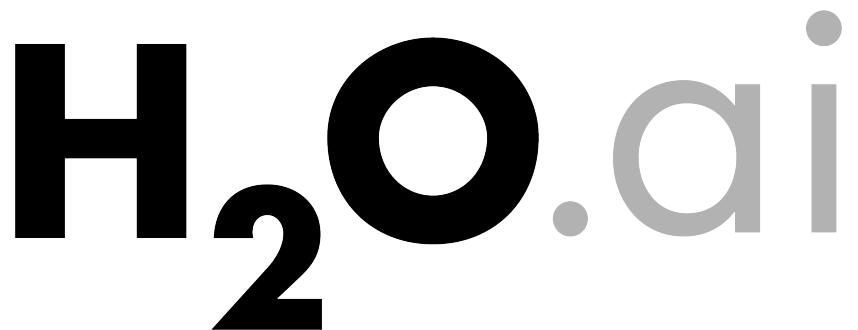


# Introduction to Machine Learning with H<sub>2</sub>O

- Introduction (30 mins)
  - From Engineering to Data Science – Why?
  - Machine Learning Basics
  - H<sub>2</sub>O Machine Learning Platform
- Tutorial (60 mins)
  - Regression (House Price)
  - Classification (Human Activities)
  - Clustering (Water Treatment Plant)
- Extra: More About H<sub>2</sub>O (15 mins)
  - Use Cases
  - H<sub>2</sub>O on a Multi-Node Cluster
  - Next-Gen H<sub>2</sub>O
- Q & A (15 mins)



Jo-fai (Joe) Chow  
Data Scientist at H<sub>2</sub>O.ai  
joe@h2o.ai

Download Link: [http://bit.ly/h2o\\_exeter\\_2017](http://bit.ly/h2o_exeter_2017)

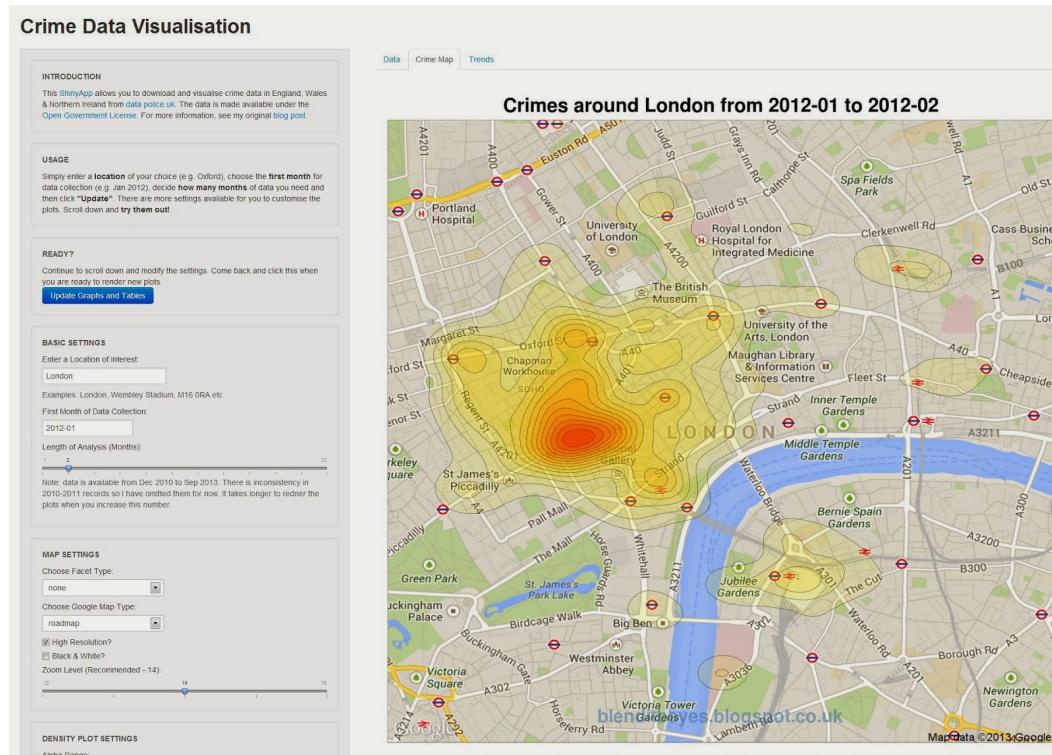
# From Engineering to Data Science

- My Story

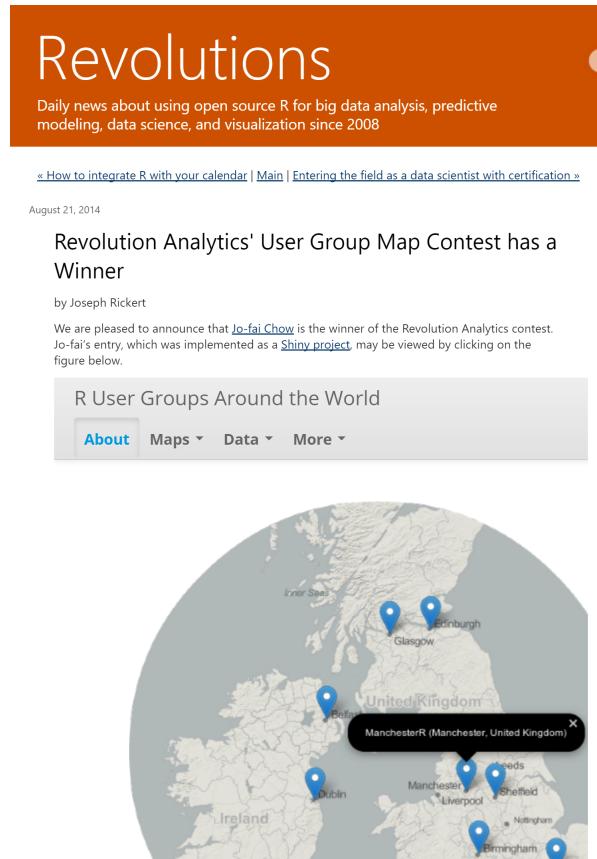
# About Me

- Civil (Water) Engineer
  - 2005 – 2010
    - Consultant (SEAMS, UK)
      - Utilities / Asset Management
      - Constrained Optimization
  - 2010 – 2014
    - STREAM EngD (Exeter)
      - Infrastructure Design Optimisation
      - Machine Learning + Water Engineering
      - **Discovered H<sub>2</sub>O in 2014**
- Data Scientist
  - 2015 – 2016
    - Virgin Media (UK)
    - Domino Data Lab (Silicon Valley)
  - 2016 – Present
    - H<sub>2</sub>O.ai (Silicon Valley)
  - How?
    - [bit.ly/joe\\_kaggle\\_story](http://bit.ly/joe_kaggle_story)

# About Me – I ❤️ DataViz



My First Data Viz & Shiny App Experience  
[CrimeMap \(2013\)](#)



Revolution Analytics' Data Viz Contest  
[RUGSMAPS \(2014\)](#)



Jo-fai (Joe) Chow  
@matlabulous

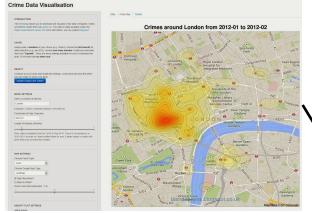
Thank you very much @RevolutionR  
@revodavid @RevoJoe #iloveR  
bit.ly/rugsmaps #Shiny #rMaps



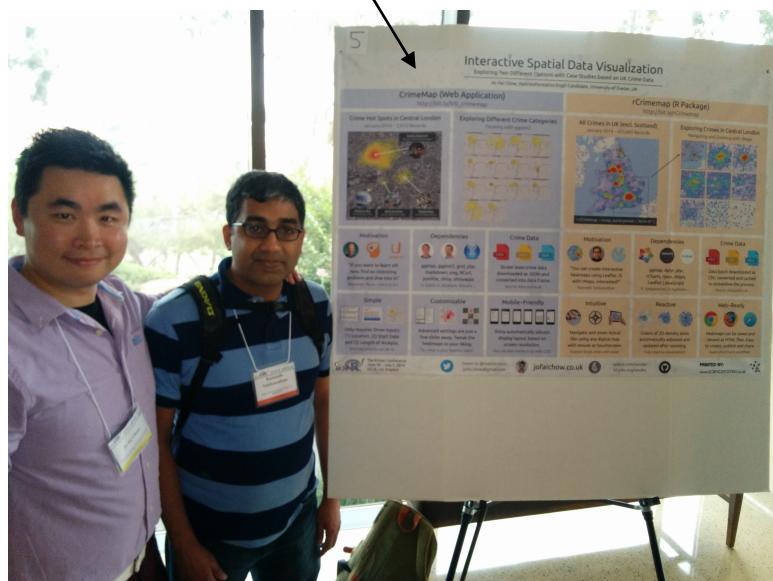
RETWEETS  
3  

1:25 AM - 29 Aug 2014

# useR! 2014



CrimeMap -> poster



H2O.ai @h2oai · Following

#linus for #rstat #user2014 #JohnChambers

H2O.ai @h2oai · Following

#SparkSummit ~~~> #useR2014  
#h2o

Some Promising Projects

1. Rcpp, Rcpp++: Interface to C++ and programming in C++ [CRAN, for Rcpp]
2. LLVM for R: Compiling toolkit for R [omegahat.org for RLLVM, RLLVMCompile]
3. h2o: Interface and Java-based computations for big data [CRAN]

John Chambers mentioned H<sub>2</sub>O

H2O.ai @h2oai · 28 Jul 2014

Smoooooth! - if I have to explain it in one word. Oxdata made this really easy for #R users. r-bloggers.com/things-to-try-... #Thanks #JoFaiChow

Things to try after useR! – Part 1: Deep Learning wit...

Annual R User Conference 2014The useR! 2014 conference was a mind-blowing experience. Hundreds of R enthusiasts and the beautiful UCLA campus, I am rea...

r-bloggers.com

Jo-fai (Joe) Chow  
@matlalobus

Replying to @h2oai

Hi @srisatish @ArnoCandel and every1  
@hexadata thx 4 making and open-sourcing  
the powerful #H2O shd hv tried it during (not  
after) #user2014

LIKES  
2

1:41 PM - 28 Jul 2014

H<sub>2</sub>O.ai

# About Me – I ❤️ Kaggle

The screenshot shows a blog post on the Domino Data Lab website. The header features a dark blue background with abstract white shapes. The title 'How to use R, H2O, and Domino for a Kaggle competition' is centered above a guest post by Jo-Fai Chow. Below the title, there's a note about a sample project being available on Domino, followed by a list of three tutorials. The introduction section discusses the purpose of the post as a sequel to a previous one. The footer contains links to the Domino App Site, Twitter, and email.

19 Sep 2014 •

Like 0 Tweet 21 g+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.

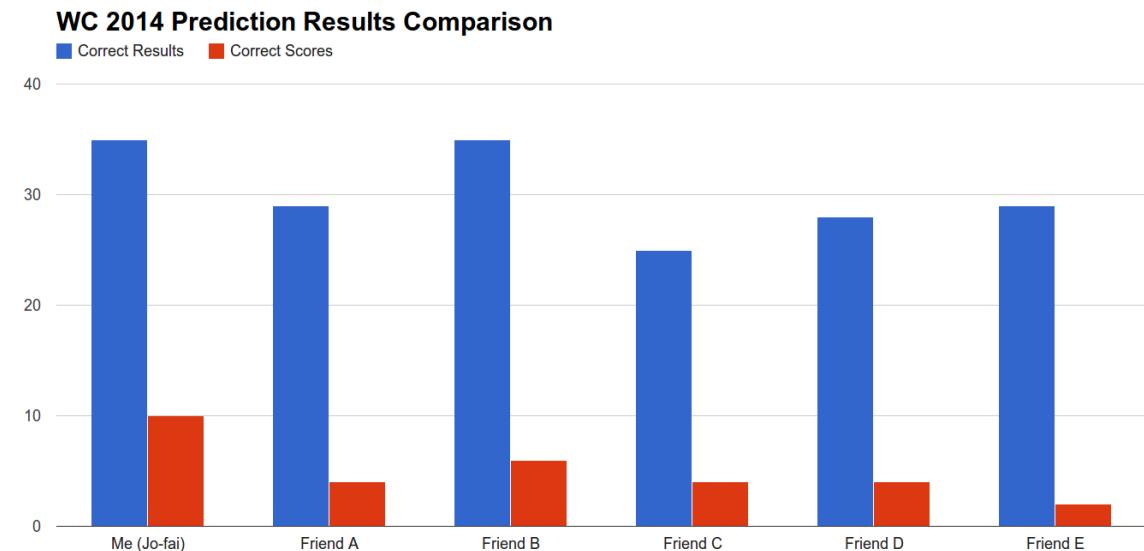
R + H<sub>2</sub>O + Domino for Kaggle  
[Guest Blog Post for Domino & H<sub>2</sub>O \(2014\)](#)

- The Long Story
  - [bit.ly/joe\\_kaggle\\_story](http://bit.ly/joe_kaggle_story)

# World Cup 2014

- Machine Learning vs My Friends

- Team performance data from web
- Simple machine learning models
- Objective: Predict Correct Score
- Me : 10 out of 64 (15.6%)
- Friends' Avg : 4 out of 64 (6.3%)



# From Engineering to Data Science

- Knowledge Driven Decision Making
  - Write code to deal with data based on knowledge and assumptions.
  - Strong domain knowledge.
- Data Driven Decision Making
  - Define a problem and collect data related to the problem.
  - Write code to show data to computers and let computers (algorithms) find patterns in data.
  - Domain knowledge is helpful but not necessary.





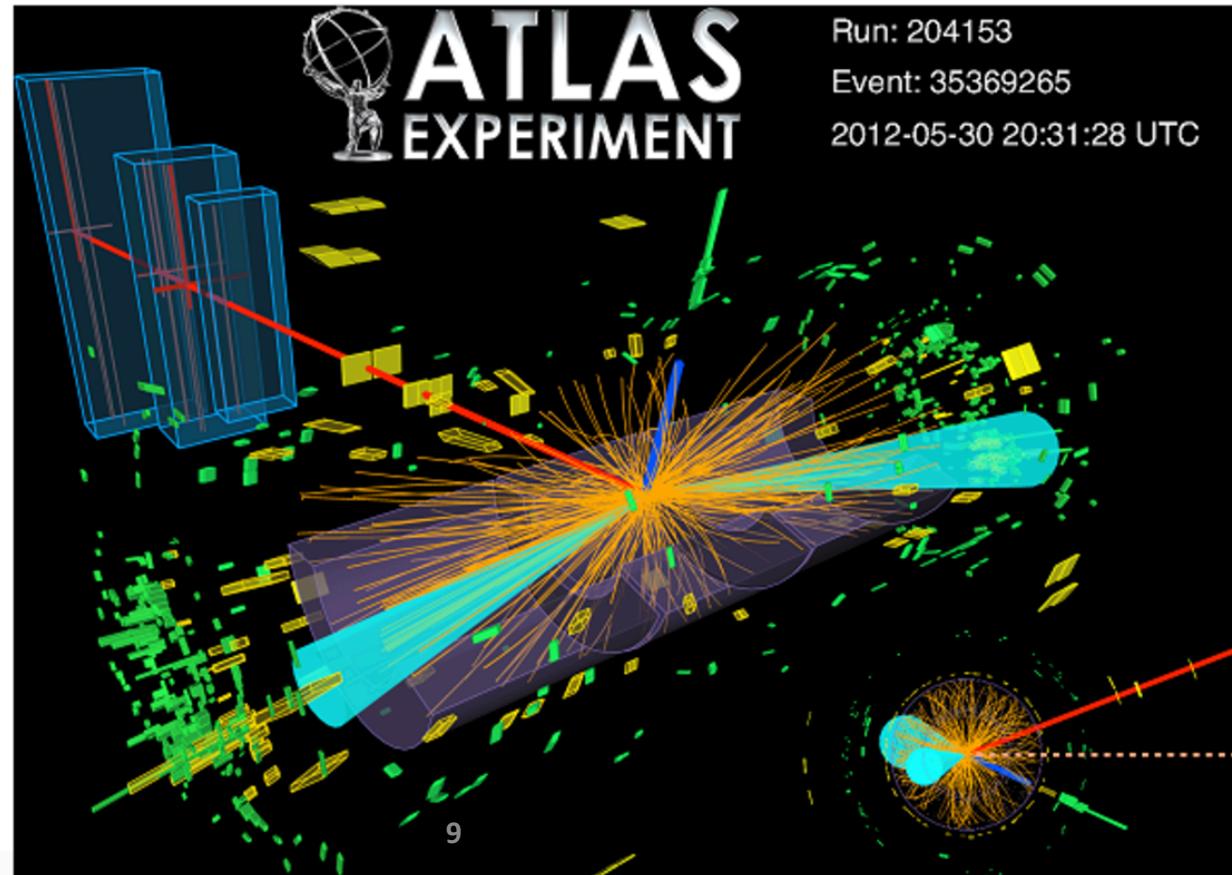
## Higgs Boson Machine Learning Challenge

Use the ATLAS experiment to identify the Higgs boson

\$13,000 · 1,785 teams · 3 years ago

[Overview](#)[Data](#)[Discussion](#)[Leaderboard](#)[Rules](#)[Team](#)[My Submissions](#)[Late Submission](#)[Overview](#)

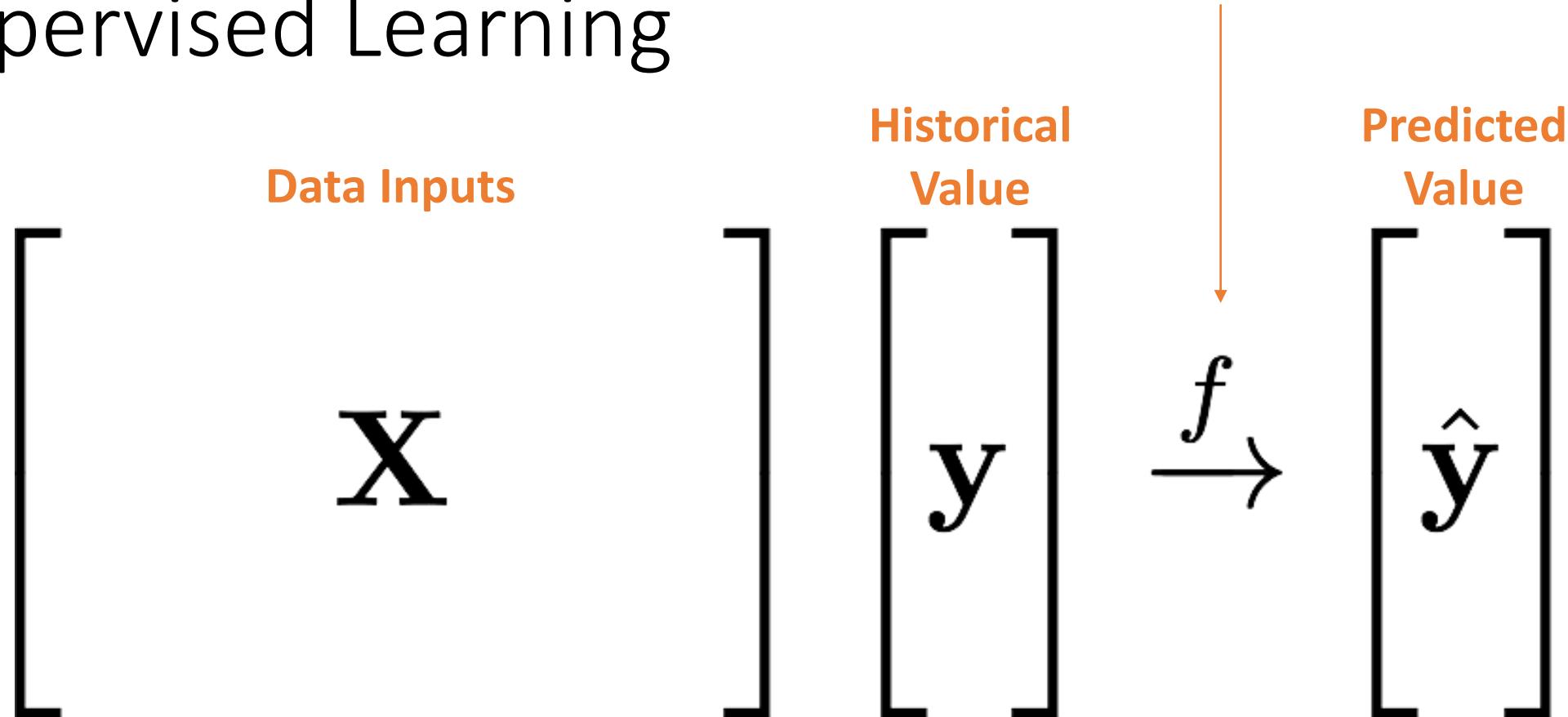
<https://www.kaggle.com/c/higgs-boson>

[Description](#)[Evaluation](#)[Prizes](#)[About The Sponsors](#)[Timeline](#)[Winners](#)

# What is Machine Learning?

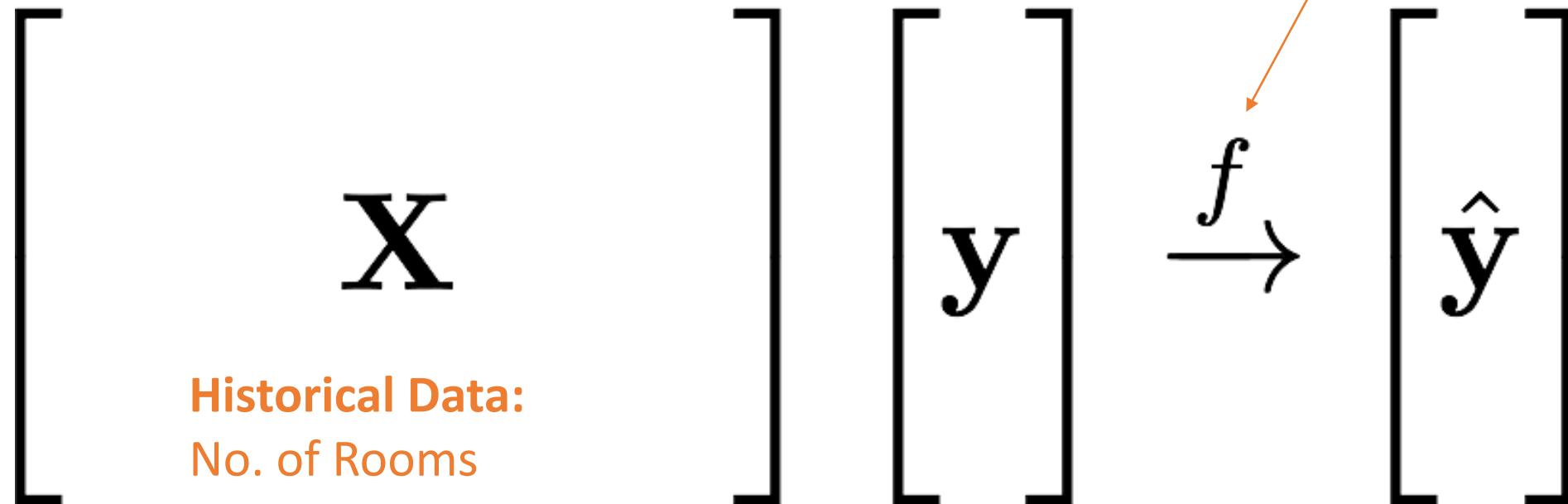
- Supervised Learning
- Unsupervised Learning

# Supervised Learning



# Supervised Learning Example

Machine Learning:  
Learn Patterns  
from Data



Target:  
House Value

Predicted Value  
(for evaluation)

# Supervised Learning Example

[

**X**

New Data:

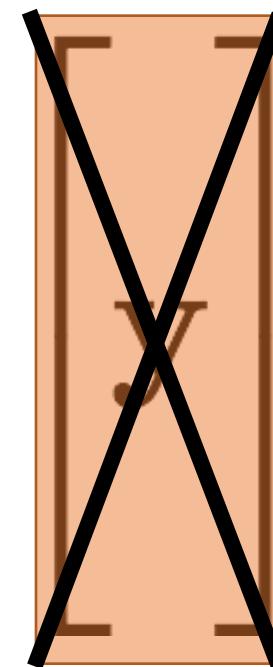
No. of Rooms

Crime Rate

Pupil-Teacher Ratio

...

