

Lecture 18: Neural Network and Backpropagation II

[SCS4049-02] Machine Learning and Data Science

Seongsik Park (s.park@dgu.edu)

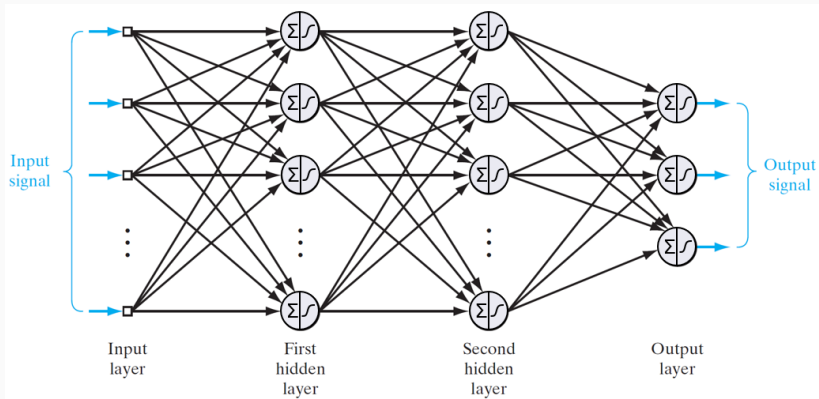
AI Department, Dongguk University

Tentative schedule

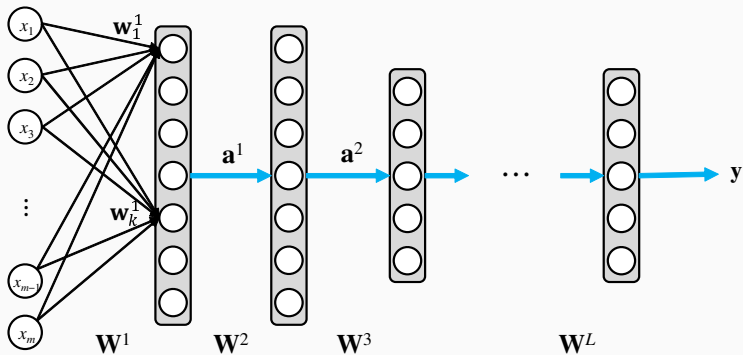
week	topic	date (수 / 월)
1	Machine Learning Introduction & Basic Mathematics	09.02 / 09.07
2	Python Practice I & Regression	09.09 / 09.14
3	AI Department Seminar I & Clustering I	09.16 / 09.21
4	Clustering II & Classification I	09.23 / 09.28
5	Classification II	(추석) / 10.05
6	Python Practice II & Support Vector Machine I	10.07 / 10.12
7	Support Vector Machine II & Decision Tree and Ensemble Learning	10.14 / 10.19
8	Mid-Term Practice & Mid-Term Exam	10.21 / 10.26
9	휴강 & Dimensional Reduction I	10.28 / 11.02
10	Dimensional Reduction II & Neural networks and Back Propagation I	11.04 / 11.09
11	Neural networks and Back Propagation II & Convolutional Neural Network	11.11 / 11.16
12	Model Optimization	11.18 / 11.23
13	Recurrent Neural network	11.25 / 11.30
14	Autoencoders	12.02 / 12.07
15	Final exam	(휴강) / 12.14

Feed-Forward Propagation and Cost Function

Deep feedforward networks



Goal of deep feedforward networks



Goal of deep feedforward networks

Let $\mathbf{W}^l = \begin{bmatrix} \mathbf{w}_1^l & \mathbf{w}_2^l & \cdots & \mathbf{w}_{m_l}^l \end{bmatrix}$ and φ^l be the activation function for layer l , then

