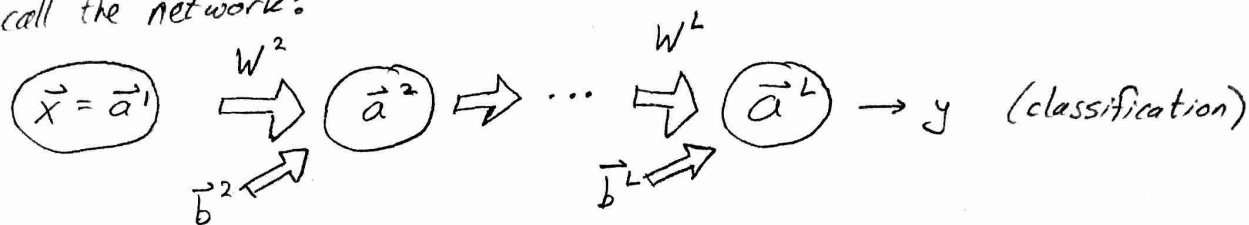


# Backpropagation ("learning") algorithm derivation

11

Recall the network:



We will have the network "learn" by algorithmically adjusting the weights and biases between each layer to make the output closer to the correct classification. Since the output is  $y = \arg \max_i \vec{a}^L$ , we want to make  $a_i^L$  closer to its maximum value 1, for  $i = \text{index of correct classification}$ .

We will do this by minimizing the "error" \* of the last layer of activations  $\vec{a}^L$  from the optimal output  $\hat{e}_i$  (again  $i$  is index of correct classification). The heuristic chosen for error largely determines how effective the learning is; here, we will use Euclidean square distance, in accordance with Nielson's first two chapters:

\* Sometimes also referred to as "loss" or "cost"

$$C = \frac{1}{2} \|\vec{y} - \vec{a}^L\|^2$$

where  $\vec{y} = \hat{e}_i$ . (We use  $C$  to stand for "cost" in accordance with Nielson, though we will often refer to it as "error".)

So we give the network an "example"  $\vec{x}$  with a known classification  $y$  (and corresponding correct last layer  $\vec{y}$ ).

Then adjust weights and biases a bit in the direction that decreases  $C$  (i.e. opposite the direction of the derivative of  $C$  w.r.t the weights and w.r.t the biases).

Thus we need to compute the derivative of  $C$  w.r.t each weight  $w$  and bias  $b$  in the network. It turns out that it's much easier to do this by taking many derivatives at once in parallel - i.e. gradients and derivative matrices. So we will approach it this way.

2

Before we get to the derivation, we have to settle two items: first, I will use a non-standard notation for derivatives, since I believe it makes it more clear what type of derivative it is once you get used to it (e.g. gradient of  $\ell$  w.r.t the biases  $\vec{b}$  in a particular layer). Refer to the Appendix A for the description of this notation.

Second, we must recall a few rules of derivatives and observe some properties specific to our purposes:

### Proposition 1 (Chain Rule)

For a function  $f = f(\vec{x})$ , where each component of  $\vec{x}$  itself depends on a vector  $\vec{u}$  (i.e.  $x_i = x_i(u_1, \dots, u_k)$   $1 \leq i \leq n$ ):

$$\vec{d}_{\vec{u}} f = D_{\vec{u}} \vec{x} \cdot \vec{d}_{\vec{x}} f$$

Recall the reason behind this is that:

$$\vec{d}_{\vec{u}} f = \begin{bmatrix} \partial_{u_1} f \\ \vdots \\ \partial_{u_k} f \end{bmatrix} \quad \text{where} \quad \partial_{u_i} f = \partial_{x_1} f \cdot \partial_{u_i} x_1 + \dots + \partial_{x_n} f \cdot \partial_{u_i} x_n \\ = \vec{d}_{\vec{x}} f \cdot \vec{\partial}_{u_i} \vec{x}$$

Thus,

$$\vec{d}_{\vec{u}} f = \begin{bmatrix} \vec{d}_{\vec{x}} f \cdot \vec{\partial}_{u_1} \vec{x} \\ \vdots \\ \vec{d}_{\vec{x}} f \cdot \vec{\partial}_{u_k} \vec{x} \end{bmatrix} = \begin{bmatrix} \vec{\partial}_{u_1} \vec{x} \\ \vdots \\ \vec{\partial}_{u_k} \vec{x} \end{bmatrix} \cdot \vec{d}_{\vec{x}} f = D_{\vec{u}} \vec{x} \cdot \vec{d}_{\vec{x}} f \quad \square$$

