

# Optimizing AI and Machine Learning Solutions

---

*Your ultimate guide to building  
high-impact ML/AI solutions*

---

**Mirza Rahim Baig**



[www.bpbonline.com](http://www.bpbonline.com)

First Edition 2024

Copyright © BPB Publications, India

ISBN: 978-93-55519-818

*All Rights Reserved.* No part of this publication may be reproduced, distributed or transmitted in any form or by any means or stored in a database or retrieval system, without the prior written permission of the publisher with the exception to the program listings which may be entered, stored and executed in a computer system, but they can not be reproduced by the means of publication, photocopy, recording, or by any electronic and mechanical means.

### **LIMITS OF LIABILITY AND DISCLAIMER OF WARRANTY**

The information contained in this book is true to correct and the best of author's and publisher's knowledge. The author has made every effort to ensure the accuracy of these publications, but publisher cannot be held responsible for any loss or damage arising from any information in this book.

All trademarks referred to in the book are acknowledged as properties of their respective owners but BPB Publications cannot guarantee the accuracy of this information.

To View Complete  
BPB Publications Catalogue  
Scan the QR Code:



**Dedicated to**

*My lovely wife*

## About the Author

Since 2009, **Mirza Rahim Baig** has been exploring, practicing, learning, and teaching all things machine learning / artificial intelligence. He is a seasoned data science expert and renowned thought leader. Rahim is adept at solving complex business problems using AI/ML in a career spanning multiple domains and geographies. Numerous job titles aside, his focus has always been using data science to solve business problems and create high impact.

Rahim is also an author, a speaker and educator. In addition to this book, he has authored two well received books *The Deep Learning Workshop* and *Data Science for Marketing Analytics*. In all his titles, Rahim focuses on application and outcomes, and shares best practices derived from his experience. A key feature of his work is the approach of focusing not merely on model building but on the entire end-to-end solution building process. The explanations are easy to follow, even for non-technical learners, thanks to his style developed over years of teaching.

Rahim is a renowned subject matter expert in data science topics, having worked with most of the popular ed-tech platforms in designing courses and delivering content for master's level programs. He is a visiting faculty at NMIMS for MBA programs. In addition to having been published multiple times, Rahim speaks on various platforms about data science.

## About the Reviewer

**Sujit Pal** is a Technology Research Director at Elsevier Health Markets, where he works as an applied researcher at the intersection of search and Machine Learning. His interests are in Information Retrieval, Natural Language Processing, Knowledge Graphs and Machine Learning, and more recently, Large Language Models and Generative AI. He has co-authored two books on Deep Learning and one on Graph-based Natural Language Processing.

## Acknowledgement

In writing this book, my aim has been to benefit others with the little knowledge I have. However, it would be insincere to claim that all the insights and knowledge contained in these chapters originated solely from me. A lot comes from my incredible mentors and excellent teachers. What you find here are also their tremendous efforts reaching you through my words.

I would like to acknowledge the achievements of the open-source community. In the past few years, this community has democratized artificial intelligence and machine learning like never before, making state-of-the-art models accessible to anyone with an internet connection. We are truly witnessing a revolution and I am genuinely grateful for it.

I am also grateful to the team at BPB Publications for their guidance and patience throughout the journey of writing this book. Working with some extremely competent editors, reviewers, and technical experts was an enriching experience.

Finally, I would like to thank all the readers who have taken an interest in my book and for their support in making it a reality. Your encouragement has been invaluable.

# Preface

Building High-impact Machine Learning / Artificial Intelligence solutions is intended for several types of readers:

- Data science developers (data analysts/scientists) who will create the data science solutions using data,
- Data science professionals who manage data science teams or data science projects,
- Aspiring data scientists who want to learn more than just the techniques,
- Anyone interested in learning the hows of building great solutions and the whys in making the multiple choices in data science solution building.

A key difference between this book and the others in the market is that this book provides a holistic view on data science problem solving. Throughout, the book keeps reminding that business impact and utility are paramount and reminds the learner to be pragmatic in their approach. The book approaches data science solution building using a principled framework. Machine learning model building is one of the many steps in the framework, and this book treats it accordingly. The book will teach what optimization means at each step. Whether it is problem formulation or hyper-parameter tuning for deep learning models, this book provides frameworks and case studies, with extensive hands-on wherever possible.

This book is not focused on the algorithms for building solutions and is not a replacement for algorithm focused books. The book will however provide you enough intuition, functional information, and references to help you make even complex deep learning models. Using this book you will create complex models without implementing the detailed mathematical operations yourself. Instead, you will use abstractions and helpful libraries wherever possible to create clean and simple code. In the later chapters, the book will even teach you how to use transfer learning for using readymade ML and AI models as a part of your solutions. You do not need extensive coding experience for machine learning, the book will guide you step-by-step and provide code and explanations.

