

Daniel Voigt Godoy

Deep Learning with PyTorch Step-by-Step



A Beginner's Guide



Deep Learning with PyTorch

Step-by-Step

A Beginner's Guide

Daniel Voigt Godoy

Version 1.1.1

Deep Learning with PyTorch Step-by-Step: A Beginner's Guide

by Daniel Voigt Godoy

Copyright © 2020-2022 by Daniel Voigt Godoy. All rights reserved.

May 2021: First Edition

Revision History for the First Edition:

- 2021-05-18: v1.0
- 2021-12-15: v1.1
- 2022-02-12: v1.1.1

This book is for sale at <http://leanpub.com/pytorch>. For more information, please send an email to contact@dvgodoy.com

Although the author has used his best efforts to ensure that the information and instructions contained in this book are accurate, under no circumstances shall the author be liable for any loss, damage, liability, or expense incurred or suffered as a consequence, directly or indirectly, of the use and/or application of any of the contents of this book. Any action you take upon the information in this book is strictly at your own risk. If any code samples or other technology this book contains or describes is subject to open source licenses or the intellectual property rights of others, it is your responsibility to ensure that your use thereof complies with such licenses and/or rights. The author does not have any control over and does not assume any responsibility for third-party websites or their content. All trademarks are the property of their respective owners. Screenshots are used for illustrative purposes only.

No part of this book may be reproduced or transmitted in any form or by any means (electronic, mechanical, photocopying, recording, or otherwise), or by any information storage and retrieval system without the prior written permission of the copyright owner, except where permitted by law. Please purchase only authorized electronic editions. Your support of the author's rights is appreciated.

"What I cannot create, I do not understand."

Richard P. Feynman

Table of Contents

Preface.....	xxiii
Acknowledgements.....	xxiv
About the Author.....	xxv
Frequently Asked Questions (FAQ)	1
Why PyTorch?	1
Why This Book?.....	2
Who Should Read This Book?.....	3
What Do I Need to Know?	4
How to Read This Book	4
What's Next?.....	7
Setup Guide.....	8
Official Repository.....	8
Environment	8
Google Colab.....	8
Binder	9
Local Installation.....	10
1. Anaconda	11
2. Conda (Virtual) Environments	11
3. PyTorch.....	13
4. TensorBoard.....	15
5. GraphViz and Torchviz (optional)	16
6. Git.....	17
7. Jupyter.....	19
Moving On.....	19
Part I: Fundamentals	21
Chapter 0: Visualizing Gradient Descent.....	22
Spoilers	22
Jupyter Notebook	22
Imports	23
Visualizing Gradient Descent.....	23
Model	24
Data Generation	25
Synthetic Data Generation.....	25

Train-Validation-Test Split	27
Step 0 - Random Initialization	28
Step 1 - Compute Model's Predictions	29
Step 2 - Compute the Loss	30
Loss Surface	32
Cross-Sections	36
Step 3 - Compute the Gradients	37
Visualizing Gradients	39
Backpropagation	40
Step 4 - Update the Parameters	41
Learning Rate	43
Low Learning Rate	44
High Learning Rate	46
Very High Learning Rate	47
"Bad" Feature	48
Scaling / Standardizing / Normalizing	51
Step 5 - Rinse and Repeat!	55
The Path of Gradient Descent	56
Recap	58
Chapter 1: A Simple Regression Problem	60
Spoilers	60
Jupyter Notebook	60
Imports	61
A Simple Regression Problem	61
Data Generation	62
Synthetic Data Generation	62
Gradient Descent	63
Step 0 - Random Initialization	64
Step 1 - Compute Model's Predictions	64
Step 2 - Compute the Loss	64
Step 3 - Compute the Gradients	65
Step 4 - Update the Parameters	66
Step 5 - Rinse and Repeat!	67
Linear Regression in Numpy	67
PyTorch	71
Tensor	71

Loading Data, Devices, and CUDA	76
Creating Parameters	81
Autograd	85
backward	85
grad	87
zero_	88
Updating Parameters	89
no_grad	92
Dynamic Computation Graph	92
Optimizer	96
step / zero_grad	97
Loss	99
Model	103
Parameters	105
state_dict	106
Device	107
Forward Pass	107
train	109
Nested Models	109
Sequential Models	112
Layers	113
Putting It All Together	115
Data Preparation	116
Model Configuration	117
Model Training	118
Recap	121
Chapter 2: Rethinking the Training Loop	123
Spoilers	123
Jupyter Notebook	123
Imports	123
Rethinking the Training Loop	124
Training Step	130
Dataset	134
TensorDataset	136
DataLoader	136
Mini-Batch Inner Loop	142

Random Split.....	145
Evaluation	147
Plotting Losses.....	151
TensorBoard	152
Running It Inside a Notebook.....	152
Running It Separately (Local Installation)	154
Running It Separately (Binder)	155
SummaryWriter	155
add_graph.....	157
add_scalars	158
Saving and Loading Models	164
Model State.....	164
Saving.....	164
Resuming Training.....	165
Deploying / Making Predictions	168
Setting the Model's Mode	169
Putting It All Together.....	170
Recap.....	173
Chapter 2.1: Going Classy.....	175
Spoilers	175
Jupyter Notebook	175
Imports	175
Going Classy.....	176
The Class.....	176
The Constructor	177
Arguments	177
Placeholders	178
Variables	180
Functions.....	180
Training Methods.....	187
Saving and Loading Models.....	191
Visualization Methods.....	192
The Full Code	193
Classy Pipeline	194
Model Training	197
Making Predictions	199

Checkpointing	199
Resuming Training	200
Putting It All Together	202
Recap	204
Chapter 3: A Simple Classification Problem	206
Spoilers	206
Jupyter Notebook	206
Imports	206
A Simple Classification Problem	207
Data Generation	208
Data Preparation	209
Model	210
Logits	211
Probabilities	212
Odds Ratio	212
Log Odds Ratio	214
From Logits to Probabilities	215
Sigmoid	217
Logistic Regression	218
Loss	221
BCELoss	223
BCEWithLogitsLoss	225
Imbalanced Dataset	228
Model Configuration	231
Model Training	232
Decision Boundary	236
Classification Threshold	241
Confusion Matrix	243
Metrics	245
True and False Positive Rates	245
Precision and Recall	248
Accuracy	249
Trade-offs and Curves	250
Low Threshold	250
High Threshold	252
ROC and PR Curves	253

The Precision Quirk	255
Best and Worst Curves	256
Comparing Models	257
Putting It All Together	259
Recap	261
Part II: Computer Vision	264
Chapter 4: Classifying Images	265
Spoilers	265
Jupyter Notebook	265
Imports	265
Classifying Images	266
Data Generation	267
Shape (NCHW vs NHWC)	271
Torchvision	274
Datasets	274
Models	274
Transforms	274
Transforms on Images	278
Transforms on Tensor	278
Normalize Transform	279
Composing Transforms	281
Data Preparation	283
Dataset Transforms	283
SubsetRandomSampler	285
Data Augmentation Transforms	288
WeightedRandomSampler	289
Seeds and more (seeds)	293
Putting It Together	295
Pixels as Features	296
Shallow Model	298
Notation	299
Model Configuration	300
Model Training	301
Deep-ish Model	301
Model Configuration	304
Model Training	304

Show Me the Math!	306
Show Me the Code!	308
Weights as Pixels	311
Activation Functions	312
Sigmoid	312
Hyperbolic Tangent (TanH)	314
Rectified Linear Unit (ReLU)	315
Leaky ReLU	317
Parametric ReLU (PReLU)	319
Deep Model	320
Model Configuration	321
Model Training	322
Show Me the Math Again!	324
Putting It All Together	326
Recap	330
Bonus Chapter: Feature Space	331
Two-Dimensional Feature Space	331
Transformations	332
A Two-Dimensional Model	333
Decision Boundary, Activation Style!	335
More Functions, More Boundaries	338
More Layers, More Boundaries	340
More Dimensions, More Boundaries	341
Recap	343
Chapter 5: Convolutions	344
Spoilers	344
Jupyter Notebook	344
Imports	344
Convolutions	345
Filter / Kernel	345
Convolving	347
Moving Around	348
Shape	351
Convolving in PyTorch	352
Striding	356
Padding	358

A REAL Filter	362
Pooling	364
Flattening	366
Dimensions	367
Typical Architecture	367
LeNet-5	368
A Multiclass Classification Problem	371
Data Generation	371
Data Preparation	372
Loss	375
Logits	375
Softmax	375
LogSoftmax	378
Negative Log-Likelihood Loss	378
Cross-Entropy Loss	382
Classification Losses Showdown!	384
Model Configuration	384
Model Training	387
Visualizing Filters and More!	388
Visualizing Filters	391
Hooks	394
Visualizing Feature Maps	403
Visualizing Classifier Layers	406
Accuracy	407
Loader Apply	409
Putting It All Together	410
Recap	414
Chapter 6: Rock, Paper, Scissors	416
Spoilers	416
Jupyter Notebook	416
Imports	416
Rock, Paper, Scissors... ..	417
Rock Paper Scissors Dataset	418
Data Preparation	419
ImageFolder	419
Standardization	420

The Real Datasets	424
Three-Channel Convolutions	425
Fancier Model	428
Dropout	431
Two-Dimensional Dropout	437
Model Configuration	438
Optimizer	438
Learning Rate	439
Model Training	439
Accuracy	440
Regularizing Effect	440
Visualizing Filters	442
Learning Rates	444
Finding LR	446
Adaptive Learning Rate	454
Moving Average (MA)	454
EWMA	455
EWMA Meets Gradients	461
Adam	462
Visualizing Adapted Gradients	463
Stochastic Gradient Descent (SGD)	470
Momentum	471
Nesterov	474
Flavors of SGD	475
Learning Rate Schedulers	478
Epoch Schedulers	479
Validation Loss Scheduler	480
Schedulers in StepByStep — Part I	482
Mini-Batch Schedulers	485
Schedulers in StepByStep — Part II	487
Scheduler Paths	489
Adaptive vs Cycling	492
Putting It All Together	492
Recap	495
Chapter 7: Transfer Learning	498
Spoilers	498

