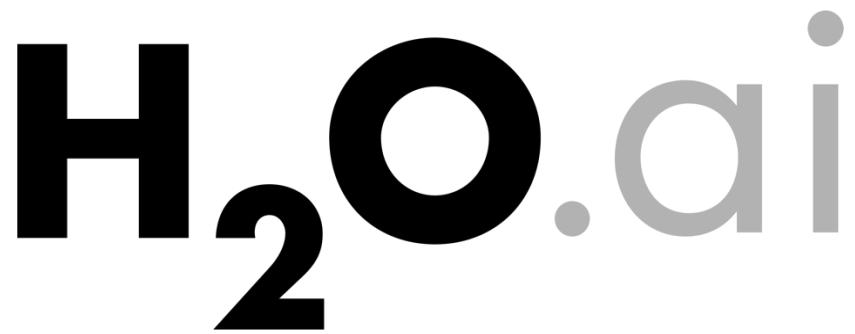


# 使用 H<sub>2</sub>O 进行机器学习

- 概观
  - 公司 / 人员
- 开源的H2O机器学习平台 – H2O-3
  - 概观
  - 信用卡风险示例
- 商业的H2O机器学习平台 - Driverless AI
  - 概观
  - 信用卡风险示例
- Q & A 提问



Jo-fai (Joe) Chow 周祖輝  
数据科学家  
[joe@h2o.ai](mailto:joe@h2o.ai)

# 公司简介

|                    |  |
|--------------------|--|
| <b>Founded 成立</b>  | 2011年获得风险投资支持，在2012年首次亮相   |
| <b>Products 产品</b> | <ul style="list-style-type: none"><li>• H<sub>2</sub>O Open Source Machine Learning Platform 开源的H2O机器学习平台</li><li>• Sparkling Water (H<sub>2</sub>O + Spark)</li><li>• Enterprise Steam</li><li>• Driverless AI (商业的H2O机器学习平台)</li></ul> |
| <b>Mission 任务</b>  | Operationalize Data Science, and provide a platform for users to build beautiful data products<br>实现数据科学的操作化，为用户构建美丽的数据产品提供平台  |
| <b>Team 团队</b>     | 75 人 <ul style="list-style-type: none"><li>• 分布式系统工程师</li><li>• 机器学习专家</li><li>• 世界级的可视化设计师</li></ul>  |
| <b>HQ 总部</b>       | Mountain View, California 加州山景城  |



# Scientific Advisory Council

## 科学咨询委员会



### 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*



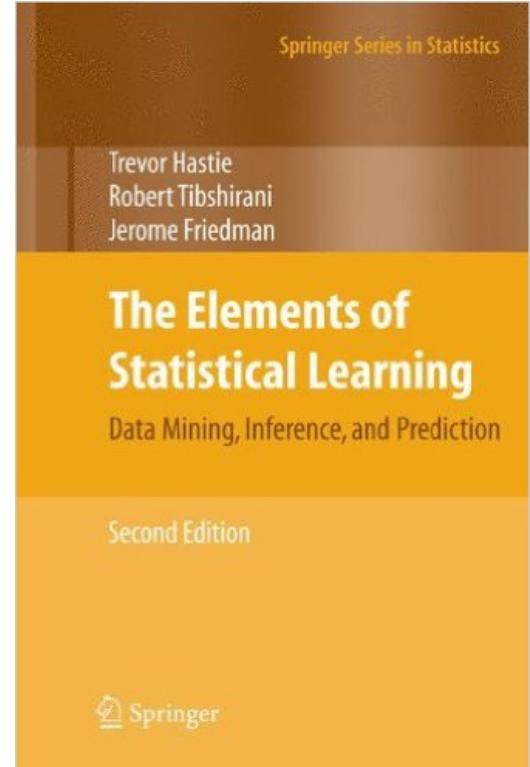
### Dr. Robert Tibshirani

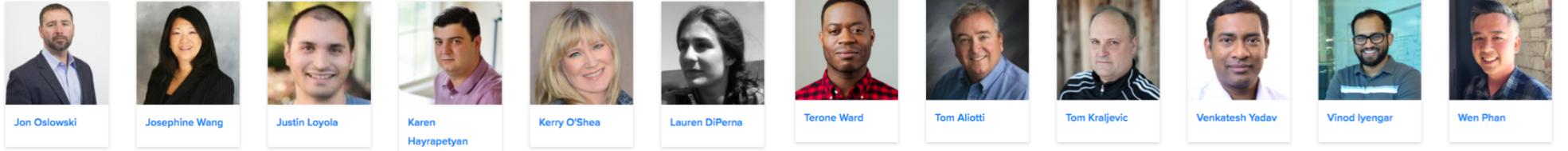
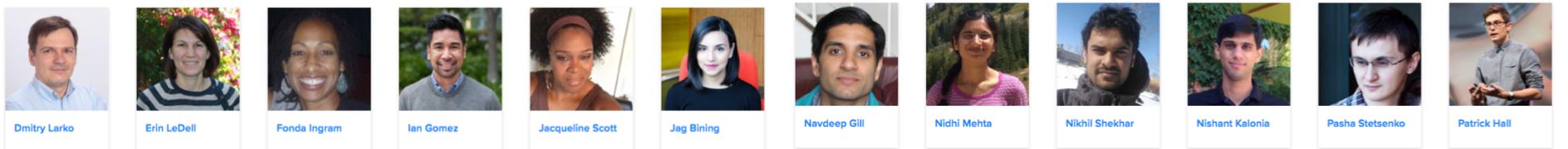
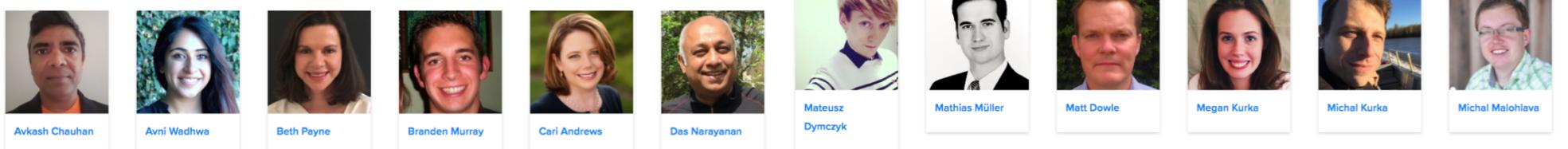
- Professor of Statistics and Health Research and Policy, Stanford University
- PhD in Statistics, Stanford University
- Co-author, *The Elements of Statistical Learning: Prediction, Inference and Data Mining*
- Author, *Regression Shrinkage and Selection via the Lasso*
- Co-author, *An Introduction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



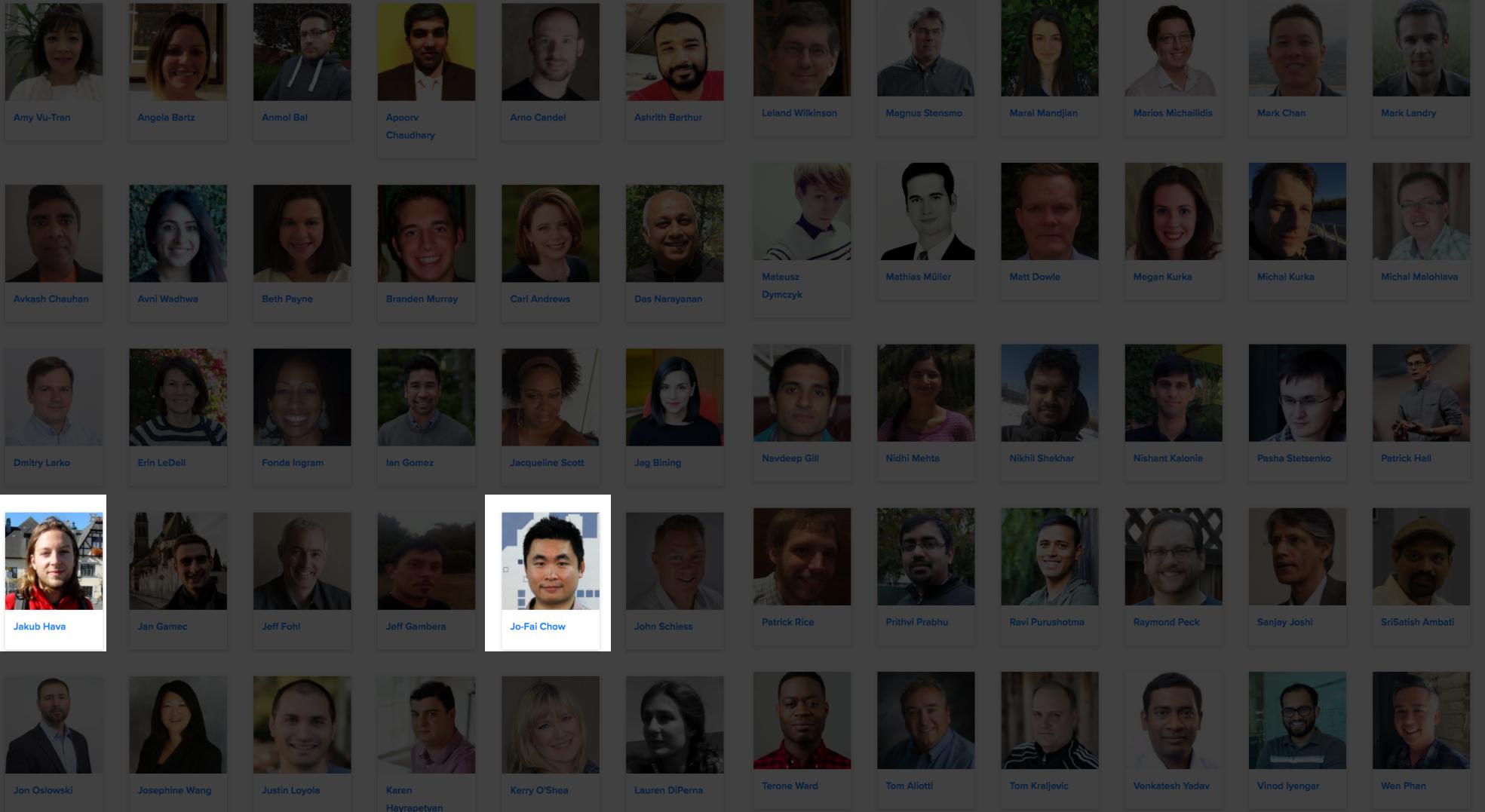
### Dr. Steven Boyd

- 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 团队

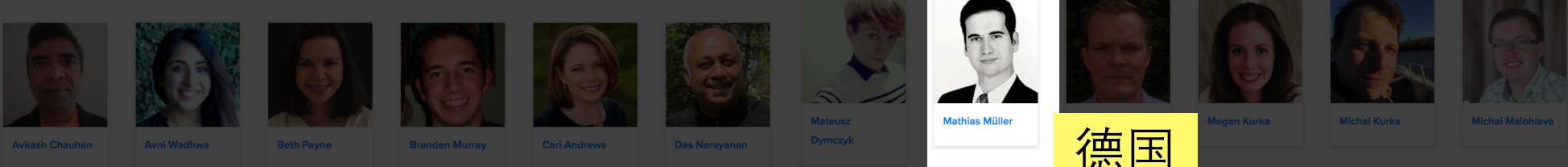


# H<sub>2</sub>O 欧洲团队(2016)

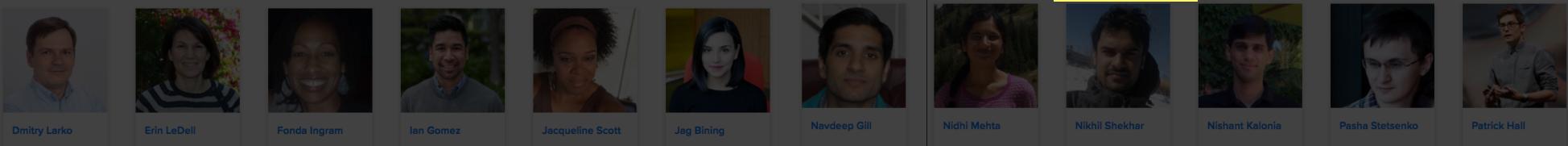




英国



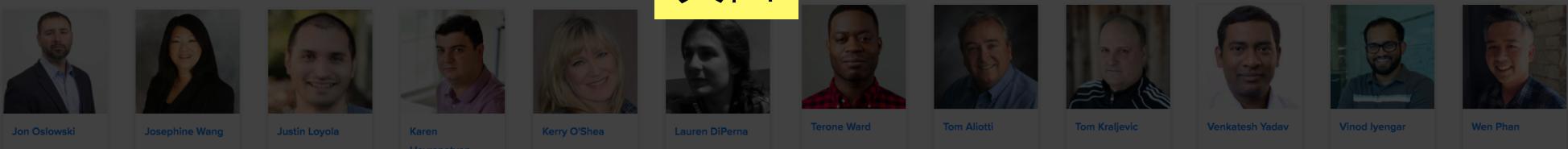
德国



捷克



英国



H<sub>2</sub>O 欧洲团队(2017)



Wendy Wong

# What is Joe's role at H<sub>2</sub>O.ai 我的职责



- Data Scientist (数据科学家) / Sales Engineer (销售) / Meetup Organiser (Meetup组织者) (on paper)
- Unofficial Photographer of H2O.ai SWAG (非官方摄影师) (the travelling data scientist)

# #AroundTheWorldWithH2Oai

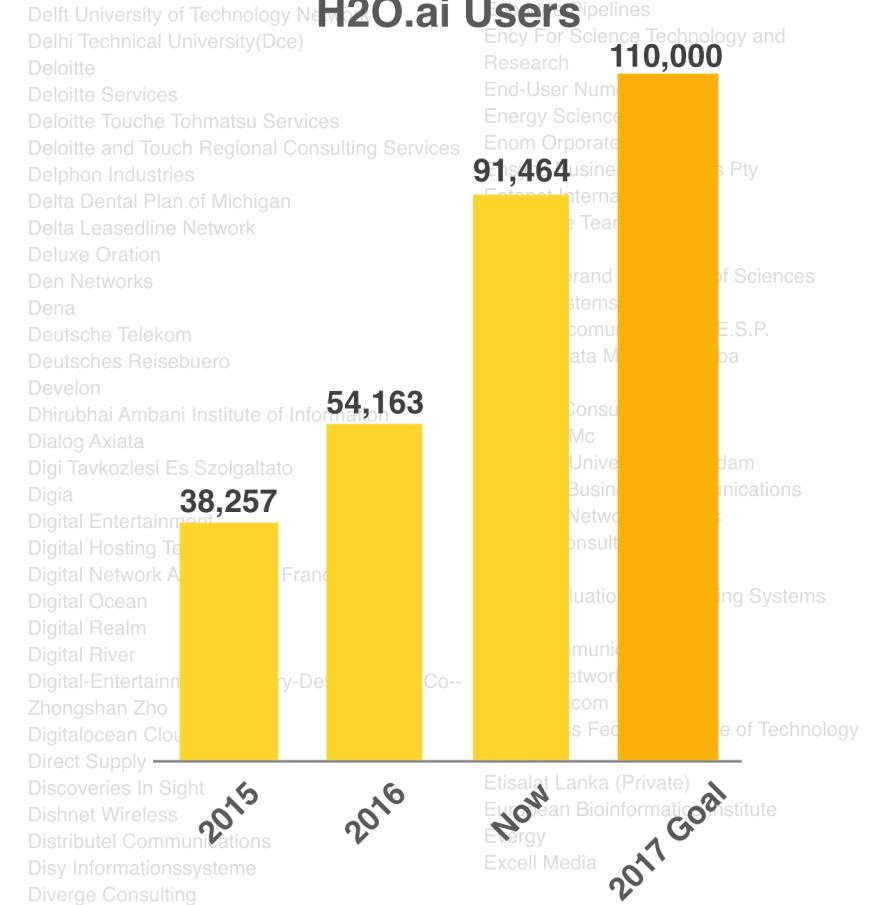
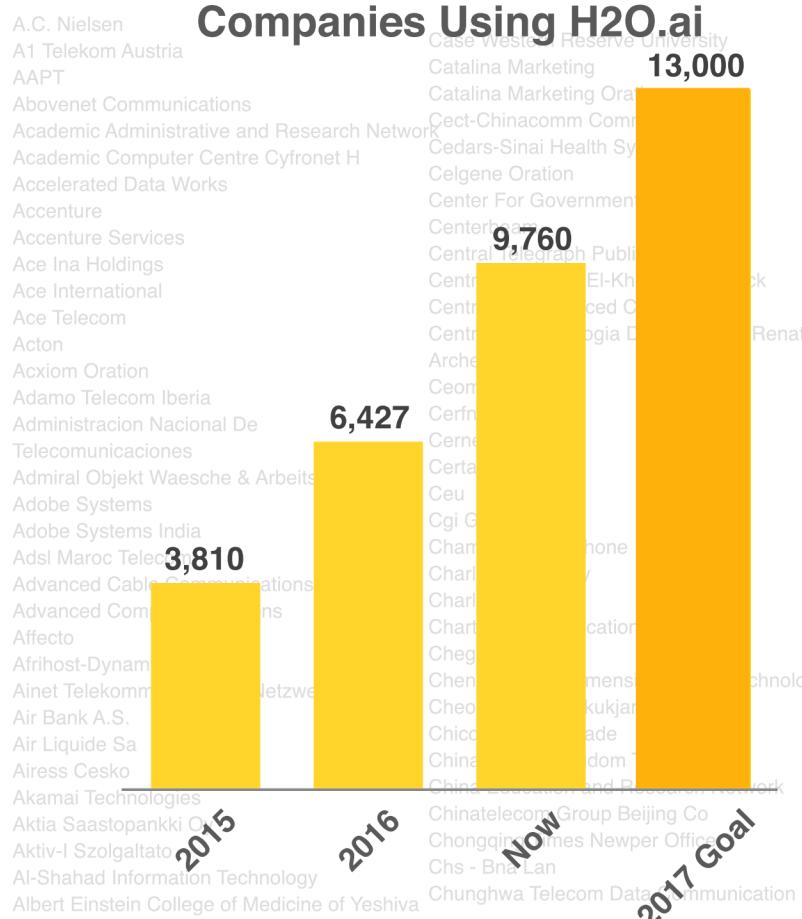
| Month | Joe's H <sub>2</sub> O Events in 2017                          |
|-------|--|
| Jan   | <b>London</b>  |
| Feb   | <b>London</b> , Warsaw, Oxford                                 |
| Mar   | Bay Area, <b>London</b> , Cologne, Barcelona, Madrid, Vienna   |
| Apr   | Amsterdam, Rotterdam, Poznan, <b>London</b>                    |
| May   | Belgrade, Hamburg, Berlin                                      |
| Jun   | Amsterdam, Stockholm, Budapest, <b>London</b> , Munich, Prague |
| Jul   | Berlin, Brussels   |
| Aug   |  |
| Sep   | <b>London</b> , Dublin   |
| Oct   | Exeter, Munich, Dublin, The Hague, Amsterdam, Frankfurt        |
| Nov   | Munich, Zurich, <b>London</b> , Glasgow                        |
| Dec   | Bay Area, <b>London</b>  |

Meetups,  
Workshops &  
Conferences in  
25+ Cities



Jo-fai (Joe) Chow •  
Unofficial Photographer of H2O.ai Stress Ball

# H2O Community & Fortune 100 customers



## Select Reference Customers:

**“Overall customer satisfaction is very high.” - Gartner**



# 客戶



“Overall customer satisfaction is very high.” – Gartner  
(整体客户满意度非常高)

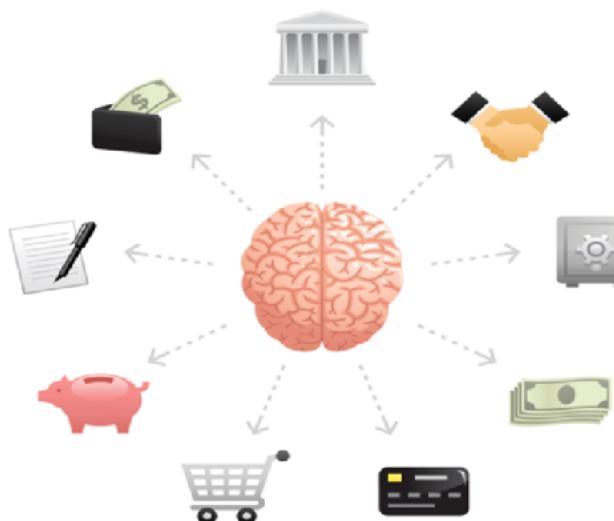
# AI in Financial Services (金融服务 AI 应用)

## Wholesale / Commercial Banking

- Know Your Customers (KYC)
- Anti-Money Laundering (AML)

## Retail Banking

- Deposit Fraud
- Customer Churn Prediction
- Auto-Loan



Today's Use Case Examples

## IT Infrastructure

- Security Cyberlake
- DoS Detection and Protection
- Master Data Management

## Card/Payments Business

- Transaction Frauds
- Real-time Targeting
- **Credit Risk Scoring**
- In-Context Promotion



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

### *Setting the scene*

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

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



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

Laura Needham partner, PwC UK



Follow

Exciting night at this year's @WAI\_News Awards: PwC wins 2017 Audit Innovation of the Year! [pwc.to/Gla17](http://pwc.to/Gla17) #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>

# Community Expansion 社区扩展

# 15 Meetups a month



The Meetup logo consists of the word "meetup" in a bold, red, cursive sans-serif font.