But in a special case (which will come up frequently in this derivation), we can simplify this substantially:

### Proposition 2 ("Parallel" Chain Rule)

Given the same conditions as Prop. 1, except  $\vec{x}$  and  $\vec{u}$  have same length, and each  $x_i$  depends only on the corresponding  $u_i$ , i.e.  $x_i = x_i(u_i)$ :

$$\vec{d}_{\vec{u}} f = \vec{d}_{\vec{x}} f \odot \vec{d}_{\vec{u}} \vec{x}$$

(recall  $\odot$  is the Hadamard Product, and  $\vec{d}_{\vec{u}} \vec{x} = (\partial_{u_1} x_1, \dots, \partial_{u_n} x_n)$ )

pf

Since  $x_i = x_i(u_i)$ , then  $\partial_{u_j} x_i = 0$  for  $i \neq j$ . Thus:

$$\vec{\partial}_{u_i} \vec{x} = \begin{bmatrix} 0 \\ \vdots \\ \partial_{u_i} x_i \\ \vdots \\ 0 \end{bmatrix}$$

and putting together all these columns  $1 \leq i \leq n$ , we form a diagonal matrix:

$$D_{\vec{u}} \vec{x} = \begin{bmatrix} \partial_{u_1} x_1 & & 0 \\ & \ddots & \\ 0 & & \partial_{u_n} x_n \end{bmatrix}$$

So by Prop. 1:

$$\begin{aligned} \vec{d}_{\vec{u}} f &= D_{\vec{u}} \vec{x} \cdot \vec{d}_{\vec{x}} f = \begin{bmatrix} \partial_{u_1} x_1 & & 0 \\ & \ddots & \\ 0 & & \partial_{u_n} x_n \end{bmatrix} \begin{bmatrix} \partial_{x_1} f \\ \vdots \\ \partial_{x_n} f \end{bmatrix} = \begin{bmatrix} \partial_{x_1} f \cdot \partial_{u_1} x_1 \\ \vdots \\ \partial_{x_n} f \cdot \partial_{u_n} x_n \end{bmatrix} \\ &= \vec{d}_{\vec{x}} f \odot \vec{d}_{\vec{u}} \vec{x} \quad \square \end{aligned}$$

In fact, we can extend this statement more generally:

### Proposition 2E (Extended)

Given the same conditions as Prop. 2, except  $\vec{x}$  is a f'n of a matrix  $U$  of values:

$$U = \begin{bmatrix} u_{11} & \cdots & u_{1n} \\ \vdots & \ddots & \vdots \\ u_{m1} & \cdots & u_{mn} \end{bmatrix}$$

where we denote rows by  $\vec{r}_i$ ,  $1 \leq i \leq m$ , and cols  $\vec{c}_j$ ,  $1 \leq j \leq n$ .

If each  $x_i$  depends only on the corresponding row  $\vec{r}_i$ , i.e.  $x_i = x_i(\vec{r}_i)$ , then:

$$\vec{d}_{\vec{c}_j} f = \vec{d}_{\vec{x}} f \odot \vec{d}_{\vec{c}_j} \vec{x}$$

The proof is based on the same idea as in Prop. 2.  
(See next page)

14

Pf

Note that for any component  $x$  and matrix entry  $u$ ,  $\partial_u x \neq 0$  only when  $u$  is in the same row as  $x$ . Furthermore, if we restrict to only derivatives w.r.t. a single column  $\vec{c}$ , then  $u$  must also be in the column  $\vec{c}$ .

