

Which machine learning algorithm should I use?

By [Hui Li](#) on [Subconscious Musings](#) | April 12, 2017

[Ad](#)

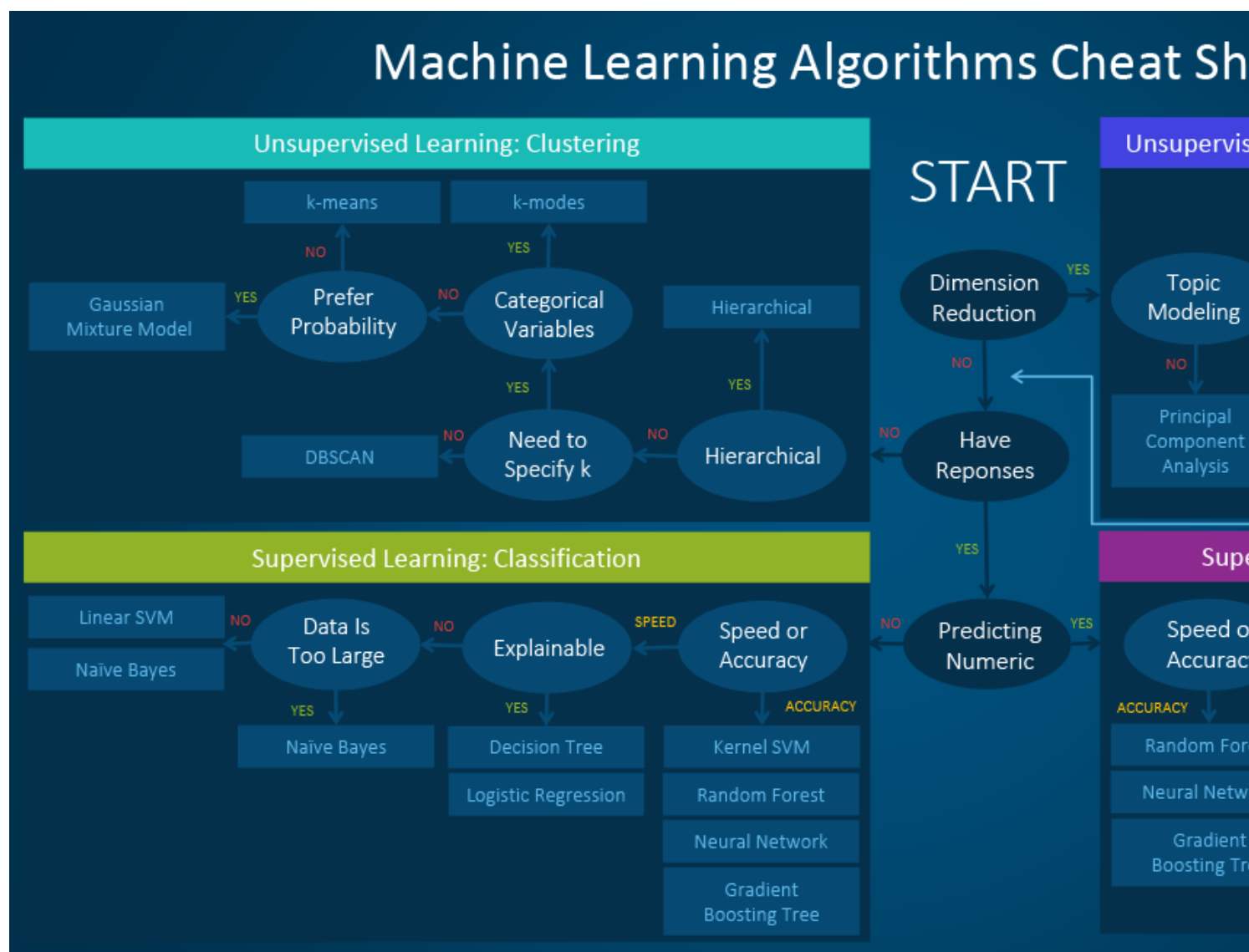
This resource is designed primarily for beginner to intermediate data scientists or analysts who are applying [machine learning](#) algorithms to address the problems of their interest.

