

# ALGORITHM SELECTION

Machine Learning

## ALGORITHM SELECTION CRITERIA

### CATEGORIZE PROBLEM

- Input category (supervised | unsupervised | reinforcement learning)
- Output category (regression | classification | clustering | anomaly detection)

### UNDERSTAND CONSTRAINTS

- Data storage capacity
- Speed of prediction
- Speed of learning

### UNDERSTAND DATA

- Sample set size
- Data types (categorical | continuous)
- Missing data sensitivity
- Outliers sensitivity
- Features multi-collinearity

### BUSINESS CRITERIA

- Meeting business goals
- Pre-processing requirements
- Model accuracy
- Model scalability
- Model explain-ability
- Model complexity
  - # of features
  - Feature engineering needs

- computational overhead
- # of hyper-parameters

### SELECTION QUESTIONS

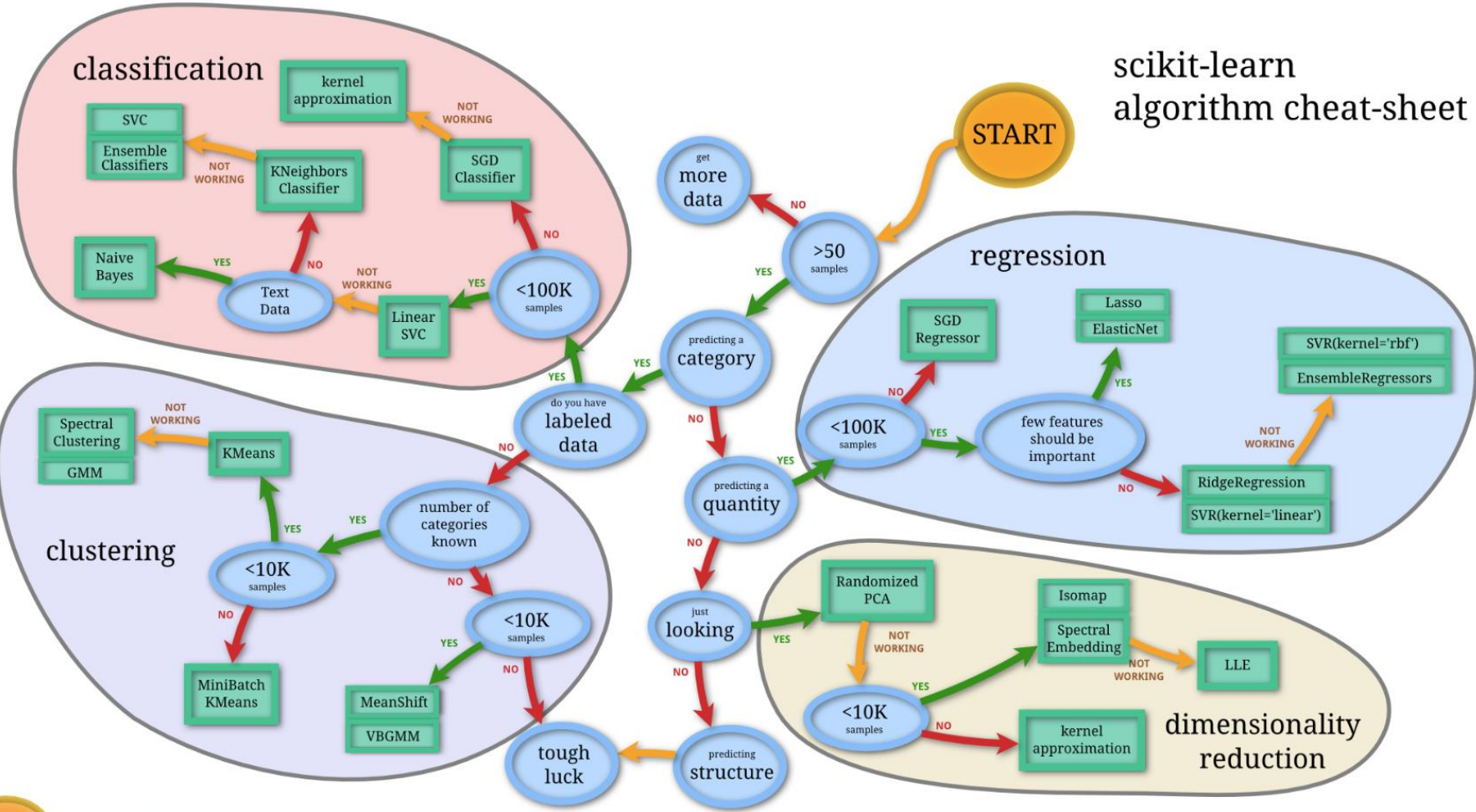
- Number of training examples, (how large is the training dataset)
  - If small: high bias/low variance classifiers (e.g., Naïve Bayes), less likely to over-fit
  - If large: low bias/high variance classifiers (e.g., kNN or Logistic regression)
- Dimensionality of the feature space
- Is the problem linearly separable?
- Are features independent?
- Are features expected to be linearly dependent with the target variable?
- Is overfitting expected to be a problem/
- System requirements: speed, performance, memory usage
- Does it requires variables to be normally distributed?
- Does it suffers multi-collinearity issue?
- Does it do as well with categorical as continuous variables?
- Does it calculate CI without CV?
- Does it conduct variables selection without stepwise?
- Does it apply to sparse data?

