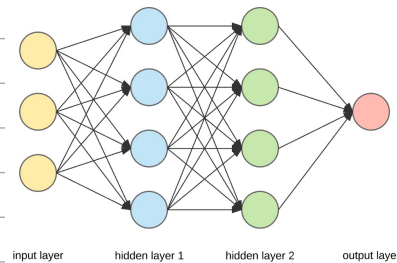


13 Neural Network < 뉴런 닮은 >

[input → hidden layer → output] 구조의 함수 근사 모델



13.1 Overview of NN Architectures

* Architecture :

$$x \xrightarrow{W^{(1)}, \sigma} h_1 \xrightarrow{W^{(2)}, \sigma} h_2 \xrightarrow{\dots W^{(L)}} \hat{y}$$

- 각 layer : $h^{(l)} = \sigma(W^{(l)}h^{(l-1)} + b^{(l)})$
- σ : non-linear function
- $W^{(l)}$: weight matrix, parameter
- $b^{(l)}$: bias vector

• Activation function

e.g)

Sigmoid $\sigma(x) = \frac{1}{1+e^{-x}}$

tanh $\frac{e^x - e^{-x}}{e^x + e^{-x}}$

ReLU $\max(0, x)$

Leaky ReLU $\max(0.01x, x)$

- image/text 분류
 - voice 인식
 - 예측 모델링
- (수학적...)

입력공간 \mathbb{R}^d 을 출력공간 \mathbb{R}^k 에 mapping 하는 가중치 벡터 θ 로 매개변수화된 함수. 입력을 층별로 출력으로 변환.

층 0 + 1 의 값은 이전 층 0 의 값의 함수

→ linear / Nonlinear

→ layer를 여러개 쌓으면 복잡한 nonlinear function 근사가능해짐

• 학습은 loss function ↓ : backpropagation + GD

구조 변형

- MLP (Fully connected) : 모든 node 연결
- CNN : spatial structure (공간구조) 유지 (image)
- RNN : Recurrent : 순서 정보 (시계열, text)
- ResNet, Transformer 더 깊은 attention 도입...

MLP (Fully connected) / CNN (공간

판도리...

CNN : Convolutional NN

image data에 구조 사용.

∴ 이미지의 spatial structure (공간구조) 활용 가능

Fully connected layer 사용

- 32 x 32 x 3 이미지 → 3072 개 node
- 다음 layer 가 1000 개 node 가지면? → 3072000 개

↓ parameter 너무 많음

$$Z_{i,j} = \sum X_{i+m,j+n} \cdot K_{m,n}$$

X : 입력 image

K : 필터

Z : 출력 feature map

작은 kernel 이 입력 이미지를 일정간격으로 슬라이딩하여 local feature 를 추출

* RNN : Recurrent NN

순서가 중요한 data (시계열 / text / audio)

구조) 이전의 hidden state 를 현재로 전달.

$$h_t = \sigma(W_{hh} h_{t-1} + W_{xh} x_t + b) \dots$$

$$\text{output}) y_t = W_{hy} h_t$$

hidden state 가 순서를 따라 time step 마다 반복 Prob) 긴 문장 등. long-term dependency 기억 어려움 → vanishing / exploding gradient

sol) LSTM, GRU,

NN : 고정된 vector 를 입력하고 (모든) node 를 완전 연결 & 공간/순서 무시

13.1 Overview of NN Architectures

- Linear Layer (Affine Transformation + Activation)

$$z^{(l+1)} = \phi(A^l z^{(l)} + b^l) \quad A^l: \text{weights matrix} \\ b^l: \text{bias vector}$$

- Convolutional layer:

is organized in a tensor

image dimension

$$w_l \times h_l \times c_l \rightarrow w_{l+1} \times h_{l+1} \times c_{l+1}$$

width height channel

- Attention layer:

act on a sequence of n data vectors of dimension d .

13.2 Optimization of Neural Network

loss:

$$\hat{R}(\theta) = \frac{1}{n} \sum_{i=1}^n \text{loss}(y_i, f_{\theta}(x_i))$$

- Classification 이면 cross-entropy,

regression 이면 MSE 사용

목표 = $\hat{R}(\theta)$ 값을 줄인다.

* Problem:

Neural Network의 최적화는 non-linear & non-convex 문제

여러개의 local minima가 존재 가능

→ 해결방법

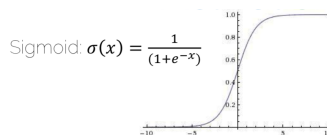
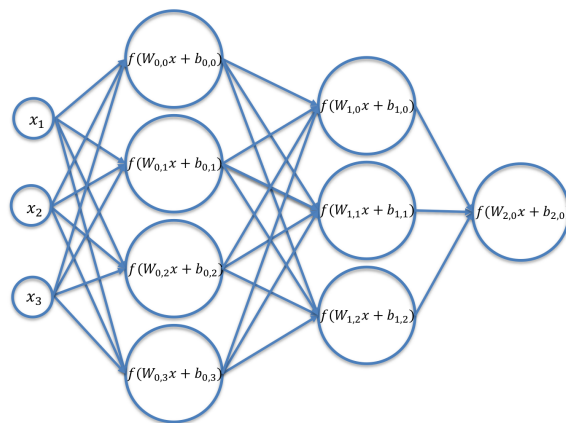
$$\text{Gradient Descent} : \theta^{(t+1)} = \theta^{(t)} - \eta \cdot \nabla_{\theta} \hat{R}(\theta^{(t)})$$

