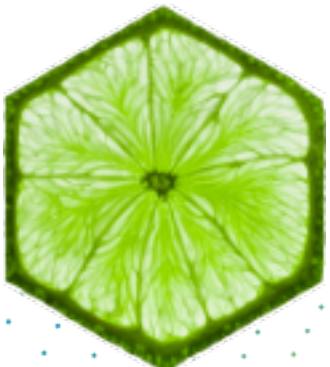


# Automatic and Interpretable Machine Learning in R with H<sub>2</sub>O and LIME



Jo-fai (Joe) Chow

Data Science Evangelist /  
Community Manager

joe@h2o.ai

@matlabulous

Download → [https://bit.ly/  
\*\*joe\\_eRum\\_2018\*\*](https://bit.ly/joe_eRum_2018)



# Why?

- Most users/organizations can benefit from automatic machine learning pipelines.
  - Eliminate time wasted on human errors, debugging etc.
- Model interpretations is crucial for those who must explain their models to regulators or customers.

# You will learn ...

- How to build high quality H<sub>2</sub>O models (almost) automatically.
- How to explain predictions from complex H<sub>2</sub>O models with LIME.
- **Bonus:** A real use-case that led to multimillion-dollar baseball decisions earlier this year.

# About Me



- **Before H<sub>2</sub>O**

- Water Engineer / Researcher / Matlab Fan Boy @matlabulous
- Discovered R, Python, H<sub>2</sub>O ... never look back again
- Data Scientist at Virgin Media (UK), Domino Data Lab (US)

- **At H<sub>2</sub>O ...**

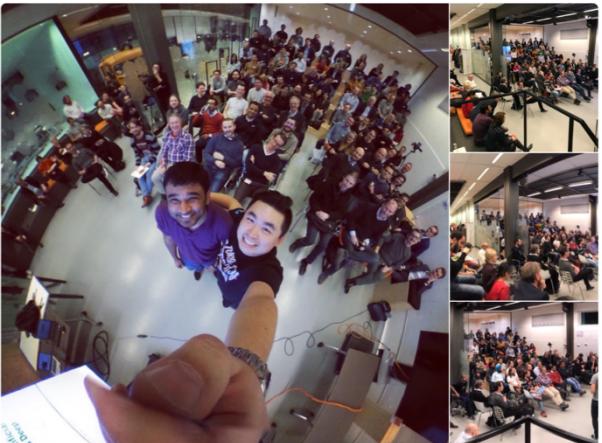
- Data Scientist / Evangelist /
- Sales Engineer / Solution Architect /
- Community Manager  
(The harsh reality of startup life)
- H<sub>2</sub>O SWAG Photographer (#AroundTheWorldWithH2Oai)
- H<sub>2</sub>O SWAG Distributor  
(Love H<sub>2</sub>O? Come get some stickers!)

# What I really do ...



Jo-fai (Joe) Chow  
@matlabulous

Thanks @ingnl for hosting @h2oai #meetup in #Amsterdam last week. Tremendous turnout and great discussions.  
#AroundTheWorldWithH2Oai #360Selfie 🇳🇱  
cc @fishnets88



7:15 AM - 26 Feb 2018 from Amsterdam, The Netherlands



Jo-fai (Joe) Chow  
@matlabulous

Another #FullHouse @h2oai #LondonAI #meetup tonight. Thanks @MSFTReactor for the amazing venue and food! #OpenSource #Community #MVPBuzz  
#AroundTheWorldWithH2Oai #360Selfie 🇬🇧  
cc our guest speakers @SKREDDY99  
@cheukting\_ho & Josh Warwick



7:15 PM - 12 Mar 2018 from London, England



Jo-fai (Joe) Chow  
@matlabulous

Awesome #KNIMESummit2018  
#KNIMESpringSummit in #Berlin. @knime  
@Kuriooos Marten here is our #360Selfie cc  
@h2oai #AroundTheWorldWithH2Oai 🇩🇪  
#OpenSource #MachineLearning  
#Community 💪

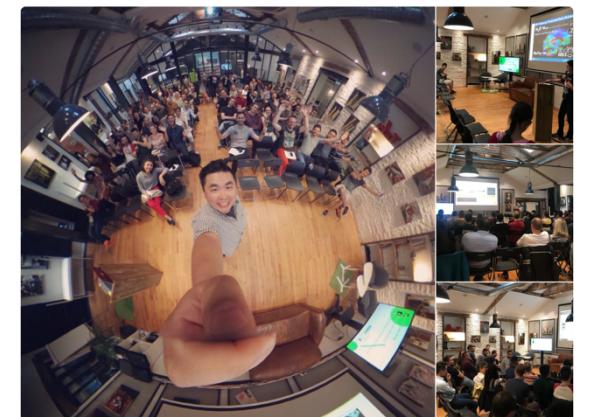


1:54 PM - 7 Mar 2018 from Hotel Berlin



Jo-fai (Joe) Chow  
@matlabulous

Merci beaucoup Alexia, Samia & Aurelie from @tlse\_dasci. We had our very first @h2oai #meetup in #Toulouse tonight. Fantastic crowd and awesome @HarryCoworking venue. We hope to see you all again in the future. Here is our #360selfie 📸  
#AroundTheWorldWithH2Oai 🇫🇷



10:35 PM - 23 Apr 2018 from Toulouse, France

Reminder: #360Selfie

# Agenda

Time	Topics / Tasks
1:30 – 1:45 pm	Install <b>h2o</b> , <b>lime</b> , <b>mlbench</b> from CRAN slides/code: <a href="https://bit.ly/joe_eRum_2018">bit.ly/joe_eRum_2018</a>
1:45 – 2:00 pm	Introduction ( $H_2O$ , AutoML, LIME)
2:00 – 2:30 pm	Regression Example
2:30 – 3:00 pm	Classification Example
3:00 – 3:30 pm	☕️🍰🍪
3:30 – 3:45 pm	Quick Recap
3:45 – 4:15 pm	Real Use-Case: Moneyball
4:15 – 4:30 pm	Other $H_2O$ News + Q & A





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



Time	Topics / Tasks
1:30 – 1:45 pm	Install <b>h2o</b> , <b>lime</b> , <b>mlbench</b> from CRAN slides/code: <a href="https://bit.ly/joe_eRum_2018">bit.ly/joe_eRum_2018</a>

## LIME

**Reference:** <https://github.com/thomasp85/lime>

```
# Install 'lime' from CRAN
install.packages('lime')
```

## H2O

**Reference:** <https://www.h2o.ai/download/>

```
# Install 'h2o' from CRAN
install.packages('h2o')
```

... and **mlbench** for datasets



Time	Topics / Tasks
1:30 – 1:45 pm	Install h2o, lime, mlbench from CRAN slides/code: <a href="https://bit.ly/joe_eRum_2018">bit.ly/joe_eRum_2018</a>
1:45 – 2:00 pm	<b>Introduction (H<sub>2</sub>O, AutoML, LIME)</b>
2:00 – 2:30 pm	Regression Example
2:30 – 3:00 pm	Classification Example
3:00 – 3:30 pm	
3:30 – 3:45 pm	Quick Recap
3:45 – 4:15 pm	Real Use-Case: Moneyball
4:15 – 4:30 pm	Other H <sub>2</sub> O News + Q & A

# Have you seen Avengers: Infinity War?

Do you know all the characters in the movie? (No spoilers - I promise)

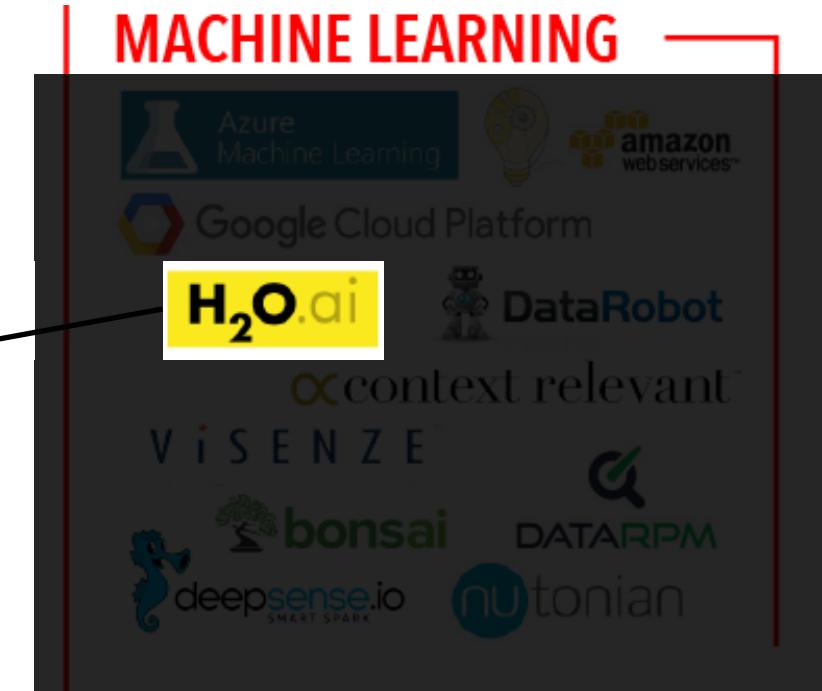
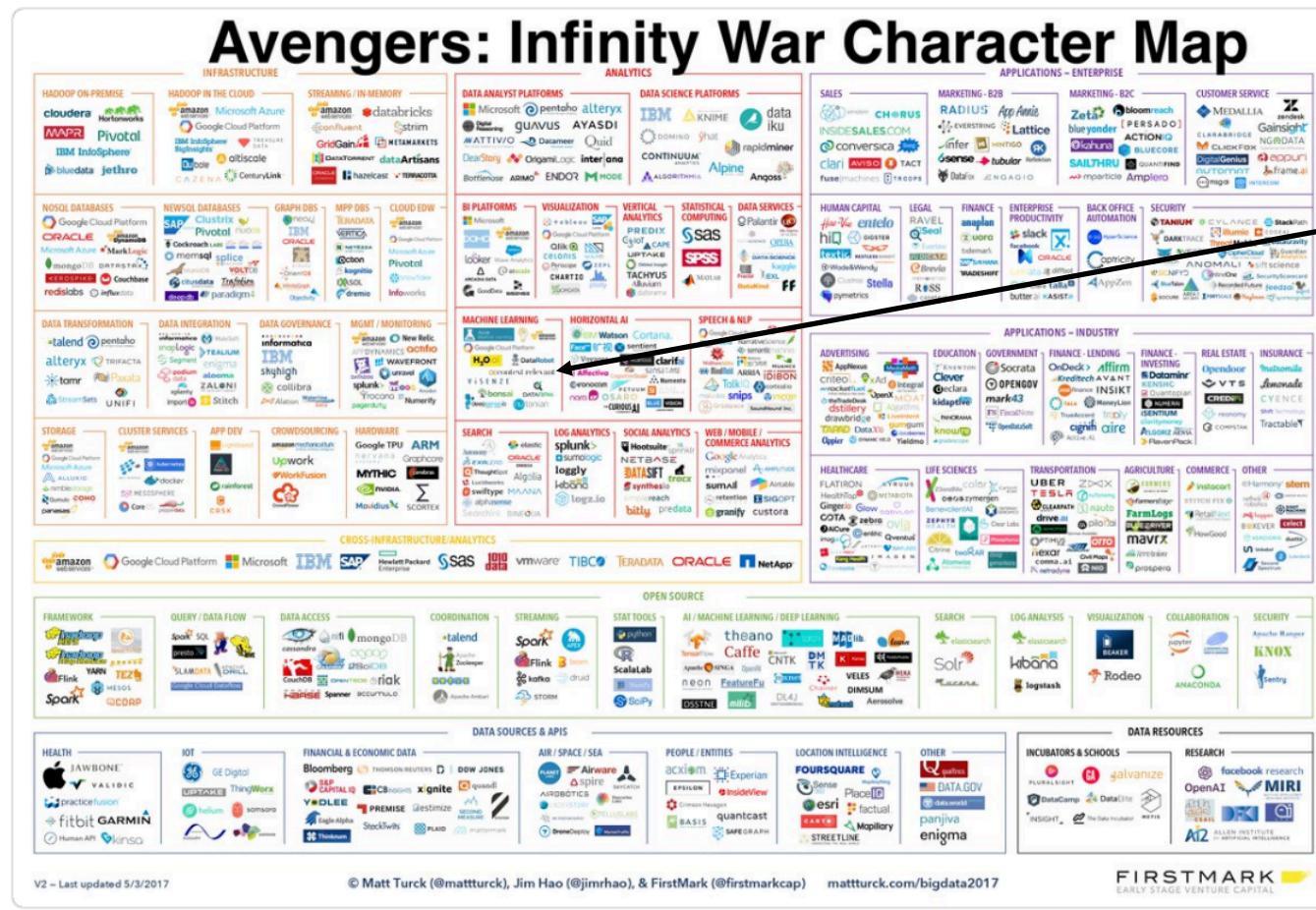


Vicki Boykis  
@vboykis

Follow



I made a guide for anyone who was as confused by all the characters in Infinity War as I was.



# Gartner names H2O as Leader with the most completeness of vision

- H2O.ai recognized as a **technology leader with most completeness of vision**
- H2O.ai was recognized for the mindshare, partner network and status as a **quasi-industry standard** for machine learning and AI.
- **H2O customers gave the highest overall score** among all the vendors for sales relationship and account management, customer support (onboarding, troubleshooting, etc.) and overall service and support.

Figure 1. Magic Quadrant for Data Science and Machine-Learning Platforms



Source: Gartner (February 2018)

As of January 2018

© Gartner, Inc

# Platforms with H<sub>2</sub>O integration



srisatish  
@srisatish

Following

Replying to @BobMuenchen @knime @h2oai

@KNIME gained the ability to run @H2O.ai algorithms, so these two may be viewed as complementary, not competitors  
#Ecosystem #OpenSource

3:32 PM - 2 Mar 2018



H<sub>2</sub>O + KNIME Talk  
at KNIME Summit  
Mar 2017

1:54 PM - 7 Mar 2018 from Hotel Berlin

Figure 1. Magic Quadrant for Data Science and Machine-Learning Platforms



Source: Gartner (February 2018)

© Gartner, Inc

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

# Company Overview

<b>Founded</b>	2012, Series C in Nov, 2017
<b>Products</b>	<ul style="list-style-type: none"><li>• Driverless AI – Automated Machine Learning</li><li>• H<sub>2</sub>O Open Source Machine Learning</li><li>• Sparkling Water</li></ul>
<b>Mission</b>	Democratize AI. Do Good
<b>Team</b>	<p>~100 employees</p> <ul style="list-style-type: none"><li>• Distributed Systems Engineers doing Machine Learning</li><li>• World-class visualization designers</li></ul>
<b>Offices</b>	Mountain View, London, Prague



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



In-Memory, Distributed  
Machine Learning Algorithms  
with H<sub>2</sub>O Flow GUI



H2O AI Open Source Engine  
Integration with Spark



Lightning Fast machine  
learning on GPUs

DRIVERLESSAI

Automatic feature  
engineering, machine  
learning and interpretability

# Steam

Secure multi-tenant H<sub>2</sub>O clusters



## This Workshop

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



In-Memory, Distributed  
Machine Learning Algorithms  
with H<sub>2</sub>O Flow GUI



H2O AI Open Source Engine  
Integration with Spark



Lightning Fast machine  
learning on GPUs

DRIVERLESSAI

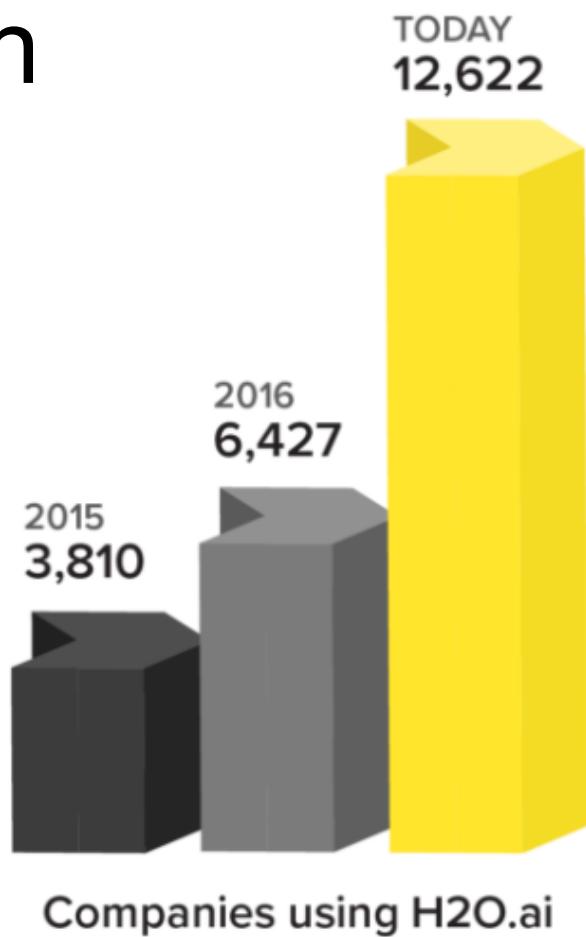
Automatic feature  
engineering, machine  
learning and interpretability

# Steam

Secure multi-tenant H<sub>2</sub>O clusters



# Worldwide Community Adoption



\* DATA FROM GOOGLE ANALYTICS EMBEDDED IN THE END USER PRODUCT

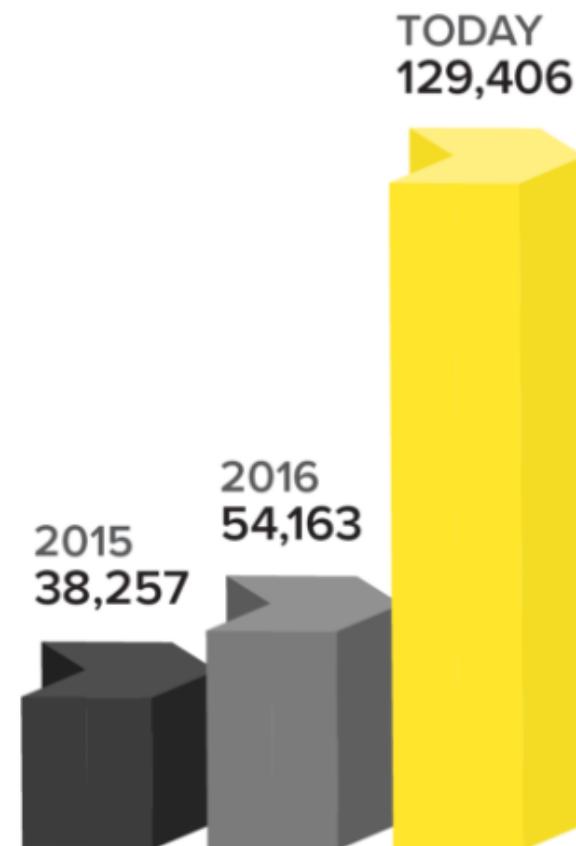
222 OF FORTUNE  
THE 500



8 OF TOP 10  
BANKS

7 OF TOP 10  
INSURANCE COMPANIES

4 OF TOP 10  
HEALTHCARE COMPANIES



H2O.ai **H<sub>2</sub>O.ai**

# H2O.ai Solution Leadership Across Verticals



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

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



**Steph Locke**  
@SteffLocke

Following

My #rstats #datascience goto   
IO: odbc readxl httr  
EDA: DataExplorer  
Prep: tidyverse  
Sampling: rsample modelr  
Feature Engineering: recipes  
Modelling: glmnet **h2o** FFTrees  
Evaluation: broom yardstick  
Deployment: sqlrutils AzureML opencpu  
Monitoring: flexdashboard  
Docs: rmarkdown

4:29 PM - 28 Apr 2018

143 Retweets 591 Likes



10

143

591



**Szilard**  
@DataScienceLA

Following

Friday fun: what's your favorite gradient boosting machine (GBM) library?