$$\mathbf{y}^1 = \varphi^1(\mathbf{W}^1 \mathbf{x}) \quad (1)$$

$$\mathbf{y}^2 = \varphi^2(\mathbf{W}^2 \mathbf{y}^1) \quad (2)$$

$$= \varphi^2(\mathbf{W}^2 \varphi^1(\mathbf{W}^1 \mathbf{x})) \quad (3)$$

$$\mathbf{y}^3 = \varphi^3(\mathbf{W}^3 \mathbf{y}^2) \quad (4)$$

$$= \varphi^3(\mathbf{W}^3 \varphi^2(\mathbf{W}^2 \mathbf{y}^1)) \quad (5)$$

$$= \varphi^3(\mathbf{W}^3 \varphi^2(\mathbf{W}^2 \varphi^1(\mathbf{W}^1 \mathbf{x}))) \quad (6)$$

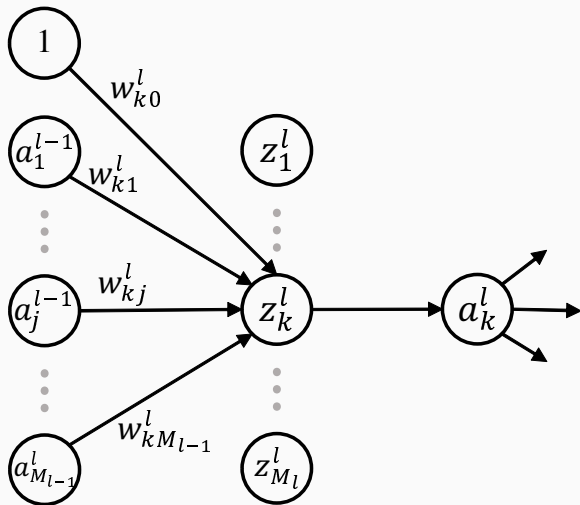
...

$$\mathbf{y} = \varphi^L(\mathbf{W}^L \varphi^{L-1}(\mathbf{W}^{L-1} \varphi^{L-2}(\cdots \varphi^2(\mathbf{W}^2 \varphi^1(\mathbf{W}^1 \mathbf{x})) \cdots))) \quad (7)$$

Very much complex nonlinear function

\implies require universal function approximator

Feedforward propagation



Feedforward propagation

For l -th hidden layer

$$z_k^l = \sum_{j=0}^{M_{l-1}} w_{kj}^l a_j^{l-1} \quad a_k^l = \sigma^l(z_k^l) \quad (8)$$

$$\mathbf{z}^l = \mathbf{W}^l \mathbf{a}^{l-1} \quad \mathbf{a}^l = \boldsymbol{\sigma}^l(\mathbf{z}^l) \quad (9)$$

- z_k^l : layer l 의 unit k 로 들어오는 activation들의 weighted sum
- w_{kj}^l : layer $(l-1)$ 의 unit j 에서 layer l 의 unit k 로 가는 가중치
- w_{k0}^l : layer l 의 unit k 에 대한 bias
- a_j^{l-1} : layer $(l-1)$ 의 unit j 에서 나오는 activation 값
- a_0^l : layer l 의 bias에 대응 (= 1)
- a_k^l : layer l 의 unit k 의 activation 값
- σ^l : layer l 의 activation function
- M_{l-1} : layer $(l-1)$ 의 hidden unit 수

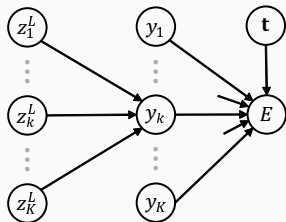
Cost function

Cost function: cross-entropy

$$E = -\frac{1}{N} \sum_{n=1}^N \sum_{j=1}^K t_j \log y_j \quad (10)$$

$$y_k = \frac{\exp(z_k)}{\sum_{j=1}^K \exp(z_j)} \quad (11)$$

- t_j : \mathbf{x} 의 class label (target value)
- K : 출력층의 unit 수 = class의 수, $M_L = K$



Stochastic gradient and backpropagation

SGD weight update rule

$$w_{ij} \leftarrow w_{ij} - \eta \frac{\partial E}{\partial w_{ij}} \quad (12)$$

