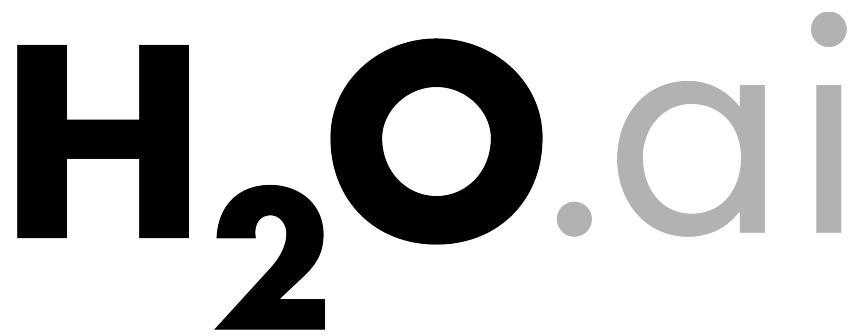


Introduction to Scalable Machine Learning with H₂O



Jo-fai (Joe) Chow

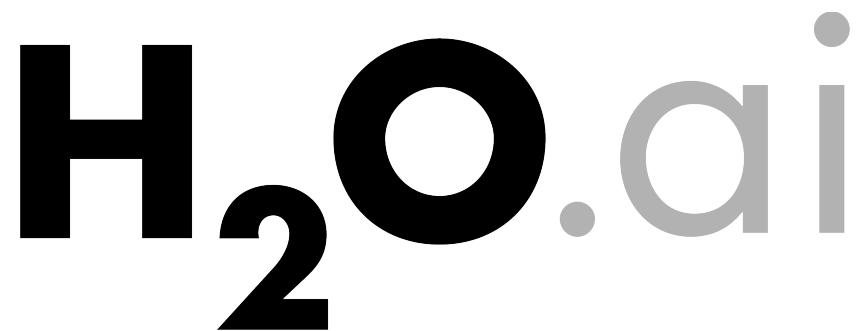
Data Scientist

joe@h2o.ai

@matlabulous

Scotland Data Science and Technology Meetup
30th November, 2017

Introduction to Scalable Machine Learning with H₂O



Jo-fai (Joe) Chow

Data Scientist

joe@h2o.ai

@matlabulous

All slides, data and code examples

http://bit.ly/h2o_meetups

Agenda

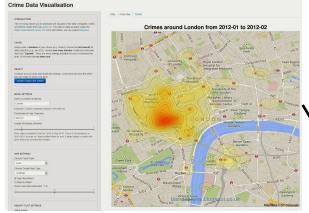
- Overview
 - Company / People
- Open Source H₂O Machine Learning Platform
 - Overview
 - Demos
 - H₂O on Hadoop
 - Automatic Machine Learning
- Q&A



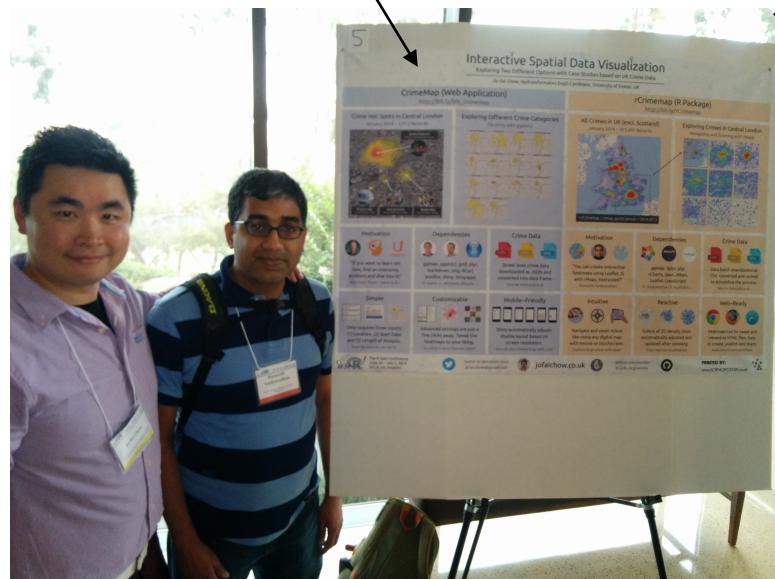
About Me

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

useR! 2014 Conference in US



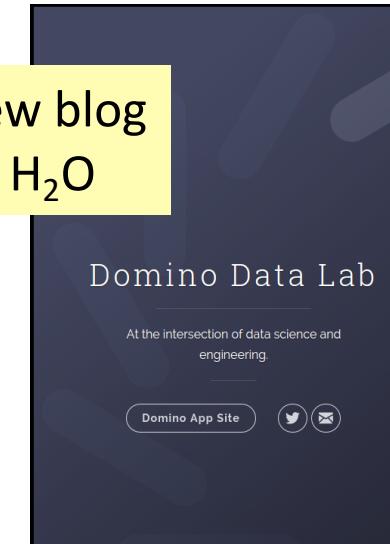
CrimeMap -> poster



John Chambers mentioned H₂O

... so I gave it a try

and wrote a few blog posts about H₂O



19 Sep 2014 · 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.

The Long Story:
bit.ly/joe_kaggle_story

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 (H₂O + Spark)• Enterprise Steam• Driverless AI
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">• Distributed Systems Engineers doing Machine Learning• World-class visualization designers
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*