58% xgboost

16% lightgbm

24% h2o

2% spark mllib

127 votes • Final results

11:21 PM - 11 May 2018

9 Retweets 9 Likes



3

9

9



Tweet your reply



**Arno Candel** @ArnoCandel · May 12

Replies to @DataScienceLA

Did you know? H2O-3 has XGBoost integration (incl. support for GPU and distributed mode) with standalone Java scoring (MOJO) - train from Flow, R or Python. H2O AutoML and Driverless AI use XGBoost too, and @h2oai contributes to XGBoost in collaboration with @nvidia #h2o4gpu

10

11

11



# Our Mission: Make Machine Learning Accessible to Everyone



Complexity is your enemy. Any fool can make something complicated. It is hard to keep things simple.

— *Richard Branson* —

AZ QUOTES

# 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*



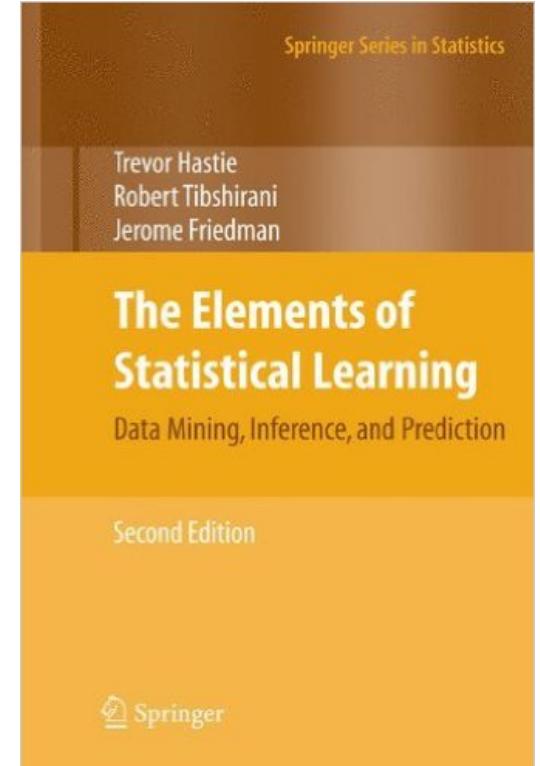
## Dr. Robert Tibshirani

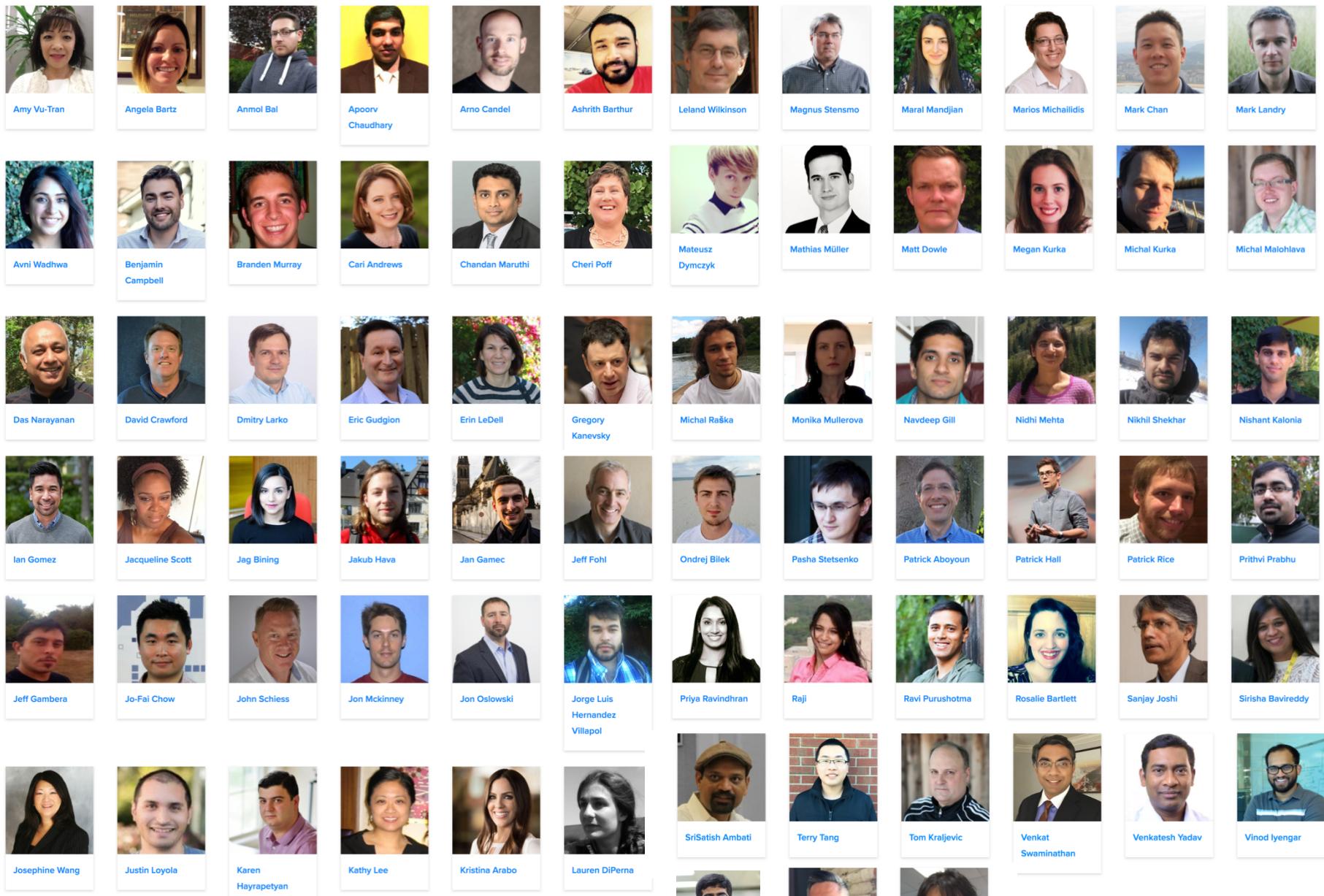
- Professor of Statistics and Health Research and Policy, Stanford University
- PhD in Statistics, Stanford University
- Co-author, *The Elements of Statistical Learning: Prediction, Inference and Data Mining*
- Author, *Regression Shrinkage and Selection via the Lasso*
- Co-author, *An Introduction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



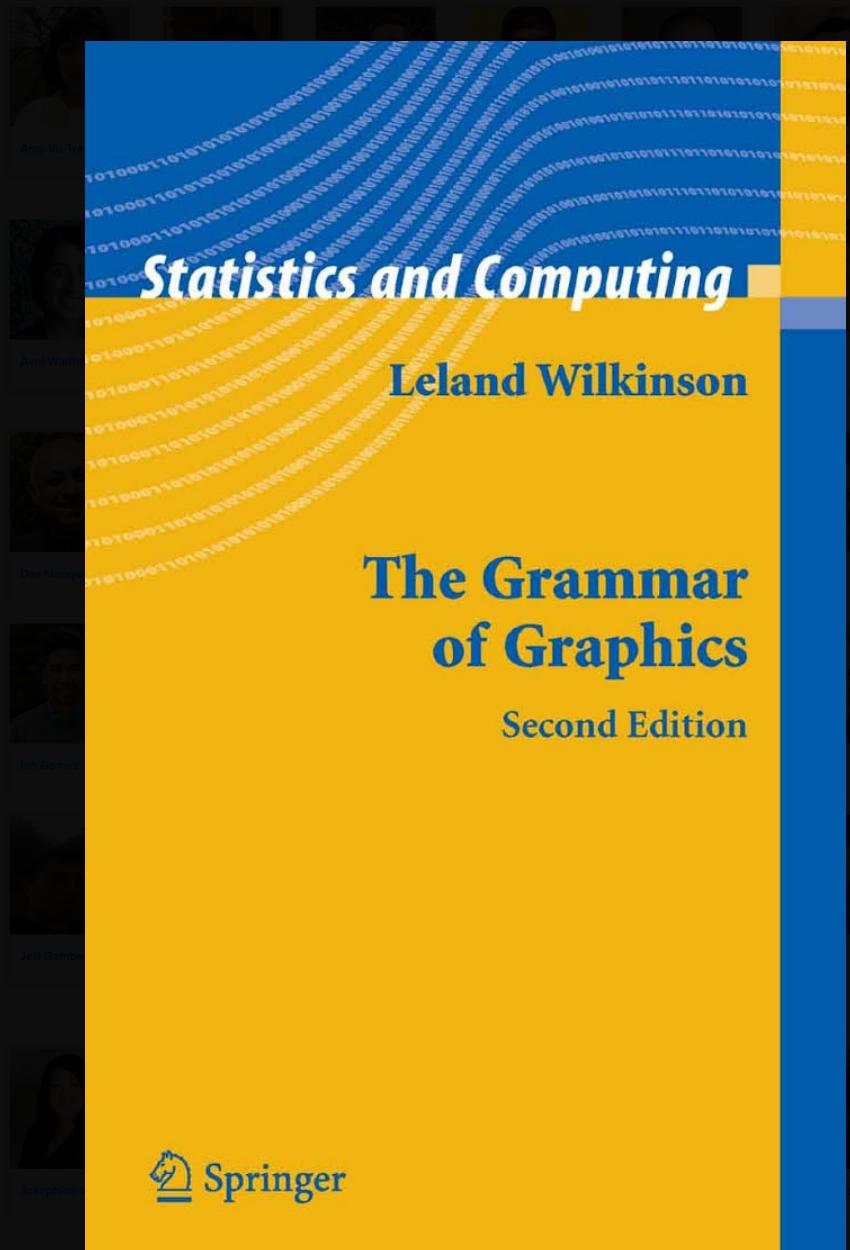
## Dr. Steven Boyd

- 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*





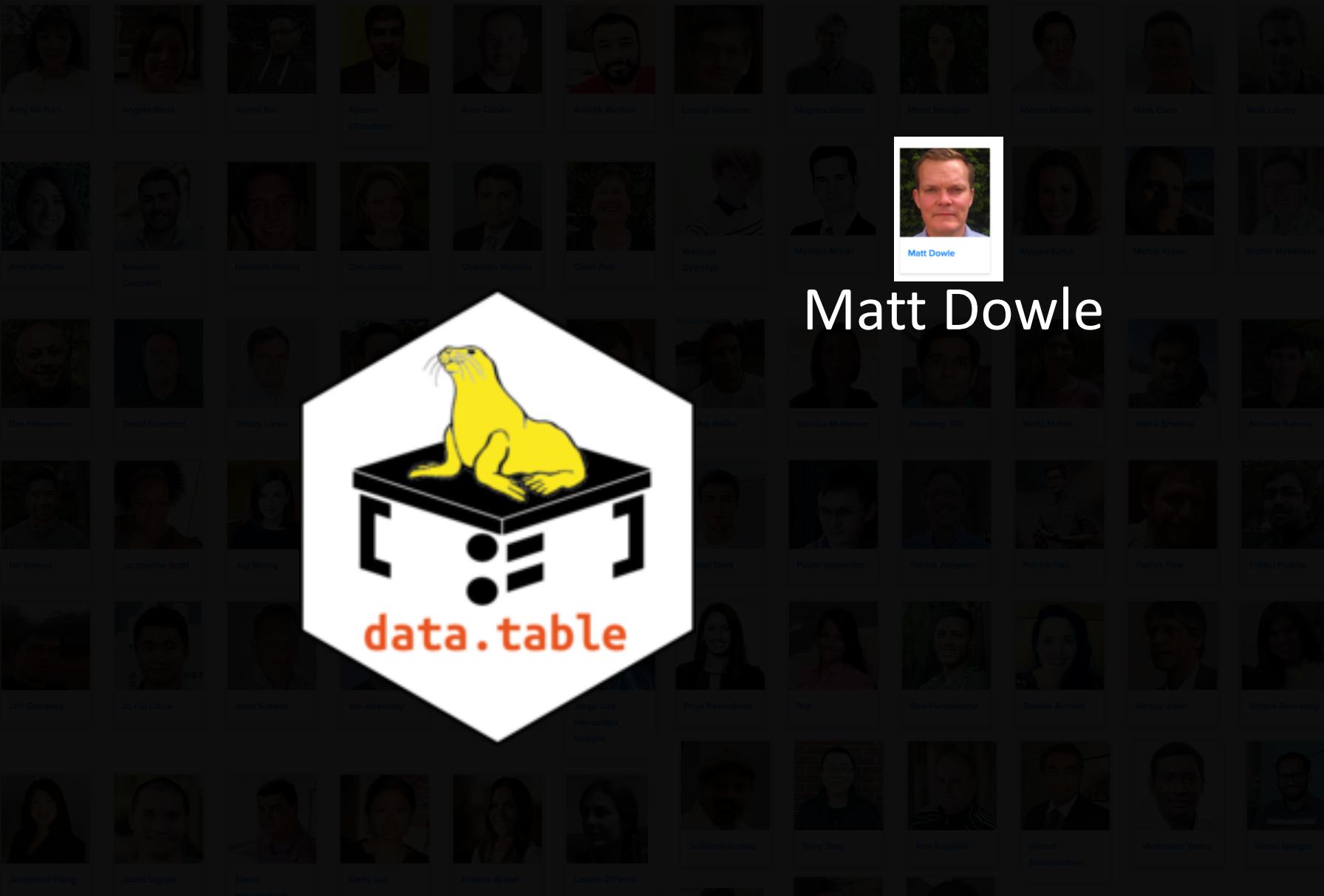
# H<sub>2</sub>O Team



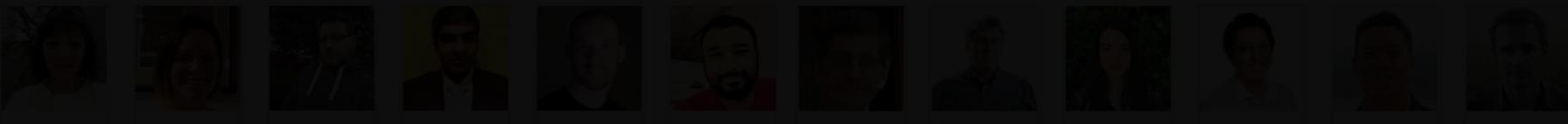
Leland Wilkinson

## Origin of R Package `ggplot2`





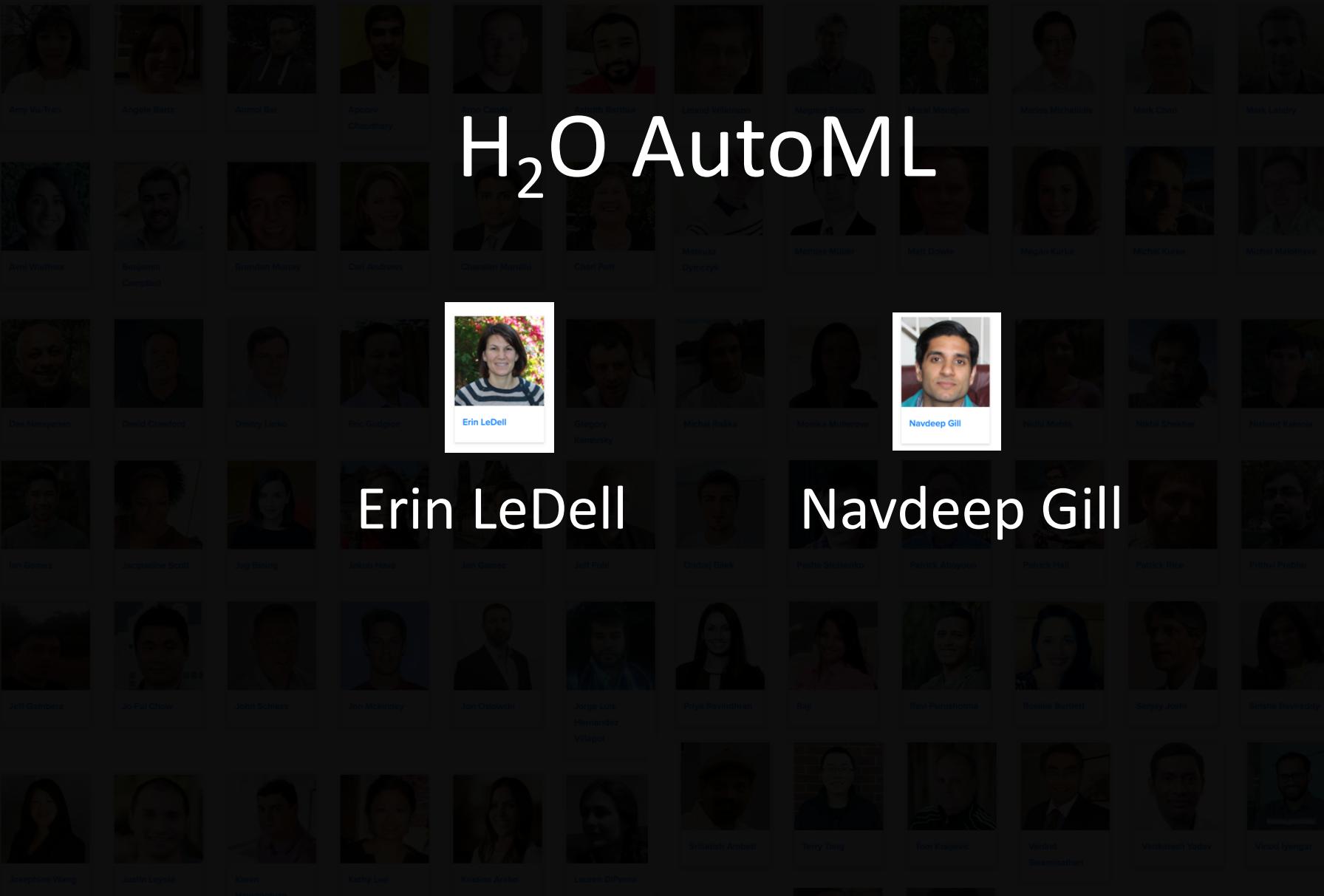
H<sub>2</sub>O Team



# Erin LeDell, Chief ML Scientist Women in ML/DS & R-Ladies Global



H<sub>2</sub>O Team



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

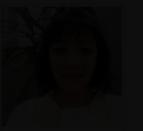
H<sub>2</sub>O Team



About 80,000 Kagglers

H<sub>2</sub>O Team

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



Amy Vu-Tran



Angela Bartz



48th



Arno Candel



Ashwin Berthar



Leland Wilkinson

Magnus Stensmo  
Maral Mandjari

Marios Michailidis

Mark Chan  
Mark Landry

Avni Wedhwa



Benjamin Campbell



Branden Murray



Cari Andrews



Chandan Manohar



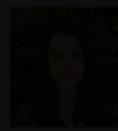
Chen Poff



Mateusz Dymczyk



Mathias Müller



Matt Dowde



Megan Kurka



Michael Kurka



Michal Malofitava



Das Narayanan



David Crawford



Dmitry Larko



Eric Gudgion



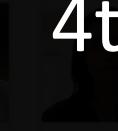
Erin LeDell



Gregory Kanovsky



Michal Rabka



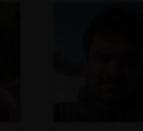
Monika Mullerova



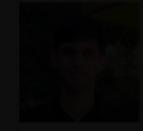
Navdeep Gill



Nidhi Mehta



Nichil Shukhar



Ian Gomez



Jacqueline Scott



Jag Birring



Jakub Horec



Jan Gamec



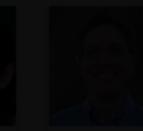
Jeff Follé



Ondrej Blatik



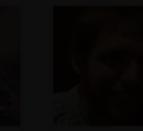
Pasha Biesenski



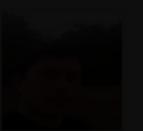
Patrick Abeyoum



Patrick Hall



Patrick Rice



Jeff Gambetta



Jo-Fai Chow



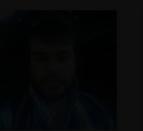
John Seeger



Jozef Luptak



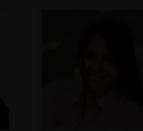
Kristina Arbov



Lauren DiPerna



SriSatish Ambati



Terry Tang



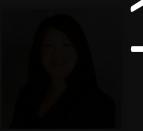
Tom Kraljevic



Venkat Swaminathan



Venkatesh Yadav



Josephine Wang



Justin Loyola



Karen Heynepeyan



Kathy Lee



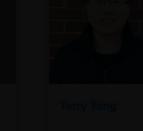
Kristina Arbov



Lauren DiPerna



SriSatish Ambati



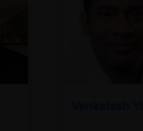
Terry Tang



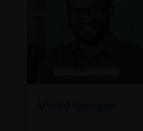
Tom Kraljevic



Venkat Swaminathan



Venkatesh Yadav

H<sub>2</sub>O Team

13th

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

Hoping to get closer to them at some point ...

