

Venelin Valkov

# Get SH\*T Done with PyTorch



Solve Real-World Machine Learning Problems

# Get SH\*T Done with PyTorch

Solve Real-World Machine Learning Problems

Venelin Valkov

This book is for sale at <http://leanpub.com/getting-things-done-with-pytorch>

This version was published on 2020-06-27



This is a [Leanpub](#) book. Leanpub empowers authors and publishers with the Lean Publishing process. [Lean Publishing](#) is the act of publishing an in-progress ebook using lightweight tools and many iterations to get reader feedback, pivot until you have the right book and build traction once you do.

© 2020 Venelin Valkov

# Contents

1.	<b>Getting Started with PyTorch</b>	1
	PyTorch & NumPy	1
	Tensors	3
	Running on GPU	5
	Common Issues	6
	Conclusion	6
	References	7
2.	<b>Build Your First Neural Network with PyTorch</b>	8
	Data	9
	Data Preprocessing	10
	Building a Neural Network	12
	Training	17
	Saving the model	21
	Evaluation	21
	Conclusion	23
	References	24
3.	<b>Transfer Learning for Image Classification using Torchvision</b>	25
	Recognizing traffic signs	26
	Building a dataset	29
	Using a pre-trained model	33
	Adding class “unknown”	46
	Summary	53
	References	54
4.	<b>Time Series Forecasting with LSTMs for Daily Coronavirus Cases</b>	55
	Novel Coronavirus (COVID-19)	55
	Daily Cases Dataset	57
	Data exploration	57
	Preprocessing	60
	Building a model	63
	Training	64
	Predicting daily cases	67

## CONTENTS

Use all data for training . . . . .	69
Predicting future cases . . . . .	70
Conclusion . . . . .	73
References . . . . .	74
<b>5. Time Series Anomaly Detection using LSTM Autoencoders . . . . .</b>	<b>75</b>
Data . . . . .	75
Exploratory Data Analysis . . . . .	77
LSTM Autoencoder . . . . .	79
Anomaly Detection in ECG Data . . . . .	80
Training . . . . .	84
Saving the model . . . . .	86
Choosing a threshold . . . . .	87
Evaluation . . . . .	88
Summary . . . . .	91
References . . . . .	91
<b>6. Face Detection on Custom Dataset with Detectron2 . . . . .</b>	<b>92</b>
Detectron 2 . . . . .	93
Face Detection Data . . . . .	96
Data Preprocessing . . . . .	96
Face Detection with Detectron 2 . . . . .	101
Evaluating Object Detection Models . . . . .	105
Finding Faces in Images . . . . .	106
Conclusion . . . . .	110
References . . . . .	111
<b>7. Create Dataset for Sentiment Analysis by Scraping Google Play App Reviews . . . . .</b>	<b>112</b>
Setup . . . . .	112
The Goal of the Dataset . . . . .	113
Scraping App Information . . . . .	114
Scraping App Reviews . . . . .	118
Summary . . . . .	120
References . . . . .	120
<b>8. Sentiment Analysis with BERT and Transformers by Hugging Face . . . . .</b>	<b>121</b>
What is BERT? . . . . .	121
Setup . . . . .	123
Data Exploration . . . . .	124
Data Preprocessing . . . . .	127
Sentiment Classification with BERT and Hugging Face . . . . .	135
Evaluation . . . . .	143
Summary . . . . .	149
References . . . . .	150

## CONTENTS

<b>9. Deploy BERT for Sentiment Analysis as REST API using FastAPI . . . . .</b>	<b>151</b>
Project setup . . . . .	151
Building a skeleton REST API . . . . .	152
Adding our model . . . . .	153
Putting everything together . . . . .	157
Testing the API . . . . .	157
Summary . . . . .	158
References . . . . .	159
<b>10. Object Detection on Custom Dataset with YOLO (v5) . . . . .</b>	<b>160</b>
Prerequisites . . . . .	160
Build a dataset . . . . .	161
Fine-tuning YOLO v5 . . . . .	170
Evaluation . . . . .	172
Making predictions . . . . .	173
Summary . . . . .	175
References . . . . .	175

# 1. Getting Started with PyTorch

PyTorch<sup>1</sup> is:

An open source machine learning framework that accelerates the path from research prototyping to production deployment.

In my humble opinion, PyTorch is the *sweet*<sup>2</sup> way to solve Machine Learning problems, in the real world! The vast community allows you to work state-of-the-art models and deploy them to production in no time (relatively speaking). Let's get started!

```
1 In [0]: !pip install -q -U torch watermark
2
3 In [0]: %load_ext watermark
4     %watermark -v -p numpy,torch
5
6 Out[0]: CPython 3.6.9
7         IPython 5.5.0
8
9         numpy 1.17.5
10        torch 1.4.0
```

## PyTorch □ NumPy

Do you know NumPy? If you do, learning PyTorch will be a breeze! If you don't, prepare to learn the skills that will guide you on your journey Machine Learning Mastery!

Let's start with something simple:

---

<sup>1</sup><https://pytorch.org/>

<sup>2</sup><https://www.youtube.com/watch?v=-h4spfuMGDI>

```
1 In [0]: import torch
2     import numpy as np
3
4 In [0]: a = np.array([1, 2])
5     b = np.array([8, 9])
6
7     c = a + b
8
9
10 Out[0]: array([ 9, 11])
```

Adding the same arrays with PyTorch looks like this:

```
1 In [0]: a = torch.tensor([1, 2])
2     b = torch.tensor([8, 9])
3
4     c = a + b
5
6
7 Out[0]: tensor([ 9, 11])
```

Fortunately, you can go from NumPy to PyTorch:

```
1 In [0]: a = torch.tensor([1, 2])
2
3     a.numpy()
4
5 Out[0]: array([1, 2])
```

and vice versa:

```
1 In [0]: a = np.array([1, 2])
2     torch.from_numpy(a)
3
4 Out[0]: tensor([1, 2])
```

The good news is that the conversions incur almost no cost on the performance of your app. The NumPy and PyTorch store data in memory in the same way. That is, PyTorch is reusing the work done by NumPy.

## Tensors

Tensors are just n-dimensional number (including booleans) containers. You can find the complete list of supported data types at [PyTorch's Tensor Docs<sup>3</sup>](#).

So, how can you create a Tensor (try to ignore that I've already shown you how to do it)?

```
1 In [0]: torch.tensor([[1, 2], [2, 1]])  
2  
3 Out[0]: tensor([[1, 2],  
4                  [2, 1]])
```

You can create a tensor from floats:

```
1 In [0]: torch.FloatTensor([[1, 2], [2, 1]])  
2  
3 Out[0]: tensor([[1., 2.],  
4                  [2., 1.]])
```

Or define the type like so:

```
1 In [0]: torch.tensor([[1, 2], [2, 1]], dtype=torch.bool)  
2  
3 Out[0]: tensor([[True, True],  
4                  [True, True]])
```

You can use a wide range of factory methods to create Tensors without manually specifying each number. For example, you can create a matrix with random numbers like this:

```
1 In [0]: torch.rand(3, 2)  
2  
3 Out[0]: tensor([[0.6686, 0.7622],  
4                  [0.0341, 0.5835],  
5                  [0.2423, 0.0651]])
```

Or one full of ones:

---

<sup>3</sup><https://pytorch.org/docs/stable/tensors.html>

```
1 In [0]: torch.ones(3, 2)
2
3 Out[0]: tensor([[1., 1.],
4                  [1., 1.],
5                  [1., 1.]])
```

PyTorch has a variety of useful operations:

```
1 In [0]: x = torch.tensor([[2, 3], [1, 2]])
2     print(x)
3     print(f'sum: {x.sum()}')
4
5 Out[0]: tensor([[2, 3],
6                  [1, 2]])
7     sum: 8
```

Get the transpose of a 2-D tensor:

```
1 In [0]: x.t()
2
3 Out[0]: tensor([[2, 1],
4                  [3, 2]])
```

Get the shape of each dimension:

```
1 In [0]: x.size()
2
3 Out[0]: torch.Size([2, 2])
```

Generally, performing some operation creates a new Tensor:

```
1 In [0]: y = torch.tensor([[2, 2], [5, 1]])
2     z = x.add(y)
3     z
4
5 Out[0]: tensor([[4, 5],
6                  [6, 3]])
```

But you can do it in-place:

```

1 In [0]: x.add_(y)
2         x
3
4 Out[0]: tensor([[4, 5],
5                 [6, 3]])

```

Almost all operations have an in-place version - the name of the operation, followed by an underscore.

## Running on GPU

At this point, you might be like: “Why do I need PyTorch at all? All of this is perfectly doable with NumPy?”. PyTorch has three major superpowers: - you can run your operations on the GPU(s) (or something else) - [Autograd: automatic differentiation](#)<sup>4</sup> - A set of tools to build Neural Networks. Including several additional packages for working [with text](#)<sup>5</sup> or [images](#)<sup>6</sup>.

Doing your Deep Learning computations on the GPU speeds up your experiment by a lot! And PyTorch makes it ridiculously easy to do it. Let’s start by checking if GPU is available:

```

1 In [0]: device = torch.device("cuda") if torch.cuda.is_available() else torch.device\
2 ("cpu")
3         device
4
5 Out[0]: device(type='cuda')

```

Good, we have a [CUDA](#)<sup>7</sup>-enabled GPU device on our hands. Let’s store a Tensor on it:

```

1 In [0]: x = torch.tensor([[2, 3], [1, 2]])
2         x.to(device)
3
4 Out[0]: tensor([[2, 3],
5                  [1, 2]]), device='cuda:0'

```

Notice that our Tensor is now on device `cuda:0`. What can we do with it? Pretty much everything as before:

---

<sup>4</sup>[https://pytorch.org/tutorials/beginner/blitz/autograd\\_tutorial.html](https://pytorch.org/tutorials/beginner/blitz/autograd_tutorial.html)

<sup>5</sup><https://github.com/pytorch/text>

<sup>6</sup><https://github.com/pytorch/vision>

<sup>7</sup><https://en.wikipedia.org/wiki/CUDA>

```

1 In [0]: x = x.to(device)
2
3     y = torch.tensor([[2, 2], [5, 1]])
4     y = y.to(device)
5
6     x.add(y)
7
8 Out[0]: tensor([[4, 5],
9                  [6, 3]], device='cuda:0')

```

## Common Issues

I got to be honest with you. You will fuck up, multiple times, before understanding how this whole thing works out. That's alright!

However, there are a couple of things you can do that might minimize the frustrations along your journey:

- Doing operations between GPU and CPU Tensors is **not allowed**
- Size mismatch between Tensors occurs often and is (almost every time) easy to fix:

In [0]: a = torch.ones(2, 2) b = torch.ones(1, 3) a \* b

Out[0]: ----- RuntimeError Traceback (most recent call last) <ipython-input-21-b2a4e8765762> in <module>() 1 a = torch.ones(2, 2) 2 b = torch.ones(1, 3) --> 3 a \* b

```

1      RuntimeError: The size of tensor a (2) must match the size of tensor b (3) at \
2 non-singleton dimension 1

```

PyTorch is very descriptive in this case. When doing more complex stuff, you would want to check the shape of your Tensors obsessively, after every operation. Just print the size!

- Running out of GPU memory: You might be leaking memory or too large of a dataset/model. Faster/better GPU always helps. But remember, you can solve really large problems with a single powerful GPU these days. Think carefully if that is not enough for you - why that is?

## Conclusion

Welcome to the dark side! You might've been working with Keras, TensorFlow, or another Deep Learning framework, until recently. Almost every framework is great, but PyTorch has really solid roots. Easy to use and understand, allows for fast experimentation and standard debugging tools apply! Enjoy!

## References

- “PyTorch: A Modern Library for Machine Learning” with Adam Paszke<sup>8</sup>
- Recitation 1 | Your First Deep Learning Code<sup>9</sup>
- What is PyTorch?<sup>10</sup>

---

<sup>8</sup><https://www.youtube.com/watch?v=5bSAipCNqXo>

<sup>9</sup>[https://www.youtube.com/watch?v=KrCp\\_yPVOxs](https://www.youtube.com/watch?v=KrCp_yPVOxs)

<sup>10</sup>[https://pytorch.org/tutorials/beginner/blitz/tensor\\_tutorial.html](https://pytorch.org/tutorials/beginner/blitz/tensor_tutorial.html)

# 2. Build Your First Neural Network with PyTorch

In this tutorial, you'll build your first Neural Network using PyTorch. You'll use it to predict whether or not it is going to rain tomorrow using real weather information.

- Run the complete notebook in your browser (Google Colab)<sup>1</sup>
- Read the Getting Things Done with Pytorch book<sup>2</sup>

You'll learn how to:

- Preprocess CSV files and convert the data to Tensors
- Build your own Neural Network model with PyTorch
- Use a loss function and an optimizer to train your model
- Evaluate your model and learn about the perils of imbalanced classification

```
1 %reload_ext watermark
2 %watermark -v -p numpy,pandas,torch
```

```
1 CPython 3.6.9
2 IPython 5.5.0
3
4 numpy 1.17.5
5 pandas 0.25.3
6 torch 1.4.0
```

---

<sup>1</sup>[https://colab.research.google.com/drive/1lDXVkd7GC8jK\\_nGmOMKeDywXse-DY-u](https://colab.research.google.com/drive/1lDXVkd7GC8jK_nGmOMKeDywXse-DY-u)

<sup>2</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>

```

1 import torch
2
3 import os
4 import numpy as np
5 import pandas as pd
6 from tqdm import tqdm
7 import seaborn as sns
8 from pylab import rcParams
9 import matplotlib.pyplot as plt
10 from matplotlib import rc
11 from sklearn.model_selection import train_test_split
12 from sklearn.metrics import confusion_matrix, classification_report
13
14 from torch import nn, optim
15
16 import torch.nn.functional as F
17
18 %matplotlib inline
19 %config InlineBackend.figure_format='retina'
20
21 sns.set(style='whitegrid', palette='muted', font_scale=1.2)
22
23 HAPPY_COLORS_PALETTE = ["#01BEFE", "#FFDD00", "#FF7D00", "#FF006D", "#93D30C", "#8F0\
24 OFF"]
25
26 sns.set_palette(sns.color_palette(HAPPY_COLORS_PALETTE))
27
28 rcParams['figure.figsize'] = 12, 8
29
30 RANDOM_SEED = 42
31 np.random.seed(RANDOM_SEED)
32 torch.manual_seed(RANDOM_SEED)

```

## Data

Our dataset contains daily weather information from multiple Australian weather stations. We're about to answer a simple question. *Will it rain tomorrow?*

The data is hosted on [Kaggle](#)<sup>3</sup> and created by [Joe Young](#)<sup>4</sup>. I've uploaded the dataset to Google Drive. Let's get it:

---

<sup>3</sup><https://www.kaggle.com/jsphyg/weather-dataset-rattle-package>

<sup>4</sup><https://www.kaggle.com/jsphyg>

```
1 !gdown --id 1Q1wUptbNDYdfizk5abhmoFxIQiX19Tn7
```

And load it into a data frame:

```
1 df = pd.read_csv('weatherAUS.csv')
```

We have a large set of features/columns here. You might also notice some *Nans*. Let's have a look at the overall dataset size:

```
1 df.shape
```

```
1 (142193, 24)
```

Looks like we have plenty of data. But we got to do something about those missing values.

## Data Preprocessing

We'll start by simplifying the problem by removing most of the data. We'll use only 4 columns for predicting whether or not it is going to rain tomorrow:

```
1 cols = ['Rainfall', 'Humidity3pm', 'Pressure9am', 'RainToday', 'RainTomorrow']
2
3 df = df[cols]
```

Neural Networks don't work with much else than numbers. We'll convert *yes* and *no* to 1 and 0, respectively:

```
1 df['RainToday'].replace({'No': 0, 'Yes': 1}, inplace = True)
2 df['RainTomorrow'].replace({'No': 0, 'Yes': 1}, inplace = True)
```

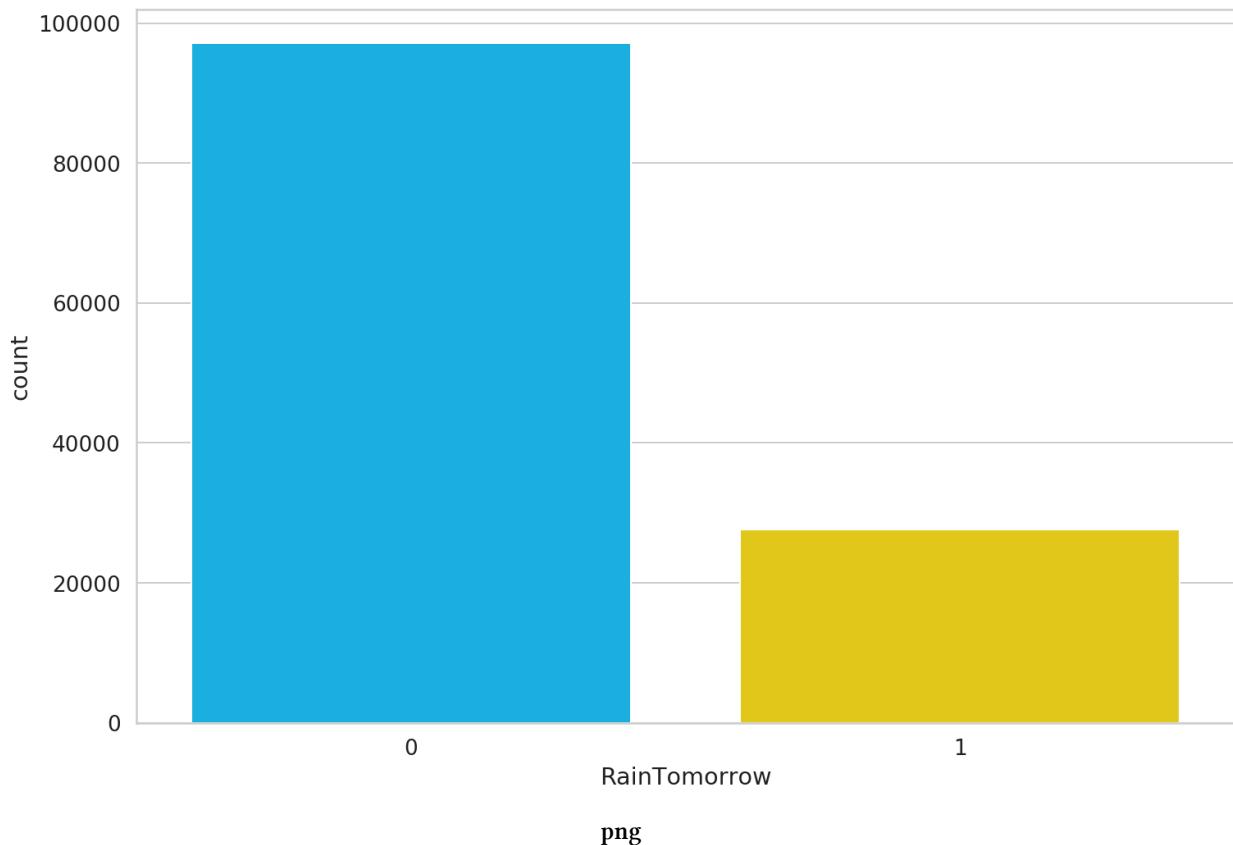
Let's drop the rows with missing values. There are better ways to do this, but we'll keep it simple:

```
1 df = df.dropna(how='any')
```

Finally, we have a dataset we can work with.

One important question we should answer is - *How balanced is our dataset?*. Or *How many times did it rain or not rain tomorrow?*:

```
1 sns.countplot(df.RainTomorrow);
```



```
1 df.RainTomorrow.value_counts() / df.shape[0]
```

```
1 0    0.778762
2 1    0.221238
3 Name: RainTomorrow, dtype: float64
```

Things are not looking good. About 78% of the data points have a non-rainy day for tomorrow. This means that a model that predicts there will be no rain tomorrow will be correct about 78% of the time.

You can read and apply the [Practical Guide to Handling Imbalanced Datasets<sup>5</sup>](#) if you want to mitigate this issue. Here, we'll just hope for the best.

The final step is to split the data into train and test sets:

---

<sup>5</sup><https://www.curiously.com/posts/practical-guide-to-handling-imbalanced-datasets/>

```

1 X = df[['Rainfall', 'Humidity3pm', 'RainToday', 'Pressure9am']]
2 y = df[['RainTomorrow']]
3
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state\
5 e=RANDOM_SEED)

```

And convert all of it to Tensors (so we can use it with PyTorch):

```

1 X_train = torch.from_numpy(X_train.to_numpy()).float()
2 y_train = torch.squeeze(torch.from_numpy(y_train.to_numpy()).float())
3
4 X_test = torch.from_numpy(X_test.to_numpy()).float()
5 y_test = torch.squeeze(torch.from_numpy(y_test.to_numpy()).float())
6
7 print(X_train.shape, y_train.shape)
8 print(X_test.shape, y_test.shape)

1 torch.Size([99751, 4]) torch.Size([99751])
2 torch.Size([24938, 4]) torch.Size([24938])

```

## Building a Neural Network

We'll build a simple Neural Network (NN) that tries to predict if it will rain tomorrow.

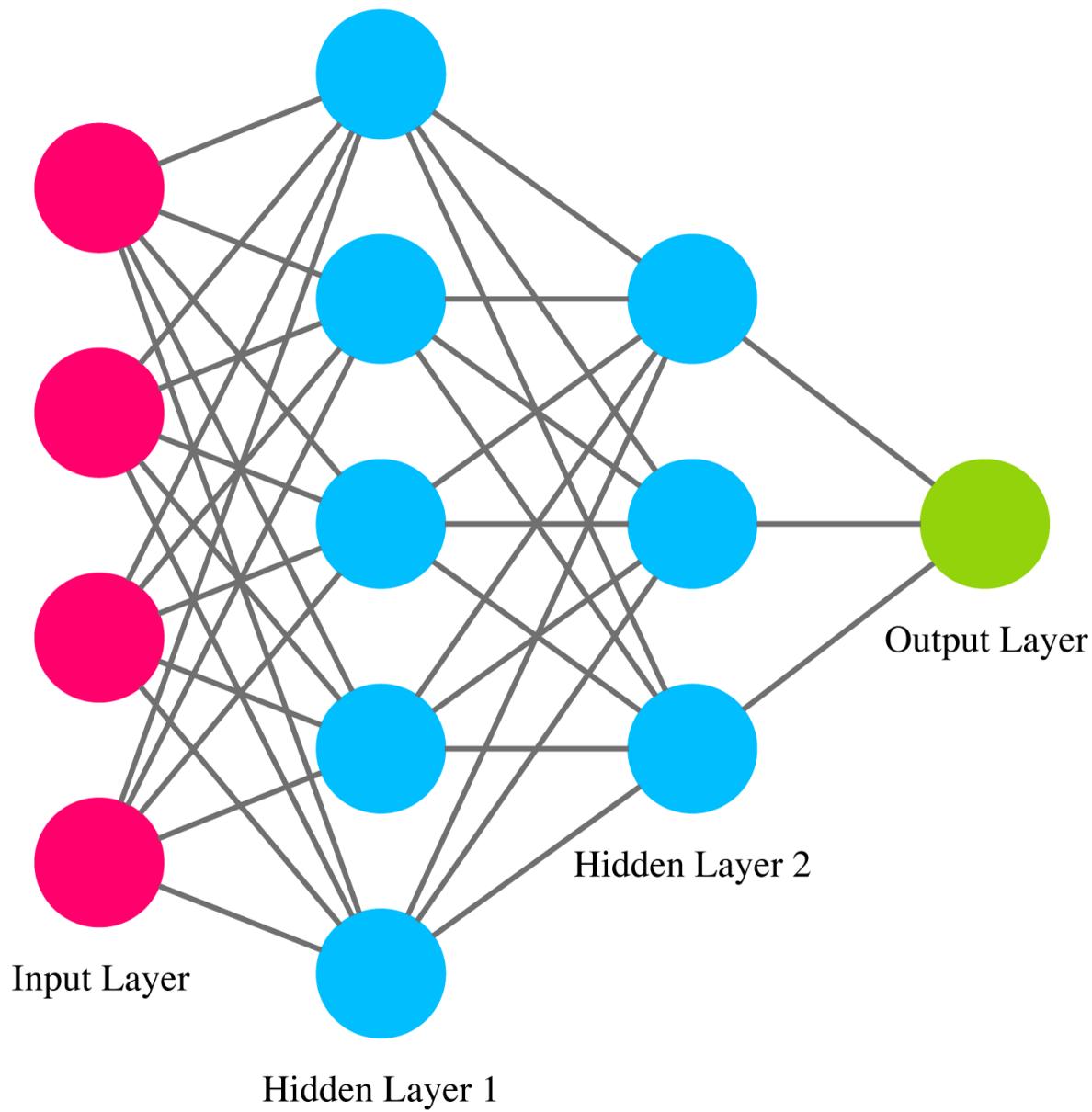
Our input contains data from the four columns: Rainfall, Humidity3pm, RainToday, Pressure9am. We'll create an appropriate input layer for that.

The output will be a number between 0 and 1, representing how likely (our model thinks) it is going to rain tomorrow. The prediction will be given to us by the final (output) layer of the network.

We'll have two (hidden) layers between the input and output layers. The parameters (neurons) of those layers will decide the final output. All layers will be fully-connected.

One easy way to build the NN with PyTorch is to create a class that inherits from `torch.nn.Module`:

```
1 class Net(nn.Module):  
2  
3     def __init__(self, n_features):  
4         super(Net, self).__init__()  
5         self.fc1 = nn.Linear(n_features, 5)  
6         self.fc2 = nn.Linear(5, 3)  
7         self.fc3 = nn.Linear(3, 1)  
8  
9     def forward(self, x):  
10        x = F.relu(self.fc1(x))  
11        x = F.relu(self.fc2(x))  
12        return torch.sigmoid(self.fc3(x))  
  
1 net = Net(X_train.shape[1])  
2  
3 ann_viz(net, view=False)
```



svg

We start by creating the layers of our model in the constructor. The `forward()` method is where the magic happens. It accepts the input `x` and allows it to flow through each layer.

There is a corresponding backward pass (defined for you by PyTorch) that allows the model to learn from the errors that it is currently making.

## Activation Functions

You might notice the calls to `F.relu` and `torch.sigmoid`. Why do we need those?

One of the cool *features* of Neural Networks is that they can approximate non-linear functions. In fact, [it is proven that they can approximate any function<sup>6</sup>](#).

Good luck approximating non-linear functions by stacking linear layers. Activation functions allow you to break from the linear world and learn (hopefully) more. You'll usually find them applied to an output of some layer.

Those functions must be hard to define, right?

### ReLU

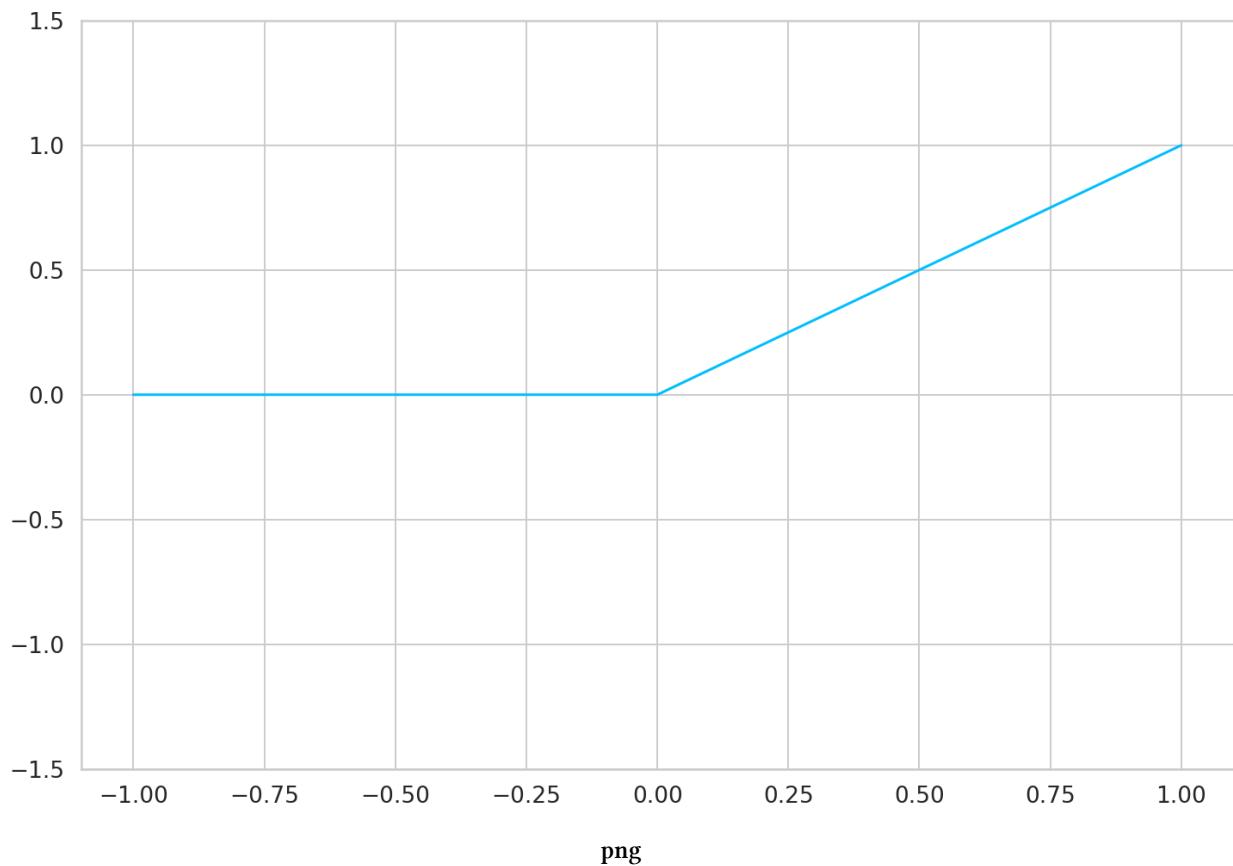
Not at all, let start with the ReLU definition (one of the most widely used activation function):

$$\text{ReLU}(x) = \max(0, x)$$

Easy peasy, the result is the maximum value of zero and the input.

```
1 ax = plt.gca()
2
3 plt.plot(
4     np.linspace(-1, 1, 5),
5     F.relu(torch.linspace(-1, 1, steps=5)).numpy()
6 )
7 ax.set_ylim([-1.5, 1.5]);
```

<sup>6</sup>[https://en.wikipedia.org/wiki/Universal\\_approximation\\_theorem](https://en.wikipedia.org/wiki/Universal_approximation_theorem)



## Sigmoid

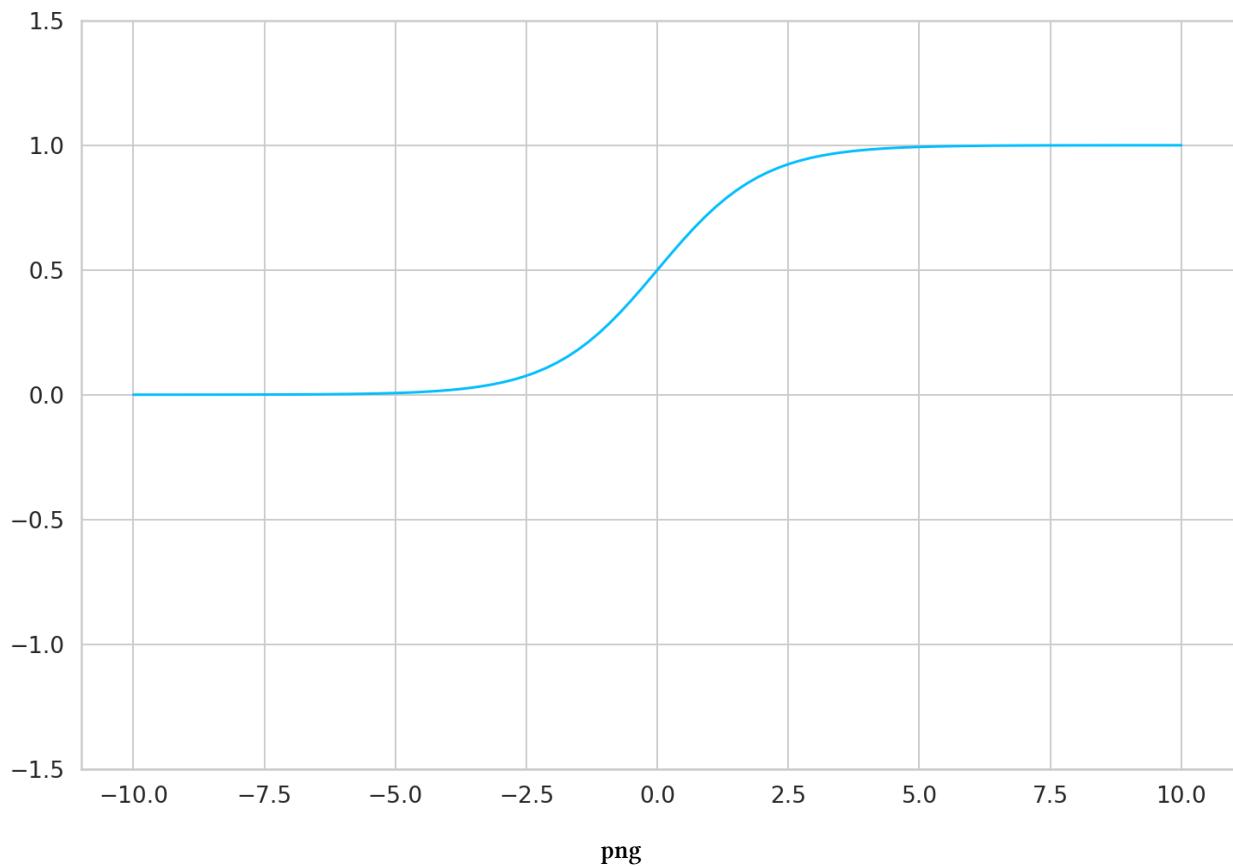
The sigmoid is useful when you need to make a binary decision/classification (answering with a *yes* or a *no*).

It is defined as:

$$\text{Sigmoid}(x) = \frac{1}{1 + e^{-x}}$$

The sigmoid squishes the input values between 0 and 1. But in a super kind of way:

```
1 ax = plt.gca()
2
3 plt.plot(
4     np.linspace(-10, 10, 100),
5     torch.sigmoid(torch.linspace(-10, 10, steps=100)).numpy()
6 )
7 ax.set_ylim([-1.5, 1.5]);
```



## Training

With the model in place, we need to find parameters that predict will it rain tomorrow. First, we need something to tell us how good we're currently doing:

```
1 criterion = nn.BCELoss()
```

The [BCELoss](#)<sup>7</sup> is a loss function that measures the difference between the two binary vectors. In our case, the predictions of our model and the real values. It expects the values to be outputted by the sigmoid function. The closer this value gets to 0, the better your model should be.

But how do we find parameters that minimize the loss function?

## Optimization

Imagine that each parameter of our NN is a knob. The optimizer's job is to find the perfect positions for each knob so that the loss gets close to 0.

<sup>7</sup><https://pytorch.org/docs/stable/nn.html#bceloss>

Real-world models can contain millions or even billions of parameters. With so many knobs to turn, it would be nice to have an efficient optimizer that quickly finds solutions.

Contrary to what you might believe, optimization in Deep Learning is just satisfying. In practice, you're content with good enough parameter values.

While there are tons of optimizers you can choose from, [Adam](#)<sup>8</sup> is a safe first choice. PyTorch has a well-debugged implementation you can use:

```
1 optimizer = optim.Adam(net.parameters(), lr=0.001)
```

Naturally, the optimizer requires the parameters. The second argument `lr` is *learning rate*. It is a tradeoff between how good parameters you're going to find and how fast you'll get there. Finding good values for this can be black magic and a lot of brute-force “experimentation”.

## Doing it on the GPU

Doing massively parallel computations on GPUs is one of the enablers for modern Deep Learning. You'll need nVIDIA GPU for that.

PyTorch makes it really easy to transfer all the computation to your GPU:

```
1 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")

1 X_train = X_train.to(device)
2 y_train = y_train.to(device)
3
4 X_test = X_test.to(device)
5 y_test = y_test.to(device)

1 net = net.to(device)
2
3 criterion = criterion.to(device)
```

We start by checking whether or not a CUDA device is available. Then, we transfer all training and test data to that device. Finally, we move our model and loss function.

## Finding Good Parameters

Having a loss function is great, but tracking the accuracy of our model is something easier to understand, for us mere mortals. Here's the definition for our accuracy:

---

<sup>8</sup><https://pytorch.org/docs/stable/optim.html#torch.optim.Adam>

```
1 def calculate_accuracy(y_true, y_pred):  
2     predicted = y_pred.ge(.5).view(-1)  
3     return (y_true == predicted).sum().float() / len(y_true)
```

We convert every value below 0.5 to 0. Otherwise, we set it to 1. Finally, we calculate the percentage of correct values.

With all the pieces of the puzzle in place, we can start training our model:

```
1 def round_tensor(t, decimal_places=3):  
2     return round(t.item(), decimal_places)  
3  
4 for epoch in range(1000):  
5  
6     y_pred = net(X_train)  
7  
8     y_pred = torch.squeeze(y_pred)  
9     train_loss = criterion(y_pred, y_train)  
10  
11    if epoch % 100 == 0:  
12        train_acc = calculate_accuracy(y_train, y_pred)  
13  
14        y_test_pred = net(X_test)  
15        y_test_pred = torch.squeeze(y_test_pred)  
16  
17        test_loss = criterion(y_test_pred, y_test)  
18  
19        test_acc = calculate_accuracy(y_test, y_test_pred)  
20        print(  
21            f'''epoch {epoch}  
22            Train set - loss: {round_tensor(train_loss)}, accuracy: {round_tensor(train_acc)}  
23            Test set - loss: {round_tensor(test_loss)}, accuracy: {round_tensor(test_acc)}  
24            '''  
25  
26        optimizer.zero_grad()  
27  
28        train_loss.backward()  
29  
30        optimizer.step()
```

```
1 epoch 0
2 Train set - loss: 2.513, accuracy: 0.779
3 Test set - loss: 2.517, accuracy: 0.778
4
5 epoch 100
6 Train set - loss: 0.457, accuracy: 0.792
7 Test set - loss: 0.458, accuracy: 0.793
8
9 epoch 200
10 Train set - loss: 0.435, accuracy: 0.801
11 Test set - loss: 0.436, accuracy: 0.8
12
13 epoch 300
14 Train set - loss: 0.421, accuracy: 0.814
15 Test set - loss: 0.421, accuracy: 0.815
16
17 epoch 400
18 Train set - loss: 0.412, accuracy: 0.826
19 Test set - loss: 0.413, accuracy: 0.827
20
21 epoch 500
22 Train set - loss: 0.408, accuracy: 0.831
23 Test set - loss: 0.408, accuracy: 0.832
24
25 epoch 600
26 Train set - loss: 0.406, accuracy: 0.833
27 Test set - loss: 0.406, accuracy: 0.835
28
29 epoch 700
30 Train set - loss: 0.405, accuracy: 0.834
31 Test set - loss: 0.405, accuracy: 0.835
32
33 epoch 800
34 Train set - loss: 0.404, accuracy: 0.834
35 Test set - loss: 0.404, accuracy: 0.835
36
37 epoch 900
38 Train set - loss: 0.404, accuracy: 0.834
39 Test set - loss: 0.404, accuracy: 0.836
```

During the training, we show our model the data for 10,000 times. Each time we measure the loss, propagate the errors through our model and asking the optimizer to find better parameters.

The `zero_grad()` method clears up the accumulated gradients, which the optimizer uses to find

better parameters.

What about that accuracy? 83.6% accuracy on the test set sounds reasonable, right? Well, I am about to disappoint you. But first, let's learn how to save and load our trained models.

## Saving the model

Training a good model can take a lot of time. And I mean weeks, months or even years. So, let's make sure that you know how you can save your precious work. Saving is easy:

```
1 MODEL_PATH = 'model.pth'  
2  
3 torch.save(net, MODEL_PATH)
```

Restoring your model is easy too:

```
1 net = torch.load(MODEL_PATH)
```

## Evaluation

Wouldn't it be perfect to know what kinds of errors your model makes? Of course, that's impossible. But you can get an estimate.