]



Target:  
Unknown

Patterns Learned  
from Historical Data

$$f \rightarrow$$

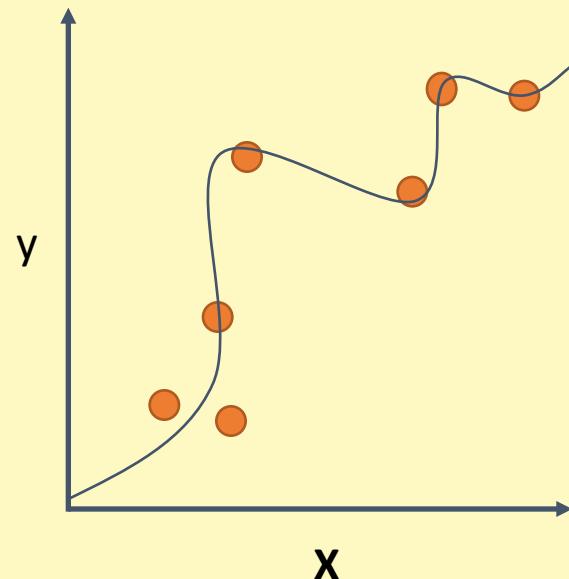
[  $\hat{y}$  ]

Predicted Value  
(for decision making)

# Supervised Learning – You Already Have Target Data

**Regression:**

**How much will a customer spend?**

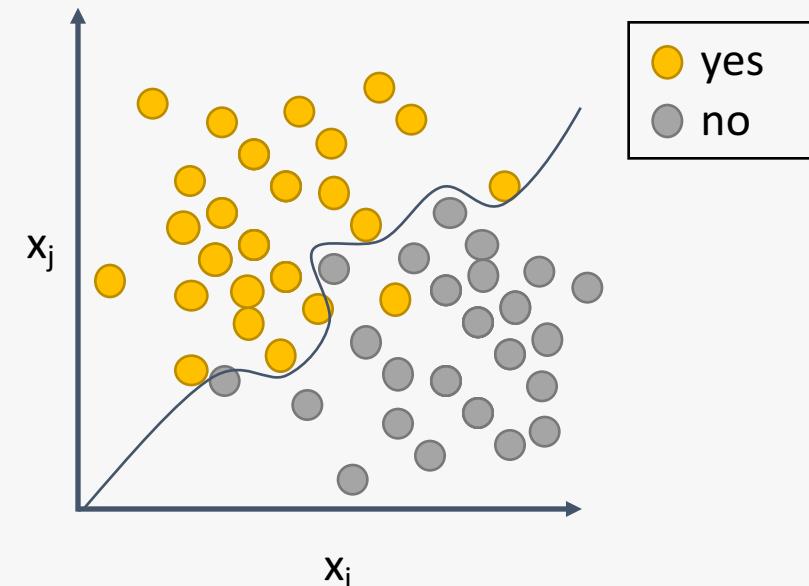


**H<sub>2</sub>O algos:**

**Penalized Linear Models**  
**Random Forest**  
**Gradient Boosting**  
**Neural Networks**  
**Stacked Ensembles**

**Classification:**

**Will a customer make a purchase? Yes or No**

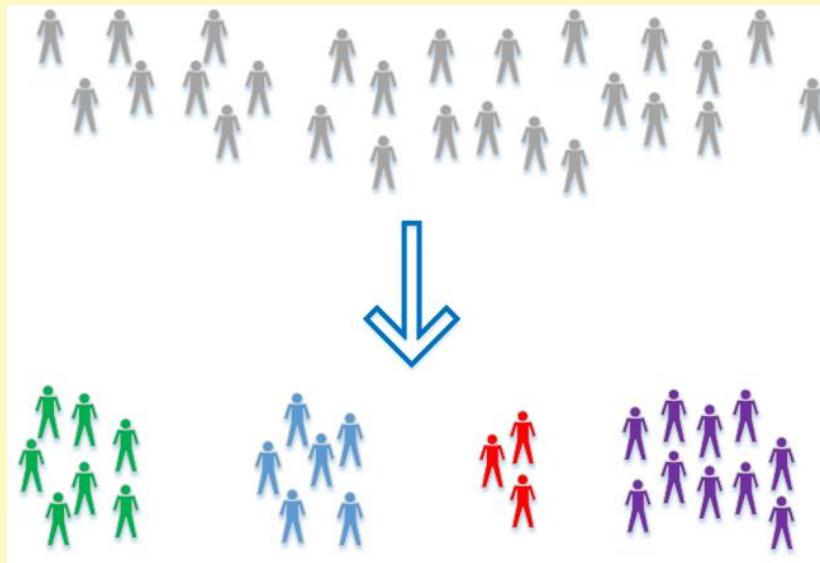


**H<sub>2</sub>O algos:**

**Penalized Linear Models**  
**Naïve Bayes**  
**Random Forest**  
**Gradient Boosting**  
**Neural Networks**  
**Stacked Ensembles**

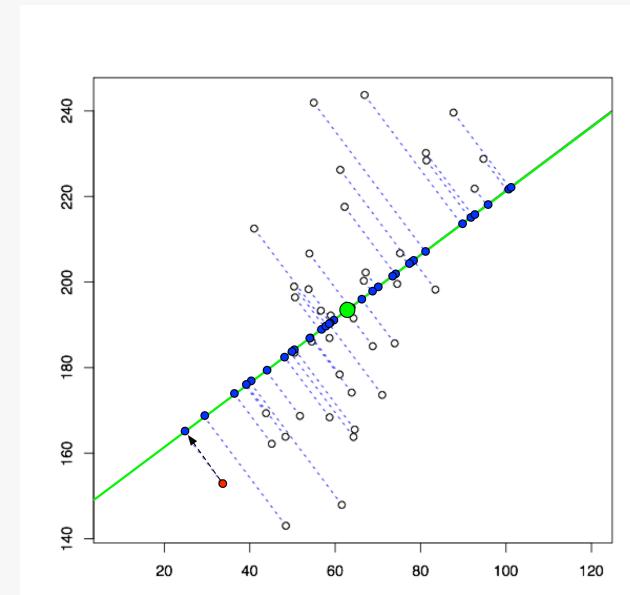
# Unsupervised Learning – Discover Hidden Patterns

**Clustering:**  
**Customer Segmentation**



**H<sub>2</sub>O algos:**  
**K-Means**  
**Generalised Low Rank Model**

**Dimensionality Reduction:**  
**Linear Transformation of Variables**



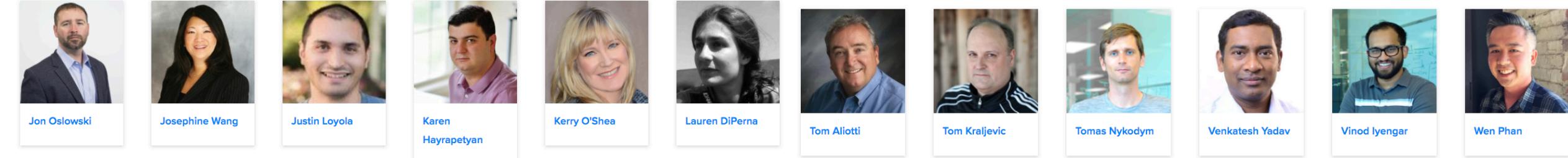
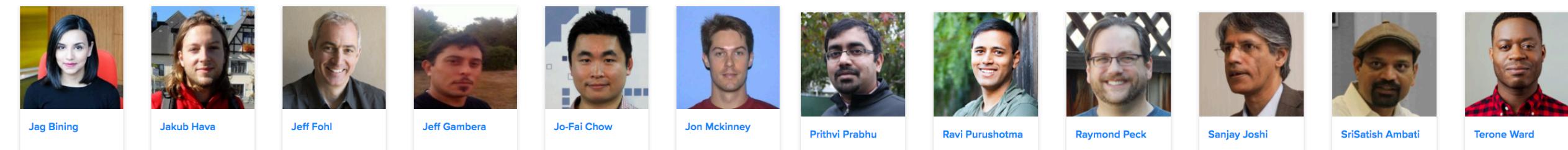
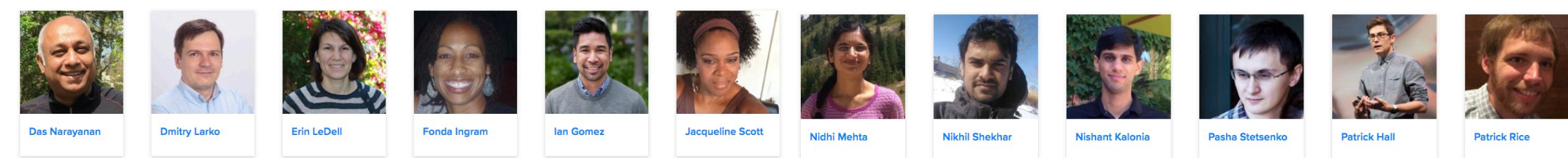
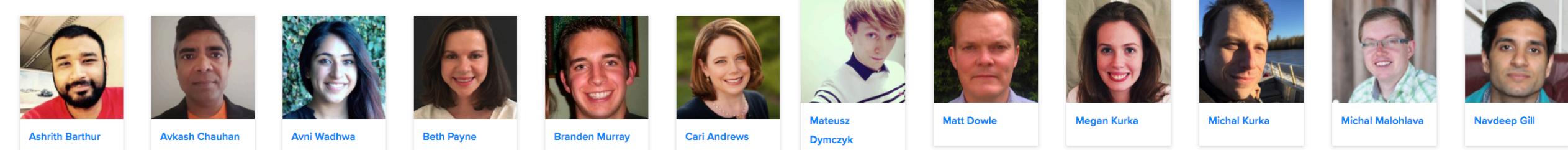
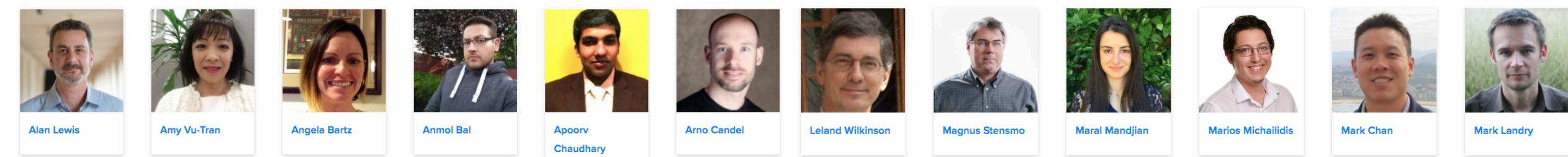
**H<sub>2</sub>O algos:**  
**Principal Component Analysis**  
**Generalised Low Rank Model**

# About H<sub>2</sub>O.ai

# Company Overview

Founded	2011 Venture-backed, debuted in 2012
Products	<ul style="list-style-type: none"><li>• <b>H<sub>2</sub>O Open Source In-Memory AI Prediction Engine</b></li><li>• Sparkling Water (H<sub>2</sub>O + Spark)</li><li>• Deep Water (H<sub>2</sub>O + Other Deep Learning Frameworks)</li><li>• Driverless AI (Next-Gen H<sub>2</sub>O)</li></ul>
Mission	Operationalize Data Science, and provide a platform for users to build beautiful data products
Team	<p>75+ employees</p> <ul style="list-style-type: none"><li>• Distributed Systems Engineers doing Machine Learning</li><li>• World-class visualization designers</li></ul>
Headquarters	Mountain View, CA





# 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*



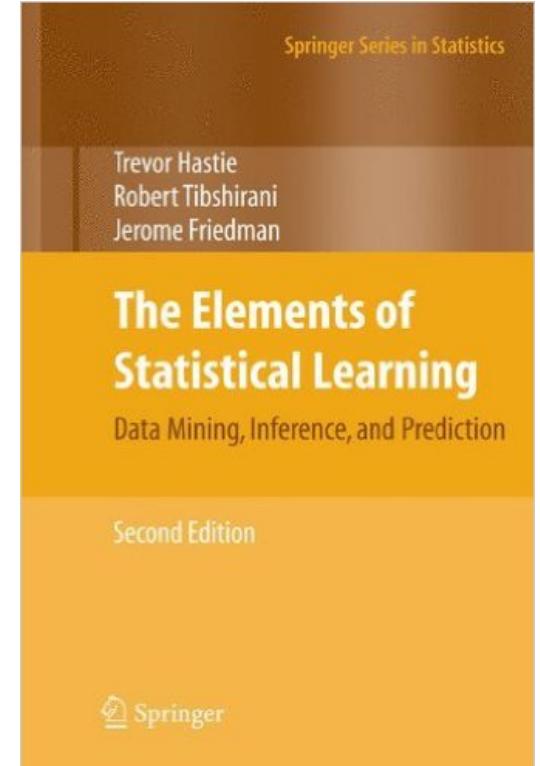
## Dr. Robert Tibshirani

- 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*

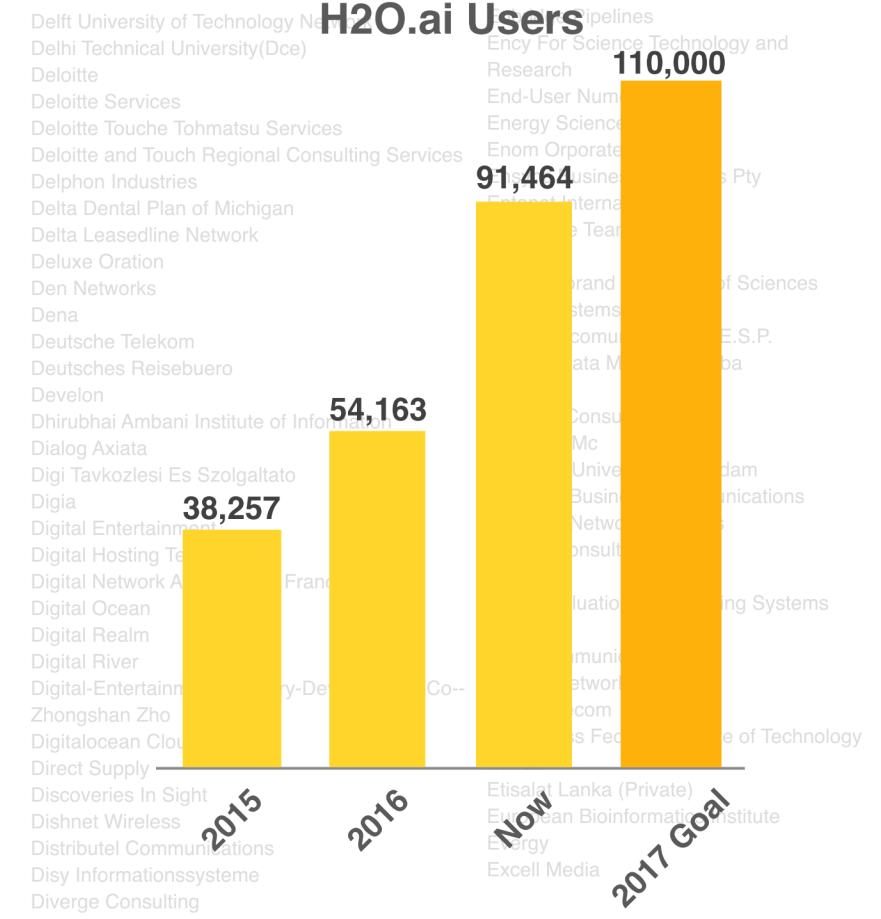
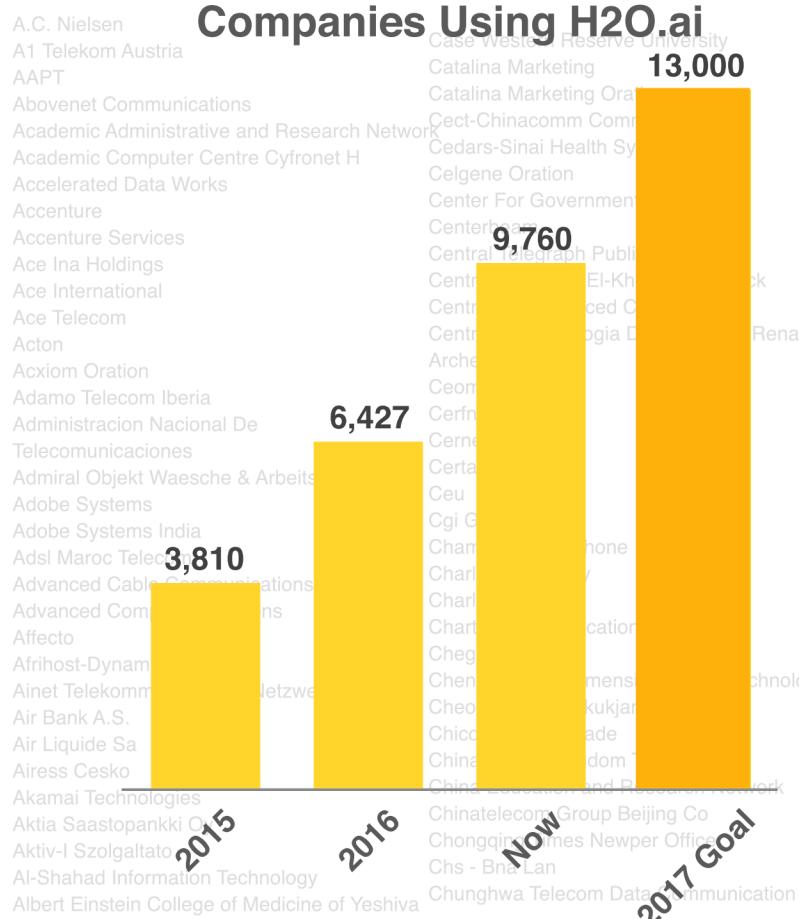


## Dr. Steven Boyd

- 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*



# H2O Community & Fortune 100 customers



## Select Reference Customers:

**"Overall customer satisfaction is very high." - Gartner**



## *Harnessing the power of AI to transform the detection of fraud and error*

### *Setting the scene*

PwC has invested significantly in pioneering the use of artificial intelligence for the audit and has partnered with H2O.ai, a leading Silicon Valley-based AI company.

Following 18 months of development, the first outcome of this partnership is PwC's GL.ai, the first module of PwC's Audit.ai - a revolutionary bot that does what humans can't. Its AI analyses billions of different data points in seconds and applies judgement to detect anomalies in general ledger transactions.



*“The reason this is such a brilliant tool is the ability to look at different risks in context at the same time. For example, it would be uneconomical for an auditor to look at every single user’s pattern of activity and decide what was unusual. With GL.ai, the algorithms do it for us.”*

Laura Needham partner, PwC UK

<http://www.pwc.com/gx/en/about/stories-from-across-the-world/harnessing-the-power-of-ai-to-transform-the-detection-of-fraud-and-error.html>

