

A concise overview of classification and clustering methods

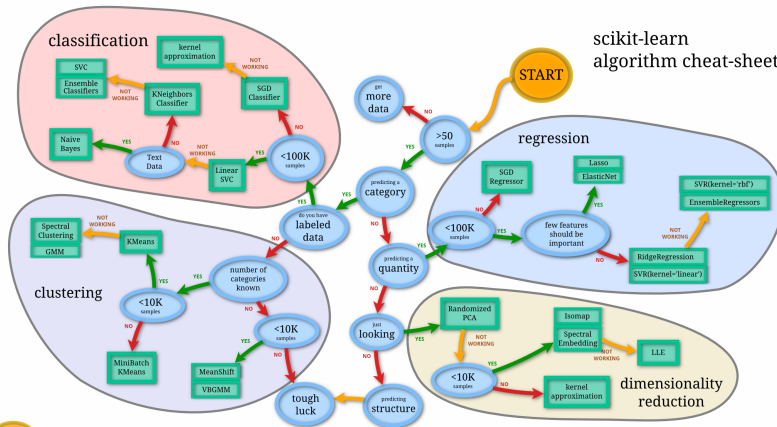
Maxime Sangnier

GDR IAMAT - May 31, 2022

Sorbonne Université, CNRS, LPSM, Paris, France

A scikit-learn map

scikit-learn algorithm cheat-sheet



Source: <https://scikit-learn.org/>

Mathematical setting

Vectors and matrices

- Each experiment is defined by initial conditions $\rightarrow x$ and a result $\rightarrow y$.
- Each observation is defined by an individual $\rightarrow x$ and a feature $\rightarrow y$.
- For computational purposes, $x \in \mathbb{R}^d$ and $y \in \mathbb{R}$.
- Many repetitions of the experiment/observation provide $(x_1, y_1), \dots, (x_n, y_n)$.

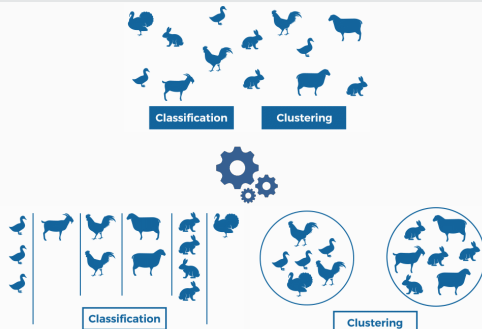
- Data matrix $\mathbf{X} = \begin{pmatrix} x_1 \text{ in row} \\ \vdots \\ x_n \text{ in row} \end{pmatrix}$ (size $n \times d$) and vector of outputs $\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$.

Numerical formulation

Goal

1. **Classification/Regression** Given x_{new} , “predict” *i.e.* approximate y_{new} : $y_{new} \approx f(x_{new})$.
2. **Clustering** When y_i 's are *not* observed, gather x_i 's that are *similar*: $\{X_1, X_3, X_4, \dots\} - \{X_2, X_7, X_{11}, \dots\} - \{X_5, X_{13}, X_{22}, \dots\}$.

Classification and clustering are roughly similar except that no information concerning the group is available for clustering.



Randomness

- (x_i, y_i) is a realization of a random pair (X_i, Y_i) .
- Why randomness? Because y_i cannot be computed only based on x_i :
 - y_i is noisy and does not reflect actually our desire.
 - x_i is incomplete and a perfect description is out of reach.

Statistics

- All (X_i, Y_i) 's are independent and have the same distribution.
- We are interested in the expected behavior: overall true now and in the future.

1. Classification/Regression

$$\mathbb{P}(f(X_{new}) \text{ is close to } Y_{new}) = \mathbb{E}[\text{similarity}(f(X_{new}), Y_{new})].$$

2. Clustering

$$\begin{aligned} & \mathbb{P}(X_{new} \text{ is similar to } X_{other} \text{ when } X_{new} \text{ and } X_{other} \text{ fall in the same group}) \\ &= \mathbb{E}[\text{similarity}(X_{new}, X_{other}) \text{ when } X_{new} \text{ and } X_{other} \text{ fall in the same group}]. \end{aligned}$$

- The distribution of (X_i, Y_i) 's is unknown but. . .
- By the law of large numbers,

$$\frac{1}{n} \sum_{i=1}^n \text{similarity}(f(X_i), Y_i) \xrightarrow[n \rightarrow \infty]{a.s.} \mathbb{E}[\text{similarity}(f(X_{new}), Y_{new})].$$

Supervised learning

Two methodologies

- We observe inputs and outputs: $(X_1, Y_1), \dots, (X_n, Y_n)$.
- Let ℓ be a loss (*i.e.* a dissimilarity) function.
- **Goal** Find f such that $\mathbb{E}[\ell(f(X_{new}), Y_{new})]$ is minimal.
 - **Classification** $\ell(f(X_{new}), Y_{new}) = \begin{cases} 1 & \text{if } f(X_{new}) \neq Y_{new} \\ 0 & \text{if } f(X_{new}) = Y_{new} \end{cases}$.
 - **(L^2) Regression** $\ell(f(X_{new}), Y_{new}) = (Y_{new} - f(X_{new}))^2$.