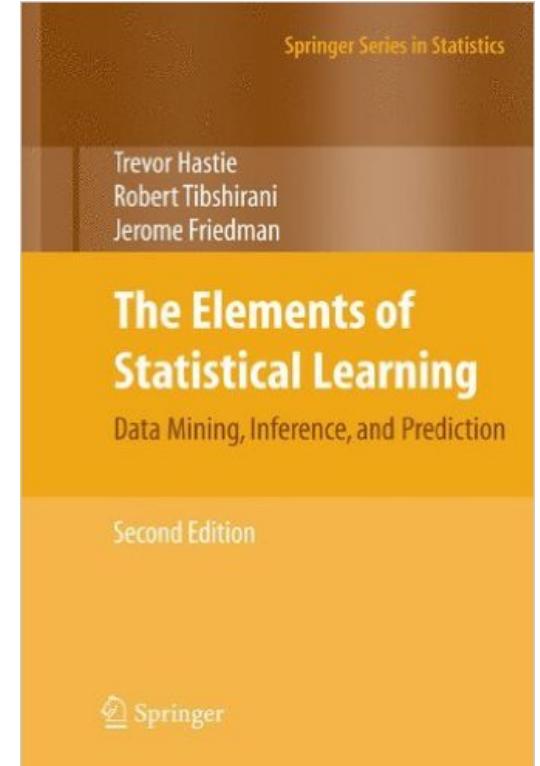
Dr. Robert Tibshirani

- 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*

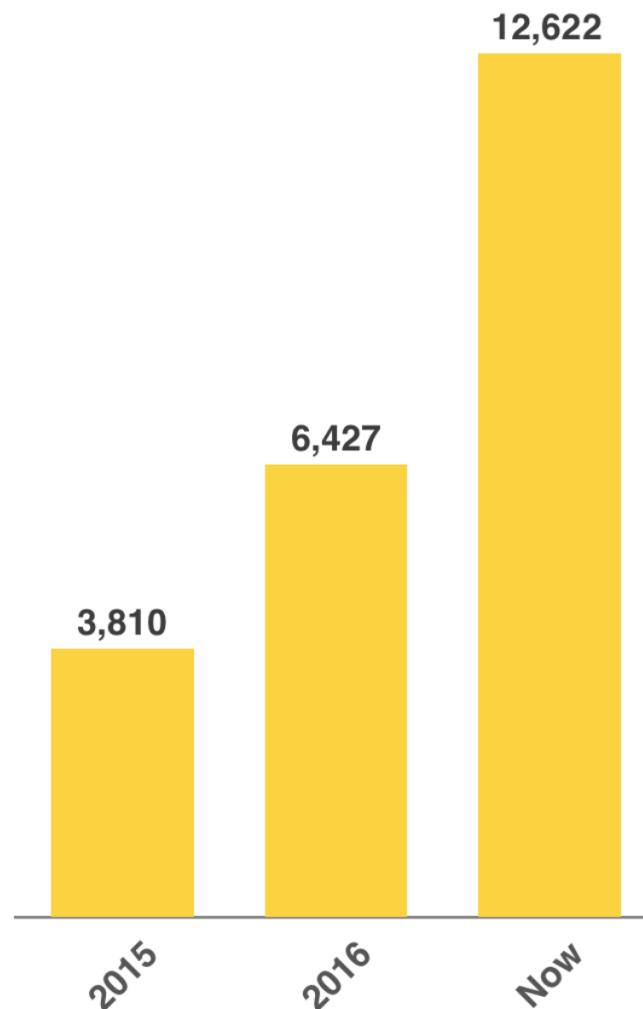


Dr. Steven Boyd

- 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*



Companies Using H2O.ai



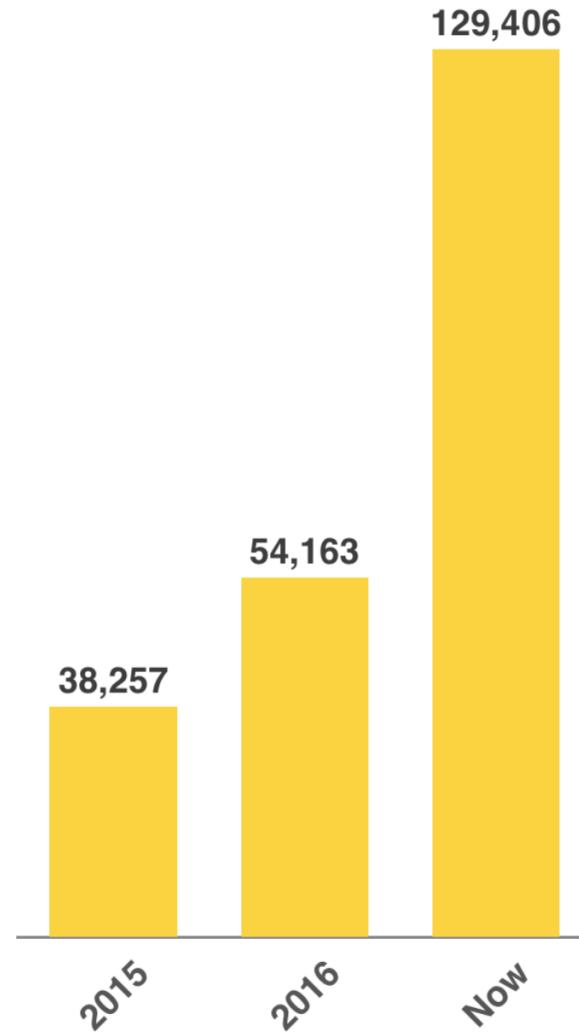
222 OF **FORTUNE
500**
 **H₂O**

**8 OF TOP 10
BANKS**

**7 OF TOP 10
INSURANCE COMPANIES**

**4 OF TOP 10
HEALTHCARE COMPANIES**

H2O.ai Users



Community Expansion

15 Meetups a month



66,843
members 33
interested 50
Meetups 45
cities 18
countries

Find out more: www.h2o.ai/community/

Select 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 made a significant investment in pioneering artificial intelligence (AI) for the audit. For the past 18 months, we have partnered with H₂O.ai – a leading Silicon Valley company – to build a revolutionary bot that uses AI and machine learning to 'x-ray' a business, analysing billions of data points in milliseconds, seeing what humans can't, and applying judgement to detect anomalies in the general ledger. Called GL.ai, it is the first module of PwC's Audit.ai.

GL.ai was named the 'Audit Innovation of the Year' by the International Accounting Bulletin in October 2017.



"The reason this is such a brilliant tool is its 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 to decide what's unusual. With GL.ai, the algorithms do it for us."

Laura Needham partner, PwC UK



Exciting night at this year's @WAI_News Awards: PwC wins 2017 Audit Innovation of the Year! pwc.to/Glai17 #taandiab17



10:15 PM - 4 Oct 2017

<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>

Partners



NVIDIA's (NASDAQ: NVDA) invention of the GPU in 1999 sparked the growth of the PC gaming market, redefined modern computer graphics, and revolutionized parallel computing. More recently, GPU deep learning ignited modern AI -- the next era of computing -- with the GPU acting as the brain of computers, robots, and self-driving cars that can perceive and understand the world.

[Website](#)



Founded in 1975, Microsoft (Nasdaq "MSFT") is the worldwide leader in software, services, devices and solutions that help people and businesses realize their full potential. You can launch Sparkling Water on Azure HDInsights with just a few clicks to build your data science pipeline on the cloud.

[Website](#) | [Documentation](#)



Watson Data Platform enables AI-powered decision making by simplifying and automating the development and operationalization of new insights. It enables unprecedented levels of collaboration amongst data savvy professionals.



Cloudera delivers a modern platform for data management and analytics, helping businesses solve their most challenging problems with data.



Anaconda is the leading open data science platform powered by Python. Continuum Analytics is the creator and driving force behind Anaconda. We put superpowers into the hands of people who are changing the world.



Databricks provides a just-in-time data platform, to simplify data integration, real-time experimentation, and robust deployment of production applications.



Hortonworks drives actionable intelligence with Connected Data Platforms that maximize the value of all data—data-in-motion and data-at-rest.



Kensu's mission is to lift Data Science to the Enterprise level focusing on the production environment and the maintenance across time.

[Website](#)



SigOpt is the optimization platform that amplifies your research. SigOpt takes any research pipeline and tunes it, right in place.

[Website](#)



Nimbix is the leading provider of purpose-built cloud computing for machine learning, AI and HPC applications. Powered by JARVICE™, the Nimbix Cloud provides high-performance software as a service, dramatically speeding up data processing for Energy, Life Sciences, Manufacturing, Media and Analytics applications.

[Website](#)



DataScience.com pairs data expertise with powerful tools to help businesses unlock the value in their data. Enabling data science for every business.

[Website](#)



MapD makes queries faster and creates a fluid and immersive data exploration experience that removes the disconnect between an analyst and their data. Making extracting insight from data effortless and lightning fast.

[Website](#)



The MapR Converged Data Platform integrates enormous power of Hadoop and Spark with global event streaming, real-time database capabilities, and enterprise storage.

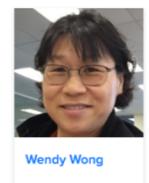
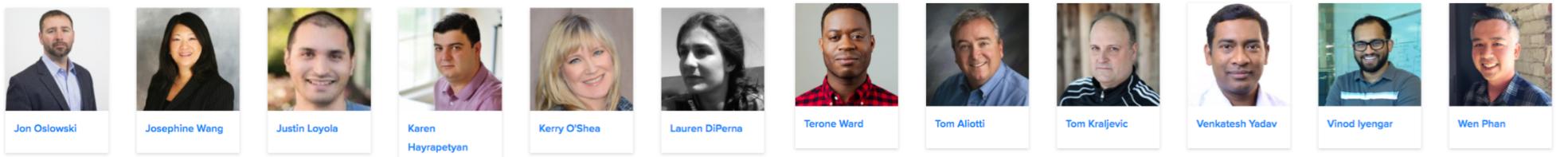
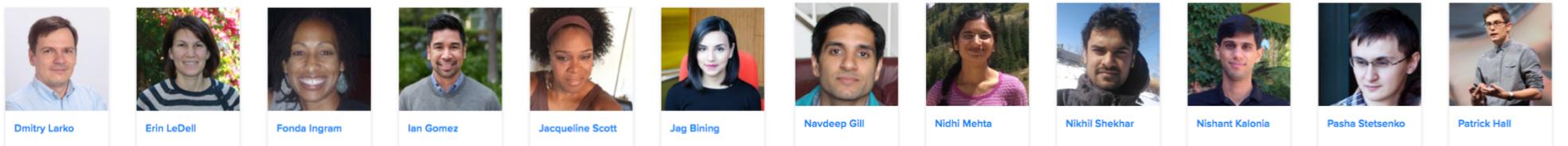
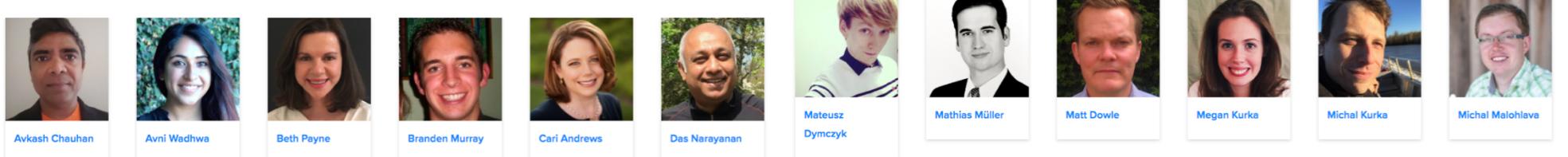


Splunk Inc. provides the leading platform for Operational Intelligence. Customers use Splunk to search, monitor, analyze and visualize machine data.

[Website](#)



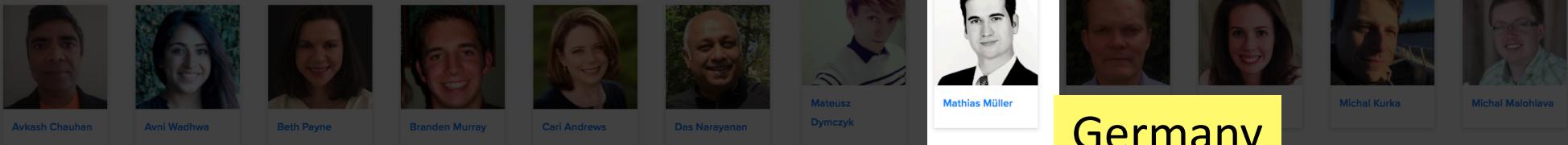
Minio is an object storage server built for cloud application developers and devops.



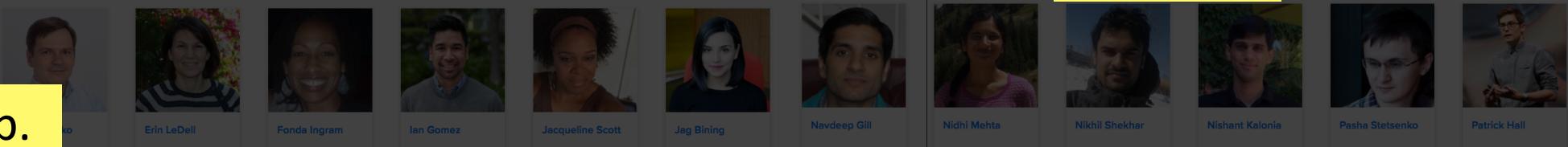
H₂O Team



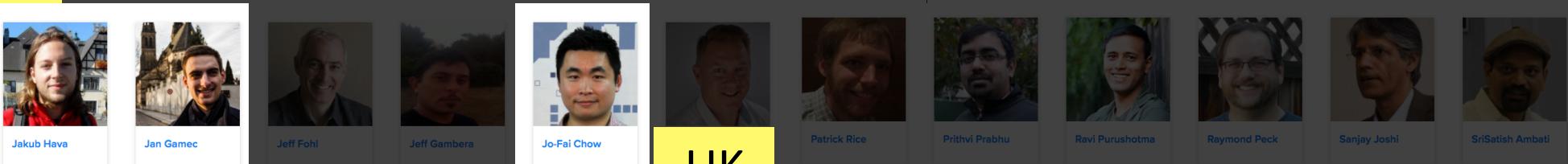
UK



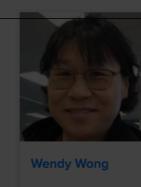
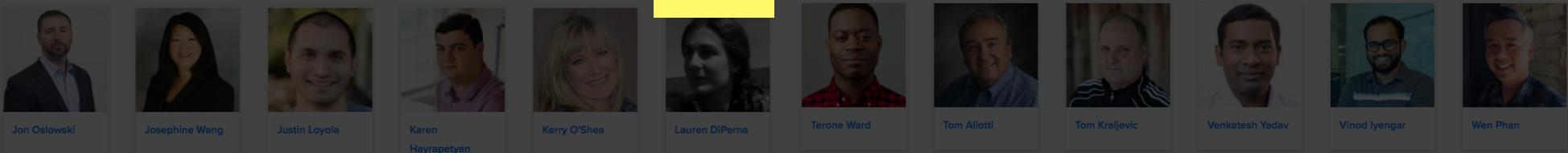
Germany



Czech Rep.



