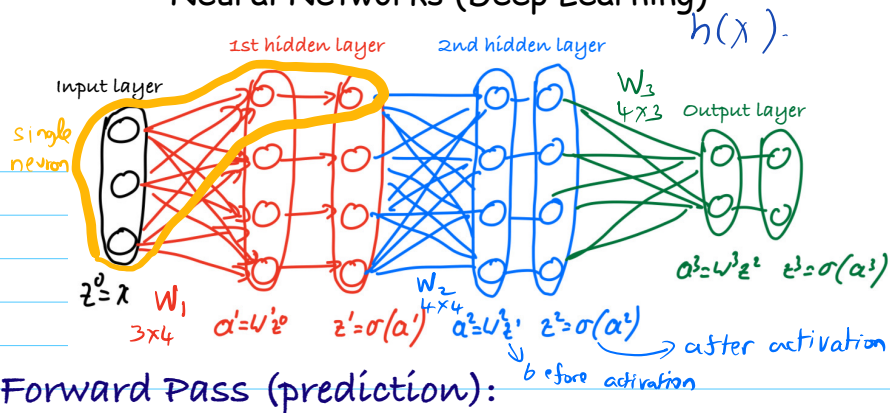


Neural Networks (Deep Learning)



Forward Pass (prediction):

$$z_0 = x \quad (\text{Input Layer})$$

FOR $i = 1$ to L

$$a_i = W_i z_{i-1} + b_i$$

$$z_i = \sigma_i(a_i)$$

END

Return z_L

Typically σ_i is different from others

As usual bias can be aborted in each layer by inserting 1 to each layer as extra feature

To learn the parameters of the NN, we use gradient descent!

How easy is it to compute gradients?

warmup with 1 dimensional example without bias:

$$L(h(x), y) = \frac{1}{2} (h(x) - y)^2$$

$x \in \mathbb{R}$

$w_1, w_2, w_3 \in \mathbb{R}$
 $b_1, b_2, b_3 \in \mathbb{R}$

$$h(x) = w_3 \left(\underbrace{\underbrace{\underbrace{w_2 \sigma(w_1 x + b_1) + b_2}_{z_1}}_{z_2}}_{z_3} \right) + b_3$$

L' derivative of Loss
 σ' derivative of σ

gradients wrt.

w_3

$$\frac{\partial L(h(x), y)}{\partial w_3} = L'(h(x), y) \times \frac{\partial z_3}{\partial w_3} = L'(h(x), y) \times z_2$$

b_3

$$\frac{\partial L(h(x), y)}{\partial b_3} = L'(h(x), y) \times \frac{\partial z_3}{\partial b_3} = L'(h(x), y)$$

Already computed

w_2, b_2

$$\frac{\partial L(h(x), y)}{\partial w_2} = L'(h(x), y) \times \frac{\partial z_3}{\partial w_2} = L'(h(x), y) \sigma'(a_2) \frac{\partial a_2}{\partial w_2}$$

$$= L'(h(x), y) \sigma'(a_2) z_1$$

$$\frac{\partial L(h(x), y)}{\partial b_2} = L'(h(x), y) \sigma'(a_2)$$

Backward pass (gradient step)

elementwise product

$$\vec{a} \circ \vec{f} = \begin{bmatrix} a_1 f_1 \\ a_2 f_2 \\ a_3 f_3 \\ \vdots \end{bmatrix}$$

$$\vec{\delta}_L = \nabla \ell(z_L, y) \circ \sigma'_L(a_L)$$

For $j = L-1$ to 1

$$W_j = W_j - \eta \vec{\delta}_j z_{j-1}^T$$

$$b_j = b_j - \eta \vec{\delta}_j$$

$$\vec{\delta}_{j-1} = \sigma'(a_{j-1}) \circ W_j^T \vec{\delta}_j$$

End

1. Compute a's and z's from forward pass (store them)
2. Run backward pass using a's and z's from forward pass (computes gradients)

To minimize training loss using gradient descent:

$$\text{Training loss} = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i) \quad \text{each one using backprop}$$

$$\nabla \text{Training loss} = \frac{1}{n} \sum_{i=1}^n \nabla \ell(h(x_i), y_i)$$

this can be expensive when n is large

Stochastic Gradient Descent:

Approximate $\nabla \text{Training loss}$ using 1 or $m \ll n$ samples
 $\nabla \approx \nabla \ell(h(x_i), y_i) \quad (x_i, y_i) \sim \mathcal{D}$

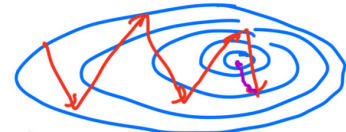
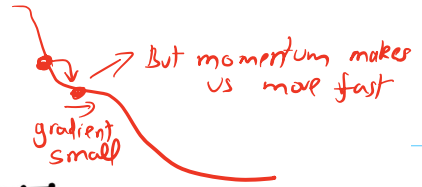
For every layer k , $W_k \leftarrow W_k - \eta \nabla_{W_k} \ell(h(x_i), y_i)$

One update using one or m samples randomly drawn from \mathcal{D}

Important optimization tricks:

- Use Momentum
- Reduce stepsize during optimization
- Use mini-batch (SGD) with $m=64$ inputs
- scale features to be within $[0,1]$
- de-correlate features (PCA)

$$\begin{aligned} G &\leftarrow \beta G + \frac{\partial L}{\partial w} \\ W &\leftarrow W - \alpha G \end{aligned}$$



large step-size
small step-size

Batch-Norm: Re-normalize activations based on current mini-batch.

$$a_i \leftarrow \gamma_i \frac{a_i - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}} + \beta_i$$

↑
parameters that are learned

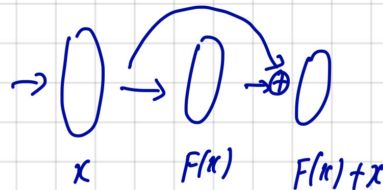
μ_i = mean activation in mini-batch

σ_i^2 = variance in mini-batch

Why? 1. reduces internal covariate shift

2. smooths out loss

Residual Connection:

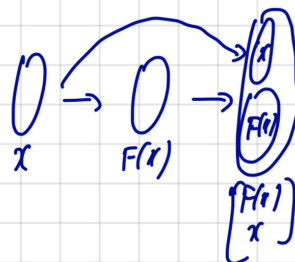


Add input to output activations after each layer*

Why? 1. vanishing gradients problem

2. $\phi(x) = x$ is hard to learn

Dense Connection:



Concatenate input to output activations after each layer*