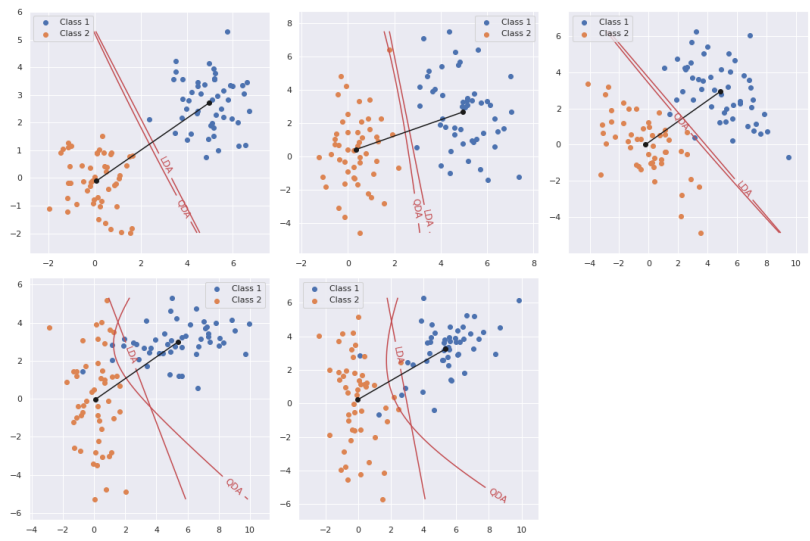
Plug-in estimator

- Bayes minimizer: $f^*(x) = L(\mathbb{E}[Y_{new} | X_{new} = x])$, where L is known.
- Estimator: $\hat{f}(x) = L(\hat{\mathbb{E}}[Y_{new} | X_{new} = x])$.
- Needs statistics.
- Often:
 - We have to model the data distribution.
 - We resort to simple and non-robust estimators.
 - We are quite limited (the Bayes estimator is not always explicit).
- **Examples** Parametric models (LDA, QDA, Logistic Regression, Linear Regression), kernel methods (k -Nearest Neighbors, Trees, Random Forests).

Empirical Risk Minimization

- By the law of large numbers: $\frac{1}{n} \sum_{i=1}^n \ell(f(X_i), Y_i) \approx \mathbb{E}[\ell(f(X_{new}), Y_{new})]$.
- Estimator: \hat{f} such that $\frac{1}{n} \sum_{i=1}^n \ell(\hat{f}(X_i), Y_i)$ is minimal.
- Needs optimization tools.
- Sometimes:
 - We have to design efficient optimization algorithms (convex optimization).
 - We have to settle for a good algorithm applied in an inadequate setting (non-convex optimization).
- **Examples** Boosting, Support Vector Machines, Neural Networks.

Discriminant analysis



The method

- For binary or multiclass classification only.
- Assumption on data:

$$\begin{cases} X_{new} \mid Y_{new} = \bullet \sim \mathcal{N}(\mu, \Sigma), & \mathbb{P}(Y_{new} = \bullet) = \pi, \\ X_{new} \mid Y_{new} = \bullet \sim \mathcal{N}(\mu, \Sigma), & \mathbb{P}(Y_{new} = \bullet) = \pi, \end{cases}$$

for unknown $\mu, \mu, \Sigma, \Sigma, \pi, \pi$.

- The Bayes classifier draws a parabolic frontier:

$$f^*(x) = \begin{cases} \bullet & \text{if } \frac{1}{2}x^\top (\Sigma^{-1} - \Sigma^{-1})x + (\Sigma^{-1}\mu - \Sigma^{-1}\mu)^\top x + \dots \\ & + \log\left(\frac{\pi}{\pi}\right) \geq 0 \\ \bullet & \text{otherwise.} \end{cases}$$

- Plug-in estimation via Maximum Likelihood: closed-form expressions for $\hat{\mu}, \hat{\mu}, \hat{\Sigma}, \hat{\Sigma}, \hat{\pi}, \hat{\pi}$.

Discriminant analysis

What if $\Sigma = \Sigma$?