Using just accuracy wouldn't be a good way to do it. Recall that our data contains mostly no rain examples.

One way to delve a bit deeper into your model performance is to assess the precision and recall for each class. In our case, that will be *no rain* and *rain*:

```
1 classes = ['No rain', 'Raining']  
2  
3 y_pred = net(X_test)  
4  
5 y_pred = y_pred.ge(.5).view(-1).cpu()  
6 y_test = y_test.cpu()  
7  
8 print(classification_report(y_test, y_pred, target_names=classes))
```

		precision	recall	f1-score	support
1					
2					
3	No rain	0.85	0.96	0.90	19413
4	Raining	0.74	0.40	0.52	5525
5					
6	accuracy			0.84	24938
7	macro avg	0.80	0.68	0.71	24938
8	weighted avg	0.83	0.84	0.82	24938

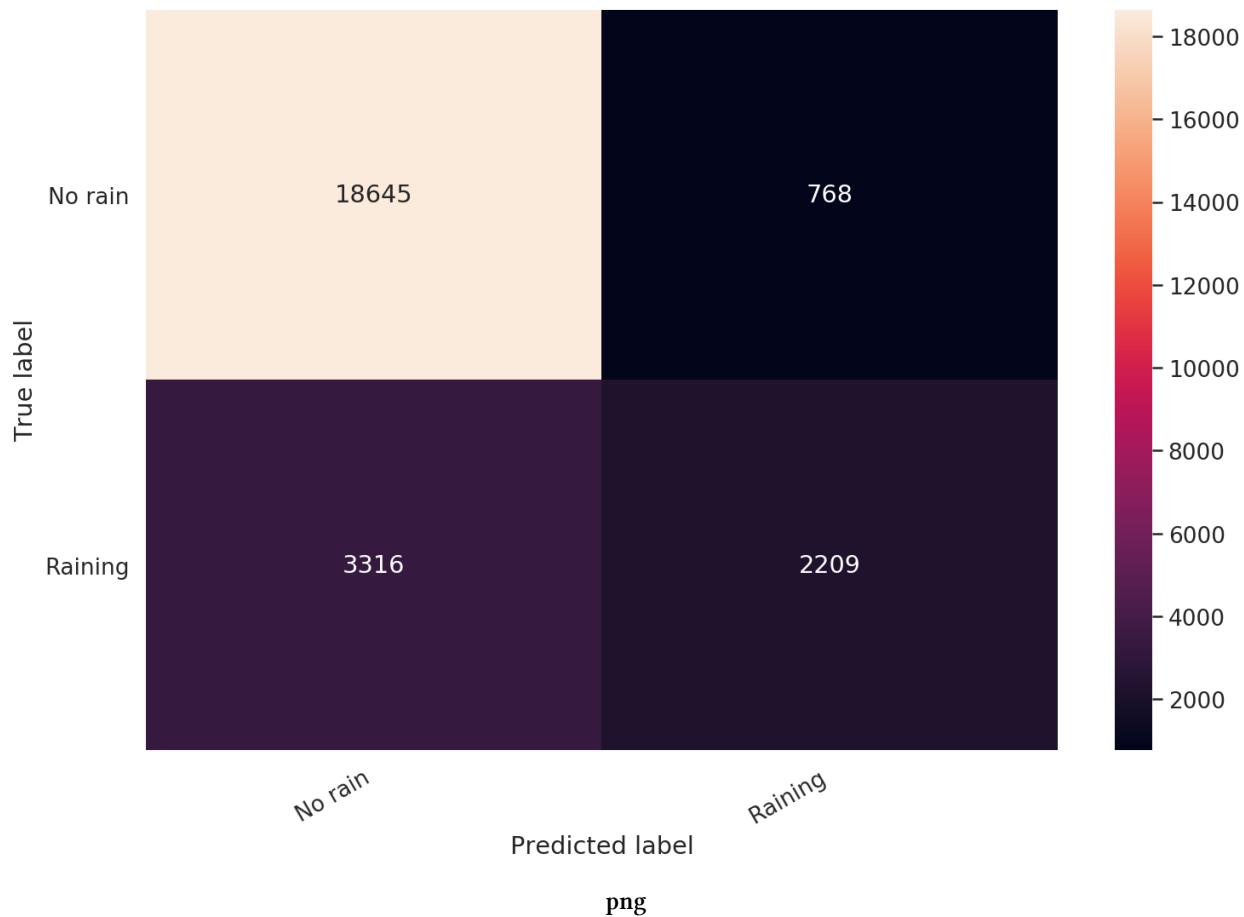
A maximum precision of 1 indicates that the model is perfect at identifying only relevant examples. A maximum recall of 1 indicates that our model can find all relevant examples in the dataset for this class.

You can see that our model is doing good when it comes to the *No rain* class. We have so many examples. Unfortunately, we can't really trust predictions of the *Raining* class.

One of the best things about binary classification is that you can have a good look at a simple confusion matrix:

```

1 cm = confusion_matrix(y_test, y_pred)
2 df_cm = pd.DataFrame(cm, index=classes, columns=classes)
3
4 hmap = sns.heatmap(df_cm, annot=True, fmt="d")
5 hmap.yaxis.set_ticklabels(hmap.yaxis.get_ticklabels(), rotation=0, ha='right')
6 hmap.xaxis.set_ticklabels(hmap.xaxis.get_ticklabels(), rotation=30, ha='right')
7 plt.ylabel('True label')
8 plt.xlabel('Predicted label');
```



You can clearly see that our model shouldn't be trusted when it says it's going to rain.

## Conclusion

Well done! You now have a Neural Network that can predict the weather. Well, sort of. Building well-performing models is hard, really hard. But there are tricks you'll pick up along the way and (hopefully) get better at your craft!

- Run the complete notebook in your browser (Google Colab)<sup>9</sup>
- Read the Getting Things Done with Pytorch book<sup>10</sup>

You learned how to:

- Preprocess CSV files and convert the data to Tensors

<sup>9</sup>[https://colab.research.google.com/drive/1lDXVkd7GC8jK\\_nGmOMKeDywXse-DY-u](https://colab.research.google.com/drive/1lDXVkd7GC8jK_nGmOMKeDywXse-DY-u)

<sup>10</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>

- Build your own Neural Network model with PyTorch
- Use a loss function and an optimizer to train your model
- Evaluate your model and learn about the perils of imbalanced classification

## References

- Precision and Recall<sup>11</sup>
- Beyond Accuracy: Precision and Recall<sup>12</sup>

<sup>11</sup>[https://en.wikipedia.org/wiki/Precision\\_and\\_recall](https://en.wikipedia.org/wiki/Precision_and_recall)

<sup>12</sup><https://towardsdatascience.com/beyond-accuracy-precision-and-recall-3da06bea9f6c>

# 3. Transfer Learning for Image Classification using Torchvision

TL;DR Learn how to use Transfer Learning to classify traffic sign images. You'll build a dataset of images in a format suitable for working with Torchvision. Get predictions on images from the wild (downloaded from the Internet).

In this tutorial, you'll learn how to fine-tune a pre-trained model for classifying raw pixels of traffic signs.

- Read the tutorial<sup>1</sup>
- Run the notebook in your browser (Google Colab)<sup>2</sup>
- Read the Getting Things Done with Pytorch book<sup>3</sup>

Here's what we'll go over:

- Overview of the traffic sign image dataset
- Build a dataset
- Use a pre-trained model from Torchvision
- Add a new *unknown* class and re-train the model

Will this model be ready for the real world?

```
1 import torch, torchvision  
2  
3 from pathlib import Path  
4 import numpy as np  
5 import cv2  
6 import pandas as pd  
7 from tqdm import tqdm  
8 import PIL.Image as Image  
9 import seaborn as sns  
10 from pylab import rcParams  
11 import matplotlib.pyplot as plt  
12 from matplotlib import rc
```

---

<sup>1</sup><https://www.curiously.com/posts/transfer-learning-for-image-classification-using-torchvision-pytorch-and-python/>

<sup>2</sup>[https://colab.research.google.com/drive/1Lk5R4pECDxDhd1uXcv26YRQ02fb\\_mrL9?usp=sharing](https://colab.research.google.com/drive/1Lk5R4pECDxDhd1uXcv26YRQ02fb_mrL9?usp=sharing)

<sup>3</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>

```

13 from matplotlib.ticker import MaxNLocator
14 from torch.optim import lr_scheduler
15 from sklearn.model_selection import train_test_split
16 from sklearn.metrics import confusion_matrix, classification_report
17 from glob import glob
18 import shutil
19 from collections import defaultdict
20
21 from torch import nn, optim
22
23 import torch.nn.functional as F
24 import torchvision.transforms as T
25 from torchvision.datasets import ImageFolder
26 from torch.utils.data import DataLoader
27 from torchvision import models
28
29 %matplotlib inline
30 %config InlineBackend.figure_format='retina'
31
32 sns.set(style='whitegrid', palette='muted', font_scale=1.2)
33
34 HAPPY_COLORS_PALETTE = ["#01BEFE", "#FFDD00", "#FF7D00", "#FF006D", "#ADFF02", "#8F0\OFF"]
35
36 sns.set_palette(sns.color_palette(HAPPY_COLORS_PALETTE))
37
38 rcParams['figure.figsize'] = 12, 8
39
40 RANDOM_SEED = 42
41 np.random.seed(RANDOM_SEED)
42 torch.manual_seed(RANDOM_SEED)
43
44 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

```

## Recognizing traffic signs

German Traffic Sign Recognition Benchmark (GTSRB)<sup>4</sup> contains more than 50,000 annotated images of 40+ traffic signs. Given an image, you'll have to recognize the traffic sign on it.

---

<sup>4</sup><http://benchmark.ini.rub.de/?section=gtsrb&subsection=dataset>

```

1 !wget https://sid.elda.dk/public/archives/daaeac0d7ce1152aea9b61d9f1e19370/GTSRB_Final\\
2 a1_Training_Images.zip
3 !unzip -qq GTSRB_Final_Training_Images.zip

```

## Exploration

Let's start by getting a feel of the data. The images for each traffic sign are stored in a separate directory. How many do we have?

```

1 train_folders = sorted(glob('GTSRB/Final_Training/Images/*'))
2 len(train_folders)

```

```
1      43
```

We'll create 3 helper functions that use OpenCV and Torchvision to load and show images:

```

1 def load_image(img_path, resize=True):
2     img = cv2.cvtColor(cv2.imread(img_path), cv2.COLOR_BGR2RGB)
3
4     if resize:
5         img = cv2.resize(img, (64, 64), interpolation = cv2.INTER_AREA)
6
7     return img
8
9 def show_image(img_path):
10    img = load_image(img_path)
11    plt.imshow(img)
12    plt.axis('off')
13
14 def show_sign_grid(image_paths):
15    images = [load_image(img) for img in image_paths]
16    images = torch.as_tensor(images)
17    images = images.permute(0, 3, 1, 2)
18    grid_img = torchvision.utils.make_grid(images, nrow=11)
19    plt.figure(figsize=(24, 12))
20    plt.imshow(grid_img.permute(1, 2, 0))
21    plt.axis('off');

```

Let's have a look at some examples for each traffic sign:

```
1 sample_images = [np.random.choice(glob(f'{tf}/*ppm')) for tf in train_folders]
2 show_sign_grid(sample_images)
```



png

And here is a single sign:

```
1 img_path = glob(f'{train_folders[16]}/*ppm')[1]
2
3 show_image(img_path)
```



png

## Building a dataset

To keep things simple, we'll focus on classifying some of the most used traffic signs:

```

1 class_names = ['priority_road', 'give_way', 'stop', 'no_entry']
2
3 class_indices = [12, 13, 14, 17]

```

We'll copy the images files to a new directory, so it's easier to use the Torchvision's dataset helpers. Let's start with the directories for each class:

```

1 !rm -rf data
2
3 DATA_DIR = Path('data')
4
5 DATASETS = ['train', 'val', 'test']
6
7 for ds in DATASETS:
8     for cls in class_names:
9         (DATA_DIR / ds / cls).mkdir(parents=True, exist_ok=True)

```

We'll reserve 80% of the images for training, 10% for validation, and 10% test for each class. We'll copy each image to the correct dataset directory:

```

1 for i, cls_index in enumerate(class_indices):
2     image_paths = np.array(glob(f'{train_folders[cls_index]}/*.{ppm}'))
3     class_name = class_names[i]
4     print(f'{class_name}: {len(image_paths)}')
5     np.random.shuffle(image_paths)
6
7     ds_split = np.split(
8         image_paths,
9         indices_or_sections=[int(.8*len(image_paths)), int(.9*len(image_paths))])
10    )
11
12    dataset_data = zip(DATASETS, ds_split)
13
14    for ds, images in dataset_data:
15        for img_path in images:
16            shutil.copy(img_path, f'{DATA_DIR}/{ds}/{class_name}/')

```

```

1 priority_road: 2100
2 give_way: 2160
3 stop: 780
4 no_entry: 1110

```

We have some class imbalance, but it is not that bad. We'll ignore it.

We'll apply some image augmentation techniques to artificially increase the size of our training dataset:

```

1 mean_nums = [0.485, 0.456, 0.406]
2 std_nums = [0.229, 0.224, 0.225]
3
4 transforms = { 'train': T.Compose([
5     T.RandomResizedCrop(size=256),
6     T.RandomRotation(degrees=15),
7     T.RandomHorizontalFlip(),
8     T.ToTensor(),
9     T.Normalize(mean_nums, std_nums)
10]), 'val': T.Compose([
11     T.Resize(size=256),
12     T.CenterCrop(size=224),
13     T.ToTensor(),
14     T.Normalize(mean_nums, std_nums)
15]), 'test': T.Compose([
16     T.Resize(size=256),
17     T.CenterCrop(size=224),
18     T.ToTensor(),
19     T.Normalize(mean_nums, std_nums)
20]), }
21

```

We apply some random resizing, rotation, and horizontal flips. Finally, we normalize the tensors using preset values for each channel. This is a [requirement of the pre-trained models<sup>5</sup>](#) in Torchvision.

We'll create a PyTorch dataset for each image dataset folder and data loaders for easier training:

---

<sup>5</sup><https://pytorch.org/docs/stable/torchvision/models.html>

```

1 image_datasets = {
2     d: ImageFolder(f'{DATA_DIR}/{d}', transforms[d]) for d in DATASETS
3 }
4
5 data_loaders = {
6     d: DataLoader(image_datasets[d], batch_size=4, shuffle=True, num_workers=4)
7     for d in DATASETS
8 }
```

We'll also store the number of examples in each dataset and class names for later:

```

1 dataset_sizes = {d: len(image_datasets[d]) for d in DATASETS}
2 class_names = image_datasets['train'].classes
3
4 dataset_sizes
5
6 {'test': 615, 'train': 4920, 'val': 615}
```

Let's have a look at some example images with applied transformations. We also need to reverse the normalization and reorder the color channels to get correct image data:

```

1 def imshow(inp, title=None):
2     inp = inp.numpy().transpose((1, 2, 0))
3     mean = np.array([mean_nums])
4     std = np.array([std_nums])
5     inp = std * inp + mean
6     inp = np.clip(inp, 0, 1)
7     plt.imshow(inp)
8     if title is not None:
9         plt.title(title)
10    plt.axis('off')
11
12 inputs, classes = next(iter(data_loaders['train']))
13 out = torchvision.utils.make_grid(inputs)
14
15 imshow(out, title=[class_names[x] for x in classes])
```



## Using a pre-trained model:

Our model will receive raw image pixels and try to classify them into one of four traffic signs. How hard can it be? Try to build a model from scratch.

Here, we'll use [Transfer Learning](#)<sup>6</sup> to copy the architecture of the very popular [ResNet](#)<sup>7</sup> model. On top of that, we'll use the learned weights of the model from training on the [ImageNet dataset](#)<sup>8</sup>. All of this is made easy to use by Torchvision:

```

1 def create_model(n_classes):
2     model = models.resnet34(pretrained=True)
3
4     n_features = model.fc.in_features
5     model.fc = nn.Linear(n_features, n_classes)
6
7     return model.to(device)

```

We reuse almost everything except the change of the output layer. This is needed because the number of classes in our dataset is different than ImageNet.

Let's create an instance of our model:

```
1 base_model = create_model(len(class_names))
```

## Training

We'll write 3 helper functions to encapsulate the training and evaluation logic. Let's start with `train_epoch`:

<sup>6</sup>[https://en.wikipedia.org/wiki/Transfer\\_learning](https://en.wikipedia.org/wiki/Transfer_learning)

<sup>7</sup><https://arxiv.org/abs/1512.03385>

<sup>8</sup><http://www.image-net.org/>

```
1 def train_epoch(
2     model,
3     data_loader,
4     loss_fn,
5     optimizer,
6     device,
7     scheduler,
8     n_examples
9 ):
10    model = model.train()
11
12    losses = []
13    correct_predictions = 0
14
15    for inputs, labels in data_loader:
16        inputs = inputs.to(device)
17        labels = labels.to(device)
18
19        outputs = model(inputs)
20
21        _, preds = torch.max(outputs, dim=1)
22        loss = loss_fn(outputs, labels)
23
24        correct_predictions += torch.sum(preds == labels)
25        losses.append(loss.item())
26
27        loss.backward()
28        optimizer.step()
29        optimizer.zero_grad()
30
31    scheduler.step()
32
33    return correct_predictions.double() / n_examples, np.mean(losses)
```

We start by turning our model into train mode and go over the data. After getting the predictions, we get the class with maximum probability along with the loss, so we can calculate the epoch loss and accuracy.

Note that we're also using a learning rate scheduler (more on that later).

```

1 def eval_model(model, data_loader, loss_fn, device, n_examples):
2     model = model.eval()
3
4     losses = []
5     correct_predictions = 0
6
7     with torch.no_grad():
8         for inputs, labels in data_loader:
9             inputs = inputs.to(device)
10            labels = labels.to(device)
11
12            outputs = model(inputs)
13
14            _, preds = torch.max(outputs, dim=1)
15
16            loss = loss_fn(outputs, labels)
17
18            correct_predictions += torch.sum(preds == labels)
19            losses.append(loss.item())
20
21    return correct_predictions.double() / n_examples, np.mean(losses)

```

The evaluation of the model is pretty similar, except that we don't do any gradient calculations.

Let's put everything together:

```

1 def train_model(model, data_loaders, dataset_sizes, device, n_epochs=3):
2     optimizer = optim.SGD(model.parameters(), lr=0.001, momentum=0.9)
3     scheduler = lr_scheduler.StepLR(optimizer, step_size=7, gamma=0.1)
4     loss_fn = nn.CrossEntropyLoss().to(device)
5
6     history = defaultdict(list)
7     best_accuracy = 0
8
9     for epoch in range(n_epochs):
10
11         print(f'Epoch {epoch + 1}/{n_epochs}')
12         print('-' * 10)
13
14         train_acc, train_loss = train_epoch(
15             model,
16             data_loaders['train'],
17             loss_fn,

```

```
18     optimizer,
19     device,
20     scheduler,
21     dataset_sizes['train']
22 )
23
24     print(f'Train loss {train_loss} accuracy {train_acc}')
25
26     val_acc, val_loss = eval_model(
27         model,
28         data_loaders['val'],
29         loss_fn,
30         device,
31         dataset_sizes['val']
32     )
33
34     print(f'Val    loss {val_loss} accuracy {val_acc}')
35     print()
36
37     history['train_acc'].append(train_acc)
38     history['train_loss'].append(train_loss)
39     history['val_acc'].append(val_acc)
40     history['val_loss'].append(val_loss)
41
42     if val_acc > best_accuracy:
43         torch.save(model.state_dict(), 'best_model_state.bin')
44         best_accuracy = val_acc
45
46     print(f'Best val accuracy: {best_accuracy}')
47
48     model.load_state_dict(torch.load('best_model_state.bin'))
49
50     return model, history
```

We do a lot of string formatting and recording of the training history. The hard stuff gets delegated to the previous helper functions. We also want the best model, so the weights of the most accurate model(s) get stored during the training.

Let's train our first model:

```
1 %%time
2
3 base_model, history = train_model(base_model, data_loaders, dataset_sizes, device)

1 Epoch 1/3
2 -----
3 Train loss 0.31827690804876935 accuracy 0.8859756097560976
4 Val    loss 0.0012465072916699694 accuracy 1.0
5
6 Epoch 2/3
7 -----
8 Train loss 0.12230596961529275 accuracy 0.9615853658536585
9 Val    loss 0.0007955377752130681 accuracy 1.0
10
11 Epoch 3/3
12 -----
13 Train loss 0.07771141678094864 accuracy 0.9745934959349594
14 Val    loss 0.0025791768387877366 accuracy 0.9983739837398374
15
16 Best val accuracy: 1.0
17 CPU times: user 2min 24s, sys: 48.2 s, total: 3min 12s
18 Wall time: 3min 21s
```

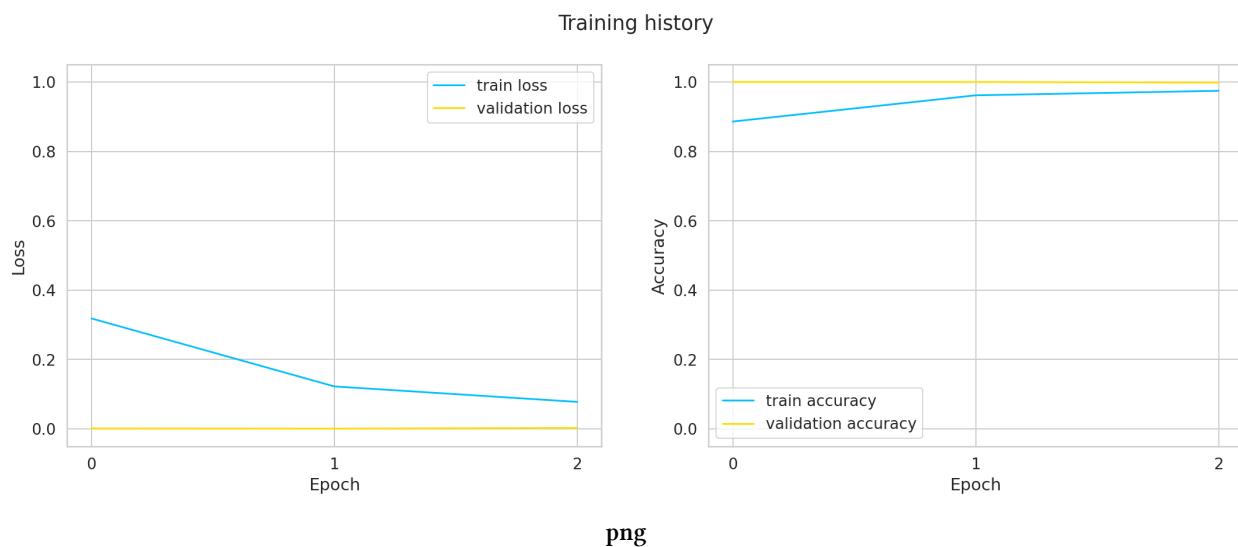
Here's a little helper function that visualizes the training history for us:

```
1 def plot_training_history(history):
2     fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(18, 6))
3
4     ax1.plot(history['train_loss'], label='train loss')
5     ax1.plot(history['val_loss'], label='validation loss')
6
7     ax1.xaxis.set_major_locator(MaxNLocator(integer=True))
8     ax1.set_ylim([-0.05, 1.05])
9     ax1.legend()
10    ax1.set_ylabel('Loss')
11    ax1.set_xlabel('Epoch')
12
13    ax2.plot(history['train_acc'], label='train accuracy')
14    ax2.plot(history['val_acc'], label='validation accuracy')
15
16    ax2.xaxis.set_major_locator(MaxNLocator(integer=True))
17    ax2.set_ylim([-0.05, 1.05])
```

```

18     ax2.legend()
19
20     ax2.set_ylabel('Accuracy')
21     ax2.set_xlabel('Epoch')
22
23     fig.suptitle('Training history')
24
25 plot_training_history(history)

```



The pre-trained model is so good that we get very high accuracy and low loss after 3 epochs. Unfortunately, our validation set is too small to get some meaningful metrics from it.

## Evaluation

Let's see some predictions on traffic signs from the test set:

```

1 def show_predictions(model, class_names, n_images=6):
2     model = model.eval()
3     images_handled = 0
4     plt.figure()
5
6     with torch.no_grad():
7         for i, (inputs, labels) in enumerate(data_loaders['test']):
8             inputs = inputs.to(device)
9             labels = labels.to(device)
10
11             outputs = model(inputs)

```

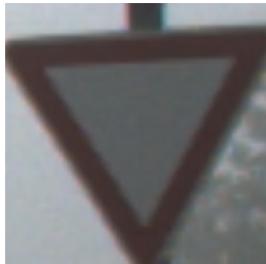
```
12     _, preds = torch.max(outputs, 1)
13
14     for j in range(inputs.shape[0]):
15         images_handeled += 1
16         ax = plt.subplot(2, n_images//2, images_handeled)
17         ax.set_title(f'predicted: {class_names[preds[j]]}')
18         imshow(inputs.cpu().data[j])
19         ax.axis('off')
20
21     if images_handeled == n_images:
22         return
```

1 show\_predictions(base\_model, class\_names, n\_images=8)

predicted: priority\_road



predicted: give\_way



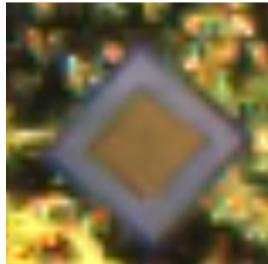
predicted: give\_way



predicted: give\_way



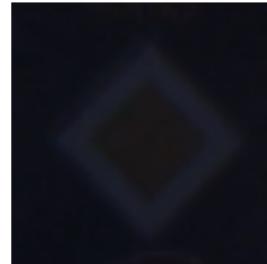
predicted: priority\_road



predicted: stop



predicted: priority\_road



predicted: stop



png

Very good! Even the almost not visible *priority road* sign is classified correctly. Let's dive a bit deeper. We'll start by getting the predictions from our model:

```

1 def get_predictions(model, data_loader):
2     model = model.eval()
3     predictions = []
4     real_values = []
5     with torch.no_grad():
6         for inputs, labels in data_loader:
7             inputs = inputs.to(device)
8             labels = labels.to(device)
9
10        outputs = model(inputs)
11        _, preds = torch.max(outputs, 1)
12        predictions.extend(preds)
13        real_values.extend(labels)
14    predictions = torch.as_tensor(predictions).cpu()
15    real_values = torch.as_tensor(real_values).cpu()
16    return predictions, real_values

1 y_pred, y_test = get_predictions(base_model, data_loaders['test'])

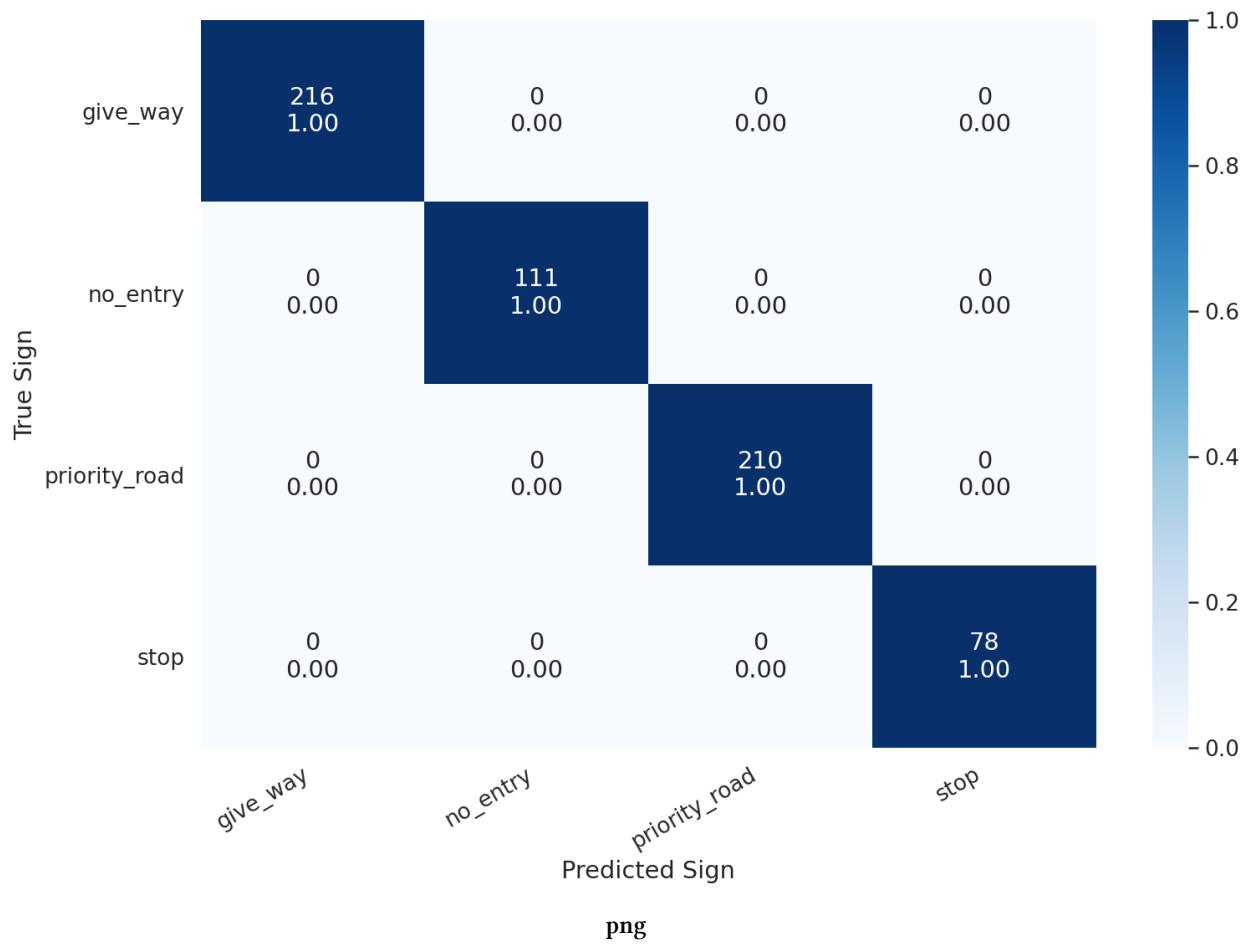
1 print(classification_report(y_test, y_pred, target_names=class_names))

```

	precision	recall	f1-score	support
give_way	1.00	1.00	1.00	216
no_entry	1.00	1.00	1.00	111
priority_road	1.00	1.00	1.00	210
stop	1.00	1.00	1.00	78
accuracy			1.00	615
macro avg	1.00	1.00	1.00	615
weighted avg	1.00	1.00	1.00	615

The classification report shows us that our model is perfect, not something you see every day! Does this thing make any mistakes?

```
1 def show_confusion_matrix(confusion_matrix, class_names):
2
3     cm = confusion_matrix.copy()
4
5     cell_counts = cm.flatten()
6
7     cm_row_norm = cm / cm.sum(axis=1)[:, np.newaxis]
8
9     row_percentages = ["{0:.2f}".format(value) for value in cm_row_norm.flatten()]
10
11    cell_labels = [f"{cnt}\n{per}" for cnt, per in zip(cell_counts, row_percentages)]
12    cell_labels = np.asarray(cell_labels).reshape(cm.shape[0], cm.shape[1])
13
14    df_cm = pd.DataFrame(cm_row_norm, index=class_names, columns=class_names)
15
16    hmap = sns.heatmap(df_cm, annot=cell_labels, fmt="", cmap="Blues")
17    hmap.yaxis.set_ticklabels(hmap.yaxis.get_ticklabels(), rotation=0, ha='right')
18    hmap.xaxis.set_ticklabels(hmap.xaxis.get_ticklabels(), rotation=30, ha='right')
19    plt.ylabel('True Sign')
20    plt.xlabel('Predicted Sign');
21
22
23 cm = confusion_matrix(y_test, y_pred)
24 show_confusion_matrix(cm, class_names)
```



No, no mistakes here!

## Classifying unseen images

Ok, but how good our model will be when confronted with a real-world image? Let's check it out:

```
1 !gdown --id 19Qz3a610u_QSHsLeTznx8LtDBu4tbqHr
```

```
1 show_image('stop-sign.jpg')
```



png

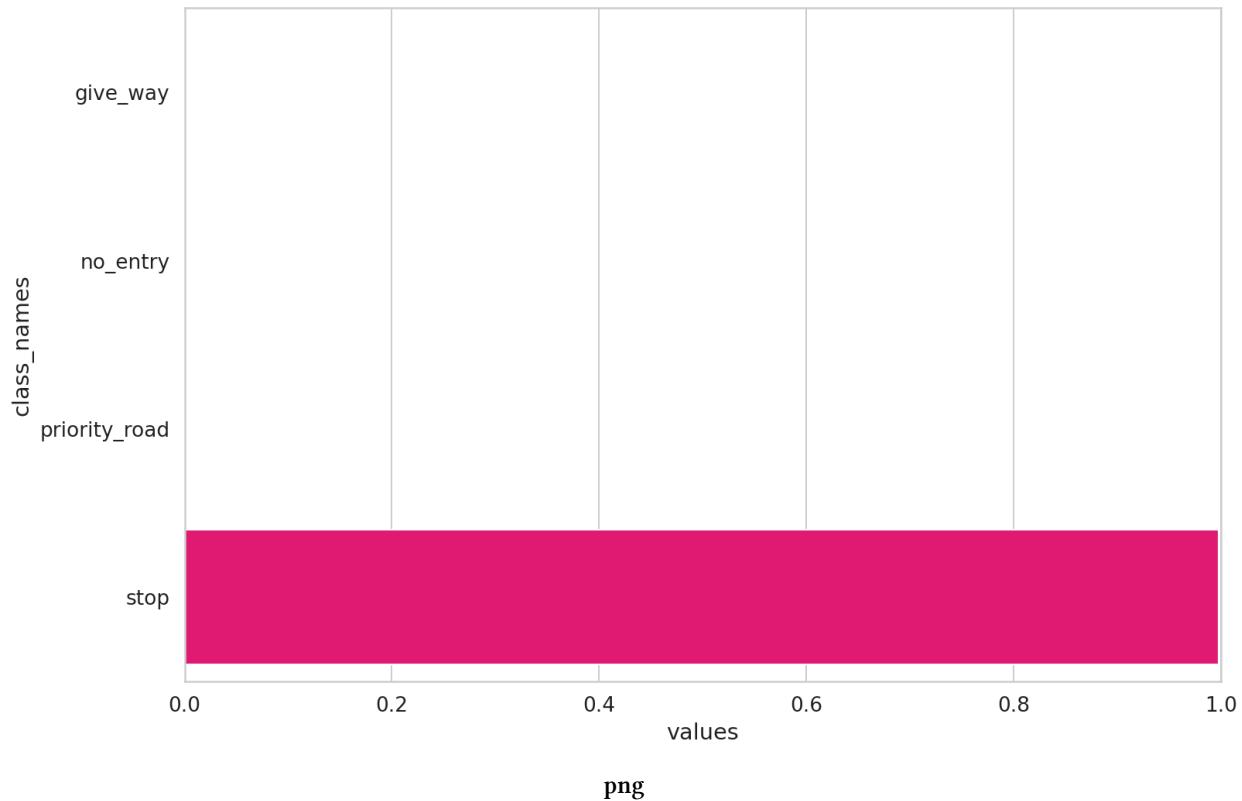
For this, we'll have a look at the confidence for each class. Let's get this from our model:

```
1 def predict_proba(model, image_path):  
2     img = Image.open(image_path)  
3     img = img.convert('RGB')  
4     img = transforms['test'](img).unsqueeze(0)  
5  
6     pred = model(img.to(device))  
7     pred = F.softmax(pred, dim=1)  
8     return pred.detach().cpu().numpy().flatten()  
  
1 pred = predict_proba(base_model, 'stop-sign.jpg')  
2 pred
```

```
1 array([1.1296713e-03, 1.9811286e-04, 3.4486805e-04, 9.9832731e-01],
2      dtype=float32)
```

This is a bit hard to understand. Let's plot it:

```
1 def show_prediction_confidence(prediction, class_names):
2     pred_df = pd.DataFrame({
3         'class_names': class_names,
4         'values': prediction
5     })
6     sns.barplot(x='values', y='class_names', data=pred_df, orient='h')
7     plt.xlim([0, 1]);
8
9 show_prediction_confidence(pred, class_names)
```



Again, our model is performing very well! Really confident in the correct traffic sign!

## Classifying unknown traffic sign

The last challenge for our model is a traffic sign that it hasn't seen before:

```
1 !gdown --id 1F61-iNh1Jk-yKZRGcu6S9P29HxDfxF0u
```

```
1 show_image('unknown-sign.jpg')
```



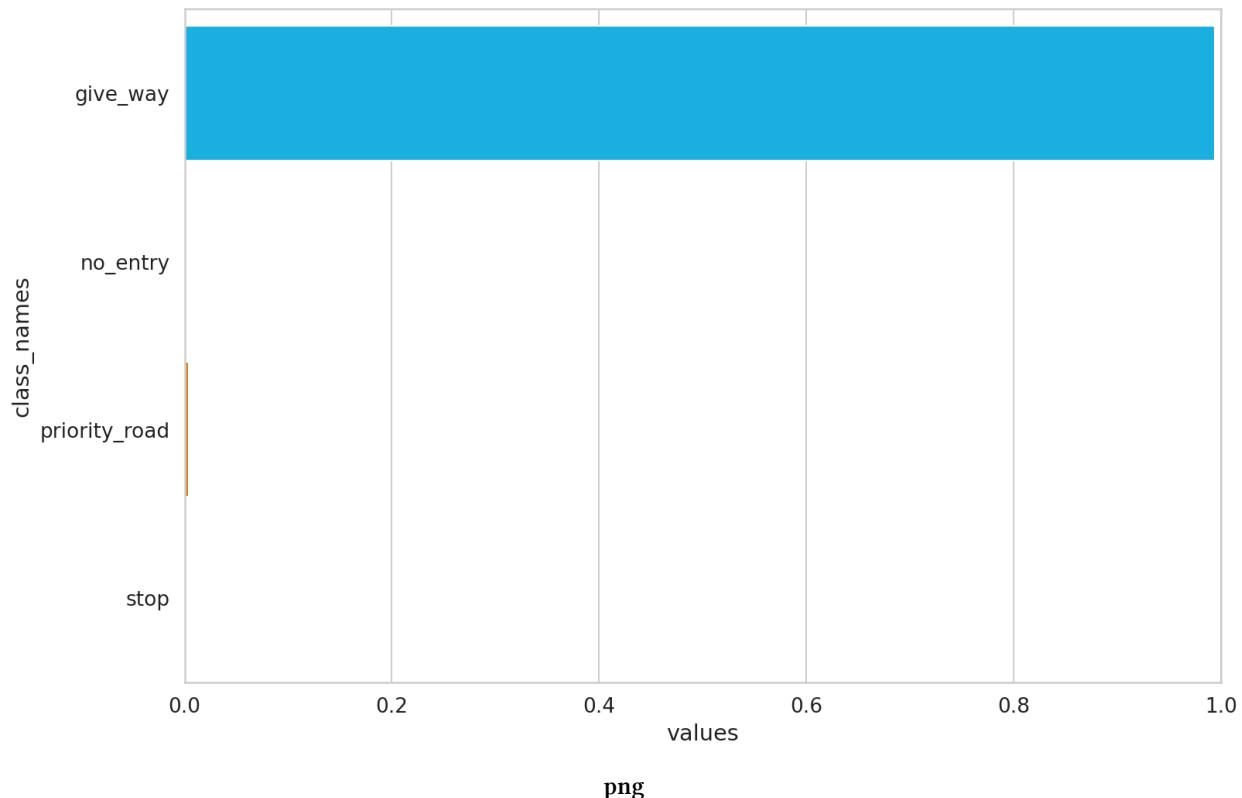
png

Let's get the predictions:

```
1 pred = predict_proba(base_model, 'unknown-sign.jpg')
2 pred
```

```
1 array([9.9413127e-01, 1.1861280e-06, 3.9936006e-03, 1.8739274e-03],
2      dtype=float32)
```

```
1 show_prediction_confidence(pred, class_names)
```



Our model is very certain (more than 95% confidence) that this is a *give way* sign. This is obviously wrong. How can you make your model see this?

