

CEE 697M: Probabilistic Machine Learning

M2 Deep Neural Networks: Neural Networks for Structured Data I

Jimi Oke



Wed, Apr 5, 2023

- 1 Introduction
- 2 Activation functions
- 3 ANN operations
- 4 Backpropagation
- 5 Summary

Notes

Notes

Neural networks

Consider the linear model:

$$f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{w}^\top \mathbf{x} + \mathbf{b} \tag{1}$$

We can increase the flexibility of the model via a basis function expansion (feature extractor) $\phi(\mathbf{x})$:

$$f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}\phi(\mathbf{x}) + \mathbf{b} \tag{2}$$

If further parameterize $\phi(\mathbf{x})$ by $\boldsymbol{\theta}_2$ for better fitting, we have:

$$f(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{W}\phi(\mathbf{x}; \boldsymbol{\theta}_2) + \mathbf{b} \tag{3}$$

To even further increase complexity, we can recursively fit more feature extractors $f_\ell(\mathbf{x}; \boldsymbol{\theta}_\ell)$:

$$f(\mathbf{x}; \boldsymbol{\theta}) = f_L(f_{L-1}(\cdots f_1(\mathbf{x}; \boldsymbol{\theta}_1)) \cdots) \tag{4}$$

- Each ℓ can be considered a layer in a **feedforward neural network** (FFNN) of L layers.
- Also known as a **multilayer perception** (MLP)
 - When L is large, this is termed a **deep neural network** (DNN)

Notes

Biological neuron

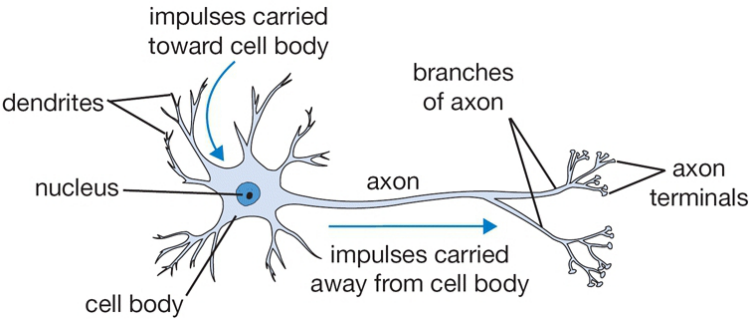


Figure: Biological neuron (Source: <https://cs231n.github.io/neural-networks-1/>)

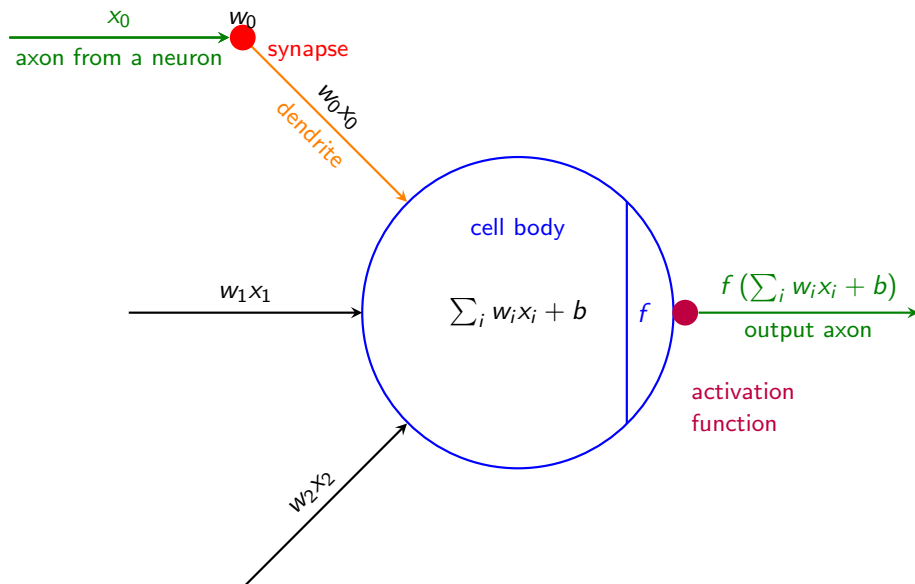
- ~86 billion neurons are found in the human nervous system
- These neurons are connected by 10^{14} to 10^{15} synapses
- Each neuron receives input signals from its dendrites and outputs signals along a single axon
- The axon in turn connects to other neurons via synapses

Notes

[Artificial] Neural networks (ANNs) are modeled as connected layers of neuron in an acyclic graph (no loops).