- The Bayes classifier draws a hyperplane frontier:

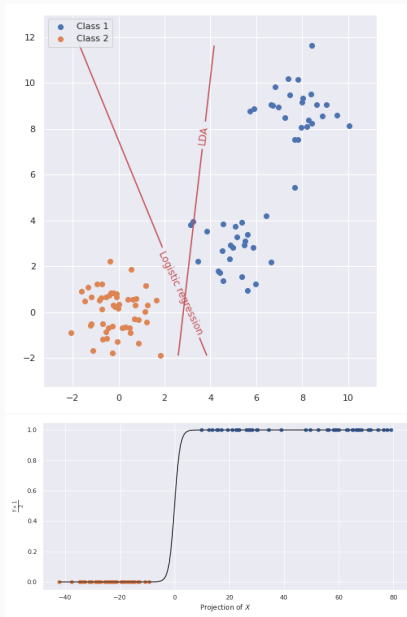
$$f^*(x) = \begin{cases} \bullet & \text{if } (\mu - \mu)^\top \Sigma^{-1} x + \dots + \log\left(\frac{\pi}{\pi}\right) \geq 0 \\ \bullet & \text{otherwise.} \end{cases}$$

- $\frac{\pi}{\pi}$ translates the hyperplane frontier.
- If $\pi = \pi$, the Bayes classifier is a minimum-Mahalanobis-distance-to-center classifier.

Pros and cons

- It is computationally tractable (OK in high dimension and with large amount of data).
- It is very restrictive (categorical data...).
- It is not robust.
- It is not very expressive.

Logistic Regression



The method

- For binary or multiclass classification only.
- Assumption on data:

$$\log \left(\frac{\mathbb{P}(Y_{new} = \bullet \mid X_{new})}{\mathbb{P}(Y_{new} = \bullet \mid X_{new})} \right) = w^\top X_{new} + b,$$

for unknown w, b .

- The Bayes classifier draws a hyperplane frontier:

$$f^*(x) = \begin{cases} \bullet & \text{if } w^\top x + b \geq 0 \\ \bullet & \text{if } w^\top x + b < 0 \end{cases}.$$

- Plug-in estimation via Regularized Maximum Likelihood: (\hat{w}, \hat{b}) such that

$$\frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-Y_i(\hat{w}^\top X_i + \hat{b})} \right) + \frac{\lambda}{2} \|\hat{w}\|^2$$

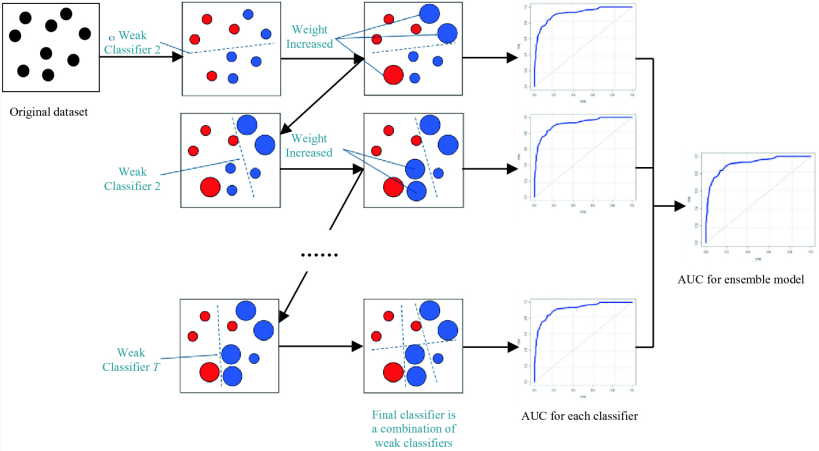
is minimal.

- Parameter: regularization coefficient $\lambda > 0$.

Pros and cons

- It is robust.
- It is computationally tractable (OK in high dimension and with large amount of data).
- It is not very expressive.

Gradient boosting



Source: <https://datascience.eu>

Gradient boosting

The method

- For binary classification and regression.
- No assumption on data.
- Combination of simple estimators (weak learners):

$$\hat{f}(x) = \hat{f}_1(x) + \dots + \hat{f}_T(x) \quad \text{such that} \quad \frac{1}{n} \sum_{i=1}^n \ell(\hat{f}(X_i), Y_i)$$

is almost minimal.

- Iterative procedure: \hat{f}_{t+1} is such that

$$\frac{1}{n} \sum_{i=1}^n \ell(\hat{f}_1(X_i) + \dots + \hat{f}_t(X_i) + \hat{f}_{t+1}(X_i), Y_i)$$

is almost minimal.