A typical question asked by a beginner, when facing a wide variety of machine learning algorithms, The answer to the question varies depending on many factors, including:

- The size, quality, and nature of data.
- The available computational time.
- The urgency of the task.
- What you want to do with the data.

Even an experienced data scientist cannot tell which algorithm will perform the best before trying and advocating a one and done approach, but we do hope to provide some guidance on which algorithm clear factors.

The machine learning algorithm cheat sheet



The **machine learning algorithm cheat sheet** helps you to choose from a variety of machine learning algorithms. This article walks you through the process of how to choose an appropriate algorithm for your specific problems.

Since the cheat sheet is designed for beginner data scientists and analysts, we will make some simplifications about the algorithms.

The algorithms recommended here result from compiled feedback and tips from several data science experts and developers. There are several issues on which we have not reached an agreement and we will highlight the commonality and reconcile the difference.

Additional algorithms will be added in later as our library grows to encompass a more complete set

How to use the cheat sheet

Read the path and algorithm labels on the chart as "If *<path label>* then use *<algorithm>*." For exam

- If you want to perform dimension reduction then use principal component analysis.
- If you need a numeric prediction quickly, use decision trees or logistic regression.
- If you need a hierarchical result, use hierarchical clustering.

Sometimes more than one branch will apply, and other times none of them will be a perfect match. paths are intended to be rule-of-thumb recommendations, so some of the recommendations are not talked with said that the only sure way to find the very best algorithm is to try all of them.

Types of machine learning algorithms

This section provides an overview of the most popular types of machine learning. If you're familiar move on to discussing specific algorithms, you can skip this section and go to "When to use specific

Supervised learning

Supervised learning algorithms make predictions based on a set of examples. For example, historical the future prices. With supervised learning, you have an input variable that consists of labeled training variable. You use an algorithm to analyze the training data to learn the function that maps the input function maps new, unknown examples by generalizing from the training data to anticipate results in

- **Classification:** When the data are being used to predict a categorical variable, supervised learning. This is the case when assigning a label or indicator, either dog or cat to an image. When there are two categories, the problems are called binary classification. When there are more than two categories, the problems are called multi-classification.
- **Regression:** When predicting continuous values, the problems become a regression problem.
- **Forecasting:** This is the process of making predictions about the future based on the past and is commonly used to analyze trends. A common example might be estimation of the next year's sales based on current year and previous years.

Semi-supervised learning

The challenge with supervised learning is that labeling data can be expensive and time consuming. Adding unlabeled examples to enhance supervised learning. Because the machine is not fully supervised, it is semi-supervised. With semi-supervised learning, you use unlabeled examples with a small amount of labeled data to improve learning accuracy.

Unsupervised learning

When performing unsupervised learning, the machine is presented with totally unlabeled data. It is the machine's job to find patterns that underlies the data, such as a clustering structure, a low-dimensional manifold, or a sparse representation.

- **Clustering:** Grouping a set of data examples so that examples in one group (or one cluster) share some criteria) than those in other groups. This is often used to segment the whole dataset into groups. A model is performed in each group to help users to find intrinsic patterns.
- **Dimension reduction:** Reducing the number of variables under consideration. In many applications, high dimensional features and some features are redundant or irrelevant to the task. Reducing the number of features reveals the true, latent relationship.

Reinforcement learning

Reinforcement learning analyzes and optimizes the behavior of an agent based on the feedback from different scenarios to discover which actions yield the greatest reward, rather than being told which actions are correct. The delayed reward distinguishes reinforcement learning from other techniques.

Considerations when choosing an algorithm

When choosing an algorithm, always take these aspects into account: accuracy, training time and cost. Accuracy is the most important, while beginners tend to focus on algorithms they know best.