# Community Expansion



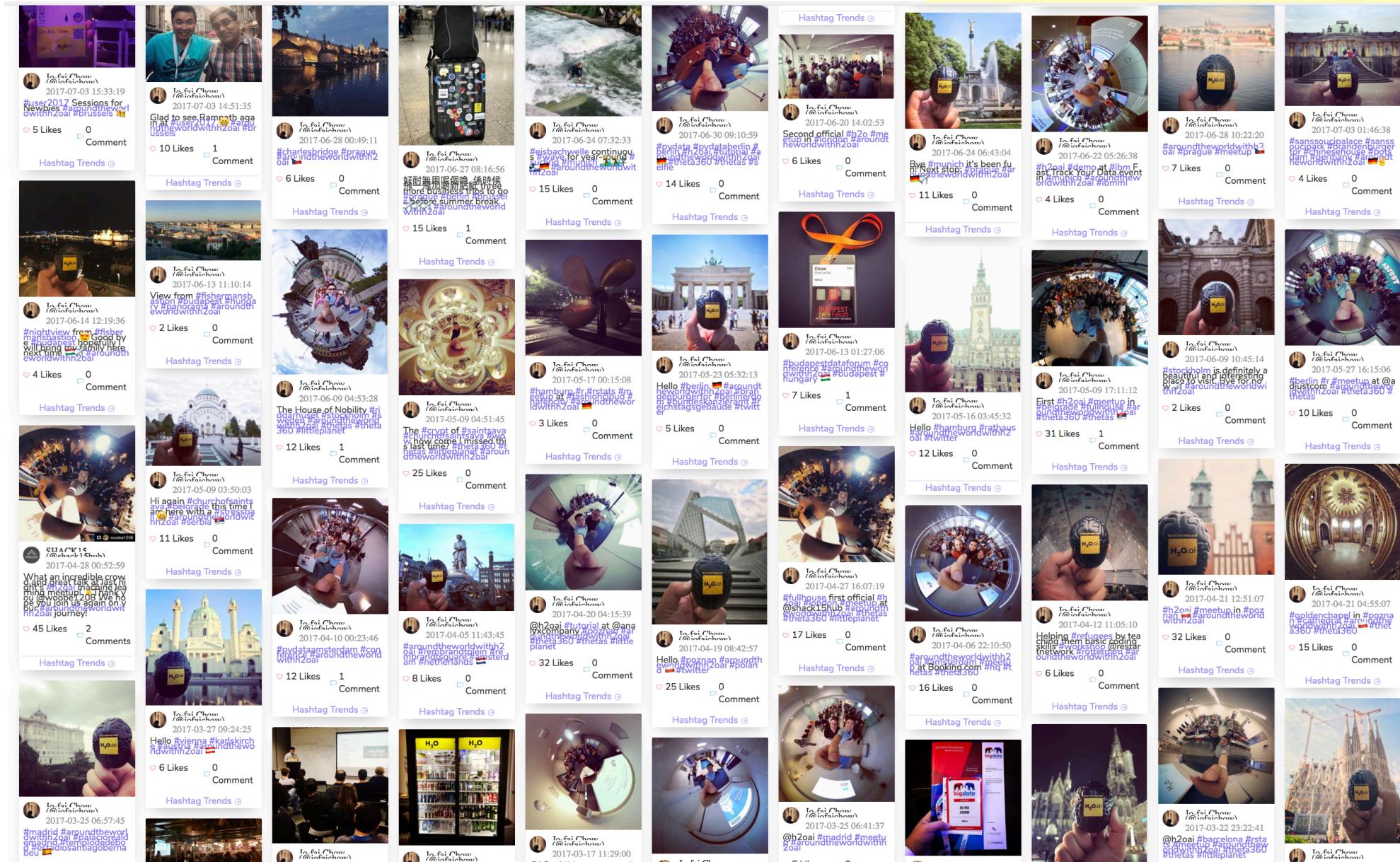
H<sub>2</sub>O.ai

56,536 members    32 interested    50 Meetups    44 cities    18 countries

15 Meetups a month

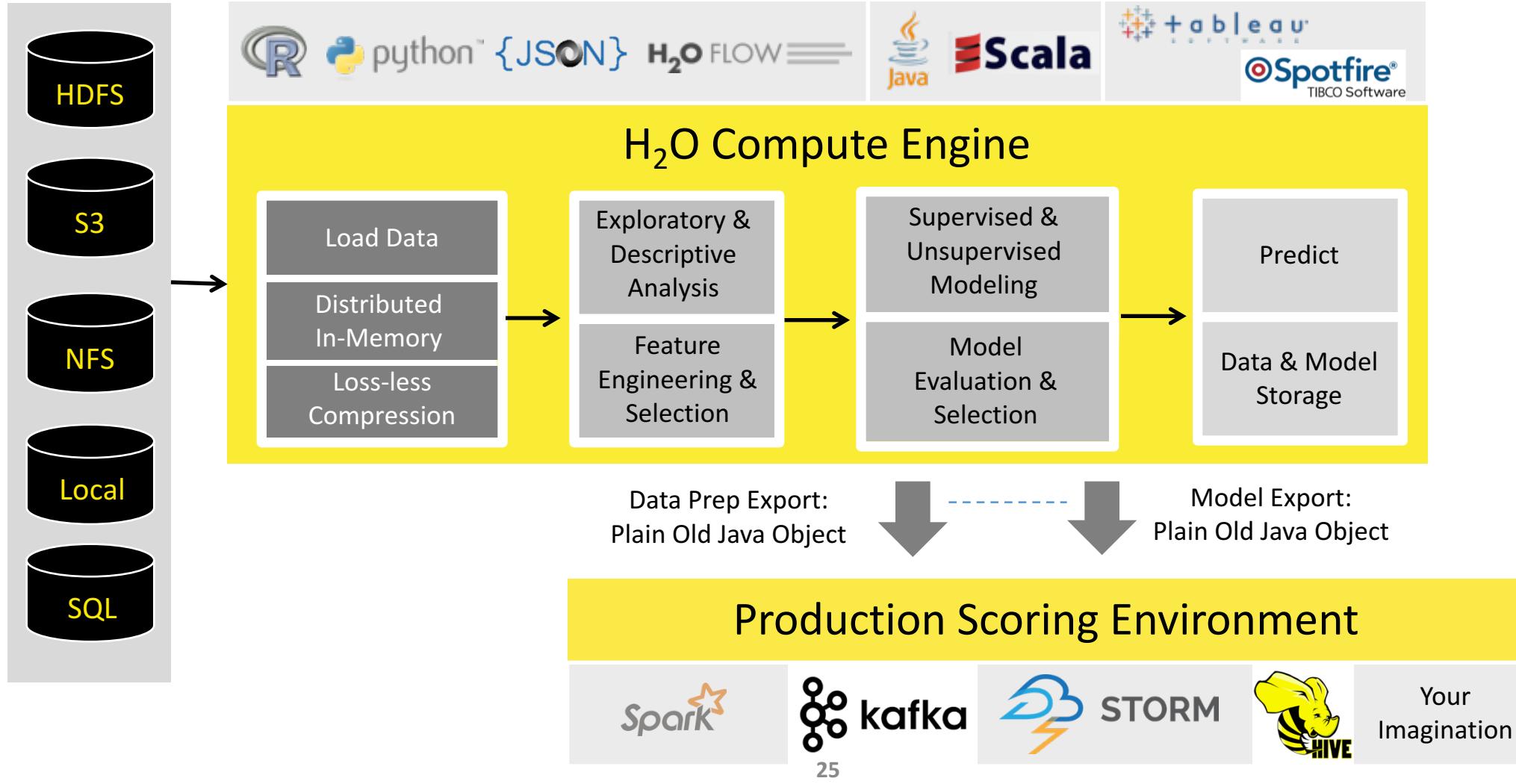
Oct 17 – Amsterdam  
Nov 17 – London

# #AroundTheWorldWithH2Oai



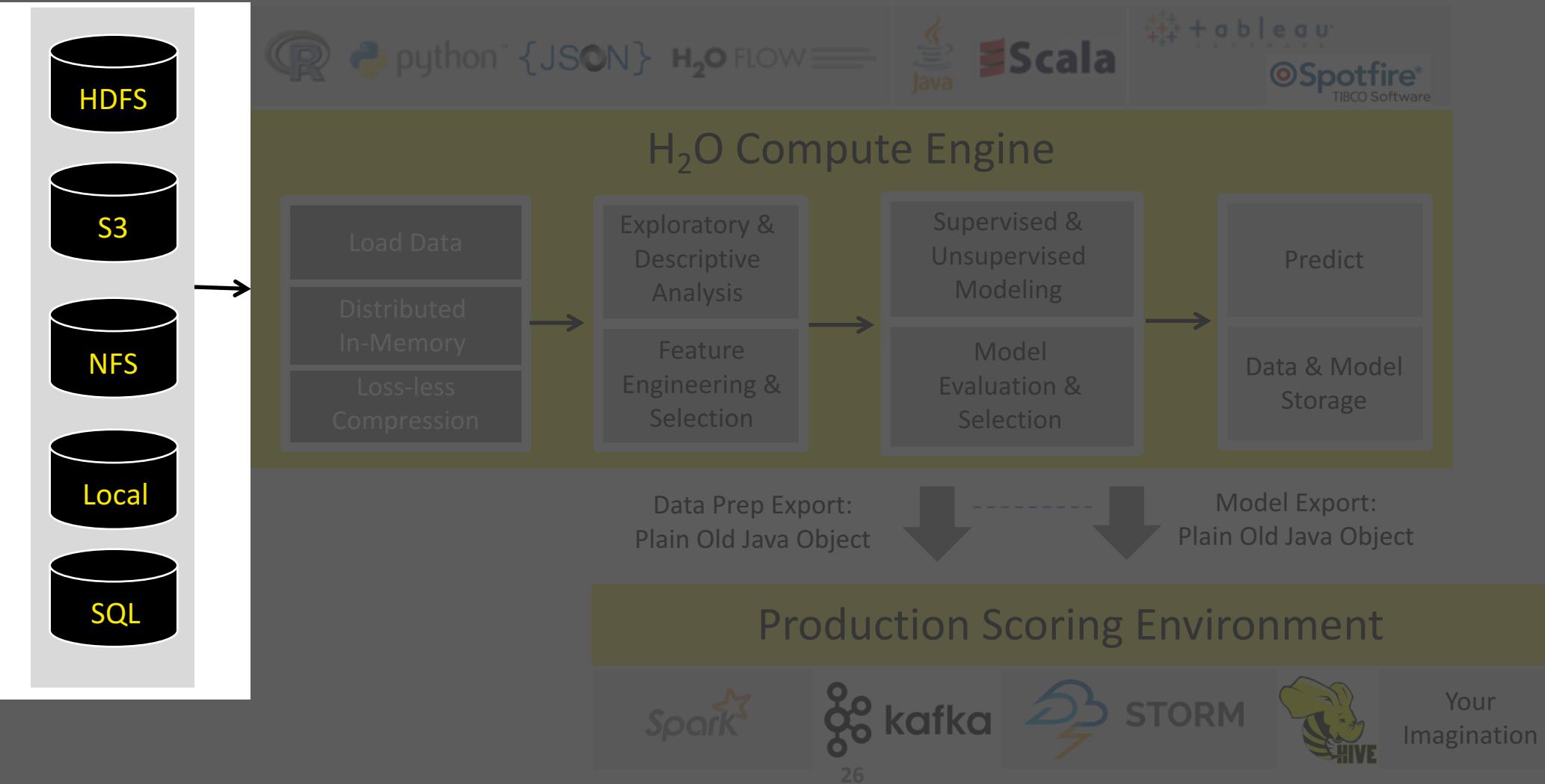
# H<sub>2</sub>O Machine Learning Platform

# High Level Architecture



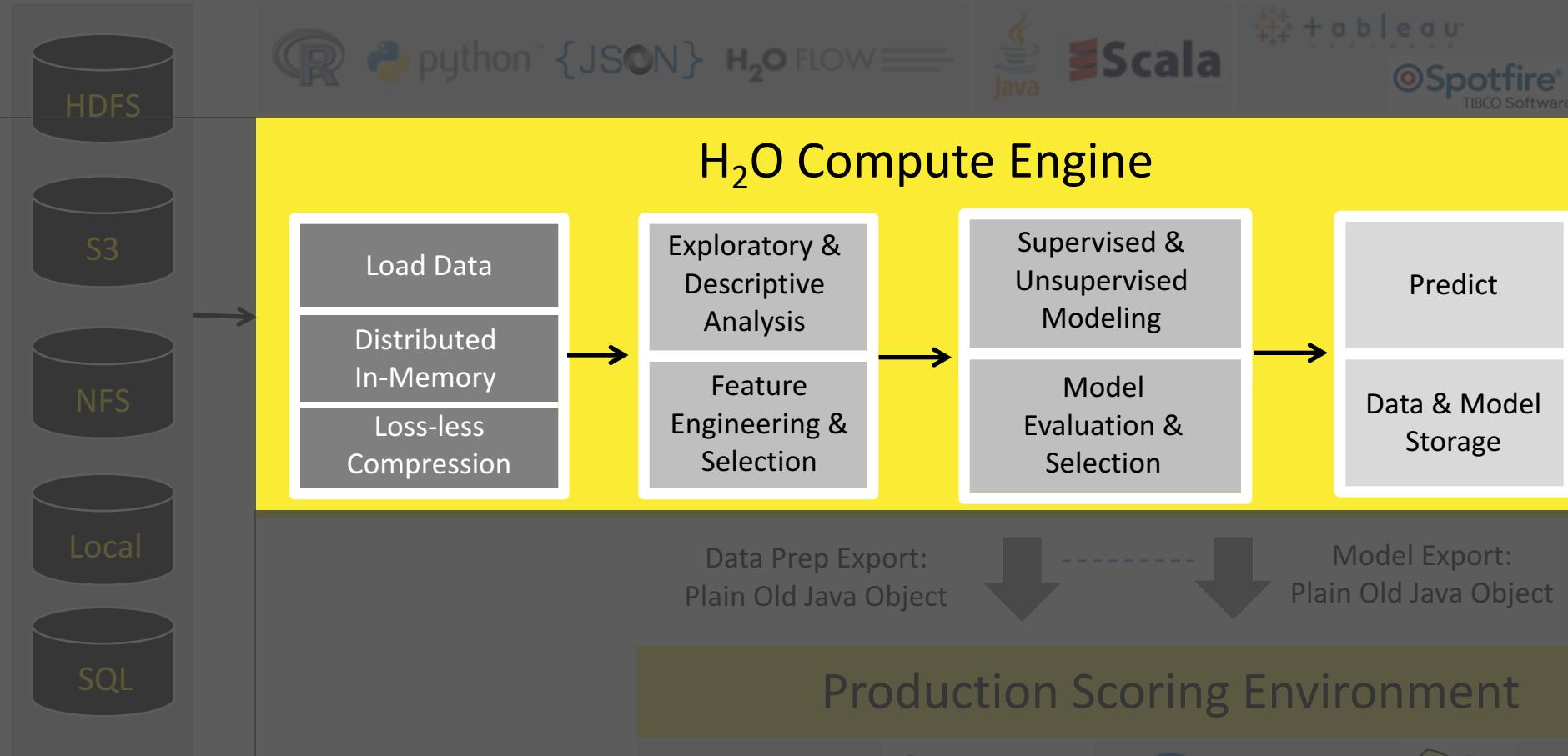
# High Level Architecture

Import Data from  
Multiple Sources



# High Level Architecture

Fast, Scalable & Distributed  
Compute Engine Written in  
Java



# 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

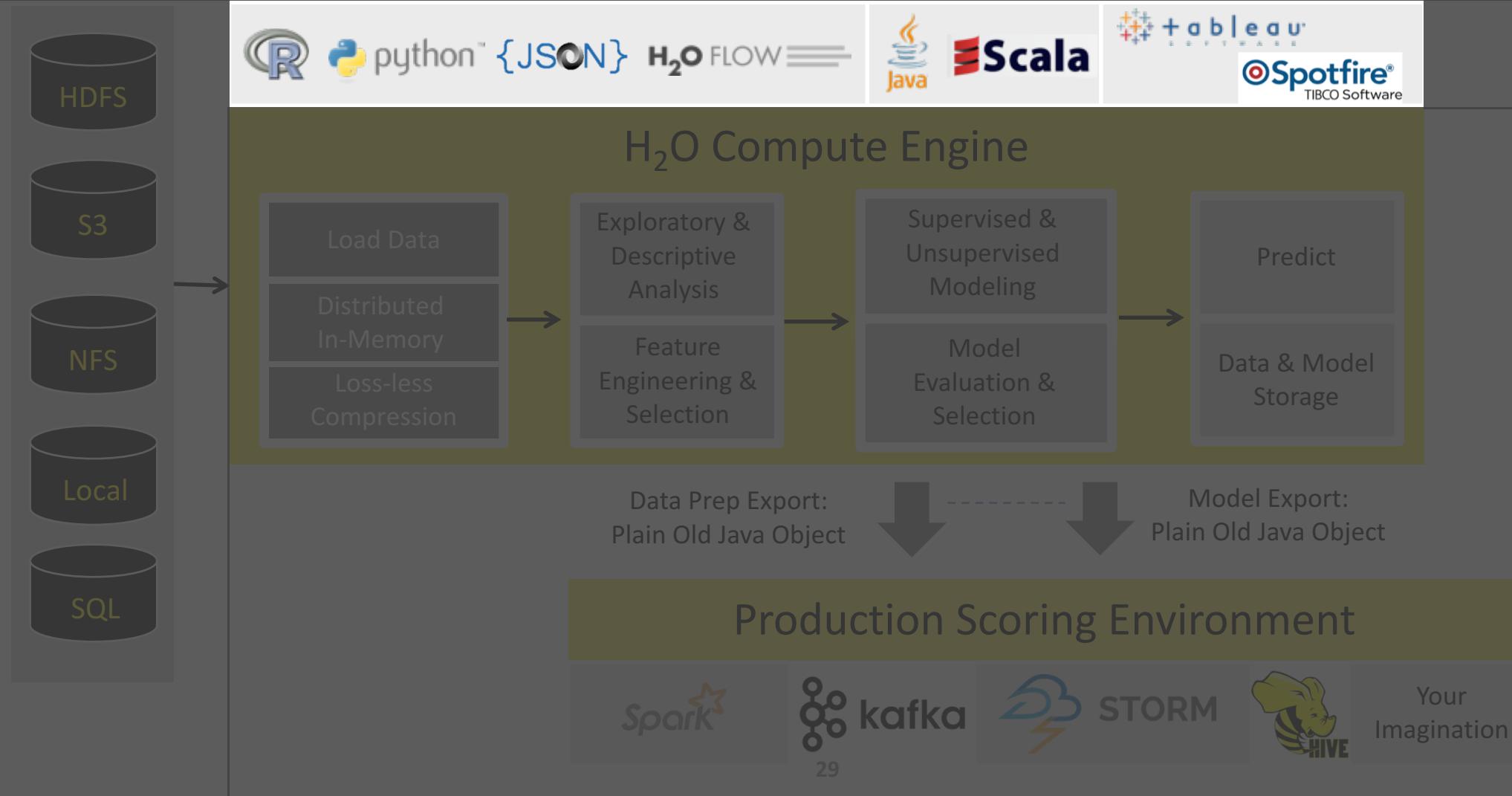
### Dimensionality Reduction

- **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

# High Level Architecture



# H<sub>2</sub>O Flow (Web)

The screenshot shows the H2O Flow (Web) interface running in a web browser. The title bar reads "H2O Flow" and the address bar shows "localhost:54321/flow/index.html". The top navigation bar includes "Flow", "Cell", "Data", "Model" (which is highlighted in yellow), "Score", "Admin", and "Help". A toolbar below the navigation bar contains various icons for file operations like opening, saving, and deleting.

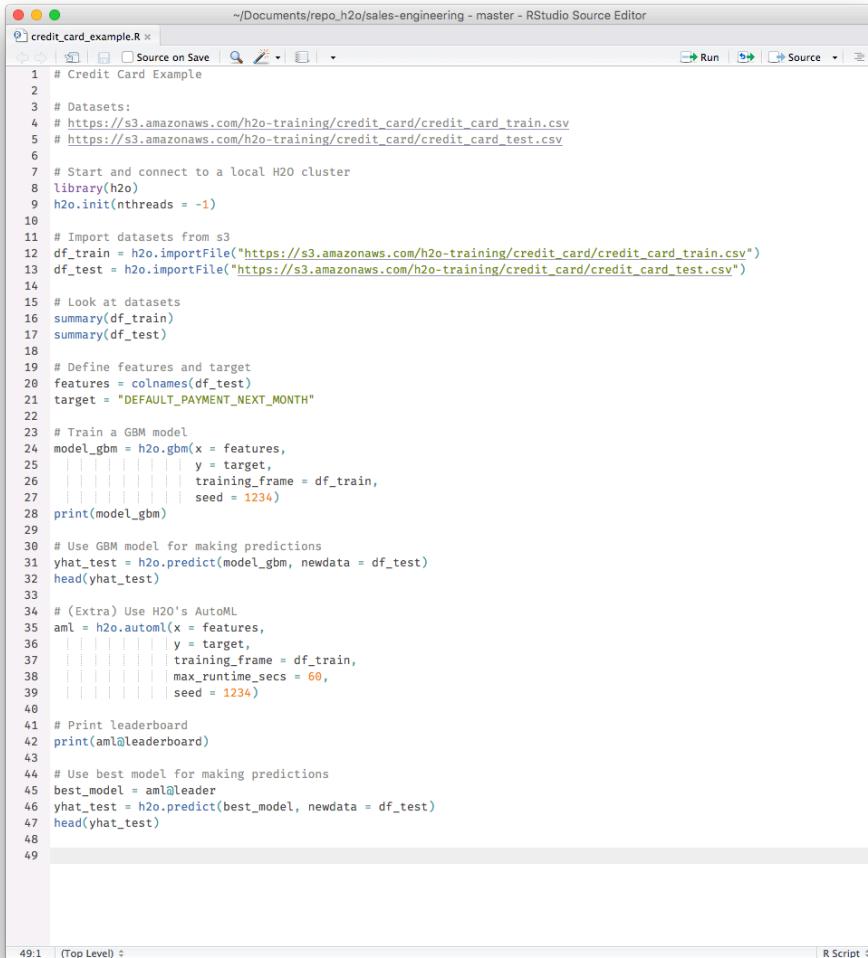
The main workspace is titled "Untitled Flow" and contains a single step labeled "assist". To the left of the workspace is a sidebar titled "Assistance" which lists various H2O routines with their descriptions:

Routine	Description
<code>importFiles</code>	Import file(s) into H2O
<code>getFrames</code>	Get a list of frames in H2O
<code>splitFrame</code>	Split a frame into two or more
<code>mergeFrames</code>	Merge two frames into one
<code>getModels</code>	Get a list of models in H2O
<code>getGrids</code>	Get a list of grid search results
<code>getPredictions</code>	Get a list of predictions in H2O
<code>getJobs</code>	Get a list of jobs running in H2O
<code>buildModel</code>	Build a model
<code>runAutoML</code>	Automatically train and tune
<code>importModel</code>	Import a saved model
<code>predict</code>	Make a prediction

A context menu is open over the "assist" step, showing options such as Aggregator..., Deep Learning..., Distributed Random Forest..., Gradient Boosting Machine..., Generalized Linear Modeling..., Generalized Low Rank Modeling..., K-means..., Naive Bayes..., Principal Components Analysis..., Stacked Ensemble..., Word2Vec..., and XGBoost... .

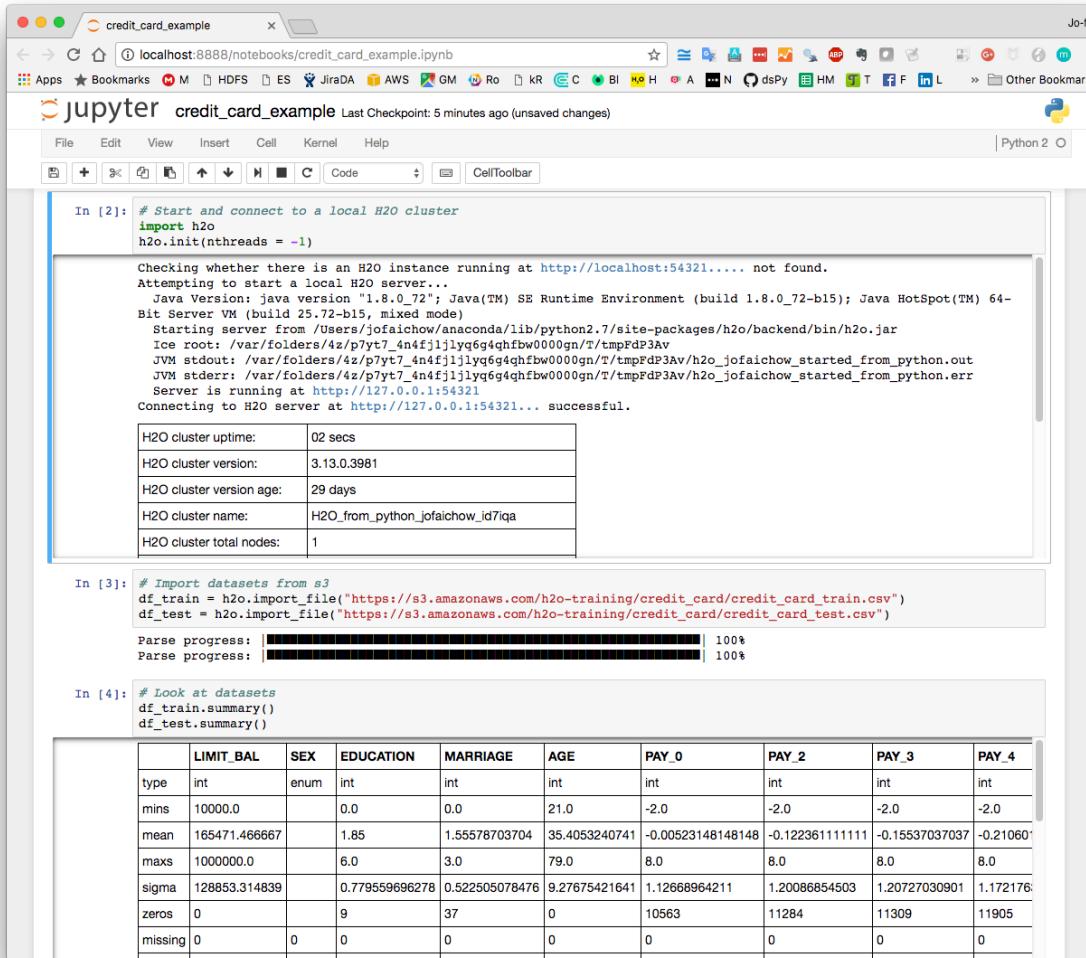
The right side of the interface features a "HELP" panel with sections for "Using Flow for the first time?", "Quickstart Videos", "Or, view example Flows to explore and learn H2O.", "STAR H2O ON GITHUB!", "GENERAL" (with links to Flow Web UI, Importing Data, Building Models, Making Predictions, Using Flows, and Troubleshooting Flow), and "EXAMPLES" (describing Flow packs and providing a link to Browse installed packs...). The bottom right corner shows "Connections: 0" and the H2O logo.

