

CSCI 8960: Privacy-Preserving Data Analysis

PPML Summer School

Jaewoo Lee
jaewoo.lee@uga.edu

July 11, 2023

Department of Computer Science



UNIVERSITY OF
GEORGIA

Introduction

An Era of Deep Learning

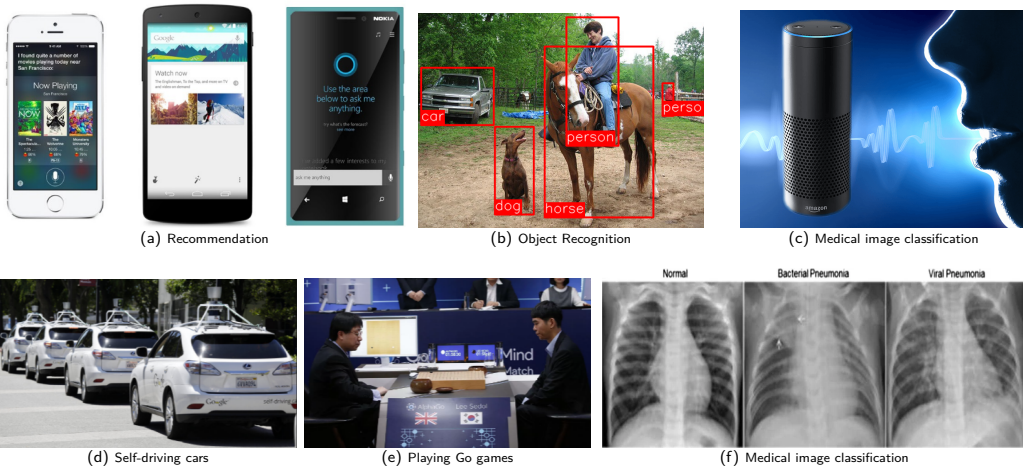
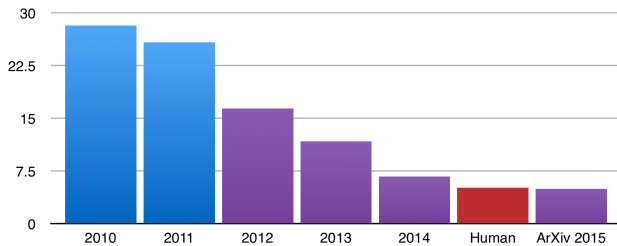


Fig. 1. Applications of deep learning

ILSVRC top-5 error on ImageNet

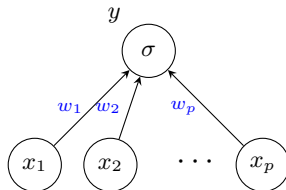


- Logistic regression

$$y = \phi(\mathbf{w}^T \mathbf{x} + w_0) = \phi\left(\sum_{i=1}^p w_i x_i + w_0\right),$$

where

- ▶ $\mathbf{x} = (x_1, x_2, \dots, x_p)$ feature vector
- ▶ $\mathbf{w} = (w_1, w_2, \dots, w_p)$ weight vector (or parameter vector)
- ▶ w_0 bias (or intercept) term



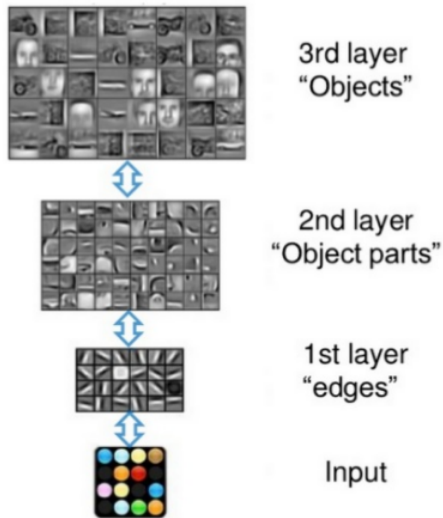
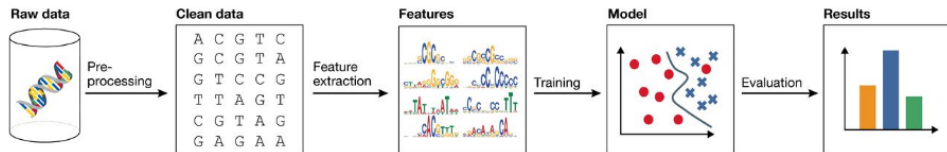


Fig. 3. Multiple layers of information processing



- 1 **Dataset:** an *i.i.d.* sample drawn from the population
 - $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
- 2 **Task:** predict Y given X
 - regression or classification?
- 3 **Model:** assumption on the (distribution) of data
 - linear, non-linear, parametric, non-parametric
- 4 **Loss**(a.k.a. score/cost): measure of how good a model is
 - $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$
- 5 **Training:** normally involve an optimization
 - (stochastic) gradient descent, Newton's method
- 6 **Evaluation:** validate the learning outcome

MLP

Multilayer Perceptrons

- Multiple *layers* of units
- Connecting units in one layer to those in the subsequent layer
- Feed-forward neural network
- No cycle
- All input units are connected to all output units (fully connected layer).

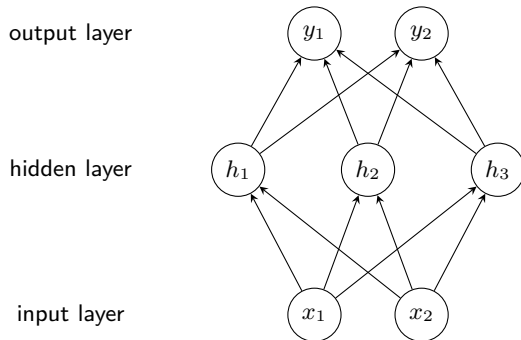


Fig. 4. Simple MLP

Use machine learning algorithms as a black box!