Jupyter Notebook	498
Imports	499
Transfer Learning	499
ImageNet	500
ImageNet Large Scale Visual Recognition Challenge (ILSVRC)	500
ILSVRC-2012	501
AlexNet (SuperVision Team)	501
ILSVRC-2014	501
VGG	502
Inception (GoogLeNet Team)	502
ILSVRC-2015	502
ResNet (MSRA Team)	503
Comparing Architectures	503
Transfer Learning in Practice	505
Pre-Trained Model	505
Adaptive Pooling	507
Loading Weights	508
Model Freezing	509
Top of the Model	510
Model Configuration	513
Data Preparation	513
Model Training	515
Generating a Dataset of Features	516
Top Model	519
Auxiliary Classifiers (Side-Heads)	521
1x1 Convolutions	524
Inception Modules	527
Batch Normalization	532
Running Statistics	535
Evaluation Phase	541
Momentum	542
BatchNorm2d	544
Other Normalizations	545
Small Summary	545
Residual Connections	546
Learning the Identity	546

The Power of Shortcuts	550
Residual Blocks	551
Putting It All Together	554
Fine-Tuning	555
Feature Extraction	556
Recap	558
Extra Chapter: Vanishing and Exploding Gradients	561
Spoilers	561
Jupyter Notebook	561
Imports	561
Vanishing and Exploding Gradients	562
Vanishing Gradients	562
Ball Dataset and Block Model	563
Weights, Activations, and Gradients	565
Initialization Schemes	567
Batch Normalization	570
Exploding Gradients	571
Data Generation & Preparation	571
Model Configuration & Training	572
Gradient Clipping	574
Value Clipping	575
Norm Clipping (or Gradient Scaling)	576
Model Configuration & Training	580
Clipping with Hooks	583
Recap	584
Part III: Sequences	586
Chapter 8: Sequences	587
Spoilers	587
Jupyter Notebook	587
Imports	588
Sequences	588
Data Generation	589
Recurrent Neural Networks (RNNs)	591
RNN Cell	594
RNN Layer	601
Shapes	604

Stacked RNN	607
Bidirectional RNN	611
Square Model	615
Data Generation	615
Data Preparation	616
Model Configuration	616
Model Training	618
Visualizing the Model	619
Transformed Inputs	619
Hidden States	620
The Journey of a Hidden State	622
Can We Do Better?	624
Gated Recurrent Units (GRUs)	625
GRU Cell	626
GRU Layer	634
Square Model II — The Quickening	635
Model Configuration & Training	636
Visualizing the Model	637
Hidden States	637
The Journey of a Gated Hidden State	638
Can We Do Better?	640
Long Short-Term Memory (LSTM)	640
LSTM Cell	641
LSTM Layer	649
Square Model III — The Sorcerer	650
Model Configuration & Training	651
Visualizing the Hidden States	652
Variable-Length Sequences	653
Padding	654
Packing	657
Unpacking (to padded)	661
Packing (from padded)	663
Variable-Length Dataset	664
Data Preparation	664
Collate Function	666
Square Model IV — Packed	667

Model Configuration & Training.....	669
1D Convolutions	670
Shapes	671
Multiple Features or Channels	672
Dilation	674
Data Preparation	676
Model Configuration & Training.....	676
Visualizing the Model.....	678
Putting It All Together.....	679
Fixed-Length Dataset.....	679
Variable-Length Dataset	680
There Can Be Only ONE ... Model	681
Model Configuration & Training.....	682
Recap.....	683
Chapter 9 — Part I: Sequence-to-Sequence.....	686
Spoilers.....	686
Jupyter Notebook	686
Imports	686
Sequence-to-Sequence.....	687
Data Generation.....	687
Encoder-Decoder Architecture	689
Encoder.....	689
Decoder.....	691
Teacher Forcing.....	696
Encoder + Decoder	698
Data Preparation	701
Model Configuration & Training.....	703
Visualizing Predictions	704
Can We Do Better?	704
Attention	705
"Values".....	708
"Keys" and "Queries"	708
Computing the Context Vector.....	710
Scoring Method.....	714
Attention Scores.....	716
Scaled Dot Product	717

Attention Mechanism	723
Source Mask	726
Decoder	728
Encoder + Decoder + Attention	730
Model Configuration & Training	732
Visualizing Predictions	733
Visualizing Attention	734
Multi-Headed Attention	735
Chapter 9 – Part II: Sequence-to-Sequence	740
Spoilers	740
Self-Attention	740
Encoder	741
Cross-Attention	746
Decoder	748
Subsequent Inputs and Teacher Forcing	750
Attention Scores	751
Target Mask (Training)	752
Target Mask (Evaluation/Prediction)	754
Encoder + Decoder + Self-Attention	758
Model Configuration & Training	762
Visualizing Predictions	763
Sequential No More	764
Positional Encoding (PE)	765
Encoder + Decoder + PE	778
Model Configuration & Training	780
Visualizing Predictions	781
Visualizing Attention	782
Putting It All Together	784
Data Preparation	784
Model Assembly	785
Encoder + Decoder + Positional Encoding	787
Self-Attention "Layers"	788
Attention Heads	790
Model Configuration & Training	792
Recap	793
Chapter 10: Transform and Roll Out	796

Spoilers	796
Jupyter Notebook	796
Imports	796
Transform and Roll Out	797
Narrow Attention	797
Chunking	798
Multi-Headed Attention	801
Stacking Encoders and Decoders	807
Wrapping "Sub-Layers"	808
Transformer Encoder	811
Transformer Decoder	816
Layer Normalization	821
Batch vs Layer	826
Our Seq2Seq Problem	828
Projections or Embeddings	829
The Transformer	831
Data Preparation	834
Model Configuration & Training	835
Visualizing Predictions	838
The PyTorch Transformer	838
Model Configuration & Training	844
Visualizing Predictions	845
Vision Transformer	846
Data Generation & Preparation	846
Patches	849
Rearranging	849
Embeddings	851
Special Classifier Token	853
The Model	857
Model Configuration & Training	859
Putting It All Together	861
Data Preparation	861
Model Assembly	861
1. Encoder-Decoder	863
2. Encoder	866
3. Decoder	867

4. Positional Encoding.....	868
5. Encoder "Layer"	869
6. Decoder "Layer"	870
7. "Sub-Layer" Wrapper.....	871
8. Multi-Headed Attention.....	873
Model Configuration & Training.....	875
Recap.....	876
Part IV: Natural Language Processing.....	879
Chapter 11: Down the Yellow Brick Rabbit Hole	880
Spoilers	880
Jupyter Notebook	880
Additional Setup.....	881
Imports	881
"Down the Yellow Brick Rabbit Hole"	883
Building a Dataset	883
Sentence Tokenization	885
HuggingFace's Dataset	891
Loading a Dataset.....	892
Attributes	893
Methods.....	894
Word Tokenization.....	896
Vocabulary.....	900
HuggingFace's Tokenizer	906
Before Word Embeddings	914
One-Hot Encoding (OHE).....	914
Bag-of-Words (BoW)	915
Language Models	916
N-grams.....	918
Continuous Bag-of-Words (CBoW)	919
Word Embeddings.....	919
Word2Vec	919
What Is an Embedding Anyway?.....	924
Pre-trained Word2Vec	927
Global Vectors (GloVe)	928
Using Word Embeddings	931
Vocabulary Coverage	931

Tokenizer	934
Special Tokens' Embeddings	935
Model I – GloVe + Classifier	937
Data Preparation	937
Pre-trained PyTorch Embeddings	939
Model Configuration & Training	941
Model II – GloVe + Transformer	942
Visualizing Attention	945
Contextual Word Embeddings	948
ELMo	949
BERT	957
Document Embeddings	959
Model III – Preprocessed Embeddings	962
Data Preparation	962
Model Configuration & Training	964
BERT	965
Tokenization	968
Input Embeddings	970
Pre-training Tasks	975
Masked Language Model (MLM)	975
Next Sentence Prediction (NSP)	978
Outputs	979
Model IV – Classifying Using BERT	984
Data Preparation	986
Model Configuration & Training	988
Fine-Tuning with HuggingFace	989
Sequence Classification (or Regression)	989
Tokenized Dataset	992
Trainer	994
Predictions	999
Pipelines	1001
More Pipelines	1002
GPT-2	1004
Putting It All Together	1008
Data Preparation	1008
"Packed" Dataset	1009

Model Configuration & Training 1012

Generating Text..... 1014

Recap..... 1016

Thank You!..... 1018

Preface

If you're reading this, I probably don't need to tell you that deep learning is amazing and PyTorch is cool, right?

But I will tell you, briefly, how this book came to be. In 2016, I started teaching a class on machine learning with Apache Spark and, a couple of years later, another class on the fundamentals of machine learning.

At some point, I tried to find a blog post that would visually explain, in a clear and concise manner, the concepts behind binary cross-entropy so that I could show it to my students. Since I could not find any that fit my purpose, I decided to write one myself. Although I thought of it as a fairly basic topic, it turned out to be my most popular blog post^[1]! My readers have welcomed the simple, straightforward, and conversational way I explained the topic.

Then, in 2019, I used the same approach for writing another blog post: "Understanding PyTorch with an example: a step-by-step tutorial."^[2] Once again, I was amazed by the reaction from the readers!

It was their positive feedback that motivated me to write this book to help beginners start their journey into deep learning and PyTorch. I hope you enjoy reading this book as much as I enjoyed writing it.

[1] <https://bit.ly/2UW5iTg>

[2] <https://bit.ly/2TpzwxR>

Acknowledgements

First and foremost, I'd like to thank YOU, my reader, for making this book possible. If it weren't for the amazing feedback I got from the thousands of readers of my blog post about PyTorch, I would have never mustered the strength to start *and finish* such a major undertaking as writing a 1,000-page book!

