

Overview of Machine Learning and H2O.ai



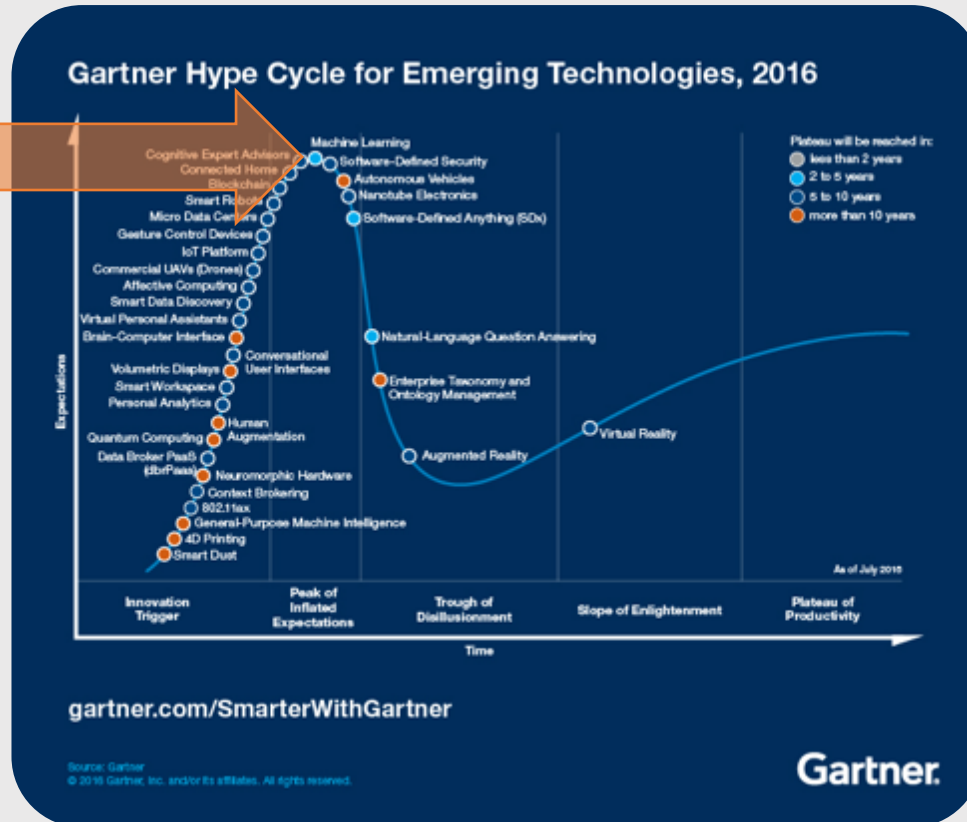
Machine Learning Overview

What is machine learning?

A field of study that gives computers the ability to learn without being explicitly programmed.

-- Arthur Samuel, 1959





Why now?

- Data, computers, and algorithms are commodities
- Unstructured data
- Increasing competition in business

Estimating a model for inference

What happened? Why?