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

Find out more: [www.h2o.ai/community/](http://www.h2o.ai/community/)



# London Artificial Intelligence & Deep Learning PRO

H2O Artificial Intelligence and Machine Learning -  
39 groups

Location

London, United Kingdom

Members

4,184

4100+ Members



Organizers  
Ian Gomez and 2 others

Schedule

...



Our group

Meetups

Members

Photos

Discussions

More

Next Meetup

12  
DEC

Tuesday, December 12, 2017, 6:00 PM

Interpretable Machine Learning, Tweet Classifier, H2O World Highlights and More



Hosted by Jo-fai Chow

Dear All, We are thrilled to see the growth of this meetup group (4100+ members right now). Let's end the year with one more exciting meetup. This time we will have speakers from Aviva, Theodo and Barclays. Many thanks to our friends from Moody's Analytics, we have a super cool venue right in the heart of Canary Wharf. They also very kindly provide food and drinks for the event. Agenda (T.B.C.): - Doors open at 6 for pizzas and drinks as usual. -

Next London Meetup: 12 Dec at Moody's London HQ

Edit



Moody's Analytics

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



# H<sub>2</sub>O WORLD 2017

LEARNING IS FUN

REGISTER NOW

Space is limited!

Dec 4 - 5, 2017

Mountain View, CA  
Computer History Museum

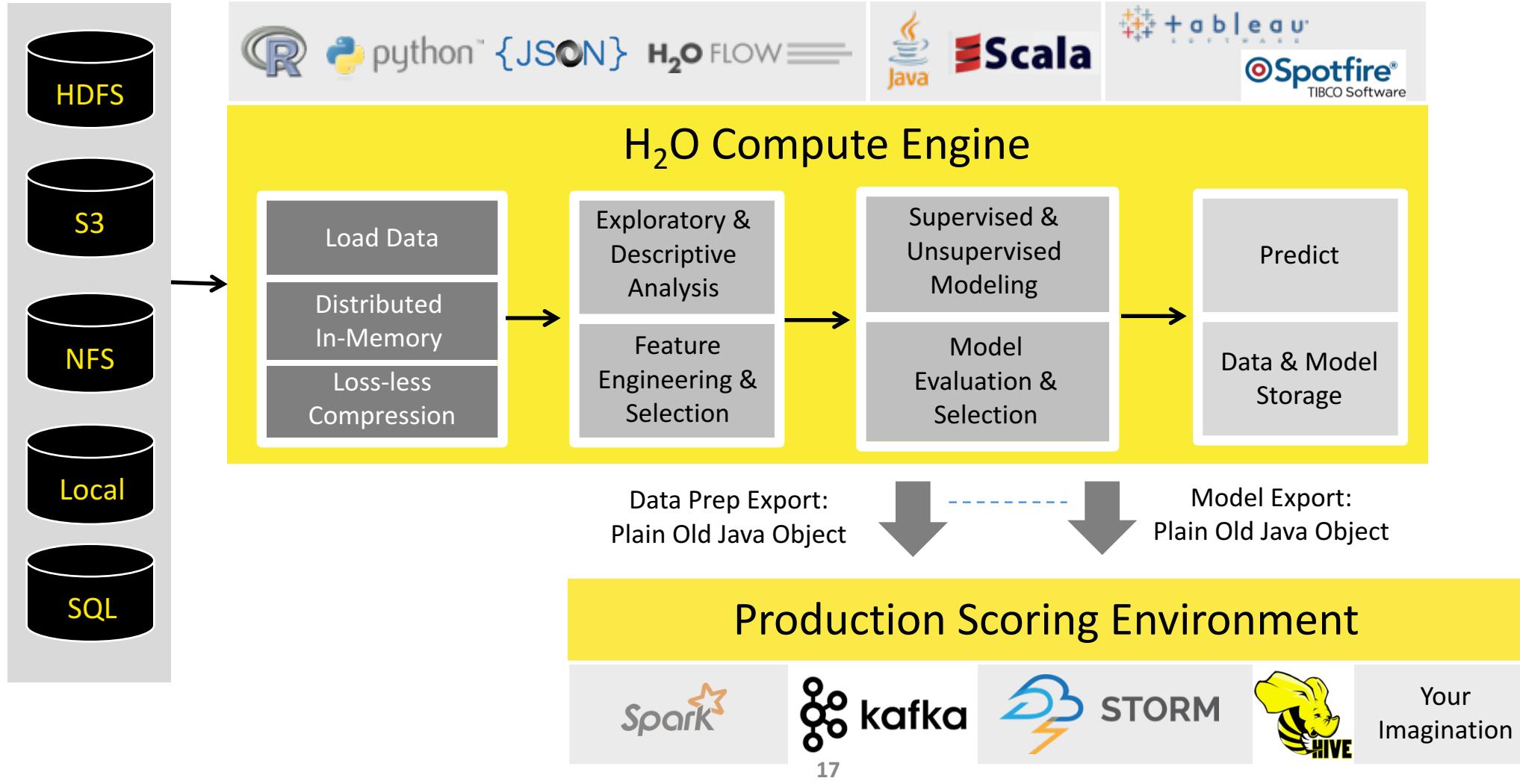
H2O is back with its flagship event, H2O World 2017.

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

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

开源的 H<sub>2</sub>O 机器学习平台

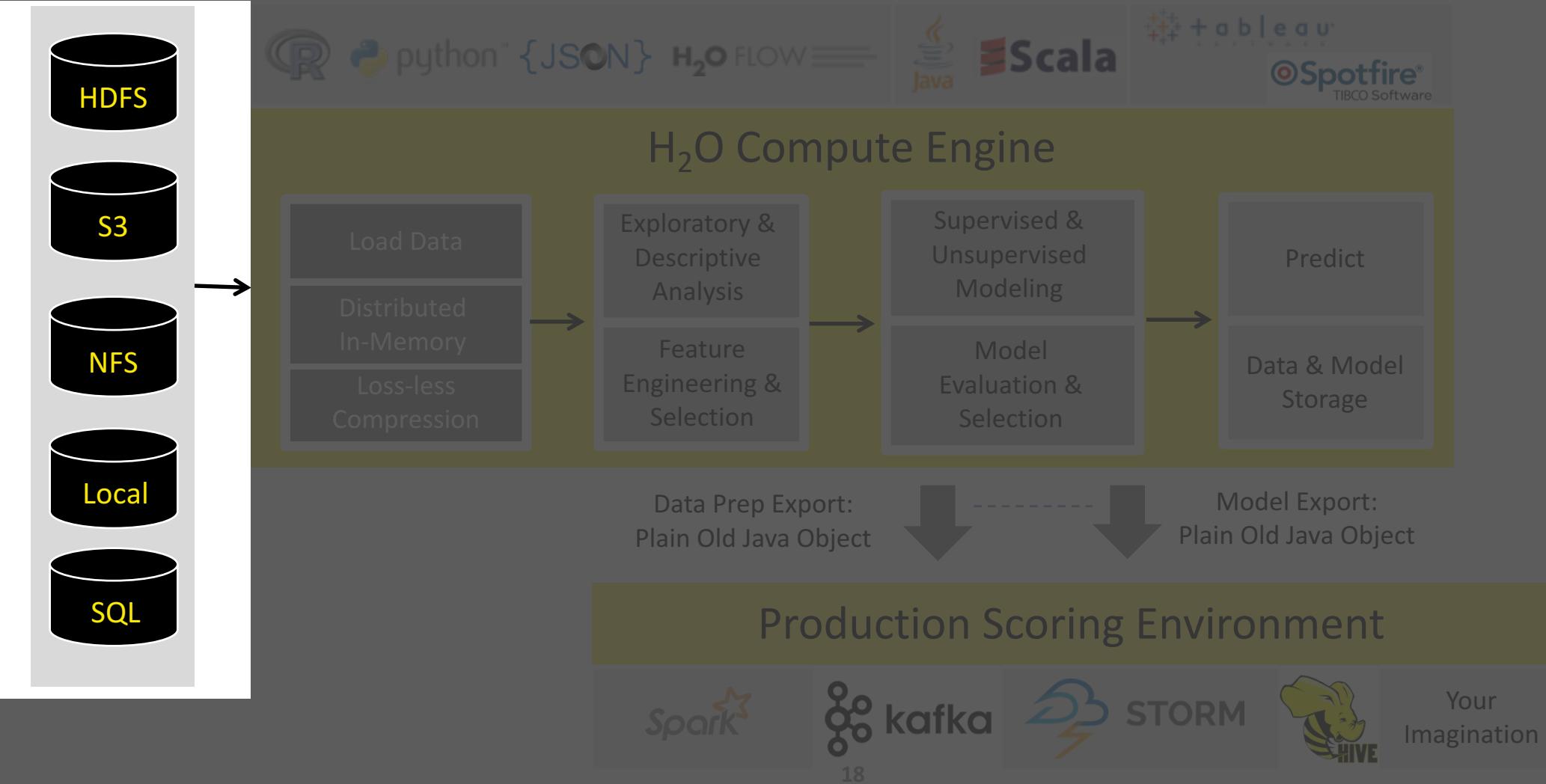
# High Level Architecture 架构



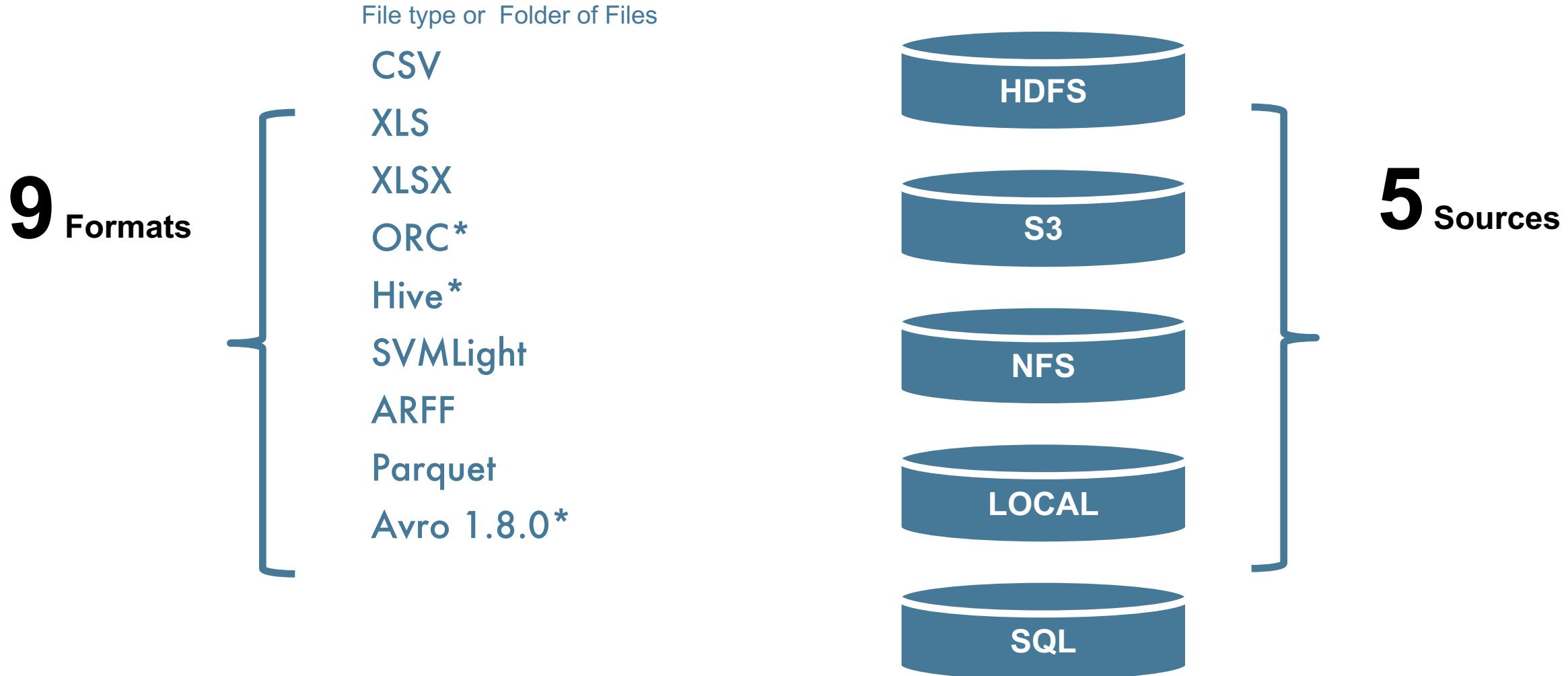
从多个来源导入数据

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

# High Level Architecture



# Supported Formats & Data Sources 支持的格式和数据源

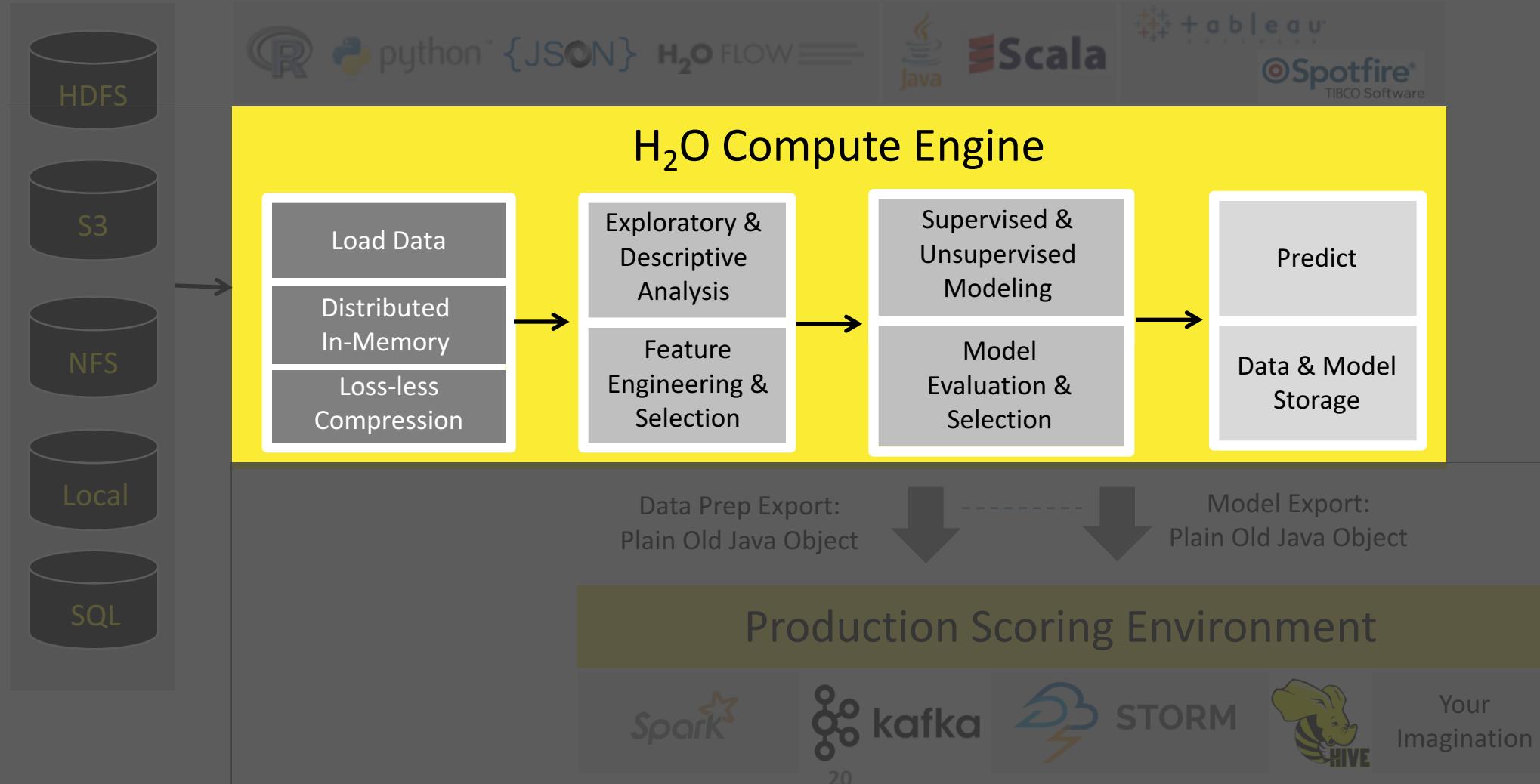


\* 1. only if H2O is running as a Hadoop job

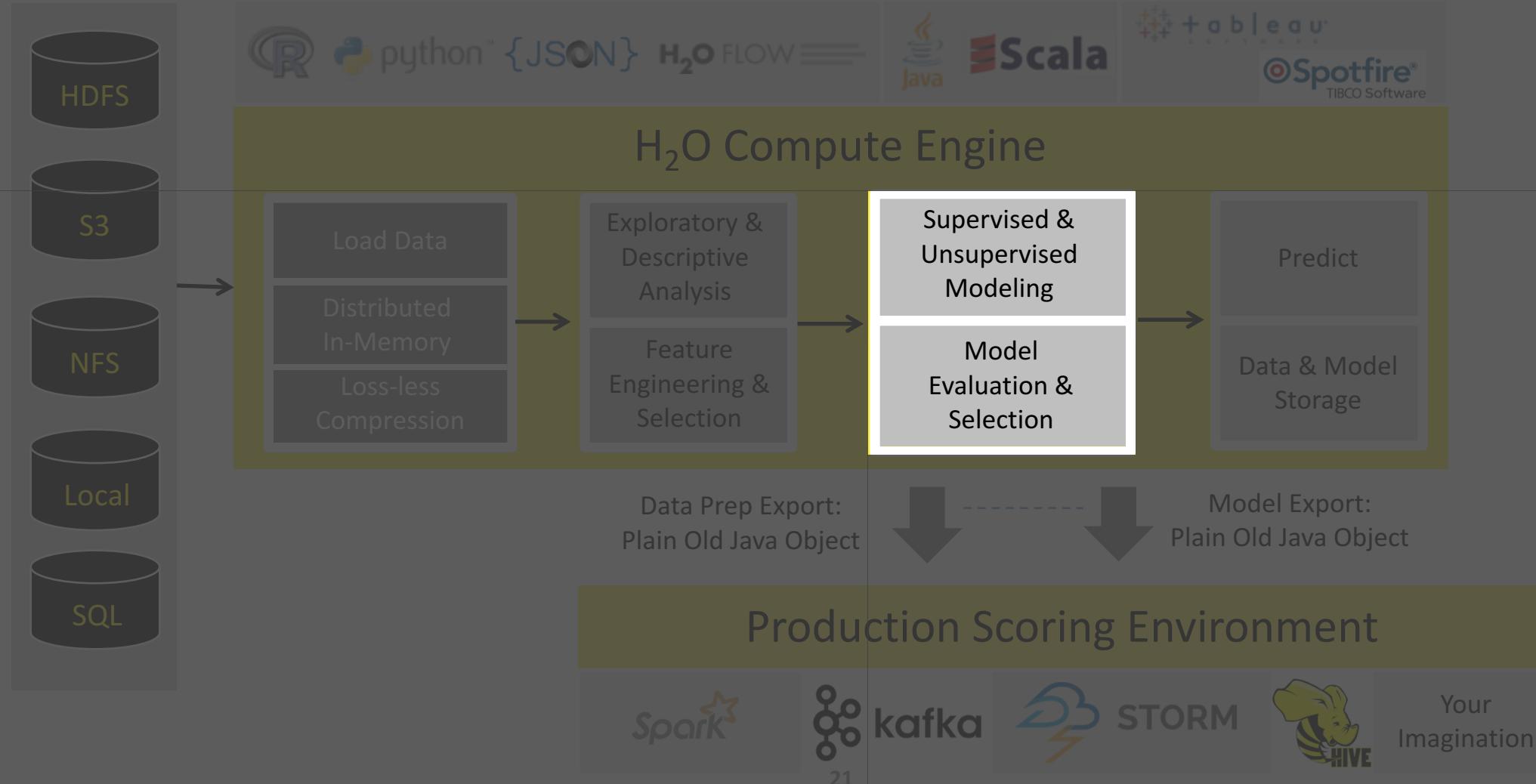
\* 2. Hive files that are saved in ORC format