I'd like to thank my good friends Jesús Martínez-Blanco (who managed to read absolutely *everything* that I wrote), Jakub Cieslik, Hannah Berscheid, Mihail Vieru, Ramona Theresa Steck, Mehdi Belayet Lincon, and António Góis for helping me out and dedicating a good chunk of their time to reading, proofing, and suggesting improvements to my drafts. I'm forever grateful for your support! I'd also like to thank my friend José Luis Lopez Pino for the initial push I needed to actually *start* writing this book.

Many thanks to my friends José Quesada and David Anderson for taking me as a student at the Data Science Retreat in 2015 and, later on, for inviting me to be a teacher there. That was the starting point of my career both as a data scientist and as teacher.

I'd also like to thank the PyTorch developers for developing such an amazing framework, and the teams from Leanpub and Towards Data Science for making it incredibly easy for content creators like me to share their work with the community.

Finally, I'd like to thank my wife, Jerusa, for always being supportive throughout the entire writing of this book, and for taking the time to read *every* single page in it :-)

About the Author



Daniel is a data scientist, developer, writer, and teacher. He has been teaching machine learning and distributed computing technologies at Data Science Retreat, the longest-running Berlin-based bootcamp, since 2016, helping more than 150 students advance their careers.

Daniel is also the main contributor of two Python packages: [HandySpark](#) and [DeepReplay](#).

His professional background includes 20 years of experience working for companies in several industries: banking, government, fintech, retail, and mobility.

Frequently Asked Questions (FAQ)

Why PyTorch?

First, coding in PyTorch is **fun** :-). Really, there is something to it that makes it very enjoyable to write code in. Some say it is because it is very **pythonic**, or maybe there is something else, who knows? I hope that, by the end of this book, you feel like that too!

Second, maybe there are even some *unexpected benefits* to your health—check Andrej Karpathy’s [tweet](#)^[3] about it!

Jokes aside, PyTorch is the **fastest-growing**^[4] framework for developing deep learning models and it has a **huge ecosystem**.^[5] That is, there are many *tools* and *libraries* developed on top of PyTorch. It is the **preferred framework**^[6] in academia already and is making its way in the industry.

Several companies are already **powered by PyTorch**,^[7] to name a few:

- **Facebook:** The company is the original developer of PyTorch, released in October 2016.
- **Tesla:** Watch Andrej Karpathy (AI director at Tesla) speak about "how Tesla is using PyTorch to develop full self-driving capabilities for its vehicles" in [this video](#).^[8]
- **OpenAI:** In January 2020, OpenAI decided to standardize its deep learning framework on PyTorch ([source](#)^[9]).
- **fastai:** fastai is a [library](#)^[10] built on top of PyTorch to simplify model training and is used in its "[Practical Deep Learning for Coders](#)"^[11] course. The fastai library is deeply connected to PyTorch and "you can't become really proficient at using fastai if you don't know PyTorch well, too."^[12]
- **Uber:** The company is a significant contributor to PyTorch's ecosystem, having developed libraries like [Pyro](#)^[13] (probabilistic programming) and [Horovod](#)^[14] (a distributed training framework).
- **Airbnb:** PyTorch sits at the core of the company's dialog assistant for customer service.^{(source}^{[15])}

This book **aims to get you started with PyTorch** while giving you a **solid understanding of how it works**.

Why This Book?

If you're looking for a book where you can **learn about deep learning and PyTorch** without having to spend hours deciphering cryptic text and code, and one that's **easy and enjoyable to read**, this is it :-)

The book covers from the **basics of gradient descent** all the way up to **fine-tuning large NLP models** (BERT and GPT-2) using HuggingFace. It is divided into four parts:

- **Part I:** Fundamentals (gradient descent, training linear and logistic regressions in PyTorch)
- **Part II:** Computer Vision (deeper models and activation functions, convolutions, transfer learning, initialization schemes)
- **Part III:** Sequences (RNN, GRU, LSTM, seq2seq models, attention, self-attention, Transformers)
- **Part IV:** Natural Language Processing (tokenization, embeddings, contextual word embeddings, ELMo, BERT, GPT-2)

This is **not** a typical book: most tutorials *start* with some nice and pretty *image classification problem* to illustrate how to use PyTorch. It may seem cool, but I believe it **distracts** you from the **main goal**: learning **how PyTorch works**. In this book, I present a **structured, incremental, and from-first-principles** approach to learn PyTorch.

Moreover, this is **not** a **formal book** in any way: I am writing this book **as if I were having a conversation with you**, the reader. I will ask you **questions** (and give you answers shortly afterward), and I will also make (silly) **jokes**.

My job here is to make you **understand** the topic, so I will **avoid fancy mathematical notation** as much as possible and **spell it out in plain English**.

In this book, I will **guide** you through the **development** of many models in PyTorch, showing you why PyTorch makes it much **easier** and more **intuitive** to build models in Python: *autograd, dynamic computation graph, model classes*, and much, much more.

We will build, **step-by-step**, not only the models themselves but also your **understanding** as I show you both the **reasoning** behind the code and **how to avoid** some **common pitfalls** and **errors** along the way.

There is yet another advantage of **focusing on the basics**: this book is likely to have a **longer shelf life**. It is fairly common for technical books, especially those focusing on cutting-edge technology, to become outdated quickly. Hopefully, this is not going to be the case here, since the **underlying mechanics are not changing, and neither are the concepts**. It is expected that some syntax changes over time, but I do not see backward compatibility breaking changes coming anytime soon.



One more thing: If you hadn't noticed already, I **really** like to make use of **visual cues**, that is, **bold** and *italic* highlights. I firmly believe this helps the reader to **grasp** the **key ideas** I am trying to convey in a sentence more easily. You can find more on that in the section "**How to Read This Book**."

Who Should Read This Book?

I wrote this book for **beginners in general**—not only PyTorch beginners. Every now and then, I will spend some time explaining some **fundamental concepts** that I believe are **essential** to have a proper **understanding of what's going on in the code**.

The best example is **gradient descent**, which most people are familiar with at some level. Maybe you know its general idea, perhaps you've seen it in Andrew Ng's Machine Learning course, or maybe you've even **computed some partial derivatives yourself**!

In real life, the **mechanics** of gradient descent will be **handled automatically by PyTorch** (uh, spoiler alert!). But, I will walk you through it anyway (unless you choose to skip Chapter 0 altogether, of course), because lots of **elements in the code**, as well as **choices of hyper-parameters** (learning rate, mini-batch size, etc.), can be much more easily understood if you know **where they come from**.

Maybe you already know some of these concepts well: If this is the case, you can simply *skip* them, since I've made these explanations as *independent* as possible from the rest of the content.

But **I want to make sure everyone is on the same page**, so, if you have just heard about a given concept or if you are unsure if you have entirely understood it, these explanations are for you.

What Do I Need to Know?

This is a book for beginners, so I am assuming as **little prior knowledge** as possible; as mentioned in the previous section, I will take the time to explain fundamental concepts whenever needed.

That being said, this is what I expect from you, the reader:

- to be able to code in **Python** (if you are familiar with object-oriented programming [OOP], even better)
- to be able to work with PyData stack (**numpy**, **matplotlib**, and **pandas**) and **Jupyter notebooks**
- to be familiar with some basic concepts used in **machine learning**, like:
 - supervised learning: regression and classification
 - loss functions for regression and classification (mean squared error, cross-entropy, etc.)
 - training-validation-test split
 - underfitting and overfitting (bias-variance trade-off)

Even so, I am still briefly touching on **some** of these topics, but I need to draw a line somewhere; otherwise, this book would be gigantic!

How to Read This Book

Since this book is a **beginner's guide**, it is meant to be read **sequentially**, as ideas and concepts are progressively built. The same holds true for the **code** inside the book—you should be able to *reproduce* all outputs, provided you execute the chunks of code in the same order as they are introduced.

This book is **visually** different than other books: As I've mentioned already in the "**Why This Book?**" section, I **really** like to make use of **visual cues**. Although this is not, *strictly speaking*, a **convention**, this is how you can interpret those cues:

- I use **bold** to highlight what I believe to be the **most relevant words** in a sentence or paragraph, while *italicized* words are considered *important* too (not important enough to be bold, though)
- *Variables, coefficients, and parameters* in general, are *italicized*

- Classes and methods are written in a monospaced font, and they link to [PyTorch](#)^[16] documentation the first time they are introduced, so you can easily follow it (unlike other links in this book, links to documentation are *numerous* and thus *not* included in the footnotes)
- Every **code cell** is followed by *another* cell showing the corresponding **outputs** (if any)
- **All code** presented in the book is available at its **official repository** on GitHub:

<https://github.com/dvgodoy/PyTorchStepByStep>

Code cells with **titles** are an important piece of the workflow:

Title Goes Here

```
1 # Whatever is being done here is going to impact OTHER code
2 # cells. Besides, most cells have COMMENTS explaining what
3 # is happening
4 x = [1., 2., 3.]
5 print(x)
```

If there is any output to the code cell, titled or not, there *will* be another code cell depicting the corresponding **output** so you can *check* if you successfully reproduced it or not.

Output

```
[1.0, 2.0, 3.0]
```

Some code cells **do not** have titles—running them does not affect the workflow:

```
# Those cells illustrate HOW TO CODE something, but they are
# NOT part of the main workflow
dummy = ['a', 'b', 'c']
print(dummy[::-1])
```

But even these cells have their outputs shown!

```
['c', 'b', 'a']
```

I use asides to communicate a variety of things, according to the corresponding icon:



WARNING

Potential **problems** or things to **look out** for.



TIP

Key aspects I really want you to **remember**.



INFORMATION

Important information to **pay attention** to.



IMPORTANT

Really important information to **pay attention** to.



TECHNICAL

Technical aspects of **parameterization** or **inner workings of algorithms**.