Assumptions, parsimony,
interpretation

Linear models, statistics

Models tend to be static

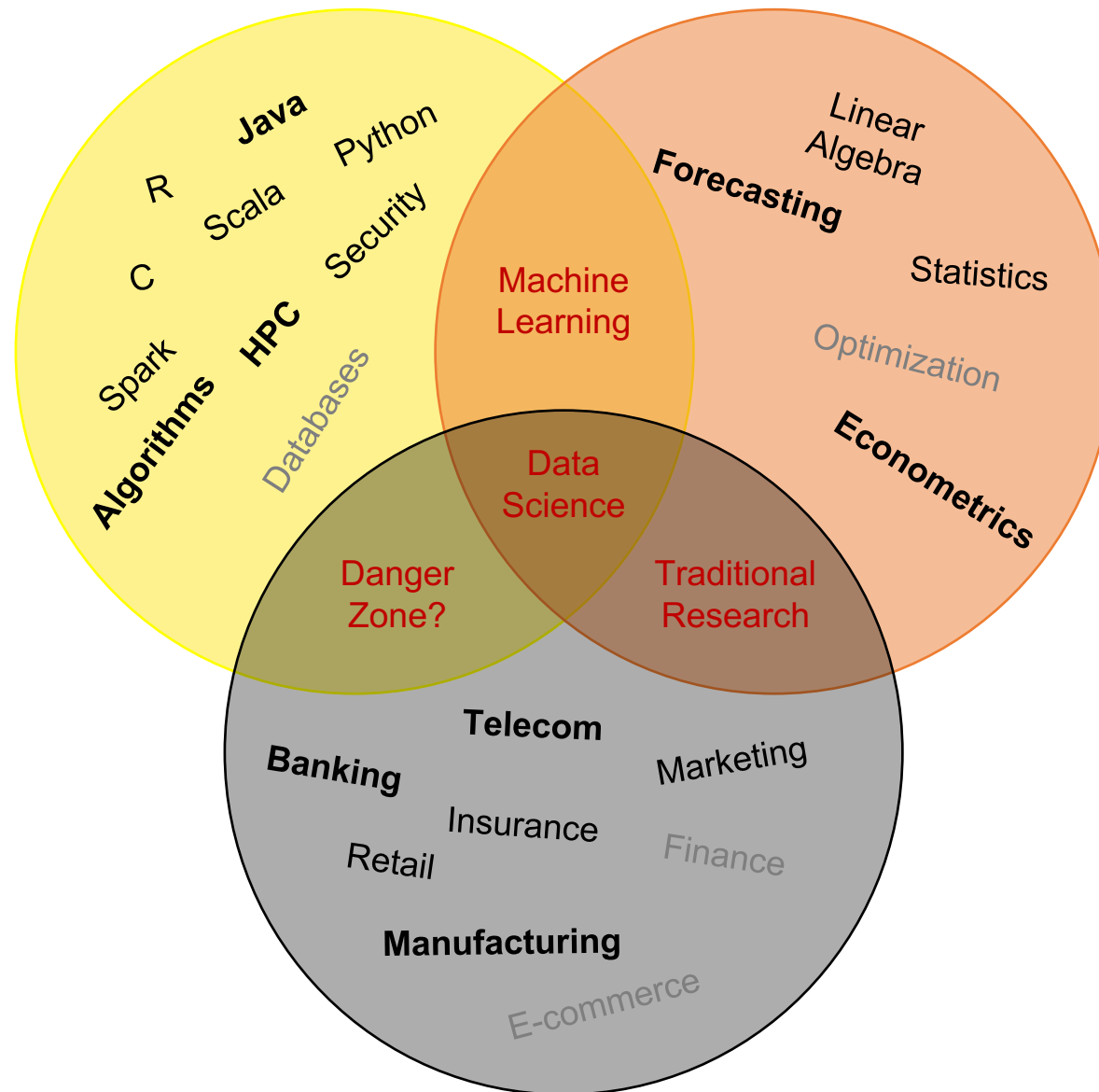
Training a model for prediction

What will happen?

Predictive accuracy,
production deployment

Machine learning

Many models can evolve
elegantly



1. There is no perfect language.



If someone claims to have the perfect programming language, he is either a fool or a salesman or both.

-- Bjarne Stroustrup

2. There is no perfect algorithm.



Algorithms that search for an extremum of a cost function perform exactly the same when averaged over all possible cost functions.

-- D.H. Wolpert

3. Doing things right is always hard.



Developing and deploying ML systems is relatively fast and cheap, but maintaining them over time is difficult and expensive.

-- Google,
*Hidden Technical Debt in
Machine Learning Systems*

H₂O.ai Overview

Company Overview

Founded	2011 Venture-backed, debuted in 2012
Products	<ul style="list-style-type: none">• H2O: In-Memory AI Prediction Engine• Sparkling Water: Spark Integration• Steam: Deployment engine• Deep Water: Deep Learning
Mission	Operationalize Data Science, and provide a platform for users to build beautiful data products
Team	70 employees <ul style="list-style-type: none">• Distributed Systems Engineers doing Machine Learning• World-class visualization designers
Headquarters	Mountain View, CA



H₂O.ai

H2O.ai Offers AI Open Source Platform

Product Suite to Operationalize Data Science

100% Open Source



In-Memory, Distributed
Machine Learning
Algorithms with Speed
and Accuracy

The logo for Deep Water, featuring the text "Deep Water" in black on a light yellow background.