# H<sub>2</sub>O Team at eRum 2018

Erin LeDell – Ask her about AutoML



Erin LeDell

Jakub (or Kuba) Hava – Ask him  
about Sparkling Water (H<sub>2</sub>O + Spark)



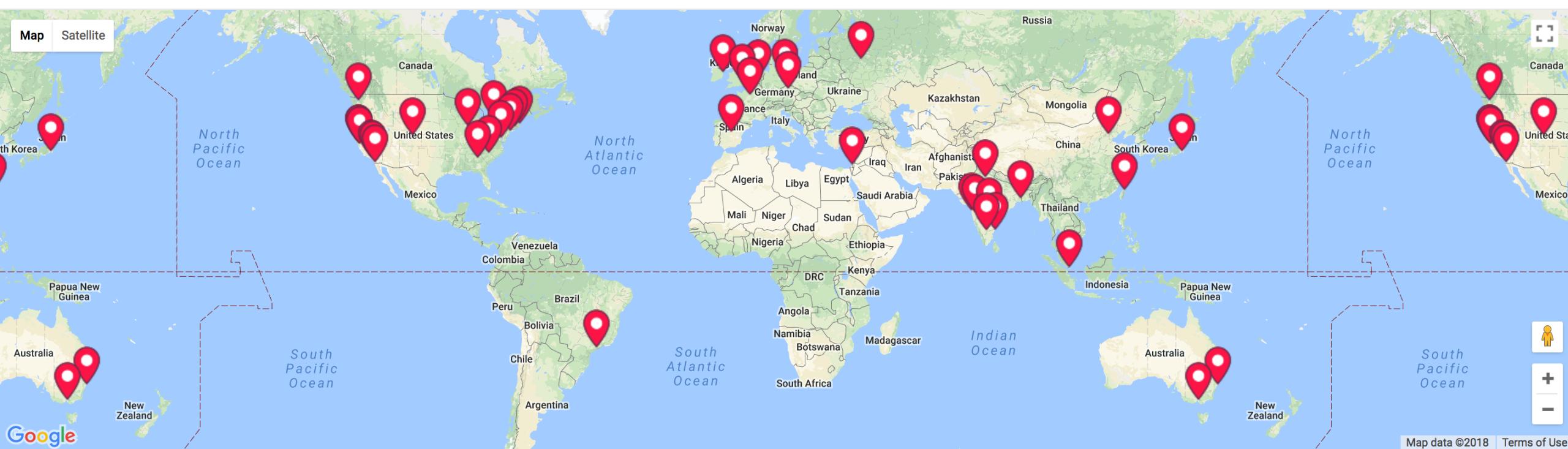
Jakub Hava



Jo-Fai Chow

Joe Chow – Happy to trade rare H<sub>2</sub>O swag for beers  
... or talk about H<sub>2</sub>O-3, Driverless AI, R, Shiny ...

H<sub>2</sub>O Team



# H2O Artificial Intelligence and Machine Learning

Members  
**78,356**

Groups  
**39**

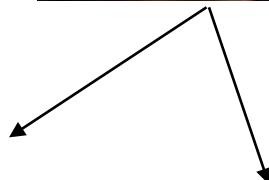
# Countries **18**

<https://www.meetup.com/pro/h2oai/>

# We sponsor meetups



Source: <https://www.maddycoupons.in/blog/yummy-yummy-pizza/>



**R-Ladies London**  
Location  
London, United Kingdom  
Members  
1,040  
Organizers  
Chin Tan and 7 others



**Artificial Intelligence (AI) Club for Gender Minorities!**  
London, United Kingdom · 489 members · Public group

Organized by  
Chin Tan and 5 others

Share: [Facebook](#) [Twitter](#) [LinkedIn](#)

**London Data Science Workshop (formerly London Kaggle Meetup)**  
Location  
London, United Kingdom  
Members  
2,783  
Organizers  
Alex Glaser and 6 others



**Women in Kaggle**  
Location  
London, United Kingdom

Members  
261  
Organizers  
Julia MacMillan and 3 others

You're a member [Join](#) [Share](#)

... and more

# We encourage diversity

Meetups	Female Speaker	Female Speaker Ratio
London Dec 2017	Kasia Kulma	1/3
Amsterdam Feb 2018	Andreea Bejinaru	1/2
London Mar 2018	Cheuk Ting Ho	1/3

Since Dec 2017 = 3/8 = 37.5%

Encourage your friends/colleagues to give a talk.

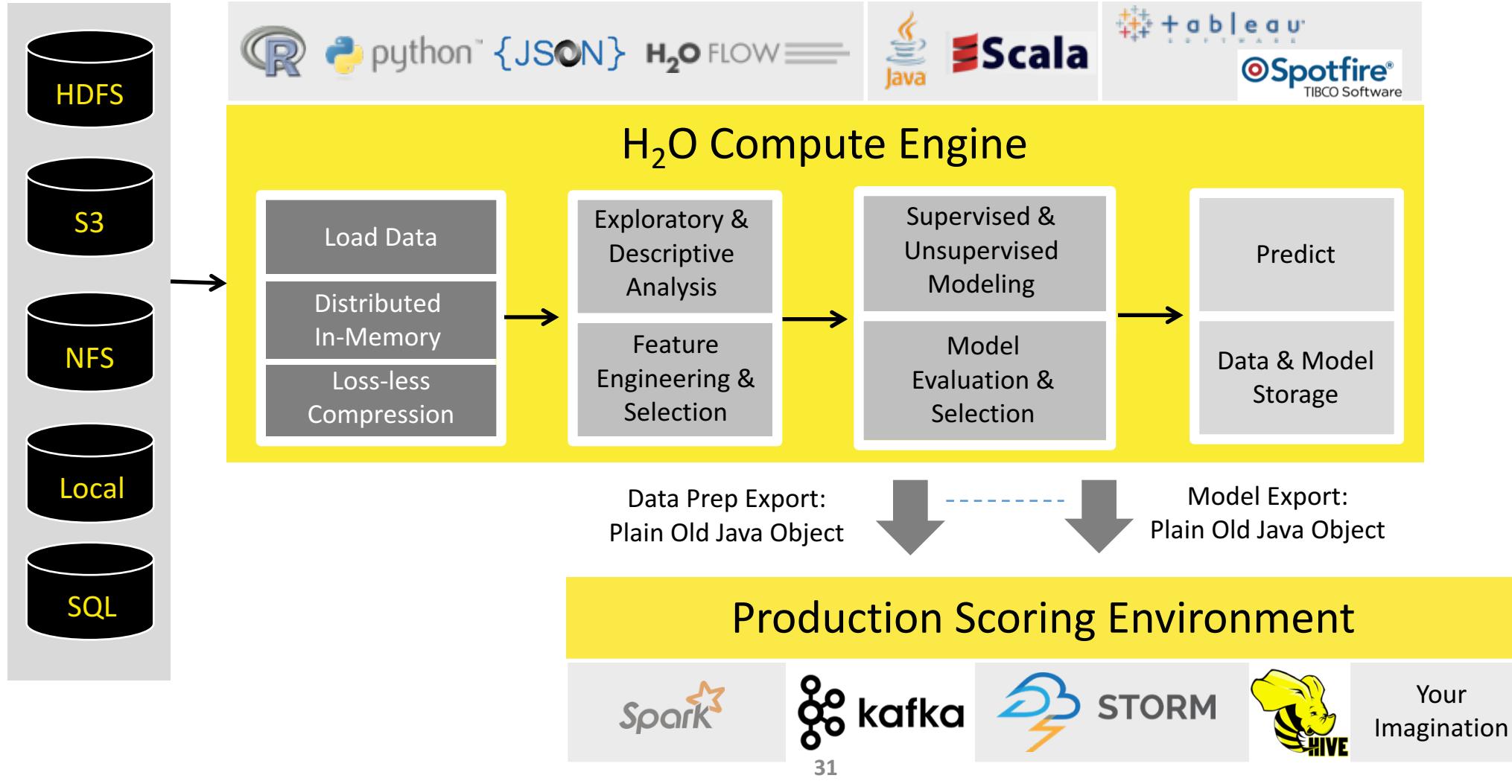
joe@h2o.ai

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

Automatic Machine Learning with H<sub>2</sub>O

<http://docs.h2o.ai/h2o/latest-stable/h2o-docs/automl.html>

# High Level Architecture

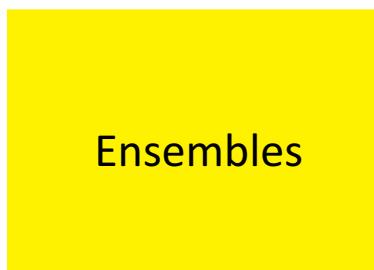


# H<sub>2</sub>O-3 Algorithms Overview

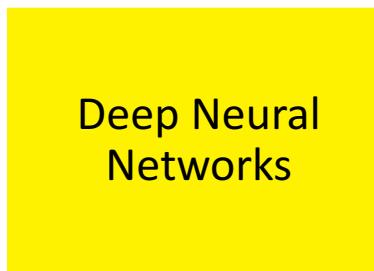
## Supervised Learning



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



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

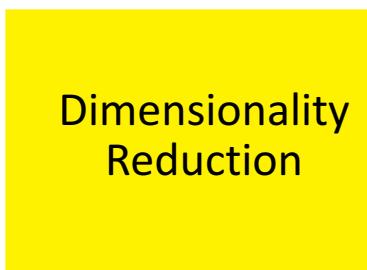


- **Deep learning:** Create multi-layer feed forward neural networks starting with an input layer followed by multiple layers of nonlinear transformations

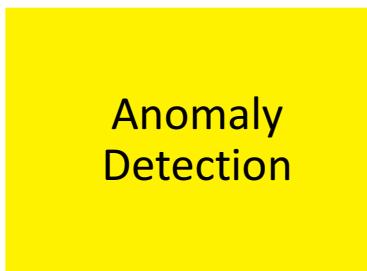
## Unsupervised Learning



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

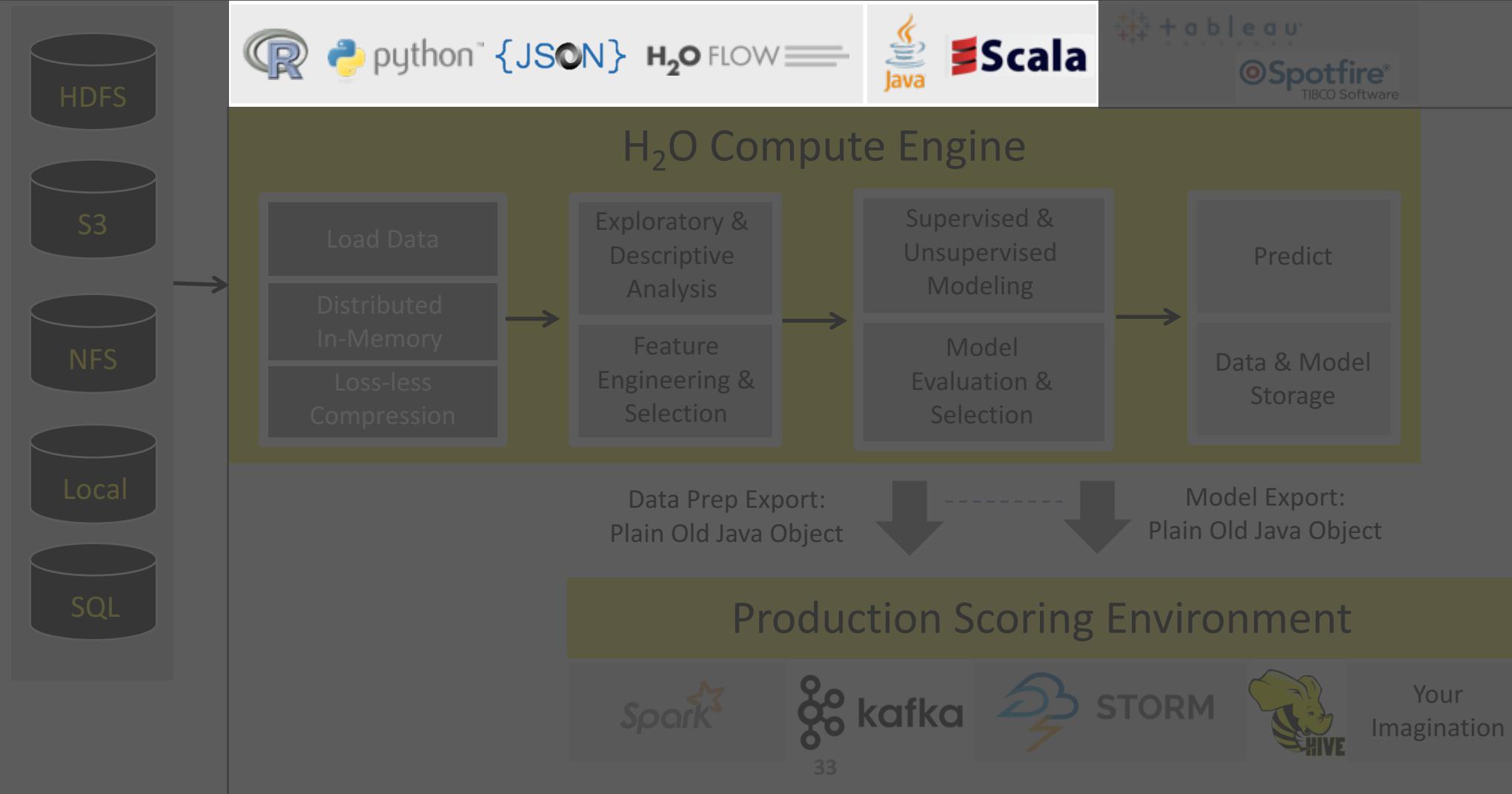


- **Principal Component Analysis:** Linearly transforms correlated variables to independent components
- **Generalized Low Rank Models:** extend the idea of PCA to handle arbitrary data consisting of numerical, Boolean, categorical, and missing data



- **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 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 machine learning models and related routines. The main workspace on the left is titled "Untitled Flow" and contains a single step labeled "assist". Below this is a table titled "Assistance" with a list of routines and their descriptions. The right side of the interface features a sidebar with sections for "OUTLINE", "FLOWS", "CLIPS", and "HELP" (which is also highlighted in yellow). The "HELP" section contains links for "Using Flow for the first time?", "Quickstart Videos", and "view example Flows". It also includes a "STAR H2O ON GITHUB!" button and a "GENERAL" section with a bulleted list of links. At the bottom, there's an "EXAMPLES" section with a paragraph about Flow packs and a "Browse installed packs..." link.

Model

- 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...

ROUTINE

Routine	Description
<code>importFiles</code>	Import file(s) into H <sub>2</sub> O
<code>getFrames</code>	Get a list of frames in H <sub>2</sub> O
<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 H <sub>2</sub> O
<code>getGrids</code>	Get a list of grid search results
<code>getPredictions</code>	Get a list of predictions in H <sub>2</sub> O
<code>getJobs</code>	Get a list of jobs running in H <sub>2</sub> O
<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

ASSIST

?

### Assistance

Routine	Description
<code>importFiles</code>	Import file(s) into H <sub>2</sub> O
<code>getFrames</code>	Get a list of frames in H <sub>2</sub> O
<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 H <sub>2</sub> O
<code>getGrids</code>	Get a list of grid search results
<code>getPredictions</code>	Get a list of predictions in H <sub>2</sub> O
<code>getJobs</code>	Get a list of jobs running in H <sub>2</sub> O
<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

OUTLINE FLOWS CLIPS HELP

?

### Help

Using Flow for the first time?

Quickstart Videos

Or, view example Flows to explore and learn H<sub>2</sub>O.

STAR H2O ON GITHUB!

Star 2,387

GENERAL

- Flow Web UI ...
- ... Importing Data
- ... Building Models
- ... Making Predictions
- ... Using Flows
- ... Troubleshooting Flow

EXAMPLES

Flow packs are a great way to explore and learn H<sub>2</sub>O. Try out these Flows and run them in your browser.

Browse installed packs...

localhost:54321/flow/index.html#

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

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

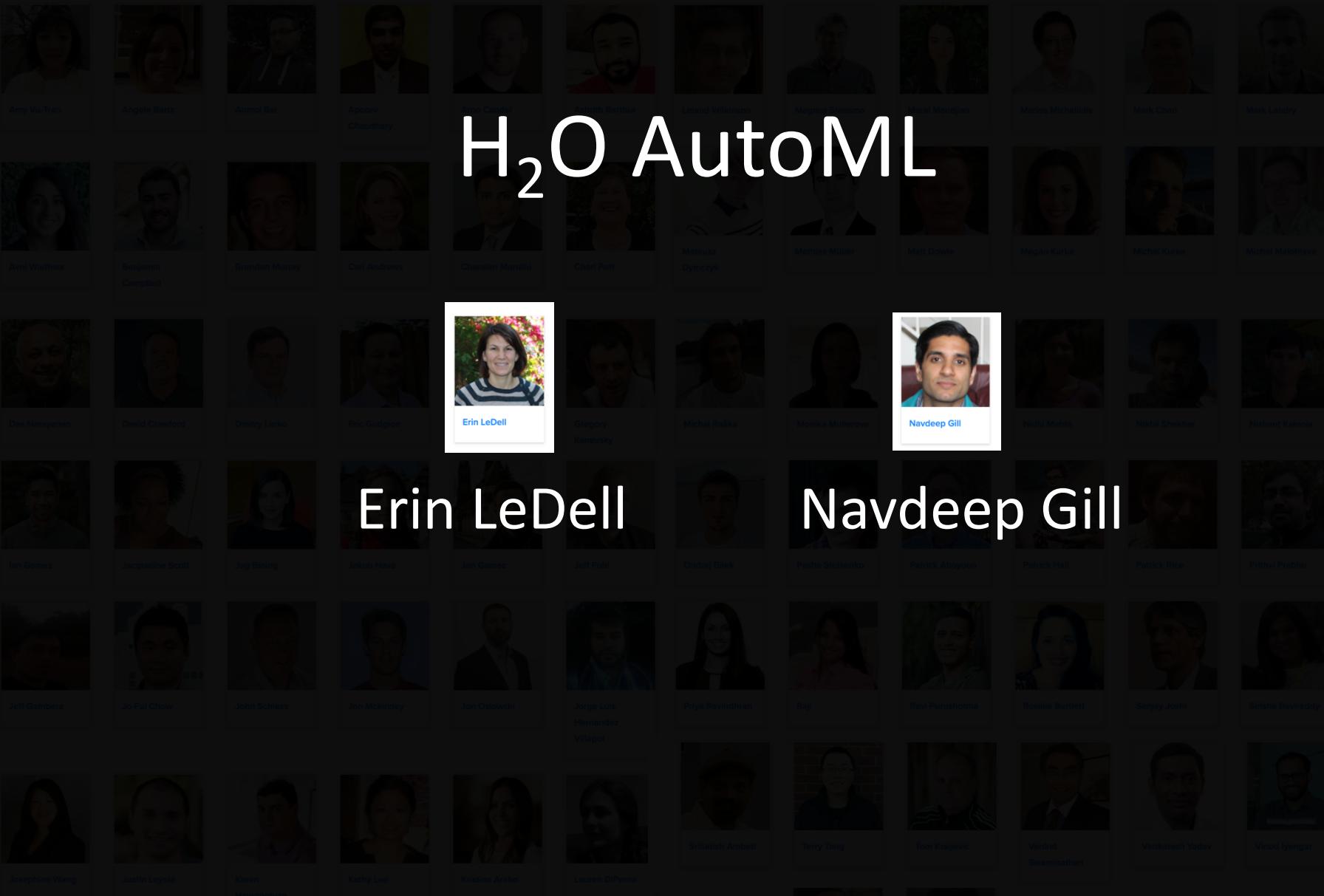
The image shows two side-by-side screenshots illustrating the use of H<sub>2</sub>O with R and Python.