UK



Wendy Wong

H₂O Team in Europe

What is Joe's role at H₂O.ai?



- Data Scientist / Sales Engineer / Speaker / Meetup Organiser / Community Growth Hacker (on paper)
- Unofficial Photographer of H₂O.ai SWAG (the travelling data scientist)
- H₂O.ai SWAG EMEA Distributor (please help yourself)

#AroundTheWorldWithH2Oai

Month Joe's H₂O Events in 2017

Jan	London
Feb	London, Warsaw, Oxford
Mar	Bay Area, London, Cologne, Barcelona, Madrid, Vienna
Apr	Amsterdam, Rotterdam, Poznan, London
May	Belgrade, Hamburg, Berlin
Jun	Amsterdam, Stockholm, Budapest, London, Munich, Prague
Jul	Berlin, Brussels
Aug	☀️
Sep	London, Dublin
Oct	Exeter, Munich, Dublin, The Hague, Amsterdam, Frankfurt
Nov	Munich, Zurich, London, Glasgow (TODAY)
Dec	Bay Area (Next Week), London, 🎄

25+ Cities in 2017



Jo-fai (Joe) Chow
@matlabulous

Minor #stickers upgrade after a few conferences. Ready for @h2oai first #meetup in #Glasgow and then #H2OWorld in US #letsgo #AroundTheWorldWithH2Oai ✈️🇺🇸



12:09 PM - 29 Nov 2017 from Heathrow Airport



London Artificial Intelligence & Deep Learning

PRO

H2O Artificial Intelligence and Machine Learning -
39 groups

Location

London, United Kingdom

Members

4,354

4000+ Members

**Organizers**

Ian Gomez and 2 others

[Schedule](#)

...

[Our group](#)[Meetups](#)[Members](#)[Photos](#)[Discussions](#)[More](#)**Next Meetup****Next London Meetup: 12 Dec**[See all](#)**12
DEC**

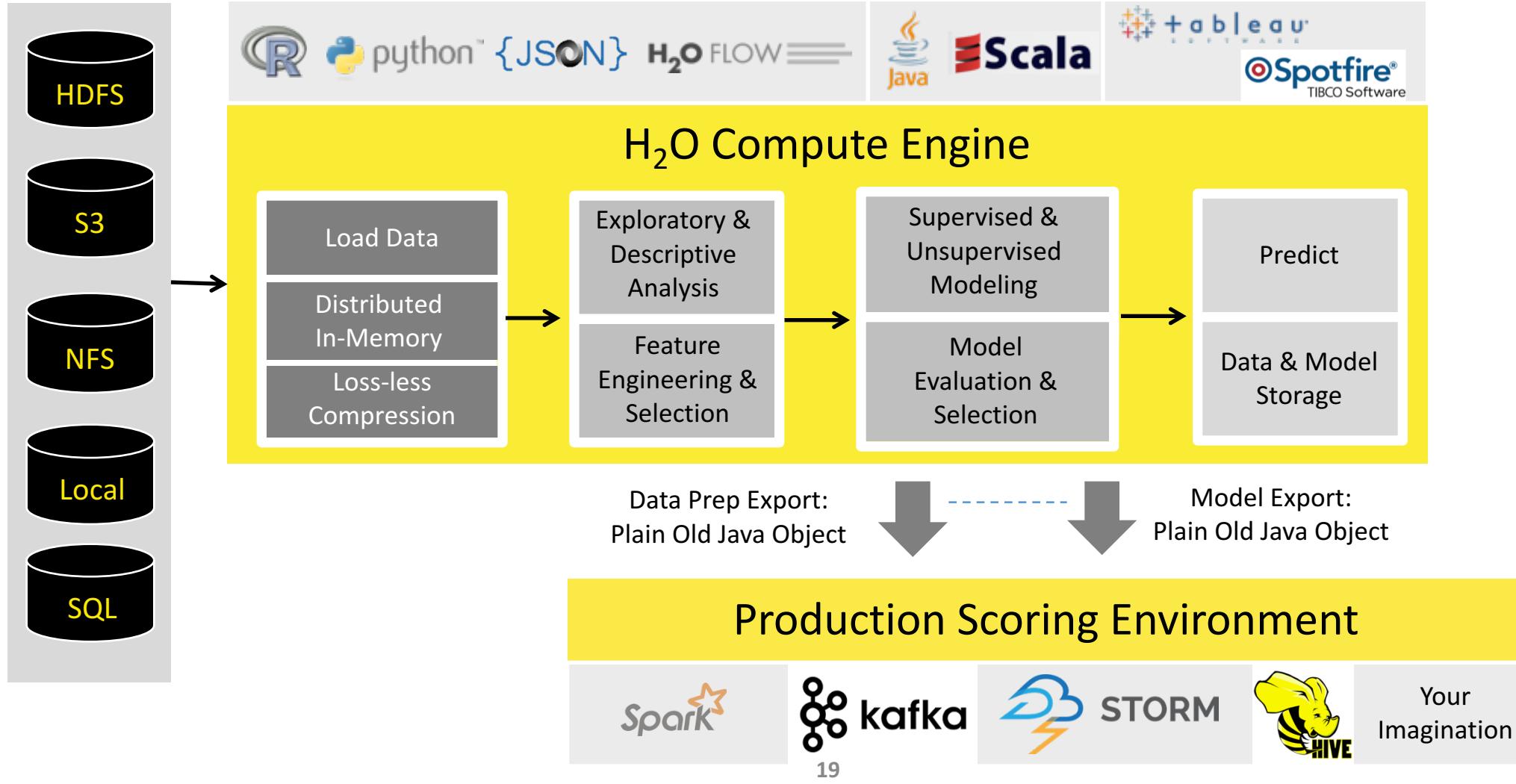
Tuesday, December 12, 2017, 6:00 PM

Interpretable Machine Learning, Tweet Classifier, H2O World Highlights and MoreHosted by [Jo-fai Chow](#)[Edit](#)**Moody's Analytics**

