# H2O Introductory Training

Erin LeDell Ph.D.

Machine Learning Scientist, H2O.ai
May 3, 2016

H2O.ai

- Statistician & Machine Learning Scientist at H2O.ai in Mountain View, California, USA

- Ph.D. in Biostatistics with Designated Emphasis in Computational Science and Engineering from UC Berkeley (focus on Machine Learning)

- Worked as a data scientist at several startups

- Written a handful of machine learning R packages
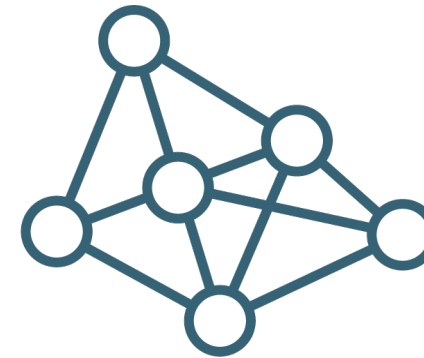
What is Data Science?

Data Science Tools?

What is Machine Learning?

What is Deep Learning?

What is Ensemble Learning?

R & Python Code Tutorial

H₂O.ai

Data Processing

- Identify a data task or prediction problem
- Collect relevant data

Problem Formulation

- Clean, transform, filter, aggregate, impute
- Convert into X and Y

Machine Learning

- Train models
- Evaluate models

H2O.ai

The number of data scientists has doubled over the last 4 years.

The top five skills listed by data scientists:

1. Data Analysis
2. R
3. Python
4. Data Mining
5. Machine Learning

H₂O.ai

2013 was the year of the data science "language wars."

H₂O.ai

In 2016, we have evolved beyond this…

We are too busy doing actual data science!

H2O.ai

We are headed toward language agnostic data science, where friendly APIs connect to powerful data processing engines.

"Field of study that gives computers the ability to learn without being explicitly programmed."

— Arthur Samuel, 1959

Unlike rules-based systems which require a human expert to hard-code domain knowledge directly into the system, a machine learning algorithm learns how to make decisions from the data alone.

H₂O.ai

## Regression

- Predict a real-valued response (e.g. viral load, price)
- Gaussian, Gamma, Poisson, etc. distributed response
- Evaluate with MSE or R^2
- Super Learner finds the optimal combination of a combination of a collection of base learning algorithms.

## Classification

- Multi-class or binary classification
- Ranking (e.g. Google Search results order)
- Evaluate with Classification Error or AUC

## Clustering

- Unsupervised learning (no training labels)
- Partition the data; identify clusters or sub-populations
- Evaluate with AIC, BIC or Total Sum of Squares

H2O.ai

| Train | Validation | Test |
|-------|------------|------|

- If you plan on doing any model tuning, you should split your dataset into three parts:  Train, Validation and Test
- There is no general rule for how you should partition the data and it will depend on how strong the signal in your data is, but an example could be:
  50% Train, 25% Validation and 25% Test
- The validation set is used strictly for model tuning (via validation of models with different parameters) and the test set is used to make a final estimate of the generalization error.

**H₂O**.ai

- K-fold Cross-validation (CV) is used to evaluate the performance of machine learning algorithms.

- CV will give you the most "mileage" on your training data.

- Performance metrics are averaged across k folds.

"A branch of machine learning based on a set of algorithms that attempt to model high-level abstractions in data by using model architectures, composed of multiple non-linear transformations."

— Wikipedia (2016)

- Deep neural networks have more than one hidden layer in their architecture. That's why they are called "deep" neural networks.
- Very useful for complex input data such as images, video, audio.

- Deep learning architectures, specifically artificial neural networks (ANNs) have been around since 1980.

- However, there were breakthroughs in training techniques that lead to their recent resurgence in the mid 2000's.

- Combined with modern computing power, they are quite effective.

"Ensemble methods use multiple learning algorithms to obtain better predictive performance that could be obtained from any of the constituent learning algorithms."

— Wikipedia (2016)

- Random Forests and Gradient Boosting Machines (GBM) are both ensembles of decision trees.

- Stacking, or Super Learning, is technique for combining various learners into a single, powerful learner using a second-level metalearning algorithm.

H₂O.ai

"Even after the observation of the frequent or constant conjunction of objects, we have no reason to draw any inference concerning any object beyond those of which we have had experience."

— David Hume (1711-1776)

- No general purpose algorithm to solve all problems.

- No right answer on optimal data preparation.

- Some algorithms may have such strong biases that they can only learn certain kinds of functions.

**H₂O**.ai

The "Intro to H2O" tutorial introduces five popular supervised machine learning algorithms in the context of a binary classification problem.

The training module demonstrates how to train models and evaluating model performance on a test set.

- Generalized Linear Model (GLM)
- Random Forest (RF)
- Gradient Boosting Machine (GBM)
- Deep Learning (DL)
- Naive Bayes (NB)

H₂O.ai

```
> print(gbm_gridperf)
H2O Grid Details
================

Grid ID: gbm_grid2
Used hyper parameters:
  - sample_rate
  - max_depth
  - learn_rate
  - col_sample_rate
Number of models: 72
Number of failed models: 0

Hyper-Parameter Search Summary: ordered by decreasing auc
  sample_rate max_depth learn_rate col_sample_rate        model_ids              auc
1           1         3       0.19               1 gbm_grid2_model_38 0.685166598389755
2         0.9         3       0.15               1 gbm_grid2_model_53 0.684956999713052
3         0.8         5       0.06               1 gbm_grid2_model_22 0.684843506375254
4         0.6         4       0.07               1  gbm_grid2_model_4 0.684327718715252
5        0.95         4       0.13               1 gbm_grid2_model_48 0.684042497773235
```

The second training module demonstrates how to find the best set of model parameters for each model using Grid Search.

**H2O**.ai

- H2O Online Training (free):  http://learn.h2o.ai

- H2O Slidedecks:  http://www.slideshare.net/0xdata

- H2O Video Presentations:  https://www.youtube.com/user/0xdata

- H2O Community Events & Meetups:  http://h2o.ai/events

- Machine Learning & Data Science courses:  http://coursebuffet.com