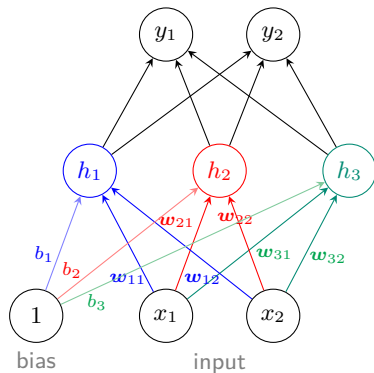
- Pre-activation

- ▶ $z_1 = w_{11}x_1 + w_{12}x_2$
- ▶ $z_2 = w_{21}x_1 + w_{22}x_2$
- ▶ $z_3 = w_{31}x_1 + w_{32}x_2$

- Matrix notation

$$\begin{pmatrix} z_1 \\ z_2 \\ z_3 \end{pmatrix} = \begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix}$$

- ▶ $\mathbf{h} = \phi(\underbrace{\mathbf{W}\mathbf{x} + \mathbf{b}}_{\text{pre-activation}})$



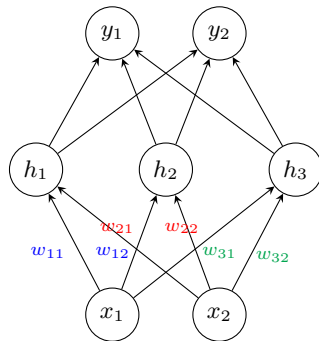
- $h_i = \phi(\mathbf{w}_i^T \mathbf{x} + b)$
- Matrix notation

$$\mathbf{z} = \underbrace{\mathbf{W}\mathbf{x} + \mathbf{b}}_{\text{pre-activation}}$$

$$\mathbf{h} = (h_1, h_2, h_3) = \phi(\mathbf{z}),$$

where $\mathbf{W} \in \mathbb{R}^{m \times n}$.

- ▶ \mathbf{W} : weight matrix
- ▶ \mathbf{b} : bias term
- ▶ \mathbf{z} : pre-activation
- ▶ ϕ : activation function



What if we have two hidden layers?

- $h_i = \phi(\mathbf{w}_i^T \mathbf{x} + b)$

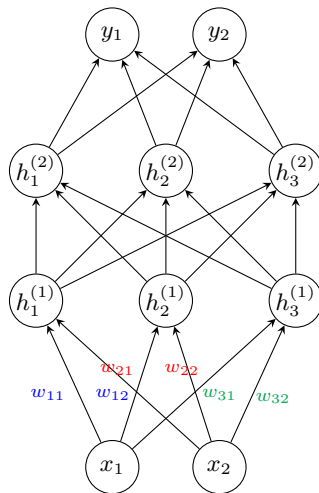
- Matrix notation

$$\mathbf{h}^{(1)} = \phi(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)})$$

$$\mathbf{h}^{(2)} = \phi(\mathbf{W}^{(2)} \mathbf{h}^{(1)} + \mathbf{b}^{(2)})$$

- Hierarchy

$$\mathbf{h}^{(2)} = \phi(\mathbf{W}^{(2)} \phi(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$$



Activation functions

- Determines what (or whether) to pass on from each neuron to the next layer
- Creates **Non-linearity** in the network

Activation function	Equation	Example	1D Graph
Unit step (Heaviside)	$\phi(z) = \begin{cases} 0, & z < 0, \\ 0.5, & z = 0, \\ 1, & z > 0, \end{cases}$	Perceptron variant	
Sign (Signum)	$\phi(z) = \begin{cases} -1, & z < 0, \\ 0, & z = 0, \\ 1, & z > 0, \end{cases}$	Perceptron variant	
Linear	$\phi(z) = z$	Adaline, linear regression	
Piece-wise linear	$\phi(z) = \begin{cases} 1, & z \geq \frac{1}{2}, \\ z + \frac{1}{2}, & -\frac{1}{2} < z < \frac{1}{2}, \\ 0, & z \leq -\frac{1}{2}, \end{cases}$	Support vector machine	
Logistic (sigmoid)	$\phi(z) = \frac{1}{1 + e^{-z}}$	Logistic regression, Multi-layer NN	
Hyperbolic tangent	$\phi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$	Multi-layer NN	

Is deep learning difficult?

- Deep models:

$$\mathbf{y} = \phi(\mathbf{W}^{(2)} \phi(\mathbf{W}^{(1)} \mathbf{x}))$$

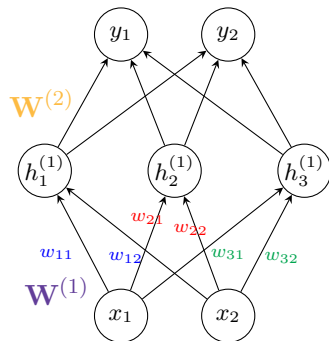
- Classical Models:

$$\mathbf{y} = \phi(\mathbf{W}\mathbf{x})$$

- Loss: $L(W) = \frac{1}{2} \|\hat{Y}(W, X) - Y\|_F^2$

- Simplistic case: $\phi(x) = x$.

▶ $W^{(1)} = w_1, W^{(2)} = w_2, X = Y = 1$



$$L_{\text{deep}}(W) = (w_1 w_2 - 1)^2$$

VS

$$L_{\text{lin}}(W) = (w_1 - 1)^2 + (w_2 - 1)^2$$

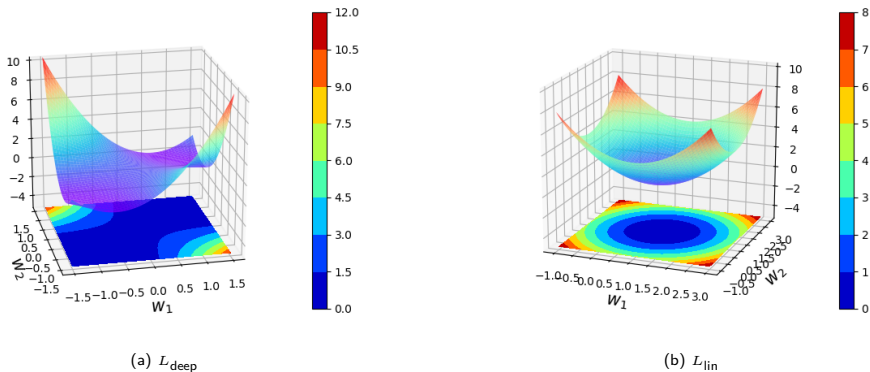
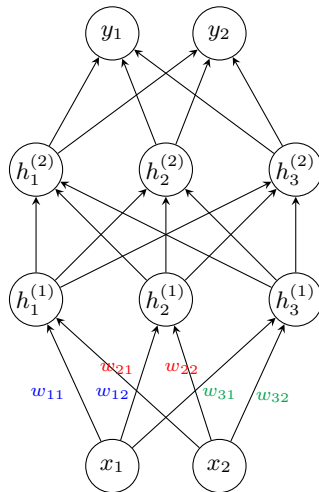


Fig. 5. Landscape of loss functions $L(W)$



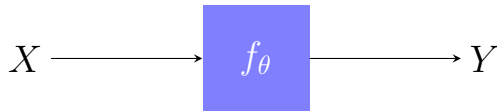
The depth of NN creates non-convex loss surfaces.