- ANNs are organized into layers of neurons (or “units”)
- Fully-connected layers are common
- The basic ANN architecture with multiple hidden layers is called the **multilayer perceptron (MLP)**
 - An ANN with only one hidden layer is called the **single layer perceptron**
 - N -layer neural network (number of hidden layers + output layer)
- The output neurons have no activation function. Instead, they perform a final transformation of outputs from the penultimate layer

Notes



Notes

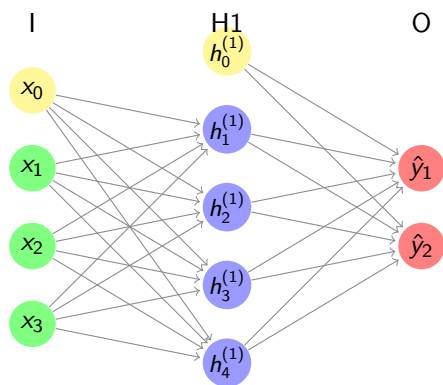
Computational neuron model (cont.)

- x_i : signals traveling along axons (inputs)
- w_i : measure of synaptic strength, which is learned;
 - $w_i > 0 \rightarrow$ excitory influence
 - $w_i < 0 \rightarrow$ inhibitory influence
- Dendrites carry signals $w_i x_i$ to the cell body, where they are summed.
- If the final sum $w_i x_i + b > t$ where t is a threshold¹, the neuron sends a spike along its axon (i.e. fires)
- Computationally, the firing rate of a neuron is represented by an **activation function f**
- The output of a neuron is also called the **activation**

¹intercept b is referred to as the “bias” in ML literature

Notes

Two-layer neural network (with bias neurons)

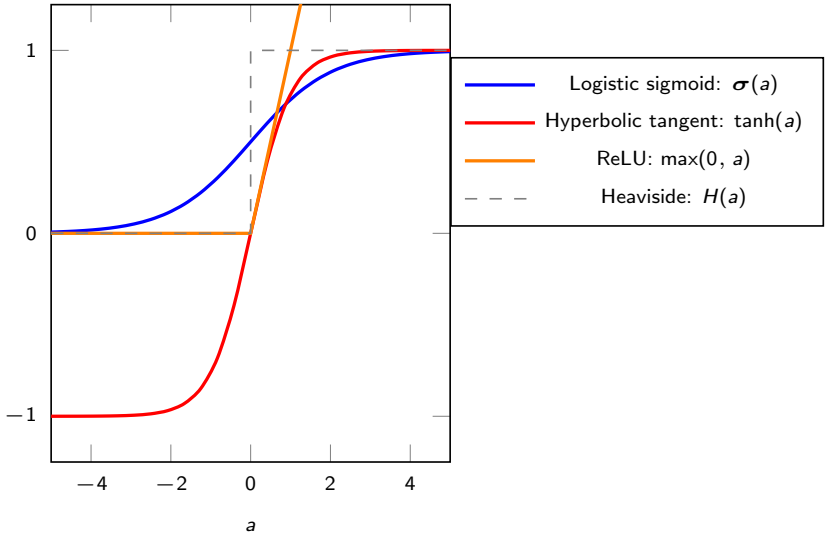


- **Layers:** 2 (input layer not counted); **Hidden layers:** 1
- **Neurons:** 7 (inputs not counted)
- **Learnable parameters:** $(4 \times 4) + (5 \times 2)$; total = 26

Notes

[illegible]

Notes

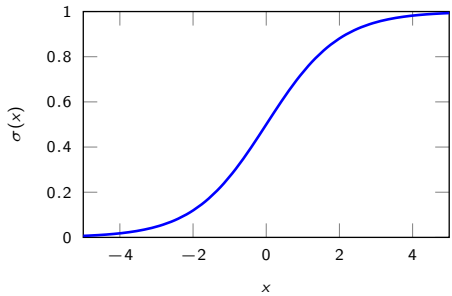


Notes

- The form of the logistic sigmoid function is given by:

$$\sigma(x) = \frac{1}{1 + e^{-x}} \tag{8}$$

- It transforms a real-valued input in the interval $[0, 1]$.



