

```
In [1]: try:
import torch as t
import torch.nn as tnn
except ImportError:
    print("Colab users: pytorch comes preinstalled. Select Change Ru")
    print("Local users: Please install pytorch for your hardware using instr")
    print("ACG users: Please follow instructions here: https://vikasdhiman.i")

    raise
```

```
In [2]: def wget(url, filename):
    """
    Download files using requests package.
    Better than wget command line because this is cross platform.
    """
    try:
        import requests
    except ImportError:
        import subprocess
        subprocess.call("pip install requests".split())
        import requests
    r = requests.get(url)
    with open(filename, 'wb') as fd:
        for chunk in r.iter_content():
            fd.write(chunk)
```

C /

```
In [3]: ## Doing it the Pytorch way without using our custom feature extraction
DEVICE='cuda:0'
DTYPE=t.float32
import torch
import torch.nn
import torch.optim
import torchvision
from torchvision.transforms import ToTensor
from torch.utils.data import DataLoader

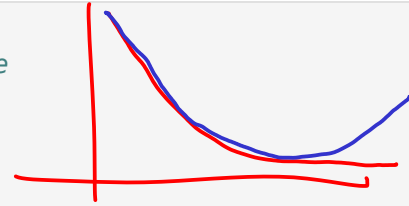
torch.manual_seed(17)
These are special classes that are subclassed from DataLoader/Dataset
# Getting the dataset, the Pytorch way
all_training_data = torchvision.datasets.MNIST(
    root="data", where to store data downloaded data
    train=True, train or test
    download=True, complain or download if the dataset is not in the root
    transform=ToTensor()
)
what to do with data directory
ToTensor() : initializes your data as tensor object
test_data = torchvision.datasets.MNIST(
    root="data",
    train=False, get test data
    download=True,
    transform=ToTensor()
)
```

$28 \times 28$   
=

all the training  
~~80%, 20%~~  $\sqrt{70\%}$   $\sqrt{10\%}$

```
In [4]: training_data, validation_data = torch.utils.data.random_split(all_training
```

```
In [8]: # Hyper parameters
learning_rate = 1e-3 # controls how fast the
batch_size = 64.
epochs = 5
momentum = 0.9
```



avoid overfitting →

```
training_dataloader = DataLoader(training_data, shuffle=True, batch_size=batch_size)
validation_dataloader = DataLoader(validation_data, batch_size=batch_size)
test_dataloader = DataLoader(test_data, batch_size=batch_size)
```

```
loss = torch.nn.CrossEntropyLoss()
```

# TODO:

# Define model = ?

```
class MLPNetwork(torch.nn.Module):
```

```
def __init__(self, hidden_size=10, nclasses=10, input_size=28*28):
    super().__init__()
    self.layers = torch.nn.ModuleList([torch.nn.Flatten(),
        torch.nn.Linear(input_size, hidden_size),
        torch.nn.ReLU(),
        torch.nn.Linear(hidden_size, nclasses)])
```

```
def forward(self, x):
    for l in self.layers:
        xnext = l(x) # call the layers in sequence
        x = xnext
    return x
```

```
model = MLPNetwork()
```

# alternatively you can also

# hidden\_size=10

# nclasses=10

# input\_size=28\*28

```
# model = torch.nn.Sequential(torch.nn.Flatten(),
#     torch.nn.Linear(input_size, hidden_size),
#     torch.nn.ReLU(),
#     torch.nn.Linear(hidden_size, nclasses))
```

# Define optimizer

```
optimizer = torch.optim.SGD(model.parameters(), lr=learning_rate, momentum=
```

```
def loss_and_accuracy(model, loss, validation_dataloader, device=DEVICE):
```

# Validation loop

```
validation_size = len(validation_dataloader.dataset)
```

```
num_batches = len(validation_dataloader)
```

```
test_loss, correct = 0, 0
```

```
with torch.no_grad():
```

```
for X, y in validation_dataloader:
```

```
    X = X.to(device)
```

```
    y = y.to(device)
```

```
    pred = model(X)
```

nn = neural network

self.l1 = torch.nn.Linear

← when the weights are initialize

input

output

x SGD → Stochastic Gradient descent

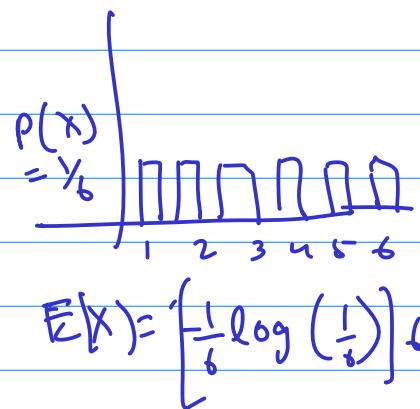
-- setattr --  
 checks for instance of torch.nn.Module

