

✓ Convolutional Neural Networks

We now apply the MLP to MNIST (handwritten digits). First, we use densely connected networks, as done with non-image data.

Then, we look into using convolutional layers designed for images. Note that because MNIST is an easy data set to classify, the overall performances may be similar.

But there are many benefits of using CNN in images, over densely-connected networks, such as spatial understanding, less parameters, non-diminishing gradients, and others.

```
import numpy as np
import torch
import torch.nn as nn
import torch.optim as optim

# PyTorch TensorBoard support
# from torch.utils.tensorboard import SummaryWriter
# import torchvision
# import torchvision.transforms as transforms

from datetime import datetime

import torchvision
import torchvision.transforms as transforms

from torchvision.datasets import MNIST
import matplotlib.pyplot as plt
%matplotlib inline

from torch.utils.data import random_split
from torch.utils.data import DataLoader
import torch.nn.functional as F

from PIL import Image
# import torchvision.transforms as T

# load the dataset
from torchvision import datasets
mnist_dataset = datasets.FashionMNIST(root = 'data/', download=True, train = True, transform = transforms.ToTensor())
print(mnist_dataset)
```

 Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/train-images-idx3-ubyte.gz>
 Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/train-images-idx3-ubyte.gz> to data/FashionMNIST/raw/train-im-
 100%|██████████| 26421880/26421880 [00:00<00:00, 94011246.06it/s]
 Extracting data/FashionMNIST/raw/train-images-idx3-ubyte.gz to data/FashionMNIST/raw

Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/train-labels-idx1-ubyte.gz>
 Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/train-labels-idx1-ubyte.gz> to data/FashionMNIST/raw/train-lal-
 100%|██████████| 29515/29515 [00:00<00:00, 5672419.47it/s]Extracting data/FashionMNIST/raw/train-labels-idx1-ubyte.gz to data/Fashi

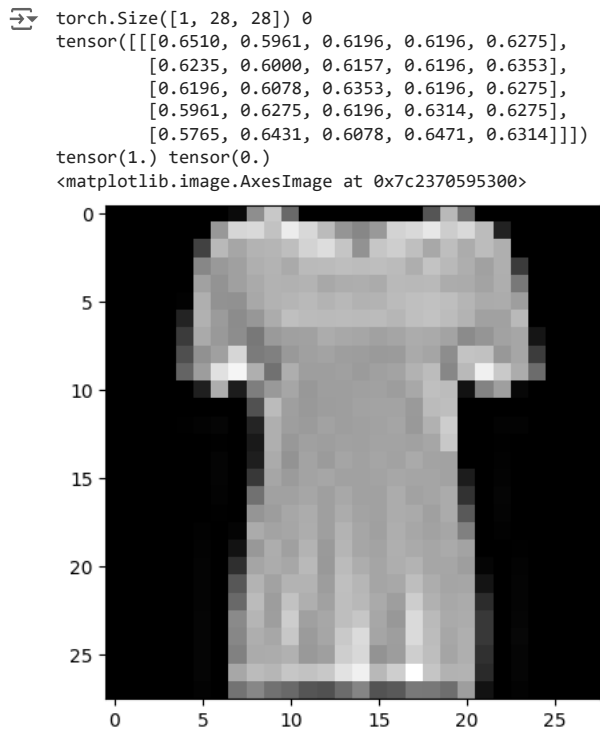
Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/t10k-images-idx3-ubyte.gz>
 Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/t10k-images-idx3-ubyte.gz> to data/FashionMNIST/raw/t10k-imag-
 100%|██████████| 4422102/4422102 [00:00<00:00, 22031922.56it/s]
 Extracting data/FashionMNIST/raw/t10k-images-idx3-ubyte.gz to data/FashionMNIST/raw

Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/t10k-labels-idx1-ubyte.gz>
 Downloading <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/t10k-labels-idx1-ubyte.gz> to data/FashionMNIST/raw/t10k-labe-
 100%|██████████| 5148/5148 [00:00<00:00, 15412046.39it/s]Extracting data/FashionMNIST/raw/t10k-labels-idx1-ubyte.gz to data/FashionM

