

Unsupervised Learning	Category	Approach	Application	Description	Hyperparameters & Regularization	Strength	Weaknesses	Performance Metrics
	Hard Clustering	K-Means-Clustering	Clustering Analysis (grouping instances to clusters)	<ul style="list-style-type: none"> <li>- randomly assign initial centroids spaced out on data</li> <li>- assign points to their nearest centroids (euclidean &amp; manhattan)</li> <li>- calculate center for all clusters</li> <li>- make center new centroid</li> <li>- re-assign data to clusters and repeat until convergence</li> </ul>	- number of clusters	- good for large datasets- easy cluster assignment	- only works for spherical data- only works for numerical features	- Inertia (mean squared distance of all centroids to instances in cluster) - Elbow-Method (Inertia vs. number of clusters; "joint" is optimal) - Silhouette-Coefficient
	Hard & Soft Clustering	Gaussian Mixture Models	Density Estimation	<ul style="list-style-type: none"> <li>- assume number of clusters, mu and sigma for each cluster</li> <li>- Expectation:               <ul style="list-style-type: none"> <li>* Compute membership probability for each data record</li> <li>* Calculate expression for responsibility for likelihood of all data points</li> </ul> </li> <li>- Maximization:               <ul style="list-style-type: none"> <li>* Calculate total likelihood for all clusters</li> <li>* Update mu and sigma for all clusters based on likelihood</li> </ul> </li> <li>- Repeat until convergence</li> <li>- Assign instances to cluster with highest membership probability</li> </ul>	<ul style="list-style-type: none"> <li>- number of clusters</li> <li>- initial mu and sigma</li> </ul>	<ul style="list-style-type: none"> <li>- hard and soft clustering</li> <li>- shows probabilities for data points</li> <li>- good for more complex datasets (ellipsoids)</li> </ul>	<ul style="list-style-type: none"> <li>- slow and computationally expensive</li> <li>- can get stuck in local optima</li> <li>- only numerical features</li> </ul>	- Bayesian Information Criterion - Akaike Information Criterion
	Hard Clustering	Hierarchical Clustering	Clustering Analysis (grouping instances to clusters)	<ul style="list-style-type: none"> <li>- Bottom-up or top-down approach; dendrogram-representation:               <ul style="list-style-type: none"> <li>- calculate proximity matrix from one instance to the other</li> </ul> </li> <li>- Bottom-Up (agglomerative):               <ul style="list-style-type: none"> <li>* initiate every data point as a single cluster</li> <li>* step-by-step, add closest datapoints to clusters and clusters to clusters (single linkage, average linkage or complete linkage)</li> <li>* repeat until all data points are within one cluster (universe cluster)</li> </ul> </li> <li>- Top-Down (divisive):               <ul style="list-style-type: none"> <li>* all instances start in universe cluster</li> <li>* split clusters based on farthest distances</li> <li>* repeat until all data points are individual clusters</li> </ul> </li> <li>- find longest vertical line from one cluster to the other (biggest distance)</li> </ul>	- None	<ul style="list-style-type: none"> <li>- non-parametric</li> <li>- easy to implement</li> <li>- informative cluster hierarchy</li> </ul>	<ul style="list-style-type: none"> <li>- very sensitive to outliers</li> <li>- not feasible for large datasets</li> <li>- clusters cannot be un-made</li> </ul>	
	Hard Clustering	DBSCAN	Outlier Detection	<ul style="list-style-type: none"> <li>- set MinPTS and Surface Radius</li> <li>- start at any data record and check if it forms a cluster; if it does, make it a center point</li> <li>- check for all other points in the cluster if they are either center points or border points (don't form their own cluster)</li> <li>- start at record that was not assigned before and restart</li> <li>- non-assigned data records are outliers</li> </ul>	- MinPTS and Surface Radius	<ul style="list-style-type: none"> <li>- handles outliers well</li> <li>- outlier detection</li> <li>- doesn't need number of clusters</li> <li>- can handle arbitrarily shaped datasets</li> </ul>	<ul style="list-style-type: none"> <li>- hard to setup minPTS and radius (domain expert needed)</li> <li>- needs clusters of same density</li> </ul>	