$$U = \begin{array}{|c|c|c|} \hline u_{11} & u_{1j} & \vdots \\ \hline \vdots & \vdots & \\ \hline u_{mj} & u_{mn} & \\ \hline \end{array} \quad \vec{x} = \begin{bmatrix} x_1 = x_1(u_{1j}) \\ \vdots \\ x_m = x_m(u_{mj}) \end{bmatrix}$$

col j

Only one element is both in the row  $\vec{r}_i$  corresponding to  $x_i$  and in the column  $\vec{c}_j$  — namely,  $u_{ij}$ . Thus, for each  $x_i$ , only one derivative in  $\vec{d}_{\vec{c}_j} x_i$  will be nonzero;  $\partial_{u_{ij}} x_i$ :

$$\vec{d}_{\vec{c}_j} x_i = \begin{bmatrix} 0 \\ \vdots \\ \partial_{u_{ij}} x_i \\ \vdots \\ 0 \end{bmatrix}$$

So again, when we put these column vectors together for  $1 \leq i \leq m$ , we get a diagonal matrix:

$$D_{\vec{c}_j} \vec{x} = \begin{bmatrix} \partial_{u_{ij}} x_1 & & 0 \\ & \ddots & \\ 0 & & \partial_{u_{mj}} x_m \end{bmatrix}$$

And by Prop. 1:

$$\vec{d}_{\vec{c}_j} f = D_{\vec{c}_j} \vec{x} \cdot \vec{d}_{\vec{x}} f = \begin{bmatrix} \partial_{u_{ij}} x_1 & & 0 \\ & \ddots & \\ 0 & & \partial_{u_{mj}} x_m \end{bmatrix} \begin{bmatrix} \partial_{x_1} f \\ \vdots \\ \partial_{x_m} f \end{bmatrix} = \begin{bmatrix} \partial_{x_1} f \cdot \partial_{u_{ij}} x_1 \\ \vdots \\ \partial_{x_m} f \cdot \partial_{u_{mj}} x_m \end{bmatrix}$$

$$= \vec{d}_{\vec{x}} f \odot \vec{d}_{\vec{c}_j} \vec{x}$$

□

With these propositions established, we can begin the derivation.

So we wish to find  $\partial_b C$  and  $\partial_w C$  for each bias and weight in the network. The expression will depend on which layer those biases and weights correspond to, so we'll start in the last layer (this turns out to be easiest). Also we'll find the entire vector  $\vec{\partial_b} C$  and matrix  $D_w C$  at once in a much cleaner way than finding each component individually. (Until further notice I will use the symbols  $\vec{b}$ ,  $W$ ,  $\vec{a}$ ,  $\vec{z}$ , etc. all referencing those values in the last layer  $L$ .)

Let's begin with something more straight-forward:

$$\vec{\partial_{\vec{a}}} C$$

$$\vec{\partial_{\vec{a}}} C = \vec{\partial_{\vec{a}}} \left( \frac{1}{2} \|\vec{y} - \vec{a}\|^2 \right)$$

If we let  $\vec{\delta} = \vec{y} - \vec{a}$ , then  $\vec{\delta}$  satisfies properties of Prop. 2:

$$\delta_i = \delta_i(a_i) = y_i - a_i$$

So:

$$= \vec{\partial_{\vec{\delta}}} \left( \frac{1}{2} \|\vec{\delta}\|^2 \right) \odot \vec{\partial_{\vec{a}}} \vec{\delta}$$

By converting the first term of this product into a dot product  $\vec{\delta} \cdot \vec{\delta}$  and applying the product rule, we get just  $\vec{\delta}$ . For the second term:

$$\vec{\partial_{\vec{a}}} \vec{\delta} = \vec{\partial_{\vec{a}}} (\vec{y} - \vec{a}) = -\vec{1}$$

(I informally use  $\vec{1}$  to denote the vector of 1's.)

Together, we get:

$$\begin{aligned} \vec{\partial_{\vec{a}}} C &= \vec{\delta} \odot -\vec{1} \\ &= (\vec{y} - \vec{a}) \odot -\vec{1} \\ &= \vec{a} - \vec{y} \end{aligned}$$

