

---

+ Code

+ Text

---

Make sure you fill in any place that says `YOUR CODE HERE`.

---

## ▼ Homework 8

*This* is a Python Notebook homework. It consists of various types of cells:

- Text: you can read them :-)
- Code: you should run them, as they may set up the problems that you are asked to solve.
- **Solution:** These are cells where you should enter a solution. You will see a marker in these cells that indicates where your work should be inserted.

```
# YOUR CODE HERE
```

- Test: These cells contains some tests, and are worth some points. You should run the cells as a way to debug your code, and to see if you understood the question, and whether the output of your code is produced in the correct format. The notebook contains both the tests you see, and some secret ones that you cannot see. This prevents you from using the simple trick of hard-coding the desired output.

## Questions

There are several groups of questions in this notebook:

- Implementation of `TunedSin`
- Implementation of a model using `TunedSin` for MNIST
- Fitting of above model

There are other pieces of text called "exercises", but you only have to do those that are explicitly marked with a place in the code for you to write the answer.

## Working on Your Notebook

To work on your notebook:

- Click on *File > Save a copy in Drive* : this will create a copy of this file in your Google Drive; you will find the notebook in your *Colab Notebooks* folder.
- Work on that notebook. Check that the runtime has GPUs (Runtime > Change Runtime Type, and check that GPU is selected).

## IMPORTANT

Please, in the cell below, define `colab = True`, so that you are able to run the problems on Google Colab. The only reason I define it `False` here is that, when autograding, I do not want to run all the model fitting. If you leave `colab` to be `False`, the notebook will not do anything interesting.

## Submitting Your Notebook

Submit your work as follows:

- Download the notebook from Colab, clicking on "File > Download .ipynb".
- Upload the resulting file to [this Google form](#).
- **Deadline:** [see home page](#)

You can submit multiple times, and the last submission before the deadline will be used to assign you a grade.

```
1 # Change this to True, or else, this notebook will not do anything interesting.
2 colab = True
```

In the previous chapters, we have implemented a machine-learning framework from scratch, on the basis of a symbolic representation for expressions, and automated gradient computation. It is now time for us to experiment with a serious machine-learning framework: we will tackle the problem of image recognition using [PyTorch](#). We will see that PyTorch works in a way that is fundamentally similar to our homegrown framework. On the other hand, PyTorch is vastly more sophisticated in a number of aspects, including:

- *Tensors*. Our variables represented one floating point number. Many machine-learning problems have as input vectors or matrices of numbers; for example, an image is commonly modeled as a 3D-matrix of floating point numbers with dimensions  $w \times h \times c$ , where  $w$  and  $h$  are the image width and height, and where  $c$  is the number of color channels ( $c = 1$  for black and white images, and  $c = 3$  for RGB color images). PyTorch's unit of processing consists in [tensors](#), which represent such multi-dimensional matrices of numbers.
- *Layers and functions*. PyTorch contains [pre-defined layers and functions](#) (corresponding to the expressions of our framework) that many basic and useful machine-learning building blocks, such as fully-connected neural-network layers, convolutional neural-network layers, recurrent networks, and more.
- *GPU optimizations*. PyTorch can exploit [GPUs](#) that support the [CUDA](#) processing language; training or using neural networks on GPUs is often about a hundred times faster than on regular GPUs. Platforms such as [Google Colab](#) provide runtimes with GPUs.
- *Batches*. Our framework processed input one point at a time. PyTorch can process *batches* of data, so that it can train on hundreds of images in the same step. The use of batches makes

machine learning more efficient, as well as better behaved. We will discuss later how the batch size is chosen.

- *Algorithms for tuning the step size.* Finding a rough step size is always a bit of guesswork, and PyTorch provides [many useful algorithms](#) to help tune the learning step size.

## [PyTorch](#), TensorFlow, and scikit-learn

PyTorch and [TensorFlow](#) are likely the two most commonly used deep-learning machine learning frameworks available today. Another widely used machine-learning framework is [scikit-learn](#), pronounced similarly to *psychic learn*.

scikit-learn provides a greater variety of learning algorithms than PyTorch or TensorFlow, including decision trees, SVMs, and more. On the other hand, scikit-learn does not provide the same optimizations as PyTorch and TensorFlow to work efficiently on vast amounts of data, and neither does scikit-learn provide support for GPU acceleration.

PyTorch and TensorFlow are much closer to each other, and often the decision of using one of them rather than the other stems mainly from personal preference. The main reason why we chose PyTorch for this chapter is that its gradient-computation methods are very close to the ones we developed for our homebrew framework.

### ▼ A flat neural network for MNIST

The first problem we will tackle will be the one of *digit recognition*. We will use a dataset made of single digits, each drawn on a  $28 \times 28$  canvas. This dataset, called the [MNIST dataset](#), and the dataset has become the "hello world" of deep machine learning. To use it, we first load PyTorch. If using Google Colab, make sure you have selected a GPU backend.

```
1 import numpy as np
2 import matplotlib
3 import matplotlib.pyplot as plt
4 import random
5 import torch
6 import torch.nn as nn
7 import torch.nn.functional as F
8 import torch.optim as optim
9 from torchvision import datasets, transforms

1 # We use the CUDA device to get GPU acceleration.
2 USE_CUDA = colab
3 if USE_CUDA:
4     device = torch.device("cuda")
```

```

5     kwargs = {'num_workers': 1, 'pin_memory': True}
6 else:
7     device = torch.device("cpu")
8     kwargs = {}

