

Introduction to H₂O Driverless AI

Use Cases and Live Demos



 FinTech
Scotland 2018

In association with VisitScotland Business Events

Jo-fai (Joe) Chow

Data Science Evangelist /
Community Manager

joe@h2o.ai

@matlabulous

More Info → [https://bit.ly/
h2o_meetups](https://bit.ly/h2o_meetups)

About Me



Glasgow 2017



Glasgow 2018

- **Before H₂O**

- Water Engineer / EngD Researcher / Matlab Fan Boy
(wonder why @matlabulous?)
- Discovered R, Python, H₂O ... never look back again
- Data Scientist at Virgin Media (UK), Domino Data Lab (US)

- **At H₂O ...**

- Data Scientist / Evangelist /
 - Sales Engineer / Solution Architect /
 - Community Manager
- ... The harsh reality of startup life ...

Reminder: #360Selfie

H₂O.ai

H2O.ai Overview

Company	Founded in Silicon Valley in 2012 Funded: \$75M Investors: Wells Fargo, NVIDIA, Nexus Ventures, Paxion Ventures
Products	<ul style="list-style-type: none">• H2O Open Source Machine Learning (14,000 organizations)• H2O Driverless AI – Automatic Machine Learning
Leadership	Leader in Gartner MQ Machine Learning and Data Science Platform
Team	120 AI experts (Kaggle Grandmasters, Distributed Computing, Visualization)
Global	Mountain View, London, Prague, India



Scientific Advisory Council



Dr. Trevor Hastie

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



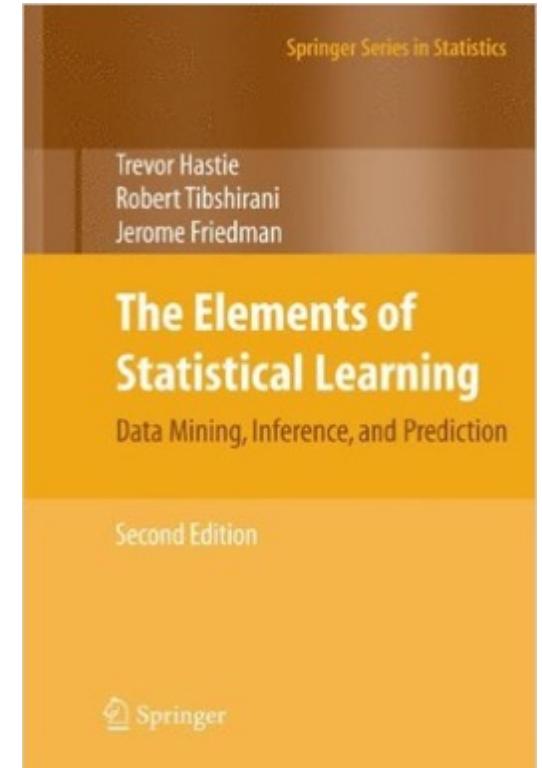
Dr. Robert Tibshirani

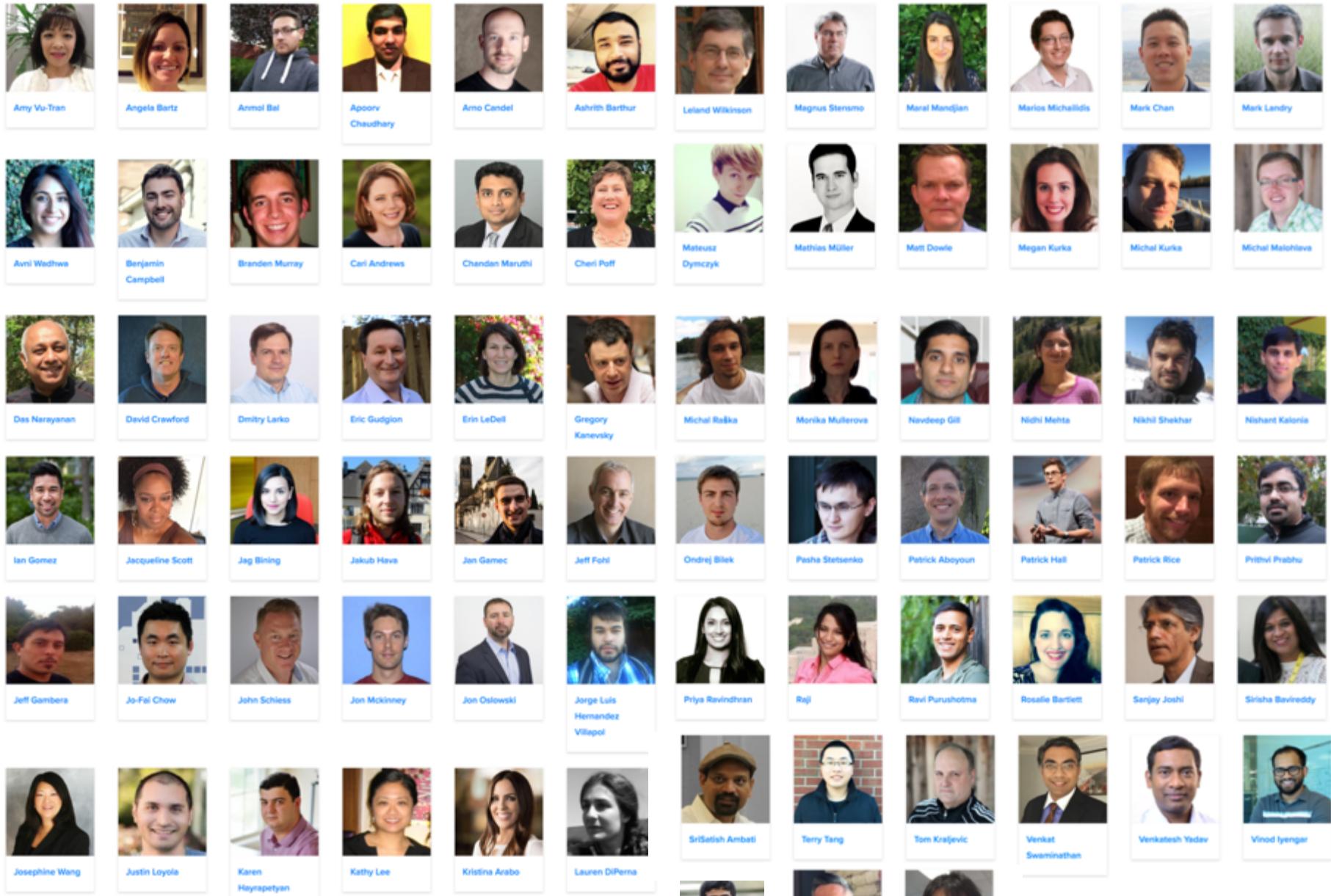
- Professor of Statistics and Health Research and Policy, Stanford University
- PhD in Statistics, Stanford University
- Co-author, *The Elements of Statistical Learning: Prediction, Inference and Data Mining*
- Author, *Regression Shrinkage and Selection via the Lasso*
- Co-author, *An Introduction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



