

Machine Perception - Pen & Paper Solutions

Backpropagation

1 Multivariable chain rule

The gradient $\nabla_{\mathbf{x}} \mathbf{f}$ of a function $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{f} : \mathbf{x} \mapsto y$ is defined as $\nabla_{\mathbf{x}} \mathbf{f} = [\frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \dots, \frac{\partial y}{\partial x_n}]^T \in \mathbb{R}^n$. The Jacobian $J(g)_{\mathbf{x}}$ of a function $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $\mathbf{g} : \mathbf{x} \mapsto \mathbf{y}$ is defined as a matrix of size $\mathbb{R}^{m \times n}$, where each element $J(g)_{\mathbf{x},ij}$ is defined as $J(g)_{\mathbf{x},ij} = \frac{\partial y_i}{\partial x_j}$. Therefore, the *transposed* gradient can be regarded as a special case of the Jacobian, where $m = 1$. The chain-rule in one dimension, for functions $f, g : \mathbb{R} \rightarrow \mathbb{R}$, $g(x) = y$, $f(y) = z$, is expressed as the following:

$$\frac{\partial f(g(x))}{\partial x} = \frac{\partial f(g(x))}{\partial g(x)} \frac{\partial g(x)}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \quad (1)$$

In higher dimensions, where $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $\mathbf{x} \mapsto \mathbf{y}$ and $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^l$, $\mathbf{y} \mapsto \mathbf{z}$, the chain rule states the following:

$$\frac{\partial z_k}{\partial x_j} = \sum_{i=1}^m \frac{\partial z_k}{\partial y_i} \frac{\partial y_i}{\partial x_j} \quad (2)$$

1.1 Chain rule for Jacobians

Given the definition above, show that $\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \frac{\partial \mathbf{z}}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$. What are the dimensions?

Hint: First derive $\frac{\partial z_k}{\partial \mathbf{x}}$

We have $\frac{\partial z_k}{\partial \mathbf{x}} = [\frac{\partial z_k}{\partial x_1}, \frac{\partial z_k}{\partial x_2}, \dots, \frac{\partial z_k}{\partial x_n}] \in \mathbb{R}^n$ and $\frac{\partial z_k}{\partial x_j} = \sum_{i=1}^m \frac{\partial z_k}{\partial y_i} \frac{\partial y_i}{\partial x_j} = \frac{\partial z_k}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial x_j}$, hence:

$$\frac{\partial z_k}{\partial \mathbf{x}} = [\frac{\partial z_k}{\partial x_1}, \frac{\partial z_k}{\partial x_2}, \dots, \frac{\partial z_k}{\partial x_n}] = [\frac{\partial z_k}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial x_1}, \frac{\partial z_k}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial x_2}, \dots, \frac{\partial z_k}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial x_n}] = [\frac{\partial z_k}{\partial \mathbf{y}} (\frac{\partial \mathbf{y}}{\partial \mathbf{x}})_{:,1}, \frac{\partial z_k}{\partial \mathbf{y}} (\frac{\partial \mathbf{y}}{\partial \mathbf{x}})_{:,2}, \dots, \frac{\partial z_k}{\partial \mathbf{y}} (\frac{\partial \mathbf{y}}{\partial \mathbf{x}})_{:,n}] = \frac{\partial z_k}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$$

Similarly, we have for $\frac{\partial \mathbf{z}}{\partial \mathbf{x}}$ using the above results:

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial z_1}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \\ \frac{\partial z_2}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial z_l}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \end{bmatrix} = \frac{\partial \mathbf{z}}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \in \mathbb{R}^{l \times n}$$

In other words, the chain rule has the same form for higher dimensional functions as in the one-dimensional case.