(cf) 보통 SGD 또는 Adam 등 변형된 방법 사용

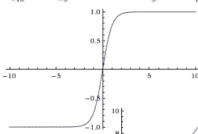
* 실전에서는 convex가 아니지만 NN 구조와 초기화가 잘 설계되면 global minima에 수렴 가능 (empirical risk가 0이 되는 경우)

eg) 2-layer NN:

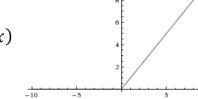
$$f_{\theta}(x) = \sigma(A_1 \cdot \text{ReLU}(A_0 x + b_0))$$



tanh: $\tanh(x)$

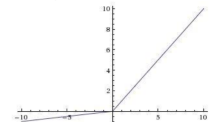


ReLU: $\max(0, x)$



Prof. Dai

Leaky ReLU $\max(0.1x, x)$



Parametric ReLU $\max(\alpha x, x)$

Maxout $\max(w_1^T x + b_1, w_2^T x + b_2)$

ELU $f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{if } x \leq 0 \end{cases}$

Prof. Dai

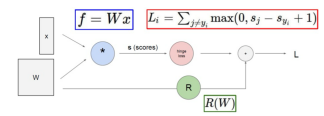
13.3 Automatic Differentiation & Backpropagation

$\nabla_{\theta} \hat{R}(\theta)$ 를 최소화할 때 (minimizing empirical risk) 4가지 방법 있음

방법	설명	한계
by hand	수학적 직접 유도 (기울기 직접 계산 등)	비현실적
Symbolic diff	Maple, Mathematica 같은 소프트웨어 사용	속도 저감
Numeric diff	미분 $\frac{f(x+h)-f(x-h)}{2h}$ 으로 계산	부정확
Auto diff	연산 그래프 기반 (PyTorch, Tensorflow, JAX)	정확, 효율

2.2 Backpropagation

- Calculate gradients numerically if function becomes too complex for analytical solutions
- Analytically solutions would at least partly be possible, but can't e.g. be parallelized
- Possible problems and solution:
 - Solution for multiple outputs in a compute node: take average
 - Solution for loops: No loops in computational graphs
 - Values of forward pass required for backward pass: Cache results of forward pass to reach faster runtime in backward pass
- Computational graph for Linear activation with hinge loss and regularization:



* Backpropagation

1) feedforward

$$z_1 = x_1 w_1 + x_2 w_2 = 0.5 \cdot 0.7 + 0.3 \cdot 0.4 = 0.47$$

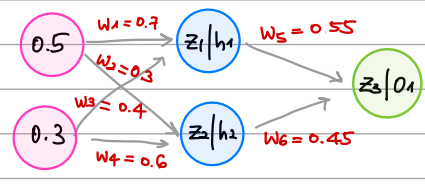
$$h_1 = \text{sigmoid}(z_1) = \frac{1}{1 + e^{-0.47}} = 0.615$$

$$z_2 = x_1 w_3 + x_2 w_4 = 0.5 \cdot 0.3 + 0.3 \cdot 0.6 = 0.33$$

$$h_2 = \text{sigmoid}(z_2) = \frac{1}{1 + e^{-0.33}} = 0.582$$

$$z_3 = h_1 w_5 + h_2 w_6 = 0.615 \cdot 0.55 + 0.582 \cdot 0.45 = 0.6$$

$$o_1 = \text{sigmoid}(z_3) = \frac{1}{1 + e^{-0.6}} = 0.645$$



2) $L = \frac{1}{n} \sum_i (y_i - \hat{y}_i)^2 \Rightarrow n=1: L = \frac{1}{1} (1 - 0.645)^2 = 0.126$

df...

$$s(x) = \frac{1}{1 + e^{-x}} \rightarrow \frac{e^{-x}}{(1 + e^{-x})^2} \approx s(x)(1 - s(x))$$

3) Chain rule로 계산

$$L = \frac{1}{n} \sum_i (y_i - \hat{y}_i)^2 \rightarrow n=1: L = (y - o_1)^2 \quad \frac{\partial L}{\partial o_1} = -2(y - o_1) = -2(1 - 0.645) = -0.71$$

뉴턴 Sigmoid의 정답: $o_1 = \text{sigmoid}(z_3)$ $\frac{\partial o_1}{\partial z_3} = s(z_3)(1 - s(z_3)) = o_1(1 - o_1) = 0.229$

$$\frac{\partial L}{\partial w_5} = \frac{\partial z_3}{\partial w_5} \cdot \frac{\partial o_1}{\partial z_3} \cdot \frac{\partial L}{\partial o_1}$$

$$z_3 = h_1 w_5 + h_2 w_6 \quad \frac{\partial z_3}{\partial w_5} = h_1 = 0.615$$