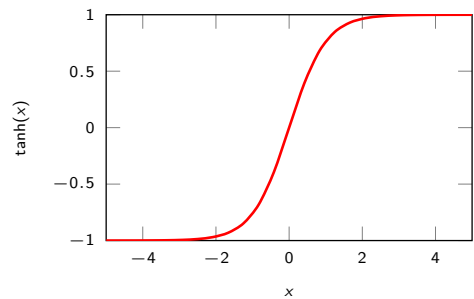
- Historically, it was used as it nicely represents the firing rate
- Recently, it has been superseded by the hyperbolic tangent due to its (a) gradient saturation and (b) non-zero-centeredness.

Hyperbolic tangent (tanh)

- The hyperbolic tangent function is given by:

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = 2\sigma(2x) - 1 \quad (9)$$

- It transforms a real-valued input in the interval $[-1, 1]$.



- Preferred to sigmoid activation function due to its zero-centeredness.

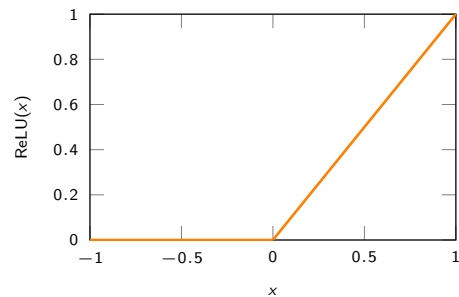
Notes

Rectified linear unit (ReLU)

- The ReLU is given by

$$\text{ReLU}(x) = \max(0, x) \quad (10)$$

- Performs a simple thresholding of input at 0.



- Demonstrates faster convergence than $\sigma(x)$ and $\tanh(x)$
- Popular for deep convolutional networks (several hidden layers)
- Neurons can be fragile, however, requiring care in selection of learning rate

Notes

The sigmoid **activation (output)** of a neuron is denoted:

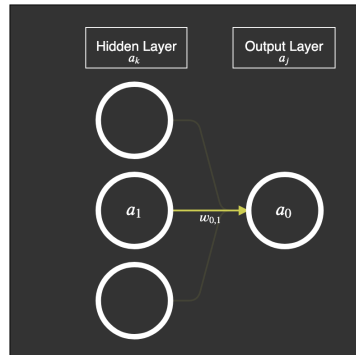
$$\sigma(w_0 a_m + w_1 a_2 + \cdots + w_{m-1} a_{m-1} + b) = \sigma\left(\sum w_i a_i + b\right) = \text{new neuron} \quad (11)$$

Further, we denote each activation as $a_{\text{neuron}}^{(\text{layer})}$, e.g.

- $a_4^{(1)}$: fourth neuron in first layer (layers are counted from first hidden layer)

Weights are denoted as $w_{\text{to,from}}$, e.g.

- $w_{2,3}^2$: from the third neuron in the layer 1 to the second neuron in layer 2
- The superscript is not often used, as it is clear from the context which layer we are dealing with



Notes

Given the activation vector (D neurons) in the zeroth (input) layer:

$$\mathbf{x} \in \mathbb{R}^D = \mathbf{a}^{(0)} = \begin{bmatrix} a_1^0 \\ a_2^0 \\ \vdots \\ a_D^0 \end{bmatrix} \quad (12)$$

Then the activations in the next layer (M neurons) are given by:

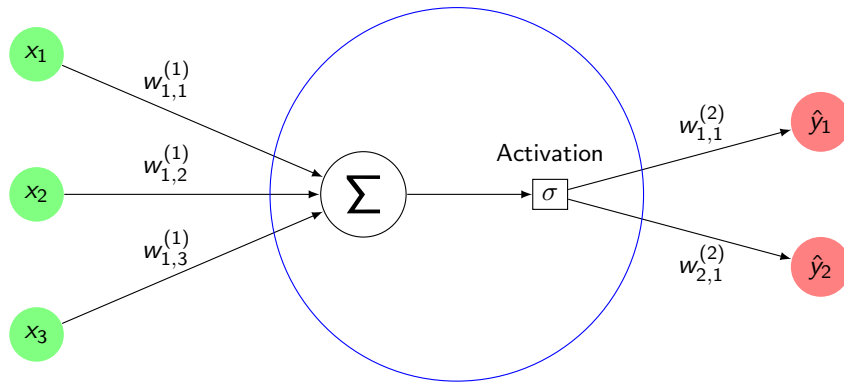
$$\mathbf{a}^{(1)} = \sigma(\mathbf{W}\mathbf{a}^{(0)} + \mathbf{b}) = \sigma\left(\begin{bmatrix} w_{1,1} & w_{1,2} & \cdots & w_{1,D} \\ w_{2,1} & w_{2,2} & \cdots & w_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ w_{M,1} & w_{M,2} & \cdots & w_{M,D} \end{bmatrix} \begin{bmatrix} a_1^0 \\ a_2^0 \\ \vdots \\ a_D^0 \end{bmatrix} + \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_M \end{bmatrix}\right) \quad (13)$$

Example: If Layer 1 had only two neurons, then the weight matrix \mathbf{W} would have only 2 rows.

Notes

Example: MLP with two outputs

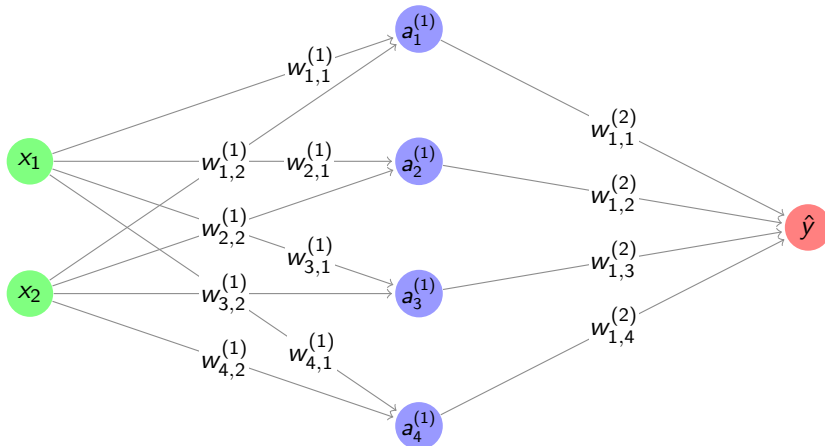
This simple MLP has 2 layers (1 hidden, one outer), and



Notes

Example: 2-layer regression MLP

Two-layer MLP for regression



Notes

Given an observation x_{nd} with $d = 1, \dots, D$ features, these equations describe the output from a 2-layer network:

$$a_m^{(1)} = \sigma \left(\sum_{d=1}^D w_{1,d}^{(1)} x_{nd} + b_d^{(1)} \right)$$

$$y_i(x_i) = \sum_{m=1}^M w_{1,m}^{(2)} a_m^{(1)} + b^{(2)}$$

- D is number of input neurons
- M is number of hidden neurons
- Total number of learnable parameters: $M(D + 1)$ weights and $(D + 1)$ biases
- Linear/identity activation is used in output

Notes

Given K output neurons and N observations (where f_k is the output), we can compute the loss (cost) functions C as follows.

For regression:

$$\mathcal{L} = \sum_{k=1}^K \sum_{n=1}^N (y_{nk} - f_k(x_n))^2 \tag{14}$$

Thus, we can write, where $K = 1$ (univariate output):

$$\mathcal{L} = \sum_{n=1}^N (y_n - \hat{y}_n)^2 \tag{15}$$

For classification, we use the cross-entropy (deviance) given K classes:

$$\mathcal{L} = - \sum_{n=1}^N \sum_{k=1}^K y_{nk} \log f_k(x_n) \tag{16}$$

Notes

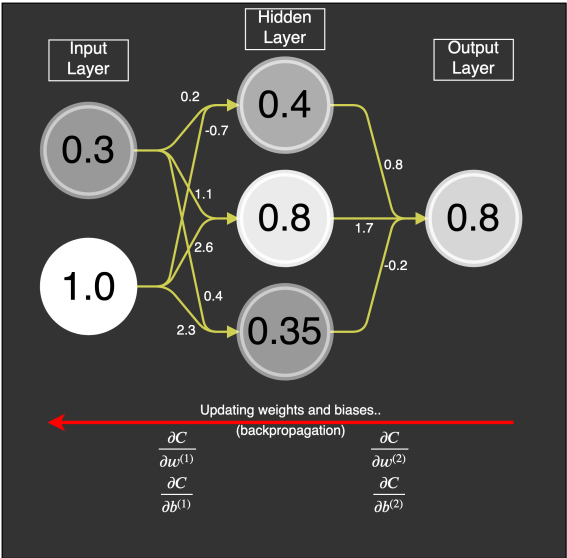
- A neural network is trained or fitted by *learning* the optimal values of the weights (and biases).
- This learning is done via optimization (e.g. gradient descent)
- Gradient descent update:

$$w^{\text{new}} = w^{\text{old}} - \eta \frac{\partial C}{\partial w^{\text{old}}} \tag{17}$$

where:

- η is the learning rate
- C is the cost function (e.g. residual sum of squares)
- w the weight
- In neural networks, the gradients are computed via **backpropagation**

Notes



Notes

Training procedure

- Fix initial weights and perform a forward sweep/pass through the network computing the activations a (outputs) of each layer l as:

$$\mathbf{a}^{(l)} = \sigma(\mathbf{W}\mathbf{a}^{(l-1)} + \mathbf{b}^{(l)}) \quad (18)$$

- At the output layer, we compute the cost function C (what we want to minimize)
- Then, we *backpropagate* the errors through each layer in order to compute the gradients for the weight updates:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(L)}} \frac{\partial \mathbf{a}^{(L)}}{\partial \mathbf{z}^{(L)}} \frac{\partial \mathbf{z}^{(L)}}{\partial \mathbf{w}^{(L)}} \quad (19)$$

where L is the last layer and $\mathbf{z} = \mathbf{W}\mathbf{a}^{(L-1)} + \mathbf{b}$

- Repeat the forward and backward passes until cost is sufficiently minimized

Notes

Equation summary: outer layer (regression case)

At the outer layer L (without indexing by neuron):

$$\mathbf{z}^{(L)} = \mathbf{w}^{(L)} \times \mathbf{a}^{(L-1)} + \mathbf{b}^{(L)} \quad (20)$$

$$\mathbf{a}^{(L)} = \mathbf{z}^{(L)} \quad (21)$$

$$C = (\mathbf{a}^{(L)} - \mathbf{y})^2 \quad (22)$$

The gradient of the cost function with respect to $\mathbf{w}^{(L)}$ is:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(L)}} \frac{\partial \mathbf{a}^{(L)}}{\partial \mathbf{z}^{(L)}} \frac{\partial \mathbf{z}^{(L)}}{\partial \mathbf{w}^{(L)}} = 2 (\mathbf{a}^{(L)} - \mathbf{y}) \sigma'(\mathbf{z}^{(L)}) \mathbf{a}^{(L-1)} \quad (23)$$

Thus, we see that this gradient depends on the activation from the previous layer $\mathbf{a}^{(L-1)}$. Also wrt to the bias:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(L)}} \frac{\partial \mathbf{a}^{(L)}}{\partial \mathbf{z}^{(L)}} \frac{\partial \mathbf{z}^{(L)}}{\partial \mathbf{b}^{(L)}} = 2 (\mathbf{a}^{(L)} - \mathbf{y}) \sigma'(\mathbf{z}^{(L)}) (1) \quad (24)$$

Notes

We can then update the weights for the last layer for the next iteration $r + 1$:

$$w^{(L),r+1} = w^{(L),r} - \eta \frac{\partial \mathcal{L}}{\partial w^{(L)}} \quad (25)$$

$$b^{(L),r+1} = b^{(L),r} - \eta \frac{\partial \mathcal{L}}{\partial b^{(L)}} \quad (26)$$

To update the weights for layer $L - 1$, we need to find the gradients $\frac{\partial \mathcal{L}}{\partial w^{(L-1)}}$ and $\frac{\partial \mathcal{L}}{\partial b^{(L-1)}}$.