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

# High Level Architecture

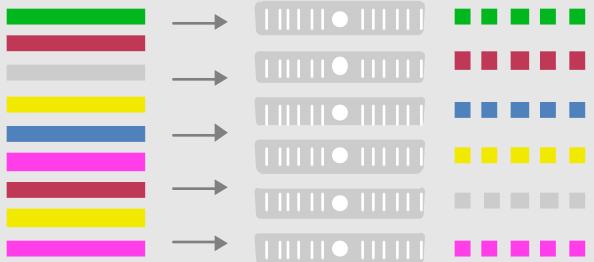


# High Level Architecture

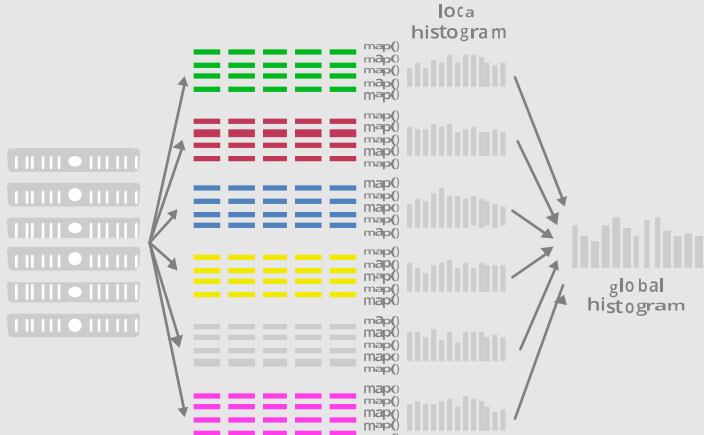


# Distributed Algorithms 分布式算法

## 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**

# Algorithms Overview 算法概述

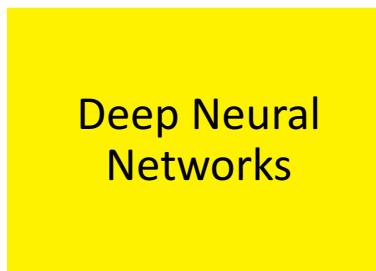
## Supervised Learning



- **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

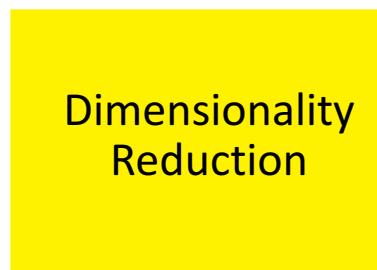


- **Deep learning:** Create multi-layer feed forward neural networks starting with an input layer followed by multiple layers of nonlinear transformations

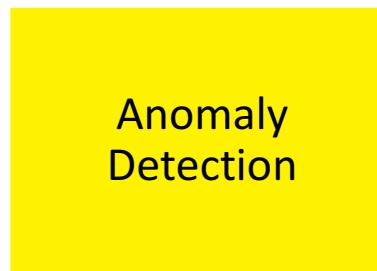
## Unsupervised Learning



- **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

# 自动机器学习

## AutoML: Automatic Machine Learning

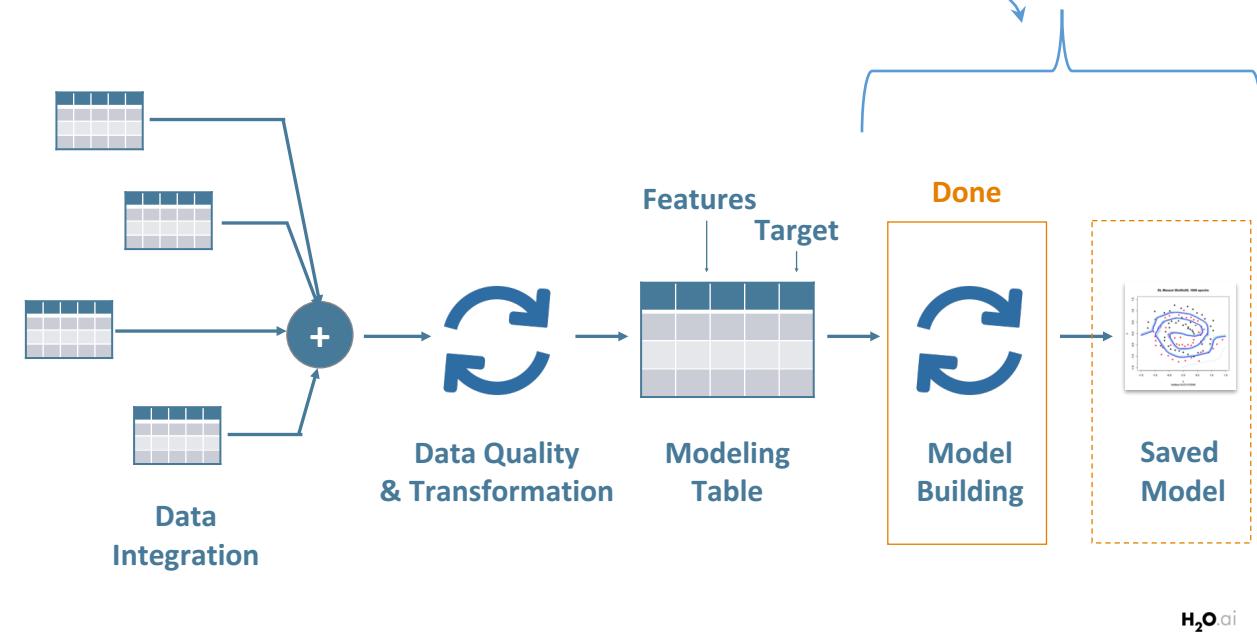
In recent years, the demand for machine learning experts has outpaced the supply, despite the surge of people entering the field. To address this gap, there have been big strides in the development of user-friendly machine learning software that can be used by non-experts. The first steps toward simplifying 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.



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.

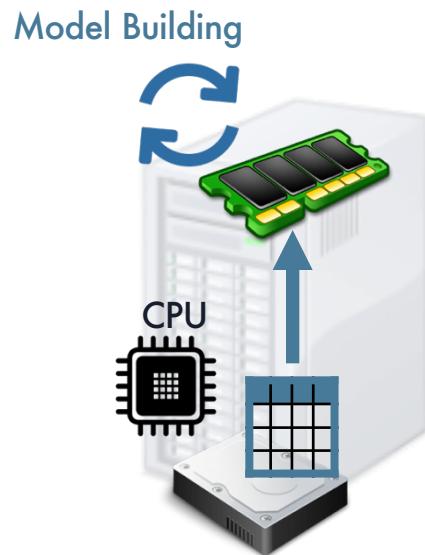
Here is an example leaderboard for a binary classification task:

| 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 |

# H<sub>2</sub>O Core 核心



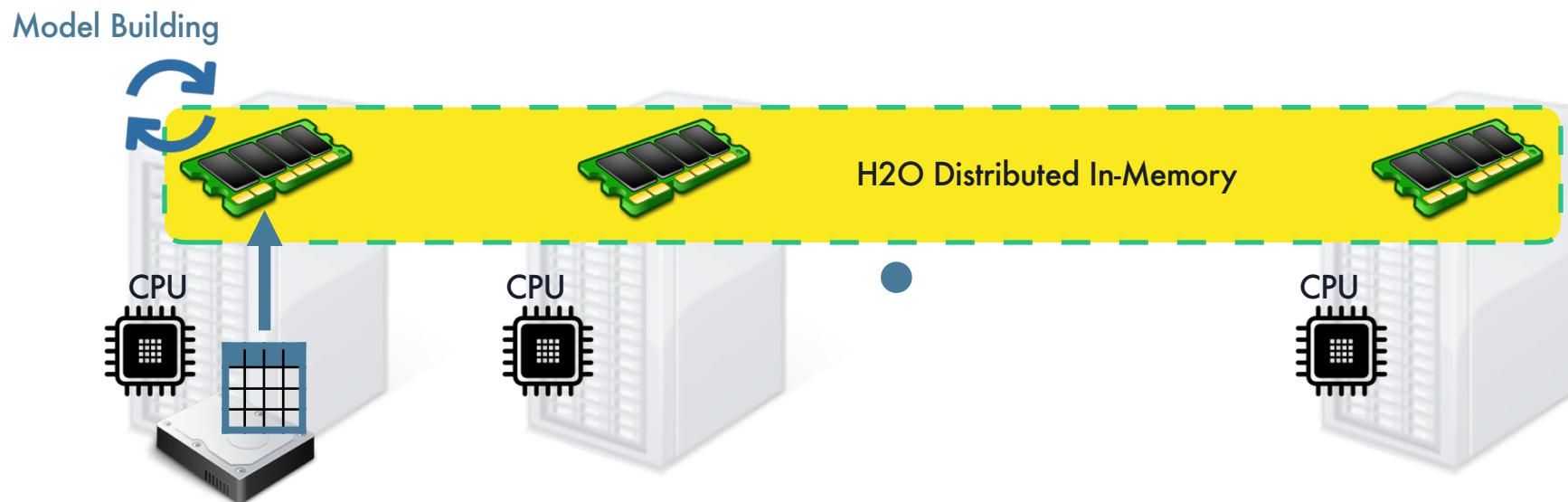
# H<sub>2</sub>O Core 核心



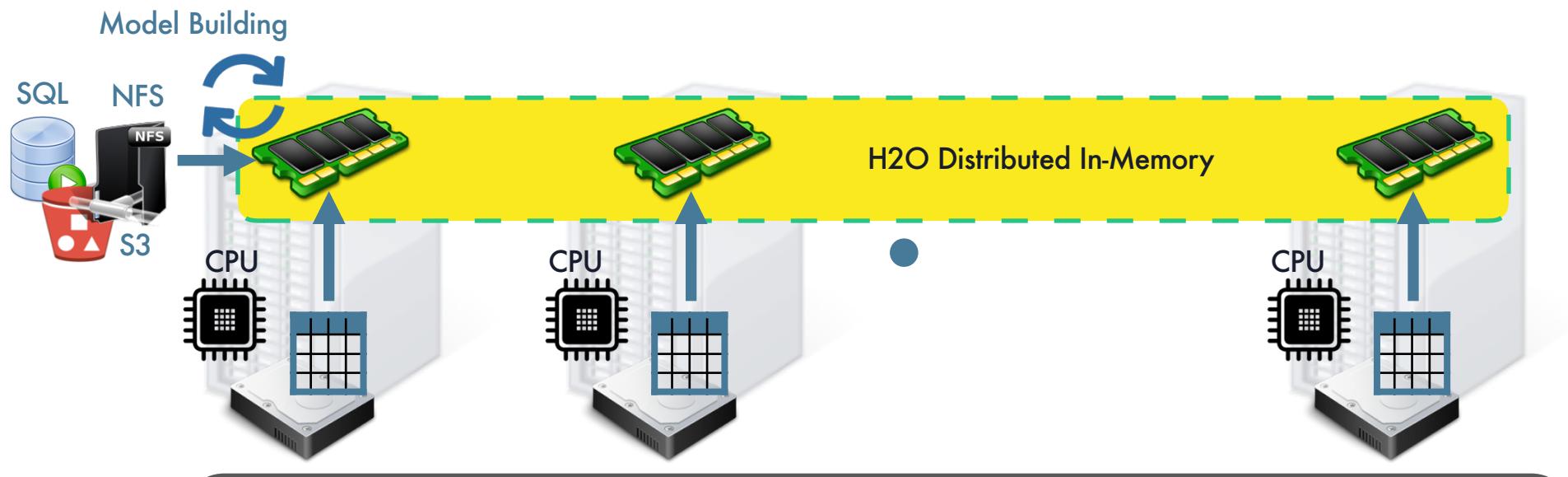
# H<sub>2</sub>O Core 核心



# H<sub>2</sub>O Core 核心



# H<sub>2</sub>O Core 核心



YARN

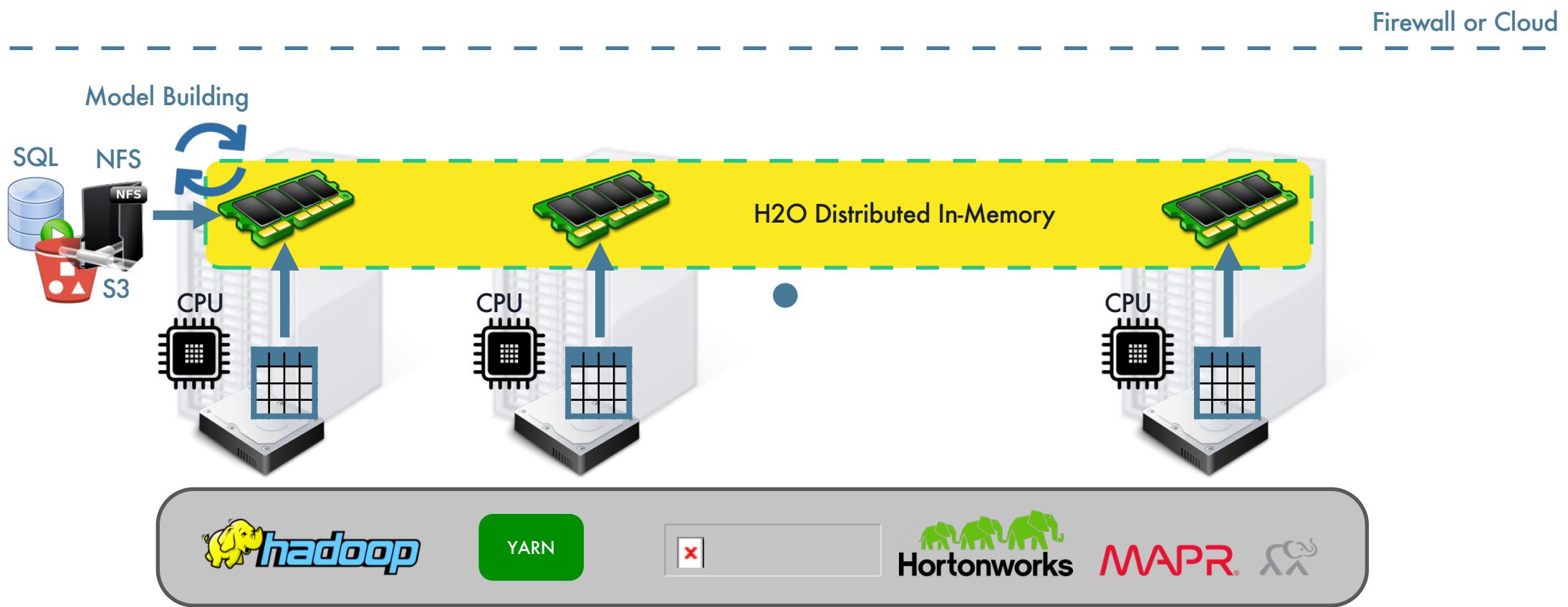
cloudera



MAPR



# H<sub>2</sub>O Core 核心

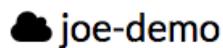


## Untitled Flow



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%           |

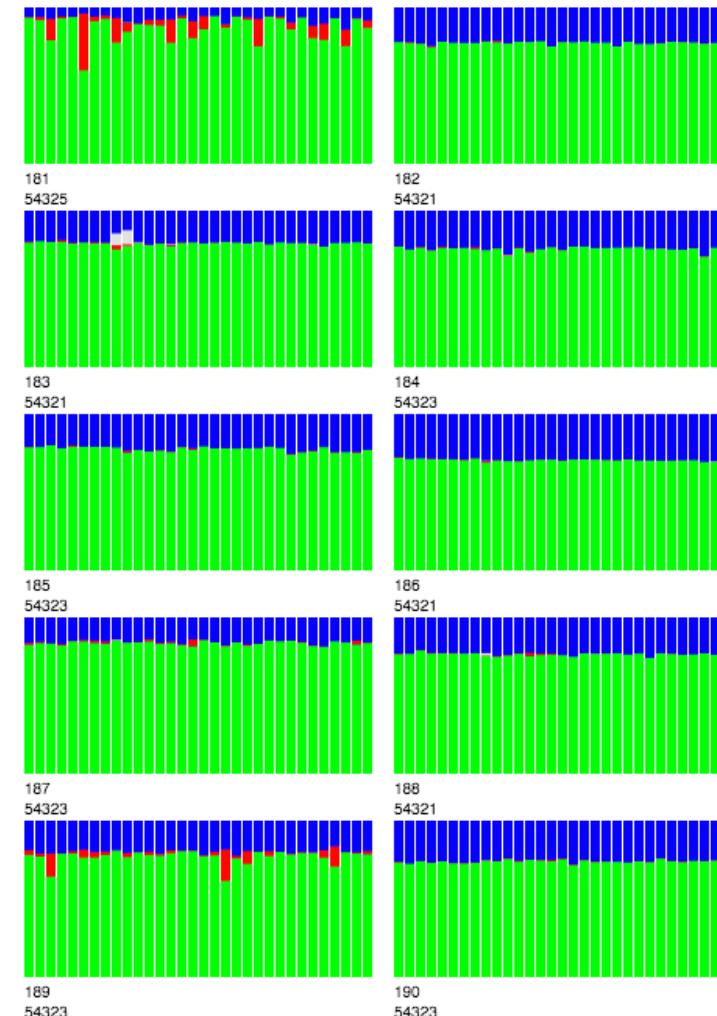
10 x 32 = 320  
CPU核

10 x 29.6 = 296  
GB 记忆体

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

# H<sub>2</sub>O Water Meter (CPU 功率监视器)

10 x 32 = 320 Cores



## Legend

Each bar represents one CPU.

