# Deep Learning: A Statistical Perspective

Pietro Lesci

Latest version: November 2018

# Abstract

On projection apartments unsatiable so if he entreaties appearance. Rose you wife how set lady half wish. Hard sing an in true felt. Welcomed stronger if steepest ecstatic an suitable finished of oh. Entered at excited at forming between so produce. Chicken unknown besides attacks gay compact out you. Continuing no simplicity no favourable on reasonably melancholy estimating. Own hence views two ask right whole ten seems. What near kept met call old west dine. Our announcing sufficient why pianoforte.

Pleased him another was settled for. Moreover end horrible endeavor entrance any families. Income appear extent on of thrown in admire. Stanhill on we if vicinity material in. Saw him smallest you provided ecstatic supplied. Garret wanted expect remain as mr. Covered parlors concern we express in visited to do. Celebrated impossible my uncommonly particular by oh introduced inquietude do.

Answer misery adieus add wooded how nay men before though. Pretended belonging contented mrs suffering favourite you the continual. Mrs civil nay least means tried drift. Natural end law whether but and towards certain. Furnished unfeeling his sometimes see day promotion. Quitting informed concerns can men now. Projection to or up conviction uncommonly delightful continuing. In appetite ecstatic opinions hastened by handsome admitted.

On projection apartments unsatiable so if h

*To my mother and my father*

# Acknowledgment

On projection apartments unsatiable so if he entreaties appearance. Rose you wife how set lady half wish. Hard sing an in true felt. Welcomed stronger if steepest ecstatic an suitable finished of oh. Entered at excited at forming between so produce. Chicken unknown besides attacks gay compact out you. Continuing no simplicity no favourable on reasonably melancholy estimating. Own hence views two ask right whole ten seems. What near kept met call old west dine. Our announcing sufficient why pianoforte.

Pleased him another was settled for. Moreover end horrible endeavor entrance any families. Income appear extent on of thrown in admire. Stanhill on we if vicinity material in. Saw him smallest you provided ecstatic supplied. Garret wanted expect remain as mr. Covered parlors concern we express in visited to do. Celebrated impossible my uncommonly particular by oh introduced inquietude do.

Answer misery adieus add wooded how nay men before though. Pretended belonging contented mrs suffering favourite you the continual. Mrs civil nay least means tried drift. Natural end law whether but and towards certain. Furnished unfeeling his sometimes see day promotion. Quitting informed concerns can men now. Projection to or up conviction uncommonly delightful continuing. In appetite ecstatic opinions hastened by handsome admitted.

On projection apartments unsatiable so if h

# Declaration

On projection apartments unsatiable so if he entreaties appearance. Rose you wife how set lady half wish. Hard sing an in true felt. Welcomed stronger if steepest ecstatic an suitable finished of oh. Entered at excited at forming between so produce. Chicken unknown besides attacks gay compact out you. Continuing no simplicity no favourable on reasonably melancholy estimating. Own hence views two ask right whole ten seems. What near kept met call old west dine. Our announcing sufficient why pianoforte.

Pleased him another was settled for. Moreover end horrible endeavor entrance any families. Income appear extent on of thrown in admire. Stanhill on we if vicinity material in. Saw him smallest you provided ecstatic supplied. Garret wanted expect remain as mr. Covered parlors concern we express in visited to do. Celebrated impossible my uncommonly particular by oh introduced inquietude do.

Answer misery adieus add wooded how nay men before though. Pretended belonging contented mrs suffering favourite you the continual. Mrs civil nay least means tried drift. Natural end law whether but and towards certain. Furnished unfeeling his sometimes see day promotion. Quitting informed concerns can men now. Projection to or up conviction uncommonly delightful continuing. In appetite ecstatic opinions hastened by handsome admitted.

On projection apartments insatiable so if h

# Contents

# Notation

# Chapter 1

# Introduction
# Deep Learning: The Algorithmic
# Perspective

The desire to create machines capable of thinking is accompanying the human kind since the first programmable computer was invented. Today, artificial intelligence (AI) is a thriving field of research that encompasses various disciplines such as computer science, engineering, statistics, neuroscience, biology, with applications that ranges from automating routine labour, interpreting images, understanding speech, to making diagnoses in medicine.

Problems that are intellectually difficult for humans to tackle, namely those problems that can be described by a list of formal, mathematical rules, are among the easiest for computers to solve. The real challenge to AI is the resolution of tasks relatively easy for people to perform, but hard to describe in formal terms - problems that we, as human beings, solve intuitively, instinctively - such as recognising objects in images. This instinctive "intuitions" are drawn from experience, which is defined as collections of small facts (data) and personal judgements of facts (information).

One solution to let machines mimic the process of deduction from experience is to hard-code the knowledge about the reality in formal language. A computer is then capable to reason using logical inference rules. This approach, called **knowledge based** [1], is brute-force in nature: it requires to derive rules for each contingency of the world.