QUESTION AND ANSWER

Asking myself **questions** (pretending to be you, the reader) and answering them, either in the same block or shortly after.



DISCUSSION

Really brief discussion on a concept or topic.



LATER

Important topics that will be covered in more detail later.



SILLY

Jokes, puns, memes, quotes from movies.

What's Next?

It's time to **set up** an environment for your learning journey using the **Setup Guide**.

- [3] <https://bit.ly/2MQoYRo>
- [4] <https://bit.ly/37uZgLB>
- [5] <https://pytorch.org/ecosystem/>
- [6] <https://bit.ly/2MTN0Lh>
- [7] <https://bit.ly/2UFHFve>
- [8] <https://bit.ly/2XXJkyo>
- [9] <https://openai.com/blog/openai-pytorch/>
- [10] <https://docs.fast.ai/>
- [11] <https://course.fast.ai/>
- [12] <https://course.fast.ai/>
- [13] <http://pyro.ai/>
- [14] <https://github.com/horovod/horovod>
- [15] <https://bit.ly/30CPhm5>
- [16] <https://bit.ly/3cT1aH2>

Setup Guide

Official Repository

This book's official repository is available on GitHub:

<https://github.com/dvgodoy/PyTorchStepByStep>

It contains **one Jupyter notebook** for every **chapter** in this book. Each notebook contains **all the code shown** in its corresponding chapter, and you should be able to **run its cells in sequence** to get the **same outputs**, as shown in the book. I strongly believe that being able to **reproduce the results** brings **confidence** to the reader.



Even though I did my best to ensure the **reproducibility** of the results, you may **still** find some minor differences in your outputs (especially during model training). Unfortunately, completely reproducible results are not guaranteed across PyTorch releases, and results may not be reproducible between CPU and GPU executions, even when using identical seeds.^[17]

Environment

There are **three options** for you to run the Jupyter notebooks:

- Google Colab (<https://colab.research.google.com>)
- Binder (<https://mybinder.org>)
- Local Installation

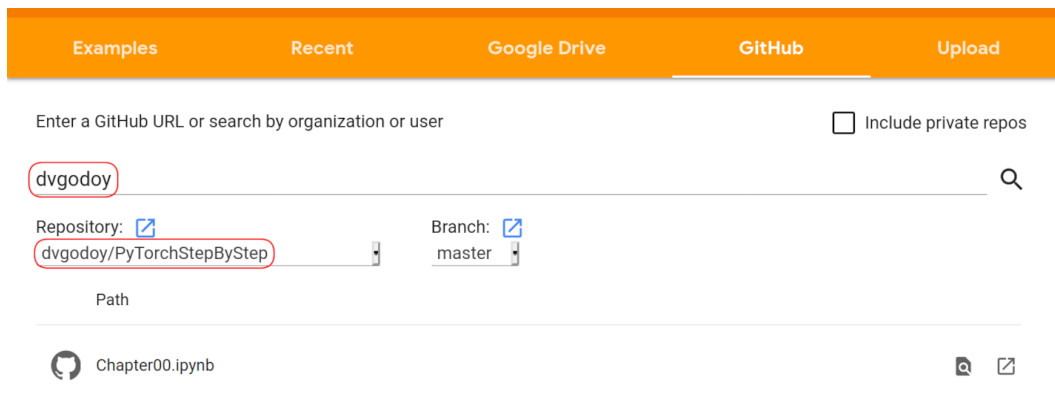
Let's briefly explore the **pros** and **cons** of each of these options.

Google Colab

Google Colab *"allows you to write and execute Python in your browser, with zero configuration required, free access to GPUs and easy sharing."*^[18]

You can easily **load notebooks directly from GitHub** using Colab's special URL (<https://colab.research.google.com/github/>). Just type in the GitHub's user or organization (like mine, dvgodoy), and it will show you a list of all its public repositories (like this book's, PyTorchStepByStep).

After choosing a repository, it will list the available notebooks and corresponding links to open them in a new browser tab.



The screenshot shows the Google Colab interface. At the top, there are tabs for 'Examples', 'Recent', 'Google Drive', 'GitHub' (which is selected), and 'Upload'. Below the tabs, there is a search bar with the text 'Enter a GitHub URL or search by organization or user'. To the right of the search bar is a checkbox labeled 'Include private repos'. Below the search bar, the text 'dvgodoy' is entered and highlighted with a red box. To the right of the search bar is a magnifying glass icon. Below the search bar, there are two dropdown menus. The first dropdown menu is labeled 'Repository:' and has 'dvgodoy/PyTorchStepByStep' selected, highlighted with a red box. The second dropdown menu is labeled 'Branch:' and has 'master' selected. Below the dropdown menus, there is a text input field labeled 'Path' which is empty. Below the 'Path' field, there is a list of notebooks. The first notebook is 'Chapter00.ipynb' with a Jupyter logo icon to its left. To the right of the notebook name are two icons: a magnifying glass and a share icon.

Figure S.1 - Google Colab's special URL

You also get access to a **GPU**, which is very useful to train deep learning models **faster**. More important, if you **make changes** to the notebook, Google Colab will **keep them**. The whole setup is very convenient; the only **cons** I can think of are:

- You need to be **logged in** to a Google account.
- You need to (re)install Python packages that are not part of Google Colab's default configuration.

Binder

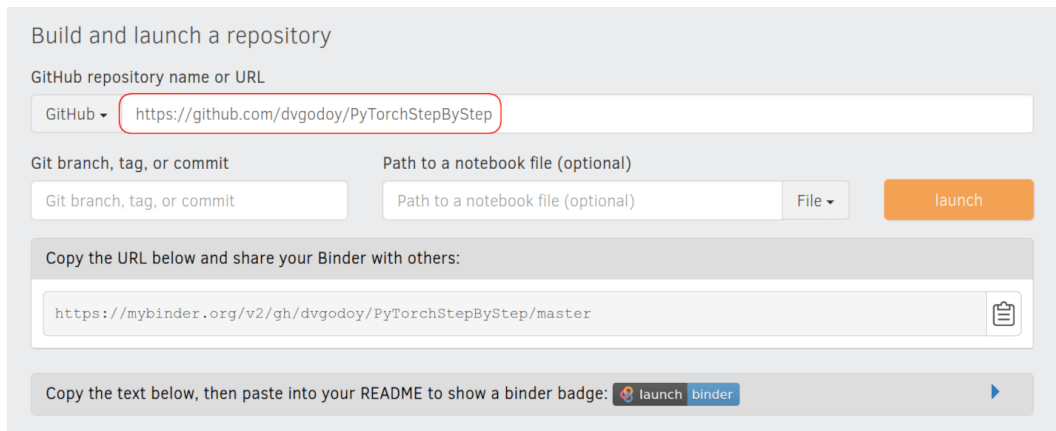
Binder "*allows you to create custom computing environments that can be shared and used by many remote users.*"^[12]

You can also **load notebooks directly from GitHub**, but the process is slightly different. Binder will create something like a *virtual machine* (technically, it is a container, but let's leave it at that), clone the repository, and start Jupyter. This allows you to have access to **Jupyter's home page** in your browser, just like you would if you were running it locally, but everything is running in a JupyterHub server on their end.

Just go to Binder's site (<https://mybinder.org/>) and type in the URL to the GitHub repository you want to explore (for instance, <https://github.com/dvgodoy/PyTorchStepByStep>) and click on **Launch**. It will take a couple of minutes to build the image and open Jupyter's home page.

You can also **launch Binder** for this book's repository directly using the following

link: <https://mybinder.org/v2/gh/dvgodoy/PyTorchStepByStep/master>.



The screenshot shows the Binder website interface. At the top, it says "Build and launch a repository". Below this, there's a section for "GitHub repository name or URL" with a dropdown menu set to "GitHub" and a text input field containing "https://github.com/dvgodoy/PyTorchStepByStep". To the right of this is a "launch" button. Below the input field, there are two more input fields: "Git branch, tag, or commit" and "Path to a notebook file (optional)". To the right of the second input field is a "File" dropdown menu and another "launch" button. Below these fields, there's a section titled "Copy the URL below and share your Binder with others:" with a text input field containing the URL "https://mybinder.org/v2/gh/dvgodoy/PyTorchStepByStep/master" and a copy icon. At the bottom, there's a section titled "Copy the text below, then paste into your README to show a binder badge:" with a small "launch binder" badge and a right arrow.

Figure S.2 - Binder's page

Binder is very convenient since it **does not require a prior setup** of any kind. Any Python packages needed to successfully run the environment are likely installed during launch (if provided by the author of the repository).

On the other hand, it may **take time** to start, and it **does not keep your changes** after your session expires (so, make sure you **download** any notebooks you modify).

Local Installation

This option will give you more **flexibility**, but it will require **more effort to set up**. I encourage you to try setting up your own environment. It may seem daunting at first, but you can surely accomplish it by following **seven easy steps**:

Checklist

- ☐ 1. Install **Anaconda**.
- ☐ 2. Create and activate a **virtual environment**.
- ☐ 3. Install **PyTorch** package.
- ☐ 4. Install **TensorBoard** package.
- ☐ 5. Install **GraphViz** software and **TorchViz** package (**optional**).
- ☐ 6. Install **git** and **clone** the repository.
- ☐ 7. Start **Jupyter** notebook.

1. Anaconda

If you don't have **Anaconda's Individual Edition**^[20] installed yet, this would be a good time to do it. It is a convenient way to start since it contains most of the Python libraries a data scientist will ever need to develop and train models.

Please follow the **installation instructions** for your OS:

- Windows (<https://docs.anaconda.com/anaconda/install/windows/>)
- macOS (<https://docs.anaconda.com/anaconda/install/mac-os/>)
- Linux (<https://docs.anaconda.com/anaconda/install/linux/>)



Make sure you choose **Python 3.X** version since Python 2 was discontinued in January 2020.

After installing Anaconda, it is time to create an **environment**.