## Adding class “unknown”

While there are a variety of ways to handle this situation (one described in this paper: [A Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks](#)<sup>9</sup>), we’ll do something simpler.

We’ll get the indices of all traffic signs that weren’t included in our original dataset:

---

<sup>9</sup><https://arxiv.org/pdf/1610.02136.pdf>

```

1 unknown_indices = [
2     i for i, f in enumerate(train_folders) \
3     if i not in class_indices
4 ]
5
6 len(unknown_indices)

```

1 39

We'll create a new folder for the unknown class and copy some of the images there:

```

1 for ds in DATASETS:
2     (DATA_DIR / ds / 'unknown').mkdir(parents=True, exist_ok=True)
3
4 for ui in unknown_indices:
5     image_paths = np.array(glob(f'{train_folders[ui]}/*.ppm'))
6     image_paths = np.random.choice(image_paths, 50)
7
8     ds_split = np.split(
9         image_paths,
10        indices_or_sections=[int(.8*len(image_paths)), int(.9*len(image_paths))])
11    )
12
13 dataset_data = zip(DATASETS, ds_split)
14
15 for ds, images in dataset_data:
16     for img_path in images:
17         shutil.copy(img_path, f'{DATA_DIR}/{ds}/unknown/')

```

The next steps are identical to what we've already done:

```

1 image_datasets = {
2     d: ImageFolder(f'{DATA_DIR}/{d}', transforms[d]) for d in DATASETS
3 }
4
5 data_loaders = {
6     d: DataLoader(image_datasets[d], batch_size=4, shuffle=True, num_workers=4)
7     for d in DATASETS
8 }
9
10 dataset_sizes = {d: len(image_datasets[d]) for d in DATASETS}

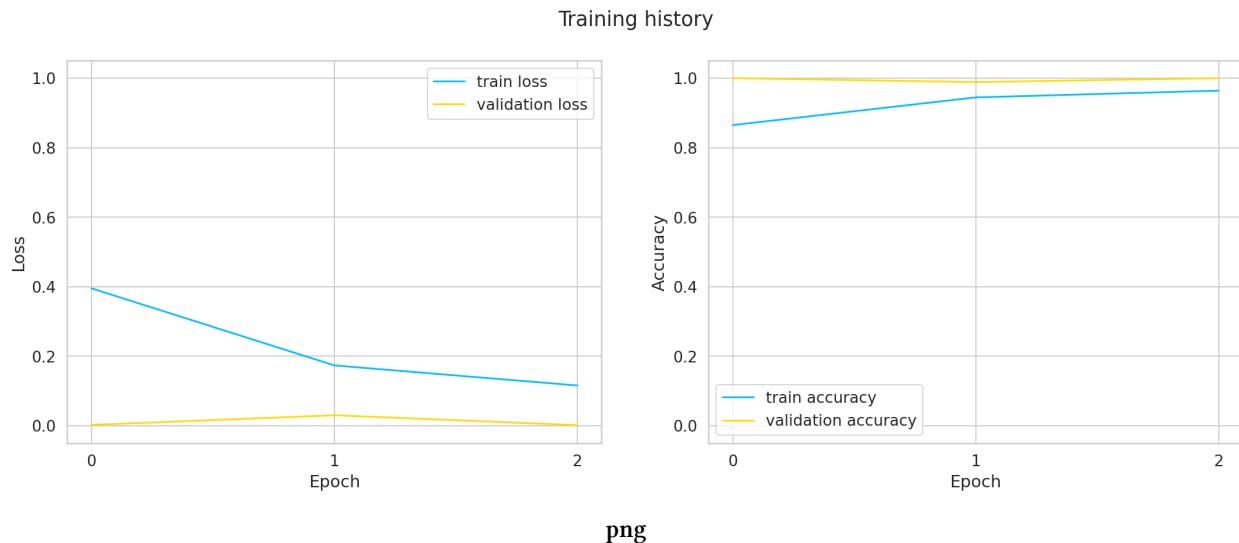
```

```
11 class_names = image_datasets['train'].classes
12
13 dataset_sizes
1
1     {'test': 784, 'train': 5704, 'val': 794}

1 %%time
2
3 enhanced_model = create_model(len(class_names))
4 enhanced_model, history = train_model(enhanced_model, data_loaders, dataset_sizes,\device)
5

1 Epoch 1/3
2 -----
3 Train loss 0.39523224640235327 accuracy 0.8650070126227208
4 Val    loss 0.002290595416447625 accuracy 1.0
5
6 Epoch 2/3
7 -----
8 Train loss 0.173455789528505 accuracy 0.9446002805049089
9 Val    loss 0.030148923471944415 accuracy 0.9886649874055415
10
11 Epoch 3/3
12 -----
13 Train loss 0.11575758963990512 accuracy 0.9640603085553997
14 Val    loss 0.0014996432778823317 accuracy 1.0
15
16 Best val accuracy: 1.0
17 CPU times: user 2min 47s, sys: 56.2 s, total: 3min 44s
18 Wall time: 3min 53s

1 plot_training_history(history)
```



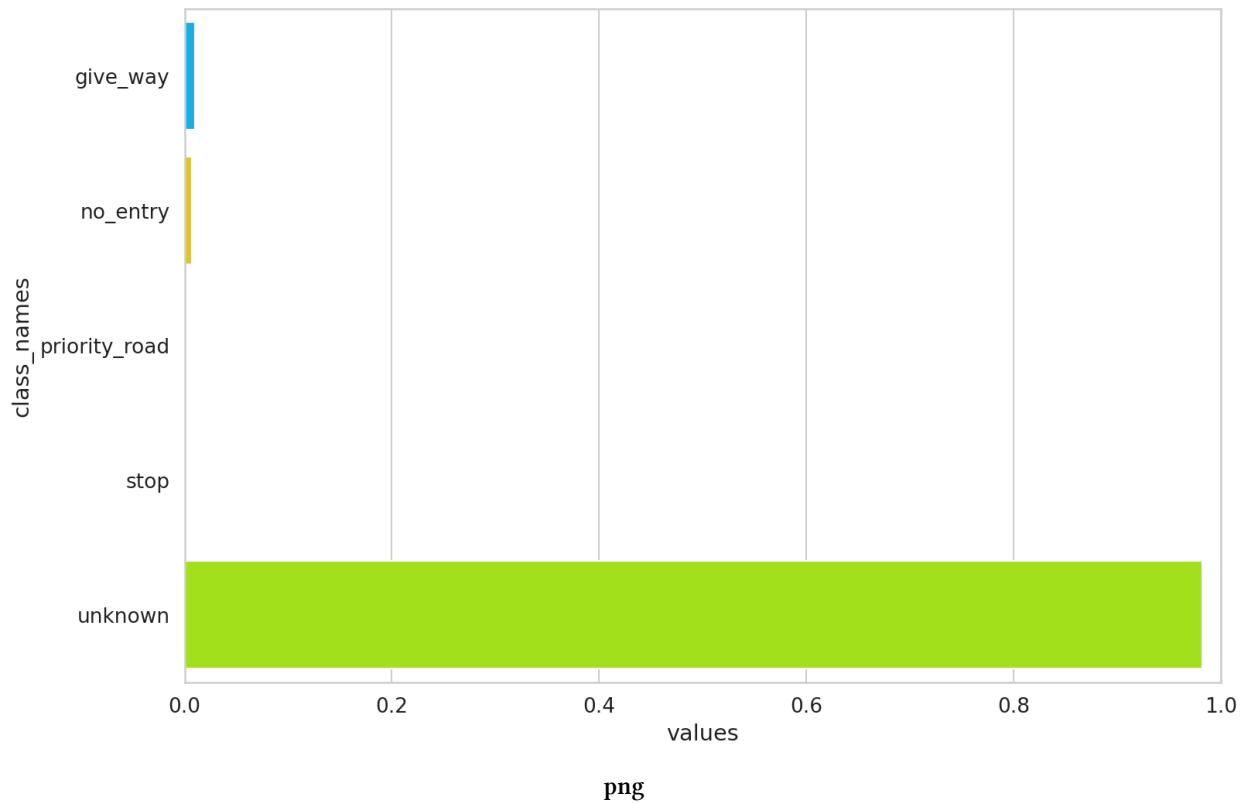
Again, our model is learning very quickly. Let's have a look at the sample image again:

```
1 show_image('unknown-sign.jpg')
```



png

```
1 pred = predict_proba(enhanced_model, 'unknown-sign.jpg')
2 show_prediction_confidence(pred, class_names)
```



Great, the model doesn't give much weight to any of the known classes. It doesn't magically know that this is a two-way sign, but recognizes it as unknown.

Let's have a look at some examples of our new dataset:

```
1 show_predictions(enhanced_model, class_names, n_images=8)
```



Let's get an overview of the new model's performance:

```

1 y_pred, y_test = get_predictions(enhanced_model, data_loaders['test'])

1 print(classification_report(y_test, y_pred, target_names=class_names))

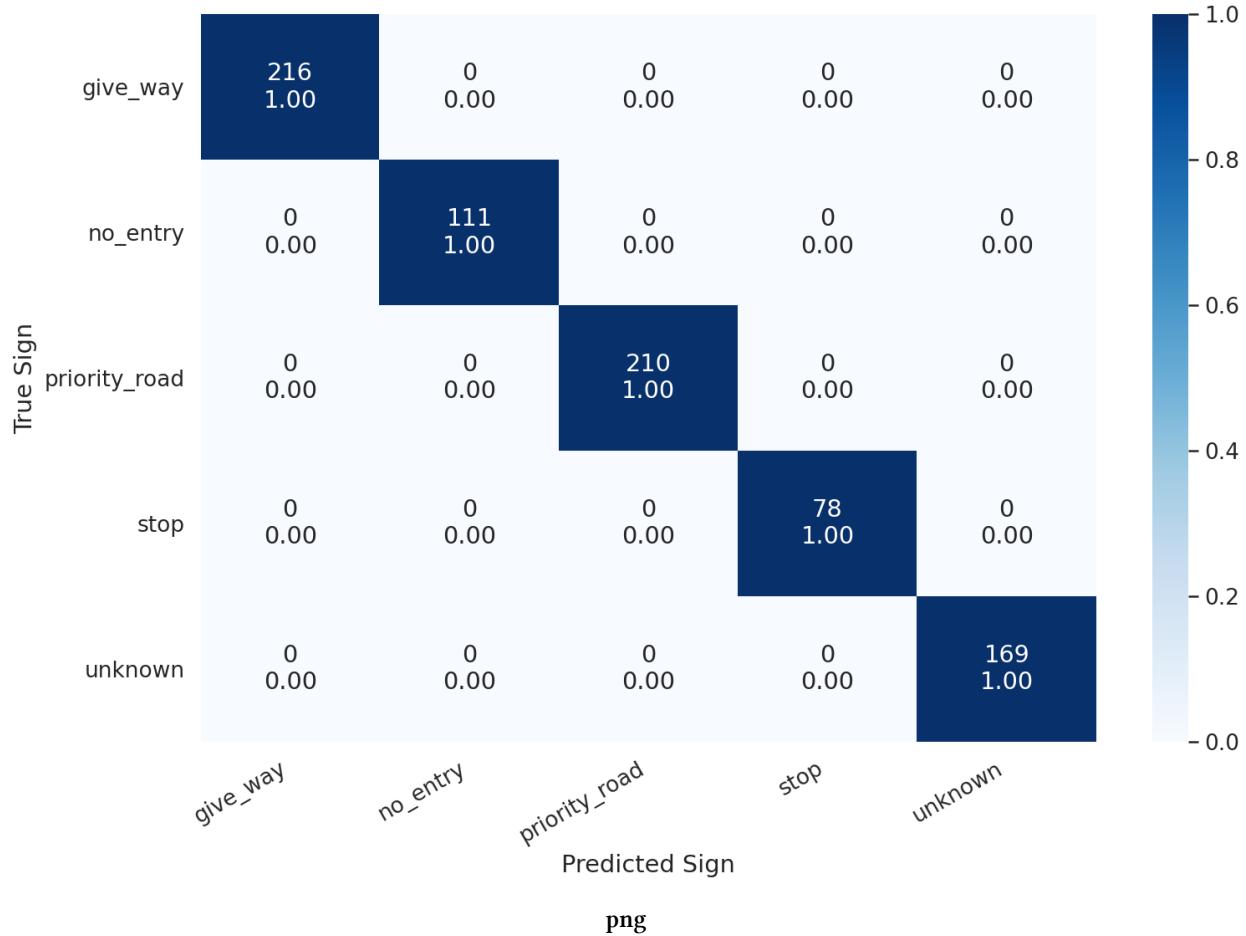
1          precision    recall  f1-score   support
2
3      give_way    1.00    1.00    1.00     216
4      no_entry    1.00    1.00    1.00     111
5  priority_road    1.00    1.00    1.00     210
6       stop    1.00    1.00    1.00      78
7     unknown    1.00    1.00    1.00     169
8
9      accuracy                           1.00      784
10     macro avg    1.00    1.00    1.00      784
11    weighted avg    1.00    1.00    1.00      784

```

```

1 cm = confusion_matrix(y_test, y_pred)
2 show_confusion_matrix(cm, class_names)

```



Our model is still perfect. Go ahead, try it on more images!

## Summary

Good job! You trained two different models for classifying traffic signs from raw pixels. You also built a dataset that is compatible with Torchvision.

- Read the tutorial<sup>10</sup>
- Run the notebook in your browser (Google Colab)<sup>11</sup>
- Read the Getting Things Done with Pytorch book<sup>12</sup>

<sup>10</sup><https://www.curiously.com/posts/transfer-learning-for-image-classification-using-torchvision-pytorch-and-python/>

<sup>11</sup>[https://colab.research.google.com/drive/1Lk5R4pECDxDhd1uXcv26YRQ02fb\\_mrL9?usp=sharing](https://colab.research.google.com/drive/1Lk5R4pECDxDhd1uXcv26YRQ02fb_mrL9?usp=sharing)

<sup>12</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>

Here's what you've learned:

- Overview of the traffic sign image dataset
- Build a dataset
- Use a pre-trained model from Torchvision
- Add a new *unknown* class and re-train the model

Can you use transfer learning for other tasks? How do you do it? Let me know in the comments below.

## References

- ResNet: the intuition behind it<sup>13</sup>
- Understanding ResNet Intuitively<sup>14</sup>
- Conv Nets: A Modular Perspective<sup>15</sup>
- An intuitive guide to Convolutional Neural Networks<sup>16</sup>
- A friendly introduction to Convolutional Neural Networks and Image Recognition<sup>17</sup>
- A Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks<sup>18</sup>
- How to Train an Image Classifier in PyTorch and use it to Perform Basic Inference on Single Images<sup>19</sup>
- Transfer Learning with Convolutional Neural Networks in PyTorch<sup>20</sup>
- Image Classification with Transfer Learning and PyTorch<sup>21</sup>

<sup>13</sup><https://wiseodd.github.io/techblog/2016/10/13/residual-net/>

<sup>14</sup><https://mca.ai/understanding-resnet-intuitively/>

<sup>15</sup><https://colah.github.io/posts/2014-07-Conv-Nets-Modular/>

<sup>16</sup><https://www.freecodecamp.org/news/an-intuitive-guide-to-convolutional-neural-networks-260c2de0a050/>

<sup>17</sup><https://youtu.be/2-Ol7ZB0MmU?t=721>

<sup>18</sup><https://arxiv.org/pdf/1610.02136.pdf>

<sup>19</sup><https://towardsdatascience.com/how-to-train-an-image-classifier-in-pytorch-and-use-it-to-perform-basic-inference-on-single-images-99465a1e9bf5>

<sup>20</sup><https://towardsdatascience.com/transfer-learning-with-convolutional-neural-networks-in-pytorch-dd09190245ce>

<sup>21</sup><https://stackabuse.com/image-classification-with-transfer-learning-and-pytorch/>

# 4. Time Series Forecasting with LSTMs for Daily Coronavirus Cases

This tutorial is NOT trying to build a model that predicts the Covid-19 outbreak/pandemic in the best way possible. This is an example of how you can use Recurrent Neural Networks on some real-world Time Series data with PyTorch. Hopefully, there are much better models that predict the number of daily confirmed cases.

Time series data captures a series of data points recorded at (usually) regular intervals. Some common examples include daily weather temperature, stock prices, and the number of sales a company makes.

Many classical methods (e.g. ARIMA) try to deal with Time Series data with varying success (not to say they are bad at it). In the last couple of years, **Long Short Term Memory Networks (LSTM)**<sup>1</sup> models have become a very useful method when dealing with those types of data.

Recurrent Neural Networks (LSTMs are one type of those) are very good at processing sequences of data. They can “recall” patterns in the data that are very far into the past (or future). In this tutorial, you’re going to learn how to use LSTMs to predict future Coronavirus cases based on real-world data.

- Run the complete notebook in your browser (Google Colab)<sup>2</sup>
- Read the Getting Things Done with Pytorch book<sup>3</sup>

## Novel Coronavirus (COVID-19)

The novel Coronavirus (Covid-19) has spread around the world very rapidly. At the time of this writing, [Worldometers.info](https://www.worldometers.info)<sup>4</sup> shows that there are more than 95,488 confirmed cases in more than 84 countries.

The top 4 worst-affected (by far) are China (the source of the virus), South Korea, Italy, and Iran. Unfortunately, many cases are currently not reported due to:

- A person can get infected without even knowing (asymptomatic)
- Incorrect data reporting
- Not enough test kits
- The symptoms look a lot like the common flu

---

<sup>1</sup>[https://en.wikipedia.org/wiki/Long\\_short-term\\_memory](https://en.wikipedia.org/wiki/Long_short-term_memory)

<sup>2</sup><https://colab.research.google.com/drive/1nQYJq1f7f4R0yeZOzQ9rBKgk00AfLoS0>

<sup>3</sup><https://github.com/cvutkor/Getting-Things-Done-with-Pytorch>

<sup>4</sup><https://www.worldometers.info/coronavirus/>

## How dangerous is this virus?

Except for the common statistics you might see cited on the news, there are some good and some bad news:

- More than 80% of the confirmed cases recover without any need of medical attention
- [3.4% Mortality Rate estimate by the World Health Organization \(WHO\) as of March 3<sup>5</sup>](https://www.worldometers.info/coronavirus/coronavirus-death-rate/#who-03-03-20)
- The reproductive number which represents the average number of people to which a single infected person will transmit the virus is between 1.4 and 2.5 ([WHO's estimated on Jan. 23<sup>6</sup>](https://www.worldometers.info/coronavirus/#repro))

The last one is really scary. It sounds like we can witness some crazy exponential growth if appropriate measures are not put in place.

Let's get started!

```

1 %reload_ext watermark
2 %watermark -v -p numpy,pandas,torch

1 CPython 3.6.9
2 IPython 5.5.0
3
4 numpy 1.17.5
5 pandas 0.25.3
6 torch 1.4.0

1 import torch
2
3 import os
4 import numpy as np
5 import pandas as pd
6 from tqdm import tqdm
7 import seaborn as sns
8 from pylab import rcParams
9 import matplotlib.pyplot as plt
10 from matplotlib import rc
11 from sklearn.preprocessing import MinMaxScaler
12 from pandas.plotting import register_matplotlib_converters
13 from torch import nn, optim
14
```

<sup>5</sup><https://www.worldometers.info/coronavirus/coronavirus-death-rate/#who-03-03-20>

<sup>6</sup><https://www.worldometers.info/coronavirus/#repro>

```

15 %matplotlib inline
16 %config InlineBackend.figure_format='retina'
17
18 sns.set(style='whitegrid', palette='muted', font_scale=1.2)
19
20 HAPPY_COLORS_PALETTE = ['#01BEFE', '#FFDD00', '#FF7D00', '#FF006D', '#93D30C', '#8F0\OFF']
21
22
23 sns.set_palette(sns.color_palette(HAPPY_COLORS_PALETTE))
24
25 rcParams['figure.figsize'] = 14, 10
26 register_matplotlib_converters()
27
28 RANDOM_SEED = 42
29 np.random.seed(RANDOM_SEED)
30 torch.manual_seed(RANDOM_SEED)

1 <torch._C.Generator at 0x7faeaa744d30>

```

## Daily Cases Dataset

The data is provided by the Johns Hopkins University Center for Systems Science and Engineering (JHU CSSE) and contains the number of reported daily cases by country. [The dataset is available on GitHub<sup>7</sup>](#) and is updated regularly.

We're going to take the Time Series data only for confirmed cases (number of deaths and recovered cases are also available):

```

1 # !wget https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_\
2 19_data/csse_covid_19_time_series/time_series_19-covid-Confirmed.csv

```

Or you can take the same dataset that I've used for this tutorial (the data snapshot is from 3 March 2020):

```

1 !gdown --id 1AsfdLrGESCQnRW5rbMz56A1KBc3Fe5aV

```

## Data exploration

Let's load the data and have a peek:

---

<sup>7</sup><https://github.com/CSSEGISandData/COVID-19>

```
1 df = pd.read_csv('time_series_19-covid-Confirmed.csv')
```

Two things to note here:

- The data contains a province, country, latitude, and longitude. We won't be needing those.
- The number of cases is cumulative. We'll undo the accumulation.

Let's start by getting rid of the first four columns:

```
1 df = df.iloc[:, 4:]
```

Let's check for missing values:

```
1 df.isnull().sum().sum()
```

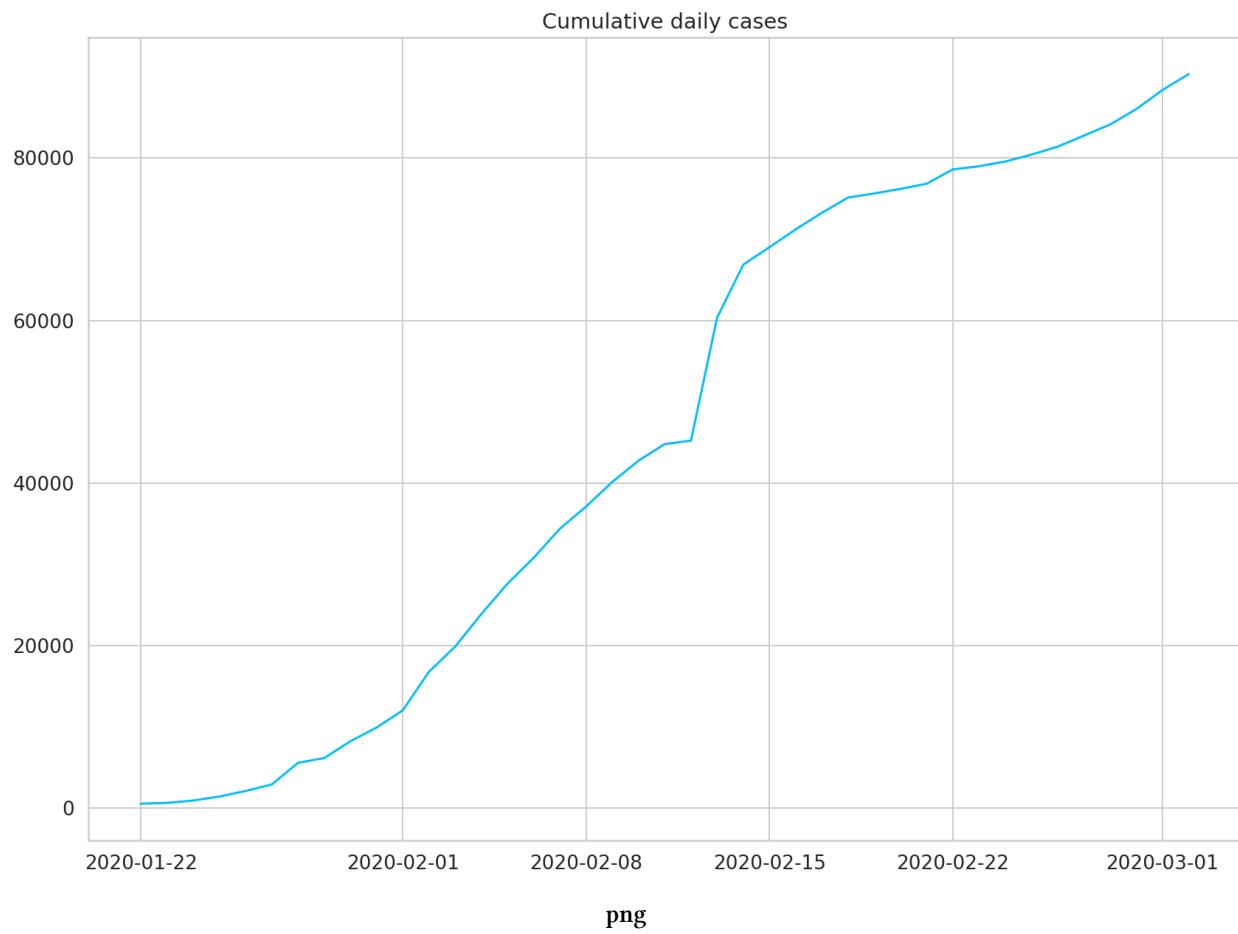
```
1 0
```

Everything seems to be in place. Let's sum all rows, so we get the cumulative daily cases:

```
1 daily_cases = df.sum(axis=0)
2 daily_cases.index = pd.to_datetime(daily_cases.index)
3 daily_cases.head()
```

```
1 2020-01-22      555
2 2020-01-23      653
3 2020-01-24      941
4 2020-01-25     1434
5 2020-01-26     2118
6 dtype: int64
```

```
1 plt.plot(daily_cases)
2 plt.title("Cumulative daily cases");
```



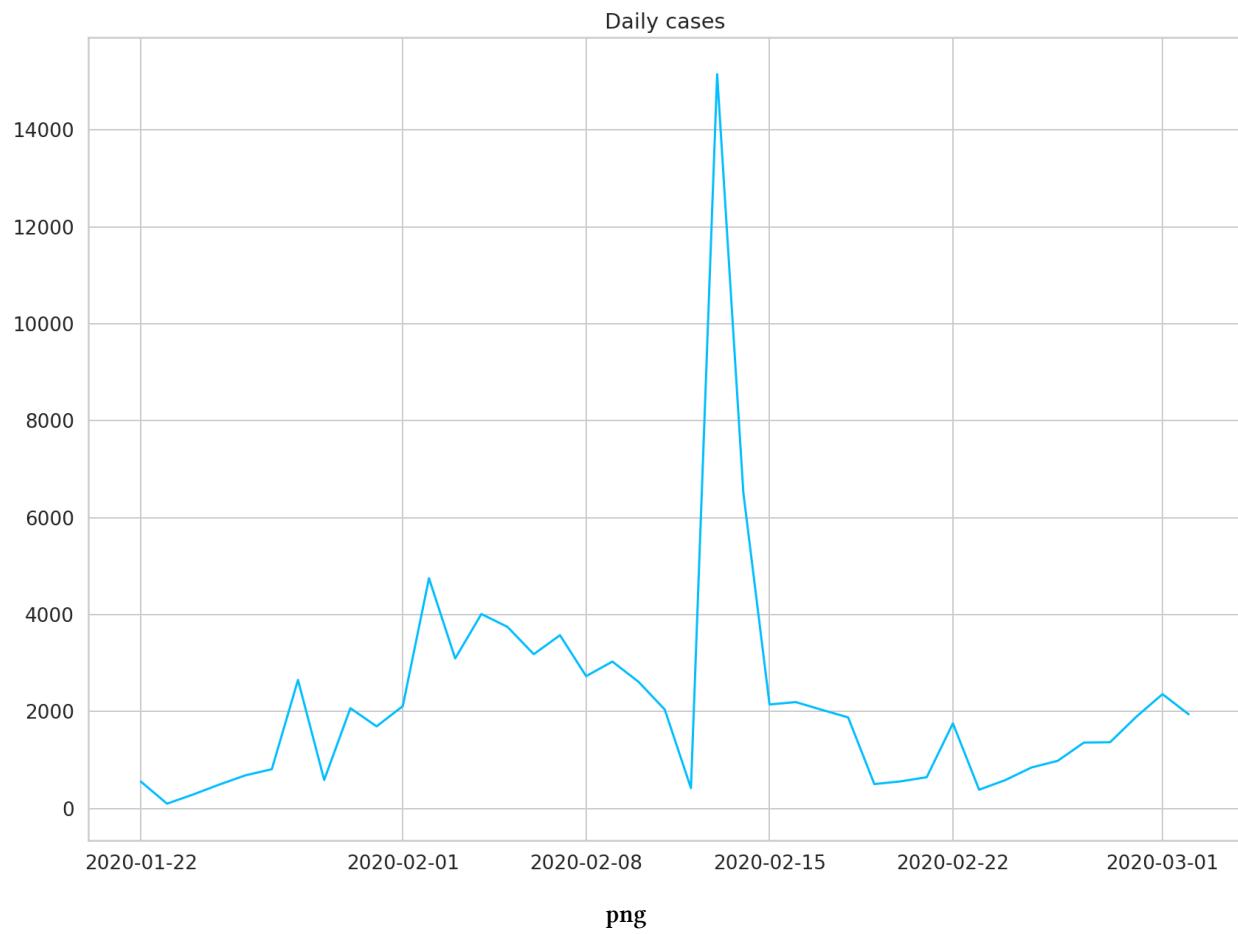
png

We'll undo the accumulation by subtracting the current value from the previous. We'll preserve the first value of the sequence:

```
1 daily_cases = daily_cases.diff().fillna(daily_cases[0]).astype(np.int64)
2 daily_cases.head()
```

```
1 2020-01-22    555
2 2020-01-23     98
3 2020-01-24   288
4 2020-01-25   493
5 2020-01-26   684
6 dtype: int64
```

```
1 plt.plot(daily_cases)
2 plt.title("Daily cases");
```



The huge spike (in the middle) is mostly due to a change of criteria for testing patients in China. This will certainly be a challenge for our model.

Let's check the amount of data we have:

```
1 daily_cases.shape
```

```
1 (41, )
```

Unfortunately, we have data for only 41 days. Let's see what we can do with it!

## Preprocessing

We'll reserve the first 27 days for training and use the rest for testing:

```

1 test_data_size = 14
2
3 train_data = daily_cases[:-test_data_size]
4 test_data = daily_cases[-test_data_size:]
5
6 train_data.shape
1 (27, )

```

We have to scale the data (values will be between 0 and 1) if we want to increase the training speed and performance of the model. We'll use the `MinMaxScaler` from scikit-learn:

```

1 scaler = MinMaxScaler()
2
3 scaler = scaler.fit(np.expand_dims(train_data, axis=1))
4
5 train_data = scaler.transform(np.expand_dims(train_data, axis=1))
6
7 test_data = scaler.transform(np.expand_dims(test_data, axis=1))

```

Currently, we have a big sequence of daily cases. We'll convert it into smaller ones:

```

1 def create_sequences(data, seq_length):
2     xs = []
3     ys = []
4
5     for i in range(len(data)-seq_length-1):
6         x = data[i:(i+seq_length)]
7         y = data[i+seq_length]
8         xs.append(x)
9         ys.append(y)
10
11    return np.array(xs), np.array(ys)

```

```
1 seq_length = 5
2 X_train, y_train = create_sequences(train_data, seq_length)
3 X_test, y_test = create_sequences(test_data, seq_length)
4
5 X_train = torch.from_numpy(X_train).float()
6 y_train = torch.from_numpy(y_train).float()
7
8 X_test = torch.from_numpy(X_test).float()
9 y_test = torch.from_numpy(y_test).float()
```

Each training example contains a sequence of 5 data points of history and a label for the real value that our model needs to predict. Let's dive in:

```
1 X_train.shape
```

```
1 torch.Size([21, 5, 1])
```

```
1 X_train[:2]
```

```
1 tensor([[[0.0304],
2           [0.0000],
3           [0.0126],
4           [0.0262],
5           [0.0389]],
6
7           [[0.0000],
8           [0.0126],
9           [0.0262],
10          [0.0389],
11          [0.0472]]])
```

```
1 y_train.shape
```

```
1 torch.Size([21, 1])
```

```
1 y_train[:2]
```

```

1 tensor([[0.0472,
2             [0.1696]]))

1 train_data[:10]

1 array([[0.03036545],
2        [0.          ],
3        [0.01262458],
4        [0.02624585],
5        [0.03893688],
6        [0.04724252],
7        [0.16963455],
8        [0.03255814],
9        [0.13089701],
10       [0.10598007]])

```

## Building a model

We'll encapsulate the complexity of our model into a class that extends from `torch.nn.Module`:

```

1 class CoronaVirusPredictor(nn.Module):
2
3     def __init__(self, n_features, n_hidden, seq_len, n_layers=2):
4         super(CoronaVirusPredictor, self).__init__()
5
6         self.n_hidden = n_hidden
7         self.seq_len = seq_len
8         self.n_layers = n_layers
9
10        self.lstm = nn.LSTM(
11            input_size=n_features,
12            hidden_size=n_hidden,
13            num_layers=n_layers,
14            dropout=0.5
15        )
16
17        self.linear = nn.Linear(in_features=n_hidden, out_features=1)
18
19    def reset_hidden_state(self):

```

```

20     self.hidden = (
21         torch.zeros(self.n_layers, self.seq_len, self.n_hidden),
22         torch.zeros(self.n_layers, self.seq_len, self.n_hidden)
23     )
24
25     def forward(self, sequences):
26         lstm_out, self.hidden = self.lstm(
27             sequences.view(len(sequences), self.seq_len, -1),
28             self.hidden
29         )
30         last_time_step = \
31             lstm_out.view(self.seq_len, len(sequences), self.n_hidden)[-1]
32         y_pred = self.linear(last_time_step)
33         return y_pred

```

Our CoronaVirusPredictor contains 3 methods:

- constructor - initialize all helper data and create the layers
- reset\_hidden\_state - we'll use a stateless LSTM, so we need to reset the state after each example
- forward - get the sequences, pass all of them through the LSTM layer, at once. We take the output of the last time step and pass it through our linear layer to get the prediction.

## Training

Let's build a helper function for the training of our model (we'll reuse it later):

```

1  def train_model(
2      model,
3      train_data,
4      train_labels,
5      test_data=None,
6      test_labels=None
7  ):
8      loss_fn = torch.nn.MSELoss(reduction='sum')
9
10     optimiser = torch.optim.Adam(model.parameters(), lr=1e-3)
11     num_epochs = 60
12
13     train_hist = np.zeros(num_epochs)
14     test_hist = np.zeros(num_epochs)

```

```

15
16     for t in range(num_epochs):
17         model.reset_hidden_state()
18
19         y_pred = model(X_train)
20
21         loss = loss_fn(y_pred.float(), y_train)
22
23         if test_data is not None:
24             with torch.no_grad():
25                 y_test_pred = model(X_test)
26                 test_loss = loss_fn(y_test_pred.float(), y_test)
27                 test_hist[t] = test_loss.item()
28
29             if t % 10 == 0:
30                 print(f'Epoch {t} train loss: {loss.item()} test loss: {test_loss.item()}')
31         elif t % 10 == 0:
32             print(f'Epoch {t} train loss: {loss.item()}')
33
34         train_hist[t] = loss.item()
35
36         optimiser.zero_grad()
37
38         loss.backward()
39
40         optimiser.step()
41
42     return model.eval(), train_hist, test_hist

```

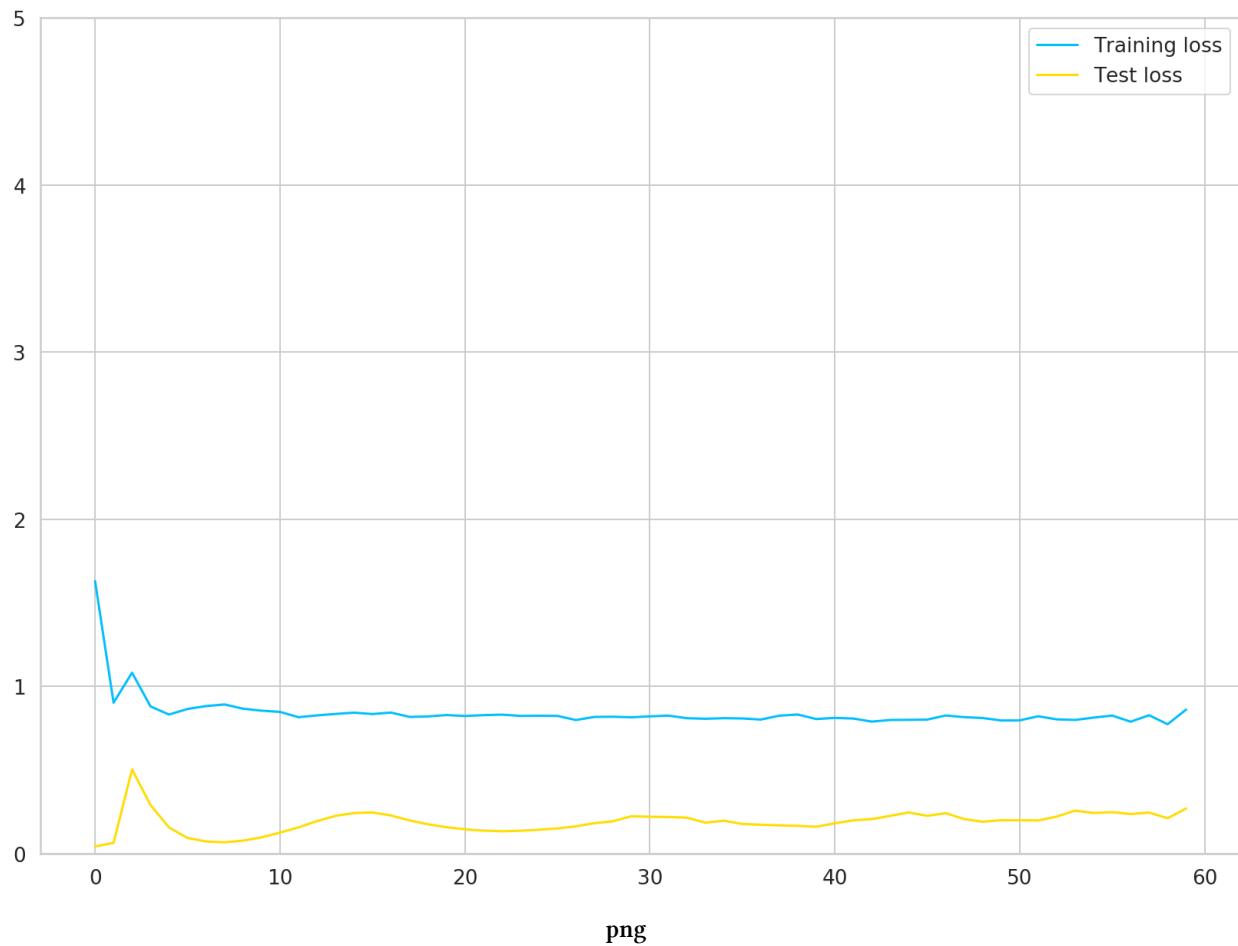
Note that the hidden state is reset at the start of each epoch. We don't use batches of data our model sees every example at once. We'll use mean squared error to measure our training and test error. We'll record both.

Let's create an instance of our model and train it:

```
1 model = CoronaVirusPredictor(  
2     n_features=1,  
3     n_hidden=512,  
4     seq_len=seq_length,  
5     n_layers=2  
6 )  
7 model, train_hist, test_hist = train_model(  
8     model,  
9     X_train,  
10    y_train,  
11    X_test,  
12    y_test  
13 )  
  
1 Epoch 0 train loss: 1.6297188997268677 test loss: 0.041186608374118805  
2 Epoch 10 train loss: 0.8466923832893372 test loss: 0.12416432797908783  
3 Epoch 20 train loss: 0.8219934105873108 test loss: 0.1438201516866684  
4 Epoch 30 train loss: 0.8200693726539612 test loss: 0.2190694659948349  
5 Epoch 40 train loss: 0.810839056968689 test loss: 0.1797715127468109  
6 Epoch 50 train loss: 0.795730471611023 test loss: 0.19855864346027374
```

Let's have a look at the train and test loss:

```
1 plt.plot(train_hist, label="Training loss")  
2 plt.plot(test_hist, label="Test loss")  
3 plt.ylim((0, 5))  
4 plt.legend();
```



Our model's performance doesn't improve after 15 epochs or so. Recall that we have very little data. Maybe we shouldn't trust our model that much?