1 Canada Square, Canary Wharf, E14 5AB · London

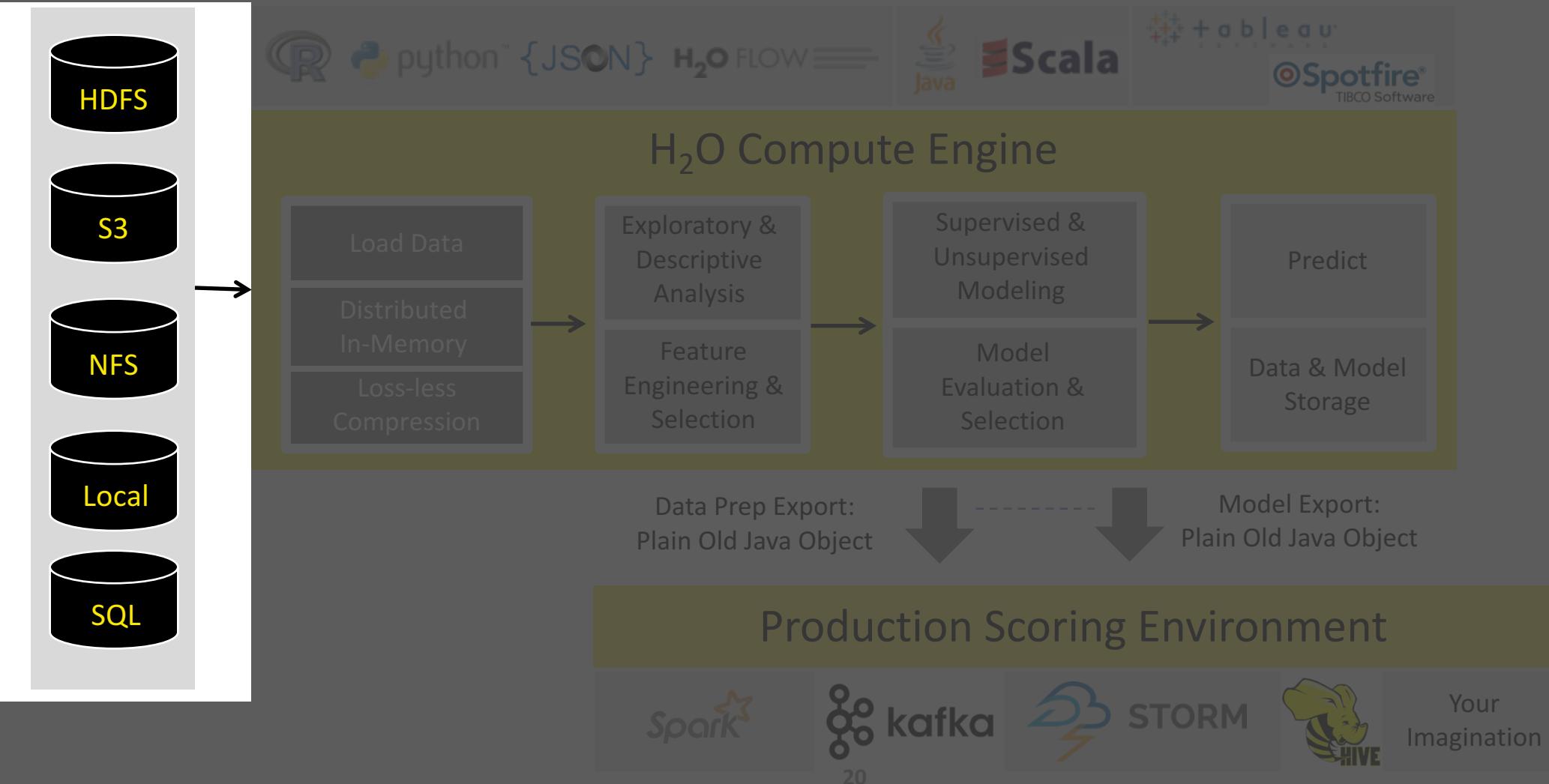
H₂O Machine Learning Platform

High Level Architecture

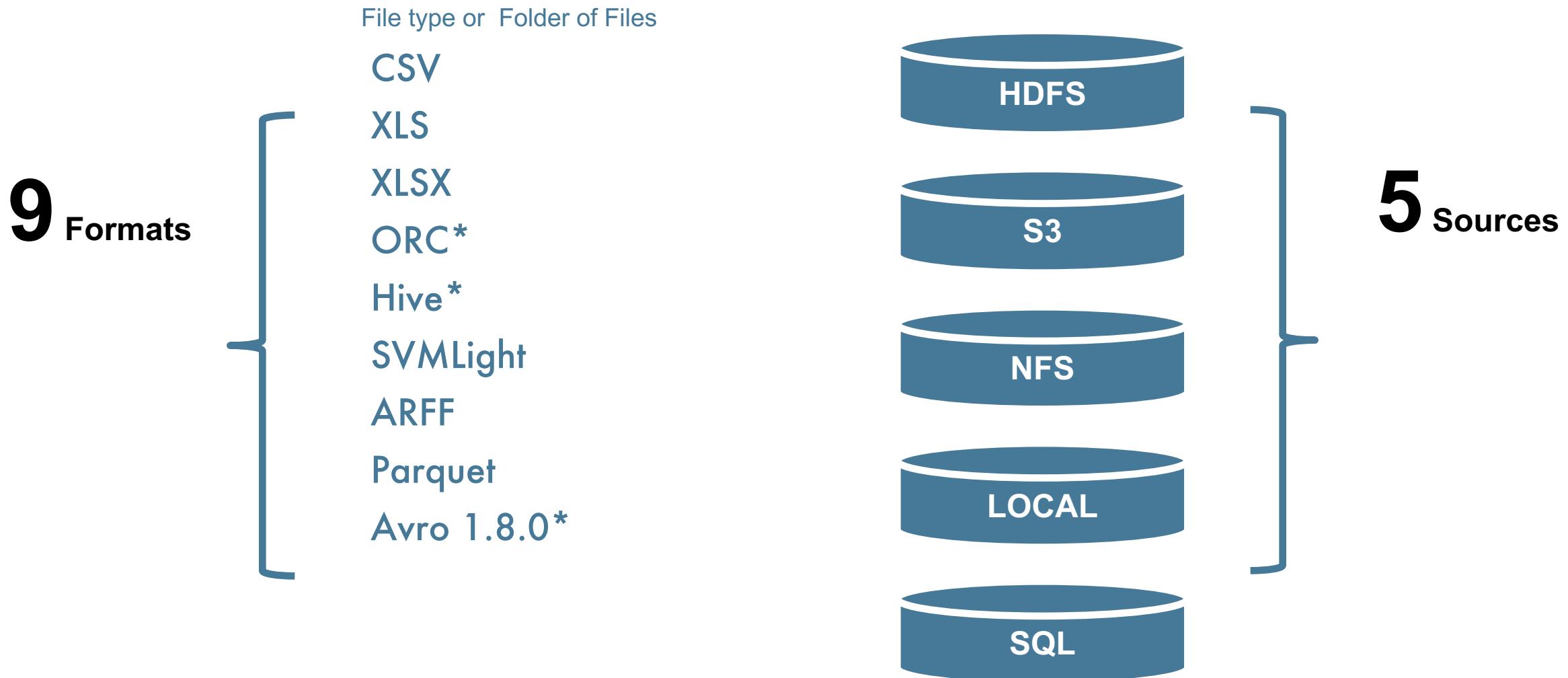


High Level Architecture

Import Data from
Multiple Sources



Supported Formats & Data Sources



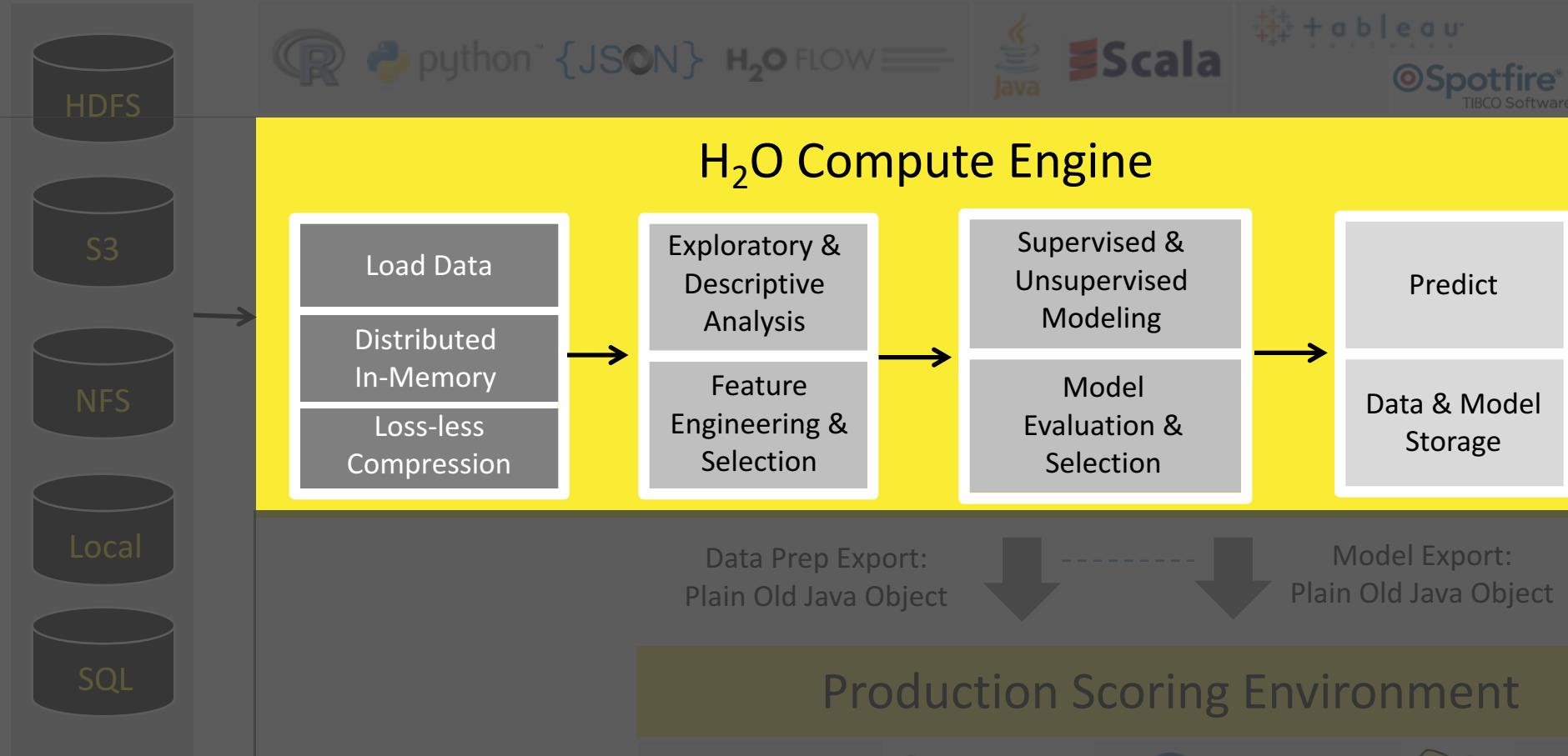
* 1. only if H2O is running as a Hadoop job

* 2. Hive files that are saved in ORC format

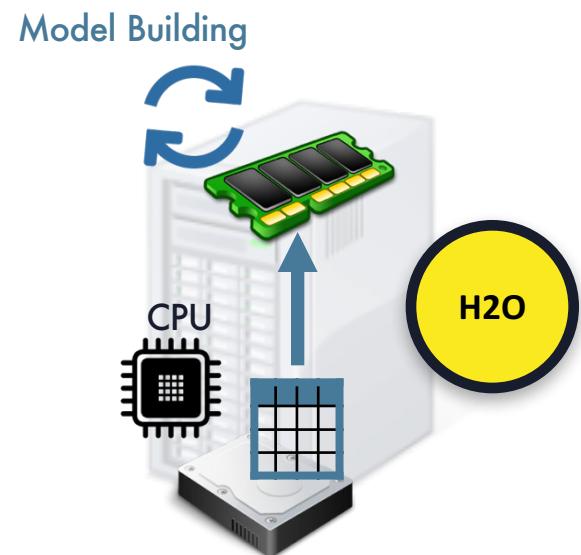
* 3. without multi-file parsing or column type modification

