[skutil Docs](#)

Java

[POJO and MOJO Model Javadoc](#)
[H2O Core Javadoc](#)
[H2O Algorithms Javadoc](#)

Scala

Sparkling Water API	2.0	1.6
Sparkling Water Scaladoc	2.0	1.6
H2O Scaladoc	2.11	2.10

Tutorials, Examples, & Presentations

Tutorials and Blogs

[H2O Tutorials HTML | PDF](#)
[H2O Blogs](#)
[H2O University](#)

Use Case Examples

Chicago crime prediction	R	Python	ScalaSW	PySW
Airlines delays prediction	R	Python	ScalaSW	PySW
Lending Club loan prediction	R	Python	ScalaSW	PySW
Ham or Spam	R	Python	ScalaSW	PySW
Prediction with prostate dataset	R	Python	ScalaSW	PySW

Presentations

[H2O Meetups](#)
[H2O World 2014 Videos](#)
[H2O World 2015 Videos](#)
[Open Tour Chicago Videos](#)
[Open Tour NYC Videos](#)
[Open Tour Dallas Videos](#)

New Developments

Deep Water

H₂O.ai Caffe  mxnet  TensorFlow

Deep Water

Next-Gen Distributed Deep Learning with H₂O

One Interface - GPU Enabled - Significant Performance Gains

Inherits All H₂O Properties in Scalability, Ease of Use and Deployment



H₂O integrates with existing **GPU** backends
for **significant performance gains**



Convolutional Neural Networks enabling
Image, video, speech recognition

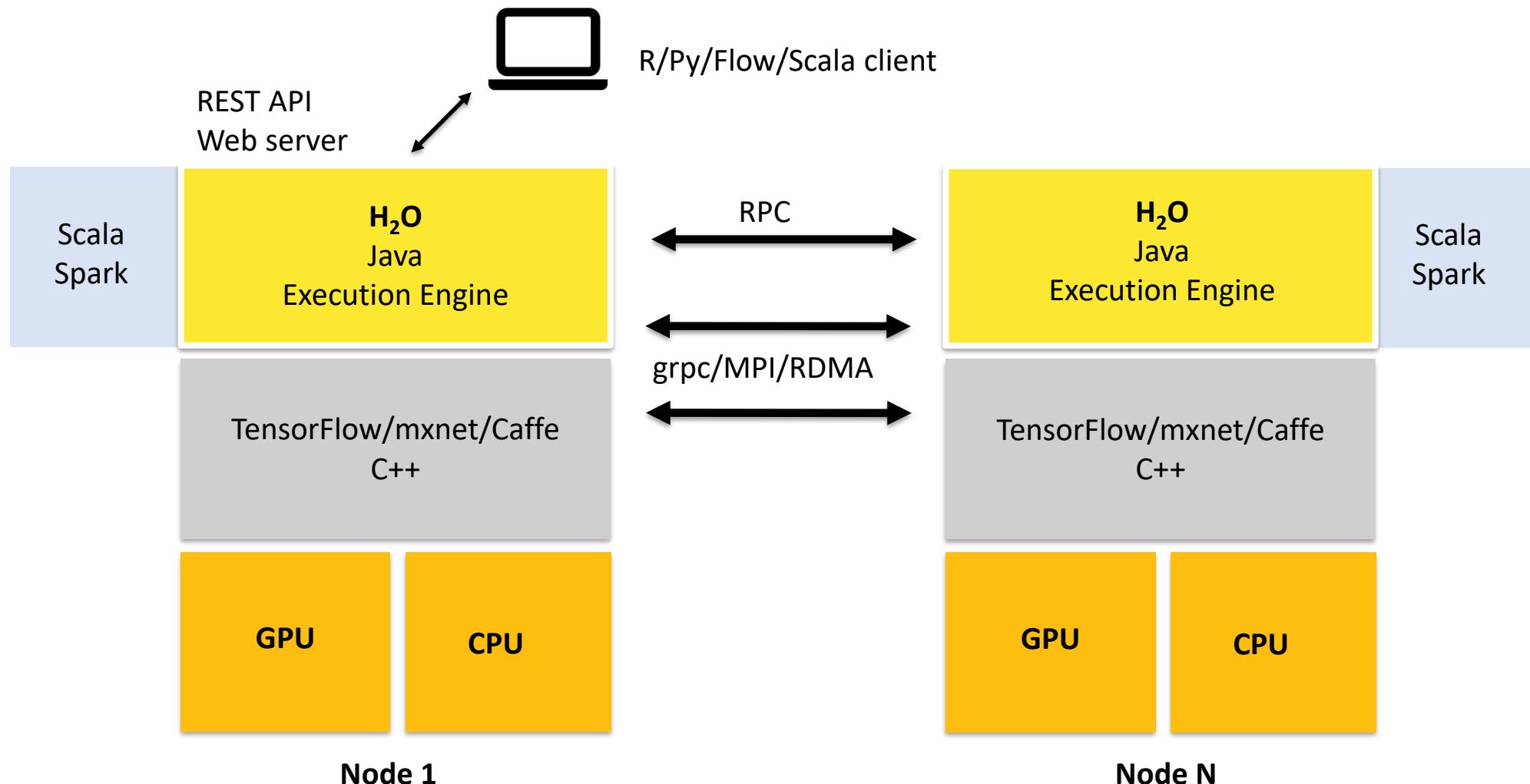


Recurrent Neural Networks
enabling **natural language processing, sequences, time series**, and more



