

Machine Learning with H₂O

H₂O.ai

H2O Algorithms

Un-supervised

H2O Algorithms

Supervised Learning

Statistical
Analysis

- **Generalized Linear Models:** Binomial, Gaussian, Gamma, Poisson, and Tweedie
- **Naive Bayes:** Binary Text Classification

Ensembles

- **Distributed Random Forest:** Classification or Regression Models
- **Gradient Boosting Machine:** Ensembles of shallow decision trees with increasing refined approximations

Deep Neural Networks

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

H2O Algorithms

Unsupervised Learning

Clustering

- **K-means:** Partition observations into k clusters of the same spatial size. Categorical features are one hot encoded.
- **Archetypes [GLRM]:** Partition observations into k archetypes.

Dimensionality Reduction

- **Principal Component Analysis:** Linearly transforms correlated variables to independent components
- **Generalized Low Rank Model:** Approximates data set as a product of two low dimensional factors. Extends PCA to handle sparse data, categorical data, and adds regularization.

Anomaly Detection

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

Unsupervised Learning Models

- No target!
- Ways to characterize or simplify data
- Uses
 - Clustering
 - Variable reduction
 - Compression
 - Anomaly detection

Unsupervised Learning:

K-MEANS CLUSTERING

Clustering

- We are going to try and “bundle” customers/companies/transactions based on “likeness” of certain features.
- These “clusters” will help us simplify and potentially “segment” observations.
 - Segment customers
 - Categorize events

K-Means in H2O

CS buildModel "kmeans"
12ms

Build a Model

Select an algorithm: K-means

PARAMETERS

model_id	kmeans-bdc9ae94-5337-460d-9692-7460bdbb429a	Destination id for this model; auto-generated if not specified.
training_frame	(Choose...)	Id of the training data frame (Not required, to allow initial validation of model parameters).
validation_frame	(Choose...)	Id of the validation data frame.
nfolds	0	Number of folds for N-fold cross-validation (0 to disable or >= 2).
ignored_columns	Search...	

All None ◀ Previous 100 ▶ Next 100

Only show columns with more than 0 % missing values.

ADVANCED

ignore_const_cols	<input checked="" type="checkbox"/>	Ignore constant columns.
k	1	The max. number of clusters. If estimate_k is disabled, the model will find k centroids, otherwise it will find up to k centroids.
estimate_k	<input type="checkbox"/>	Whether to estimate the number of clusters (<=k) iteratively and deterministically.
max_iterations	10	Maximum training iterations (if estimate_k is enabled, then this is for each inner Lloyds iteration)
standardize	<input checked="" type="checkbox"/>	Standardize columns before computing distances
init	Furthest	Initialization mode

EXPERT

fold_assignment	AUTO	Cross-validation fold assignment scheme, if fold_column is not specified. The 'Stratified' option will stratify the folds based on the response variable, for classification problems.
fold_column	(Choose...)	Column with cross-validation fold index assignment per observation.
score_each_iteration	<input type="checkbox"/>	Whether to score during each iteration of model training.
seed	-1	RNG Seed
max_runtime_secs	0	Maximum allowed runtime in seconds for model training. Use 0 to disable.
categorical_encoding	AUTO	Encoding scheme for categorical features

keep_cross_validation_predictions
 keep_cross_validation_fold_assignment
user_points (Choose...) This option allows you to specify a dataframe, where each row represents an initial cluster center. The user-specified points must have the same number of columns as the training observations. The number of rows must equal the number of clusters