- What functions can be expressed by my model?
 - ▶ A NN can be viewed as a function (or a mapping).
 - ▶ $f_{\theta} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$
 - ▶ $f_{\theta} : \mathbf{x} \mapsto \mathbf{y}$
- Model parameter
 - ▶ $\theta = \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(2)}, \mathbf{W}^{(3)}, \mathbf{b}^{(3)}\}$





- We want learn a map $f : \mathcal{X} \rightarrow \mathcal{Y}$.



- ▶ input $\{(x_i, y_i)\}$, $y_i = f$ evaluated at x_i
 - ▶ output $\{\hat{y}_i = f_\theta(x_i)\}$
 - ▶ parameter θ : weight vectors
- Training: optimize θ to approximate the mapping f

What to optimize?

- Define a distribution over a variable to predict $p(y|x; \theta)$
- Maximum Likelihood Principle

$$L(\theta) = \mathbb{E}_{x, y \sim p_{\text{data}}} [\log p_{\text{model}}(y | x)]$$

- Example: when $p_{\text{model}}(y | x) = \mathcal{N}(y; f(x; \theta), \mathbf{I})$, MLE reduces to

$$L(\theta) = \frac{1}{2} \mathbb{E}_{x, y \sim p_{\text{data}}} [\|y - f(x; \theta)\|^2].$$

x_1	x_2	y
0.1	0.5	1
0.2	0.7	0
0.9	0.1	0

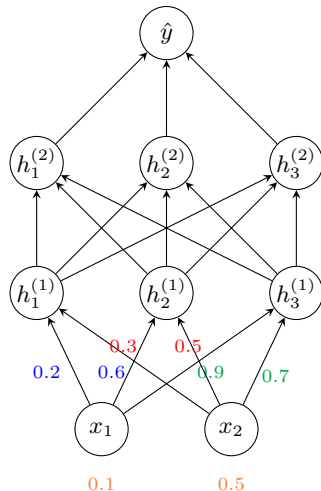
Table 1. Input data

- Assume $\phi(z) = \max(0, z)$ (ReLU activation function)

$$z_1^{(1)} = 0.2 \cdot 0.1 + 0.6 \cdot 0.5 = 0.32$$

$$z_2^{(1)} = 0.3 \cdot 0.1 + 0.5 \cdot 0.5 = 0.28$$

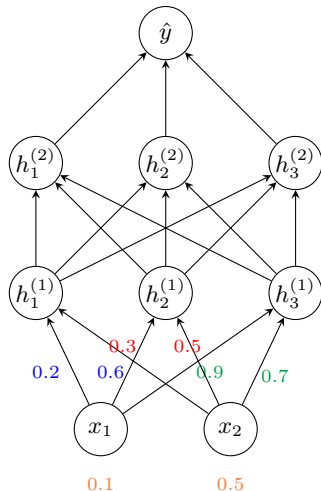
$$z_3^{(1)} = 0.9 \cdot 0.1 + 0.7 \cdot 0.5 = 0.44$$



x_1	x_2	\hat{y}	y
0.1	0.5	0.8	1
0.2	0.7	0.3	0
0.9	0.1	0.1	0

Table 2. After the forward step

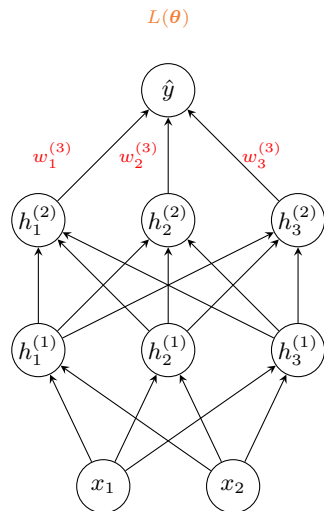
- For each $\mathbf{x} \in D$, compute the *loss*
- Average the losses



$$L(\theta) = \frac{1}{3} \left\{ (0.8 - 1)^2 + (0.3 - 0)^2 + (0.1 - 0)^2 \right\}$$

- Backpropagating *error correction* info.
 - ▶ How should we modify $\theta_l = \{\mathbf{w}^{(l)}, \mathbf{b}^{(l)}\}$?
 - ▶ It is given by the gradient: $\frac{\partial L(\theta)}{\partial \theta_l}$.
 - ▶ Let's start with $\frac{\partial L(\theta)}{\partial \theta_3}$.

$$\hat{y} = (\mathbf{w}^{(3)})^\top \mathbf{h}^{(2)} = w_1^{(3)} h_1^{(2)} + w_2^{(3)} h_2^{(2)} + w_3^{(3)} h_3^{(2)}$$

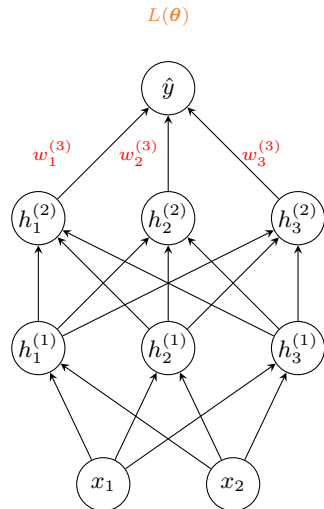


- Recall that $L(\theta) = \frac{1}{2}(y - \hat{y})^2$.

$$\hat{y} = (\mathbf{w}^{(3)})^\top \mathbf{h}^{(2)} = w_1^{(3)}h_1^{(2)} + w_2^{(3)}h_2^{(2)} + w_3^{(3)}h_3^{(2)}$$

- By the chain rule, we have

$$\frac{\partial L}{\partial \theta_3} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \theta_3}.$$





- How about $\frac{\partial L}{\partial \theta_2}$?