Blue: idle time

Green: user time

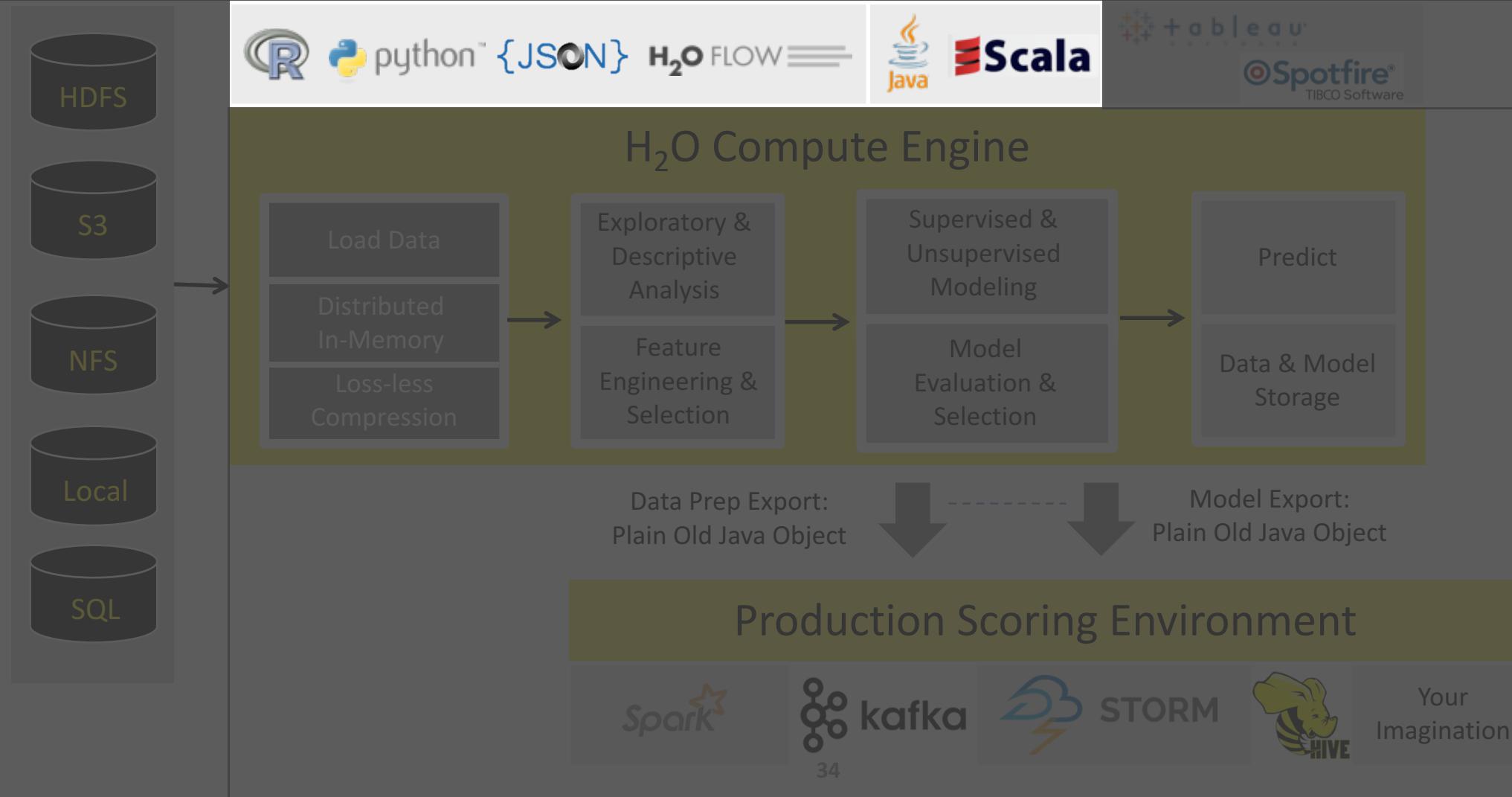
Red: system time

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

多介面

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

# High Level Architecture



# H<sub>2</sub>O Flow (网页界面)

The screenshot shows the H2O Flow web interface running in a browser window titled "H2O Flow". The URL in the address bar is "localhost:54321/flow/index.html". The browser's top bar includes standard icons for back, forward, and search, along with a user profile "Jo-fai". The main interface has a header with tabs for "Flow", "Cell", "Data", "Model" (which is highlighted in yellow), "Score", "Admin", and "Help". Below the header is a toolbar with various icons for file operations like import, export, and model building. A sidebar on the left is titled "Untitled Flow" and contains a "CS" section with a single entry labeled "assist". To the right of the sidebar is a large panel titled "Assistance" which lists various H2O routines with their descriptions. A context menu is open over the "Model" tab, listing options such as Aggregator..., Deep Learning..., Distributed Random Forest..., Gradient Boosting Machine..., Generalized Linear Modeling..., Generalized Low Rank Modeling..., K-means..., Naive Bayes..., Principal Components Analysis..., Stacked Ensemble..., Word2Vec..., XGBoost..., List All Models, List Grid Search Results, Import Model..., Export Model..., and Run AutoML... . On the far right, there is a "HELP" panel with sections for "Using Flow for the first time?", "Quickstart Videos", "view example Flows", "STAR H2O ON GITHUB!", "GENERAL" (with links to Flow Web UI, Importing Data, Building Models, Making Predictions, Using Flows, and Troubleshooting Flow), and "EXAMPLES" (with a note about Flow packs and a link to Browse installed packs...). The bottom right corner of the interface shows "Connections: 0" and the H2O logo.

# H<sub>2</sub>O - R / Python 介面

~/Documents/repo\_h2o/sales-engineering - master - RStudio Source Editor

```

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
49:1 (Top Level) +

```

R Script

credit\_card\_example x

localhost:8888/notebooks/credit\_card\_example.ipynb

Jupyter credit\_card\_example Last Checkpoint: 5 minutes ago (unsaved changes)

In [2]: # Start and connect to a local H2O cluster  
import h2o  
h2o.init(nthreads = -1)  
Checking whether there is an H2O instance running at http://localhost:54321.... not found.  
Attempting to start a local H2O server...  
Java Version: java version "1.8.0\_72"; Java(TM) SE Runtime Environment (build 1.8.0\_72-b15); Java HotSpot(TM) 64-Bit Server VM (build 25.72-b15, mixed mode)  
Starting server from /Users/jofaichow/anaconda/lib/python2.7/site-packages/h2o/backend/bin/h2o.jar  
Ice root: /var/folders/4z/p7yt7\_4n4fjijiy6g4qfbw000gn/T/tmpPdP3Av  
JVM stdout: /var/folders/4z/p7yt7\_4n4fjijiy6g4qfbw000gn/T/tmpPdP3Av/h2o\_jofaichow\_started\_from\_python.out  
JVM stderr: /var/folders/4z/p7yt7\_4n4fjijiy6g4qfbw000gn/T/tmpPdP3Av/h2o\_jofaichow\_started\_from\_python.err  
Server is running at http://127.0.0.1:54321  
Connecting to H2O server at http://127.0.0.1:54321... successful.

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

In [3]: # Import datasets from s3  
df\_train = h2o.import\_file("https://s3.amazonaws.com/h2o-training/credit\_card/credit\_card\_train.csv")  
df\_test = h2o.import\_file("https://s3.amazonaws.com/h2o-training/credit\_card/credit\_card\_test.csv")  
Parse progress: |██████████| 100%  
Parse progress: |██████████| 100%

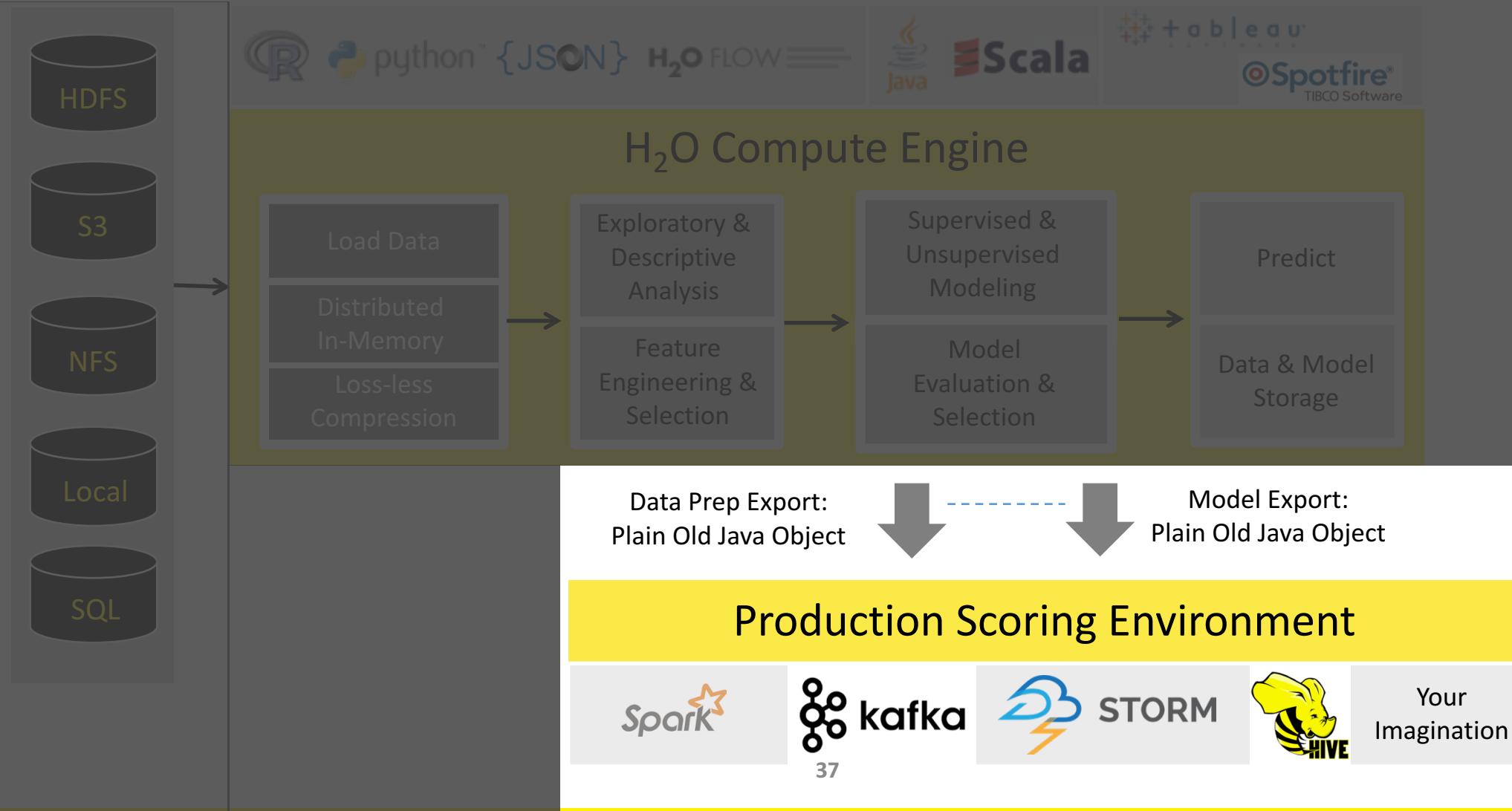
In [4]: # Look at datasets  
df\_train.summary()  
df\_test.summary()

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

导出生产的独立模型

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

# High Level Architecture



## H<sub>2</sub>O Open Source Software Documentation | H<sub>2</sub>O Commercial Software Documentation

### Open Source Software Documentation

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

#### Getting Started & User Guides

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

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

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

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

##### Sparkling Water

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

[Quick Start Video - Scala](#)

[Download Sparkling Water](#)

##### Steam

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

[Download Steam](#)

##### Deep Water (preview)

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

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

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

[H<sub>2</sub>O4GPU Readme](#)  
[Open Source License \(Apache V2\)](#)

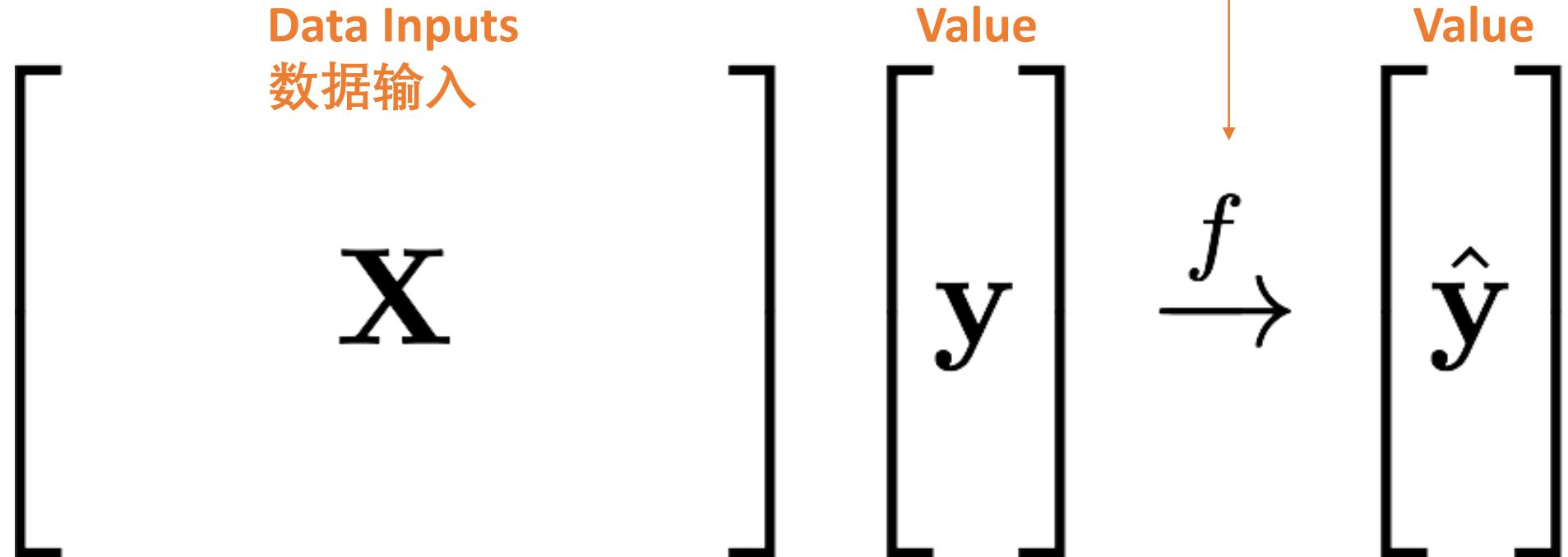
[Download H<sub>2</sub>O4GPU](#)

# Credit Card Demo

信用卡风险演示

# Learning from Data

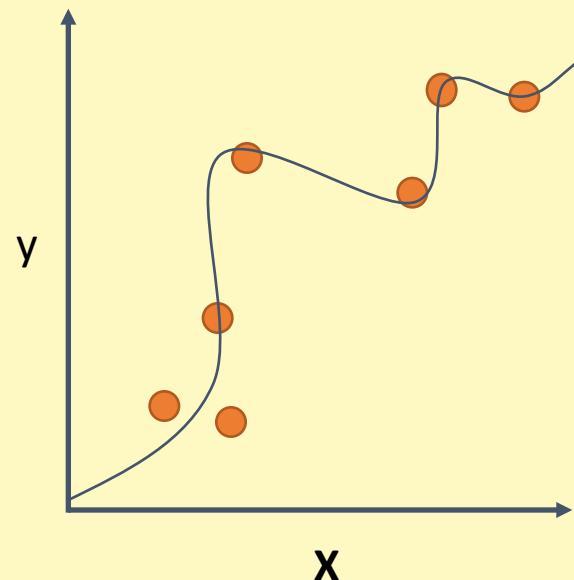
## 从数据中学习



# Supervised Learning 监督学习

Regression:

How much will a customer spend?

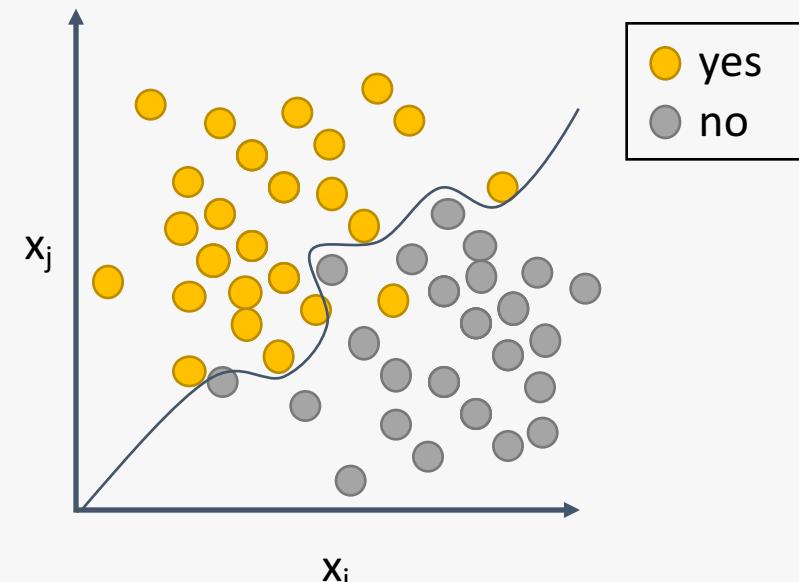


H<sub>2</sub>O algos:

Penalized Linear Models  
Random Forest  
Gradient Boosting  
Neural Networks  
Stacked Ensembles

Classification:

Will a customer make a purchase? Yes or No



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

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

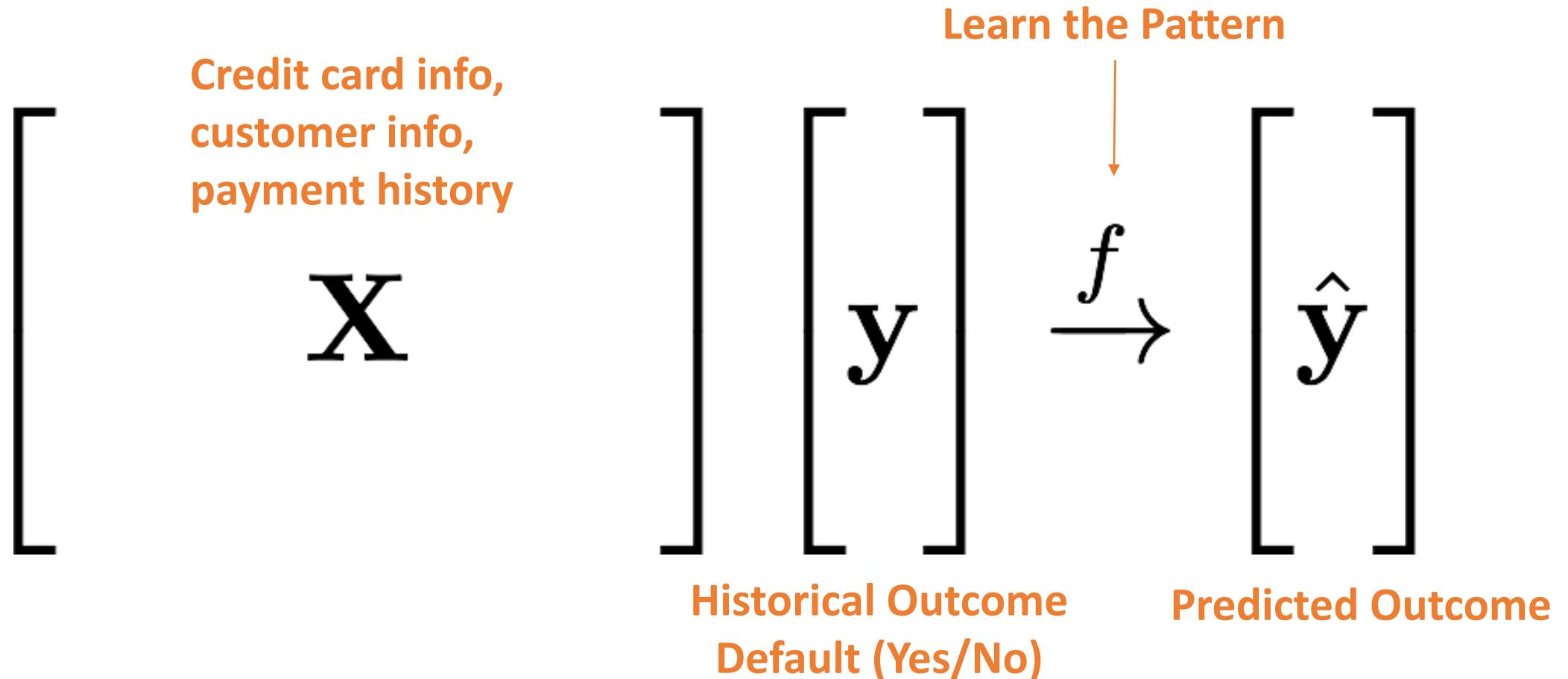
# Demo Introduction

/ **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



credit\_card\_train.csv

|    | A         | B      | C         | D        | E   | F     | G     | H     | I     | J     | K     | L         | M         | N         | O         | P         | Q         | R        | S        | T        | U        | V        | W        | X                          | Y   |
|----|-----------|--------|-----------|----------|-----|-------|-------|-------|-------|-------|-------|-----------|-----------|-----------|-----------|-----------|-----------|----------|----------|----------|----------|----------|----------|----------------------------|-----|
| 1  | LIMIT_BAL | SEX    | EDUCATION | MARRIAGE | AGE | PAY_0 | PAY_2 | PAY_3 | PAY_4 | PAY_5 | PAY_6 | BILL_AMT1 | BILL_AMT2 | BILL_AMT3 | BILL_AMT4 | BILL_AMT5 | BILL_AMT6 | PAY_AMT1 | PAY_AMT2 | PAY_AMT3 | PAY_AMT4 | PAY_AMT5 | PAY_AMT6 | DEFAULT_PAYMENT_NEXT_MONTH |     |
| 2  | 20000     | Female | 2         | 1        | 24  | 2     | 2     | -1    | -1    | -2    | -2    | 3913      | 3102      | 689       | 0         | 0         | 0         | 0        | 689      | 0        | 0        | 0        | 0        | 0                          | Yes |
| 3  | 120000    | Female | 2         | 2        | 26  | -1    | 2     | 0     | 0     | 0     | 2     | 2682      | 1725      | 2682      | 3272      | 3455      | 3261      | 0        | 1000     | 1000     | 1000     | 0        | 0        | 2000                       | Yes |
| 4  | 90000     | Female | 2         | 2        | 34  | 0     | 0     | 0     | 0     | 0     | 0     | 29239     | 14027     | 13559     | 14331     | 14948     | 15549     | 1518     | 1500     | 1000     | 1000     | 1000     | 5000     | No                         |     |
| 5  | 50000     | Female | 2         | 1        | 37  | 0     | 0     | 0     | 0     | 0     | 0     | 46990     | 48233     | 49291     | 28314     | 28959     | 29547     | 2000     | 2019     | 1200     | 1100     | 1069     | 1000     | No                         |     |
| 6  | 50000     | Male   | 2         | 1        | 57  | -1    | 0     | -1    | 0     | 0     | 0     | 8617      | 5670      | 35835     | 20940     | 19146     | 19131     | 2000     | 36681    | 10000    | 9000     | 689      | 679      | No                         |     |
| 7  | 100000    | Female | 2         | 2        | 23  | 0     | -1    | -1    | 0     | 0     | -1    | 11876     | 380       | 601       | 221       | -159      | 567       | 380      | 601      | 0        | 581      | 1687     | 1542     | No                         |     |
| 8  | 140000    | Female | 3         | 1        | 28  | 0     | 0     | 2     | 0     | 0     | 0     | 11285     | 14096     | 12108     | 12211     | 11793     | 3719      | 3329     | 0        | 432      | 1000     | 1000     | 1000     | No                         |     |
| 9  | 20000     | Male   | 3         | 2        | 35  | -2    | -2    | -2    | -2    | -1    | -1    | 0         | 0         | 0         | 0         | 13007     | 13912     | 0        | 0        | 0        | 13007    | 1122     | 0        | No                         |     |
| 10 | 200000    | Female | 3         | 2        | 34  | 0     | 0     | 2     | 0     | 0     | -1    | 11073     | 9787      | 5535      | 2513      | 1828      | 3731      | 2306     | 12       | 50       | 300      | 3738     | 66       | No                         |     |
| 11 | 630000    | Female | 2         | 2        | 41  | -1    | 0     | -1    | -1    | -1    | -1    | 12137     | 6500      | 6500      | 6500      | 6500      | 2870      | 1000     | 6500     | 6500     | 6500     | 2870     | 0        | No                         |     |
| 12 | 70000     | Male   | 2         | 2        | 30  | 1     | 2     | 2     | 0     | 0     | 2     | 65802     | 67369     | 65701     | 66782     | 36137     | 36894     | 3200     | 0        | 3000     | 3000     | 1500     | 0        | Yes                        |     |
| 13 | 250000    | Male   | 1         | 2        | 20  | 0     | 0     | 0     | 0     | 0     | 0     | 70007     | 67000     | 67000     | 67000     | 67000     | 55512     | 3000     | 3000     | 3000     | 3000     | 3000     | 3000     | No                         |     |
| 14 | 50000     | Female | 3         | 2        | 41  | -1    | 0     | -1    | -1    | -1    | -1    | 12137     | 6500      | 6500      | 6500      | 6500      | 2870      | 1000     | 6500     | 6500     | 6500     | 2870     | 0        | No                         |     |
| 15 | 20000     | Male   | 1         | 2        | 20  | 0     | 0     | 0     | 0     | 0     | 0     | 70007     | 67000     | 67000     | 67000     | 67000     | 30211     | 0        | 1500     | 1100     | 1200     | 1300     | 1100     | No                         |     |
| 16 | 320000    | Male   | 1         | 2        | 36  | 0     | 0     | 0     | 0     | 0     | 0     | 19104     | 3200      | 0         | 1500      | 0         | 19104     | 3200     | 0        | 1500     | 0        | 1650     | 0        | Yes                        |     |
| 17 | 360000    | Female | 1         | 2        | 36  | 0     | 0     | 0     | 0     | 0     | 0     | 195599    | 10358     | 10000     | 75940     | 20000     | 195599    | 50000    | 0        | 0        | 0        | 0        | 0        | No                         |     |
| 18 | 180000    | Female | 1         | 2        | 30  | 0     | 0     | 0     | 0     | 0     | 0     | 930       | 3000      | 1537      | 1000      | 2000      | 930       | 33764    | 0        | 0        | 0        | 0        | 0        | No                         |     |
| 19 | 130000    | Female | 3         | 2        | 30  | 0     | 0     | 0     | 0     | 0     | 0     | 46012     | 316       | 316       | 316       | 0         | 46012     | 2007     | 3582     | 0        | 3601     | 0        | 1820     | Yes                        |     |
| 20 | 120000    | Female | 2         | 2        | 30  | 0     | 0     | 0     | 0     | 0     | 0     | 0         | 0         | 0         | 0         | 0         | 0         | 0        | 19428    | 1473     | 560      | 0        | 0        | 1128                       |     |
| 21 | 70000     | Female | 2         | 2        | 45  | 0     | 0     | 0     | 0     | 0     | 0     | 0         | 0         | 0         | 0         | 0         | 0         | 0        | 0        | 0        | 0        | 0        | No       |                            |     |
| 22 | 450000    | Female | 1         | 2        | 23  | 0     | 0     | 0     | -1    | 0     | 0     | 4744      | 7070      | 0         | 5398      | 6360      | 8292      | 3        | 1200     | 2045     | 2000     | 2000     | 2000     | No                         |     |
| 23 | 90000     | Male   | 1         | 2        | 23  | 0     | 0     | 0     | -1    | 0     | 0     | 4744      | 7070      | 0         | 5398      | 6360      | 8292      | 3        | 1200     | 2045     | 2000     | 2000     | 2000     | No                         |     |
| 24 | 50000     | Male   | 3         | 2        | 23  | 0     | 0     | 0     | 0     | 0     | 0     | 47620     | 41810     | 36023     | 28967     | 29829     | 30046     | 1        | 1432     | 1062     | 997      | 997      | 997      | No                         |     |
| 25 | 60000     | Male   | 1         | 2        | 27  | 1     | -2    | -1    | -1    | -1    | -1    | -109      | -425      | 259       | -57       | 127       | -189      | 1        | 500      | 0        | 1000     | 0        | 1000     | Yes                        |     |
| 26 | 50000     | Female | 3         | 2        | 30  | 0     | 0     | 0     | 0     | 0     | 0     | 22541     | 16138     | 17163     | 17878     | 18931     | 19617     | 1        | 150      | 0        | 0        | 0        | 0        | No                         |     |
| 27 | 50000     | Female | 3         | 1        | 47  | -1    | -1    | -1    | -1    | -1    | -1    | 650       | 3415      | 3416      | 2040      | 30430     | 257       | 1        | 3043     | 0        | 0        | 0        | 0        | No                         |     |
| 28 | 50000     | Male   | 1         | 2        | 26  | 0     | 0     | 0     | 0     | 0     | 0     | 15329     | 16575     | 17496     | 17907     | 18375     | 11400     | 1        | 100      | 0        | 0        | 0        | 0        | No                         |     |
| 29 | 230000    | Female | 1         | 2        | 27  | -1    | -1    | -1    | -1    | -1    | -1    | 16646     | 17265     | 13266     | 15339     | 14307     | 36923     | 1        | 1200     | 2045     | 2000     | 2000     | 2000     | No                         |     |
| 30 | 100000    | Male   | 1         | 2        | 32  | 0     | 0     | 0     | 0     | 0     | 0     | 93036     | 84071     | 82880     | 80958     | 78703     | 75589     | 3023     | 3511     | 3302     | 320      | 320      | 320      | No                         |     |
| 31 | 50000     | Female | 2         | 2        | 54  | -2    | -2    | -2    | -2    | -2    | -2    | 10929     | 4152      | 22722     | 7521      | 71439     | 8981      | 4152     | 22827    | 7521     | 7143     | 7143     | 7143     | No                         |     |
| 32 | 500000    | Male   | 1         | 1        | 58  | -2    | -2    | -2    | -2    | -2    | -2    | 13709     | 5006      | 31130     | 3180      | 0         | 5293      | 5006     | 31178    | 3180     | 0        | 1000     | 1000     | No                         |     |
| 33 | 160000    | Male   | 1         | 2        | 30  | -1    | -1    | -2    | -2    | -2    | -1    | 30265     | -131      | -527      | -923      | -1488     | -1884     | 131      | 396      | 396      | 396      | 56       | 56       | No                         |     |
| 34 | 280000    | Male   | 2         | 1        | 40  | 0     | 0     | 0     | 0     | 0     | 0     | 186503    | 181328    | 180422    | 170410    | 173901    | 177413    | 8026     | 8060     | 6300     | 6300     | 640      | 640      | No                         |     |
| 35 | 60000     | Female | 2         | 2        | 22  | 0     | 0     | 0     | 0     | 0     | -1    | 15054     | 9806      | 11068     | 6026      | -28335    | 18660     | 1500     | 1518     | 2043     | 2043     | 2043     | 2043     | No                         |     |
| 36 | 50000     | Male   | 1         | 2        | 25  | 1     | -1    | -1    | -2    | -2    | -2    | 0         | 780       | 0         | 0         | 0         | 0         | 780      | 0        | 0        | 0        | 0        | 0        | Yes                        |     |
| 37 | 280000    | Male   | 1         | 2        | 31  | -1    | -1    | 2     | -1    | 0     | -1    | 498       | 9075      | 4641      | 9976      | 17976     | 9477      | 9075     | 0        | 9976     | 800      | 800      | 800      | No                         |     |
| 38 | 360000    | Male   | 1         | 2        | 33  | 0     | 0     | 0     | 0     | 0     | 0     | 218668    | 221296    | 206895    | 628699    | 195969    | 179224    | 10000    | 7000     | 6000     | 6000     | 18884    | 18884    | No                         |     |
| 39 | 70000     | Female | 1         | 2        | 25  | 0     | 0     | 0     | 0     | 0     | 0     | 67521     | 66999     | 63949     | 63699     | 64718     | 65970     | 3000     | 4500     | 4042     | 4042     | 250      | 250      | No                         |     |
| 40 | 10000     | Male   | 2         | 2        | 22  | 0     | 0     | 0     | 0     | 0     | 0     | 1977      | 2184      | 6002      | 2575      | 2620      | 2454      | 1500     | 2022     | 1000     | 1000     | 1000     | 1000     | No                         |     |

**Data Inputs (x):**  
credit card limit balance,  
customer info (sex, education, marriage, age),  
payment history (six months)

**Learn the Pattern**

**Historical Outcome (y):**  
Default Credit Card Payment after six months?  
(Yes / No)

default\_payment\_test\_data

credit\_card\_test.csv

Home Insert Page Layout Formulas Data Review View

Cut Copy Format

Calibri (Body) 12 A A Wrap Text General Conditional Formatting

Merge & Center Format as Table

|    | A         | B      | C         | D        | E   | F     | G     | H     | I     | J     | K     | L         | M         | N         | O         | P         | Q         | R        | S        | T        | U        | V        | W        | X | Y | Z | AA |
|----|-----------|--------|-----------|----------|-----|-------|-------|-------|-------|-------|-------|-----------|-----------|-----------|-----------|-----------|-----------|----------|----------|----------|----------|----------|----------|---|---|---|----|
| 1  | LIMIT_BAL | SEX    | EDUCATION | MARRIAGE | AGE | PAY_0 | PAY_2 | PAY_3 | PAY_4 | PAY_5 | PAY_6 | BILL_AMT1 | BILL_AMT2 | BILL_AMT3 | BILL_AMT4 | BILL_AMT5 | BILL_AMT6 | PAY_AMT1 | PAY_AMT2 | PAY_AMT3 | PAY_AMT4 | PAY_AMT5 | PAY_AMT6 |   |   |   |    |
| 2  | 50000     | Male   | 1         | 2        | 37  | 0     | 0     | 0     | 0     | 0     | 0     | 64400     | 57069     | 57608     | 19394     | 19619     | 20024     | 2500     | 1815     | 657      | 1000     | 1000     | 800      |   |   |   |    |
| 3  | 500000    | Male   | 1         | 2        | 29  | 0     | 0     | 0     | 0     | 0     | 0     | 367965    | 412023    | 445007    | 542653    | 483003    | 473944    | 55000    | 40000    | 38000    | 20239    | 13750    | 13770    |   |   |   |    |
| 4  | 260000    | Female | 1         | 2        | 51  | -1    | -1    | -1    | -1    | -1    | 2     | 12261     | 21670     | 9966      | 8517      | 22287     | 13668     | 21818    | 9966     | 8583     | 22301    | 0        | 3640     |   |   |   |    |
| 5  | 50000     | Male   | 2         | 2        | 33  | 2     | 0     | 0     | 0     | 0     | 0     | 30518     | 29618     | 22102     | 22734     | 23217     | 23680     | 1718     | 1500     | 1000     | 1000     | 1000     | 716      |   |   |   |    |
| 6  | 150000    | Female | 5         | 2        | 46  | 0     | 0     | -1    | 0     | 0     | -2    | 4463      | 3034      | 1170      | 1170      | 0         | 1013      | 1170     | 0        | 0        | 0        | 0        | 0        |   |   |   |    |
| 7  | 20000     | Male   | 1         | 2        | 24  | 0     | 0     | 0     | 0     | 0     | 0     | 17447     | 18479     | 19476     | 19865     | 20480     | 20063     | 1318     | 1315     | 704      | 928      | 912      | 1069     |   |   |   |    |
| 8  | 130000    | Female | 2         | 1        | 51  | -1    | -1    | -2    | -2    | -1    | -1    | 99        | 0         | 0         | 0         | 0         | 0         | 0        | 0        | 0        | 2353     | 0        | 0        |   |   |   |    |
| 9  | 320000    | Male   | 2         | 2        | 29  | 2     | 2     | 2     | 2     | 2     | 2     | 58267     | 59246     | 60184     | 58622     | 62307     | 63526     | 2500     | 2500     | 0        | 4800     | 2400     | 1600     |   |   |   |    |
| 10 | 50000     | Male   | 3         | 2        | 25  | -1    | 0     | 0     | 0     | 0     | 0     | 42838     | 37225     | 36087     | 9636      | 9590      | 10030     | 1759     | 1779     | 320      | 500      | 1000     | 1000     |   |   |   |    |
| 11 | 130000    | Female | 1         | 1        | 35  | 0     | 0     | 0     | -1    | -1    | -1    | 81313     | 117866    | 17740     | 1330      | 7095      | 1190      | 40000    | 5000     | 1330     | 7095     | 1190     | 2090     |   |   |   |    |
| 12 | 20000     | Male   | 3         | 2        |     |       |       |       |       |       |       |           |           |           |           |           |           | 0        | 1651     | 1000     | 2000     | 0        | 1500     |   |   |   |    |
| 13 | 100000    | Female | 1         | 2        |     |       |       |       |       |       |       |           |           |           |           |           |           | 7555     | 0        | 0        | 0        | 0        | 0        |   |   |   |    |
| 14 | 400000    | Male   | 2         | 1        |     |       |       |       |       |       |       |           |           |           |           |           |           | 9677     | 11867    | 7839     | 14837    | 7959     | 5712     |   |   |   |    |
| 15 | 180000    | Male   | 1         | 1        |     |       |       |       |       |       |       |           |           |           |           |           |           | 4655     | 2690     | 2067     | 2142     | 2217     | 1000     |   |   |   |    |
| 16 | 260000    | Female | 1         | 1        |     |       |       |       |       |       |       |           |           |           |           |           |           | 0        | 22500    | 0        | 969      | 1000     | 0        |   |   |   |    |
| 17 | 140000    | Male   | 2         | 1        |     |       |       |       |       |       |       |           |           |           |           |           |           | 3455     |          |          |          |          | 2602     |   |   |   |    |
| 18 | 210000    | Male   | 3         | 1        |     |       |       |       |       |       |       |           |           |           |           |           |           | 10478    |          |          |          |          | 10478    |   |   |   |    |
| 19 | 370000    | Male   | 1         | 2        |     |       |       |       |       |       |       |           |           |           |           |           |           | 15383    |          |          |          |          | 4699     |   |   |   |    |
| 20 | 50000     | Female | 1         | 2        | 24  | 1     | -2    | -2    | -2    | -2    | -2    | -709      | -709      | -709      | -2898     | -3272     | -3272     | 0        |          |          |          |          | 0        |   |   |   |    |
| 21 | 180000    | Female | 1         | 2        | 29  | -1    | -1    | -1    | -2    | -1    | 0     | 11386     | 199       | 0         | 0         | 17227     | 17042     | 199      |          |          |          |          | 5114     |   |   |   |    |
| 22 | 120000    | Male   | 2         | 2        | 26  | 0     | 0     | 0     | 0     | 0     | 0     | 107314    | 110578    | 113736    | 116000    | 119131    | 122135    | 5000     |          |          |          |          | 5000     |   |   |   |    |
| 23 | 470000    | Male   | 2         | 2        | 27  | 2     | 2     | 2     | 2     | 0     | 0     | 296573    | 303320    | 307843    | 479978    | 305145    | 309959    | 13000    | 11001    | 0        | 10484    | 10838    | 10367    |   |   |   |    |
| 24 | 50000     | Male   | 2         | 2        | 23  | 2     | 0     | 0     | 0     | 0     | 0     | 49758     | 48456     | 44116     | 21247     | 20066     | 18858     | 2401     | 2254     | 2004     | 704      | 707      | 1004     |   |   |   |    |
| 25 | 20000     | Male   | 2         | 2        | 23  | 1     | 2     | 0     | 0     | 2     | 0     | 20235     | 17132     | 16856     | 16875     | 13454     | 10104     | 0        | 1200     | 1000     | 0        | 1000     | 10000    |   |   |   |    |
| 26 | 60000     | Female | 1         | 2        | 28  | 1     | 2     | 2     | -2    | -2    | -1    | 21501     | 20650     | 0         | 0         | 0         | 2285      | 0        | 0        | 0        | 0        | 0        | 2285     | 0 |   |   |    |
| 27 | 250000    | Female | 2         | 1        | 75  | 0     | -1    | -1    | -1    | -1    | -1    | 52874     | 1631      | 1536      | 1010      | 5572      | 794       | 1631     | 1536     | 1010     | 5572     | 794      | 1184     |   |   |   |    |
| 28 | 30000     | Male   | 2         | 2        | 28  | 0     | 0     | 0     | 0     | 0     | 0     | 29242     | 29507     | 29155     | 25255     | 22001     | 0         | 5006     | 1244     | 851      | 955      | 0        | 0        |   |   |   |    |
| 29 | 100000    | Female | 3         | 1        | 43  | 0     | 0     | -2    | -2    | -2    | -2    | 62170     | 0         | 0         | 0         | 0         | 0         | 0        | 0        | 0        | 0        | 0        | 0        |   |   |   |    |
| 30 | 50000     | Female | 1         | 2        | 26  | -1    | -1    | -1    | -1    | -1    | -1    | 1156      | 316       | 316       | 316       | 316       | 316       | 316      | 316      | 316      | 316      | 316      | 316      |   |   |   |    |
| 31 | 110000    | Female | 2         | 2        | 36  | 0     | 0     | 0     | 0     | 0     | 0     | 47819     | 48947     | 50330     | 50894     | 52175     | 53652     | 2200     | 2500     | 2000     | 2100     | 2500     | 2200     |   |   |   |    |
| 32 | 180000    | Male   | 2         | 2        | 29  | 1     | -2    | -2    | -2    | -2    | -2    | -2        | -2        | -2        | -2        | -2        | 0         | 0        | 0        | 0        | 0        | 0        |          |   |   |   |    |
| 33 | 110000    | Male   | 2         | 2        | 29  | 1     | 2     | 2     | 0     | 0     | 0     | 58362     | 56598     | 51908     | 48647     | 47862     | 47969     | 2500     | 0        | 2000     | 2000     | 1854     | 2000     |   |   |   |    |
| 34 | 20000     | Male   | 1         | 2        | 27  | 0     | 0     | 0     | 0     | 0     | 0     | 20571     | 19089     | 19658     | 19453     | 19108     | 18868     | 1323     | 1600     | 830      | 700      | 674      | 376      |   |   |   |    |
| 35 | 140000    | Female | 2         | 2        | 29  | 0     | 0     | 0     | 0     | 0     | 0     | 20110     | 17102     | 18862     | 19996     | 21214     | 21085     | 3000     | 3000     | 3000     | 3500     | 2000     | 2000     |   |   |   |    |
| 36 | 60000     | Female | 1         | 2        | 23  | 1     | 2     | 2     | 2     | 2     | 2     | 29332     | 28577     | 30805     | 31601     | 32349     | 32965     | 0        | 2709     | 1600     | 1400     | 1300     | 1200     |   |   |   |    |
| 37 | 230000    | Female | 2         | 2        | 27  | 1     | 2     | 0     | 0     | 0     | 0     | 13668     | 12647     | 13135     | 10596     | 9218      | 5068      | 0        | 1064     | 423      | 313      | 1000     | 4641     |   |   |   |    |
| 38 | 70000     | Male   | 1         | 2        | 27  | 0     | 0     | 0     | 0     | 0     | 0     | 70119     | 68536     | 66601     | 29401     | 28949     | 29795     | 3600     | 1646     | 600      | 28468    | 1327     | 1000     |   |   |   |    |
| 39 | 90000     | Male   | 3         | 1        | 48  | 1     | 2     | 2     | 2     | 2     | 2     | 77604     | 73317     | 71334     | 67009     | 63228     | 59378     | 1700     | 4000     | 1600     | 1600     | 1500     | 4086     |   |   |   |    |
| 40 | 30000     | Female | 2         | 1        | 43  | 2     | 2     | 2     | 2     | 2     | 2     | 28702     | 26622     | 24022     | 24268     | 20850     | 10622     | 1200     | 1608     | 0        | 0        | 0        | 800      |   |   |   |    |

