

Deep Learning on Microcontrollers

*Learn how to develop embedded
AI applications using TinyML*

Atul Krishna Gupta
Dr. Siva Prasad Nandyala



www.bpbonline.com

Copyright © 2023 BPB Online

All rights reserved. No part of this book may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, without the prior written permission of the publisher, except in the case of brief quotations embedded in critical articles or reviews.

Every effort has been made in the preparation of this book to ensure the accuracy of the information presented. However, the information contained in this book is sold without warranty, either express or implied. Neither the author, nor BPB Online or its dealers and distributors, will be held liable for any damages caused or alleged to have been caused directly or indirectly by this book.

BPB Online has endeavored to provide trademark information about all of the companies and products mentioned in this book by the appropriate use of capitals. However, BPB Online cannot guarantee the accuracy of this information.

First published: 2023

Published by BPB Online

WeWork

119 Marylebone Road

London NW1 5PU

UK | UAE | INDIA | SINGAPORE

ISBN 978-93-55518-057

www.bpbonline.com

Dedicated to

*My Parents Mr. Keshava Kumar Gupta and
Late. Smt. Lakshmi Devi Gupta and to
my wife Richa Gupta and to
my children Ananya Gupta and Avi Gupta*

— Atul Krishna Gupta

*My Parents Late Mr. Koti Nagaiah Nandyala and
Smt. Durgamba Nandyala and to
my brother Sambasiva Rao Nandyala and to
my wife Sandya Kemisetti*

— Dr. Sivaprasad Nandyala

About the Authors

- **Atul Krishna Gupta** has held many positions as Research & Development Executive in companies such as Syntiant, Macom, Inphi (now Marvell) and Gennum (now Semtech). He has over 25 years of experience in delivering all aspects of systems from IC design to software support. He has made contributions to various forums such as IEEE, SMPTE and OIF. Two technical Emmy Awards were granted to two companies for the technical work he led in the past. He was awarded with the Employee of the year award and Excellence in R&D award at Gennum. Atul holds over 20 patents. Currently, his research interests are in the field of Battery Management Systems (BMS) where he is finding ways to use AI to make Electrical Vehicles (EV) safer and last longer. He has received his B.Tech degree in Electrical Engineering from Indian Institute of Technology, Kanpur, India and MS degree in Electrical and Computer Engineering from University of Calgary, Canada.
- **Dr. Sivaprasad Nandyala** worked in Eaton Research Labs as Lead Engineer (Data Science) at Eaton India Innovation Center, Pune, India. Prior to Eaton, he worked in companies like Tata Elxsi, Wipro Technologies, Analog Devices & Ikanos Communications in multiple technology areas. Dr. Nandyala has over 35+ research publications, 1 patent grant and 6 patents under review. He obtained his Ph.D. in Speech Processing from NIT Warangal, India. He was an ERASMUS MUNDUS scholarship holder from the European government for his Postdoctoral Research at Politecnico di Milano (POLIMI), Italy.

About the Reviewer

Dr. Sanjay Boddhu is an experienced Research and Engineering leader with expertise in leading and mentoring geographically distributed teams, in the domains of Computer Vision, Image Processing, Natural Language Processing, Predictive Analytics, and Modelling. He is skilled in using various Machine Learning Ops approaches to design and develop real-world applications in Cloud and at Edge.

Acknowledgements

- **Atul Krishna Gupta:** I am incredibly grateful to Kurt Busch for providing me the opportunity to learn and contribute to the emerging field of Artificial Intelligence (AI). I would like to thank Dr. Jeremy Holleman for their insightful conversation over various topics related to neural networks. I would also like to thank Dr. Stephen Bailey for help in the firmware of the TinyML board. I would not have been able to showcase the TinyML board without guidance from Mallik Maturi and Poupak Khodabandeh.

I would like to offer a special thanks to Zack Shelby and Aurelien Lequertier for making available their Machine Learning (ML) platform to the developer community free of cost. The platform enables zero code deployment of production grade AI deployment.

I would like to thank my wife and children for their patience and support to finish the book.

- **Dr. Sivaprasad Nandyala:** The writing of a book is never a solo effort, and this “**Deep Learning on Microcontrollers**” book is no exception. Before anything else, I want to express my deepest gratitude to Mr. Atul Krishna Gupta who had faith in me from the beginning in writing this book. In addition, I would like to thank Dr. Sanjay Boddhu for his feedback and suggestions in reviewing the book. I want to thank the TinyML community for all the great things they have done for the field. My heartfelt thanks go to my family and friends for their unwavering support and understanding during this hard journey.

Preface

As the title of the book suggests, this book is intended to enable readers from different backgrounds to make a tangible AI application, which can be deployed on the edge on off-the-shelf platforms such as Arduino or TinyML board. The focus of this book is on the practical aspects of AI deployment. The journey of AI deployment from demo quality to production grade is not easy. We have taken a realistic example to show the pitfalls and given ideas on how to overcome the roadblocks.

While the focus of the book is on the practical side, the book also provides a good academic background as well. The field of AI is evolving and it is not practical to have one comprehensive book on all the topics, but we have given insight into some of the advanced topics of the AI field.

Deployment of AI on the edge will require some hardware. For cost effective deployment, it is expected that companies will develop their unique hardware. However, for getting started, there are several hardware boards available from websites such as Digikey or Amazon. Readers can buy this type of hardware in the range of \$35-\$100.

This book is divided into **9 chapters**. Each chapter description is listed as follows.

Chapter 1: Introduction to AI – will show a continuum of traditional code-based solution and Artificial Intelligence based solution. It will show where an AI based solution will be suitable and how to approach the solution.

Chapter 2: Traditional ML Lifecycle – will cover how machine learning is different from classical methods, introduction to traditional ML life cycle, performance metrics, and the basics of deep learning (DL) and different DL algorithms. It also covers transfer learning. We will discuss several tools, libraries, and frameworks for developing and deploying ML models on various embedded devices and microcontrollers. We also cover the differences between learning and inference, ML model deployment and inferencing on different hardware platforms and their comparison at various deployment levels.

