# H2O.ai Algorithms

# 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 Learning**

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

**Term Embeddings**
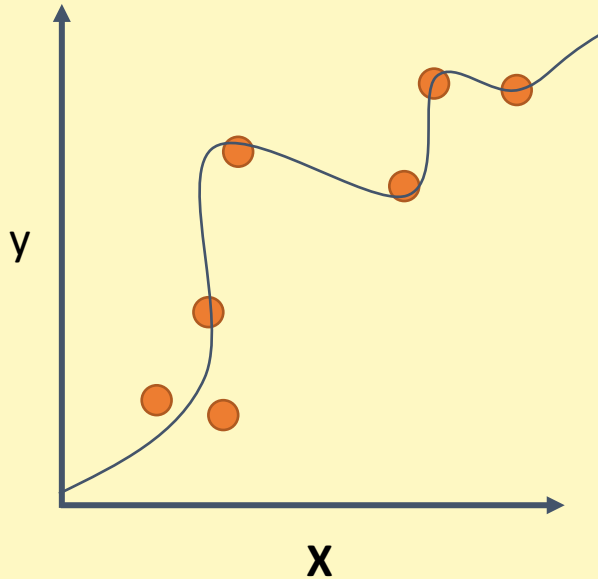
- **Autoencoders**: Find outliers using a nonlinear dimensionality reduction technique
- **Word2vec:** Generate context-sensitive numerical representations of a large text corpus