```

Let us load the dataset itself, and display its first digit.

```

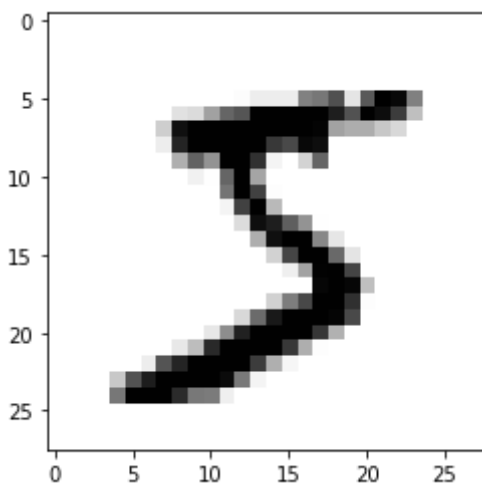
1 if colab:
2     mnist = datasets.MNIST('../data', train=True, download=True,
3                             transform=transforms.ToTensor())
4     digit_idx = 0
5     print("This digit is a:", mnist[digit_idx][1])
6     plt.imshow(mnist[digit_idx][0].reshape(28,28), vmin=0., vmax=1., cmap='binary')

```

```

0%|          | 0/9912422 [00:00<?, ?it/s]Downloading http://yann.lecun.com/exd
9920512it [00:00, 21683140.56it/s]
Extracting ../data/MNIST/raw/train-images-idx3-ubyte.gz to ../data/MNIST/raw
32768it [00:00, 310988.67it/s]
0it [00:00, ?it/s]Downloading http://yann.lecun.com/exdb/mnist/train-labels-idx1
Extracting ../data/MNIST/raw/train-labels-idx1-ubyte.gz to ../data/MNIST/raw
Downloading http://yann.lecun.com/exdb/mnist/t10k-images-idx3-ubyte.gz to ../dat
1654784it [00:00, 5313160.94it/s]
8192it [00:00, 130169.75it/s]
Extracting ../data/MNIST/raw/t10k-images-idx3-ubyte.gz to ../data/MNIST/raw
Downloading http://yann.lecun.com/exdb/mnist/t10k-labels-idx1-ubyte.gz to ../dat
Extracting ../data/MNIST/raw/t10k-labels-idx1-ubyte.gz to ../data/MNIST/raw
Processing...
Done!
This digit is a: 5

```



## ▼ Defining the neural network

We will be using a [fully connected neural net](#), consisting of three layers. The first layer has  $28 \times 28$  inputs, one for each pixel of the image, and `layer1_size` outputs (and so, `layer1_size` neurons); the second layer has `layer1_size` inputs and `layer2_size` outputs (and thus, neurons), and the final layer has `layer2_size` inputs and 10 outputs, one output for each of the 10 digit classes.

The first two layers of the net consist of [ReLU units](#). A ReLU unit with  $n$  inputs is defined by  $n + 1$  weights  $b, w_1, \dots, w_n$ , and when it is fed inputs  $x_1, \dots, x_n$ , it produces output

$$\max \left( 0, b + \sum_{i=1}^n w_i x_i \right) .$$

The weight  $b$  is called the *bias* of the unit. We will be using not one neuron at a time, but one whole layer of neurons. If the layer contains  $m$  neurons, and we write  $y_j$  for its  $j$ -th output, with  $1 \leq j \leq m$ , we have:

$$y_j = \max \left( 0, b_j + \sum_{i=1}^n w_{ji} x_i \right) ,$$

or denoting by  $b = [b_j]_{1 \leq j \leq m}$  the bias vector, and by  $W = [w_{ji}]_{1 \leq j \leq m, 1 \leq i \leq n}$  the weight matrix, by:

$$y = \max (0, b + Wx) ,$$

where  $x$  is the input vector  $x = [x_1, \dots, x_n]$ , and  $y$  is the output vector  $[y_1, \dots, y_m]$ . We see thus that at the heart of a layer of ReLU units is a matrix multiplication, followed by a max operator.

The matrix multiplication is performed by a PyTorch [linear layer](#), and the maximum operator via the [Pytorch relu function](#).

The last layer of our neural net is a *softmax* layer, as is common in classification problems. Indicate by  $y_0, \dots, y_9$  be the 10 outputs of the 10 linear neurons, computed via  $b + Wx$ . These  $y_0, \dots, y_9$  are general real numbers, not restricted to a particular range. As our aim is to identify which one among the 10 digits is the digit in the input, we prefer the output of the overall layer to be a vector  $[p_0, \dots, p_9]$  of probabilities for each digit  $0, \dots, 9$ , with  $0 \leq p_i \leq 1$  for all  $0 \leq i \leq 9$ , and  $\sum_{i=0}^9 p_i = 1$ . We convert from linear outputs to probabilities using the [softmax](#) transformation:

$$p_i = \frac{e^{y_i}}{\sum_{j=0}^9 e^{y_j}} .$$

It turns out that is numerically more stable to output not these probabilities directly, but their logarithm. We output:

$$q_i = \log p_i = (\log e^{y_i}) - \log \sum_{j=0}^9 e^{y_j} = y_i - \log \sum_{j=0}^9 e^{y_j} .$$