## Predicting daily cases

Our model can (due to the way we've trained it) predict only a single day in the future. We'll employ a simple strategy to overcome this limitation. Use predicted values as input for predicting the next days:

```

1 with torch.no_grad():
2     test_seq = X_test[:1]
3     preds = []
4     for _ in range(len(X_test)):
5         y_test_pred = model(test_seq)
6         pred = torch.flatten(y_test_pred).item()
7         preds.append(pred)
8         new_seq = test_seq.numpy().flatten()
9         new_seq = np.append(new_seq, [pred])
10        new_seq = new_seq[1:]
11        test_seq = torch.as_tensor(new_seq).view(1, seq_length, 1).float()

```

We have to reverse the scaling of the test data and the model predictions:

```

1 true_cases = scaler.inverse_transform(
2     np.expand_dims(y_test.flatten().numpy(), axis=0)
3 ).flatten()
4
5 predicted_cases = scaler.inverse_transform(
6     np.expand_dims(preds, axis=0)
7 ).flatten()

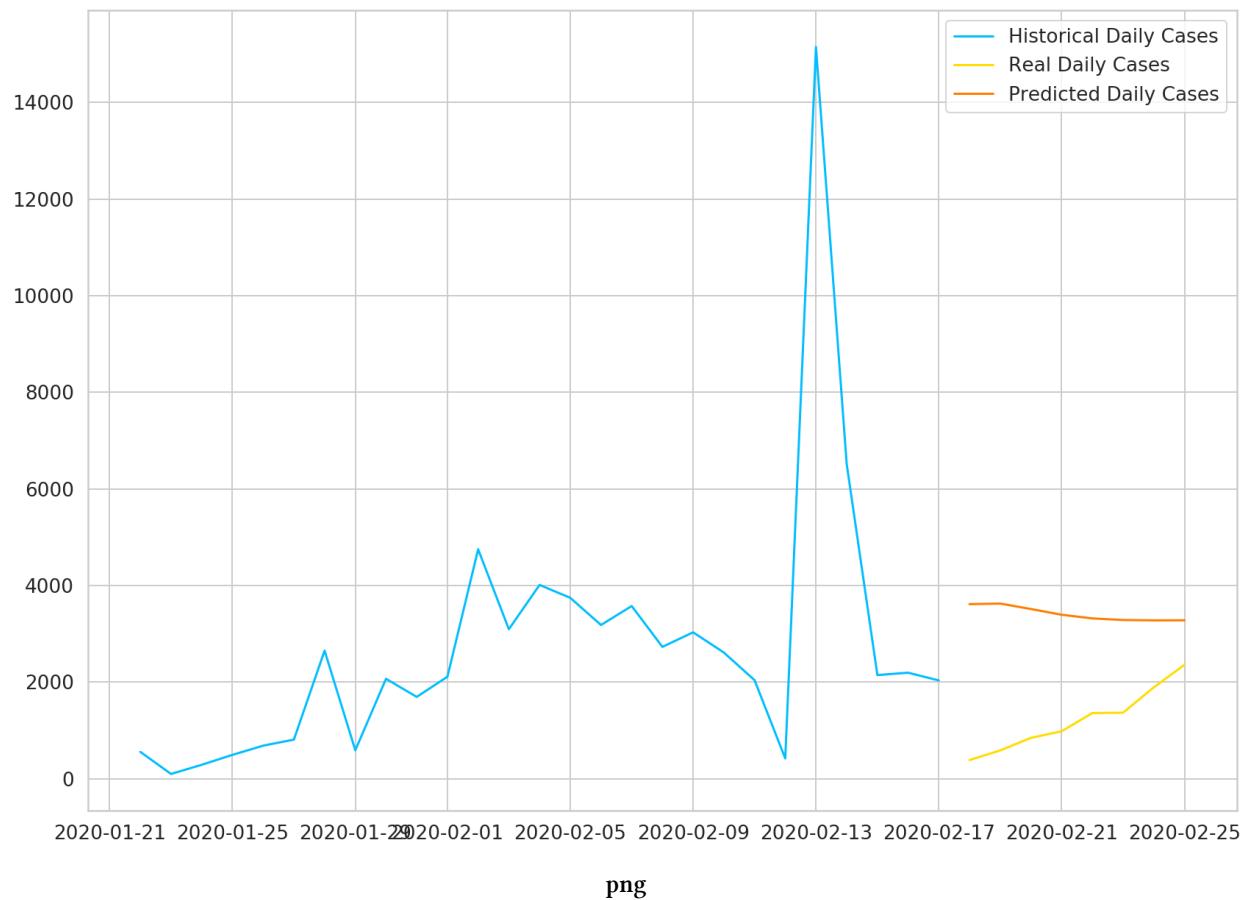
```

Let's look at the results:

```

1 plt.plot(
2     daily_cases.index[:len(train_data)],
3     scaler.inverse_transform(train_data).flatten(),
4     label='Historical Daily Cases'
5 )
6
7 plt.plot(
8     daily_cases.index[len(train_data):len(train_data) + len(true_cases)],
9     true_cases,
10    label='Real Daily Cases'
11 )
12
13 plt.plot(
14     daily_cases.index[len(train_data):len(train_data) + len(true_cases)],
15     predicted_cases,
16     label='Predicted Daily Cases'
17 )
18
19 plt.legend();

```



As expected, our model doesn't perform very well. That said, the predictions seem to be in the right ballpark (probably due to using the last data point as a strong predictor for the next).

## Use all data for training

Now, we'll use all available data to train the same model:

```

1  scaler = MinMaxScaler()
2
3  scaler = scaler.fit(np.expand_dims(daily_cases, axis=1))
4
5  all_data = scaler.transform(np.expand_dims(daily_cases, axis=1))
6
7  all_data.shape

```

```
1 (41, 1)
```

The preprocessing and training steps are the same:

```
1 X_all, y_all = create_sequences(all_data, seq_length)
2
3 X_all = torch.from_numpy(X_all).float()
4 y_all = torch.from_numpy(y_all).float()
5
6 model = CoronaVirusPredictor(
7     n_features=1,
8     n_hidden=512,
9     seq_len=seq_length,
10    n_layers=2
11 )
12 model, train_hist, _ = train_model(model, X_all, y_all)

1 Epoch 0 train loss: 1.9441421031951904
2 Epoch 10 train loss: 0.8385428786277771
3 Epoch 20 train loss: 0.8256545066833496
4 Epoch 30 train loss: 0.8023681640625
5 Epoch 40 train loss: 0.8125611543655396
6 Epoch 50 train loss: 0.8225002884864807
```

## Predicting future cases

We'll use our "fully trained" model to predict the confirmed cases for 12 days into the future:

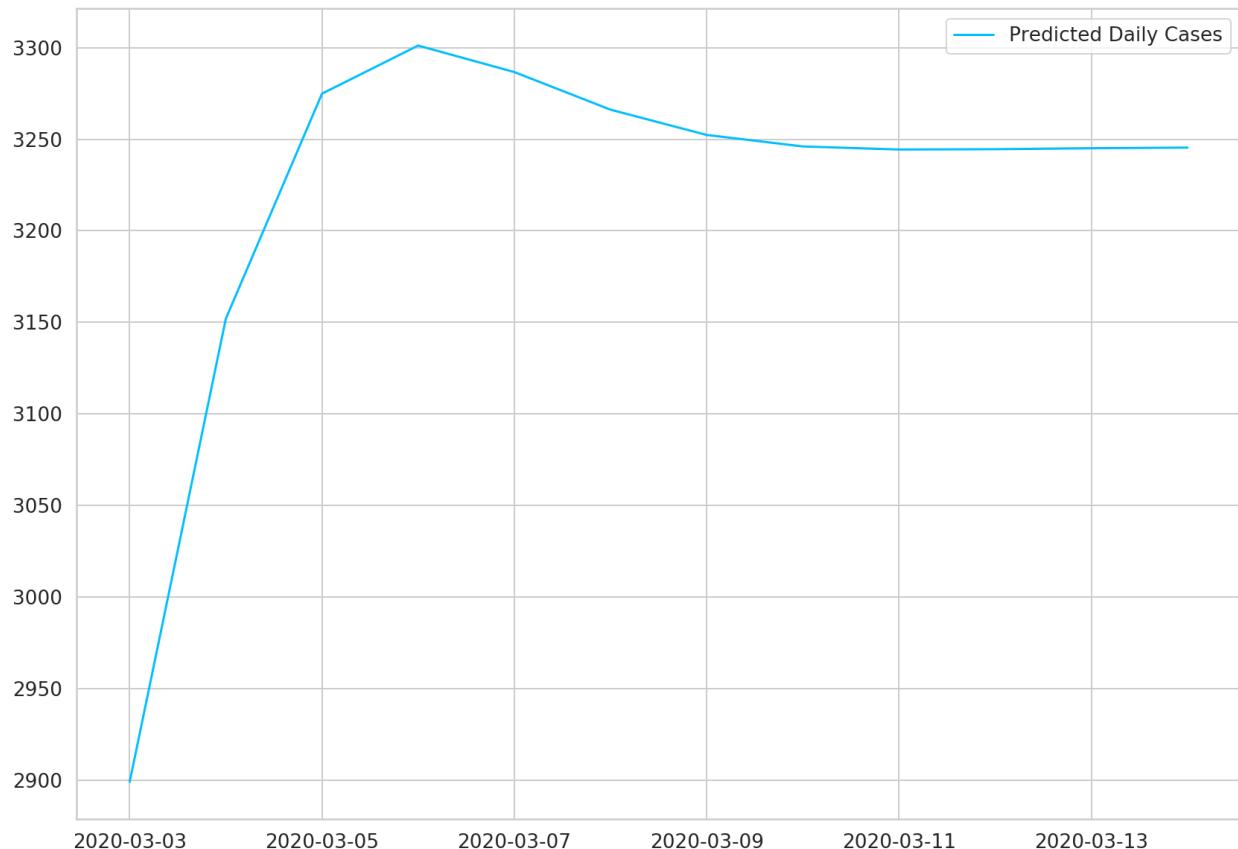
```
1 DAYS_TO_PREDICT = 12
2
3 with torch.no_grad():
4     test_seq = X_all[:1]
5     preds = []
6     for _ in range(DAYS_TO_PREDICT):
7         y_test_pred = model(test_seq)
8         pred = torch.flatten(y_test_pred).item()
9         preds.append(pred)
10        new_seq = test_seq.numpy().flatten()
11        new_seq = np.append(new_seq, [pred])
12        new_seq = new_seq[1:]
13        test_seq = torch.as_tensor(new_seq).view(1, seq_length, 1).float()
```

As before, we'll inverse the scalar transformation:

```
1 predicted_cases = scaler.inverse_transform(  
2     np.expand_dims(preds, axis=0)  
3 ).flatten()
```

To create a cool chart with the historical and predicted cases, we need to extend the date index of our data frame:

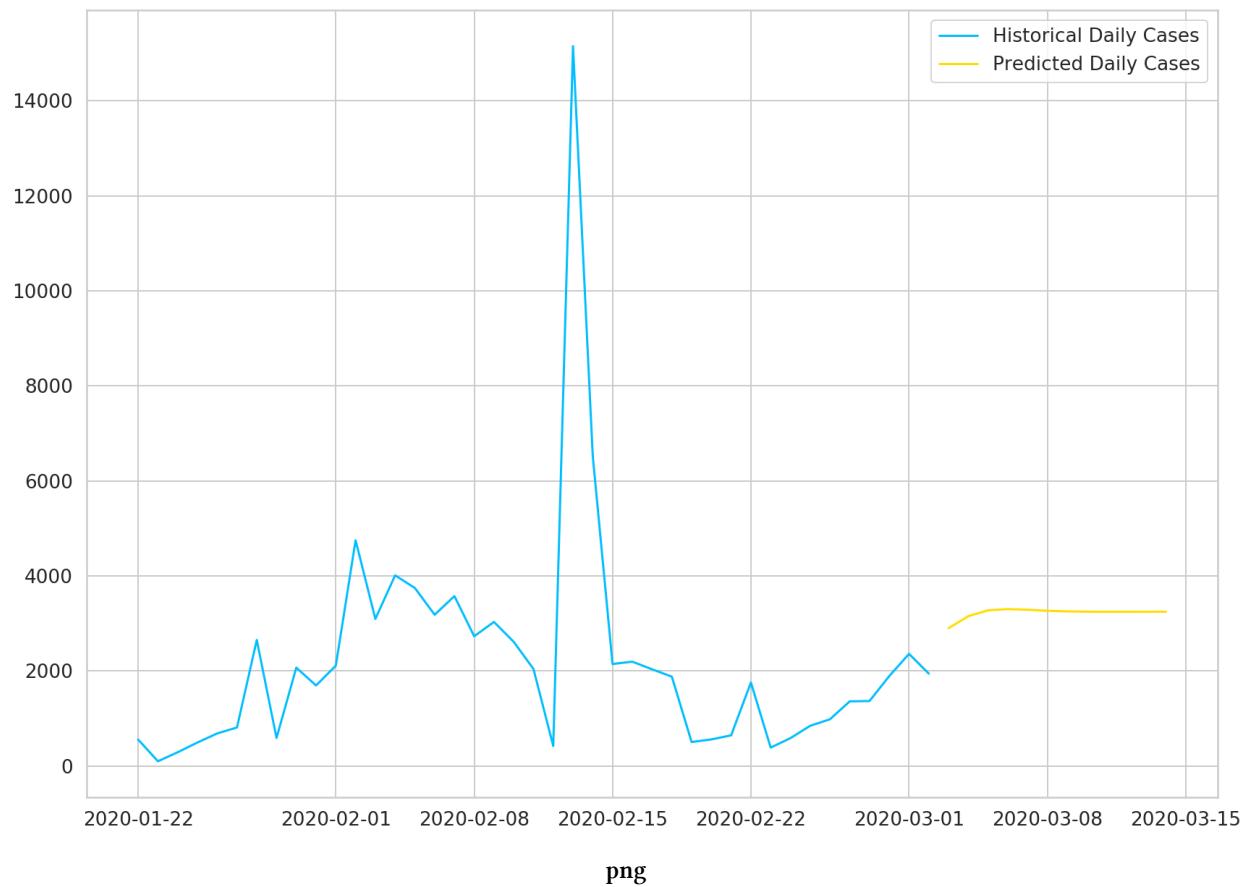
```
1 daily_cases.index[-1]  
  
1 Timestamp('2020-03-02 00:00:00')  
  
1 predicted_index = pd.date_range(  
2     start=daily_cases.index[-1],  
3     periods=DAYS_TO_PREDICT + 1,  
4     closed='right'  
5 )  
6  
7 predicted_cases = pd.Series(  
8     data=predicted_cases,  
9     index=predicted_index  
10 )  
11  
12 plt.plot(predicted_cases, label='Predicted Daily Cases')  
13 plt.legend();
```



png

Now we can use all the data to plot the results:

```
1 plt.plot(daily_cases, label='Historical Daily Cases')
2 plt.plot(predicted_cases, label='Predicted Daily Cases')
3 plt.legend();
```



Our model thinks that things will level off. Note that the more you go into the future, the more you shouldn't trust your model predictions.

## Conclusion

Well done! You learned how to use PyTorch to create a Recurrent Neural Network that works with Time Series data. The model performance is not that great, but this is expected, given the small amounts of data.

- Run the complete notebook in your browser (Google Colab)<sup>8</sup>
- Read the Getting Things Done with Pytorch book<sup>9</sup>

The problem of predicting daily Covid-19 cases is a hard one. We're amidst an outbreak, and there's more to be done. Hopefully, everything will be back to normal after some time.

<sup>8</sup><https://colab.research.google.com/drive/1nQYJq1f7f4R0yeZOzQ9rBKgk00AfLoS0>

<sup>9</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>

## References

- Sequence Models PyTorch Tutorial<sup>10</sup>
- LSTM for time series prediction<sup>11</sup>
- Time Series Prediction using LSTM with PyTorch in Python<sup>12</sup>
- Stateful LSTM in Keras<sup>13</sup>
- LSTMs for Time Series in PyTorch<sup>14</sup>
- Novel Coronavirus (COVID-19) Cases, provided by JHU CSSE<sup>15</sup>
- covid-19-analysis<sup>16</sup>
- How does Coronavirus compare to Ebola, SARS, etc?<sup>17</sup>
- Worldometer COVID-19 Coronavirus Outbreak<sup>18</sup>
- How contagious is the Wuhan Coronavirus? (Ro)<sup>19</sup>
- Systemic Risk of Pandemic via Novel Pathogens - Coronavirus: A Note<sup>20</sup>
- Statistical Consequences of Fat Tails: Real World Preasymptotics, Epistemology, and Applications<sup>21</sup>

<sup>10</sup>[https://pytorch.org/tutorials/beginner/nlp/sequence\\_models\\_tutorial.html](https://pytorch.org/tutorials/beginner/nlp/sequence_models_tutorial.html)

<sup>11</sup><https://towardsdatascience.com/lstm-for-time-series-prediction-de8aeb26f2ca>

<sup>12</sup><https://stackabuse.com/time-series-prediction-using-lstm-with-pytorch-in-python/>

<sup>13</sup><https://philipperemy.github.io/keras-stateful-lstm/>

<sup>14</sup><https://www.jessicayung.com/lstms-for-time-series-in-pytorch/>

<sup>15</sup><https://github.com/CSSEGISandData/COVID-19>

<sup>16</sup><https://github.com/AaronWard/covid-19-analysis>

<sup>17</sup><https://www.youtube.com/watch?v=6dDD2tHWnU>

<sup>18</sup><https://www.worldometers.info/coronavirus/>

<sup>19</sup><https://www.worldometers.info/coronavirus/#repro>

<sup>20</sup>[https://www.academia.edu/41743064/Systemic\\_Risk\\_of\\_Pandemic\\_via\\_Novel\\_Pathogens\\_-\\_Coronavirus\\_A\\_Note](https://www.academia.edu/41743064/Systemic_Risk_of_Pandemic_via_Novel_Pathogens_-_Coronavirus_A_Note)

<sup>21</sup><https://www.researchers.one/article/2020-01-21>

# 5. Time Series Anomaly Detection using LSTM Autoencoders

TL;DR Use real-world Electrocardiogram (ECG) data to detect anomalies in a patient heartbeat. We'll build an LSTM Autoencoder, train it on a set of normal heartbeats and classify unseen examples as normal or anomalies

In this tutorial, you'll learn how to detect anomalies in Time Series data using an LSTM Autoencoder. You're going to use real-world ECG data from a single patient with heart disease to detect abnormal hearbeats.

- Run the complete notebook in your browser (Google Colab)<sup>1</sup>
- Read the Getting Things Done with Pytorch book<sup>2</sup>

By the end of this tutorial, you'll learn how to:

- Prepare a dataset for Anomaly Detection from Time Series Data
- Build an LSTM Autoencoder with PyTorch
- Train and evaluate your model
- Choose a threshold for anomaly detection
- Classify unseen examples as normal or anomaly

## Data

The [dataset<sup>3</sup>](#) contains 5,000 Time Series examples (obtained with ECG) with 140 timesteps. Each sequence corresponds to a single heartbeat from a single patient with congestive heart failure.

An electrocardiogram (ECG or EKG) is a test that checks how your heart is functioning by measuring the electrical activity of the heart. With each heart beat, an electrical impulse (or wave) travels through your heart. This wave causes the muscle to squeeze and pump blood from the heart. [Source<sup>4</sup>](#)

We have 5 types of hearbeats (classes):

<sup>1</sup>[https://colab.research.google.com/drive/1\\_J2MrBSvsJfOcVmYAN2-WSp36BtsFZCa](https://colab.research.google.com/drive/1_J2MrBSvsJfOcVmYAN2-WSp36BtsFZCa)

<sup>2</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>

<sup>3</sup><http://timeseriesclassification.com/description.php?Dataset=ECG5000>

<sup>4</sup><https://www.heartandstroke.ca/heart/tests/electrocardiogram>

- Normal (N)
- R-on-T Premature Ventricular Contraction (R-on-T PVC)
- Premature Ventricular Contraction (PVC)
- Supra-ventricular Premature or Ectopic Beat (SP or EB)
- Unclassified Beat (UB).

Assuming a healthy heart and a typical rate of 70 to 75 beats per minute, each cardiac cycle, or heartbeat, takes about 0.8 seconds to complete the cycle. Frequency: 60–100 per minute (Humans) Duration: 0.6–1 second (Humans) [Source<sup>5</sup>](#)

The dataset is available on my Google Drive. Let's get it:

```
1 !gdown --id 16M1eqoIr1vYx1Gk4GKnGmrsCPuWkkpT
1 !unzip -qq ECG5000.zip
1 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
```

The data comes in multiple formats. We'll load the `arff` files into Pandas data frames:

```
1 with open('ECG5000_TRAIN.arff') as f:
2     train = a2p.load(f)
3
4 with open('ECG5000_TEST.arff') as f:
5     test = a2p.load(f)
```

We'll combine the training and test data into a single data frame. This will give us more data to train our Autoencoder. We'll also shuffle it:

```
1 df = train.append(test)
2 df = df.sample(frac=1.0)
3 df.shape
1 (5000, 141)
```

We have 5,000 examples. Each row represents a single heartbeat record. Let's name the possible classes:

---

<sup>5</sup>[https://en.wikipedia.org/wiki/Cardiac\\_cycle](https://en.wikipedia.org/wiki/Cardiac_cycle)

```
1 CLASS_NORMAL = 1
2
3 class_names = ['Normal', 'R on T', 'PVC', 'SP', 'UB']
```

Next, we'll rename the last column to target, so its easier to reference it:

```
1 new_columns = list(df.columns)
2 new_columns[-1] = 'target'
3 df.columns = new_columns
```

## Exploratory Data Analysis

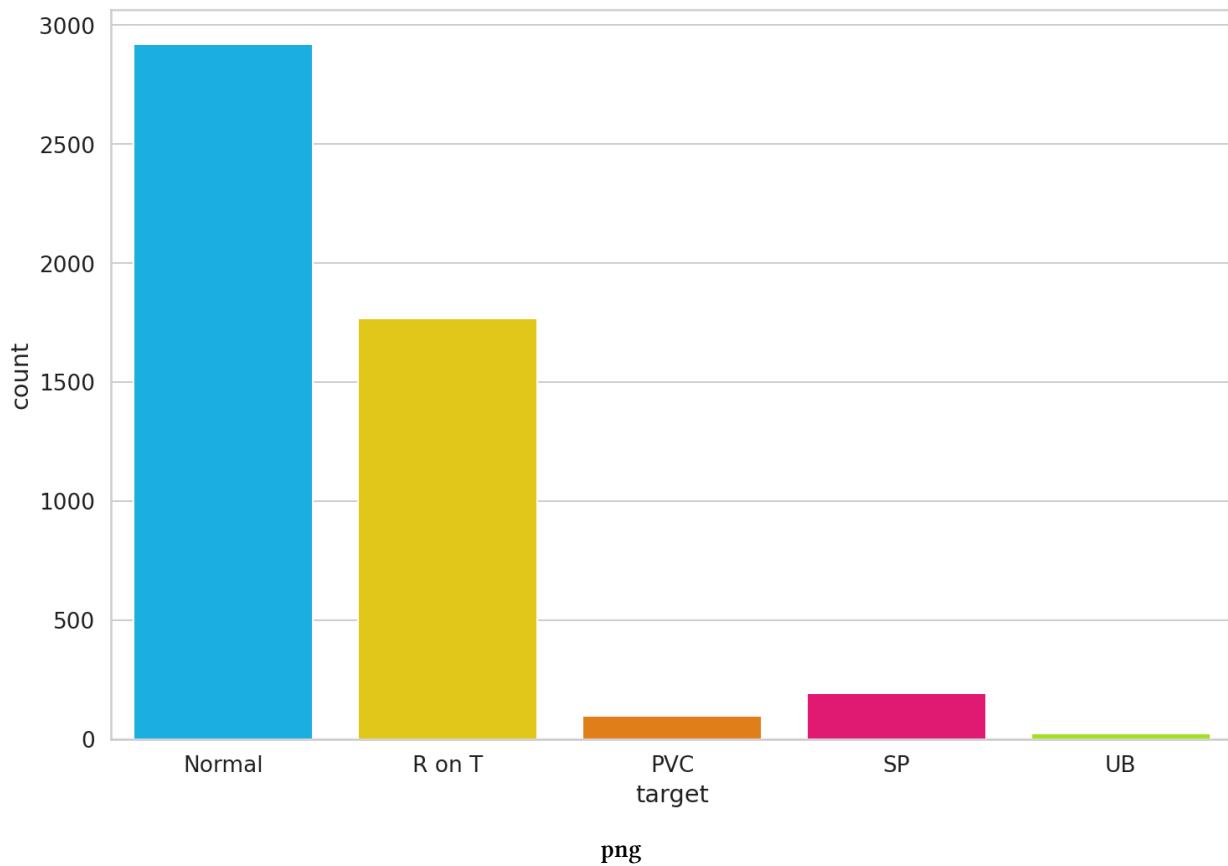
Let's check how many examples for each heartbeat class do we have:

```
1 df.target.value_counts()
```

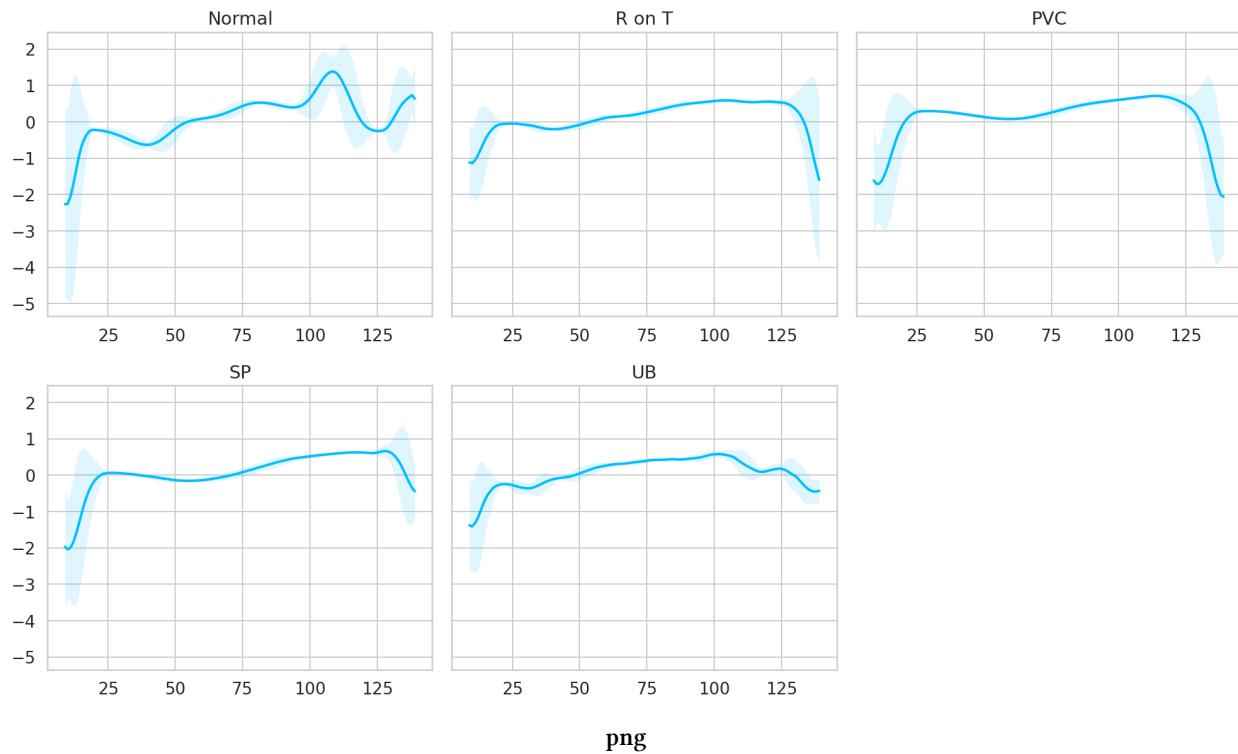
```
1 1    2919
2 2    1767
3 4     194
4 3      96
5 5      24
6 Name: target, dtype: int64
```

Let's plot the results:



The normal class, has by far, the most examples. This is great because we'll use it to train our model.

Let's have a look at an averaged (smoothed out with one standard deviation on top and bottom of it) Time Series for each class:



It is very good that the normal class has a distinctly different pattern than all other classes. Maybe our model will be able to detect anomalies?

## LSTM Autoencoder

The [Autoencoder's<sup>6</sup>](#) job is to get some input data, pass it through the model, and obtain a reconstruction of the input. The reconstruction should match the input as much as possible. The trick is to use a small number of parameters, so your model learns a compressed representation of the data.

In a sense, Autoencoders try to learn only the most important features (compressed version) of the data. Here, we'll have a look at how to feed Time Series data to an Autoencoder. We'll use a couple of LSTM layers (hence the LSTM Autoencoder) to capture the temporal dependencies of the data.

To classify a sequence as normal or an anomaly, we'll pick a threshold above which a heartbeat is considered abnormal.

## Reconstruction Loss

When training an Autoencoder, the objective is to reconstruct the input as best as possible. This is done by minimizing a loss function (just like in supervised learning). This function is known as

<sup>6</sup><https://en.wikipedia.org/wiki/Autoencoder>

*reconstruction loss.* Cross-entropy loss and Mean squared error are common examples.

## Anomaly Detection in ECG Data

We'll use normal heartbeats as training data for our model and record the *reconstruction loss*. But first, we need to prepare the data:

### Data Preprocessing

Let's get all normal heartbeats and drop the target (class) column:

```
1 normal_df = df[df.target == str(CLASS_NORMAL)].drop(labels='target', axis=1)
2 normal_df.shape

1 (2919, 140)
```

We'll merge all other classes and mark them as anomalies:

```
1 anomaly_df = df[df.target != str(CLASS_NORMAL)].drop(labels='target', axis=1)
2 anomaly_df.shape

1 (2081, 140)
```

We'll split the normal examples into train, validation and test sets:

```
1 train_df, val_df = train_test_split(
2     normal_df,
3     test_size=0.15,
4     random_state=RANDOM_SEED
5 )
6
7 val_df, test_df = train_test_split(
8     val_df,
9     test_size=0.33,
10    random_state=RANDOM_SEED
11 )
```

We need to convert our examples into tensors, so we can use them to train our Autoencoder. Let's write a helper function for that:

```

1 def create_dataset(df):
2
3     sequences = df.astype(np.float32).to_numpy().tolist()
4
5     dataset = [torch.tensor(s).unsqueeze(1).float() for s in sequences]
6
7     n_seq, seq_len, n_features = torch.stack(dataset).shape
8
9     return dataset, seq_len, n_features

```

Each Time Series will be converted to a 2D Tensor in the shape *sequence length x number of features* (140x1 in our case).

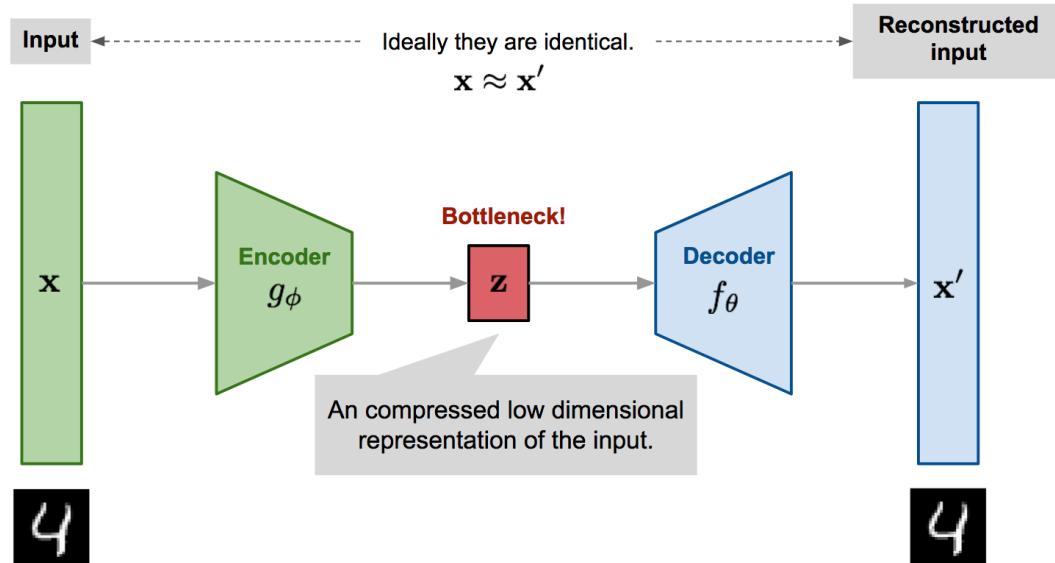
Let's create some datasets:

```

1 train_dataset, seq_len, n_features = create_dataset(train_df)
2 val_dataset, _, _ = create_dataset(val_df)
3 test_normal_dataset, _, _ = create_dataset(test_df)
4 test_anomaly_dataset, _, _ = create_dataset(anomaly_df)

```

## LSTM Autoencoder



*Sample Autoencoder Architecture* [Image Source<sup>7</sup>](#)

The general Autoencoder architecture consists of two components. An *Encoder* that compresses the input and a *Decoder* that tries to reconstruct it.

<sup>7</sup><https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html>

We'll use the LSTM Autoencoder from this [GitHub repo](#)<sup>8</sup> with some small tweaks. Our model's job is to reconstruct Time Series data. Let's start with the *Encoder*:

```

1  class Encoder(nn.Module):
2
3      def __init__(self, seq_len, n_features, embedding_dim=64):
4          super(Encoder, self).__init__()
5
6          self.seq_len, self.n_features = seq_len, n_features
7          self.embedding_dim, self.hidden_dim = embedding_dim, 2 * embedding_dim
8
9          self.rnn1 = nn.LSTM(
10              input_size=n_features,
11              hidden_size=self.hidden_dim,
12              num_layers=1,
13              batch_first=True
14          )
15
16          self.rnn2 = nn.LSTM(
17              input_size=self.hidden_dim,
18              hidden_size=embedding_dim,
19              num_layers=1,
20              batch_first=True
21          )
22
23      def forward(self, x):
24          x = x.reshape((1, self.seq_len, self.n_features))
25
26          x, (_, _) = self.rnn1(x)
27          x, (hidden_n, _) = self.rnn2(x)
28
29          return hidden_n.reshape((self.n_features, self.embedding_dim))

```

The *Encoder* uses two LSTM layers to compress the Time Series data input.

Next, we'll decode the compressed representation using a *Decoder*:

---

<sup>8</sup><https://github.com/shobrook/sequitur>

```

1  class Decoder(nn.Module):
2
3      def __init__(self, seq_len, input_dim=64, n_features=1):
4          super(Decoder, self).__init__()
5
6          self.seq_len, self.input_dim = seq_len, input_dim
7          self.hidden_dim, self.n_features = 2 * input_dim, n_features
8
9          self.rnn1 = nn.LSTM(
10              input_size=input_dim,
11              hidden_size=input_dim,
12              num_layers=1,
13              batch_first=True
14      )
15
16      self.rnn2 = nn.LSTM(
17          input_size=input_dim,
18          hidden_size=self.hidden_dim,
19          num_layers=1,
20          batch_first=True
21      )
22
23      self.output_layer = nn.Linear(self.hidden_dim, n_features)
24
25  def forward(self, x):
26      x = x.repeat(self.seq_len, self.n_features)
27      x = x.reshape((self.n_features, self.seq_len, self.input_dim))
28
29      x, (hidden_n, cell_n) = self.rnn1(x)
30      x, (hidden_n, cell_n) = self.rnn2(x)
31      x = x.reshape((self.seq_len, self.hidden_dim))
32
33  return self.output_layer(x)

```

Our Decoder contains two LSTM layers and an output layer that gives the final reconstruction.

Time to wrap everything into an easy to use module:

```

1 class RecurrentAutoencoder(nn.Module):
2
3     def __init__(self, seq_len, n_features, embedding_dim=64):
4         super(RecurrentAutoencoder, self).__init__()
5
6         self.encoder = Encoder(seq_len, n_features, embedding_dim).to(device)
7         self.decoder = Decoder(seq_len, embedding_dim, n_features).to(device)
8
9     def forward(self, x):
10        x = self.encoder(x)
11        x = self.decoder(x)
12
13    return x

```

Our Autoencoder passes the input through the Encoder and Decoder. Let's create an instance of it:

```

1 model = RecurrentAutoencoder(seq_len, n_features, 128)
2 model = model.to(device)

```

## Training

Let's write a helper function for our training process:

```

1 def train_model(model, train_dataset, val_dataset, n_epochs):
2     optimizer = torch.optim.Adam(model.parameters(), lr=1e-3)
3     criterion = nn.L1Loss(reduction='sum').to(device)
4     history = dict(train=[], val[])
5
6     best_model_wts = copy.deepcopy(model.state_dict())
7     best_loss = 10000.0
8
9     for epoch in range(1, n_epochs + 1):
10        model = model.train()
11
12        train_losses = []
13        for seq_true in train_dataset:
14            optimizer.zero_grad()
15
16            seq_true = seq_true.to(device)
17            seq_pred = model(seq_true)

```

```

18
19     loss = criterion(seq_pred, seq_true)
20
21     loss.backward()
22     optimizer.step()
23
24     train_losses.append(loss.item())
25
26 val_losses = []
27 model = model.eval()
28 with torch.no_grad():
29     for seq_true in val_dataset:
30
31         seq_true = seq_true.to(device)
32         seq_pred = model(seq_true)
33
34         loss = criterion(seq_pred, seq_true)
35         val_losses.append(loss.item())
36
37 train_loss = np.mean(train_losses)
38 val_loss = np.mean(val_losses)
39
40 history['train'].append(train_loss)
41 history['val'].append(val_loss)
42
43 if val_loss < best_loss:
44     best_loss = val_loss
45     best_model_wts = copy.deepcopy(model.state_dict())
46
47 print(f'Epoch {epoch}: train loss {train_loss} val loss {val_loss}')
48
49 model.load_state_dict(best_model_wts)
50 return model.eval(), history

```

At each epoch, the training process feeds our model with all training examples and evaluates the performance on the validation set. Note that we're using a batch size of 1 (our model sees only 1 sequence at a time). We also record the training and validation set losses during the process.