# Using H<sub>2</sub>O with R and Python



The screenshot shows the RStudio Source Editor window with the file `credit_card_example.R` open. The code implements a machine learning pipeline using the H2O library in R. It starts by importing datasets from S3, then connects to a local H2O cluster. The code defines features and target variables, trains a GBM model, and prints the results. It also uses AutoML to find the best model and prints the leaderboard.

```
~/Documents/repo_h2o/sales-engineering - master - RStudio Source Editor
credit_card_example.R
Source on Save Run Source Cell Toolbar
1 # Credit Card Example
2
3 # Datasets:
4 # https://s3.amazonaws.com/h2o-training/credit_card/credit_card_train.csv
5 # https://s3.amazonaws.com/h2o-training/credit_card/credit_card_test.csv
6
7 # Start and connect to a local H2O cluster
8 library(h2o)
9 h2o.init(nthreads = -1)
10
11 # Import datasets from s3
12 df_train = h2o.importFile("https://s3.amazonaws.com/h2o-training/credit_card/credit_card_train.csv")
13 df_test = h2o.importFile("https://s3.amazonaws.com/h2o-training/credit_card/credit_card_test.csv")
14
15 # Look at datasets
16 summary(df_train)
17 summary(df_test)
18
19 # Define features and target
20 features = colnames(df_test)
21 target = "DEFAULT_PAYMENT_NEXT_MONTH"
22
23 # Train a GBM model
24 model_gbm = h2o.gbm(x = features,
25                      y = target,
26                      training_frame = df_train,
27                      seed = 1234)
28 print(model_gbm)
29
30 # Use GBM model for making predictions
31 yhat_test = h2o.predict(model_gbm, newdata = df_test)
32 head(yhat_test)
33
34 # (Extra) Use H2O's AutoML
35 aml = h2o.automl(x = features,
36                   y = target,
37                   training_frame = df_train,
38                   max_runtime_secs = 60,
39                   seed = 1234)
40
41 # Print leaderboard
42 print(aml@leaderboard)
43
44 # Use best model for making predictions
45 best_model = aml@leader
46 yhat_test = h2o.predict(best_model, newdata = df_test)
47 head(yhat_test)
48
49
```



The screenshot shows a Jupyter Notebook interface with the notebook `credit_card_example.ipynb` open. The notebook contains Python code for connecting to a local H2O cluster, importing datasets, and summarizing them. The output cell for the connection attempt shows the process of starting the H2O server and connecting successfully. Subsequent cells show the import of datasets and their summaries, including a detailed table of the dataset's numerical and categorical columns and their statistics.

```
In [2]: # Start and connect to a local H2O cluster
import h2o
h2o.init(nthreads = -1)

Checking whether there is an H2O instance running at http://localhost:54321.... not found.
Attempting to start a local H2O server...
Java Version: java version "1.8.0_72"; Java(TM) SE Runtime Environment (build 1.8.0_72-b15); Java HotSpot(TM) 64-Bit Server VM (build 25.72-b15, mixed mode)
Starting server from /Users/jofaichow/anaconda/lib/python2.7/site-packages/h2o/backend/bin/h2o.jar
Ice root: /var/folders/4z/p7yt7_4n4fjijiy6g4qfbw000gn/T/tmpPdP3Av
JVM stdout: /var/folders/4z/p7yt7_4n4fjijiy6g4qfbw000gn/T/tmpPdP3Av/h2o_jofaichow_started_from_python.out
JVM stderr: /var/folders/4z/p7yt7_4n4fjijiy6g4qfbw000gn/T/tmpPdP3Av/h2o_jofaichow_started_from_python.err
Server is running at http://127.0.0.1:54321
Connecting to H2O server at http://127.0.0.1:54321... successful.

H2O cluster uptime: 02 secs
H2O cluster version: 3.13.0.3981
H2O cluster version age: 29 days
H2O cluster name: H2O_from_python_jofaichow_id7qa
H2O cluster total nodes: 1

In [3]: # Import datasets from s3
df_train = h2o.import_file("https://s3.amazonaws.com/h2o-training/credit_card/credit_card_train.csv")
df_test = h2o.import_file("https://s3.amazonaws.com/h2o-training/credit_card/credit_card_test.csv")

Parse progress: |██████████| 100%
Parse progress: |██████████| 100%

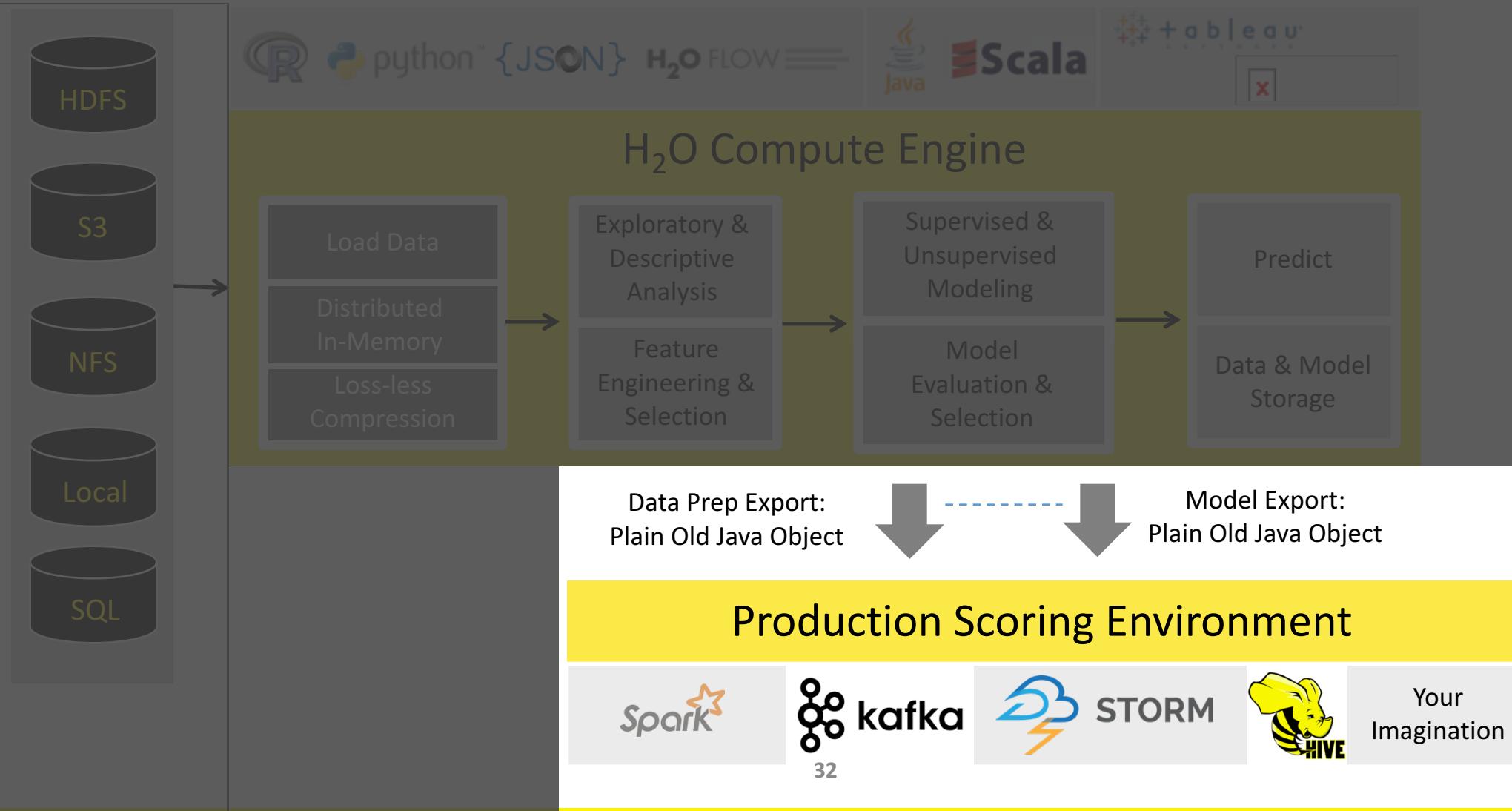
In [4]: # Look at datasets
df_train.summary()
df_test.summary()


```

	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_0	PAY_2	PAY_3	PAY_4
type	int	enum	int	int	int	int	int	int	int
mins	10000.0		0.0	0.0	21.0	-2.0	-2.0	-2.0	-2.0
mean	165471.466667		1.85	1.55578703704	35.4053240741	-0.00523148148148	-0.122361111111	-0.15537037037	-0.210601
maxs	1000000.0		6.0	3.0	79.0	8.0	8.0	8.0	8.0
sigma	128853.314839		0.779559696278	0.522505078476	9.27675421641	1.12668964211	1.20086854503	1.20727030901	1.172176
zeros	0		9	37	0	10563	11284	11309	11905
missing	0		0	0	0	0	0	0	0

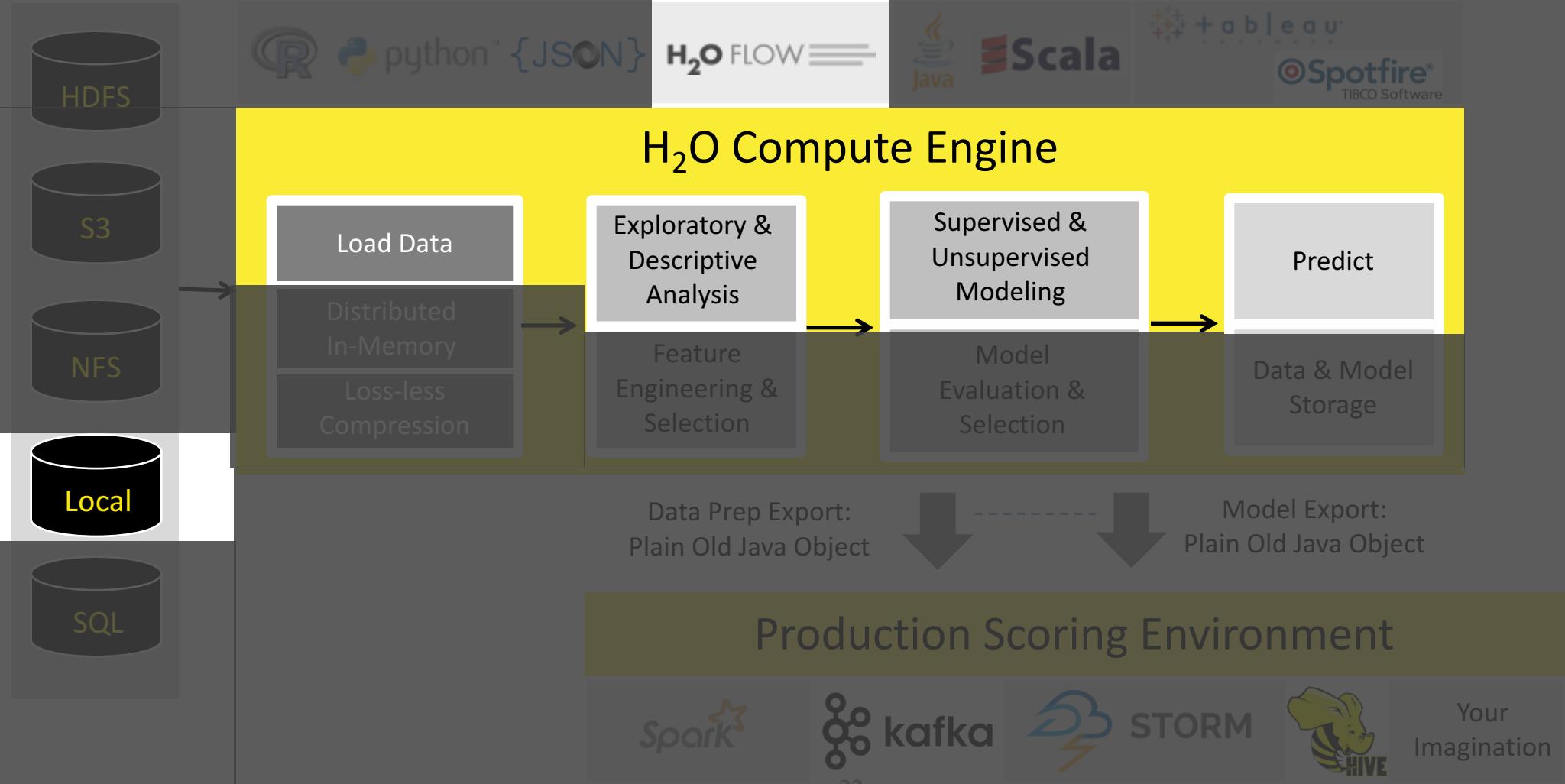
# High Level Architecture

Export Standalone Models  
for Production



# High Level Architecture

This Workshop



## Getting Started & User Guides

### H<sub>2</sub>O

[What is H<sub>2</sub>O?](#)  
[H<sub>2</sub>O User Guide](#) (Main docs)  
[H<sub>2</sub>O Book \(O'Reilly\)](#)  
[Recent Changes](#)  
[Open Source License \(Apache V2\)](#)

[Quick Start Video - Flow Web UI](#)  
[Quick Start Video - R](#)  
[Quick Start Video - Python](#)

[Download H<sub>2</sub>O](#)

### Sparkling Water

[What is Sparkling Water?](#)  
[Sparkling Water Booklet](#)  
[PySparkling Readme](#) 2.0 | 2.1 | 2.2  
[RSparkling Readme](#)  
[Open Source License \(Apache V2\)](#)

[Quick Start Video - Scala](#)

[Download Sparkling Water](#)

### Steam

[What is Steam?](#)  
[Steam User Guide](#)  
[Recent Changes](#)  
[Open Source License \(AGPL\)](#)

[Download Steam](#)

### Deep Water (preview)

[Deep Water Readme](#)  
[Deep Water Booklet](#)  
[Deep Water AMI Guide](#)  
[Deep Water Docker Image](#)  
[Open Source License \(Apache V2\)](#)

[Launch Deep Water AMI  
\(choose p2.xlarge\)](#)

### Q & A

[FAQ](#)  
[Issue Tracking \(JIRA\)](#)  
[Stack Overflow](#)  
[h2ostream Google Group](#)  
[Gitter](#)  
[Cross Validated](#)

**For Supported Enterprise Customers**  
[Enterprise Support Web | Email](#)

## Algorithms

### Supervised Learning

Generalized Linear Modeling (GLM)	<a href="#">Tutorial</a>	<a href="#">Booklet</a>	<a href="#">Reference</a>	<a href="#">Tuning</a>
Gradient Boosting Machine (GBM)	<a href="#">Tutorial</a>	<a href="#">Booklet</a>	<a href="#">Reference</a>	<a href="#">Tuning</a>
Deep Learning	<a href="#">Tutorial</a>	<a href="#">Booklet</a>	<a href="#">Reference</a>	<a href="#">Tuning</a>
Distributed Random Forest	<a href="#">Tutorial</a>	<a href="#">Booklet</a>	<a href="#">Reference</a>	<a href="#">Tuning</a>
Naive Bayes	<a href="#">Tutorial</a>	<a href="#">Booklet</a>	<a href="#">Reference</a>	<a href="#">Tuning</a>
Stacked Ensembles	<a href="#">Tutorial</a>	<a href="#">Booklet</a>	<a href="#">Reference</a>	<a href="#">Tuning</a>
XGBoost	<a href="#">Tutorial</a>	<a href="#">Booklet</a>	<a href="#">Reference</a>	<a href="#">Tuning</a>

### Unsupervised Learning

Generalized Low Rank Models (GLRM)	<a href="#">Tutorial</a>	<a href="#">Reference</a>
K-Means Clustering	<a href="#">Tutorial</a>	<a href="#">Reference</a>
Principal Components Analysis (PCA)	<a href="#">Tutorial</a>	<a href="#">Reference</a>

### Miscellaneous

Word2vec	<a href="#">Tutorial</a>	<a href="#">Reference</a>
----------	--------------------------	---------------------------

TOP

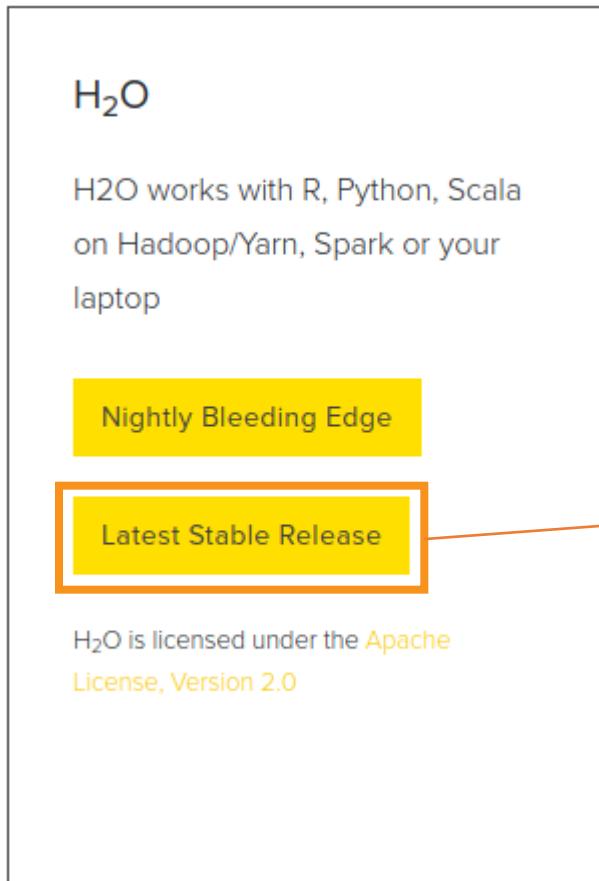
# H<sub>2</sub>O Tutorials

- Quick Start
- Regression
- Classification
- Clustering

# H<sub>2</sub>O Quick Start

# www.h2o.ai/download

Prerequisite: Java version 7 or 8  
(Note: Java version 9 is not yet supported)



**H<sub>2</sub>O**  
Version 3.14.0.3

Fast Scalable Machine Learning API  
For Smarter Applications

DOWNLOAD AND RUN    INSTALL IN R    INSTALL IN PYTHON    INSTALL ON HADOOP    USE FROM MAVEN

**DOWNLOAD H<sub>2</sub>O**

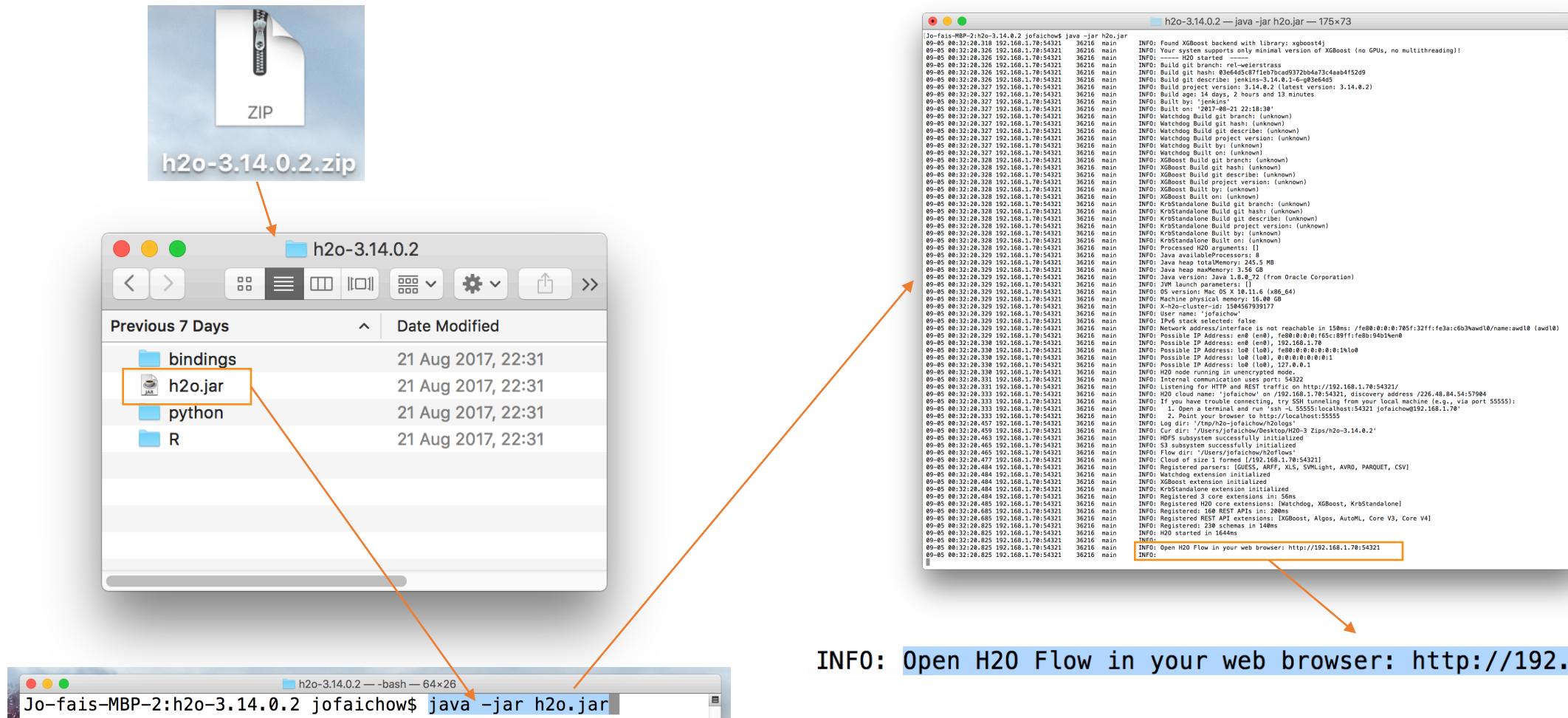
Get started with H<sub>2</sub>O in 3 easy steps

1. Download H<sub>2</sub>O. This is a zip file that contains everything you need to get started.
2. From your terminal, run:

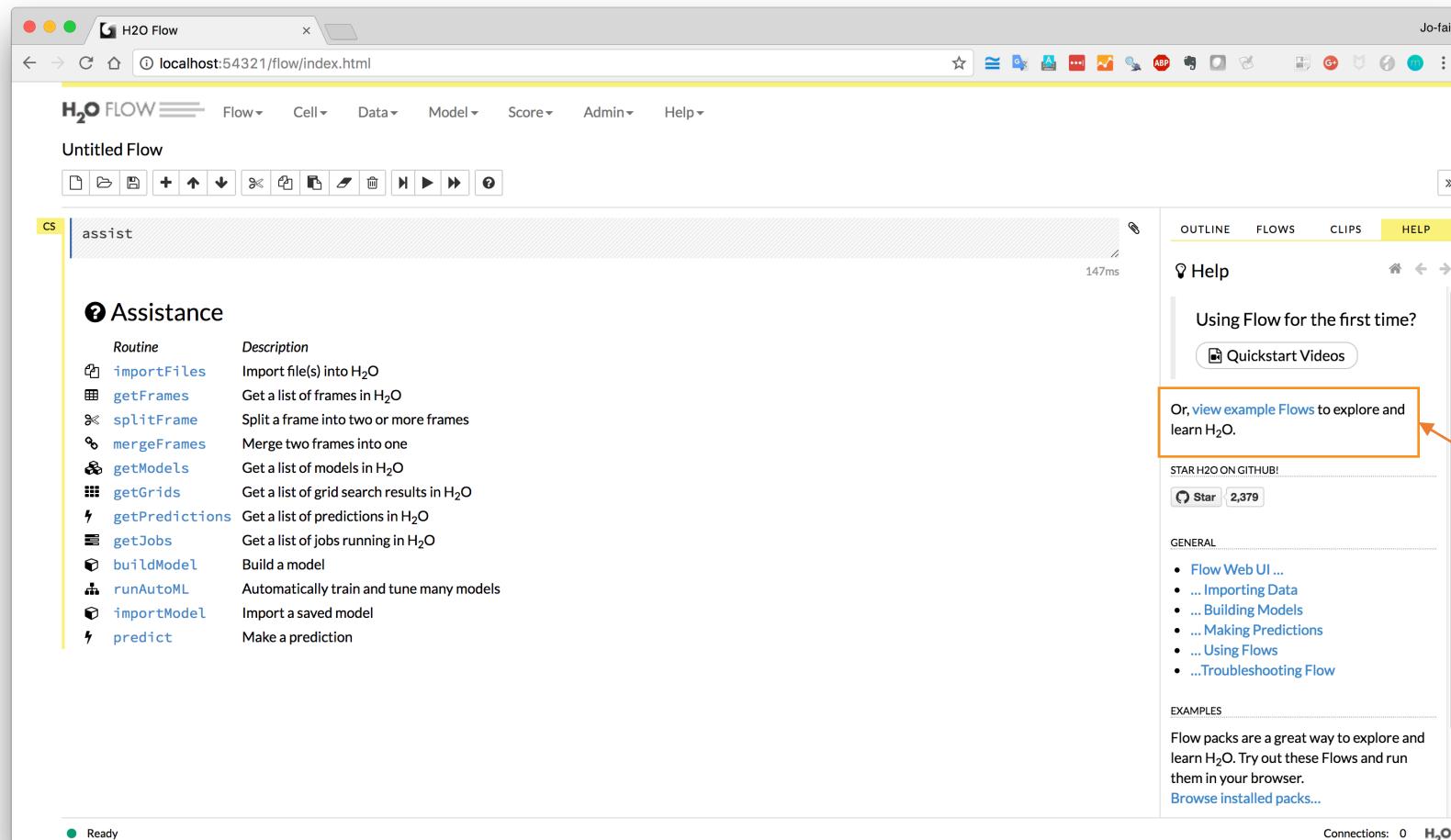
```
cd ~/Downloads  
unzip h2o-3.14.0.3.zip  
cd h2o-3.14.0.3  
java -jar h2o.jar
```