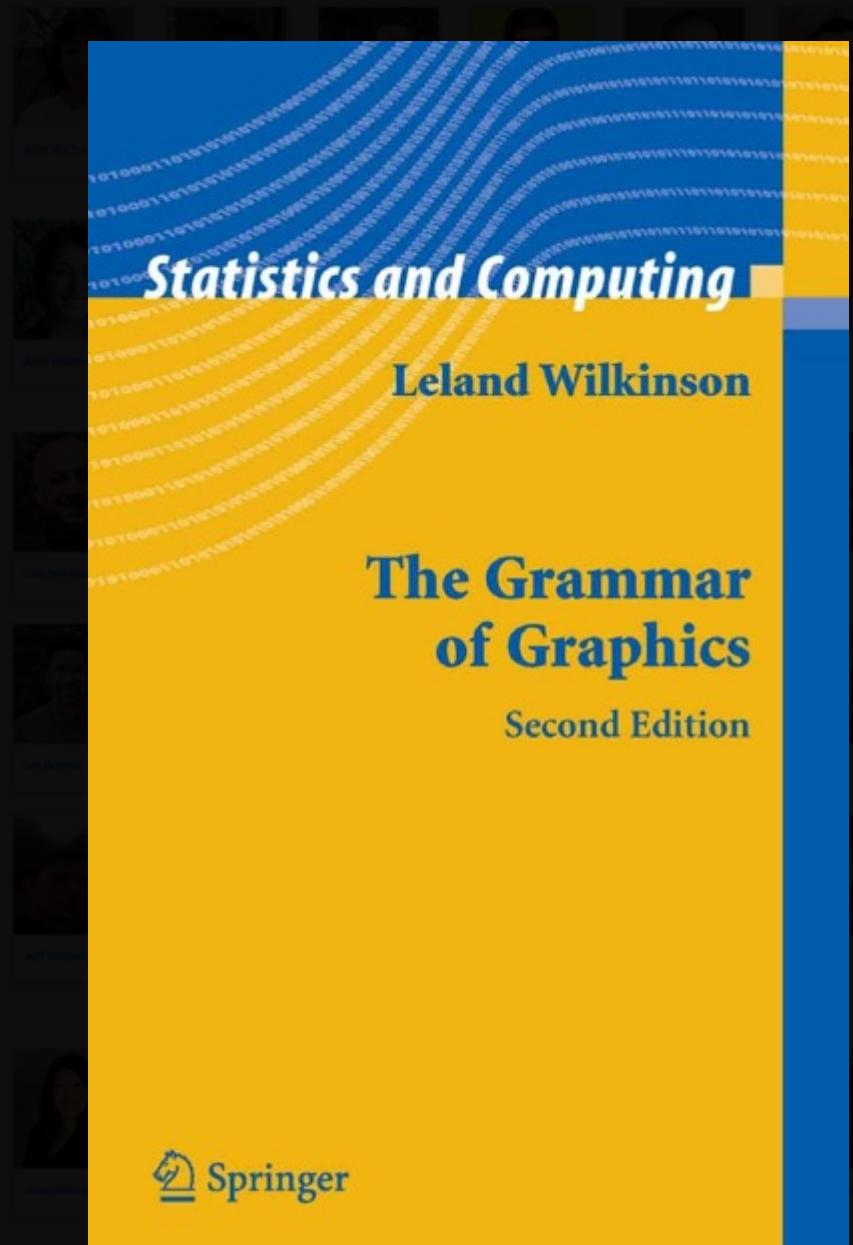
Dr. Steven Boyd

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





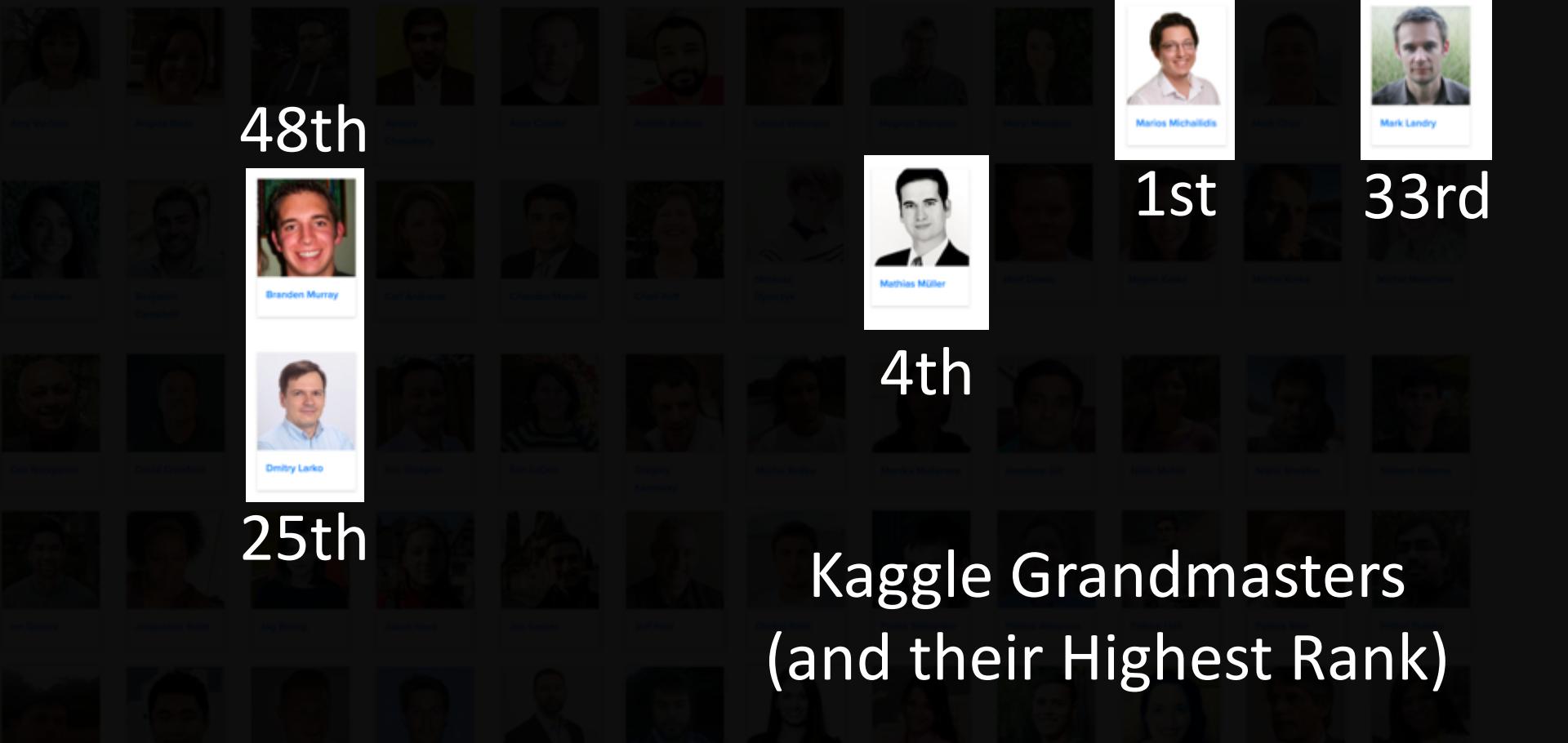
H₂O Team



Leland Wilkinson

Origin of R Package `ggplot2`





 113
Grandmasters

 980
Masters

 3,339
Experts

 46,135
Contributors

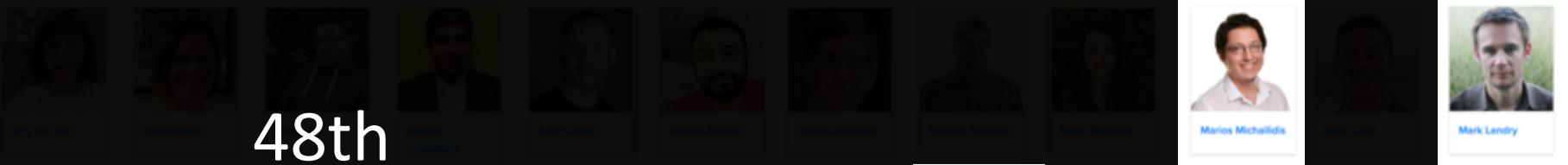
 33,242
Novices

About 80,000 Kagglers

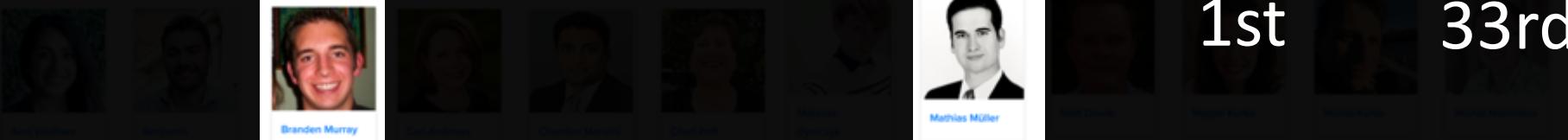
H₂O Team

13th

H₂O.ai



48th



Branden Murray

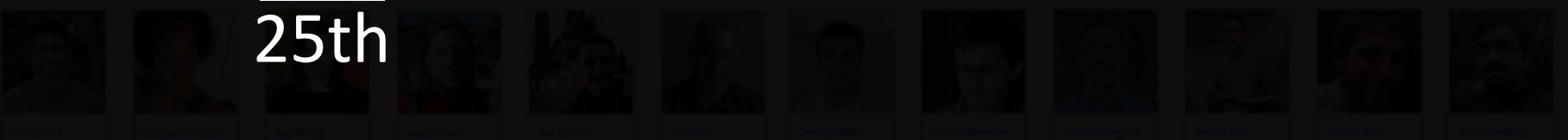
1st

33rd

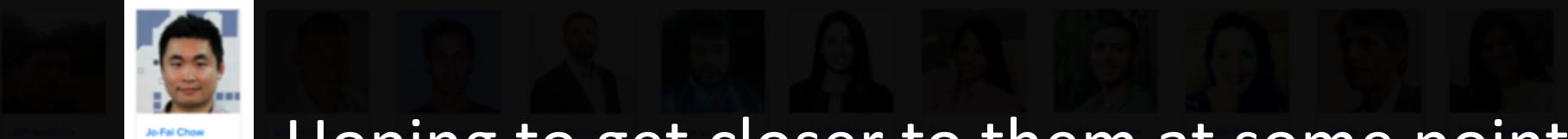


Dmitry Larko

4th

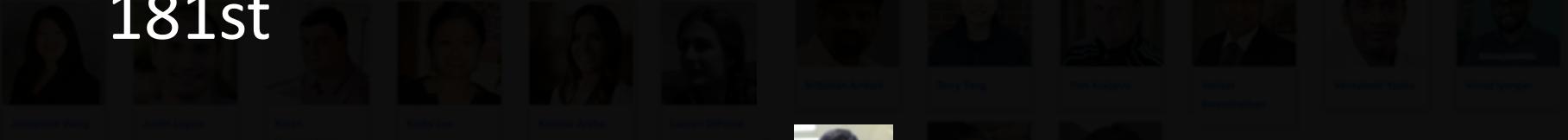


25th



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

181st



13th

H₂O Team

H₂O.ai

H2O.ai Product Suite



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

- 100% open source – Apache V2 licensed
- Built for data scientists – interface using R, Python on H2O Flow (interactive notebook interface)
- Enterprise Support subscriptions

- Enterprise software
- Built for domain users, analysts & data scientists – GUI based interface for end-to-end data science
- Fully automated machine learning from ingest to deployment
- User licenses on a per seat basis (annual subscription)

Worldwide Recognition in the H2O.ai Community

Open source
community

222 OF FORTUNE
THE 500
 H₂O