High Level Architecture

Fast, Scalable & Distributed
Compute Engine Written in
Java



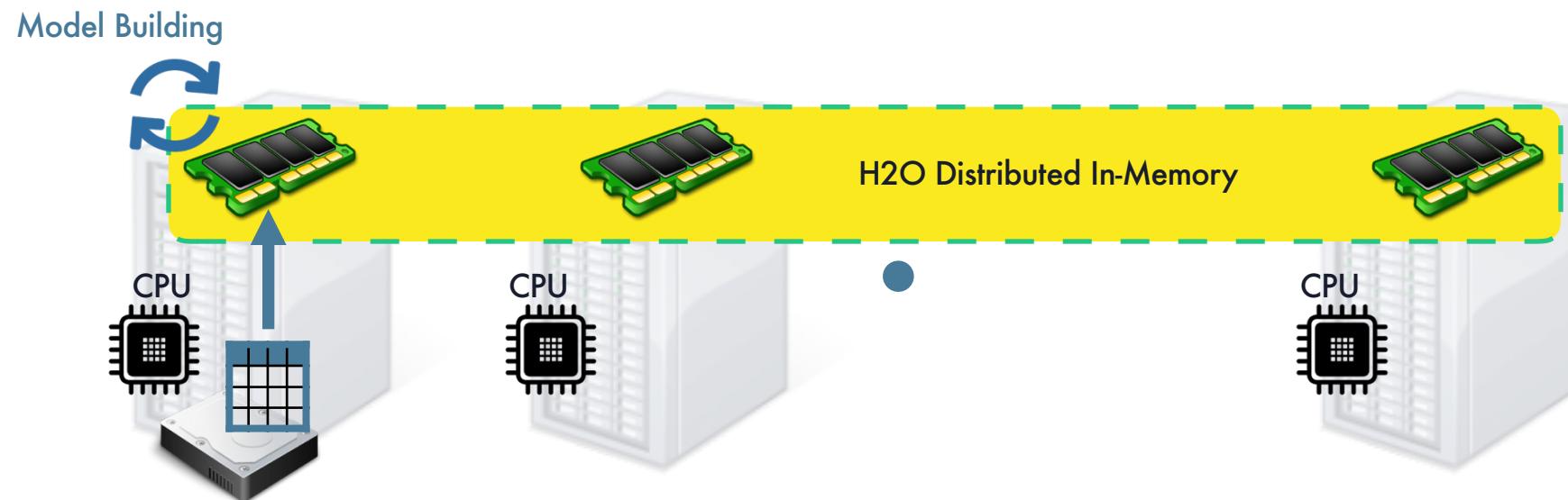
H₂O Core



H₂O Core

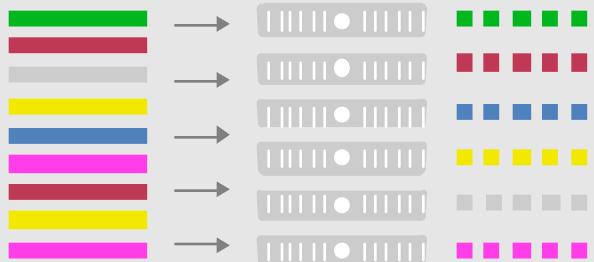


H₂O Core



Distributed Algorithms

Foundation for Distributed Algorithms



Parallel Parse into **Distributed Rows**



Fine Grain Map Reduce Illustration: Scalable
Distributed Histogram Calculation for GBM

Advantageous Foundation

- Foundation for In-Memory Distributed Algorithm Calculation - **Distributed Data Frames** and **columnar compression**
- All algorithms are distributed in H₂O: GBM, GLM, DRF, Deep Learning and more. Fine-grained map-reduce iterations.
- **Only enterprise-grade, open-source distributed algorithms in the market**

User Benefits

- “Out-of-box” functionalities for all algorithms (**NO MORE SCRIPTING**) and uniform interface across all languages: R, Python, Java
- **Designed for all sizes of data sets, especially large data**
- **Highly optimized Java code for model exports**
- **In-house expertise for all algorithms**