- Similar to gradient descent:

$$\hat{f}_{t+1}(X_i) \approx - \frac{\eta}{\text{normalization}} \frac{\partial \ell}{\partial x}(\hat{f}_1(X_i) + \dots + \hat{f}_t(X_i), Y_i).$$

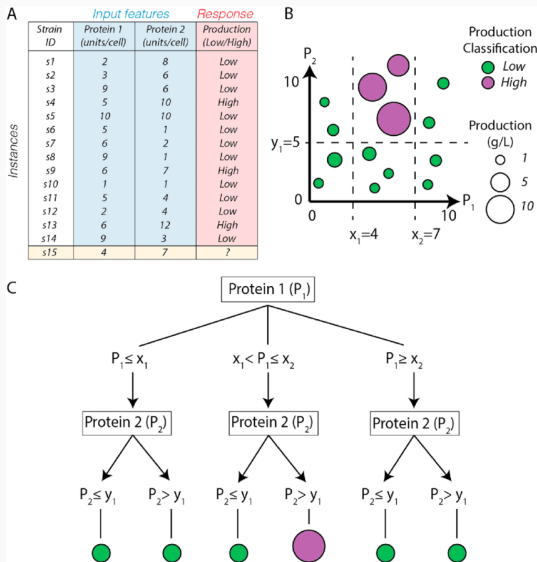
The method

- The case of L^2 regression: $\hat{f}_{t+1}(X_i) \approx Y_i - (\hat{f}_1(X_i) + \dots + \hat{f}_t(X_i))$.
- Parameters: number of simple estimators T , shrinkage coefficient $\eta \in]0, 1]$.

Pros and cons

- It is robust and efficient.
- It is fairly computationally tractable.
- It is very expressive.
- In practice, it has more than two parameters.

Decision tree



Source: <https://www.researchgate.net>

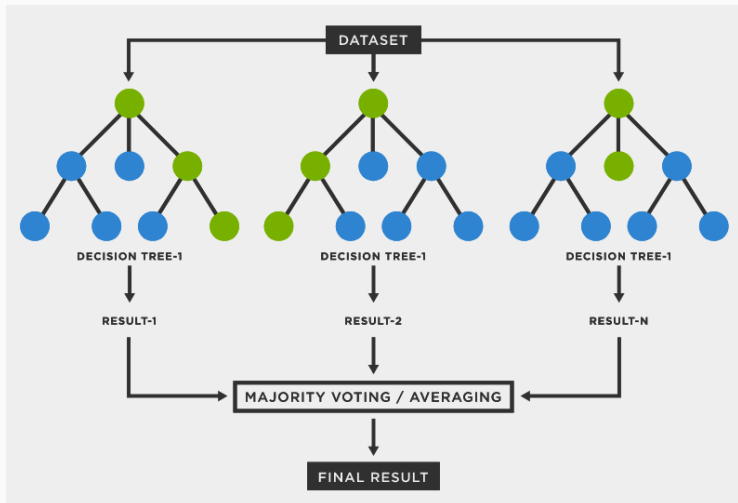
The method

- For multiclass classification and regression.
- No assumption on data.
- Data dependent partitioning of the input space: $x \in \text{Cell}_1$ or $\text{Cell}_2 \dots$
- Partition implemented as a binary tree.
- Piece-wise constant estimator: $\hat{f}(x) = \text{Value}_k$ for $x \in \text{Cell}_k$.
- Cells and values are determined in order to maximize the output *homogeneity*.
- Parameters: size of the tree (*i.e.* number of cells).

Pros and cons

- It is computationally tractable.
- It is very expressive.
- It is prone to overfitting.
- It is used in gradient boosting with few cells.

Random forest



Source: <https://www.tibco.com>

The method

- For multiclass classification and regression.
- No assumption on data.
- Combination of decision trees:

$$\hat{f}(x) = \frac{\text{tree}_1(x) + \dots + \text{tree}_T(x)}{T},$$

such that trees are roughly independent.

- Trees are learned on bootstrap samples.
 - Construction of partitions is perturbed by noise.
- Parameters: number of trees T , size of trees, level of noise in partitions building.

Pros and cons

- It is robust and efficient.
- It is fairly computationally tractable.
- It is very expressive.
- Construction is parallelizable.
- In practice, it has many parameters.

Feature importance and selection

- It is mainly model-dependent.
- In-model importance and post-selection:
 - Feature weights for linear models.
 - Contribution to the score for decision trees.
- Iterative procedure: score improvement when adding a feature.
- In-model importance and in-selection: sparse regularization for linear models.