Start with simple like logistic regression to set a baseline and only make it more complicated if you need to. At that point, tree ensembles, and in particular Random Forests since they are easy to tune, might be the right way to go. If you feel there is still room for improvement, try GBDT or get even fancier and go for Deep learning

## REFERENCES

- [Choose the right machine learning algorithm](#)
- [Choosing right estimator](#)
- [ML algorithms pros and cons](#)
- [Modern machine learning algorithms: strengths and weaknesses](#)
- [ML algorithms cheat sheets](#)

# scikit-learn algorithm cheat-sheet



Algorithm	Applications	Advantages	Disadvantages
<b>Principal Component Analysis (PCA)</b>	<ul style="list-style-type: none"> <li>Dimensionality reduction</li> <li>Data pre-processing method for other ML algorithms sensitive to feature multi-collinearity</li> </ul>	<ul style="list-style-type: none"> <li>Only the data pre-processing option of highly correlated features to prevent overfitting for other ML algorithms sensitive to multi-collinearity</li> <li>Visual representation of highly dimensional features</li> </ul>	<ul style="list-style-type: none"> <li>Sometimes is difficult to interpret low dimensional representation of variables (principal components) as linear combination of features</li> </ul>
<b>K-Means Method</b>	<p>Finding natural grouping of observations:</p> <ul style="list-style-type: none"> <li>Customer segmentation</li> <li>Grouping similar items in e-commerce</li> <li>Social network analysis</li> </ul>	<ul style="list-style-type: none"> <li>Flexible enough to get reasonable results if you pre-process your data and engineer useful features</li> <li>Fast and simple algorithm</li> <li>Can process large datasets</li> <li>Often terminates at a local optimum</li> <li>Applicable only when mean is defined</li> </ul>	<ul style="list-style-type: none"> <li>Not applicable for categorical data</li> <li>Unable to handle noisy data</li> <li>Need to define number of clusters in advance</li> <li>Not suitable to discover clusters with non-convex shapes</li> <li>Poor performance for not globular clusters</li> </ul>
<b>K-nearest Neighbour (k-NN) Classifier</b>		<ul style="list-style-type: none"> <li>Nonparametric</li> <li>Zero cost in the learning process</li> <li>Classifying any data whenever finding similarity measures of any given instance intuitive approach</li> <li>Robust to outliers on the predictors</li> </ul>	<ul style="list-style-type: none"> <li>Hard to interpret the result</li> <li>Performs poorly for high-dimensional data, use regression or ensembles</li> <li>Lack of explicit model training</li> <li>Susceptible to correlated inputs and irrelevant features</li> <li>Very difficult in handling data of mixed types</li> <li>Memory intensive as they need to save each training observation</li> </ul>
<b>Affinity Propagation</b>		<ul style="list-style-type: none"> <li>Used for smaller and uneven classes patterns</li> <li>No need to specify number of clusters</li> </ul>	<ul style="list-style-type: none"> <li>Slow and memory intensive, so it does not scale well with datasets</li> <li>Poor performance for non-globular clusters</li> </ul>
<b>Hierarchical or Agglomerative Clustering</b>		<ul style="list-style-type: none"> <li>Clusters are not assume to be globular</li> <li>Scales well with dataset</li> </ul>	<ul style="list-style-type: none"> <li>Need to choose number of the clusters</li> </ul>
<b>DBSCAN Clustering</b>		<ul style="list-style-type: none"> <li>Density based clustering</li> <li>Not assume globular clusters</li> <li>Scalable with dataset</li> </ul>	<ul style="list-style-type: none"> <li>Difficult to tune sensitive hype-parameters</li> </ul>
<b>Linear Regression (Regularized)</b>	<ul style="list-style-type: none"> <li>Time to go one location to another</li> <li>Predicting sales of particular product next month</li> <li>Predict real estate prices</li> </ul>	<ul style="list-style-type: none"> <li>Easy to understand – clearly seen what the biggest drivers of the model are</li> <li>Can model only linear relationships between variables</li> </ul>	<ul style="list-style-type: none"> <li>Not flexible enough to capture more complex patterns and adding the right interaction terms or polynomials can be tricky and time consuming</li> </ul>