▶ Again, by the chain rule, we have

$$\frac{\partial L}{\partial \theta_3} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \theta_3}$$

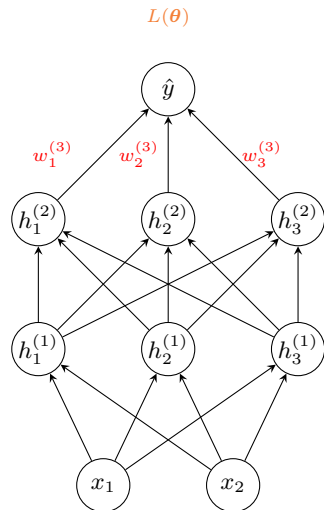
$$\frac{\partial L}{\partial \theta_2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}^{(2)}} \frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{z}^{(2)}} \frac{\partial \mathbf{z}^{(2)}}{\partial \theta_2} .$$

From the forward phase, we have

$$\hat{y} = (\mathbf{w}^{(3)})^\top \mathbf{h}^{(2)}$$

$$\mathbf{h}^{(2)} = \phi(\mathbf{z}^{(2)})$$

$$\mathbf{z}^{(2)} = \mathbf{W}^{(2)} \mathbf{h}^{(1)} .$$



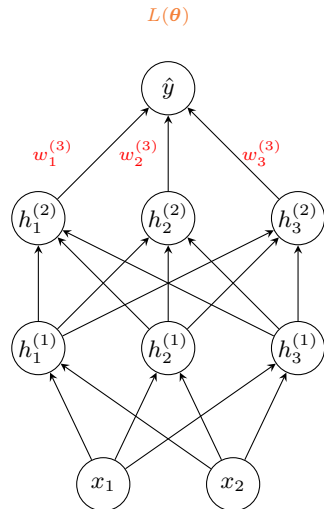
- Now we compute $\frac{\partial L}{\partial \theta_1}$.

► By the chain rule, we have

$$\frac{\partial L}{\partial \theta_3} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \theta_3}$$

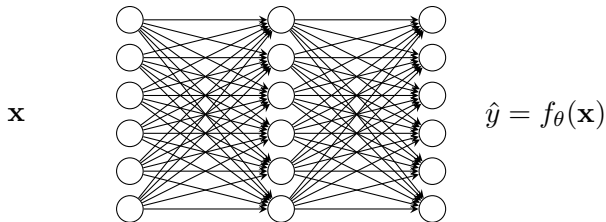
$$\frac{\partial L}{\partial \theta_2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}^{(2)}} \underbrace{\frac{\partial \mathbf{h}^{(2)}}{\partial \mathbf{z}^{(2)}} \frac{\partial \mathbf{z}^{(2)}}{\partial \theta_2}}_{\frac{\partial L}{\partial \mathbf{z}^{(2)}}}$$

$$\begin{aligned} \frac{\partial L}{\partial \theta_1} &= \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{z}_2} \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} \frac{\partial \mathbf{z}_1}{\partial \theta_1} \\ &= \frac{\partial L}{\partial \mathbf{z}_2} \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} \frac{\partial \mathbf{z}_1}{\partial \theta_1} \end{aligned}$$



CNN

- MLPs are *universal* function approximators.
 - ▶ can learn a complex mapping $f : \mathcal{X} \rightarrow \mathcal{Y}$ between X and Y



- Given input feature *vector* \mathbf{x} , f_{θ} returns desired output.
 - ▶ $f_{\theta} : \mathbf{x} \mapsto \hat{\mathbf{y}}$
 - ▶ *fully connected*
 - ▶ Forward step : matrix multiplication
 - $\mathbf{h}^{(1)} = \phi(\mathbf{W}_1^T \mathbf{x})$
 - $\mathbf{h}^{(2)} = \phi(\mathbf{W}_2^T \mathbf{h}^{(1)})$

MLP can be too complex!

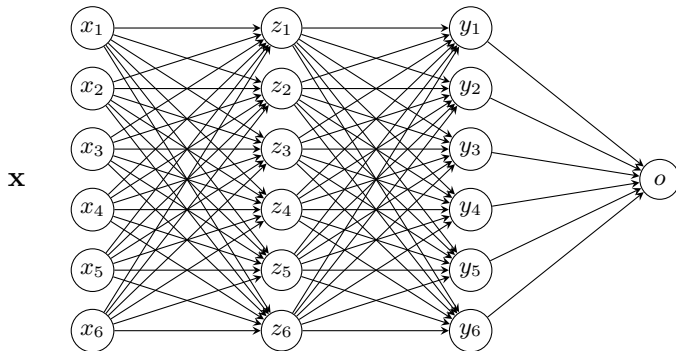
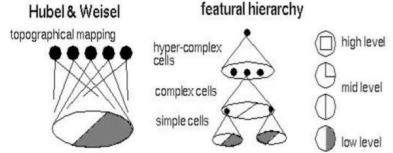
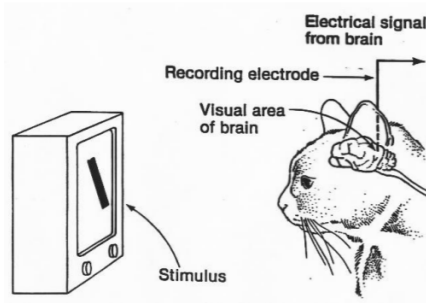


Fig. 6. A fully connected network with 2 hidden layers

- How *many* parameters to learn?
 - ▶ If we have n units in each layers, what is the total number of parameters?
 - ▶ high complexity may lead to *overfitting* and local optima
 - ▶ requires careful *initialization*

Hubel and Wiesel 1959



- measured neural responses in cat's striate cortex
- suggested a hierarchy of *feature detectors* in the visual cortex
- higher level features responding to patterns of activations in lower level cells

- Computer vision needs some robustness:
 - ▶ *translation invariance*

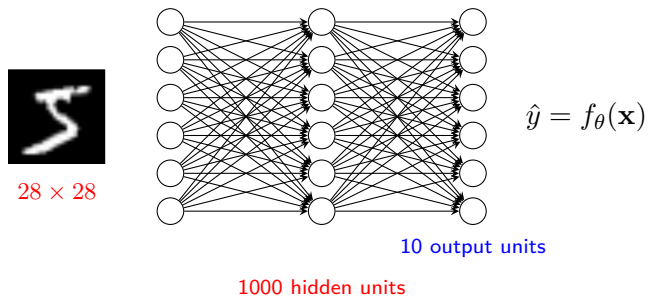


Cat

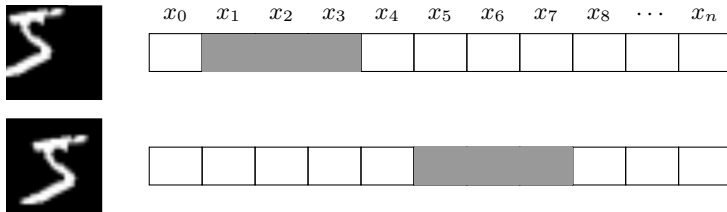


Cat

- need a network that recognizes the cat *regardless of its location*
- location might not be an important factor