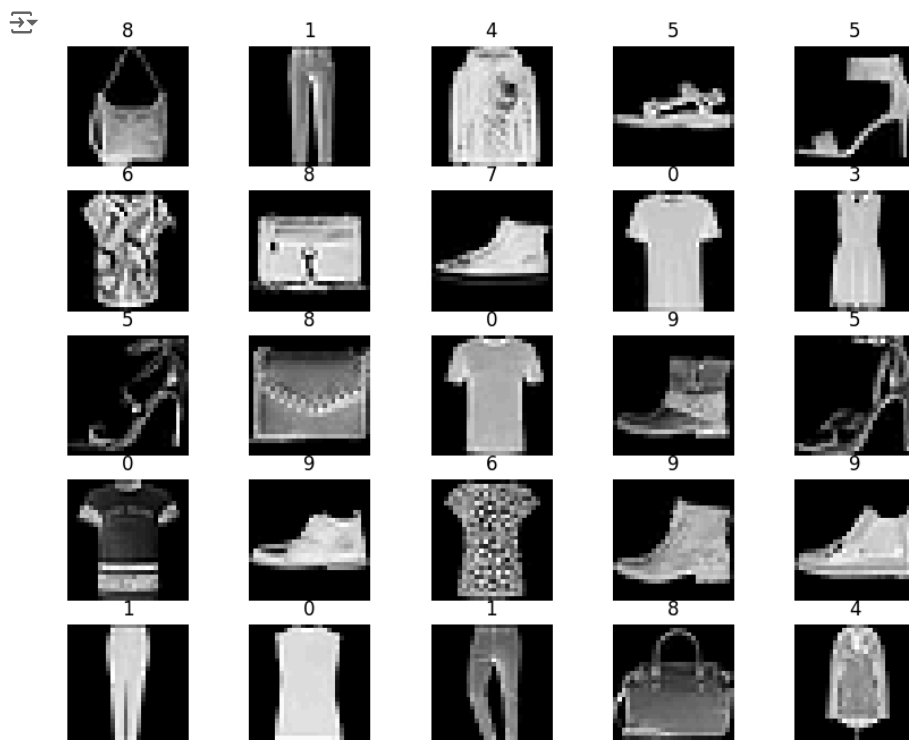
```
Dataset FashionMNIST
  Number of datapoints: 60000
  Root location: data/
  Split: Train
  StandardTransform
  Transform: ToTensor()
```

```
# mnist_dataset has 'images as tensors' so that they can't be displayed directly
sampleTensor, label = mnist_dataset[10]
print(sampleTensor.shape, label)
tpil = transforms.ToPILImage() # using the __call__ to
image = tpil(sampleTensor)
image.show()

# The image is now convert to a 28 X 28 tensor.
# The first dimension is used to keep track of the color channels.
# Since images in the MNIST dataset are grayscale, there's just one channel.
# The values range from 0 to 1, with 0 representing black, 1 white and the values between different shades of grey.
print(sampleTensor[:,10:15,10:15])
print(torch.max(sampleTensor), torch.min(sampleTensor))
plt.imshow(sampleTensor[0,:,:], cmap = 'gray')
```

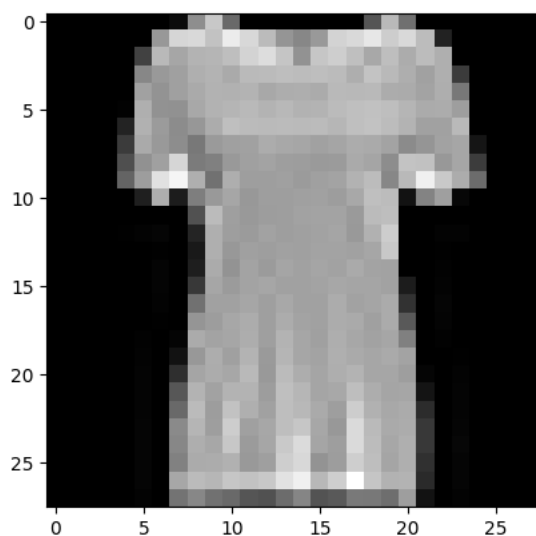


```
# Print multiple images at once
figure = plt.figure(figsize=(10, 8))
cols, rows = 5, 5
for i in range(1, cols * rows + 1):
    sample_idx = torch.randint(len(mnist_dataset), size=(1,)).item()
    img, label = mnist_dataset[sample_idx]
    figure.add_subplot(rows, cols, i)
    plt.title(label)
    plt.axis("off")
    plt.imshow(img.squeeze(), cmap="gray")
plt.show()
```



```
# The image is now convert to a 28 X 28 tensor.
# The first dimension is used to keep track of the color channels.
# Since images in the MNIST dataset are grayscale, there's just one channel.
# The values range from 0 to 1, with 0 representing black, 1 white and the values between different shades of grey.
print(sampleTensor[:,10:15,10:15])
print(torch.max(sampleTensor), torch.min(sampleTensor))
plt.imshow(sampleTensor[0,:,:], cmap = 'gray')
```

```
tensor([[[[0.6510, 0.5961, 0.6196, 0.6196, 0.6275],
          [0.6235, 0.6000, 0.6157, 0.6196, 0.6353],
          [0.6196, 0.6078, 0.6353, 0.6196, 0.6275],
          [0.5961, 0.6275, 0.6196, 0.6314, 0.6275],
          [0.5765, 0.6431, 0.6078, 0.6471, 0.6314]]]])
tensor(1.) tensor(0.)
<matplotlib.image.AxesImage at 0x7c236ddb3790>
```



✓ Training and validation data

While building a ML/DP models, it is common to split the dataset into 3 parts:

- Training set - to train the model, compute the loss and adjust the weights of the model using gradient descent.
- Validation set - to evaluate the training model, adjusting the hyperparameters and pick the best version of the model.

- Test set - to final check the model predictions on the new unseen data to evaluate how well the model is performing.