Algorithm	Applications	Advantages	Disadvantages
<b>Linear Regression (Regularized)</b>	<ul style="list-style-type: none"> <li>• Predict stock price movements</li> <li>• Predict student scores</li> <li>• Impact of blood alcohol content on coordination</li> <li>• Predict monthly gift card sales and improve yearly revenue projections</li> </ul>	<ul style="list-style-type: none"> <li>• Model can be easily updated with new data using stochastic gradient descent</li> </ul>	<ul style="list-style-type: none"> <li>• Perform poorly if there are non-linear relationships</li> <li>• Tendency for model to over-fit, risk can be reduced by regularization (Lasso, Ridge, Elastic-Net)</li> <li>• Sensitive to outliers</li> <li>• Unstable in case features are redundant, multi-collinear</li> </ul>
<b>Naïve Bayes</b>	<ul style="list-style-type: none"> <li>• Sentiment analysis and text classification</li> <li>• Recommendation systems like Netflix, Amazon</li> <li>• Spam detectors</li> <li>• Face recognition</li> </ul>	<ul style="list-style-type: none"> <li>• Easy to train as it converge quicker than discriminative models (ex. Logistic regression) under conditional independence assumption, so you need less training data</li> <li>• Scale well with growing dataset</li> <li>• Along with the simplicity outperform highly sophisticated classification methods (in case of not proper tuning)</li> <li>• Good choice if CPU and memory resources are a limiting factor</li> <li>• Can easily obtain the probability for a prediction</li> <li>• Relatively simple and straightforward to use</li> <li>• Can deal with some noise and missing data</li> <li>• Nonparametric – non distribution requirements</li> <li>• Good for few categorical variables</li> </ul>	<ul style="list-style-type: none"> <li>• Can handle multiple classes</li> <li>• Prone to bias when increasing the number of training sets</li> <li>• Assumes all features are independent and equally important, which is unlikely in real-world cases</li> <li>• Sensitive to how the input data is prepared</li> <li>• Can't learn interactions between features</li> <li>• Suffers multi-collinearity, but even if conditional independence assumption is not met it perform pretty well in practice</li> <li>• Often beaten by more sophisticated models properly tuned</li> </ul>