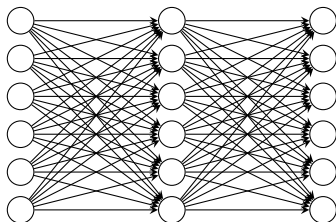
- How many parameters do we need for the first layer?
 - ▶ Parameters: $280 \times 280 \times 1000 = 78.4M$
- What if the input is slightly shifted?



- A unit's activation is dependent on the outputs of *all* nodes in the previous layer.



28 × 28

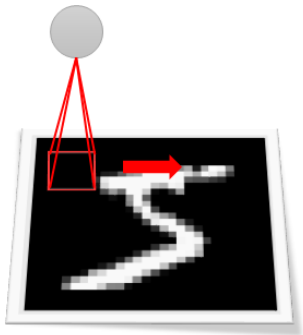


10 output units

1000 hidden units

$$\hat{y} = 5$$

- Each unit in the hidden layer is affected by the *entire* inputs.
- How can we detect the same feature at different positions?
- expect low-level features to be *local*
- expect high-level features to be *coarser*



- Let's have a feature detector
- An interesting pattern can be anywhere in the image.
 - ▶ Define a kernel (or filter) \mathbf{W} that search for a specific pattern.
 - ▶ The detector slides \mathbf{W} over the given image.
 - ▶ Inspect $\langle \mathbf{X}[\cdot, \cdot], \mathbf{W} \rangle$ and determine if the region contains the pattern
- Let's have multiple of them!

What is special about $\langle \mathbf{X}[\cdot, \cdot], \mathbf{W} \rangle$ operation?

- Suppose we have a 1D image $\mathbf{f} = (f_1, f_2, \dots, f_n)$ and a 1D kernel $\mathbf{g} = (g_1, \dots, g_m)$.
- $\langle \mathbf{X}[\cdot, \cdot], \mathbf{W} \rangle$ in 1D corresponds to

$$h[i] = (f * g)[i] = \sum_{j=1}^m g[j] \cdot f[i - j + m/2].$$

$f =$

9	0	2	1	0	9
---	---	---	---	---	---

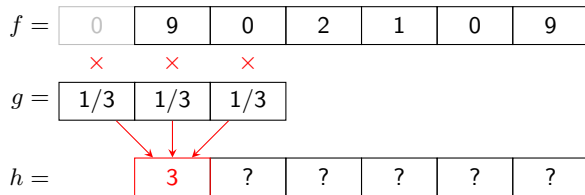
$g =$

1/3	1/3	1/3
-----	-----	-----

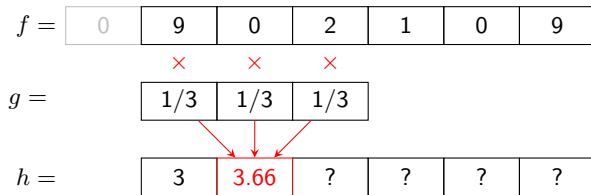
$h =$

?	?	?	?	?	?
---	---	---	---	---	---

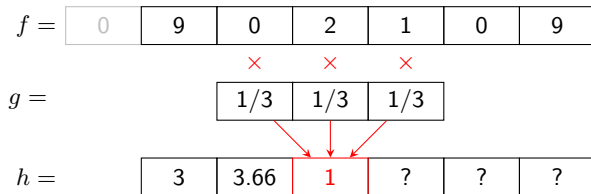
1D Convolution




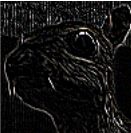
1D Convolution




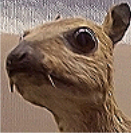
1D Convolution





- This can be easily extended to 2D.
- But, is h useful?


$$* \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix} =$$


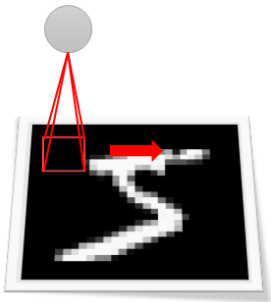
(a) Edge detection


$$* \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix} =$$


(b) Sharpen


$$* \frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} =$$


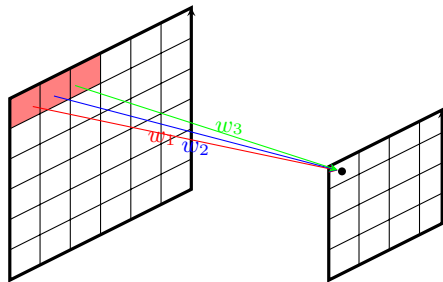
(c) Blur



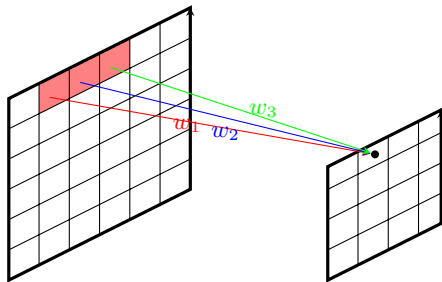
- But how do we know which filter to apply?

$$\mathbf{W} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix}$$

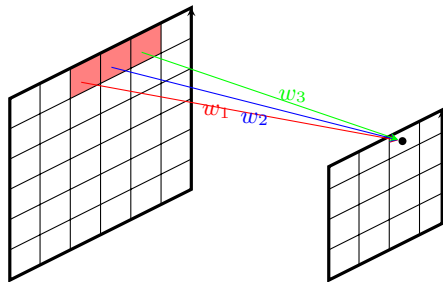
- ▶ Let our NN *learn* it from the data.
- ▶ \mathbf{W} are *shared* across different locations.