8 OF TOP 10
BANKS

7 OF TOP 10
INSURANCE COMPANIES

4 OF TOP 10
HEALTHCARE COMPANIES

Paying Customers



"H2O.ai's reference customers gave it the highest overall score for sales relationship and overall service and support" - Gartner MQ 2018

H₂O.ai

Partner Ecosystem

SCAN
computers



PENGUIN
COMPUTING

AMAX

SUPERMICRO



Google
Cloud Platform

Microsoft Azure



IBM Cloud

NVIDIA

IBM

TRACE3
World Wide Technology, Inc.

TechData

accenture

Capgemini

snowflake

Hortonworks

cloudera

MINIO

MAPR

HW Vendors

Cloud Providers

Strategic
Partners

Value Added
Resellers

System
Integrators

Data Stores

H2O.ai

H2O.ai is a **Leader** in the 2018 Gartner Data Science and Machine Learning Platforms Magic Quadrant

- Technology leader with most completeness of vision
- Recognized for the mindshare, partner network and status as a **quasi-industry standard** for machine learning and AI
- H2O.ai 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



Get the
Gartner
Magic
Quadrant
[here](#)

Why Driverless AI?

Why Driverless AI for **Enterprise AI** adoption

Data is a Team Sport

~100

Data science experts in the world

Lack of AI Talent

Months to
Insight

Time for a data scientist to
build a model

Time to Insights slow

"US alone faces a shortage of 190,000
people with analytical expertise."

