

# Optimization Problems: Fully-connected Networks

Chih-Jen Lin  
National Taiwan University

Last updated: April 23, 2021

# Multi-class Classification I

- Our training set includes  $(\mathbf{y}^i, \mathbf{x}^i)$ ,  $i = 1, \dots, l$ .
- $\mathbf{x}^i \in R^{n_1}$  is the feature vector.
- $\mathbf{y}^i \in R^K$  is the label vector.
- As label is now a vector, we change (label, instance) from

$$(y_i, \mathbf{x}_i) \text{ to } (\mathbf{y}^i, \mathbf{x}^i)$$

- $K$ : # of classes
- If  $\mathbf{x}^i$  is in class  $k$ , then

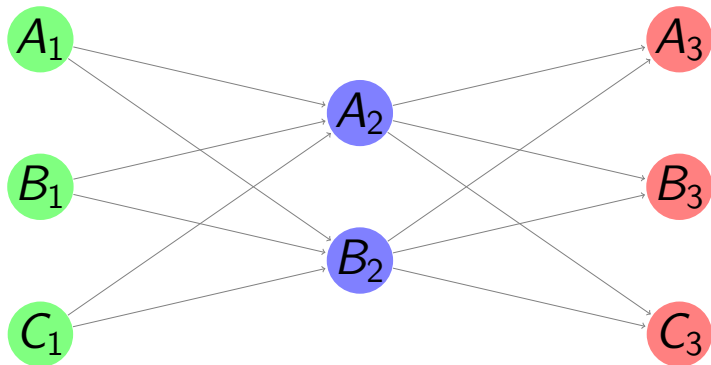
$$\mathbf{y}^i = \underbrace{[0, \dots, 0]}_{k-1}, 1, 0, \dots, 0]^T \in R^K$$

# Multi-class Classification II

- A neural network maps each feature vector to one of the class labels by the connection of nodes.

# Fully-connected Networks

- Between two layers a weight matrix maps inputs (the previous layer) to outputs (the next layer).



# Operations Between Two Layers I

- The weight matrix  $W^m$  at the  $m$ th layer is

$$W^m = \begin{bmatrix} W_{11}^m & W_{12}^m & \cdots & W_{1n_m}^m \\ W_{21}^m & W_{22}^m & \cdots & W_{2n_m}^m \\ \vdots & \vdots & \vdots & \vdots \\ W_{n_{m+1}1}^m & W_{n_{m+1}2}^m & \cdots & W_{n_{m+1}n_m}^m \end{bmatrix}_{n_{m+1} \times n_m}$$

- $n_m$  : # input features at layer  $m$
- $n_{m+1}$  : # output features at layer  $m$ , or # input features at layer  $m + 1$
- $L$ : number of layers

# Operations Between Two Layers II

- $n_1 = \#$  of features,  $n_{L+1} = \#$  of classes
- Let  $\mathbf{z}^m$  be the input of the  $m$ th layer,  $\mathbf{z}^1 = \mathbf{x}$  and  $\mathbf{z}^{L+1}$  be the output
- From  $m$ th layer to  $(m + 1)$ th layer

$$\mathbf{s}^m = \mathbf{W}^m \mathbf{z}^m,$$
$$z_j^{m+1} = \sigma(s_j^m), \quad j = 1, \dots, n_{m+1},$$

$\sigma(\cdot)$  is the activation function.

# Operations Between Two Layers III

- Usually people do a bias term

$$\begin{bmatrix} b_1^m \\ b_2^m \\ \vdots \\ b_{n_{m+1}}^m \end{bmatrix}_{n_{m+1} \times 1},$$

so that

$$\mathbf{s}^m = \mathbf{W}^m \mathbf{z}^m + \mathbf{b}^m$$

# Operations Between Two Layers IV

- Activation function is usually an

$$R \rightarrow R$$

non-linear transformation.

- There are various reasons of using an activation function. An important one is to introduce the non-linearity.



# Operations Between Two Layers V

- If without an activation function, all

$$W^L \dots W^2 W^1$$

becomes a single matrix and we end up with having only a linear mapping from the input feature to the output layer

# Operations Between Two Layers VI

- We collect **all variables**:

$$\theta = \begin{bmatrix} \text{vec}(W^1) \\ \mathbf{b}^1 \\ \vdots \\ \text{vec}(W^L) \\ \mathbf{b}^L \end{bmatrix} \in R^n$$

$n$  : total # variables =  $(n_1 + 1)n_2 + \dots + (n_L + 1)n_{L+1}$

- The  $\text{vec}(\cdot)$  operator stacks columns of a matrix to a vector

# Optimization Problem I

- We solve the following optimization problem,

$$\min_{\theta} f(\theta), \quad \text{where}$$

$$f(\theta) = \frac{1}{2} \theta^T \theta + C \sum_{i=1}^l \xi(z^{L+1,i}(\theta); \mathbf{y}^i, \mathbf{x}^i).$$

C: regularization parameter

- $\mathbf{z}^{L+1}(\theta) \in R^{n_{L+1}}$ : last-layer output vector of  $\mathbf{x}$ .  
 $\xi(\mathbf{z}^{L+1}; \mathbf{y}, \mathbf{x})$ : loss function. Example:

$$\xi(\mathbf{z}^{L+1}; \mathbf{y}, \mathbf{x}) = \|\mathbf{z}^{L+1} - \mathbf{y}\|^2$$

# Optimization Problem II

- The formulation is **same as linear classification**
- However, the loss function is **more complicated**
- Further, it's **non-convex**
- Note that in the earlier discussion we consider a single instance
- In the training process we actually have for  $i = 1, \dots, l$ ,

$$\begin{aligned} \mathbf{s}^{m,i} &= W^m \mathbf{z}^{m,i}, \\ z_j^{m+1,i} &= \sigma(s_j^{m,i}), \quad j = 1, \dots, n_{m+1}, \end{aligned}$$

This makes the training more complicated