# Supervised Learning
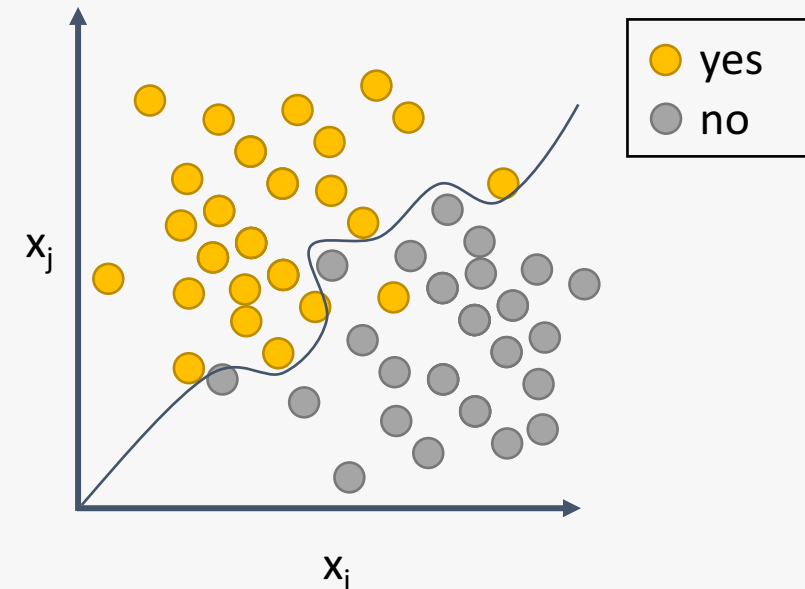
**Regression:**

**How much will a customers spend?**



**H₂O algos:**

**Penalized Linear Models**

**Random Forest**

**Gradient Boosting**

**Neural Networks**

**Stacked Ensembles**

**Classification:**

**Will a customer make a purchase? Yes or No**
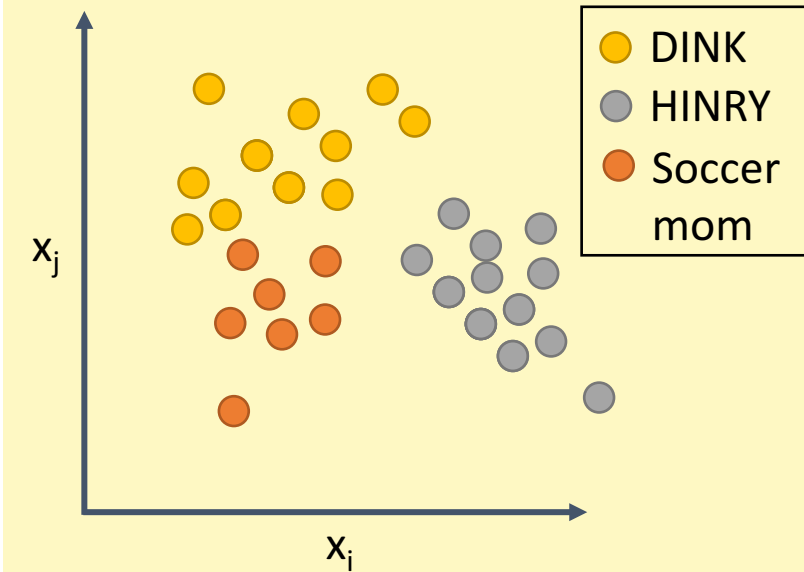


**H₂O algos:**

**Penalized Linear Models**

**Naïve Bayes**

**Random Forest**

**Gradient Boosting**

**Neural Networks**

**Stacked Ensembles**

H₂O.ai

# Unsupervised Learning

## Clustering:

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



- $\circ$ DINK
- $\circ$ HINRY
- $\circ$ Soccer mom

$x_j$

$x_i$

**H₂O algos:**

k – means

## Feature extraction:

Grouping columns – Create a small number of new representative dimensions



$x_j$

$x_i$

$PC_1 = -0.3\ x_i - 0.4\ x_i$

**H₂O algos:**

**Principal components**
**Generalized low rank models**
**Autoencoders**
**Word2Vec**

## Anomaly detection:

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



Fraudster

$x_j$

Weirdo

Billionaire

$x_i$

**H₂O algos:**

**Principal components**
**Generalized low rank models**
**Autoencoders**

H₂O.ai

# H₂O.ai

| | Usage | Recommendations | Problems |
|---|---|---|---|
| **Penalized Linear Models** | • Regression<br>• Classification | • Creates interpretable models with super-fast training time<br>• Nonlinear and interaction terms to be specified manually<br>• Can extrapolate beyond training data domain<br>• Select the correct target distribution<br>• Few hyperparameters to tune | • NAs<br>• Outliers/influential points<br>• Strongly correlated inputs<br>• Rare categorical levels in new data |
| **Naïve Bayes** | • Classification | • Nonlinear and interaction terms should be specified by users | • Linear independence assumption<br>• Often less accurate than more sophisticated classifiers<br>• Rare categorical levels in new data |
| **Random Forest** | • Regression<br>• Classification | • Builds accurate models without overfitting<br>• Few hyperparameters to tune<br>• Requires less data prep<br>• Great for implicitly modeling interactions | • Difficulty extrapolating beyond training data domain<br>• Can be difficult to interpret<br>• Rare categorical levels in new data |
| **Gradient Boosting Machines** | • Regression<br>• Classification | • Builds accurate models without overfitting (often more accurate than random forest)<br>• Requires less data prep<br>• Great for implicitly modeling interactions | • Many hyperparameters<br>• Difficulty extrapolating beyond training data domain<br>• Can be difficult to interpret<br>• Rare categorical levels in new data |
| **Neural Networks** (Deep learning & MLP) | • Regression<br>• Classification | • Great for modeling interactions in fully connected topologies<br>• Can extrapolate beyond training data domain<br>• Deep learning architectures best-suited for pattern recognition in images, videos, and sound | • NAs<br>• Overfitting<br>• Outliers/influential points<br>• Long training times<br>• Difficult to interpret<br>• Many hyperparameters<br>• Strongly correlated inputs<br>• Rare categorical levels in new data |

| | Usage | Recommendations | Problems | |
|---|---|---|---|---|
| **_k_ - means** | • Clustering | • Great for creating Gaussian, non-overlapping, roughly equally sized clusters<br>• The number of clusters can be unknown | • NAs<br>• Outliers/influential points<br>• Strongly correlated inputs<br>• Cluster labels sensitive to initialization<br>• Curse of dimensionality | |
| **Principal Components Analysis** | • Feature extraction<br>• Dimension reduction<br>• Anomaly detection | • Great for extracting a number <= $N$ of linear, orthogonal features from i.i.d. numeric data<br>• Great for plotting extracted features in a reduced-dimensional space to analyze data structure, e.g. clusters, hierarchy, sparsity, outliers | • NAs<br>• Outliers/influential points<br>• Categorical inputs | |
| **Generalized Low Rank Models** | • Feature extraction<br>• Dimension reduction<br>• Anomaly detection<br>• Matrix completion<br>• Recommender Systems | • Great for extracting linear features from mixed data<br>• Great for plotting extracted features in a reduced-dimensional space to analyze data structure, e.g. clusters, hierarchy, sparsity, outliers<br>• Great for imputing NAs<br>• Great for creating recommendations | • Outliers/influential points | |
| **Autoencoders (Neural Networks)** | • Feature extraction<br>• Dimension reduction<br>• Anomaly detection | • Great for extracting a number of nonlinear features from mixed data<br>• Great for plotting extracted features in a reduced dimensional space to analyze structure, e.g. clusters, hierarchy, sparsity, outliers | • NAs<br>• Overtraining<br>• Outliers/influential points<br>• Long training times | • Many hyperparameters<br>• Strongly correlated inputs<br>• Rare categorical levels in new data |
| **Word2Vec** | • Highly representative feature extraction from text | • Great for extracting highly representative, context sensitive term embeddings (e.g. numerical vectors) from text<br>• Great for text preprocessing prior to further supervised or unsupervised analysis | • Many Hyperparameters<br>• Overtraining<br>• Specifying term weightings prior to training | • Long training times |