State-of-the-art
Deep Learning on GPUs
with TensorFlow, MXNet
or Caffe with the ease of
use of H2O



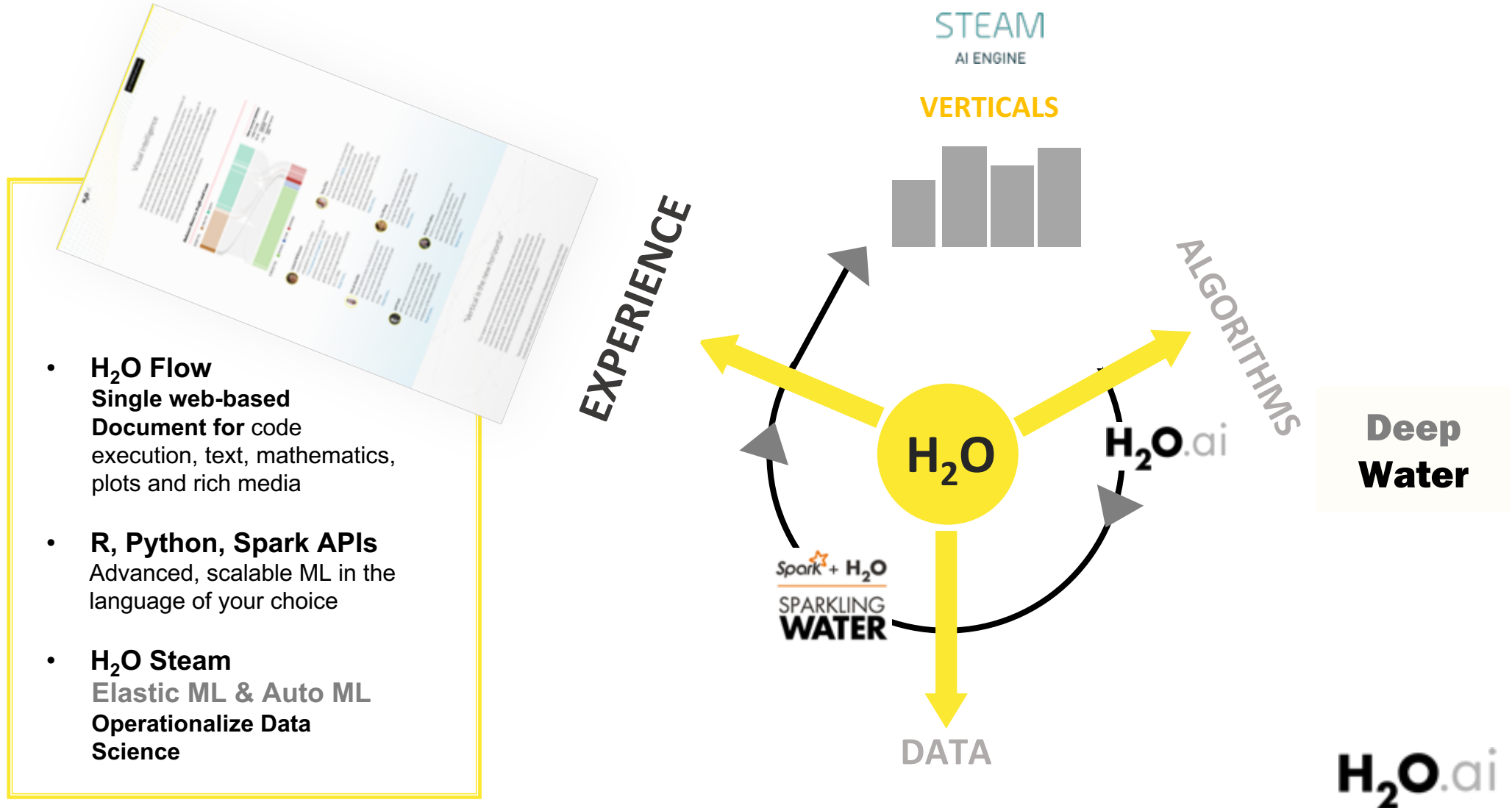
H2O Integration with
Spark. Best Machine
Learning on Spark.