The definition of our net is done in the class `FlatNet` below. We define the three linear layers in the `__init__` method: these are the layers that contain trainable weights. The computation of the output of the net happens in the `forward` method. The input  $x$  is a tensor, that is, a matrix with size

$[28 * 28, \kappa]$ , where  $\kappa$  is the batch size: the net will process the inputs for  $\kappa$  digits at a time. To compute the net output, we apply the transformation for each layer.

```

1 class FlatNet(nn.Module):
2     """This is a 3-layer NN, flat (not convolutional)."""
3
4     def __init__(self, layer1_size=128, layer2_size=128):
5         super().__init__()
6         self.layer1 = nn.Linear(28 * 28, layer1_size)
7         self.layer2 = nn.Linear(layer1_size, layer2_size)
8         self.layer3 = nn.Linear(layer2_size, 10)
9
10    def forward(self, x):
11        """x is the input to the NN. Computes the output."""
12        x = F.relu(self.layer1(x))
13        x = F.relu(self.layer2(x))
14        return F.log_softmax(self.layer3(x), dim=1)

```

## ▼ Training and testing

### Training and testing data

Once the network is defined, we can write two functions, for training and testing the network. The two functions are very similar, except that they work on different portions of the data, and that only the training function actually trains the network.

Separating the dataset into a training and a testing portion may seem wasteful: why not train from all the data, rather than from only a portion of it? The reason for using separate training and testing data is that we want to be able to check for *overtraining*. If a neural network has many degrees of freedom (many trainable weights) compared to how much training data we have, it is possible that the network will learn the particular examples on which we train it, rather than understand the general rules that makes a "2" a "2" and a "7" a "7" digit. We want the network to understand the rules, not to memorize the training examples: our goal is to use the network to classify *new* data that it has not seen yet, new digits written by people in the future, not the same old digits we have in the training set.

To check whether the network is overtraining, we will measure the network classification accuracy on testing data that has not been used to train the network. We split the overall data (actually, MNIST comes pre-split) in two sets: a training set consisting of 80% of the data, and a testing set consisting of 20% of the data.

## ▼ The model and optimizer

The `model` is an instance of our `FlatNet` net.

The `optimizer` is in charge of optimizing the model: precisely, the optimizer decides the size of the learning step, and it applies the update to all the model parameters. The optimizer is created by giving its initializer the parameters of the model, and a learning rate (a base learn step size). There are [many optimizers available in PyTorch](#), and the [SGD optimizer](#) is the simplest of them all. It implements standard gradient descent, with a bit of "[momentum](#)". A more sophisticated optimizer, such as [Adam](#), could be used as a drop-in replacement for the SGD one; you are encouraged to experiment with this. Note that when you change the optimizer, you may have to give a different base learning rate, as each optimizer may interpret this rate in a different way.

```
1 model = FlatNet(layer1_size=128, layer2_size=128).to(device)
2
3 LEARNING_RATE = 0.01
4 MOMENTUM = 0.5
5 optimizer = optim.SGD(model.parameters(), lr=LEARNING_RATE, momentum=MOMENTUM)
```

## ▼ The train and test functions

The `train` and `test` functions correspond closely to the code we used to train our homebrew models. The `train` function iterates over the training data, one batch at a time. The `data.to(device)` and `target.to(device)` expressions move the data (the images to be classified) and the target (the correct digits) to the `device`, which is a GPU (if you wish to run this on the CPU rather than the GPU, change `USE_CUDA` to `False`).

First, we zero the gradient, to start the computation of a new gradient, exactly as we did in our homegrown framework.

Next, the images in the batch are flattened, that is, transformed from a square grid of pixels to a linear vector of  $28 * 28$  pixels, and they are fed to the model, which is an instance of our `FlatNet` class.

We then compute the loss via the [log-likelihood](#) function. The details of why we use this loss are not essential; suffices to say, the log likelihood measures the difference between our output digit probability distribution, and the true distribution, which is a distribution with a 1 for the correct digit and 0 for all other digits.

Finally, the loss gradient is propagated, via the `backward()` method call, and the parameters are updated, in the `optimizer.step()` method call.

```
1 def train(model, device, train_loader, optimizer, epoch, batch_size,
2           flatten=False, log_interval=100):
3     model.train()
4     correct = 0.
```

```

5     for batch_idx, (data, target) in enumerate(train_loader):
6         data, target = data.to(device), target.to(device)
7         optimizer.zero_grad()
8         output = model(data.view(-1, 28 * 28)) if flatten else model(data)
9         loss = F.nll_loss(output, target)
10        loss.backward()
11        optimizer.step()
12        # Get the prediction.
13        pred = output.max(1, keepdim=True)[1]
14        correct += pred.eq(target.view_as(pred)).sum().item()
15        if (batch_idx + 1) % log_interval == 0:
16            print('Train Epoch: {} [{}/{} ({:.0f}%)]\tLoss: {:.6f} \tAccuracy:{:.5f}
17                  epoch, batch_idx * len(data), len(train_loader.dataset),
18                  100. * batch_idx / len(train_loader), loss.item(),
19                  100. * correct / (log_interval * batch_size))
20        correct = 0.

1 def test(model, device, test_loader, flatten=False):
2     model.eval()
3     test_loss = 0
4     correct = 0.
5     with torch.no_grad():
6         for data, target in test_loader:
7             data, target = data.to(device), target.to(device)
8             output = model(data.view(-1, 28 * 28)) if flatten else model(data)
9             test_loss += F.nll_loss(output, target).item() # sum up batch loss
10            pred = output.max(1, keepdim=True)[1] # get the index of the max log-p
11            correct += pred.eq(target.view_as(pred)).sum().item()
12
13    test_loss /= len(test_loader.dataset)
14    print('\nTest set: Average loss: {:.4f}, Accuracy: {}/{} ({:.5f}%)\n'.format(
15          100. * test_loss, correct, len(test_loader.dataset),
16          100. * correct / len(test_loader.dataset)))

```

## ▼ Let's give it a spin

We are finally ready to give our neural network a spin.