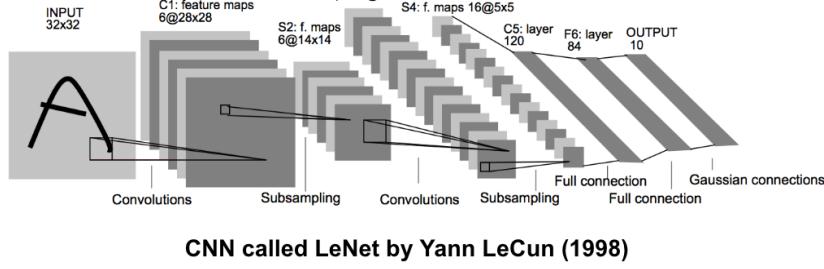
Hybrid Neural Network Architectures
enabling **speech to text translation, image captioning, scene parsing** and more

Deep Water Architecture

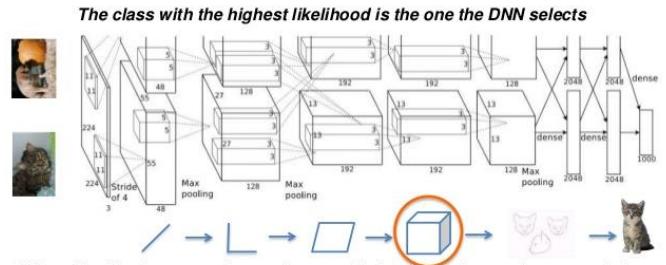


Available Networks in Deep Water

- LeNet
- AlexNet
- VGGNet
- Inception (GoogLeNet)
- ResNet (Deep Residual Learning)
- Build Your Own

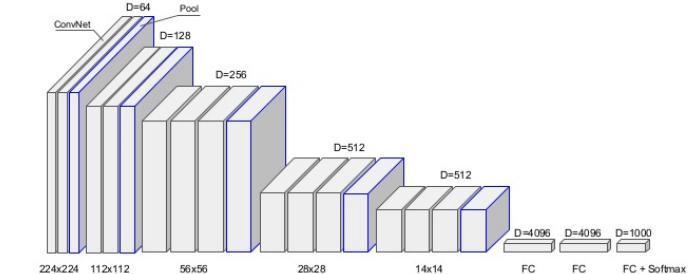


AlexNet (Krizhevsky et al. 2012)

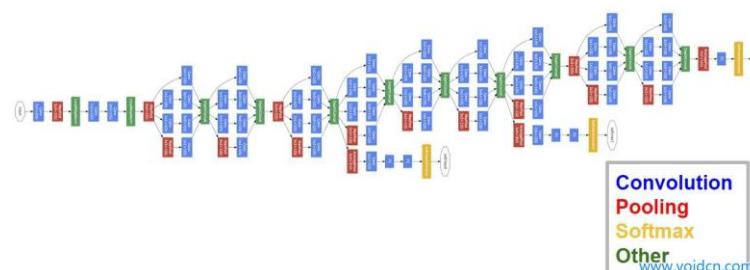


When AlexNet is processing an image, this is what is happening at each layer.

Classical CNN topology - VGGNet (2013)

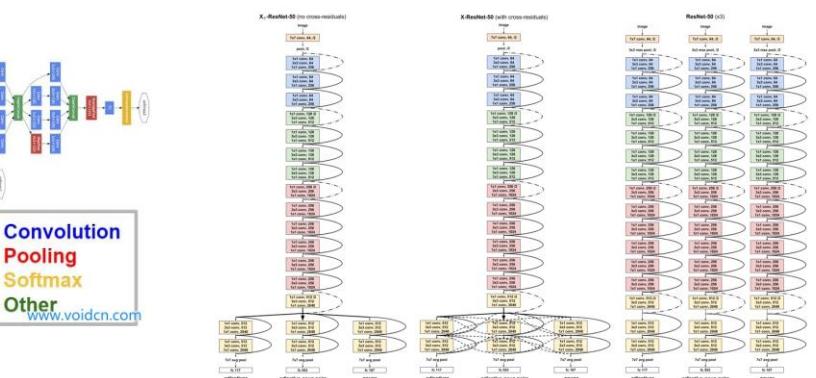


GoogLeNet



30

ResNet



Deep Water H2O and TensorFlow Demo



All None

Only show columns with more than % missing values.

epochs 500

How many times the dataset should be iterated (streamed), can be fractional.

ignore_const_cols

Ignore constant columns.

network lenet



Network architecture.

activation

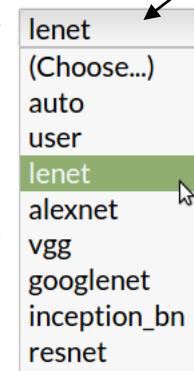
Activation function. Only used if no user-defined network architecture file is provided, and only for problem_type=dataset.

hidden

Hidden layer sizes (e.g. [200, 200]). Only used if no user-defined network architecture file is provided, and only for problem_type=dataset.

problem_type

Problem type, auto-detected by default. If set to image, the H2OFrame must contain a string column containing the path (URI or URL) to the images in the first column. If set to text, the H2OFrame must contain a string column containing the text in the first column. If set to dataset, Deep Water behaves just like any other H2O Model and builds a model on the provided H2OFrame (non-String columns).



Example: Deep Water + H₂O Flow Choosing different network structures

ADVANCED

GRID ?

checkpoint

Model checkpoint to resume training with.

autoencoder

Auto-Encoder.

balance_classes

Balance training data class counts via over/under-sampling (for imbalanced data).

fold_column

Column with cross-validation fold index assignment per observation.

offset_column

Offset column. This will be added to the combination of columns before applying the link function.