Multiclass problems

- Natively handled by some methods: k -Nearest neighbors, trees, random forests, neural networks.
- Others are mainly for binary classification: linear models, boosting, Support Vector Machines.
- Four strategies:
 - Binary encoding of the class number + $\lceil \log_2(\# \text{ classes}) \rceil$ classifiers.
 - One versus One + $\binom{\# \text{ classes}}{2}$ classifiers.
 - One versus Rest + $\# \text{ classes}$ classifiers.
 - Hierarchical One versus Rest + $\# \text{ classes} - 1$ classifiers.

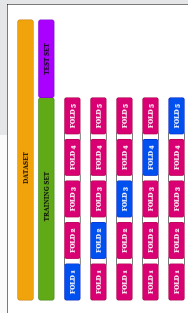
Evaluation and selection

Metrics

- It depends on what is important for you but is not necessarily reflected in the loss ℓ .
- Accuracy, balanced accuracy, top- k accuracy.
- Area under the ROC curve, F1 score.
- Mean squared error, R^2 score.

Generalization and model selection

- Cross-validation.
- Grid or random search.
- Regularization path.



Clustering

Three methodologies

- We observe only inputs: X_1, \dots, X_n .
- **Goal** Find a partition of the space $\{A, \bar{A}\}$ such that

$$\begin{cases} X_{new} \in A, X_{other} \in A \iff X_{new} \text{ and } X_{other} \text{ are similar;} \\ X_{new} \in \bar{A}, X_{other} \in \bar{A} \iff X_{new} \text{ and } X_{other} \text{ are similar.} \end{cases}$$

The latent variable model

- Partial observation in the classification setting: $(X_1, Y_1), \dots, (X_n, Y_n)$,
 $Y_i = \bullet$ or $\color{red}\bullet$.
- Expected clustering (\approx Bayes classifier):

$$A = \{x : \mathbb{P}(Y_{new} = \bullet \mid X_{new} = x) \geq \mathbb{P}(Y_{new} = \color{red}\bullet \mid X_{new} = x)\},$$

\bar{A} is the rest.

- Estimator:

$$\hat{A} = \left\{ x : \hat{\mathbb{P}}(Y_{new} = \bullet \mid X_{new} = x) \geq \hat{\mathbb{P}}(Y_{new} = \color{red}\bullet \mid X_{new} = x) \right\}.$$

The latent variable model

- Needs statistics.
- In practice:
 - We have to model the data distribution.
 - We resort to simple and non-robust estimators.
 - We are very limited.
- **Examples** Soft k -Means.

Empirical Risk Minimization

- For a dissimilarity function ℓ , find $\{A, A\}$ such that

$$\mathbb{E} [\ell(X_{new}, X_{other}) \mathbb{1}_{X_{new} \in A, X_{other} \in A \text{ or } X_{new} \in A, X_{other} \in A}]$$

is minimal.

- By the law of large numbers:

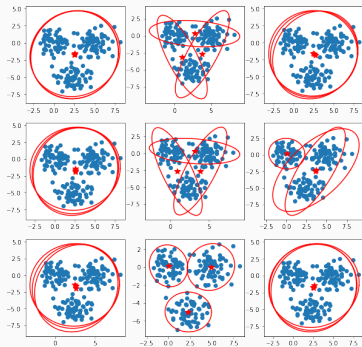
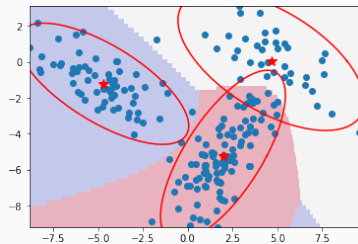
$$\frac{1}{2} \sum_{1 \leq i, j \leq n} \ell(X_i, X_j) \left[\frac{\mathbb{1}_{X_i \in A, X_j \in A}}{\#X_\ell \in A} + \frac{\mathbb{1}_{X_i \in A, X_j \in A}}{\#X_\ell \in A} \right] \approx \mathbb{E} [\ell(X_{new}, X_{other}) \mathbb{1} \dots].$$

- Estimator: $\{\hat{A}, \hat{A}\}$ such that $\frac{1}{n} \sum_{i=1}^n \ell(X_{new}, X_{other}) \mathbb{1} \dots$ is minimal.
- Needs optimization tricks.
- In practice:
 - We have to change the optimization problem because it is not tractable (NP-hard).
 - We have to design a simple and non-optimal iterative procedure because the optimization problem it is not tractable (NP-hard).
- **Examples** k -means, spectral clustering, hierarchical clustering.

Density-based approaches

- Implicitly estimate the density of X_{new} .
- Detect the modes and set automatically the number of groups.
- Needs intuition and geometry.
- In practice:
 - We have only an intuitive understanding.
 - We have to set a *density* parameter (neighborhood size, bandwidth).
- **Examples** DBSCAN, OPTICS, Mean shift.