Linear  $W_1, b_1$   
 ↓  
 ReLU  
 ↓  
 Linear  $W_2, b_2$

Entropy (Shannon's): measure of randomness

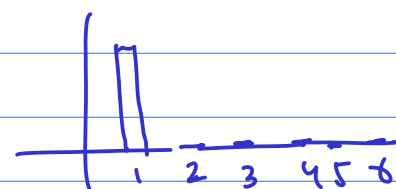
$$E[X] = \sum_{x \in X} -p(x=x) \log_2(p(x=x))$$

bits  $\leftarrow$  bits required to capture information in  $X$



Cross entropy

$$H[Y, \hat{Y}] = \sum_{y \in Y} - \overset{\text{true}}{p(y)} \log \left( \overset{\text{predicted}}{p(\hat{y})} \right)$$



$$E[X] = -1 \log(1).$$

Cross entropy is a measure of information provided by  $Y$  about  $\hat{Y}$  and vice versa

$$\begin{aligned} &+ -0 \log(0) \\ &+ -0 \log(0) \\ &+ \dots \\ &\text{5 times} \\ &= 0 \end{aligned}$$

- is maximum when the probabilities mismatch

- is minimum when the prob. match

$$H[Y, \hat{Y}] = -0 \log 0$$

$$+ \text{ } \boxed{-1 \log(0)}$$

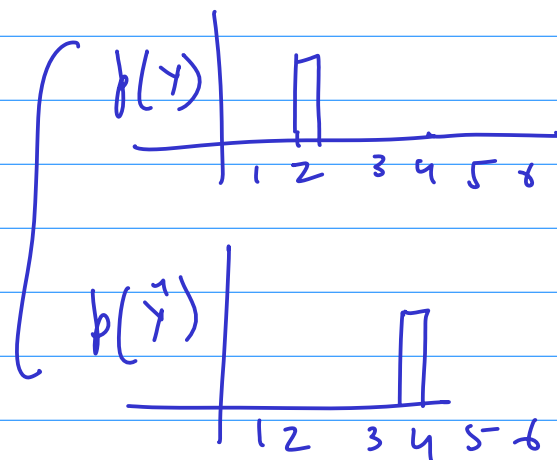
$$+ -0 \log 0$$

$$+ \dots -0 \log 1$$

+

+

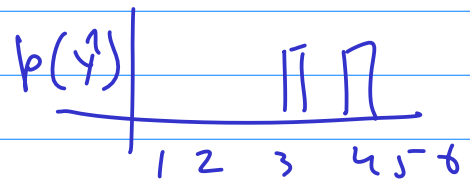
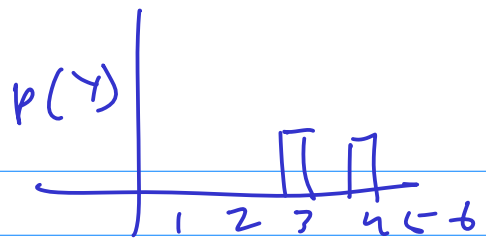
=  $\infty$



$$\log(0^+) \rightarrow -\infty$$

$$H[Y, \hat{Y}] = -0 \log(0) - 1 \log(1) - 0 \log(0) - \dots$$

$$= 0$$



$$H[Y, \hat{Y}] = -0.5 \log(0.5) - 0.5 \log(0.5) = -\log(0.5) \quad \log(1) = 0$$

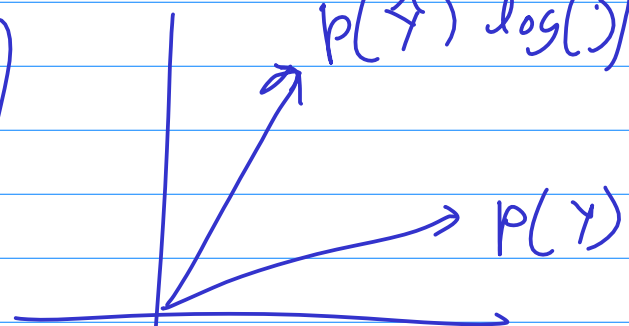
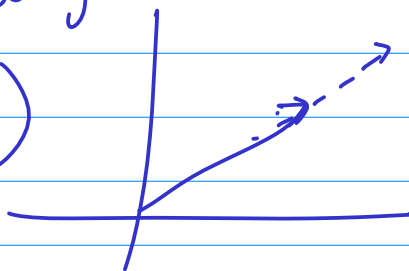
$$\underline{p}(Y) = \begin{bmatrix} p(Y=1) \\ p(Y=2) \\ \vdots \\ p(Y=6) \end{bmatrix}_{6 \times 1}$$