Regression	Category	Approach	Application	Description	Hyperparameters & Regularization	Strength	Weaknesses	Performance Metrics
	Regression	(Multi-) Linear Regression	Predicting continuous value	<ul style="list-style-type: none"> <li>- needs standard deviation of the data, since RSS is MLE for normal distribution</li> <li>- needs to pass test for linearity</li> <li>- learns linear relation of data to target variable and makes predictions for target variable based on new feature values</li> </ul>	<ul style="list-style-type: none"> <li>- lambda-Parameter for weights</li> <li>- Drastic: Lasso Regression: add <math>\lambda \sum \text{coefficients}</math> to cost function; cancels out small weights</li> <li>- Less drastic: Ridge regression; adds <math>\lambda \sum \text{coefficients}^2</math> to the cost function, dampening large weights</li> </ul>	<ul style="list-style-type: none"> <li>- weight parameters are transparent and easy to understand for linear models</li> </ul>	<ul style="list-style-type: none"> <li>- prone to overfitting, always needs regularization</li> <li>- susceptible to outliers</li> </ul>	<ul style="list-style-type: none"> <li>- mean squared error (susceptible to outliers)</li> <li>- root mean squared error (RMSE – more interpretable like standard dev.)</li> <li>- mean absolute error</li> <li>- <math>R^2</math> (how much of variance in dependent variable stems from independent variables; 0 to 1, where 1 is optimal)</li> </ul>
	Regression	Polynomial Regression	Predicting continuous value	-applicable to non-linear data; rest: see above	- see above	- works with non-linear data	- see above	- see above
	Regression	Quantile Regression	Predicting continuous value	- separate dependent variable into different segments (quantiles) to train the model individually for each segment	- see above	- works with multimodal or skewed data	- see above	- see above
	Regression	Support Vector Machines	Predicting continuous value	<ul style="list-style-type: none"> <li>- learns relation of data to target variable by fitting as many points as possible on the hyperplane</li> <li>- tries to minimize weights instead of minimizing error!</li> <li>- all points have to be within certain margin (E)</li> <li>- derived function is used to make predictions for new instances</li> <li>- works like the "inverse" of the large margin classifier</li> <li>- for non-linear problems, soft margins are introduced with slack variables ksi</li> </ul>	<ul style="list-style-type: none"> <li>- Parameter C: Increase to decrease regularization; decrease to increase regularization (allows more point on margin)</li> <li>- Slack Variables ksi allow for points outside the margin (soft margin)</li> </ul>	- Not given in script	- Not given in script	- Not given in script
	Regression	Regression Trees	Predicting target variable	<ul style="list-style-type: none"> <li>- Build the tree:               <ul style="list-style-type: none"> <li>- calculate standard deviation of target variable</li> <li>- calculate standard deviation of all manifestations of feature variable and weight it with its probability</li> <li>- calculate reduction in standard deviation of target variable by subtracting standard deviation of all feature variables</li> <li>- select the feature that causes the highest reduction in standard deviation as the root/ next node</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>- Pruning the tree: continue as long as coefficient of variation (<math>\sigma / \mu</math>) stays under a certain value</li> </ul>	<ul style="list-style-type: none"> <li>- little assumptions on the data</li> <li>- allows regression based on non-numerical data (e.g. categorical data)</li> </ul>	- needs pruning or would overfit	- Not given in script

Category	Approach	Application	Description	Hyperparameters & Regularization	Strength	Weaknesses	Performance Metrics