McKinsey&Company



Black box models

Lack of Trust in AI

Driverless AI Delivers **AI for Enterprise**

Talent

Top 10

Data Science Experts

Kaggle Grandmasters

Time

Months
to Hours

GPU accelerated ML
Automatic Pipelines

Time to Insight

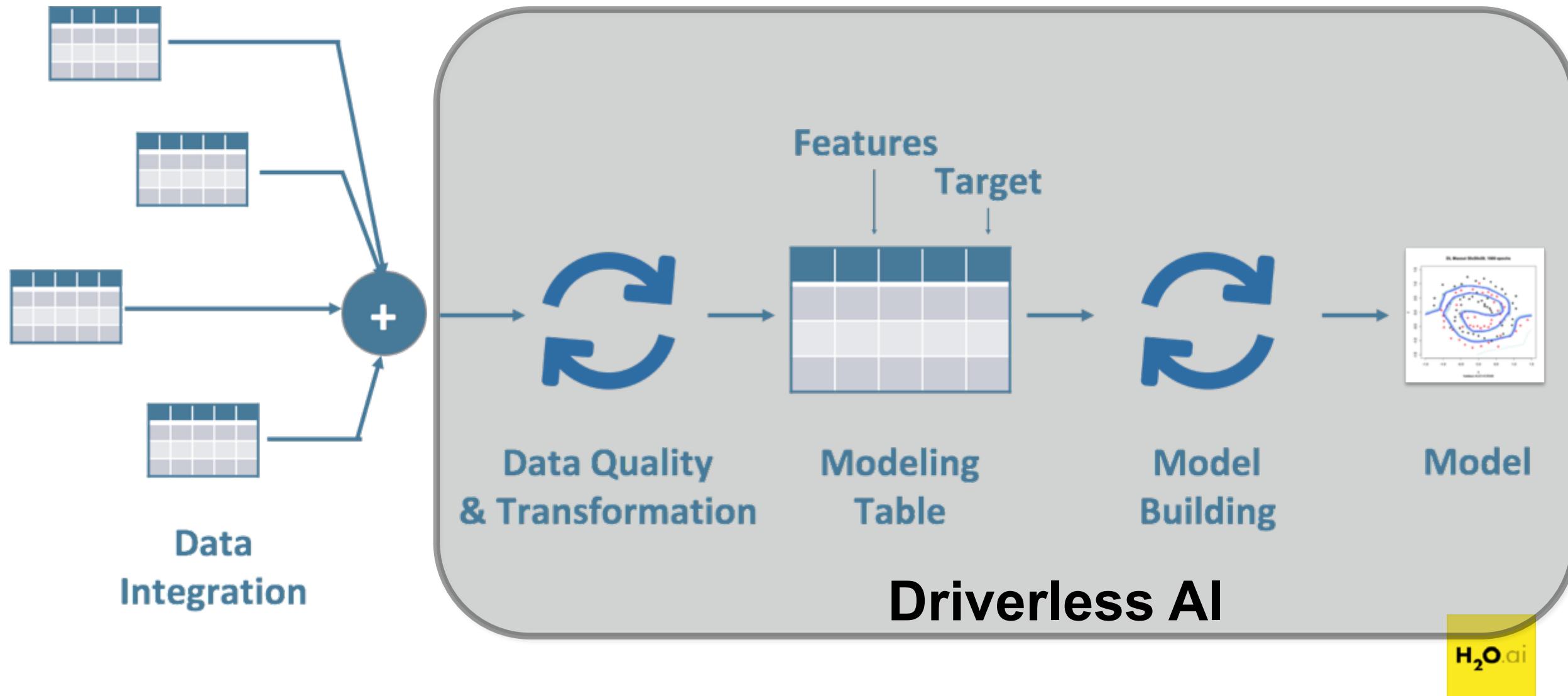
Trust

MLI
Auto Doc

Auto Visualization

Explainability & Transparency

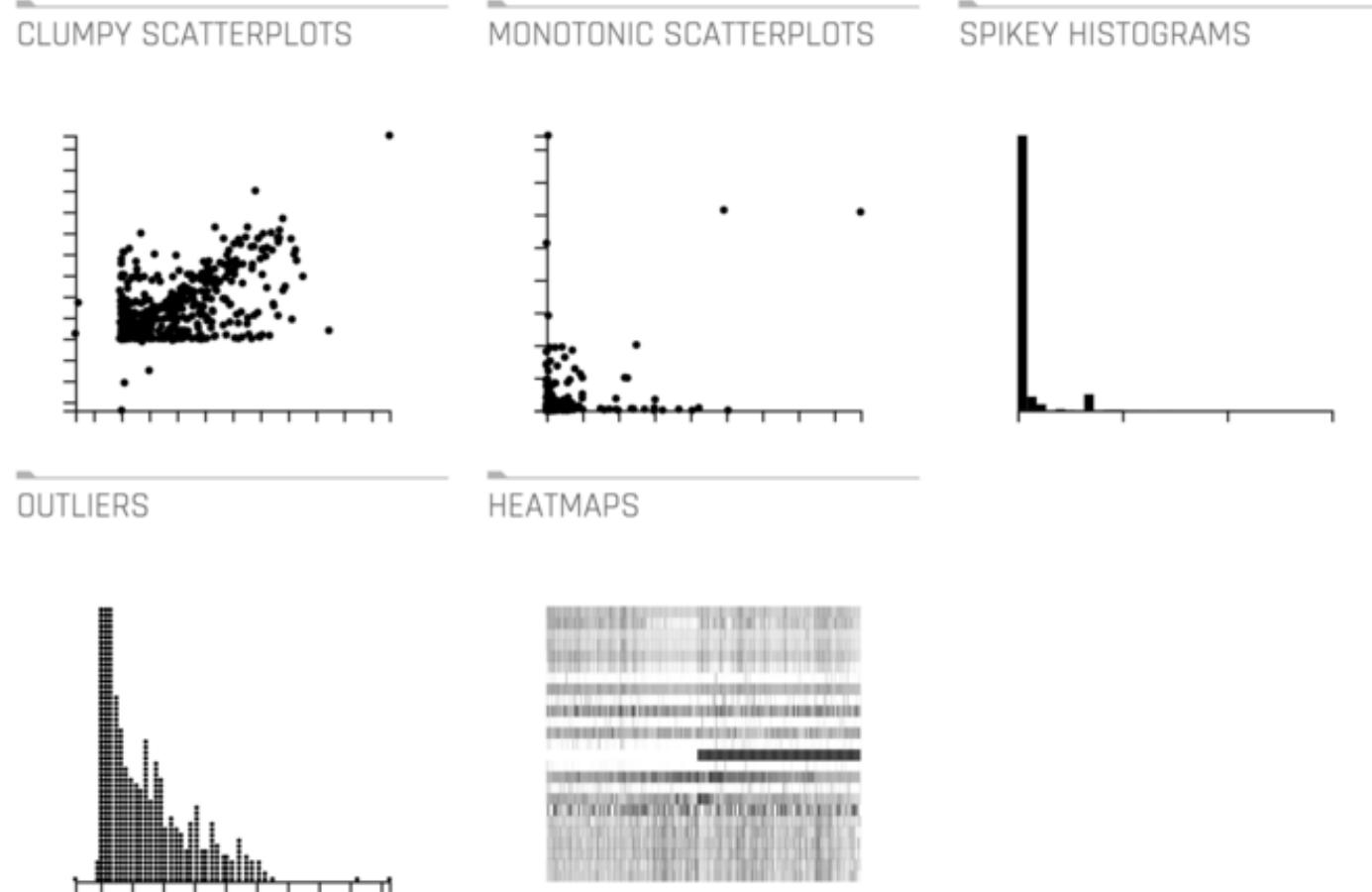
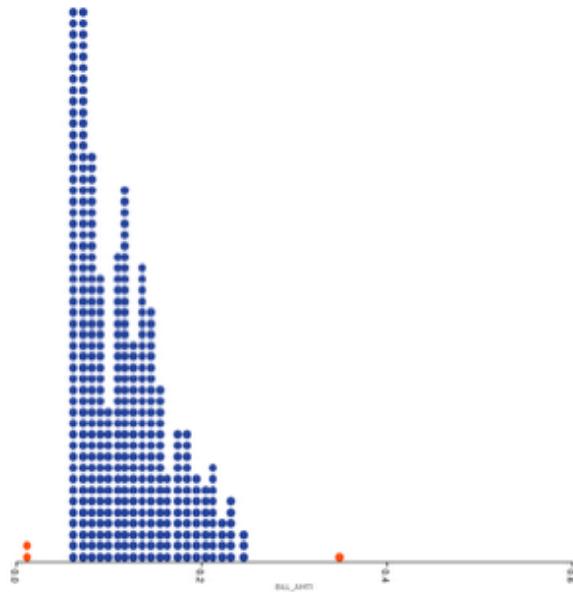
Driverless AI: Automates Data Science and ML Workflows