Flow ▾ Cell ▾ Data ▾ Model ▾ Score ▾ Admin ▾ Help ▾

Deep Water H2O and TensorFlow Demo



Choosing different backends (TensorFlow, MXNet, Caffe)

score_training_samples	10000	Number of training set samples for scoring (0 for all).	<input type="checkbox"/>
score_validation_samples	0	Number of validation set samples for scoring (0 for all).	<input type="checkbox"/>
score_duty_cycle	1	Maximum duty cycle fraction for scoring (lower: more training, higher: more scoring).	<input type="checkbox"/>
stopping_rounds	5	Early stopping based on convergence of stopping_metric. Stop if simple moving average of length k of the stopping_metric does not improve for k:=stopping_rounds scoring events (0 to disable)	<input type="checkbox"/>
stopping_metric	AUTO	Metric to use for early stopping (AUTO: logloss for classification, deviance for regression)	<input type="checkbox"/>
stopping_tolerance	0	Relative tolerance for metric-based stopping criterion (stop if relative improvement is not at least this much)	<input type="checkbox"/>
max_runtime_secs	0	Maximum allowed runtime in seconds for model training. Use 0 to disable.	<input type="checkbox"/>
backend	tensorflow ▾	Deep Learning Backend.	<input type="checkbox"/>
image_shape	28,28	Width and height of image.	<input type="checkbox"/>
channels	3	Number of (color) channels.	<input type="checkbox"/>
network_definition_file		Path of file containing network definition (graph, architecture).	<input type="checkbox"/>
network_parameters_file		Path of file containing network (initial) parameters (weights, biases).	<input type="checkbox"/>
mean_image_file		Path of file containing the mean image data for data normalization.	<input type="checkbox"/>
export_native_parameters_prefix		Path (prefix) where to export the native model parameters after every iteration.	<input type="checkbox"/>
input_dropout_ratio	0	Input layer dropout ratio (can improve generalization, try 0.1 or 0.2).	<input type="checkbox"/>
hidden_dropout_ratios		Hidden layer dropout ratios (can improve generalization), specify one value per hidden layer, defaults to 0.5.	<input type="checkbox"/>

H_2O + xgboost

Brand-new: H2O XGBoost Integration (Gradient Boosting)

Why XGBoost?

Competitive **accuracy** and **speed** (great for Kaggle)

GPU support (for small/medium data)

Efficient on **sparse** data

Why integrate into H2O?

Ease of use (**Flow** GUI, R/Py APIs)

Real-time model status (var imp, metrics)

Efficient **data preprocessing** (sparse, categorical)

Integration into **H2O ecosystem** (modeling, deployment, support)

Live Demo of GPU Gradient Boosting in H2O

Build a Model

Select an algorithm: XGBoost

booster gbtree

reg_lambda 1

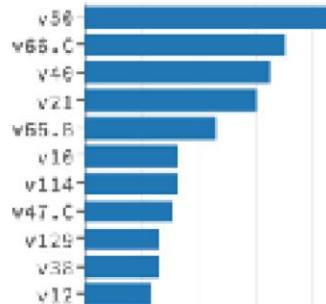
reg_alpha (Choose...)

auto

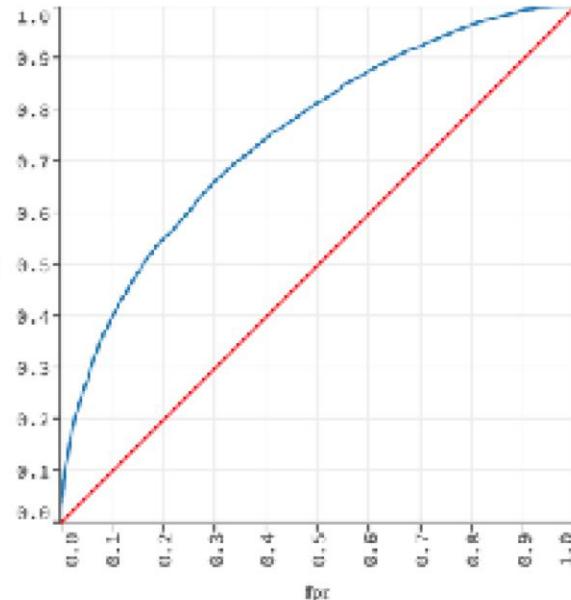
backend gpu

cpu

VARIABLE IMPORTANCES



* ROC CURVE - VALIDATION METRICS , AUC = 0.750331



```
+-----+  
| NVIDIA-SMI 375.28      Driver Version: 375.28 |  
+-----+  
| GPU  Name     Persistence-M| Bus-Id     Disp.A  Volatile Uncorr. ECC |  
| Fan  Temp  Perf  Pwr.Usage/Cap| Memory-Usage  GPU-Util  Compute M. |  
|-----+-----+-----+-----+-----+-----+-----+  
|  0  GeForce GTX 1080     Off  | 0000:02:00.0  On   | N/A  |  
| 27%   43C    P2    83W / 188W | 2848MiB / 8145MiB | 94%  | Default |  
+-----+
```

```
gbm = h2o.get_model("h2o-gbm")
xgb = h2o.get_model("h2o-xgboost")
print("H2O GBM:      Validation AUC=%r" % gbm.auc(valid=True))
print("H2O XGBoost: Validation AUC=%r" % xgb.auc(valid=True))

H2O GBM:      Validation AUC=0.7517126133633125
H2O XGBoost: Validation AUC=0.750331001716736
```

Stacked Ensembles

Stacked Ensembles in H2O



https://github.com/h2oai/h2o-meetups/blob/master/2017_02_23_Metis_SF_Sacked_Eensemles_Deep_Water/stacked_ensembles_in_h2o_feb2017.pdf

February 2017

Erin LeDell Ph.D.
Machine Learning Scientist

H₂O.ai

Available in both Python and R

Model Stacking

```
In [23]: # Define a list of models to be stacked
# i.e. best model from each grid
all_ids = [best_gbm_model_id, best_drf_model_id, best_dnn_model_id]

In [24]: # Set up Stacked Ensemble
ensemble = H2OStackedEnsembleEstimator(model_id = "my_ensemble",
                                         base_models = all_ids)

In [25]: # use .train to start model stacking
# GLM as the default metalearner
ensemble.train(x = features,
                y = 'quality',
                training_frame = wine_train)

stackedensemble Model Build progress: |██████████| 100%
```