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \nabla_{\mathbf{W}} E \quad (13)$$

Credit assignment

- Weight 값들이 잘못 되었기 때문에 Cost 가 발생
⇒ 따라서 모든 weight에게 그 weight가 cost 발생에 기여한 만큼씩 gradient 를 배분해 주는 것
- 문제는 weight 들이 여러 층에 분산되어 있으며, 각 층은 비선형 함수이므로 배분 방법이 단순하지 않음 ⇒ backpropagation algorithm

Backpropagation

출력층부터 입력층까지 역순으로 모든 layer에서 weight w_{ij}^l 을 업데이트

$$w_{ij} \leftarrow w_{ij} - \eta \frac{\partial E}{\partial w_{ij}} \quad (14)$$

Local gradient (or error of neuron)의 정의

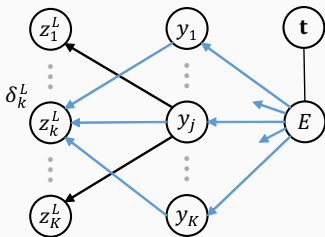
$$\delta_j^l \triangleq \frac{\partial E}{\partial z_j^l} \quad (15)$$

- Layer l 의 unit k 에서 error, 또는 하위층으로 분배할 gradient
- Activation a_j^l 를 기준으로 생각할 수 있으나, nonlinear 함수를 포함하므로 다루기 어려움

출력층에서의 local gradient δ_k^L 계산

$$\delta_k^L = \frac{\partial E}{\partial z_k^L} = \sum_j \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial z_k^L} \quad (16)$$

- Softmax 함수는 $z_1^L, z_2^L, \dots, z_k^L$ 에 의존
- 따라서 z_k^L 의 변동은 y_1, y_2, \dots, y_K 에 모두 영향



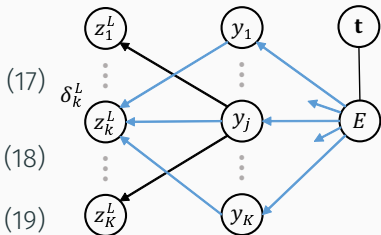
Backpropagation: output layer 2/7

$\frac{\partial y_j}{\partial z_k^L}$ 의 계산: $k = j$ 인 경우

$$\frac{\partial y_k}{\partial z_k^L} = \frac{\partial \exp(z_k)}{\partial z_k^L \sum_i \exp(z_i)}$$

$$= \frac{\exp(z_k)}{\sum_i \exp(z_i)} \left(1 - \frac{\exp(z_k)}{\sum_i \exp(z_i)} \right) \quad (18)$$

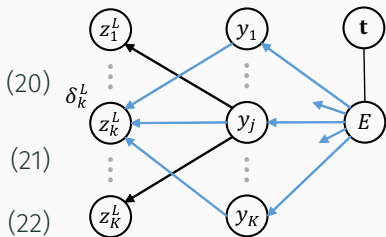
$$= y_k(1 - y_k)$$



Backpropagation: output layer 3/7

$\frac{\partial y_j}{\partial z_k^L}$ 의 계산: $k \neq j$ 인 경우

$$\begin{aligned} \frac{\partial y_j}{\partial z_k^L} &= \frac{\partial}{\partial z_k^L} \frac{\exp(z_j)}{\sum_i \exp(z_i)} \\ &= -\frac{\exp(z_k)}{\sum_i \exp(z_i)} \frac{\exp(z_j)}{\sum_i \exp(z_i)} \\ &= -y_k y_j \end{aligned}$$