### Epochs

The training occurs in *epochs*, where an epoch, in machine-learning parlance, consists in feeding to the network the entire training data. The data is obtained via training and testing loaders: these loaders shuffle the original data so that it is presented in a different order in each epoch, and then return the data in batches of a chosen size.

Shuffling the data at each epoch helps the network train uniformly over the whole dataset.

### Learning in batches



We do not feed the learning examples (the digits) to our neural network one by one: rather, we feed them to the network in batches.

One reason for this is efficiency. The computations in a batch are performed in a GPU, and a GPU is able to perform many operations in parallel. If we fed the data one by one, we would be able to exploit the parallel processing capabilities of a GPU only in part, as the size of a single example would limit the amount of available parallelism. For instance, in a digit, we would be limited to performing  $28 \cdot 28 \cdot l_1$  operations in parallel, where  $l_1$  is the number of layers in the first network. If we batch  $\kappa$  examples together, this number grows by a factor of  $\kappa$ .

The other reason, however, is more fundamental, and it has to do with the very process for machine learning. Ideally, we would like to tune the neural network so as to decrease its loss (its error) not on one example only, but on all training data at once: after all, we want it to recognize all digits.

Denoting by  $L_x$  the loss for a particular input  $x$ , we would like to compute the loss  $L_T = \sum_{x \in T} L_x$  over the whole training set  $T$ . We would then train the model parameters (the weights)  $\theta_1, \dots, \theta_m$  according to the complete gradient

$$\nabla_{\theta} L_T = \left( \frac{\partial L_T}{\partial \theta_1}, \dots, \frac{\partial L_T}{\partial \theta_m} \right)$$

To do this, we can call the `zero_gradient()` method at the beginning of a training epoch. Then, we propagate each training batch through the network, and backpropagate the loss; this accumulates the gradient, and allows the computation of  $g_T = \nabla_{\theta} L_T$ .

However, it turns out it is often better to compute the gradient for each batch, and apply it, rather than add the batch gradients into the total gradient. That is, if  $L_B = \sum_{x \in B} L_x$  is the sum for a batch of training data  $B$ , after processing  $B$ , we train the network with respect to the batch gradient  $g_B = \nabla_{\theta} L_B$  only. The batch gradient  $g_B$  is an approximation of the full gradient  $g_T$ , in the same way in which the average height of 10 students in a class provides an approximation for the average height of all students in the class. Saying this in an equivalent way, we can write  $g_B = g_T + \eta$ , where  $\eta$  is noise. The larger the batch size, the smaller the noise  $\eta$  that enters the gradient at each training step.

The process of using the gradient, estimated from a subset of the data, and thus affected by noise, is known as [stochastic gradient descent](#). This is an accurate yet magnificently grand-sounding name for the simple process of estimating the gradient, and then following the estimate. The next time your read only part of the assigned reading, you can tell the professor you are doing stochastic learning.

It turns out that a small amount of training noise is actually beneficial, for two reasons. First, some amount of noise is often helpful in learning, as it can be used to break symmetries, avoid getting stuck in small local minima, and similar. Second, even the complete gradient  $g_T$  is an approximation for the true (and uncomputable) gradient  $g_U$  computed on the universe  $U$  of possible input data. Training on a noisy version of  $g_T$  helps avoid overtraining to  $T$ .

In practice, as for many things in machine learning, the choice of batch size is guided by experimentation. We use here a batch size of 100; in a 100-input batch, each digit will be represented 10 times on average, so this is a good compromise between an example-specific gradient, and the gradient on all training data.

## Measuring the performance

After each epoch of training, we feed to the classifier the testing data, and we print the resulting accuracy.

```
1 # How big are the batches for training and testing.
2 TRAIN_BATCH_SIZE = 100
3 TEST_BATCH_SIZE = 100
4
5 # Loads the datasets.
6 if colab:
7     train_loader = torch.utils.data.DataLoader(
8         datasets.MNIST('../data', train=True, download=True,
9             transform=transforms.ToTensor()),
10         batch_size=TRAIN_BATCH_SIZE, shuffle=True, **kwargs)
11     test_loader = torch.utils.data.DataLoader(
12         datasets.MNIST('../data', train=False, transform=transforms.ToTensor()),
13         batch_size=TEST_BATCH_SIZE, shuffle=True, **kwargs)

1 NUM_EPOCHS = 2
2 torch.manual_seed(1) # Init random number generator.
3 if colab:
4     for epoch in range(1, NUM_EPOCHS + 1):
5         train(model, device, train_loader, optimizer, epoch, TRAIN_BATCH_SIZE, flat
6         test(model, device, test_loader, flatten=True)
7
```



## ▼ Convolutional neural network models

Once we have created our first neural network, it is a fairly simple matter to replace it with a more sophisticated one. We will see how to replace our "flat" neural network with a [convolutional](#) one.