Comparison of Model Performance on Test Data

```
In [26]: print('Best GBM model from Grid (MSE) : ', best_gbm_from_rand_grid.model_performance(wine_test).mse())
print('Best DRF model from Grid (MSE) : ', best_drf_from_rand_grid.model_performance(wine_test).mse())
print('Best DNN model from Grid (MSE) : ', best_dnn_from_rand_grid.model_performance(wine_test).mse())
print('Stacked Ensembles (MSE) : ', ensemble.model_performance(wine_test).mse())

Best GBM model from Grid (MSE) : 0.4013942890547201
Best DRF model from Grid (MSE) : 0.478156285687009
Best DNN model from Grid (MSE) : 0.489784141303471
Stacked Ensembles (MSE) : 0.39965430199959595
```

© 2017 GitHub, Inc. Terms Privacy Security Status Help

Contact GitHub API Training Shop Blog About

Model Stacking

```
In [20]: # Define a list of models to be stacked
# i.e. best model from each grid
all_ids = list(best_gbm_model_id, best_drf_model_id, best_dnn_model_id)

In [21]: # Stack models
# GLM as the default metalearner
ensemble = h2o.stackedEnsemble(x = features,
                                y = 'quality',
                                training_frame = wine_train,
                                model_id = "my_ensemble",
                                base_models = all_ids)
```

|=====| 100%

Comparison of Model Performance on Test Data

```
In [22]: cat('Best GBM model from Grid (MSE) : ', h2o.performance(best_gbm_from_rand_grid, wine_test)$MSE,
      "\n")
cat('Best DRF model from Grid (MSE) : ', h2o.performance(best_drf_from_rand_grid, wine_test)$MSE,
      "\n")
cat('Best DNN model from Grid (MSE) : ', h2o.performance(best_dnn_from_rand_grid, wine_test)$MSE,
      "\n")
cat('Stacked Ensembles (MSE) : ', h2o.performance(ensemble, wine_test)$MSE, "\n")

Best GBM model from Grid (MSE) : 0.4013943
Best DRF model from Grid (MSE) : 0.4781568
Best DNN model from Grid (MSE) : 0.5543555
Stacked Ensembles (MSE) : 0.3989076
```

© 2017 GitHub, Inc. Terms Privacy Security Status Help

Contact GitHub API Training Shop Blog About

https://github.com/woobe/odsc_h2o_machine_learning

Automatic Machine Learning (AutoML)

H2O AutoML

- AutoML stands for “Automatic Machine Learning”
- The idea here is to remove most (or all) of the parameters from the algorithm, as well as automatically generate derived features that will aid in learning.
- Single algorithms are tuned automatically using a carefully constructed random grid search.
- Optionally, a Stacked Ensemble can be constructed.

Public code coming soon!

Model Interpretation

Ideas on interpreting machine learning

Mix-and-match approaches for visualizing data and interpreting machine learning models and results.

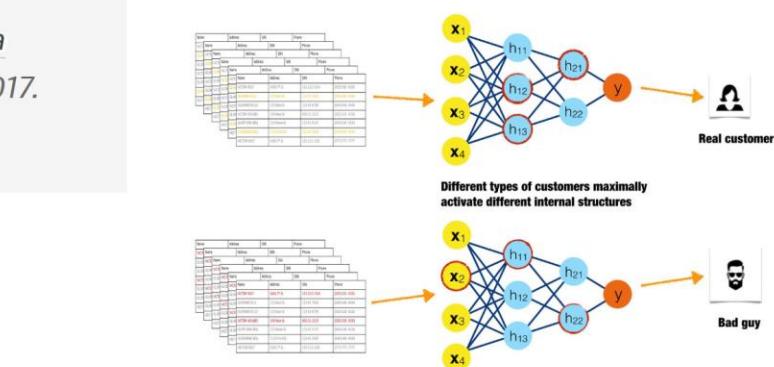
By Patrick Hall, Wen Phan and SriSatish Ambati, March 15, 2017

For more on advances in machine learning, prediction, and technology, check out the [Data science and advanced analytics sessions at Strata + Hadoop World London, May 22-25, 2017](#). Early price ends April 7.

You've probably heard by now that machine learning algorithms can use big data to predict whether a donor will give to a charity, whether an infant in a NICU will develop sepsis, whether a customer will respond to an ad, and on and on. Machine learning can even [drive cars](#) and [predict elections](#).

... Err, wait. Can it? I believe it can, but these recent high-profile hiccups should leave everyone who works with data (big or not) and machine learning algorithms asking themselves some very hard questions: do I understand my data? Do I understand the model and answers my machine learning algorithm is giving me? And do I trust these answers?

Unfortunately, the complexity that bestows the extraordinary predictive



Inputs activating different neurons in a neural network.
(source: Image courtesy of Patrick Hall and the h2o.ai team, used with permission)