1.2 Chain rule for graphs

We saw the multivariable chain rule defined for higher dimensional functions. In term of graphs, f, g can be visualized as follows:

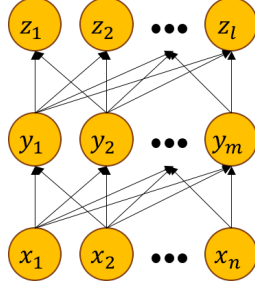
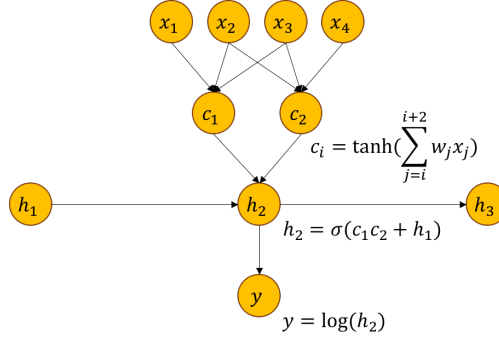


Figure 1: Visualization of the functions f and g as graphs.

Therefore, for a given z_k and x_i , one can imagine the multivariable chain rule as the sum of univariable chain rules over all possible child nodes of x_i that connects it to z_k . Neural networks can be regarded as graphs, where a node outputs values that are a function of its inputs. The input is defined by the edges that connect nodes to each other. Using this, it is possible to define complicated architectures that are still fully differentiable, hence can be optimized.

Let $\sigma = \frac{1}{1+\exp(-x)}$ be the sigmoid activation and $\tanh = \frac{\exp(x)-\exp(-x)}{\exp(x)+\exp(-x)}$. Given the following graph:



a) Derive $\frac{\partial y}{\partial x_2}$

We define $(X)_{:,i}$ as selecting the i -th column of the matrix X . Using the chain rule, we get:

$$\frac{\partial y}{\partial x_2} = \frac{\partial y}{\partial h_2} \left(\frac{\partial h_2}{\partial c_1} \frac{\partial c_1}{\partial x_2} + \frac{\partial h_2}{\partial c_2} \frac{\partial c_2}{\partial x_2} \right)$$

Where each partial derivative has the following values:

$$\begin{aligned} \frac{\partial}{\partial x} \tanh(x) &= 1 - \tanh^2(x), & \frac{\partial}{\partial x} \sigma(x) &= \sigma(x)(1 - \sigma(x)), & \frac{\partial}{\partial x} \log(x) &= \frac{1}{x} \\ \frac{\partial y}{\partial h_2} &= \frac{1}{h_2}, & \frac{\partial h_2}{\partial c_1} &= h_2(1 - h_2)c_2, & \frac{\partial h_2}{\partial c_2} &= h_2(1 - h_2)c_1, & \frac{\partial c_1}{\partial x_2} &= (1 - c_1^2)w_2, & \frac{\partial c_2}{\partial x_2} &= (1 - c_2^2)w_2 \end{aligned}$$

Putting the partial derivatives together, we get:

$$\frac{\partial y}{\partial x_2} = \frac{1}{h_2} \left(h_2(1 - h_2)c_2(1 - c_1^2)w_2 + h_2(1 - h_2)c_1(1 - c_2^2)w_2 \right) = (1 - h_2)w_2(c_1 + c_2 - c_2c_1^2 - c_1c_2^2)$$

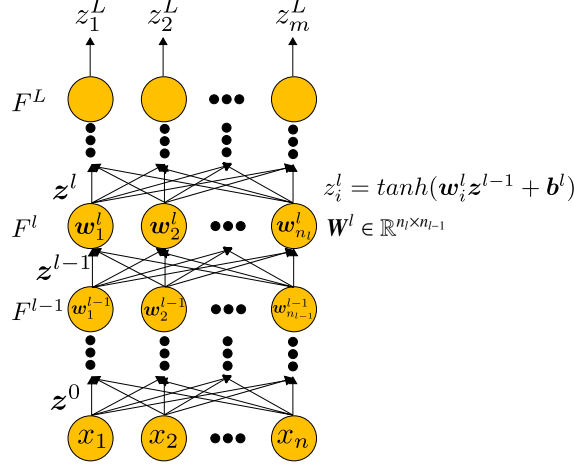


Figure 2: Visualization of multiplayer perceptron with L layers.