The logo for Steam, featuring the word "Steam" in orange.

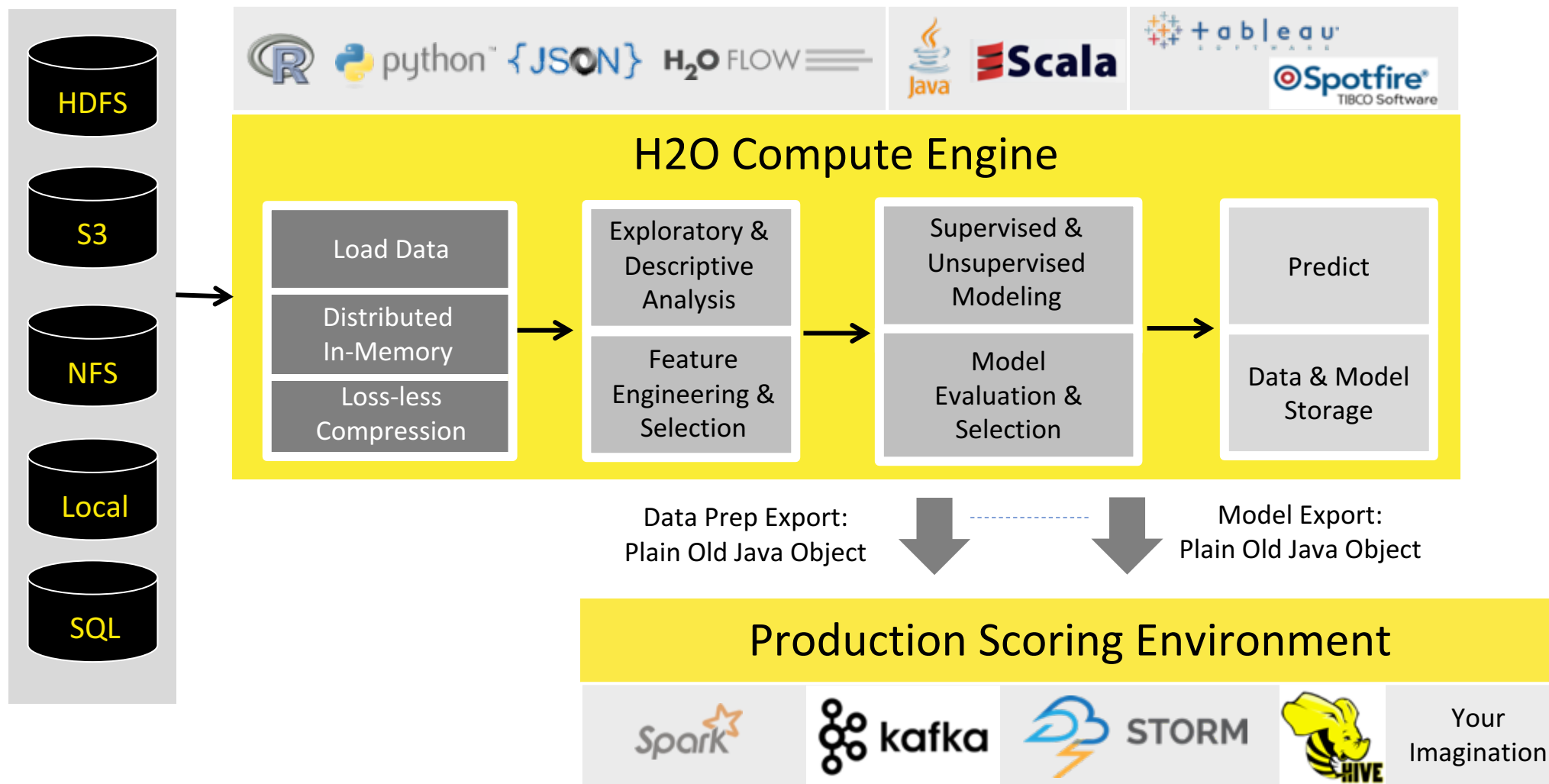
Operationalize and
Streamline Model
Building, Training and
Deployment Automatically
and Elastically