$$\underline{p}(\hat{Y}) = \begin{bmatrix} p(\hat{Y}=1) \\ \vdots \\ p(\hat{Y}=6) \end{bmatrix}_{6 \times 1} \rightarrow \log(p(\hat{Y})) = \begin{bmatrix} \log(p(\hat{Y}=1)) \\ \vdots \\ \log(p(\hat{Y}=6)) \end{bmatrix}$$

$$H[Y, \hat{Y}] = -\underline{p}(Y)^T \underline{\log(p(\hat{Y}))}$$

minimum

when  $p(Y)$  is aligned  
with  $\underline{\log(p(\hat{Y}))}$



$$\vec{a} \cdot \vec{b} \geq |\vec{a}| \cdot |\vec{b}| \cos \theta$$

max when  $\theta = 0$

Classification → Binary: (two classes)  
Hinge loss

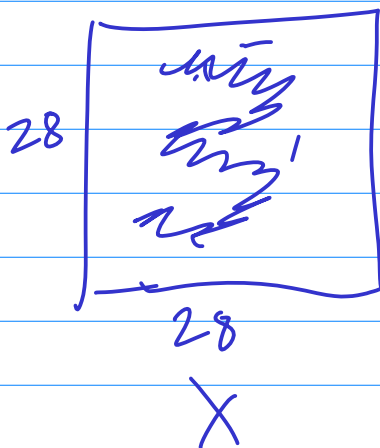
→ Multi-class

Cross entropy loss

MNIST - 10 digits

$\left. \begin{array}{l} Y=0 \\ Y=1 \\ Y=2 \\ \vdots \\ Y=9 \end{array} \right\}$  10 possible classes