**Left Screenshot (RStudio):** A screenshot of the RStudio Source Editor window titled "credit\_card\_example.R". The code is an R script demonstrating how to use the H<sub>2</sub>O library to train a GBM model on a credit card dataset. It includes importing datasets from S3, initializing H<sub>2</sub>O, defining features and target, training a GBM model, making predictions, and printing the leaderboard. The code is annotated with comments explaining the steps.

```
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" running on "localhost:8888". The notebook displays the output of the R script. It shows the H<sub>2</sub>O cluster starting up, including logs and a table of cluster statistics. Below this, it shows the import of datasets and their summaries. The summary table for the "df\_train" dataset is shown below:

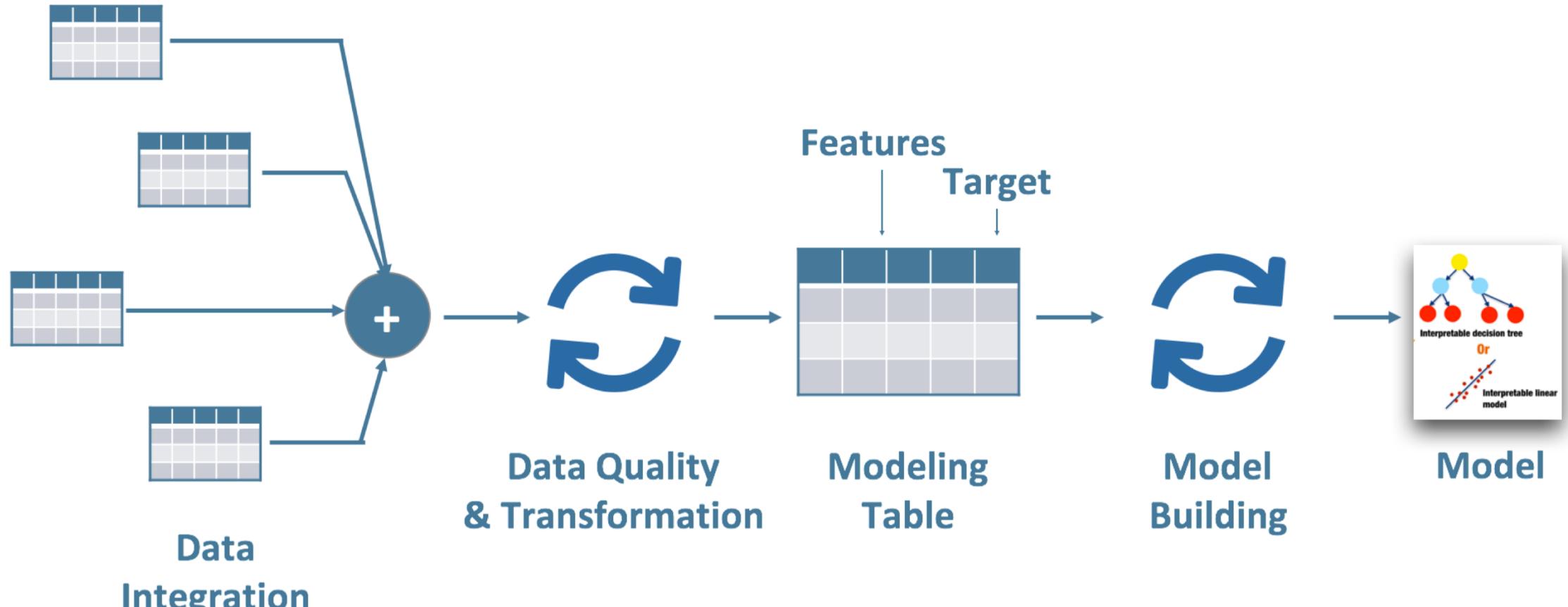
	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_0	PAY_2	PAY_3	PAY_4
type	int	enum	int	int	int	int	int	int	int
mins	10000.0		0.0	0.0	21.0	-2.0	-2.0	-2.0	-2.0
mean	165471.466667		1.85	1.55578703704	35.4053240741	-0.00523148148148	-0.122361111111	-0.15537037037	-0.210601
maxs	1000000.0		6.0	3.0	79.0	8.0	8.0	8.0	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



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

H<sub>2</sub>O Team

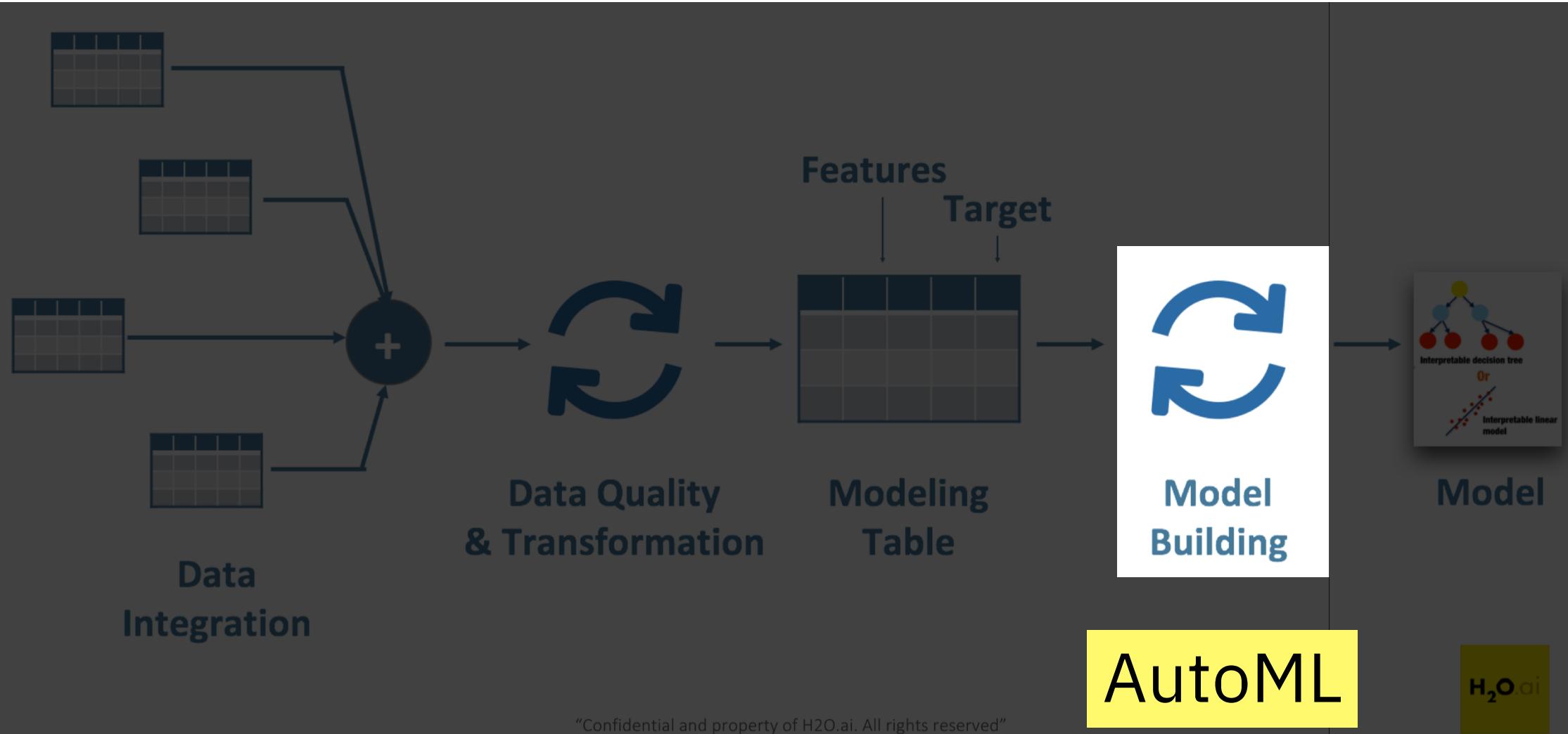
# Typical Enterprise Machine Learning Workflow



“Confidential and property of H2O.ai. All rights reserved”



# Typical Enterprise Machine Learning Workflow



"Confidential and property of H2O.ai. All rights reserved"

# 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 or limit on the number of total models trained.

In both the R and Python API, AutoML uses the same data-related arguments, `x`, `y`, `training_frame`, `validation_frame`, as the other H2O algorithms. Most of the time, all you'll need to do is specify the data arguments. You can then configure values for `max_runtime_secs` and/or `max_models` to set explicit time or number-of-model limits on your run.

## Required Parameters

### Required Data Parameters

- `y`: This argument is the name (or index) of the response column.
- `training_frame`: Specifies the training set.

### Required Stopping Parameters

One of the following stopping strategies (time or number-of-model based) must be specified. When both options are set, then the AutoML run will stop as soon as it hits one of either of these limits.

- `max_runtime_secs`: This argument controls how long the AutoML run will execute for. This defaults to 3600 seconds (1 hour).
- `max_models`: Specify the maximum number of models to build in an AutoML run, excluding the Stacked Ensemble models. Defaults to `NULL/None`.

## AutoML Output

The AutoML object includes a “leaderboard” of models that were trained in the process, including the 5-fold cross-validated model performance (by default). The number of folds used in the model evaluation process can be adjusted using the `n_folds` parameter. If the user would like to score the models on a specific dataset, they can specify the `leaderboard_frame` argument, and then the leaderboard will show scores on that dataset instead.

The models are 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_AllModels_0_AutoML_20171121_012135	0.788321	0.554019
StackedEnsemble_BestOfFamily_0_AutoML_20171121_012135	0.783099	0.559286
GBM_grid_0_AutoML_20171121_012135_model_1	0.780554	0.560248
GBM_grid_0_AutoML_20171121_012135_model_0	0.779713	0.562142
GBM_grid_0_AutoML_20171121_012135_model_2	0.776206	0.564970
GBM_grid_0_AutoML_20171121_012135_model_3	0.771026	0.570270
DRF_0_AutoML_20171121_012135	0.734653	0.601520
XRT_0_AutoML_20171121_012135	0.730457	0.611706
GBM_grid_0_AutoML_20171121_012135_model_4	0.727098	0.666513
GLM_grid_0_AutoML_20171121_012135_model_0	0.685211	0.635138

# About Machine Learning Interpretability

LIME (Local Interpretable Model-Agnostic Explanations)

... and more

# Acknowledgement

- **Marco Tulio Ribeiro:** Original LIME Framework and Python package 
- **Thomas Lin Pedersen:** LIME R package 
- **Matt Dancho:** LIME + H2O AutoML example + LIME R package improvement 
- **Kasia Kulma:** LIME + H2O AutoML example 

# Why Should I Trust Your Model?



System that performs behaviour but you don't know how it works

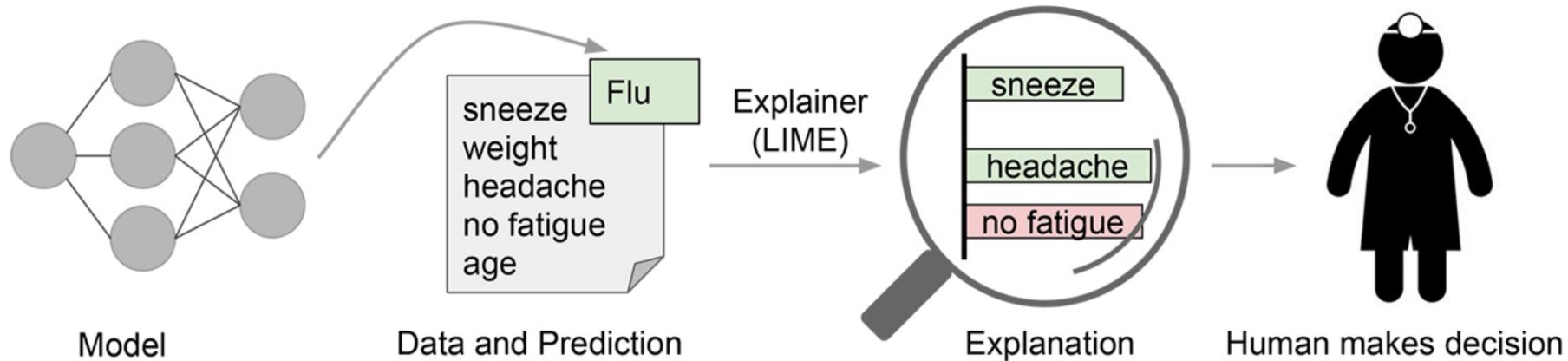


Figure 1. Explaining individual predictions to a human decision-maker. Source: Marco Tulio Ribeiro.

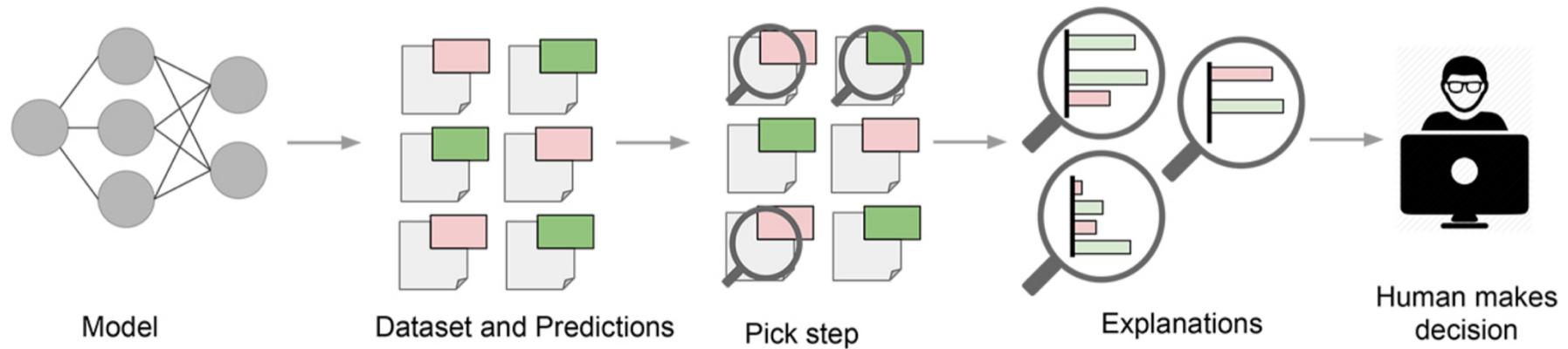


Figure 2. Explaining a model to a human decision-maker. Source: Marco Tulio Ribeiro.

# Local Interpretable Model-Agnostic Explanations

## LIME - How does it work?

### Theory

- LIME approximates model locally as logistic or linear model
- Repeats process many times
- Outputs features that are most important to local models

### Outcome

- Approximate reasoning
- Complex models can be interpreted
  - Neural nets, Random Forest, Ensembles etc.

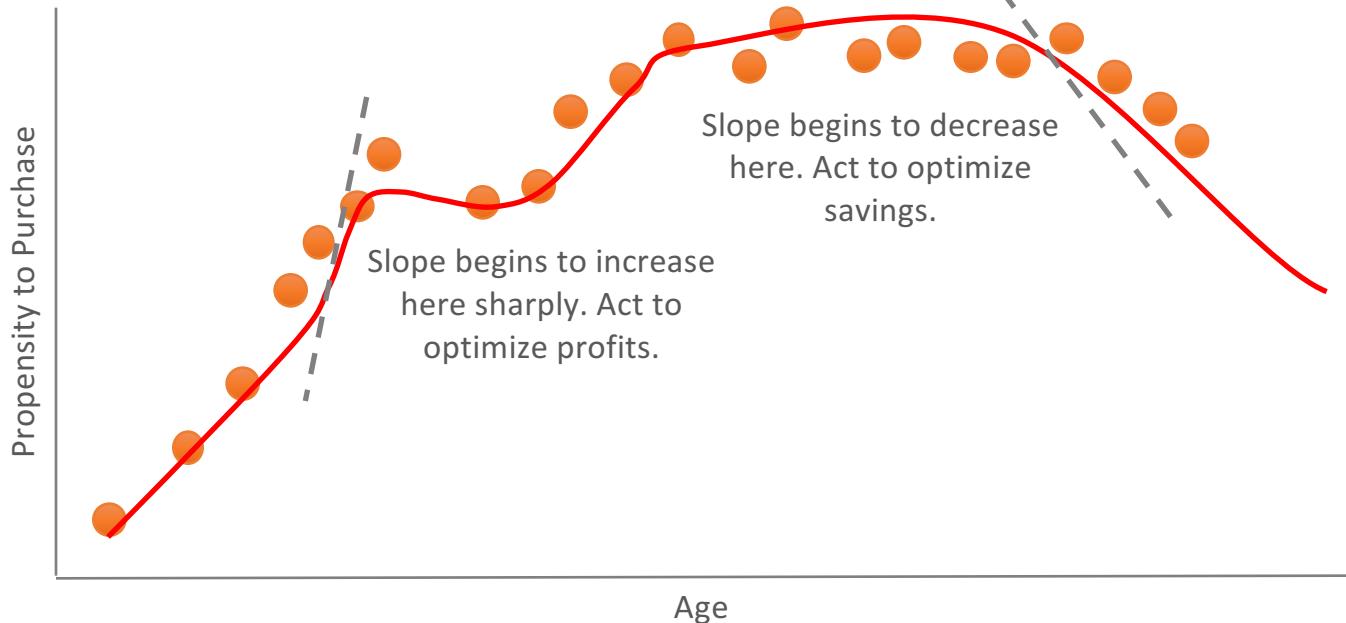
## Linear Models

*Exact explanations for approximate models.*



## Machine Learning

*Approximate explanations for exact models.*