Automatic Visualization

H2O.ai

Automatic Scagnostics and other visualizations to generate the most relevant visualizations for each dataset



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

H₂O.ai

Secret Sauce: 1) Grandmaster Feature Engineering

VARIABLE IMPORTANCE

632_WoE:v56:v79.0	1.00
441_WoE:v13:v66.0	0.25
197_v50	0.20
94_CVTE:v66.0	0.18
556_NumToCatWoE:v50.0	0.13
621_NumToCatWoE:v114:v49:v72:v73.0	0.12
588_InteractionMul:v114:v50	0.11
592_TruncSVD:v12:v50:v85:v90.1	0.09
417_ClusterDist7:v104:v114:v50.1	0.06
535_NumToCatWoE:v114:v25:v39:v49:v95.0	0.05
607_NumToCatTE:v114:v46:v62:v65:v98.0	0.05
616_TruncSVD:v12:v15:v18:v50.1	0.04
438_CVTE:v107:v127:v22:v37:v5.0	0.04
430_ClusterDist7:v49:v50.1	0.04

Numerical/Categorical Interactions, Target Encoding, Clustering, Dimensionality Reduction, Weight of Evidence, etc.

VARIABLE IMPORTANCE

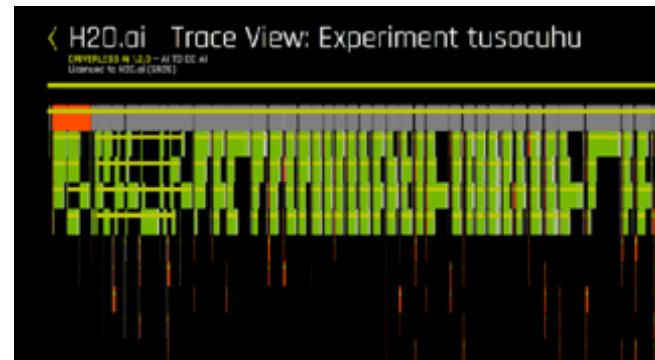
11_EWMA(0.05)(0)TargetLags:Dept:Store.44:45:51:52:53:105	1.00
18_TargetLog:Dept:Store.52	0.32
22_ClusterDlst4:Dept.3	0.20
13_LagsMax:Dept:Store.Weekly_Sales.44:45:51:52:53:105	0.15
22_ClusterDlst4:Dept.0	0.15
33_Freq:Dept:Store	0.15
18_TargetLog:Dept:Store.44	0.14
22_ClusterDlst4:Dept.1	0.11
18_TargetLog:Dept:Store.45	0.10
2_Freq:Store	0.08
22_ClusterDlst4:Dept.2	0.05
31_EWMA(0.05)(0)TargetLogs:Dept.44:45:51:52:53:105	0.05
25_TargetLog:Store.44	0.05
13_LagsMedian:Dept:Store.Weekly_Sales.44:45:51:52:53:1...	0.04

Time-Series: Lags and historical aggregates with causality constraints

Secret Sauce: 2) Grandmaster Pipeline Tuning + Validation

Example: Driverless AI BNP Paribas on 3-GPU workstation

Recipe: AutoDL (171 iterations, 12 individuals)
Validation scheme: stratified, 1 internal holdout
Feature engineering: 18923 features tested (344 selected)
Timing:
Data preparation: 8.44 secs
Model parameter tuning: 403.98 secs (19 models trained)
Feature engineering: 15424.53 secs (1008 models trained)
Final model training: 1935.21 secs (26 models trained)
Validation score: LOGLOSS = 0.47811 +/- 0.0023019 (baseline)
Validation score: LOGLOSS = 0.43681 +/- 0.0037107 (final model)
Test score: LOGLOSS = N/A (no target)

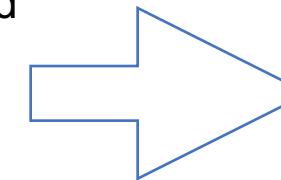


MTV

evolutionary strategies

massively parallel processing
(multi-CPU, multi-GPU)

19,000 features tested

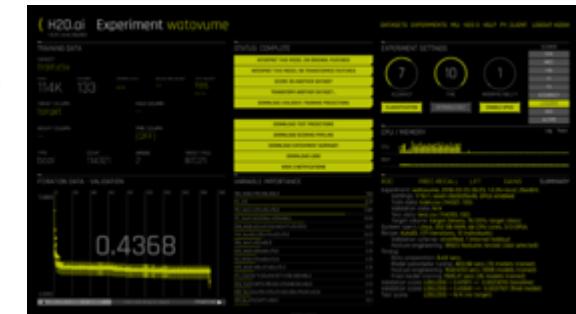
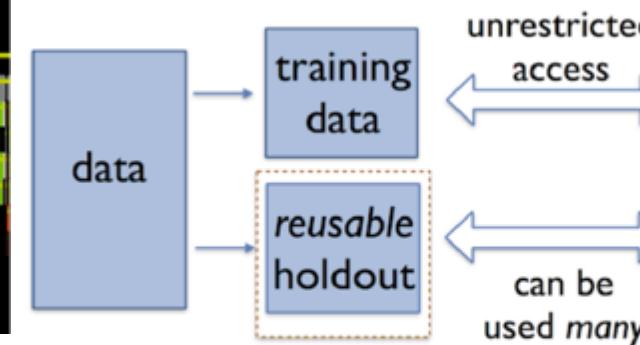


1 final optimal
scoring pipeline

1,000 models trained

reliable generalization estimates (overfitting avoidance)

Reusable holdout method



essentially as good as
using fresh data each time!