for every image



→ label  $Y=3$

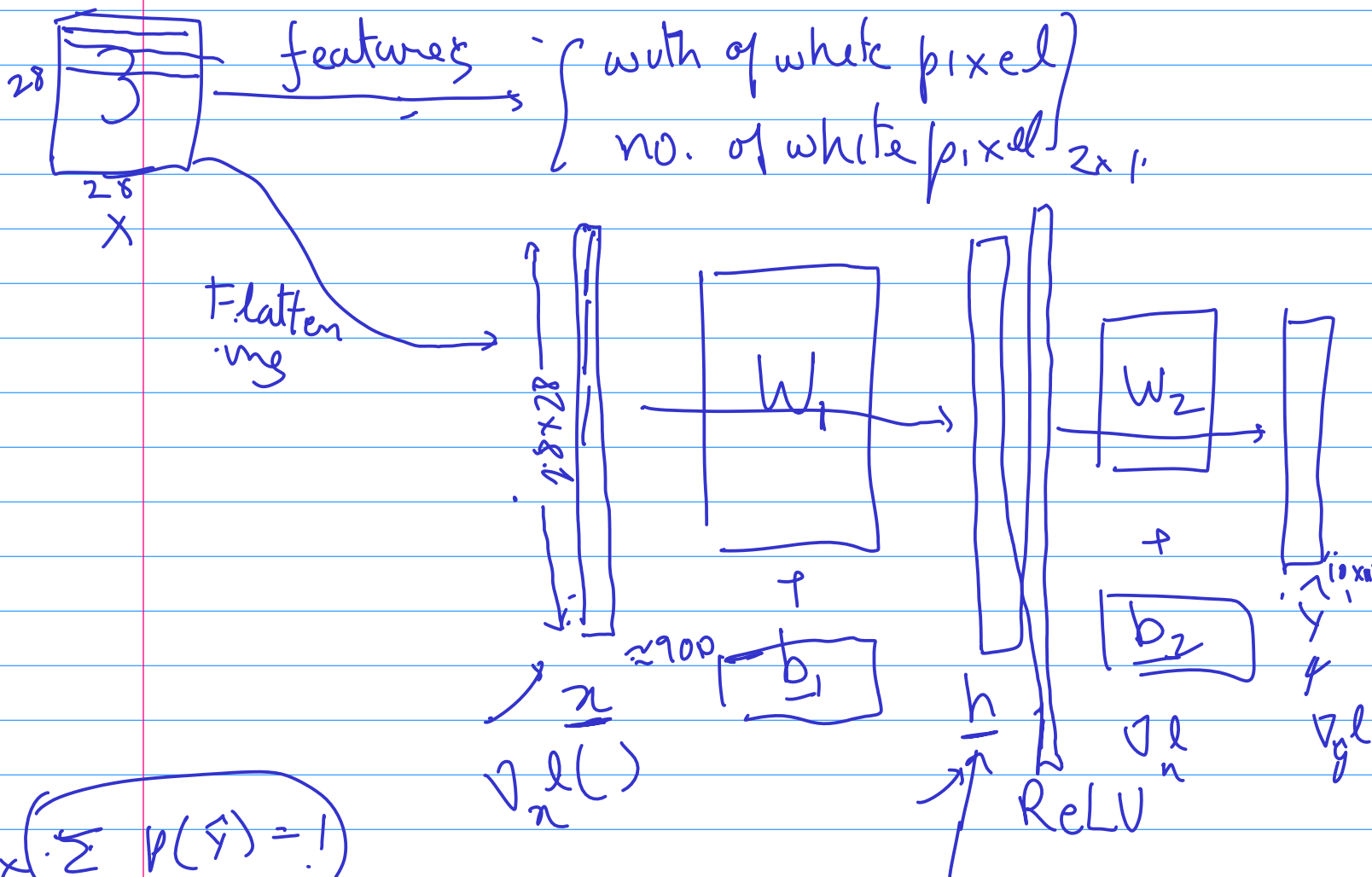
↓ convert label to probability

$\left[ \begin{array}{l} P(Y=0) = 0 \\ \vdots \\ P(Y=3) = 1 \\ \vdots \\ P(Y=9) = 0 \end{array} \right]$  One-hot vector

model

$$f(x; \theta) \rightarrow \hat{y} = f(x; \theta)$$

# MLP : Multi layer perceptron



$$\sum_{y \in Y} P(y) = 1$$

Freedom to choose

$$h \in \mathbb{R}^{5 \times 1}$$

$$P(\hat{y}) = \text{MLP}(x; W_2, b_2, W_1, b_1) = \underbrace{b_2 + W_2}_{\text{Linear}} \left( \underbrace{\text{ReLU}(W_1 x + b_1)}_{h \in \mathbb{R}^{5 \times 1}} \right)$$

Linear

Linear

4 bytes

$$x \in [0, 1]$$

$$P(\hat{y}) \in \mathbb{R}^{10 \times 1}$$

$$\mathbb{R} = \text{float32} =$$

$$\begin{bmatrix} W_1 \in \mathbb{R}^{h \times 784} & b_1 \in \mathbb{R}^{h \times 1} \\ W_2 \in \mathbb{R}^{10 \times h} & b_2 \in \mathbb{R}^{10 \times 1} \end{bmatrix}$$

Total bytes

$$= 5 \times 784 \times 4 + 5 \times 4 + 10 \times 5 \times 4 + 10 \times 4$$

$$\underline{z} \in \underbrace{(-\infty, \infty)}_{\mathbb{R}}^{n \times 1} \xrightarrow{\text{softmax}} \underbrace{(0, 1)}_{\substack{p(\underline{z}) \\ \sum p(\underline{z}) = 1}}^{n \times 1}$$

Softmax

$$p(z_i) = \frac{\exp(z_i)}{\sum_{i=1}^n \exp(z_i)}$$

$$\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{bmatrix} \xrightarrow{\text{softmax}} \underbrace{\begin{bmatrix} \frac{\exp(z_1)}{\sum \exp(z_i)} \\ \frac{\exp(z_2)}{\sum \exp(z_i)} \\ \vdots \\ \vdots \end{bmatrix}}_{p(\underline{z})}$$

$$\exp(\underline{z}) = \begin{bmatrix} \exp(z_1) \\ \vdots \\ \vdots \\ \exp(z_n) \end{bmatrix}$$

$$\boxed{\text{softmax}(\underline{z}) = \exp(\underline{z}) / \exp(\underline{z})^T \mathbb{1}_{n \times 1}}$$

$$\begin{aligned}
 & p(\hat{y}) = \text{Softmax}(\text{MLP}(\underline{x};)) \in \mathbb{R}^{10 \times 1} \\
 & p(y) = \text{one-hot vector} \\
 & \rightarrow \text{Loss}(y, \hat{y}) = - \sum_{y \in Y} p(y) \log(\hat{y}) \\
 & \quad \text{if correct label} = 3 \\
 & \quad p(y=3) = 1 \qquad \qquad \qquad = - \log(p(\hat{y}=3))
 \end{aligned}$$