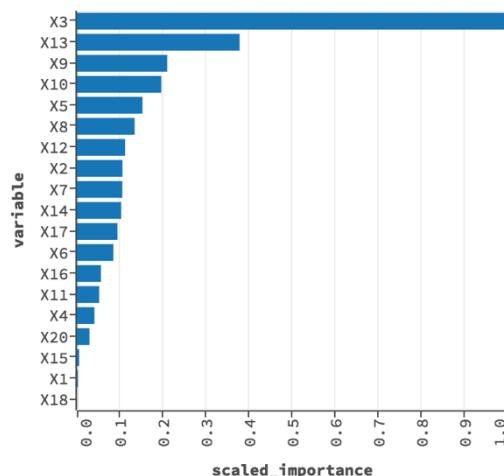


... there are more techniques!

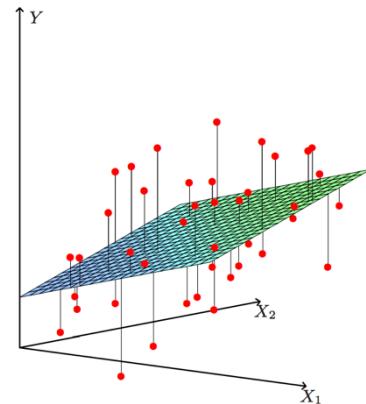
# Machine Learning Interpretability

$$\begin{bmatrix} & \\ & \mathbf{X} & \\ & & \end{bmatrix} \quad \begin{bmatrix} & \\ & \hat{\mathbf{y}} & \\ & & \end{bmatrix}$$

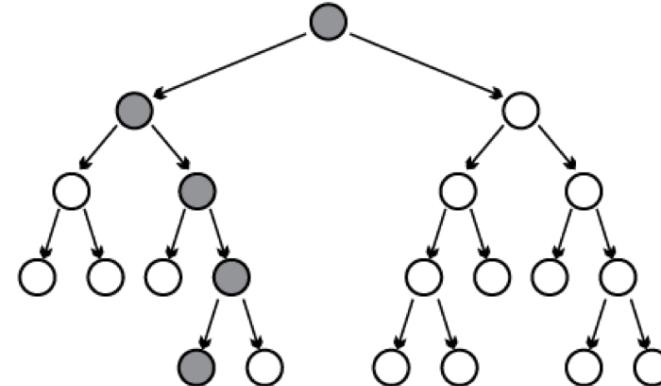
Variable Importance



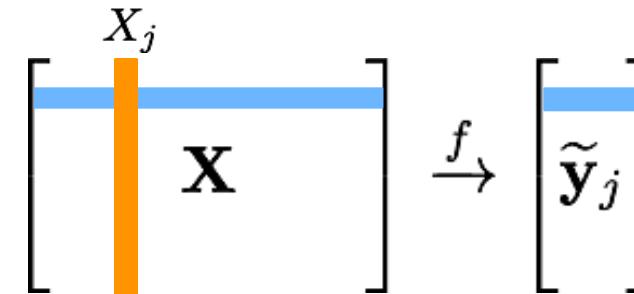
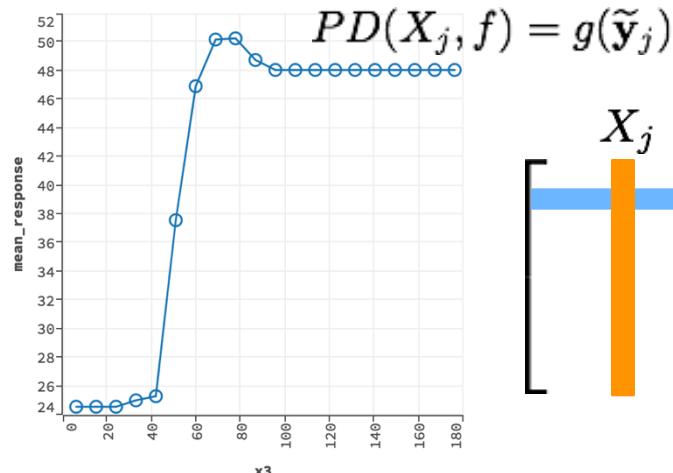
Local Models



Surrogate Model



Partial Dependence



## Workshop day (May 14, 2018 -- Monday)

Room	N15 Aud A (150 seats)	N15 Aud B (150 seats)	N15 103 (75 seats)	N15 101 (50 seats)	N15 106 (40 seats)	N15 203 (25 seats)	N15 202 (20 seats)
8:00							
8:30							
9:00							
9:30	Efficient R programming	DALEX: Descriptive mAchine Learning EXplanations	Clean R code - how to write it and what will the benefits be	Building an Interpretable NLP model to classify tweets	Geocomputation with R	Graphs: A datastructure to query	Forwards Package Development Workshop for Women
10:00							
10:30							
11:00							
11:30	Efficient R programming	DALEX: Descriptive mAchine Learning EXplanations	Clean R code - how to write it and what will the benefits be	Building an Interpretable NLP model to classify tweets	Geocomputation with R	Graphs: A datastructure to query	Forwards Package Development Workshop for Women
12:00							
12:30							
13:00							
13:30							
14:00	Deep Learning with Keras for R	Automatic and Interpretable Machine Learning in R with H2O and LIME	The beauty of data manipulation with data.table	Building a package that lasts	Building a pipeline for reproducible data screening and quality control	Plotting spatial data in R	Forwards Package Development Workshop for Women
14:30							
15:00							
15:30							
16:00	Deep Learning with Keras for R	Automatic and Interpretable Machine Learning in R with H2O and LIME	The beauty of data manipulation with data.table	Building a package that lasts	Building a pipeline for reproducible data screening and quality control	Plotting spatial data in R	Forwards Package Development Workshop for Women
16:30							
...							

This Workshop

Walk from CEU to Akvárium Klub

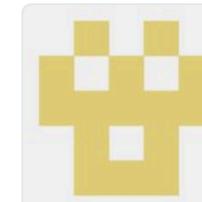


Erin LeDell  
@ledell

Following

List of #MachineLearning model  
interpretability pkgs in R:  
lime, ShapleyR, live, xgboostExplainer,  
breakDown, DALEX #eRum2018

[github.com/thomasp85/lime](https://github.com/thomasp85/lime)  
[github.com/redichh/Shaple...](https://github.com/redichh/Shaple...)  
[mi2datalab.github.io/live/](https://mi2datalab.github.io/live/)  
[github.com/AppliedDataSci ...](https://github.com/AppliedDataSci ...)  
[pbiecek.github.io/breakDown/](https://pbiecek.github.io/breakDown/)  
[pbiecek.github.io/DALEX/](https://pbiecek.github.io/DALEX/)



AppliedDataSciencePartners/xgboostExplainer  
xgboostExplainer - An R package that makes xgboost  
models fully interpretable  
[github.com](https://github.com)

8:34 AM - 14 May 2018

4 Retweets 4 Likes



1 4 4 4



Tweet your reply

Erin LeDell @ledell · 10m

One more (for the randomForest package):  
[mi2datalab.github.io/randomForestEx...](https://mi2datalab.github.io/randomForestEx...)

# Ideas on interpreting machine learning

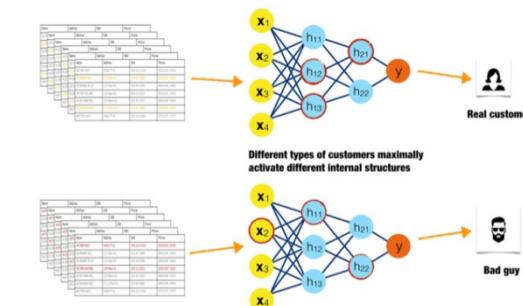
Mix-and-match approaches for visualizing data and interpreting machine learning models and results.

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

*Check out the "Data Science & Machine Learning" sessions at the Strata Data Conference in London, May 21-24, 2018.*

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

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



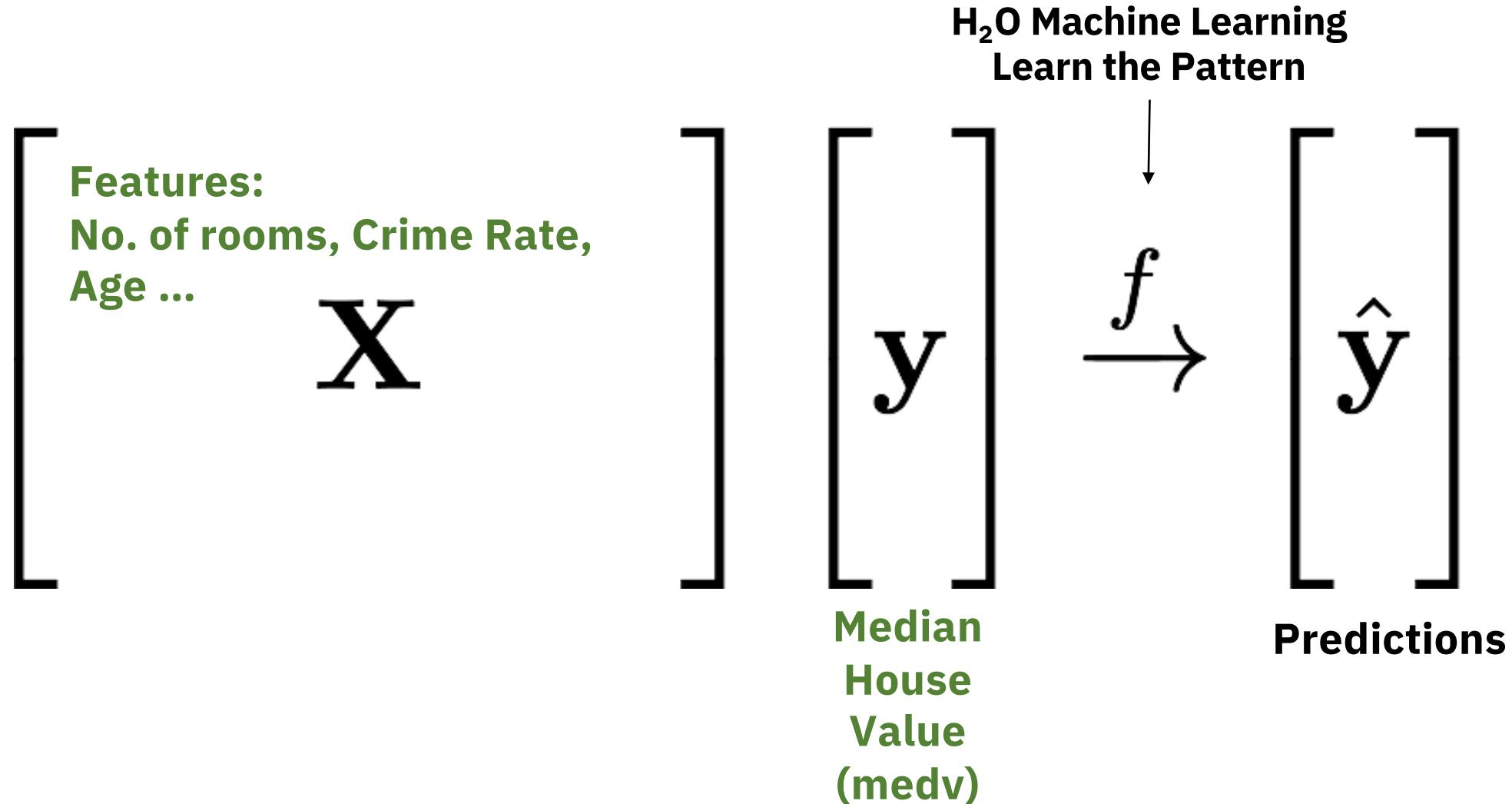
Inputs activating different neurons in a neural network.  
(source: Image courtesy of Patrick Hall and the h2o.ai team, used with permission)

<https://www.oreilly.com/ideas/ideas-on-interpreting-machine-learning>



Time	Topics / Tasks
1:30 – 1:45 pm	Install h2o, lime, mlbench from CRAN slides/code: <a href="https://bit.ly/joe_eRum_2018">bit.ly/joe_eRum_2018</a>
1:45 – 2:00 pm	Introduction (H <sub>2</sub> O, AutoML, LIME)
2:00 – 2:30 pm	Regression Example <code>\examples\regression_...Rmd</code>
2:30 – 3:00 pm	Classification Example
3:00 – 3:30 pm	☕️🍰🍪
3:30 – 3:45 pm	Quick Recap
3:45 – 4:15 pm	Real Use-Case: Moneyball
4:15 – 4:30 pm	Other H <sub>2</sub> O News + Q & A

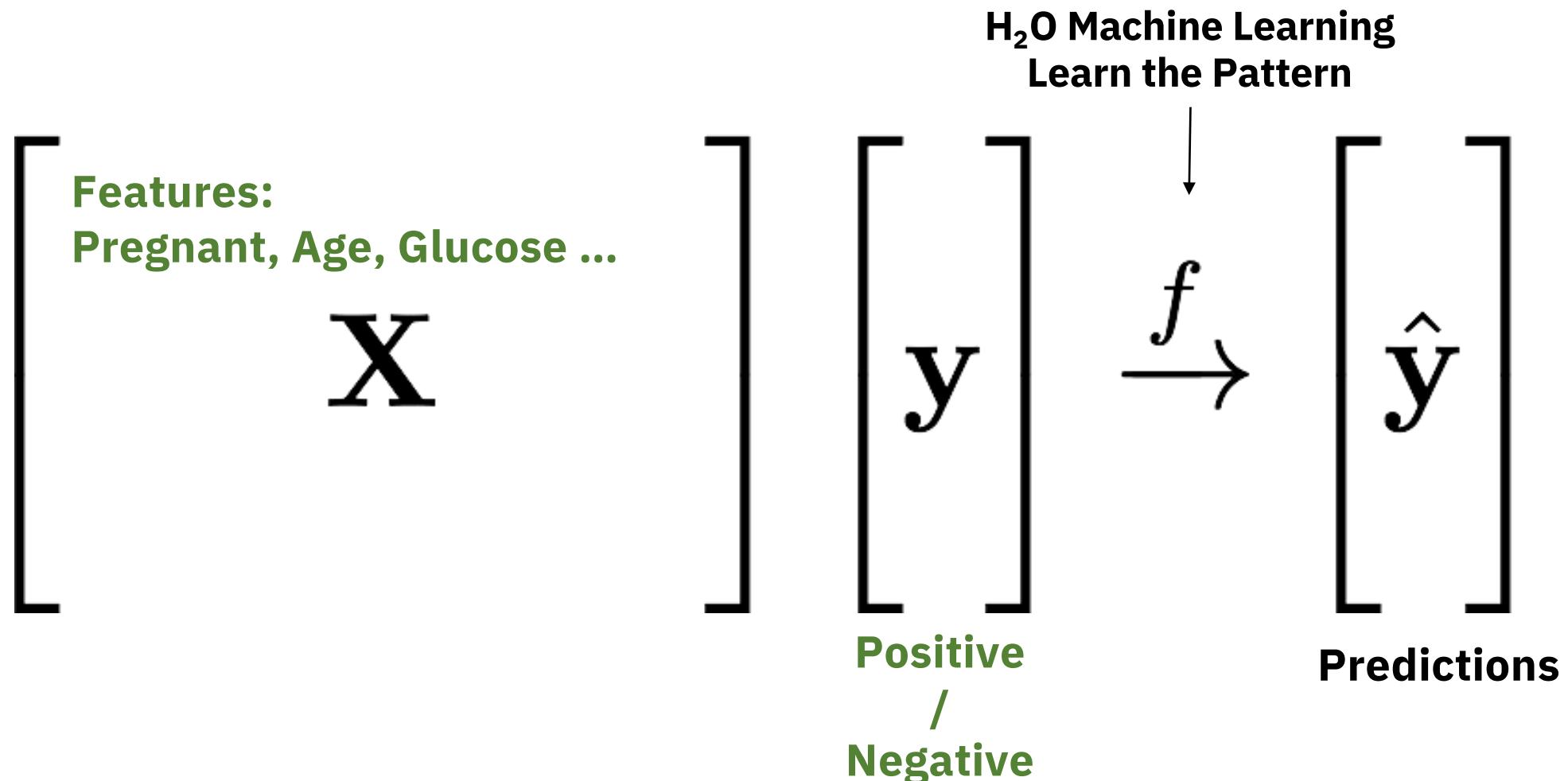
# Learning from **Boston Housing** Data





Time	Topics / Tasks
1:30 – 1:45 pm	Install h2o, lime, mlbench from CRAN slides/code: <a href="https://bit.ly/joe_eRum_2018">bit.ly/joe_eRum_2018</a>
1:45 – 2:00 pm	Introduction (H <sub>2</sub> O, AutoML, LIME)
2:00 – 2:30 pm	Regression Example
2:30 – 3:00 pm	<b>Classification Example</b> <code>\examples\classification_...Rmd</code>
3:00 – 3:30 pm	☕️🍰🍪
3:30 – 3:45 pm	Quick Recap
3:45 – 4:15 pm	Real Use-Case: Moneyball
4:15 – 4:30 pm	Other H <sub>2</sub> O News + Q & A

# Learning from Diabetes Data





@h2oai @matlabulous  
#eRum2018 #AutoML #LIME

Please come back at 3:30pm

Late to the party? Download → [bit.ly/joe\\_eRum\\_2018](https://bit.ly/joe_eRum_2018)