Fact 1:  $\vec{\partial_{\vec{a}}} C = \vec{a} - \vec{y}$

6

Recalling that  $\vec{a} = \vec{\sigma}(\vec{z})$ , we can take this to the next step:

$$\begin{aligned} \vec{d}_{\vec{z}} C \\ C = C(\vec{a}(\vec{z})) \end{aligned}$$

Note that by definition,  $\vec{a}$  also satisfies the conditions for Prop. 2, since  $a_i = \sigma(z_i)$ . Thus, we can also parallelize this derivative:

$$\vec{d}_{\vec{z}} C = \vec{d}_{\vec{a}} C \odot \vec{d}_{\vec{z}} \vec{a}$$

Expanding the second term, we get:

$$\vec{d}_{\vec{z}} \vec{a} = \vec{d}_{\vec{z}} \vec{\sigma}(\vec{z}) = \begin{bmatrix} \partial_{z_1} \sigma(z_1) \\ \vdots \\ \partial_{z_n} \sigma(z_n) \end{bmatrix} = \begin{bmatrix} \sigma'(z_1) \\ \vdots \\ \sigma'(z_n) \end{bmatrix}$$

We will label this vector  $\vec{\sigma}'(\vec{z})$ . So we have our next result:

$$\text{Fact 2: } \vec{d}_{\vec{z}} C = \vec{d}_{\vec{a}} C \odot \vec{\sigma}'(\vec{z})$$

Armed with these two facts, we can go for the derivative w.r.t the biases, recalling  $\vec{z} = W\vec{a}' + \vec{b}$  (here I use  $\vec{a}'$  to denote activations in the previous layer, as in the Feedforward derivation).

$$\vec{d}_{\vec{b}} C$$

$$C = C(\vec{z}(\vec{b}))$$

Again, by the parallel nature of the network,  $\vec{z}$  satisfies the conditions of Prop. 2:  $z_i = z_i(b_i) = \sum w a' + b_i$ . So:

$$\vec{d}_{\vec{b}} C = \vec{d}_{\vec{z}} C \odot \vec{d}_{\vec{b}} \vec{z}$$

Examining the second term (since we have the first from Fact 2):

$$\vec{d}_{\vec{b}} \vec{z} = \vec{d}_{\vec{b}} (W\vec{a}' + \vec{b}) = \vec{1}$$

Thus:

$$\text{Fact 3: } \vec{d}_{\vec{b}} C = \vec{d}_{\vec{z}} C$$

And at last, we can go for the derivative w.r.t weights.

7

$D_w C$

We'll tackle the matrix by examining just one column at a time:

$$D_w C = \left[ \vec{d}_{\vec{w}_1} C, \dots, \vec{d}_{\vec{w}_n} C \right]$$

Let's try the general column  $\vec{d}_{\vec{w}_i} C$  using the same techniques:

$$C = C(\vec{z}(w))$$

But this is not quite like the other cases. We cannot parallelize w/ Prop. 2 because each  $z_i = \sum w a' + b_i$  depends on multiple  $w$ 's. Which  $w$ 's though? This sum is a dot product over the  $i$ th row of  $W$  and  $\vec{a}'$ . So each  $z_i$  depends on row  $i$  of  $W$ .

These are the conditions for Prop. 2E, and since we're differentiating w.r.t. a column of  $W$ , we can apply Prop. 2E to use the same technique:

$$\vec{d}_{\vec{w}_i} C = \vec{d}_{\vec{z}} C \odot \vec{d}_{\vec{w}_i} \vec{z}$$

We compute the second term:

$$\begin{aligned} \vec{d}_{\vec{w}_i} \vec{z} &= \vec{d}_{\vec{w}_i} (W \vec{a}' + \vec{b}) \\ &= \vec{d}_{\vec{w}_i} \left( \sum_k a'_k \vec{w}_k + \vec{b} \right) \end{aligned}$$

Differentiating columns  $\vec{w}_k$  for  $k \neq i$  will give  $\vec{0}$ , and likewise  $\vec{b}$  will give  $\vec{0}$ , as none of those columns depend on  $\vec{w}_i$ . However, differentiating the  $\vec{w}_i$  term will give us back the coefficient:

$$= a'_i \vec{I}$$

Putting this together, we have:

Fact 4:  $\vec{d}_{\vec{w}_i} C = a'_i \vec{d}_{\vec{z}} C$

8

To get the full matrix of derivatives, we put together all the columns:

$$D_w C = \begin{bmatrix} \vec{d}_{\vec{w}}, C, \dots, \vec{d}_{\vec{w}_n} C \end{bmatrix}$$

$$= \begin{bmatrix} a'_1 \vec{d}_{\vec{z}} C, \dots, a'_n \vec{d}_{\vec{z}} C \end{bmatrix}$$