Algorithm	Applications	Advantages	Disadvantages
<b>Support Vector Machines (SVM)</b>	<ul style="list-style-type: none"> <li>Used for pattern recognition and classification with two classes</li> <li>Popular in text classification problems</li> <li>Detecting persons with common diseases such as diabetes</li> <li>Hand-written character recognition</li> <li>Text categorization – news articles by topics</li> <li>Stock market price prediction</li> </ul>	<ul style="list-style-type: none"> <li>Can utilize predictive power of linear combination of inputs</li> <li>Good prediction in a variety of situations</li> <li>Low generalization error</li> <li>Easy to interpret results</li> <li>Use non-linear kernels to model complex decision boundaries</li> <li>Robust to overfitting especially in high dimensional space</li> <li>Linear kernel is similar to Logistic regression which can also be used with different kernels</li> <li>Good in a high-dimensional space (e.g text classification)</li> <li>High accuracy</li> <li>Good theoretical guarantees regarding overfitting</li> <li>No distribution requirement</li> <li>Compute hinge loss</li> <li>Flexible selection of kernel for non-linear correlation</li> <li>Not suffer multi-collinearity</li> </ul>	<ul style="list-style-type: none"> <li>Very black box</li> <li>Sensitive to tuning parameters and kernel choice</li> <li>Testing data should be near to training data</li> <li>Hard to interpret</li> <li>Can be inefficient to train as it is memory-intensive for large datasets</li> <li>Not for most „industry scale“ applications (anything beyond the toy or lab problem)</li> <li>Weak in natural handling of mixed data types and scalability to large datasets</li> <li>Currently in industry Random Forest are preferred over SVM's</li> </ul>
<b>Stochastic Gradient Descent</b>	<p>Large scale and sparse machine learning problems</p> <ul style="list-style-type: none"> <li>Text classification</li> <li>Natural Language Processing</li> </ul>	<ul style="list-style-type: none"> <li>Simple yet very efficient approach to discriminative learning</li> <li>Used for large scale learning in order of magnitude <math>10^5</math> for training examples and features</li> <li>Easy implementation (lots of opportunities for code tuning)</li> </ul>	<ul style="list-style-type: none"> <li>Only linear classifier comparable with SVC using linear kernel or Logistic Regression</li> <li>Sensitive to feature scaling</li> </ul>

Algorithm	Applications	Advantages	Disadvantages
<b>Decision Trees</b>	<ul style="list-style-type: none"> <li>Used very rarely, but basis of very efficient ensemble methods, like random forest or gradient tree boosting</li> </ul> <p>Excellent tools for helping to choose between several courses of action</p> <ul style="list-style-type: none"> <li>Investment decisions</li> <li>Customer churn</li> <li>Bank loan defaulters</li> <li>Build vs. buy decisions</li> <li>Sales lead qualifications</li> </ul>	<ul style="list-style-type: none"> <li>Easily handle feature interactions and complex non-linear relationships due to branching algorithm based on maximizing information gain</li> <li>Non-parametric – no distribution requirements</li> <li>Less sensitive for outliers</li> <li>Easily interpretable</li> <li>Can handle missing data</li> <li>Able to handle both numerical and categorical data</li> <li>Performs well with large datasets</li> <li>Heuristic</li> <li>No suffer multi-collinearity</li> <li>Handle very well large dimensional spaces as well as large number of training examples</li> </ul>	<ul style="list-style-type: none"> <li>Can't work on (linear) combinations of features</li> <li>Relatively less predictive in many situations</li> <li>Practical decision-tree learning algorithms cannot guarantee to return the globally optimal decision tree</li> <li>Can lead to overfitting, but ensemble methods prevent overfitting by feature or training sample randomization</li> <li>Not often used its own for prediction because it's also often too simple and not powerful enough for complex data</li> <li>Good for only few categorical variables</li> <li>Can easily over-fit, this can be alleviated by using ensembles like RF or GBM</li> </ul>
<b>Bagging</b>		<ul style="list-style-type: none"> <li>Helps reduce variance and avoid overfitting</li> </ul>	<ul style="list-style-type: none"> <li>Loss of interpretability of the model</li> <li>Problem with high bias if not modelled properly</li> <li>Computationally expensive</li> </ul>
<b>Random Forest</b>	<ul style="list-style-type: none"> <li>Genetics</li> <li>Predict patients for high risks</li> <li>Predict part failures in manufacturing</li> <li>Predict loan defaulters</li> </ul>	<ul style="list-style-type: none"> <li>Can solve regression and classification problems with large datasets</li> <li>Highly scalable to any number of dimensions with quite acceptable performance</li> <li>A sort of „wisdom of the crowd“, tends to result in very high quality models out of the box, without extra need of testing and tuning</li> <li>Train each tree independently, using a random sample for the data, so the trained model is more robust than single decision tree, and less likely to over-fit</li> </ul>	<ul style="list-style-type: none"> <li>Can be slow to output predictions relative to other algorithms</li> <li>Not easy to understand predictions</li> <li>Learning can be slow</li> <li>Not possible to iteratively improve</li> </ul>