Chapter 3: TinyML Hardware and Software Platforms – will cover CPUs, GPUs, Raspberry Pi boards, TPUs, and Data Center Servers. We will also look at TinyML compatible microcontrollers and Raspberry Pi boards. We then focus on TinyML's hardware boards and software platforms for machine learning. We will discuss important software platforms, data engineering, and model compression frameworks.

Chapter 4: End-to-End TinyML Deployment Phases – will discuss embedded machine learning's (EML) basics, characteristics, and examples. Next, we will also explore EML's building blocks, pros and cons, and how to run an ML model on microcontrollers. We will discuss Edge Impulse and Arduino IDE platforms, their pros and cons, and how to use different hardware boards with them. Data collection from sensors and the different platforms will be covered. We will cover data engineering, model training with Edge Impulse, optimization, and inferencing for model deployment on TinyML hardware platforms.

Chapter 5: Real World Use Cases – will cover various use cases of the TinyML deployment. The chapter categorizes these deployments in seven categories. However, many applications overlap multiple categories. These applications just show the tip of the iceberg because we just got started. Over the next few decades, we are expecting an explosion of TinyML deployment. These examples are provided just to ignite the creativity of the reader, so that they can lead innovation and deploy AI solutions which do not exist today.

Chapter 6: Practical Experiments with TinyML – will utilize Arduino IDE for TinyML hardware experimentation. We will collect sensor data using the TinyML board, clean the data for the practical experiment (Air Gesture Digit Recognition), upload it to the Edge Impulse platform, train and test the model with Nano RP2040 board sensor data. Finally, we will download the Edge Impulse inference model and test it on the RP2040 using Arduino IDE to evaluate performance.

Chapter 7: Advance Implementation with TinyML Board – will deep dive on specific hardware accelerator chips, which provide AI specific computation at a fraction of cost and power relative to microcontroller-based architecture. The development boards are readily available on these hardware accelerator chips where readers can deploy an AI solution. The power of the entire solution can run from batteries months to years. The chapter describes the entire flow of deployment in a few easy steps on the readily available Edge Impulse software platform.

Chapter 8: Continuous Improvement – will cover topics in improving the accuracy of the AI solution. AI is a data driven flow where the accuracy depends on the data. This chapter takes a deeper dive into a keyword detection application to demonstrate how to curate the data and improve the performance to take the solution from demo to production quality.

Chapter 9: Conclusion – will provide the conclusion of various aspects learned in the earlier chapters. This is an introduction book on AI and there are many topics which will require many more books. Some of those topics are mentioned in this chapter to ensure that the reader knows there is more to AI than what is covered in the book.

Code Bundle and Coloured Images

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

<https://rebrand.ly/yt0v6ae>

The code bundle for the book is also hosted on GitHub at **<https://github.com/bpbpublications/Deep-Learning-on-Microcontrollers>**. 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

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 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 for more details.

At **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** 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**. 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**.

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

1. Introduction to AI.....	1
Introduction.....	1
Structure.....	2
Objectives.....	2
Artificial Intelligence.....	3
Continuum of code writing and artificial intelligence.....	3
<i>Exercise</i>	3
Changing the paradigm.....	5
Neural Network.....	7
Machine Learning.....	11
Intelligent IoT System vs. Cloud based IoT system.....	13
<i>Arduino Nano 33 BLE Sense board</i>	14
<i>Limited compute resources</i>	15
<i>Battery power limits</i>	15
<i>TinyML and Nicla Voice board</i>	16
<i>>10x parameters</i>	18
<i>>200x Power advantage</i>	18
<i>>20x Throughput</i>	18
TinyML Ecosystem	19
Key applications for Intelligent IoT systems	19
<i>Smart agriculture</i>	20
<i>Smart appliances</i>	20
<i>Smart cities</i>	20
<i>Smart health</i>	21
<i>Smart homes</i>	21
<i>Smart industry</i>	21
Conclusion.....	21
Key facts.....	22
Questions	22

References	23
2. Traditional ML Lifecycle	25
Introduction.....	25
Structure.....	26
Objectives.....	26
Traditional methods	26
Machine learning landscape	27
<i>Supervised learning</i>	29
<i>Unsupervised learning</i>	29
<i>Reinforcement Learning (RL)</i>	30
ML Performance Metrics	30
<i>Confusion matrix</i>	30
Basics of DL and different DL algorithms.....	32
Transfer Learning.....	35
Tools and Different ML, DL frameworks	35
<i>Python</i>	36
<i>Jupyter Notebooks</i>	36
<i>Google Colaboratory</i>	36
<i>TensorFlow (TF), TFLite and TensorFlow Lite Micro</i>	36
<i>TensorFlow Lite</i>	37
<i>TensorFlow Lite Micro</i>	38
<i>AI Model Efficiency Toolkit (AIMET)</i>	38
<i>Convolutional Architecture for Fast Feature Embedding (Caffe)</i>	39
<i>CoreML</i>	40
<i>Open Neural Network Exchange (ONNX)</i>	41
<i>Open Visual Inference and Neural network Optimization (OpenVINO)</i>	41
<i>Pytorch and PyTorch Mobile</i>	42
Embedded Machine Learning (EML).....	42
Difference between Learning and Inference.....	43
ML model deployment and inferencing on different platforms.....	44
Conclusion.....	46
Key facts.....	47

Questions	47
References	48
3. TinyML Hardware and Software Platforms.....	51
Introduction.....	51
Structure.....	52
Objectives.....	52
Servers at Data Centers: CPUs, GPUs and TPUs	52
Mobile CPU, Raspberry Pi board and its types.....	53
Microcontrollers and Microcontroller with AI accelerator	55
TinyML Hardware Boards.....	57
<i>Arduino and Arduino Nano 33 BLE</i>	59
<i>Arduino Nicla Sense ME</i>	60
<i>Adafruit Feather</i>	61
<i>SparkFun Edge</i>	62
<i>NVIDIA Jetson Nano</i>	62
<i>Google Coral Edge TPU</i>	63
<i>Qualcomm QCS605</i>	64
<i>NXP i.MX 8M</i>	65
<i>STMicroelectronics STM32L4</i>	66
<i>Intel Curie</i>	67
<i>Syntiant TinyML</i>	68
TinyML Software Suites.....	69
<i>TensorFlow Lite Micro (Google)</i>	70
<i>uTensor (ARM)</i>	71
<i>Arduino Create</i>	72
<i>EloquentML</i>	72
<i>EdgeML (Microsoft)</i>	72
<i>EON Compiler (Edge Impulse)</i>	73
<i>STM32Cube.AI and NanoEdge AI Studio (STMicroelectronics)</i>	73
<i>PYNQ</i>	74
<i>OpenMV</i>	76
<i>SensiML</i>	76