Quite often, validation and test sets are interchanged (i.e., the validation set is used to final check the model predictions...). Read carefully of the setup.

Following adapted from [Kaggle notebook](#)

```
train_data, validation_data = random_split(mnist_dataset, [50000, 10000])
## Print the length of train and validation datasets
print("length of Train Datasets: ", len(train_data))
print("length of Validation Datasets: ", len(validation_data))
```

```
batch_size = 128
train_loader = DataLoader(train_data, batch_size, shuffle = True)
val_loader = DataLoader(validation_data, batch_size, shuffle = False)
## MNIST data from pytorch already provides held-out test set!
```

```
↗ length of Train Datasets: 50000
length of Validation Datasets: 10000
```

✓ Multi-class Logistic Regression (a building block of DNN)

Since nn.Linear expects the each training example to a vector, each 1 X 28 X 28 image tensor needs to be flattened out into a vector of size 784(28 X 28), before being passed into the model.

The output for each image is vector of size 10, with each element of the vector signifying the probability a particular target label(i.e 0 to 9). The predicted label for an image is simply the one with the highest probability.

```
## Basic set up for a logistic regression model (won't be used in practice or for training)
input_size = 28 * 28
num_classes = 10

# we gradually build on this inherited class from pytorch
model = nn.Linear(input_size, num_classes)
```

We define the class with multiple methods so that we can train, evaluate, and do many other routine tasks with the model.

Particularly, we are looking at multi-class logistic regression (a generalization of one-class logistic regression) using the softmax function (more about this in a few cells down)

```
# Slowly build the model, first with basic
class MnistModel(nn.Module):
    def __init__(self):
        super(MnistModel, self).__init__()
        self.linear = nn.Linear(input_size, num_classes)
        self.conv1 = nn.Conv2d(1, 32, kernel_size=3, stride=1, padding=1)
        self.conv2 = nn.Conv2d(32, 64, kernel_size=3, stride=1, padding=1)
        self.fc1 = nn.Linear(64 * 7 * 7, 128)
        self.fc2 = nn.Linear(128, num_classes)

    def forward(self, xb):
        xb = xb.view(-1, 1, 28, 28)
        xb = F.relu(self.conv1(xb))
        xb = F.max_pool2d(xb, 2, 2)
        xb = F.relu(self.conv2(xb))
        xb = F.max_pool2d(xb, 2, 2)
        xb = xb.view(-1, 64 * 7 * 7)
        xb = F.relu(self.fc1(xb))
        xb = self.fc2(xb)
        return xb

model = MnistModel()
print(model.linear.weight.shape, model.linear.bias.shape)
list(model.parameters())
```



```

[-0.0158, -0.0015, -0.0101, ..., 0.0008, 0.0075, -0.0055],
...,
[ 0.0022, -0.0045, 0.0069, ..., 0.0064, -0.0156, 0.0169],
[ 0.0129, 0.0011, 0.0044, ..., 0.0149, -0.0014, -0.0006],
[ 0.0108, -0.0004, 0.0161, ..., 0.0162, 0.0040, -0.0161]],
requires_grad=True),
Parameter containing:
tensor([-2.1060e-03, -1.2702e-02, -1.3244e-02, -5.6491e-03, -6.3562e-03,
-3.9946e-03, -3.5581e-03, 1.4547e-02, 5.0350e-04, 9.4306e-03,
-1.4695e-02, 5.5750e-05, -1.5179e-02, 6.3846e-04, -1.4080e-02,
-1.2257e-02, 3.5683e-04, -7.5767e-03, 1.0268e-02, -3.1599e-03,
5.0761e-03, 9.6564e-03, 2.2310e-03, 6.1720e-03, -1.0855e-02,
-7.0898e-03, -1.2005e-02, 8.0798e-05, 4.6886e-03, 5.2460e-03,
-4.8769e-03, -9.1819e-03, -1.6731e-02, 1.7213e-04, 8.2179e-03,
-4.8937e-03, 1.0858e-02, 1.3604e-02, -1.1970e-02, 2.2889e-03,
-1.3540e-02, 9.4843e-03, -8.3981e-04, 1.0741e-02, 8.4427e-03,
1.0565e-02, -5.1802e-03, 1.2463e-02, 7.6866e-03, 1.6010e-02,
-7.9018e-04, -1.1453e-02, 4.0929e-04, -1.5332e-02, 9.4306e-03,
-9.4779e-04, -2.1443e-03, -1.4661e-02, 1.3889e-02, 8.0632e-03,
1.6024e-02, 1.1728e-02, -8.8501e-03, -1.0342e-03, 5.1759e-03,
9.1075e-03, -1.2967e-02, 1.9553e-03, 6.7015e-03, -9.6814e-03,
1.0687e-03, 1.2383e-02, 1.7190e-02, -1.4136e-02, -1.3058e-02,
-1.7080e-02, 4.7423e-03, 1.2991e-02, 1.1205e-02, -1.5805e-02,
-6.2784e-04, 1.0194e-03, -7.2938e-03, -1.7486e-02, 1.2272e-02,
1.3408e-02, 1.0831e-02, -1.1701e-03, 1.5427e-02, 7.7649e-03,
-1.1368e-02, 5.6955e-03, 1.3745e-03, 1.6602e-02, -8.0642e-03,
5.9271e-03, 1.4633e-02, 8.9540e-03, 1.1271e-02, -7.3087e-03,
-1.3241e-02, -1.3205e-02, -1.2171e-02, -7.1846e-03, -1.0315e-02,
-6.1198e-04, -1.1278e-02, -2.1239e-03, -1.0396e-02, -2.7074e-03,
-5.1589e-03, -1.5386e-02, 8.8836e-03, 9.2723e-03, 1.0651e-02,
2.9225e-03, -1.9699e-03, -4.1419e-03, -7.3009e-03, 3.0299e-03,
-1.3176e-02, -1.2748e-02, -1.0689e-02, 1.2017e-02, -1.3192e-02,
1.0155e-02, -1.2415e-02, 8.0570e-03], requires_grad=True),
Parameter containing:
tensor([[ 0.0851, 0.0652, 0.0340, ..., 0.0829, 0.0004, -0.0340],
[-0.0841, 0.0163, 0.0538, ..., 0.0671, -0.0283, 0.0261],
[-0.0268, -0.0836, 0.0749, ..., -0.0454, -0.0210, 0.0152],
...,
[-0.0567, -0.0220, -0.0403, ..., -0.0603, 0.0035, 0.0535],
[ 0.0012, -0.0780, 0.0722, ..., 0.0012, -0.0282, 0.0682],
[-0.0546, -0.0737, -0.0824, ..., 0.0641, -0.0454, -0.0519]],
requires_grad=True),
Parameter containing:
tensor([ 0.0240, -0.0164, 0.0058, 0.0495, 0.0820, -0.0305, 0.0130, 0.0554,
0.0321, 0.0831], requires_grad=True)]