Algorithm	Applications	Advantages	Disadvantages
<b>Random Forest</b>		<ul style="list-style-type: none"> <li>• 2 parameters: number of trees and number of features to be selected at each node</li> <li>• Good for parallel or distributed computing</li> <li>• Lower classification error and better scores than decision trees</li> <li>• Perform as well as or better than SVMs, but far easier for humans to understand</li> <li>• Good with uneven data sets with missing variables</li> <li>• It helps identify significant features (feature importance)</li> <li>• Train faster than SVMs</li> </ul>	
<b>AdaBoost</b>	<p>Basic and more complex recognition problems</p> <ul style="list-style-type: none"> <li>• Biology</li> <li>• Computer vision</li> <li>• Speech processing</li> </ul>	<ul style="list-style-type: none"> <li>• Easily interpretable</li> <li>• Can be manually tuned</li> <li>• Can achieve similar classification results with much less tweaking of parameters or settings</li> <li>• Less sensitive to overfitting</li> </ul>	<ul style="list-style-type: none"> <li>• Can be sensitive to noisy data and outliers</li> </ul>
<b>Gradient Boosting</b>	<ul style="list-style-type: none"> <li>•</li> </ul>	<ul style="list-style-type: none"> <li>• High performing</li> <li>• Build trees one at a time, each new tree corrects some errors made by the previous trees, the model becomes even more expressive</li> <li>• Usually performs better than Random forests, but it is harder to tune</li> </ul>	<ul style="list-style-type: none"> <li>• A small change in the feature set or training set can create radical changes in the model</li> <li>• Not easy to understand predictions</li> <li>• 3 parameters – number of trees, depth, of trees, and learning rate: trees are generally shallow</li> <li>• The hyper parameters are hardest to tune and more prone to overfitting. RFs can almost work „out of box“</li> <li>• Training takes slower since trees are built sequentially</li> </ul>
<b>Linear Discriminant Analysis</b>	<ul style="list-style-type: none"> <li>•</li> </ul>	<ul style="list-style-type: none"> <li>• Compute the addition of multivariate distribution</li> <li>• Compute confidence intervals</li> </ul>	<ul style="list-style-type: none"> <li>• Require normal distribution</li> <li>• Not good for few categorical variables</li> <li>• Suffers multi-collinearity</li> </ul>



