

Table of Contents

 Learn

 Colab

 Notebook

 GitHub

[Learn the Basics](#) || [Quickstart](#) || [Tensors](#) || [Datasets & DataLoaders](#) || [Transforms](#) || **Build Model** || [Autograd](#) || [Optimization](#) || [Save & Load Model](#)

Build the Neural Network

Created On: Feb 09, 2021 | Last Updated: Jan 24, 2025 | Last Verified: Not Verified

Neural networks comprise of layers/modules that perform operations on data. The `torch.nn` namespace provides all the building blocks you need to build your own neural network. Every module in PyTorch subclasses the `nn.Module`. A neural network is a module itself that consists of other modules (layers). This nested structure allows for building and managing complex architectures easily.

In the following sections, we'll build a neural network to classify images in the FashionMNIST dataset.

```
import os
import torch
from torch import nn
from torch.utils.data import DataLoader
from torchvision import datasets, transforms
```

Get Device for Training

We want to be able to train our model on an `accelerator` such as CUDA, MPS, MTIA, or XPU. If the current accelerator is available, we will use it. Otherwise, we use the CPU.

```
device = torch.accelerator.current_accelerator().type if torch.accelerator.is_available() else "cpu"
print(f"Using {device} device")
```

Out: Using cuda device

Define the Class

We define our neural network by subclassing `nn.Module` , and initialize the neural network layers in `__init__` . Every `nn.Module` subclass implements the operations on input data in the `forward` method.

```
class NeuralNetwork(nn.Module):
    def __init__(self):
        super().__init__()
        self.flatten = nn.Flatten()
        self.linear_relu_stack = nn.Sequential(
            nn.Linear(28*28, 512),
            nn.ReLU(),
            nn.Linear(512, 512),
            nn.ReLU(),
            nn.Linear(512, 10),
        )

    def forward(self, x):
        x = self.flatten(x)
        logits = self.linear_relu_stack(x)
        return logits
```

We create an instance of `NeuralNetwork` , and move it to the `device` , and print its structure.

```
model = NeuralNetwork().to(device)
print(model)
```

```
Out:  NeuralNetwork(
      (flatten): Flatten(start_dim=1, end_dim=-1)
      (linear_relu_stack): Sequential(
        (0): Linear(in_features=784, out_features=512, bias=True)
        (1): ReLU()
        (2): Linear(in_features=512, out_features=512, bias=True)
        (3): ReLU()
        (4): Linear(in_features=512, out_features=10, bias=True)
      )
    )
```

To use the model, we pass it the input data. This executes the model’s `forward` , **along with some background operations**. Do not call `model.forward()` directly!

Calling the model on the input returns a 2-dimensional tensor with `dim=0` corresponding to each output of 10 raw predicted values for each class, and `dim=1` corresponding to the individual values of each output. We get the prediction probabilities by passing it through an instance of the `nn.Softmax` module.

```
X = torch.rand(1, 28, 28, device=device)
logits = model(X)
pred_probab = nn.Softmax(dim=1)(logits)
y_pred = pred_probab.argmax(1)
print(f"Predicted class: {y_pred}")
```

```
Out:  Predicted class: tensor([7], device='cuda:0')
```

Model Layers

Let’s break down the layers in the FashionMNIST model. To illustrate it, we will take a sample minibatch of 3 images of size 28x28 and see what happens to it as we pass it through the network.

```
input_image = torch.rand(3,28,28)
print(input_image.size())
```

```
Out:  torch.Size([3, 28, 28])
```

nn.Flatten

We initialize the **nn.Flatten** layer to convert each 2D 28x28 image into a contiguous array of 784 pixel values (the minibatch dimension (at `dim=0`) is maintained).

```
flatten = nn.Flatten()
flat_image = flatten(input_image)
print(flat_image.size())
```

```
Out:  torch.Size([3, 784])
```

nn.Linear

The **linear layer** is a module that applies a linear transformation on the input using its stored weights and biases.

```
layer1 = nn.Linear(in_features=28*28, out_features=20)
hidden1 = layer1(flat_image)
print(hidden1.size())
```

```
Out:  torch.Size([3, 20])
```

nn.ReLU

Non-linear activations are what create the complex mappings between the model’s inputs and outputs. They are applied after linear transformations to introduce *nonlinearity*, helping neural networks learn a wide variety of phenomena.

In this model, we use **nn.ReLU** between our linear layers, but **there’s other activations to introduce non-linearity in your model**.

```
print(f"Before ReLU: {hidden1}\n\n")
hidden1 = nn.ReLU()(hidden1)
print(f"After ReLU: {hidden1}")
```