```

```
# Always check the dimensions and sample data/image
```

```
for images, labels in train_loader:
```

```
    outputs = model(images)
```

```
    break
```

```
print('Outputs shape: ', outputs.shape) # torch.Size([128, 10])
```

```
print('Sample outputs: \n', outputs[:2].data) # example outputs
```



```
Outputs shape: torch.Size([128, 10])
```

```
Sample outputs:
```

```
tensor([[ 0.0274, 0.0073, -0.0123, -0.0052, -0.0017, -0.1133, 0.0490, 0.0152,
0.0529, 0.0978],
[ 0.0246, 0.0168, -0.0234, -0.0085, -0.0274, -0.1173, 0.0724, 0.0114,
0.0458, 0.0686]])
```

✓ Softmax function

The softmax formula is as follows:

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

Mathematical definition of the softmax function

where all the z_i values are the elements of the input vector and can take any real value. The term on the bottom of the formula is the normalization term which ensures that all the output values of the function will sum to 1, thus constituting a valid probability distribution.

The softmax function is a function that turns a vector of K real values into a vector of K real values that sum to 1. The input values can be positive, negative, zero, or greater than one, but the softmax transforms them into values between 0 and 1, so that they can be interpreted as probabilities. If one of the inputs is small or negative, the softmax turns it into a small probability, and if an input is large, then it turns it into a large probability, but it will always remain between 0 and 1.

The softmax function is sometimes called the softargmax function, or multi-class logistic regression. This is because the softmax is a generalization of logistic regression that can be used for multi-class classification, and its formula is very similar to the sigmoid function which is used for logistic regression. The softmax function can be used in a classifier only when the classes are mutually exclusive.

Many multi-layer neural networks end in a penultimate layer which outputs real-valued scores that are not conveniently scaled and which may be difficult to work with. Here the softmax is very useful because it converts the scores to a normalized probability distribution, which can be displayed to a user or used as input to other systems. For this reason it is usual to append a softmax function as the final layer of the neural network.

Softmax Activation Function

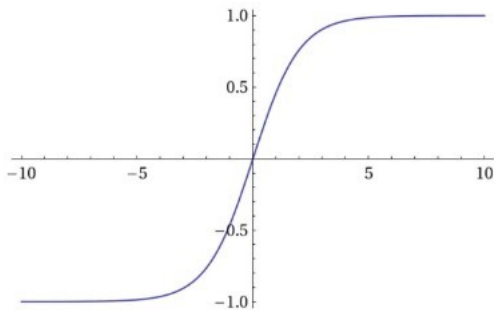


Image from <https://insideaiml.com/blog/SoftMaxActivation-Function-1034>

```
## Apply softmax for each output row
probs = F.softmax(outputs, dim = 1)

## checking at sample probabilities
print("Sample probabilities:\n", probs[:2].data)

# print(preds)
# print("\n")
# print(max_probs)
```

→ Sample probabilities:

```
tensor([[0.1014, 0.0994, 0.0975, 0.0982, 0.0985, 0.0881, 0.1037, 0.1002, 0.1041,
         0.1088],
        [0.1017, 0.1009, 0.0969, 0.0984, 0.0966, 0.0883, 0.1067, 0.1004, 0.1039,
         0.1063]])
```

✓ Evaluation Metric and Loss Function

Here we evaluate our model by finding the percentage of labels that were predicted correctly i.e. the accuracy of the predictions. We can simply find the label with maximum value (before OR after the softmax layer).