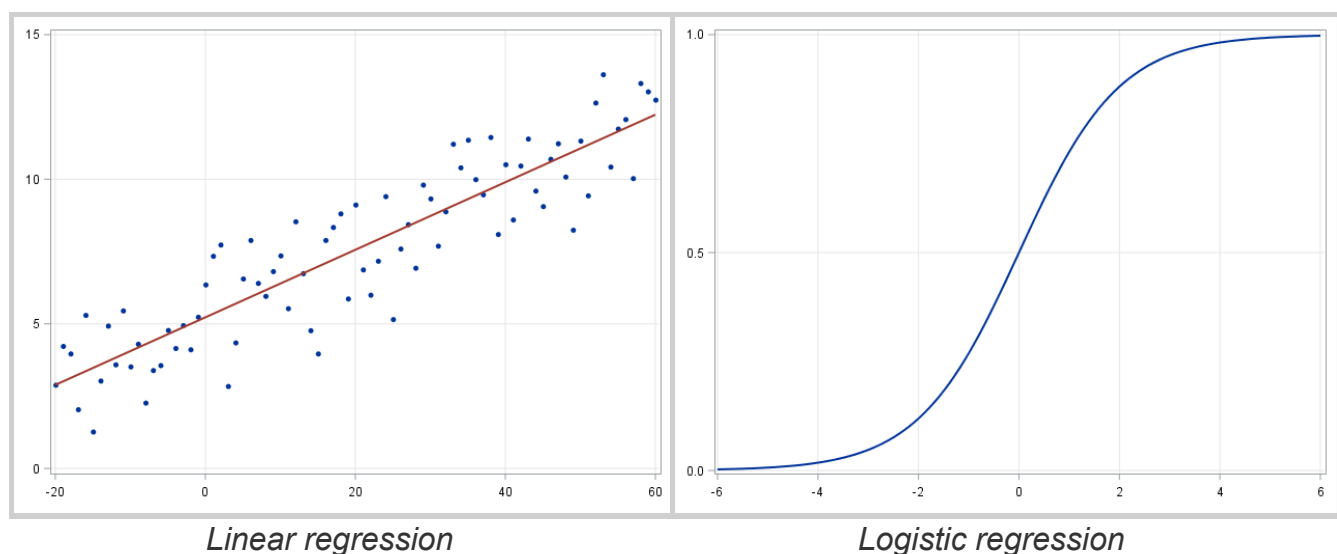
When presented with a dataset, the first thing to consider is how to obtain results, no matter what the results are. Beginners tend to choose algorithms that are easy to implement and can obtain results quickly. This is the first step in the process. Once you obtain some results and become familiar with the data, you can then try more sophisticated algorithms to strengthen your understanding of the data, hence further improving the results.

Even in this stage, the best algorithms might not be the methods that have achieved the highest performance. Usually, this usually requires careful tuning and extensive training to obtain its best achievable performance.

When to use specific algorithms

Looking more closely at individual algorithms can help you understand what they provide and how they provide more details and give additional tips for when to use specific algorithms, in alignment with the goals of the problem.

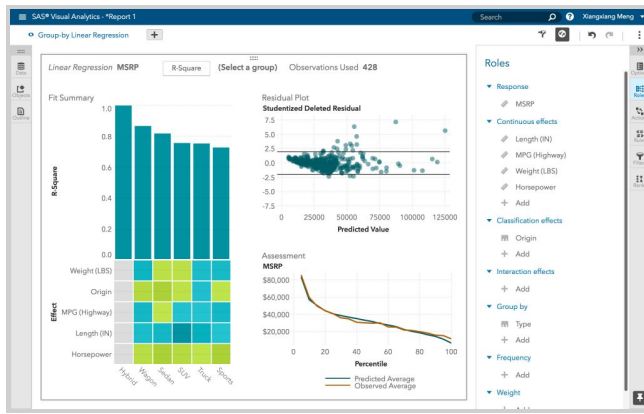
Linear regression and Logistic regression