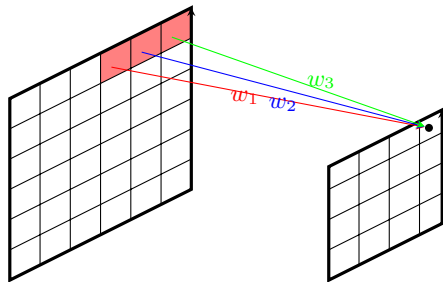
- $h_1 = w_1x_1 + w_2x_2 + w_3x_3$



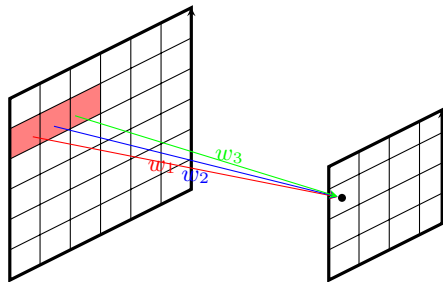
- $h_2 = w_1x_2 + w_2x_3 + w_3x_4$



- $h_3 = w_1x_3 + w_2x_4 + w_3x_5$



- $h_4 = w_1x_4 + w_2x_5 + w_3x_6$



- $h_5 = w_1x_7 + w_2x_8 + w_3x_9$

1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved
Feature

1	1 _{x1}	1 _{x0}	0 _{x1}	0
0	1 _{x0}	1 _{x1}	1 _{x0}	0
0	0 _{x1}	1 _{x0}	1 _{x1}	1
0	0	1	1	0
0	1	1	0	0

Image

4	3	

Convolved
Feature

1	1	1	0	0
0 _{x1}	1 _{x0}	1 _{x1}	1	0
0 _{x0}	0 _{x1}	1 _{x0}	1	1
0 _{x1}	0 _{x0}	1 _{x1}	1	0
0	1	1	0	0

Image

4	3	4
2		

Convolved
Feature

1	1	1	0	0
0	1 _{x1}	1 _{x0}	1 _{x1}	0
0	0 _{x0}	1 _{x1}	1 _{x0}	1
0	0 _{x1}	1 _{x0}	1 _{x1}	0
0	1	1	0	0

Image

4	3	4
2	4	

Convolved
Feature

1	1	1	0	0
0	1	1 _{x1}	1 _{x0}	0 _{x1}
0	0	1 _{x0}	1 _{x1}	1 _{x0}
0	0	1 _{x1}	1 _{x0}	0 _{x1}
0	1	1	0	0

Image

4	3	4
2	4	3

Convolved
Feature

1	1	1	0	0
0	1	1	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0 _{x0}	0 _{x1}	1 _{x0}	1	0
0 _{x1}	1 _{x0}	1 _{x1}	0	0

Image

4	3	4
2	4	3
2		

Convolved
Feature

1	1	1	0	0
0	1	1	1	0
0	0 _{x1}	1 _{x0}	1 _{x1}	1
0	0 _{x0}	1 _{x1}	1 _{x0}	0
0	1 _{x1}	1 _{x0}	0 _{x1}	0

Image

4	3	4
2	4	3
2	3	

Convolved
Feature

1	1	1	0	0
0	1	1	1	0
0	0	1 _{x1}	1 _{x0}	1 _{x1}
0	0	1 _{x0}	1 _{x1}	0 _{x0}
0	1	1 _{x1}	0 _{x0}	0 _{x1}

Image

4	3	4
2	4	3
2	3	4

Convolved
Feature

What does a Convolution do?

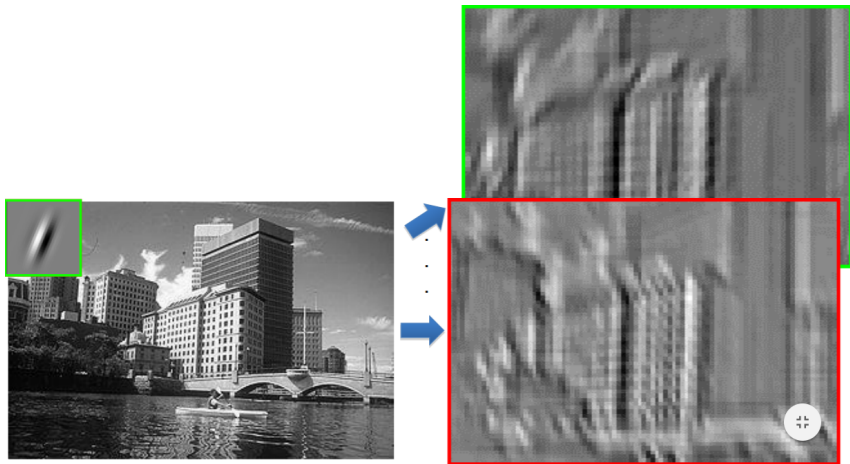


Fig. 8. Image credit: I. Kokkinos



Output Size

In our example, we have

- a 5×5 image,
- a Kernel (or filter) of size 3×3 , and
- Stride of 1.

This gives an output of size

$$\frac{N - K}{S} + 1.$$

- Suppose we have
 - ▶ 7×7 image
 - ▶ 3×3 kernel
 - ▶ stride of 3
- What is the output size?

Zero Padding



0	0	0	0	0	0	0	0	0	0
0	77	80	82	78	70	82	82	140	0
0	83	78	80	83	82	77	94	151	0
0	87	82	81	80	74	75	112	152	0
0	87	87	85	77	66	99	151	167	0
0	84	79	77	78	76	107	162	160	0
0	86	72	70	72	81	151	166	151	0
0	78	72	73	73	107	166	170	148	0
0	76	76	77	84	147	180	168	142	0
0	0	0	0	0	0	0	0	0	0

- image size = 8×8
- kernel = 3×3 with stride of 1
- What is the size of output? $\frac{10 - 3}{1} + 1 = 8$