NOTE that while accuracy is a great way to evaluate the model, it can't be used as a loss function for optimizing our model using gradient descent, because it does not take into account the actual probabilities predicted by the model, so it can't provide sufficient feedback for incremental improvements.

Due to this reason accuracy is a great evaluation metric (and human-understandable) for classification metric, but not a good loss function. A commonly used loss function for classification problems is the Cross Entropy (implemented directly, no extra coding required).

```
# accuracy calculation
def accuracy(outputs, labels):
    _, preds = torch.max(outputs, dim = 1)
    return(torch.tensor(torch.sum(preds == labels).item() / len(preds)))

print("Accuracy: ", accuracy(outputs, labels))
print("\n")
loss_fn = F.cross_entropy
print("Loss Function: ", loss_fn)
print("\n")
## Loss for the current batch
loss = loss_fn(outputs, labels)
print(loss)
```

→ Accuracy: tensor(0.1172)

Loss Function: <function cross_entropy at 0x7c23d33d5fc0>

```
tensor(2.2925, grad_fn=<NllLossBackward0>)
```

✓ Cross-Entropy

Cross-entropy is commonly used to quantify the difference between two probabilities distribution. Usually the "True" distribution is expressed in terms of a one-hot distribution.

Read more on:

- https://en.wikipedia.org/wiki/Cross_entropy
- <https://machinelearningmastery.com/cross-entropy-for-machine-learning/>
- <https://stackoverflow.com/questions/41990250/what-is-cross-entropy>

```

# We put all of the above:
class MnistModel(nn.Module):
    def __init__(self):
        super().__init__()
        self.linear = nn.Linear(input_size, num_classes)

    def forward(self, xb):
        xb = xb.reshape(-1, 784)
        out = self.linear(xb)
        return(out)

# We add extra methods
def training_step(self, batch):
    # when training, we compute the cross entropy, which help us update weights
    images, labels = batch
    out = self(images) ## Generate predictions
    loss = F.cross_entropy(out, labels) ## Calculate the loss
    return(loss)

def validation_step(self, batch):
    images, labels = batch
    out = self(images) ## Generate predictions
    loss = F.cross_entropy(out, labels) ## Calculate the loss
    # in validation, we want to also look at the accuracy
    # ideally, we would like to save the model when the accuracy is the highest.
    acc = accuracy(out, labels) ## calculate metrics/accuracy
    return({'val_loss':loss, 'val_acc': acc})

def validation_epoch_end(self, outputs):
    # at the end of epoch (after running through all the batches)
    batch_losses = [x['val_loss'] for x in outputs]
    epoch_loss = torch.stack(batch_losses).mean()
    batch_accs = [x['val_acc'] for x in outputs]
    epoch_acc = torch.stack(batch_accs).mean()
    return({'val_loss': epoch_loss.item(), 'val_acc' : epoch_acc.item()})

def epoch_end(self, epoch,result):
    # log epoch, loss, metrics
    print("Epoch [{}], val_loss: {:.4f}, val_acc: {:.4f}".format(epoch, result['val_loss'], result['val_acc']))

# we instantiate the model
model = MnistModel()

# a simple helper function to evaluate
def evaluate(model, data_loader):
    # for batch in data_loader, run validation_step
    outputs = [model.validation_step(batch) for batch in data_loader]
    return(model.validation_epoch_end(outputs))

# actually training
def fit(epochs, lr, model, train_loader, val_loader, opt_func = torch.optim.SGD):
    history = []
    optimizer = opt_func(model.parameters(), lr)
    for epoch in range(epochs):
        ## Training Phase
        for batch in train_loader:
            loss = model.training_step(batch)
            loss.backward() ## backpropagation starts at the loss and goes through all layers to model inputs
            optimizer.step() ## the optimizer iterate over all parameters (tensors); use their stored grad to update their values
            optimizer.zero_grad() ## reset gradients

        ## Validation phase
        result = evaluate(model, val_loader)
        model.epoch_end(epoch, result)
        history.append(result)
    return(history)

# test the functions, with a randomly initialized model (weights are random, e.g., untrained)
result0 = evaluate(model, val_loader)
result0

→ {'val_loss': 2.315847635269165, 'val_acc': 0.09038765728473663}

# let's train for 10 epochs
history1 = fit(10, 0.001, model, train_loader, val_loader)

→ Epoch [0], val_loss: 1.7198, val_acc: 0.6012
Epoch [1], val_loss: 1.4279, val_acc: 0.6617
Epoch [2], val_loss: 1.2590, val_acc: 0.6678
Epoch [3], val_loss: 1.1499, val_acc: 0.6743
Epoch [4], val_loss: 1.0736, val_acc: 0.6821
Epoch [5], val_loss: 1.0168, val_acc: 0.6911

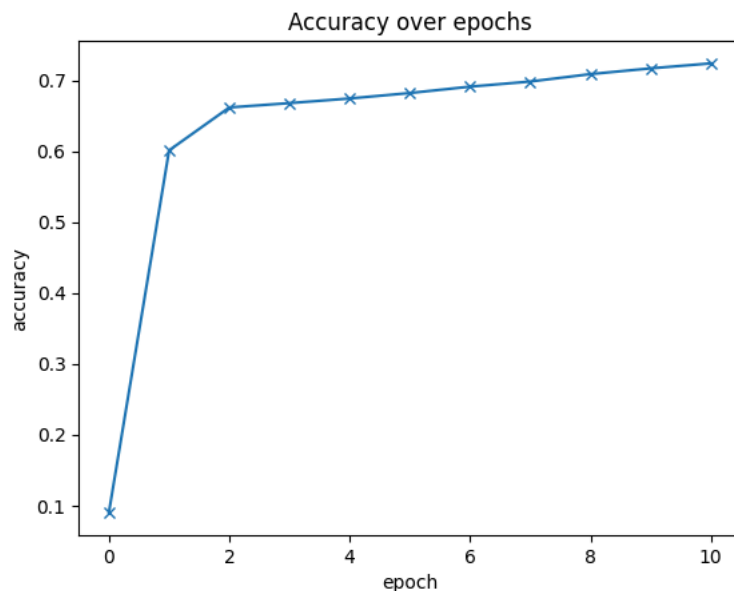
```



```
Epoch [6], val_loss: 0.9724, val_acc: 0.6983
Epoch [7], val_loss: 0.9368, val_acc: 0.7088
Epoch [8], val_loss: 0.9072, val_acc: 0.7169
Epoch [9], val_loss: 0.8821, val_acc: 0.7239
```