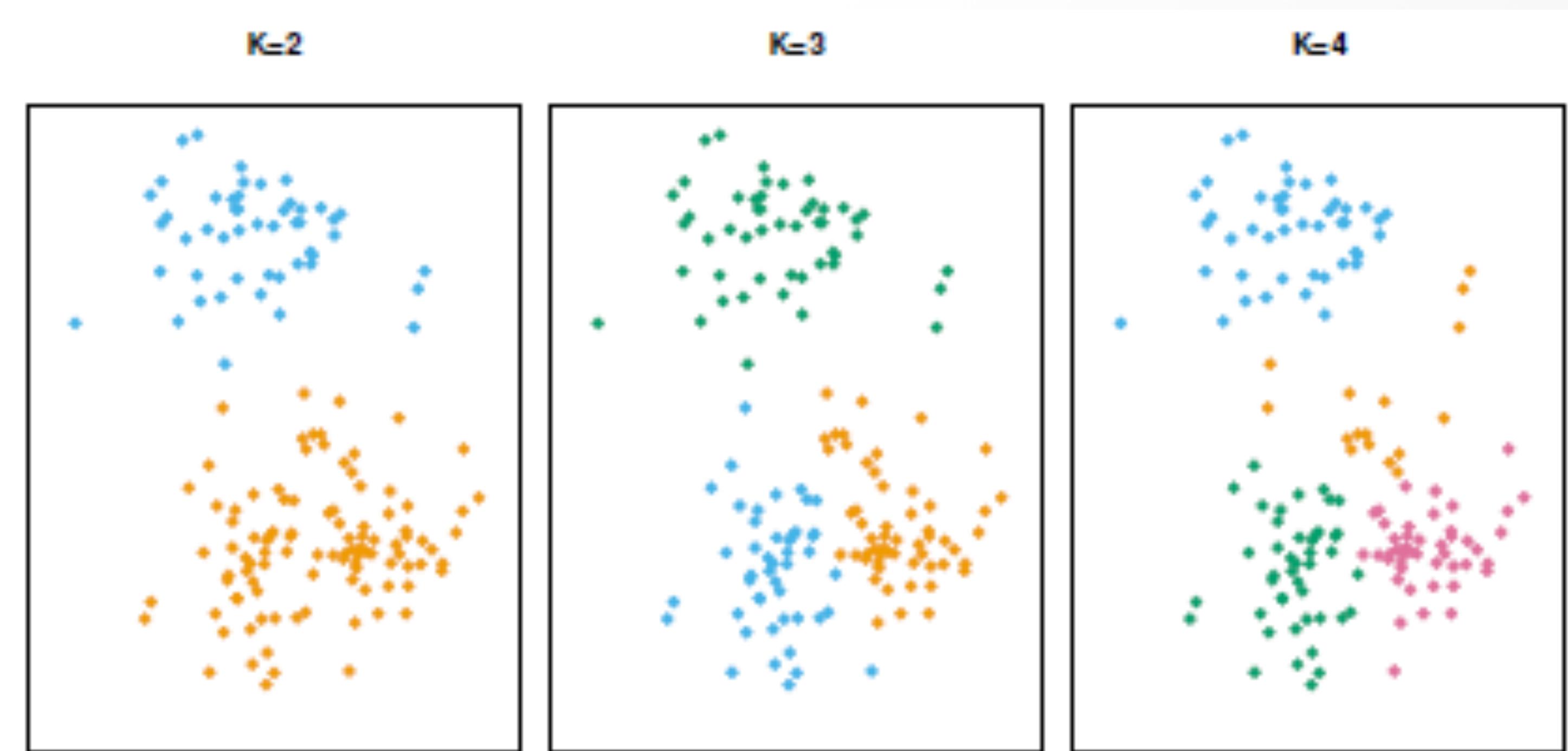
Build Model

Clustering

Objective: Minimize the within cluster variation for a given k

How to choose k?

- Choose several and evaluate performance
- Often business rules determine feasible k



What is Principal Components Analysis

- Another method of data simplification
- Can I reduce the width (dimensionality) of the dataset by creating one or more new PCA variables that do a “good job” of representing the original signal in the data?
- Uses:
 - Variable reduction=> Fewer predictors
 - Feature extraction=> Anomaly detection

Principal Component Analysis in H2O

CS buildModel "pca" 9ms

Build a Model

Select an algorithm: Principal Components Analysis

PARAMETERS

model_id	pca-5385d795-5fb4-466f-b884-c703a6a89424	Destination id for this model; auto-generated if not specified.
training_frame	(Choose...)	Id of the training data frame (Not required, to allow initial validation of model parameters).
validation_frame	(Choose...)	Id of the validation data frame.
ignored_columns	Search...	