In our flat neural network, the  $28 \times 28$  image is flattened into an array of  $28 \cdot 28 = 784$  pixels, and these 784 pixels are then fed to the neural network. Thus, all spacial information in the image -- the information on which pixel is close to which other pixel -- is lost. The network will learn the weights for pixels 23, 457, 537, and so on, but will have no idea of where they are in the image! Clearly, the network would be implementing a very different method of image processing than the one we perform in our brains!

[Convolutional networks](#) do a much better job of exploiting spacial information. They examine the image through a sliding window of a given size, in our case,  $3 \times 3$  pixels, and compute intermediate results on the basis of such sliding window. Sliding a  $3 \times 3$  window over a  $28 \times 28$  image with padding of 1 again generates a  $28 \times 28$  grid of results; you can [find here](#) a good visualization of how convolution works, and what padding does. The original pixels of the  $28 \times 28$  image had one value per pixel, the brightness value. The resulting  $28 \times 28$  grid of results has more than one value per cell: it now has as many values as there are neurons in the  $3 \times 3$  convolution. We use multiple neurons so that each of the neurons can learn to look for a specific aspect in the  $3 \times 3$  grid: one neuron might learn to look for vertical lines, one for horizontal lines, and so forth. A good discussion with examples is provided in [Figure 3 of this well-known paper](#).

After the convolutional step, we *pool* the results, computing the maximum in every  $2 \times 2$  square, and reducing the image to a  $14 \times 14$  image, as a consequence. The pooling is used to establish that features are present, while "forgetting" the details of their spacial placement.

We repeat the convolution and pooling a second time, obtaining a  $7 \times 7$  grid of results containing 20 values per cell. These results are then flattened into a  $7 \times 7 \times 20$  flat array, and this flat array is passed through two flat layers of neural nets.

This takes a lot of words to say, but only a few lines of code to do. We use the [Conv2d](#) layer for convolution, and the [MaxPool2d](#) layer for the pooling. Rather than writing `__init__` and `forward` methods, we simply list our layers as input to the `nn.Sequential` model: this is a convenient shorthand for defining a model.

```
1 conv1_neurons = 20
2 conv2_size = 2
3 conv2_neurons = 40
4
5 flat1_in = (28 // 4) * (28 // 4) * conv2_neurons
6 flat1_out = 256
7 flat2_out = 128
```

```

8
9 conv_model = nn.Sequential(
10     nn.Conv2d(1, conv1_neurons, 3, padding=1),
11     nn.ReLU(),
12     nn.MaxPool2d(2),
13     nn.Conv2d(conv1_neurons, conv2_neurons, 3, padding=1),
14     nn.ReLU(),
15     nn.MaxPool2d(2),
16     nn.Flatten(),
17     nn.Linear(flat1_in, flat1_out),
18     nn.ReLU(),
19     nn.Linear(flat1_out, flat2_out),
20     nn.ReLU(),
21     nn.Linear(flat2_out, 10),
22     nn.LogSoftmax(dim=1)
23 )
24 if colab:
25     gpu_conv_model = conv_model.to(device)

```

Each optimizer is defined on the basis of the model parameters, so when we change the model, we need to define a new optimizer too -- it would not work to use the optimizer for the previous linear model to train this convolutional model!

```

1 if colab:
2     conv_optimizer = optim.SGD(gpu_conv_model.parameters(), lr=0.01, momentum=0.5)

```

We can now train our convolutional model. This is not a small net, and it will train much faster on a GPU-enabled machine (if you are on Google Colab, you can select a non-GPU kernel and a GPU kernel to see the difference for yourself).

```

1 if colab:
2     for epoch in range(1, NUM_EPOCHS + 1):
3         train(gpu_conv_model, device, train_loader, conv_optimizer, epoch, TRAIN_BATCH_SIZE)
4         test(gpu_conv_model, device, test_loader)

```



Train Epoch: 1 [9900/60000 (16%) ]	Loss: 2.278747	Accuracy:20.28000
Train Epoch: 1 [19900/60000 (33%) ]	Loss: 2.243984	Accuracy:39.68000
Train Epoch: 1 [29900/60000 (50%) ]	Loss: 1.915723	Accuracy:47.92000
Train Epoch: 1 [39900/60000 (66%) ]	Loss: 0.814999	Accuracy:70.07000
Train Epoch: 1 [49900/60000 (83%) ]	Loss: 0.505966	Accuracy:80.91000
Train Epoch: 1 [59900/60000 (100%) ]	Loss: 0.502425	Accuracy:85.41000

## ▼ Defining new layers

Most ML tutorials at this point go on to more sophisticated neural architectures or more challenging problems; [a number of very good tutorials](#) is available on the PyTorch site, and more are available if you search. We will do something different here: we go back to basics, to the mechanism that makes ML work. In the previous chapter, we have seen how to implement automatic gradient computation for simple expressions; it is now time to see how to do the same in PyTorch: this will make the relationship between our simple autograd setup and PyTorch clear.

The task we will set ourselves to do is to define a new layer, which we call the `sinLayer` layer. The `sinLayer` layer takes as input a vector  $x = [x_1, \dots, x_n]$ , and internally has a set of weights  $[w_1, \dots, w_n]$ . The output of the layer is

$$\sin(wx) = [\sin(w_1 x_1), \dots, \sin(w_n x_n)] .$$

There is no particular reason to think this layer should work especially well. Indeed, we can think of many reasons why it should work poorly, including the fact that normally, the activation functions used in machine learning have uniform behavior when the input goes to infinity, not an oscillatory behavior. But we will give it a try nevertheless: the goal is to learn implementing a new layer, rather than obtaining the ultimate neural network for MNIST.