Soft k -means



The method

- Partial observation in the classification setting: $(X_1, Y_1), \dots, (X_n, Y_n)$.
- Assumption on data:

$$\begin{cases} X_{new} \mid Y_{new} = \bullet \sim \mathcal{N}(\mu, \Sigma), & \mathbb{P}(Y_{new} = \bullet) = \pi, \\ X_{new} \mid Y_{new} = \bullet \sim \mathcal{N}(\mu, \Sigma), & \mathbb{P}(Y_{new} = \bullet) = \pi, \end{cases}$$

for unknown $\mu, \mu, \Sigma, \Sigma, \pi, \pi$.

- A posteriori distribution:

$$\mathbb{P}(Y_{new} = \bullet \mid X_{new} = x) = \frac{\pi \varphi(x)}{\pi \varphi(x) + \pi \varphi(x)} = p(x).$$

- Plug-in estimation via Maximum Likelihood: $\hat{\pi}, \hat{\pi}, \hat{\varphi}, \hat{\varphi}$ such that

$$\sum_{i=1}^n \log(\hat{\pi} \hat{\varphi}(X_i) + \hat{\pi} \hat{\varphi}(X_i))$$

is maximal.

The method

- Proxy maximization:

$$\mathbb{E} \left[\sum_{i=1}^n \log \left(\text{density}_{(X_1, Z_{t,1})}(X_i, Z_{t,i}) \right) \mid X_1, \dots, X_n \right]$$

$$= \sum_{i=1}^n [\hat{\rho}_t(X_i) \log(\hat{\pi}_{t+1} \hat{\varphi}_{t+1}(X_i)) + (1 - \hat{\rho}_t(X_i)) \log(\hat{\pi}_{t+1} \hat{\varphi}_{t+1}(X_i))],$$

where $Z_{t,i} \mid X_i = x$ has the distribution of $Y_{new} \mid X_{new} = x$ estimated with $\hat{\pi}_t, \hat{\pi}_{t+1}, \hat{\varphi}_t, \hat{\varphi}_{t+1}$.

- Estimator: $\hat{A} = \{x : \hat{\rho}_t(x) \geq \frac{1}{2}\}$.

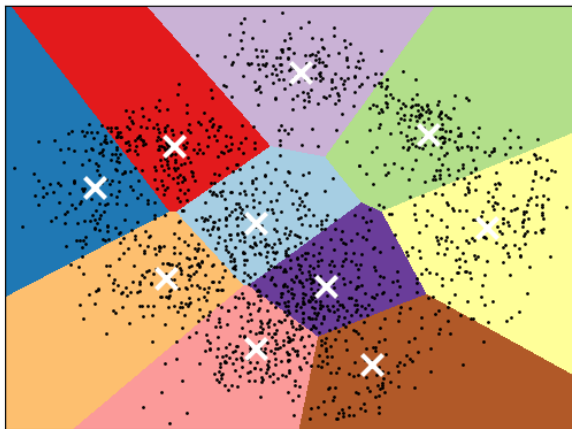
Iterative algorithm

1. Compute $\hat{\rho}_t(X_1), \dots, \hat{\rho}_t(X_n)$ with $\hat{\pi}_t, \hat{\pi}_t, \hat{\varphi}_t, \hat{\varphi}_t$.
2. $\hat{\pi}_{t+1} = \frac{1}{n} \sum_{i=1}^n \hat{\rho}_t(X_i)$.
3. $\hat{\mu}_{t+1} =$ empirical mean weighted by $\hat{\rho}_t(X_1), \dots, \hat{\rho}_t(X_n)$.
4. $\hat{\Sigma}_{t+1} =$ empirical covariance centered at $\hat{\mu}_{t+1}$ and weighted by $\hat{\rho}_t(X_1), \dots, \hat{\rho}_t(X_n)$.
5. Same for $\hat{\pi}_{t+1}, \hat{\mu}_{t+1}, \hat{\Sigma}_{t+1}$ with $(1 - \hat{\rho}_t(X_1)), \dots, (1 - \hat{\rho}_t(X_n))$.

Pros and cons

- It is a very simple and cheap iterative algorithm.
- It is suboptimal (does not necessarily return the MLEs).
- It is very sensitive to initialization $\hat{\mu}_0$ and $\hat{\rho}_0$ (in practice, it is initialized with k -means++ output).