Data Inputs (x) only – no outcome data  
(as if this is the data from new customers)

We use this dataset to make prediction

Make Predictions

Probability of Default Payments  
(Decision Makers to Take Actions)

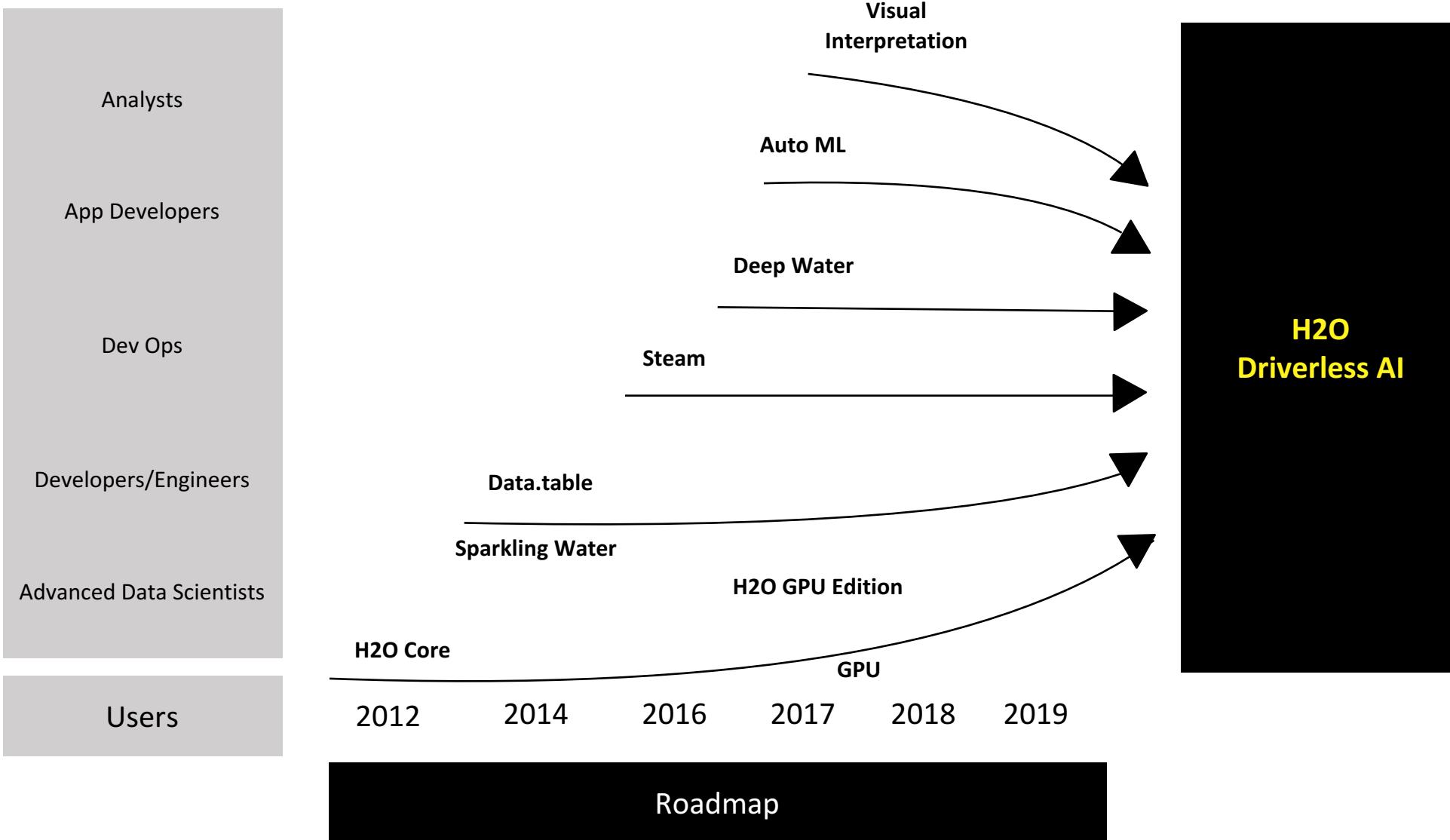
# Live Demo

- Import and explore credit card dataset
- Train a gradient boosting model and make predictions
- Train many models using AutoML and make predictions using the best model

# Driverless AI

Using AI to do AI 用人工智能制造人工智能

# H<sub>2</sub>O AI Platform Timeline



# Kaggle Masters at H<sub>2</sub>O.ai

**Dmitry Larko**

Sr. Data Scientist at H2O.ai  
San Francisco Bay Area, CA, United States  
Joined 5 years ago · last seen in the past day

[Twitter](#) [LinkedIn](#) <https://h2o.ai>



[Home](#) [Competitions \(34\)](#) [Discussion \(32\)](#) [Organizations \(2\)](#) [Followers \(4\)](#) [Contact User](#) [Follow User](#)

**Branden Murray**

Data Scientist at h2o.ai  
California, United States  
Joined 3 years ago · last seen in the past day

[GitHub](#) [Twitter](#) [LinkedIn](#)



[Home](#) [Competitions \(24\)](#) [Kernels \(5\)](#) [Discussion \(250\)](#) [Organizations \(1\)](#) [Followers \(11\)](#) [Contact User](#) [Follow User](#)

**mlandy**

Mountain View, CA, United States  
Joined 5 years ago · last seen a day ago

[GitHub](#) [Twitter](#) [LinkedIn](#)



[Home](#) [Competitions \(66\)](#) [Kernels \(12\)](#) [Discussion \(123\)](#) [Organizations \(1\)](#) [Followers \(17\)](#) [Contact User](#) [Follow User](#)

**Μαριος Μιχαηλιδης Kazanova**

Data Scientist at H2O ai  
Volos, Greece  
Joined 4 years ago · last seen in the past day  
[GitHub](#) [LinkedIn](#) <http://www.kazanovaforanalytics.com/>



Competitions Grandmaster

[Home](#) [Competitions \(93\)](#) [Kernels \(5\)](#) [Discussion \(468\)](#) [Organizations \(1\)](#) [Followers \(72\)](#) [Contact User](#) [Follow User](#)

Competitions Grandmaster

Kernels Contributor

Discussion Master

| Current Rank   | Highest Rank |
|----------------|--------------|
| 2<br>of 58,377 | 1            |
| 23             | 23           |
| 21             |              |
| 0              | 0            |
| 0              | 0            |

Unranked

Rank  
2  
of 30,803

**Faron**

Data Scientist at H2O.ai  
Deutschland  
Joined 3 years ago · last seen a day ago  
[GitHub](#) [Twitter](#) [LinkedIn](#) <http://kagglenizer.com/>



Competitions Grandmaster

Followers 402

[Home](#) [Competitions \(26\)](#) [Kernels \(9\)](#) [Discussion \(377\)](#) [Organizations \(1\)](#) [Followers \(402\)](#) [Contact User](#) [Follow User](#)

Competitions Grandmaster

Kernels Expert

Discussion Master

| Current Rank     | Highest Rank |
|------------------|--------------|
| 5<br>of 66,219   | 4            |
| 24<br>of 108,021 | 8            |
| 12               | 6            |

Current Rank  
24  
of 108,021

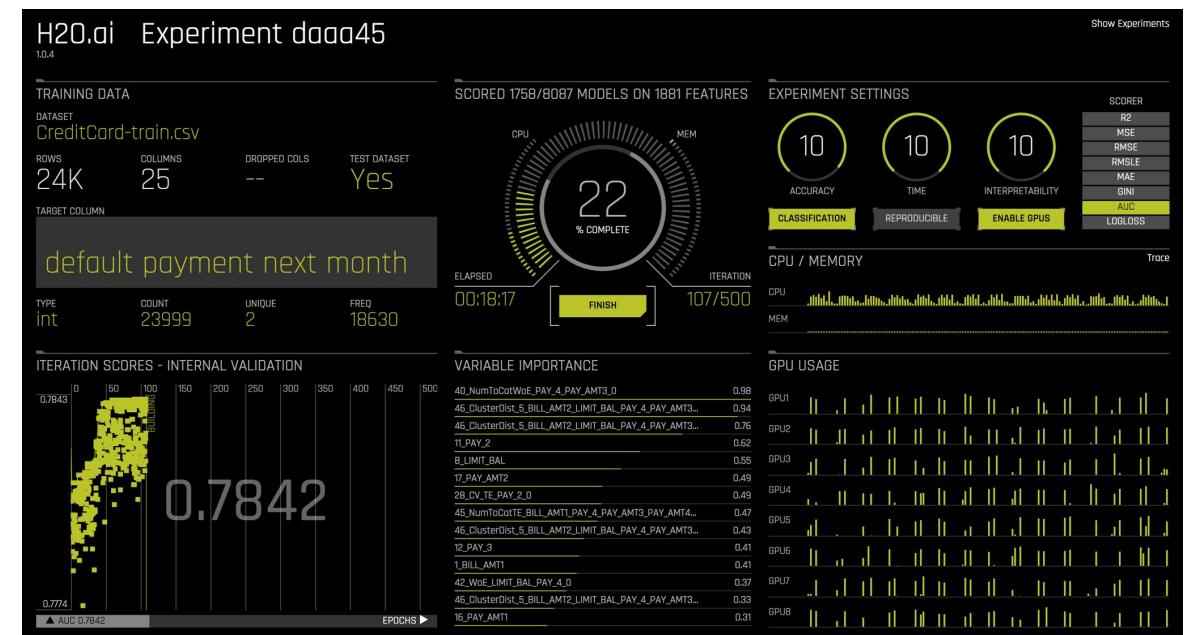
Highest Rank  
8

Current Rank  
12  
of 40,552

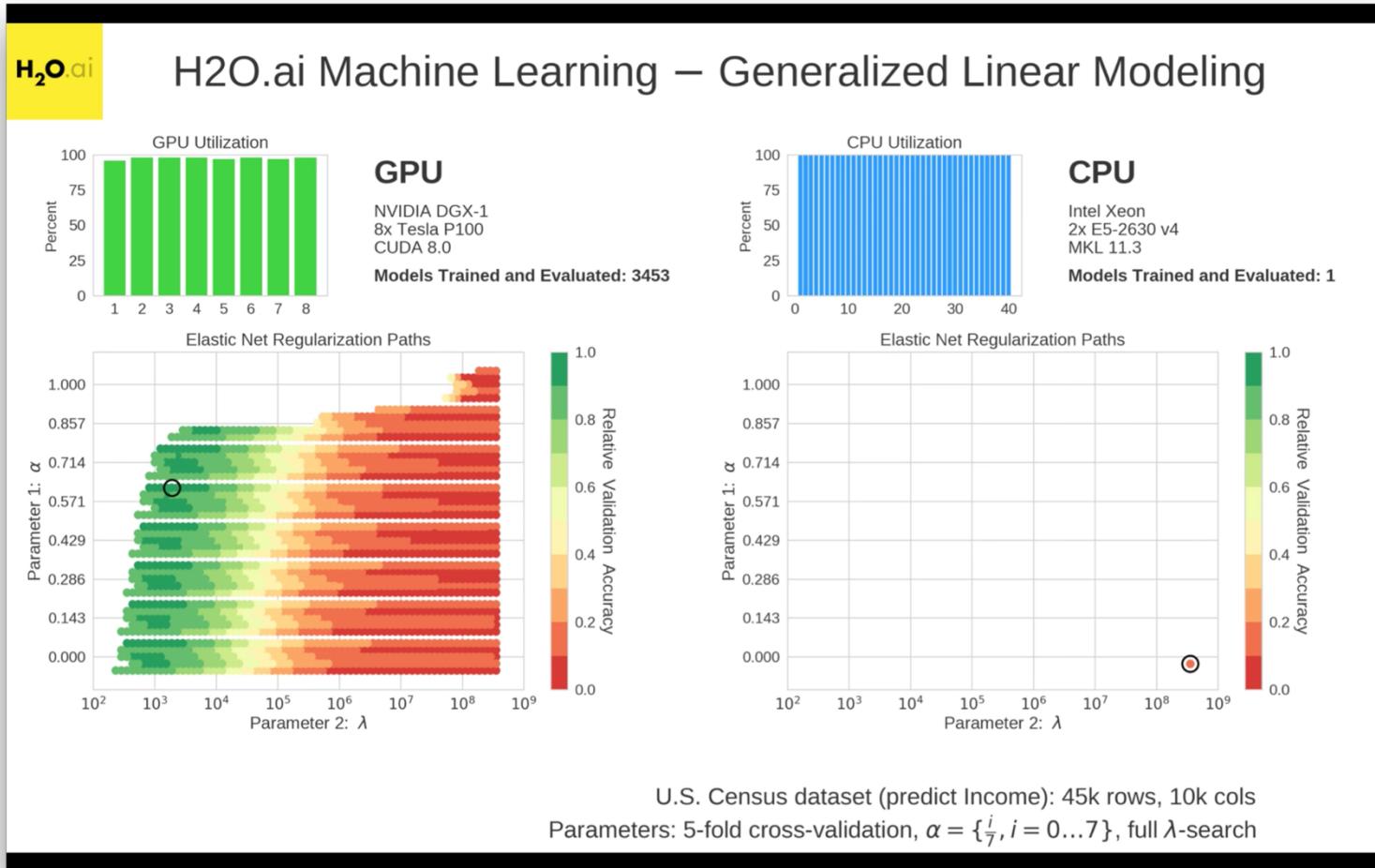
Highest Rank  
6

# The Product: Driverless AI

- Kaggle Grandmasters in a Box
  - A solution to the shortage of data scientists
- Fast
  - H<sub>2</sub>O Algorithms optimised for GPUs
- Accurate
  - Auto Feature Engineering
  - Auto Model Tuning / Selection / Ensemble
- Interpretable
  - Machine Learning Interpretation
- Production Ready
  - Export pipeline as a standalone Python package

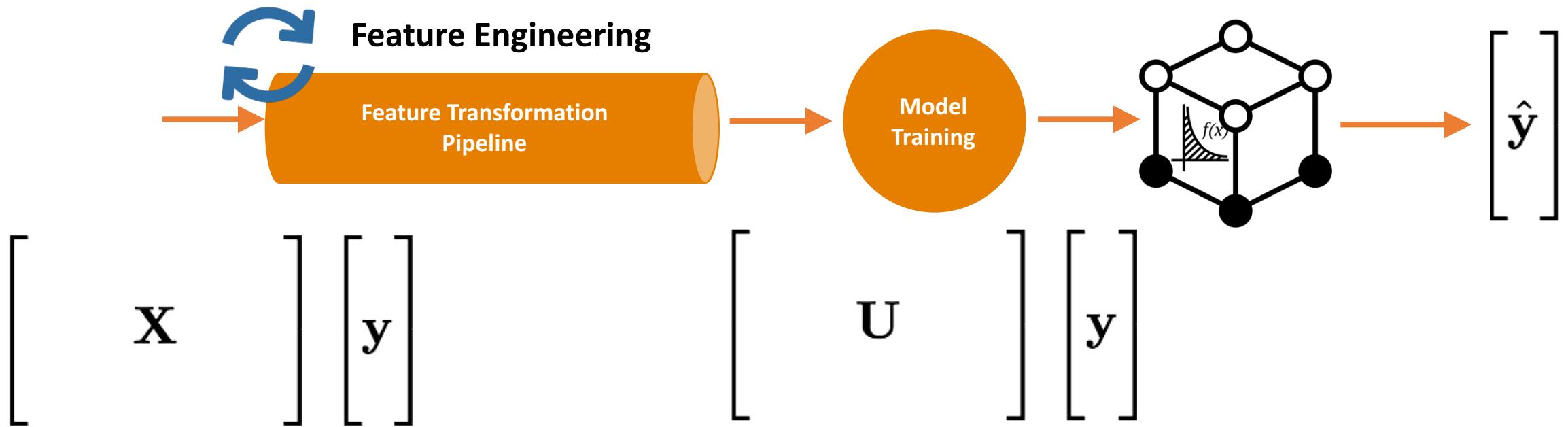


# Fast: H<sub>2</sub>O Algorithms on GPUs

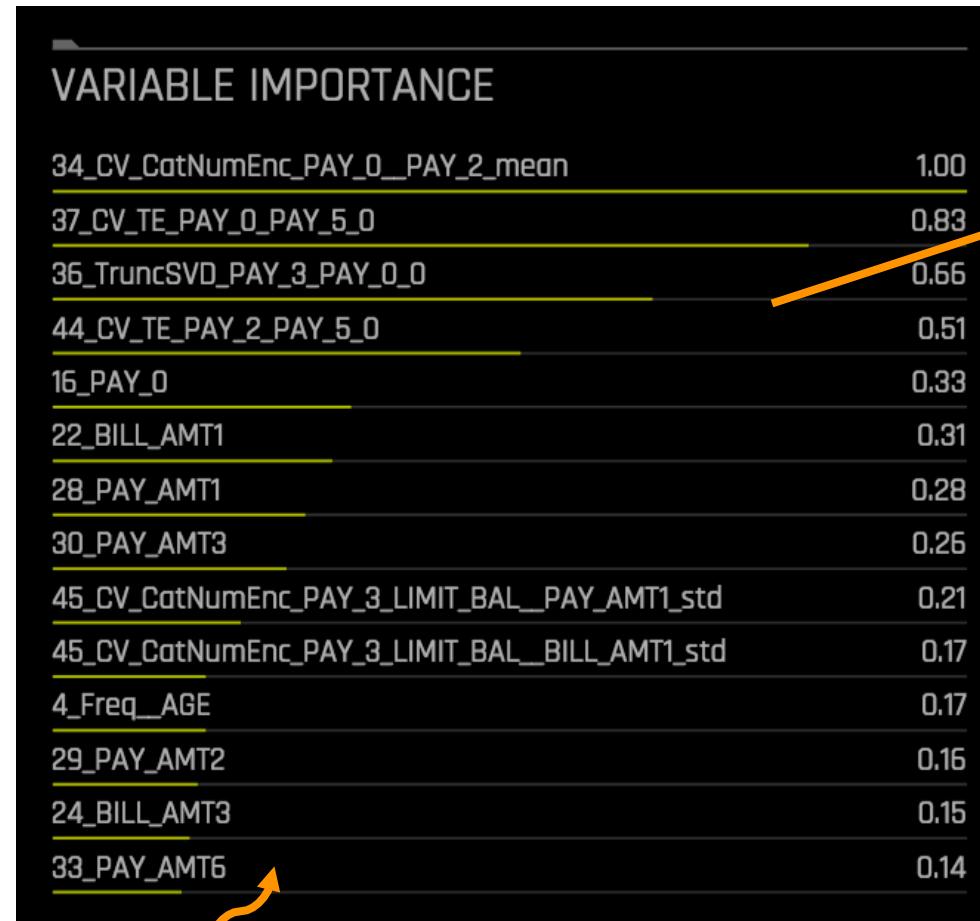


3400+ models  
were trained on  
GPUs by the time  
the first model  
completed training  
on CPUs