Using the chain rule again, we write:

$$\frac{\partial \mathcal{L}}{\partial w^{(L-1)}} = \frac{\partial \mathcal{L}}{\partial a^{(L-1)}} \frac{\partial a^{(L-1)}}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial w^{(L-1)}} \quad (27)$$

$$\frac{\partial \mathcal{L}}{\partial b^{(L-1)}} = \frac{\partial \mathcal{L}}{\partial a^{(L-1)}} \frac{\partial a^{(L-1)}}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial b^{(L-1)}} \quad (28)$$

Notes

But we recall that \mathcal{L} is not *explicitly* dependent on $a^{(L-1)}$ as $C = (a^{(L)} - y)^2$. However, it is *implicitly* dependent, since

$$C \propto a^{(L)}, \quad (29)$$

$$a^{(L)} \propto z^{(L)} \quad (30)$$

and

$$z^{(L)} \propto a^{(L-1)} \quad (31)$$

So, we use the chain rule to expand $\frac{\partial \mathcal{L}}{\partial a^{(L-1)}}$ as follows:

$$\frac{\partial \mathcal{L}}{\partial a^{(L-1)}} = \frac{\partial \mathcal{L}}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial z^{(L)}}{\partial a^{(L-1)}} \quad (32)$$

Notes

Backward pass (cont.)

We can then expand the cost function gradient wrt to weights for layer $L - 1$ as:

$$\frac{\partial \mathcal{L}}{\partial w^{(L-1)}} = \frac{\partial \mathcal{L}}{\partial a^{(L-1)}} \frac{\partial a^{(L-1)}}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial w^{(L-1)}} \quad (33)$$

$$= \frac{\partial \mathcal{L}}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial z^{(L)}}{\partial a^{(L-1)}} \frac{\partial a^{(L-1)}}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial w^{(L-1)}} \quad (34)$$

$$\frac{\partial \mathcal{L}}{\partial b^{(L-1)}} = \frac{\partial \mathcal{L}}{\partial a^{(L-1)}} \frac{\partial a^{(L-1)}}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial b^{(L-1)}} \quad (35)$$

$$= \frac{\partial \mathcal{L}}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial z^{(L)}}{\partial a^{(L-1)}} \frac{\partial a^{(L-1)}}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial b^{(L-1)}} \quad (36)$$

Once these gradients are computed, we update the weights for the $(r + 1)$ th iteration using:

$$w^{(L-1),r+1} = w^{(L-1),r} - \eta \frac{\partial \mathcal{L}}{\partial w^{(L-1)}} \quad (37)$$

$$b^{(L-1),r+1} = b^{(L-1),r} - \eta \frac{\partial \mathcal{L}}{\partial b^{(L-1)}} \quad (38)$$

Notes

Summary: forward pass

- ① ($r = 0$): Initialize weights and biases: $w^{(l),0}$, $b^{(l),0}$
- ② Perform forward pass to compute activations:

$$z^{(l),0} = w^{(l),0} \times a^{(l-1),0} + b^{(l),0} \quad (39)$$

$$a^{(l)} = \sigma(z^{(l),0}) \quad (40)$$

At output layer:

$$z^{(L),0} = w^{(L),0} \times a^{(L-1),0} + b^{(L),0} \quad (41)$$

$$a^{(L),0} = \sigma(z^{(L),0}) \quad (42)$$

$$C = (a^{(L),0} - y)^2 \quad (43)$$

Notes

③ Backward pass, outer layer (L):

① Compute gradients:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(L)}} \frac{\partial \mathbf{a}^{(L)}}{\partial \mathbf{z}^{(L)}} \frac{\partial \mathbf{z}^{(L)}}{\partial \mathbf{w}^{(L)}} \quad (44)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(L)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(L)}} \frac{\partial \mathbf{a}^{(L)}}{\partial \mathbf{z}^{(L)}} \frac{\partial \mathbf{z}^{(L)}}{\partial \mathbf{b}^{(L)}} \quad (45)$$