$$= \frac{-\log(\exp(\hat{y}_3))}{\sum_i \exp(\hat{y}_i)}$$

Gradient descent

$$\mathcal{D}_T = \{(\underline{x}_i, y_i) \dots$$

$$l(y_i, \hat{y}_i) \quad \hat{y}_i = f(\underline{x}_i; \underline{w})$$

$$\underset{\text{1-step}}{L(\mathcal{D}, \underline{w})} = \sum_{i \in \mathcal{D}} l(y_i, \hat{y}_i) \quad \left. \vphantom{\sum_{i \in \mathcal{D}}} \right\} \text{Loss over the entire dataset}$$

$$\underbrace{w_{t+1} = w_t - \alpha \nabla_w \sum_{i \in \mathcal{D}} l(y_i, \hat{y}_i)}_{\text{not feasible for big data sets}} \quad \times$$



$$L \rightarrow w_{t+1} = w_t - \alpha \sum_{i \in D} \left( \nabla_w l(y_i, \hat{y}_i; w_{t,0}) \right) \quad \text{GD}$$

1 epoch  
for  $i$  in  $D$ :

$$w_{t+1,i} = w_{t,i} - \alpha \nabla_w l(y_i, \hat{y}_i; w_{t,i})$$

one sample at a time

Stochastic GD

Batch SGD

batch size

1 epoch

$$\text{for } b \text{ in } \text{range}(D / \text{batch-size})$$

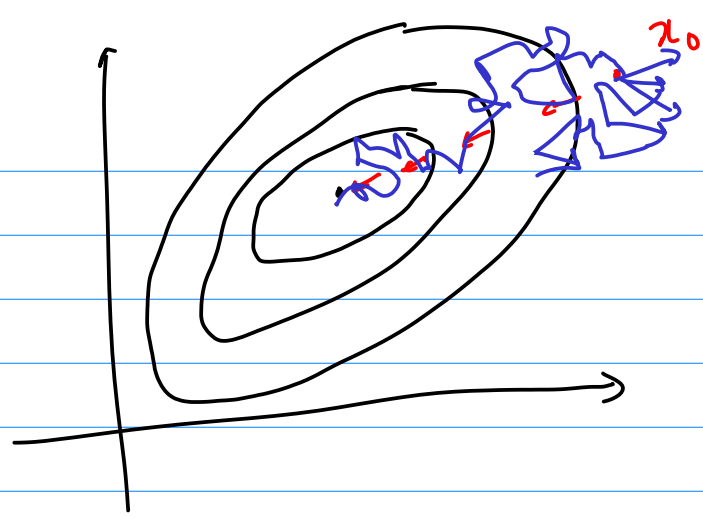
$$w_{t+1,b} = w_{t,b} - \alpha \nabla_w \sum_{\text{batchsize}} l(y_b, \hat{y}_b)$$

one batch at time

SGD with momentum (Next time)

Algo on,

$$L(D, w_t)$$



- GD
- SGD
- Batch SGD

Use previous directions to move towards minimum faster  $\equiv$  SGD with momentum

for  $i$  in range(D):

$$w_{t+1,i} = w_{t,i} - \alpha \nabla_w l(y_i, \hat{y}_i; w_{t,i})$$

SGD with momentum  $\lambda \approx 0.9$   $\lambda = 0.99$

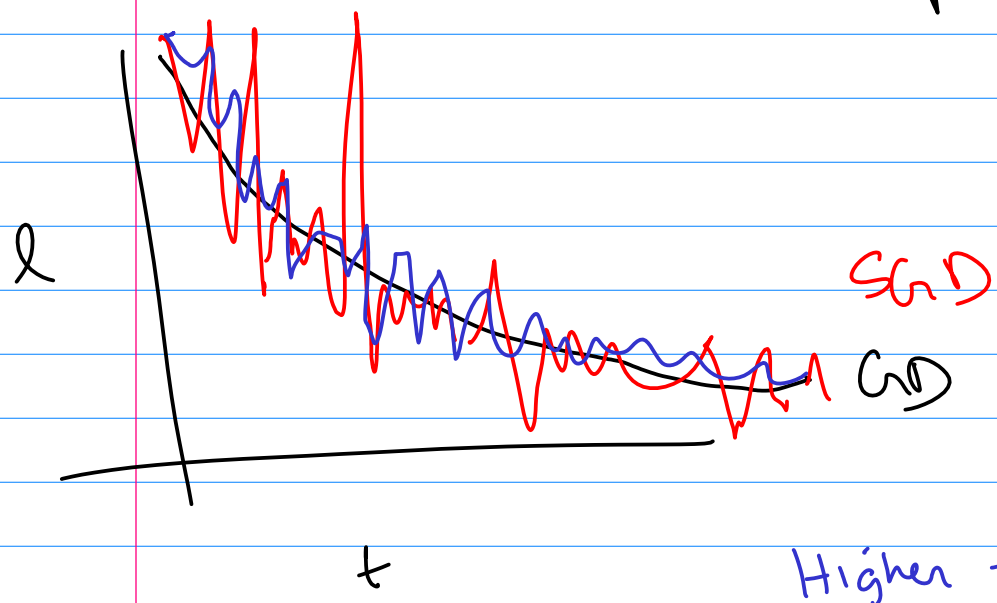
for  $i$  in range(D):