H₂O.ai Now Focused On Experience

Beyond Algorithms and Data



High Level Architecture



Intro to Machine Learning

Algos

Algorithms on H2O

Supervised Learning

Statistical Analysis

- **Penalized Linear Models:** Super-fast, super-scalable, and interpretable
- **Naïve Bayes:** Straightforward linear classifier

Decision Tree Ensembles

- **Distributed Random Forest:** Easy-to-use tree-bagging ensembles
- **Gradient Boosting Machine:** Highly tunable tree-boosting ensembles

Stacking

- **Stacked Ensemble:** Combine multiple types of models for better predictions

Neural Networks

Multilayer Perceptron

- **Deep neural networks:** Multi-layer feed-forward neural networks for standard data mining tasks

Deep Learning

- **Convolutional neural networks:** Sophisticated architectures for pattern recognition in images, sound, and text

Unsupervised Learning

Clustering

- **K-means:** Partitions observations into similar groups; automatically detects number of groups

Dimensionality Reduction

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

Aggregator

- **Aggregator:** Efficient, advanced sampling that creates smaller data sets from larger data sets

Anomaly Detection

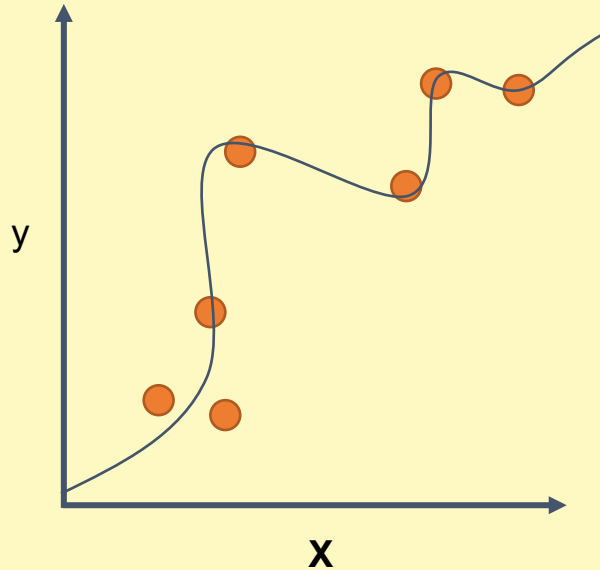
- **Autoencoders:** Find outliers using a nonlinear dimensionality reduction technique

Term Embeddings

- **Word2vec:** Generate context-sensitive numerical representations of a large text corpus

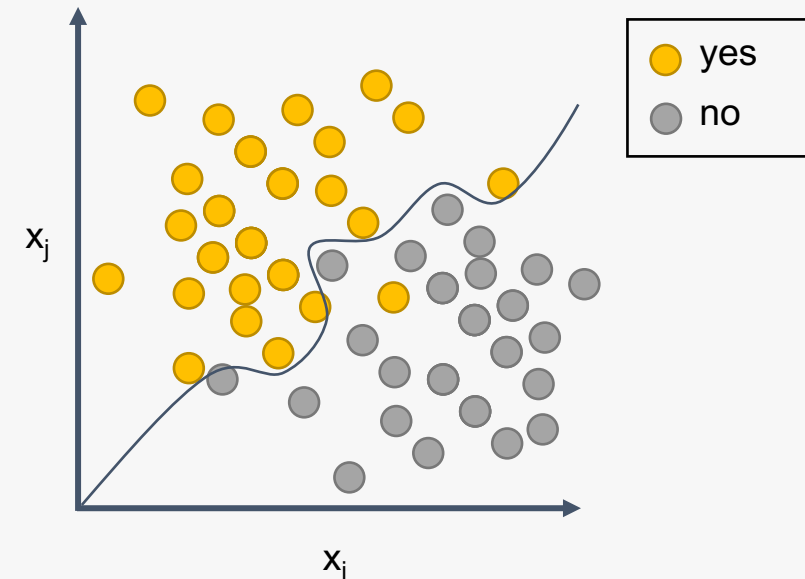
Supervised Learning

Regression:
How much will a customer spend?