## ▼ Defining an autograd function

At the core of the `sinLayer` layer is the  $\sin(wx)$  function. We could simply implement the function using the `*` and `torch.sin` functions of PyTorch, but rather than doing that, we will instead define a custom function called `TunedSin`: in this way, we will get a close look at how the automatic gradient computation works in PyTorch.

To implement `TunedSin`, we need to implement a subclass of [torch.autograd.function.Function](#). The class has two static methods: a `forward` method, which computes  $\sin(wx)$  from  $x$  and  $w$ , and a `backward` method, that propagates the gradient.

We wrote the `forward` method already. The method takes three arguments: a context `ctx`, and  $x$  and  $w$ . The context is a place we can use to store values that we wish to have available during the computation of the gradient. We store in the context  $x$ ,  $w$ , and also their product  $xw$ ; of course, the product can be reconstructed from  $x$  and  $w$ , but if we need it, there is no point computing it both in the forward and in the backward pass. By storing it, we are trading a little bit of memory in exchange

of a little bit of efficiency. Aside from storing  $x$ ,  $w$ , and  $xw$  in the context, the forward method simply returns  $\sin(wx)$ .

We leave the interesting bits of the `backward` method for you to write. The method receives as input the context, and the gradient of the loss with respect to its output, that is, indicating with  $L$  the loss and with  $y = \sin(wx)$  the output of the forward pass, it receives

$$\text{grad\_output} = \frac{\partial L}{\partial y} .$$

Corresponding to the two inputs  $x$  and  $w$ , the `backward` method must return the two gradients  $\partial L/\partial x$  and  $\partial L/\partial w$ . These gradients are computed using:

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial}{\partial x} \sin(wx) = \frac{\partial L}{\partial y} w \cos(wx) ,$$

and similarly,

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} x \cos(wx) .$$

The `backward` method first retrieves the values of  $x$ ,  $w$ , and  $xw$  from the context. Then, the method must compute and return  $\partial L/\partial x$  and  $\partial L/\partial w$  as above, and return them. We leave it to you

```

1 from torch.autograd.function import Function
2
3 class TunedSin(Function):
4     """Implements sin(w * x), where w is a tunable weight vector."""
5
6     @staticmethod
7     def forward(ctx, x, w):
8         """
9         @param ctx: context, we can use to store x in it.
10        @param x: the tensor to propagate; it has shape b * n, where b is the
11        batch size, and n is the number of inputs.
12        @param w: weight vector, of the same shape as x.
13        @returns: sin(w * x)
14        """
15        # We form the element-by-element product of x and w.
16        xw = x * w
17        # We save, in case it is useful for backward propagation,
18        # x, w, and the product xw.
19        ctx.save_for_backward(x, w, xw)
20        # We compute the result.
21        return torch.sin(xw)
22
23     @staticmethod
24     def backward(ctx, grad_output):
25         """
26        @param ctx: the context, from which we can retrieve x and w.
27        @param grad_output: the gradient of the output.
28        @returns a pair (gx, gw) consisting of the gradient with

```

```

28     returns a pair (gx, gw), consisting of the gradient with
29     respect to x, and the gradient with respect to w.
30     """
31     # We retrieve x, w, xw from the context.
32     x, w, xw = ctx.saved_tensors
33     # Now we must compute gx and gw.
34     # YOUR CODE HERE
35     gx = grad_output*w*torch.cos(xw)
36     gw = grad_output*x*torch.cos(xw)
37     return gx, gw

```

Let us test this function in the small, so we can check that we have implemented it correctly before moving on. Let us load our classical test library.

```

1 try:
2     from nose.tools import assert_equal, assert_almost_equal
3     from nose.tools import assert_true, assert_false
4     from nose.tools import assert_not_equal
5 except:
6     !pip install nose
7     from nose.tools import assert_equal, assert_almost_equal
8     from nose.tools import assert_true, assert_false
9     from nose.tools import assert_not_equal

```

```

[ ] Collecting nose
    Downloading https://files.pythonhosted.org/packages/15/d8/dd071918c040f50falcfd
    |████████████████████████████████████████| 163kB 3.3MB/s
Installing collected packages: nose
Successfully installed nose-1.3.7

```

Let us define a shorthand function to compare tensors.

```

1 # Let's define a function to compare tensors.
2 def assert_tensors_equal(x, y, places=3):
3     v = torch.sum(torch.abs(x - y)).item()
4     assert_almost_equal(v, 0., places=places)

```

The forward and backward method take a *context* as their first argument. Rather than producing a proper context, we create a *mock* of a context. In computer science, a *mock* is a stand-in for the real thing, that is used in testing. For instance, a database mock is something that (for small amounts of test data) behaves like a database, but is simpler to create, and is suited to be used in tests. Mocks are instrumental in testing complex code, and are easier to create in "duck typing" languages such as Python, where you just need to create an object with the required behavior and methods.

```

1 # We also need to define a context mock for our tests
2 class MockContext(object):

```

```

3
4     def __init__(self):
5         self.saved_tensors = None
6
7     def save_for_backward(self, *tensors):
8         self.saved_tensors = tensors