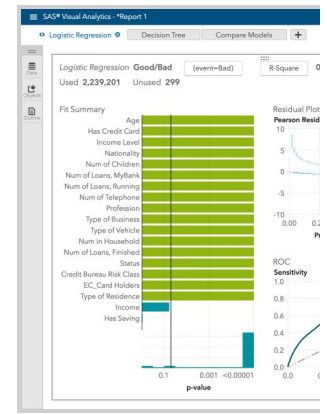
Linear regression is an approach for modeling the relationship between a continuous dependent variable y and predictors X . The relationship between y and X can be linearly modeled as $y = \beta^T X + \epsilon$. Given the data, the parameter vector β can be learnt.

If the dependent variable is not continuous but categorical, linear regression can be transformed to use a link function. Logistic regression is a simple, fast yet powerful classification algorithm. Here we discuss a dependent variable y only takes binary values $\{y_i \in (-1, 1)\}_{i=1}^N$ (it which can be easily extended to multi-class problems).

In logistic regression we use a different hypothesis class to try to predict the probability that a given instance belongs to the "1" class versus the probability that it belongs to the "-1" class. Specifically, we will try to learn a function of x and $p(y_i = -1|x_i) = 1 - \sigma(\beta^T x_i)$. Here $\sigma(x) = \frac{1}{1 + \exp(-x)}$ is a sigmoid function. Given the training data, the parameter vector β can be learnt by maximizing the log-likelihood of β given the data set.



Group By Linear Regression



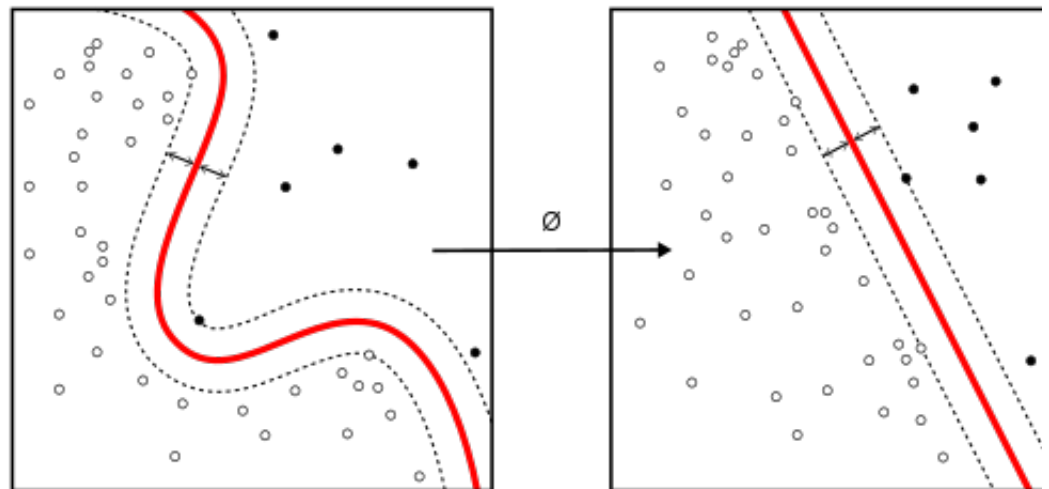
Logistic Regression i

Linear SVM and kernel SVM

Kernel tricks are used to map a non-linearly separable functions into a higher dimension linearly separable function. The SVM training algorithm finds the classifier represented by the normal vector w and bias b . The hyperplane (boundary) separates different classes by as wide a margin as possible. The problem is a constrained optimization problem:

$$\begin{aligned} & \underset{w}{\text{minimize}} && ||w|| \\ & \text{subject to} && y_i(w^T X_i - b) \geq 1, \quad i = 1, \dots, n. \end{aligned}$$

A support vector machine (SVM) training algorithm finds the classifier represented by the normal vector v . This hyperplane (boundary) separates different classes by as wide a margin as possible. The problem is a constrained optimization problem:

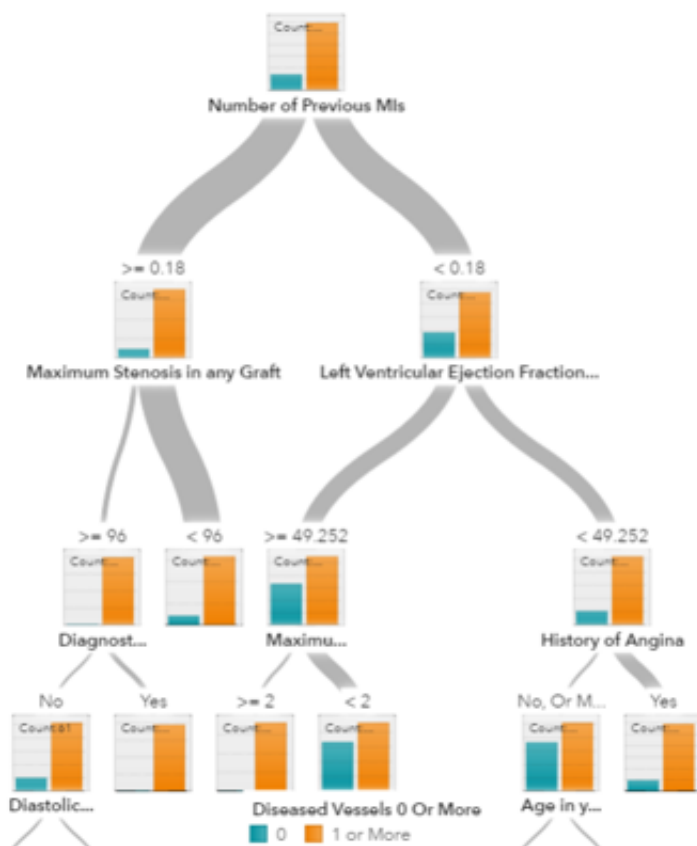


Kernel tricks are used to map a non-linearly separable functions into a higher dimension linearly separable function.

When the classes are not linearly separable, a kernel trick can be used to map a non-linearly separable space to a high dimension linearly separable space.

When most dependent variables are numeric, logistic regression and SVM should be the first try for easy to implement, their parameters easy to tune, and the performances are also pretty good. So for beginners.