With this book, you will be able to approach a ML/AI solution in a systematic way, optimizing each step of the solution building process. You will thus be able to create high impact, high value machine learning/artificial intelligence solutions – which is the ultimate goal of every data science professional.

**Chapter 1: Optimizing a Machine Learning /Artificial Intelligence Solution** – introduces some core concepts and sets up the foundation for our journey together in this book. It provides an overview of machine learning, followed by addressing the various practical challenges in machine learning. It introduces some key ideas which will be expanded on in the later chapters. It is crucial to distinguish between simply making a model and carefully designing an end-to-end solution to the business problem. A framework to approach such end-to-end solutions is introduced and the chapter will point to the chapters in the book that help you optimize the solution at each step highlighted in the framework, to ultimately develop a truly optimized machine learning / artificial intelligence solution.

**Chapter 2: ML Problem Formulation: Setting the Right Objective** – discusses perhaps the most important step in any data science project that involves machine learning problem formulation. A problem can be formulated in various ways employing different solutions: machine learning or not. Further, there are multiple possible ML approaches. Making these decisions is not a trivial matter. The chapter also highlights the important of aligning the model objective with the business objective for achieving maximum impact.

**Chapter 3: Data Collection and Pre-processing** – addresses a basic reality in the ML / AI space – that the sophistication level of the technique is irrelevant if the input data is poor. The chapter explains the various factors that influence the data collection process and how you should tackle associated problems. A principled approach is introduced that helps you optimize the entire pre-processing process as a pipeline, with the help of a case study.

**Chapter 4: Model Evaluation and Debugging** – teaches model evaluation in the context of solving business problems. Learn various considerations that make a model better than another. The chapter introduces metrics for model evaluation, classification, and regression problems. Learn a principled way to diagnose and debug models to identify the cause of poor performance, along with approaches to handle underfitting and overfitting.

**Chapter 5: Imbalanced Machine Learning** – shows how to tackle one of the most common and challenging problems in machine learning / artificial intelligence. The chapter explains how imbalance impacts model evaluation and teaches how to diagnose issues due to imbalance. Apart from introducing better metrics for evaluation, it also introduces multiple approaches to tackle imbalance. Learn using several types of under/over sampling techniques to balance the dataset better for a more optimized machine learning solution.

**Chapter 6: Hyper-parameter Tuning** – is dedicated to the process of finding the best set of *settings* for the model, i.e., hyper-parameters, to get the best generalized performance. The chapter covers the key considerations in defining *best* and proceeds to demonstrate multiple automated approaches to searching for the best hyper-parameters – from exhaustive methods like grid search to approaches like Randomized search.



---

**Chapter 7: Parameter Optimization Algorithms** – provides a good understanding of mathematical optimization, the core of machine learning/artificial intelligence models. The chapter contextualizes optimization for machine learning and discusses aspects that make optimization for machine learning unique. It then teaches a general, principled approach to solving optimization problems for classical machine learning and explains the variations developed for deep neural networks. The chapter helps you make important choices in optimization, like loss functions and optimization algorithms, to maximize the performance of the model.

**Chapter 8: Optimizing Deep Learning Models** – focuses on optimizing neural networks. Like with classical machine learning, hyper-parameter tuning is a key step for optimizing deep neural networks as well. The chapter focusses on the simple fully connected network architecture for these concepts and the accompanying case studies. In addition to automated hyper-parameter tuning for deep neural networks, the chapter will introduce several techniques to improve model performance, like data addition, data augmentation, ensembling, regularization, and injecting your domain knowledge via manual features.

**Chapter 9: Optimizing Image Models** – dives deeper by teaching optimization of deep learning models for image tasks. This chapter teaches concepts, architectures and parameters that are unique to Image models and then demonstrates how to optimize these using TensorFlow 2.0. The chapter discusses the design approach used in various popular model approaches and teaches you to implement your own versions of VGG and ResNet models. It teaches how to effectively use regularization using multiple approaches - dropout, data augmentation, and the powerful batch normalization technique, ending with some general guidelines for designing deep learning models for image tasks.

**Chapter 10: Optimizing Natural Language Processing Models** – continues the deep learning deep dive by focusing on NLP models. The chapter introduces the peculiarities of natural language data and teaches handling them by appropriate text pre-processing and data representation (using embeddings). The chapter details the recurrent and architecture and the transformer architecture with hands on case studies. In addition to hyper-parameter tuning for NLP models, the chapter shares tricks like using 1D convolutions and pre-trained embeddings (a prelude to transfer learning).

