

A Brief History of Machine Learning

Norbert Cristian BERECHKI

April 2018

Abstract

One of the most actual and interesting topics of the new millennium is how to make machines learn and help our society. In this paper we will present a brief history of machine learning and how it came to be the hot topic that it is today. We will also put much emphasis on Deep Learning and its potential.

1 Motivation

In today's world, Machine Learning is an indispensable part of our life. Many of the technologies that are built in this time are based on Machine Learning. Moreover, it creates new jobs and automates many simple-to-do tasks. So, it looks more and more that people will stop doing mundane tasks and start concentrating on complex ones. Consequently, it is natural to investigate how it came to the capabilities that ML-tools provide and the direction the domain is heading towards.

This paper consists mostly of 2 chapters. The first one presents a more general view of how ML tools developed and the most important milestones that occurred. The second chapter puts emphasis upon Deep Learning(DL), and why, even though neural networks were invented a long time ago, it became so popular only in the past two to three decades.

2 A broader view of ML

This section comprises of more subsections that describe parts of the Machine Learning.

2.1 Statistical Learning

We begin with this topic to emphasize its importance to the field in discussion. It is fair to say that without statistics there would be no Machine Learning nor there would be any advanced AI algorithm that we have today. However, we will only concentrate upon the part which includes *learning*.

Interesting enough, the first tabulating machine made by Herman Hollerith was created for the government to offer support for some statistical computations related to the population (1890). Statistical Learning is preoccupied with the employing of statistics into modeling data. Most of the state-of-the-art techniques use at some degree statistics, but in this part we will deal only with the ones that make much use of them.

The revolutionary work made by Lagrange and Gauss (in the 19th century) lead to the discovery of the *linear least square method* which was a sort of **regression** method. By the time of approximately 1950, more and more linear classifications and regression methods were developed. However only after 1970 most people started using these methods in practical applications such as predictive analysis, anomaly detection and classification. This is mostly due to the rise of the more powerful computers which could support complex computations and had much more memory and to the rise of non-linear classification and regression techniques. Not to forget the importance of the creation of the R language which boosted the field of statistics and provided very good (and optimized) tools for the development of statistical models and inference models. [5]

2.2 Unsupervised Learning

Methods provided by this topic are widely used in subfields such as Data Science, Data Analysis and Business Intelligence. Unsupervised learning implies finding structure and ordering data where there is unlabeled data. So you will not be able to know for certain how well your model works.

The most trivial example for this is **clustering**. The apparition of clustering techniques (including here the most known clustering algorithm: **k-means**) is a little bit debatable many people referring to it as clustering or using k-means but not publishing it. Although, the first one to publish work related to the k-means algorithm is Stuart Lloyd in 1957. In the early 2000s the rise of data mining led to the invention of yet another clustering algorithm, namely the hierarchical clustering algorithm. Although it is not as used as other, more advanced algorithms, the hierarchical clustering algo had much popularity until approximately 2010s.[1]

Other areas where unsupervised learning shows its prowess are the fields of Anomaly Detection, Structural Prediction and feature reduction (mainly PCA). The PCA method was created in 1901 by Karl Pearson (as [2] states) and since it has been used in many fields. However, the years 2000-2013 have known a large increase in papers upon PCA, mainly the year of 2013 where most of the papers on PCA were published, as stated in [1].

2.3 Supervised Learning

When taling about supervised Learning, it is usually implied that for a problem of regression or classification we must train a model, using some already known examples, to give us the desired output.

One of the most common set of techniques is the one regarding **regression analysis**. Here [3] it states that: 'Regression analysis is a set of statistical processes for estimating the relationships among variables. More specifically, regression analysis helps one understand how the typical value of the dependent variable (or 'criterion variable') changes when any one of the independent variables is varied, while the other independent variables are held fixed.'

In this article [6] there is much being said about the early years of regression. We find out that Karl Pearson, who was at the time the first professor of Statistics in Britain chose Galton over Auguste Bravaris as the first one to offer a reasonable definition to regression. He

accounts that the idea came to Galton when he was sheltering from a rainstorm. At that time Galton worked on eugenics and tried to see how geniuses can be bred. He was interested in the behaviour of outliers and if they will regress back to the mean. Methods for regression evolved over time by the works of Pearson and later Fisher. By the 1950s, regression was used in many branches of our society including economics.

3 Deep Learning

3.1 Artificial Neural Networks

The roots of neural networks lies in the Hebbian Learning Theory which was proposed in the late 1940s. After that period, researchers tried to bring the hebbian learning model into the computational world. The first major breakthrough was made by Rosenblatt, namely the discovery of the **Perceptron** which could classify linearly separable classes. For instance, take the logical AND operation. You can draw a line to separate the output(0's and 1's) made by the AND function when inserting a 2-input value (Figure 1). In spite of this big breakthrough, much research did not follow until later on. This was due to mostly two things: firstly, the perceptron could not classify more complex functions such as the logical XOR, secondly, the computing power was weak compared to these days.

A big step futher was taken in 1975 by Werbos when he proposed a variant of the backpropagation algorithm, which was simply gradient descent applied recursively on each node, now having a network of perceptrons. Because of this, the network was able to learn the XOR function. G. Hinton had also worked to demonstrate the generalization of backpropagation on multi-layer networks. That was the spark that triggered the Deep Learning Revolution. To give some clarification, by Deep Learning we understand the methods to train "deep models", such as neural networks with many layers (more than 10).

3.2 Deeper Networks

This idea to stack layers in order to get deeper network stroke a revolution in Machine Learning. However, an important discovery was the Neocognitron (created by Kunihiro Fukushima in 1980) which was mostly a neural network that was classifying handwritten digits. This had a great impact on the work of Yann LeCun, called the "Father of Convolutional Nets". One of the first important works that LeCun did, was during his PhD. and was on the backpropagation algorithm. He continued his work later on at Bell Labs (where, from his saying, he got his own cluster to experiment on). There, he created the first Convolutional Neural Network called LeNet 1, which did handwritten digit classification. He improved his model until iteration 5 (LeNet 5). After that he went on working on the famous DjVu image compression system. His groundbreaking work enabled researchers to invest and develop a stil-growing set of tools related to Convolutional Neural Networks.

Other architectures began to rise, including: the Cresceptron(which had the job of detecting objects in a 3D space) and the SOFT (a weightless self-organizing feature extraction neural network). Although there were many issues related to deep learning, research stil continued. The issues were connected to the time of the training phase(for most models days or weeks were needed to train on even a small amount of data) and to the lack in labeled data. So in this time

other methods such as SVMs(Support Vector Machines) and Gabor-Filter-based models were used for computer vision tasks and Hidden Markov Models for sequence modeling tasks.

4 Conclusion and Fun Facts

[4]

References

- [1] URL: <https://f1000research.com/articles/6-2012>.
- [2] URL: https://en.wikipedia.org/wiki/Principal_component_analysis.
- [3] URL: https://en.wikipedia.org/wiki/Regression_analysis.
- [4] URL: <http://www.mellanox.com/interconnected-planet/related-docs/machine-learning-infographic.pdf>.
- [5] Gareth James et al. *An Introduction to Statistical Learning*. Springer, 2013. ISBN: 978-1-4614-7137-0.
- [6] T.J. Barnes. "A history of regression: actors, networks, machines, and numbers". In: *Environment and Planning* 30 (1998), pp. 203–223. DOI: <http://journals.sagepub.com/doi/pdf/10.1068/a300203>.