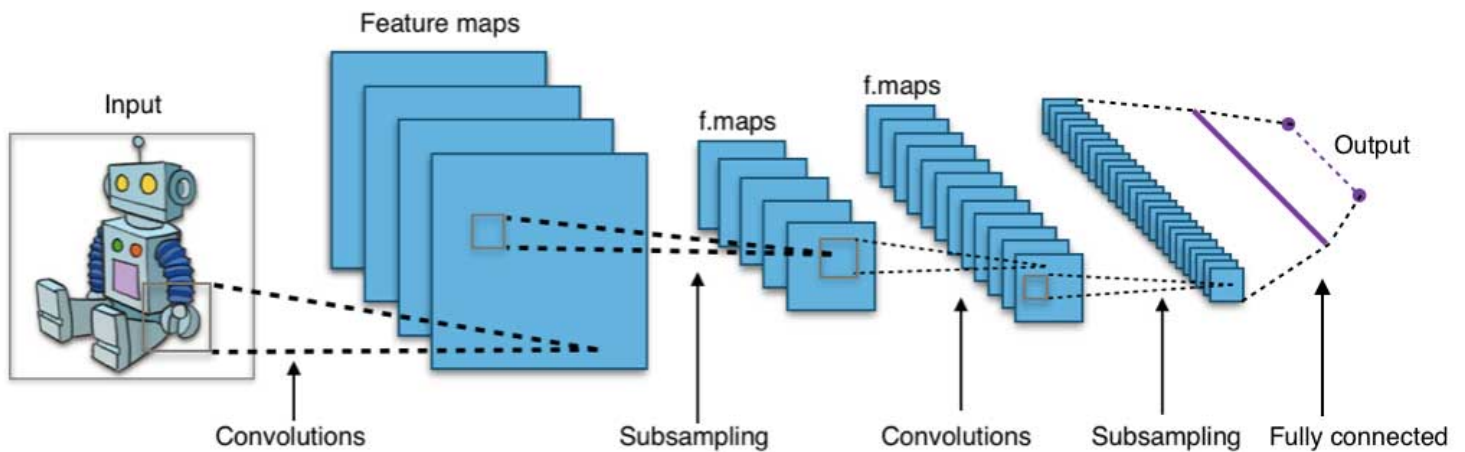
Trees and ensemble trees



A decision tree for prediction model

Decision trees, random forest and gradient boosting are all algorithms based on decision trees. They all do the same thing – subdivide the feature space into regions with mostly the same class. However, they tend to over fit data when we exhaust the branches and create too many regions. Random Forest and gradient boosting are two popular ways to use tree algorithms to achieve good results and avoid the over-fitting problem.

Neural networks and deep learning



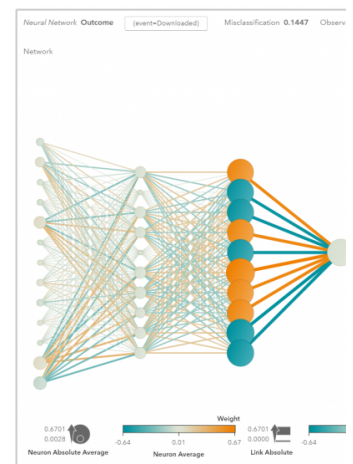
A convolution neural network architecture (image source: wikipedia creative commons)

Neural networks flourished in the mid-1980s due to their parallel and distributed processing ability. They were impeded by the ineffectiveness of the back-propagation training algorithm that is widely used to train neural networks. Support vector machines (SVM) and other simpler models, which can be easily trained for many problems, gradually replaced neural networks in machine learning.

In recent years, new and improved training techniques such as unsupervised pre-training and layer-wise training have led to a resurgence of interest in neural networks. Increasingly powerful computational capabilities, such as GPU and massively parallel processing (MPP), have also spurred the revived adoption of neural networks. The use of neural networks has given rise to the invention of models with thousands of layers.

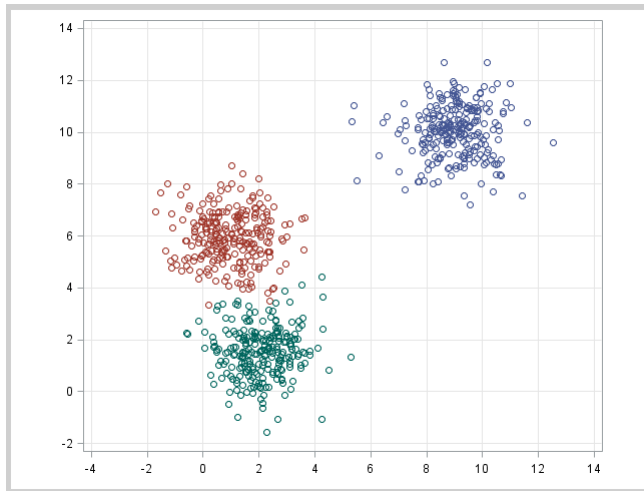
In other words, shallow neural networks have evolved into **deep learning** neural networks. Deep neural networks have been very successful for supervised learning. When used for speech and image recognition, deep learning performs as well as, or even better than, humans. Applied to unsupervised learning tasks, such as **feature extraction**, deep learning also extracts features from raw images or speech with much less human intervention.

A neural network consists of three parts: input layer, hidden layers and output layer. The training samples define the input and output layers. When the output layer is a categorical variable, then the neural network is a way to address classification problems. When the output layer is a continuous variable, then the network can be used to do regression. When the output layer is the same as the input layer, the network can be used to extract intrinsic features. The number of hidden layers defines the model's capacity.