$$v_{t+1,i} = \lambda v_{t,i} + (1 - \lambda) \nabla_w l(y_i, \hat{y}_i; w_{t,i})$$

acceleration

$$w_{t+1,i} = w_{t,i} - \alpha v_{t+1,i}$$

position      velocity



SGD with momentum  $\lambda \in [0, 1]$

Higher the  $\lambda$ , harder it is to change the velocity

$\lambda$  is called the momentum factor

```
test_loss += loss(pred, y).item()
correct += (pred.argmax(dim=-1) == y).type(DTYPE).sum().item()
```

```
test_loss /= num_batches
correct /= validation_size
return test_loss, correct
```

```
def train(model, loss, training_dataloader, validation_dataloader, device=DE
    model.to(device)
    train_losses = []
    valid_losses = []
    for t in range(epochs):
        # Train loop
        training_size = len(training_dataloader.dataset)
        for batch, (X, y) in enumerate(training_dataloader):
            X = X.to(device)
            y = y.to(device)
            # Compute prediction and loss
            pred = model(X)
            loss_t = loss(pred, y)

            # Backpropagation
            optimizer.zero_grad()
            loss_t.backward()
            optimizer.step()

            if batch % 100 == 0:
                loss_t, current = loss_t.item(), (batch + 1) * len(X)
                print(f"loss: {loss_t:>7f} [{current:>5d}/{training_size:>5d}]")
                train_losses.append(loss_t)
                valid_loss, correct = loss_and_accuracy(model, loss, validation_dataloader)
                valid_losses.append(valid_loss)
                print(f"Validation Error: \n Accuracy: {(100*correct):>0.1f}%")
    return model, train_losses, valid_losses
```

64x28x28

one batch at a time

similar to GD

specific to SGD

64x1 → one-hot vector 64x10

SGD update  
 $w_{t+1} = w_t - \alpha \nabla_x$

```
trained_model, train_losses, valid_losses = train(model, loss, training_dataloader, validation_dataloader, device=device)

test_loss, correct = loss_and_accuracy(model, loss, test_dataloader)
print(f"Test Error: \n Accuracy: {(100*correct):>0.1f}%, Avg loss: {test_loss:>7f}")
```