**Chapter 11: Transfer Learning** – teaches you how to stand on the shoulders of giants and benefit from the excellent work done by the ML/AI community. The chapter explains the different types of transfer learning and how they are helpful in modern machine learning/artificial intelligence, highlighting both the benefits and limitations of it. The chapter demonstrates the usage of a SOTA model out of the box and eventually how to fine tune it for better performance. Transfer learning is slightly different for image processing and NLP tasks. The chapter will employ multiple SOTA models for both image and NLP tasks, employing popular libraries.

# Code Bundle and Coloured Images

Please follow the link to download the *Code Bundle* and the *Coloured Images* of the book:

**<https://rebrand.ly/zce3wew>**

The code bundle for the book is also hosted on GitHub at

**<https://github.com/bpbpublications/Optimizing-AI-and-Machine-Learning-Solutions>**.

In case there's an update to the code, it will be updated on the existing GitHub repository.

We have code bundles from our rich catalogue of books and videos available at **<https://github.com/bpbpublications>**. Check them out!

## Errata

We take immense pride in our work at BPB Publications and follow best practices to ensure the accuracy of our content to provide with an indulging reading experience to our subscribers. Our readers are our mirrors, and we use their inputs to reflect and improve upon human errors, if any, that may have occurred during the publishing processes involved. To let us maintain the quality and help us reach out to any readers who might be having difficulties due to any unforeseen errors, please write to us at :

**[errata@bpbonline.com](mailto:errata@bpbonline.com)**

Your support, suggestions and feedbacks are highly appreciated by the BPB Publications' Family.

Did you know that BPB offers eBook versions of every book published, with PDF and ePub files available? You can upgrade to the eBook version at [www.bpbonline.com](http://www.bpbonline.com) and as a print book customer, you are entitled to a discount on the eBook copy. Get in touch with us at :

**[business@bpbonline.com](mailto:business@bpbonline.com)** for more details.

At **[www.bpbonline.com](http://www.bpbonline.com)**, you can also read a collection of free technical articles, sign up for a range of free newsletters, and receive exclusive discounts and offers on BPB books and eBooks.

## Piracy

If you come across any illegal copies of our works in any form on the internet, we would be grateful if you would provide us with the location address or website name. Please contact us at [business@bpbonline.com](mailto:business@bpbonline.com) with a link to the material.

## If you are interested in becoming an author

If there is a topic that you have expertise in, and you are interested in either writing or contributing to a book, please visit [www.bpbonline.com](http://www.bpbonline.com). We have worked with thousands of developers and tech professionals, just like you, to help them share their insights with the global tech community. You can make a general application, apply for a specific hot topic that we are recruiting an author for, or submit your own idea.

## Reviews

Please leave a review. Once you have read and used this book, why not leave a review on the site that you purchased it from? Potential readers can then see and use your unbiased opinion to make purchase decisions. We at BPB can understand what you think about our products, and our authors can see your feedback on their book. Thank you!

For more information about BPB, please visit [www.bpbonline.com](http://www.bpbonline.com).

## Join our book's Discord space

Join the book's Discord Workspace for Latest updates, Offers, Tech happenings around the world, New Release and Sessions with the Authors:

<https://discord.bpbonline.com>



# Table of Contents

<b>1. Optimizing a Machine Learning /Artificial Intelligence Solution .....</b>	<b>1</b>
Introduction .....	1
Structure .....	1
Objectives .....	2
Case study: Text deduplication for online fashion.....	2
Understanding machine learning.....	3
Machine learning styles.....	4
<i>Supervised machine learning</i> .....	5
<i>Unsupervised machine learning</i> .....	6
<i>Reinforcement learning</i> .....	6
<i>Choosing the ML style</i> .....	7
Challenges in ML / AI .....	7
<i>Poor formulation</i> .....	8
<i>Invalid/poor assumptions</i> .....	8
<i>Data availability and hygiene</i> .....	9
<i>Representative data (lack of)</i> .....	9
<i>Model scalability</i> .....	10
<i>Infeasible consumption</i> .....	10
<i>Misalignment with business outcomes</i> .....	11
ML / AI models vs. end-to-end solutions .....	12
CRISP-DM framework .....	13
Optimization at each step of solution development.....	14
<i>Business understanding</i> .....	15
<i>Data understanding</i> .....	16
<i>Data preparation</i> .....	17
<i>Model building</i> .....	18
<i>Evaluation</i> .....	19
<i>Deployment</i> .....	20
Conclusion .....	21