Accuracy

- Automatic feature engineering to increase accuracy - AlphaGo for AI
- Automatic Kaggle Grandmaster recipes in a box for solving wide variety of use-cases
- Automatic machine learning to find and tune the right ensemble of models

Driverless AI: top 5% in Amazon Kaggle competition

Driverless AI produces feature engineering pipeline (“more columns”) for downstream use



Amazon.com - Employee Access

Predict an employee's access needs
\$5,000 - 1,687 teams - 4 years ago

Driverless AI: 80th place (out of 1687 - top 5%)

Driverless AI: Top-10 in BNP Paribas Kaggle competition



single run, fully automated: 2h on DGX Station! 6h on PC

BNP Paribas Cardif Claims Management

Can you accelerate BNP Paribas Cardif's claims management process?

\$30,000 - 2,926 teams - 2 years ago

Submission and Description
[sub.csv](#)
2 minutes ago by Arna Candel
94098f7/10/1 cv-0.4354 finished after 172 iter

Private Score: 0.42945
Public Score: 0.43156

#	In the money	Team Name	Kernel	Team Members	Score (P)	Entries	Last
1	-	Dexter's Lab			0.42037	198	2y
2	-	escalated chi			0.42079	162	2y
3	-	Exploding Kittens			0.42182	524	2y
4	-	Brandon Nickel (utility)			0.42259	254	2y
5	-	the flying bumble brothers			0.42450	264	2y
6	-	n_m			0.42535	4	2y
7	-	PAFY			0.42557	310	2y
8	-	KAIME			0.42688	121	2y
9	-	Jack (Dapper)			0.42744	22	2y
10	+1	Ondrej & Bohdjan			0.43000	192	2y
11	+1	Li-Der			0.43006	56	2y
12	+2	BIGMEIPRS			0.43089	338	2y
13	-	x2Red			0.43187	55	2y
14	-	Franchies			0.43146	134	2y
15	+1	Aina			0.43168	55	2y
16	+1	maze numbers			0.43362	164	2y
17	-	BB-B2			0.43313	129	2y
18	+3	no one			0.43367	68	2y

Driverless AI: 10th place in private LB at Kaggle (out of 2926)

2 months for Grandmasters — 2 hours for Driverless AI

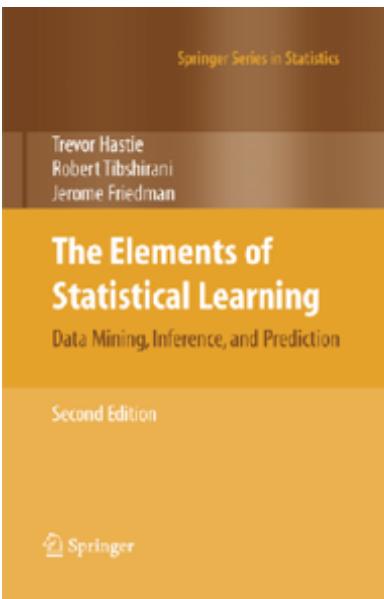
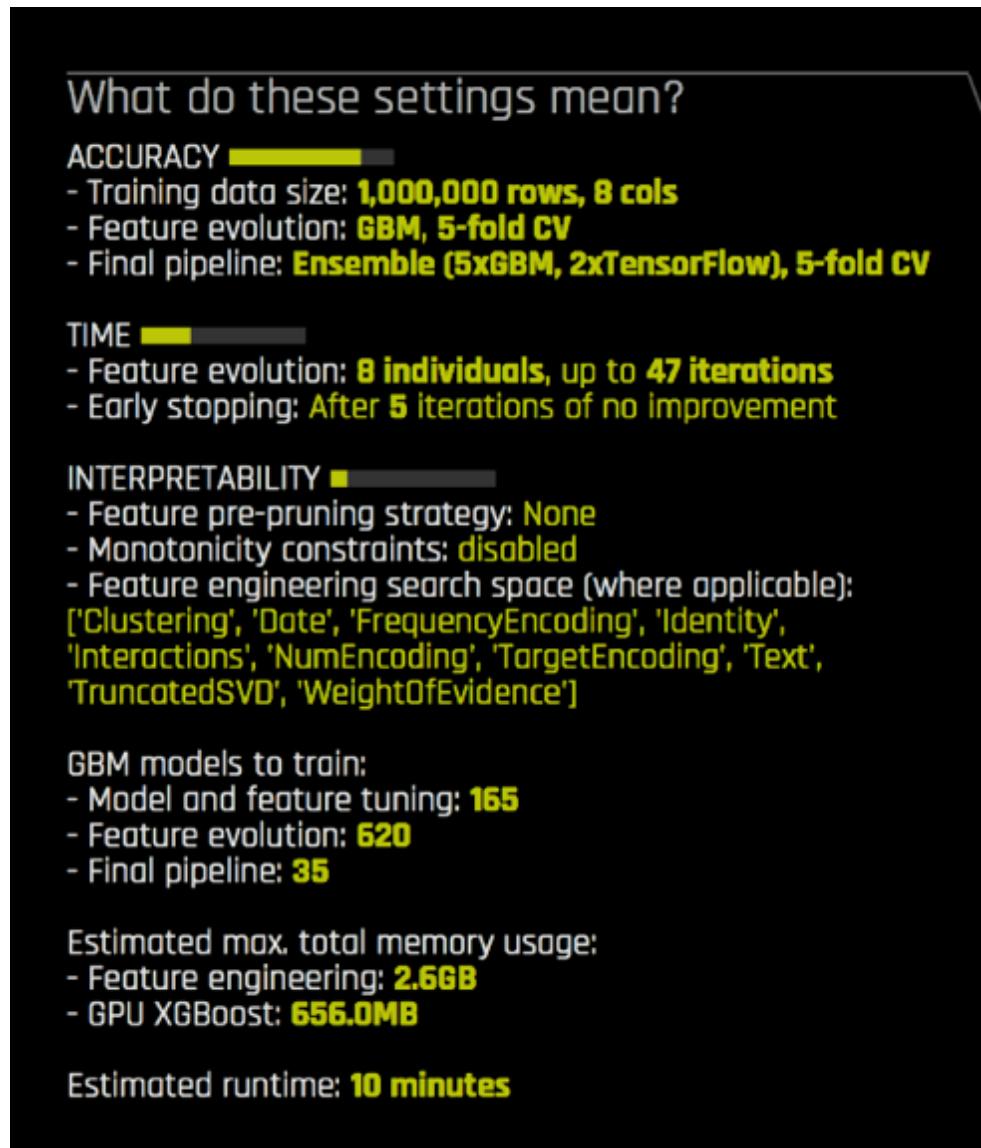