Backpropagation: output layer 4/7

From definition of local gradient

$$\delta_j^l = \frac{\partial E}{\partial z_j^l} = - \sum_j t_j \frac{\partial}{\partial z_k} \log y_j = - \sum_j t_j \frac{1}{y_j} \frac{\partial y_j}{\partial z_k} \quad (23)$$

$$= - \frac{t_k}{y_k} \frac{\partial y_j}{\partial z_k} - \sum_{j \neq k} \frac{t_j}{y_j} \frac{\partial y_j}{\partial z_k} \quad (24)$$

$$= - \frac{t_k}{y_k} y_k (1 - y_k) - \sum_{j \neq k} \frac{t_j}{y_j} (-y_k y_j) \quad (25)$$

$$= -t_k + t_k y_k + \sum_{j \neq k} t_j y_k \quad (26)$$

$$= -t_k + \sum_{j=1}^K t_j y_k \quad (27)$$

$$= -(t_k - y_k) \quad (28)$$

Backpropagation: output layer 5/7

In matrix form

$$\delta^L = \frac{\partial E}{\partial \mathbf{z}^L} = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{z}^L} \right)^T \left(\frac{\partial E}{\partial \mathbf{y}} \right) \quad (29)$$

$$= \begin{bmatrix} \frac{\partial y_1}{\partial z_1} & \cdots & \frac{\partial y_j}{\partial z_1} & \cdots & \frac{\partial y_K}{\partial z_1} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial z_k} & \cdots & \frac{\partial y_j}{\partial z_k} & \cdots & \frac{\partial y_K}{\partial z_k} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial z_K} & \cdots & \frac{\partial y_j}{\partial z_K} & \cdots & \frac{\partial y_K}{\partial z_K} \end{bmatrix} \begin{bmatrix} \frac{\partial E}{\partial y_1} \\ \vdots \\ \frac{\partial E}{\partial y_j} \\ \vdots \\ \frac{\partial E}{\partial y_K} \end{bmatrix} \quad (30)$$

$$= \begin{bmatrix} \sum_j \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial z_1} \\ \vdots \\ \sum_j \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial z_r} \\ \vdots \\ \sum_j \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial z_K} \end{bmatrix} = \begin{bmatrix} -(t_1 - y_1) \\ \vdots \\ -(t_k - y_k) \\ \vdots \\ -(t_K - y_K) \end{bmatrix} \quad (31)$$

With vector form

$$\mathbf{y} = \varphi(\mathbf{z}^L) \quad (32)$$

$$\delta^L = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{z}^L} \right)^T \left(\frac{\partial E}{\partial \mathbf{y}} \right) \quad (33)$$

$$= (\varphi'(\mathbf{z}^L))^T \nabla_{\mathbf{y}} E \quad (34)$$

$$= -(\mathbf{t} - \mathbf{y}) \quad (35)$$

Backpropagation: output layer 7/7

$\frac{\partial E}{\partial w_{kj}^L}$ 의 계산

$$z_k^L = \sum_{j=0}^{M_{L-1}} w_{kj}^L a_j^{L-1} \quad (36)$$

$$\frac{\partial z_k^L}{\partial w_{kj}^L} = a_j^{L-1} \quad (37)$$

$$\frac{\partial E}{\partial w_{kj}^L} = \frac{\partial E}{\partial z_k^L} \frac{\partial z_k^L}{\partial w_{kj}^L} = -(t_k - y_k) a_j^{L-1} = \delta_k^L a_j^{L-1} \quad (38)$$

그러므로 weight update rule은

$$w_{kj}^L \leftarrow w_{kj}^L - \eta \frac{\partial E}{\partial w_{kj}^L} \quad (39)$$

$$= w_{kj}^L - \eta \delta_k^L a_j^{L-1} \quad (40)$$

$$= w_{kj}^L + \eta (t_k - y_k) a_j^{L-1} \quad (41)$$

$$(42)$$

Backpropagation: local gradient propagation

From δ_j^{l+1} to δ_k^l

$$\delta_k^l = \frac{\partial E}{\partial z_k^l} = \sum_j \frac{\partial E}{\partial z_j^{l+1}} \frac{\partial z_j^{l+1}}{\partial z_k^l} = \sum_j \frac{\partial z_j^{l+1}}{\partial z_k^l} \delta_j^{l+1} \quad (43)$$

$$z_j^{l+1} = \sum_{i=1} w_{ji}^{l+1} a_i^l = \sum_{i=0} w_{ji}^{l+1} \sigma^l(z_i^l) \quad (44)$$

$$\frac{\partial z_j^{l+1}}{\partial z_k^l} = w_{jk}^{l+1} \sigma'^l(z_k^l) \quad (45)$$

