

# Tutorial on basic neural network

## ① Supervised Learning

- Learn a function from a family  $\mathcal{F} = \{f_\theta : \theta \in \Theta\}$ .

$$f_\theta : X \mapsto Y \quad \text{e.g. } X = \mathbb{R}^d, Y = \{0, 1\} / \mathbb{R}^{d_2} / \dots$$

- Loss function  $L(f_\theta(x), y)$ :

$$\begin{aligned} \text{E.g. } L(f_\theta(x), y) &= -y \log f_\theta(x) - (1-y) \log (1-f_\theta(x)) \\ L(f_\theta(x), y) &= \|y - f_\theta(x)\|_2^2. \end{aligned}$$

- ERM:

$$\theta^* = \underset{\theta \in \Theta}{\operatorname{argmin}} \sum_{i=1}^n L(f_\theta(x_i), y_i).$$

---

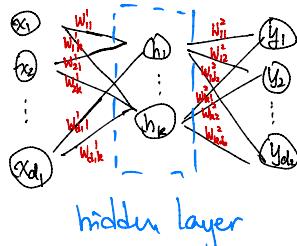
$$\begin{cases} \mathcal{F} = \text{all the linear functions, i.e. } \{a^T x + b : a \in \mathbb{R}^d, b \in \mathbb{R}\} \\ L(f_\theta(x), y) = \|y - f_\theta(x)\|_2^2 \end{cases}$$

→ Linear regression

## ② Neural Networks - Multi-layer Perceptron (MLP)

- What is the  $\mathcal{F}$ ?

$$f_\theta : \mathbb{R}^{d_1} \mapsto \mathbb{R}^{d_2}$$



For some nonlinear function  $\sigma(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ .

$$h_1 = \sigma\left(\sum_{i=1}^{d_1} w_{1i} x_i + b_1\right), \dots, h_k = \sigma\left(\sum_{i=1}^{d_1} w_{ki} x_i + b_k\right),$$

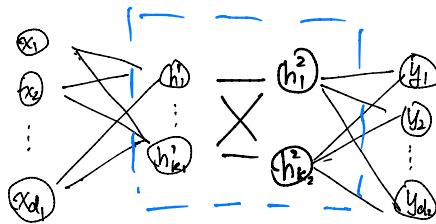
$$y_1 = \sum_{j=1}^k w_{j1}^2 h_j + b_1^2, \dots, y_{d_2} = \sum_{j=1}^k w_{jd_2}^2 h_j + b_{d_2}^2$$

$$h = \sigma(W^1 x + b^1), \quad y = W^2 h + b^2.$$

Popular choice of  $\sigma$  { Sigmoid:  $\sigma(x) = \frac{1}{1+e^{-x}}$   
ReLU:  $\sigma(x) = \max(0, x)$ .

$$\theta = (W^1, W^2, b^1, b^2)$$

We can have more than one hidden layer:



- If we are dealing with a  $d_2$ -class classification problem, we can apply a softmax function to the output  $y_1, \dots, y_{d_2}$ , i.e.

$$p_i = \frac{\exp(y_i)}{\sum_{j=1}^{d_2} \exp(y_j)}$$

Then,  $(p_1, \dots, p_{d_2})$  satisfies.

- $p_i \geq 0$  for all  $i \in \{1, \dots, d_2\}$ .
- $\sum_{i=1}^{d_2} p_i = 1$

and therefore is a valid probability distribution.

What are the choices of  $L$ ?

- Cross entropy loss.  
(Maximum Likelihood).

$$L = \begin{cases} - \sum_{i=1}^n (y_i \log p_i + (1-y_i) \log (1-p_i)) & \text{(Binary)} \\ - \sum_{i=1}^n \sum_{k=1}^K y_{ik} \log p_{ik} & \text{(General)} \end{cases}$$

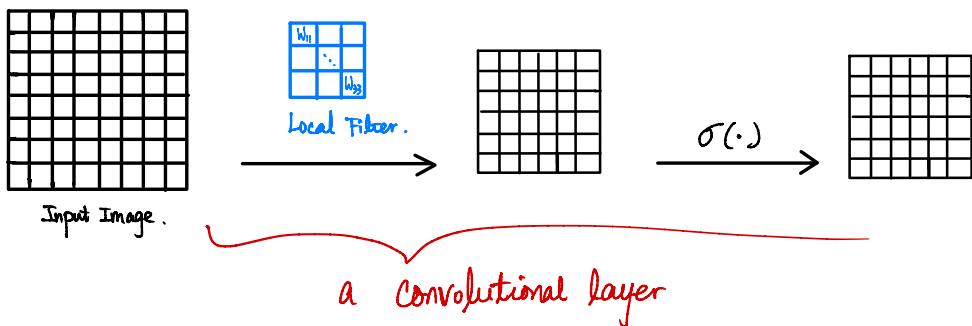
- Hinge Loss:  $L = \sum_{i=1}^n \max(0, y_i - (1-2y_i) \cdot p_i)$ .
- $L_1$  Loss,  $L_2$  Loss, ...

### ③ Neural net - Convolution Neural Net. (CNN).

Consider image input, with resolution  $100 \times 100$  pixels.  
If we use MLP with one hidden layer of dimension 20.  
What is the number of parameter to learn?

- Slow Training
- Overfitting

For image data, local information is some time more useful.  
Instead of training the full weight matrix.  $\mathbb{R}^{10000} \rightarrow \mathbb{R}^{20}$ ,  
we consider the local linear filter as follows.



Following that, one can.

- apply another conv layer directly; or
- flatten the output to a vector and apply a fully connect layer.
- ...

There are small tricks to finetune the conv layer.

Zero padding, striding, ...

#### ④ Train neural nets

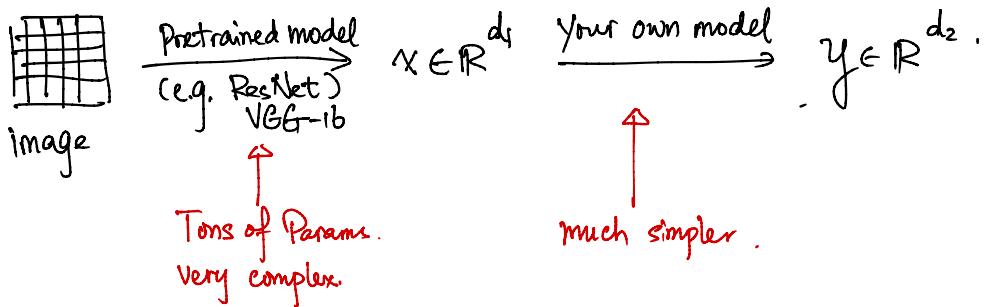
Target:  $\min_{\theta \in \Theta} L(f_\theta(x), y)$ ,  $\theta = ((W^1, b^1), (W^2, b^2), \dots)$ .

Method: Gradient descent or other optimization algorithms.

Implementation: PyTorch / Tensorflow.

They have great tutorials with codes.

Use pretrained model as feature extractors?



#### Why is neural nets popular?

- Universal Approximation: MLP with a single hidden layer can approximate any continuous function in any compact set.
- Good empirical performance in many applications.

#### What are the shortcomings?

- Poor interpretability
- High computational costs