Classification	Logistic Regression	Binary or categorical classification	<ul style="list-style-type: none"> <li>- based on bernoulli distribution</li> <li>- uses coefficients to determine a probability</li> <li>- decision for classification is done based on threshold (e.g. 90%)</li> </ul>	<ul style="list-style-type: none"> <li>- lambda or inverse (C)</li> <li>- Lasso Regression (drastic)</li> <li>- Ridge Regression (less drastic)</li> </ul>	<ul style="list-style-type: none"> <li>- takes binomious categorical or discrete data</li> <li>- works with big samples</li> </ul>	<ul style="list-style-type: none"> <li>- not feasible for non-linearly seperable datasets</li> </ul>	Confusion Matrix: <ul style="list-style-type: none"> <li>- Precision <math>(TP / (TP + FP))</math></li> <li>- Recall <math>(TP / (TP + FN))</math></li> <li>- F1-Score <math>[2 * Precision * Recall / (Precision + Recall)]</math></li> <li>- Accuracy <math>[(TP + TN) / (TP + TN + FP + FN)]</math></li> </ul>
Classification	Support Vector Machines	Categorical or discrete classification	Large margin classifier: Maximizes the width between support vectors of different classes: <ul style="list-style-type: none"> <li>- Hard Margin Problem (linearly seperable data): All data points have to be classified in the correct class; no room for error (overfitting?)</li> <li>- Soft Margin Problem (nonlinearly seperable data): Model allows certain errors within the classification, making the model more robust (less overfitting) w. slack-variable zeta and parameter C</li> <li>- Dimensionality: To make non-linearly separable data linearly separable, one can increase its dimension → curse of dimensionality (overfitting and cost for calculation)</li> <li>- Kernel-Trick: To avoid this, one can use the kernel trick, because it only transfers the relevant datapoints more efficiently w. dot products               <ul style="list-style-type: none"> <li>→ Polynomial Kernel</li> <li>→ RBF (Gaussian) Kernel</li> </ul> </li> </ul>	<ul style="list-style-type: none"> <li>- Regularization Parameter C (increasing → less regularization)</li> <li>- for RBF (Gaussian) Kernel: gamma (increasing → less regularization)</li> </ul>	<ul style="list-style-type: none"> <li>- high discrimination power</li> <li>- takes continuous and categorical data</li> <li>- works with linear and nonlinear problems</li> <li>- memory efficient (only needs support vectors)</li> <li>- good generalization</li> </ul>	<ul style="list-style-type: none"> <li>- needs proper mathematical understanding</li> <li>- harder to understand than more simple models</li> </ul>	Confusion Matrix: <ul style="list-style-type: none"> <li>- Precision <math>(TP / (TP + FP))</math></li> <li>- Recall <math>(TP / (TP + FN))</math></li> <li>- F1-Score <math>[2 * Precision * Recall / (Precision + Recall)]</math></li> <li>- Accuracy <math>[(TP + TN) / (TP + TN + FP + FN)]</math></li> </ul>
Classification	Decision Trees	Non-parametric classification for all kinds of data	Asking binary question: Yes or no – trees are built on roots (initial question) nodes (decision points) and leafs (answers) Splits are made by the highest gain of information for the target variable: <ul style="list-style-type: none"> <li>- ID 3: "Information Gain":               <ul style="list-style-type: none"> <li>→ Gain <math>(f) = Entropy (dataset) - weighted Entropy (f) (Dataset)</math></li> <li>- ID 4.5: "Split Info":</li> </ul> </li> <li>- Avoids completely homogenous datasets (poor generalization!) and can use continuous and discrete features, while handling missing values and decision tree pruning.</li> <li>- CART</li> <li>- Reduction in Impurity by using the Gini Index, implements post-pruning</li> </ul>	Pruning: Pre-Pruning: Only split a node if a certain condition is met (max depth, min members in a node ...) Post-Pruning: Remove parts of a tree that hold little information with a leaf that holds the label of the most frequent feature in the leaf (implemented in CART)	<ul style="list-style-type: none"> <li>- non-parametric model</li> <li>- easy to interpret</li> <li>- uses discrete and numerical features</li> </ul>	<ul style="list-style-type: none"> <li>- always overfits if not pruned (recreates the input data by design)</li> <li>- very susceptible to outliers</li> </ul>	Confusion Matrix and:  ROC (Receiver Operating Characteristic Curve): <ul style="list-style-type: none"> <li>- shows performance at all classification thresholds</li> <li>- Recall on abzissa, False Positive Rate on ordinate:</li> <li>→ Recall = <math>TP / (TP + FN)</math>, FPR = <math>FP / (FP + TN)</math></li> </ul> AOC: Area under ROC (from 0 to 1)