그러므로

$$\delta_k^l = \sigma'^l(z_k^l) \sum_j w_{jk}^{l+1} \delta_j^{l+1} \quad (46)$$

Backpropagation: weight in layer l

$\frac{\partial E}{\partial w_{kj}^l}$ 의 계산

$$z_k^l = \sum_{j=0}^{M_{l-1}} w_{kj}^l a_j^{l-1} \quad (47)$$

$$\frac{\partial z_k^l}{\partial w_{kj}^l} = a_j^{l-1} \quad (48)$$

$$\frac{\partial E}{\partial w_{kj}^l} = \frac{\partial E}{\partial z_k^l} \frac{\partial z_k^l}{\partial w_{kj}^l} = \delta_k^l a_j^{l-1} \quad (49)$$

그러므로 weight update rule은

$$w_{kj}^l \leftarrow w_{kj}^l - \eta \frac{\partial E}{\partial w_{kj}^l} \quad (50)$$

$$= w_{kj}^l - \eta \delta_k^l a_j^{l-1} \quad (51)$$

Backpropagation: summary

1. Input \mathbf{x} : 입력층에서 activation \mathbf{a}^1 을 계산
2. Feedforward: 각 층 $l = 2, 3, \dots, L$ 에 대해서 다음을 계산

$$\mathbf{z}^l = \mathbf{W}^l \mathbf{a}^{l-1} \quad \mathbf{a}^l = \sigma^l(\mathbf{z}^l) \quad (52)$$

3. Output error δ^L : 출력층에서 비용함수를 이용해 local gradient 계산

$$\delta^L = -(\mathbf{t} - \mathbf{y}) \quad (53)$$

4. Backpropagate the error: 출력층부터 거꾸로 다음 식을 계산

$$\delta_k^l = \sigma^{l'}(z_k^l) \sum_j w_{jk}^{l+1} \delta_{jk}^{l+1} \quad (54)$$

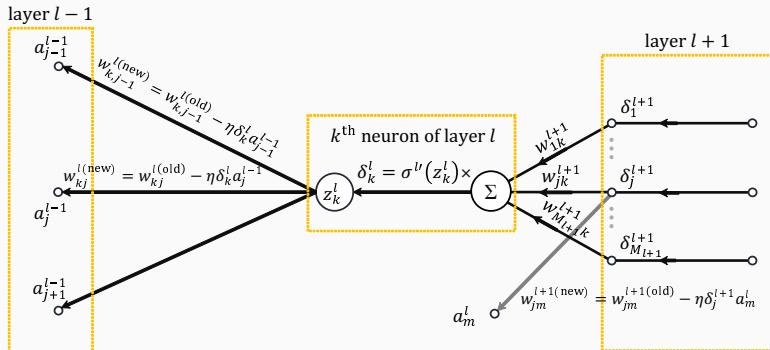
5. Output the gradient: 각 weight에 대해 gradient 계산

$$\frac{\partial E}{\partial w_{kj}^l} = \delta_k^l a_j^{l-1} \quad (55)$$

6. Update weight:

$$w_{kj}^l \leftarrow w_{kj}^l - \eta \delta_k^l a_j^{l-1} \quad (56)$$

Backpropagation: summary



Appendix

Reference and further reading

- “Chap 10 | Introduction to Artificial Neural Network with Keras” and “ Chap 11 | Training Deep Neural Networks” of A. Geron, Hands-On Machine Learning with Scikit-Learn, Keras & TensorFlow
- “Lecture 9 | Neural Networks and Deep Learning” of Kwang Il Kim, Machine Learning (2019)