---

<b>2. ML Problem Formulation: Setting the Right Objective .....</b>	<b>23</b>
Introduction .....	23
Structure .....	23
Objectives .....	24
Identifying a Machine Learning problem.....	24
Choosing the ML style.....	26
Problem 1: Fraud detection.....	27
<i>Approach 1: Supervised classification .....</i>	<i>27</i>
<i>Approach 2: Unsupervised clustering.....</i>	<i>28</i>
Problem 2: Predicting high-selling products.....	29
<i>Approach 1: Supervised Regression .....</i>	<i>29</i>
<i>Approach 2: Supervised Classification .....</i>	<i>29</i>
Problem 3: Question de-duplication for an e-commerce website .....	30
<i>Approach 1: Supervised classification using Deep Learning .....</i>	<i>30</i>
<i>Approach 2: Unsupervised clustering of questions.....</i>	<i>31</i>
Blending ML styles .....	31
<i>For better accuracy .....</i>	<i>31</i>
<i>For dimensionality reduction .....</i>	<i>32</i>
<i>For better data representation .....</i>	<i>33</i>
Choosing the right dependent variable .....	33
Business objective vs. ML objective.....	35
Conclusion .....	36
<b>3. Data Collection and Pre-processing.....</b>	<b>37</b>
Introduction .....	37
Structure .....	37
Objectives .....	38
Building a machine learning solution.....	38
<i>The data collection process .....</i>	<i>39</i>
<i>The nature of the data.....</i>	<i>40</i>
<i>Domain-specific aspects.....</i>	<i>40</i>
<i>Influence of the task at hand.....</i>	<i>40</i>

---

The case study .....	40
The pre-processing process.....	42
<i>Step 1: Gather data + basic checks</i> .....	43
<i>Exercise 3.1: Loading and exploring bank marketing dataset</i> .....	44
<i>Exercise 3.2: Fixing formats and identifying missing values</i> .....	47
<i>Step 2: Separate into train and test datasets</i> .....	48
<i>Exercise 3.3: Splitting the data into train and test sets</i> .....	49
<i>Step 3: Outlier treatment</i> .....	51
<i>Percentile threshold approach</i> .....	52
<i>Z score approach</i> .....	52
<i>Box plots/Tukey criterion</i> .....	53
<i>Exercise 3.4: Outlier detection and treatment</i> .....	54
<i>Step 4: Missing value treatment</i> .....	57
<i>Step 5: Categorical feature handling</i> .....	59
<i>Solving our case study</i> .....	60
<i>Exercise 3.5: Handling remaining categorical features</i> .....	66
<i>Step 6: Transformation of numeric features</i> .....	69
<i>Exercise 3.6: Transforming numerical variables</i> .....	70
Case study steps review .....	73
Conclusion .....	74
<b>4. Model Evaluation and Debugging</b> .....	<b>75</b>
Introduction .....	75
Structure .....	76
Objectives .....	76
Ad click prediction case study .....	76
<i>Exercise 4.1: Load and prepare the data</i> .....	77
Model evaluation considerations.....	78
<i>Exercise 4.2: Building the competing models</i> .....	80
Model evaluation metrics .....	82
<i>Metrics for classification problems</i> .....	82

---

<i>Accuracy</i> .....	83
<i>Confusion matrix</i> .....	83
<i>Exercise 4.3: Confusion matrix for click prediction models</i> .....	85
<i>Precision, Recall, and F1 Score</i> .....	87
<i>Area Under ROC curve (AUC)</i> .....	90
<i>Area under PR curve</i> .....	93
<i>Metrics for regression</i> .....	95
Model evaluation schemes .....	96
<i>Exercise 4.4: Model evaluation using cross-validation</i> .....	98
Model debugging.....	100
<i>Overfitting and underfitting</i> .....	101
<i>Overfitting</i> .....	101
<i>Underfitting</i> .....	102
<i>Validation curves and goodness of fit</i> .....	103
<i>Exercise 4.5: Validation curves for Ad click prediction</i> .....	105
<i>Handling overfitting</i> .....	107
<i>Feature selection</i> .....	107
<i>Hyper-parameter control</i> .....	108
<i>Regularization</i> .....	109
<i>Data augmentation</i> .....	110
<i>Optimal training (early stopping)</i> .....	111
Conclusion .....	112
<b>5. Imbalanced Machine Learning</b> .....	<b>113</b>
Introduction .....	113
Structure .....	114
Objectives .....	115
Understanding the business problem .....	115
<i>Exercise 5.1: Understanding and staging the data</i> .....	116
Difficulty in model evaluation .....	118
<i>Exercise 5.2: Accuracy for imbalanced classes</i> .....	119