Correlation Graphs

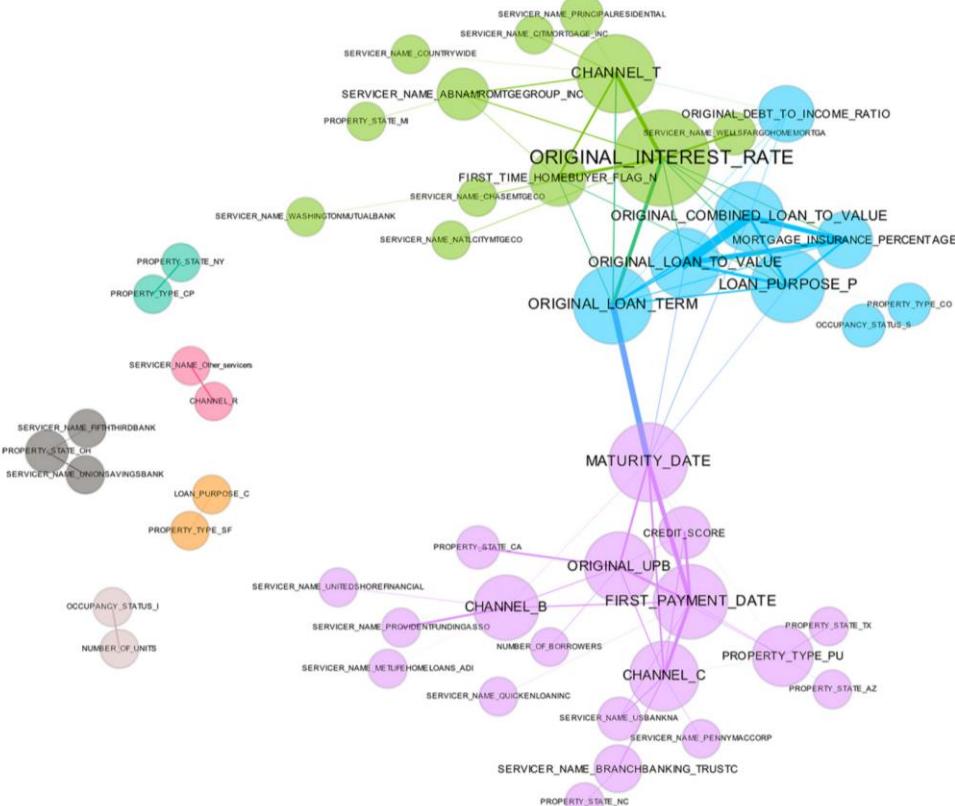


Figure 3. A correlation graph representing loans made by a large financial firm. Figure courtesy of Patrick Hall and the H2O.ai team.

Partial dependence plots

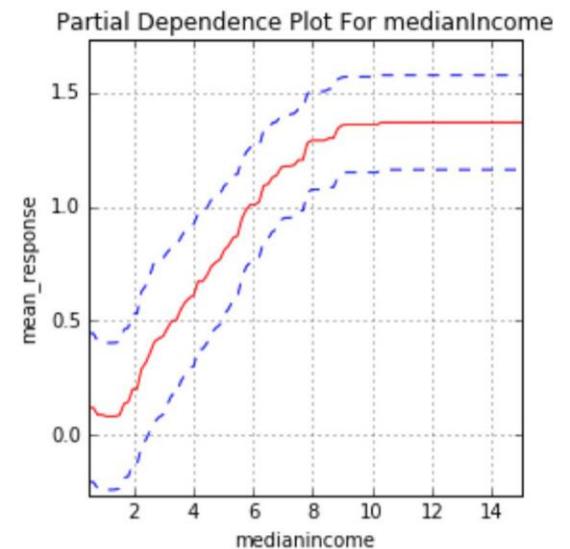
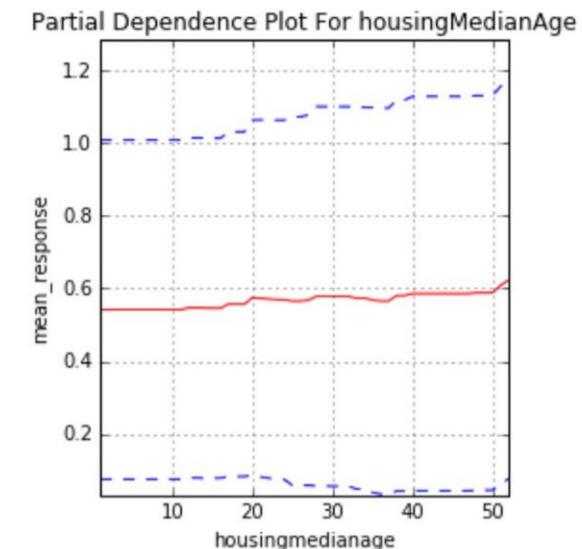


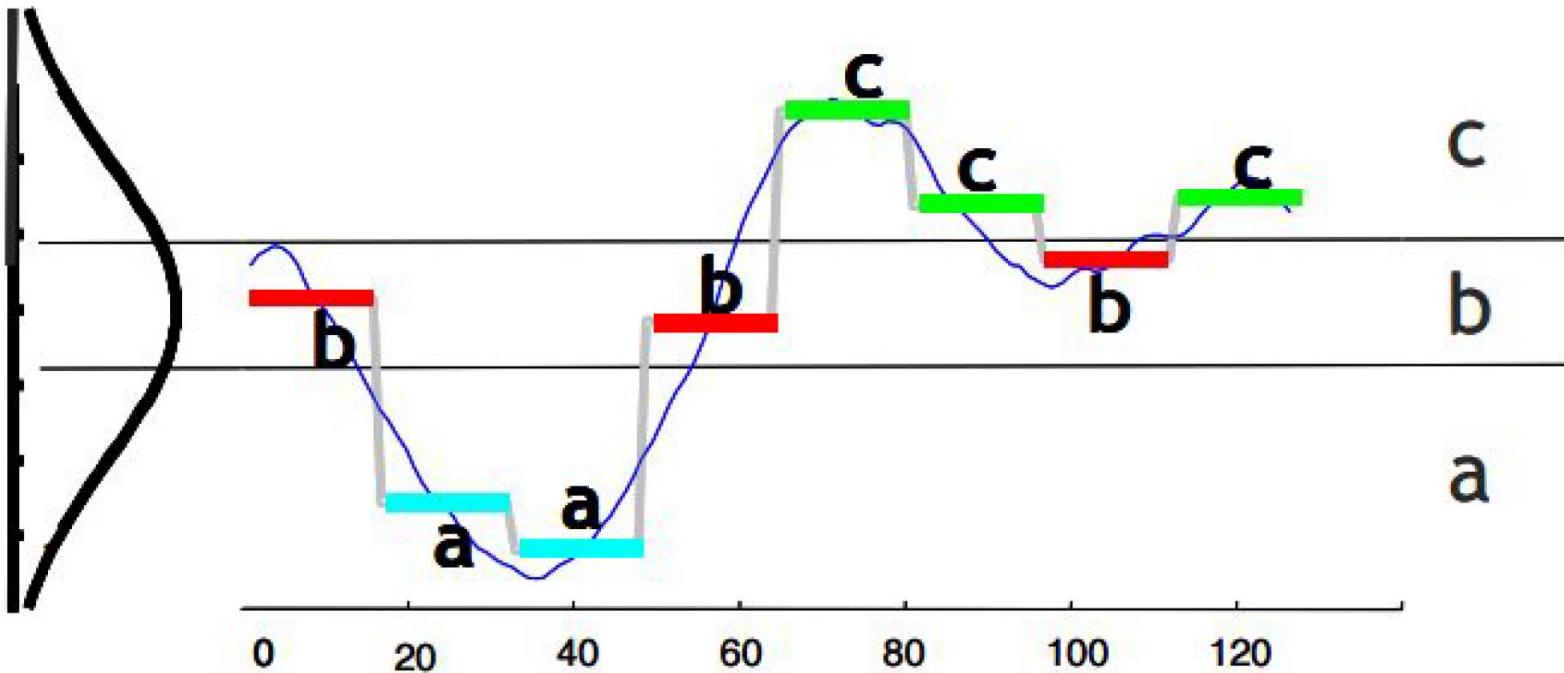
Figure 5. One-dimensional partial dependence plots from a gradient boosted tree ensemble model of the well-known California housing data set. Image courtesy Patrick Hall and the H2O.ai team.