Time	Topics / Tasks
1:30 – 1:45 pm	Install h2o, lime, mlbench from CRAN slides/code: <a href="https://bit.ly/joe_eRum_2018">bit.ly/joe_eRum_2018</a>
1:45 – 2:00 pm	Introduction (H <sub>2</sub> O, AutoML, LIME)
2:00 – 2:30 pm	Regression Example
2:30 – 3:00 pm	Classification Example
3:00 – 3:30 pm	
<b>3:30 – 3:45 pm</b>	<b>Quick Recap</b>
3:45 – 4:15 pm	Real Use-Case: Moneyball
4:15 – 4:30 pm	Other H <sub>2</sub> O News + Q & A



# Why?

- Most users/organizations can benefit from automatic machine learning pipelines.
  - Eliminate time wasted on human errors, debugging etc.
- Model interpretations is crucial for those who must explain their models to regulators or customers.

# You will learn ...

- How to build high quality H<sub>2</sub>O models (almost) automatically.
- How to explain predictions from complex H<sub>2</sub>O models with LIME.
- **Bonus:** A real use-case that led to multimillion-dollar baseball decisions earlier this year.

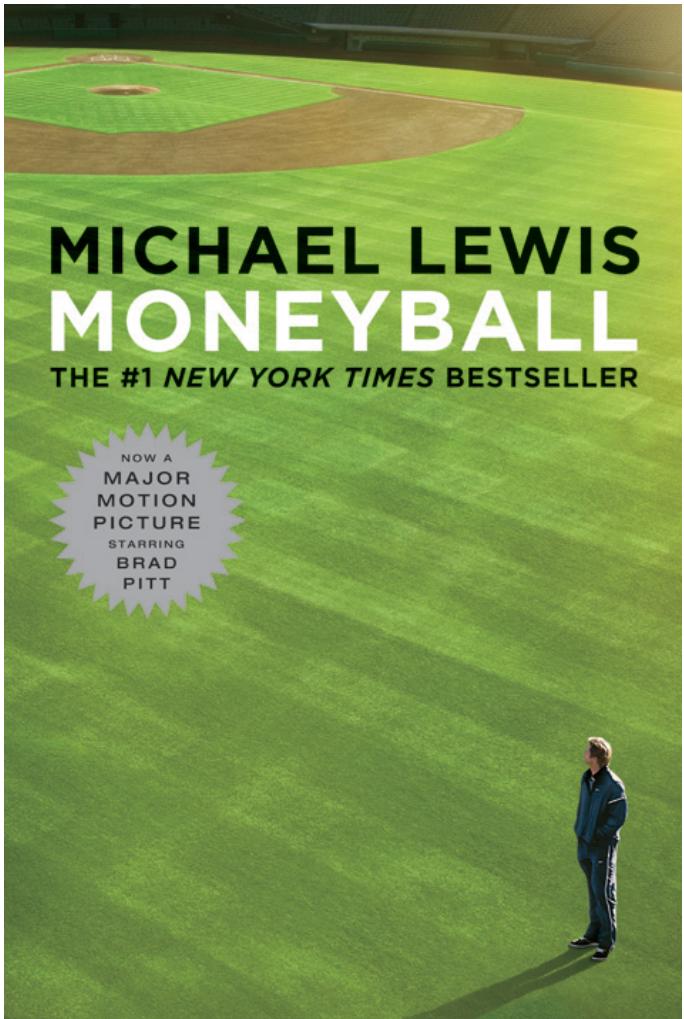


Time	Topics / Tasks
1:30 – 1:45 pm	Install h2o, lime, mlbench from CRAN slides/code: <a href="https://bit.ly/joe_eRum_2018">bit.ly/joe_eRum_2018</a>
1:45 – 2:00 pm	Introduction (H <sub>2</sub> O, AutoML, LIME)
2:00 – 2:30 pm	Regression Example
2:30 – 3:00 pm	Classification Example
3:00 – 3:30 pm	
3:30 – 3:45 pm	Quick Recap
<b>3:45 – 4:15 pm</b>	<b>Real Use-Case: Moneyball</b>
4:15 – 4:30 pm	Other H <sub>2</sub> O News + Q & A

# Making Multimillion-Dollar Decisions with H<sub>2</sub>O AutoML, LIME and Shiny

My journey to a real Moneyball application

# About Moneyball



Billy Beane

Peter Brand  
(based on Paul DePodesta)

# Ari Kaplan – the Real “Moneyball” Guy

- The real characters in the movie (Billy Beane and Paul DePodesta) did not want to work with Hollywood.
- The filmmaker interviewed Ari instead and created the Paul DePodesta character based on Ari’s real-life story.
- Ari happens to work at Aginity so we have a real “Moneyball” guy for this demo.



# A Proof-of-Concept Demo for IBM Think Conference



**Moneyball** [Demo](#)

- [Introduction](#)
- [Results \(Pitching\)](#)
- [Results \(Batting\)](#)
- [About Us](#)
- [YouTube](#)

**Hit a Home Run Making Baseball Decisions Using Artificial Intelligence and Machine Learning**

Thursday, 1:30 PM - 2:10 PM | Session ID: 3456A  
Mandalay Bay South, Level 2 | Breakers C

**IBM + aginity + H<sub>2</sub>O.ai**

Join Ari Kaplan, a real "MoneyBall" and well known around Major League Baseball, Joe Chow, a H2O data scientist, and David Kearns from IBM's Analytics Ecosystem team for this fun, interactive session where you will have the chance to see where artificial intelligence meets business intelligence. Ari and Joe will briefly present the latest machine learning technologies and concepts powering today's baseball decisions, including Hortonworks Data Platform, Spark, Aginity Amp, H2O.ai, IBM Data Science Experience and more. You will then step up to the plate as general manager to see how your player decisions would stack up under World Series pressure. Are you ready to play ball?

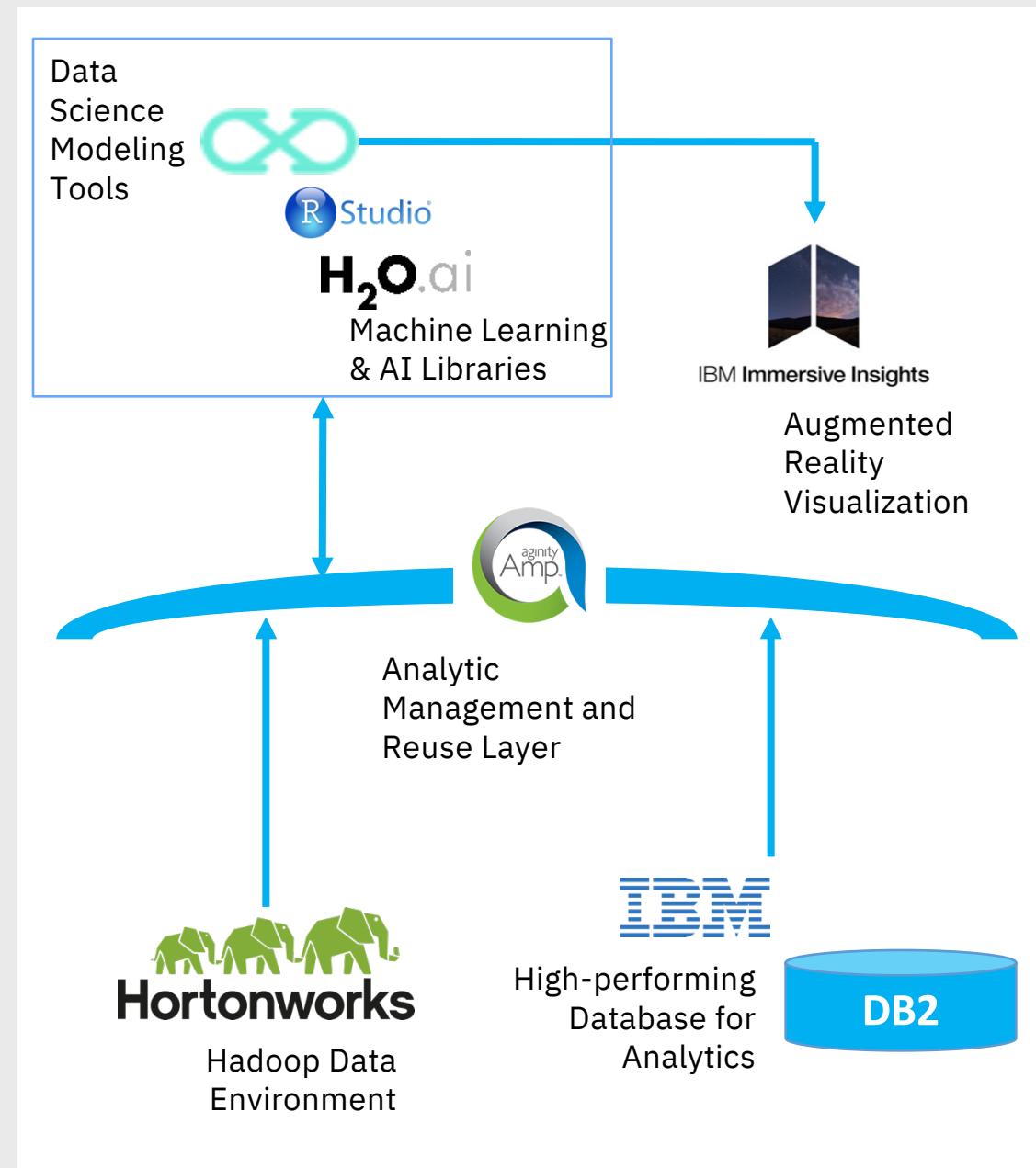
**Speakers:**

- Ari Kaplan, Aginity
- Jo-fai Chow, H2O.ai
- David Kearns, IBM

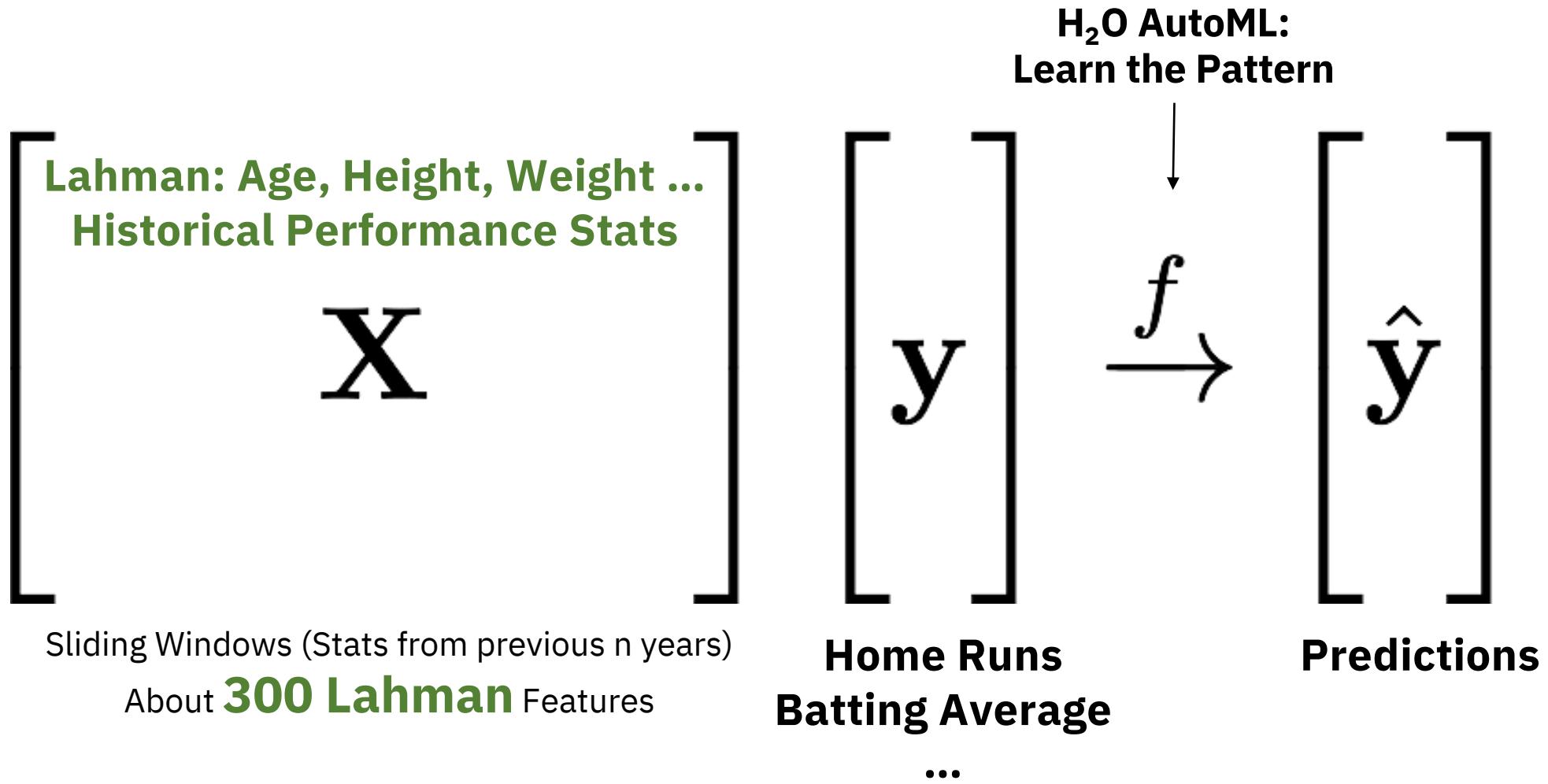
# Enterprise Solution

## The Workflow

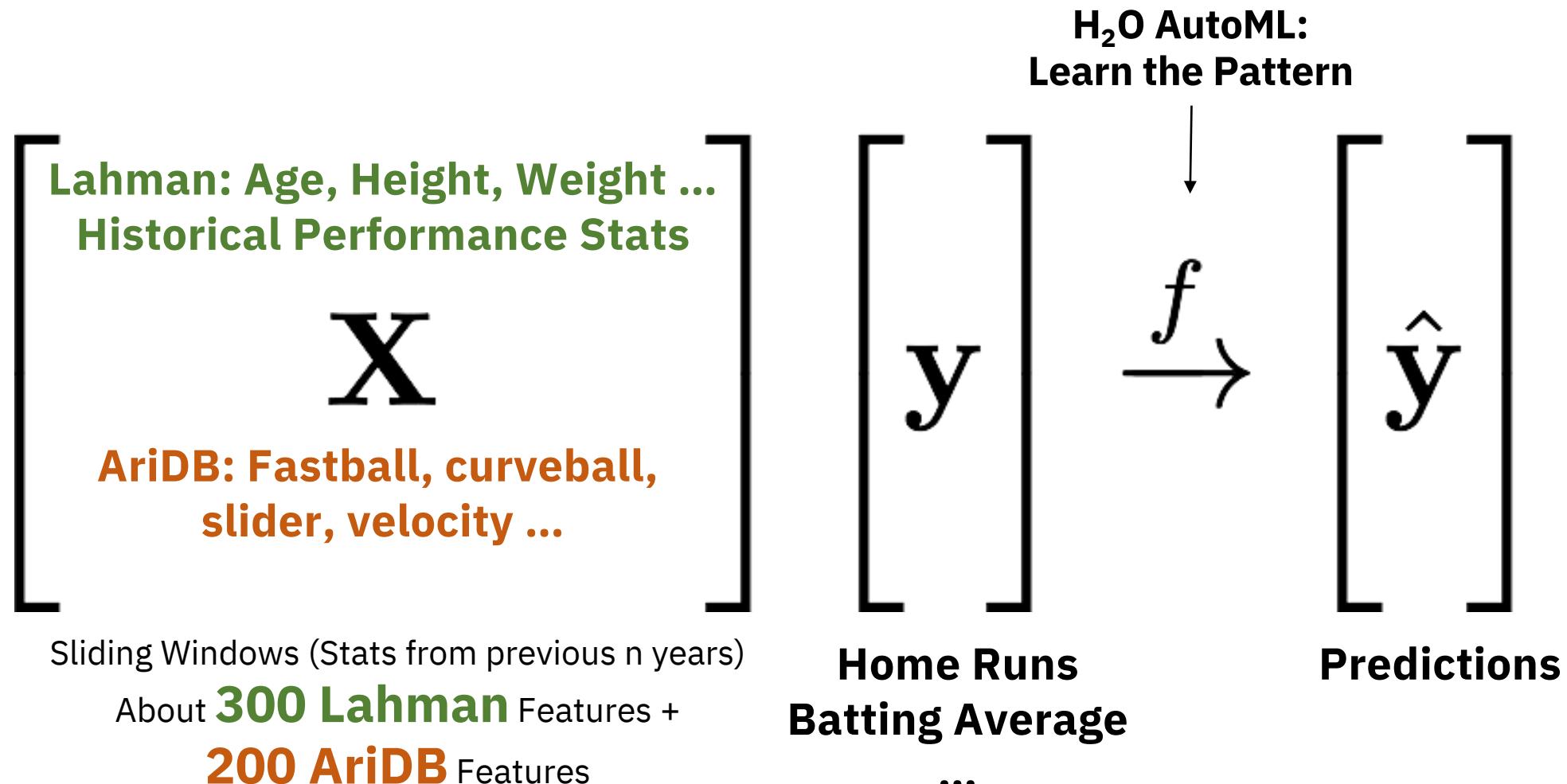
1. Data loaded into the databases
2. Connected diverse data sources to Amp
3. Amp used to create derived attributes and publish them and data to DSX and H<sub>2</sub>O
4. DSX and H<sub>2</sub>O to build and tweak statistical and machine learning models
5. Visualizations tested in Immersive Insights
6. Steps 4 and 5 repeated to get settled data
7. Statistical and machine learning models saved in Amp
8. Data exported to Immersive Insights for final visualizations



# Approach One: Learning from **Lahman** only

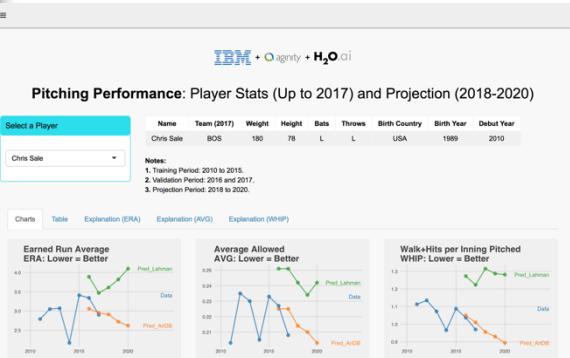
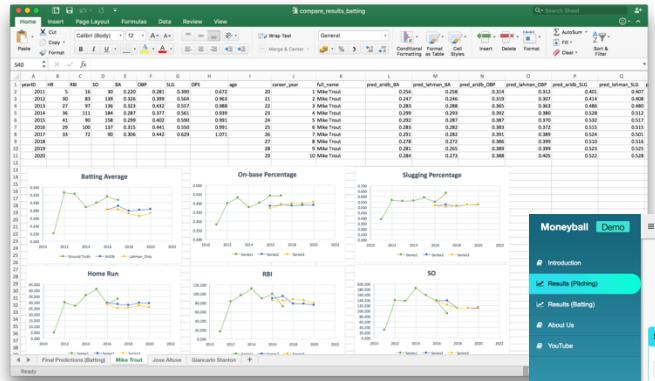


# Approach Two: Learning from **Lahman** & **AriDB**



# Timeline

- **March 19** – AutoML Predictions finalized. Initial presentation in Excel.
- **March 20** – Version 1 of Shiny app. Ari used to app to validate some players he had in mind and recommended one player to his team.
- **March 21** – Multimillion-dollar contract finalized.
- **March 22** – Moneyball presentation at IBM Think



# Presentation Shiny App

**IBM + aginity + H<sub>2</sub>O.ai**

### Pitching Performance: Player Stats (Up to 2017) and Projection (2018-2020)

Select a Player

Name	Team (2017)	Weight	Height	Bats	Throws	Birth Country	Birth Year	Debut Year
Chris Sale	BOS	180	78	L	L	USA	1989	2010

Notes:

1. Training Period: 2010 to 2015.
2. Validation Period: 2016 and 2017.
3. Projection Period: 2018 to 2020.

Charts    Table    Explanation (ERA)    Explanation (AVG)    Explanation (WHIP)

**Green: Predictions based on Lahman only**

**Orange: Predictions based on AriDB + Lahman**

**IBM + aginity + H<sub>2</sub>O.ai**

### Pitching Performance: Player Stats (Up to 2017) and Projection (2018-2020)

Select a Player

Name	Team (2017)	Weight	Height	Bats	Throws	Birth Country	Birth Year	Debut Year