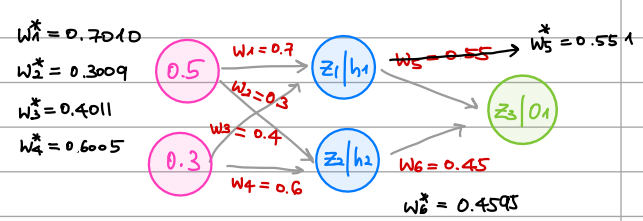
$$\rightarrow \frac{\partial L}{\partial w_5} = \frac{\partial L}{\partial o_1} \cdot \frac{\partial o_1}{\partial z_3} \cdot \frac{\partial z_3}{\partial w_5} = -0.71 \cdot 0.229 \cdot 0.615 = -0.01$$

SGD: $\theta^{(t+1)} = \theta^t - \eta \cdot \nabla \hat{R}(\theta)$

$$w^* = w^5 - \eta \cdot \frac{\partial L}{\partial w_5}$$

$$= 0.55 - 0.1 \cdot (-0.01) = 0.551$$

이리까지



* NN의 오차가 줄었는지 확인.

$$z_1 = x_1 w_1 + x_2 w_3 = 0.5 \cdot 0.7010 + 0.3 \cdot 0.4011 = 0.4708$$

$$h_1 = \frac{1}{1 + e^{-0.4708}} = 0.6156$$

$$z_2 = x_1 w_2 + x_2 w_4 = 0.5 \cdot 0.3009 + 0.3 \cdot 0.6005 = 0.3306$$

$$h_2 = \frac{1}{1 + e^{-0.3306}} = 0.5819$$

$$z_3 = h_1 w_5 + h_2 w_6 = 0.6156 \cdot 0.551 + 0.5819 \cdot 0.4595 = 0.6066$$

$$o_1 = \frac{1}{1 + e^{-0.6066}} = 0.6472$$

$L = \frac{1}{n} \sum_i (y_i - \hat{y}_i)^2 \rightarrow n=1: L = (y - o_1)^2$

$$= (1 - 0.6472)^2 = 0.1245$$

13.4 Vanishing Gradients

↳ 아주 깊은 Network 에서 gradient 가 사라지게 작거나 커짐 → 학습이 안 됨

Vanishing problem
Exploding problem

• Example : $loss(\theta_1, \dots, \theta_L) = (y - x \cdot \theta_1 \cdot \theta_2 \cdot \theta_3 \dots \theta_L)^2$
 L 이 커질수록 (많아질수록) $y - x \prod_{i=1}^L \theta_i$ 형태 $\rightarrow 0$ 또는 무한대

• Solution - Gradient 의 분포를 조절해서 학습을 안정화 시켜자

1) Normalization layers

- Batch Norm : 각 배치마다 $mean = 0, Var = 1$ Regularization (정규화)
- Layer Norm : Layer 단위로 정규화
- Channel Normalization : 각 채널마다 독립적으로 정규화

$$CN(z) = \frac{z - mean(z)}{\sqrt{var(z) + \epsilon}} \gamma + \beta$$

γ, β 는 trainable 하거나 고정된 채널별 scale or shift parameter
 ϵ 는 numerical stability 사용

• If $CN(\theta) = \frac{\theta}{|\theta|}$ 이면 $\rightarrow loss(\theta_1, \dots, \theta_L) = \left(y - x \cdot \theta_1 \frac{\theta_2}{|\theta_2|} \dots \frac{\theta_L}{|\theta_L|}\right)^2$
~~~~~  
벡터  $\theta$  를 norm 1 로  
정규화  
  
 $\underbrace{\theta_2 \dots \theta_L}_{1 \dots 1}$

• Normalization makes  
the deeper layers obsolete  
 $\theta_2 \sim \theta_L$  까지를 규제해버려서

2) Residual Connections (Skip Connection)

각 layer 들을 identity map 에 더 근사시켜자. = 입력을 다음층에도 직접 더하자

$$z_{l+1} = ReLU(A_l z_l + b_l)$$

$$z_{l+1} = z_l + ReLU(A_l z_l + b_l)$$

• ResNet 에서 사용, 이미지 복원의 U-Net 구조에서도 사용

$$loss(\theta_1, \theta_2, \dots, \theta_L) = (y - x \theta_1 (1 + \theta_2) \dots (1 + \theta_L))^2$$

$\theta_2, \dots, \theta_L$  이 1 에 비해  $(1 + \theta_i)$  term 에서 작기때문에  
gradient 가 너무 작거나 너무 크기 많으며 vanishing problem 해소

3) Good initialization

적절한 초기화(initialization)안으로 vanishing problem 해결가능

e.g) 우리  $loss( )$  에서  $\theta_2, \dots, \theta_L$  를 1 근방값으로 초기화하면 residual 구조와 유사  $\approx$  해소