Algorithm	Applications	Advantages	Disadvantages	
Logistic Regression	<ul style="list-style-type: none"> <li>Predicting the customer churn</li> <li>Credit scoring</li> <li>Fraud detection</li> <li>Measuring effectiveness of marketing campaigns</li> </ul>	<ul style="list-style-type: none"> <li>Probabilistic interpretation</li> <li>Easy understand feature importance</li> <li>Provides confidence intervals</li> <li>Quickly update classification model to incorporate new data using stochastic gradient descent</li> <li>Output can be interpreted as probability, which can be used for ranking instead of classification</li> <li>Can easily „feature engineering“ most non-linear features into linear ones</li> <li>Robust to noise</li> <li>Can use L1 and L2 to avoid overfitting (and for feature selection)</li> <li>Efficient, can be distributed (ADMM)</li> <li>No need to worry about features being correlated</li> </ul>	<ul style="list-style-type: none"> <li>Does not handle the missing value of continuous variables</li> <li>Suffers multi-collinearity</li> <li>Sensitive to extreme values of continuous variables</li> <li>Sometimes too simple to capture complex relationships between variables</li> <li>Tendency for model to over-fit</li> <li>No distribution requirements</li> <li>Variable selection</li> <li>Features are expected to be roughly linear</li> <li>Cannot handle categorical (binary) variables well</li> <li>Underperform when there are multiple or non-linear decision boundaries</li> </ul> <p>Don't use when:</p> <ul style="list-style-type: none"> <li>If the variables are normally distributed and the categorical variables all have 5+ categories: use Linear discriminant analysis</li> <li>If the correlations are mostly nonlinear: use VSM</li> <li>If sparsity and multi-collinearity are a concern: Adaptive Lasso with Ridge (weights) + Lasso</li> </ul>	
		Lasso (L1)	<ul style="list-style-type: none"> <li>No distribution requirement</li> </ul>	<p>Lasso (L1)</p> <ul style="list-style-type: none"> <li>Variable selection</li> <li>Suffers multi-collinearity</li> </ul>
		Ridge (L2)	<ul style="list-style-type: none"> <li>No distribution requirement</li> </ul>	<p>Ridge (L2)</p> <ul style="list-style-type: none"> <li>No variable selection</li> <li>No suffer multi-collinearity</li> </ul>

Algorithm	Applications	Advantages	Disadvantages
<b>Neural Networks and Deep Learning</b>	<ul style="list-style-type: none"> <li>• Object recognition</li> <li>• Feature extractions</li> <li>• Computer vision</li> <li>• Speech recognition</li> </ul>	<ul style="list-style-type: none"> <li>• Used for classification and regression problems</li> <li>• Good prediction generally</li> <li>• Some tolerance to correlated inputs</li> <li>• Incorporating the predictive power of different combinations of inputs</li> <li>• Can handle extremely complex tasks – no other algorithms comes close to image recognition</li> <li>• Good to model the non-linear data with large number of input features</li> <li>• Widely used in industry</li> <li>• Many open source implementations</li> <li>• No need for complex feature engineering usually performed by hidden layers</li> <li>• Can learn extremely complex patterns due to hidden layers modelling intermediary representation of features</li> <li>• Good in image classification, video, audio, text</li> <li>• Performs better than any other algorithms</li> <li>• Can be easily updated with new data using batch propagation</li> <li>• Flexible as their architecture can be adapted for many types of problems</li> </ul>	<ul style="list-style-type: none"> <li>• Not a general purpose technique for classification as it needs very large amount of data to train</li> <li>• Not robust to outliers</li> <li>• Susceptible to irrelevant features</li> <li>• Difficult in dealing with big data with complex model</li> <li>• Only for numerical inputs, vectors with constant number of values, and datasets with non-missing data</li> <li>• „black box, the classification boundaries are hard to understand intuitively</li> <li>• Computational expensive</li> <li>• The trained model depends crucially on initial parameters</li> <li>• Difficult to troubleshoot when they don't work as expected</li> <li>• Not sure if they will generalize well to data not in training set</li> <li>• Multi-layer neural networks are usually hard to train, and require tuning lots of parameters what needs more expertise in architecture and hyper-parameters</li> <li>• Not probabilistic, unlike their more statistical or Bayesian counterparts. The continuous number output (eg. A score) can be difficult to translate that into a probability</li> <li>• Computationally and memory expensive</li> <li>• Usually outperform by ensembles for classical machine learning problems</li> </ul>