

# Mathematical Engineering of Deep Learning

Book Draft

Benoit Liquet, Sarat Moka and Yoni Nazarathy

February 28, 2024

# Contents

<b>Preface - DRAFT</b>	<b>3</b>
<b>1 Introduction - DRAFT</b>	<b>1</b>
1.1 The Age of Deep Learning . . . . .	1
1.2 A Taste of Tasks and Architectures . . . . .	7
1.3 Key Ingredients of Deep Learning . . . . .	12
1.4 DATA, Data, data! . . . . .	17
1.5 Deep Learning as a Mathematical Engineering Discipline . . . . .	20
1.6 Notation and Mathematical Background . . . . .	23
Notes and References . . . . .	25
<b>2 Principles of Machine Learning - DRAFT</b>	<b>27</b>
2.1 Key Activities of Machine Learning . . . . .	27
2.2 Supervised Learning . . . . .	32
2.3 Linear Models at Our Core . . . . .	39
2.4 Iterative Optimization Based Learning . . . . .	48
2.5 Generalization, Regularization, and Validation . . . . .	52
2.6 A Taste of Unsupervised Learning . . . . .	62
Notes and References . . . . .	72
<b>3 Simple Neural Networks - DRAFT</b>	<b>75</b>
3.1 Logistic Regression in Statistics . . . . .	75
3.2 Logistic Regression as a Shallow Neural Network . . . . .	82
3.3 Multi-class Problems with Softmax . . . . .	86
3.4 Beyond Linear Decision Boundaries . . . . .	95
3.5 Shallow Autoencoders . . . . .	99
Notes and References . . . . .	111
<b>4 Optimization Algorithms - DRAFT</b>	<b>113</b>
4.1 Formulation of Optimization . . . . .	113
4.2 Optimization in the Context of Deep Learning . . . . .	120
4.3 Adaptive Optimization with ADAM . . . . .	128
4.4 Automatic Differentiation . . . . .	135
4.5 Additional Techniques for First-Order Methods . . . . .	143
4.6 Concepts of Second-Order Methods . . . . .	152
Notes and References . . . . .	164
<b>5 Feedforward Deep Networks - DRAFT</b>	<b>167</b>
5.1 The General Fully Connected Architecture . . . . .	167
5.2 The Expressive Power of Neural Networks . . . . .	173
5.3 Activation Function Alternatives . . . . .	180
5.4 The Backpropagation Algorithm . . . . .	184
5.5 Weight Initialization . . . . .	192

## Contents

5.6	Batch Normalization . . . . .	194
5.7	Mitigating Overfitting with Dropout and Regularization . . . . .	197
	Notes and References . . . . .	203
<b>6</b>	<b>Convolutional Neural Networks - DRAFT</b>	<b>205</b>
6.1	Overview of Convolutional Neural Networks . . . . .	205
6.2	The Convolution Operation . . . . .	209
6.3	Building a Convolutional Layer . . . . .	216
6.4	Building a Convolutional Neural Network . . . . .	226
6.5	Inception, ResNets, and Other Landmark Architectures . . . . .	236
6.6	Beyond Classification . . . . .	240
	Notes and References . . . . .	247
<b>7</b>	<b>Sequence Models - DRAFT</b>	<b>249</b>
7.1	Overview of Models and Activities for Sequence Data . . . . .	249
7.2	Basic Recurrent Neural Networks . . . . .	255
7.3	Generalizations and Modifications to RNNs . . . . .	265
7.4	Encoders Decoders and the Attention Mechanism . . . . .	271
7.5	Transformers . . . . .	279
	Notes and References . . . . .	294
<b>8</b>	<b>Specialized Architectures and Paradigms - DRAFT</b>	<b>297</b>
8.1	Generative Modelling Principles . . . . .	297
8.2	Diffusion Models . . . . .	306
8.3	Generative Adversarial Networks . . . . .	315
8.4	Reinforcement Learning . . . . .	328
8.5	Graph Neural Networks . . . . .	338
	Notes and References . . . . .	353
	<b>Epilogue - DRAFT</b>	<b>355</b>
<b>A</b>	<b>Some Multivariable Calculus - DRAFT</b>	<b>357</b>
A.1	Vectors and Functions in $\mathbb{R}^n$ . . . . .	357
A.2	Derivatives . . . . .	359
A.3	The Multivariable Chain Rule . . . . .	362
A.4	Taylor's Theorem . . . . .	364
<b>B</b>	<b>Cross Entropy and Other Expectations with Logarithms - DRAFT</b>	<b>367</b>
B.1	Divergences and Entropies . . . . .	367
B.2	Computations for Multivariate Normal Distributions . . . . .	369
	<b>Bibliography</b>	<b>399</b>
	<b>Index</b>	<b>401</b>

## Epilogue - DRAFT

Our story was about the **mathematical engineering of deep learning**. Our goal was to describe deep learning ideas in simple mathematical terms. Our goal was not to study implementation of deep learning; it was not to discuss the history and evolution of deep learning; and it was not to dive into subtle mathematical properties of deep learning. We simply wanted to present a basic **mathematical** description, empowering the reader with an understanding of key concepts and terminology. Mathematics is a language of choice.

We focused on the most popular and successful **deep learning** architectures and ideas that emerged over recent years. Somewhat anti-climatically we claim that the popularity and success of these ideas is due to their practical applicability, and not so much due to mathematical elegance. There are many other variants that we did not present here which are interesting and elegant yet have not been as popular from a practical perspective. With this we note that the aspect of **engineering** focusing on the empirical evaluation of architectures was not discussed and studied in the book at all.

Take as an example the *transformer architecture* studied in Section 7.5. This architecture has been pivotal in *large language models*. Indeed, in the same years that we worked on writing this book, 2021–2023, large language models, almost exclusively powered by the transformer architecture, have risen in popularity. Yet it is fair to say that the transformer architecture is quite arbitrary. If a couple of years prior to the development of this architecture, published in 2017 with [410], we the authors would have been presented with a transformer, without empirical trials and experimentation results, we would have no proof that transformers work so well.

It is also important to note that the pace and unpredictability of deep learning developments moves fast. By now, large language models have effectively beaten the Turing test, [38], a goal which seemed yet unattainable in the days when we conceived this book in late 2020. So our humble claim is that while **mathematical engineering** is important, in its own right, without computers, GPUs, software, data, and experimentation, it is void of substance. Nevertheless, we do believe that our presentation approach is succinct and unique, and given that the ideas that we present were previously shown to be winning ideas, the knowledge that you gained by reading this book will be beneficial.

Finally we close by mentioning that while this is a mathematical book, one cannot ignore the vast area of ethical issues associated with deep learning and artificial intelligence. Now, as we are in the third decade of the twenty first century, artificial intelligence is at the center of discussions associated with politics, freedom, social justice, violence, equity, and many other domains. Since this book is not about applications, we as authors had the luxury of ignoring the many ethical issues associated with deep learning in our exposition. Nevertheless, any practitioner using deep learning should at onset make sure to consider what defines responsible use and what not. We certainly want the technology to be used for purposes that do good rather than bad.