3. Point your browser to <http://localhost:54321>

# Install and Start H<sub>2</sub>O Flow (Web Interface)

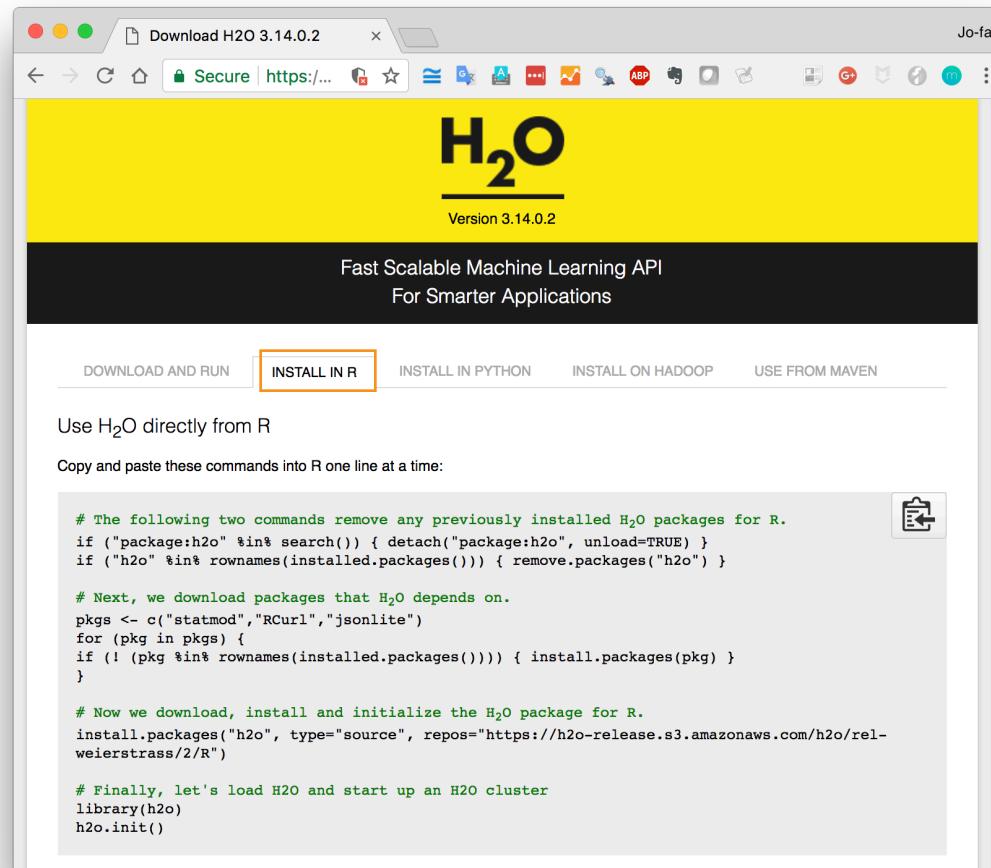


# H<sub>2</sub>O Flow (Web Interface)



More Examples

# Install and Start H<sub>2</sub>O in R / Python



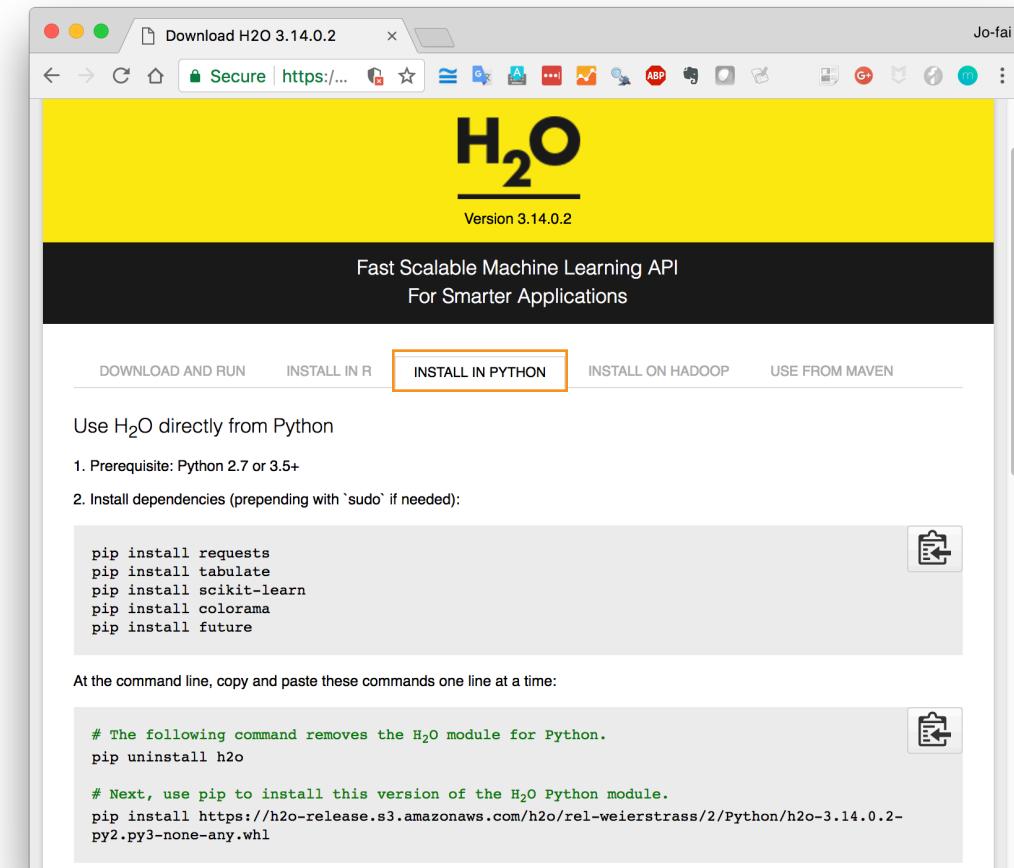
The screenshot shows the H2O download page for version 3.14.0.2. The top navigation bar includes links for 'Download H2O 3.14.0.2', 'Secure | https://...', and 'Jo-fai'. Below the header, the H2O logo and version information are displayed. A yellow banner reads 'Fast Scalable Machine Learning API For Smarter Applications'. Below the banner, there are four buttons: 'DOWNLOAD AND RUN' (disabled), 'INSTALL IN R' (highlighted in orange), 'INSTALL IN PYTHON', 'INSTALL ON HADOOP', and 'USE FROM MAVEN'. A section titled 'Use H2O directly from R' contains R code for package installation and cluster initialization. A clipboard icon is located next to the code block.

```
# The following two commands remove any previously installed H2O packages for R.
if ("package:h2o" %in% search()) { detach("package:h2o", unload=TRUE) }
if ("h2o" %in% rownames(installed.packages())) { remove.packages("h2o") }

# Next, we download packages that H2O depends on.
pkgs <- c("statmod", "RCurl", "jsonlite")
for (pkg in pkgs) {
  if (! (pkg %in% rownames(installed.packages()))) { install.packages(pkg) }
}

# Now we download, install and initialize the H2O package for R.
install.packages("h2o", type="source", repos="https://h2o-release.s3.amazonaws.com/h2o/rel-weierstrass/2/R")

# Finally, let's load H2O and start up an H2O cluster
library(h2o)
h2o.init()
```



The screenshot shows the H2O download page for version 3.14.0.2. The top navigation bar includes links for 'Download H2O 3.14.0.2', 'Secure | https://...', and 'Jo-fai'. Below the header, the H2O logo and version information are displayed. A yellow banner reads 'Fast Scalable Machine Learning API For Smarter Applications'. Below the banner, there are five buttons: 'DOWNLOAD AND RUN' (disabled), 'INSTALL IN R' (disabled), 'INSTALL IN PYTHON' (highlighted in orange), 'INSTALL ON HADOOP', and 'USE FROM MAVEN'. A section titled 'Use H2O directly from Python' contains instructions and command-line code. It lists prerequisites (Python 2.7 or 3.5) and dependencies (requests, tabulate, scikit-learn, colorama, future). A clipboard icon is located next to the dependency list. Another section provides Python installation commands at the command line.

1. Prerequisite: Python 2.7 or 3.5+
2. Install dependencies (prepending with `sudo` if needed):

```
pip install requests
pip install tabulate
pip install scikit-learn
pip install colorama
pip install future
```

At the command line, copy and paste these commands one line at a time:

```
# The following command removes the H2O module for Python.
pip uninstall h2o

# Next, use pip to install this version of the H2O Python module.
pip install https://h2o-release.s3.amazonaws.com/h2o/rel-weierstrass/2/Python/h2o-3.14.0.2-py2.py3-none-any.whl
```

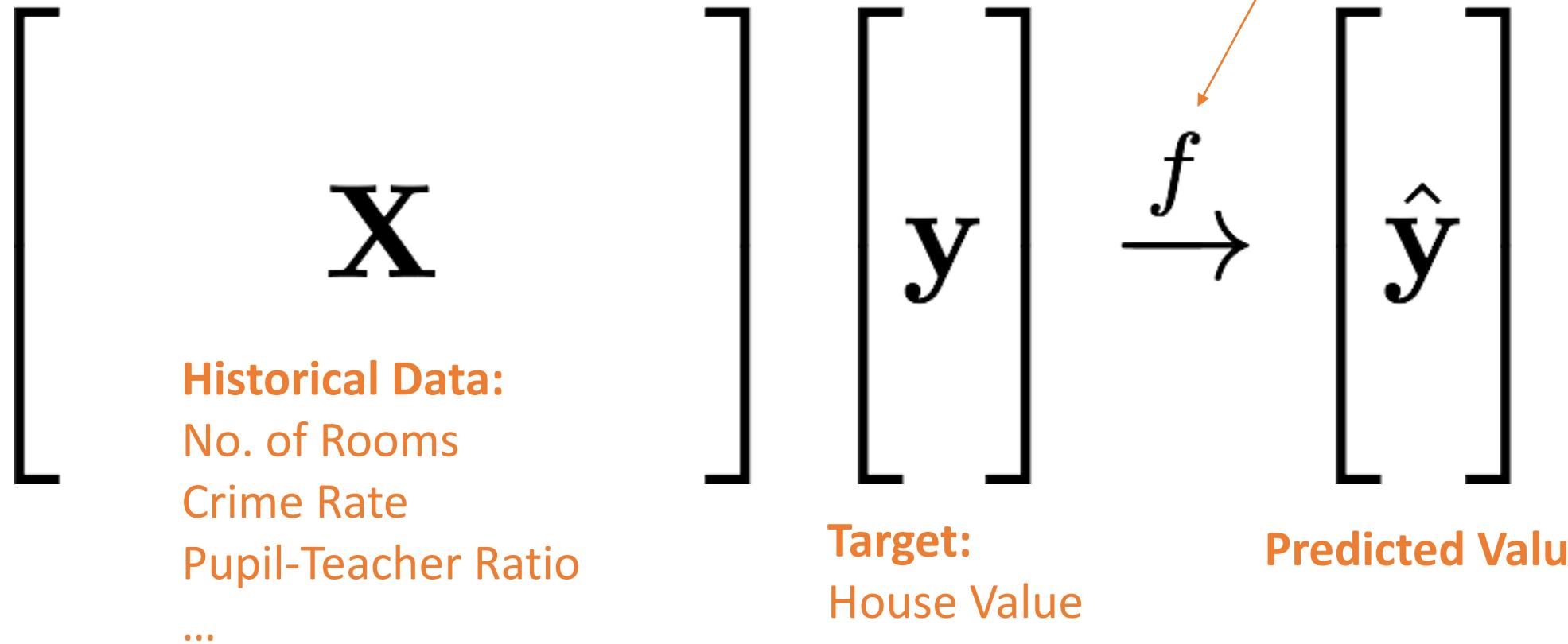
# Tutorial 1 - Regression

- **Data:** Boston Housing (1978)
- **Source:** <http://www.cs.toronto.edu/~delve/data/boston/bostonDetail.html>



# Supervised Learning Example

Machine Learning:  
Learn Patterns  
from Data



1. CRIM - per capita crime rate by town
2. ZN - proportion of residential land zoned for lots over 25,000 sq.ft.
3. INDUS - proportion of non-retail business acres per town.
4. CHAS - Charles River dummy variable (1 if tract bounds river; 0 otherwise)
5. NOX - nitric oxides concentration (parts per 10 million)
6. RM - average number of rooms per dwelling
7. AGE - proportion of owner-occupied units built prior to 1940
8. DIS - weighted distances to five Boston employment centres
9. RAD - index of accessibility to radial highways
10. TAX - full-value property-tax rate per \$10,000
11. PTRATIO - pupil-teacher ratio by town
12. B -  $1000(Bk - 0.63)^2$  where Bk is the proportion of blacks by town
13. LSTAT - % lower status of the population
14. MEDV - Median value of owner-occupied homes in \$1000's

house\_price\_train

Search Sheet

Home Insert Page Layout Formulas Data Review View

Cut Copy Paste Format

Calibri (Body) 12 A A Wrap Text General

Merge & Center Conditional Formatting Format as Table Cell Styles Insert Delete Format

AutoSum Fill Clear Sort & Filter

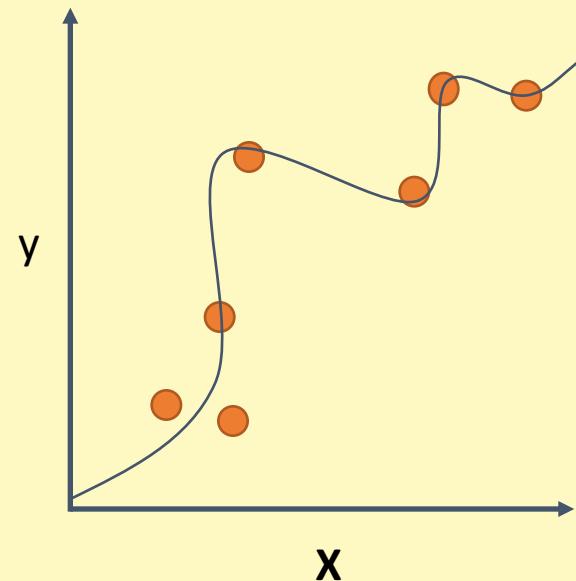
N1 fx medv

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1	crim	zn	indus	chas	nox	rm	age	dis	rad	tax	ptratio	b	lstat	medv	
2	0.02731	0	7.07	0	0.469	6.421	78.9	4.9671	2	242	17.8	396.9	9.14	21.6	
3	0.02729	0	7.07	0	0.469	7.185	61.1	4.9671	2	242	17.8	392.83	4.03	34.7	
4	0.03237	0	2.18	0	0.458	6.998	45.8	6.0622	3	222	18.7	394.63	2.94	33.4	
5	0.06905	0	2.18	0	0.458	7.147	54.2	6.0622	3	222	18.7	396.9	5.33	36.2	
6	0.08829	12.5	7.87	0	0.524	6.012	66.6	5.5605	5	311	15.2	395.6	12.43	22.9	
7	0.14455	12.5	7.87	0	0.524	6.172	96.1	5.9505	5	311	15.2	396.9	19.15	27.1	
8	0.21124	12.5	7.87	0	0.524	5.631	100	6.0821	5	311	15.2	386.63	29.93	16.5	
9	0.17004	12.5	7.87	0	0.524	6.004	85.9	6.5921	5	311	15.2	386.71	17.1	18.9	
10	0.22489	12.5	7.87	0	0.524	6.377	94.3	6.3467	5	311	15.2	392.52	20.45	15	
11	0.11747	12.5	7.87	0	0.524	6.009	82.9	6.2267	5	311	15.2	396.9	13.27	18.9	
12	0.09378	12.5	7.87	0	0.524	5.889	39	5.4509	5	311	15.2	390.5	15.71	21.7	
13	0.62976	0	8.14	0	0.538	5.949	61.8	4.7075	4	307	21	396.9	8.26	20.4	
14	0.63796	0	8.14	0	0.538	6.096	84.5	4.4619	4	307	21	380.02	10.26	18.2	
15	0.62739	0	8.14	0	0.538	5.834	56.5	4.4986	4	307	21	395.62	8.47	19.9	
16	1.05393	0	8.14	0	0.538	5.935	29.3	4.4986	4	307	21	386.85	6.58	23.1	
17	0.7842	0	8.14	0	0.538	5.99	81.7	4.2579	4	307	21	386.75	14.67	17.5	
18	0.80271	0	8.14	0	0.538	5.456	36.6	3.7965	4	307	21	288.99	11.69	20.2	
19	0.7258	0	8.14	0	0.538	5.727	69.5	3.7965	4	307	21	390.95	11.28	18.2	
20	1.25179	0	8.14	0	0.538	5.57	98.1	3.7979	4	307	21	376.57	21.02	13.6	
21	0.85204	0	8.14	0	0.538	5.965	89.2	4.0123	4	307	21	392.53	13.83	19.6	
22	0.75026	0	8.14	0	0.538	5.924	94.1	4.3996	4	307	21	394.33	16.3	15.6	
23	0.84054	0	8.14	0	0.538	5.599	85.7	4.4546	4	307	21	303.42	16.51	13.9	
24	0.67191	0	8.14	0	0.538	5.813	90.3	4.682	4	307	21	376.88	14.81	16.6	
25	0.95577	0	8.14	0	0.538	6.047	88.8	4.4534	4	307	21	306.38	17.28	14.8	
26	0.77299	0	8.14	0	0.538	6.495	94.4	4.4547	4	307	21	387.94	12.8	18.4	
27	1.00245	0	8.14	0	0.538	6.674	87.3	4.239	4	307	21	380.23	11.98	21	

# Supervised Learning – You Already Have Target Data

**Regression:**

**How much will a customers spend?**

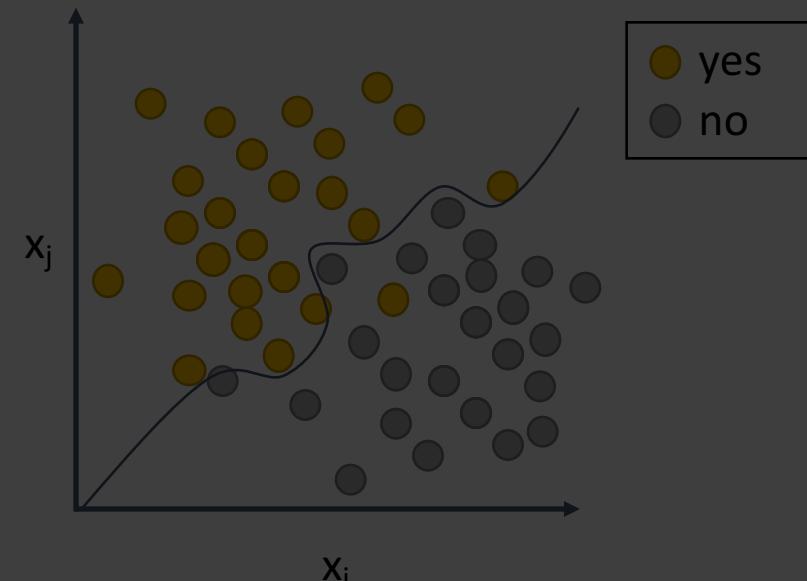


**H<sub>2</sub>O algos:**

**Penalized Linear Models**  
**Random Forest**  
**Gradient Boosting**  
**Neural Networks**  
**Stacked Ensembles**

**Classification:**

**Will a customer make a purchase? Yes or No**



**H<sub>2</sub>O algos:**

**Penalized Linear Models**  
**Naïve Bayes**  
**Random Forest**  
**Gradient Boosting**  
**Neural Networks**  
**Stacked Ensembles**

# Regression – Key Steps

- Import Data (CSV Files)
  - ./data/regression/ ...
  - Have a quick look
- Train a Random Forest Model
  - Look at variable importance
- Make Predictions
  - Compare with ground truth
- Build Partial Dependence Plots
  - Explain variable



H2O Flow Jo-fai

localhost:54321/flow/index.html

Untitled Flow

Data ▾ Model ▾ Score ▾ Admin ▾ Help ▾

Import Files...  
Upload File...  
Split Frame...  
Merge Frames...  
List All Frames  
Impute...

cs | assist

67ms

?

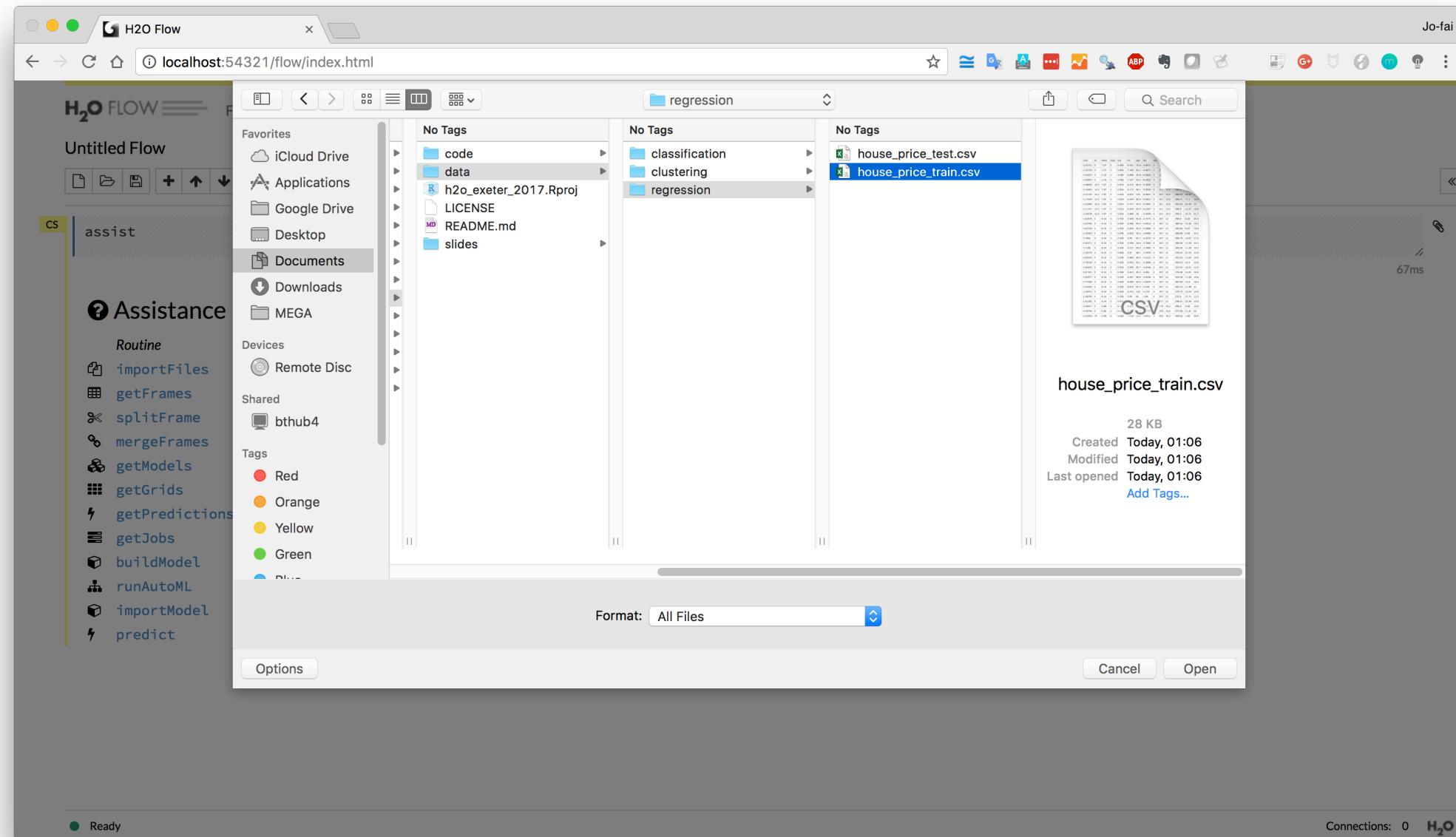
### Assistance

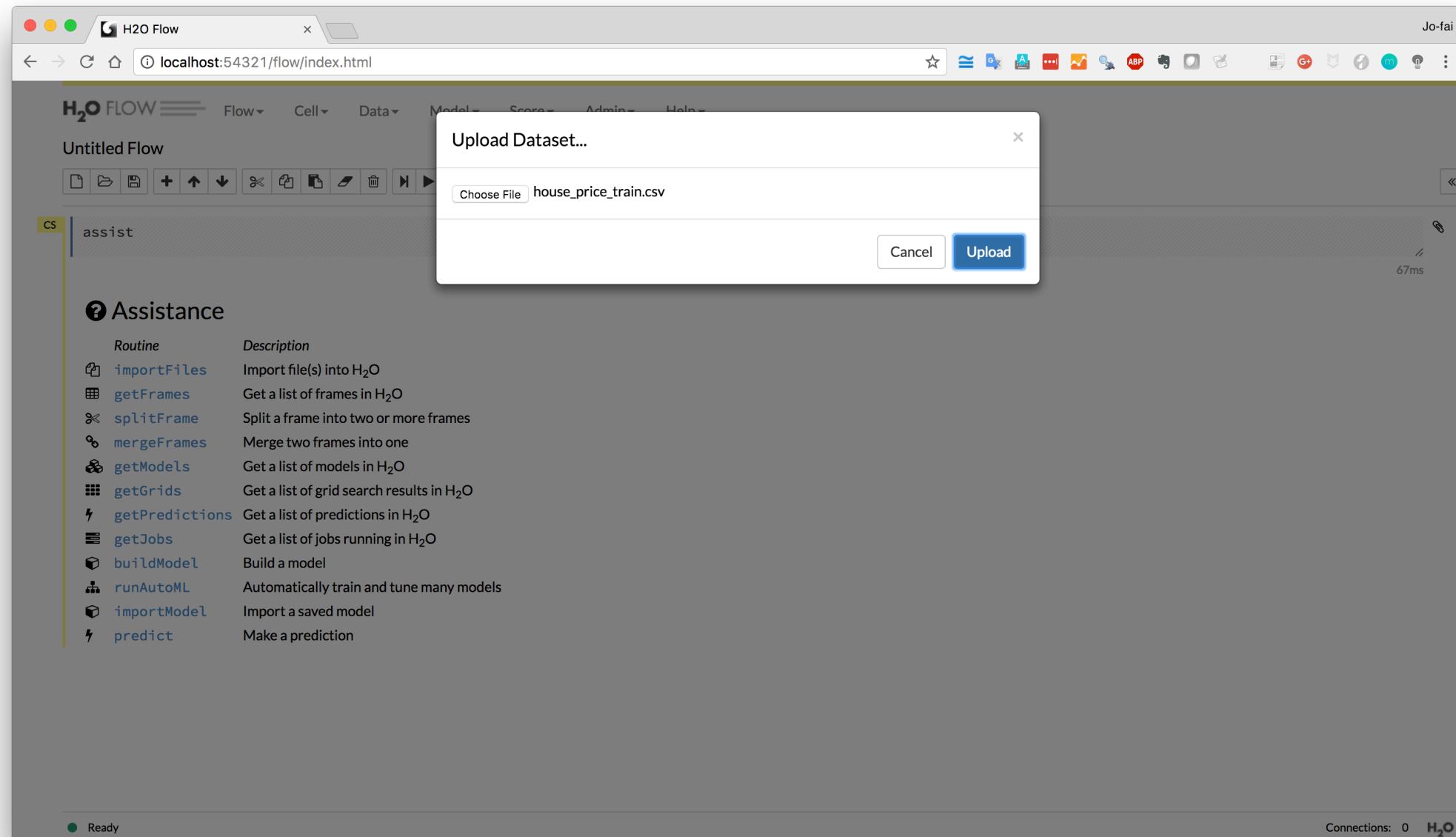
Routine	Description
importFiles	Import file(s) into H2O
getFrames	Get a list of frames in H2O
splitFrame	Split a frame into two or more frames
mergeFrames	Merge two frames into one
getModels	Get a list of models in H2O
getGrids	Get a list of grid search results in H2O
getPredictions	Get a list of predictions in H2O
getJobs	Get a list of jobs running in H2O
buildModel	Build a model
runAutoML	Automatically train and tune many models
importModel	Import a saved model
predict	Make a prediction

localhost:54321/flow/index.html#

Connections: 0 H2O

The screenshot shows the H2O Flow web application. At the top, there's a navigation bar with tabs for Data, Model, Score, Admin, and Help. The Data tab is currently active and has a yellow background. A dropdown menu is open under the Data tab, listing several options: Import Files..., Upload File..., Split Frame..., Merge Frames..., List All Frames, and Impute... . Below the menu, the main workspace is visible, showing a sidebar with a table of assistance routines and a central area with a 'assist' frame. The bottom of the screen shows the browser's address bar with 'localhost:54321/flow/index.html#' and a status bar indicating 'Connections: 0' and the H2O logo.