2. Conda (Virtual) Environments

Virtual environments are a convenient way to isolate Python installations associated with different projects.



"What is an environment?"

It is pretty much a **replication of Python itself and some (or all) of its libraries**, so, effectively, you'll end up with multiple Python installations on your computer.



"Why can't I just use one single Python installation for everything?"

With so many independently developed Python **libraries**, each having many different **versions** and each version having various **dependencies** (on other libraries), **things can get out of hand** real quick.

It is beyond the scope of this guide to debate these issues, but take my word for it (or Google it!)—you'll benefit a great deal if you pick up the habit of **creating a different environment for every project you start working on**.



"How do I create an environment?"

First, you need to choose a **name** for your environment :-) Let's call ours

pytorchbook (or anything else you find easy to remember). Then, you need to open a **terminal** (in Ubuntu) or **Anaconda Prompt** (in Windows or macOS) and type the following command:

```
$ conda create -n pytorchbook anaconda
```

The command above creates a Conda environment named **pytorchbook** and includes **all Anaconda packages** in it (time to get a coffee, it will take a while...). If you want to learn more about creating and using Conda environments, please check Anaconda's "[Managing Environments](#)"^[21] user guide.

Did it finish creating the environment? Good! It is time to **activate it**, meaning, making **that Python installation** the one to be used now. In the same terminal (or Anaconda prompt), just type:

```
$ conda activate pytorchbook
```

Your prompt should look like this (if you're using Linux):

```
(pytorchbook)$
```

or like this (if you're using Windows):

```
(pytorchbook)C:\>
```

Done! You are using a **brand new Conda environment** now. You'll need to **activate it** every time you open a new terminal, or, if you're a Windows or macOS user, you can open the corresponding Anaconda prompt (it will show up as **Anaconda Prompt (pytorchbook)**, in our case), which will have it activated from the start.



IMPORTANT: From now on, I am assuming you'll activate the **pytorchbook** environment every time you open a terminal or Anaconda prompt. Further installation steps **must** be executed inside the environment.

3. PyTorch

PyTorch is the coolest **deep learning framework**, just in case you skipped the introduction.

It is *"an open source machine learning framework that accelerates the path from research prototyping to production deployment."*^[22] Sounds good, right? Well, I probably don't have to convince you at this point :-)

It is time to install the star of the show :-) We can go straight to the **Start Locally** (<https://pytorch.org/get-started/locally/>) section of PyTorch's website, and it will automatically select the options that best suit your local environment, and it will show you the **command to run**.

START LOCALLY

Select your preferences and run the install command. Stable represents the most currently tested and supported version of PyTorch. This should be suitable for many users. Preview is available if you want the latest, not fully tested and supported, 1.5 builds that are generated nightly. Please ensure that you have met the prerequisites below (e.g., numpy), depending on your package manager. Anaconda is our recommended package manager since it installs all dependencies. You can also [install previous versions of PyTorch](#). Note that LibTorch is only available for C++.

PyTorch Build	Stable (1.5.1)		Preview (Nightly)	
Your OS	Linux		Mac	Windows
Package	Conda	Pip	LibTorch	Source
Language	Python		C++ / Java	
CUDA	9.2	10.1	10.2	None
Run this Command:	<code>conda install pytorch torchvision cudatoolkit=10.2 -c pytorch</code>			

Figure S.3 - PyTorch's Start Locally page

Some of these options are given:

- PyTorch Build: Always select a **Stable** version.
- Package: I am assuming you're using **Conda**.
- Language: Obviously, **Python**.

So, two options remain: **Your OS** and **CUDA**.



"What is CUDA?" you ask.

Using GPU / CUDA

CUDA "is a parallel computing platform and programming model developed by NVIDIA for general computing on graphical processing units (GPUs)."^[23]

If you have a **GPU** in your computer (likely a GeForce graphics card), you can leverage its power to train deep learning models **much faster** than using a CPU. In this case, you should choose a PyTorch installation that includes CUDA support.

This is not enough, though: If you haven't done so yet, you need to install up-to-date drivers, the CUDA Toolkit, and the CUDA Deep Neural Network library (cuDNN). Unfortunately, more detailed installation instructions for CUDA are outside the scope of this book.

The **advantage** of using a GPU is that it allows you to **iterate faster** and **experiment with more-complex models and a more extensive range of hyper-parameters**.

In my case, I use **Linux**, and I have a **GPU** with CUDA version 10.2 installed. So I would run the following command in the **terminal** (after activating the environment):

```
(pytorchbook)$ conda install pytorch torchvision\
cudatoolkit=10.2 -c pytorch
```

Using CPU

If you **do not** have a **GPU**, you should choose **None** for CUDA.



"Would I be able to run the code **without** a GPU?" you ask.

Sure! The code and the examples in this book were designed to allow **all readers** to follow them promptly. Some examples may demand a bit more computing power, but we are talking about a **couple of minutes** in a CPU, not hours. If you do not have a GPU, **don't worry!** Besides, you can always use Google Colab if you need to use a GPU for a while!

If I had a **Windows** computer, and **no GPU**, I would have to run the following command in the **Anaconda prompt (pytorchbook)**:

```
(pytorchbook) C:\> conda install pytorch torchvision cpuonly\
-c pytorch
```

Installing CUDA

CUDA: Installing drivers for a GeForce graphics card, NVIDIA's cuDNN, and CUDA Toolkit can be challenging and is highly dependent on which model you own and which OS you use.

For installing GeForce's drivers, go to GeForce's website (<https://www.geforce.com/drivers>), select your OS and the model of your graphics card, and follow the installation instructions.

For installing NVIDIA's CUDA Deep Neural Network library (cuDNN), you need to register at <https://developer.nvidia.com/cudnn>.

For installing CUDA Toolkit (<https://developer.nvidia.com/cuda-toolkit>), please follow instructions for your OS and choose a local installer or executable file.

macOS: If you're a macOS user, please beware that PyTorch's binaries **DO NOT** support **CUDA**, meaning you'll need to install PyTorch **from source** if you want to use your GPU. This is a somewhat **complicated** process (as described in <https://github.com/pytorch/pytorch#from-source>), so, if you don't feel like doing it, you can choose to proceed **without CUDA**, and you'll still be able to execute the code in this book promptly.

4. TensorBoard

TensorBoard is TensorFlow's **visualization toolkit**, and *"provides the visualization and tooling needed for machine learning experimentation."*^[24]

TensorBoard is a powerful tool, and we can use it even if we are developing models in PyTorch. Luckily, you don't need to install the whole TensorFlow to get it; you can easily **install TensorBoard alone** using **Conda**. You just need to run this command in your **terminal** or **Anaconda prompt** (again, after activating the environment):

```
(pytorchbook)$ conda install -c conda-forge tensorboard
```


5. GraphViz and Torchviz (optional)



This step is optional, mostly because the installation of GraphViz can sometimes be *challenging* (especially on Windows). If for any reason you do not succeed in installing it correctly, or if you decide to skip this installation step, you will still be **able to execute the code in this book** (except for a couple of cells that generate images of a model's structure in the "Dynamic Computation Graph" section of Chapter 1).

GraphViz is an open source graph visualization software. It is "*a way of representing structural information as diagrams of abstract graphs and networks.*"^[25]

We need to install GraphViz to use **TorchViz**, a neat package that allows us to visualize a model's structure. Please check the **installation instructions** for your OS at <https://www.graphviz.org/download/>.



If you are using Windows, please use the **GraphViz's Windows Package** installer at <https://graphviz.gitlab.io/pages/Download/windows/graphviz-2.38.msi>.



You also need to add GraphViz to the PATH (environment variable) in Windows. Most likely, you can find the GraphViz executable file at C:\ProgramFiles(x86)\Graphviz2.38\bin. Once you find it, you need to set or change the PATH accordingly, adding GraphViz's location to it.

For more details on how to do that, please refer to "[How to Add to Windows PATH Environment Variable.](#)"^[26]

For additional information, you can also check the "[How to Install Graphviz Software](#)"^[27] guide.

After installing GraphViz, you can install the **torchviz**^[28] package. This package is **not** part of [Anaconda Distribution Repository](#)^[29] and is only available at **PyPI**^[30], the Python Package Index, so we need to `pip install` it.

Once again, open a **terminal** or **Anaconda prompt** and run this command (just once more: after activating the environment):

```
(pytorchbook)$ pip install torchviz
```

To check your GraphViz / TorchViz installation, you can try the Python code below:

```
(pytorchbook)$ python

Python 3.7.5 (default, Oct 25 2019, 15:51:11)
[GCC 7.3.0] :: Anaconda, Inc. on linux
Type "help", "copyright", "credits" or "license" for more
information.
>>> import torch
>>> from torchviz import make_dot
>>> v = torch.tensor(1.0, requires_grad=True)
>>> make_dot(v)
```

If everything is **working correctly**, you should see something like this:

Output

```
<graphviz.dot.Digraph object at 0x7ff540c56f50>
```

If you get an **error** of any kind (the one below is pretty common), it means there is still some kind of **installation issue** with GraphViz.

Output

```
ExecutableNotFound: failed to execute ['dot', '-Tsvg'], make
sure the Graphviz executables are on your systems' PATH
```

6. Git

It is way beyond this guide's scope to introduce you to version control and its most popular tool: git. If you are familiar with it already, great, you can skip this section altogether!

Otherwise, I'd recommend you to learn more about it; it will **definitely** be useful for you later down the line. In the meantime, I will show you the bare minimum so you can use git to **clone the repository** containing all code used in this book and get your **own, local copy** of it to modify and experiment with as you please.

First, you need to install it. So, head to its downloads page (<https://git-scm.com/downloads>) and follow instructions for your OS. Once the installation is complete, please open a **new terminal** or **Anaconda prompt** (it's OK to close the previous one). In the new terminal or Anaconda prompt, you should be able to **run git commands**.

To clone this book's repository, you only need to run:

```
(pytorchbook)$ git clone https://github.com/dvgodoy/\nPyTorchStepByStep.git
```

The command above will create a PyTorchStepByStep folder that contains a local copy of everything available on GitHub's repository.

conda install vs pip install

Although they may seem equivalent at first sight, you should **prefer conda install** over **pip install** when working with Anaconda and its virtual environments.

This is because **conda install** is sensitive to the active virtual environment: The package will be installed only for that environment. If you use **pip install**, and **pip** itself is not installed in the active environment, it will fall back to the **global pip**, and you definitely **do not** want that.

Why not? Remember the problem with **dependencies** I mentioned in the virtual environment section? That's why! The **conda** installer assumes it handles all packages that are part of its repository and keeps track of the complicated network of dependencies among them (to learn more about this, check this [link](#)^[31]).

To learn more about the differences between **conda** and **pip**, read "[Understanding Conda and Pip](#)."^[32]

As a rule, first try to **conda install** a given package and, only if it does not exist there, fall back to **pip install**, as we did with **torchviz**.

7. Jupyter

After cloning the repository, navigate to the PyTorchStepByStep folder and, **once inside it, start Jupyter** on your terminal or Anaconda prompt:

```
(pytorchbook)$ jupyter notebook
```

This will open your browser, and you will see **Jupyter's home page** containing the repository's notebooks and code.

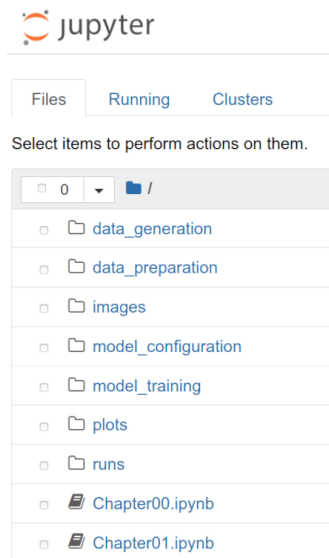


Figure S.4 - Running Jupyter

Moving On

Regardless of which of the three environments you chose, now you are ready to move on and tackle the development of your first PyTorch model, **step-by-step!**

[17] <https://pytorch.org/docs/stable/notes/randomness.html>

[18] <https://colab.research.google.com/notebooks/intro.ipynb>

[19] <https://mybinder.readthedocs.io/en/latest/>

[20] <https://www.anaconda.com/products/individual>

[21] <https://bit.ly/2MVk0CM>

[22] <https://pytorch.org/>

[23] <https://developer.nvidia.com/cuda-zone>

[24] <https://www.tensorflow.org/tensorboard>

[25] <https://www.graphviz.org/>

[26] <https://bit.ly/3flwYA5>

- [27] <https://bit.ly/30Ayct3>
- [28] <https://github.com/szagoruyko/pytorchviz>
- [29] <https://docs.anaconda.com/anaconda/packages/pkg-docs/>
- [30] <https://pypi.org/>
- [31] <https://bit.ly/37onBTt>
- [32] <https://bit.ly/2AAh8J5>

Part I

Fundamentals

Chapter 0

Visualizing Gradient Descent

Spoilers

In this chapter, we will:

- define a **simple linear regression model**
- walk through **every step of gradient descent**: initializing parameters, performing a forward pass, computing errors and loss, computing gradients, and updating parameters
- understand **gradients** using **equations**, **code**, and **geometry**
- understand the difference between **batch**, **mini-batch**, and **stochastic** gradient descent
- visualize the **effects on the loss** of using different **learning rates**
- understand the importance of **standardizing / scaling features**
- and much, much more!

There is **no** actual PyTorch code in this chapter... it is *Numpy* all along because our focus here is to understand, inside and out, how gradient descent works. PyTorch will be introduced in the next chapter.

Jupyter Notebook

The Jupyter notebook corresponding to [Chapter 0^{\[33\]}](#) is part of the official *Deep Learning with PyTorch Step-by-Step* repository on GitHub. You can also run it directly in [Google Colab^{\[34\]}](#).

If you're using a *local installation*, open your terminal or Anaconda prompt and navigate to the PyTorchStepByStep folder you cloned from GitHub. Then, *activate* the pytorchbook environment and run `jupyter notebook`:

```
$ conda activate pytorchbook
```

```
(pytorchbook)$ jupyter notebook
```

If you're using Jupyter's default settings, [this link](#) should open Chapter 0's

notebook. If not, just click on Chapter00.ipynb on your Jupyter's home page.

Imports

For the sake of organization, all libraries needed throughout the code used in any given chapter are imported at its very beginning. For this chapter, we'll need the following imports:

```
import numpy as np
from sklearn.linear_model import LinearRegression
from sklearn.preprocessing import StandardScaler
```

Visualizing Gradient Descent



According to [Wikipedia^{\[35\]}](#): "**Gradient descent** is a first-order iterative optimization algorithm for finding a local minimum of a differentiable function."

But I would go with: "**Gradient descent** is an iterative technique commonly used in machine learning and deep learning to find the best possible set of parameters / coefficients for a given model, data points, and loss function, starting from an initial, and usually random, guess."



"Why **visualizing** gradient descent?"

I believe *the way gradient descent is usually explained lacks intuition*. Students and beginners are left with a *bunch of equations and rules of thumb*—**this is not the way one should learn such a fundamental topic**.

If you **really understand** how gradient descent works, you will also understand how the **characteristics of your data** and your **choice of hyper-parameters** (mini-batch size and learning rate, for instance) have an **impact** on how **well** and how **fast** the model training is going to be.

By *really understanding*, I do not mean working through the equations manually: this does not develop intuition either. I mean **visualizing** the effects of different settings; I mean **telling a story** to illustrate the concept. That's how you **develop intuition**.

That being said, I'll cover the **five basic steps** you'd need to go through to use gradient descent. I'll show you the corresponding *Numpy* code while explaining lots of **fundamental concepts** along the way.

But first, we need some **data** to work with. Instead of using some *external dataset*, let's

- define which **model** we want to train to better understand gradient descent; and
- **generate synthetic data** for that model.

Model

The model must be **simple** and **familiar**, so you can focus on the **inner workings** of gradient descent.

So, I will stick with a model as simple as it can be: a **linear regression with a single feature, x** !

$$y = b + wx + \epsilon$$

Equation 0.1 - Simple linear regression model

In this model, we use a **feature (x)** to try to predict the value of a **label (y)**. There are three elements in our model:

- **parameter b** , the *bias* (or *intercept*), which tells us the expected average value of y when x is zero
- **parameter w** , the *weight* (or *slope*), which tells us how much y increases, on average, if we increase x by one unit
- and that **last term** (why does it *always* have to be a Greek letter?), *epsilon*, which is there to account for the inherent **noise**; that is, the **error** we cannot get rid of

We can also conceive the very same model structure in a less abstract way:

salary = minimum wage + increase per year * years of experience + noise

And to make it even more concrete, let's say that the **minimum wage** is **\$1,000** (whatever the currency or time frame, this is not important). So, if you have **no experience**, your salary is going to be the **minimum wage** (parameter b).

Also, let's say that, **on average**, you get a **\$2,000 increase** (parameter w) for every year of experience you have. So, if you have **two years of experience**, you are expected to earn a salary of **\$5,000**. But your actual salary is **\$5,600** (lucky you!). Since the model cannot account for those **extra \$600**, your extra money is, technically speaking, **noise**.

Data Generation

We know our model already. In order to generate **synthetic data** for it, we need to pick values for its **parameters**. I chose $b = 1$ and $w = 2$ (as in thousands of dollars) from the example above.

First, let's generate our **feature (x)**: We use *Numpy's* `rand()` method to randomly generate 100 (N) points between 0 and 1.

Then, we plug our **feature (x)** and our **parameters b and w** into our **equation** to compute our **labels (y)**. But we need to add some **Gaussian noise**^[36] (**epsilon**) as well; otherwise, our synthetic dataset would be a perfectly straight line. We can generate noise using *Numpy's* `randn()` method, which draws samples from a normal distribution (of mean 0 and variance 1), and then multiply it by a **factor** to adjust for the **level of noise**. Since I don't want to add too much noise, I picked 0.1 as my factor.

Synthetic Data Generation

Data Generation

```
1 true_b = 1
2 true_w = 2
3 N = 100
4
5 # Data Generation
6 np.random.seed(42)
7 x = np.random.rand(N, 1)
8 epsilon = (.1 * np.random.randn(N, 1))
9 y = true_b + true_w * x + epsilon
```

Did you notice the `np.random.seed(42)` at line 6? This line of code is actually more important than it looks. It guarantees that, every time we run this code, the **same random numbers will be generated**.



"Wait; what?! Aren't the numbers supposed to be **random**? How could they possibly be the **same** numbers?" you ask, perhaps even a bit annoyed by this.

(Not So) Random Numbers

Well, you know, random numbers are not **quite** random... They are really **pseudo-random**, which means *Numpy's* number generator spits out a **sequence of numbers** that **looks like it's random**. But it is not, really.

The **good** thing about this behavior is that we can tell the generator to **start a particular sequence of pseudo-random numbers**. To some extent, it works as if we tell the generator: *"please generate sequence #42,"* and it will spill out a sequence of numbers. That number, 42, which works like the *index* of the sequence, is called a **seed**. Every time we give it the **same seed**, it generates the **same numbers**.

This means we have the **best of both worlds**: On the one hand, we do **generate** a sequence of numbers that, for all intents and purposes, is **considered to be random**; on the other hand, we have the **power to reproduce any given sequence**. I cannot stress enough how convenient that is for **debugging** purposes!

Moreover, you can guarantee that **other people will be able to reproduce your results**. Imagine how annoying it would be to run the code in this book and get different outputs every time, having to wonder if there is anything wrong with it. But since I've set a seed, you and I can achieve the very same outputs, even if it involved generating random data!

Next, let's **split** our synthetic data into **train** and **validation** sets, shuffling the array of indices and using the first 80 shuffled points for training.