Time Series in H₂O

ISAX

- Time series data compression algorithm, implemented on H2O's distributed architecture
- Find groups of similar time series patterns in one billion
 - <http://blog.h2o.ai/2016/11/indexing-1-billion-time-series-with-h2o-and-isax/>
- ISAX applied to fields such as IoT, Finance, Bioinformatics, Image/Sound processing.
 - <http://www.cs.ucr.edu/~eamonn/SAX.htm>
- Compress your time series data and run:
 - Clustering
 - Classification
 - Anomaly Detection
 - Predictive analytics

ISAX



Blog Post

<http://blog.h2o.ai/2016/11/indexing-1-billion-time-series-with-h2o-and-isax/>

Indexing 1 Billion Time Series with H2O and ISax

At H2O, we have recently debuted a new feature called ISax that works on time series data in an H2O Dataframe. ISax stands for Indexable Symbolic Aggregate APPROXimation, which means it can represent complex time series patterns using a symbolic notation and thereby reducing the dimensionality of your data. From there you can run H2O's ML algos or use the index for searching or data analysis. ISax has many uses in a variety of fields including finance, biology and cybersecurity.

Today in this blog we will use H2O to create an ISax index for analytical purposes. We will generate 1 Billion time series of 256 steps on an integer U(100,100) distribution. Once we have the index we'll show how you can search for similar patterns using the index.

We'll show you the steps and you can run along, assuming you have enough hardware and patience. In this example we are using a 9 machine cluster, each with 32 cores and 256GB RAM. We'll create a 1B row synthetic data set and form random walks for more interesting time series patterns. We'll run ISax and perform the search, the whole process takes ~30 minutes with our cluster.

Raw H2O Frame Creation

In the typical use case, H2O users would be importing time series data from disk. H2O can read from local filesystems, NFS, or distributed systems like Hadoop. H2O cluster file reads are parallelized across the nodes for speed. In our case we'll be generating a 256 columns, 1B row frame. By the way H2O Dataframes scales better by increasing rows instead of columns. Each row will be an individual time series. The ISax algo assumes the time series data is row based.

```
rawdf = h2o.create_frame(cols=256, rows=1000000000, real_fraction=0.0, integer_fraction=1.0, missing_fraction=0.0)

In [4]: print(datetime.datetime.now())
        rawdf = h2o.create.frame(cols<256, rows=1000000000, real_fraction=0.0, integer_fraction=1.0, missing_fraction=0.0)
        print(datetime.datetime.now())
2016-11-08 13:11:00.456099
Create Frame progress: [██████████] 100%
2016-11-08 13:15:40.852390
```

Random Walk

Here we do a row wise cumulative sum to simulate random walks. The .head call triggers the execution graph so we can do a time measurement.

```
tsdf = rawdf.cumsum(axis=1)
print tsdf.head()

In [7]: print(datetime.datetime.now())
        tsdf = rawdf.cumsum(axis=1)
        print tsdf.head()
        print(datetime.datetime.now())
2016-11-08 13:44:42.466140
          C1   C2   C3   C4   C5   C6   C7   C8   C9
-3  -40  -104  -128  -178  -180  -202  -138  -229
-23  -54  -86  -170  -154  -235  -294  -312  -357
29   19   -2  -86  -24  -11  -3  58  78
88  123   80  -20  36  -63  -25  52  1
-63  -116  -194  -181  -167  -188  -214  -188  -101
-9  -18  -34  -48  -46  -45  -12  -58  -21
100  118  216  228  258  345  279  262  313
-1  34   59  -26  -39  4   84  113  36
```

Lets take a quick peek at our time series

```
tsdf[0:2,:].transpose().as_data_frame(use_pandas=True).plot()

In [8]: tsdf[0:2,:].transpose().as_data_frame(use_pandas=True).plot()
Out[8]: <matplotlib.axes._subplots.AxesSubplot at 0x7f59bc14ad90>
```