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

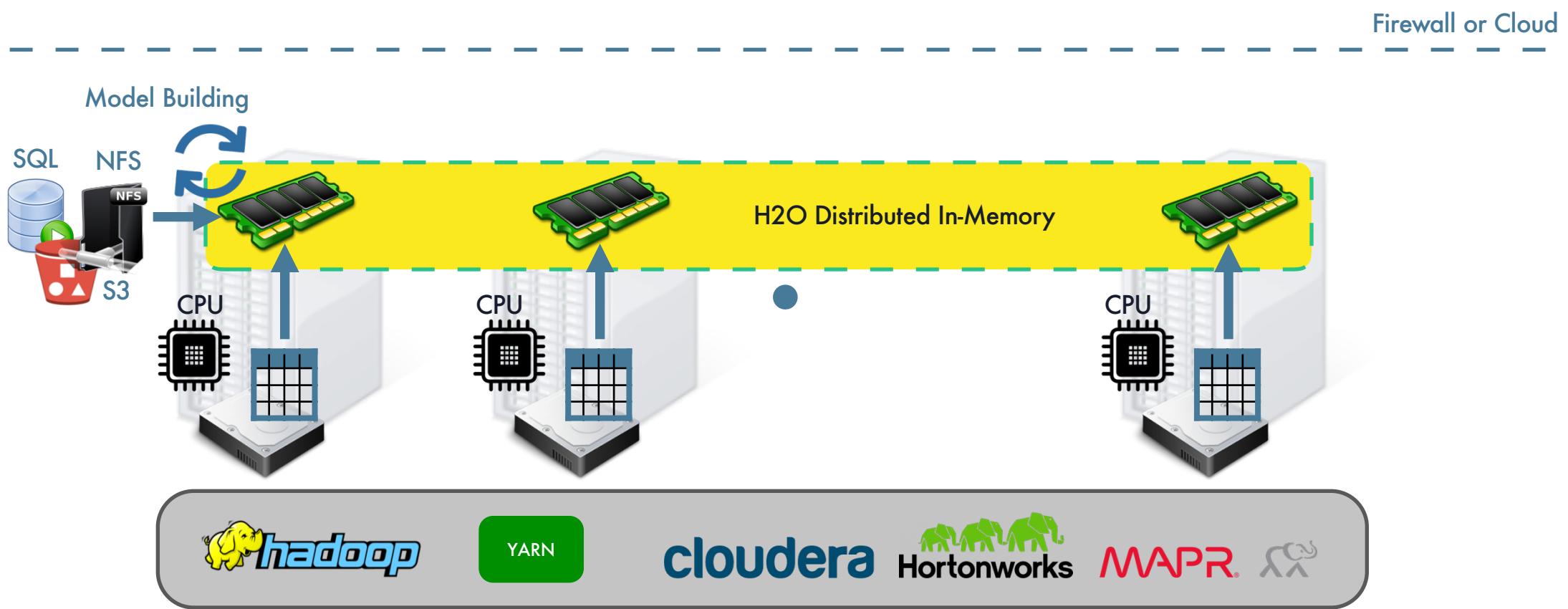
Dimensionality Reduction

- **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

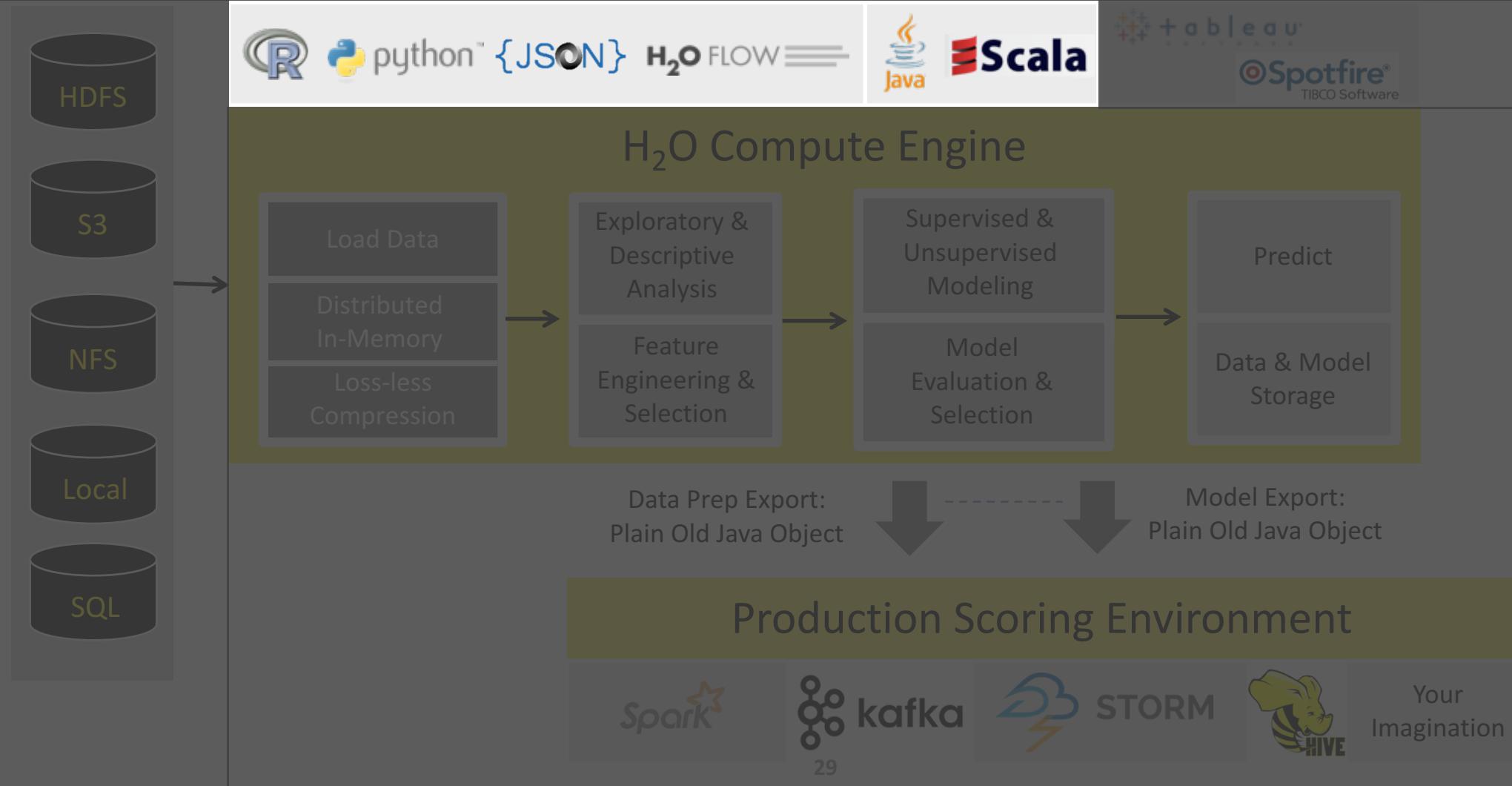
Anomaly Detection

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

H₂O Core



High Level Architecture



H₂O Flow (Web) – Today's Demos

The screenshot shows the H2O Flow (Web) interface running in a browser window. 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 context menu is open under the "Model" dropdown, listing various modeling routines: 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..., and Run AutoML... . To the left, there's a sidebar titled "Assistance" with a table of routines and their descriptions. The main workspace shows a single step named "assist". On the right, there's a "Help" panel with sections for "Using Flow for the first time?", "Quickstart Videos", "View example Flows", "Star" button (2,387), "GENERAL" section with links to Flow Web UI, Importing Data, Building Models, Making Predictions, Using Flows, and Troubleshooting Flow, and an "EXAMPLES" section encouraging users to try out flows.

Using H₂O with R and Python

The image shows two side-by-side screenshots illustrating the use of H₂O with R and Python.

Left Screenshot (RStudio): A screenshot of the RStudio Source Editor window titled "credit_card_example.R". The code is an R script for a credit card example, demonstrating the import of datasets from S3, training a GBM model, and using AutoML. The code includes imports for h2o and h2o.automl, and various training parameters like training_frame, seed, and max_runtime_secs.

```
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
```

Right Screenshot (Jupyter Notebook): A screenshot of a Jupyter notebook titled "credit_card_example". The notebook shows the execution of an R script. In cell [2], the H₂O cluster is started, and its status is displayed in a table:

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

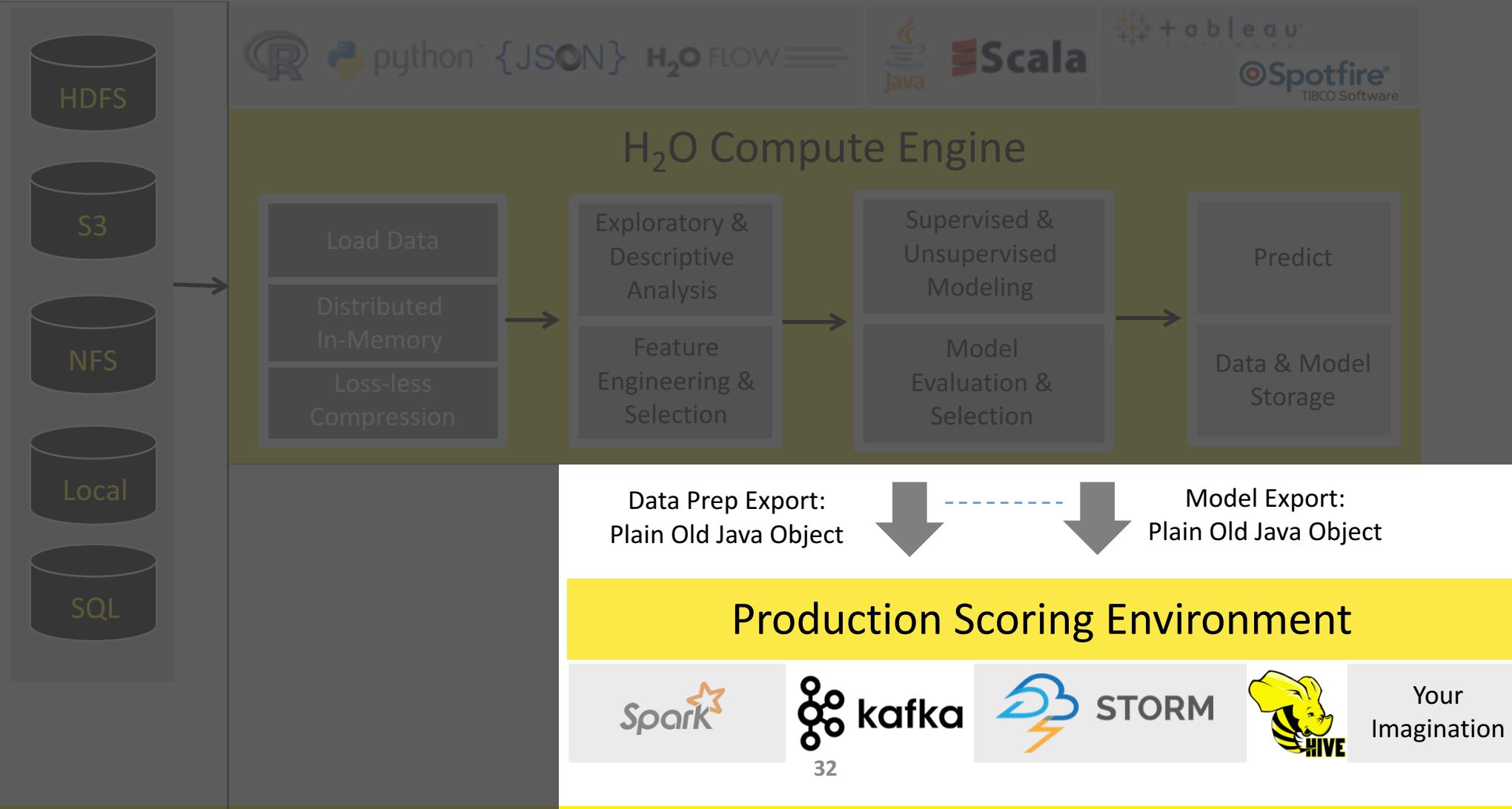
In cell [3], datasets are imported from S3, and the progress bar shows 100% completion for both operations.

In cell [4], the datasets are summarized, and the resulting table is shown:

	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



My Training Materials

- H₂O + Python
 - Workshop at PyData Conferences
 - bit.ly/joe_h2o_tutorials
- H₂O + R
 - github.com/woobe/h2o_training_2017_10

The screenshot shows a video player interface for a PyData Berlin 2017 talk. The video frame displays a man speaking at a podium with a laptop. The title bar of the video player reads "PyData Berlin 2017". Below the video frame, the title of the talk is "Jo-fai Chow - Introduction to Machine Learning with H2O and Python". The video has 665 views and was published on 26 Jul 2017. The video content itself shows a Jupyter Notebook cell for "GBM with CV, Early Stopping and Random Grid Search". The code in the cell includes:

```
In [26]: # define the criteria for random grid search
search_criteria = {"strategy": "randomdiscrete",
                   "max_models": 9,
                   "seed": 1234}

In [27]: # define the range of hyper-parameters for grid search
hyper_params = {"col_sample_rate": [0.7, 0.8, 0.9],
                 "col_sample_rate": [0.7, 0.8, 0.9],
                 "max_depth": [3, 5, 7]}

In [28]: # Set up GBM grid search
gbm_rand_grid = H2OGridSearch(
    H2OGradientBoostingEstimator(
        model_id = "gbm_rand_grid",
        seed = 1234,
        max_models = 9,
        nfolds = 5,
        stopping_metric = "mse",
        stopping_rounds = 15,
        score_tree_interval = 1),
    search_criteria = search_criteria, # full grid search
    hyper_params = hyper_params)

In [29]: # Use .train() to start the grid search
gbm_rand_grid.train(x = features,
                     y = target,
                     training_frame = wine_train)
```

A yellow callout box highlights the text "Only search for 9 combinations". Another yellow callout box highlights the text "Expand Search Space". The video player also shows the "Founding Sponsor" logo for ANACONDA.

H₂O Documentation

Getting Started & User Guides | Q & A | Algorithms | Languages | Tutorials, Examples, & Presentations | API & Developer Docs | For the Enterprise

Getting Started & User Guides

 Open Source |  Commercial

H₂O

What is H₂O?
[H₂O User Guide](#) (Main docs)
H₂O 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 H₂O](#)

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](#)

Driverless AI

What is Driverless AI?
Driverless AI User Guide [HTML](#) [PDF](#)
Driverless AI Booklet
MLI with Driverless AI Booklet

Driverless AI Webinars

[Download Driverless AI](#)

H₂O4GPU (alpha)

[H₂O4GPU Readme](#)
[Open Source License \(Apache V2\)](#)

[Download H₂O4GPU](#)