---

Evaluation metrics for imbalanced classes.....	123
<i>F1 score</i> .....	123
<i>Precision-Recall Curve</i> .....	125
<i>Summary of model evaluation</i> .....	126
Handling imbalance .....	127
<i>Adjusting the cost function</i> .....	127
<i>Exercise 5.3: Adjusting class_weight</i> .....	129
<i>Data balancing</i> .....	130
Undersampling.....	131
<i>Exercise 5.4: Random undersampling of insurance data</i> .....	131
<i>Model evaluation with re-balancing</i> .....	133
<i>Using Pipelines for evaluating re-balanced data</i> .....	135
<i>Undersampling methods: Tomek Links</i> .....	136
<i>Exercise 5.5: Tomek Links applied to insurance data</i> .....	138
<i>Undersampling methods: edited nearest neighbour</i> .....	139
<i>Exercise 5.6: Edited nearest neighbour undersampling</i> .....	140
Oversampling .....	142
<i>Synthetic Minority Oversampling Technique</i> .....	144
<i>Exercise 5.7: Synthetic Data Generation with SMOTE</i> .....	145
<i>ADASYN - Adaptive Synthesis</i> .....	146
<i>Exercise 5.8: Synthetic Data Generation with ADASYN</i> .....	147
Combination of under and oversampling.....	148
<i>Exercise 5.9: Using a combination of over and undersampling</i> .....	149
Comparing the approaches.....	150
Conclusion .....	152
<b>6. Hyper-parameter Tuning.....</b>	<b>155</b>
Introduction .....	155
Structure .....	156
Objectives .....	156
Parameters and hyper-parameters .....	157
Importance of hyper-parameters .....	159



---

Hyper-parameter tuning.....	163
Manual looping over hyper-parameters.....	164
<i>Perform your own manual hyper-parameter tuning</i> .....	164
GridSearchCV - Grid Search with Cross-Validation .....	166
<i>Using GridSearchCV for hyper-parameter tuning</i> .....	168
<i>GridSearchCV with multiple hyper-parameters</i> .....	170
<i>The effect of each hyper-parameter</i> .....	172
Limitations of exhaustive methods (looping, GridSearchCV) .....	175
RandomizedSearchCV.....	178
Grid Search versus Randomized Search.....	181
Choosing a hyper-parameter search approach .....	182
Choosing the 'best' model .....	183
Training on the entire data and making predictions.....	185
The complete process .....	186
Conclusion .....	187
<b>7. Parameter Optimization Algorithms.....</b>	<b>189</b>
Introduction .....	189
Structure .....	189
Objectives .....	190
Optimization in machine learning .....	190
Fundamentals of mathematical optimization.....	191
<i>Objective/loss functions</i> .....	191
<i>Local and global optima</i> .....	192
<i>Convex vs. non-convex loss functions</i> .....	193
Numerical optimization for machine learning.....	193
<i>Regression</i> .....	194
<i>Classification</i> .....	194
<i>Regularization</i> .....	195
A general solving approach.....	196
<i>Gradient descent</i> .....	197
<i>Stopping condition</i> .....	199

---

Faster gradient descent variants .....	199
<i>Momentum</i> .....	199
<i>AdaGrad optimizer</i> .....	201
<i>RMSProp optimizer</i> .....	202
<i>Adam optimizer</i> .....	202
<i>Choosing the optimization approach</i> .....	203
Conclusion .....	204
<b>8. Optimizing Deep Learning Models.....</b>	<b>205</b>
Introduction .....	205
Structure .....	206
Objectives .....	207
The image processing case study.....	207
<i>The overall modeling approach</i> .....	210
Building a baseline model.....	211
Hyper-parameters in neural networks .....	216
<i>Activation function</i> .....	216
<i>Learning rate</i> .....	226
<i>Tuning the learning rate</i> .....	228
<i>Network structure</i> .....	232
<i>Optimizer</i> .....	235
<i>Other hyper-parameters</i> .....	236
Coarse plus fine tuning .....	237
Practical tips for tuning Neural Networks .....	237
Beyond hyper-parameters: other approaches to boost performance .....	238
<i>Regularization</i> .....	238
<i>More data</i> .....	239
<i>Data augmentation</i> .....	239
<i>Ensembles</i> .....	241
<i>Inject Domain knowledge through features</i> .....	241
Conclusion .....	241
Additional reads.....	242