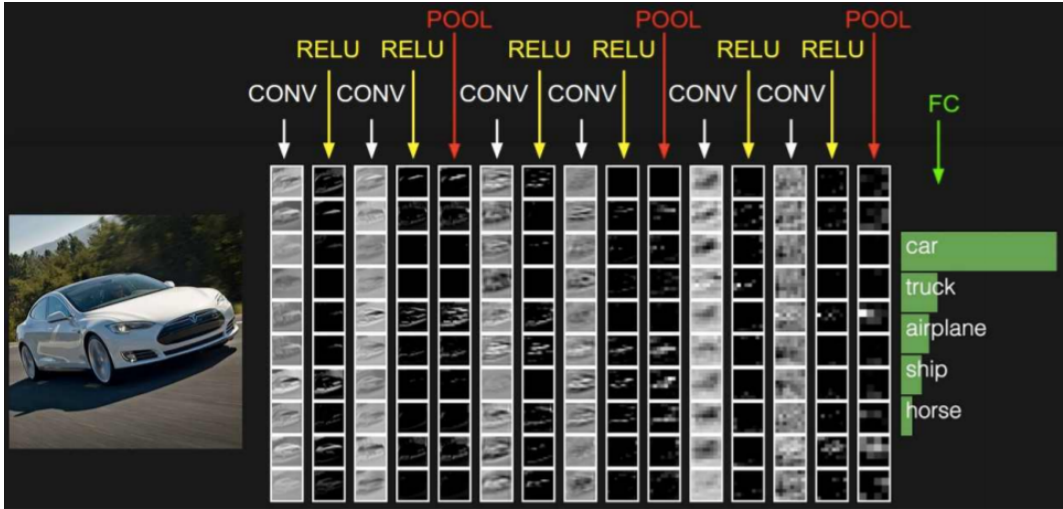
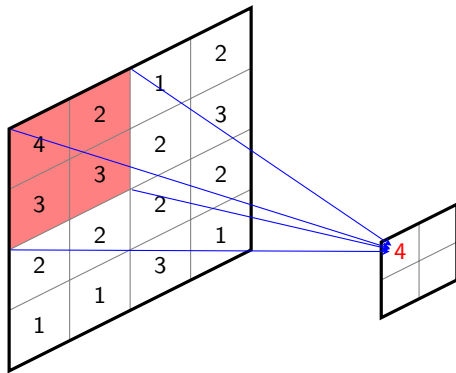
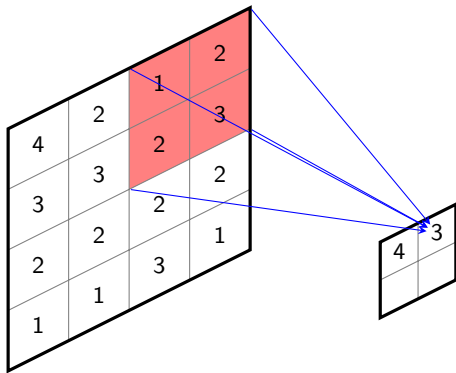


Fig. 9. Image source: Andrej Karpathy

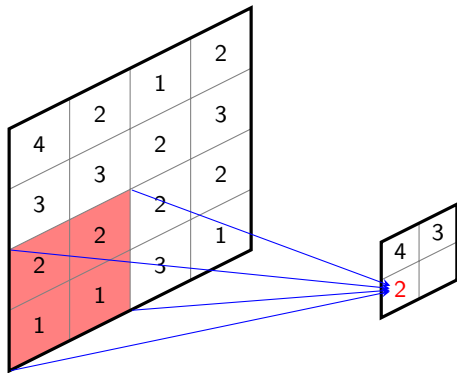
Pooling Layer



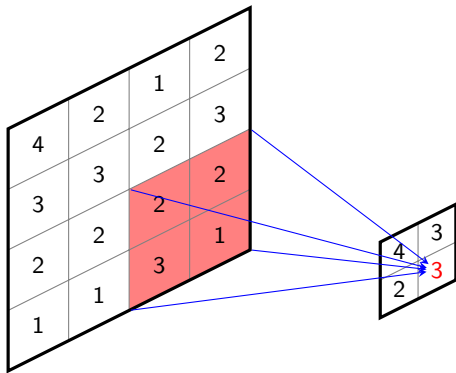
Pooling Layer



Pooling Layer



Pooling Layer



- What does the pooling layer do?
 - ▶ result in a smaller representation (shrunk image)
 - ▶ target feature values becomes less sensitive to small change in the image
 - ▶ stride > 1 is called “down sampling”
- Max pooling
- Average pooling
- L_p pooling ($p = 1 \Rightarrow$ average, $p = \infty \Rightarrow$ max pooling)

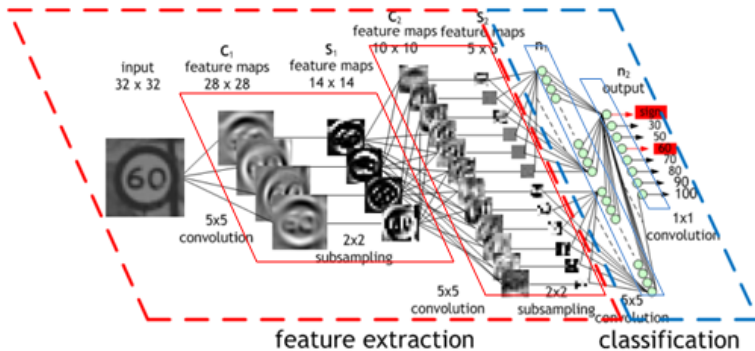


Fig. 10. Image source: NVIDIA

- The input to the next layer is *down sampled* image.
- Higher-layer filters see a larger region of the input than the equal-sized filters in the lower layers.

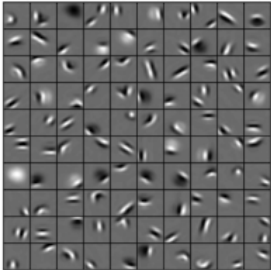


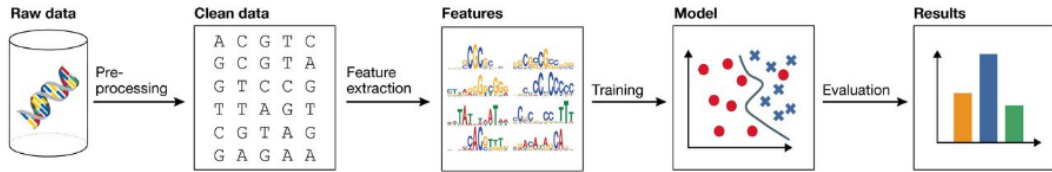
Fig. 11. Image source: Honglak Lee et al. 2011

PyTorch

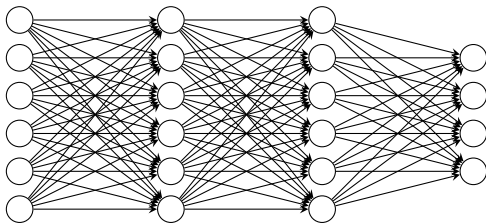
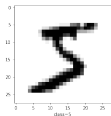
What is PyTorch?