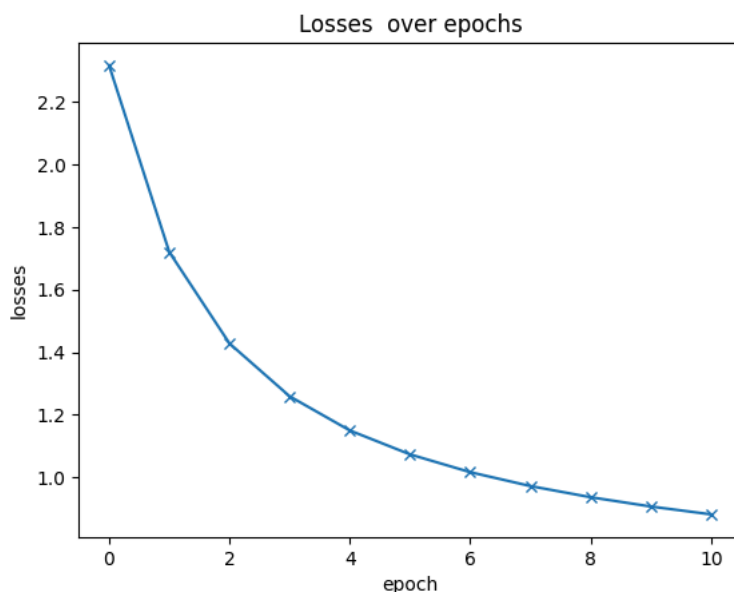
```
# we combine the first result (no training) and the training results of 5 epoches
# plotting accuracy
history = [result0] + history1
accuracies = [result['val_acc'] for result in history]
plt.plot(accuracies, '-x')
plt.xlabel('epoch')
plt.ylabel('accuracy')
plt.title('Accuracy over epochs')
```

↗ Text(0.5, 1.0, 'Accuracy over epochs')



```
# plotting losses
history = [result0] + history1
losses = [result['val_loss'] for result in history]
plt.plot(losses, '-x')
plt.xlabel('epoch')
plt.ylabel('losses')
plt.title('Losses over epochs')
```

↗ Text(0.5, 1.0, 'Losses over epochs')

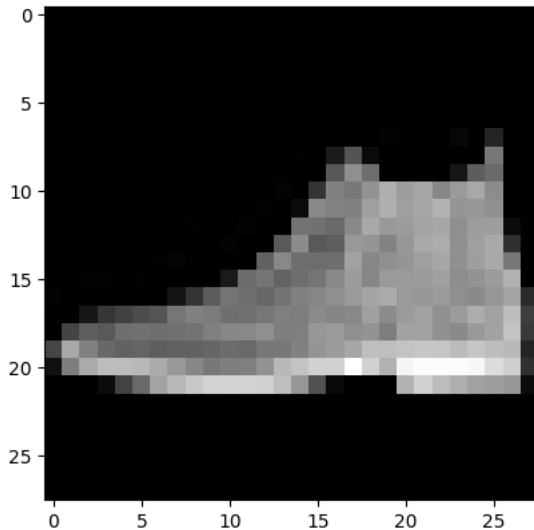


✓ Final check using the (held-out) test dataset.