---

<b>9. Optimizing Image Models .....</b>	<b>243</b>
Introduction .....	243
Structure .....	244
Objectives .....	245
Applications of image processing.....	245
Case study: Auto image classification.....	247
Feature detection using convolutions .....	251
Traditional Convolutional Neural Networks.....	253
Parameter regularization .....	258
<i>Spatial Dropout</i> .....	259
<i>Exercise 9.1: Convolution model with dropout</i> .....	260
<i>Batch normalization</i> .....	262
<i>Data augmentation</i> .....	265
<i>Exercise 9.2: Conv model with data augmentation</i> .....	268
Popular architectures.....	270
<i>VGG design approach</i> .....	271
<i>Exercise 9.3: Creating VGG16 from scratch</i> .....	272
<i>ResNet design approach</i> .....	275
Guidelines for designing image models.....	280
Conclusion .....	281
<b>10. Optimizing Natural Language Processing Models .....</b>	<b>283</b>
Introduction .....	283
Structure .....	284
Objectives .....	284
Key tasks in NLP .....	285
Importance of sequence processing.....	286
<i>A text classification case study</i> .....	287
Text pre-processing .....	289
Text representation.....	292
<i>Count-based features</i> .....	292
<i>Text embeddings</i> .....	292

---

<i>Embedding layer in Keras</i> .....	293
Architectures for NLP .....	294
<i>Recurrent architectures</i> .....	295
<i>Long Short Term Memory</i> .....	300
<i>Gated recurrent units</i> .....	303
<i>Bi-directional recurrent architectures</i> .....	304
<i>Recurrent architectures with attention</i> .....	305
<i>Transformer architecture</i> .....	306
Tuning network hyper-parameters .....	309
<i>Exercise 10.1: Hyper-parameter tuning of GRU model</i> .....	310
Using convolutions for NLP .....	313
<i>Exercise 10.2: Using 1D convolutions and RNNs</i> .....	314
Using pre-trained embeddings .....	315
<i>Benefits</i> .....	316
<i>Drawbacks</i> .....	316
<i>Exercise- Text classification with pre-trained embeddings</i> .....	319
<i>Notable pre-trained embeddings</i> .....	321
Conclusion .....	323
<b>11. Transfer Learning</b> .....	<b>325</b>
Introduction .....	325
Structure .....	325
Objectives .....	326
Motivation for transfer learning .....	326
Using a SOTA image classification model .....	327
What is transfer learning? .....	330
<i>Applications of transfer learning</i> .....	331
<i>Ready to use pre-trained models</i> .....	331
<i>Fine tuning models</i> .....	332
<i>Data representation/feature extraction</i> .....	332
<i>Benefits of transfer learning</i> .....	332
<i>Limitations of transfer learning</i> .....	333

---

<i>Sources for pre-trained models</i> .....	334
<i>Keras Applications</i> .....	334
<i>Tensorflow Hub</i> .....	335
<i>Hugging Face</i> .....	336
<i>Kaggle</i> .....	337
<i>GitHub</i> .....	338
<i>Word embeddings for text</i> .....	339
General transfer learning workflow .....	339
Transfer learning for images .....	340
<i>ImageNet database</i> .....	341
<i>Choosing a Transfer Learning approach</i> .....	341
<i>Fine tuning approaches</i> .....	342
<i>Fine tuning a SOTA image model</i> .....	345
Transfer learning for text .....	350
<i>Applications of transfer learning for text</i> .....	350
<i>Prediction using pre-trained models</i> .....	350
<i>Feature extraction</i> .....	351
<i>Fine tuning</i> .....	351
<i>Hugging face for transfer learning</i> .....	351
<i>Sentiment analysis</i> .....	354
<i>Exercise 11.1: RoBERTa for tweet sentiment classification</i> .....	354
<i>Text generation</i> .....	355
<i>Zero-shot classification</i> .....	357
<i>Translation</i> .....	358
<i>Fine tuning a SOTA NLP model</i> .....	359
<i>Exercise 11.2: Fine tuning a BERT model</i> .....	360
Conclusion .....	363
<b>Index</b> .....	<b>365-370</b>



# CHAPTER 1

# Optimizing a Machine Learning / Artificial Intelligence Solution

## Introduction

This chapter will provide an overview of **Machine Learning (ML)**, followed by addressing the various practical challenges in machine learning. This chapter will introduce some key ideas which will be expanded on in the later chapters. We will make the crucial distinction between simply making a model and carefully designing an end-to-end solution to the business problem. We will learn about a framework to approach such end-to-end solutions learn what it means to optimize at each step, and ultimately develop a truly optimized machine learning/ artificial intelligence solution.