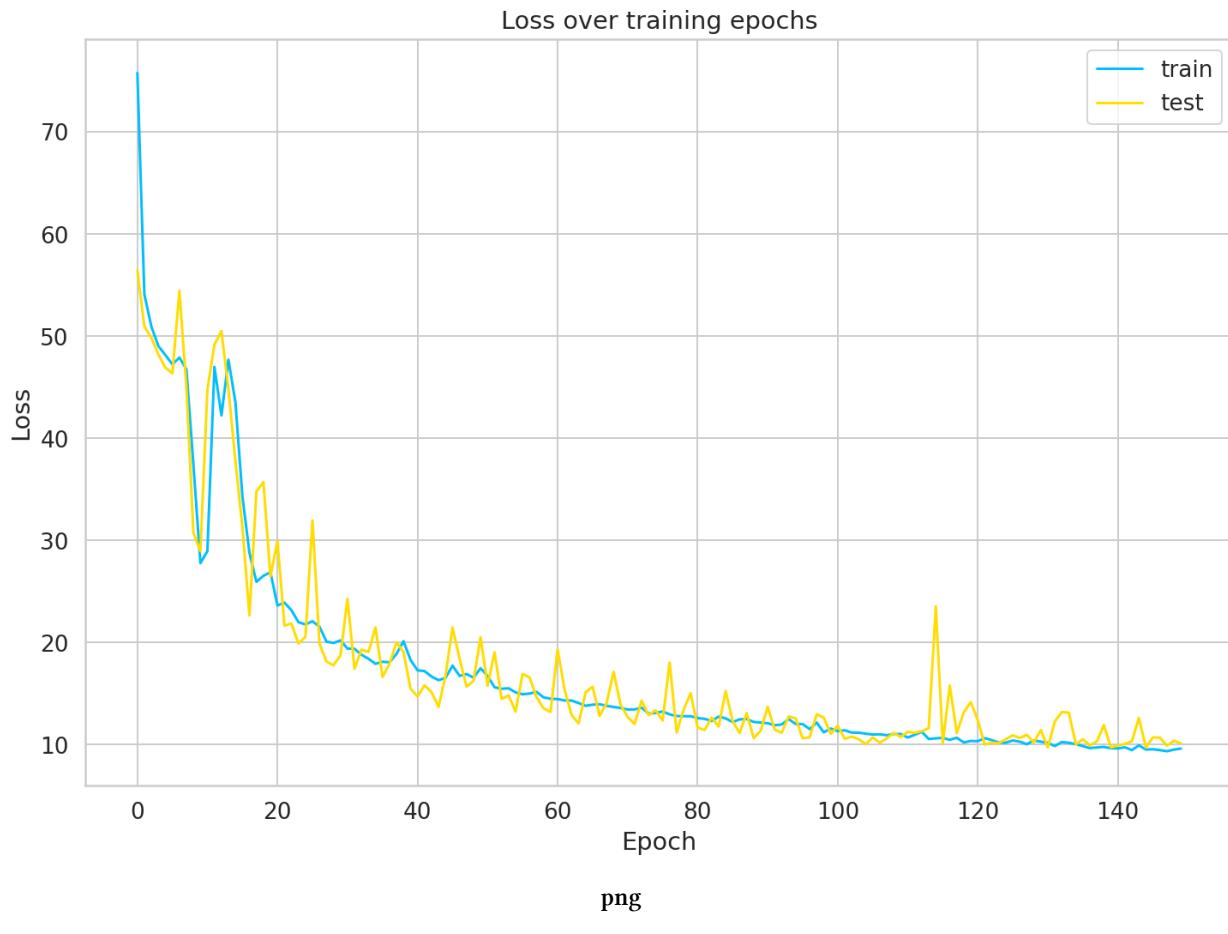
Note that we're minimizing the [L1Loss](#)<sup>9</sup>, which measures the MAE (mean absolute error). Why? The reconstructions seem to be better than with MSE (mean squared error).

We'll get the version of the model with the smallest validation error. Let's do some training:

---

<sup>9</sup><https://pytorch.org/docs/stable/nn.html#l1loss>

```
1 model, history = train_model(  
2     model,  
3     train_dataset,  
4     val_dataset,  
5     n_epochs=150  
6 )
```



Our model converged quite well. Seems like we might've needed a larger validation set to smoothen the results, but that'll do for now.

## Saving the model

Let's store the model for later use:

```

1 MODEL_PATH = 'model.pth'
2
3 torch.save(model, MODEL_PATH)

```

Uncomment the next lines, if you want to download and load the pre-trained model:

```

1 # !gdown --id 1jEYx5wGsb7Ix8cZAw315p5p0wHs3_I9A
2 # model = torch.load('model.pth')
3 # model = model.to(device)

```

## Choosing a threshold

With our model at hand, we can have a look at the reconstruction error on the training set. Let's start by writing a helper function to get predictions from our model:

```

1 def predict(model, dataset):
2     predictions, losses = [], []
3     criterion = nn.L1Loss(reduction='sum').to(device)
4     with torch.no_grad():
5         model = model.eval()
6         for seq_true in dataset:
7             seq_true = seq_true.to(device)
8             seq_pred = model(seq_true)
9
10            loss = criterion(seq_pred, seq_true)
11
12            predictions.append(seq_pred.cpu().numpy().flatten())
13            losses.append(loss.item())
14
15    return predictions, losses

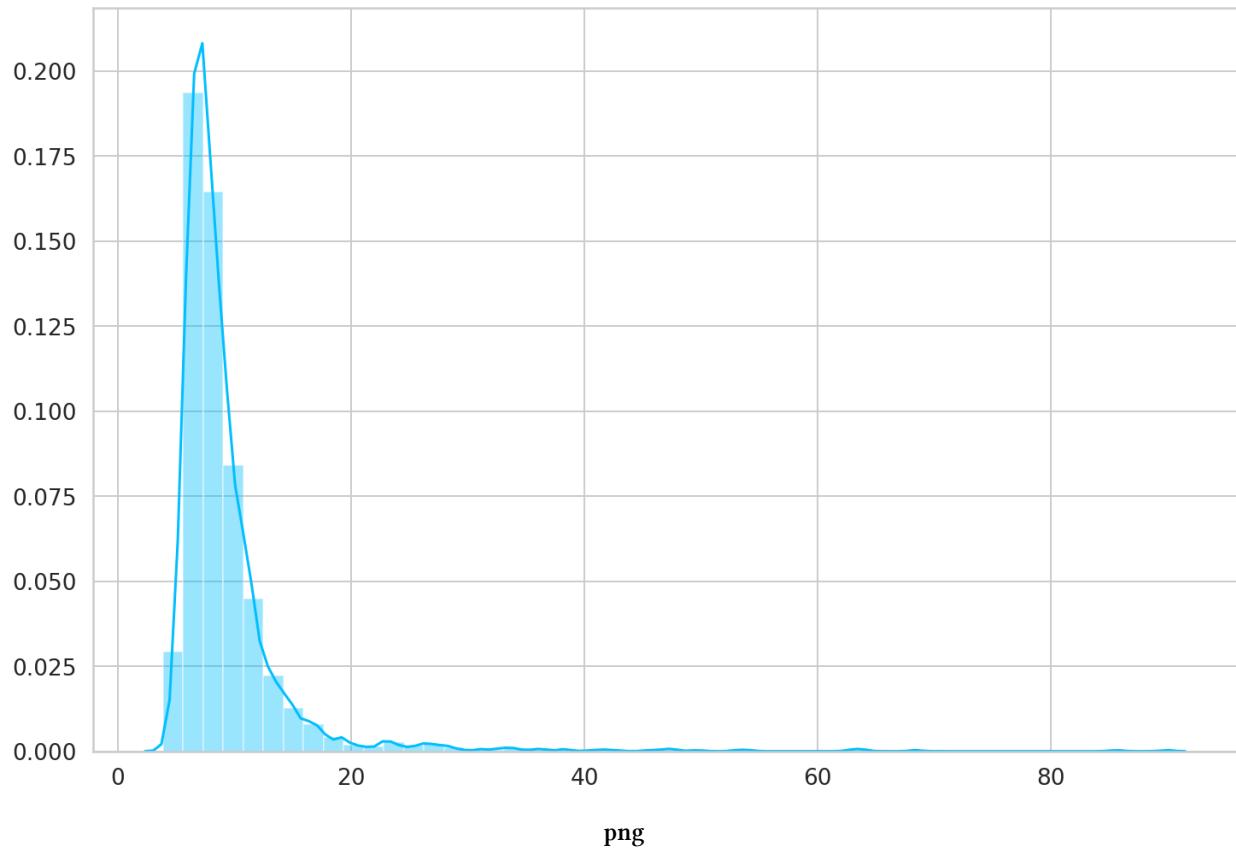
```

Our function goes through each example in the dataset and records the predictions and losses. Let's get the losses and have a look at them:

```

1 _, losses = predict(model, train_dataset)
2
3 sns.distplot(losses, bins=50, kde=True);

```



```
1 THRESHOLD = 26
```

## Evaluation

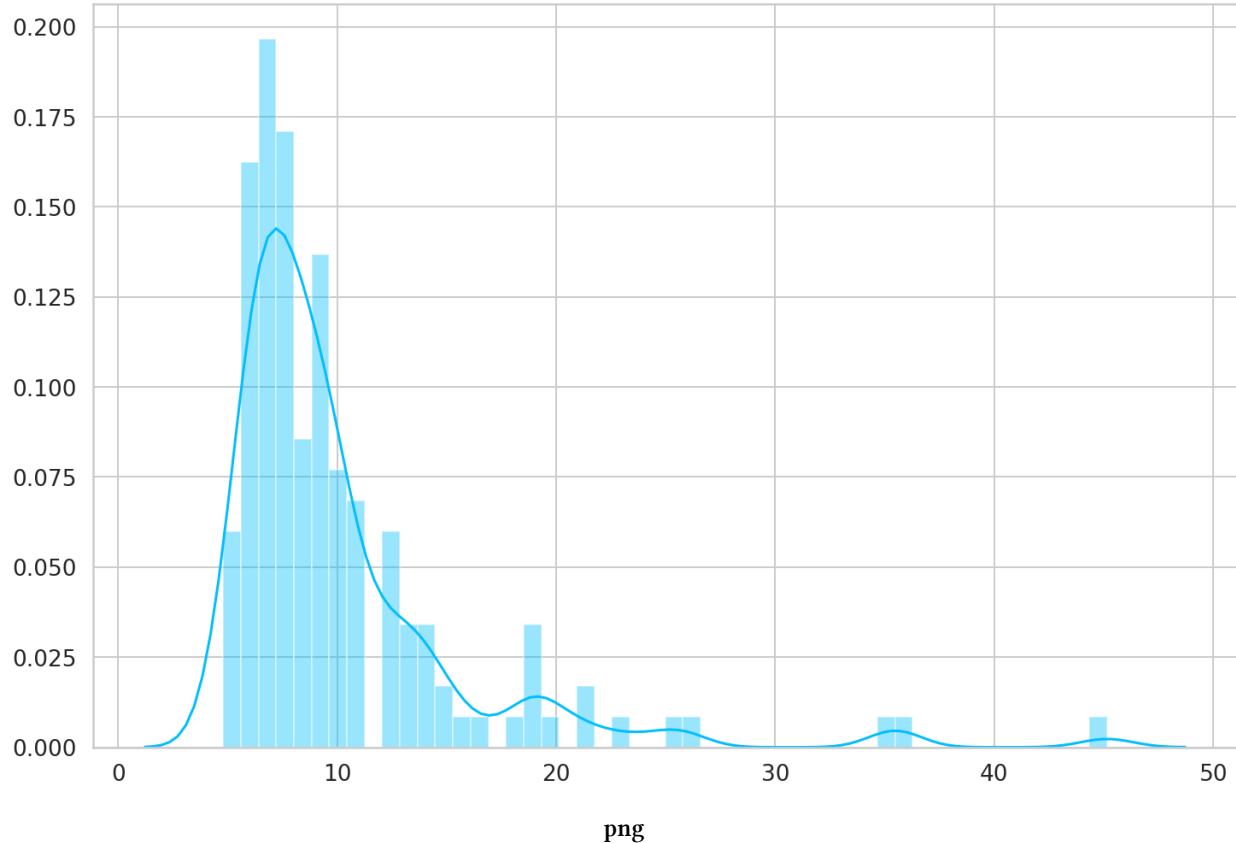
Using the threshold, we can turn the problem into a simple binary classification task:

- If the reconstruction loss for an example is below the threshold, we'll classify it as a *normal* heartbeat
- Alternatively, if the loss is higher than the threshold, we'll classify it as an anomaly

### Normal heartbeats

Let's check how well our model does on normal heartbeats. We'll use the normal heartbeats from the test set (our model haven't seen those):

```
1 predictions, pred_losses = predict(model, test_normal_dataset)
2 sns.distplot(pred_losses, bins=50, kde=True);
```



We'll count the correct predictions:

```
1 correct = sum(1 <= THRESHOLD for 1 in pred_losses)
2 print(f'Correct normal predictions: {correct}/{len(test_normal_dataset)}')

1 Correct normal predictions: 142/145
```

## Anomalies

We'll do the same with the anomaly examples, but their number is much higher. We'll get a subset that has the same size as the normal heartbeats:

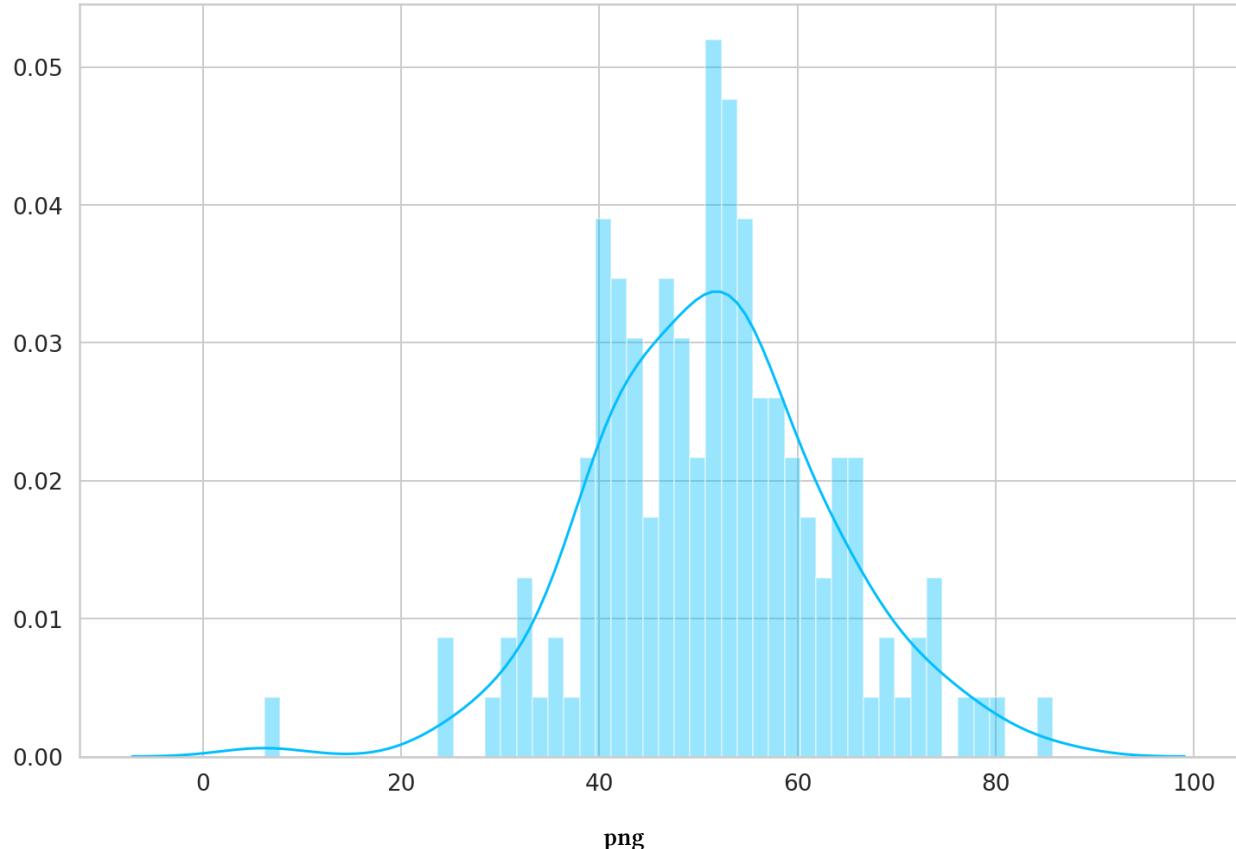
```
1 anomaly_dataset = test_anomaly_dataset[:len(test_normal_dataset)]
```

Now we can take the predictions of our model for the subset of anomalies:

```

1 predictions, pred_losses = predict(model, anomaly_dataset)
2 sns.distplot(pred_losses, bins=50, kde=True);

```



Finally, we can count the number of examples above the threshold (considered as anomalies):

```

1 correct = sum(1 > THRESHOLD for l in pred_losses)
2 print(f'Correct anomaly predictions: {correct}/{len(anomaly_dataset)}')

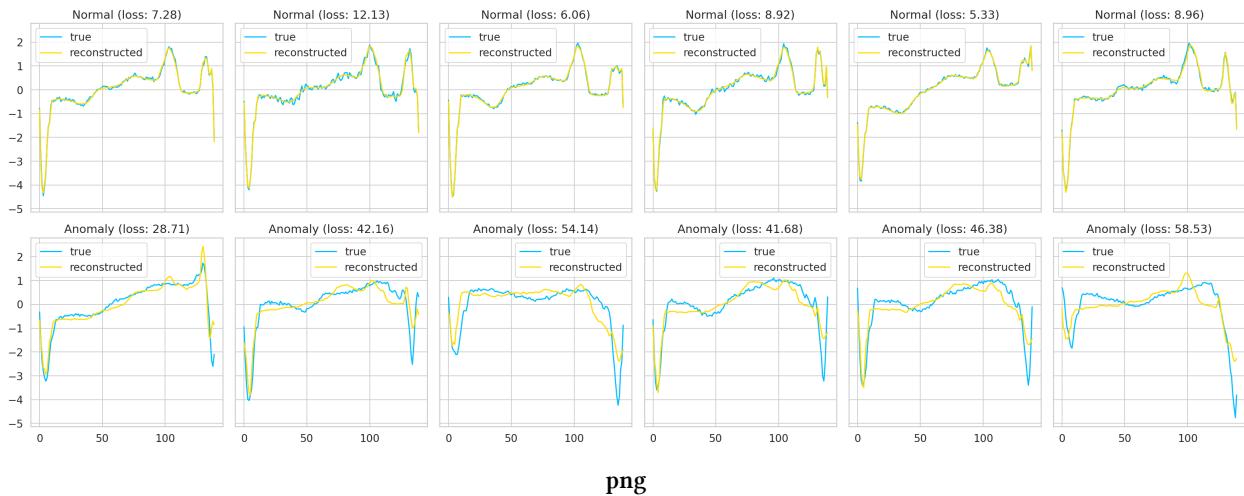
1 Correct anomaly predictions: 142/145

```

We have very good results. In the real world, you can tweak the threshold depending on what kind of errors you want to tolerate. In this case, you might want to have more false positives (normal heartbeats considered as anomalies) than false negatives (anomalies considered as normal).

## Looking at Examples

We can overlay the real and reconstructed Time Series values to see how close they are. We'll do it for some normal and anomaly cases:



## Summary

In this tutorial, you learned how to create an LSTM Autoencoder with PyTorch and use it to detect heartbeat anomalies in ECG data.

- Run the complete notebook in your browser (Google Colab)<sup>10</sup>
- Read the Getting Things Done with Pytorch book<sup>11</sup>

You learned how to:

- Prepare a dataset for Anomaly Detection from Time Series Data
- Build an LSTM Autoencoder with PyTorch
- Train and evaluate your model
- Choose a threshold for anomaly detection
- Classify unseen examples as normal or anomaly

While our Time Series data is univariate (we have only 1 feature), the code should work for multivariate datasets (multiple features) with little or no modification. Feel free to try it!

## References

- Sequitur - Recurrent Autoencoder (RAE)<sup>12</sup>
- Towards Never-Ending Learning from Time Series Streams<sup>13</sup>
- LSTM Autoencoder for Anomaly Detection<sup>14</sup>

<sup>10</sup>[https://colab.research.google.com/drive/1\\_J2MrBSvsJfOcVmYAN2-WSp36BtsFZCa](https://colab.research.google.com/drive/1_J2MrBSvsJfOcVmYAN2-WSp36BtsFZCa)

<sup>11</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>

<sup>12</sup><https://github.com/shobrook/sequitur>

<sup>13</sup><https://www.cs.ucr.edu/~eamonn/neverending.pdf>

<sup>14</sup><https://towardsdatascience.com/lstm-autoencoder-for-anomaly-detection-e1f4f2ee7ccf>

# 6. Face Detection on Custom Dataset with Detectron2

TL;DR Learn how to prepare a custom Face Detection dataset for Detectron2 and PyTorch. Fine-tune a pre-trained model to find face boundaries in images.

Face detection is the task of finding (boundaries of) faces in images. This is useful for

- security systems (the first step in recognizing a person)
- autofocus and smile detection for making great photos
- detecting age, race, and emotional state for marking (yep, we already live in that world)

Historically, this was a really tough problem to solve. Tons of manual feature engineering, novel algorithms and methods were developed to improve the state-of-the-art.

These days, face detection models are included in almost every computer vision package/framework. Some of the best-performing ones use Deep Learning methods. OpenCV, for example, provides a variety of tools like the [Cascade Classifier<sup>1</sup>](#).

- Run the complete notebook in your browser (Google Colab)<sup>2</sup>
- Read the [Getting Things Done with Pytorch](#) book<sup>3</sup>

In this guide, you'll learn how to:

- prepare a custom dataset for face detection with Detectron2
- use (close to) state-of-the-art models for object detection to find faces in images
- You can extend this work for face recognition.

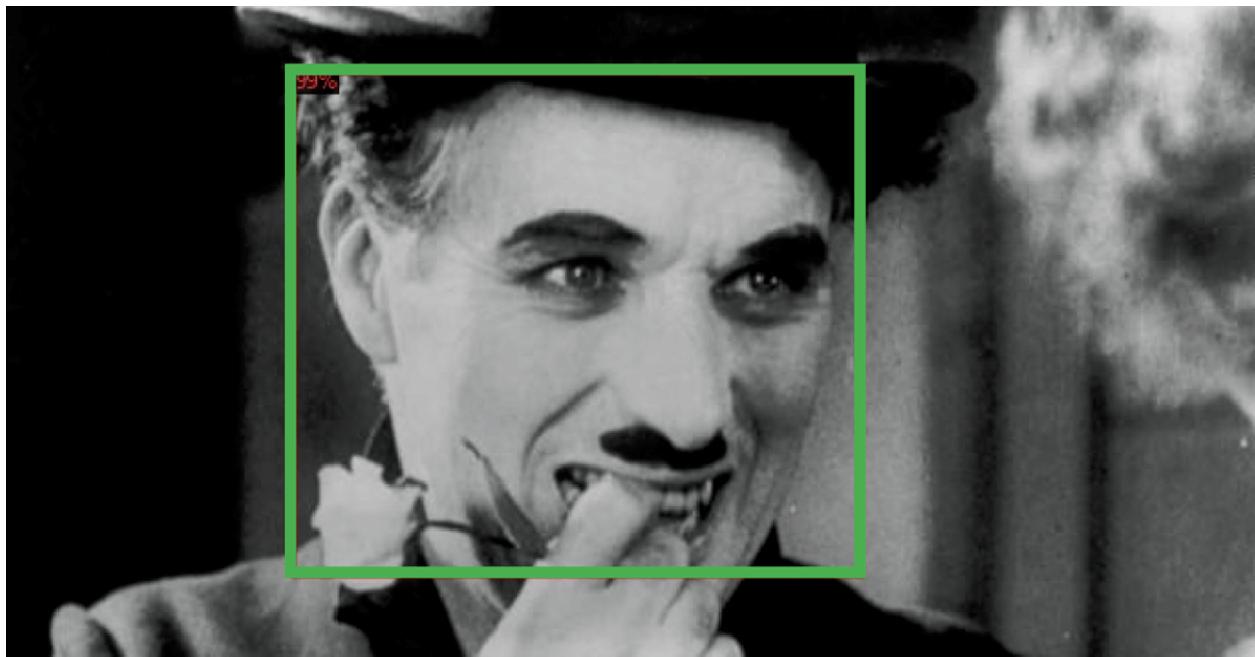
Here's an example of what you'll get at the end of this guide:

---

<sup>1</sup>[https://docs.opencv.org/4.2.0/db/d28/tutorial\\_cascade\\_classifier.html](https://docs.opencv.org/4.2.0/db/d28/tutorial_cascade_classifier.html)

<sup>2</sup><https://colab.research.google.com/drive/1Jk4-qX9zdYGsBrTnh2vF52CV9ucuqpjk>

<sup>3</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>



png

## Detectron 2



# Detectron2

png

Detectron2<sup>4</sup> is a framework for building state-of-the-art object detection and image segmentation models. It is developed by the Facebook Research team. Detectron2 is a complete rewrite of the [first version](#)<sup>5</sup>.

Under the hood, Detectron2 uses PyTorch (compatible with the latest version(s)) and allows for [blazing fast training](#)<sup>6</sup>. You can learn more at [introductory blog post](#)<sup>7</sup> by Facebook Research.

The real power of Detectron2 lies in the HUGE amount of pre-trained models available at the [Model Zoo](#)<sup>8</sup>. But what good that would it be if you can't fine-tune those on your own datasets? Fortunately,

<sup>4</sup><https://github.com/facebookresearch/detectron2>

<sup>5</sup><https://github.com/facebookresearch/Detectron>

<sup>6</sup><https://detectron2.readthedocs.io/notes/benchmarks.html>

<sup>7</sup><https://ai.facebook.com/blog/-/detectron2-a-pytorch-based-modular-object-detection-library-/>

<sup>8</sup>[https://github.com/facebookresearch/detectron2/blob/master/MODEL\\_ZOO.md](https://github.com/facebookresearch/detectron2/blob/master/MODEL_ZOO.md)

that's super easy! We'll see how it is done in this guide.

## Installing Detectron2

At the time of this writing, Detectron2 is still in an alpha stage. While there is an official release, we'll clone and compile from the master branch. This should equal version 0.1.

Let's start by installing some requirements:

```
1 !pip install -q cython pyyaml==5.1
2 !pip install -q -U 'git+https://github.com/cocodataset/cocoapi.git#subdirectory=Pyth\
3 onAPI'
```

And download, compile, and install the Detectron2 package:

```
1 !git clone https://github.com/facebookresearch/detectron2 detectron2_repo
2 !pip install -q -e detectron2_repo
```

At this point, you'll need to restart the notebook runtime to continue!

```
1 !pip install -q -U watermark

1 %reload_ext watermark
2 %watermark -v -p numpy,pandas,pycocotools,torch,torchvision,detectron2

1 CPython 3.6.9
2 IPython 5.5.0
3
4 numpy 1.17.5
5 pandas 0.25.3
6 pycocotools 2.0
7 torch 1.4.0
8 torchvision 0.5.0
9 detectron2 0.1
```

```
1 import torch, torchvision
2 import detectron2
3 from detectron2.utils.logger import setup_logger
4 setup_logger()
5
6 import glob
7
8 import os
9 import ntpath
10 import numpy as np
11 import cv2
12 import random
13 import itertools
14 import pandas as pd
15 from tqdm import tqdm
16 import urllib
17 import json
18 import PIL.Image as Image
19
20 from detectron2 import model_zoo
21 from detectron2.engine import DefaultPredictor, DefaultTrainer
22 from detectron2.config import get_cfg
23 from detectron2.utils.visualizer import Visualizer, ColorMode
24 from detectron2.data import DatasetCatalog, MetadataCatalog, build_detection_test_lo\
ader
25
26 from detectron2.evaluation import COCOEvaluator, inference_on_dataset
27 from detectron2.structures import BoxMode
28
29 import seaborn as sns
30 from pylab import rcParams
31 import matplotlib.pyplot as plt
32 from matplotlib import rc
33
34 %matplotlib inline
35 %config InlineBackend.figure_format='retina'
36
37 sns.set(style='whitegrid', palette='muted', font_scale=1.2)
38
39 HAPPY_COLORS_PALETTE = ["#01BEFE", "#FFDD00", "#FF7D00", "#FF006D", "#ADFF02", "#8F0\\
OFF"]
40
41 sns.set_palette(sns.color_palette(HAPPY_COLORS_PALETTE))
42
43
```

```

44 rcParams['figure.figsize'] = 12, 8
45
46 RANDOM_SEED = 42
47 np.random.seed(RANDOM_SEED)
48 torch.manual_seed(RANDOM_SEED)

```

## Face Detection Data

Our dataset is provided by [Dataturks<sup>9</sup>](#), and it is hosted on [Kaggle<sup>10</sup>](#). Here's an excerpt from the description:

Faces in images marked with bounding boxes. Have around 500 images with around 1100 faces manually tagged via bounding box.

I've downloaded the JSON file containing the annotations and uploaded it to Google Drive. Let's get it:

```
1 !gdown --id 1K79wJgmPTWamqb040p2GxW0SW9oxw8KS
```

Let's load the file into a Pandas dataframe:

```
1 faces_df = pd.read_json('face_detection.json', lines=True)
```

Each line contains a single face annotation. Note that multiple lines might point to a single image (e.g. multiple faces per image).

## Data Preprocessing

The dataset contains only image URLs and annotations. We'll have to download the images. We'll also normalize the annotations, so it's easier to use them with Detectron2 later on:

---

<sup>9</sup><https://dataturks.com/>

<sup>10</sup><https://www.kaggle.com/dataturks/face-detection-in-images>

```

1 os.makedirs("faces", exist_ok=True)
2
3 dataset = []
4
5 for index, row in tqdm(faces_df.iterrows(), total=faces_df.shape[0]):
6     img = urllib.request.urlopen(row["content"])
7     img = Image.open(img)
8     img = img.convert('RGB')
9
10    image_name = f'face_{index}.jpeg'
11
12    img.save(f'faces/{image_name}', "JPEG")
13
14    annotations = row['annotation']
15    for an in annotations:
16
17        data = {}
18
19        width = an['imageWidth']
20        height = an['imageHeight']
21        points = an['points']
22
23        data['file_name'] = image_name
24        data['width'] = width
25        data['height'] = height
26
27        data["x_min"] = int(round(points[0]["x"] * width))
28        data["y_min"] = int(round(points[0]["y"] * height))
29        data["x_max"] = int(round(points[1]["x"] * width))
30        data["y_max"] = int(round(points[1]["y"] * height))
31
32        data['class_name'] = 'face'
33
34    dataset.append(data)

```

Let's put the data into a dataframe so we can have a better look:

```

1 df = pd.DataFrame(dataset)

1 print(df.file_name.unique().shape[0], df.shape[0])

```

```
1 409 1132
```

We have a total of 409 images (a lot less than the promised 500) and 1132 annotations. Let's save them to the disk (so you might reuse them):

```
1 df.to_csv('annotations.csv', header=True, index=None)
```

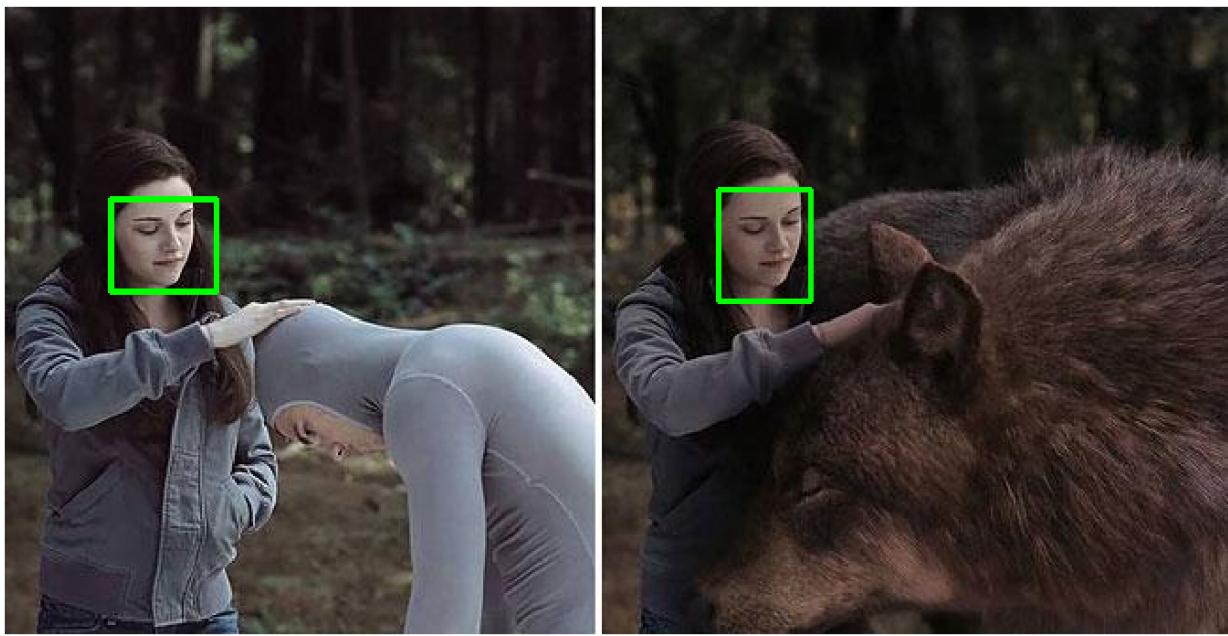
## Data Exploration

Let's see some sample annotated data. We'll use OpenCV to load an image, add the bounding boxes, and resize it. We'll define a helper function to do it all:

```
1 def annotate_image(annotations, resize=True):
2     file_name = annotations.file_name.to_numpy()[0]
3     img = cv2.cvtColor(cv2.imread(f'faces/{file_name}'), cv2.COLOR_BGR2RGB)
4
5     for i, a in annotations.iterrows():
6         cv2.rectangle(img, (a.x_min, a.y_min), (a.x_max, a.y_max), (0, 255, 0), 2)
7
8     if not resize:
9         return img
10
11    return cv2.resize(img, (384, 384), interpolation = cv2.INTER_AREA)
```

Let's start by showing some annotated images:

```
1 img_df = df[df.file_name == df.file_name.unique()[0]]
2 img = annotate_image(img_df, resize=False)
3
4 plt.imshow(img)
5 plt.axis('off');
```



png

```
1 img_df = df[df.file_name == df.file_name.unique()[1]]  
2 img = annotate_image(img_df, resize=False)  
3  
4 plt.imshow(img)  
5 plt.axis('off');
```



png

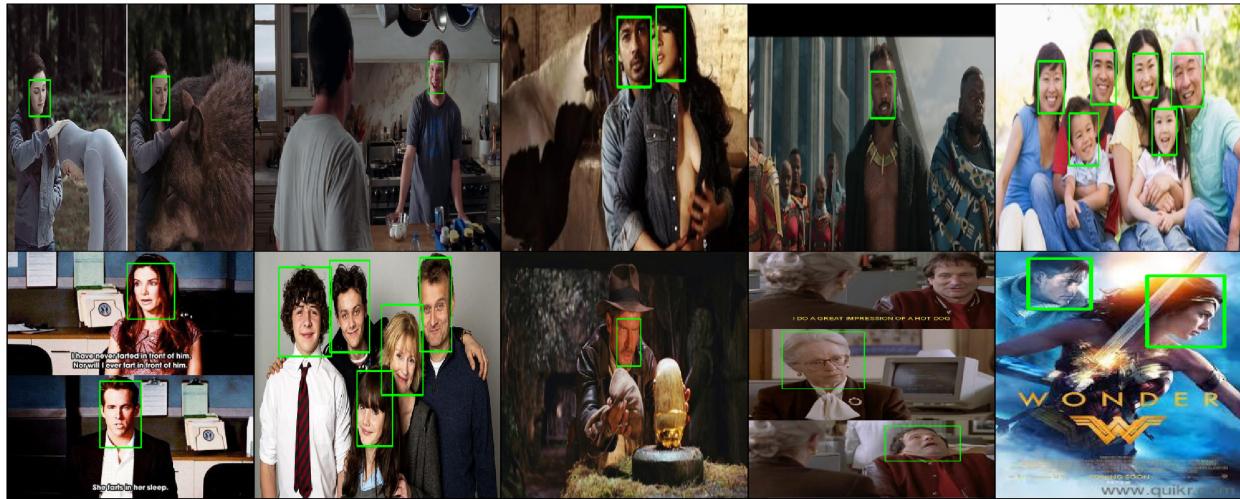
Those are good ones, the annotations are clearly visible. We can use torchvision to create a grid of images. Note that the images are in various sizes, so we'll resize them:

```
1 sample_images = [annotate_image(df[df.file_name == f]) for f in df.file_name.unique()\n2 )[:10]]\n3 sample_images = torch.as_tensor(sample_images)\n\n1 sample_images.shape\n\n1 torch.Size([10, 384, 384, 3])\n\n1 sample_images = sample_images.permute(0, 3, 1, 2)\n\n1 sample_images.shape\n\n1 torch.Size([10, 3, 384, 384])
```

```

1 plt.figure(figsize=(24, 12))
2 grid_img = torchvision.utils.make_grid(sample_images, nrow=5)
3
4 plt.imshow(grid_img.permute(1, 2, 0))
5 plt.axis('off');

```



png

You can clearly see that some annotations are missing (column 4). That's real life data for you, sometimes you have to deal with it in some way.

## Face Detection with Detectron 2

It is time to go through the steps of fine-tuning a model using a custom dataset. But first, let's save 5% of the data for testing:

```

1 df = pd.read_csv('annotations.csv')
2
3 IMAGES_PATH = f'faces'
4
5 unique_files = df.file_name.unique()
6
7 train_files = set(np.random.choice(unique_files, int(len(unique_files) * 0.95), replace=False))
8 train_df = df[df.file_name.isin(train_files)]
9 test_df = df[~df.file_name.isin(train_files)]

```

The classical `train_test_split` won't work here, cause we want a split amongst the file names.

The next parts are written in a bit more generic way. Obviously, we have a single class - face. But adding more should be as simple as adding more annotations to the dataframe:

```
1 classes = df.class_name.unique().tolist()
```

Next, we'll write a function that converts our dataset into a format that is used by Detectron2:

```
1 def create_dataset_dicts(df, classes):
2     dataset_dicts = []
3     for image_id, img_name in enumerate(df.file_name.unique()):
4
5         record = {}
6
7         image_df = df[df.file_name == img_name]
8
9         file_path = f'{IMAGES_PATH}/{img_name}'
10        record["file_name"] = file_path
11        record["image_id"] = image_id
12        record["height"] = int(image_df.iloc[0].height)
13        record["width"] = int(image_df.iloc[0].width)
14
15        objs = []
16        for _, row in image_df.iterrows():
17
18            xmin = int(row.x_min)
19            ymin = int(row.y_min)
20            xmax = int(row.x_max)
21            ymax = int(row.y_max)
22
23            poly = [
24                (xmin, ymin), (xmax, ymin),
25                (xmax, ymax), (xmin, ymax)
26            ]
27            poly = list(itertools.chain.from_iterable(poly))
28
29            obj = {
30                "bbox": [xmin, ymin, xmax, ymax],
31                "bbox_mode": BoxMode.XYXY_ABS,
32                "segmentation": [poly],
33                "category_id": classes.index(row.class_name),
34                "iscrowd": 0
35            }
36            objs.append(obj)
```

```

37
38     record["annotations"] = objs
39     dataset_dicts.append(record)
40
41     return dataset_dicts

```

We convert every annotation row to a single record with a list of annotations. You might also notice that we're building a polygon that is of the exact same shape as the bounding box. This is required for the image segmentation models in Detectron2.

You'll have to register your dataset into the dataset and metadata catalogues:

```

1 for d in ["train", "val"]:
2     DatasetCatalog.register("faces_" + d, lambda d=d: create_dataset_dicts(train_df if \
3         d == "train" else test_df, classes))
4     MetadataCatalog.get("faces_" + d).set(thing_classes=classes)
5
6 statement_metadata = MetadataCatalog.get("faces_train")

```

Unfortunately, evaluator for the test set is not included by default. We can easily fix that by writing our own trainer:

```

1 class CocoTrainer(DefaultTrainer):
2
3     @classmethod
4     def build_evaluator(cls, cfg, dataset_name, output_folder=None):
5
6         if output_folder is None:
7             os.makedirs("coco_eval", exist_ok=True)
8             output_folder = "coco_eval"
9
10    return COCOEvaluator(dataset_name, cfg, False, output_folder)

```

The evaluation results will be stored in the `coco_eval` folder if no folder is provided.

Fine-tuning a Detectron2 model is nothing like writing PyTorch code. We'll load a configuration file, change a few values, and start the training process. But hey, it really helps if you know what you're doing ☺

For this tutorial, we'll use the Mask R-CNN X101-FPN model. It is pre-trained on the [COCO dataset<sup>11</sup>](#) and achieves very good performance. The downside is that it is slow to train.