<i>Neuton TinyML</i>	77
<i>Metavision Intelligence Suite 3.0 (Vision applications)</i>	78
Data Engineering Frameworks.....	78
<i>Edge Impulse</i>	78
<i>SensiML</i>	79
<i>Qeexo AutoML</i>	81
TinyML Model Compression Frameworks.....	81
<i>Quantization</i>	82
<i>Pruning</i>	83
<i>Low ranked approximation</i>	83
<i>Knowledge distillation</i>	83
<i>TensorFlow Lite</i>	84
<i>STM32 X-CUBE-AI</i>	85
<i>QKeras</i>	86
<i>Qualcomm AIMET</i>	86
<i>Microsoft NNI</i>	87
<i>CMix-NN</i>	89
<i>OmniML</i>	89
Conclusion.....	90
Key facts.....	90
Questions	91
References	92
4. End-to-End TinyML Deployment Phases.....	93
Introduction.....	93
Structure.....	94
Objectives.....	95
Understanding Embedded ML.....	95
Introduction to Edge-impulse and Arduino IDE	99
<i>Edge-impulse</i>	99
<i>Arduino Integrated Development Environment (IDE)</i>	109
<i>Arduino Driver Installation</i>	111
Data collection from multiple sensors.....	115

<i>Data collection from an Arduino board</i>	116
<i>Data collection from Syntiant board</i>	117
Data engineering steps for TinyML.....	117
<i>Cleaning</i>	118
<i>Organizing</i>	119
<i>Transformation</i>	119
Model Training in TinyML software platforms.....	120
<i>EON Compiler (Edge Impulse)</i>	120
Model Compression.....	122
<i>Pruning</i>	122
<i>Knowledge distillation</i>	123
Model conversion.....	124
<i>Quantization</i>	124
Inferencing/Prediction of results with test data.....	126
<i>Model Deployment in TinyML Hardware board</i>	128
Conclusion.....	130
Key facts.....	130
Questions.....	131
References.....	132
5. Real World Use Cases	133
Introduction.....	133
Structure.....	133
Objectives.....	135
Smart agriculture.....	135
<i>Agriculture video analytics</i>	135
<i>Crop intruder detection</i>	136
<i>Crop yield prediction and improvement</i>	136
<i>Agribots</i>	136
<i>Insect detection and pesticide reduction</i>	137
<i>Weedicides elimination</i>	137
<i>Acoustic insect detection</i>	138
<i>Animal husbandry</i>	139

Smart appliances.....	139
<i>Vision AI for appliances.....</i>	139
<i>Audio AI for appliances.....</i>	141
<i>Sensors based AI for appliances.....</i>	142
Smart cities.....	142
<i>Safe and secure city.....</i>	142
<i>City maintenance.....</i>	143
<i>Parking enforcement systems</i>	144
<i>Traffic management.....</i>	144
<i>Maintaining bridges.....</i>	144
<i>Non-Smoking enforcement</i>	145
Smart health.....	146
<i>Cataract detection</i>	146
<i>Fall detection.....</i>	147
<i>Cough detection</i>	148
<i>Boxing Moves Detector</i>	148
<i>Mosquito detection.....</i>	149
<i>Snoring and sleep apnea detection.....</i>	150
Smart home.....	151
<i>Person detection at the door.....</i>	151
<i>Glassbreak detection.....</i>	151
<i>Smart baby monitoring.....</i>	152
<i>Voice recognition for home automation</i>	153
Smart industry	153
<i>Railway track defect detection</i>	153
<i>Telecom towers defect detection</i>	154
<i>Defect detection in components.....</i>	155
Smart automotive	156
<i>Drowsy driver alert.....</i>	156
<i>Advance collision detection</i>	156
Conclusion.....	157
Key facts.....	157

Questions	158
References	158
6. Practical Experiments with TinyML.....	161
Introduction.....	161
Structure.....	162
Objectives.....	162
Introduction to Nano RP2040 TinyML board	163
<i>Setting up Arduino IDE and testing the Nano RP2040 Board</i>	<i>163</i>
<i>High level steps involved in the air gesture digit recognition in</i> <i> Edge Impulse platform.....</i>	<i>165</i>
Data collection for the air gesture digit recognition.....	166
<i>Loading the dataset in Edge Impulse Platform.....</i>	<i>170</i>
<i>Setting up the development framework and design of</i> <i> neural network classifier.....</i>	<i>174</i>
Model training in Edge Impulse platform	177
Model testing with the collected data.....	183
Model deployment in Nano RP2040 board	184
Inferencing/Prediction of results with RP2040.....	188
Conclusion.....	192
Key facts.....	193
Questions	193
References	193
7. Advance Implementation with TinyML Board	195
Introduction.....	195
Structure.....	195
Objectives.....	196
NDP101 Architecture	196
NDP120 Architecture	199
Practical implementation and deployment.....	199
<i>Creating a project.....</i>	<i>199</i>
<i>Uploading Data</i>	<i>201</i>
<i>Impulse Design.....</i>	<i>203</i>