This fits the definition of the outer product:  $\vec{u} \otimes \vec{v} = \vec{u} \vec{v}^t = \begin{bmatrix} v_1 \vec{u}, \dots, v_n \vec{u} \end{bmatrix}$

$$= \vec{d}_{\vec{z}} C \otimes \vec{a}'$$

Fact 5:  $D_w C = \vec{d}_{\vec{z}} C \otimes \vec{a}'$

Remember that these all apply to the last layer  $L$ . Here we will recap the formulas, putting back in the usual layer superscript notation:

$$1) \vec{d}_{\vec{a}^L} C = \vec{a}^L - \vec{y}$$

$$2) \vec{d}_{\vec{z}^L} C = \vec{d}_{\vec{a}^L} C \odot \vec{\sigma}'(\vec{z}^L)$$

$$3) \vec{d}_{\vec{b}^L} C = \vec{d}_{\vec{z}^L} C$$

$$4) \vec{d}_{\vec{w}_i^L} C = a_i^{L-1} \vec{d}_{\vec{z}^L} C$$

$$5) D_w C = \vec{d}_{\vec{z}^L} C \otimes \vec{a}^{L-1}$$

Notice that each term (except the first) can be expressed in terms of a previous term. Also note the relationship between terms from this layer  $L$  and the activations from the previous layer  $L-1$ . This suggests that there may be a recurrence relation in general for all layers  $l$ , and perhaps we'll have the same dependence on earlier terms with a different "base" term  $\vec{d}_{\vec{a}^L} C$ . In fact, this is the case, and we will prove it by induction backwards on  $l$ , from  $L$  down to 1. We use the last layer  $L$  and these five facts as our base case, and extend to the general layer  $l$  next.



9

First, we must derive the expression for the new "base" term  $\vec{d}_{\vec{a}} C$  for the second to last layer  $L-1$  to see the pattern, and then we can inductively prove the general case for this term.

Throughout this section, we will adopt the notation using primes to indicate layer numbers. For instance, no primes ( $\vec{a}$ ) denotes the current layer, one prime ( $\vec{a}'$ ) denotes the previous layer, two primes ( $\vec{a}''$ ) denotes two layers before the current layer, and so on. This keeps superscripts clean.

In the following analysis, we treat the last layer as the "current" one.

$\vec{d}_{\vec{a}} C$

$$C = C(\vec{z}(\vec{a}')) \quad \text{recall } \vec{z} = W\vec{a}' + \vec{b}$$

$$\vec{d}_{\vec{a}} C = D_{\vec{a}} \vec{z} \cdot \vec{d}_{\vec{z}} C \quad \text{by Prop. 1 (Chain Rule)}$$

Expanding the first term:

$$D_{\vec{a}} \vec{z} = D_{\vec{a}} (W\vec{a}' + \vec{b}) = W^t$$

The fact that this is transposed  $W$  is not totally obvious. It's a consequence of how we define the  $D$  operator in this context, which is:

$$D_{\vec{u}} \vec{x} = \begin{bmatrix} \vec{d}_{\vec{u}} x_1, \dots, \vec{d}_{\vec{u}} x_n \end{bmatrix} = \begin{bmatrix} \partial_{u_1} x_1 & \dots & \partial_{u_1} x_n \\ \vdots & \ddots & \vdots \\ \partial_{u_m} x_1 & \dots & \partial_{u_m} x_n \end{bmatrix}$$

Working this out substituting  $\vec{a}'$  for  $\vec{u}$  and  $W\vec{a}' + \vec{b}$  for  $\vec{x}$ , we get  $W^t$ . Thus:

$$\vec{d}_{\vec{a}} C = W^t \vec{d}_{\vec{z}} C \quad \text{where } \vec{d}_{\vec{z}} C \text{ is as in Fact 2.}$$