Let's load the config file and the pre-trained model weights:

---

<sup>11</sup><http://cocodataset.org/#home>

```

1  cfg = get_cfg()
2
3  cfg.merge_from_file(
4      model_zoo.get_config_file(
5          "COCO-InstanceSegmentation/mask_rcnn_X_101_32x8d_FPN_3x.yaml"
6      )
7  )
8
9  cfg.MODEL.WEIGHTS = model_zoo.get_checkpoint_url(
10     "COCO-InstanceSegmentation/mask_rcnn_X_101_32x8d_FPN_3x.yaml"
11 )

```

Specify the datasets (we registered those) we'll use for training and evaluation:

```

1  cfg.DATASETS.TRAIN = ("faces_train",)
2  cfg.DATASETS.TEST = ("faces_val",)
3  cfg.DATALOADER.NUM_WORKERS = 4

```

And for the optimizer, we'll do a bit of magic to converge to something nice:

```

1  cfg.SOLVER.IMS_PER_BATCH = 4
2  cfg.SOLVER.BASE_LR = 0.001
3  cfg.SOLVER.WARMUP_ITERS = 1000
4  cfg.SOLVER.MAX_ITER = 1500
5  cfg.SOLVER.STEPS = (1000, 1500)
6  cfg.SOLVER.GAMMA = 0.05

```

Except for the standard stuff (batch size, max number of iterations, and learning rate) we have a couple of interesting params:

- WARMUP\_ITERS - the learning rate starts from 0 and goes to the preset one for this number of iterations
- STEPS - the checkpoints (number of iterations) at which the learning rate will be reduced by GAMMA

Finally, we'll specify the number of classes and the period at which we'll evaluate on the test set:

```

1 cfg.MODEL.ROI_HEADS.BATCH_SIZE_PER_IMAGE = 64
2 cfg.MODEL.ROI_HEADS.NUM_CLASSES = len(classes)
3
4 cfg.TEST.EVAL_PERIOD = 500

```

Time to train, using our custom trainer:

```

1 os.makedirs(cfg.OUTPUT_DIR, exist_ok=True)
2
3 trainer = CocoTrainer(cfg)
4 trainer.resume_or_load(resume=False)
5 trainer.train()

```

## Evaluating Object Detection Models

Evaluating object detection models is a bit different when compared to evaluating standard classification or regression models.

The main metric you need to know about is IoU (intersection over union). It measures the overlap between two boundaries - the predicted and ground truth one. It can get values between 0 and 1.

$$\text{IoU} = \frac{\text{area of overlap}}{\text{area of union}}$$

Using IoU, one can define a threshold (e.g. >0.5) to classify whether a prediction is a true positive (TP) or a false positive (FP).

Now you can calculate average precision (AP) by taking the area under the precision-recall curve.

Now AP@X (e.g. AP50) is just AP at some IoU threshold. This should give you a working understanding of how object detection models are evaluated.

I suggest you read the [mAP \(mean Average Precision\) for Object Detection<sup>12</sup>](#) tutorial by Jonathan Hui if you want to learn more on the topic.

I've prepared a pre-trained model for you, so you don't have to wait for the training to complete. Let's download it:

```

1 !gdown --id 18Ev2bpdKsBaDufhVKf0cT6RmM3FjW3nL
2 !mv face_detector.pth output/model_final.pth

```

We can start making predictions by loading the model and setting a minimum threshold of 85% certainty at which we'll consider the predictions as correct:

---

<sup>12</sup>[https://medium.com/@jonathan\\_hui/map-mean-average-precision-for-object-detection-45c121a31173](https://medium.com/@jonathan_hui/map-mean-average-precision-for-object-detection-45c121a31173)

```

1 cfg.MODEL.WEIGHTS = os.path.join(cfg.OUTPUT_DIR, "model_final.pth")
2 cfg.MODEL.ROI_HEADS.SCORE_THRESH_TEST = 0.85
3 predictor = DefaultPredictor(cfg)

```

Let's run the evaluator with the trained model:

```

1 evaluator = COCOEvaluator("faces_val", cfg, False, output_dir="./output/")
2 val_loader = build_detection_test_loader(cfg, "faces_val")
3 inference_on_dataset(trainer.model, val_loader, evaluator)

```

## Finding Faces in Images

Next, let's create a folder and save all images with predicted annotations in the test set:

```

1 os.makedirs("annotated_results", exist_ok=True)
2
3 test_image_paths = test_df.file_name.unique()

1 for clothing_image in test_image_paths:
2     file_path = f'{IMAGES_PATH}/{clothing_image}'
3     im = cv2.imread(file_path)
4     outputs = predictor(im)
5     v = Visualizer(
6         im[:, :, ::-1],
7         metadata=statement_metadata,
8         scale=1.,
9         instance_mode=ColorMode.IMAGE
10    )
11    instances = outputs["instances"].to("cpu")
12    instances.remove('pred_masks')
13    v = v.draw_instance_predictions(instances)
14    result = v.get_image()[:, :, ::-1]
15    file_name = ntpath.basename(clothing_image)
16    write_res = cv2.imwrite(f'annotated_results/{file_name}', result)

```

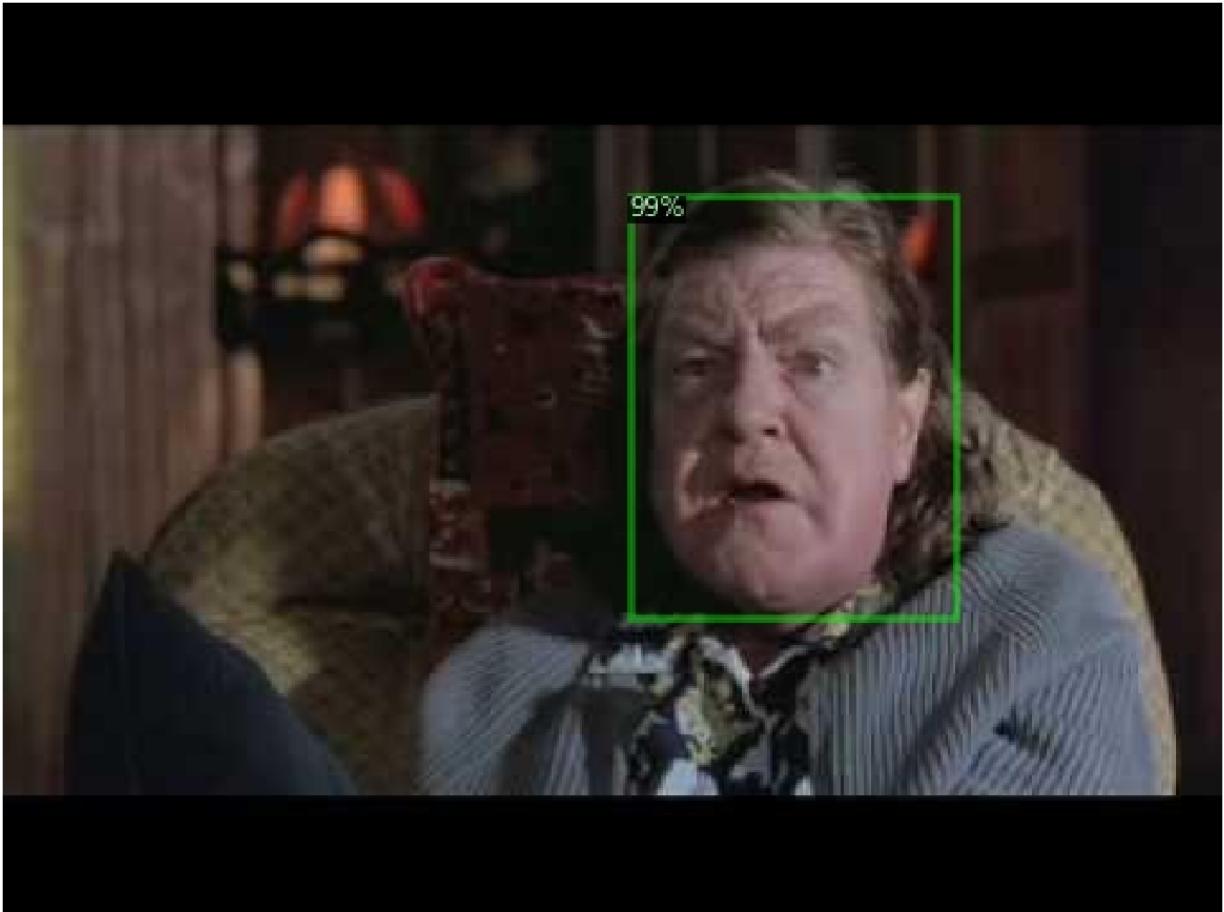
Let's have a look:

```

1 annotated_images = [f'annotated_results/{f}' for f in test_df.file_name.unique()]

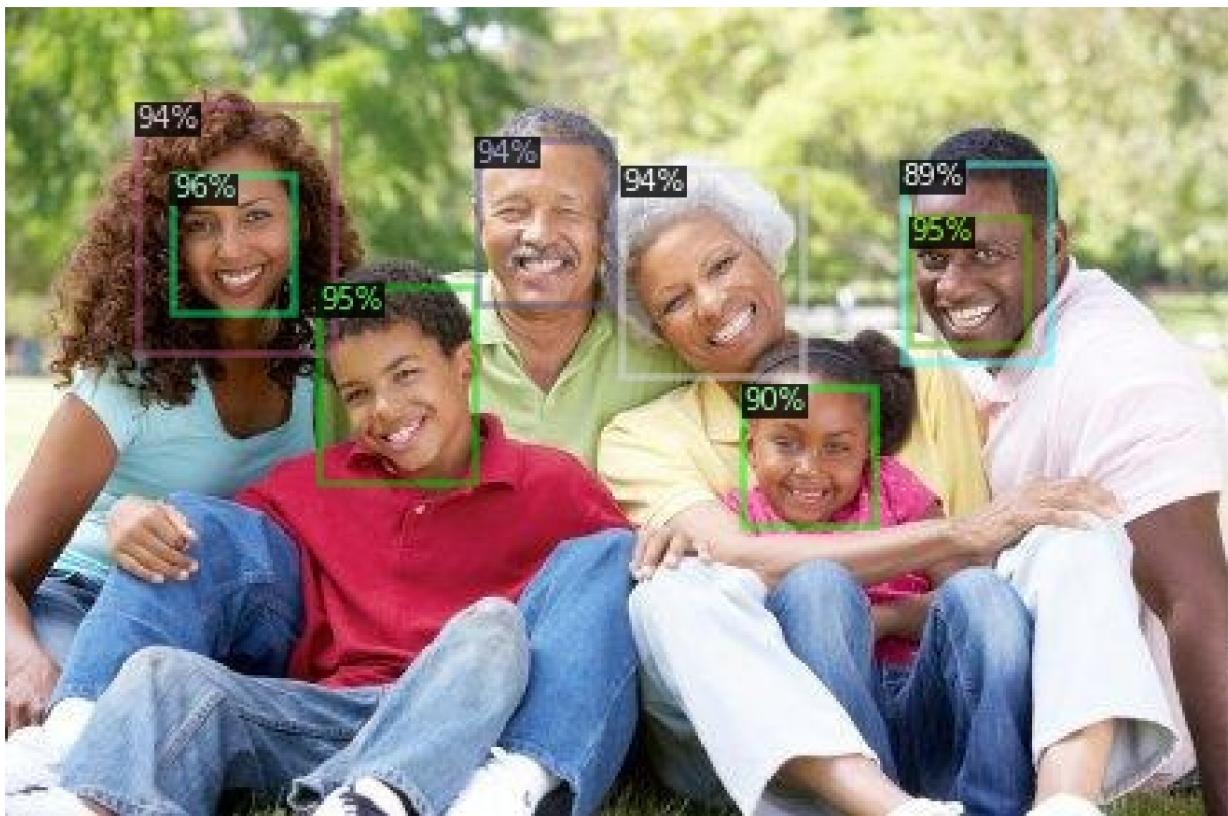
```

```
1 img = cv2.cvtColor(cv2.imread(annotated_images[0]), cv2.COLOR_BGR2RGB)
2
3 plt.imshow(img)
4 plt.axis('off');
```



png

```
1 img = cv2.cvtColor(cv2.imread(annotated_images[1]), cv2.COLOR_BGR2RGB)
2
3 plt.imshow(img)
4 plt.axis('off');
```



png

```
1 img = cv2.cvtColor(cv2.imread(annotated_images[3]), cv2.COLOR_BGR2RGB)
2
3 plt.imshow(img)
4 plt.axis('off');
```



png

```
1 img = cv2.cvtColor(cv2.imread(annotated_images[4]), cv2.COLOR_BGR2RGB)
2
3 plt.imshow(img)
4 plt.axis('off');
```



Not bad. Not bad at all. I suggest you explore more images on your own, too!

Note that some faces have multiple bounding boxes (on the second image) with different degrees of certainty. Maybe training the model longer will help? How about adding more or augmenting the existing data?

## Conclusion

Congratulations! You now know the basics of Detectron2 for object detection! You might be surprised by the results, given the small dataset we have. That's the power of large pre-trained models for you ☺

- Run the complete notebook in your browser (Google Colab)<sup>13</sup>
- Read the **Getting Things Done with Pytorch** book<sup>14</sup>

You learned how to:

- prepare a custom dataset for face detection with Detectron2
- use (close to) state-of-the-art models for object detection to find faces in images
- You can extend this work for face recognition.

<sup>13</sup><https://colab.research.google.com/drive/1Jk4-qX9zdYGsBrTnh2vF52CV9ucuqpjk>

<sup>14</sup><https://github.com/curiousily/Getting-Things-Done-with-Pytorch>

## References

- Face Detection in Images<sup>15</sup>
- Detectron2 on GitHub<sup>16</sup>
- mAP (mean Average Precision) for Object Detection<sup>17</sup>

<sup>15</sup><https://www.kaggle.com/dataturks/face-detection-in-images>

<sup>16</sup><https://github.com/facebookresearch/detectron2/>

<sup>17</sup>[https://medium.com/@jonathan\\_hui/map-mean-average-precision-for-object-detection-45c121a31173](https://medium.com/@jonathan_hui/map-mean-average-precision-for-object-detection-45c121a31173)

# 7. Create Dataset for Sentiment Analysis by Scraping Google Play App Reviews

TL;DR In this tutorial, you'll learn how to create a dataset for Sentiment Analysis by scraping user reviews for Android apps. You'll convert the app and review information into Data Frames and save that to CSV files.

- Run the notebook in your browser (Google Colab)<sup>1</sup>
- Read the Getting Things Done with Pytorch book<sup>2</sup>

You'll learn how to:

- Set a goal and inclusion criteria for your dataset
- Get real-world user reviews by scraping Google Play
- Use Pandas to convert and save the dataset into CSV files

## Setup

Let's install the required packages and setup the imports:

```
1 %watermark -v -p pandas,matplotlib,seaborn,google_play_scraper  
  
1 CPython 3.6.9  
2 IPython 5.5.0  
3  
4 pandas 1.0.3  
5 matplotlib 3.2.1  
6 seaborn 0.10.0  
7 google_play_scraper 0.0.2.3
```

---

<sup>1</sup><https://colab.research.google.com/drive/1GDJIpz7BXw55jl9wTOMQDool9m8DIOyp>

<sup>2</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>

```

1 import json
2 import pandas as pd
3 from tqdm import tqdm
4
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7
8 from pygments import highlight
9 from pygments.lexers import JsonLexer
10 from pygments.formatters import TerminalFormatter
11
12 from google_play_scraper import Sort, reviews, app
13
14 %matplotlib inline
15 %config InlineBackend.figure_format='retina'
16
17 sns.set(style='whitegrid', palette='muted', font_scale=1.2)

```

## The Goal of the Dataset

You want to get feedback for your app. Both negative and positive are good. But the negative one can reveal critical features that are missing or downtime of your service (when it is much more frequent).

Lucky for us, Google Play has plenty of apps, reviews, and scores. We can scrape app info and reviews using the [google-play-scraper<sup>3</sup>](#) package.

You can choose plenty of apps to analyze. But different app categories contain different audiences, domain-specific quirks, and more. We'll start simple.

We want apps that have been around some time, so opinion is collected organically. We want to mitigate advertising strategies as much as possible. Apps are constantly being updated, so the time of the review is an important factor.

Ideally, you would want to collect every possible review and work with that. However, in the real world data is often limited (too large, inaccessible, etc). So, we'll do the best we can.

Let's choose some apps that fit the criteria from the *Productivity* category. We'll use [AppAnnie<sup>4</sup>](#) to select some of the top US apps:

---

<sup>3</sup><https://github.com/JoMingyu/google-play-scraper>

<sup>4</sup>[https://www.appannie.com/apps/google-play/top-chart/?country=US&category=29&device=&date=2020-04-05&feed=All&rank\\_=sorting\\_type=rank&page\\_number=1&page\\_size=100&table\\_selections=](https://www.appannie.com/apps/google-play/top-chart/?country=US&category=29&device=&date=2020-04-05&feed=All&rank_=sorting_type=rank&page_number=1&page_size=100&table_selections=)

```

1 app_packages = [
2     'com.anydo',
3     'com.todoist',
4     'com.ticktick.task',
5     'com.habitrpg.android.habitica',
6     'cc.forestapp',
7     'com.oristats.habitbull',
8     'com.levor.liferpgtasks',
9     'com.habitnow',
10    'com.microsoft.todos',
11    'prox.lab.calclock',
12    'com.gmail.jmartindev.timetune',
13    'com.artfulagenda.app',
14    'com.tasks.android',
15    'com.appgenix.bizcal',
16    'com.appxy.planner'
17 ]

```

## Scraping App Information

Let's scrape the info for each app:

```

1 app_infos = []
2
3 for ap in tqdm(app_packages):
4     info = app(ap, lang='en', country='us')
5     del info['comments']
6     app_infos.append(info)
7
8
9 100%|████████████████| 15/15 [00:02<00:00,  6.34it/s]

```

We got the info for all 15 apps. Let's write a helper function that prints JSON objects a bit better:

```
1 def print_json(json_object):
2     json_str = json.dumps(
3         json_object,
4         indent=2,
5         sort_keys=True,
6         default=str
7     )
8     print(highlight(json_str, JsonLexer(), TerminalFormatter()))
```

Here is a sample app information from the list:

```
1 print_json(app_infos[0])
```

```
1 {  
2     "adSupported": null,  
3     "androidVersion": "Varies",  
4     "androidVersionText": "Varies with device",  
5     "appId": "com.anydo",  
6     "containsAds": null,  
7     "contentRating": "Everyone",  
8     "contentRatingDescription": null,  
9     "currency": "USD",  
10    "description": "<b>\ud83c\udfc6 Editor's Choice by Google</b>\r\n\r\nAny.do is a T\\  
11 o Do List, Calendar, Planner...",  
12    "descriptionHTML": "<b>\ud83c\udfc6 Editor's Choice by Google</b><br><br>Any.d\\  
13 o is a To Do List, Calendar, Planner...",  
14    "developer": "Any.do Calendar & To-Do List",  
15    "developerAddress": "Any.do Inc.\n\n6 Agripas Street, Tel Aviv\n6249106 ISRAEL",  
16    "developerEmail": "feedback+androidtodo@any.do",  
17    "developerId": "5304780265295461149",  
18    "developerInternalID": "5304780265295461149",  
19    "developerWebsite": "https://www.any.do",  
20    "free": true,  
21    "genre": "Productivity",  
22    "genreId": "PRODUCTIVITY",  
23    "headerImage": "https://lh3.googleusercontent.com/dZknnlk1LM8fYS3wj0vVH0mWKOGH1HAe\\  
24 691Yuh7LAeBj6a730A1CQqZnXxjNahAYUFFw",  
25    "histogram": [27291, 9246, 13735, 29904, 262997],  
26    "icon": "https://lh3.googleusercontent.com/zgOLUXCHkF91H8xuMTMLT17smwgLPwSBjU1KVWF\\  
27 -cZRFj1v-Uvtman7DiHEii54fbEE",  
28    "installs": "10,000,000+",  
29    "minInstalls": 10000000,
```

```

30  "offersIAP": true,
31  "price": 0,
32  "privacyPolicy": "https://www.any.do/privacy",
33  "ratings": 343174,
34  "recentChanges": "Faster and smoother for better user experience!",
35  "recentChangesHTML": "Faster and smoother for better user experience!",
36  "released": "Nov 10, 2011",
37  "reviews": 122170,
38  "score": 4.43388,
39  "screenshots": [
40      "https://lh3.googleusercontent.com/C-L3_FPM1KVrZItAORaszhnQz1zMyXcqF_-oGaabHm_On\wUW1jz02BXBVSKi0HRUTQ",
41      "https://lh3.googleusercontent.com/uAP6G5ANQcgVs4Uj6yrcsAo40UhejTJRVCX0xnAVA5Efi\OtAnrOYyL1SUhj1rv",
42      "https://lh3.googleusercontent.com/AI5mLFu0Ats10km2F09_IwJXNy_1q1_X6Ua3EVMZNedp0\dsDToDRaWQ1UDvI6mb1-I0",
43      "https://lh3.googleusercontent.com/bYCAAn3mjgB4ugSY0PL-PCcMBfbvXCSFkzL-pLSI1bZ8sQ\ByQPerHboPQ2fA126K4LDtU",
44      "https://lh3.googleusercontent.com/u-dX4lpTepsVXs33ds4xxYpApuGS4JBAEb0UsvY_fPbpt\xnF0QxaKNW0-tJVXaP8a1E",
45      "https://lh3.googleusercontent.com/qvUz_9IXHQd6FSLUALZo8NKLx-s4uDGYElPOGRsU28TCE\ficQc0BoNRloRRLqUkH2A",
46      "https://lh3.googleusercontent.com/tEyGs6MG1Y97ccLc4c_HxV9xNOpSvwQyHz6uGAezkVtxm\1ydAaTj5EZSUgqlg69qrrk",
47      "https://lh3.googleusercontent.com/StN0i2BskOs6HCfaP00DMBOCQMCAg3okWVI_S1FJtMytw\bgNMBnD5i9hbSqdN1Gxffffmn",
48      "https://lh3.googleusercontent.com/GRKqWfo-PLzCKwpgZ8fej4PGsUp1q9eM5a3LQeiYCOW-K\UpCOIHXOp3mteZWbJ-pz4My",
49      "https://lh3.googleusercontent.com/pFQQ_qi8u92duWCNXpEcNKpH21VpD_hFd5f-U1TP_f6wf\t3YyYLMzwLitxt-UI6G8vs",
50      "https://lh3.googleusercontent.com/AoeCU6bT1x0eHRvJwvQy0SKJ31oSayox959qMNVaSzz3u\N9bvk1cGek5zyRDe1BdtA",
51      "https://lh3.googleusercontent.com/vICme1f4J9vFt8wY3xBY-LshGgYyvSbsa4TLJyEtNsyoa\1UI0i9oMQVq8oJ41_yR1Aw",
52      "https://lh3.googleusercontent.com/7sn9m__iVM-peiG6_jkKBuE-QVH_xDaycF_oR1XJ1wcAC\45ybNZ_Exor09ENOJ41Q2U",
53      "https://lh3.googleusercontent.com/9I_m2ZXgPtU4Po4cw_cyIaEpZxynxQ1n3YkhFgakATfb\u63a8_f8vGQDxKOHYITzew"
54  ],
55  "size": "Varies with device",
56  "summary": "Task Manager \u2705 Organizer \ud83d\udcc5 Agenda \ud83d\udcdd Daily R\eminders \ud83d\udd14 All-in-One Simple App.",
57  "summaryHTML": "Task Manager \u2705 Organizer \ud83d\udcc5 Agenda \ud83d\udcdd Dai\ud83d\udcbb"
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72

```

```

73  "ly Reminders \ud83d\udd14 All-in-One Simple App.",
74  "title": "Any.do: To do list, Calendar, Planner & Reminders",
75  "updated": 1586258773,
76  "url": "https://play.google.com/store/apps/details?id=com.anydo&hl=en&gl=us",
77  "version": "Varies with device",
78  "video": "https://www.youtube.com/embed/2nk11LD0x6o?ps=play&vq=large&rel=0&autohid\
79 e=1&showinfo=0",
80  "videoImage": "https://i.ytimg.com/vi/2nk11LD0x6o/hqdefault.jpg"
81 }

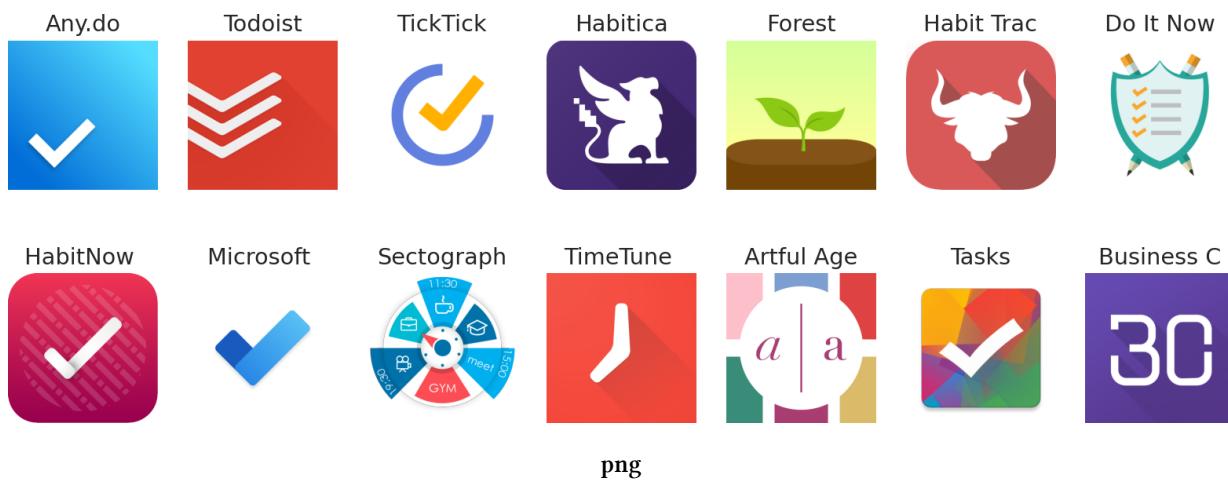
```

This contains lots of information including the number of ratings, number of reviews and number of ratings for each score (1 to 5). Let's ignore all of that and have a look at their beautiful icons:

```

1 def format_title(title):
2     sep_index = title.find(':') if title.find(':') != -1 else title.find('-')
3     if sep_index != -1:
4         title = title[:sep_index]
5     return title[:10]
6
7 fig, axs = plt.subplots(2, len(app_infos) // 2, figsize=(14, 5))
8
9 for i, ax in enumerate(axs.flat):
10    ai = app_infos[i]
11    img = plt.imread(ai['icon'])
12    ax.imshow(img)
13    ax.set_title(format_title(ai['title']))
14    ax.axis('off')

```



We'll store the app information for later by converting the JSON objects into a Pandas dataframe and saving the result into a CSV file:

```

1 app_infos_df = pd.DataFrame(app_infos)
2 app_infos_df.to_csv('apps.csv', index=None, header=True)

```

## Scraping App Reviews

In an ideal world, we would get all the reviews. But there are lots of them and we're scraping the data. That wouldn't be very polite. What should we do?

We want:

- Balanced dataset - roughly the same number of reviews for each score (1-5)
- A representative sample of the reviews for each app

We can satisfy the first requirement by using the scraping package option to filter the review score. For the second, we'll sort the reviews by their helpfulness, which are the reviews that Google Play thinks are most important. Just in case, we'll get a subset from the newest, too:

```

1 app_reviews = []
2
3 for ap in tqdm(app_packages):
4     for score in list(range(1, 6)):
5         for sort_order in [Sort.MOST_RELEVANT, Sort.NEWEST]:
6             rvs, _ = reviews(
7                 ap,
8                 lang='en',
9                 country='us',
10                sort=sort_order,
11                count= 200 if score == 3 else 100,
12                filter_score_with=score
13            )
14            for r in rvs:
15                r['sortOrder'] = 'most_relevant' if sort_order == Sort.MOST_RELEVANT else '\n\
16 newest'
17                r['appId'] = ap
18            app_reviews.extend(rvs)

1 100%|██████████████| 15/15 [00:45<00:00,  3.01s/it]

```

Note that we're adding the app id and sort order to each review. Here's an example for one:

```

1 print_json(app_reviews[0])

1 {
2     "appId": "com.anydo",
3     "at": "2020-04-05 22:25:57",
4     "content": "Update: After getting a response from the developer I would change my \
5 rating to 0 stars if possible. These guys hide behind confusing and opaque terms and\
6 refuse to budge at all. I'm so annoyed that my money has been lost to them! Really \
7 terrible customer experience. Original: Be very careful when signing up for a free t\
8 rial of this app. If you happen to go over they automatically charge you for a full \
9 years subscription and refuse to refund. Terrible customer experience and the app is\
10 just OK.",
11     "repliedAt": "2020-04-07 14:09:03",
12     "replyContent": "Our policy and TOS are completely transparent and can be found in\
13 the Help Center and our main page. In addition, a payment can only be made upon the\
14 user's authorization via the app and Google Play. We provide users with a full 7 da\
15 ys trial to test the app with an additional 48 hours for a refund, along with priori\
16 ty support for all issues.",
17     "reviewCreatedVersion": "4.17.0.3",
18     "score": 1,
19     "sortOrder": "most_relevant",
20     "thumbsUpCount": 37,
21     "userImage": "https://lh3.googleusercontent.com/a-/AOh14GiHdfNEu1DwwcJ6yNyju8Yvn4J\
22 wjpzuXvD74aVmDA",
23     "userName": "Andrew Thomas"
24 }
```

repliedAt and replyContent contain the developer response to the review. Of course, they can be missing.

How many app reviews did we get?

```
1 len(app_reviews)
```

```
1 15750
```

Let's save the reviews to a CSV file:

```

1 app_reviews_df = pd.DataFrame(app_reviews)
2 app_reviews_df.to_csv('reviews.csv', index=None, header=True)
```

## Summary

Well done! You now have a dataset with more than 15k user reviews from 15 productivity apps. Of course, you can go crazy and get much much more.

- Run the notebook in your browser (Google Colab)<sup>5</sup>
- Read the Getting Things Done with Pytorch book<sup>6</sup>

You learned how to:

- Set goals and expectations for your dataset
- Scrape Google Play app information
- Scrape user reviews for Google Play apps
- Save the dataset to CSV files

Next, we're going to use the reviews for sentiment analysis with BERT. But first, we'll have to do some text preprocessing!

## References

- Google Play Scraper for Python<sup>7</sup>

<sup>5</sup><https://colab.research.google.com/drive/1GDJIpz7BXw55jl9wTOMQDool9m8DIOyp>

<sup>6</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>

<sup>7</sup><https://github.com/JoMingyu/google-play-scraper>

# 8. Sentiment Analysis with BERT and Transformers by Hugging Face

TL;DR In this tutorial, you'll learn how to fine-tune BERT for sentiment analysis. You'll do the required text preprocessing (special tokens, padding, and attention masks) and build a Sentiment Classifier using the amazing Transformers library by Hugging Face!

- Run the notebook in your browser (Google Colab)<sup>1</sup>
- Read the Getting Things Done with Pytorch book<sup>2</sup>

You'll learn how to:

- Intuitively understand what BERT is
- Preprocess text data for BERT and build PyTorch Dataset (tokenization, attention masks, and padding)
- Use Transfer Learning to build Sentiment Classifier using the Transformers library by Hugging Face
- Evaluate the model on test data
- Predict sentiment on raw text

Let's get started!

## What is BERT?

BERT (introduced in [this paper<sup>3</sup>](#)) stands for Bidirectional Encoder Representations from Transformers. If you don't know what most of that means - you've come to the right place! Let's unpack the main ideas:

- Bidirectional - to understand the text you're looking you'll have to look back (at the previous words) and forward (at the next words)
- Transformers - The [Attention Is All You Need<sup>4</sup>](#) paper presented the Transformer model. The Transformer reads entire sequences of tokens at once. In a sense, the model is non-directional, while LSTMs read sequentially (left-to-right or right-to-left). The attention mechanism allows for learning contextual relations between words (e.g. his in a sentence refers to Jim).

---

<sup>1</sup><https://colab.research.google.com/drive/1PHv-IRLPcv7oTcIGbsgZHqrB5LPvB7S>

<sup>2</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>

<sup>3</sup><https://arxiv.org/abs/1810.04805>

<sup>4</sup><https://arxiv.org/abs/1706.03762>

- (Pre-trained) contextualized word embeddings - [The ELMO paper<sup>5</sup>](#) introduced a way to encode words based on their meaning/context. Nails has multiple meanings - fingernails and metal nails.

BERT was trained by masking 15% of the tokens with the goal to guess them. An additional objective was to predict the next sentence. Let's look at examples of these tasks:

## **Masked Language Modeling (Masked LM)**

The objective of this task is to guess the masked tokens. Let's look at an example, and try to not make it harder than it has to be:

That's [mask] she [mask] -> That's what she said

## **Next Sentence Prediction (NSP)**

Given a pair of two sentences, the task is to say whether or not the second follows the first (binary classification). Let's continue with the example:

*Input* = [CLS] That's [mask] she [mask]. [SEP] Hahaha, nice! [SEP]

*Label* = *IsNext*

*Input* = [CLS] That's [mask] she [mask]. [SEP] Dwight, you ignorant [mask]! [SEP]

*Label* = *NotNext*

The training corpus was comprised of two entries: [Toronto Book Corpus<sup>6</sup>](#) (800M words) and English Wikipedia (2,500M words). While the original Transformer has an encoder (for reading the input) and a decoder (that makes the prediction), BERT uses only the decoder.

BERT is simply a pre-trained stack of Transformer Encoders. How many Encoders? We have two versions - with 12 (BERT base) and 24 (BERT Large).

## **Is This Thing Useful in Practice?**

The BERT paper was released along with [the source code<sup>7</sup>](#) and pre-trained models.

The best part is that you can do Transfer Learning (thanks to the ideas from OpenAI Transformer) with BERT for many NLP tasks - Classification, Question Answering, Entity Recognition, etc. You can train with small amounts of data and achieve great performance!

---

<sup>5</sup><https://arxiv.org/abs/1802.05365v2>

<sup>6</sup><https://arxiv.org/abs/1506.06724>

<sup>7</sup><https://github.com/google-research/bert>

## Setup

We'll need the [Transformers library](#)<sup>8</sup> by Hugging Face:

```
1 !pip install -qq transformers

1 %reload_ext watermark
2 %watermark -v -p numpy,pandas,torch,transformers

1 CPython 3.6.9
2 IPython 5.5.0
3
4 numpy 1.18.2
5 pandas 1.0.3
6 torch 1.4.0
7 transformers 2.8.0

1 import transformers
2 from transformers import BertModel, BertTokenizer, AdamW, get_linear_schedule_with_w\
3 armup
4 import torch
5
6 import numpy as np
7 import pandas as pd
8 import seaborn as sns
9 from pylab import rcParams
10 import matplotlib.pyplot as plt
11 from matplotlib import rc
12 from sklearn.model_selection import train_test_split
13 from sklearn.metrics import confusion_matrix, classification_report
14 from collections import defaultdict
15 from textwrap import wrap
16
17 from torch import nn, optim
18 from torch.utils.data import Dataset, DataLoader
19
20 %matplotlib inline
21 %config InlineBackend.figure_format='retina'
```

<sup>8</sup><https://huggingface.co/transformers/>

```

22
23 sns.set(style='whitegrid', palette='muted', font_scale=1.2)
24
25 HAPPY_COLORS_PALETTE = ["#01BEFE", "#FFDD00", "#FF7D00", "#FF006D", "#ADFF02", "#8F0\OFF"]
26
27 sns.set_palette(sns.color_palette(HAPPY_COLORS_PALETTE))
28
29
30 rcParams['figure.figsize'] = 12, 8
31
32 RANDOM_SEED = 42
33 np.random.seed(RANDOM_SEED)
34 torch.manual_seed(RANDOM_SEED)
35 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")

```

## Data Exploration

We'll load the Google Play app reviews dataset, that we've put together in the previous part:

```

1 !gdown --id 1S6qMioqPJjyBLpLVz4gmRTnJHnjitnuV
2 !gdown --id 1zdmewp7ayS4js4VtrJEHzAheSW-5NBZv

```

```

1 df = pd.read_csv("reviews.csv")
2 df.shape

1 (15746, 11)

```

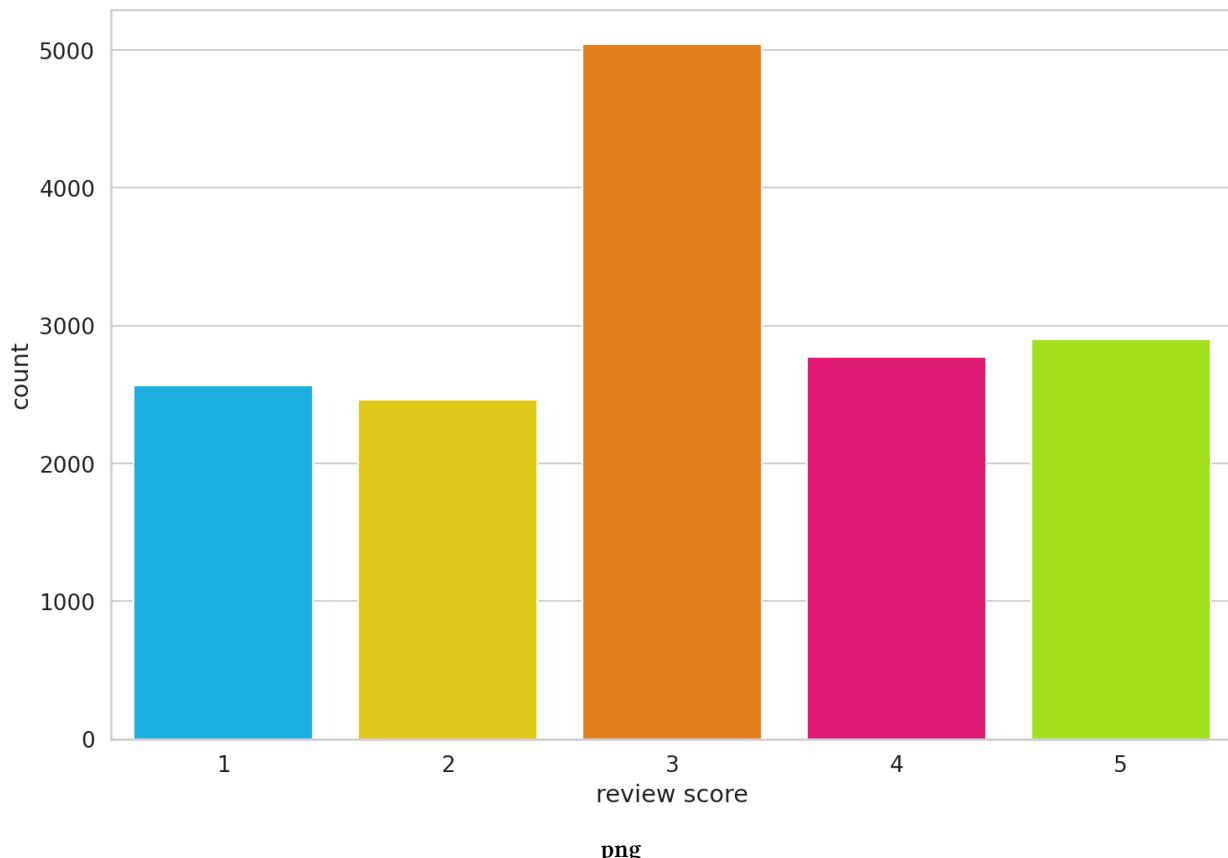
We have about 16k examples. Let's check for missing values:

```
1 df.info()
```

```
1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 15746 entries, 0 to 15745
3 Data columns (total 11 columns):
4 #   Column           Non-Null Count Dtype
5 ---  -----
6 0   userName        15746 non-null  object
7 1   userImage       15746 non-null  object
8 2   content         15746 non-null  object
9 3   score           15746 non-null  int64
10 4   thumbsUpCount  15746 non-null  int64
11 5   reviewCreatedVersion 13533 non-null  object
12 6   at              15746 non-null  object
13 7   replyContent    7367 non-null   object
14 8   repliedAt       7367 non-null   object
15 9   sortOrder       15746 non-null  object
16 10  appId          15746 non-null  object
17 dtypes: int64(2), object(9)
18 memory usage: 1.3+ MB
```

Great, no missing values in the score and review texts! Do we have class imbalance?

```
1 sns.countplot(df.score)
2 plt.xlabel('review score');
```

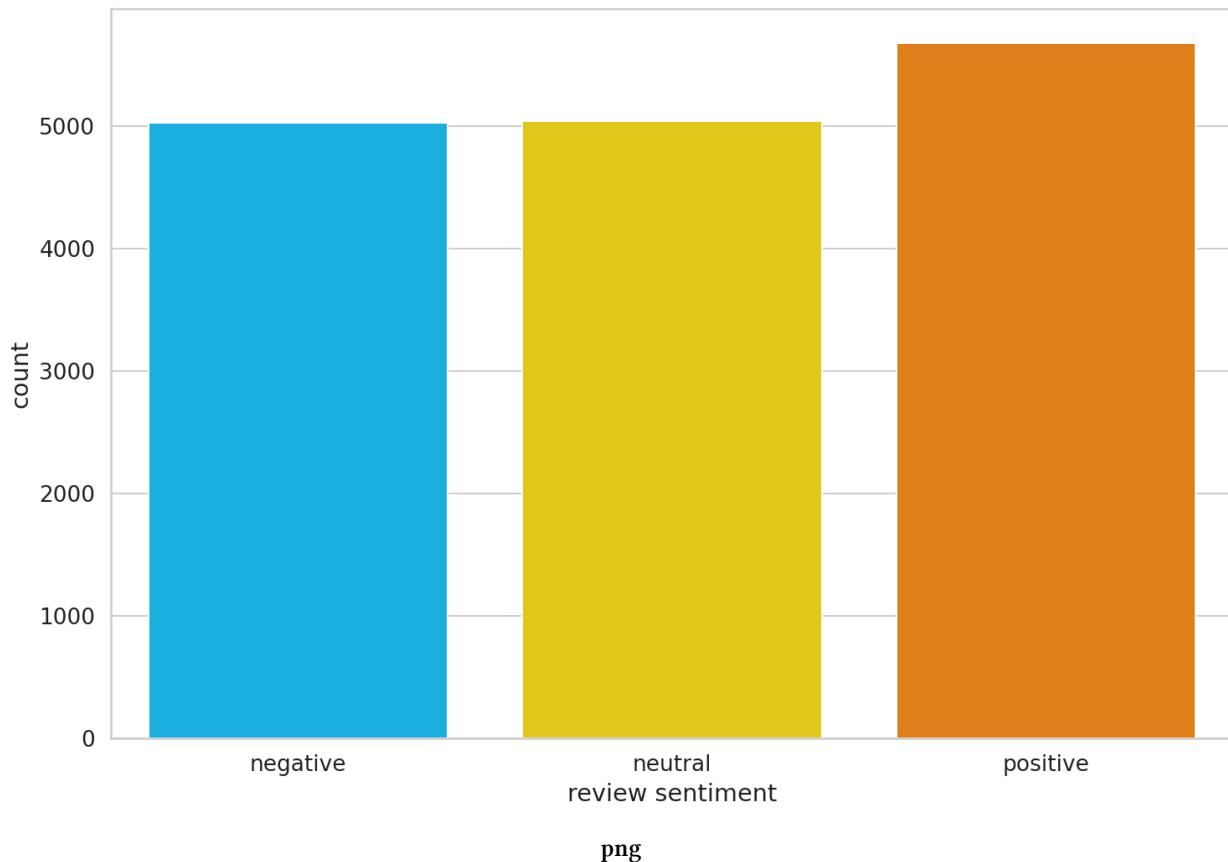


That's hugely imbalanced, but it's okay. We're going to convert the dataset into negative, neutral and positive sentiment:

```
1 def to_sentiment(rating):
2     rating = int(rating)
3     if rating <= 2:
4         return 0
5     elif rating == 3:
6         return 1
7     else:
8         return 2
9
10 df['sentiment'] = df.score.apply(to_sentiment)

1 class_names = ['negative', 'neutral', 'positive']
```

```
1 ax = sns.countplot(df.sentiment)
2 plt.xlabel('review sentiment')
3 ax.set_xticklabels(class_names);
```



The balance was (mostly) restored.