- Python implementation of deep learning tools
 - ▶ A collection of classes
 - ▶ <https://pytorch.org>
- Developed by Meta



- 1 Dataset:** an *i.i.d.* sample drawn from the population
 - $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
- 2 Task:** predict Y given X
 - regression or classification?
- 3 Model:** assumption on the (distribution) of data
 - linear, non-linear, parametric, non-parametric
- 4 Loss(a.k.a. score/cost):** measure of how good a model is
 - $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$
- 5 Training:** normally involve an optimization
 - (stochastic) gradient descent, Newton's method
- 6 Evaluation:** validate the learning outcome



$$\hat{y} = f_{\theta}(\mathbf{x})$$

- Input: MNIST dataset
 - ▶ 28×28 gray scale image
 - ▶ vectorize the image to $\mathbf{x} \in \mathbb{R}^{784}$
- Architecture
 - ▶ 1st hidden layer: 512 units
 - ▶ 2nd hidden layer: 512 units
 - ▶ Output layer: 10 units (10 classes)

Building a Neural Network (1)



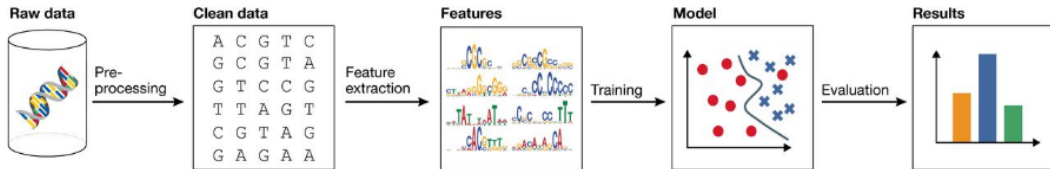
```
1 import torch.nn as nn
2
3 # Get cpu or gpu device for training.
4 device = "cuda" if torch.cuda.is_available() else "cpu"
5 print(f"Using {device} device")
6
7 # Define model
8 class NeuralNetwork(nn.Module):
9     def __init__(self):
10         super(NeuralNetwork, self).__init__()
11         self.flatten = nn.Flatten()
12         self.linear_relu_stack = nn.Sequential(
13             nn.Linear(28*28, 512),
14             nn.ReLU(),
15             nn.Linear(512, 512),
16             nn.ReLU(),
17             nn.Linear(512, 10)
18         )
19
20     def forward(self, x):
21         x = self.flatten(x)
22         logits = self.linear_relu_stack(x)
23         return logits
24
25 model = NeuralNetwork().to(device)
26 print(model)
```

```
1 import torch.nn.functional as F
2
3 # Define model
4 class NeuralNetwork(nn.Module):
5     def __init__(self):
6         super(NeuralNetwork, self).__init__()
7         self.flatten = nn.Flatten()
8         self.linear1 = nn.Linear(28*28, 512)
9         self.linear2 = nn.Linear(512, 512)
10        self.classifier = nn.Linear(512, 10)
11
12    def forward(self, x):          # define how the input tensor is processed
13        x = self.flatten(x)
14        out = F.relu(self.linear1(x))
15        out = F.relu(self.linear2(out))
16        logits = self.classifier(out)
17
18        return logits
19
20 model = NeuralNetwork().to(device)
21 print(model)
```

Defining Datasets

```
1 import torch
2 from torch import nn
3 from torch.utils.data import DataLoader
4 from torchvision import datasets
5 from torchvision.transforms import ToTensor
6
7 # Download training data from open datasets.
8 training_data = datasets.FashionMNIST(
9     root="data", # path to the directory containing the dataset
10    train=True,
11    download=True, # If True, downloads the dataset from the internet
12    transform=ToTensor(),
13 )
14
15 # Download test data from open datasets.
16 test_data = datasets.FashionMNIST(
17    root="data",
18    train=False,
19    download=True,
20    transform=ToTensor(),
21 )
```

```
1  batch_size = 64
2
3  # Create data loaders.
4  train_dataloader = DataLoader(training_data, batch_size=batch_size)
5  test_dataloader = DataLoader(test_data, batch_size=batch_size)
6
7  for X, y in test_dataloader:
8      print(f"Shape of X [N, C, H, W]: {X.shape}")
9      print(f"Shape of y: {y.shape} {y.dtype}")
10     break
```

```
1 loss_fn = nn.CrossEntropyLoss()
2 optimizer = torch.optim.SGD(model.parameters(), lr=1e-3)
```

- `nn.MSELoss`
- `nn.NLLLoss`
- `nn.L1Loss`
- ...
- See [here](#) for other available functions.

```
1 def train(dataloader, model, loss_fn, optimizer):
2     size = len(dataloader.dataset)
3     model.train()
4
5     for batch, (X, y) in enumerate(dataloader):
6         X, y = X.to(device), y.to(device)
7
8         # Compute prediction error
9         pred = model(X)
10        loss = loss_fn(pred, y)
11
12        # Backpropagation
13        optimizer.zero_grad()
14        loss.backward()
15        optimizer.step()
16
17        if batch % 100 == 0:
18            loss, current = loss.item(), batch * len(X)
19            print(f"loss: {loss:>7f}    [{current:>5d}/{size:>5d}]")
```

```
1 def test(dataloader, model, loss_fn):
2     size = len(dataloader.dataset)
3     num_batches = len(dataloader)
4     model.eval()
5     test_loss, correct = 0, 0
6     with torch.no_grad():
7         for X, y in dataloader:
8             X, y = X.to(device), y.to(device)
9             pred = model(X)
10            test_loss += loss_fn(pred, y).item()
11            correct += (pred.argmax(1) == y).type(torch.float).sum().item()
12    test_loss /= num_batches
13    correct /= size
14    print(f"Test Error: \n Accuracy: {(100*correct):>0.1f}%, Avg loss: {test_loss:>8f} \n")
```

```
1 epochs = 5
2
3 for t in range(epochs):
4     print(f"Epoch {t+1}\n-----")
5     train(train_dataloader, model, loss_fn, optimizer)
6     test(test_dataloader, model, loss_fn)
7 print("Done!")
```