## Structure

In this chapter, we will cover the following topics:

- Case study
- Understanding machine learning
- Machine learning styles
- Challenges in ML / AI
  - Poor formulation
  - Invalid assumptions

- Data availability and hygiene
- Representative data (lack of)
- Model scalability
- Infeasible consumption
- Misalignment with business objectives
- ML/ AI models vs. end-to-end solutions
- CRISP-DM framework for solution development
- Optimization at each step of solution development
  - Business understanding
  - Data understanding
  - Data preparation
  - Model building
  - Evaluation
  - Deployment
- Conclusion

## Objectives

In this chapter, we will take a good, holistic look at the field of machine learning. This chapter will introduce some key ideas which will be expanded on in the later chapters. We will make the crucial distinction between simply making a model and carefully designing an end-to-end solution to the business problem. We will learn about a framework to approach such end-to-end solutions and learn what it means to optimize at each step, and ultimately develop a truly optimized machine learning/artificial intelligence solution. The various examples and case studies in this chapter will make the ideas concrete.

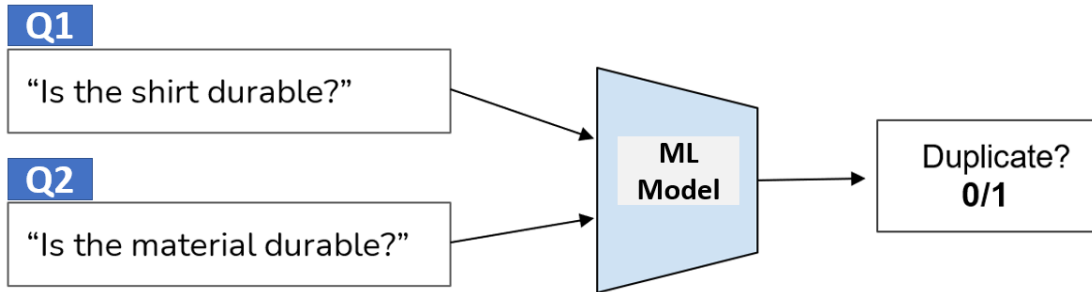
Consider this chapter as the gateway – where you get an overview of the steps in creating high-impact, optimized machine learning/artificial intelligence solutions. Each of the steps/ideas we discuss in this chapter will be dealt with in detail in the chapters that follow.

## Case study: Text deduplication for online fashion

Consider a data scientist working at the online fashion giant Azra Inc., to make the product detail page most helpful to the shopper. The product page contains detailed information about the product, including ratings, reviews, and questions that users ask about the product. The user questions section is of particular concern. There is a severe duplication



of questions. Users are asking the same question with minor variations in language. For example, *Is the material durable?* can be considered a duplicate of *Is the shirt durable?* Due to this, a few common questions are suppressing the visibility of other useful questions and their answers, withholding useful information from users, and affecting product sales. The task for the data scientist is to use their ML / AI expertise to identify duplicate questions as shown in *Figure 1.1*:



*Figure 1.1: Deduplication using supervised classification*

The data scientist formulates this as a supervised classification problem, as illustrated in *Figure 1.1*, using a deep learning model for text classification. This makes intuitive sense as we expect deep learning methods to shine in such situations. Using the latest transformer architecture should solve this, right? Unfortunately, in this case, the project was stopped after about 3 months of effort. The reason was a lack of sufficient labeled data.

For the text deduplication task, using a transformer architecture would require at least a few thousand pairs of questions labeled. The problem is labeling tens of thousands of question pairs. Manual labeling would take time and solid guidelines for the labelers so that their labels are in agreement. This is an expensive and time-consuming approach. This logical approach failed because of a presumption of data availability. The solution that eventually worked used an unsupervised clustering approach. The lesson is that improper problem formulation and presumptions can spell disaster for a ML / AI solution.

For success in ML / AI solutions, there are various considerations and decisions at various steps that need to be optimized; this is why data science projects fail. Before we discuss those, let us take a step back. Let us establish the understanding of machine learning that we will employ throughout this book. It is imperative that we take a holistic look at what ML is and more importantly, what it is good for, and how to make it work.

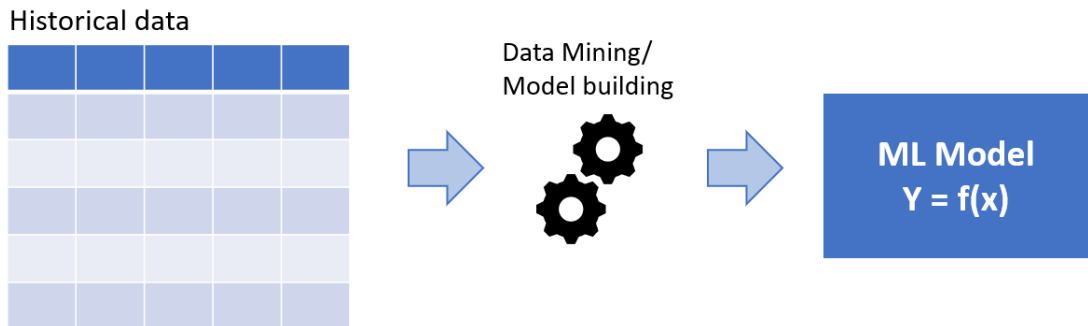
## Understanding machine learning

Our modern, data-driven world seeks to make decisions based on data insights to increasingly employ machines to perform repetitive tasks, drive cars, diagnose patients, allocate ads, recommend connections and songs, summarize news, etc. If data is oil in this new world, machine learning is the closest thing we have to the engine of this machinery.