```

With this, we can finally test our `TunedSin` function.

```

1 # Let's build two tensors x and w of the same size.
2 x = torch.tensor([1., 2.])
3 w = torch.tensor([2., 3.])
4
5 # Let's test the forward method (this should work!).
6 ctx = MockContext()
7 y = TunedSin.forward(ctx, x, w)
8 yy = torch.tensor([ 0.9093, -0.2794])
9 assert_tensors_equal(y, yy)
10
11 # And the backpropagation.
12 grad_output = torch.tensor([1.2, 1.6])
13 gx, gw = TunedSin.backward(ctx, grad_output)
14 assert_tensors_equal(gx, torch.tensor([-0.9988,  4.6088]))
15 assert_tensors_equal(gw, torch.tensor([-0.4994,  3.0725]))
16
17 # Once more.
18 x = torch.tensor([-1., 0.1])
19 w = torch.tensor([0.2, 2.1])
20 y = TunedSin.forward(ctx, x, w)
21 grad_output = torch.tensor([-1., 3.4])
22 gx, gw = TunedSin.backward(ctx, grad_output)
23 assert_tensors_equal(gx, torch.tensor([-0.1960,  6.9831]))
24 assert_tensors_equal(gw, torch.tensor([0.9801, 0.3325]))
25
26 # Oh, let's check that when we backpropagate 0 we get 0.
27 gx, gw = TunedSin.backward(ctx, torch.tensor([0., 0.]))
28 assert_tensors_equal(gx, torch.tensor([0., 0.]))
29 assert_tensors_equal(gw, torch.tensor([0., 0.]))
30

```

## ▼ Defining the `SinLayer` layer

We are now ready to define our new layer. Layers in PyTorch are called modules, and are subclasses of the `nn.Module` class. Models usually have learnable parameters, which are declared as subclasses of the [Parameter class](#). When we will build an optimizer later on, we will call the [parameter\(\). method](#) on the overall module; that method call will collect all instances of



Parameter from all the sub-models of the model (all the layers of the network), and return it as the

```

1 from torch.nn.parameter import Parameter
2
3 class SinLayer(nn.Module):
4
5     def __init__(self, in_features):
6         """Initializes the layer, which is done to process in_features."""
7         super().__init__() # Let's not forget to initialize the Module.
8         self.in_features = in_features
9         # We would expect a size (in_features,) rather than (1, in_features),
10        # but the first 1, is for the batch size. In this way, when we pass
11        # to the layer a batch of size (b, in_features), the weights self.w
12        # will be "broadcast" from shape (1, in_features) to shape
13        # (b, in_features). We fill the weights with zeros now, but we
14        # then overwrite them with random values.
15        self.w = Parameter(torch.zeros(1, in_features))
16        # Let's overwrite it with random values. In a tensor, self.w.data is
17        # the actual data (as a numpy array).
18        self.w.data.normal_()
19
20    def forward(self, x):
21        """Propagates the (batch of) values x."""
22        # We just apply our TunedSin function.
23        return TunedSin.apply(x, self.w)
24
25    def extra_repr(self):
26        return "in_features=out_features=%d" % self.in_features

```

We can now define the overall model. We will define it using the `nn.Sequential` shorthand, as we did for the convolutional networks. You need to build a sequence of layers like this:

- A Linear layer, with as input as many pixels as a MNIST image, and as output,  $n$  outputs.
- A SinLayer layer, with  $n$  inputs and outputs.
- A Linear layer, with  $n$  inputs and outputs.
- A SinLayer layer, with  $n$  inputs and outputs.
- A Linear layer, with  $n$  inputs and 10 outputs.
- A final LogSoftmax layer; this is given for you.

We will test later that the model is correctly built. We will have you experiment with  $n = 64, 128, 256, 512$ . We call  $n$  the `LAYER_SIZE` in the code.

```

1 #@title Choose your layer size
2 LAYER_SIZE = 512 #@param {type:"integer"}

```

**Choose your layer size**

**LAYER\_SIZE: 512**

---

```

1 sin_model = nn.Sequential(
2     # YOUR CODE HERE
3     nn.Linear(784, LAYER_SIZE),
4     SinLayer(LAYER_SIZE),
5     nn.Linear(LAYER_SIZE, LAYER_SIZE),
6     SinLayer(LAYER_SIZE),
7     nn.Linear(LAYER_SIZE, 10),
8     nn.LogSoftmax(dim=1)
9 ).to(device)

```

Let us check that you built the model correctly. If this check does not pass, there is no point in your trying to train the model.

```

1 # Let us check that you built the model correctly.
2 correct_fmt = [
3     "Linear(in_features=784, out_features={sz}, bias=True)",
4     "SinLayer(in_features=out_features={sz})",
5     "Linear(in_features={sz}, out_features={sz}, bias=True)",
6     "SinLayer(in_features=out_features={sz})",
7     "Linear(in_features={sz}, out_features=10, bias=True)",
8     "LogSoftmax()" ]
9 correct_layer_types = [l.format(sz=LAYER_SIZE) for l in correct_fmt]
10 print("Correct layer types:")
11 for c in correct_layer_types:
12     print("\t", c)
13 print("Actual layer types:")
14 for c in sin_model.children():
15     print("\t", c)
16 declared_layer_types = [repr(layer) for layer in sin_model.children()]
17 assert correct_layer_types == declared_layer_types, "Wrong layer types"

```

☞ Correct layer types:

```

Linear(in_features=784, out_features=512, bias=True)
SinLayer(in_features=out_features=512)
Linear(in_features=512, out_features=512, bias=True)
SinLayer(in_features=out_features=512)
Linear(in_features=512, out_features=10, bias=True)
LogSoftmax()

```

Actual layer types:

```

Linear(in_features=784, out_features=512, bias=True)
SinLayer(in_features=out_features=512)
Linear(in_features=512, out_features=512, bias=True)
SinLayer(in_features=out_features=512)
Linear(in_features=512, out_features=10, bias=True)
LogSoftmax()

```

You need to define now an optimizer for your model.

```

1 # Now define an optimizer.
2 sin_optimizer = None # You need to define sin_optimizer in the next line.