We will first load the test dataset (from MNIST) and individually check the prediction made by the model. And then, we will put through all images in the test dataset to obtain the final accuracy

```
# Testing with individual images
## Define the test dataset
test_dataset = datasets.FashionMNIST(root = 'data/', train = False, transform = transforms.ToTensor())
print("Length of Test Datasets: ", len(test_dataset))
img, label = test_dataset[0]
plt.imshow(img[0], cmap = 'gray')
print("Shape: ", img.shape)
print('Label: ', label)
```

```
→ Length of Test Datasets: 10000
Shape: torch.Size([1, 28, 28])
Label: 9
```



```
def predict_image(img, model):
    xb = img.unsqueeze(0)
    yb = model(xb)
    _, preds = torch.max(yb, dim = 1)
    return(preds[0].item())
```

```
img, label = test_dataset[0]
print('Label:', label, ', Predicted :', predict_image(img, model))
```

```
→ Label: 9 , Predicted : 9
```

```
# the final check on the test dataset (not used in any training)
test_loader = DataLoader(test_dataset, batch_size = 256, shuffle = False)
result = evaluate(model, test_loader)
result
```

```
→ {'val_loss': 0.903281033039093, 'val_acc': 0.7079101800918579}
```

✓ Convolutional Neural Network (CNN)

So far we treated the MNIST data by flattening each image into a vector. However, there's a lot of information embedded in spatial information. In order to fully 'understand' the image, we need to consider its 2 or more dimensions. Convolutional layers help us in this regard. In most of cases, CNN outperforms densely connected networks and is the most popular architecture for imaging analysis.

CNN is the main force behind revolutionizing the AI or deep learning in the recent decade. Deep neural networks using CNN has shown unprecedented performances when they were first introduced at many competitions (e.g., the ImageNet) by large margins. For imaging analysis, CNN remains the mainstay.

Looking ahead, there are more recent architectures such as the transformer and the diffusion model. We won't be covering them in this course ;)

Convolutional layer is implemented in pytorch as **nn.Conv2d**. As you can see, it is essentially a drop in replacement for **nn.Linear** and other classes.

The explanation for the pytorch class **nn.Conv2d**.

in_channels (int) — Number of channels in the input image, 1 for a grayscale image

out_channels (int) — Number of channels produced by the convolution

kernel_size (int or tuple) — Size of the convolving kernel

stride (int or tuple, optional) — Stride of the convolution. Default: 1

padding (int or tuple, optional) — Zero-padding added to both sides of the input. Default: 0

padding_mode (string, optional) — 'zeros', 'reflect', 'replicate' or 'circular'. Default: 'zeros'

dilation (int or tuple, optional) — Spacing between kernel elements. Default: 1

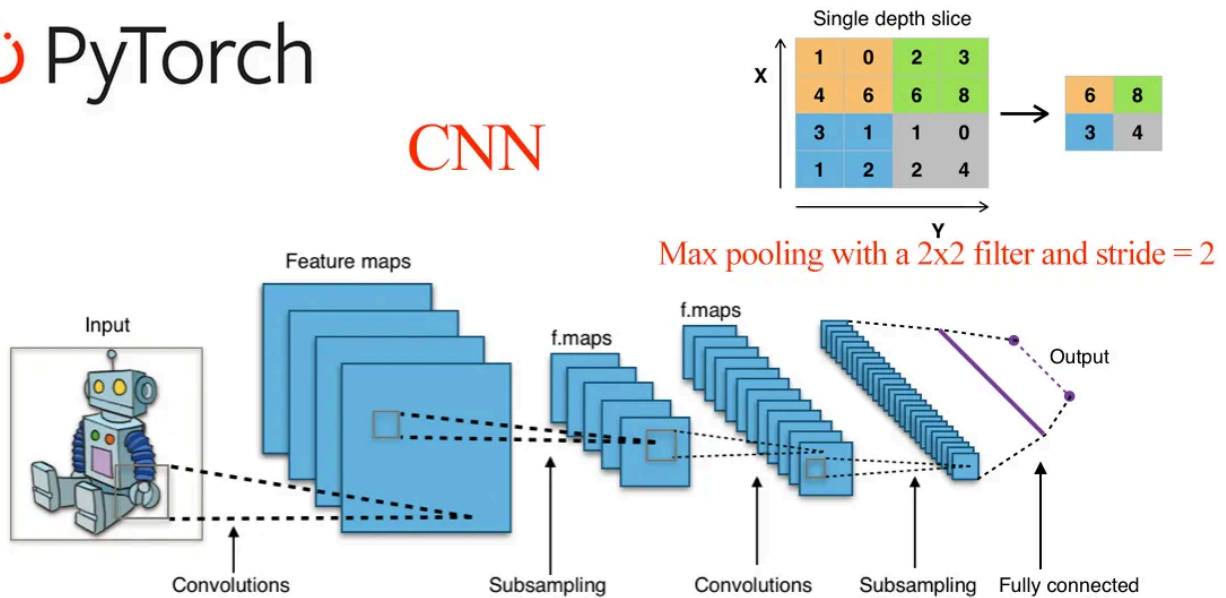
groups (int, optional) — Number of blocked connections from input channels to output channels. Default: 1

bias (bool, optional) — If True, adds a learnable bias to the output. Default: True