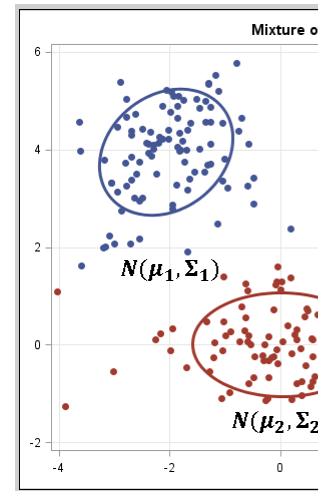


A neural network

k-means/k-modes, GMM (Gaussian mixture model) clustering



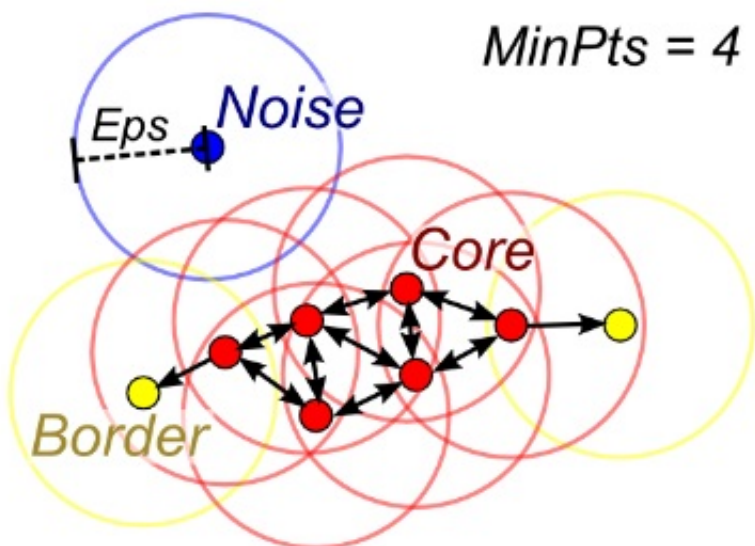
K Means Clustering



Gaussian Mixture Model

Kmeans/k-modes, GMM clustering aims to partition n observations into k clusters. K-means define to be and only to be associated to one cluster. GMM, however define a soft assignment for each sample with a probability to be associated with each cluster. Both algorithms are simple and fast enough for clustering when k is given.

DBSCAN

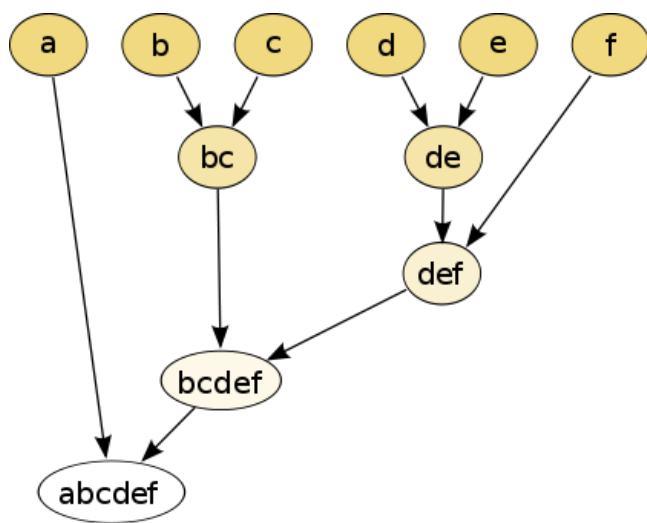


A DBSCAN illustration (image source: [Wikipedia](#))

When the number of clusters k is not given, DBSCAN (density-based spatial clustering) can be used.

density diffusion.