<i>Epochs Setting</i>	209
<i>Learning rate setting</i>	209
<i>Validation data set setting</i>	209
<i>Auto balance setting</i>	210
<i>Data augmentation</i>	210
<i>Neural network architecture</i>	211
<i>Neural network training</i>	212
<i>Model testing</i>	216
<i>Deployment</i>	217
Conclusion	222
Key facts.....	222
Questions	222
References	223
8. Continuous Improvement	225
Introduction.....	225
Structure.....	225
Objectives.....	226
Expectation gap.....	226
Unique issues about audio application	226
<i>Raw neural network output and softmax transformation</i>	228
Handling anomalous behavior during target classifier testing	229
<i>Method 1: Running window averaging</i>	230
<i>Method 2: Enriching target classifier</i>	231
<i>Method 3: Enriching open set classifier</i>	232
False Acceptance Rate testing	238
Optimization of window size in running window averaging	240
Phrase recognition constraints to improve system level performance	246
FRR testing under noisy conditions.....	249
Improving FRR performance under noisy conditions	251
Data collection for continuous improvement.....	255
Conclusion	255
Key facts.....	255

Questions	256
References	256
9. Conclusion.....	257
Introduction.....	257
Structure.....	257
Objectives.....	258
Review of material covered in this book.....	258
<i>Chapter 1.....</i>	<i>258</i>
<i>Chapter 2.....</i>	<i>258</i>
<i>Chapter 3.....</i>	<i>259</i>
<i>Chapter 4.....</i>	<i>259</i>
<i>Chapter 5.....</i>	<i>259</i>
<i>Chapter 6.....</i>	<i>259</i>
<i>Chapter 7.....</i>	<i>260</i>
<i>Chapter 8.....</i>	<i>260</i>
Advanced topics	260
<i>Different types of neural networks</i>	<i>261</i>
<i>Neural network optimization.....</i>	<i>262</i>
<i>Zero-shot, One-shot or Few-shot learning.....</i>	<i>263</i>
<i>Federated learning</i>	<i>264</i>
<i>Transfer learning.....</i>	<i>267</i>
<i>Tuning pretrained networks</i>	<i>268</i>
<i>MLOps.....</i>	<i>270</i>
Key facts.....	270
Questions	270
References	271
Index	273-280

CHAPTER 1

Introduction to AI

Introduction

Artificial Intelligence (AI) today has touched all our lives without us realizing it. You may have used **Siri**, **Alexa** or **Google Assistant**. Did you ever wonder how it understands speech? During international trips, one often sees facial recognition in action, in airports. It used to take a lot of time for airline agents to check passports, but now, it is as easy as simply walking through a door. This door opens only when facial recognition is done and matches the passport. When people are sick and cannot type on their phone, all they need to do is use the voice assistant that comes along with all smartphones nowadays.

AI has become an integral part of many organizational ecosystems, not just at the consumer levels, which has brought forth many benefits such as increasing efficiency, and automating multiple tasks, while reducing installation and setup costs. For example, most machines which were installed in the last several decades, have an analog display. To replace all the monitors with digital meters is a nebulous task. Now, an image detector is placed on top of these displays, which can recognize the position of the needle and interpret the measurement in digital form. The information can be sent to the master control room via wireless protocol, such as

Wi-Fi, Bluetooth, **Long Range Radio (LoRa)** or even **Narrow Band IOT (NB-IOT)**. Refer to *Figure 1.1* for an illustration of the same:

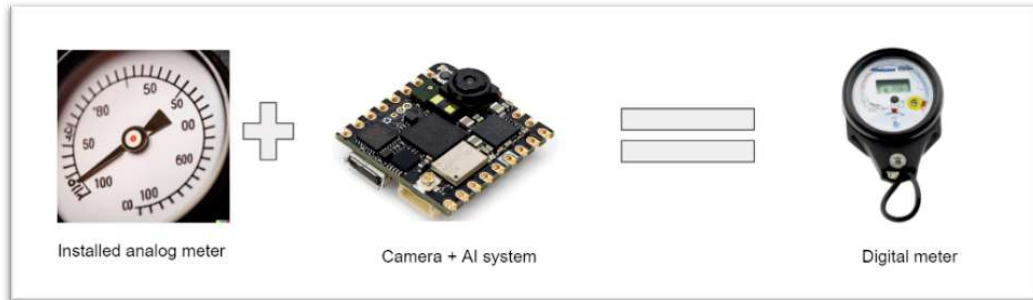


Figure 1.1: Transforming analog to digital with AI

Structure

In this chapter, the following topics will be covered:

- Artificial Intelligence
- Continuum of code writing and artificial intelligence
- Changing the paradigm
- Neural Network
- Machine Learning
- Intelligent IoT System vs. Cloud based IoT system
 - Arduino Nano 33 BLE Sense board
 - TinyML and Nicla Voice board
- TinyML Ecosystem
- Key applications for intelligent IoT systems

Objectives

By the end of this chapter, the reader will be able to relate what Artificial Intelligence has to offer. They will be familiar with the common lingo needed to bring an idea utilizing AI to a real system. Readers who are already doing firmware and software programming for the IoT devices can relate how their work will change when they plan to apply AI in their system. A concrete example is presented which shows where traditional methods will reach their limitations and AI deployment will be an easier path.

Artificial Intelligence

The intelligence demonstrated by computer systems is termed as artificial intelligence (Reference AI), as compared to natural intelligence demonstrated by living beings. The term **intelligence** could be controversial because as of today, the demonstrated capability of machines is still nowhere close to human intelligence. For this book, we will use the work Artificial Intelligence in the context of computers solving unique problems, which otherwise is not practical to solve with traditional code writing.