# Accurate: Auto Feature Engineering



# Accurate: Auto Feature Engineering

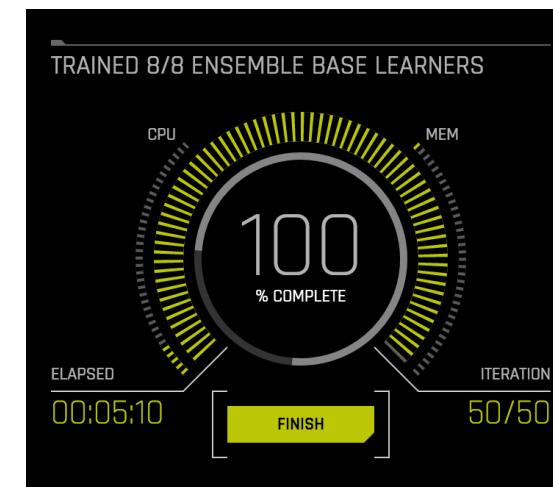
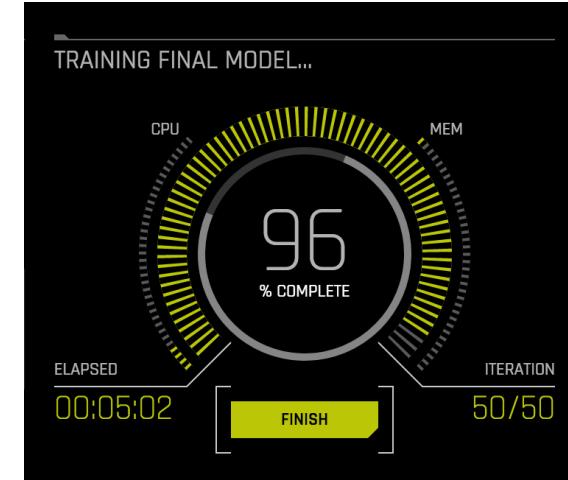
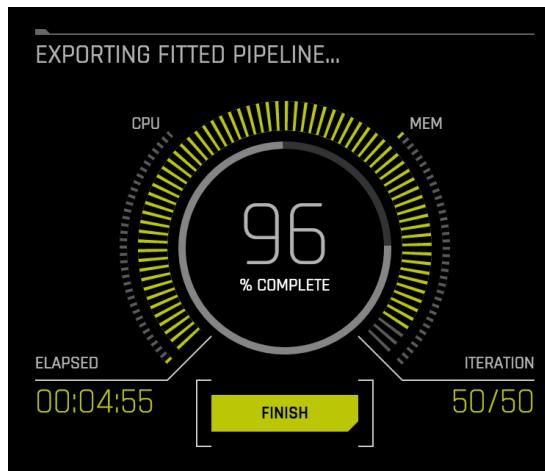
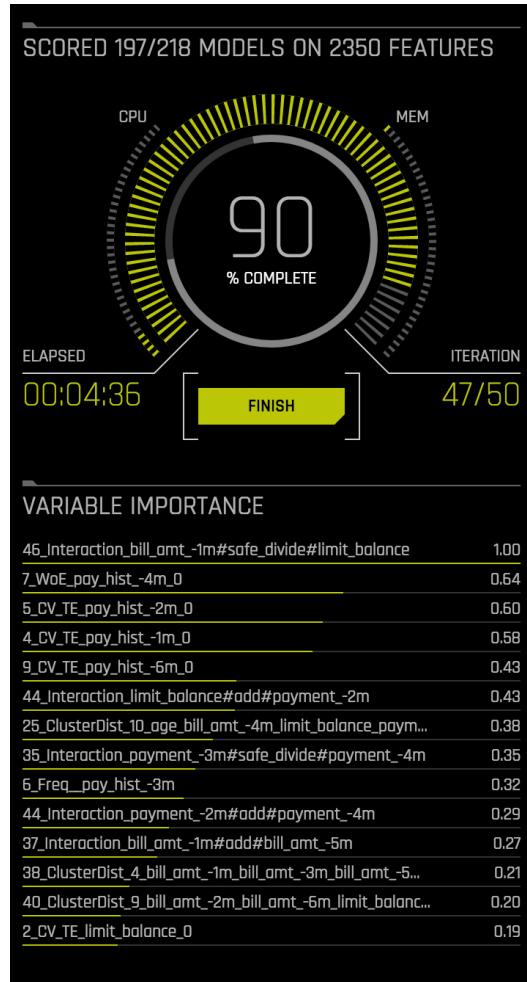


Original Features

## Feature Transformations

- Cross Validation  
Categorical Encoding
- Frequency Encoding
- Cross Validation Target  
Encoding
- Truncated SVD and More

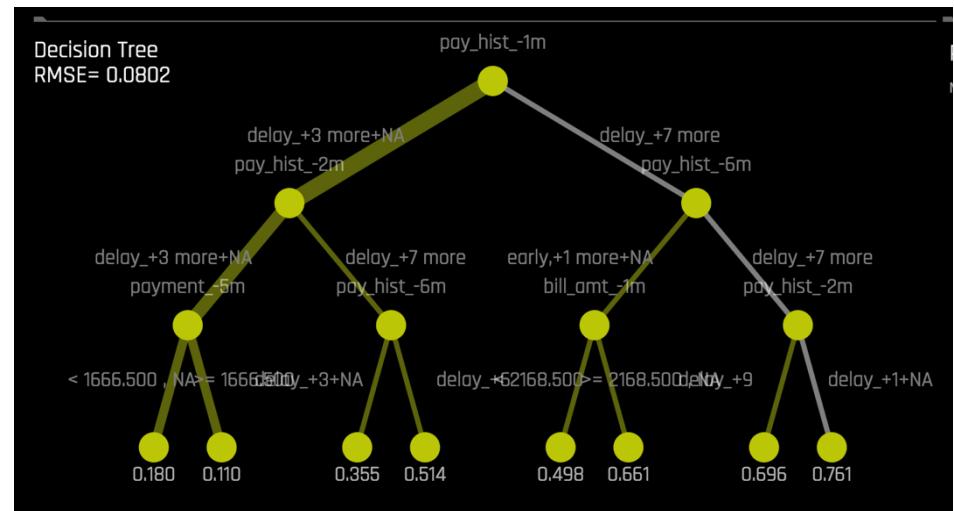
# Accurate: Auto Model Tuning/Selection/Ensemble



# Interpretable: Explainable Machine Learning

The screenshot shows a web browser window for the O'ReILLY Ideas Learning Platform. The URL is <https://www.oreilly.com/ideas/ideas-on-interpreting-machine-learning>. The page title is "Ideas on interpreting machine learning". Below the title, it says "Mix-and-match approaches for visualizing data and interpreting machine learning models and results." It was written by Patrick Hall, Wen Phan, and SriSatish Ambati on March 15, 2017. A sidebar on the left lists categories: AI, BUSINESS, DATA, DESIGN, ECONOMY, OPERATIONS, SECURITY, SOFTWARE, and SEE ALL. The "DATA SCIENCE" category is highlighted. At the bottom, there's a callout for Strata Data Conference sessions and two small diagrams illustrating neural network structures.

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



| k-LIME Local Attributions              | Variable | with value       | is associated with default |
|--|----------|------------------|----------------------------|
| <b>Top Positive Local Attributions</b> |          |                  |                            |
| pay_hist_-1m                           | delay_2m | increase of 0.35 |                            |
| pay_hist_-4m                           | delay_7m | increase of 0.1  |                            |
| pay_hist_-2m                           | delay_2m | increase of 0.06 |                            |
| pay_hist_-5m                           | delay_7m | increase of 0.05 |                            |
| age                                    | 26.000   | increase of 0.02 |                            |
| pay_hist_-3m                           | delay_7m | increase of 0.01 |                            |
| sex                                    | male     | increase of 0.01 |                            |
| bill_amt_-6m                           | 2400.000 | increase of 0    |                            |
| bill_amt_-3m                           | 2400.000 | increase of 0    |                            |

# External Review

App Dev • Machine Learning & AI • Database • Analytics & Big Data • Cloud • Open Source • Insider Articles • Reviews • Resources & White Papers

**InfoWorld**  
FROM IDG

Home > Artificial Intelligence > Machine Learning

INSIDER

## Review: H2O.ai automates machine learning

Driverless AI really is able to create and train good machine learning models without requiring machine learning expertise from users



By **Martin Heller**

Contributing Editor, InfoWorld | NOV 6, 2017

### AT A GLANCE

H2O.ai Driverless AI 1.0.5

★★★★★

LEARN MORE

on H2O.ai



Machine learning, and especially deep learning, have turned out to be incredibly useful in the right hands, as well as incredibly demanding of computer hardware.

The boom in availability of high-end GPGPUs (general purpose graphics processing units), FPGAs (field-programmable gate arrays), and custom chips such as Google's Tensor Processing Unit (TPU) isn't an accident, nor is their appearance on cloud services.

### Pros

- Driverless AI is able to create and train good models without requiring user expertise
- Good integration with Nvidia GPUs (K80 and above)
- Approximate linear models help to explain important factors in a decision
- Makes quick work of generating and evaluating many models
- Generates and exports prediction pipeline for trained model

<https://www.infoworld.com/article/3236048/machine-learning/review-h2oai-automates-machine-learning.html>

H2O.ai

[www.h2o.ai/driverless-ai-download/](http://www.h2o.ai/driverless-ai-download/)

## Docker Image

## H2O Driverless AI

Thank you very much for your interest with our new and exciting Driverless AI product. This product leverages GPU Machines with a focus on Auto Feature Engineering, Model Interpretability, and Automatic Data Visualization.

[DOWNLOAD DRIVERLESS AI](#)

[DRIVERLESS AI DOCUMENTATION](#)

Don't have a registration key? Apply [here](#) to try Driverless AI.

## 30-Day Free Trial

Here are some other resources to help you get started:

[Using Driverless AI Booklet](#)

[Machine Learning Interpretability with H2O Driverless AI Booklet](#)

[Driverless AI Data Sheet](#)

[Webinars recently delivered on Driverless AI](#)

# Demo: Credit Card Default Risk

Training  
Data

Target

Simple  
Settings

#### TRAINING DATA

DATASET

credit\_card\_train.csv

ROWS

25K

COLUMNS

24

DROPPED COLS

--

TEST DATASET

Yes

#### TARGET COLUMN

default

TYPE

int

COUNT

25000

UNIQUE

2

FREQ

19470

#### EXPERIMENT SETTINGS

5

ACCURACY

5

TIME

5

INTERPRETABILITY

CLASSIFICATION

REPRODUCIBLE

ENABLE GPUs

#### SCORER

- R2
- MSE
- RMSE
- RMSLE
- MAE
- GINI
- AUC
- LOGLOSS

LAUNCH EXPERIMENT

Optional  
Test Set

Select Metric

That's it!

# H2O.ai Experiment b1569b

1.0.5

Show Experiments

## TRAINING DATA

DATASET

credit\_card\_train.csv

ROWS  
25K

COLUMNS  
24

DROPPED COLS  
--

TEST DATASET  
Yes

TARGET COLUMN

default

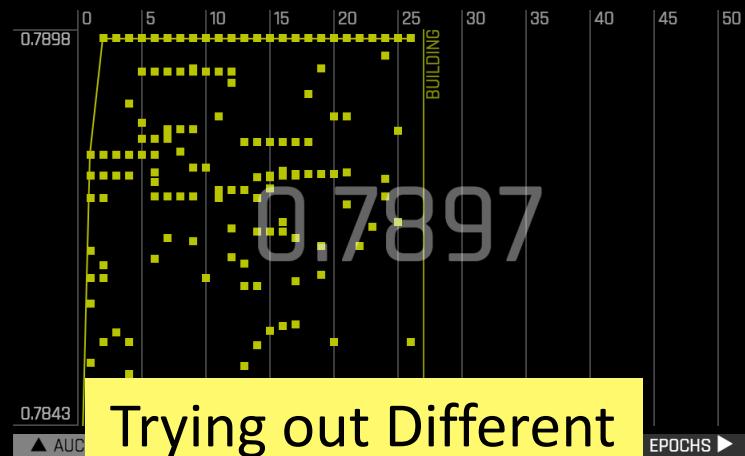
TYPE  
int

COUNT  
25000

UNIQUE  
2

FREQ  
19470

## ITERATION SCORES - INTERNAL VALIDATION



## SCORED 113/218 MODELS ON 1720 FEATURES



## EXPERIMENT SETTINGS



CLASSIFICATION

REPRODUCIBLE

ENABLE GPUs

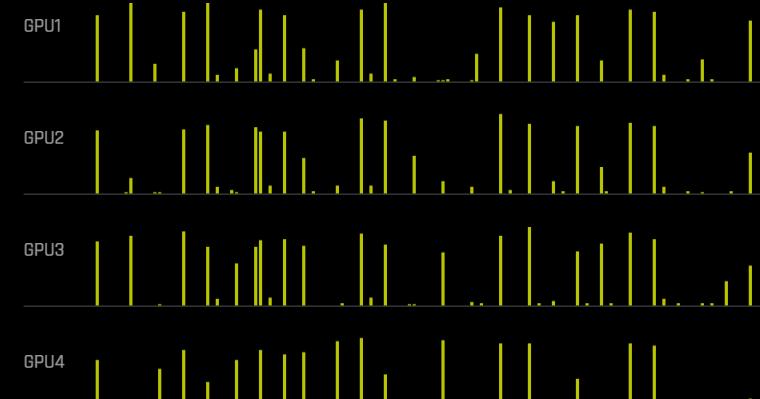
| SCORER     |
|------------|
| R2         |
| MSE        |
| RMSE       |
| RMSLE      |
| MAE        |
| GINI       |
| <b>AUC</b> |
| LOGLOSS    |

Trace

## CPU / MEMORY



## GPU USAGE



## VARIABLE IMPORTANCE

|  |      |
|--|------|
| 46_Interaction_bill_amt_-1m#safe_divide#limit_balance      | 1.00 |
| 7_WoE_pay_hist_-4m_0                                       | 0.78 |
| 5_CV_TE_pay_hist_-2m_0                                     | 0.72 |
| 4_CV_TE_pay_hist_-1m_0                                     | 0.68 |
| 25_ClusterDist_10_oge_bill_amt_-4m_limit_balance_pym...    | 0.50 |
| 44_Interaction_payment_-2m#add#payment_-4m                 | 0.42 |
| 44_Interaction_limit_balance#add#payment_-2m               | 0.40 |
| 8_WoE_pay_hist_-5m_0                                       | 0.40 |
| 6_Freq_pay_hist_-3m  | 0.37 |
| 9_CV_TE_pay_hist_-6m_0                                     | 0.36 |
| 38_ClusterDist_4_bill_amt_-1m_bill_amt_-3m_bill_amt_-5...  | 0.34 |
| 40_ClusterDist_9_bill_amt_-2m_bill_amt_-6m_limit_balanc... | 0.32 |
| 35_In  | 0.29 |
| 19_pa  | 0.25 |

Auto Feature Engineering

Using both CPUs & GPUs

# Options for Interpretation, Scoring, Logs, Pipeline

[Show Experiments](#)

## TRAINING DATA

## DATASET

credit\_card\_train.csv

## ROWS

25K

## COLUMNS

24

## DROPPED COLS

--

## TEST DATASET

Yes

## TARGET COLUMN

default

## TYPE

int

## COUNT

25000

## UNIQUE

2

## FREQ

19470

## ITERATION SCORES - INTERNAL VALIDATION



## STATUS: COMPLETE

- [INTERPRET THIS MODEL](#)
- [SCORE ON ANOTHER DATASET](#)
- [DOWNLOAD \(HOLDOUT\) TRAINING PREDICTIONS](#)
- [DOWNLOAD TEST PREDICTIONS](#)
- [DOWNLOAD TRANSFORMED TRAINING DATA](#)
- [DOWNLOAD TRANSFORMED TEST DATA](#)
- [DOWNLOAD LOGS](#)
- [DOWNLOAD SCORING PACKAGE](#)

## EXPERIMENT SETTINGS



## CLASSIFICATION

## REPRODUCIBLE

## ENABLE GPUs

| SCORER     |
|------------|
| R2         |
| MSE        |
| RMSE       |
| RMSLE      |
| MAE        |
| GINI       |
| <b>AUC</b> |
| LOGLOSS    |