Adapted from [@nutanbhogendrasharma](#)

PyTorch

CNN



We construct a fundamental CNN class.

```
class CNN(nn.Module):
    def __init__(self):
        super(CNN, self).__init__()
        self.conv1 = nn.Sequential(
            nn.Conv2d(
                in_channels=1,
                out_channels=16,
                kernel_size=5,
                stride=1,
                padding=2,
            ),
            nn.ReLU(),
            nn.MaxPool2d(kernel_size=2),
        )
        self.conv2 = nn.Sequential(
            nn.Conv2d(16, 32, 5, 1, 2),
            nn.ReLU(),
            nn.MaxPool2d(2),
        )
        # fully connected layer, output 10 classes
        self.out = nn.Linear(32 * 7 * 7, 10)

    def forward(self, x):
        x = self.conv1(x)
        x = self.conv2(x)
        # flatten the output of conv2 to (batch_size, 32 * 7 * 7)
        x = x.view(x.size(0), -1)
        output = self.out(x)
        return output, x  # return x for visualization
```

```
cnn = CNN()
print(cnn)
```

```
CNN(
  (conv1): Sequential(
    (0): Conv2d(1, 16, kernel_size=(5, 5), stride=(1, 1), padding=(2, 2))
    (1): ReLU()
    (2): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
  )
  (conv2): Sequential(
    (0): Conv2d(16, 32, kernel_size=(5, 5), stride=(1, 1), padding=(2, 2))
```

```

    (1): ReLU()
    (2): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)
  )
  (out): Linear(in_features=1568, out_features=10, bias=True)
)

loss_func = nn.CrossEntropyLoss()
loss_func

# unlike earlier example using optim.SGD, we use optim.Adam as the optimizer
# lr(Learning Rate): Rate at which our model updates the weights in the cells each time back-propagation is done.
optimizer = optim.Adam(cnn.parameters(), lr = 0.01)
optimizer

Adam (
  Parameter Group 0
    amsgrad: False
    betas: (0.9, 0.999)
    capturable: False
    differentiable: False
    eps: 1e-08
    foreach: None
    fused: None
    lr: 0.01
    maximize: False
    weight_decay: 0
)

# train_data, validation_data = random_split(mnist_dataset, [50000, 10000])
# ## Print the length of train and validation datasets
# print("length of Train Datasets: ", len(train_data))
# print("length of Validation Datasets: ", len(validation_data))

# batch_size = 128
# train_loader = DataLoader(train_data, batch_size, shuffle = True)
# val_loader = DataLoader(validation_data, batch_size, shuffle = False)
from torch.autograd import Variable

def train(num_epochs, cnn, loaders):
    cnn.train()
    optimizer = optim.Adam(cnn.parameters(), lr = 0.01)
    loss_func = nn.CrossEntropyLoss()
    # Train the model
    total_step = len(loaders)

    for epoch in range(num_epochs):
        for i, (images, labels) in enumerate(loaders):

            # gives batch data, normalize x when iterate train_loader
            b_x = Variable(images) # batch x
            b_y = Variable(labels) # batch y
            output = cnn(b_x)[0]
            loss = loss_func(output, b_y)

            # clear gradients for this training step
            optimizer.zero_grad()

            # backpropagation, compute gradients
            loss.backward()
            # apply gradients
            optimizer.step()

            if (i+1) % 100 == 0:
                print ('Epoch [{}/{}], Step [{}/{}], Loss: {:.4f}'.format(epoch + 1, num_epochs, i + 1, total_step, loss.item()))
                pass
        pass

```

```
# instiate the CNN model
cnn = CNN()
# for testing purpose, we calculate the accuracv of the initial
train(num_epochs=5, cnn=cnn, loaders=train_loader)
```

```
↗ Epoch [1/5], Step [100/391], Loss: 0.3047
Epoch [1/5], Step [200/391], Loss: 0.6164
Epoch [1/5], Step [300/391], Loss: 0.2523
Epoch [2/5], Step [100/391], Loss: 0.3769
Epoch [2/5], Step [200/391], Loss: 0.2451
Epoch [2/5], Step [300/391], Loss: 0.3340
Epoch [3/5], Step [100/391], Loss: 0.2767
Epoch [3/5], Step [200/391], Loss: 0.2540
Epoch [3/5], Step [300/391], Loss: 0.1804
Epoch [4/5], Step [100/391], Loss: 0.2824
Epoch [4/5], Step [200/391], Loss: 0.2473
Epoch [4/5], Step [300/391], Loss: 0.3704
Epoch [5/5], Step [100/391], Loss: 0.1642
Epoch [5/5], Step [200/391], Loss: 0.2867
Epoch [5/5], Step [300/391], Loss: 0.2433
```

✓ Evaluate the model on test data