Continuum of code writing and artificial intelligence

It is expected that the reader is familiar with writing computer code. It can be argued that the problem which artificial intelligence solves, can also be solved with traditional computer code writing. However, the purpose of this book is to show that sometimes, seemingly simple problems could be very difficult to solve by traditional code writing. To appreciate the value of artificial intelligence, a hypothetical problem is posed here. Let us say an image of 200x200 pixels could contain a single line or could contain a circle. To simplify, let us assume the image is only of white or black color, as shown in *Figure 1.2*:

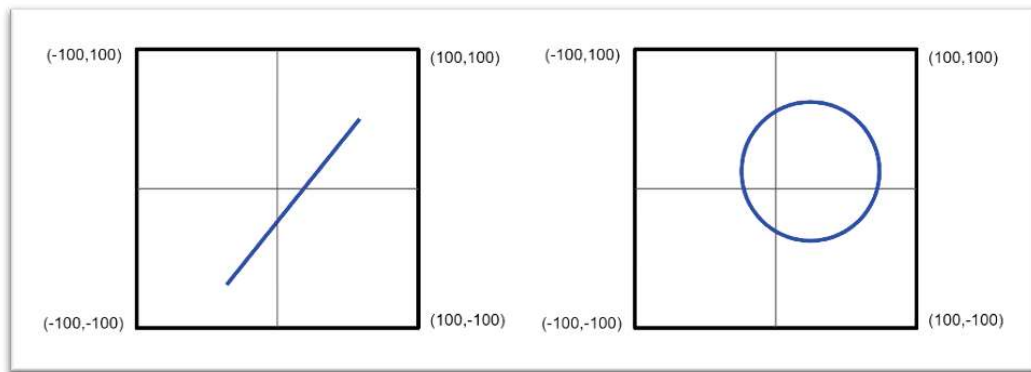


Figure 1.2: Image containing single line or circle

Exercise

Follow the given steps to perform the exercise:

- 1 Generate a 200x200 matrix with (0,0) points as origin, with value 0 (representing white space) or 1 (representing a dot in the curve).

- 2 Make multiple instances with lines of slope, ranging between ± 1 and y axis intercept of ± 50 , as shown in Figure 1.3:

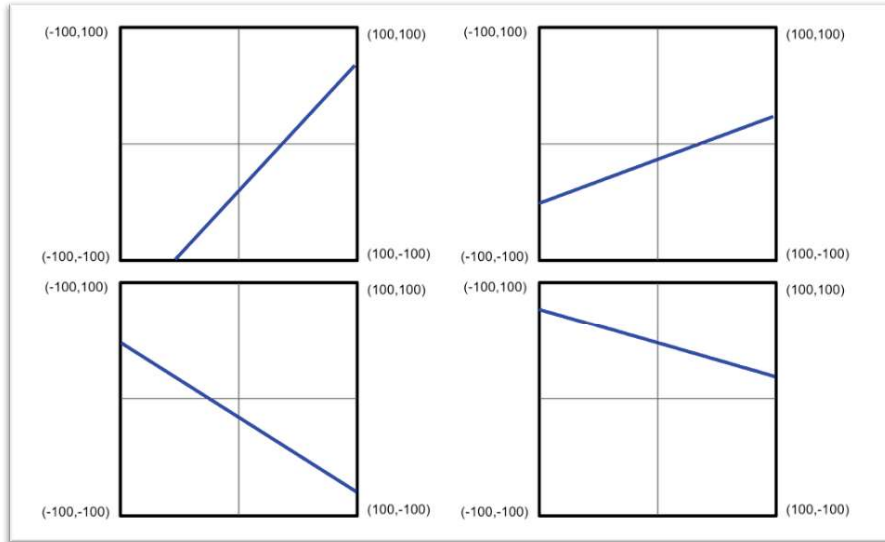


Figure 1.3: Instances with lines of different slopes and y intercepts

- 3 Similarly, make multiple instances of circles which fit completely within the image. Choose circles with radius of 10 to 100 and center within ± 50 units of (0,0) coordinate, as shown in Figure 1.4:

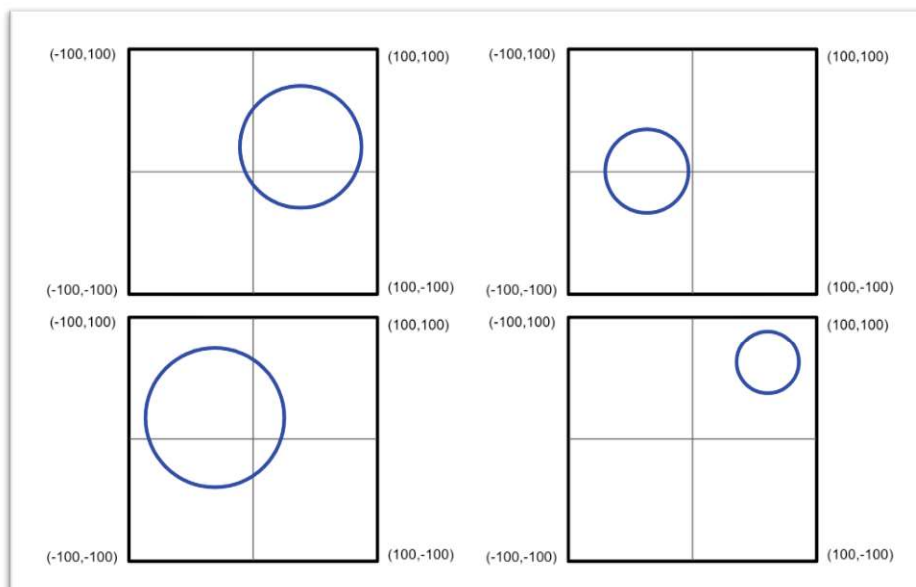


Figure 1.4: Instances with circles

- 4 Then, write a code with traditional logic and classify if this is line or circle.
- 5 Now extend the code to determine 0-9 digits in 28x28 pixels, using MNIST data Reference MNIST, as shown in *Figure 1.5*. If it takes more than a month to write a code to successfully recognize over 90% accuracy, then the user will appreciate the advances in artificial intelligence. The artificial intelligence methodology can find a solution with over 97% accuracy in much shorter coder's time. Please refer to the following figure:

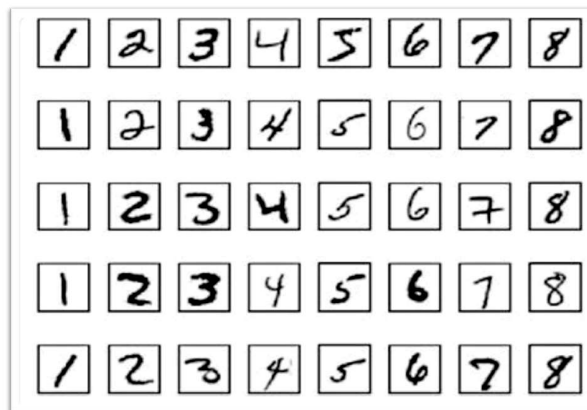


Figure 1.5: Sample images of MNIST dataset

In an artificial intelligence flow, the code is written once without analyzing a particular problem. The code uses already compiled libraries. Tensor flow library which is developed by Google is one such library. The user can scale the model with thousands to billions of variables which are also known as parameters. These variables are optimized during a process which is termed as a training aspect of machine learning. Thousands of data sets are required to train the system. Once the parameters are optimized, the test patterns are fed, and the classifications are checked. The test patterns are not part of the training set.