Classification	Ensemble Models	Non-parametric classification for all kinds of data	Meta Method: Combine several homogen or heterogen (weak) learners that are strong at specific parts of input space → combination evens out errors but overall prediction stays strong  Bagging (for homegen weak learners): <ul style="list-style-type: none"> <li>- all samples are thrown in a "Bag" and repeatedly drawn to train the weak lerners, afterwards "returned to the bag" (each sample can be drawn multiple times)</li> <li>- all weak learners vote for classification, decision by majority vote ("agg"regation step)</li> </ul> Pasting (for homogen weak learners): <ul style="list-style-type: none"> <li>- all samples are thrown in a bag and drawn only once as a subset to train the weak lerners (each sample can only be used once!)</li> </ul> Boosting: <ul style="list-style-type: none"> <li>- iterative approach → performance of all weak learners is improved by training samples step-wise</li> <li>- all samples are drawn repeatedly for training</li> <li>- misclassified instances get a bigger probability to be re-drawn</li> </ul> Adaptive Boosting: <ul style="list-style-type: none"> <li>- like boosting, but successful learners get more important vote in final classification step</li> </ul>	<ul style="list-style-type: none"> <li>- depends on the chosen models</li> </ul>	<ul style="list-style-type: none"> <li>- combines several different ML approaches and evens out their weaknesses; usually a good choice</li> </ul>	<ul style="list-style-type: none"> <li>- results might become more opaque</li> </ul>	Confusion Matrix and:  ROC (Receiver Operating Characteristic Curve): <ul style="list-style-type: none"> <li>- shows performance at all classification thresholds</li> <li>- Recall on abzissa, False Positive Rate on ordinate:</li> <li>→ Recall = <math>TP / (TP + FN)</math>, FPR = <math>FP / (FP + TN)</math></li> </ul> AOC: Area under ROC (from 0 to 1)
Classification	Random Forest	Non-parametric classification for all kinds of data	Forest creation and voting: <ul style="list-style-type: none"> <li>- build trees by choosing subset of training samples and usually subset of training features</li> <li>- trees vote for classification</li> <li>- training on "bagging" method</li> </ul>	<ul style="list-style-type: none"> <li>- number of trees in forest</li> <li>- maximal depth of a tree</li> <li>- minimal samples per leaf</li> </ul>	<ul style="list-style-type: none"> <li>- big features spaces</li> <li>- random subset of features reduces correlation between learners; reducing variance</li> </ul>	<ul style="list-style-type: none"> <li>- ?</li> </ul>	See above

Category	Approach	Application	Description	Hyperparameters & Regularization	Strength	Weaknesses	Performance Metrics
Optimization	Genetic Algorithm	Heuristic optimization methods for hyperparameter tuning	Creation of population: <ul style="list-style-type: none"> <li>- individual problem solutions are encoded in a Chromosome with Genes (0 and 1 values)</li> </ul> Evolution of fitness function: <ul style="list-style-type: none"> <li>- putting individual solutions towards the fitness function</li> </ul> Selection of fittest: <ul style="list-style-type: none"> <li>- looking for a certain amount of highest scoring chromosomes</li> </ul> Crossover: <ul style="list-style-type: none"> <li>- re-combining chromosomes on single or dual crossover</li> </ul> Mutation: <ul style="list-style-type: none"> <li>- flipping one bit to prevent being stuck in local optima</li> </ul>	None	None	None	None