## Data Preprocessing

You might already know that Machine Learning models don't work with raw text. You need to convert text to numbers (of some sort). BERT requires even more attention (good one, right?). Here are the requirements:

- Add special tokens to separate sentences and do classification
- Pass sequences of constant length (introduce padding)
- Create array of 0s (pad token) and 1s (real token) called *attention mask*

The Transformers library provides (you've guessed it) a wide variety of Transformer models (including BERT). It works with TensorFlow and PyTorch! It also includes prebuild tokenizers that do the heavy lifting for us!

```
1 PRE_TRAINED_MODEL_NAME = 'bert-base-cased'
```

You can use a cased and uncased version of BERT and tokenizer. I've experimented with both. The cased version works better. Intuitively, that makes sense, since "BAD" might convey more sentiment than "bad".

Let's load a pre-trained `BertTokenizer`<sup>9</sup>:

```
1 tokenizer = BertTokenizer.from_pretrained(PRE_TRAINED_MODEL_NAME)
```

We'll use this text to understand the tokenization process:

```
1 sample_txt = 'When was I last outside? I am stuck at home for 2 weeks.'
```

Some basic operations can convert the text to tokens and tokens to unique integers (ids):

```
1 tokens = tokenizer.tokenize(sample_txt)
2 token_ids = tokenizer.convert_tokens_to_ids(tokens)
3
4 print(f' Sentence: {sample_txt}')
5 print(f' Tokens: {tokens}')
6 print(f' Token IDs: {token_ids}')

1 Sentence: When was I last outside? I am stuck at home for 2 weeks.
2     Tokens: ['When', 'was', 'I', 'last', 'outside', '?', 'I', 'am', 'stuck', 'at', 'h\
3 ome', 'for', '2', 'weeks', '.']
4 Token IDs: [1332, 1108, 146, 1314, 1796, 136, 146, 1821, 5342, 1120, 1313, 1111, 123\
5 , 2277, 119]
```

## Special Tokens

[SEP] - marker for ending of a sentence

```
1 tokenizer.sep_token, tokenizer.sep_token_id
```

---

<sup>9</sup>[https://huggingface.co/transformers/model\\_doc/bert.html#berttokenizer](https://huggingface.co/transformers/model_doc/bert.html#berttokenizer)

```

1 ('[SEP]', 102)

[CLS] - we must add this token to the start of each sentence, so BERT knows we're doing
classification

1 tokenizer.cls_token, tokenizer.cls_token_id

1 ('[CLS]', 101)

```

There is also a special token for padding:

```

1 tokenizer.pad_token, tokenizer.pad_token_id

1 ('[PAD]', 0)

```

BERT understands tokens that were in the training set. Everything else can be encoded using the [UNK] (unknown) token:

```

1 tokenizer.unk_token, tokenizer.unk_token_id

1 ('[UNK]', 100)

```

All of that work can be done using the `encode_plus()`<sup>10</sup> method:

```

1 encoding = tokenizer.encode_plus(
2     sample_txt,
3     max_length=32,
4     add_special_tokens=True, # Add '[CLS]' and '[SEP]'
5     return_token_type_ids=False,
6     pad_to_max_length=True,
7     return_attention_mask=True,
8     return_tensors='pt',   # Return PyTorch tensors
9 )
10
11 encoding.keys()

```

---

<sup>10</sup>[https://huggingface.co/transformers/main\\_classes/tokenizer.html#transformers.PreTrainedTokenizer.encode\\_plus](https://huggingface.co/transformers/main_classes/tokenizer.html#transformers.PreTrainedTokenizer.encode_plus)

```
1 dict_keys(['input_ids', 'attention_mask'])
```

The token ids are now stored in a Tensor and padded to a length of 32:

```
1 print(len(encoding['input_ids'][0]))  
2 encoding['input_ids'][0]
```

```

1 32
2 tensor([ 101, 1332, 1108, 146, 1314, 1796, 136, 146, 1821, 5342, 1120, 1313,
3           1111, 123, 2277, 119, 102, 0, 0, 0, 0, 0, 0, 0, 0],
4           0, 0, 0, 0, 0, 0, 0, 0])

```

The attention mask has the same length:

```
1 print(len(encoding['attention_mask'][0]))  
2 encoding['attention_mask']
```

We can inverse the tokenization to have a look at the special tokens:

```
1 tokenizer.convert_ids_to_tokens(encoding['input_ids'][0])
```

```
1 [ '[CLS]' ,  
2 'When' ,  
3 'was' ,  
4 'I' ,  
5 'last' ,  
6 'outside' ,  
7 '?' ,  
8 'I' ,  
9 'am' ,  
10 'stuck' ,  
11 'at' ,  
12 'home' ,  
13 'for' ,  
14 '2' ,  
15 'weeks'
```

```

16    '.',  

17    '[SEP]',  

18    '[PAD]',  

19    '[PAD]',  

20    '[PAD]',  

21    '[PAD]',  

22    '[PAD]',  

23    '[PAD]',  

24    '[PAD]',  

25    '[PAD]',  

26    '[PAD]',  

27    '[PAD]',  

28    '[PAD]',  

29    '[PAD]',  

30    '[PAD]',  

31    '[PAD]',  

32    '[PAD]' ]

```

## Choosing Sequence Length

BERT works with fixed-length sequences. We'll use a simple strategy to choose the max length. Let's store the token length of each review:

```

1 token_lens = []  

2  

3 for txt in df.content:  

4     tokens = tokenizer.encode(txt, max_length=512)  

5     token_lens.append(len(tokens))

```

and plot the distribution:

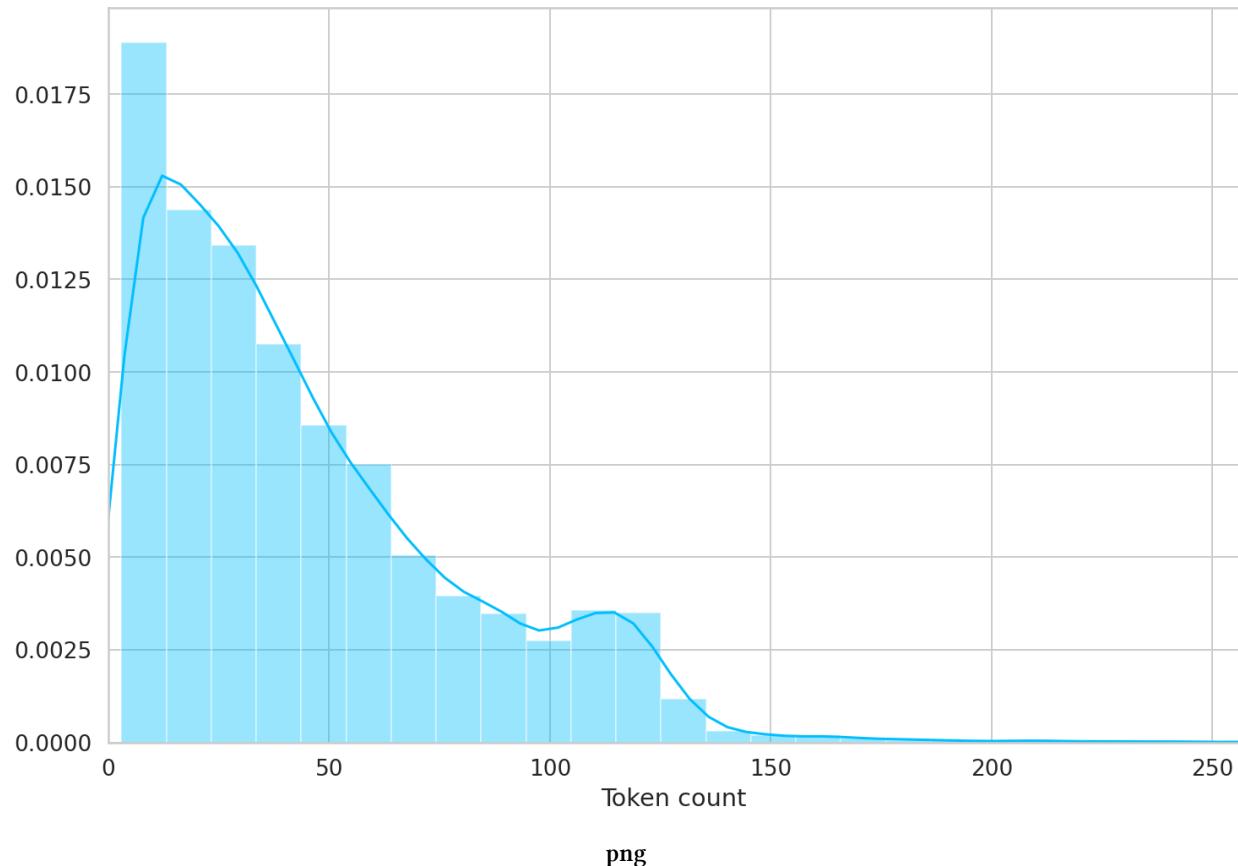
```

1 sns.distplot(token_lens)  

2 plt.xlim([0, 256]);  

3 plt.xlabel('Token count');

```



Most of the reviews seem to contain less than 128 tokens, but we'll be on the safe side and choose a maximum length of 160.

```
1 MAX_LEN = 160
```

We have all building blocks required to create a PyTorch dataset. Let's do it:

```
1 class GPReviewDataset(Dataset):
2
3     def __init__(self, reviews, targets, tokenizer, max_len):
4         self.reviews = reviews
5         self.targets = targets
6         self.tokenizer = tokenizer
7         self.max_len = max_len
8
9     def __len__(self):
10        return len(self.reviews)
11
12    def __getitem__(self, item):
```

```

13     review = str(self.reviews[item])
14     target = self.targets[item]
15
16     encoding = self.tokenizer.encode_plus(
17         review,
18         add_special_tokens=True,
19         max_length=self.max_len,
20         return_token_type_ids=False,
21         pad_to_max_length=True,
22         return_attention_mask=True,
23         return_tensors='pt',
24     )
25
26     return {
27         'review_text': review,
28         'input_ids': encoding['input_ids'].flatten(),
29         'attention_mask': encoding['attention_mask'].flatten(),
30         'targets': torch.tensor(target, dtype=torch.long)
31     }

```

The tokenizer is doing most of the heavy lifting for us. We also return the review texts, so it'll be easier to evaluate the predictions from our model. Let's split the data:

```

1 df_train, df_test = train_test_split(
2     df,
3     test_size=0.1,
4     random_state=RANDOM_SEED
5 )
6 df_val, df_test = train_test_split(
7     df_test,
8     test_size=0.5,
9     random_state=RANDOM_SEED
10 )

```

```

1 df_train.shape, df_val.shape, df_test.shape

```

```

1 ((14171, 12), (787, 12), (788, 12))

```

We also need to create a couple of data loaders. Here's a helper function to do it:

```

1 def create_data_loader(df, tokenizer, max_len, batch_size):
2     ds = GPReviewDataset(
3         reviews=df.content.to_numpy(),
4         targets=df.sentiment.to_numpy(),
5         tokenizer=tokenizer,
6         max_len=max_len
7     )
8
9     return DataLoader(
10        ds,
11        batch_size=batch_size,
12        num_workers=4
13    )

1 BATCH_SIZE = 16
2
3 train_data_loader = create_data_loader(df_train, tokenizer, MAX_LEN, BATCH_SIZE)
4 val_data_loader = create_data_loader(df_val, tokenizer, MAX_LEN, BATCH_SIZE)
5 test_data_loader = create_data_loader(df_test, tokenizer, MAX_LEN, BATCH_SIZE)

```

Let's have a look at an example batch from our training data loader:

```

1 data = next(iter(train_data_loader))
2 data.keys()

1 dict_keys(['review_text', 'input_ids', 'attention_mask', 'targets'])

1 print(data['input_ids'].shape)
2 print(data['attention_mask'].shape)
3 print(data['targets'].shape)

1 torch.Size([16, 160])
2 torch.Size([16, 160])
3 torch.Size([16])

```

## Sentiment Classification with BERT and Hugging Face

There are a lot of helpers that make using BERT easy with the Transformers library. Depending on the task you might want to use [BertForSequenceClassification<sup>11</sup>](#), [BertForQuestionAnswering<sup>12</sup>](#) or something else.

But who cares, right? We're *hardcore!* We'll use the basic [BertModel<sup>13</sup>](#) and build our sentiment classifier on top of it. Let's load the model:

```
1 bert_model = BertModel.from_pretrained(PRE_TRAINED_MODEL_NAME)
```

And try to use it on the encoding of our sample text:

```
1 last_hidden_state, pooled_output = bert_model(
2     input_ids=encoding['input_ids'],
3     attention_mask=encoding['attention_mask']
4 )
```

The `last_hidden_state` is a sequence of hidden states of the last layer of the model. Obtaining the `pooled_output` is done by applying the [BertPooler<sup>14</sup>](#) on `last_hidden_state`:

```
1 last_hidden_state.shape
```

```
1 torch.Size([1, 32, 768])
```

We have the hidden state for each of our 32 tokens (the length of our example sequence). But why 768? This is the number of hidden units in the feedforward-networks. We can verify that by checking the config:

```
1 bert_model.config.hidden_size
```

```
1 768
```

You can think of the `pooled_output` as a summary of the content, according to BERT. Albeit, you might try and do better. Let's look at the shape of the output:

---

<sup>11</sup>[https://huggingface.co/transformers/model\\_doc/bert.html#bertforsequenceclassification](https://huggingface.co/transformers/model_doc/bert.html#bertforsequenceclassification)

<sup>12</sup>[https://huggingface.co/transformers/model\\_doc/bert.html#bertforquestionanswering](https://huggingface.co/transformers/model_doc/bert.html#bertforquestionanswering)

<sup>13</sup>[https://huggingface.co/transformers/model\\_doc/bert.html#bertmodel](https://huggingface.co/transformers/model_doc/bert.html#bertmodel)

<sup>14</sup>[https://github.com/huggingface/transformers/blob/edf0582c0be87b60f94f41c659ea779876fc7be/src/transformers/modeling\\_bert.py#L426](https://github.com/huggingface/transformers/blob/edf0582c0be87b60f94f41c659ea779876fc7be/src/transformers/modeling_bert.py#L426)

```
1 pooled_output.shape
```

```
1 torch.Size([1, 768])
```

We can use all of this knowledge to create a classifier that uses the BERT model:

```
1 class SentimentClassifier(nn.Module):
2
3     def __init__(self, n_classes):
4         super(SentimentClassifier, self).__init__()
5         self.bert = BertModel.from_pretrained(PRE_TRAINED_MODEL_NAME)
6         self.drop = nn.Dropout(p=0.3)
7         self.out = nn.Linear(self.bert.config.hidden_size, n_classes)
8
9     def forward(self, input_ids, attention_mask):
10        _, pooled_output = self.bert(
11            input_ids=input_ids,
12            attention_mask=attention_mask
13        )
14        output = self.drop(pooled_output)
15        return self.out(output)
```

Our classifier delegates most of the heavy lifting to the BertModel. We use a dropout layer for some regularization and a fully-connected layer for our output. Note that we're returning the raw output of the last layer since that is required for the cross-entropy loss function in PyTorch to work.

This should work like any other PyTorch model. Let's create an instance and move it to the GPU

```
1 model = SentimentClassifier(len(class_names))
2 model = model.to(device)
```

We'll move the example batch of our training data to the GPU:

```
1 input_ids = data['input_ids'].to(device)
2 attention_mask = data['attention_mask'].to(device)
3
4 print(input_ids.shape) # batch size x seq length
5 print(attention_mask.shape) # batch size x seq length
```

```

1 torch.Size([16, 160])
2 torch.Size([16, 160])

```

To get the predicted probabilities from our trained model, we'll apply the softmax function to the outputs:

```

1 F.softmax(model(input_ids, attention_mask), dim=1)

1 tensor([[0.5879, 0.0842, 0.3279],
2         [0.4308, 0.1888, 0.3804],
3         [0.4871, 0.1766, 0.3363],
4         [0.3364, 0.0778, 0.5858],
5         [0.4025, 0.1040, 0.4935],
6         [0.3599, 0.1026, 0.5374],
7         [0.5054, 0.1552, 0.3394],
8         [0.5962, 0.1464, 0.2574],
9         [0.3274, 0.1967, 0.4759],
10        [0.3026, 0.1118, 0.5856],
11        [0.4103, 0.1571, 0.4326],
12        [0.4879, 0.2121, 0.3000],
13        [0.3811, 0.1477, 0.4712],
14        [0.3354, 0.1354, 0.5292],
15        [0.3999, 0.2822, 0.3179],
16        [0.5075, 0.1684, 0.3242]], device='cuda:0', grad_fn=<SoftmaxBackward>)

```

## Training

To reproduce the training procedure from the BERT paper, we'll use the `AdamW`<sup>15</sup> optimizer provided by Hugging Face. It corrects weight decay, so it's similar to the original paper. We'll also use a linear scheduler with no warmup steps:

---

<sup>15</sup>[https://huggingface.co/transformers/main\\_classes/optimizer\\_schedules.html#adamw](https://huggingface.co/transformers/main_classes/optimizer_schedules.html#adamw)

```

1 EPOCHS = 10
2
3 optimizer = AdamW(model.parameters(), lr=2e-5, correct_bias=False)
4 total_steps = len(train_data_loader) * EPOCHS
5
6 scheduler = get_linear_schedule_with_warmup(
7     optimizer,
8     num_warmup_steps=0,
9     num_training_steps=total_steps
10 )
11
12 loss_fn = nn.CrossEntropyLoss().to(device)

```

How do we come up with all hyperparameters? The BERT authors have some recommendations for fine-tuning:

- Batch size: 16, 32
- Learning rate (Adam): 5e-5, 3e-5, 2e-5
- Number of epochs: 2, 3, 4

We're going to ignore the number of epochs recommendation but stick with the rest. Note that increasing the batch size reduces the training time significantly, but gives you lower accuracy.

Let's continue with writing a helper function for training our model for one epoch:

```

1 def train_epoch(
2     model,
3     data_loader,
4     loss_fn,
5     optimizer,
6     device,
7     scheduler,
8     n_examples
9 ):
10    model = model.train()
11
12    losses = []
13    correct_predictions = 0
14
15    for d in data_loader:
16        input_ids = d["input_ids"].to(device)
17        attention_mask = d["attention_mask"].to(device)
18        targets = d["targets"].to(device)

```

```

19     outputs = model(
20         input_ids=input_ids,
21         attention_mask=attention_mask
22     )
23
24
25     _, preds = torch.max(outputs, dim=1)
26     loss = loss_fn(outputs, targets)
27
28     correct_predictions += torch.sum(preds == targets)
29     losses.append(loss.item())
30
31     loss.backward()
32     nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
33     optimizer.step()
34     scheduler.step()
35     optimizer.zero_grad()
36
37     return correct_predictions.double() / n_examples, np.mean(losses)

```

Training the model should look familiar, except for two things. The scheduler gets called every time a batch is fed to the model. We're avoiding exploding gradients by clipping the gradients of the model using `clipgrad_norm`<sup>16</sup>.

Let's write another one that helps us evaluate the model on a given data loader:

```

1 def eval_model(model, data_loader, loss_fn, device, n_examples):
2     model = model.eval()
3
4     losses = []
5     correct_predictions = 0
6
7     with torch.no_grad():
8         for d in data_loader:
9             input_ids = d["input_ids"].to(device)
10            attention_mask = d["attention_mask"].to(device)
11            targets = d["targets"].to(device)
12
13            outputs = model(
14                input_ids=input_ids,
15                attention_mask=attention_mask
16            )

```

<sup>16</sup><https://pytorch.org/docs/stable/nn.html#clip-grad-norm>

```

17     _, preds = torch.max(outputs, dim=1)
18
19     loss = loss_fn(outputs, targets)
20
21     correct_predictions += torch.sum(preds == targets)
22     losses.append(loss.item())
23
24 return correct_predictions.double() / n_examples, np.mean(losses)

```

Using those two, we can write our training loop. We'll also store the training history:

```

1 %%time
2
3 history = defaultdict(list)
4 best_accuracy = 0
5
6 for epoch in range(EPOCHS):
7
8     print(f'Epoch {epoch + 1}/{EPOCHS}')
9     print('-' * 10)
10
11    train_acc, train_loss = train_epoch(
12        model,
13        train_data_loader,
14        loss_fn,
15        optimizer,
16        device,
17        scheduler,
18        len(df_train)
19    )
20
21    print(f'Train loss {train_loss} accuracy {train_acc}')
22
23    val_acc, val_loss = eval_model(
24        model,
25        val_data_loader,
26        loss_fn,
27        device,
28        len(df_val)
29    )
30
31    print(f'Val    loss {val_loss} accuracy {val_acc}')
32    print()

```

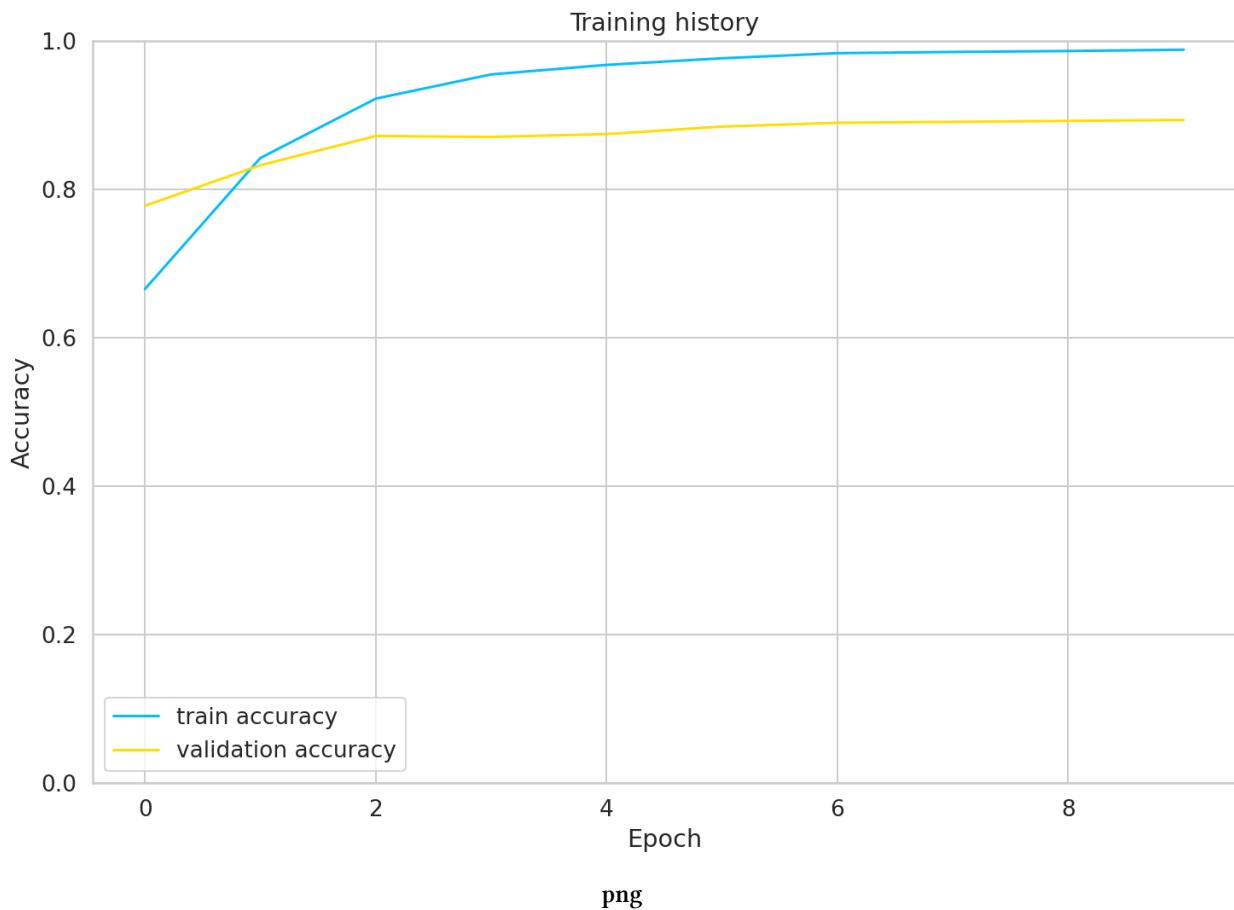
```
33
34     history['train_acc'].append(train_acc)
35     history['train_loss'].append(train_loss)
36     history['val_acc'].append(val_acc)
37     history['val_loss'].append(val_loss)
38
39     if val_acc > best_accuracy:
40         torch.save(model.state_dict(), 'best_model_state.bin')
41         best_accuracy = val_acc
42
43
44 Epoch 1/10
45 -----
46 Train loss 0.7330631300571541 accuracy 0.6653729447463129
47 Val    loss 0.5767546480894089 accuracy 0.7776365946632783
48
49 Epoch 2/10
50 -----
51 Train loss 0.4158683338330777 accuracy 0.8420012701997036
52 Val    loss 0.5365073362737894 accuracy 0.832274459974587
53
54 Epoch 3/10
55 -----
56 Train loss 0.24015077009679367 accuracy 0.922023851527768
57 Val    loss 0.5074492372572422 accuracy 0.8716645489199493
58
59 Epoch 4/10
60 -----
61 Train loss 0.16012676668187295 accuracy 0.9546962105708843
62 Val    loss 0.6009970247745514 accuracy 0.8703939008894537
63
64 Epoch 5/10
65 -----
66 Train loss 0.11209654617575301 accuracy 0.9675393409074872
67 Val    loss 0.7367783848941326 accuracy 0.8742058449809403
68
69 Epoch 6/10
70 -----
71 Train loss 0.08572274737026433 accuracy 0.9764307388328276
72 Val    loss 0.7251267762482166 accuracy 0.8843710292249047
73
74 Epoch 7/10
75 -----
```

```
33 Train loss 0.06132202987342602 accuracy 0.9833462705525369
34 Val    loss 0.7083295831084251 accuracy 0.889453621346887
35
36 Epoch 8/10
37 -----
38 Train loss 0.050604159273123096 accuracy 0.9849693035071626
39 Val    loss 0.753860274553299 accuracy 0.8907242693773825
40
41 Epoch 9/10
42 -----
43 Train loss 0.04373276197092931 accuracy 0.9862395032107826
44 Val    loss 0.7506809896230697 accuracy 0.8919949174078781
45
46 Epoch 10/10
47 -----
48 Train loss 0.03768671146314381 accuracy 0.9880036694658105
49 Val    loss 0.7431786182522774 accuracy 0.8932655654383737
50
51 CPU times: user 29min 54s, sys: 13min 28s, total: 43min 23s
52 Wall time: 43min 43s
```

Note that we're storing the state of the best model, indicated by the highest validation accuracy.

Whoo, this took some time! We can look at the training vs validation accuracy:

```
1 plt.plot(history['train_acc'], label='train accuracy')
2 plt.plot(history['val_acc'], label='validation accuracy')
3
4 plt.title('Training history')
5 plt.ylabel('Accuracy')
6 plt.xlabel('Epoch')
7 plt.legend()
8 plt.ylim([0, 1]);
```



png

The training accuracy starts to approach 100% after 10 epochs or so. You might try to fine-tune the parameters a bit more, but this will be good enough for us.

Don't want to wait? Uncomment the next cell to download my pre-trained model:

```

1 # !gdown --id 1V8itWtowCYnb2Bc9K1K9SxGff9WwmogA
2
3 # model = SentimentClassifier(len(class_names))
4 # model.load_state_dict(torch.load('best_model_state.bin'))
5 # model = model.to(device)

```

## Evaluation

So how good is our model on predicting sentiment? Let's start by calculating the accuracy on the test data:

```

1 test_acc, _ = eval_model(
2     model,
3     test_data_loader,
4     loss_fn,
5     device,
6     len(df_test)
7 )
8
9 test_acc.item()

```

```
1 0.883248730964467
```

The accuracy is about 1% lower on the test set. Our model seems to generalize well.

We'll define a helper function to get the predictions from our model:

```

1 def get_predictions(model, data_loader):
2     model = model.eval()
3
4     review_texts = []
5     predictions = []
6     prediction_probs = []
7     real_values = []
8
9     with torch.no_grad():
10         for d in data_loader:
11
12             texts = d["review_text"]
13             input_ids = d["input_ids"].to(device)
14             attention_mask = d["attention_mask"].to(device)
15             targets = d["targets"].to(device)
16
17             outputs = model(
18                 input_ids=input_ids,
19                 attention_mask=attention_mask
20             )
21             _, preds = torch.max(outputs, dim=1)
22
23             review_texts.extend(texts)
24             predictions.extend(preds)
25             prediction_probs.extend(outputs)
26             real_values.extend(targets)

```

```
27 predictions = torch.stack(predictions).cpu()
28 prediction_probs = torch.stack(prediction_probs).cpu()
29 real_values = torch.stack(real_values).cpu()
30
31 return review_texts, predictions, prediction_probs, real_values
```

This is similar to the evaluation function, except that we're storing the text of the reviews and the predicted probabilities:

```
1 y_review_texts, y_pred, y_pred_probs, y_test = get_predictions(
2     model,
3     test_data_loader
4 )
```

Let's have a look at the classification report

```
1 print(classification_report(y_test, y_pred, target_names=class_names))
```

	precision	recall	f1-score	support
negative	0.89	0.87	0.88	245
neutral	0.83	0.85	0.84	254
positive	0.92	0.93	0.92	289
accuracy			0.88	788
macro avg	0.88	0.88	0.88	788
weighted avg	0.88	0.88	0.88	788

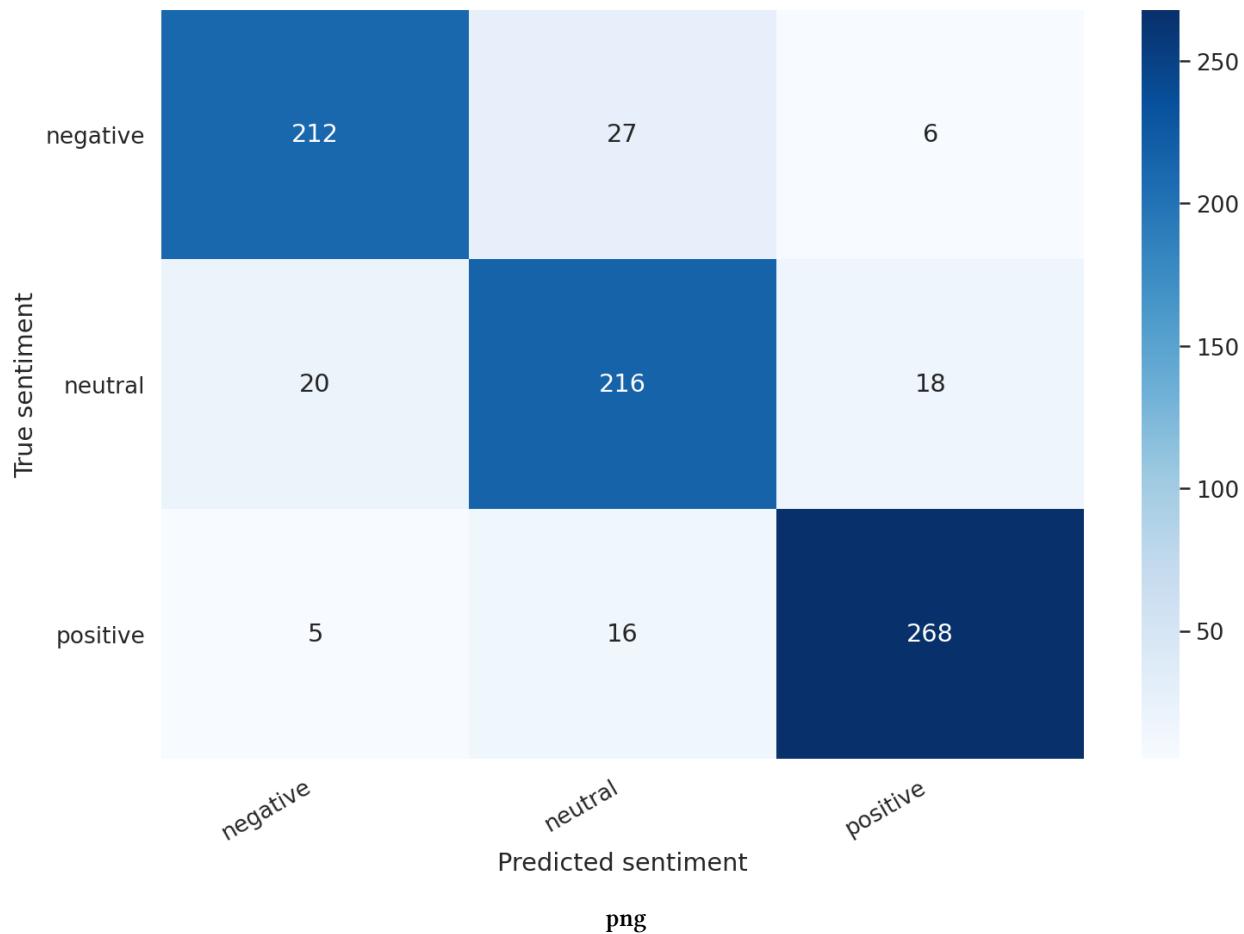
Looks like it is really hard to classify neutral (3 stars) reviews. And I can tell you from experience, looking at many reviews, those are hard to classify.

We'll continue with the confusion matrix:

```

1 def show_confusion_matrix(confusion_matrix):
2     hmap = sns.heatmap(confusion_matrix, annot=True, fmt="d", cmap="Blues")
3     hmap.yaxis.set_ticklabels(hmap.yaxis.get_ticklabels(), rotation=0, ha='right')
4     hmap.xaxis.set_ticklabels(hmap.xaxis.get_ticklabels(), rotation=30, ha='right')
5     plt.ylabel('True sentiment')
6     plt.xlabel('Predicted sentiment');
7
8 cm = confusion_matrix(y_test, y_pred)
9 df_cm = pd.DataFrame(cm, index=class_names, columns=class_names)
10 show_confusion_matrix(df_cm)

```



This confirms that our model is having difficulty classifying neutral reviews. It mistakes those for negative and positive at a roughly equal frequency.