② Update weights:

$$\mathbf{w}^{(L),1} = \mathbf{w}^{(L),0} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(L)}} \quad (46)$$

$$\mathbf{b}^{(L),1} = \mathbf{b}^{(L),0} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(L)}} \quad (47)$$

Notes

③ Backward pass, layer ($L - 1$):

③ Compute gradients:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(L-1)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(L)}} \frac{\partial \mathbf{a}^{(L)}}{\partial \mathbf{z}^{(L)}} \frac{\partial \mathbf{z}^{(L)}}{\partial \mathbf{a}^{(L-1)}} \frac{\partial \mathbf{a}^{(L-1)}}{\partial \mathbf{z}^{(L-1)}} \frac{\partial \mathbf{z}^{(L-1)}}{\partial \mathbf{w}^{(L-1)}} \quad (48)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(L-1)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(L)}} \frac{\partial \mathbf{a}^{(L)}}{\partial \mathbf{z}^{(L)}} \frac{\partial \mathbf{z}^{(L)}}{\partial \mathbf{a}^{(L-1)}} \frac{\partial \mathbf{a}^{(L-1)}}{\partial \mathbf{z}^{(L-1)}} \frac{\partial \mathbf{z}^{(L-1)}}{\partial \mathbf{b}^{(L-1)}} \quad (49)$$

④ Update weights:

$$\mathbf{w}^{(L-1),1} = \mathbf{w}^{(L-1),0} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(L-1)}} \quad (50)$$

$$\mathbf{b}^{(L-1),1} = \mathbf{b}^{(L-1),0} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(L-1)}} \quad (51)$$

Notes

Summary: backward pass—second-to-last hidden layer

③ Backward pass, layer $(L - 2)$:

⑤ Compute gradients:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(L-2)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(L)}} \frac{\partial \mathbf{a}^{(L)}}{\partial \mathbf{z}^{(L)}} \frac{\partial \mathbf{z}^{(L)}}{\partial \mathbf{a}^{(L-1)}} \frac{\partial \mathbf{a}^{(L-1)}}{\partial \mathbf{z}^{(L-1)}} \frac{\partial \mathbf{z}^{(L-1)}}{\partial \mathbf{a}^{(L-2)}} \frac{\partial \mathbf{a}^{(L-2)}}{\partial \mathbf{z}^{(L-2)}} \frac{\partial \mathbf{z}^{(L-2)}}{\partial \mathbf{w}^{(L-2)}} \quad (52)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(L-2)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(L)}} \frac{\partial \mathbf{a}^{(L)}}{\partial \mathbf{z}^{(L)}} \frac{\partial \mathbf{z}^{(L)}}{\partial \mathbf{a}^{(L-1)}} \frac{\partial \mathbf{a}^{(L-1)}}{\partial \mathbf{z}^{(L-1)}} \frac{\partial \mathbf{a}^{(L-2)}}{\partial \mathbf{z}^{(L-2)}} \frac{\partial \mathbf{z}^{(L-2)}}{\partial \mathbf{b}^{(L-2)}} \quad (53)$$

$$(54)$$

⑥ Update weights:

$$\mathbf{w}^{(L-2),1} = \mathbf{w}^{(L-2),0} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(L-2)}} \quad (55)$$

$$\mathbf{b}^{(L-2),1} = \mathbf{b}^{(L-2),0} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(L-2)}} \quad (56)$$

Notes

Summary: backward pass—first hidden layer

③ Backward pass, layer (1):

⑤ Compute gradients:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(1)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(L)}} \frac{\partial \mathbf{a}^{(L)}}{\partial \mathbf{z}^{(L)}} \frac{\partial \mathbf{z}^{(L)}}{\partial \mathbf{a}^{(L-1)}} \frac{\partial \mathbf{a}^{(L-1)}}{\partial \mathbf{z}^{(L-1)}} \cdots \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{a}^{(1)}} \frac{\partial \mathbf{a}^{(1)}}{\partial \mathbf{z}^{(1)}} \frac{\partial \mathbf{z}^{(1)}}{\partial \mathbf{w}^{(1)}} \quad (57)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(1)}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(L)}} \frac{\partial \mathbf{a}^{(L)}}{\partial \mathbf{z}^{(L)}} \frac{\partial \mathbf{z}^{(L)}}{\partial \mathbf{a}^{(L-1)}} \frac{\partial \mathbf{a}^{(L-1)}}{\partial \mathbf{z}^{(L-1)}} \cdots \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{a}^{(1)}} \frac{\partial \mathbf{a}^{(1)}}{\partial \mathbf{z}^{(1)}} \frac{\partial \mathbf{z}^{(1)}}{\partial \mathbf{b}^{(1)}} \quad (58)$$