"Why do you need to **shuffle** randomly generated data points? Aren't they random enough?"

Yes, they **are** random enough, and shuffling them is indeed redundant in this example. But it is best practice to **always shuffle** your data points before training a model to improve the performance of gradient descent.



There is one **exception** to the "always shuffle" rule, though: **time series** problems, where shuffling can lead to data leakage.

Train-Validation-Test Split

It is beyond the scope of this book to explain the reasoning behind the **train-validation-test split**, but there are two points I'd like to make:

1. The split should **always** be the **first thing** you do—no preprocessing, no transformations; **nothing happens before the split**. That's why we do this **immediately after the synthetic data generation**.
2. In this chapter we will use **only the training set**, so I did not bother to create a **test set**, but I performed a split nonetheless to **highlight point #1 :-)**

Train-Validation Split

```
1 # Shuffles the indices
2 idx = np.arange(N)
3 np.random.shuffle(idx)
4
5 # Uses first 80 random indices for train
6 train_idx = idx[:int(N*.8)]
7 # Uses the remaining indices for validation
8 val_idx = idx[int(N*.8):]
9
10 # Generates train and validation sets
11 x_train, y_train = x[train_idx], y[train_idx]
12 x_val, y_val = x[val_idx], y[val_idx]
```



"Why didn't you use `train test split()` from Scikit-Learn?" you may be asking.

That's a fair point. Later on, we will refer to the **indices of the data points** belonging to either train or validation sets, instead of the points themselves. So, I thought of using them from the very start.

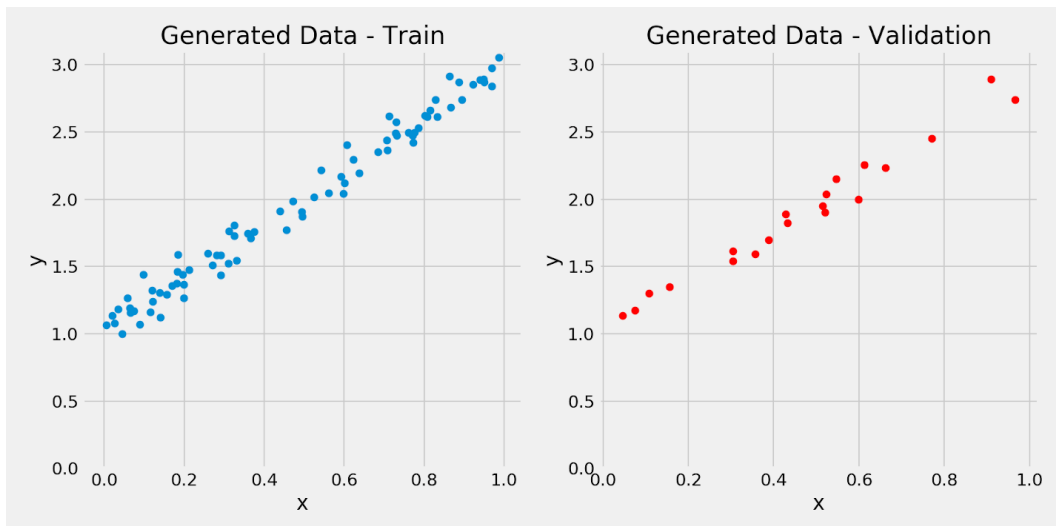


Figure 0.1 - Synthetic data: train and validation sets

We know that $b = 1$, $w = 2$, but now let's see how close we can get to the true values by using **gradient descent** and the 80 points in the **training set** (for training, $N = 80$).

Step 0 - Random Initialization

In our example, we already **know** the **true values** of the **parameters**, but this will obviously never happen in real life: If we *knew* the true values, why even bother to train a model to find them?!

OK, given that **we'll never know** the **true values** of the parameters, we need to set **initial values** for them. How do we choose them? It turns out a **random guess** is as good as any other.



Even though the initialization is **random**, there are some clever **initialization schemes** that should be used when training more-complex models. We'll get back to them (much) later.

For training a model, you need to **randomly initialize the parameters / weights** (we have only two, b and w).

Random Initialization

```
1 # Step 0 - Initializes parameters "b" and "w" randomly
2 np.random.seed(42)
3 b = np.random.randn(1)
4 w = np.random.randn(1)
5
6 print(b, w)
```

Output

```
[0.49671415] [-0.1382643]
```

Step 1 - Compute Model's Predictions

This is the **forward pass**; it simply *computes the model's predictions using the current values of the parameters / weights*. At the very beginning, we will be producing **really bad predictions**, as we started with **random values in Step 0**.

Step 1

```
1 # Step 1 - Computes our model's predicted output - forward pass
2 yhat = b + w * x_train
```

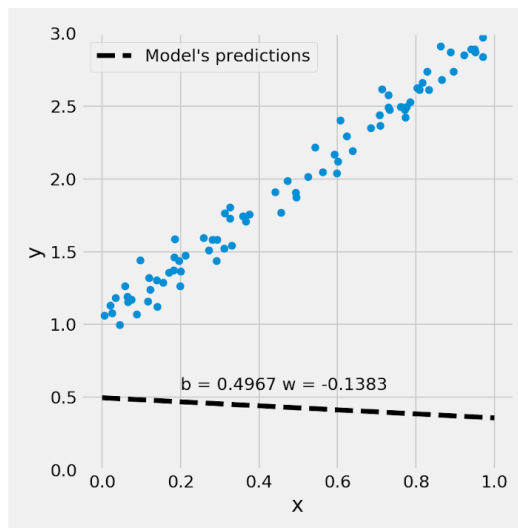


Figure 0.2 - Model's predictions (with random parameters)

Step 2 - Compute the Loss

There is a subtle but fundamental difference between **error** and **loss**.

The **error** is the difference between the **actual value (label)** and the **predicted value** computed for a single data point. So, for a given i -th point (from our dataset of N points), its error is:

$$\text{error}_i = \hat{y}_i - y_i$$

Equation 0.2 - Error

The error of the **first point** in our dataset ($i = 0$) can be represented like this:

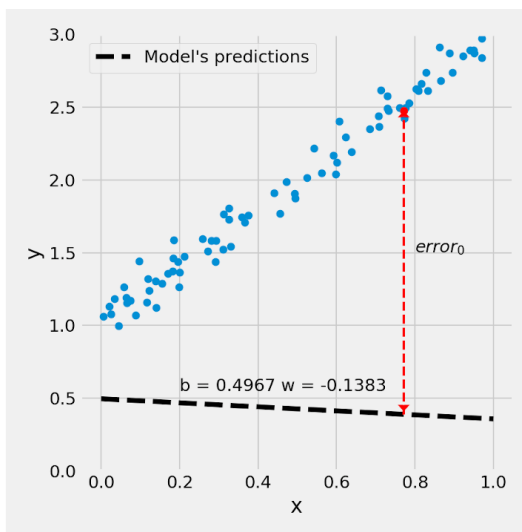


Figure 0.3 - Prediction error (for one data point)

The **loss**, on the other hand, is some sort of **aggregation of errors** for a **set of data points**.

It seems rather obvious to compute the loss for **all** (N) data points, right? Well, yes and no. Although it will surely yield a **more stable path** from the initial **random parameters** to the parameters that **minimize the loss**, it will also surely be **slow**.

This means one *needs to sacrifice (a bit of) stability for the sake of speed*. This is easily accomplished by randomly choosing (*without replacement*) a subset of n out of N data points each time we compute the loss.



Batch, Mini-batch, and Stochastic Gradient Descent

- If we use **all points** in the training set ($n = N$) to compute the loss, we are performing a **batch** gradient descent;
- If we were to use a **single point** ($n = 1$) each time, it would be a **stochastic** gradient descent;
- Anything else (n) **in between 1 and N** characterizes a **mini-batch** gradient descent;

For a regression problem, the **loss** is given by the **mean squared error (MSE)**; that is, the average of all squared errors; that is, the average of all squared differences between **labels** (y) and **predictions** ($b + wx$).

$$\begin{aligned}\text{MSE} &= \frac{1}{n} \sum_{i=1}^n \text{error}_i^2 \\ &= \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (b + wx_i - y_i)^2\end{aligned}$$

Equation 0.3 - Loss: mean squared error (MSE)

In the code below, we are using **all data points** of the training set to compute the **loss**, so $n = N = 80$, meaning we are indeed performing **batch gradient descent**.

Step 2

```
1 # Step 2 - Computing the loss
2 # We are using ALL data points, so this is BATCH gradient
3 # descent. How wrong is our model? That's the error!
4 error = (yhat - y_train)
5
6 # It is a regression, so it computes mean squared error (MSE)
7 loss = (error ** 2).mean()
8
9 print(loss)
```


Output

```
2.7421577700550976
```

Loss Surface

We have just computed the **loss** (2.74) corresponding to our **randomly initialized parameters** ($b = 0.49$ and $w = -0.13$). What if we did the same for **ALL** possible values of b and w ? Well, not *all* possible values, but *all combinations of evenly spaced values in a given range*, like:

```
# Reminder:
# true_b = 1
# true_w = 2

# we have to split the ranges in 100 evenly spaced intervals each
b_range = np.linspace(true_b - 3, true_b + 3, 101)
w_range = np.linspace(true_w - 3, true_w + 3, 101)
# meshgrid is a handy function that generates a grid of b and w
# values for all combinations
bs, ws = np.meshgrid(b_range, w_range)
bs.shape, ws.shape
```

Output

```
((101, 101), (101, 101))
```

The result of the `meshgrid()` operation was two (101, 101) matrices representing the values of each parameter inside a grid. What does one of these matrices look like?

```
bs
```

Output

```
array([[ -2.    ,  -1.94,  -1.88, ...,   3.88,   3.94,   4.    ],
       [ -2.    ,  -1.94,  -1.88, ...,   3.88,   3.94,   4.    ],
       [ -2.    ,  -1.94,  -1.88, ...,   3.88,   3.94,   4.    ],
       ...,
       [ -2.    ,  -1.94,  -1.88, ...,   3.88,   3.94,   4.    ],
       [ -2.    ,  -1.94,  -1.88, ...,   3.88,   3.94,   4.    ],
       [ -2.    ,  -1.94,  -1.88, ...,   3.88,   3.94,   4.    ]])
```

Sure, we're somewhat *cheating* here, since we *know* the **true** values of b and w , so we can choose the **perfect ranges** for the parameters. But it is for educational purposes only :-)

Next, we could use those values to compute the corresponding **predictions, errors, and losses**. Let's start by taking a **single data point** from the training set and computing the predictions for every combination in our grid:

```
dummy_x = x_train[0]
dummy_yhat = bs + ws * dummy_x
dummy_yhat.shape
```

Output

```
(101, 101)
```

Thanks to its broadcasting capabilities, *Numpy* is able to understand we want to multiply the **same x value** by **every entry** in the **ws matrix**. This operation resulted in a **grid of predictions** for that **single data point**. Now we need to do this for **every one of our 80 data points** in the training set.