Validation Error: 64/54000]  
Accuracy: 10.9%, Avg loss: 2.325170

Validation Error: 6464/54000]  
Accuracy: 13.6%, Avg loss: 2.183869

Validation Error: 2864/54000]  
Accuracy: 34.1%, Avg loss: 2.001862

Validation Error: 9264/54000]  
Accuracy: 46.8%, Avg loss: 1.787784

Validation Error: 5664/54000]  
Accuracy: 54.8%, Avg loss: 1.580385

Validation Error: 2064/54000]  
Accuracy: 62.2%, Avg loss: 1.389772

Validation Error: 8464/54000]  
Accuracy: 68.5%, Avg loss: 1.216087

Validation Error: 4864/54000]  
Accuracy: 75.8%, Avg loss: 1.061541

Validation Error: 1264/54000]  
Accuracy: 79.9%, Avg loss: 0.935546

Validation Error: 64/54000]  
Accuracy: 80.2%, Avg loss: 0.889401

Validation Error: 6464/54000]  
Accuracy: 82.1%, Avg loss: 0.803695

Validation Error: 2864/54000]  
Accuracy: 82.9%, Avg loss: 0.736784

Validation Error: 9264/54000]  
Accuracy: 83.7%, Avg loss: 0.682686

Validation Error: 5664/54000]  
Accuracy: 84.5%, Avg loss: 0.639954

Validation Error: 2064/54000]  
Accuracy: 85.1%, Avg loss: 0.606284

Validation Error: 8464/54000]  
Accuracy: 85.5%, Avg loss: 0.577947

Validation Error: 4864/54000]  
Accuracy: 86.0%, Avg loss: 0.553612

Validation Error: 1264/54000]  
Accuracy: 86.4%, Avg loss: 0.534344

Validation Error: 64/54000]  
Accuracy: 86.7%, Avg loss: 0.527205

Validation Error: 6464/54000]  
Accuracy: 86.8%, Avg loss: 0.510988

Validation Error: 2864/54000]  
Accuracy: 87.1%, Avg loss: 0.497435

Validation Error: 9264/54000]  
Accuracy: 87.0%, Avg loss: 0.485881

Validation Error: 5664/54000]  
Accuracy: 87.4%, Avg loss: 0.474297

Validation Error: 2064/54000]  
Accuracy: 87.4%, Avg loss: 0.465911

Validation Error: 8464/54000]  
Accuracy: 87.7%, Avg loss: 0.456387

Validation Error: 4864/54000]  
Accuracy: 87.9%, Avg loss: 0.449242

Validation Error: 1264/54000]  
Accuracy: 87.9%, Avg loss: 0.442204

Validation Error: 64/54000]  
Accuracy: 88.2%, Avg loss: 0.439258

Validation Error: 6464/54000]  
Accuracy: 88.1%, Avg loss: 0.433261

Validation Error: 2864/54000]  
Accuracy: 88.0%, Avg loss: 0.428486

Validation Error: 9264/54000]  
Accuracy: 88.3%, Avg loss: 0.423615

Validation Error: 5664/54000]  
Accuracy: 88.5%, Avg loss: 0.420128

Validation Error: 2064/54000]  
Accuracy: 88.5%, Avg loss: 0.414353

Validation Error: 8464/54000]  
Accuracy: 88.4%, Avg loss: 0.409636

Validation Error: 4864/54000]  
Accuracy: 88.7%, Avg loss: 0.407865

Validation Error: 1264/54000]  
Accuracy: 88.8%, Avg loss: 0.401992

Validation Error: 64/54000]  
Accuracy: 88.9%, Avg loss: 0.401415

Validation Error: 6464/54000]

Accuracy: 88.9%, Avg loss: 0.396768

Validation Error: 2864/54000]

Accuracy: 88.9%, Avg loss: 0.395058

Validation Error: 9264/54000]

Accuracy: 89.0%, Avg loss: 0.391086

Validation Error: 5664/54000]

Accuracy: 88.9%, Avg loss: 0.388569

Validation Error: 2064/54000]

Accuracy: 89.0%, Avg loss: 0.385012

Validation Error: 8464/54000]

Accuracy: 89.0%, Avg loss: 0.384191

Validation Error: 4864/54000]

Accuracy: 89.3%, Avg loss: 0.381349

Validation Error: 1264/54000]