All None

Only show columns with more than 0 % missing values.

ignore_const_cols	<input checked="" type="checkbox"/>	Ignore constant columns.
transform	NONE	Transformation of training data
pca_method	GramSVD	Method for computing PCA (Caution: GLRM is currently experimental and unstable)
k*	1	Rank of matrix approximation
max_iterations	1000	Maximum training iterations
use_all_factor_levels	<input type="checkbox"/>	Whether first factor level is included in each categorical expansion
compute_metrics	<input checked="" type="checkbox"/>	Whether to compute metrics on the training data
impute_missing	<input type="checkbox"/>	Whether to impute missing entries with the column mean
seed	-1	RNG seed for initialization

ADVANCED

score_each_iteration	<input type="checkbox"/>	Whether to score during each iteration of model training.
max_runtime_secs	0	Maximum allowed runtime in seconds for model training. Use 0 to disable.

GRID?

Build Model

Generalized Low Rank Models

- **Given:** Data table A w/ m rows and n columns
- **Find:** Compressed representation as numeric table X and Y where # cols in X = # rows in Y = small user specified K $\ll \max(m,n)$
- # cols in Y is d = (total dimension of embedded features in A) $\geq n$

$$m \left\{ \overbrace{\begin{bmatrix} A \end{bmatrix}}^n \right\} \approx m \left\{ \overbrace{\begin{bmatrix} X \end{bmatrix}}^k \left[\overbrace{\begin{bmatrix} Y \end{bmatrix}}^n \right] \right\} k$$

- Row of Y = archetypal fea
- Row of X = row of A in reduced feature space
- Can approximately reconstruct A from product XY

Generalized Low Rank Modeling in H2O

CS buildModel "glrm"

7ms

Build a Model

Select an algorithm: Generalized Low Rank Modeling

PARAMETERS

model_id glrm-767dc939-83ed-4ce7-af3f-37c084f9a76b
training_frame (Choose...)
validation_frame (Choose...)
ignored_columns Search...

Destination id for this model; auto-generated if not specified.
Id of the training data frame (Not required, to allow initial validation of model parameters).
Id of the validation data frame.

GRID?

Only show columns with more than 0 % missing values.

All None

[← Previous 100](#) [→ Next 100](#)

ignore_const_cols
loading_name
transform NONE
k* 1
loss Quadratic
loss_by_col_idx
multi_loss Categorical
period 1
regularization_x None
regularization_y None
gamma_x 0
gamma_y 0
max_iterations 1000
max_updates 2000
init_step_size 1
min_step_size 0.0001
seed -1
init PlusPlus
svd_method Randomized
user_y (Choose...)
user_x (Choose...)
expand_user_y
impute_original
recover_svd

Ignore constant columns.
Frame key to save resulting X
Transformation of training data
Rank of matrix approximation
Numeric loss function
Loss function by column index (override)
Categorical loss function
Length of period (only used with periodic loss function)
Regularization function for X matrix
Regularization function for Y matrix
Regularization weight on X matrix
Regularization weight on Y matrix
Maximum number of iterations
Maximum number of updates, defaults to 2*max_iterations
Initial step size
Minimum step size
RNG seed for initialization
Initialization mode
Method for computing SVD during initialization (Caution: Randomized is currently experimental and unstable)
User-specified initial Y
User-specified initial X
Expand categorical columns in user-specified initial Y
Reconstruct original training data by reversing transform
Recover singular values and eigenvectors of XY

ADVANCED

score_each_iteration
max_runtime_secs 0

Whether to score during each iteration of model training.
Maximum allowed runtime in seconds for model training. Use 0 to disable.

GRID?

[Build Model](#)

Why use PCA?

- Reduce storage space, e.g. 10 GB compressed to 100 MB
- Increase prediction speed, e.g. 10x speed-up with no accuracy loss
- Identify and visualize important features
- Impute missing data entries
- Supports different loss functions:
 - Quadratic (PCA), Absolute, Huber, Poisson, Hinge, Logistic, Periodic
- Supports regularization:
 - Quadratic, L2, L1, NonNegative, OneSparse, UnitOneSparse, Simplex

More Information

- Anqi Fu GLRM
 - <https://www.youtube.com/watch?v=gEZtZRANeLc>
- Madeline Udell
 - <https://www.youtube.com/watch?v=zwvzGuS82MA>
- <https://web.stanford.edu/~boyd/papers/pdf/glrm.pdf>



avkash@h2o.ai