Changing the paradigm

As you may have noticed, code writing is automated at the expense of needing a lot of data for training. As the problem becomes more convoluted, it is not easy to write a traditional code even by a seasoned engineer. Writing a software which determines a face may be trivial to a seasoned engineer. However, writing a software which determines the age of the person without obvious clues, such as facial hair, is not trivial.

If the data is governed with simple laws or rules, then it is hard to justify use of artificial intelligence flow. For complex and obscure problems, such as age recognition, artificial intelligence is well suited.

Readers may be curious to know how AI programs are written. Let us consider the program which can distinguish between lines and circles, and then make it a more realistic problem where not all the points are strictly following one line or a circle. Let us consider the input image as shown in *Figure 1.6*:

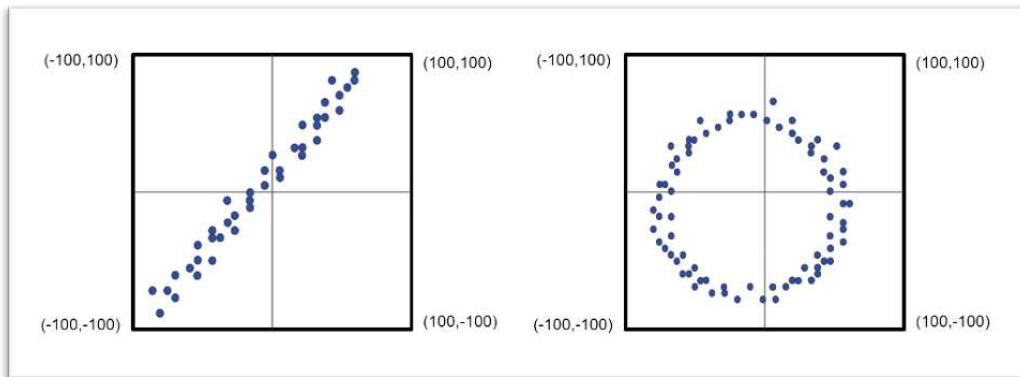


Figure 1.6: Input image where points do not follow one line or circle

As we know, it takes two points to define a line and three points to define a circle. Thus, we could use the same logic to define a line or circle, using specific chosen points. However, that will not be utilizing all the information, and results will also be different if different points are chosen, as shown in *Figure 1.7*:

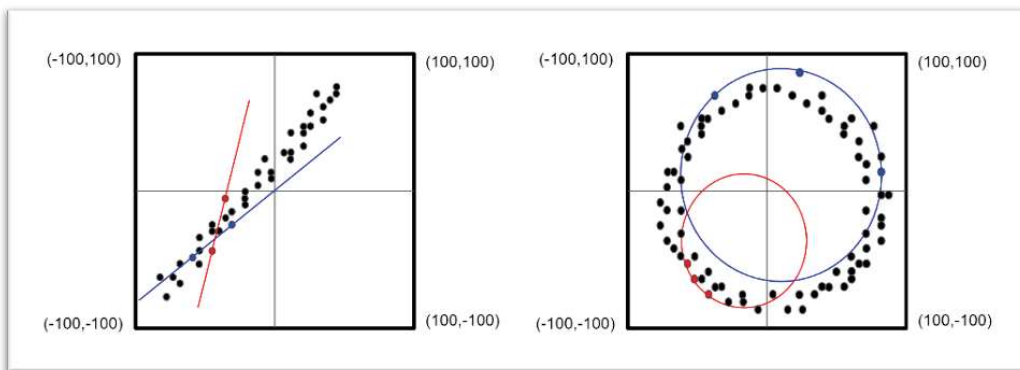


Figure 1.7: Multiple lines or circles can be estimated if only subset of the points used

For a robust solution, statistical regression methods should be used, which minimizes root mean square distance of all the points. All the data provided will be used, thus providing a robust solution. *Figure 1.8* illustrates the best fitting curve which is not dependent on few points but all the points:

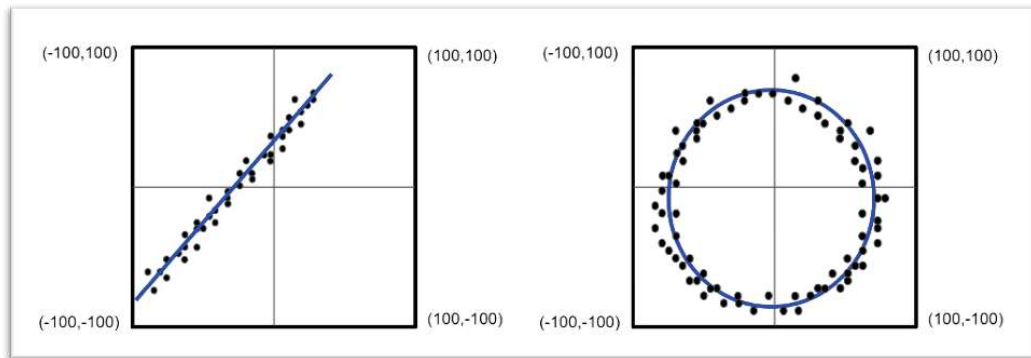


Figure 1.8: Using regression method to find best fitted line and circle

As you know, a line is specified as

$$Y = mX + c$$

where parameters, m and c are slope, and y intercept of the line respectively. Similarly, a circle can be written as

$$(X-X_0)^2 + (Y-Y_0)^2 = r^2$$

where three parameters, X_0, Y_0 and r define the circle. (X_0, Y_0) is the origin of the circle and r is the radius.

