# Artificial Intelligence: A Textbook

Charu C. Aggarwal
IBM T. J. Watson Research Center
Yorktown Heights, New York

March 23, 2021

To my wife Lata, my daughter Sayani,
and all my computer science instructors

# Contents

# Preface

"AI is likely to be either the best or the worst thing that happened to humanity."– Stephen Hawking

Artificial intelligence is a field that spans work from multiple communities, including classical logic programming, machine learning, and data mining. Since the founding of the field, there has a clear dichotomy between the *deductive reasoning* and the *inductive learning* forms of artificial intelligence. In the deductive reasoning view, one starts with various forms of domain knowledge (which are often stored as knowledge bases), and these forms of domain knowledge are used in order to make inferences. Such methods are often highly interpretable. The domain knowledge can be used in order to create hypotheses, which are then leveraged to make predictions. For example, in a chess game, the domain knowledge about the importance and position of pieces can be used to create a hypothesis about the quality of a position. This hypothesis can be used to predict moves by searching a tree of possible moves up to a specific number of moves. In learning methods, data-driven evidence is used to *learn* how to make predictions. For example, it is possible to generate data from chess games using self play, and then learn which moves are best for any particular (type of) position. Since the number of possible alternative move sequences in chess is too large to evaluate explicitly, chess programs often use various types of machine learning methods to relate typical patterns of pieces on the board to make predictions from carefully selected sequences. This approach is somewhat similar to how humans make chess moves. In the early years, deductive reasoning methods were more popular, although inductive learning methods have become increasingly popular in recent years. Many books in artificial intelligence tend to focus predominantly on deductive reasoning as a legacy from its dominance during the early years. This book has attempted to strike a balance between deductive reasoning and inductive learning.

The main disadvantage of inductive learning methods is that they are not interpretable, and they often require a lot of data. A key point is that humans do not require a lot of data to learn. For example, a child is often able to learn to recognize a truck with the use of a small number of examples. Although the best solutions to many problems in artificial intelligence integrate methods from both these areas, there is often little discussion of this type of integration. This textbook focuses on giving an integrated view of artificial intelligence, along with a discussion of the advantages of different views of artificial intelligence.

After presenting a broad overview in Chapter 1, the remaining portions of this book primarily belong to three categories:

1. *Methods based on deductive reasoning:* Chapters 2 through 5 discuss deductive reasoning methods. The primary focus areas include search and logic.

2. *Methods based on inductive learning:* Learning methods are discussed in Chapters 6 to 10. The topics covered include classification, neural networks, unsupervised learning, probabilistic graphical models, and reinforcement learning.

3. *Methods based on both reasoning and learning:* Chapter 11 to 13 discuss a number of methods that have aspects of both reasoning and learning. This include techniques like Bayesian networks, knowledge graphs, and neuro-symbolic artificial intelligence.

A number of topics of recent importance, such as transfer learning and lifelong learning, are also discussed in this book.

Throughout this book, a vector or a multidimensional data point is annotated with a bar, such as $\overline{X}$ or $\overline{y}$. A vector or multidimensional point may be denoted by either small letters or capital letters, as long as it has a bar. Vector dot products are denoted by centered dots, such as $\overline{X} \cdot \overline{Y}$. A matrix is denoted in capital letters without a bar, such as $R$. Throughout the book, the $n \times d$ matrix corresponding to the entire training data set is denoted by $D$, with $n$ data points and $d$ dimensions. The individual data points in $D$ are therefore $d$-dimensional row vectors, and are often denoted by $\overline{X}_1 \ldots \overline{X}_n$. On the other hand, vectors with one component for each data point are usually $n$-dimensional column vectors. An example is the $n$-dimensional column vector $\overline{y}$ of class variables of $n$ data points. An observed value $y_i$ is distinguished from a predicted value $\hat{y}_i$ by a circumflex at the top of the variable.