H2O Flow Jo-fai

localhost:54321/flow/index.html

## Untitled Flow

### Setup Parse

**PARSE CONFIGURATION**

Sources: house\_price\_train.csv  
ID: Key\_Frame\_house\_price\_train.hex  
Parser: CSV  
Separator: ;'44'  
Column Headers:  Auto  
 First row contains column names  
 First row contains data  
Options:  Enable single quotes as a field quotation character  
 Delete on done

**EDIT COLUMN NAMES AND TYPES**

Search by column name...

	Column Name	Type	Value 1	Value 2	Value 3	Value 4	Value 5	Value 6	Value 7	Value 8	Value 9
1	crim	Numeric	0.02731	0.02729	0.03237	0.06905	0.08829	0.14455	0.21124	0.17004	0.22489
2	zn	Numeric	0	0	0	0	12.5	12.5	12.5	12.5	12.5
3	indus	Numeric	7.07	7.07	2.18	2.18	7.87	7.87	7.87	7.87	7.87
4	chas	Numeric	0	0	0	0	0	0	0	0	0
5	nox	Numeric	0.469	0.469	0.458	0.458	0.524	0.524	0.524	0.524	0.524
6	rm	Numeric	6.421	7.185	6.998	7.147	6.012	6.172	5.631	6.004	6.377
7	age	Numeric	78.9	61.1	45.8	54.2	66.6	96.1	100	85.9	94.3
8	dis	Numeric	4.9671	4.9671	6.0622	6.0622	5.5605	5.9505	6.0821	6.5921	6.3467

Ready Connections: 0 H2O

H<sub>2</sub>O Flow Jo-fai

localhost:54321/flow/index.html

Untitled Flow

Delete on done

EDIT COLUMN NAMES AND TYPES

Search by column name...

	Column Name	Type	Value 1	Value 2	Value 3	Value 4	Value 5	Value 6	Value 7	Value 8	Value 9
1	crim	Numeric	0.02731	0.02729	0.03237	0.06905	0.08829	0.14455	0.21124	0.17004	0.22489
2	zn	Numeric	0	0	0	0	12.5	12.5	12.5	12.5	12.5
3	indus	Numeric	7.07	7.07	2.18	2.18	7.87	7.87	7.87	7.87	7.87
4	chas	Numeric	0	0	0	0	0	0	0	0	0
5	nox	Numeric	0.469	0.469	0.458	0.458	0.524	0.524	0.524	0.524	0.524
6	rm	Numeric	6.421	7.185	6.998	7.147	6.012	6.172	5.631	6.004	6.377
7	age	Numeric	78.9	61.1	45.8	54.2	66.6	96.1	100	85.9	94.3
8	dis	Numeric	4.9671	4.9671	6.0622	6.0622	5.5605	5.9505	6.0821	6.5921	6.3467
9	rad	Numeric	2	2	3	3	5	5	5	5	5
10	tax	Numeric	242	242	222	222	311	311	311	311	311
11	ptratio	Numeric	17.8	17.8	18.7	18.7	15.2	15.2	15.2	15.2	15.2
12	b	Numeric	396.9	392.83	394.63	396.9	395.6	396.9	386.63	386.71	392.52
13	lstat	Numeric	9.14	4.03	2.94	5.33	12.43	19.15	29.93	17.1	20.45
14	medv	Numeric	21.6	34.7	33.4	36.2	22.9	27.1	16.5	18.9	15

[Previous page](#) [Next page](#)

[Parse](#)

Ready Connections: 0 H<sub>2</sub>O

H2O Flow Jo-fai

localhost:54321/flow/index.html

## H2O FLOW

Untitled Flow

Flow Cell Data Model Score Admin Help

File Edit View Insert Cells Data Models Scores Admin Help

CS parseFiles

```
source_frames: ["house_price_train.csv"]
destination_frame: "Key_Frame__house_price_train.hex"
parse_type: "CSV"
separator: 44
number_columns: 14
single_quotes: false
column_names: ["crim","zn","indus","chas","nox","rm","age","dis","rad","tax","ptratio","b","lstat","medv"]
column_types:
["Numeric","Numeric","Numeric","Numeric","Numeric","Numeric","Numeric","Numeric","Numeric","Numeric","Numeric","Numeric"]
delete_on_done: true
check_header: 1
chunk_size: 4194304
```

1.1s

### Job

Run Time 00:00:00.107  
Remaining Time 00:00:00.0

Type Frame  
Key [Q Key\\_Frame\\_house\\_price\\_train.hex](#)

Description Parse  
Status DONE  
Progress 100%   
Done.

Actions [View](#)

Ready Connections: 0 H2O

H2O Flow Jo-fai

localhost:54321/flow/index.html

Untitled Flow

Flow Cell Data Model Score Admin Help

getFrameSummary "Key\_Frame\_\_house\_price\_train.hex"

86ms

Key\_Frame\_\_house\_price\_train.hex

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

Rows Columns Compressed Size

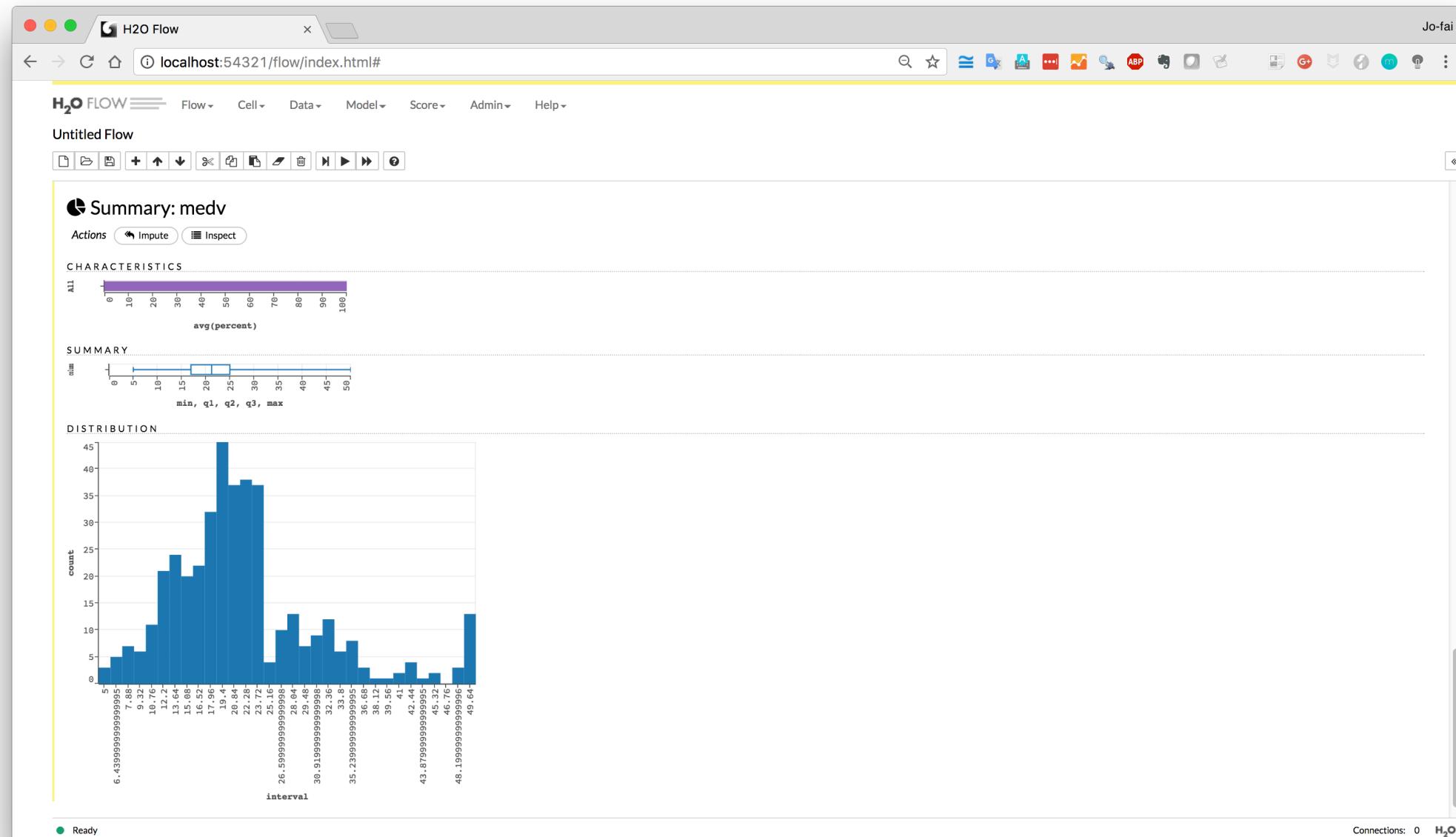
407 14 12KB

COLUMN SUMMARIES

label	type	Missing	Zeros	+Inf	-Inf	min	max	mean	sigma	cardinality	Actions
crim	real	0	0	0	0	0.0091	73.5341	3.5674	7.9480	...	
zn	real	0	301	0	0	0	100.0	10.5872	22.2598	...	
indus	real	0	0	0	0	0.4600	27.7400	11.4093	6.8145	...	
chas	int	0	379	0	0	0	1.0	0.0688	0.2534	...	Convert to enum
nox	real	0	0	0	0	0.3850	0.8710	0.5568	0.1156	...	
rm	real	0	0	0	0	3.8630	8.7250	6.2866	0.6909	...	
age	real	0	0	0	0	2.9000	100.0	69.3889	27.8179	...	
dis	real	0	0	0	0	1.1296	10.7103	3.7177	2.0152	...	
rad	int	0	0	0	0	1.0	24.0	9.8378	8.7844	...	Convert to enum
tax	int	0	0	0	0	188.0	711.0	412.3784	170.4474	...	Convert to enum
ptratio	real	0	0	0	0	12.6000	22.0	18.4474	2.1618	...	
b	real	0	0	0	0	0.3200	396.9000	354.4032	94.1752	...	
lstat	real	0	0	0	0	1.7300	37.9700	12.7920	7.0987	...	
medv	real	0	0	0	0	5.0	50.0	22.6248	9.1850	...	

Previous 20 Columns Next 20 Columns

Ready Connections: 0 H2O



H2O Flow    Flow ▾ Cell ▾

Untitled Flow

Actions: View Data Split...

Rows  
407

COLUMN SUMMARIES

label	type	Missing	Zeros	+Inf	-Inf
crim	real	0	0	0	0
zn	real	0	301	0	0
indus	real	0	0	0	0
chas	int	0	379	0	0
nox	real	0	0	0	0
rm	real	0	0	0	0
age	real	0	0	0	0
dis	real	0	0	0	0
rad	int	0	0	0	1.0
tax	int	0	0	0	188.0
ptratio	real	0	0	0	12.6000
b	real	0	0	0	0.3200
lstat	real	0	0	0	1.7300
medv	real	0	0	0	5.0

Format: All Files

Options Cancel Open

Previous 20 Columns Next 20 Columns

CHUNK COMPRESSION SUMMARY

FRAME DISTRIBUTION SUMMARY

Ready

Connections: 0 H2O

localhost:54321/flow/index.html#

regression

No Tags

house\_price\_test.csv

house\_price\_train.csv

CSV

7 KB

Created Today, 01:06

Modified Today, 01:06

Last opened Today, 01:06

Add Tags...

Delete

H2O Flow Jo-fai

localhost:54321/flow/index.html#

Untitled Flow

Type Frame  
Key Q Key\_Frame\_house\_price\_hex  
Description Parse  
Status DONE  
Progress 100% Done.  
Actions Q View

Import Files...  
Upload File...  
Split Frame...  
Merge Frames...  
List All Frames  
Impute...

getFrames 29ms

Frames

	Rows	Columns	Size
Key_Frame_house_price_hex	99	14	4KB
Key_Frame_house_price_train.hex	407	14	12KB

Type ID  
Key\_Frame\_house\_price\_hex  
Build Model... Predict... Inspect  
Key\_Frame\_house\_price\_train.hex  
Build Model... Predict... Inspect

Predict on selected frames... Delete selected frames

Connections: 0 H2O

A screenshot of the H2O Flow web application. At the top, there's a navigation bar with tabs for Data, Model, Score, Admin, and Help. The Data tab is currently active. Below the navigation is a toolbar with various icons for file operations like Import, Export, and Split. A context menu is open over the first frame, listing options such as Import Files..., Upload File..., Split Frame..., Merge Frames..., List All Frames, and Impute... The main workspace shows two frames listed under the 'Frames' section. The first frame, 'Key\_Frame\_house\_price\_hex', has 99 rows, 14 columns, and a size of 4KB. The second frame, 'Key\_Frame\_house\_price\_train.hex', has 407 rows, 14 columns, and a size of 12KB. Both frames have 'Build Model...', 'Predict...', and 'Inspect...' buttons associated with them. At the bottom, there are buttons for 'Predict on selected frames...' and 'Delete selected frames'. The status bar at the bottom shows the URL 'localhost:54321/flow/index.html#' and connection information 'Connections: 0 H2O'.

H2O Flow Jo-fai

localhost:54321/flow/index.html#

Model ▾ Score ▾ Admin ▾ Help ▾

Untitled Flow

Type Frame  
Key [Key\\_Frame\\_house\\_price\\_test.hex](#)  
Description Parse  
Status DONE  
Progress 100%  
Done.  
Actions [View](#)

getFrames

Frames

- Type ID
- [Key\\_Frame\\_house\\_price\\_test.hex](#)  
[Build Model...](#) [Predict...](#) [Inspect...](#)
- [Key\\_Frame\\_house\\_price\\_train.hex](#)  
[Build Model...](#) [Predict...](#) [Inspect...](#)

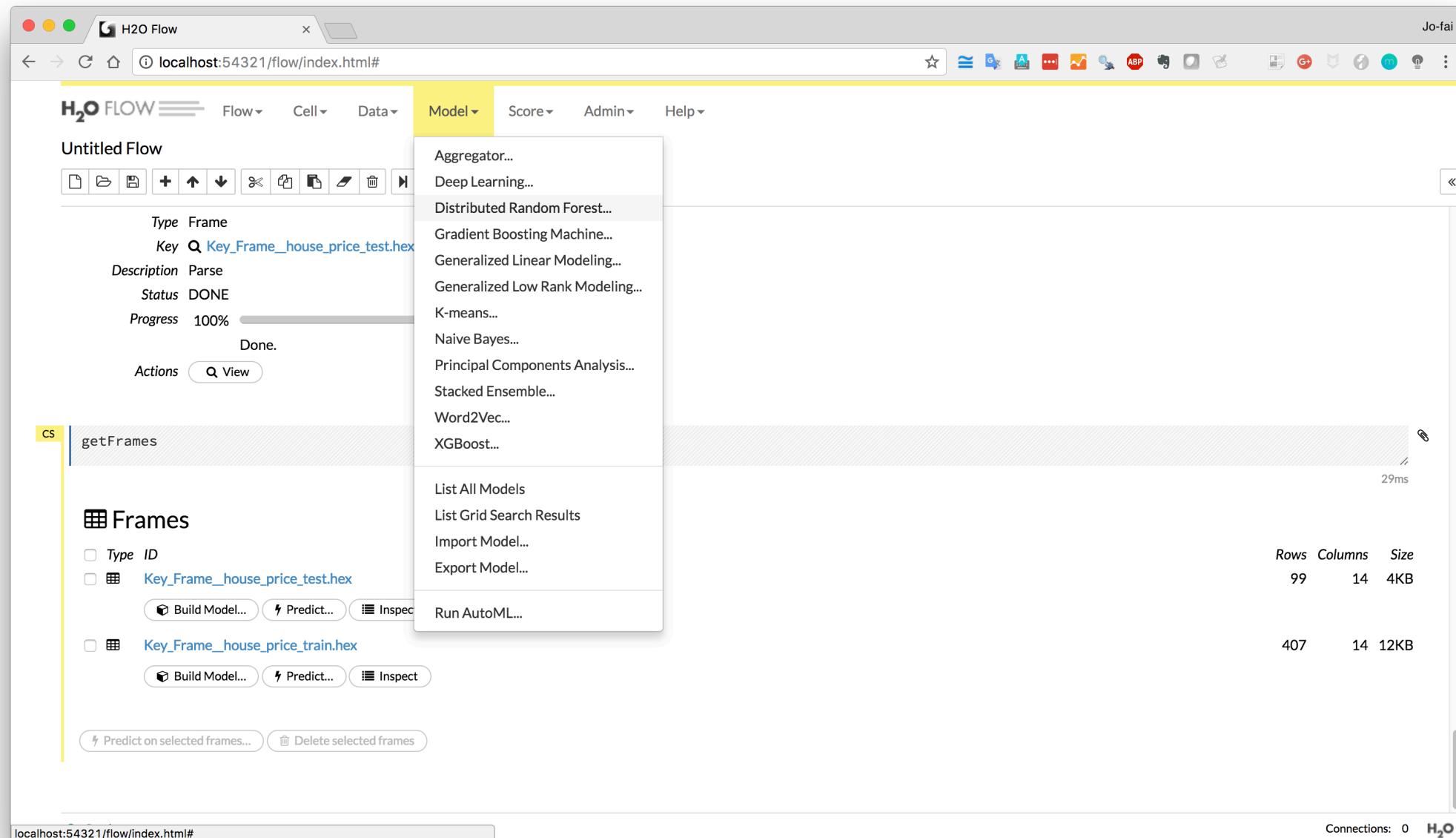
Predict on selected frames... Delete selected frames

Aggregator...  
Deep Learning...  
Distributed Random Forest...  
Gradient Boosting Machine...  
Generalized Linear Modeling...  
Generalized Low Rank Modeling...  
K-means...  
Naive Bayes...  
Principal Components Analysis...  
Stacked Ensemble...  
Word2Vec...  
XGBoost...

List All Models  
List Grid Search Results  
Import Model...  
Export Model...  
Run AutoML...

Rows	Columns	Size
99	14	4KB
407	14	12KB

Connections: 0 H2O



H2O Flow Jo-fai

localhost:54321/flow/index.html#

## H2O FLOW

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

### Untitled Flow

Build a Model

Select an algorithm: **Distributed Random Forest**

PARAMETERS

**model\_id** my\_random\_forest  
Destination id for this model; auto-generated if not specified.

**training\_frame** Key\_Frame\_house\_price\_train.hex  
Id of the training data frame (Not required, to allow initial validation of model parameters).

**validation\_frame** (Choose...)  
Id of the validation data frame.

**nfold** 0  
Number of folds for N-fold cross-validation (0 to disable or >= 2).

**response\_column** medv  
Response variable column.

**ignored\_columns** Search...  
Showing page 1 of 1.  

crim	REAL
zn	REAL
indus	REAL
chas	INT

Connections: 0 H2O

H2O Flow x localhost:54321/flow/index.html#

H2O FLOW Flow Cell Data Model Score Admin Help

Untitled Flow

Only show columns with more than  % missing values.

ignore\_const\_cols  Ignore constant columns.

ntrees 50 Number of trees.

max\_depth 20 Maximum tree depth.

min\_rows 1 Fewest allowed (weighted) observations in a leaf.

nbins 20 For numerical columns (real/int), build a histogram of (at least) this many bins, then split at the best point

seed 54321 Seed for pseudo random number generator (if applicable)

mtries -1 Number of variables randomly sampled as candidates at each split. If set to -1, defaults to  $\sqrt{p}$  for classification and  $p/3$  for regression (where  $p$  is the # of predictors)

sample\_rate 0.6320000290870667 Row sample rate per tree (from 0.0 to 1.0)

ADVANCED GRID ?

score\_each\_iteration Whether to score during each iteration of model training.

score\_tree\_interval 0 Score the model after every so many trees. Disabled if set to 0.

fold\_column (Choose...) Column with cross-validation fold index assignment per observation.

offset\_column (Choose...) Offset column. This will be added to the combination of columns before applying the link function.

weights\_column (Choose...) Column with observation weights. Giving some observation a weight of zero is equivalent to excluding it from the dataset; giving an observation a relative weight of 2 is equivalent to repeating that row twice. Negative weights are not allowed.

nbins\_top\_level 1024 For numerical columns (real/int), build a histogram of (at most) this many bins at the root level, then decrease by factor of two per level

nbins\_cats 1024 For categorical columns (factors), build a histogram of this many bins, then split at the best point. Higher values can lead to more overfitting.

r2\_stopping 1.7976931348623157e+308 r2\_stopping is no longer supported and will be ignored if set - please use stopping\_rounds, stopping\_metric and stopping\_tolerance

Connections: 0 H2O

H2O Flow x Jo-fai

localhost:54321/flow/index.html#

## H2O FLOW

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

### Untitled Flow

File ▾ + ▾ Up ▾ Down ▾ Cut ▾ Copy ▾ Paste ▾ Delete ▾ Next ▾ Previous ▾ Help ▾

histogram\_type: AUTO What type of histogram to use for finding optimal split points

categorical\_encoding: AUTO Encoding scheme for categorical features

**EXPERT**

build\_tree\_one\_node:  Run on one node only; no network overhead but fewer cpus used. Suitable for small datasets.

sample\_rate\_per\_class: A list of row sample rates per class (relative fraction for each class, from 0.0 to 1.0), for each tree

binomial\_double\_trees: For binary classification: Build 2x as many trees (one per class) - can lead to higher accuracy.

col\_sample\_rate\_change\_per\_level: 1 Relative change of the column sampling rate for every level (from 0.0 to 2.0)

calibrate\_model:  Use Platt Scaling to calculate calibrated class probabilities. Calibration can provide more accurate estimates of class probabilities.

calibration\_frame: (Choose...) Calibration frame for Platt Scaling

Build Model

Ready Connections: 0 H2O

H2O Flow Jo-fai

localhost:54321/flow/index.html#

## Untitled Flow

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

Build Model

```
CS buildModel 'drf', {"model_id":"my_random_forest","training_frame":"Key_Frame__house_price_train.hex","nfolds":0,"response_column":"medv","ignored_columns":[],"ignore_const_cols":true,"ntrees":50,"max_depth":20,"min_rows":1,"nbins":20,"seed":54321,"mtries":-1,"sample_rate":0.6320000290870667,"score_each_iteration":false,"score_tree_interval":0,"nbins_top_level":1024,"nbins_cats":1024,"r2_stopping":1.7976931348623157e+308,"stopping_rounds":0,"stopping_metric":"AUTO","stopping_tolerance":0.001,"max_runtime_secs":0,"checkpoint":"","col_sample_rate_per_tree":1,"min_split_improvement":0.00001,"histogram_type":"AUTO","categorical_encoding":"AUTO","build_tree_one_node":false,"sample_rate_per_class":[],"binomial_double_trees":false,"col_sample_rate_change_per_level":1,"calibrate_model":false}
```

1.1s

### Job

Run Time 00:00:00.654  
Remaining Time 00:00:00.0

Type Model  
Key Q my\_random\_forest  
Description DRF  
Status DONE  
Progress 100% Done.  
Actions View