We can use *Numpy's* `apply_along_axis()` to accomplish this:



Look ma, no loops!

```
all_predictions = np.apply_along_axis(
    func1d=lambda x: bs + ws * x,
    axis=1,
    arr=x_train,
)
all_predictions.shape
```

Output

```
(80, 101, 101)
```

Cool! We got **80 matrices** of shape (101, 101), **one matrix for each data point**, each matrix containing a **grid of predictions**.

The **errors** are the difference between the predictions and the labels, but we cannot perform this operation right away—we need to work a bit on our **labels (y)**, so they have the proper **shape** for it (broadcasting is good, but not *that* good):

```
all_labels = y_train.reshape(-1, 1, 1)
all_labels.shape
```

Output

```
(80, 1, 1)
```

Our **labels** turned out to be **80 matrices of shape (1, 1)**—the most boring kind of matrix—but that is enough for broadcasting to work its magic. We can compute the **errors** now:

```
all_errors = (all_predictions - all_labels)
all_errors.shape
```

Output

```
(80, 101, 101)
```

Each prediction has its own error, so we get **80 matrices** of shape (101, 101), again,

one matrix for each data point, each matrix containing a **grid of errors**.

The only step missing is to compute the **mean squared error**. First, we take the square of all errors. Then, we **average the squares over all data points**. Since our data points are in the **first dimension**, we use `axis=0` to compute this average:

```
all_losses = (all_errors ** 2).mean(axis=0)
all_losses.shape
```

Output

```
(101, 101)
```

The result is a **grid of losses**, a matrix of shape (101, 101), **each loss** corresponding to a **different combination of the parameters b and w** .

These losses are our **loss surface**, which can be visualized in a 3D plot, where the vertical axis (z) represents the loss values. If we **connect** the combinations of b and w that yield the **same loss value**, we'll get an **ellipse**. Then, we can draw this ellipse in the original $b \times w$ plane (in blue, for a loss value of 3). This is, in a nutshell, what a **contour plot** does. From now on, we'll always use the contour plot, instead of the corresponding 3D version.

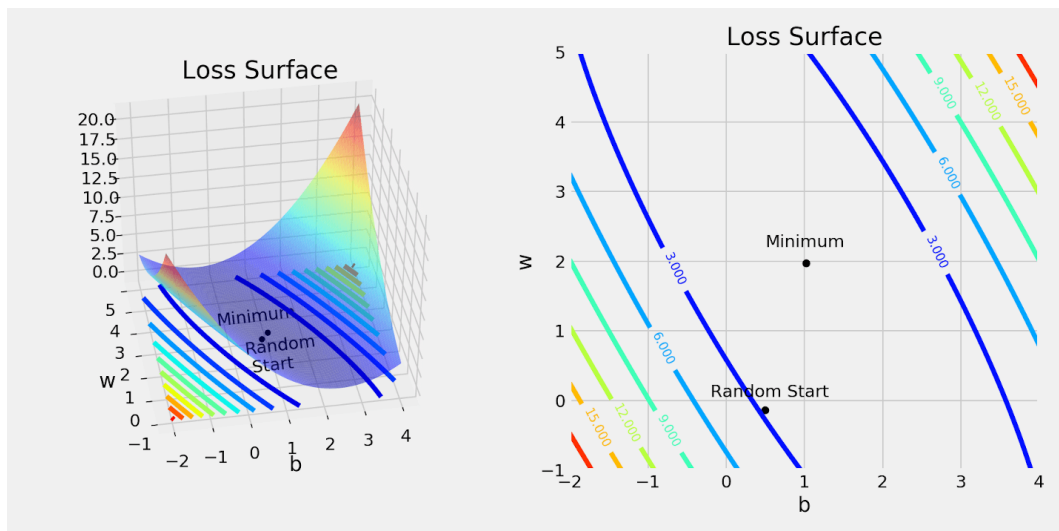


Figure 0.4 - Loss surface

In the center of the plot, where parameters (b, w) have values close to (1, 2), the loss is at its **minimum** value. This is the point we're trying to reach using gradient

descent.

In the bottom, slightly to the left, there is the **random start** point, corresponding to our randomly initialized parameters.

This is one of the nice things about tackling a simple problem like a linear regression with a single feature: We have only **two parameters**, and thus **we can compute and visualize the loss surface**.



Unfortunately, for the absolute majority of problems, **computing the loss surface is not going to be feasible**: we have to rely on gradient descent's ability to reach a point of minimum, even if it starts at some random point.

Cross-Sections

Another nice thing is that we can cut a **cross-section** in the loss surface to check what the **loss** would look like if the **other parameter were held constant**.

Let's start by making $b = 0.52$ (the value from `b_range` that is closest to our initial random value for b , 0.4967). We cut a cross-section *vertically* (the red dashed line) on our loss surface (left plot), and we get the resulting plot on the right:

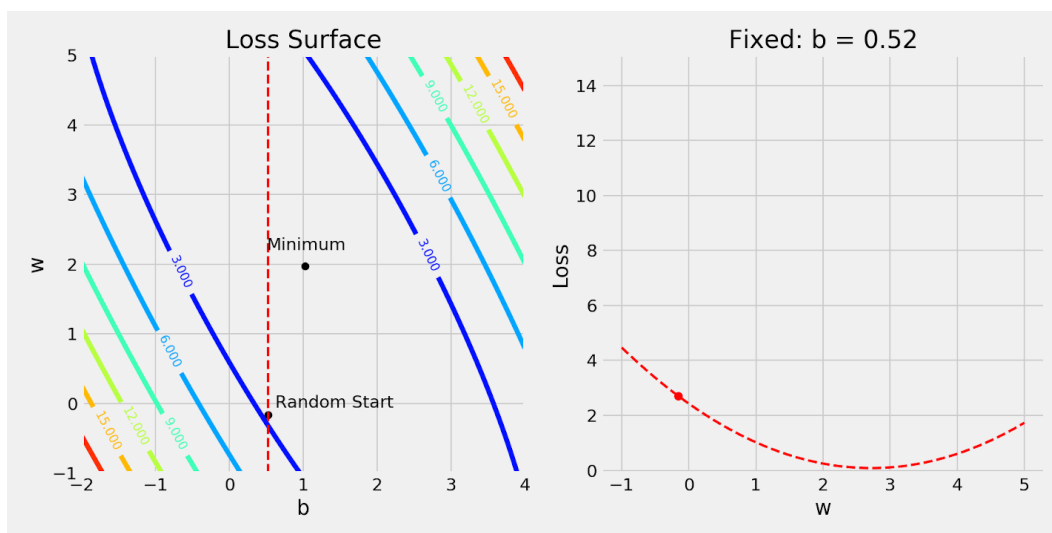
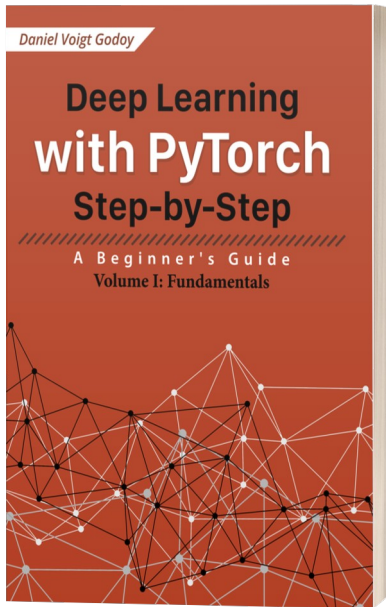


Figure 0.5 - Vertical cross-section; parameter b is fixed

What does this cross-section tell us? It tells us that, **if we keep b constant** (at 0.52), the **loss**, seen from the **perspective of parameter w** , can be minimized if **w gets increased** (up to some value between 2 and 3).



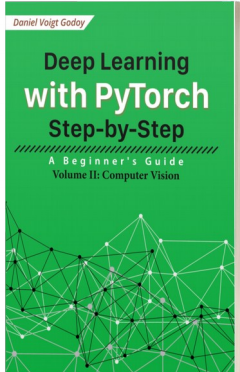
KEEP ON READING

Get the FULL version on

[Kindle](#) for \$9.95
[PDF](#) for \$14.95
[Paperback](#) for \$24.95

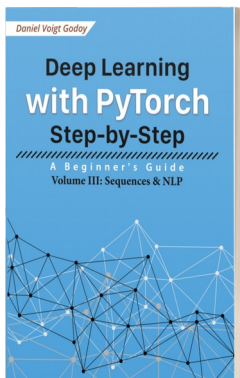
KEEP ON LEARNING

Take your skills to the next level!



Tackle Computer Vision tasks using convolutional neural networks and transfer learning in the second volume!

[Kindle](#) [PDF](#) [Paperback](#)



Discover how to handle sequences and Natural Language Processing tasks using recurrent neural networks, Transformers, and HuggingFace in the third volume!

[Kindle](#) [PDF](#) [Paperback](#)