K-means clustering on the digits dataset (PCA-reduced data)
Centroids are marked with white cross



Source: <https://scikit-learn.org>

The method

- No assumption on data.
- ERM with dissimilarity $\ell(x, x') = \|x - x'\|_2^2$: minimize

$$\begin{aligned} & \frac{1}{2} \sum_{1 \leq i, j \leq n} \|X_i - X_j\|_2^2 \left[\frac{\mathbb{1}_{X_i \in \hat{A}, X_j \in \hat{A}}}{\#X_\ell \in \hat{A}} + \frac{\mathbb{1}_{X_i \in \hat{A}, X_j \in \hat{A}^c}}{\#X_\ell \in \hat{A}^c} \right] \\ &= \sum_{i=1}^n \left[\|X_i - \hat{\mu}\|_2^2 \mathbb{1}_{X_i \in \hat{A}} + \|X_i - \hat{\mu}^c\|_2^2 \mathbb{1}_{X_i \in \hat{A}^c} \right], \end{aligned}$$

with $\hat{\mu} = \frac{1}{\#X_\ell \in \hat{A}} \sum_{i=1}^n X_i \mathbb{1}_{X_i \in \hat{A}}$.

- Alternating procedure (Lloyd's algorithm):
 1. Find $\{\hat{A}, \hat{A}^c\}$ with $\hat{\mu}$ and $\hat{\mu}^c$ fixed: Voronoi partitioning.
 2. Compute $\hat{\mu}$ and $\hat{\mu}^c$ with fixed partition $\{\hat{A}, \hat{A}^c\}$.
- Estimator: $\hat{A} = \{x : \|x - \hat{\mu}\| \leq \|x - \hat{\mu}^c\|\}$.

Iterative algorithm

1. $\hat{p}_t(X_i) = \mathbb{1}_{\|X_i - \hat{\mu}_t\| \leq \|X_i - \hat{\mu}_{t'}\|}$.
2. $\hat{\mu}_{t+1} =$ empirical mean weighted by $\hat{p}_t(X_1), \dots, \hat{p}_t(X_n)$.
3. Same for $\hat{\mu}_{t+1}$ with $(1 - \hat{p}_t(X_1)), \dots, (1 - \hat{p}_t(X_n))$.

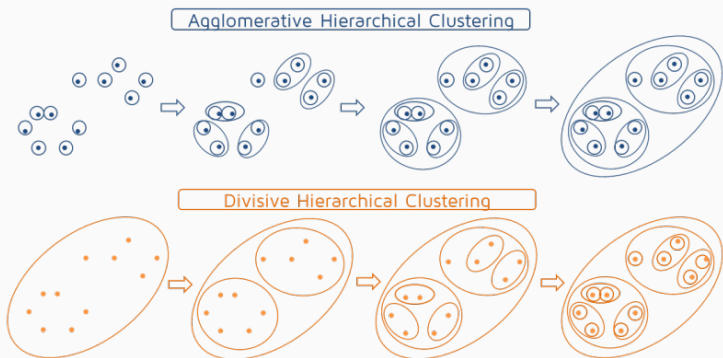
Connection with soft k-means

- Hard assignment: $\hat{p}_t(x) = \mathbb{1}_{\dots}$ instead of $\hat{\mathbb{P}}(Y_{new} = \bullet \mid X_{new} = x)$.
- No a priori in partitioning (or $\hat{\pi}_{t+1} = \frac{1}{2}$).
- No variance (or $\hat{\Sigma}_{t+1} \rightarrow 0$).

Pros and cons

- It is a simple and cheap iterative algorithm.
- It is suboptimal (does not necessarily return the optimal partition).
- It is very sensitive to initialization $\hat{\mu}_0$ and $\hat{\mu}_0$: k-means++.
- Groups are convex, not hierarchically structured.

Agglomerative clustering



Source: <https://quantdare.com>

Agglomerative clustering

The method

- No assumption on data.
- Iterative procedure:
 1. Start with n groups: a point = a group.
 2. Iteratively merge *nearest* groups.
- Famous *distances* between groups A and A :

- Single linkage:

$$d(A, A) = \min_{x \in A, x \in A} \|x - x\|.$$

- Ward's criterion:

$$d(A, A) = \text{Inertia}(A \cup A) - \text{Inertia}(A) - \text{Inertia}(A) = \frac{\text{Size}(A)\text{Size}(A)}{\text{Size}(A) + \text{Size}(A)} \|\mu - \mu\|^2.$$

Connection with k -means

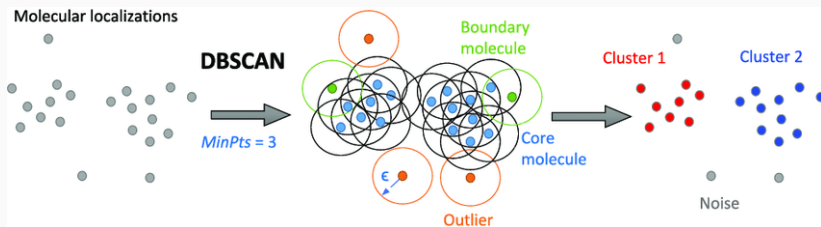
- With Ward's criterion: tend to minimize the total inertia

$$\sum_{i=1}^n \left[\|X_i - \hat{\mu}\|_2^2 \mathbb{1}_{X_i \in \hat{A}} + \|X_i - \hat{\mu}\|_2^2 \mathbb{1}_{X_i \in \hat{A}} \right]$$

with a hierarchical procedure.