Ready Connections: 0 H2O

H2O Flow Jo-fai

localhost:54321/flow/index.html#

## Untitled Flow

### Model

Model ID: my\_random\_forest  
Algorithm: Distributed Random Forest

Actions: Refresh, Predict..., Download POJO, Download Model Deployment Package (MOJO), Export, Inspect, Delete, Download Gen Model

MODEL PARAMETERS

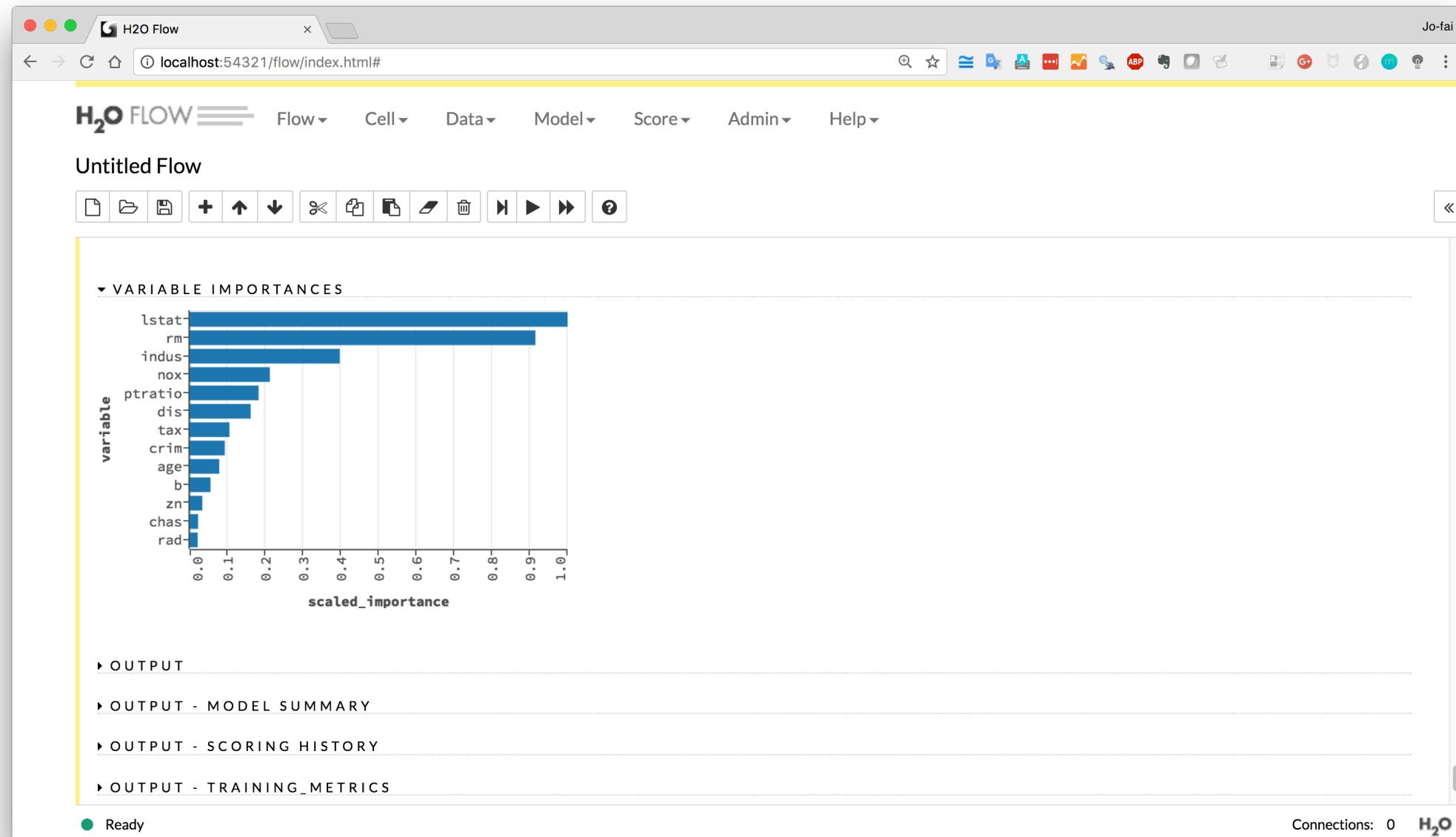
SCORING HISTORY - DEVIANCE

A line graph titled "SCORING HISTORY - DEVIANCE". The y-axis is labeled "training\_deviance" and ranges from 10 to 26. The x-axis is labeled "number\_of\_trees" and ranges from 0 to 50. The data points show a sharp initial drop in deviance from approximately 19 at 1 tree to about 12.5 at 10 trees, after which it levels off and continues to decrease very gradually towards 10.5 at 50 trees.

VARIABLE IMPORTANCES

Ready

Connections: 0 H2O



H2O Flow Jo-fai

localhost:54321/flow/index.html#

Untitled Flow

Score Admin Help

Predict... Partial Dependence Plots... List All Predictions

predict 17ms

**⚡ Predict**

Name: my\_predictions

Model: my\_random\_forest

Frame: Key\_Frame\_house\_price\_test.hex

Actions: ⚡ Predict

Ready localhost:54321/flow/index.html#

Connections: 0 H2O

The screenshot shows the H2O Flow web application. At the top, there's a navigation bar with tabs for Score, Admin, and Help. A dropdown menu under Score is open, showing options like 'Predict...', 'Partial Dependence Plots...', and 'List All Predictions'. Below the navigation is a toolbar with various icons for file operations and flow management. The main workspace is titled 'Untitled Flow' and contains a single step labeled 'predict'. This step has a yellow border and a duration of '17ms' indicated next to it. To the left of the step, there's a vertical yellow line labeled 'cs'. Below the step, there's a section titled '⚡ Predict' with four input fields: 'Name' set to 'my\_predictions', 'Model' set to 'my\_random\_forest', 'Frame' set to 'Key\_Frame\_house\_price\_test.hex', and an 'Actions' button labeled '⚡ Predict'. At the bottom of the screen, there's a status bar showing 'Ready' and the URL 'localhost:54321/flow/index.html#'. On the right side of the status bar, it says 'Connections: 0' and has the H2O logo.

H2O Flow Jo-fai

localhost:54321/flow/index.html#

H2O FLOW Flow Cell Data Model Score Admin Help

Untitled Flow

CS predict model: "my\_random\_forest", frame: "Key\_Frame\_\_house\_price\_test.hex", predictions\_frame: "my\_predictions" 53ms

**⚡ Prediction**

Actions: Inspect

PREDICTION

```
model my_random_forest
model_checksum 2387713827283948032
frame Key_Frame__house_price_test.hex
frame_checksum -3619964009561634304
description .
model_category Regression
scoring_time 1507165392633
predictions my_predictions
MSE 10.860009
RMSE 3.295453
nobs 99
r2 0.872714
mean_residual_deviance 10.860009
mae 2.415040
rmsle 0.162378
```

Combine predictions with frame

Ready Connections: 0 H2O

The screenshot shows the H2O Flow web application running in a browser. The title bar says 'H2O Flow' and 'Jo-fai'. The address bar shows 'localhost:54321/flow/index.html#'. The main menu includes 'Flow', 'Cell', 'Data', 'Model', 'Score', 'Admin', and 'Help'. Below the menu is a toolbar with various icons. The main workspace has a yellow vertical bar on the left. A code snippet 'predict model: "my\_random\_forest", frame: "Key\_Frame\_\_house\_price\_test.hex", predictions\_frame: "my\_predictions"' is shown with a timestamp '53ms'. Below it, a section titled '⚡ Prediction' displays detailed model statistics. At the bottom, there's a button 'Combine predictions with frame' and status indicators 'Ready' and 'Connections: 0'.

H2O Flow Jo-fai

localhost:54321/flow/index.html#

H2O FLOW Flow Cell Data Model Score Admin Help

Untitled Flow

combined-my\_predictions

DATA

Previous 20 Columns Next 20 Columns

Row	predict	crim	zn	indus	chas	nox	rm	age	dis	rad	tax	ptratio	b	lstat	medv
1	27.0670	0.0063	18.0	2.3100	0	0.5380	6.5750	65.2000	4.0900	1.0	296.0	15.3000	396.9000	4.9800	24.0
2	25.9460	0.0299	0	2.1800	0	0.4580	6.4300	58.7000	6.0622	3.0	222.0	18.7000	394.1200	5.2100	28.7000
3	16.9850	1.2325	0	8.1400	0	0.5380	6.1420	91.7000	3.9769	4.0	307.0	21.0	396.9000	18.7200	15.2000
4	15.5270	0.9884	0	8.1400	0	0.5380	5.8130	100.0	4.0952	4.0	307.0	21.0	394.5400	19.8800	14.5000
5	15.0130	1.1308	0	8.1400	0	0.5380	5.7130	94.1000	4.2330	4.0	307.0	21.0	360.1700	22.6000	12.7000
6	15.6990	1.1517	0	8.1400	0	0.5380	5.7010	95.0	3.7872	4.0	307.0	21.0	358.7700	18.3500	13.1000
7	21.8920	0.0801	0	5.9600	0	0.4990	5.8500	41.5000	3.9342	5.0	279.0	19.2000	396.9000	8.7700	21.0
8	20.3540	0.1751	0	5.9600	0	0.4990	5.9660	30.2000	3.8473	5.0	279.0	19.2000	393.4300	10.1300	24.7000
9	29.0590	0.0276	75.0	2.9500	0	0.4280	6.5950	21.8000	5.4011	3.0	252.0	18.3000	395.6300	4.3200	30.8000
10	20.9370	0.0887	21.0	5.6400	0	0.4390	5.9630	45.7000	6.8147	4.0	243.0	16.8000	395.5600	13.4500	19.7000
11	23.3900	0.0205	85.0	0.7400	0	0.4100	6.3830	35.7000	9.1876	2.0	313.0	17.3000	396.9000	5.7700	24.7000
12	19.1160	0.1717	25.0	5.1300	0	0.4530	5.9660	93.4000	6.8185	8.0	284.0	19.7000	378.0800	14.4400	16.0
13	23.7010	0.1103	25.0	5.1300	0	0.4530	6.4560	67.8000	7.2255	8.0	284.0	19.7000	396.9000	6.7300	22.2000
14	21.9820	0.0839	0	12.8300	0	0.4370	5.8740	36.6000	4.5026	5.0	398.0	18.7000	396.0600	9.1000	20.3000
15	23.2093	0.0355	25.0	4.8600	0	0.4260	6.1670	46.7000	5.4007	4.0	281.0	19.0	390.6400	7.5100	22.9000
16	22.9100	0.0715	0	1.1900	0	0.1100	6.1210	56.8000	3.7476	3.0	217.0	18.5000	395.1500	8.1100	22.2000

Ready Connections: 0 H2O

H2O Flow Jo-fai

localhost:54321/flow/index.html#

**Untitled Flow**

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

Rows: 99 Columns: 15 Compressed Size: 5KB

**COLUMN SUMMARIES**

label	type	Missing	Zeros	+Inf	-Inf	min	max	mean	sigma	cardinality	Actions
predict	real	0	0	0	0	8.1564	46.3040	22.7819	7.9191	...	
crim	real	0	0	0	0	0.0063	88.9762	3.8033	10.9319	...	
zn	real	0	71	0	0	0	95.0	14.5556	27.1545	...	
indus	real	0	0	0	0	0.7400	27.7400	10.0166	6.9690	...	
chas	int	0	92	0	0	0	1.0	0.0707	0.2576	...	Convert to enum
nox	real	0	0	0	0	0.3890	0.8710	0.5462	0.1174	...	
rm	real	0	0	0	0	3.5610	8.7800	6.2766	0.7525	...	
age	real	0	0	0	0	9.9000	100.0	65.2283	29.3788	...	
dis	real	0	0	0	0	1.1781	12.1265	4.1131	2.4283	...	
rad	int	0	0	0	0	1.0	24.0	8.3636	8.3207	...	Convert to enum
tax	int	0	0	0	0	187.0	711.0	391.2121	160.1494	...	Convert to enum
ptratio	real	0	0	0	0	12.6000	22.0	18.4889	2.1887	...	
b	real	0	0	0	0	6.6800	396.9000	366.0096	78.0843	...	
lstat	real	0	0	0	0	2.8800	34.4100	12.0817	7.3215	...	
medv	real	0	0	0	0	5.6000	50.0	22.1545	9.2839	...	

Previous 20 Columns Next 20 Columns

**CHUNK COMPRESSION SUMMARY**

Ready Connections: 0 H2O

H2O Flow Jo-fai

localhost:54321/flow/index.html#

H2O FLOW Flow Cell Data Model Score Admin Help

Untitled Flow

buildPartialDependence 17ms

Partial Dependence

Save Destination PDP my\_ppd as:

Model: my\_random\_forest

Frame: Key\_Frame\_house\_price\_train.hex

nbins 20

Select columns?

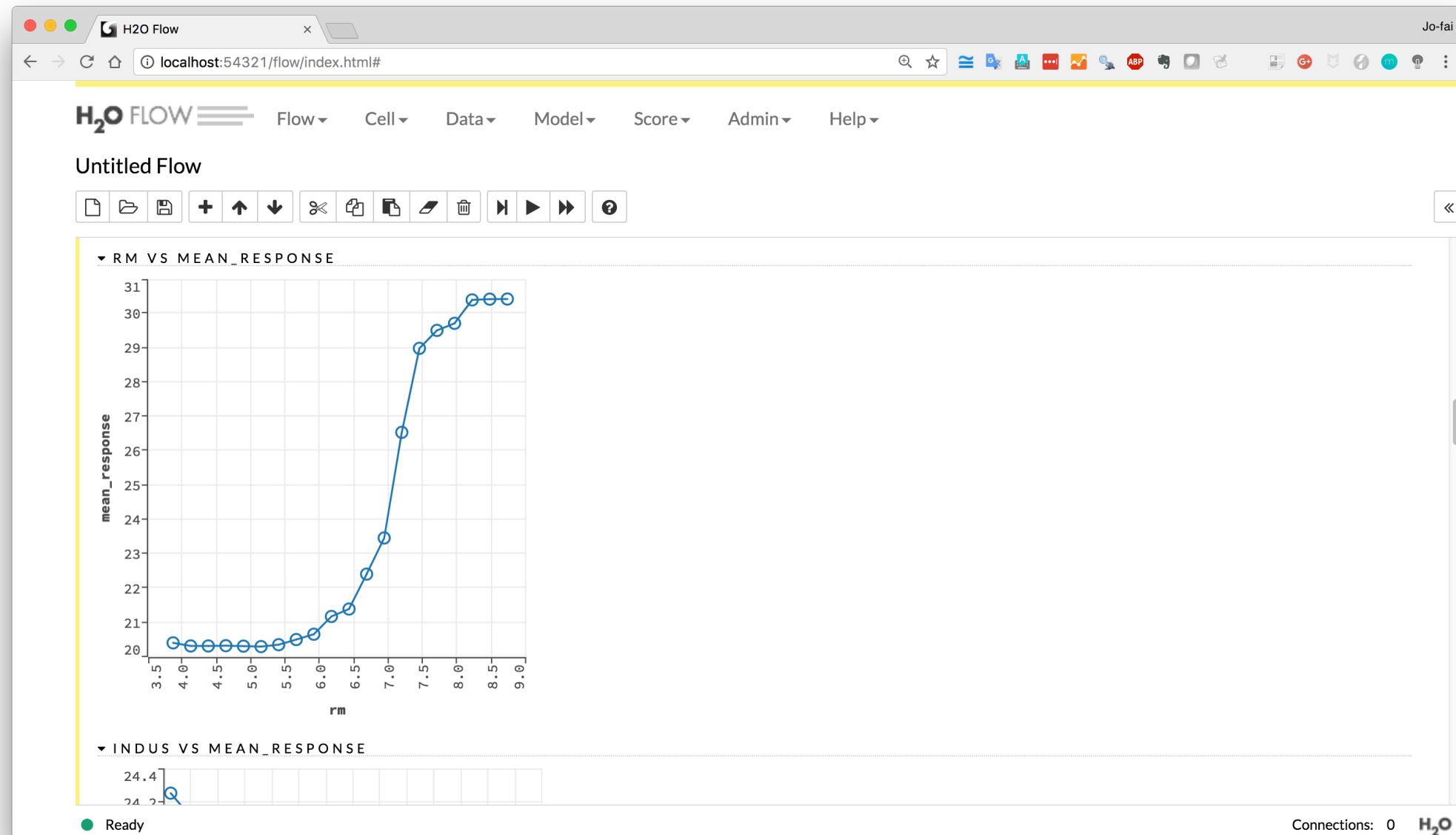
Actions: Compute

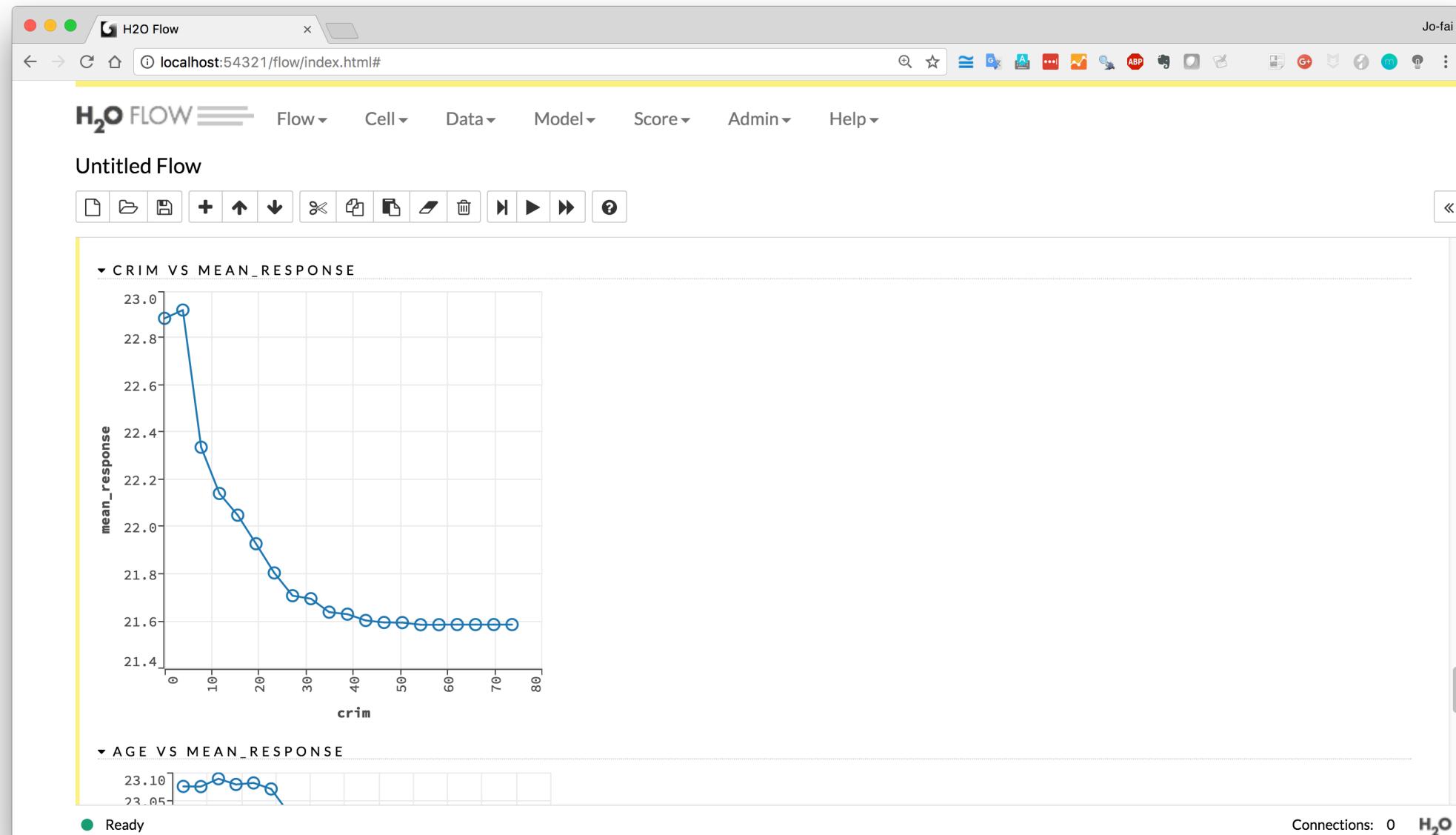
How many levels should PDP compute. More levels will make it slower.

Checking this will allow you to select custom columns for PDP. By default, the top 10 features are used. Those features are sorted by variable importance.

Ready Connections: 0 H2O

68





# Tutorial 2 - Classification

- **Data:** Human Activity Recognition Using Smartphone Sensors (2012)
- **Source:**  
<https://archive.ics.uci.edu/ml/datasets/human+activity+recognition+using+smartphones>

# Tutorial 3 - Clustering

- **Data:** Water Treatment Plant (1993)
- **Source:** <https://archive.ics.uci.edu/ml/datasets/water+treatment+plant>

# More About H<sub>2</sub>O

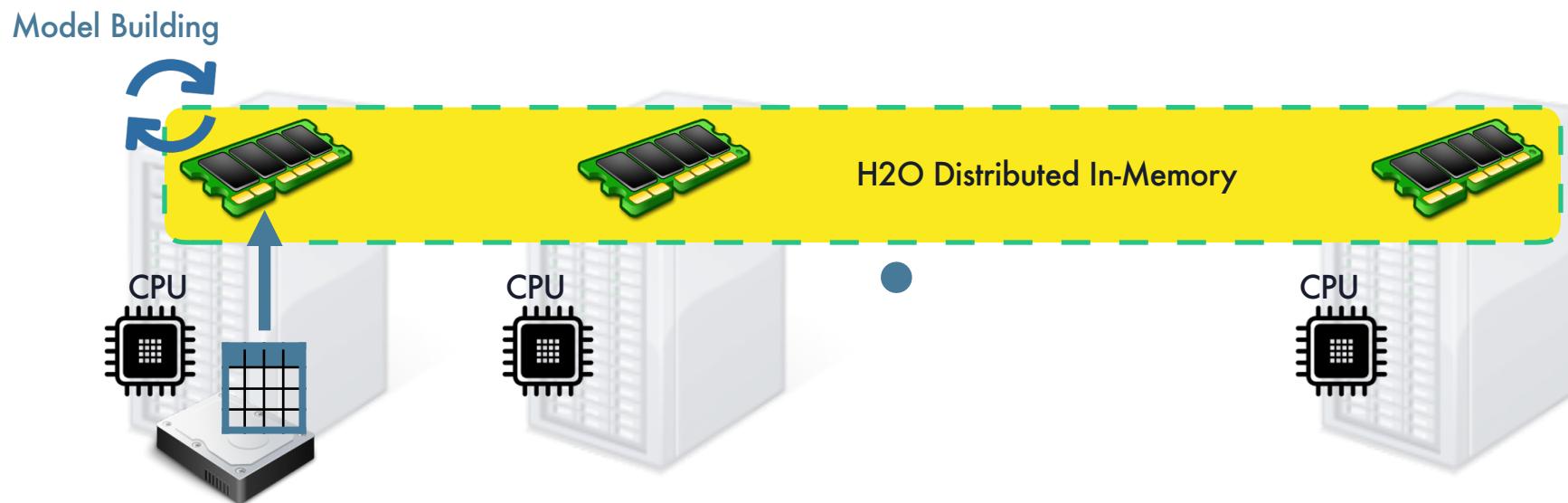
- Running H<sub>2</sub>O on Multi-Node Cluster
- Next-Gen H<sub>2</sub>O Products

# H<sub>2</sub>O on Multi-Node Cluster

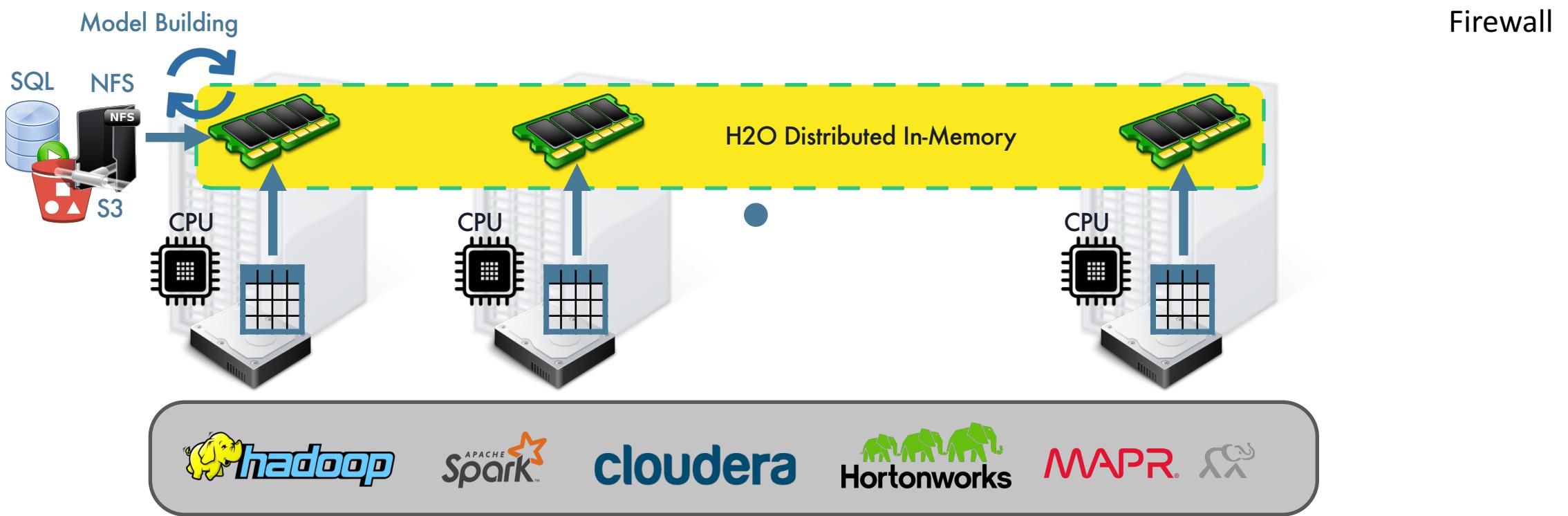
# H<sub>2</sub>O on Single Machine



# H<sub>2</sub>O on Multi-Node Cluster



# H<sub>2</sub>O with Distributed Data Storage Systems



**11M Rows****Size (Raw): 7.48 GB****Compressed: 2.00 GB ( $\approx$  27% of Raw)**

## HIGGS.hex

Actions:

View Data

Split...