URL: docs.h2o.ai

Demo: Running H₂O on Hadoop



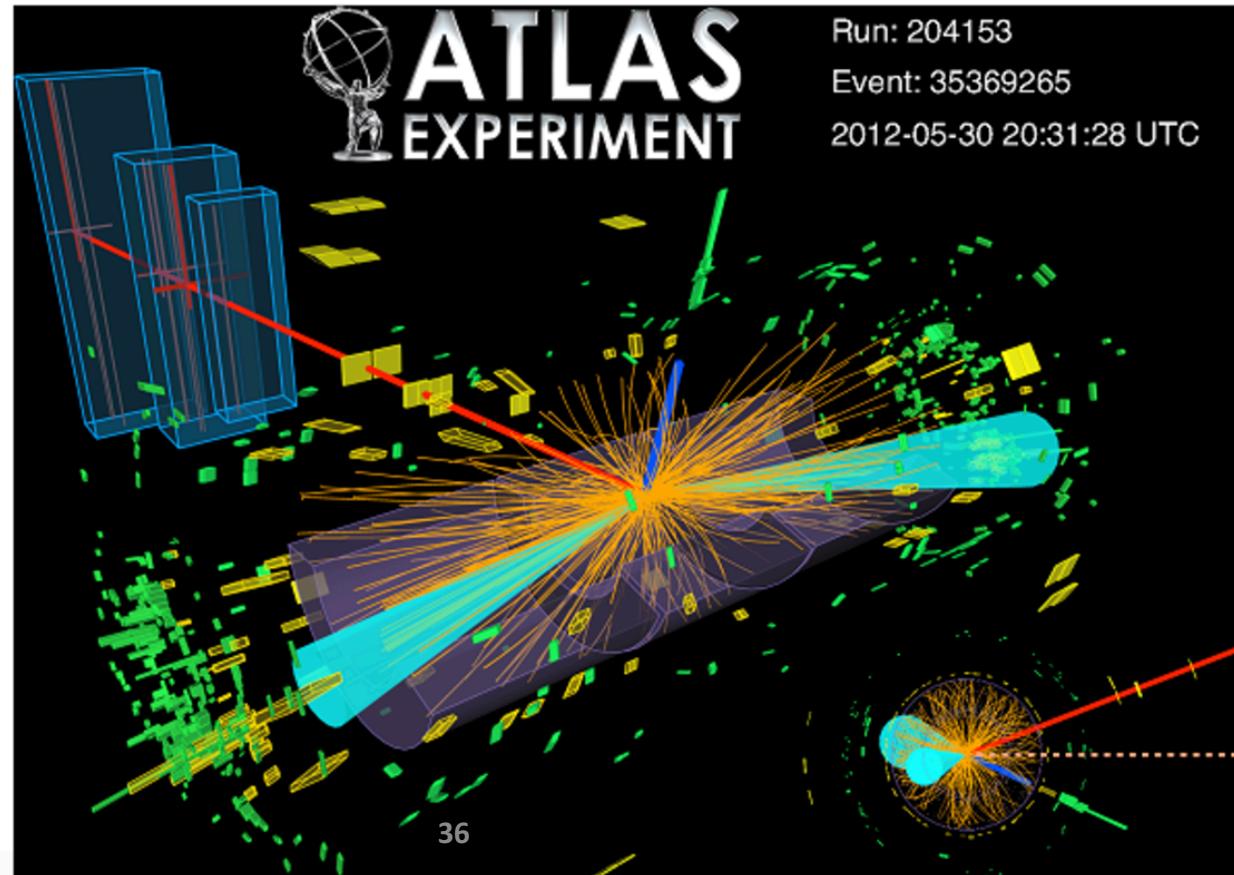
Higgs Boson Machine Learning Challenge

Use the ATLAS experiment to identify the Higgs boson

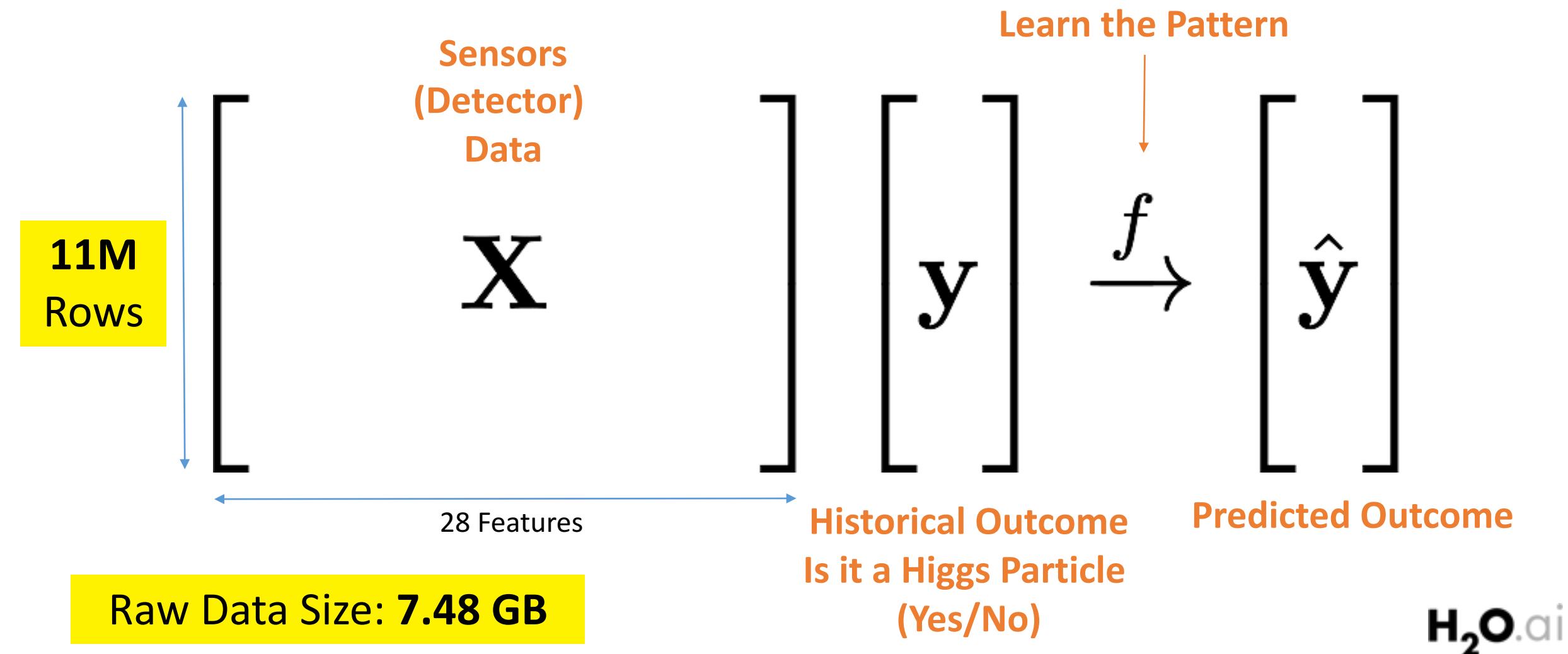
\$13,000 · 1,785 teams · 3 years ago

[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 Machine Data



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

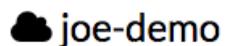
label	type	Missing	Zeros	+Inf	-Inf	min	max	mean	sigma	cardinality	Actions
C1	enum	0	5170877	0	0	0	1.0	0.5299	0.4991	2	Convert to numeric
C2	real	0	0	0	0	0.2747	12.0989	0.9915	0.5654	· ·	
C3	real	0	0	0	0	-2.4350	2.4349	-0.0	1.0088	· ·	
C4	real	0	0	0	0	-1.7425	1.7432	-0.0	1.0063	· ·	
C5	real	0	0	0	0	0.0002	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



joe-demo

10 nodes

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%

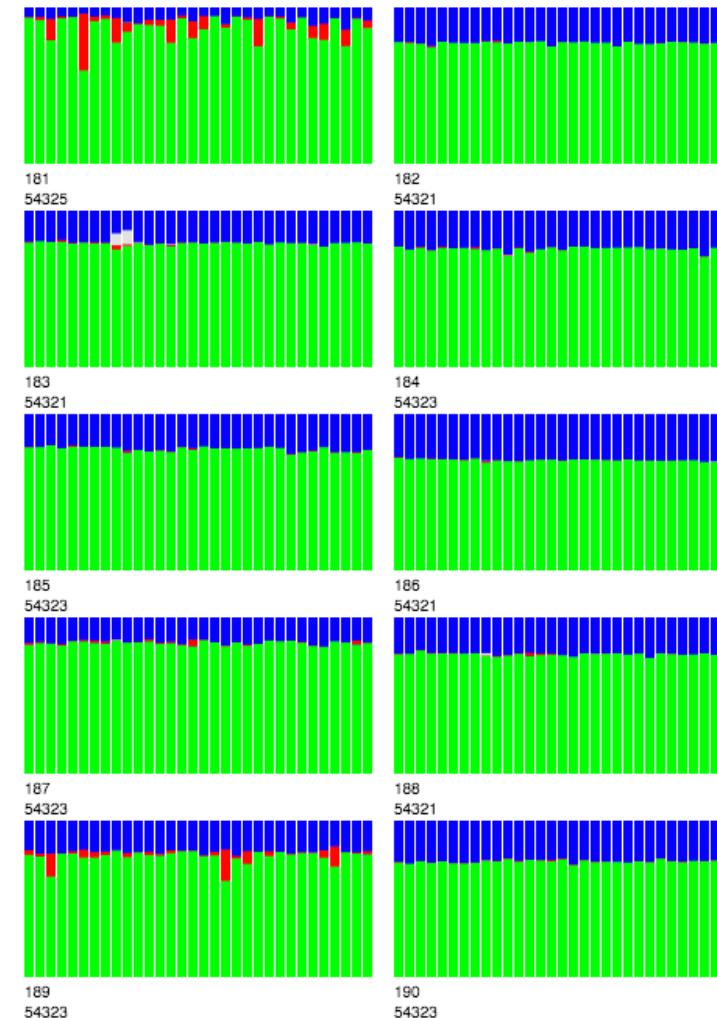
$$10 \times 32 = \\ 320 \text{ Cores}$$

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

H₂O.ai

H₂O Water Meter (CPU Monitor)

10 x 32 = 320 Cores



Legend

Each bar represents one CPU.