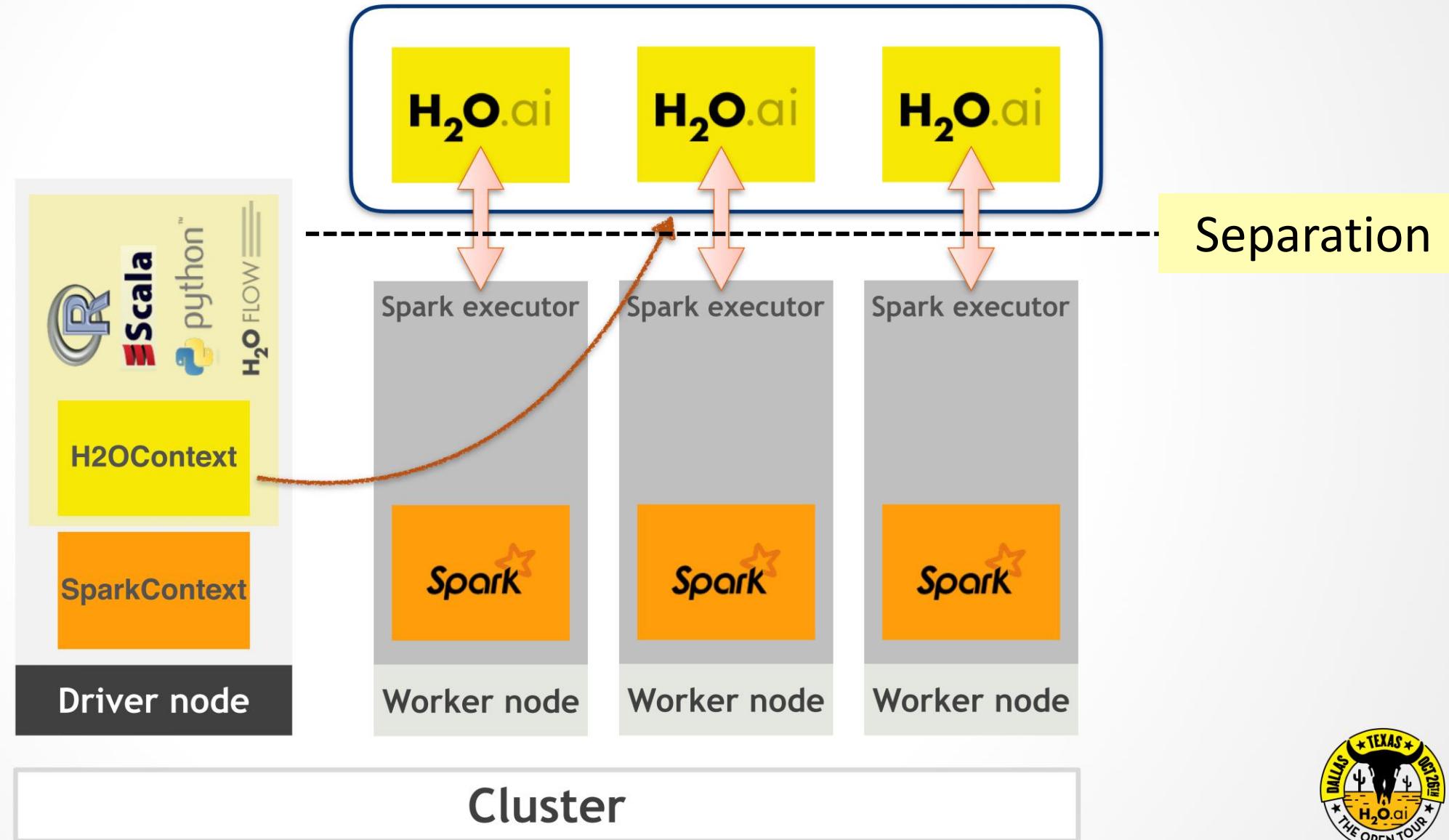
Sparkling Water

H₂O's integration with Spark – latest developments

External H₂O Backend for Sparkling Water

- High Availability Mode
 - Separating Spark and H₂O
 - while preserving same API
- Advantages
 - H₂O does not crash when Spark executor goes down
 - Better resource management since resources can be planned per tool
- Disadvantages
 - Transfer overhead between Spark and H₂O processes
 - Under measurement with cooperation of a customer
- Links
 - https://github.com/h2oai/sparkling_water/blob/master/doc/backends.md

High Availability



H₂O Community

H₂O Meetup Groups

H2O
29,796 members 21 interested **43 Meetups** (circled)
39 cities 18 countries

[Find a Meetup Group near you](#)

- www.meetup.com/topics/h2o/all/

- Need your help 😊

All H2O Meetups

Silicon Valley Big Data Science
6,663 Makers | Mountain View, CA

NYC Big Data Science
4,319 makers | New York, NY

San Francisco Big Data Science
2,784 Makers | San Francisco, CA

H2O and Data Science!
1,877 Hackers | Mountain View, CA

Dallas Big Data Science!
1,635 Data Scientists | Dallas, TX

London Artificial Intelligence & Deep Learning
914 Makers | London, United Kingdom

Amsterdam Artificial Intelligence & Deep Learning
753 Makers | Amsterdam, Netherlands

Toronto Artificial Intelligence & Deep Learning
684 Makers | Toronto, ON

Sydney Artificial Intelligence & Deep Learning
559 Makers | Sydney, Australia

Bangalore Artificial Intelligence & Deep Learning
526 Makers | Bangalore, India

Berlin Artificial Intelligence & Deep Learning
520 Makers | Berlin, Germany

Singapore Artificial Intelligence & Deep Learning
519 Makers | Singapore, Singapore

Hyderabad Artificial Intelligence & Deep Learning
507 Makers | Hyderabad, India

Madrid Artificial Intelligence & Deep Learning
500 Fabricantes | Madrid, Spain

Pune Artificial Intelligence & Deep Learning
459 Makers | Pune, India

#AroundTheWorldWithH2Oai

#aroundtheworldwithh2oai

TOP LATEST PEOPLE PHOTOS VIDEOS NEWS BROADCASTS

Search filters · Show

Who to follow · Refresh · View all

- Monica Rogati @mrogati Follow
- Julia Silge @juliasilge Follow
- David Robinson @drob Follow

Find friends

United Kingdom Trends · Change

- #WorldPhysicalActivityDay 4,215 Tweets
- #ThursdayThoughts · LIVE 20.4K Tweets
- Supreme Court Court rules children can't be taken out of school without permission
- #GenderPayGap UK companies must publish their gender pay gaps
- #ApprenticeshipLevy
- Michael Caine Sir Michael Caine explains why he backs Brexit
- #sagesummit 10.6K Tweets
- Aintrue 12.1K Tweets
- Samantha Baldwin Missing mother Samantha Baldwin and sons found safe
- Barack Obama 40.1K Tweets





Questions

Jobs

Documentation
BETA

Tags

Users

 h2o

1

1

1



?