Build Model...

Predict

Download

Export

Rows	Columns	Compressed Size
11000000	29	2GB

**▼ COLUMN SUMMARIES**

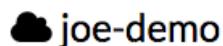
label	type	Missing	Zeros	+Inf	-Inf	min	max	mean	sigma	cardinality	Actions
C1	enum	0	5170877	0	0	0	1.0	0.5299	0.4991	2	<a href="#">Convert to numeric</a>
C2	real	0	0	0	0	0.2747	12.0989	0.9915	0.5654	..	..
C3	real	0	0	0	0	-2.4350	2.4349	-0.0	1.0088	..	..
C4	real	0	0	0	0	-1.7425	1.7432	-0.0	1.0063	..	..
C5	real	0	0	0	0	0.0002	15.3968	0.9985	0.6000	..	..
C6	real	0	0	0	0	-1.7439	1.7433	0.0	1.0063	..	..
C7	real	0	0	0	0	0.1375	9.9404	0.9909	0.4750	..	..
C8	real	0	0	0	0	-2.9697	2.9697	-0.0	1.0093	..	..
C9	real	0	0	0	0	-1.7412	1.7415	0.0	1.0059	..	..
C10	real	0	5394611	0	0	0	2.1731	1.0	1.0278	..	..
C11	real	0	0	0	0	0.1890	11.6471	0.9927	0.5000	..	..
C12	real	0	0	0	0	-2.9131	2.9132	-0.0	1.0093	..	..
C13	real	0	0	0	0	-1.7424	1.7432	-0.0	1.0062	..	..
C14	real	0	5523912	0	0	0	2.2149	1.0	1.0494	..	..
C15	real	0	0	0	0	0.2636	14.7090	0.9923	0.4877	..	..
C16	real	0	0	0	0	-2.7297	2.7300	0.0	1.0087	..	..
C17	real	0	0	0	0	-1.7421	1.7429	0.0	1.0063	..	..
C18	real	0	6265240	0	0	0	2.5482	1.0	1.1937	..	..
C19	real	0	0	0	0	0.3654	12.8826	0.9861	0.5058	..	..
C20	real	0	0	0	0	-2.4973	2.4980	-0.0	1.0077	..	..

## Untitled Flow



CS

getCloud



## CLOUD STATUS

HEALTHY	CONSENSUS	LOCKED
Version	Started	Nodes (Used / All)
3.13.0.3981	a minute ago	10 / 10

## NODES

Name	Ping	Cores	Load	My CPU %	Sys	Shut Down	Data (Used/Total)	Data (% Cached)	GC (Free / Total / Max)	Disk (Free / Max)	Disk (% Free)
✓ 172.16.2.181:54323	a few seconds ago	32	6.110	0	8	-	40.603	33.82 GB / s	29.46 GB / NaN undefined / 29.58 GB	339.08 GB / 1.70 TB	19%
✓ 172.16.2.182:54321	a few seconds ago	32	0.240	7	8	-	44.566	39.59 GB / s	29.43 GB / NaN undefined / 29.58 GB	225.64 GB / 1.70 TB	12%
✓ 172.16.2.183:54321	a few seconds ago	32	9.820	0	3	-	44.883	42.09 GB / s	29.34 GB / NaN undefined / 29.58 GB	450.18 GB / 1.70 TB	25%
✓ 172.16.2.184:54323	a few seconds ago	32	0.990	0	0	-	44.656	41.67 GB / s	29.51 GB / NaN undefined / 29.58 GB	254.96 GB / 1.70 TB	14%
✓ 172.16.2.185:54323	a few seconds ago	32	0.440	8	8	-	43.128	38.33 GB / s	29.43 GB / NaN undefined / 29.58 GB	501.02 GB / 1.70 TB	28%
✓ 172.16.2.186:54321	a few seconds ago	32	1.750	0	0	-	44.589	42.46 GB / s	29.42 GB / NaN undefined / 29.58 GB	331.27 GB / 1.70 TB	18%
✓ 172.16.2.187:54323	a few seconds ago	32	1.490	0	10	-	43.993	42.00 GB / s	29.46 GB / NaN undefined / 29.58 GB	367.40 GB / 1.70 TB	21%
✓ 172.16.2.188:54321	a few seconds ago	32	0.610	0	8	-	41.977	18.63 GB / s	28.30 GB / NaN undefined / 29.58 GB	218.27 GB / 1.70 TB	12%
✓ 172.16.2.189:54323	a few seconds ago	32	4.420	6	9	-	48.590	38.91 GB / s	29.34 GB / NaN undefined / 29.58 GB	477.97 GB / 1.70 TB	27%
✓ 172.16.2.190:54323	a few seconds ago	32	2.970	10	12	-	43.931	22.15 GB / s	29.51 GB / NaN undefined / 29.58 GB	274.50 GB / 1.70 TB	15%
✓ TOTAL	-	320	28.840	-	-	-	440.916	359.62 GB / s	293.18 GB / NaN undefined / 295.83 GB	3.36 TB / 17.04 TB	19%

$$10 \times 32 = \\ 320 \text{ Cores}$$

$$10 \times 29.6 = 296 \\ \text{GB Memory}$$



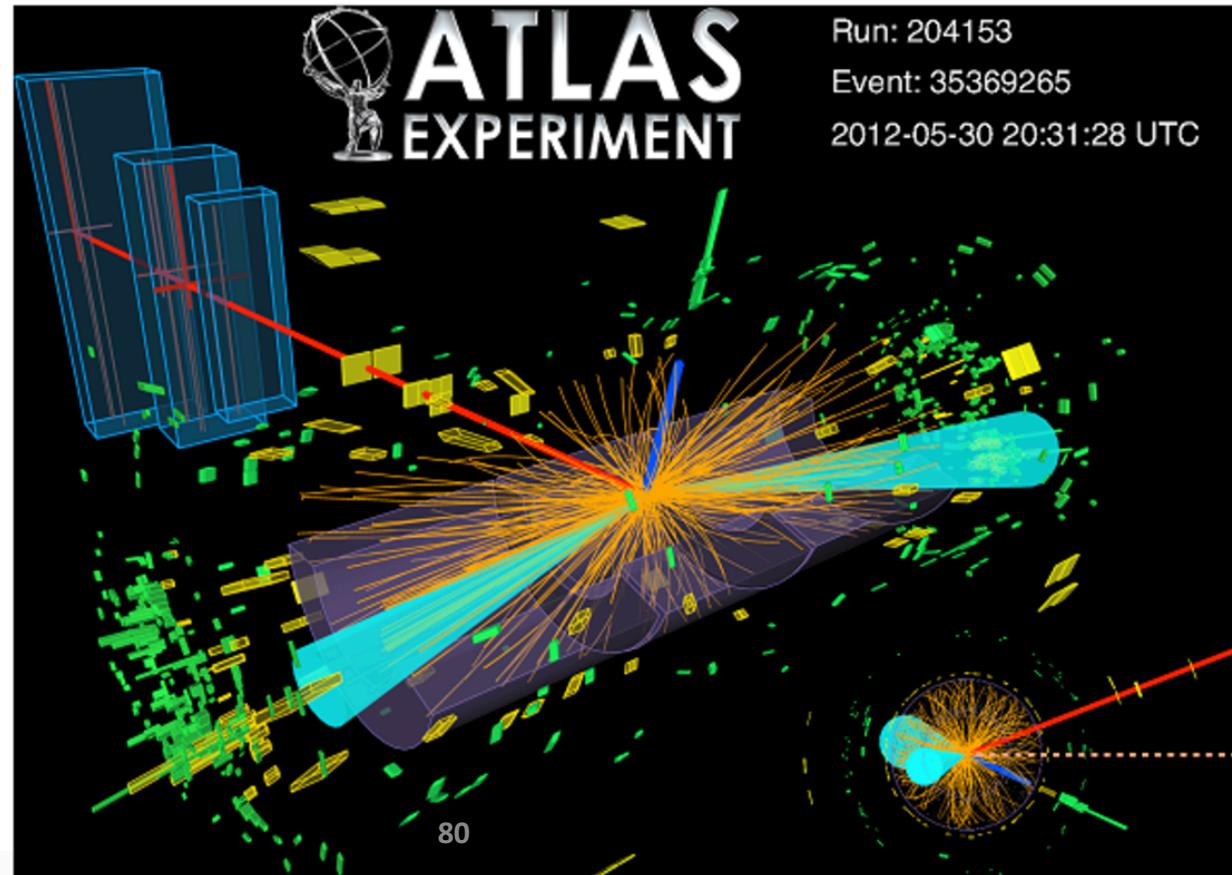
## Higgs Boson Machine Learning Challenge

Use the ATLAS experiment to identify the Higgs boson

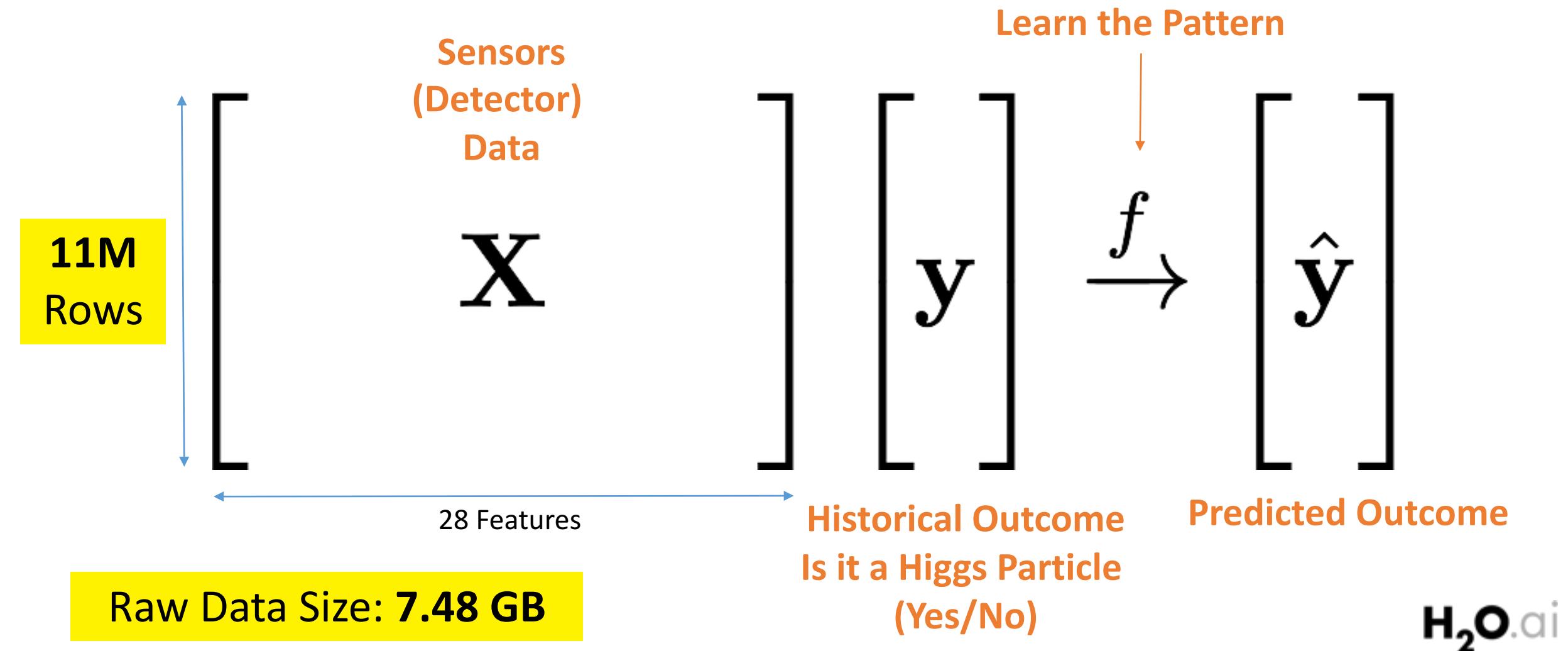
\$13,000 · 1,785 teams · 3 years ago

[Overview](#)[Data](#)[Discussion](#)[Leaderboard](#)[Rules](#)[Team](#)[My Submissions](#)[Late Submission](#)[Overview](#)

<https://www.kaggle.com/c/higgs-boson>

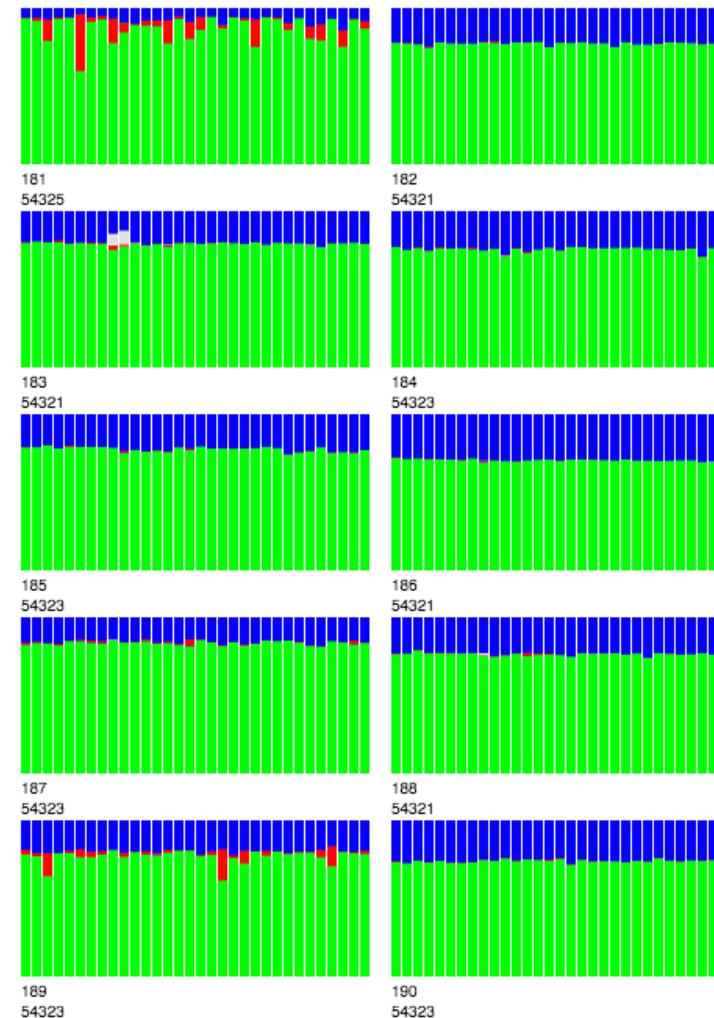
[Description](#)[Evaluation](#)[Prizes](#)[About The Sponsors](#)[Timeline](#)[Winners](#)

# Learning from Higgs Boson Dataset



# H<sub>2</sub>O Water Meter (CPU Usage Monitor)

$$10 \times 32 = 320 \text{ Cores}$$



## Legend

Each bar represents one CPU.

Blue: idle time

Green: user time

Red: system time

White: other time (e.g. i/o)

# H<sub>2</sub>O4GPU

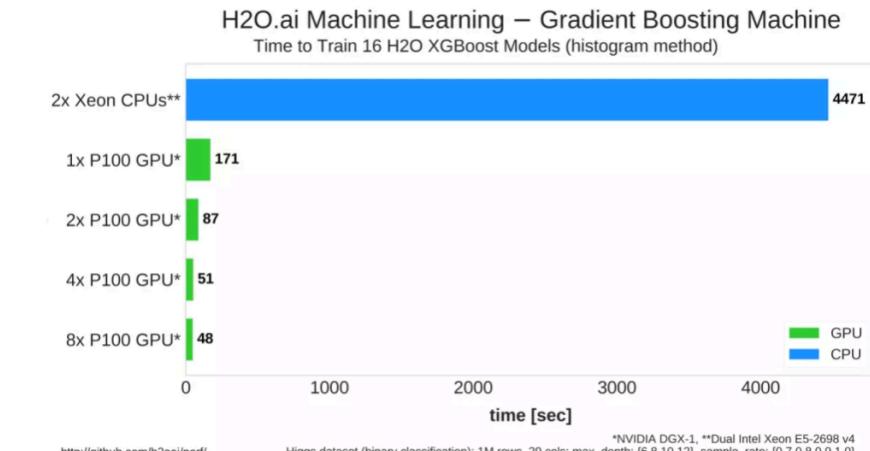
**H2O.ai Releases H2O4GPU, the Fastest Collection of GPU Algorithms on the Market, to Expedite Machine Learning in Python**

BY AVNI WADHWA ON SEPTEMBER 26, 2017 – 0 COMMENTS

H2O4GPU is an open-source collection of GPU solvers created by H2O.ai. It builds on the easy-to-use scikit-learn Python API and its well-tested CPU-based algorithms. It can be used as a drop-in replacement for scikit-learn with support for GPUs on selected (and ever-growing) algorithms. H2O4GPU inherits all the existing scikit-learn algorithms and falls back to CPU algorithms when the GPU algorithm does not support an important existing scikit-learn class option. It utilizes the efficient parallelism and high throughput of GPUs. Additionally, GPUs allow the user to complete training and inference much faster than possible on ordinary CPUs.

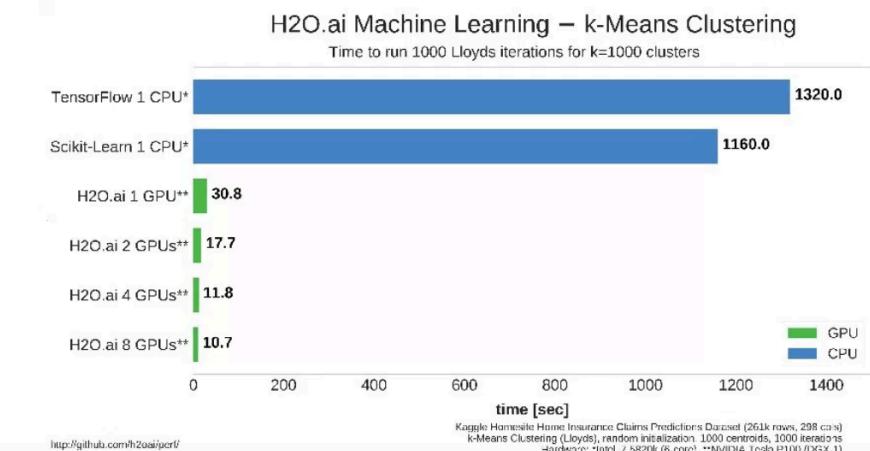
Today, select algorithms are GPU-enabled. These include Gradient Boosting Machines (GBM's), Generalized Linear Models (GLM's), and K-Means Clustering. Using H2O4GPU, users can unlock the power of GPU's through the scikit-learn API that many already use today. In addition to the scikit-learn Python API, an R API is in development.

<https://blog.h2o.ai/tag/h2o4gpu/>



#### k-Means Clustering

- Based on NVIDIA prototype of k-Means algorithm in CUDA
- Improvements to original implementation:
  - Significantly faster than scikit-learn implementation (50x) and other GPU implementations (5-10x)
  - Supports multiple GPUs



**H<sub>2</sub>O.ai**

# Driverless AI

H2O.ai and NVIDIA Bring Fast, Accurate and Interpretable Driverless AI with Automated Machine Learning and Feature Engineering

**Businesses Can Leapfrog Data Scientist Shortage and Accelerate Adoption of AI with H2O.ai Driverless AI on NVIDIA DGX Systems**

September 26, 2017 04:57 PM Eastern Daylight Time

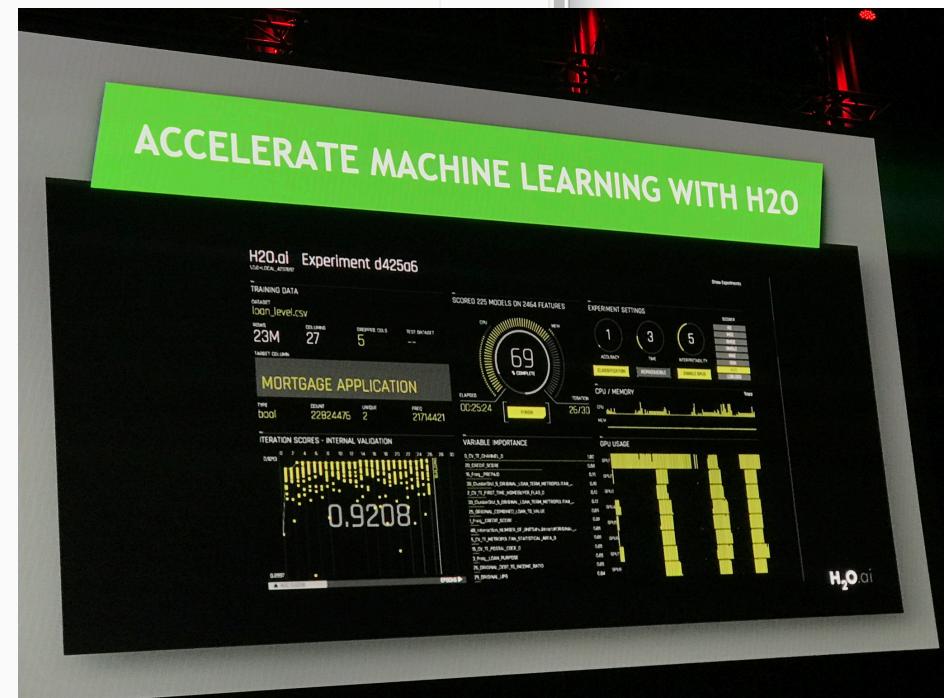
NEW YORK--(BUSINESS WIRE)--Strata Data Conference — H2O.ai today announced an offering built on NVIDIA DGX Systems to further democratize machine learning and address the growing demands placed on a limited number of trained data scientists.

"Interpretable AI builds trust in AI and automates report generation for regulatory purposes. Driverless AI on GPU Computing Platforms accelerates learning and data products with AI in enterprises."

[Tweet this](#)

Engineering and quickly develop hundreds of machine learning models to help businesses mitigate risks and maximize revenue

[View All](#)



<http://www.businesswire.com/news/home/20170926006769/en/H2O.ai-NVIDIA-Bring-Fast-Accurate-Interpretable-Driverless>

# Thank you!

- Prof. Dragan Savić
- Prof. Zoran Kapelan
- Code, Slides & Documents
  - [bit.ly/h2o\\_meetups](https://bit.ly/h2o_meetups)
  - [docs.h2o.ai](https://docs.h2o.ai)
- Contact
  - [joe@h2o.ai](mailto:joe@h2o.ai)
  - [@matlabulous](https://twitter.com/matlabulous)
  - [github.com/woobe](https://github.com/woobe)
- Please search/ask questions on  
**Stack Overflow**
  - Use the tag `h2o` (not H2 zero)