Blue: idle time

Green: user time

Red: system time

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

Demo: AutoML

Automatic Machine Learning with H₂O

AutoML: Automatic Machine Learning

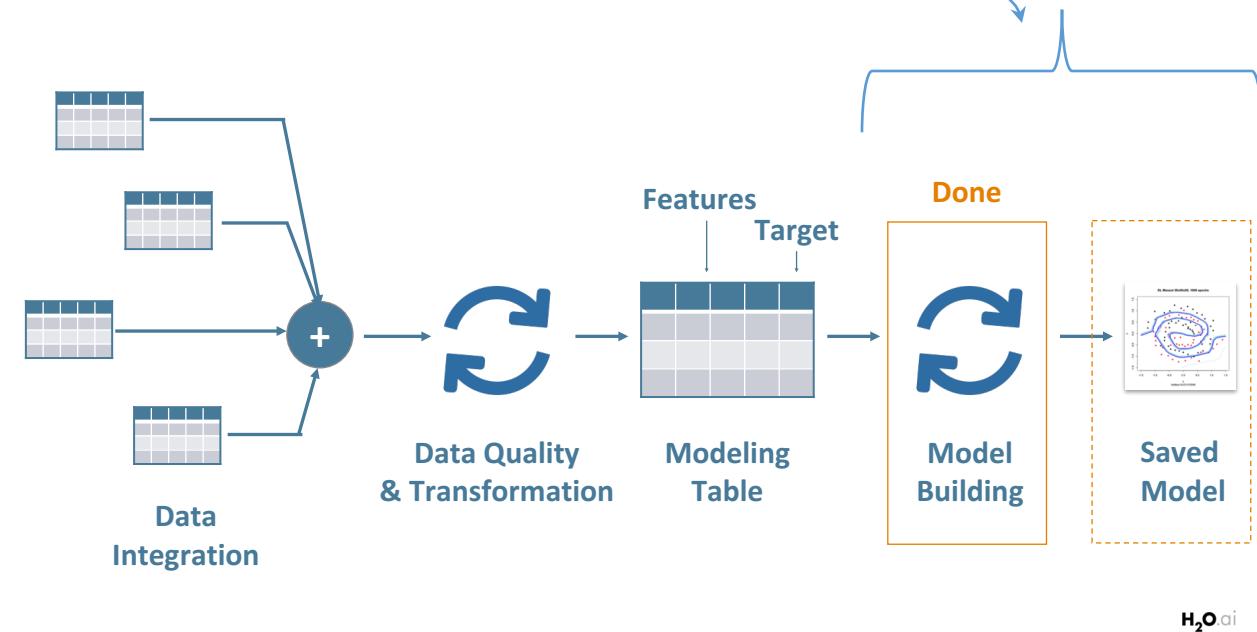
In recent years, the demand for machine learning experts has outpaced the supply, despite the surge of people entering the field. To address this gap, there have been big strides in the development of user-friendly machine learning software that can be used by non-experts. The first steps toward simplifying machine learning involved developing simple, unified interfaces to a variety of machine learning algorithms (e.g. H2O).

Although H2O has made it easy for non-experts to experiment with machine learning, there is still a fair bit of knowledge and background in data science that is required to produce high-performing machine learning models. Deep Neural Networks in particular are notoriously difficult for a non-expert to tune properly. In order for machine learning software to truly be accessible to non-experts, we have designed an easy-to-use interface which automates the process of training a large selection of candidate models. H2O's AutoML can also be a helpful tool for the advanced user, by providing a simple wrapper function that performs a large number of modeling-related tasks that would typically require many lines of code, and by freeing up their time to focus on other aspects of the data science pipeline tasks such as data-preprocessing, feature engineering and model deployment.

H2O's AutoML can be used for automating the machine learning workflow, which includes automatic training and tuning of many models within a user-specified time-limit. The user can also use a performance metric-based stopping criterion for the AutoML process rather than a specific time constraint. [Stacked Ensembles](#) will be automatically trained on the collection individual models to produce a highly predictive ensemble model which, in most cases, will be the top performing model in the AutoML Leaderboard. Stacked ensembles are not yet available for multiclass classification problems, so in that case, only singleton models will be trained.

AutoML Interface

The H2O AutoML interface is designed to have as few parameters as possible so that all the user needs to do is point to their dataset, identify the response column and optionally specify a time constraint, a maximum number of models constraint, and early stopping parameters.



AutoML Output

The AutoML object includes a “leaderboard” of models that were trained in the process, ranked by a default metric based on the problem type (the second column of the leaderboard). In binary classification problems, that metric is AUC, and in multiclass classification problems, the metric is mean per-class error. In regression problems, the default sort metric is deviance. Some additional metrics are also provided, for convenience.

Here is an example leaderboard for a binary classification task:

model_id	auc	logloss
StackedEnsemble_0_AutoML_20170605_212658	0.776164	0.564872
GBM_grid_0_AutoML_20170605_212658_model_2	0.75355	0.587546
DRF_0_AutoML_20170605_212658	0.738885	0.611997
GBM_grid_0_AutoML_20170605_212658_model_0	0.735078	0.630062
GBM_grid_0_AutoML_20170605_212658_model_1	0.730645	0.67458
XRT_0_AutoML_20170605_212658	0.728358	0.629296
GLM_grid_0_AutoML_20170605_212658_model_1	0.685216	0.635137
GLM_grid_0_AutoML_20170605_212658_model_0	0.685216	0.635137

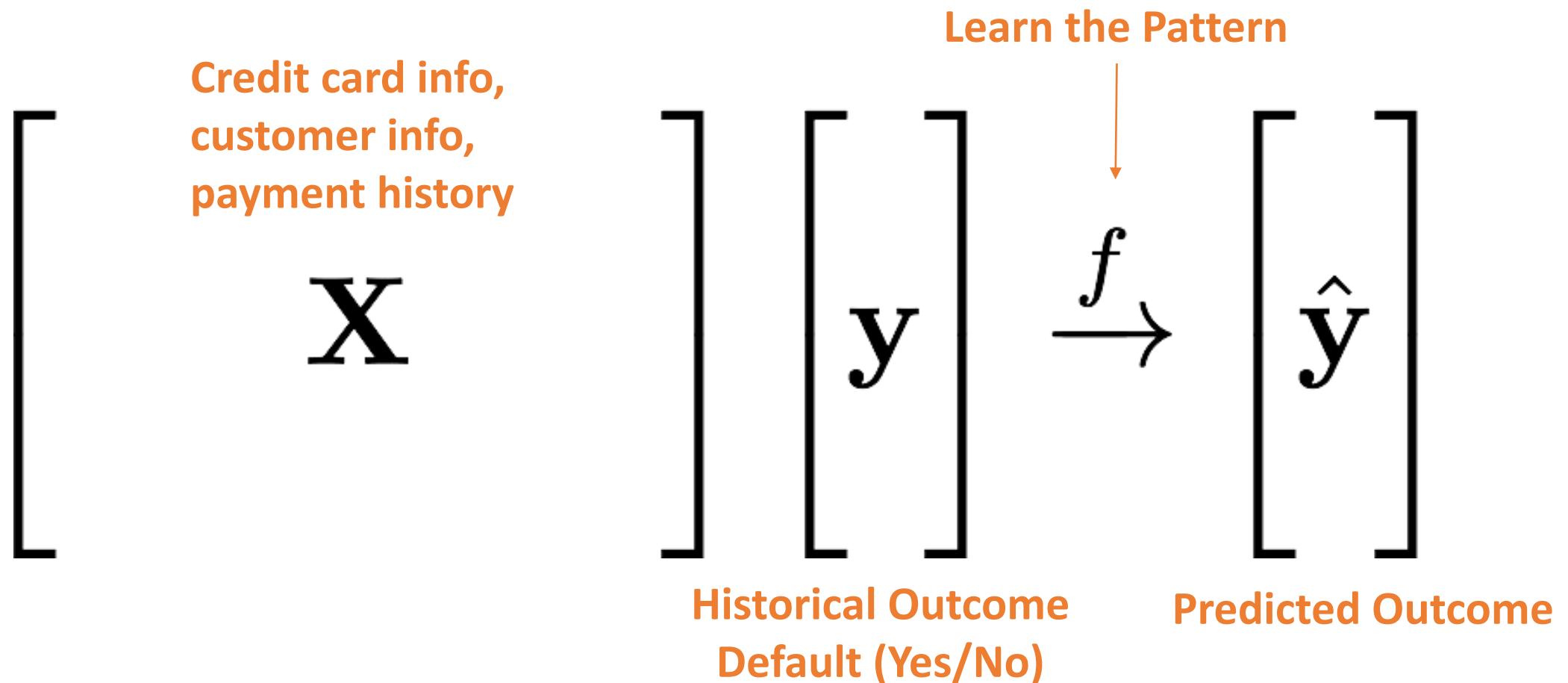
Demo Introduction

/ **Use Case:** Probability of Default for Credit Card Loans

/ **Features**

- **default.payment.next.month**: Did the next loan payment default (1=True, 0=False)
- **LIMIT_BAL**: Credit limit in (NT) dollars
- **SEX, EDUCATION, MARRIAGE, AGE**
- **PAY_0**: Was a payment received in the current month?
- **PAY_2**: Was a payment received in the 2 months ago?
...
- **BILL_AMT1**: Amount of bill statement in 1 month ago
- **BILL_AMT2**: Amount of bill statement in 2 months ago
...
- **PAY_AMT1**: Amount of previous payment 1 month ago
- **PAY_AMT2**: Amount of previous payment 2 months ago
...

Learning from Credit Card Data



Email me
joe@h2o.ai
For VIP tickets



LEARNING IS FUN

REGISTER NOW

Space is limited!

Dec 4 - 5, 2017

Mountain View, CA
Computer History Museum

We'll live-stream the
event (link T.B.C.)

H2O is back with its flagship event, H2O World 2017.

Whether you're just getting started with H2O or you're a power user looking to expand your skill set even more, join

Thank you!

- Organisers & Sponsors
 - Michael Young
 - Scotland Data Science and Technology Meetup



- 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)