

# The Most Important Machine Learn Algorithms

Pros and cons of the top 11 algorithms every machine learning engineer should know

Posted by Tobias ([https://twitter.com/semanti\\_ca](https://twitter.com/semanti_ca)) on 03-09-2018

There are plenty of machine learning algorithms. Some of them are v useful for solving some very specific problems. Others could be appl solve a wide range of tasks.

In this post, we list the algorithms most frequently used by machine engineers at [semanti.ca](https://semanti.ca) (<https://semanti.ca>) as well as those that oft [Kaggle](https://www.kaggle.com) (<https://www.kaggle.com>) competitions.

---

## Classification and Regression

### Gradient Boosting Machine

Gradient Boosting Machine is currently (as of 2018) the most popular algorithm among Kagglers and most winning teams used Gradient B Machine in their solutions.

**Pros:**

- Can handle huge datasets (millions of examples and millions of dimensions);
- Extremely accurate;
- Can be used for both classification and regression tasks;
- Lots of flexibility with the choice of loss functions: can be tailored to the characteristics of the problem.

**Cons:**

- Can be slow at training, since the trees are built sequentially;
- Model is a black box (as for all ensemble methods);
- Prone to overfitting (however hyperparameter tuning helps).

**Random Forest**

Random Forests are still frequently used by the machine learning practitioners. They provide good accuracy and speed of training and handle big datasets.

**Pros:**

- Much easier to tune than Gradient Boosting Machine (two hyperparameters vs three);
- Almost always perform as well as or better than SVMs;
- Deal well with uneven datasets that have missing variables;
- Rarely overfits.

**Cons:**

- Slow at prediction time so less appropriate for high-speed data

processing;

- If data includes categorical variables with different numbers of levels, Random Forest is biased in favor of those with more levels. Therefore, the variable importance scores from Random Forest are not reliable.

## Support Vector Machine

Support Vector Machines is a mature and well-studied machine learning algorithm, with a solid theoretical foundation. It supports kernels, so it can handle non-linearly separable classification problems.

### Pros:

- Handles both classification and regression;
- Supports very high-dimensional and sparse data (hundreds of thousands of dimensions); highly effective in text classification where high-dimensionality and sparsity are frequently observed in data.
- High accuracy, good theoretical guarantees regarding overfitting;
- Handles very well outliers and noise.
- Extremely fast at prediction time;
- Generalize well from few data.

### Cons:

- Cannot be parallelized, the whole dataset has to fit in memory;
- Not good at handling very big dataset (more than a hundred of thousands of examples);
- Sometimes hard to tune due to a wide range of possible hyperparameter values.

## LSTM Neural Network

LSTM Neural Network is a variant of the Recurrent Neural Network.

units effectively solve the problem of vanishing gradient. They are typically used for sequential classification problems: text labeling, speech recognition. The Encoder-Decoder architecture of LSTM networks allows building machine translation systems.

**Pros:**

- Classify relatively long (up to 40-50 tokens) sequences well;
- Highly accurate even with a couple of hidden layers.

**Cons:**

- Very slow at training as cannot be parallelized between multiple GPUs;
- Relatively slow at classification, since it classifies one token at a time (especially the Encoder-Decoder architecture);
- Black boxes without strong theoretical foundations.

## Convolutional Neural Network

As of 2018, one cannot imagine the image processing without Convolutional Neural Networks. Such state-of-the-art CNN architectures as AlexNet, GoogLeNet, and ResNet are all variants of a Convolutional Neural Network.

**Pros:**

- Very fast to train, as can be parallelized between multiple GPUs;
- Learn low-level and high-level features so can be used as feature extractors for other ML algorithms;
- Very accurate at image classification and many other tasks, including natural language processing.

**Cons:**

- Require a lot of training data (however, such techniques as **trans learning** can reduce the need for data);
- Sensible to parameter initialization and choice of hyperparameters
- Slow to train on CPU;
- Black boxes without strong theoretical foundations (as of 2018).

## K-Nearest Neighbors

K-Nearest Neighbors is one of the oldest and simplest algorithms of learning. In fact, there's not so much learning happening: the prediction is given as the majority (classification) or the average (regression) of the labels of the  $k$  nearest neighbors in the close neighborhood of the input, unannotated example. KNN is very accurate recently, however, it was for a long time considered as a slow algorithm. However, modern libraries made this algorithm very fast. Today it can easily compete with the state-of-the-art techniques.

### Pros:

- No assumptions about data: works well for non-linearly separable data
- Simple to implement;
- Flexible to feature/distance choices;
- Naturally handles multi-class cases;
- Does well in practice with enough representative data
- Can be used for both classification and regression.

### Cons:

- Search in a large space of examples to find nearest neighbors can be slow (mostly resolved in modern libraries);
- With a lot of training examples, the model can require a lot of memory
- Sensitive to irrelevant features and the scale of the data.

# Clustering and Dimensionality Reduction

## K-Means

The K-means algorithm remains one of the most used clustering algorithms despite its simplicity.

### Pros:

- Simple to use;
- Fast (when used with an appropriate data structure to store examples that preserve spatial locality);
- Results are simple to understand by a human.