That's a good overview of the performance of our model. But let's have a look at an example from our test data:

```

1 idx = 2
2
3 review_text = y_review_texts[idx]
4 true_sentiment = y_test[idx]
5 pred_df = pd.DataFrame({
6     'class_names': class_names,
7     'values': y_pred_probs[idx]
8 })
9
10 print("\n".join(wrap(review_text)))
11 print()
12 print(f'True sentiment: {class_names[true_sentiment]}')

I used to use Habitica, and I must say this is a great step up. I'd
like to see more social features, such as sharing tasks - only one
person has to perform said task for it to be checked off, but only
giving that person the experience and gold. Otherwise, the price for
subscription is too steep, thus resulting in a sub-perfect score. I
could easily justify $0.99/month or eternal subscription for $15. If
that price could be met, as well as fine tuning, this would be easily
worth 5 stars.

True sentiment: neutral

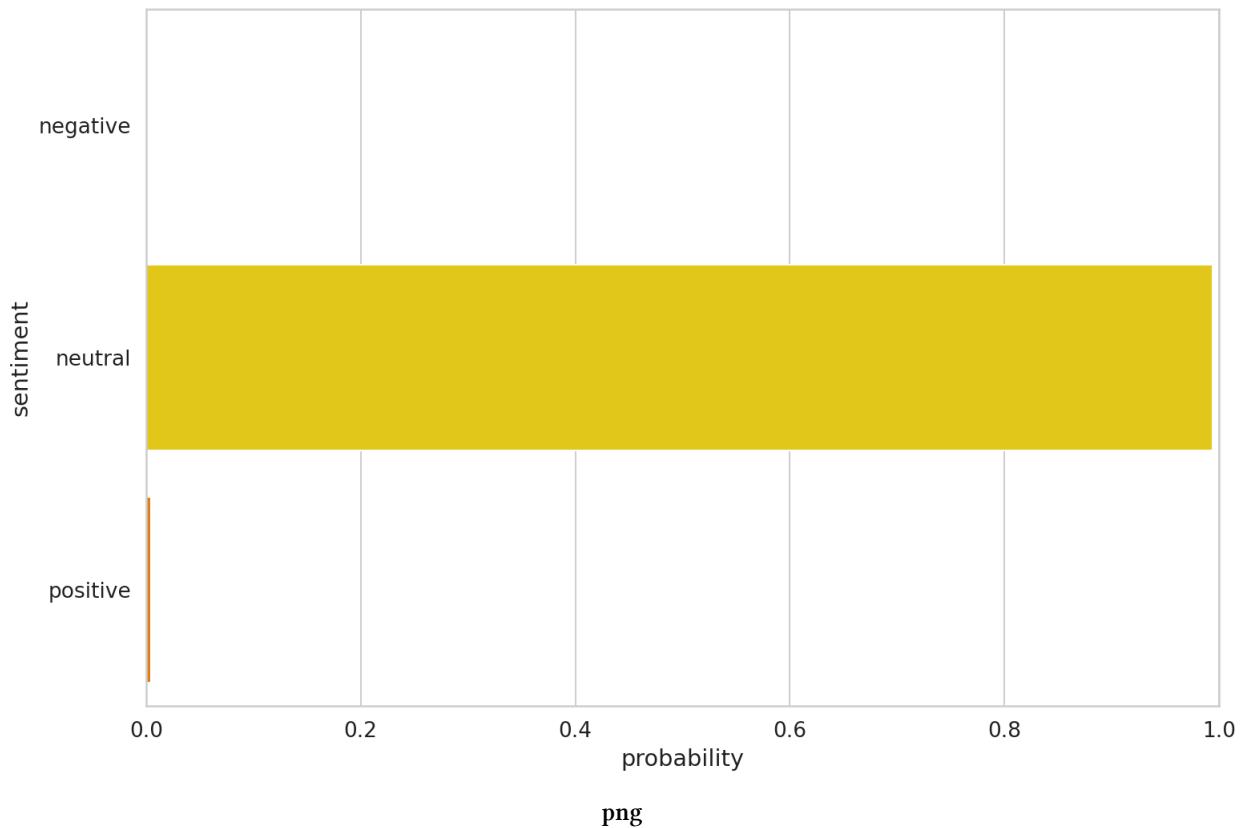
```

Now we can look at the confidence of each sentiment of our model:

```

1 sns.barplot(x='values', y='class_names', data=pred_df, orient='h')
2 plt.ylabel('sentiment')
3 plt.xlabel('probability')
4 plt.xlim([0, 1]);

```



## Predicting on Raw Text

Let's use our model to predict the sentiment of some raw text:

```
1 review_text = "I love completing my todos! Best app ever!!!"
```

We have to use the tokenizer to encode the text:

```
1 encoded_review = tokenizer.encode_plus(  
2     review_text,  
3     max_length=MAX_LEN,  
4     add_special_tokens=True,  
5     return_token_type_ids=False,  
6     pad_to_max_length=True,  
7     return_attention_mask=True,  
8     return_tensors='pt',  
9 )
```

Let's get the predictions from our model:

---

Smile, you are amazing!

```
1 input_ids = encoded_review['input_ids'].to(device)
2 attention_mask = encoded_review['attention_mask'].to(device)
3
4 output = model(input_ids, attention_mask)
5 _, prediction = torch.max(output, dim=1)
6
7 print(f'Review text: {review_text}')
8 print(f'Sentiment : {class_names[prediction]}')


1 Review text: I love completing my todos! Best app ever!!!
2 Sentiment : positive
```

## Summary

Nice job! You learned how to use BERT for sentiment analysis. You built a custom classifier using the Hugging Face library and trained it on our app reviews dataset!

- Run the notebook in your browser (Google Colab)<sup>17</sup>
- Read the Getting Things Done with Pytorch book<sup>18</sup>

You learned how to:

- Intuitively understand what BERT is
- Preprocess text data for BERT and build PyTorch Dataset (tokenization, attention masks, and padding)
- Use Transfer Learning to build Sentiment Classifier using the Transformers library by Hugging Face
- Evaluate the model on test data
- Predict sentiment on raw text

Next, we'll learn how to deploy our trained model behind a REST API and build a simple web app to access it.

<sup>17</sup><https://colab.research.google.com/drive/1PHv-IRLPCTv7oTcIGbsgZHqrB5LPvB7S>

<sup>18</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>

## References

- BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding<sup>19</sup>
- L11 Language Models - Alec Radford (OpenAI)<sup>20</sup>
- The Illustrated BERT, ELMo, and co.<sup>21</sup>
- BERT Fine-Tuning Tutorial with PyTorch<sup>22</sup>
- How to Fine-Tune BERT for Text Classification?<sup>23</sup>
- Huggingface Transformers<sup>24</sup>
- BERT Explained: State of the art language model for NLP<sup>25</sup>

<sup>19</sup><https://arxiv.org/abs/1810.04805>

<sup>20</sup><https://www.youtube.com/watch?v=BnpB3GrpsfM>

<sup>21</sup><https://jalammar.github.io/illustrated-bert/>

<sup>22</sup><https://mccormickml.com/2019/07/22/BERT-fine-tuning/>

<sup>23</sup><https://arxiv.org/pdf/1905.05583.pdf>

<sup>24</sup><https://huggingface.co/transformers/>

<sup>25</sup><https://towardsdatascience.com/bert-explained-state-of-the-art-language-model-for-nlp-f8b21a9b6270>

# 9. Deploy BERT for Sentiment Analysis as REST API using FastAPI

TL;DR Learn how to create a REST API for Sentiment Analysis using a pre-trained BERT model

- [Project on GitHub<sup>1</sup>](#)
- [Run the notebook in your browser \(Google Colab\)<sup>2</sup>](#)
- [Getting Things Done with Pytorch on GitHub<sup>3</sup>](#)

In this tutorial, you'll learn how to deploy a pre-trained BERT model as a REST API using FastAPI. Here are the steps:

- Initialize a project using Pipenv
- Create a project skeleton
- Add the pre-trained model and create an interface to abstract the inference logic
- Update the request handler function to return predictions using the model
- Start the server and send a test request

## Project setup

We'll manage our dependencies using [Pipenv<sup>4</sup>](#). Here's the complete Pipfile:

```
1 [ [source] ]
2 name = "pypi"
3 url = "https://pypi.org/simple"
4 verify_ssl = true
5
6 [dev-packages]
7 black = "==19.10b0"
8 isort = "*"
9 flake8 = "*"
10 gdown = "*"
```

---

<sup>1</sup><https://github.com/curiously/Deploy-BERT-for-Sentiment-Analysis-with-FastAPI>

<sup>2</sup>[https://colab.research.google.com/drive/154jf65arX4cHGaGXl2\\_kJ1DT8FmF4Lhf](https://colab.research.google.com/drive/154jf65arX4cHGaGXl2_kJ1DT8FmF4Lhf)

<sup>3</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>

<sup>4</sup><https://pipenv.pypa.io/en/latest/>

```

11
12 [packages]
13 fastapi = "*"
14 uvicorn = "*"
15 pydantic = "*"
16 torch = "*"
17 transformers = "*"
18
19 [requires]
20 python_version = "3.8"
21
22 [pipenv]
23 allow_prereleases = true

```

The backbone of our REST API will be:

- [FastAPI<sup>5</sup>](#) - lets you easily set up a REST API (some say it might be fast, too)
- [Uvicorn<sup>6</sup>](#) - server that lets you do async programming with Python (pretty cool)
- [Pydantic<sup>7</sup>](#) - data validation by introducing types for our request and response data.

Some tools will help us write some better code (thanks to [Momchil Hardalov<sup>8</sup>](#) for the configs):

- [Black<sup>9</sup>](#) - code formatting
- [isort<sup>10</sup>](#) - imports sorting
- [flake8<sup>11</sup>](#) - check for code style (PEP 8) compliance

## Building a skeleton REST API

Let's start by creating a skeleton structure for our project. Your directory should look like this:

---

<sup>5</sup><https://fastapi.tiangolo.com/>  
<sup>6</sup><https://www.uvicorn.org/>  
<sup>7</sup><https://pydantic-docs.helpmanual.io/>  
<sup>8</sup><https://github.com/mhardalov>  
<sup>9</sup><https://black.readthedocs.io/en/stable/>  
<sup>10</sup><https://timothycrosley.github.io/isort/>  
<sup>11</sup><https://flake8.pycqa.org/en/latest/>

```

1 .
2 └── Pipfile
3 └── Pipfile.lock
4 └── sentiment_analyzer
5     ├── api.py

```

We'll start by creating a dummy/stubbed response to test that everything is working end-to-end. Here are the contents of `api.py`:

```

1 from typing import Dict
2
3 from fastapi import Depends, FastAPI
4 from pydantic import BaseModel
5
6 app = FastAPI()
7
8
9 class SentimentRequest(BaseModel):
10     text: str
11
12
13 class SentimentResponse(BaseModel):
14
15     probabilities: Dict[str, float]
16     sentiment: str
17     confidence: float
18
19
20 @app.post("/predict", response_model=SentimentResponse)
21 def predict(request: SentimentRequest):
22     return SentimentResponse(
23         sentiment="positive",
24         confidence=0.98,
25         probabilities=dict(negative=0.005, neutral=0.015, positive=0.98)
26     )

```

Our API expects a text - the review for sentiment analysis. The response contains the sentiment, confidence (softmax output for the sentiment) and all probabilities for each sentiment.

## Adding our model

Here's the file structure of the complete project:

```

1 .
2 └── assets
3   └── model_state_dict.bin
4 └── bin
5   └── download_model
6 └── config.json
7 └── Pipfile
8 └── Pipfile.lock
9 └── sentiment_analyzer
10  ├── api.py
11  ├── classifier
12  │   ├── model.py
13  │   └── sentiment_classifier.py

```

We'll need the pre-trained model. We'll write the `download_model` script for that:

```

1 #!/usr/bin/env python
2 import gdown
3
4 gdown.download(
5     "https://drive.google.com/uc?id=1V8itWtowCYnb2Bc9K1K9SxGff9WwmogA",
6     "assets/model_state_dict.bin",
7 )

```

The model can be downloaded from my Google Drive. Let's get it:

```
1 python bin/download_model
```

Our pre-trained model is stored as a PyTorch state dict. We need to load it and use it to predict the text sentiment.

Let's start with the config file `config.json`:

```

1 {
2     "BERT_MODEL": "bert-base-cased",
3     "PRE_TRAINED_MODEL": "assets/model_state_dict.bin",
4     "CLASS_NAMES": [
5         "negative",
6         "neutral",
7         "positive"
8     ],
9     "MAX_SEQUENCE_LEN": 160
10 }

```

Next, we'll define the `sentiment_classifier.py`:

```

1 import json
2
3 from torch import nn
4 from transformers import BertModel
5
6 with open("config.json") as json_file:
7     config = json.load(json_file)
8
9
10 class SentimentClassifier(nn.Module):
11     def __init__(self, n_classes):
12         super(SentimentClassifier, self).__init__()
13         self.bert = BertModel.from_pretrained(config["BERT_MODEL"])
14         self.drop = nn.Dropout(p=0.3)
15         self.out = nn.Linear(self.bert.config.hidden_size, n_classes)
16
17     def forward(self, input_ids, attention_mask):
18         _, pooled_output = self.bert(input_ids=input_ids, attention_mask=attention_m\
19 ask)
20         output = self.drop(pooled_output)
21         return self.out(output)

```

This is the same model we've used for training. It just uses the config file.

Recall that BERT requires some special text preprocessing. We need a place to use the tokenizer from Hugging Face. We also need to do some massaging of the model outputs to convert them to our API response format.

The Model provides a nice abstraction (a Facade) to our classifier. It exposes a single predict() method and should be pretty generalizable if you want to use the same project structure as a template for your next deployment. The `model.py` file:

```

1 import json
2
3 import torch
4 import torch.nn.functional as F
5 from transformers import BertTokenizer
6
7 from .sentiment_classifier import SentimentClassifier
8
9 with open("config.json") as json_file:
10     config = json.load(json_file)
11
12

```

```
13 class Model:
14     def __init__(self):
15
16         self.device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
17
18         self.tokenizer = BertTokenizer.from_pretrained(config["BERT_MODEL"])
19
20         classifier = SentimentClassifier(len(config["CLASS_NAMES"]))
21         classifier.load_state_dict(
22             torch.load(config["PRE_TRAINED_MODEL"], map_location=self.device)
23         )
24         classifier = classifier.eval()
25         self.classifier = classifier.to(self.device)
26
27     def predict(self, text):
28         encoded_text = self.tokenizer.encode_plus(
29             text,
30             max_length=config["MAX_SEQUENCE_LEN"],
31             add_special_tokens=True,
32             return_token_type_ids=False,
33             pad_to_max_length=True,
34             return_attention_mask=True,
35             return_tensors="pt",
36         )
37         input_ids = encoded_text["input_ids"].to(self.device)
38         attention_mask = encoded_text["attention_mask"].to(self.device)
39
40         with torch.no_grad():
41             probabilities = F.softmax(self.classifier(input_ids, attention_mask), dim=1)
42             m=1)
43             confidence, predicted_class = torch.max(probabilities, dim=1)
44             predicted_class = predicted_class.cpu().item()
45             probabilities = probabilities.flatten().cpu().numpy().tolist()
46             return (
47                 config["CLASS NAMES"][predicted_class],
48                 confidence,
49                 dict(zip(config["CLASS NAMES"], probabilities)),
50             )
51
52
53 model = Model()
```

```
56 def get_model():
57     return model
```

We'll do the inference on the GPU, if one is available. We return the name of the predicted sentiment, the confidence, and the probabilities for each sentiment.

But why don't we define all that logic in our request handler function? For this tutorial, this is an example of overengineering. But in the real world, when you start testing your implementation, this will be such a nice bonus.

You see, mixing everything in the request handler logic will result in countless sleepless nights. When shit hits the fan (and it will) you'll wonder if your REST or model code is wrong. This way allows you to test them, separately.

The `get_model()` function ensures that we have a single instance of our Model (Singleton). We'll use it in our API handler.

## Putting everything together

Our request handler needs access to the model to return a prediction. We'll use the [Dependency Injection framework<sup>12</sup>](#) provided by FastAPI to inject our model. Here's the new `predict` function:

```
1 @app.post("/predict", response_model=SentimentResponse)
2 def predict(request: SentimentRequest, model: Model = Depends(get_model)):
3     sentiment, confidence, probabilities = model.predict(request.text)
4     return SentimentResponse(
5         sentiment=sentiment, confidence=confidence, probabilities=probabilities
6     )
```

The model gets injected by `Depends` and our Singleton function `get_model`. You can really appreciate the power of abstraction by looking at this!

But does it work?

## Testing the API

Let's fire up the server:

```
1 uvicorn sentiment_analyzer.api:app
```

This should take a couple of seconds to load everything and start the HTTP server.

---

<sup>12</sup><https://fastapi.tiangolo.com/tutorial/dependencies/>

```
1 http POST http://localhost:8000/predict text="This app is a total waste of time!"
```

Here's the response:

```
1 {
2     "confidence": 0.999885082244873,
3     "probabilities": {
4         "negative": 0.999885082244873,
5         "neutral": 8.876612992025912e-05,
6         "positive": 2.614063305372838e-05
7     },
8     "sentiment": "negative"
9 }
```

Let's try with a positive one:

```
1 http POST http://localhost:8000/predict text="OMG. I love how easy it is to stick to\
2 my schedule. Would recommend to everyone!"
```

```
1 {
2     "confidence": 0.999932050704956,
3     "probabilities": {
4         "negative": 1.834999602579046e-05,
5         "neutral": 4.956663542543538e-05,
6         "positive": 0.999932050704956
7     },
8     "sentiment": "positive"
9 }
```

Both results are on point. Feel free to tryout with some real reviews from the Play Store.

## Summary

You should now be a proud owner of ready to deploy (kind of) Sentiment Analysis REST API using BERT. Of course, you're missing lots of stuff to be production-ready - logging, monitoring, alerting, containerization, and much more. But hey, you did good!

- Project on GitHub<sup>13</sup>

---

<sup>13</sup><https://github.com/curiously/Deploy-BERT-for-Sentiment-Analysis-with-FastAPI>

- Run the notebook in your browser (Google Colab)<sup>14</sup>
- Getting Things Done with Pytorch on GitHub<sup>15</sup>

You learned how to:

- Initialize a project using Pipenv
- Create a project skeleton
- Add the pre-trained model and create an interface to abstract the inference logic
- Update the request handler function to return predictions using the model
- Start the server and send a test request

Go on then, deploy and make your users happy!

## References

- FastAPI Homepage<sup>16</sup>
- fastAPI ML quickstart<sup>17</sup>

---

<sup>14</sup>[https://colab.research.google.com/drive/154jf65arX4cHGaGXl2\\_kJ1DT8FmF4Lhf](https://colab.research.google.com/drive/154jf65arX4cHGaGXl2_kJ1DT8FmF4Lhf)

<sup>15</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>

<sup>16</sup><https://fastapi.tiangolo.com/>

<sup>17</sup><https://github.com/cosmic-cortex/fastAPI-ML-quickstart>

# 10. Object Detection on Custom Dataset with YOLO (v5)

TL;DR Learn how to build a custom dataset for YOLO v5 (darknet compatible) and use it to fine-tune a large object detection model. The model will be ready for real-time object detection on mobile devices.

In this tutorial, you'll learn how to fine-tune a pre-trained YOLO v5 model for detecting and classifying clothing items from images.

- Run the notebook in your browser (Google Colab)<sup>1</sup>
- Read the Getting Things Done with Pytorch book<sup>2</sup>

Here's what we'll go over:

- Install required libraries
- Build a custom dataset in YOLO/darknet format
- Learn about YOLO model family history
- Fine-tune the largest YOLO v5 model
- Evaluate the model
- Look at some predictions

How good our final model is going to be?

## Prerequisites

Let's start by installing some required libraries by the YOLOv5 project:

```
1 !pip install torch==1.5.1+cu101 torchvision==0.6.1+cu101 -f https://download.pytorch.org/whl/torch_stable.html
2 !pip install numpy==1.17
3 !pip install PyYAML==5.3.1
4 !pip install git+https://github.com/cocodataset/cocoapi.git#subdirectory=PythonAPI
```

We'll also install [Apex by NVIDIA<sup>3</sup>](#) to speed up the training of our model (this step is optional):

---

<sup>1</sup><https://colab.research.google.com/drive/1e4zvS6LyhOAayEDh3bz8MXFTJcVFSvZX?usp=sharing>

<sup>2</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>

<sup>3</sup><https://nvidia.github.io/apex/>

```

1 !git clone https://github.com/NVIDIA/apex && cd apex && pip install -v --no-cache-dir \
2 --global-option="--cpp_ext" --global-option="--cuda_ext" . --user && cd .. && rm -r \
3 rf apex

```

## Build a dataset

<YouTube youTubeId="NsxDxEJTgRw" />

The dataset contains annotations for clothing items - bounding boxes around shirts, tops, jackets, sunglasses. The dataset is from [DataTurks<sup>4</sup>](#) and is on [Kaggle<sup>5</sup>](#).

```

1 !gdown --id 1uWdQ2kn25RSQITtBHa9_zayplm27IXNC

```

The dataset contains a single JSON file with URLs to all images and bounding box data.

Let's import all required libraries:

```

1 from pathlib import Path
2 from tqdm import tqdm
3 import numpy as np
4 import json
5 import urllib
6 import PIL.Image as Image
7 import cv2
8 import torch
9 import torchvision
10 from IPython.display import display
11 from sklearn.model_selection import train_test_split
12
13 import seaborn as sns
14 from pylab import rcParams
15 import matplotlib.pyplot as plt
16 from matplotlib import rc
17
18 %matplotlib inline
19 %config InlineBackend.figure_format='retina'
20 sns.set(style='whitegrid', palette='muted', font_scale=1.2)
21 rcParams['figure.figsize'] = 16, 10
22
23 np.random.seed(42)

```

---

<sup>4</sup><https://dataturks.com/>

<sup>5</sup><https://www.kaggle.com/dataturks/clothing-item-detection-for-e-commerce>

Each line in the dataset file contains a JSON object. Let's create a list of all annotations:

```

1 clothing = []
2 with open("clothing.json") as f:
3     for line in f:
4         clothing.append(json.loads(line))

```

Here's an example annotation:

```
1 clothing[0]
```

```

1 {'annotation': [{'imageHeight': 312,
2   'imageWidth': 147,
3   'label': ['Tops'],
4   'notes': '',
5   'points': [{['x': 0.02040816326530612, 'y': 0.2532051282051282},
6     {'x': 0.9931972789115646, 'y': 0.8108974358974359}]}],
7   'content': 'http://com.dataturks.a96-i23.open.s3.amazonaws.com/2c9fafb063ad2b650163\
8 b00a1ead0017/4bb8fd9d-8d52-46c7-aa2a-9c18af10aed6__Data_xx1-top-4437-jolliy-original\
9 1-imaekasxahykhd3t.jpeg',
10  'extras': None}

```

We have the labels, image dimensions, bounding box points (normalized in 0-1 range), and an URL to the image file.

Do we have images with multiple annotations?

```

1 for c in clothing:
2     if len(c['annotation']) > 1:
3         display(c)
4
5
6
7
8
9
10

```

```

1 {'annotation': [{'imageHeight': 312,
2   'imageWidth': 265,
3   'label': ['Jackets'],
4   'notes': '',
5   'points': [{['x': 0, 'y': 0.6185897435897436},
6     {'x': 0.026415094339622643, 'y': 0.6185897435897436}]}],
7   'imageHeight': 312,
8   'imageWidth': 265,
9   'label': ['Skirts'],
10  'notes': ''}

```

```

11  'points': [{ 'x': 0.01509433962264151, 'y': 0.03205128205128205},
12  {'x': 1, 'y': 0.9839743589743589}]}],
13  'content': 'http://com.dataturks.a96-i23.open.s3.amazonaws.com/2c9fafb063ad2b650163\
14 b00a1ead0017/b3be330c-c211-45bb-b244-11aef08021c8__Data_free-sk-5108-mudrika-origin\
15 a1-imaf4fz626peqq9f.jpeg',
16  'extras': None}

```

Just a single example. We'll need to handle it, though.

Let's get all unique categories:

```

1 categories = []
2 for c in clothing:
3     for a in c['annotation']:
4         categories.extend(a['label'])
5 categories = list(set(categories))
6 categories.sort()
7 categories

```

```

1 ['Jackets',
2  'Jeans',
3  'Shirts',
4  'Shoes',
5  'Skirts',
6  'Tops',
7  'Trousers',
8  'Tshirts',
9  'sunglasses']

```

We have 9 different categories. Let's split the data into a training and validation set:

```

1 train_clothing, val_clothing = train_test_split(clothing, test_size=0.1)
2 len(train_clothing), len(val_clothing)

```

```

1 (453, 51)

```

## Sample image and annotation

Let's have a look at an image from the dataset. We'll start by downloading it:

```
1 row = train_clothing[10]
2
3 img = urllib.request.urlopen(row["content"])
4 img = Image.open(img)
5 img = img.convert('RGB')
6
7 img.save("demo_image.jpeg", "JPEG")
```

Here's how our sample annotation looks like:

```
1 row
2
3
4
5
6
7
8
9
10
```

```
{'annotation': [{ 'imageHeight': 312,
  'imageWidth': 145,
  'label': ['Tops'],
  'notes': '',
  'points': [ { 'x': 0.013793103448275862, 'y': 0.22756410256410256},
    { 'x': 1, 'y': 0.7948717948717948} ]}],
  'content': 'http://com.dataturks.a96-i23.open.s3.amazonaws.com/2c9fafb063ad2b650163\
b00a1ead0017/ec339ad6-6b73-406a-8971-f7ea35d47577__Data_s-top-203-red-srw-original-\
imaf2nfrxdzvh3k.jpeg',
  'extras': None}
```

We can use OpenCV to read the image:

```
1 img = cv2.cvtColor(cv2.imread(f'demo_image.jpeg'), cv2.COLOR_BGR2RGB)
2 img.shape
3
4
5
```

```
(312, 145, 3)
```

Let's add the bounding box on top of the image along with the label:

```
1  for a in row['annotation']:
2      for label in a['label']:
3
4          w = a['imageWidth']
5          h = a['imageHeight']
6
7          points = a['points']
8          p1, p2 = points
9
10         x1, y1 = p1['x'] * w, p1['y'] * h
11         x2, y2 = p2['x'] * w, p2['y'] * h
12
13         cv2.rectangle(
14             img,
15             (int(x1), int(y1)),
16             (int(x2), int(y2)),
17             color=(0, 255, 0),
18             thickness=2
19         )
20
21         ((label_width, label_height), _) = cv2.getTextSize(
22             label,
23             fontFace=cv2.FONT_HERSHEY_PLAIN,
24             fontScale=1.75,
25             thickness=2
26         )
27
28         cv2.rectangle(
29             img,
30             (int(x1), int(y1)),
31             (int(x1 + label_width + label_width * 0.05), int(y1 + label_height + label_height * 0.25)),
32             color=(0, 255, 0),
33             thickness=cv2.FILLED
34         )
35
36
37         cv2.putText(
38             img,
39             label,
40             org=(int(x1), int(y1 + label_height + label_height * 0.25)), # bottom left
41             fontFace=cv2.FONT_HERSHEY_PLAIN,
42             fontScale=1.75,
43             color=(255, 255, 255),
```

```
44     thickness=2  
45 )
```

The point coordinates are converted back to pixels and used to draw rectangles over the image. Here's the result:

```
1 plt.imshow(img)  
2 plt.axis('off');
```



png

## Convert to YOLO format

YOLO v5 requires the dataset to be in the *darknet format*. Here's an outline of what it looks like:

- One txt with labels file per image
- One row per object
- Each row contains: class\_index bbox\_x\_center bbox\_y\_center bbox\_width bbox\_height
- Box coordinates must be normalized between 0 and 1

Let's create a helper function that builds a dataset in the correct format for us:

```

1 def create_dataset(clothing, categories, dataset_type):
2
3     images_path = Path(f"clothing/images/{dataset_type}")
4     images_path.mkdir(parents=True, exist_ok=True)
5
6     labels_path = Path(f"clothing/labels/{dataset_type}")
7     labels_path.mkdir(parents=True, exist_ok=True)
8
9     for img_id, row in enumerate(tqdm(clothing)):
10
11         image_name = f"{img_id}.jpeg"
12
13         img = urllib.request.urlopen(row["content"])
14         img = Image.open(img)
15         img = img.convert("RGB")
16
17         img.save(str(images_path / image_name), "JPEG")
18
19         label_name = f"{img_id}.txt"
20
21         with (labels_path / label_name).open(mode="w") as label_file:
22
23             for a in row['annotation']:
24
25                 for label in a['label']:
26
27                     category_idx = categories.index(label)
28
29                     points = a['points']
30                     p1, p2 = points
31

```

```

32         x1, y1 = p1['x'], p1['y']
33         x2, y2 = p2['x'], p2['y']
34
35         bbox_width = x2 - x1
36         bbox_height = y2 - y1
37
38         label_file.write(
39             f"{category_idx} {x1 + bbox_width / 2} {y1 + bbox_height / 2} {bbox_widt\
40 h} {bbox_height}\n"
41     )

```

We'll use it to create the train and validation datasets:

```

1 create_dataset(train_clothing, categories, 'train')
2 create_dataset(val_clothing, categories, 'val')

```

Let's have a look at the file structure:

```
1 !tree clothing -L 2
```

```

1 clothing
2   └── images
3   |   ├── train
4   |   └── val
5   └── labels
6   ├── train
7   └── val
8
9 6 directories, 0 files

```

And a single annotation example:

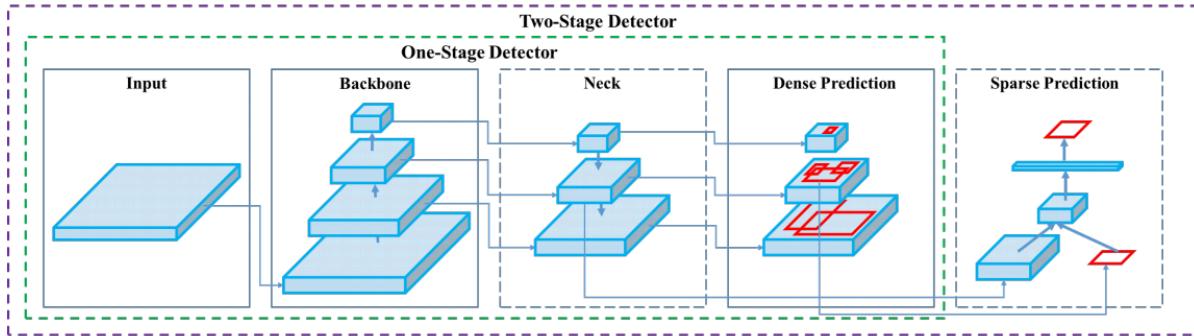
```
1 !cat clothing/labels/train/0.txt
```

```
1 4 0.525462962962963 0.5432692307692308 0.9027777777777778 0.9006410256410257
```

## Fine-tuning YOLO v5

<YouTube youTubeId="XNRzZkZ-Byg" />

The YOLO abbreviation stands for You Only Look Once. YOLO models are one stage object detectors.



*One-stage vs two-stage object detectors. Image from the YOLO v4 paper*

YOLO models are very light and fast. They are [not the most accurate object detections around<sup>6</sup>](#), though. Ultimately, those models are the choice of many (if not all) practitioners interested in [real-time object detection \(FPS >30\)<sup>7</sup>](#).

## Controversy

Joseph Redmon introduced YOLO v1 in the 2016 paper [You Only Look Once: Unified, Real-Time Object Detection<sup>8</sup>](#). The implementation uses the [Darknet Neural Networks library<sup>9</sup>](#).

He also co-authored the YOLO v2 paper in 2017 [YOLO9000: Better, Faster, Stronger<sup>10</sup>](#). A significant improvement over the first iteration with much better localization of objects.

The final iteration, from the original author, was published in the 2018 paper [YOLOv3: An Incremental Improvement<sup>11</sup>](#).

Then things got a bit wacky. Alexey Bochkovskiy published [YOLOv4: Optimal Speed and Accuracy of Object Detection<sup>12</sup>](#) on April 23, 2020. The project has an open-source repository on GitHub<sup>13</sup>.

YOLO v5 got open-sourced on [May 30, 2020<sup>14</sup>](#) by Glenn Jocher<sup>15</sup> from ultralytics. There is no published paper, but the complete project is on GitHub<sup>16</sup>.

<sup>6</sup><https://paperswithcode.com/sota/object-detection-on-coco>

<sup>7</sup><https://paperswithcode.com/sota/real-time-object-detection-on-coco>

<sup>8</sup><https://arxiv.org/pdf/1506.02640.pdf>

<sup>9</sup><https://pjreddie.com/darknet/>

<sup>10</sup><https://arxiv.org/pdf/1612.08242.pdf>

<sup>11</sup><https://arxiv.org/pdf/1804.02767.pdf>

<sup>12</sup><https://arxiv.org/abs/2004.10934>

<sup>13</sup><https://github.com/AlexeyAB/darknet>

<sup>14</sup><https://github.com/ultralytics/yolov5/commit/1e84a23f38fad9e52b59101e9f1246d93066ed1e>

<sup>15</sup><https://github.com/glenn-jocher>

<sup>16</sup><https://github.com/ultralytics/yolov5>

The community at Hacker News got into a [heated debate about the project naming<sup>17</sup>](#). Even the guys at Roboflow wrote [Responding to the Controversy about YOLOv5<sup>18</sup>](#) article about it. They also did a great comparison between YOLO v4 and v5.

My opinion? As long as you put out your work for the whole world to use/see - I don't give a flying fuck. I am not going to comment on points/arguments that are obvious.

## YOLO v5 project setup

YOLO v5 uses PyTorch, but everything is abstracted away. You need the project itself (along with the required dependencies).

Let's start by cloning the GitHub repo and checking out a specific commit (to ensure reproducibility):

```
1 !git clone https://github.com/ultralytics/yolov5
2 %cd yolov5
3 !git checkout ec72eea62bf5bb86b0272f2e65e413957533507f
```

We need two configuration files. One for the dataset and one for the model we're going to use. Let's download them:

```
1 !gdown --id 1ZycPS5Ft_0v1fgHnLsfvZPhcH6qOAqB0 -O data/clothing.yaml
2 !gdown --id 1czESPskb0WZF7_PkCcvRfTiUUJfpX12i -O models/yolov5x.yaml
```

The model config changes the number of classes to 9 (equal to the ones in our dataset). The dataset config `clothing.yaml` is a bit more complex:

```
1 train: ../clothing/images/train/
2 val: ../clothing/images/val/
3
4 nc: 9
5
6 names:
7 [
8     "Jackets",
9     "Jeans",
10    "Shirts",
11    "Shoes",
12    "Skirts",
13    "Tops",
14    "Trousers",
```

<sup>17</sup><https://news.ycombinator.com/item?id=23478151>

<sup>18</sup><https://blog.roboflow.ai/yolov4-versus-yolov5/>

```

15     "Tshirts",
16     "sunglasses",
17 ]

```

This file specifies the paths to the training and validation sets. It also gives the number of classes and their names (you should order those correctly).

## Training

Fine-tuning an existing model is very easy. We'll use the largest model **YOLOv5x** (89M parameters), which is also the most accurate.

In our case, we don't really care about speed. We just want the best accuracy you can get. The checkpoint you're going to use for a different problem(s) is contextually specific. [Take a look at the overview of the pre-trained checkpoints<sup>19</sup>](#).

To train a model on a custom dataset, we'll call the `train.py` script. We'll pass a couple of parameters:

- `img 640` - resize the images to 640x640 pixels
- `batch 4` - 4 images per batch
- `epochs 30` - train for 30 epochs
- `data ./data/clothing.yaml` - path to dataset config
- `cfg ./models/yolov5x.yaml` - model config
- `weights yolov5x.pt` - use pre-trained weights from the YOLOv5x model
- `name yolov5x_clothing` - name of our model
- `cache` - cache dataset images for faster training

```

1 !python train.py --img 640 --batch 4 --epochs 30 \
2   --data ./data/clothing.yaml --cfg ./models/yolov5x.yaml --weights yolov5x.pt \
3   --name yolov5x_clothing --cache

```

The training took around 30 minutes on Tesla P100. The best model checkpoint is saved to `weights/best_yolov5x_clothing.pt`.

## Evaluation

The project includes a great utility function `plot_results()` that allows you to evaluate your model performance on the last training run:

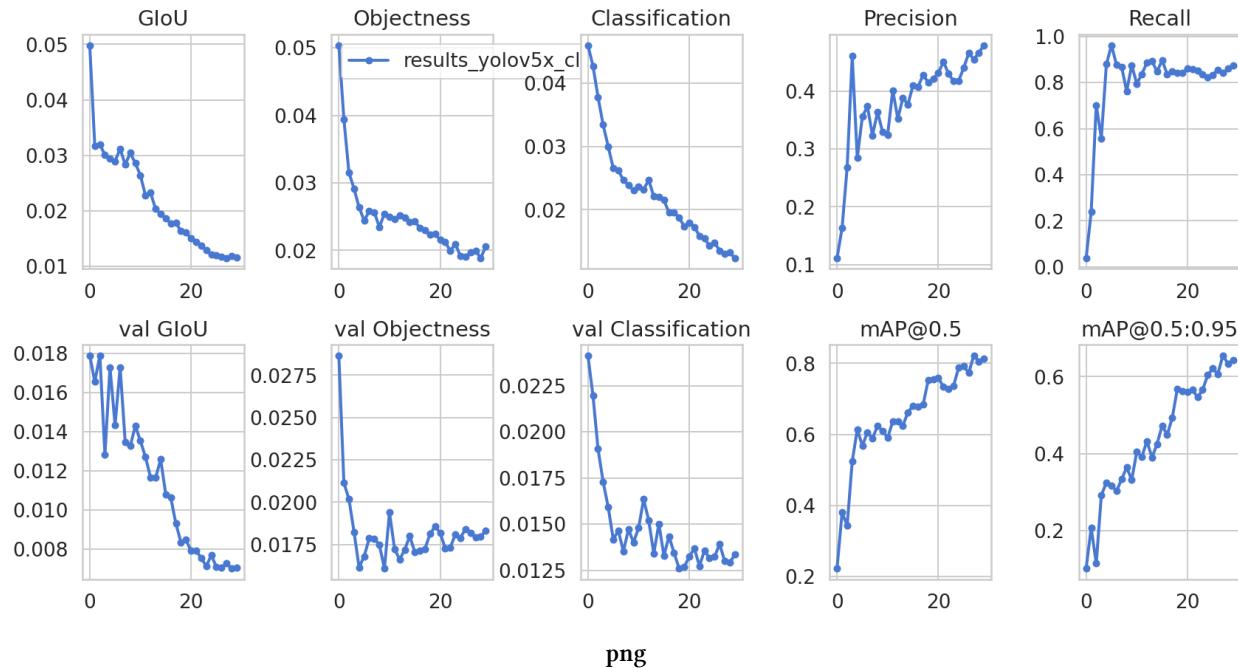
---

<sup>19</sup><https://github.com/ultralytics/yolov5/blob/f9ae460eecd30bdc43a89a37f74b9cc7b93d52f/README.md#pretrained-checkpoints>

```

1 from utils.utils import plot_results
2
3 plot_results();

```



Looks like the mean average precision (mAP) is getting better throughout the training. The model might benefit from more training, but it is good enough.

## Making predictions

Let's pick 50 images from the validation set and move them to `inference/images` to see how our model does on those:

```

1 !find ../clothing/images/val/ -maxdepth 1 -type f | head -50 | xargs cp -t "./inference/images/"
2

```

We'll use the `detect.py` script to run our model on the images. Here are the parameters we're using:

- weights `weights/best_yolov5x_clothing.pt` - checkpoint of the model
- img 640 - resize the images to 640x640 px
- conf 0.4 - take into account predictions with confidence of 0.4 or higher
- source `./inference/images/` - path to the images

```
1 !python detect.py --weights weights/best_yolov5x_clothing.pt \
2   --img 640 --conf 0.4 --source ./inference/images/
```

We'll write a helper function to show the results:

```
1 def load_image(img_path: Path, resize=True):
2     img = cv2.cvtColor(cv2.imread(str(img_path)), cv2.COLOR_BGR2RGB)
3     img = cv2.resize(img, (128, 256), interpolation = cv2.INTER_AREA)
4     return img
5
6 def show_grid(image_paths):
7     images = [load_image(img) for img in image_paths]
8     images = torch.as_tensor(images)
9     images = images.permute(0, 3, 1, 2)
10    grid_img = torchvision.utils.make_grid(images, nrow=11)
11    plt.figure(figsize=(24, 12))
12    plt.imshow(grid_img.permute(1, 2, 0))
13    plt.axis('off');
```

Here are some of the images along with the detected clothing:

```
1 img_paths = list(Path("inference/output").glob("*.jpeg"))[:22]
2 show_grid(img_paths)
```



png

To be honest with you. I am really blown away with the results!

## Summary

You now know how to create a custom dataset and fine-tune one of the YOLO v5 models on your own. Nice!

- Run the notebook in your browser (Google Colab)<sup>20</sup>
- Read the Getting Things Done with Pytorch book<sup>21</sup>

Here's what you've learned:

- Install required libraries
- Build a custom dataset in YOLO/darknet format
- Learn about YOLO model family history
- Fine-tune the largest YOLO v5 model
- Evaluate the model
- Look at some predictions

How well does your model do on your dataset? Let me know in the comments below.

In the next part, you'll learn how to deploy your model a mobile device.

## References

- Clothing Item Detection for E-Commerce dataset<sup>22</sup>
- YOLOv5 GitHub<sup>23</sup>
- YOLOv5 Train on Custom Data<sup>24</sup>
- NVIDIA Apex on GitHub<sup>25</sup>
- YOLOv4: Optimal Speed and Accuracy of Object Detection<sup>26</sup>

<sup>20</sup><https://colab.research.google.com/drive/1e4zvS6LyhOAayEDh3bz8MXFTJcVFSvZX?usp=sharing>

<sup>21</sup><https://github.com/curiously/Getting-Things-Done-with-Pytorch>

<sup>22</sup><https://www.kaggle.com/dataturks/clothing-item-detection-for-e-commerce>

<sup>23</sup><https://github.com/ultralytics/yolov5>

<sup>24</sup><https://github.com/ultralytics/yolov5/wiki/Train-Custom-Data>

<sup>25</sup><https://github.com/NVIDIA/apex>

<sup>26</sup><https://arxiv.org/pdf/2004.10934.pdf>