Hierarchical clustering



Hierarchical partitions can be visualized using a tree structure (a dendrogram). It does not need the and the partitions can be viewed at different levels of granularities (i.e., can refine/coarsen clusters

PCA, SVD and LDA

We generally do not want to feed a large number of features directly into a machine learning algorithm irrelevant or the “intrinsic” dimensionality may be smaller than the number of features. Principal component analysis (PCA), singular value decomposition (SVD), and latent Dirichlet allocation (LDA) all can be used to perform dimensionality reduction.

PCA is an unsupervised clustering method which maps the original data space into a lower dimensional space while preserving as much information as possible. The PCA basically finds a subspace that most preserves the data variance by the dominant eigenvectors of the data’s covariance matrix.

The SVD is related to PCA in the sense that SVD of the centered data matrix (features versus samples) yields singular vectors that define the same subspace as found by PCA. However, SVD is a more versatile tool than PCA. Things that PCA may not do. For example, the SVD of a user-versus-movie matrix is able to extract user and movie profiles which can be used in a recommendation system. In addition, SVD is also widely used as a tool for latent semantic analysis, in natural language processing (NLP).

A related technique in NLP is latent Dirichlet allocation (LDA). LDA is probabilistic topic model and topics in a similar way as a Gaussian mixture model (GMM) decomposes continuous data into Gaussian components. In the GMM, an LDA models discrete data (words in documents) and it constrains that the topics are Dirichlet distribution.

Conclusions

This is the work flow which is easy to follow. The takeaway messages when trying to solve a new problem are:

- Define the problem. What problems do you want to solve?
- Start simple. Be familiar with the data and the baseline results.
- Then try something more complicated.

[SAS Visual Data Mining and Machine Learning](#) provides a good platform for beginners to learn machine learning methods to their problems. [Sign up for a free trial today!](#)

Tags

machine learning algorithms

machine learning

data science basics

data science

regression

Share



ABOUT AUTHOR

Hui Li

Principal Staff Scientist, Data Science

Dr. Hui Li is a Principal Staff Scientist of Data Science Technologies at SAS. Her current research focuses on Cognitive Computing and SAS recommendation systems in SAS Viya. She received her

degree in Electrical and Computer Engineering from Duke University. Before joining SA as a research scientist and at Signal Innovation Group, Inc. as a research engineer. He machine learning for big, heterogeneous data, collaborative filtering recommendations, reinforcement learning.

8 COMMENTS

Daymond Ling on April 12, 2017 7:58 pm

Thank you for the cheat-sheet, it provides a nice taxonomy for people to understand the relationship between different machine learning algorithms. I will use it in my machine learning class to help students round out their world view.

Hui Li on April 17, 2017 9:54 am

Thank Daymond.

Let us know if you have any questions when teaching the students using the information.

[Hector Alvaro Rojas](#) on April 21, 2017 11:12 am

This is a great cheat-sheet to understand and remember the relationship between the most important machine learning algorithms. I have not seen something similar like this published online yet.

I think it could be nice to incorporate the "cost" variable, the principal's reasons why each set of examples of applications for each one. I know that this suggestion means a lot of work and so I will not expect a quick response. Anyway, it could be a nice new project to be done, don't you think so?

Congratulations for the work already done anyway!

Hui Li on April 24, 2017 11:29 am

Thanks, Hector. Incorporating the "cost" variable is a pretty wider area in machine learning.

considered as a subfield of reinforcement learning -- based on the cost (reward), the a
he/she wants to take. I considered this problem for a while and haven't found a good e
I have a time, I will write a blog specifically for the reinforcement learning.

charles on April 24, 2017 9:54 am

An excellent blog. Thank you

Hui Li on April 24, 2017 11:30 am

Thank you.

Don Maclean on April 25, 2017 11:34 am

Excellent summary but I think the target audience is a few steps beyond "beginner". I showed
study machine learning, and they were overwhelmed.

Anastassia Dr Lauterbach on April 26, 2017 6:31 am

Great blog, thank you. I will use it when talking to non tech companies about starting doing M