### Cons:

- K-means only works well when the shape of clusters are hyper-spherical or when clusters are far from one another;
- Different runs of the algorithm will most probably yield in different clusterings of the same data;
- Requires the number of clusters to be known in advance.

## Expectation Maximization

Expectation Maximization, as well as Latent Dirichlet Allocation (below), are both examples of the so-called "soft" clustering. They assign one example to multiple clusters with a probability of membership. Expectation Maximization learns Gaussian Mixture Models which solves the problem that K-Means: the clusters can be non-spherical.

### Pros:

- Simplicity and ease of implementation;
- Solutions to the M-steps often exist in the closed form;
- Can often be easily parallelized and its memory requirements are modest compared to other methods;
- The algorithm is numerically very stable;
- Soft-membership in clusters is often desirable;
- Can be naturally used for finding outliers.

**Cons:**

- Slow linear convergence;
- Can converge to local optima;
- As in K-Means, one needs to know the number of clusters;
- One needs to carefully choose the generative model.

**DBSCAN**

DBSCAN (for *Density-Based Spatial Clustering of Applications with Noise*) is another popular clustering algorithm. Its difference from K-Means lies in the fact that DBSCAN doesn't use the notion of cluster centroids. It's density-based, so the number of clusters depends on the data itself. Clusters can have absolutely any form: for DBSCAN a cluster is a subset of points that form a dense "cloud".

**Pros:**

- Doesn't need the number of clusters as input;
- Clusters can have arbitrary shapes;
- Robust to noise;
- Deterministic.

**Cons:**

- Requires connected regions of sufficiently high density;
- The maximal distance between two points that belong to one cluster to be given as input;
- Sensitive to hyperparameters;
- Sampling affects density measures.
- Datasets with varying densities are problematic.

## Latent Dirichlet Allocation

LDA is one of the most important unsupervised learning algorithms. It has been successfully applied to text analysis (topic modeling), social network analysis (community overlaps), content recommendation, genetic population analysis and many others.

In the nutshell, LDA takes a collection of documents, the number of topics and returns a distribution of topics over documents and a distribution of words over topics.

LDA can be used as a soft-clustering method: every topic can be seen as a cluster and each document can have several topics with different probabilities, thus can belong to several clusters.

### Pros:

- Excellent empirical results;
- Mitigates overfitting well;

### Cons:

- High computational complexity;
- Nondeterministic (different runs can give different results);
- Doesn't work well on short documents.



## UMAP

One important aspect of the machine learning practice is data visual and dimensionality reduction. Dimensionality reduction can be done visualize data (the human can only see up to three-dimensional data), make learning tractable and/or more accurate.

For a long Principal Component Analysis ([https://en.wikipedia.org/wiki/Principal\\_component\\_analysis](https://en.wikipedia.org/wiki/Principal_component_analysis)) or PCA was the most popular dimensionality reduction and t-SNE ([https://en.wikipedia.org/wiki/distributed\\_stochastic\\_neighbor\\_embedding](https://en.wikipedia.org/wiki/distributed_stochastic_neighbor_embedding)) was used to visualize. However, in 2018 the new algorithm, UMAP (<https://github.com/lmc/umap>) (for *Uniform Manifold Approximation and Projection*) was proposed and efficiently implemented. UMAP combines the advantage of both t-SNE: it can be used as a dimensionality reduction techniques (t-SNE) and has visualization properties similar or surpassing those of t-SNE in many tasks.

### Pros:

- Can be used for both dimensionality reduction and visualization
- Has very fast implementation in multiple programming language including Python.

### Cons:

- Has hyperparameters one has to tune to find a good visualization


These were the 11 most important machine learning algorithms, according to [semanti.ca](https://www.semanti.ca) (<https://www.semanti.ca>) engineers and scientists. Did we miss something important? Please, let us know and we will improve this to


---


Read our previous post "[Data Science Work Productivity Tips & Tricks](https://semanti.ca/blog/?data-science-work-productivity)" or [subscribe](https://semanti.ca/blog/?feed) to our RSS feed.


---


Like it? Share it!

 (<https://twitter.com/home?status=https://www.semanti.ca/blog/?the-most-important-machine-learning-algorithms>)

 (<https://www.facebook.com/sharer/sharer.php?u=https://www.semanti.ca/blog/?the-most-important-machine-learning-algorithms>)

 (<https://plus.google.com/share?url=https://www.semanti.ca/blog/?the-most-important-machine-learning-algorithms>)


 (<https://www.linkedin.com/shareArticle?mini=true&url=https://www.semanti.ca/blog/?the-most-important-machine-learning-algorithms>)

 (<http://www.reddit.com/submit?url=https://www.semanti.ca/blog/?the-most-important-machine-learning-algorithms>)



**Home (/) · Blog (/blog/) · Pricing (/pricing/) · About (/about/) ·  
Contact (/contacts/)**

 **+1 646 9050250**

 **human@semanti.ca (mailto:human@semanti.ca)**

**semanti.ca**

We build AI-powered APIs for automated extraction of data from web pages. No programming required.

**in**

**(http://www.linkedin.com**

  **/com** **(http://github.com**

**(http://github.com/semanti.ca**  
**(http://github.com/semanti.ca)**  
**/+semanti.ca/blog/feed)**