2 Backpropagation in multilayer perceptron

A multilayer perceptron, having L layers, can be defined as:

$$F = F^L \circ F^{L-1} \circ \dots \circ F^1 : \mathbb{R}^n \rightarrow \mathbb{R}^m, F(\mathbf{x}) = \mathbf{y} = \mathbf{z}^L \quad (3)$$

Where \circ is the composition operator (i.e $f(g(x)) = f \circ g(x)$), $F^l = \tanh \circ \tilde{F}^l$, where $\tilde{F}^l(\mathbf{x}) = \mathbf{W}^l \mathbf{x} + \mathbf{b}^l$ with $\mathbf{W}^l \in \mathbb{R}^{n_l \times n_{l-1}}$, $\mathbf{b}^l \in \mathbb{R}^{n_l}$, and $\tanh(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$ is the tanh activation function applied element-wise and $\mathbf{z}^0 = \mathbf{x}$. Let $\mathcal{R} : \mathbb{R}^m \rightarrow \mathbb{R}$ be a function on \mathbf{z}^L computing the loss. We define the error signal of layer l as $\delta^l := \frac{\partial \mathcal{R}}{\partial \mathbf{z}^l}$, where $\mathbf{z}^l = F^l(\mathbf{z}^{l-1})$. We left out \mathbf{z}_L from \mathcal{R} to reduce clutter, but still assume its there.

a) For a given layer l , derive δ^l as a function of error signals δ^k and Jacobians $\frac{\partial \mathbf{z}^k}{\partial \mathbf{z}^{k-1}}$

We have $\delta^l = \frac{\partial \mathcal{R}}{\partial \mathbf{z}^l}$ per definition. Expanding this yields:

$$\delta^l = \frac{\partial \mathcal{R}}{\partial \mathbf{z}^l} = \frac{\partial}{\partial \mathbf{z}^l} (\mathcal{R} \circ F^L \circ \dots \circ F^{l+1}(\mathbf{z}^l)) = \delta^L \cdot \frac{\partial \mathbf{z}^L}{\partial \mathbf{z}^{L-1}} \cdot \dots \cdot \frac{\partial \mathbf{z}^{l+1}}{\partial \mathbf{z}^l} = \delta^{l+1} \cdot \frac{\partial \mathbf{z}^{l+1}}{\partial \mathbf{z}^l}$$

b) Derive $\frac{\partial \mathcal{R}}{\partial \mathbf{w}_{ij}^l}$ and $\frac{\partial \mathcal{R}}{\partial \mathbf{b}_i^l}$ as a function of δ^l , \mathbf{W}^l , \mathbf{b}^l and \mathbf{z}^l . Which term was backpropagated from layer $l+1$ to layer l ?

Using the definition of \mathcal{R} , we get:

$$\frac{\partial \mathcal{R}}{\partial \mathbf{w}_{ij}^l} = \frac{\partial}{\partial \mathbf{w}_{ij}^l} (\mathcal{R} \circ F^L \circ \dots \circ F^{l+1}(\mathbf{z}^l)) = \frac{\partial}{\partial \mathbf{w}_{ij}^l} (\mathcal{R} \circ F^L \circ \dots \circ F^{l+1}(\tanh(\mathbf{W}^l \mathbf{x}^{l-1} + \mathbf{b}^l)))$$

Applying the chain rule, we get:

$$\frac{\partial}{\partial \mathbf{w}_{ij}^l} (\mathcal{R} \circ F^L \circ \dots \circ F^{l+1}(\mathbf{z}^l)) = \frac{\partial \mathcal{R}}{\partial \mathbf{z}^l} \frac{\partial \mathbf{z}^l}{\partial \mathbf{w}_{ij}^l} = \delta^l \frac{\partial \mathbf{z}^l}{\partial \mathbf{w}_{ij}^l}$$

Where $\frac{\partial \mathcal{R}}{\partial \mathbf{z}^l} = \delta^l$ as derived in a). To determine $\frac{\partial \mathbf{z}^l}{\partial \mathbf{w}_{ij}^l}$ let's define $\mathbf{z}^l = \tanh(\mathbf{W}^l \mathbf{z}^{l-1} + \mathbf{b}^l) =: \tanh(\mathbf{a})$. Recall that \mathbf{z}^l is a column vector, hence the derivative w.r.t. a scalar must also be a column vector. Using the chain rule we can write.

$$\frac{\partial \mathbf{z}^l}{\partial \mathbf{w}_{ij}^l} = \frac{\partial \tanh(\mathbf{a})}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \mathbf{w}_{ij}^l}$$

The first factor is the Jacobian, which in this case is a diagonal matrix, because when we derive a_j by a_i it will be zero unless $i = j$. Lets denote $\tanh'(x) = 1 - \tanh^2(x)$ and recall that \tanh and its derivative just operate element-wise for vector-valued input. It follows that

$$\left[\frac{\partial \tanh(\mathbf{a})}{\partial \mathbf{a}} \right]_{ij} = \frac{\partial a_i}{\partial a_j} = \tanh'(a_i) \quad \text{if } i = j, 0 \text{ otherwise}$$

and thus

$$\frac{\partial \tanh(\mathbf{a})}{\partial \mathbf{a}} = \begin{bmatrix} \tanh'(a_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \tanh'(a_{n_l}) \end{bmatrix} = \text{diag}(1 - \tanh^2(\mathbf{a})) \quad (4)$$

To determine $\frac{\partial \mathbf{a}}{\partial w_{ij}^l}$ let's expand \mathbf{a} again. Remember from linear algebra that a matrix-vector multiplication is nothing more than a linear combination of the matrix's column vectors with the vector's elements as weights:

$$\mathbf{a} = \mathbf{W}^l \mathbf{z}^{l-1} + \mathbf{b}^l = z_1^{l-1} \mathbf{w}_1^l + \dots + z_j^{l-1} \mathbf{w}_j^l + \dots + z_{n_{l-1}}^{l-1} \mathbf{w}_{n_{l-1}}^l + \mathbf{b}^l$$

Here, \mathbf{w}_j^l refers to the j -th column of \mathbf{W}^l . If we derive this by w_{ij}^l it is clear that the terms where \mathbf{w}_j^l does not occur just equate to 0. Hence we are left with the term $z_j^{l-1} \mathbf{w}_j^l$ which is a vector, and, again, taking the element-wise derivative of those elements that do not contain w_{ij}^l (which are all entries except the i -th row) will result in 0. We write this down as follows:

$$\frac{\partial \mathbf{a}}{\partial w_{ij}^l} = \underbrace{[0, \dots, \overbrace{z_j^{l-1}}^{i\text{-th entry}}, \dots, 0]^T}_{n_l \text{ many entries}} =: \text{vec}_{n_l}^{n_l}(z_j^{l-1}) \quad (5)$$

If you now multiply equations 4 and 5 to obtain the desired gradient you will see that lots of terms cancel out due to the sparseness of both factors. We end up with

$$\frac{\partial \mathbf{z}^l}{\partial w_{ij}^l} = \text{vec}_{n_l}^{n_l}(1 - \tanh^2(a_i)) \cdot z_j^{l-1}$$

which we can also write as

$$\left[\frac{\partial \mathbf{z}^l}{\partial w_{ij}^l} \right]_k = (1 - \tanh^2(\mathbf{a}))_k \cdot z_j^{l-1} \quad \text{if } k = i, 0 \text{ otherwise}$$

Similarly, we get for $\frac{\partial \mathcal{R}}{\partial b_i^l}$:

$$\frac{\partial \mathcal{R}}{\partial b_i^l} = \frac{\partial \mathcal{R}}{\partial \mathbf{z}^l} \frac{\partial \mathbf{z}^l}{\partial b_i^l} = \delta^l \frac{\partial \mathbf{z}^l}{\partial b_i^l}$$

with (differentiating this is very similar to the above derivations):

$$\frac{\partial \mathbf{z}^l}{\partial b_i^l} = \text{vec}_{n_l}^{n_l}(1 - \tanh^2(a_i))$$