Machine learning is the process that makes it possible to learn patterns from the provided data. The patterns learned by the machines can then be used to make some estimations/ predictions.

The outcome of the pattern learning process is often a mathematical model, capturing how the output relates to the input. The process of learning is also often referred to as model building or data mining. *Figure 1.2* illustrates this process. The historical data is input into the data mining/ model building process. The model-building process learns the patterns, which are expressed as a machine learning model. This model captures the relation between the inputs and the output and can therefore be employed to make estimations/ predictions. Depending on the technique employed, the model could be simpler and easily interpretable (e.g., a simple decision tree, or a linear regression equation), or a complex, hard-to-interpret from a deep neural network (complex series of matrix multiplications) that requires additional effort in post hoc explanations.

**Note: Data mining is a broader term that encapsulates the entire process of building a machine learning solution from data, the end-to-end process. Model building is merely one part of this process. We will discuss this in detail later in the chapter in the section titled CRISP-DM Framework.**



*Figure 1.2: Machines learning patterns from data*

## Machine learning styles

Let us now learn how to make machines learn the relevant patterns. We will have to make several decisions. For instance, we need to decide if we want to provide feedback and if yes, how to do that. We must define the kind of estimations that the machine needs to make. We need to define whether the model would be used to make predictions for the future or to uncover some patterns to aid human decision-making. Also important is to define the kind of data we input into the model. The specific solution depends on these considerations, but over the decades we have arrived at broadly three different machine learning styles, as illustrated in *Figure 1.3*:

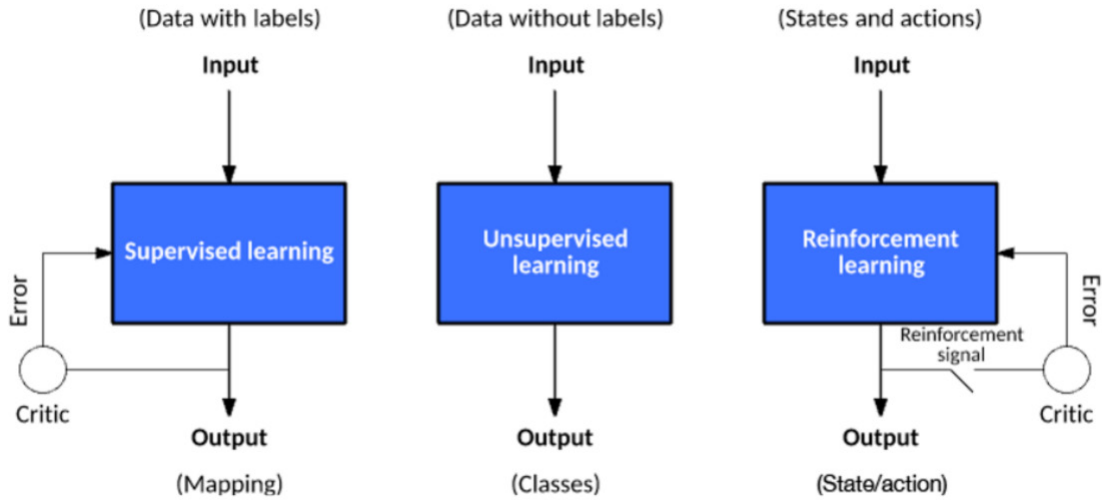


Figure 1.3: Machine learning styles

## Supervised machine learning

A key feature of supervised machine learning is that the data provided to the model contains the target as well. The input data comprises the features and the target, as shown in *Figure 1.4*. The features contain the information that will be used to predict the target. For using the model in the future, the input features will be available to us and will be used to predict the target. The machine learning process learns to predict the target using the input features, as illustrated in *Figure 1.4*:

### Supervised

Features					Target
X1	X2	...	...	Xn	Y

### Unsupervised

Features				
X1	X2	...	...	Xn

Figure 1.4: Input data for supervised vs. unsupervised ML

*Figure 1.4* illustrates how the supervision is done. The modeling technique sees the true target values (often called labels) and their predictions, compares them using a notion