H₂O algos:
Penalized Linear Models
Random Forest
Gradient Boosting
Neural Networks
Stacked Ensembles

Classification:
Will a customer make a purchase? Yes or No

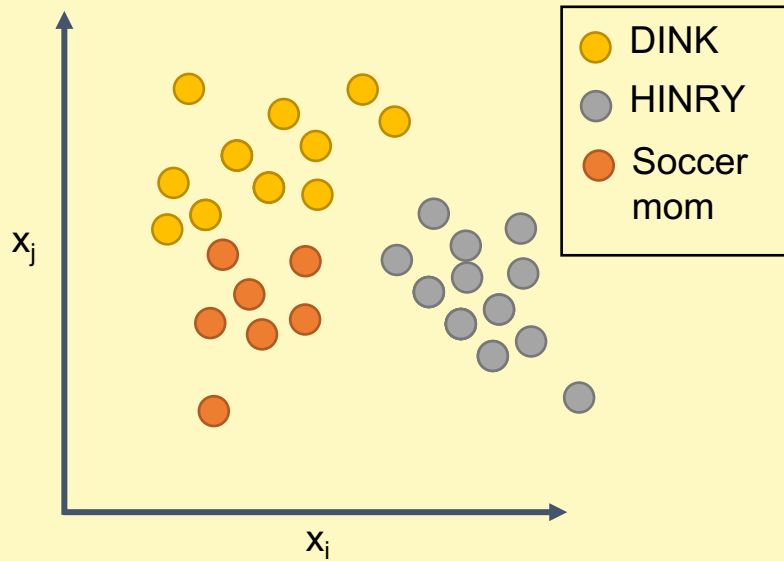


H₂O algos:
Penalized Linear Models
Naïve Bayes
Random Forest
Gradient Boosting
Neural Networks
Stacked Ensembles

Unsupervised Learning

Clustering:

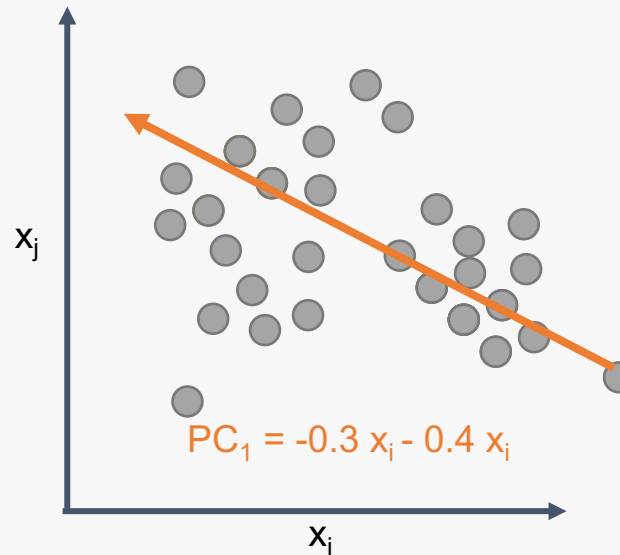
Grouping rows – e.g. creating groups of similar customers



H₂O algos:
k – means

Feature extraction:

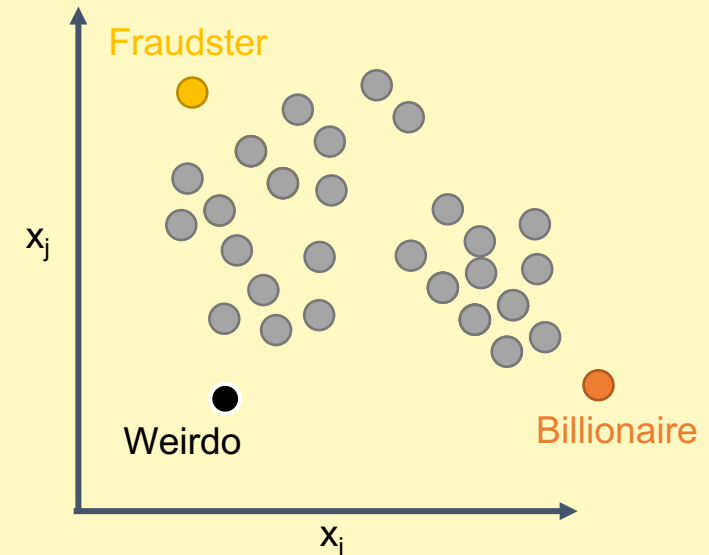
Grouping columns – Create a small number of new representative dimensions