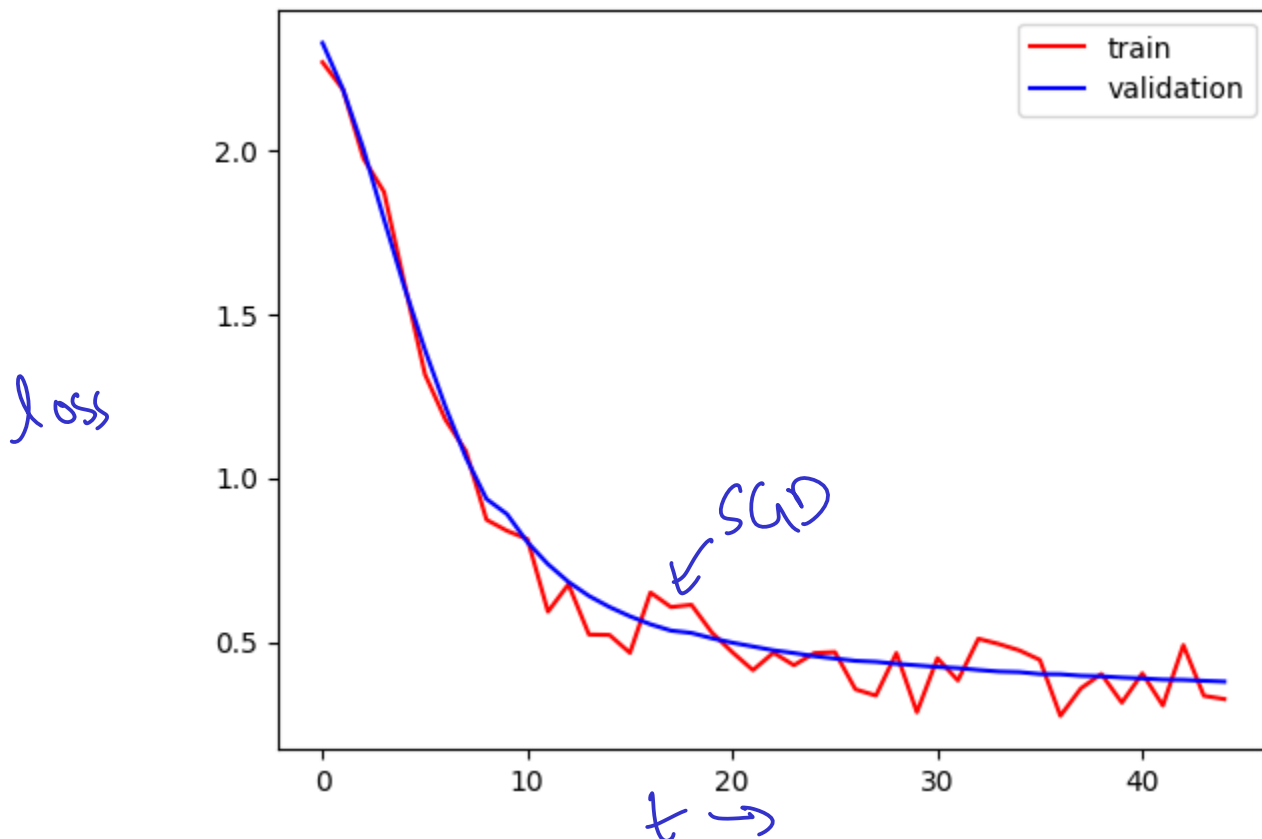
Accuracy: 89.3%, Avg loss: 0.378932

Test Error:

Accuracy: 90.0%, Avg loss: 0.351287

```
In [10]: import matplotlib.pyplot as plt
plt.plot(train_losses, 'r', label='train')
plt.plot(valid_losses, 'b', label='validation')
plt.legend()
```

Out[10]: <matplotlib.legend.Legend at 0x7f92d52b5900>



```
In [11]: torch.nn.__file__
```

```
Out[11]: '/home/vdhiman/.local/share/virtualenvs/nbgrader-notebooks-_16a_jDm/lib/python3.10/site-packages/torch/nn/__init__.py'
```

```
In [12]: X, _ = next(iter(test_dataloader))
         X.shape
```

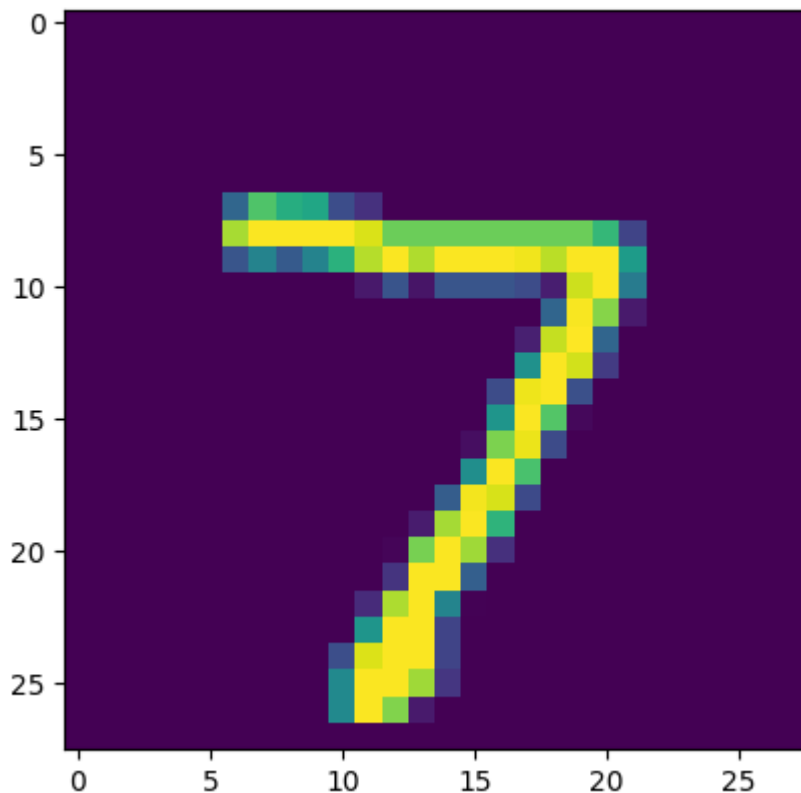
channels

```
Out[12]: torch.Size([64, 1, 28, 28])
```

Batch

```
In [13]: import matplotlib.pyplot as plt
         plt.imshow(X[0, 0])
```

```
Out[13]: <matplotlib.image.AxesImage at 0x7f92d02888b0>
```



```
In [14]: print("The predicted image label is ", model(X.to(DEVICE)).argmax(dim=-1)[0])
```

The predicted image label is 7

In [ ]:

