

Backpropagation for sequential neural networks with densely connected layers

Pawel Wocjan

February 20, 2019

Abstract

We introduce neural networks with densely connected layers and describe the backpropagation algorithm for efficiently training these networks.

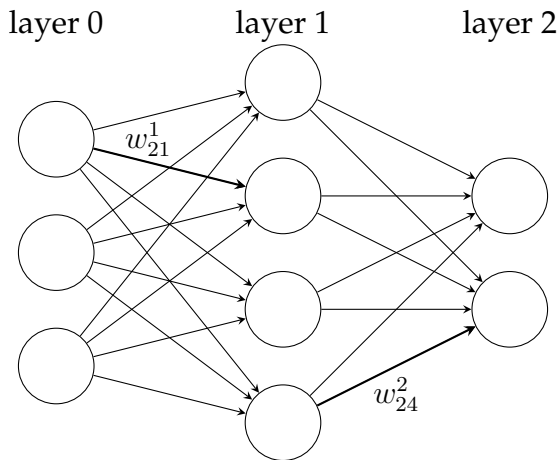
1 Forward propagation

We consider a sequential neural network consisting of the L densely connected layers of sizes n^ℓ , where $\ell \in \{0, \dots, L - 1\}$ is used to enumerate the layers.

For $\ell \geq 1$, we denote the weight of the connection from the k th neuron in the $(\ell - 1)$ th layer to the j th neuron in the ℓ th layer by

$$w_{jk}^\ell \tag{1}$$

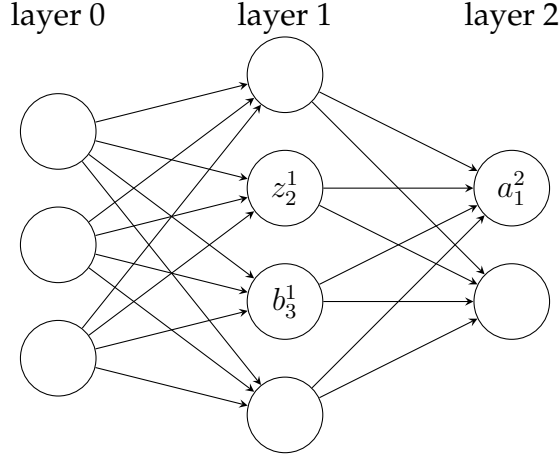
The diagram below shows two examples. The considered connections are drawn thicker.



We use a similar notation for the network's biases and activations. Let

$$b_j^\ell, \quad z_j^\ell, \quad a_j^\ell \quad (2)$$

denote the bias, weighted input, and activation of the j th neuron in the ℓ th layer, respectively.



To keep the discussion general, we use $f^{[\ell]}$ to denote the activation function of the neurons in the ℓ th layer. This is because the activation functions for different layers do not always have to be equal.¹

The forward propagation is described by the following formulas in index notation. We use $k \in \{0, \dots, n^{[0]} - 1\}$ to enumerate all neurons in the 0th layer, which is the input layer. For each neuron k of the input layer, we set its activation to

$$a_k^{[0]} = x_k \quad (3)$$

where the value x_k is the k th entry of feature vector $\mathbf{x} \in \mathbb{R}^{n^{[0]}}$ that is input to the network. (Obviously, the input neurons do not have any weights, biases, or weighted inputs.)

We use $j \in \{0, \dots, n^\ell - 1\}$ to enumerate all neurons of the ℓ layer, where $\ell \geq 1$. For each neuron of the layer ℓ , we have

$$z_j^\ell = \sum_{k=0}^{n^{\ell-1}-1} w_{jk}^\ell \cdot a_k^{\ell-1} + b_j^\ell \quad (4)$$

$$a_j^\ell = f^\ell(z_j^\ell) \quad (5)$$

where the sum is taken over all neurons in the $(\ell - 1)$ th layer.

¹We exclude the softmax activation function for now. We only consider activation functions that are applied by each neuron independently. These include sigmoid, tangent hyperbolicus, relu activation functions. We treat the softmax activation function later.

To rewrite the above formulas in a matrix form we define a weight matrix

$$W^\ell = (w_{jk}^\ell) \in \mathbb{R}^{n^{\ell-1} \times n^\ell} \quad (6)$$

for each layer ℓ . Similarly, we define the bias vector $\mathbf{b}^\ell \in \mathbb{R}^\ell$ and the activation vector $\mathbf{a}^\ell \in \mathbb{R}^\ell$ for each layer ℓ .

The forward propagation in matrix notation is as follows. For the input layer, we set

$$\mathbf{a}^0 = \mathbf{x} \quad (7)$$

For the ℓ th layer, where $\ell \geq 1$, we have

$$\mathbf{z}^\ell = W^\ell \mathbf{a}^{\ell-1} + \mathbf{b}^\ell \quad (8)$$

$$\mathbf{a}^\ell = f^\ell(\mathbf{z}^\ell) \quad (9)$$

2 Loss function

The goal of backpropagation is to compute the partial derivatives

$$\frac{\partial \mathcal{L}}{\partial w_{jk}^\ell} \quad \text{and} \quad \frac{\partial \mathcal{L}}{\partial b_j^\ell} \quad (10)$$

with respect to any weight and bias of the network. We need these partial derivatives to be able to apply stochastic gradient descent.

Let $\mathbf{x} \in \mathbb{R}^{n^0}$ be a feature vector and $\mathbf{y} \in \mathbb{R}^{n^{L-1}}$ its corresponding target vector.² For obvious reasons, we need to assume that the loss function \mathcal{L} is a function of the target vector \mathbf{y} and the output of the neural network \mathbf{a}^{L-1} that is produced when \mathbf{x} is the input. Also, we assume that the loss for a mini-batch is given by the mean of the losses for the individual examples of the mini-batch.

3 Backpropagation

Assume that a friendly demon sits in the j th neuron in layer ℓ . As the weighted input \mathbf{z}_j^ℓ is formed, the demon messes with the neuron's operation. It adds a little change $\Delta \mathbf{z}_j^\ell$ to \mathbf{z}_j^ℓ , so the neuron outputs $f^\ell(\mathbf{z}_j^\ell + \Delta \mathbf{z}_j^\ell)$ instead of $f^\ell(\mathbf{z}_j^\ell)$. This change is propagated through the subsequent layers, finally causing the loss to change by an amount

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}_j^\ell} \cdot \Delta \mathbf{z}_j^\ell. \quad (11)$$

This demon is good, and it is trying to decrease the loss, that is, to find a $\Delta \mathbf{z}_j^\ell$ that makes the loss smaller.

²For instance, this could be the one-hot-encoding vector.

- Suppose the partial derivative $\partial\mathcal{L}/\partial z_j^\ell$ is a large value (either positive or negative). Then the demon can decrease the loss substantially by choosing Δz_j^ℓ to have the opposite sign to $\partial\mathcal{L}/\partial z_j^\ell$.
- In contrast, suppose that the partial derivative $\partial\mathcal{L}/\partial z_j^\ell$ is close to zero. Then the demon cannot decrease the loss much by perturbing the weighted input z_j^ℓ . So far the helpful demon can tell, the neuron is already pretty near optimal. So there is a heuristic sense in which the partial derivative $\partial\mathcal{L}/\partial z_j^\ell$ is a measure for the error in the neuron.

We define the delta (error) δ_j^ℓ of the j th neuron in the ℓ th layer by

$$\delta_j^\ell = \frac{\partial\mathcal{L}}{\partial z_j^\ell}. \quad (12)$$

We also define δ^ℓ the vector of all errors associated with layer ℓ by

$$\delta^\ell = \nabla_{\mathbf{z}^\ell} \mathcal{L}. \quad (13)$$

We will see backpropagation gives us an efficient way of calculating δ^ℓ for every layer.

3.1 Pawel's derivation

We want to show that

$$\delta^{\ell-1} = (W^\ell)^T \delta^\ell \circ g^{\ell-1}(z^{\ell-1}) \quad (14)$$

where \circ denotes the Hadamard product and $g^{\ell-1}$ denotes the derivative of $f^{\ell-1}$. Note that $g^{\ell-1}$ is applied elementwise.

First, we have

$$\frac{\partial\mathcal{L}}{\partial a_k^{\ell-1}} = \sum_j \frac{\partial\mathcal{L}}{\partial z_j^\ell} \cdot \frac{\partial z_j^\ell}{\partial a_k^{\ell-1}} \quad (15)$$

$$= \sum_j \delta_j^\ell \cdot w_{jk}^\ell \quad (16)$$

and in matrix notation

$$\nabla_{\mathbf{a}^{\ell-1}} \mathcal{L} = (W^\ell)^T \cdot \nabla_{\mathbf{z}^\ell} \mathcal{L} \quad (17)$$

$$= (W^\ell)^T \cdot \delta^\ell. \quad (18)$$

Second, we have

$$\frac{\partial L}{\partial z_k^{\ell-1}} = \frac{\partial\mathcal{L}}{\partial a_k^{\ell-1}} \cdot \frac{\partial a_k^{\ell-1}}{\partial z_k^{\ell-1}} \quad (19)$$

and in matrix notation

$$\delta^{\ell-1} = \nabla_{\mathbf{a}^{\ell-1}} \mathcal{L} \circ g^{\ell-1}(\mathbf{z}^{\ell-1}) \quad (20)$$

3.2 An equation for the delta in the output layer

The components of δ^{L-1} are given by

$$\delta_j^{L-1} = \frac{\partial \mathcal{L}}{\partial a_j^{L-1}} \cdot \frac{\partial a_j^{L-1}}{\partial z_j^{L-1}} \quad (21)$$

$$= \frac{\partial \mathcal{L}}{\partial a_j^{L-1}} \cdot g^{L-1}(z_j^{L-1}), \quad (22)$$

where g^{L-1} denotes the derivative of the activation function f^{L-1} .

Note that everything is easily computed. In particular, we compute the weighted input z_j^{L-1} during forward propagation, and it's only a small additional cost to compute $g^{L-1}(z_j^{L-1})$. The exact form of $\partial \mathcal{L} / \partial a_j^{L-1}$ depends on the particular loss function used. However, since the loss function is known, it should be easy to compute $\partial \mathcal{L} / \partial a_j^{L-1}$.

The above equation is in index notation. However, it is easy to rewrite the equation in matrix form as follows

$$\boldsymbol{\delta}^{L-1} = \nabla_{\mathbf{a}^{L-1}} \mathcal{L} \circ g(\mathbf{z}^{L-1}), \quad (23)$$

where $\nabla_{\mathbf{a}^{L-1}} \mathcal{L}$ is defined to be the vector whose entries are the partial derivatives $\partial \mathcal{L} / \partial a_j^{L-1}$ and \circ denotes the Hadamard product³ of two vectors.

3.3 An equation for the delta in terms of the successor's delta

³The Hadamard product (also called Schur product) is the entrywise product.