Out:

```
Before ReLU: tensor([[ 0.4158, -0.0130, -0.1144,  0.3960,  0.1476, -0.0690, -0.0269,  0.2690,
                    0.1353,  0.1975,  0.4484,  0.0753,  0.4455,  0.5321, -0.1692,  0.4504,
                    0.2476, -0.1787, -0.2754,  0.2462],
                  [ 0.2326,  0.0623, -0.2984,  0.2878,  0.2767, -0.5434, -0.5051,  0.4339,
                    0.0302,  0.1634,  0.5649, -0.0055,  0.2025,  0.4473, -0.2333,  0.6611,
                    0.1883, -0.1250,  0.0820,  0.2778],
                  [ 0.3325,  0.2654,  0.1091,  0.0651,  0.3425, -0.3880, -0.0152,  0.2298,
                    0.3872,  0.0342,  0.8503,  0.0937,  0.1796,  0.5007, -0.1897,  0.4030,
                    0.1189, -0.3237,  0.2048,  0.4343]], grad_fn=<AddmmBackward0>)

After ReLU: tensor([[0.4158, 0.0000, 0.0000, 0.3960, 0.1476, 0.0000, 0.0000, 0.2690, 0.1353,
                    0.1975, 0.4484, 0.0753, 0.4455, 0.5321, 0.0000, 0.4504, 0.2476, 0.0000,
                    0.0000, 0.2462],
                  [0.2326, 0.0623, 0.0000, 0.2878, 0.2767, 0.0000, 0.0000, 0.4339, 0.0302,
                    0.1634, 0.5649, 0.0000, 0.2025, 0.4473, 0.0000, 0.6611, 0.1883, 0.0000,
                    0.0820, 0.2778],
                  [0.3325, 0.2654, 0.1091, 0.0651, 0.3425, 0.0000, 0.0000, 0.2298, 0.3872,
                    0.0342, 0.8503, 0.0937, 0.1796, 0.5007, 0.0000, 0.4030, 0.1189, 0.0000,
```

nn.Sequential

[nn.Sequential](#) is an ordered container of modules. The data is passed through all the modules in the same order as defined. You can use sequential containers to put together a quick network like `seq_modules`.

```
seq_modules = nn.Sequential(
    flatten,
    layer1,
    nn.ReLU(),
    nn.Linear(20, 10)
)
input_image = torch.rand(3,28,28)
logits = seq_modules(input_image)
```

nn.Softmax

The last linear layer of the neural network returns *logits* - raw values in `[-infy, infy]` - which are passed to the [nn.Softmax](#) module. The logits are scaled to values `[0, 1]` representing the model's predicted probabilities for each class. `dim` parameter indicates the dimension along which the values must sum to 1.

```
softmax = nn.Softmax(dim=1)
pred_probab = softmax(logits)
```

Model Parameters

Many layers inside a neural network are *parameterized*, i.e. have associated weights and biases that are optimized during training. Subclassing `nn.Module` automatically tracks all fields defined inside your model object, and makes all parameters accessible using your model's `parameters()` or `named_parameters()` methods.

In this example, we iterate over each parameter, and print its size and a preview of its values.

```
print(f"Model structure: {model}\n\n")

for name, param in model.named_parameters():
    print(f"Layer: {name} | Size: {param.size()} | Values : {param[:2]} \n")
```

Out:

```
Layer: linear_relu_stack.0.bias | Size: torch.Size([512]) | Values : tensor([-0.0155, -0.0327], device='cuda:0', grad_fn=
<SliceBackward0>)

Layer: linear_relu_stack.2.weight | Size: torch.Size([512, 512]) | Values : tensor([[ 0.0116,  0.0293, -0.0280, ...,  0.0334,
-0.0078,  0.0298],
                  [ 0.0095,  0.0038,  0.0009, ..., -0.0365, -0.0011, -0.0221]],
                  device='cuda:0', grad_fn=<SliceBackward0>)

Layer: linear_relu_stack.2.bias | Size: torch.Size([512]) | Values : tensor([ 0.0148, -0.0256], device='cuda:0', grad_fn=
<SliceBackward0>)

Layer: linear_relu_stack.4.weight | Size: torch.Size([10, 512]) | Values : tensor([[ -0.0147, -0.0229,  0.0180, ..., -0.0013,
 0.0177,  0.0070],
                  [ -0.0202, -0.0417, -0.0279, ..., -0.0441,  0.0185, -0.0268]],
                  device='cuda:0', grad_fn=<SliceBackward0>)

Layer: linear_relu_stack.4.bias | Size: torch.Size([10]) | Values : tensor([ 0.0070, -0.0411], device='cuda:0', grad_fn=
<SliceBackward0>)
```

Further Reading

- [torch.nn API](#)

Total running time of the script: (0 minutes 0.395 seconds)

Out:

```
Model structure: NeuralNetwork(
  (flatten): Flatten(start_dim=1, end_dim=-1)
  (linear_relu_stack): Sequential(
    (0): Linear(in_features=784, out_features=512, bias=True)
    (1): ReLU()
    (2): Linear(in_features=512, out_features=512, bias=True)
    (3): ReLU()
    (4): Linear(in_features=512, out_features=10, bias=True)
  )
)

Layer: linear_relu_stack.0.weight | Size: torch.Size([512, 784]) | Values : tensor([[
0.0273,  0.0296, -0.0084, ..., -0.0142,  0.0093,  0.0135],
[-0.0188, -0.0354,  0.0187, ..., -0.0106, -0.0001,  0.0115]],
device='cuda:0', grad_fn=<SliceBackward0>)

Layer: linear_relu_stack.0.bias | Size: torch.Size([512]) | Values : tensor([-0.0155,
-0.0327], device='cuda:0', grad_fn=<SliceBackward0>)
```

PyTorch

Get Started

Features

Ecosystem

Blog

Contributing

Resources

Tutorials

Docs

Discuss

Github Issues

Brand Guidelines

Stay up to date

Facebook

Twitter

YouTube

LinkedIn

PyTorch Podcasts

Spotify

Apple

Google

Amazon

Terms | Privacy

© Copyright The Linux Foundation. The PyTorch Foundation is a project of The Linux Foundation. For web site terms of use, trademark policy and other policies applicable to The PyTorch Foundation please see www.linuxfoundation.org/policies/. The PyTorch Foundation supports the PyTorch open source project, which has been established as PyTorch Project a Series of LF Projects, LLC. For policies applicable to the PyTorch Project a Series of LF Projects, LLC, please see www.lfprojects.org/policies/.