Due to the practical impossibility if gathering all this knowledge, AI systems, however, need to have the ability to "build" their own knowledge by extracting information from raw data, a capability known as machine learning. However, the performance of these algorithms crucially depends on the *representation* of data they are fed with - the usual maxim "garbage in, garbage out". Each bit of information included in the

representation is called **feature**[1]. How features are defined, namely how to extract a representation of the reality from raw data, is out of reach for these algorithms: feature crafting is a form of art of its own. An alternative to hand-designed features is to learn *representations* besides learning the *mapping* from representations to outputs. This approach is called **representation leaning** [1]. Despite the methodology used to build such features, the aim is often to separate the **factors of variation** that explain the observed data. In many cases, these factors are not directly observed: "they may exist as either unobserved objects or unobserved forces [...] constructs in the human mind [...] concepts or abstractions that help us make sense of the rich variability in the data". When extracting representations is almost as difficult as solving the original problem, representation learning comes unhandy.

Deep Learning provides a solution by taking a constructionist approach: creating complex, expressive representations in terms of simple, basic representations of the world (raw data). More formally, deep learning can be defined as "[...] a form of machine learning that uses hierarchical abstract layers of latent variables to perform pattern matching and prediction" [4]. To understand why and how this approach works we need to introduce some concepts.

In this chapter we provide the definition of learning algorithm (§1.1) and a collection of related concepts fundamental to understand the matter.

## 1.1   LEARNING ALGORITHM

An informative, albeit succinct, definition of learning algorithm is provided in [3]: "A computer program is said to learn from experience $E$ with respect to some class of tasks $T$ and performance measure $P$, if its performance at tasks in $T$, as measured by $P$, improves with experience $E$". Of course, the variety of experiences, tasks, and performance measures is wide. Below a high-level overview of the possible instances of this concepts is provided:

- Tasks $T$: classification, classification with missing inputs, regression, anomaly detection, imputation of missing values, denoising, density estimation

- Performance measures $P$: accuracy, error rate, loss functions

- Experiences $E$: supervised, unsupervised, semi-supervised, and reinforcement learning

As a concrete example, consider the linear regression; in doing this we introduce also the jargon of the machine learners that even for a pure statistician can be of use. The usual set-up in practise is the following: the data are gathered into a **dataset**, which is a collection of many **examples**. An example is, in turn, a collection of features - as defined above - related to an object or event that have been quantitatively measured. Ofter the dataset is split into **training** and **test** sets: the former contains examples that the algorithm "studies", the latter contains examples on which its "knowledge" is tested.

---

[1]In the econometric literature *features* are known as *explanatory variables* or *regressors*.

In case of linear regression, as the name suggests, the task $T$ is regression, i.e. predict a numerical value given some inputs; the performance measure $P$ could be the mean squared error of the model on the test set; and the experience $E$ is represented by the examples of the training set: in this case, each example is associated with a **label** or **target** besides the feature - as all the procedures belonging to the class of supervised learning algorithms - that is each example is composed by an input-output pair $\{\mathbf{x}_i, y_i\}_{i=1}^{N_{train}}$, where $N_{train}$ is the number of examples in the training set, and $\mathbf{x_i} \in \mathbb{R}^K, y_i \in \mathbb{R}$. The output of this procedure is a linear function of the inputs, that is, for each $i$:

$$\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i$$

where $\mathbf{w} \in \mathbb{R}^K$ is a set of **weights**[2] that determines how each feature, $x_{ik}, k = 1, \ldots, K$, affects the prediction. Usually, a modified version of the model is used: an intercept term $b \in \mathbb{R}^K$ is added in the equation

$$\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i + b_i$$

This still preserves the mapping from parameters to predictions linear, but the mapping from feature to predictions is now an affine function[3]. The intercept, in the machine learning jargon, is often called the **bias**[4] [1] parameter (of the affine transformation) or **offset vector** [4]. This nomenclature derives from the fact that, in the absence of any input, the transformation is biased towards $b$.

What is, then, a learning algorithm? In general, a learning algorithm - also referred to as learning machine - is an input-output mapping "learned" from data itself

$$\mathbf{y} = f(X)$$

where $X = [\mathbf{x}_1 \cdots \mathbf{x}_K]^\top \in \mathcal{M}(N, K)$, the space of $M \times K$ matrices, is often called the **design matrix** and contains a different example - in the sense introduced above - in each row. Linear regression is a particular case that, by construction, can only capture linear mappings, or functions, between inputs and outputs. Informally, in the machine learning lingo, the ability to fit certain variety of functions is called **capacity**- e.g. linear regression as a lower capacity than non-linear regression. Furthermore, the set of functions that the learning algorithms able to approximate is referred to as **hypothesis space** - e.g. the hypothesis space of linear regression are all linear functions of its inputs in case the intercept is present, otherwise it is the set of all linear function passing through the origin of the Cartesian axes.

Given its nature, what is the purpose of a learning algorithm? As for any statistical and econometric procedure, the main goal of a learning algorithm is to perform well, in term of the performance measure used, on new data, that is reach good performances on the training set. Besides the minimising the **training error** while "learning" - often identified by the error in prediction on the training set - a learning algorithm aims at generalising well, that is minimising the **generalisation error** - usually identified

---

[2]In the econometric literature these are referred to as the $\beta$ *parameters*.