Chris Sale	BOS	180	78	L	L	USA	1989	2010

Notes:

1. Training Period: 2010 to 2015.
2. Validation Period: 2016 and 2017.
3. Projector Period: 2018 to 2020.

Charts    Table    Explanation (ERA)    Explanation (AVG)    Explanation (WHIP)

Data	Year	ERA (Historical Data)	ERA (Predictions based on Ari_DB)	ERA (Predictions based on Lahman)	AVG (Historical Data)	AVG (Predictions based on Ari_DB)	AVG (Predictions based on Lahman)	WHIP (Historical Data)	WHIP (Predictions based on Ari_DB)	WHIP (Predictions based on Lahman)
Training	2011	2.790		0.203				1.113		
Training	2012	3.050		0.235				1.135		
Training	2013	3.070		0.230				1.073		
Training	2014	2.170		0.205				0.966		
Training	2015	3.410		0.233				1.088		
Validation	2016	3.340	3.060	0.227	0.225	0.251	1.037	1.050	1.273	
Validation	2017	2.900	2.950	0.208	0.225	0.251	0.970	1.010	1.223	
Prediction	2018		2.910	3.610	0.214	0.242	0.956	1.315		
Prediction	2019		2.720	3.620	0.210	0.234	0.950	1.287		
Prediction	2020		2.620	4.100	0.203	0.242	0.894	1.281		

**IBM + aginity + H<sub>2</sub>O.ai**

### Pitching Performance: Player Stats (Up to 2017) and Projection (2018-2020)

Select a Player

Name	Team (2017)	Weight	Height	Bats	Throws	Birth Country	Birth Year	Debut Year
Chris Sale	BOS	180	78	L	L	USA	1989	2010

Notes:

1. Training Period: 2010 to 2015.
2. Validation Period: 2016 and 2017.
3. Projector Period: 2018 to 2020.

Charts    Table    Explanation (ERA)    Explanation (AVG)    Explanation (WHIP)

**IBM + aginity + H<sub>2</sub>O.ai**

### Batting Performance: Player Stats (Up to 2017) and Projection (2018-2020)

Select a Player

Name	Team (2017)	Weight	Height	Bats	Throws	Birth Country	Birth Year	Debut Year
Giancarlo Stanton	MIA	245	78	R	R	USA	1989	2010

Notes:

1. Training Period: 2010 to 2015.
2. Validation Period: 2016 and 2017.
3. Projector Period: 2018 to 2020.

Charts (1/2)    Charts (2/2)    Table (1/2)    Table (2/2)    Exp. (BA)    Exp. (HR)    Exp. (RBI)    Exp. (OBP)    Exp. (SLG)    Exp. (SO)

# Acknowledgement



9:04 PM - 22 Mar 2018





Time	Topics / Tasks
1:30 – 1:45 pm	Install <code>h2o</code> , <code>lime</code> , <code>mlbench</code> from CRAN slides/code: <a href="https://bit.ly/joe_eRum_2018">bit.ly/joe_eRum_2018</a>
1:45 – 2:00 pm	Introduction ( $H_2O$ , AutoML, LIME)
2:00 – 2:30 pm	Regression Example
2:30 – 3:00 pm	Classification Example
3:00 – 3:30 pm	
3:30 – 3:45 pm	Quick Recap
3:45 – 4:15 pm	Real Use-Case: Moneyball
4:15 – 4:30 pm	Other $H_2O$ News + Q & A

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



In-Memory, Distributed  
Machine Learning Algorithms  
with H2O Flow GUI



H2O AI Open Source Engine  
Integration with Spark



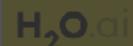
Lightning Fast machine  
learning on GPUs

DRIVERLESSAI

Automatic feature  
engineering, machine  
learning and interpretability

# Steam

Secure multi-tenant H2O clusters

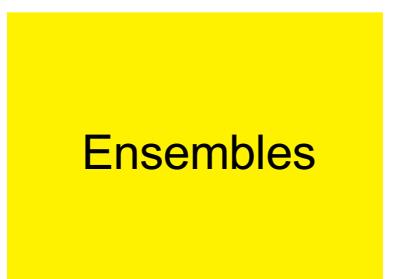


# Algorithms on H<sub>2</sub>O-3 (CPU)

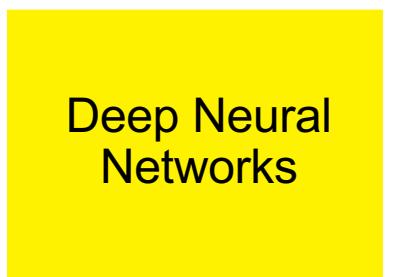
## Supervised Learning



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



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

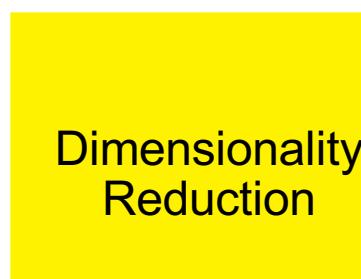


- Deep learning: Create multi-layer feed forward neural networks starting with an input layer followed by multiple layers of nonlinear transformations

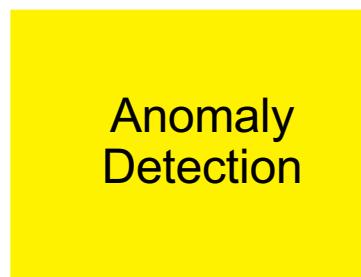
## Unsupervised Learning



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



- 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



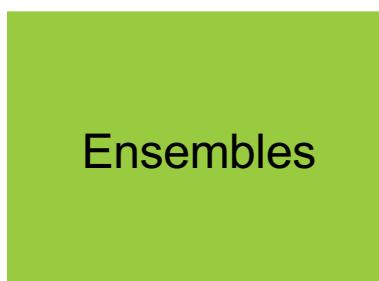
- Autoencoders: Find outliers using a nonlinear dimensionality reduction using deep learning

# Algorithms on H<sub>2</sub>O4GPU (more to come)

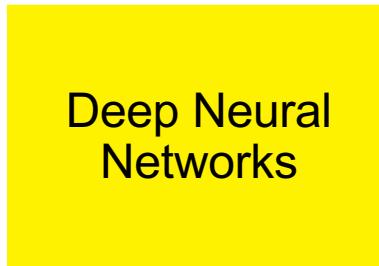
## Supervised Learning



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



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

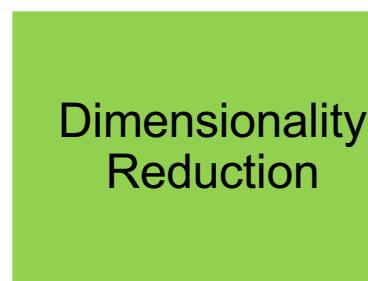


- Deep learning: Create multi-layer feed forward neural networks starting with an input layer followed by multiple layers of nonlinear transformations

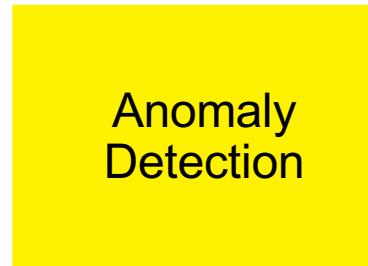
## Unsupervised Learning



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



- 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



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

# H2O4GPU now available in R

BY ERIN LEDELL ON MARCH 27, 2018 – 0 COMMENTS

In September, H2O.ai released a new open source software project for GPU machine learning called [H2O4GPU](#). The initial release (blog post [here](#)) included a Python module with a scikit-learn compatible API, which allows it to be used as a drop-in replacement for scikit-learn with support for GPUs on selected (and ever-growing) algorithms. We are proud to announce that the same collection of GPU algorithms is now available in R, and the `h2o4gpu` R package is available on [CRAN](#).



<https://github.com/h2oai/h2o4gpu>

# From Kaggle Grand Masters' Recipes to Production Ready in a Few Clicks

BY JO-FAI CHOW ON MAY 9, 2018 – 0 COMMENTS – EDIT