Interpretability

- Interpretability for debugging, not just for regulators
- Get reason codes and model interpretability in plain english
- K-Lime, LOCO, partial dependence and more



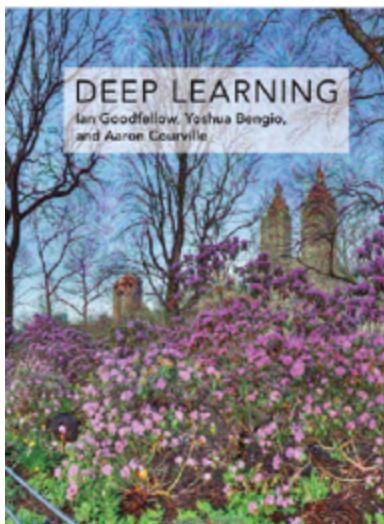
Statistical Learning vs Deep Learning - We Do Both!



GLM/CART/RF/GBM/XGBoost
K-Means/PCA/SVD

Typically better for structured data
(CSV, SQL, Transactional)

<https://web.stanford.edu/~hastie/Papers/ESLII.pdf>



TensorFlow Deep Learning

Typically better for unstructured data
(Images, Video, Audio, Text)

<http://www.deeplearningbook.org>



Live Demo


BNP PARIBAS CARDIF

BNP Paribas Cardif Claims Management

Can you accelerate BNP Paribas Cardif's claims management process?
\$30,000 · 2,926 teams · 2 years ago

Overview Data Kernels Discussion Leaderboard Rules Team My Submissions Late Submission

Overview

Description	As a global specialist in personal insurance, BNP Paribas Cardif serves 90 million clients in 36 countries across Europe, Asia and Latin America.
Evaluation	In a world shaped by the emergence of new uses and lifestyles, everything is going faster and faster. When facing unexpected events, customers expect their insurer to support them as soon as possible. However, claims management may require different levels of check before a claim can be approved and a payment can be made. With the new practices and behaviors generated by the digital economy, this process needs adaptation thanks to data science to meet the new needs and expectations of customers.
Prizes	
Timeline	
About Bnp Paribas Cardif	



In this challenge, BNP Paribas Cardif is providing an anonymized database with two categories of claims:

1. claims for which approval could be accelerated leading to faster payments
2. claims for which additional information is required before approval

Kagglers are challenged to predict the category of a claim based on features available early in the process, helping BNP Paribas Cardif accelerate its claims process and therefore provide a better service to its customers.

Driverless AI Delivers “Expert Data Scientist in a Box”

- Created and supported by world renowned AI experts
- Empowers companies to accomplish AI and ML with a single platform
- Performs the function of an expert data scientist and adds more power to both novice and expert teams
- Details and highlights insights and interpretability with easy to understand results and visualizations



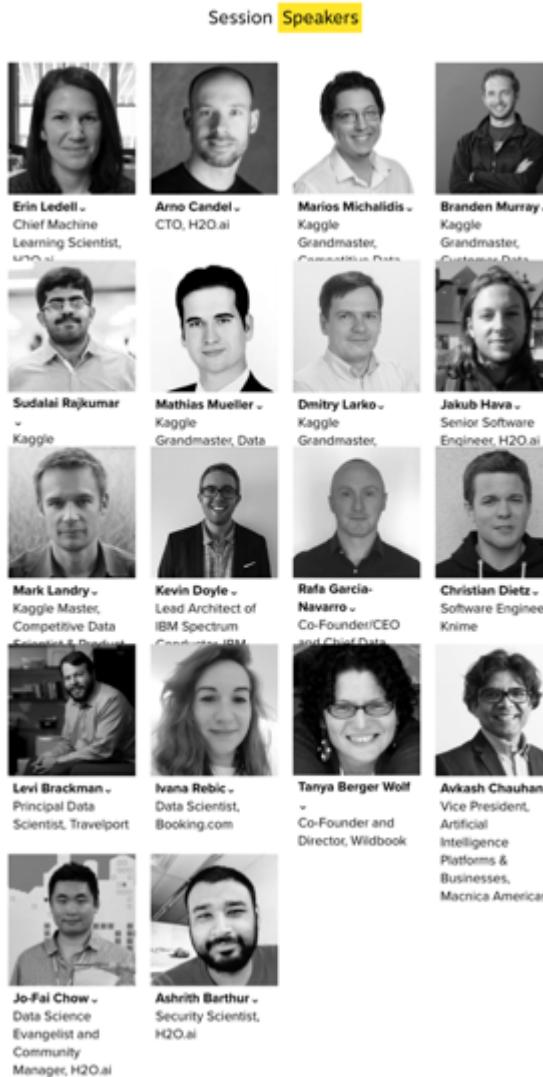
21 day free trial for [Driverless AI](#)

H₂O.ai

Our Flagship Community Event – H2O AI World is finally coming to London!



29th & 30th Oct, London



More real-world use cases + All H₂O Kaggle Grandmasters + Hands-on Training

H₂O.ai

Thanks!



- More Info, Code, and Slides
 - [bit.ly/
h2o_meetups](https://bit.ly/h2o_meetups)
- Contact
 - joe@h2o.ai
 - [@matlabulous](https://twitter.com/matlabulous)
 - github.com/woobe