H₂O algos:
Principal components
Generalized low rank models
Autoencoders
Word2Vec

Anomaly detection:

Detecting outlying rows - Finding high-value, fraudulent, or weird customers



H₂O algos:
Principal components
Generalized low rank models
Autoencoders

Penalized Linear Models

- Regression
- Classification

- Creates interpretable models with super-fast training time
- Nonlinear and interaction terms to be specified manually
- Can extrapolate beyond training data domain
- Select the correct target distribution
- Few hyperparameters to tune

- NAs
- Outliers/influential points
- Strongly correlated inputs
- Rare categorical levels in new data

Naïve Bayes

- Classification

- Nonlinear and interaction terms should be specified by users

- Linear independence assumption
- Often less accurate than more sophisticated classifiers
- Rare categorical levels in new data

Random Forest

- Regression
- Classification

- Builds accurate models without overfitting
- Few hyperparameters to tune
- Requires less data prep
- Great for implicitly modeling interactions

- Difficulty extrapolating beyond training data domain
- Can be difficult to interpret
- Rare categorical levels in new data

Gradient Boosting Machines

- Regression
- Classification

- Builds accurate models without overfitting (often more accurate than random forest)
- Requires less data prep
- Great for implicitly modeling interactions

- Many hyperparameters
- Difficulty extrapolating beyond training data domain
- Can be difficult to interpret
- Rare categorical levels in new data

Neural Networks (Deep learning & MLP)

- Regression
- Classification

- Great for modeling interactions in fully connected topologies
- Can extrapolate beyond training data domain
- Deep learning architectures best-suited for pattern recognition in images, videos, and sound

- NAs
- Overfitting
- Outliers/influential points
- Long training times
- Difficult to interpret
- Many hyperparameters
- Strongly correlated inputs
- Rare categorical levels in new data

	Usage	Recommendations	Problems
<i>k</i> - means	<ul style="list-style-type: none"> Clustering 	<ul style="list-style-type: none"> Great for creating Gaussian, non-overlapping, roughly equally sized clusters The number of clusters can be unknown 	<ul style="list-style-type: none"> NAs Outliers/influential points Strongly correlated inputs Cluster labels sensitive to initialization Curse of dimensionality
Principal Components Analysis	<ul style="list-style-type: none"> Feature extraction Dimension reduction Anomaly detection 	<ul style="list-style-type: none"> Great for extracting a number $\leq N$ of linear, orthogonal features from i.i.d. numeric data Great for plotting extracted features in a reduced-dimensional space to analyze data structure, e.g. clusters, hierarchy, sparsity, outliers 	<ul style="list-style-type: none"> NAs Outliers/influential points Categorical inputs
Generalized Low Rank Models	<ul style="list-style-type: none"> Feature extraction Dimension reduction Anomaly detection Matrix completion 	<ul style="list-style-type: none"> Great for extracting linear features from mixed data Great for plotting extracted features in a reduced-dimensional space to analyze data structure, e.g. clusters, hierarchy, sparsity, outliers Great for imputing NAs 	<ul style="list-style-type: none"> Outliers/influential points
Autoencoders (Neural Networks)	<ul style="list-style-type: none"> Feature extraction Dimension reduction Anomaly detection 	<ul style="list-style-type: none"> Great for extracting a number of nonlinear features from mixed data Great for plotting extracted features in a reduced dimensional space to analyze structure, e.g. clusters, hierarchy, sparsity, outliers 	<ul style="list-style-type: none"> NAs Overtraining Outliers/influential points Long training times Many hyperparameters Strongly correlated inputs Rare categorical levels in new data
Word2Vec	<ul style="list-style-type: none"> Highly representative feature extraction from text 	<ul style="list-style-type: none"> Great for extracting highly representative, context sensitive term embeddings (e.g. numerical vectors) from text Great for text preprocessing prior to further supervised or unsupervised analysis 	<ul style="list-style-type: none"> Many Hyperparameters Long training times Overtraining Specifying term weightings prior to training