Search

[Ask Question](#)

h2o

[search](#)

1,130 results

[relevance](#)[newest](#)[votes](#)[active](#)

3

votes

1
answer

Q: Subsetting in H2O R

I have a **h2o** object. The standard R for subset sub1<-trans[trans\$type==1,] I tried the same in **h2o**. It is not working sub1<-trans[trans\$type==1,] I also tried sub1<-h2o.exec(trans[trans\$type==1,]) note* trans is a **h2o** data Object. Any idea to do it in **h2o**? Thanks ...

[r](#) [subset](#) [h2o](#)asked Nov 28 '14 by [chee.work.stuff](#)

5

votes

4
answers

Q: h2o implementation in R

I am learning **h2o** package now, I installed **h2o** package from CRAN and couln't run this code ## To import small iris data file from H\:\\sub:`2`\\ O's package irisPath = system.file("extdata ... ", "iris.csv", package="h2o") iris.hex = h2o.importFile(localH2O, path = irisPath, key = "iris.hex") I am getting the below error, Error in h2o.importFile(localH2O, path = irisPath, key ...

[r](#) [h2o](#)asked Aug 22 '16 by [varun](#)

0

votes

1
answer

Q: MAPE metric at h2o

What is correct way to implement MAPE under **h2o** framework? I am interested to convert below function to **h2o** concept def mape(a, b): mask = a > 0 return (np.fabs(a - b)/a)[mask].mean() ...

[python](#) [python-3.x](#) [pandas](#) [dataframe](#) [h2o](#)asked Mar 29 by [SpanishBoy](#)

1

vote

Q: H2O Python module build problems

This is the error encountered when running ./gradlew build in the **h2o-3** repository. :h2o-py:buildDistFound packages: ['h2o', 'h2o.backend', 'h2o.estimators', 'h2o.grid', 'h2o.model', 'h2o.schemas ... : invalid command 'hdist_wheel': h2o-py:buildDist FAILED :h2o-py:buildDist took 0.273 secs FAILURE: Build failed with an

[Advanced Search Tips](#)results found containing
h2o**Want a [python](#) job?****Software Engineer - Back End**

Booking.com Amsterdam, Netherlands

€60K - €70K RELOCATION

[python](#) [perl](#)**Ervaren Python developer**

NRC Media Amsterdam, Netherlands

[python](#) [node.js](#)**Hot Network Questions**

What can the United Nations do after the chemical attack in Syria?

Does a LED also emit light when conducting in avalanche mode?

Can a high voltage line kill a person without touching it?

How do I compute the mean from ASCII file data in bash?

In a family business, do I refer to people by their

PyData Amsterdam 2017

Conference Schedule

[View past PyData event schedules here.](#)

Tutorial Sessions – Friday April 7, 2017

	Big Room	Small Room
8:30		Registration
9:00	From Fourier to deep convnets Ivo Everts	Pandas from the Inside / "Big Pandas" Stephen Simmons
12:00	Lunch	
13:00	Deep learning - advanced techniques Geoff French	So You Want to Be a Python Expert? James Powell
16:00	Hacking Space	Introduction to Machine Learning with H2O and Python Jo-fai Chow
18:30	PyData Amsterdam Meetup RSVP separately (Space is Limited!)	
22:00		

H₂O Tutorial
Friday – 4:00 pm
(All Materials will be available on GitHub)

General Sessions – Saturday April 8, 2017

	Big Room	Small Room
8:00		Breakfast & Registration
9:00		Opening Notes
9:15		Ethical Machine Learning: Creating Fair Models in an Unjust World Katharine Jarmul
10:00	Knowledge Repository Matthew Wardrop, Dan Frank	Creativity and AI: Deep Neural Nets "Going Wild" Roelof Pieters
10:45		Break
11:00	Successfully applying Bayesian statistics to A/B testing in your business Ruben Mak	Detecting Clickbaits using Machine Learning Abhishek Thakur
11:45	Deep Reinforcement Learning: theory, intuition, code Maxim Lapan	Making contract documents fully searchable at KPN Gianluigi Bardelloni
12:30		Lunch
13:30	Diagnosing Machine Learning Models Lucas Javier Bernardi	A Pythonic Tour of Neo4j and the Cypher Query Language Nigel Small
14:15	Training a TensorFlow model to detect lung nodules on CT scans Mark-Jan Harte	A practical guide to speed up your application with Asyncio Niels Denissen
15:00		Break
15:15	Bayesian optimization with Scikit-Optimize Gilles Louppe	Different Strategies of Scaling H2O Machine Learning on Apache Spark Jakub Hava

Lucas' Talk
Saturday – 1:30 pm

Sparkling Water Talk
Saturday – 3:15 pm

Thanks!

- Organizers & Sponsors

Booking.com



- Find us at PyData Conference
 - Live Demos

- Code, Slides & Documents

- bit.ly/h2o_meetups
- docs.h2o.ai

- Contact

- joe@h2o.ai
- [@matlabulous](https://twitter.com/matlabulous)
- github.com/woobe

- Please search/ask questions on
Stack Overflow

- Use the tag `h2o` (not H2 zero)