We can extrapolate to have several parameters to define a complex shape. The function can be defined with a set of multivariable linear equations which go through a non-linear function. We can cascade such linear and nonlinear functions to form a complex function.

To solve a generalized problem pattern recognition, the study of the biological brain has inspired a new type of processor. A biological brain contains many neuron cells which are connected to each other, making a network. It is believed that electrical signals pass through the neurons and eventually interpret a pattern. This processor is named as a neural network which indicates its origin. Let us look at the neural network and how it resembles the biological brain.

Neural Network

Most of us can guess people's age at a glance within reasonable accuracy. It comes effortlessly because this is how our brain works. Scientists got inspired from the anatomy of the biological brain, which is made of neurons. Neurons relate to multiple other neurons, which may seem a random connection. However, as the signal passes

through these neurons, living beings can make rational decisions. Refer to *Figure 1.9* for an illustration of the biological neuron:

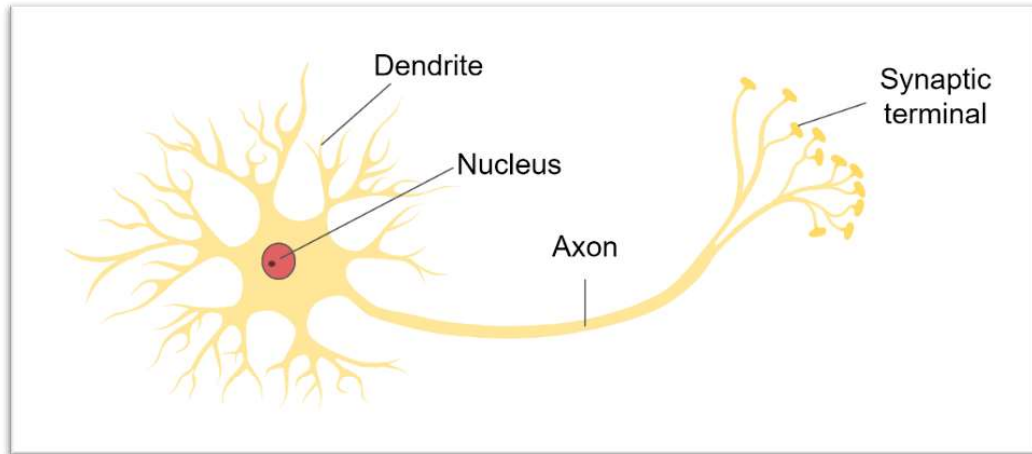


Figure 1.9: Illustration of a biological neuron

Figure 1.10 shows how multiple neurons are connected, thus forming a neural network. The bond between two neurons is called synaptic bond. It is believed that these synaptic bonds are being made over time. The synaptic bonds could have different connection strengths which would pass proportionate information. However, it may not be obvious how the simple connection between neurons suddenly would possess intelligence. A mathematical model made on similar principals are developed and that forms the basis of artificial intelligence. Please refer to the following figure:

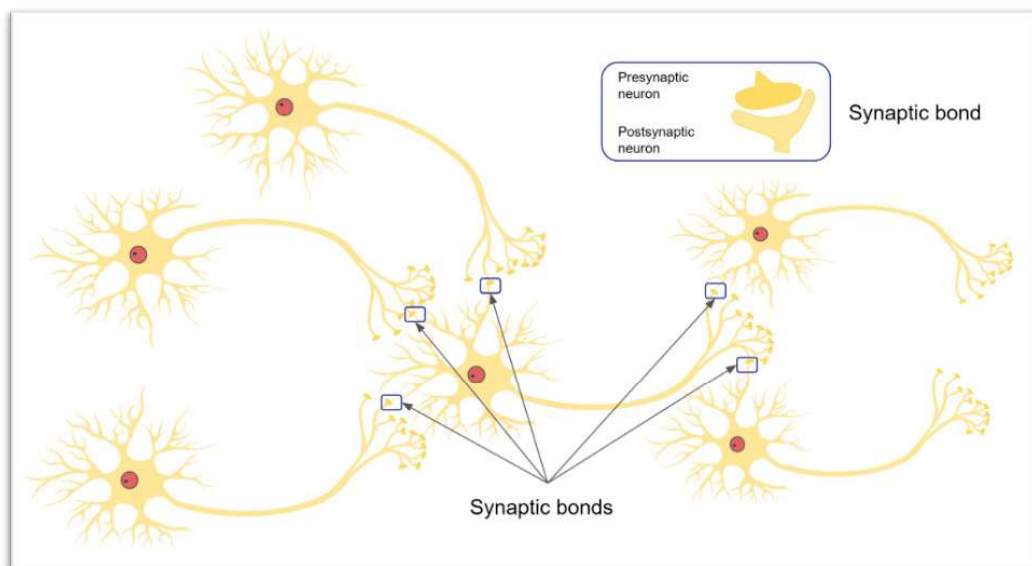


Figure 1.10: Illustration of a biological neurons forming neural networks with synaptic bonds

In a mathematical representation, the neuron simply sums the signals coming through the synaptic bonds and passes the signal to the next neuron. *Figure 1.11* draws a parallel between the biological neural network to a mathematical neural network. It shows how mathematical neurons are mimicking a biological neuron and how synapse connections are replaced by weights:

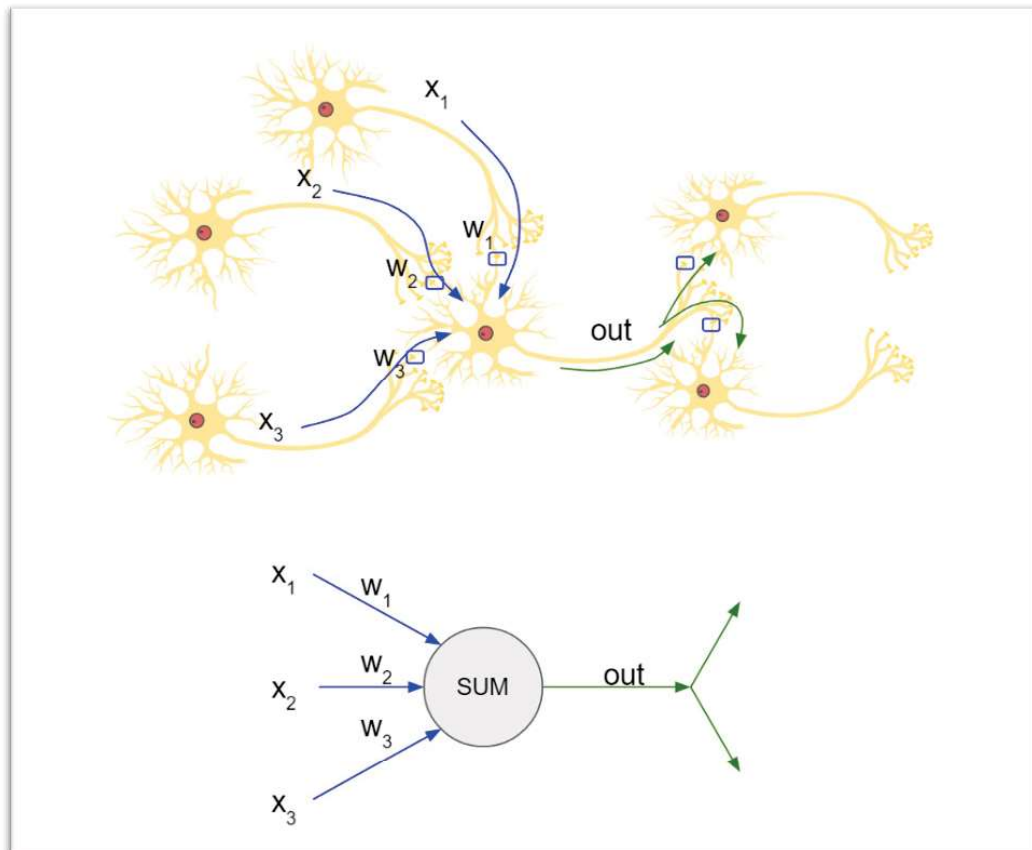


Figure 1.11: Parallel between biological neuron and mathematical neuron

After several tries over the years, many structured networks have been developed. The easiest network is defined as a fully connected neural network which is also termed as Dense Neural network. In a fully connected neural network, all inputs are applied to all neurons. The weight of the neuron is equivalent to the strength of the biological neuron. If there is no connection between two neurons, then the weight can be 0 in a mathematical neuron.

If one parameter is equated to one synapse of the human brain, then it is estimated that it will require several hundred trillion of parameters (reference Trillion). Let us approximate the number of parameters to be 1000 trillion parameters Assuming a

typical RAM of a computer is 8 Giga Bytes, a human brain is equivalent to 1,25,000 of such laptops. A typical data center can have 1 million to 10 million servers, which can be shown to have more capacity than one human brain. So, as of today, it is not impossible to mimic the human brain in a data center. It will be a while before full emulation of human brain will be economical and widely used.

However, the number of parameters used in artificial intelligence is growing at an exponential pace. A linear-log plot shows how the number of parameters has grown since 1952 to today, as shown in *Figure 1.12*:

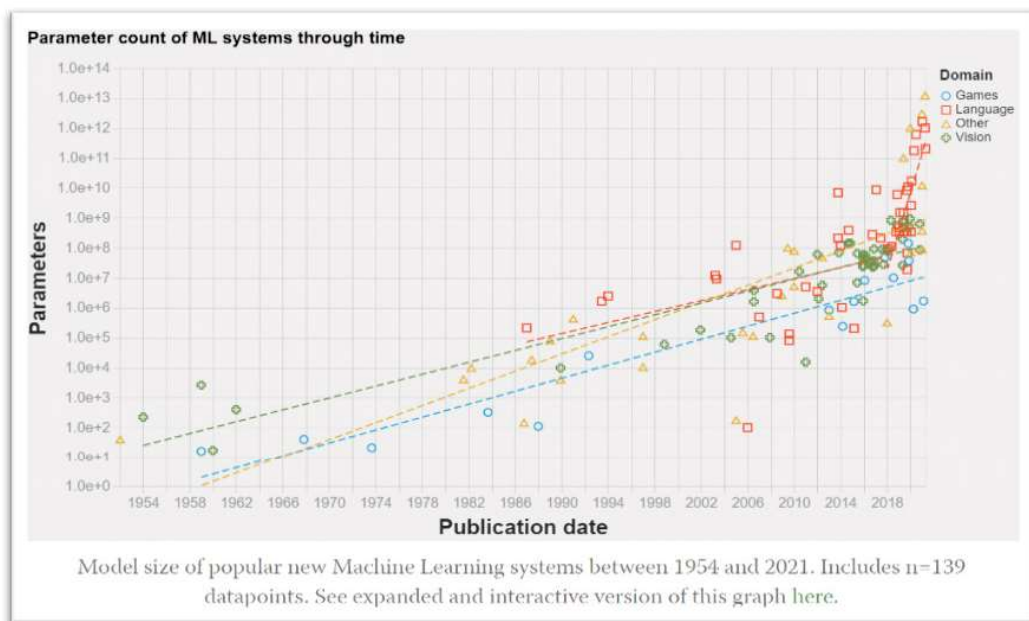


Figure 1.12: Model size of popular new Machine Learning systems between 1954 and 2021.

Even in the linear-log curve, the plot looks exponential towards the end, showing that the growth is faster than mathematical exponential growth. As of today, the highest parameter neural model in Google search shows a 175 billion parameter model, named as OpenAI LLC's GPT-3 natural language processing model (Reference GPT3). This model still is only 1/5000 of the human brain.

It is not sufficient to just have a neural network to recognize a pattern. Even in the biological world, it takes years to train the brain. Similarly, a neural network needs to be trained. As mentioned earlier, a neural network is defined with parameters. The parameters are like variables which are placeholders for a number. For different applications, the numbers will be different. The process of finding a set of these numbers comes under machine learning. Let us take a deeper look at how machines learn.