Pros and cons

- It is a simple iterative algorithm.
- It is suboptimal (does not necessarily return the optimal partition) but provides a hierarchical structure of groups.
- It is deterministic given X_1, \dots, X_n .
- Groups may be non-convex.



Source: <https://www.researchgate.net>

The method

- *Density Based Spatial Clustering and Applications with Noise.*
- No assumption on data.
- Iterative growing and birth of groups.
- Three types of points:
 1. Core points (at least m neighbors with a distance ϵ).
 2. Reachable points (non-core points in the ϵ -neighborhood of a core point).
 3. Outliers.
- Parameters: number of neighbors m and radius ϵ .
- Movie.

Connection with agglomerative clustering

- With $m = 2$, DBSCAN is similar to single linkage with a dendrogram cut at ϵ .

Pros and cons

- It is a simple iterative algorithm.
- No need to specify the number of groups.
- Groups may be non-convex.
- It is barely sensitive to initialization (for reachable points).
- It cannot detect groups with different densities: OPTICS (*Ordering Points To Identify the Clustering Structure*).

Metrics

- Global: inertia

$$I = \sum_{i=1}^n \left[\|X_i - \hat{\mu}\|_2^2 \mathbb{1}_{X_i \in \hat{A}} + \|X_i - \hat{\mu}\|_2^2 \mathbb{1}_{X_i \in \hat{A}} \right].$$

- Individual and global: silhouette coefficient

$$S = \frac{1}{n} \sum_{i=1}^n \frac{b_i - a_i}{\max(a_i, b_i)},$$

where a_i = average distance of X_i to its group and b_i = average distance of X_i to the nearest group.

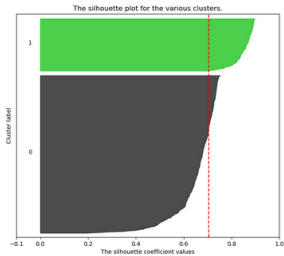
- Make sens for convex groups.

Model selection

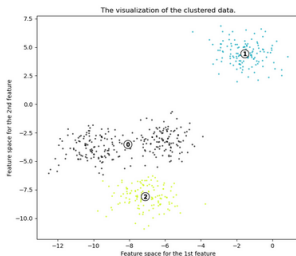
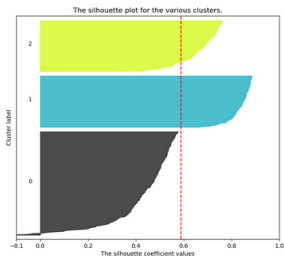
- The elbow method on inertia.
- Analyzing the silhouette coefficient: example.

Evaluation and selection

Silhouette analysis for KMeans clustering on sample data with $n_clusters = 2$

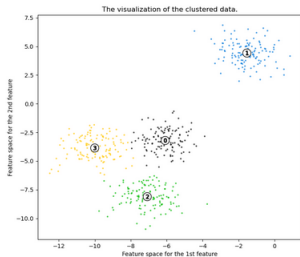
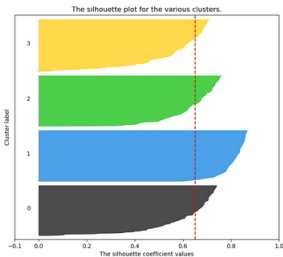


Silhouette analysis for KMeans clustering on sample data with $n_clusters = 3$

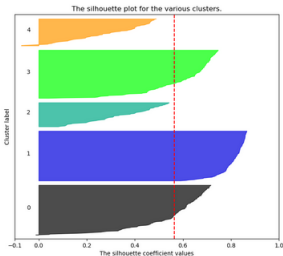


Evaluation and selection

Silhouette analysis for KMeans clustering on sample data with $n_clusters = 4$



Silhouette analysis for KMeans clustering on sample data with $n_clusters = 5$



What's next?

Other learning domains

- Neural networks.
- Dimensionality reduction.
- Data preprocessing.
- Time series.
- Reinforcement learning.
- Data generation.
- Active learning.
- Domain adaptation.
- Image and natural language processing.
- Causality.
- Visualization.
- ...