```

```

3 # You can define any optimizer you want (any algorithm you want).
4 # YOUR CODE HERE
5 # Learning rate is 0.01, Momentum is 0.5
6 sin_optimizer = optim.SGD(sin_model.parameters(), lr = 0.01, momentum=0.5)

1 # We want to make sure that this is an optimizer.
2 assert isinstance(sin_optimizer, torch.optim.Optimizer)

1 NUM_EPOCHS = 4
2 if colab:
3     torch.manual_seed(1) # Init random number generator.
4     for epoch in range(1, NUM_EPOCHS + 1):
5         train(sin_model, device, train_loader, sin_optimizer, epoch, TRAIN_BATCH_SIZE)
6         test(sin_model, device, test_loader, flatten=True)

```

```

☞ Train Epoch: 1 [9900/60000 (16%)]      Loss: 1.569339  Accuracy:53.75000
Train Epoch: 1 [19900/60000 (33%)]      Loss: 0.862156  Accuracy:74.01000
Train Epoch: 1 [29900/60000 (50%)]      Loss: 0.576350  Accuracy:82.88000
Train Epoch: 1 [39900/60000 (66%)]      Loss: 0.394682  Accuracy:85.61000
Train Epoch: 1 [49900/60000 (83%)]      Loss: 0.342849  Accuracy:87.49000
Train Epoch: 1 [59900/60000 (100%)]     Loss: 0.323784  Accuracy:88.92000

```

Test set: Average loss: 0.3787, Accuracy: 8996.0/10000 (89.96000%)

```

Train Epoch: 2 [9900/60000 (16%)]      Loss: 0.420766  Accuracy:89.84000
Train Epoch: 2 [19900/60000 (33%)]      Loss: 0.275867  Accuracy:90.50000
Train Epoch: 2 [29900/60000 (50%)]      Loss: 0.409535  Accuracy:90.46000
Train Epoch: 2 [39900/60000 (66%)]      Loss: 0.506138  Accuracy:90.87000
Train Epoch: 2 [49900/60000 (83%)]      Loss: 0.221651  Accuracy:91.47000
Train Epoch: 2 [59900/60000 (100%)]     Loss: 0.305738  Accuracy:91.62000

```

Test set: Average loss: 0.2672, Accuracy: 9223.0/10000 (92.23000%)

```

Train Epoch: 3 [9900/60000 (16%)]      Loss: 0.179633  Accuracy:92.25000
Train Epoch: 3 [19900/60000 (33%)]      Loss: 0.178755  Accuracy:92.66000
Train Epoch: 3 [29900/60000 (50%)]      Loss: 0.182310  Accuracy:92.77000
Train Epoch: 3 [39900/60000 (66%)]      Loss: 0.193280  Accuracy:92.83000
Train Epoch: 3 [49900/60000 (83%)]      Loss: 0.207884  Accuracy:93.62000
Train Epoch: 3 [59900/60000 (100%)]     Loss: 0.234229  Accuracy:92.89000

```

Test set: Average loss: 0.2137, Accuracy: 9388.0/10000 (93.88000%)

```

Train Epoch: 4 [9900/60000 (16%)]      Loss: 0.223944  Accuracy:93.61000
Train Epoch: 4 [19900/60000 (33%)]      Loss: 0.164896  Accuracy:94.03000
Train Epoch: 4 [29900/60000 (50%)]      Loss: 0.218553  Accuracy:94.14000
Train Epoch: 4 [39900/60000 (66%)]      Loss: 0.118754  Accuracy:94.20000
Train Epoch: 4 [49900/60000 (83%)]      Loss: 0.189351  Accuracy:94.89000
Train Epoch: 4 [59900/60000 (100%)]     Loss: 0.166804  Accuracy:94.03000

```

Test set: Average loss: 0.1811, Accuracy: 9475.0/10000 (94.75000%)

Assign to `BEST_LAYER_SIZE` the value of `LAYER_SIZE` among 64, 128, 256, 512 that achieves the best performance according to the above training experiments (with the given randomization seed, trained for 4 epochs).

```
1 # YOUR CODE HERE
2 BEST_LAYER_SIZE = LAYER_SIZE
3 #Best Layer Size is 512

1 assert_true(BEST_LAYER_SIZE in {64, 128, 256, 512})
```

Double-click (or enter) to edit

Double-click (or enter) to edit

In aeronautics, there is a saying: if you attach a powerful enough engine, you can make even a brick fly. This is likely a good metaphor for what we did with our `SinLayer`. Nevertheless, it is interesting to note that the `SinLayer`, on MNIST and for the brief amount of training we performed, performs a bit better than the classical ReLU neurons!