⑥ Update weights:

$$\mathbf{w}^{(1),1} = \mathbf{w}^{(1),0} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(1)}} \quad (59)$$

$$\mathbf{b}^{(1),1} = \mathbf{b}^{(1),0} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(1)}} \quad (60)$$

Notes

Summary of backpropagation

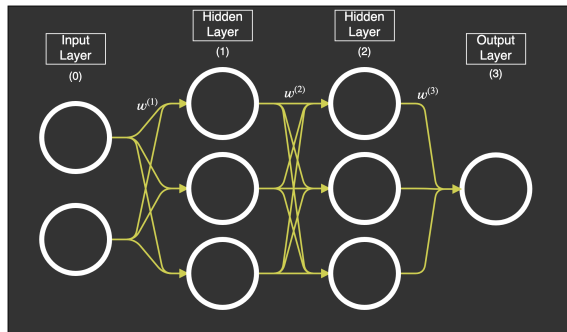
- 1 Fix initial weights $w^{(l),0}$, $b^{(l),0}$ and perform a forward sweep/pass through the network computing the activations a (outputs) of each layer l as:

$$a^{(l)} = \sigma(\mathbf{W}^{(l)}a^{(l-1)} + b^{(l)}) \quad (61)$$

- 2 At the output layer, we compute the cost function C (what we want to minimize)
- 3 Then, we *backpropagate* the errors through each layer in order to compute the gradients $\frac{\partial \mathcal{L}}{\partial w^{(l)}}$, $\frac{\partial \mathcal{L}}{\partial b^{(l)}}$ and weight updates $w^{(l),r+1}$ and $b^{(l),r+1}$
- 4 Repeat the forward and backward passes until cost is sufficiently minimized

Notes

Example: backpropagation for 3-layer network



$$\frac{\partial \mathcal{L}}{\partial w^{(3)}} = \frac{\partial \mathcal{L}}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial w^{(3)}} \quad (62)$$

$$\frac{\partial \mathcal{L}}{\partial w^{(2)}} = \frac{\partial \mathcal{L}}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial w^{(2)}} \quad (63)$$

$$\frac{\partial \mathcal{L}}{\partial w^{(1)}} = \frac{\partial \mathcal{L}}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial w^{(1)}} \quad (64)$$

Notes

Typical hyperparameter values are:

Hyperparameter	Value
# input neurons	1 per input feature
# hidden layers	Usually 1 – 5
# neurons per hidden layer	Usually 10 – 100
# output neurons	1 per prediction dimension
hidden layer activation	ReLU
output activation	None (if unbounded)
loss function	MSE or MAE/Huber

Notes

- For classification, input and hidden layers are chosen in similar fashion to the regression case
- However, the number of output neurons is given by the name of classes/labels
- The output layer activation is typically the softmax function:

$$\text{softmax}(z_k) = \frac{e^{z_k}}{\sum_{k'} e^{z_{k'}}$$

(65)

where z_k is the unnormalized log probability of each class k

- The loss function is taken as the cross entropy

Notes

The standard ANN architecture (MLP) we have studied is also called the feed-forward network.

Other architectures have been shown to give better performance for various applications:

- Recurrent neural networks (RNNs): time-series forecasting
- Convolutional neural networks (CNNs): image classification
- Long short-term memory networks (LSTMs): time-series, pattern identification, etc.

Notes

We will discuss the CNN on Wednesday, along with examples in Python.

- **PMLI**: 13.1-3
- **PML**: 8.3, 9.4
- **ESL**: 11
- **DL**: 6-8, 11, 12
- Experiment in this [playground](#)

Notes