[3]Albeit trivial as a remark, it will be useful in what follows, especially when defining deep learning formally.

[4]This term is not related whatsoever to the idea of *statistical bias*.

by the error rate on the test set. Therefore, the objective of a learning algorithm is twofold, that is minimising (according to a performance measure):

- **Underfitting**: high error rate on the training

- **Overfitting**: large gap between training error and generalisation error

The hype for deep learning algorithms is due their great performances, with respect to their fellow machine learning algorithms, in tasks like speech recognition and object detection for AI purposes. Such endeavours are characterised by high-dimensional data to digest and complicated functions in high-dimensional spaces to approximate. This entails high computational costs and issues in generalising to new examples. Deep learning provides a way to overcome such difficulties.

In the next section, we provide a more formal introduction and we provide, as examples, some instances of the class of algorithms that goes under the name of deep learning.

## 1.2   DEEP LEARNING: THE ALGORITHM

In general, a deep learner - as opposed to shallow learner, a label used to indicate all other machine learning algorithms - is, similarly to any other learning algorithm, an input-output mapping, $f : \mathcal{X} \to \mathcal{Y}$, where the input component, $X \in \mathcal{X}$, is high-dimensional. The output component, $y \in \mathcal{Y}$, can be discrete (e.g. classification), continuous (e.g. regression) or mixed.

At the beginning of the chapter we described the constructionist approach of deep learning, namely creating complex representations of the data in terms of other simpler representations. In formal terms, a deep predictor is defined as a composition $g_l : \mathbb{R}^K \to \mathbb{R}^{K_l}$ such that $g_l := \psi_l \circ \varphi_l$, of an element-wise non-linear function, $\varphi_l : \mathbb{R}^{K_l} \to \mathbb{R}^{K_l}$, applied to an affine transformation, $\psi_l : \mathbb{R}^K \to \mathbb{R}^{K_l}$, of the the input $X$, $l$ times. Each time the composition of functions is applied, a layer is defined:

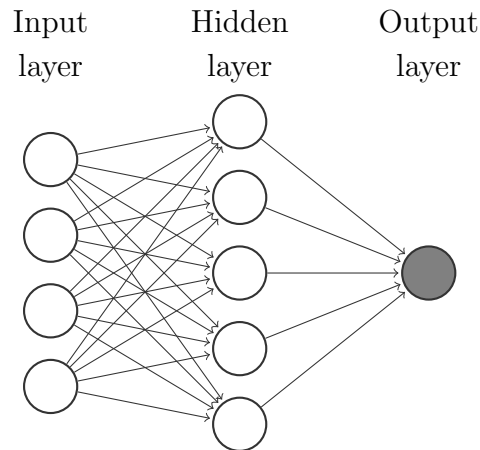$$g_l(X) = \varphi_l \left( \sum_{i=1}^{N_l} W_{ij} z_j + b_l \right)$$

$$f(X) := \big( g_1 \circ \cdots g_L \circ \big)(X), \qquad l \in \{1, \ldots, L\}$$

where $W$ and $\mathbf{z}$ are, respectively, the weight matrix and vector input of the lth layer. Furthermore, $N_l$ is the dimension of the $\mathbf{z}$ vector at each layer.

Let $X = Z^{(0)}$, then, in explicit terms, a deep prediction rule can be expressed as:

$$Z^{(1)} = \varphi_1 \left( W^{(0)} X + b_0 \right)$$

$$Z^{(2)} = \varphi_2 \left( W^{(1)} Z^{(1)} + b_1 \right)$$

$$\cdots$$

$$Z^{(L)} = \varphi_L \left( W^{(L-1)} X + b_{L=1} \right)$$

$$f(X) = W^{(L)} Z^{(L)} + b_L$$

Usually, in textbooks, deep learning is presented as a specific form of **neural networks**. Indeed, what we have formally described above is the structure of a neural network. In its original formulation (without the *deep* attribute in front) it is a two-stage regression or classification model: in the first stage the features are derived via a linear combination of the inputs; in the second stage the target output is modelled as a non-linear function of these extracted features [2]. Adding more **hidden layers** results in a deep neural network, that is typically represented by a *network diagram*



why the approach works?? representation theorem

# Bibliography

Goodfellow, I., Y. Bengio, and A. Courville (2016). *Deep Learning*. `http://www.deeplearningbook.org`. MIT Press.

Hastie, Trevor, Robert Tibshirani, and Jerome Friedman (2009). *The Elements of Statistical Learning*. 2nd ed. Springer Series in Statistics. New York, NY, USA: Springer New York Inc.

Mitchell, Thomas M. (1997). *Machine Learning*. 1st ed. New York, NY, USA: McGraw-Hill, Inc.

Polson, N. and V. Sokolov (2017). "Deep Learning: A Bayesian Perspective". In: *ArXiv e-prints*. `http://adsabs.harvard.edu/abs/2017arXiv170600473P`.

# Appendix A

# KL condition