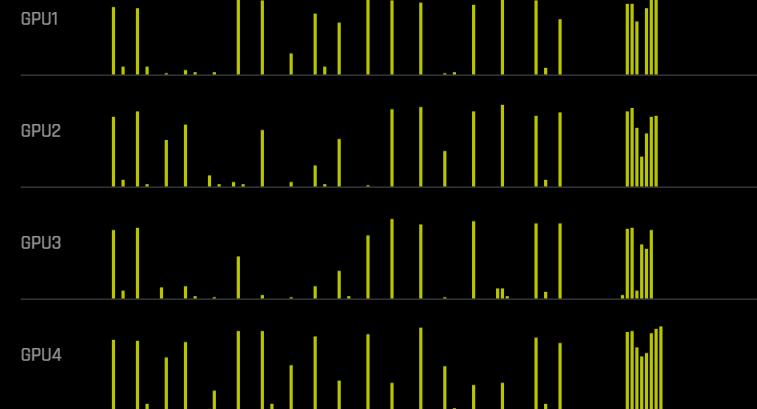
Trace

## CPU / MEMORY

## CPU

## MEM

## GPU USAGE

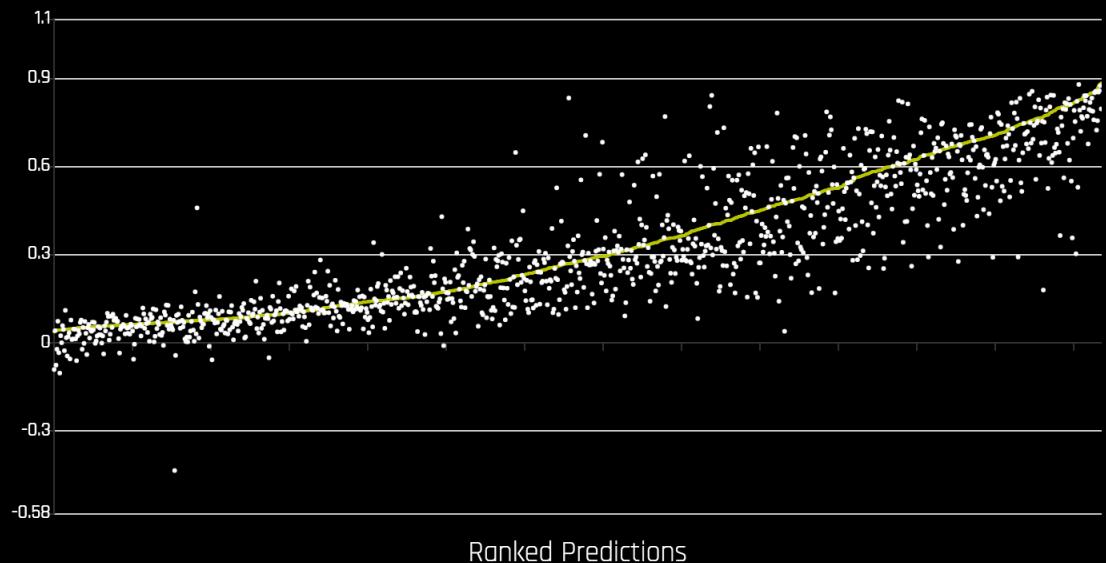


## VARIABLE IMPORTANCE

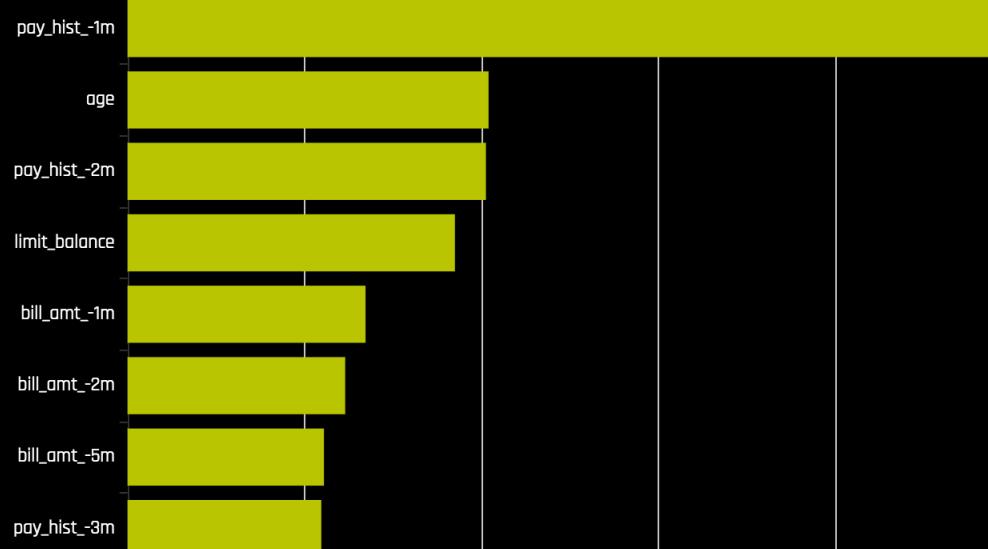
|  |      |
|--|------|
| 46_Interaction_bill_amt_-1m#safe_divide#limit_balance      | 1.00 |
| 7_WoE_pay_hist_-4m_0                                       | 0.64 |
| 5_CV_TE_pay_hist_-2m_0                                     | 0.60 |
| 4_CV_TE_pay_hist_-1m_0                                     | 0.58 |
| 9_CV_TE_pay_hist_-6m_0                                     | 0.43 |
| 44_Interaction_limit_balance#add#payment_-2m               | 0.43 |
| 25_ClusterDist_10_oge_bill_amt_-4m_limit_balance_pym...    | 0.38 |
| 35_Interaction_payment_-3m#safe_divide#payment_-4m         | 0.35 |
| 6_Freq_pay_hist_-3m  | 0.32 |
| 44_Interaction_payment_-2m#add#payment_-4m                 | 0.29 |
| 37_Interaction_bill_amt_-1m#add#bill_amt_-5m               | 0.27 |
| 38_ClusterDist_4_bill_amt_-1m_bill_amt_-3m_bill_amt_-5...  | 0.21 |
| 40_ClusterDist_9_bill_amt_-2m_bill_amt_-6m_limit_balanc... | 0.20 |
| 2_CV_TE_limit_balance_0                                    | 0.19 |

## Global Interpretable Model Explanation Plot

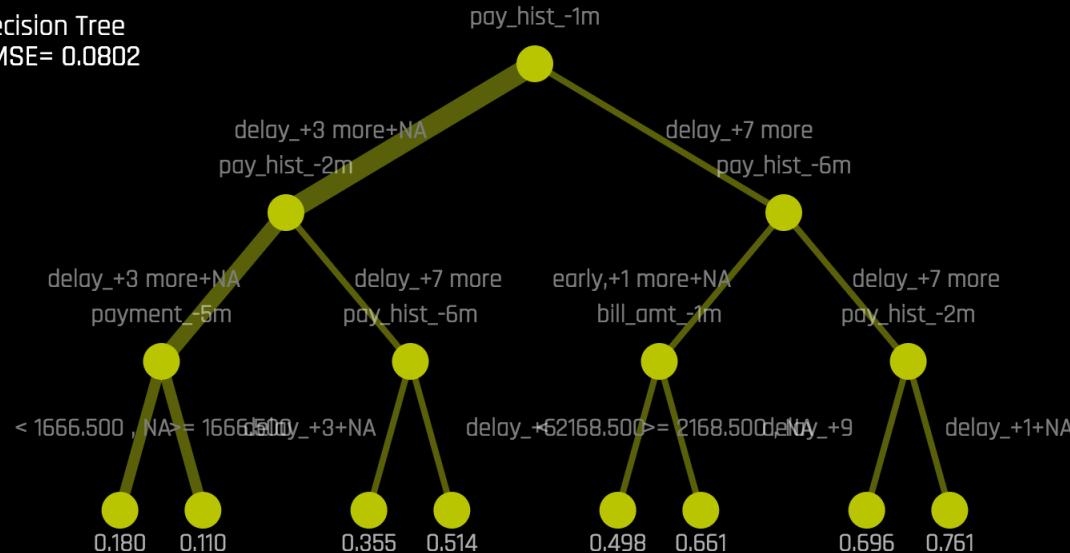
Model Prediction ● k-LIME Model Prediction ● Actual Target



## Variable Importance

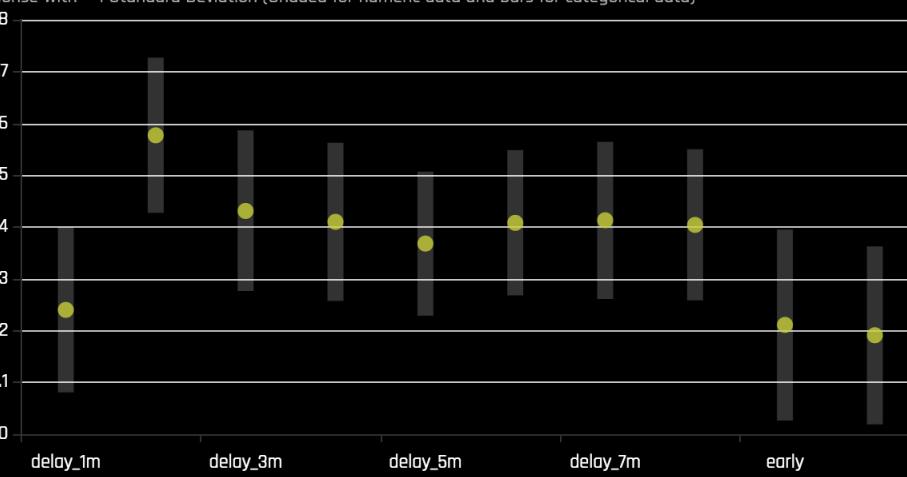


Decision Tree  
RMSE = 0.0802

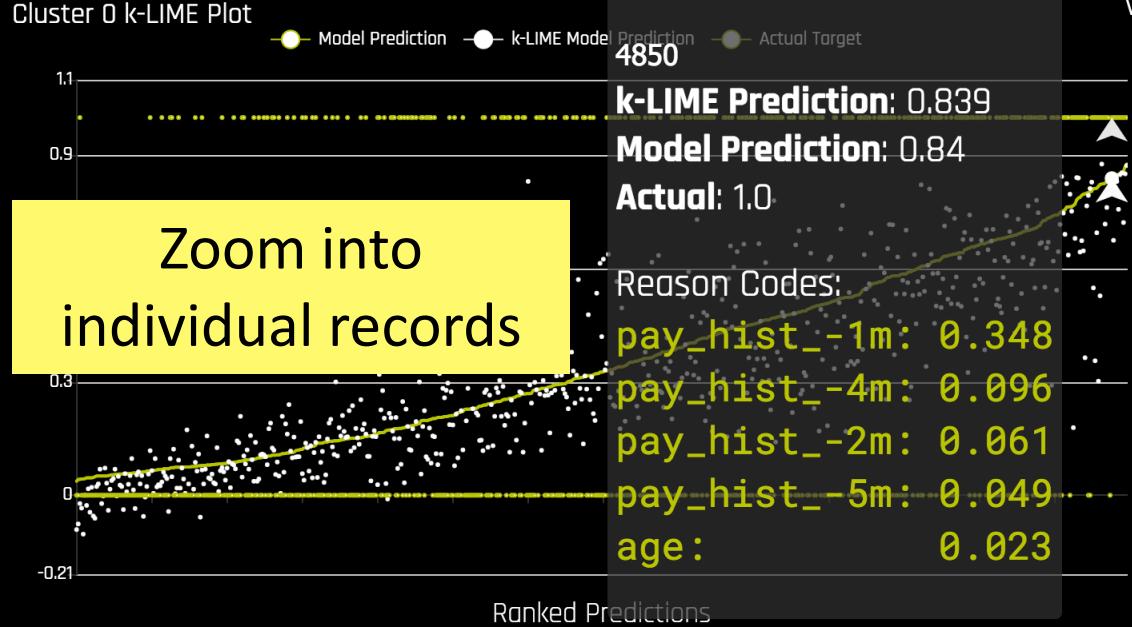


## Partial Dependence

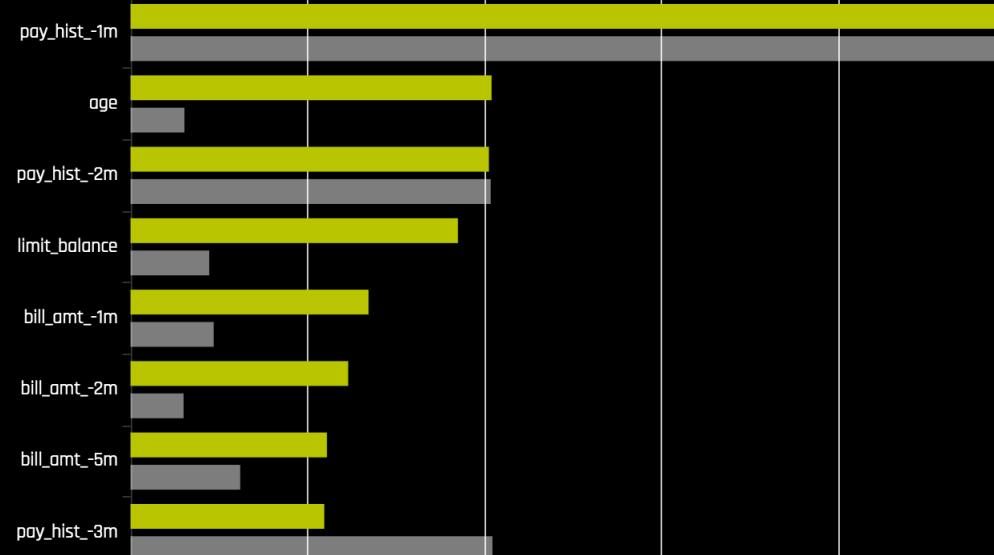
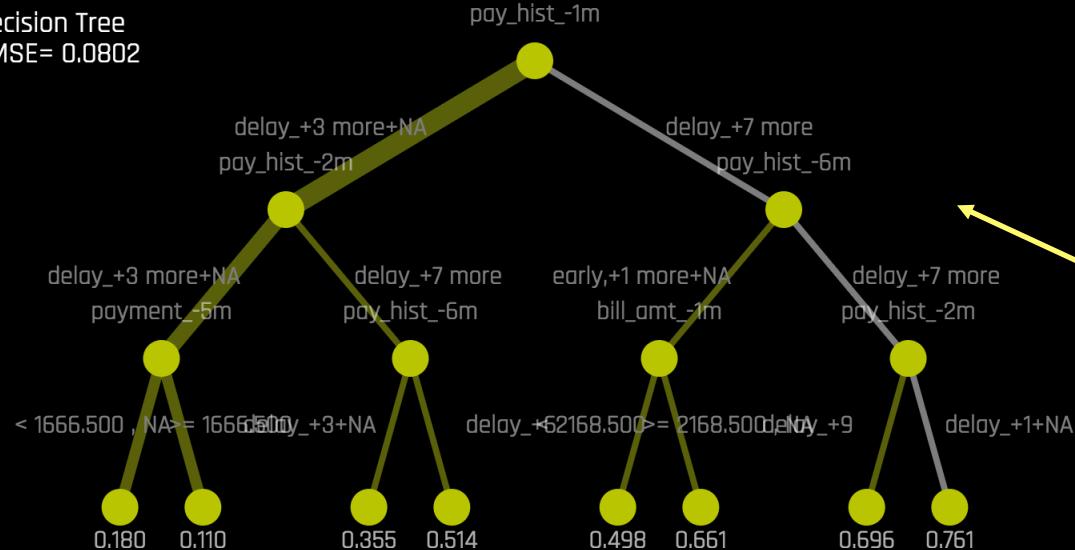
Mean Response with  $\pm 1$  Standard Deviation (Shaded for numeric data and bars for categorical data)



## Cluster 0 k-LIME Plot

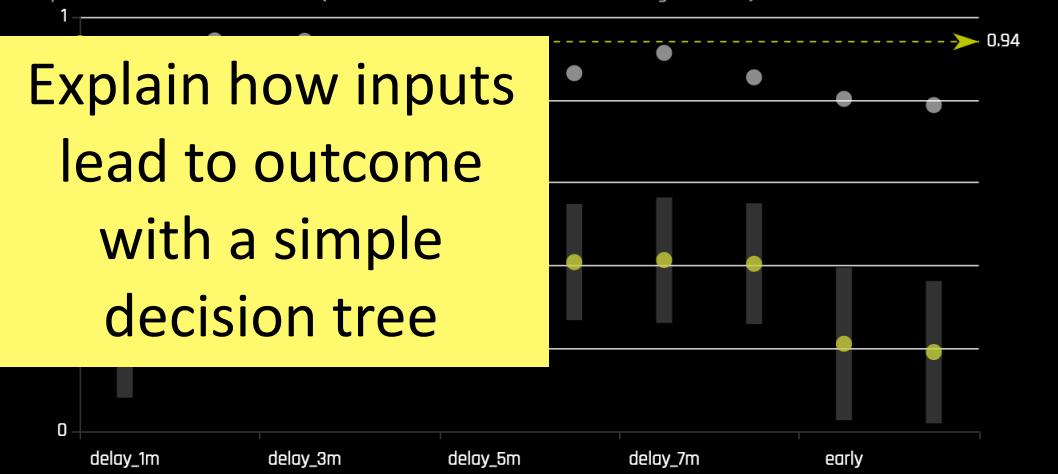


## Variable Importance

Decision Tree  
RMSE = 0.0802

## Partial Dependence

Mean Response with  $\pm 1$  Standard Deviation (Shaded for numeric data and bars for categorical data)



## Actual and Predicted Values

|                            |        |
|----------------------------|--------|
| default (Actual)           | 1      |
| Model Prediction Value     | 0.8403 |
| k-LIME Prediction Value    | 0.8392 |
| k-LIME Prediction Accuracy | 99.9%  |

## Local Reason Codes

| Top Positive Local Attributions |              |            |                               |
|---------------------------------|--------------|------------|-------------------------------|
| k-LIME Local Attributions       | Variable     | with value | is associated<br>with default |
|                                 |              |            |                               |
|                                 | pay_hist_-1m | delay_2m   | increase of 0.35              |
|                                 | pay_hist_-4m | delay_7m   | increase of 0.1               |
|                                 | pay_hist_-2m | delay_2m   | increase of 0.06              |
|                                 | pay_hist_-5m | delay_7m   | increase of 0.05              |
|                                 | age          | 26.000     | increase of 0.02              |
|                                 | pay_hist_-3m | delay_7m   | increase of 0.01              |
|                                 | sex          | male       | increase of 0.01              |
|                                 | bill_amt_-6m | 2400.000   | increase of 0                 |
|                                 | bill_amt_-3m | 2400.000   | increase of 0                 |

Turn decision tree  
into reason codes

Skipped 6 additional attributions, click to view all ...

## DOWNLOAD SCORING PACKAGE

## Download Pipeline

```
ubuntu@ip-10-10-0-226:~/scorer$ ls -l
total 86684
-rw-r--r-- 1 ubuntu ubuntu      15 Nov  8 06:56 client_requirements.txt
-rw-r--r-- 1 ubuntu ubuntu 1510256 Nov  8 06:56 datatable-0.2.2+master.301.noomp-cp36-cp36m-linux_x86_64.whl
drwxrwxr-x 6 ubuntu ubuntu    4096 Nov  8 07:01 env
-rw-r--r-- 1 ubuntu ubuntu     6738 Nov  8 06:55 example_client.py
-rw-r--r-- 1 ubuntu ubuntu    9577 Nov  8 06:55 example.py
-rw-r--r-- 1 ubuntu ubuntu   16591 Nov  8 06:56 features.txt
-rw-r--r-- 1 ubuntu ubuntu 46735145 Nov  8 06:56 h2o4gpu-0.0.4-py36-none-any.whl
-rw-r--r-- 1 ubuntu ubuntu 14778303 Nov  8 06:56 h2oaicore-1.0.5-cp36-cp36m-linux_x86_64.whl
-rw-r--r-- 1 ubuntu ubuntu   10660 Nov  8 06:56 README.txt
-rw-r--r-- 1 ubuntu ubuntu     308 Nov  8 06:56 requirements.txt
-rwxr-xr-x 1 ubuntu ubuntu     305 Nov  8 06:55 run_example.sh
-rw-r--r-- 1 ubuntu ubuntu    6219 Nov  8 06:55 run_http_client.sh
-rwxr-xr-x 1 ubuntu ubuntu     520 Nov  8 06:55 run_http_server.sh
-rwxr-xr-x 1 ubuntu ubuntu     428 Nov  8 06:55 run_tcp_client.sh
-rwxr-xr-x 1 ubuntu ubuntu     519 Nov  8 06:55 run_tcp_server.sh
-rw-r--r-- 1 ubuntu ubuntu 25622313 Nov  8 06:57 scoring_1988a3_20171108045709_7ea33-1.0.0-py3-none-any.whl
-rw-r--r-- 1 ubuntu ubuntu    2722 Nov  8 06:55 scoring.thrift
-rw-r--r-- 1 ubuntu ubuntu    9613 Nov  8 06:55 server.py
-rw-r--r-- 1 ubuntu ubuntu     29 Nov  8 06:56 server_requirements.txt
drwxrwxr-x 2 ubuntu ubuntu    4096 Nov  8 07:02 tmp
ubuntu@ip-10-10-0-226:~/scorer$ █
```

```
ubuntu@ip-10-10-0-226:~/scorer$ cat requirements.txt
numpy==1.13.1
pandas==0.19.2
scikit-learn==0.19.0
scipy==0.19.0
requests==2.13.0
pycrypto==2.6.1
hashids==1.2.0
datatable-0.2.2+master.301.noomp-cp36-cp36m-linux_x86_64.whl
h2o4gpu-0.0.4-py36-none-any.whl
h2oaicore-1.0.5-cp36-cp36m-linux_x86_64.whl
scoring_1988a3_20171108045709_7ea33-1.0.0-py3-none-any.whl
ubuntu@ip-10-10-0-226:~/scorer$
```

## Standalone Python Package

```
# Create a singleton Scorer instance.
# For optimal performance, create a Scorer instance once, and
#
scorer = Scorer()

#
# To score one row at a time, use the Scorer.score() method:
#

print('----- Score Row -----')
print(scorer.score([
    336625, # limit_balance
    'female', # sex
    'high_school', # education
    'single', # marriage
    41, # age
    'delay_2m', # pay_hist_-1m
    'delay_7m', # pay_hist_-2m
    'delay_5m', # pay_hist_-3m
    'delay_7m', # pay_hist_-4m
    'on_time', # pay_hist_-5m
    'delay_7m', # pay_hist_-6m
    780288, # bill_amt_-1m
    781180, # bill_amt_-2m
    677013, # bill_amt_-3m
    417082, # bill_amt_-4m
    -16183, # bill_amt_-5m
    -165806, # bill_amt_-6m
    684022, # payment_-1m
    450296, # payment_-2m
    125650, # payment_-3m
    109434, # payment_-4m
```

# Q & A