## Introducing Accelerated Automatic Pipelines in H2O Driverless AI

At H2O, we work really hard to make machine learning fast, accurate, and accessible to everyone. With H2O Driverless AI, users can leverage years of world-class, [Kaggle Grand Masters](#) experience and our GPU-accelerated algorithms ([H2O4GPU](#)) to produce top quality predictive models in a fully automatic and timely fashion.

In our most recent release (version 1.1), we are going one step further to streamline the deployment process with MOJO (Model ObjEcT, Optimized). Inherited from our popular H2O-3 platform, MOJO is a highly optimized, low-latency scoring engine that is easily embeddable in any Java environment. With automatic pipeline generation in Driverless AI, users can go from automatic machine learning to production ready in just a few clicks. This blog post illustrates the usage of MOJO in Driverless AI with a simple example.

### Easing the Pain Points in a Machine Learning Workflow

In a typical enterprise machine learning workflow, there are many things that could go wrong due to human errors, bad data science practices, different tools/infrastructure, incompatible code, lack of testing, versioning, communication and so on.

blog.h2o.ai

# Thanks!

- Organizers & Sponsors



- Code, Slides & Documents

- [bit.ly/joe\\_eRum\\_2018](http://bit.ly/joe_eRum_2018)
- [docs.h2o.ai](http://docs.h2o.ai)

- Contact

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

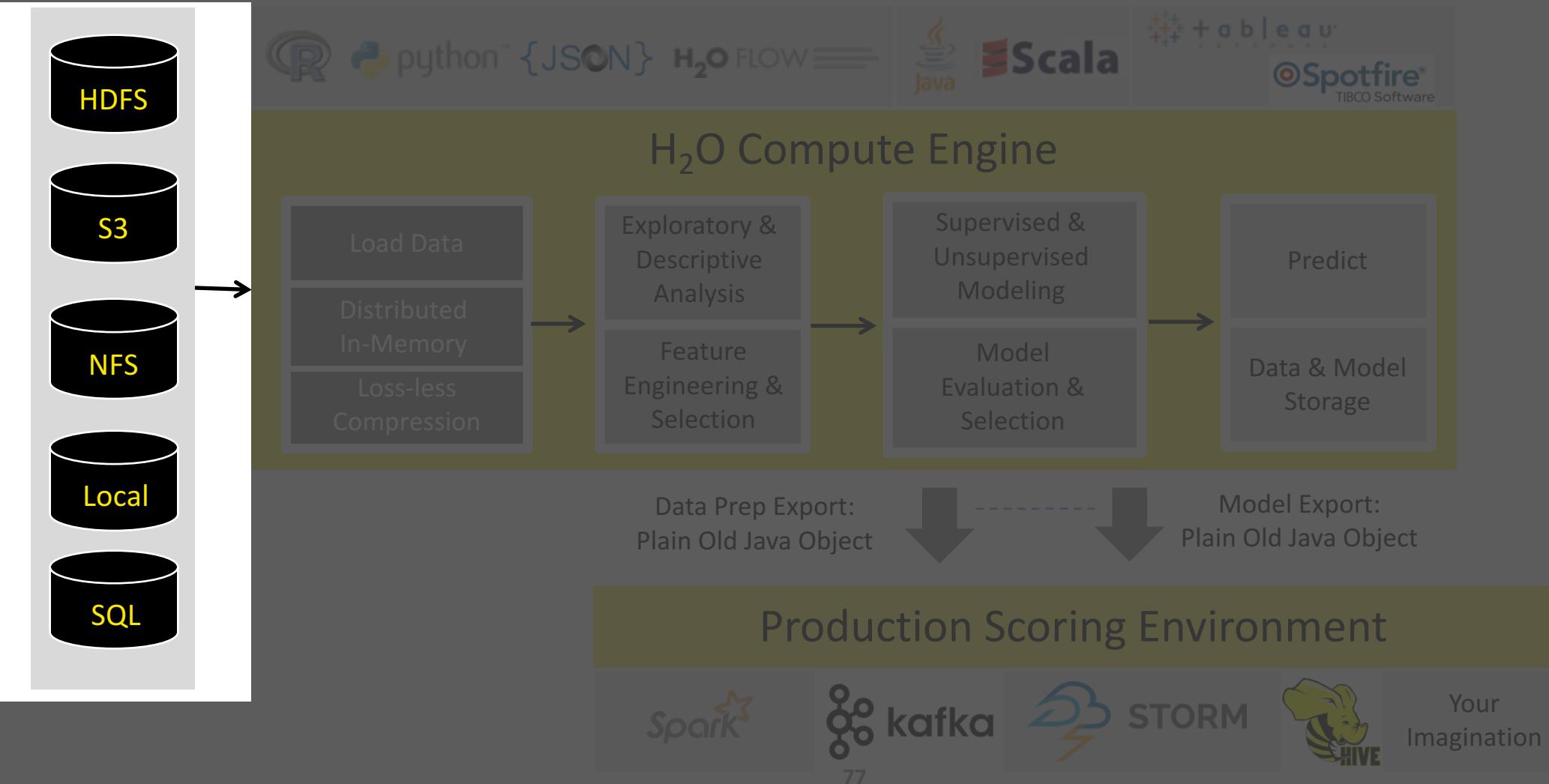
- Please search/ask questions on  
**Stack Overflow**

- Use the tag `h2o` (not h2 zero)

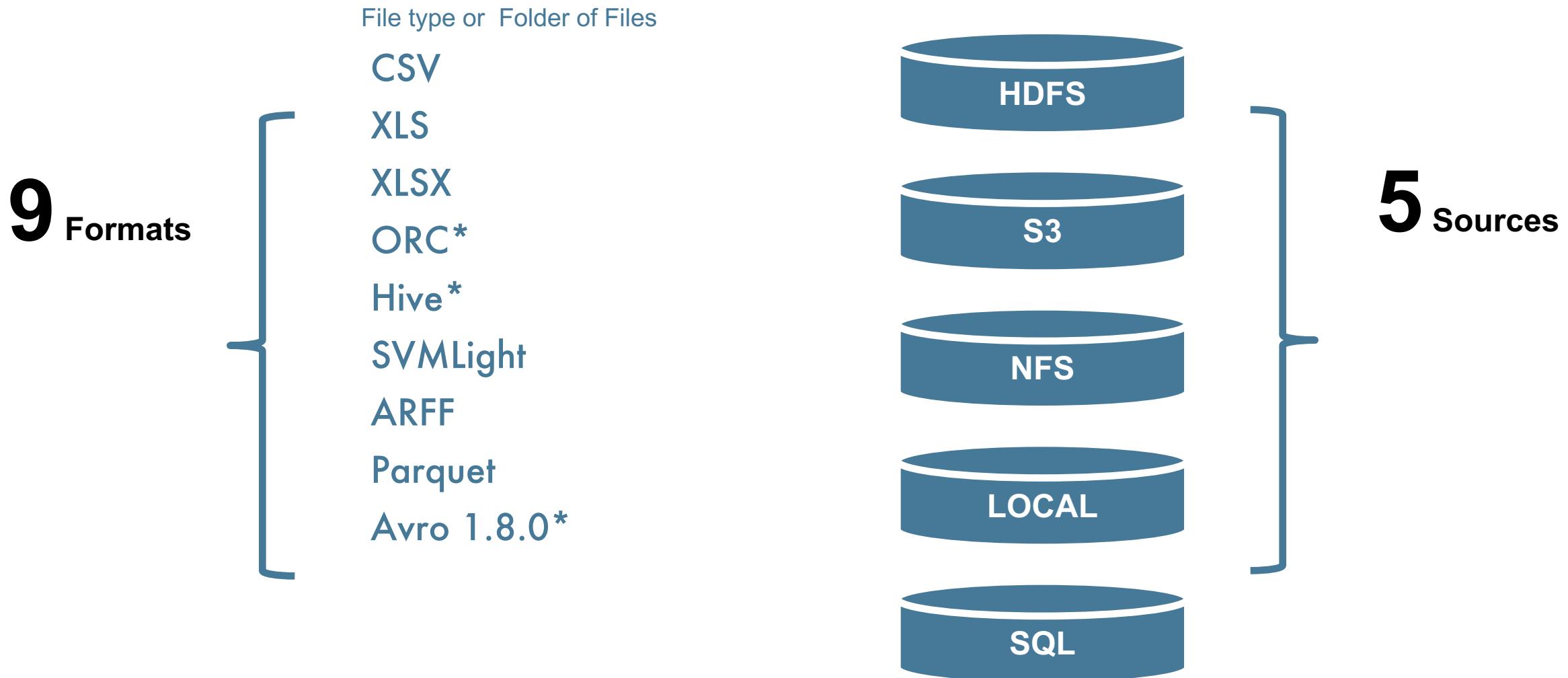
# Appendix

# High Level Architecture

Import Data from  
Multiple Sources



# Supported Formats & Data Sources



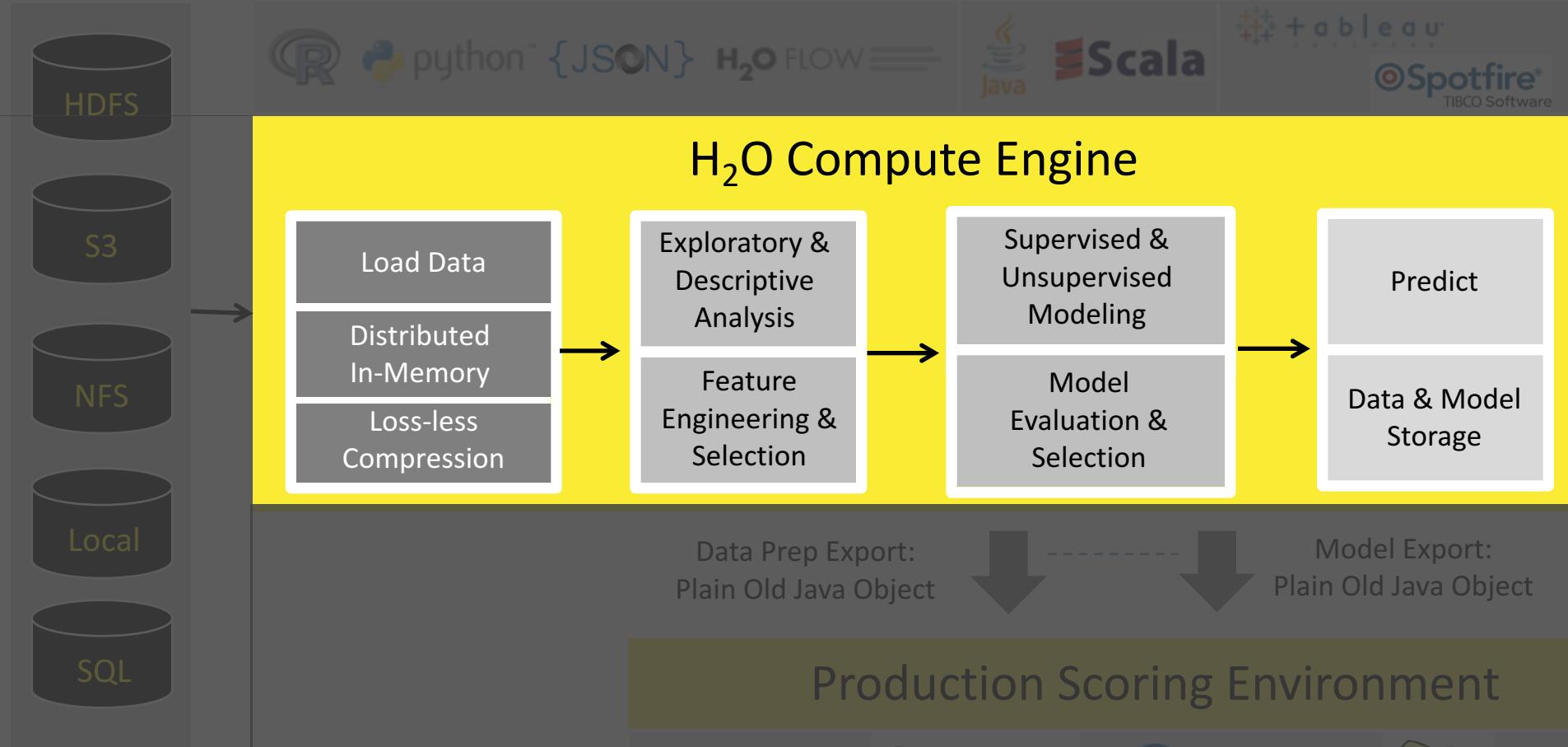
\* 1. only if H2O is running as a Hadoop job

\* 2. Hive files that are saved in ORC format

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

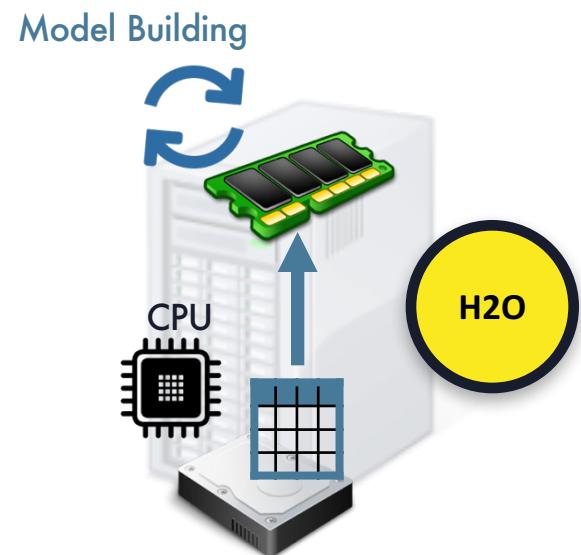
# High Level Architecture

Fast, Scalable & Distributed  
Compute Engine Written in  
Java



Your  
Imagination

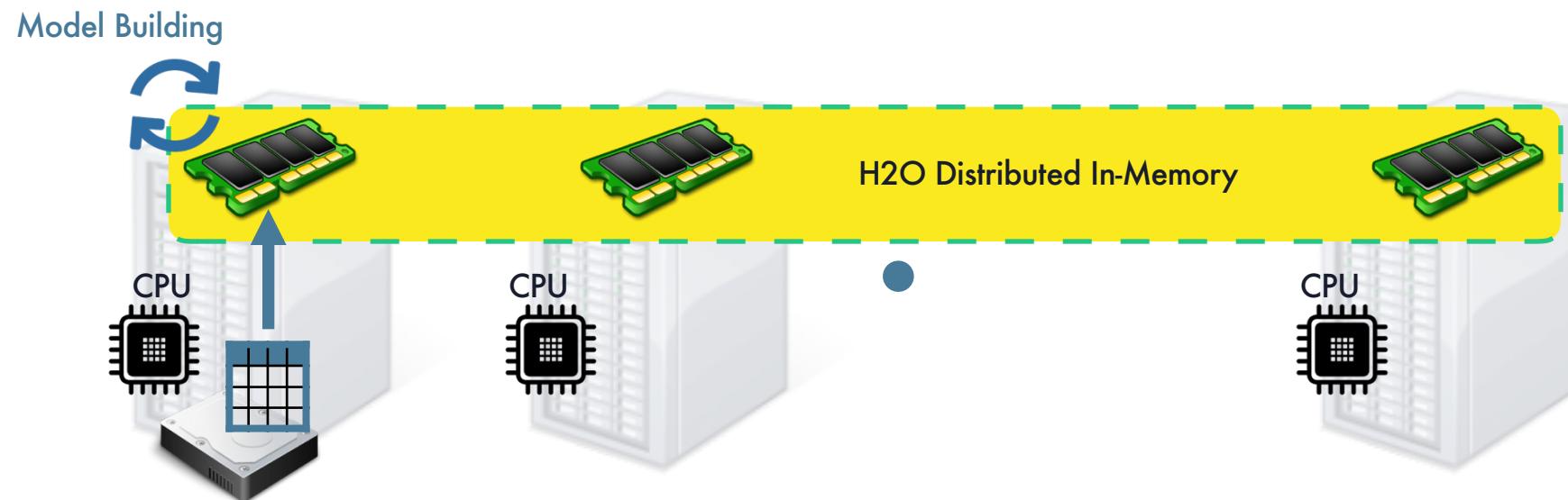
# H<sub>2</sub>O Core



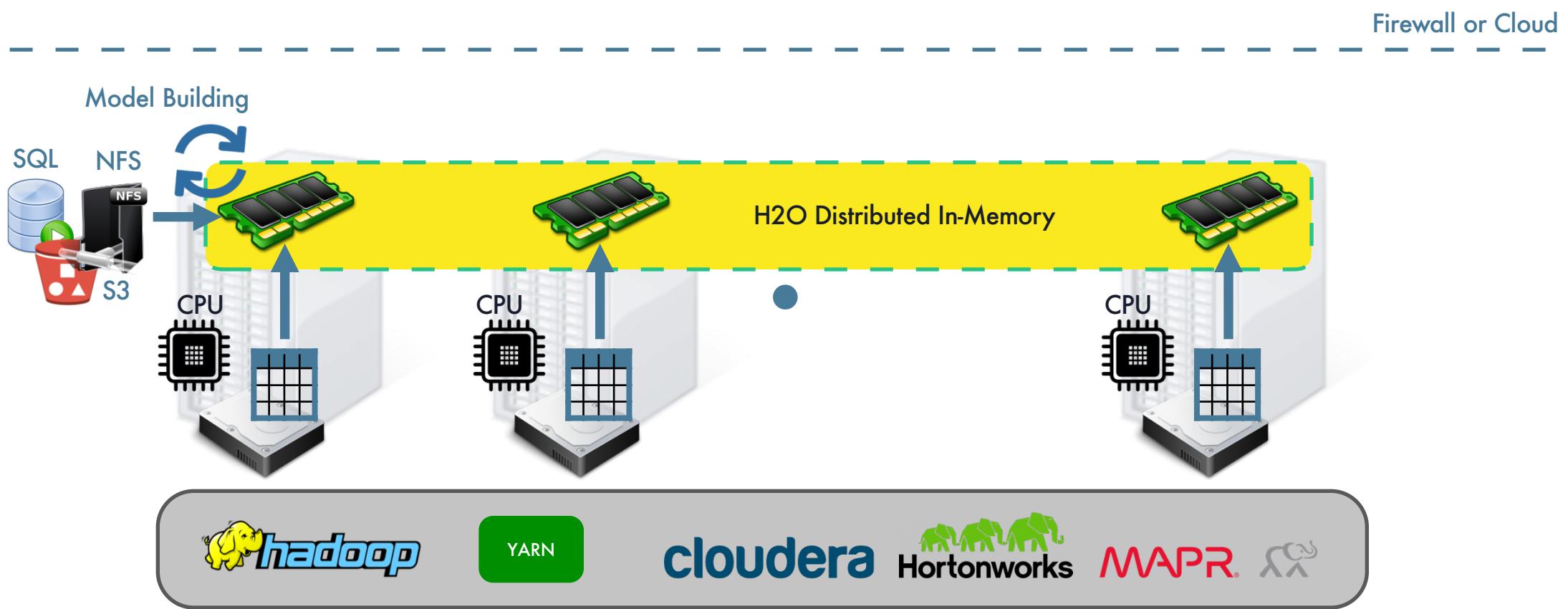
# H<sub>2</sub>O Core



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

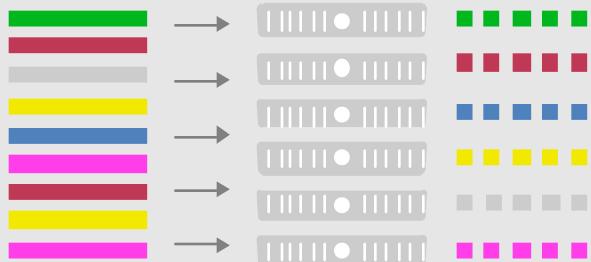


# H<sub>2</sub>O Core

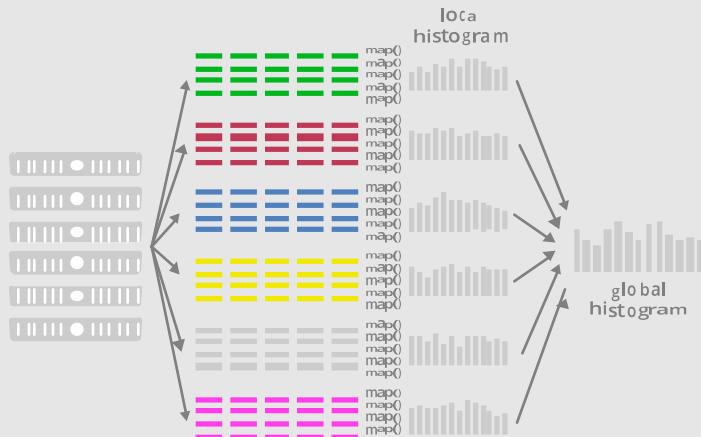


# Distributed Algorithms

## Foundation for Distributed Algorithms



Parallel Parse into **Distributed Rows**



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

## Advantageous Foundation

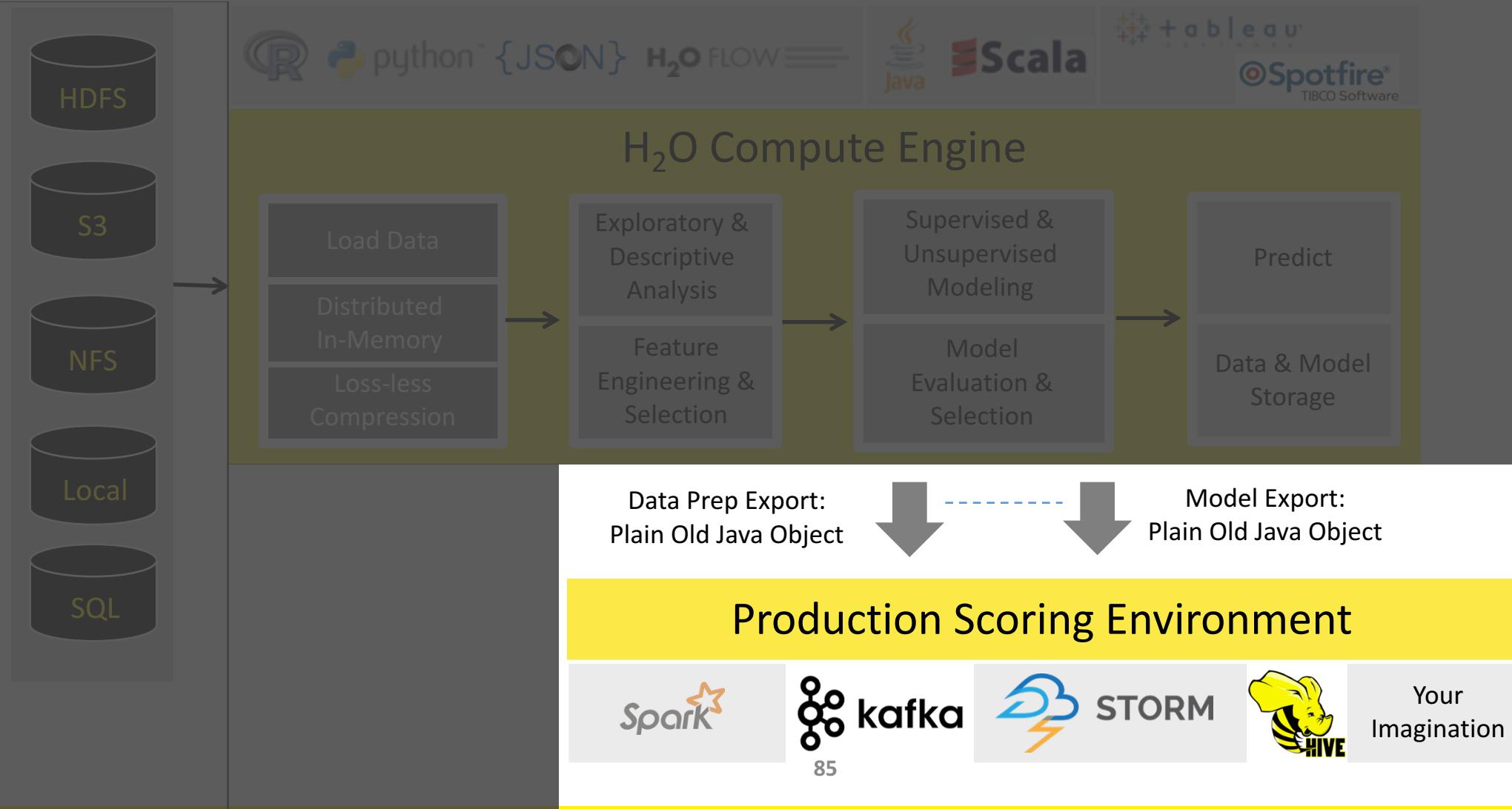
- Foundation for In-Memory Distributed Algorithm Calculation - **Distributed Data Frames** and **columnar compression**
- All algorithms are distributed in H<sub>2</sub>O: GBM, GLM, DRF, Deep Learning and more. Fine-grained map-reduce iterations.
- **Only enterprise-grade, open-source distributed algorithms 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**

# High Level Architecture

Export Standalone Models  
for Production



# H<sub>2</sub>O 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

**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 User Guide** 2.3 2.2 2.1  
Sparkling Water Booklet  
RSparkling Readme  
PySparkling User Guide 2.3 2.2 2.1  
Recent Changes 2.3 2.2 2.1  
Open Source License (Apache V2)

Quick Start Video - Scala

**Download Sparkling Water**

**Driverless AI**

What is Driverless AI?  
Driverless AI User Guide HTML PDF  
Recent Changes  
Driverless AI Booklet  
MLI with Driverless AI Booklet

Quick Start Video - Downloading Driverless AI  
Quick Start Video - Launching an Experiment  
Driverless AI Webinars

**Download Driverless AI**

**H<sub>2</sub>O4GPU (alpha)**

H<sub>2</sub>O4GPU Readme  
Open Source License (Apache V2)

**Download H<sub>2</sub>O4GPU**

**URL: docs.h2o.ai**

# Demo: $H_2O$ on a 320-Core Hadoop Cluster

(Web Interface)



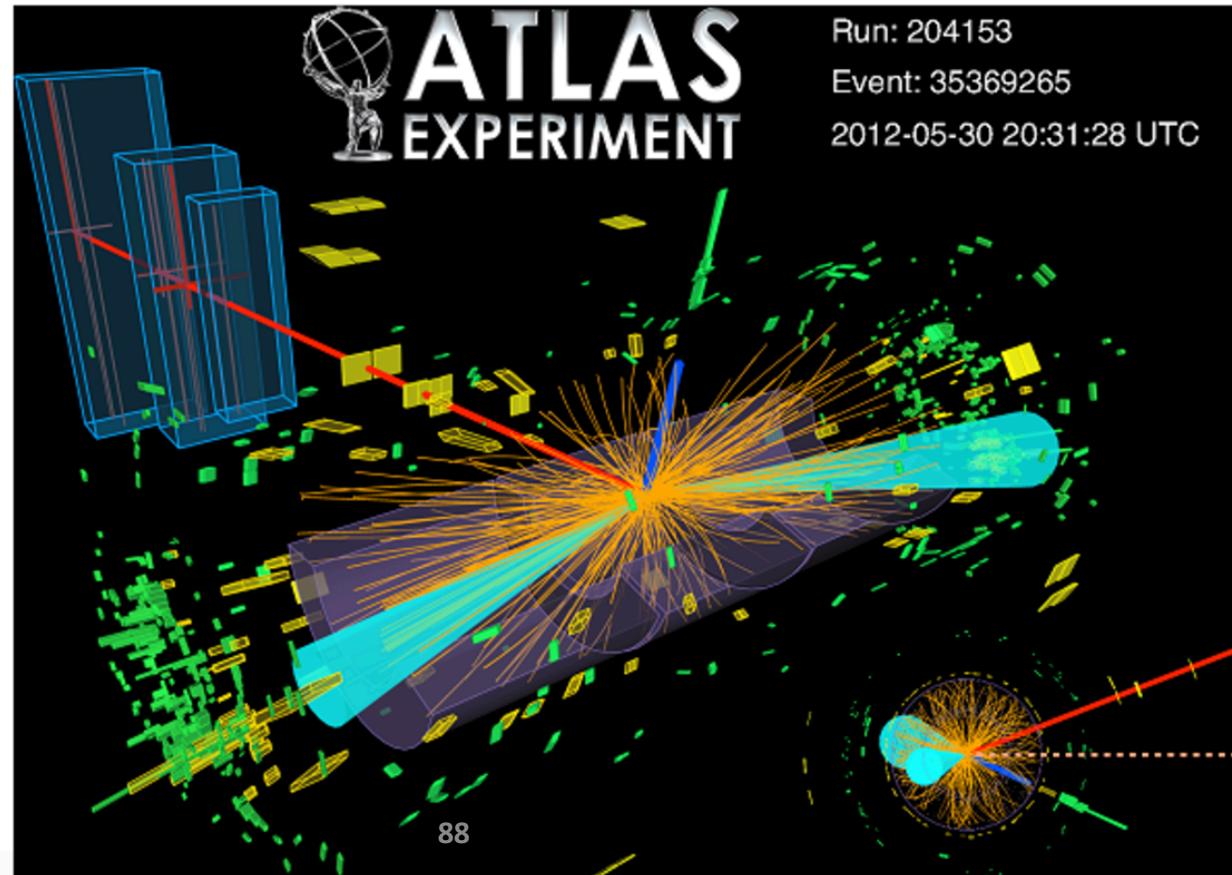
## Higgs Boson Machine Learning Challenge

Use the ATLAS experiment to identify the Higgs boson

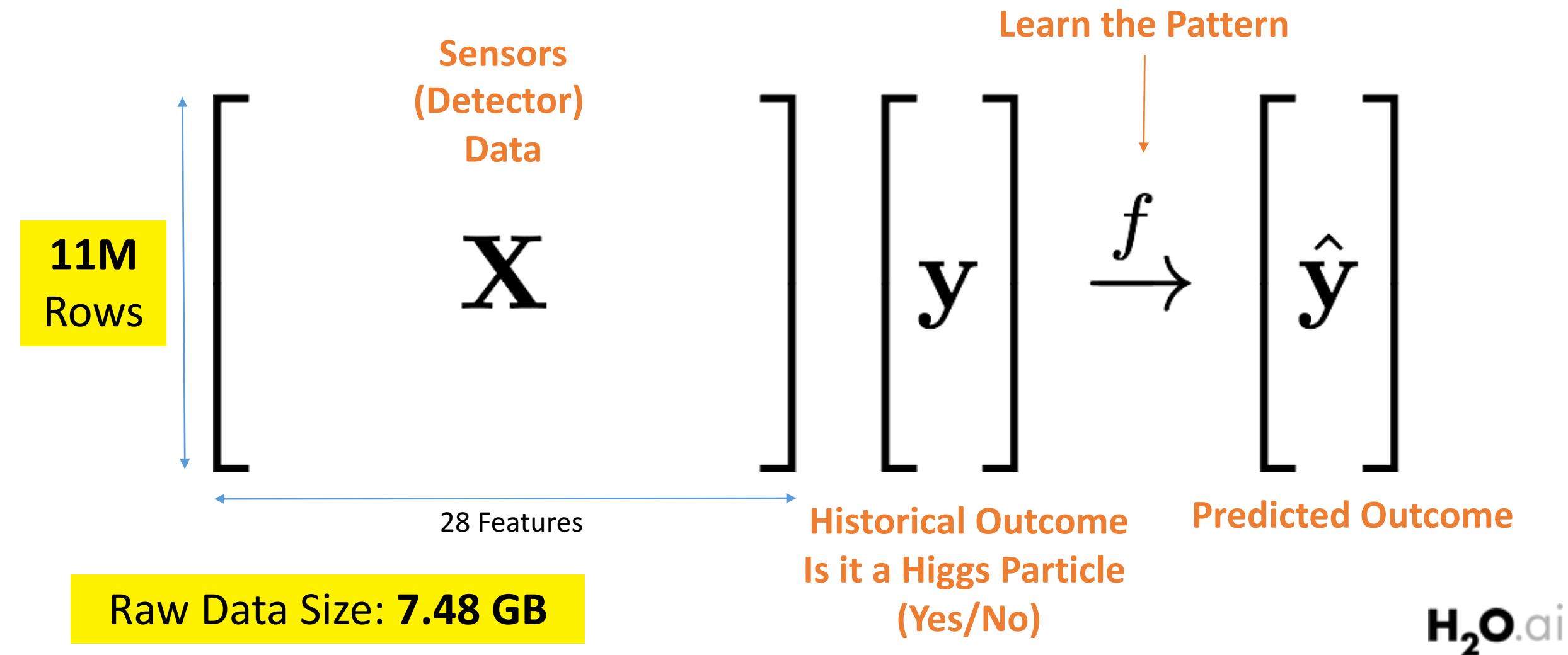
\$13,000 · 1,785 teams · 3 years ago

[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**

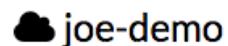
label	type	Missing	Zeros	+Inf	-Inf	min	max	mean	sigma	cardinality	Actions
C1	enum	0	5170877	0	0	0	1.0	0.5299	0.4991	2	<a href="#">Convert to numeric</a>
C2	real	0	0	0	0	0.2747	12.0989	0.9915	0.5654	..	..
C3	real	0	0	0	0	-2.4350	2.4349	-0.0	1.0088	..	..
C4	real	0	0	0	0	-1.7425	1.7432	-0.0	1.0063	..	..
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

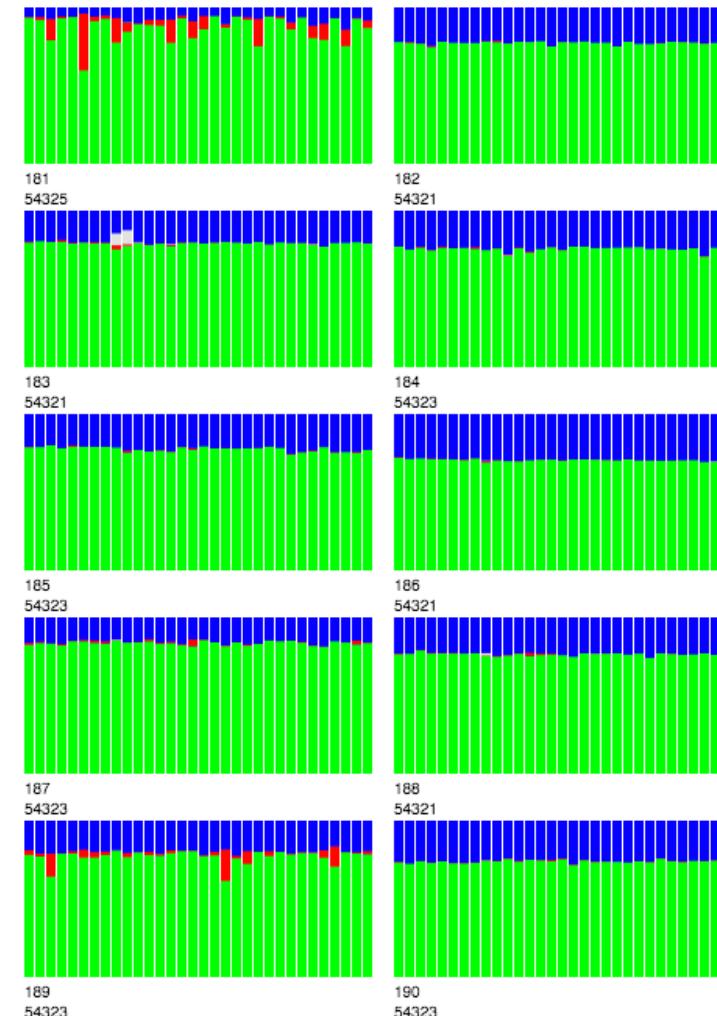
Name	Ping	Cores	Load	My CPU %	Sys	Shut Down	Data (Used/Total)	Data (% Cached)	GC (Free / Total / Max)	Disk (Free / Max)	Disk (% Free)
✓ 172.16.2.181:54323	a few seconds ago	32	6.110	0	8	-	40.603	33.82 GB / s	29.46 GB / NaN undefined / 29.58 GB	339.08 GB / 1.70 TB	19%
✓ 172.16.2.182:54321	a few seconds 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%

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

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

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