We claim this is a recurrence relation that extends to any choice of "current" layer. In fact, using this relation with layer  $L$  as our "current" layer as a base case, we now assume this holds for any "current" layer  $l \leq L$ . We now show it still holds going back one more layer, between layers  $l-1$  and  $l-2$ :

$$C = C(\vec{z}'(\vec{a}'')) \quad , \quad \vec{z}' = W'\vec{a}'' + \vec{b}'$$

$$\vec{d}_{\vec{a}''} C = D_{\vec{a}''} \vec{z}' \cdot \vec{d}_{\vec{z}'} C \quad \text{as above}$$

$$D_{\vec{a}''} \vec{z}' = D_{\vec{a}''} (W'\vec{a}'' + \vec{b}') = (W')^t$$

Thus:

$$\vec{d}_{\vec{a}''} C = (W')^t \vec{d}_{\vec{z}'} C \quad \square$$

Fact 6:  $\vec{d}_{\vec{a}'} C = W^t \vec{d}_{\vec{z}} C$

Note: This turns out to be the same argument as above, substituting one prime for no primes, and two primes for one.

Having established the base term  $\vec{d}_{\vec{a}} C$  for any layer, we can now verify the other formulas for any layer. Using Facts 2-5 at layer  $L$  as a base case, we assume they hold at any layer  $l \leq L$  and show that they hold at layer  $l-1$  (one prime):

$$\vec{d}_{\vec{z}'} C$$

$$C = C(\vec{a}'(\vec{z}')) \quad , \quad \vec{a}' = \vec{\sigma}(\vec{z}') \quad \text{where } a'_i = \sigma(z'_i)$$

$$\Rightarrow \vec{d}_{\vec{z}'} C = \vec{d}_{\vec{a}'} C \odot \vec{d}_{\vec{z}'} \vec{a}' = \vec{d}_{\vec{a}'} C \odot \vec{\sigma}'(\vec{z}') \quad \square$$

$$\vec{d}_{\vec{b}'} C$$

$$C = C(\vec{z}'(\vec{b}')) \quad , \quad \vec{z}' = W' \vec{a}'' + \vec{b}' \quad \text{where } z'_i = \sum w'_{i''} a''_{i''} + b'_{i'}$$

$$\Rightarrow \vec{d}_{\vec{b}'} C = \vec{d}_{\vec{z}'} C \odot \vec{d}_{\vec{b}'} \vec{z}' = \vec{d}_{\vec{z}'} C \odot \vec{I} = \vec{d}_{\vec{z}'} C \quad \square$$

$$D_{W'} C$$

$$C = C(\vec{z}'(W')) \quad , \quad \vec{z}' = W' \vec{a}'' + \vec{b}'$$

$$D_{W'} C = \left[ \vec{d}_{\vec{w}_1'} C, \dots, \vec{d}_{\vec{w}_n'} C \right] \quad \text{for column vectors } \vec{d}_{\vec{w}_i'} C$$

$$\vec{d}_{\vec{w}_i'} C = \vec{d}_{\vec{z}'} C \odot \vec{d}_{\vec{w}_i'} \vec{z}' \quad \text{by Prop 2E (as in Fact 4's analysis)}$$

$$= \vec{d}_{\vec{z}'} C \odot a_{i''} \vec{I}$$

$$= a_{i''} \vec{d}_{\vec{z}'} C$$

This verifies general case for Fact 4

Finally,

$$D_{W'} C = \left[ a_1'' \vec{d}_{\vec{z}'} C, \dots, a_n'' \vec{d}_{\vec{z}'} C \right]$$

$$= \vec{d}_{\vec{z}'} C \otimes \vec{a}''$$

This verifies general case for Fact 5  $\square$

Thus, we have the general case formulas at current layer  $l$  (no primes):

$$\cdot d_{\vec{a}} C = W^t \vec{d}_{\vec{z}} C \quad , \quad \text{with base case } \vec{d}_{\vec{a}^L} C = \vec{a}^L - \vec{y}$$

$$\cdot \vec{d}_{\vec{z}} C = \vec{d}_{\vec{a}} C \odot \vec{\sigma}'(\vec{z})$$

$$\cdot \vec{d}_{\vec{b}} C = \vec{d}_{\vec{z}} C$$

$$\cdot D_W C = \vec{d}_{\vec{z}} C \otimes \vec{a}'$$

So once the network has been given an example  $\vec{x}$  with its correct classification  $y$ , we feedforward  $\vec{x}$  to obtain the activations  $\vec{a}$  at each layer. Then we use them to compute the derivatives  $\vec{d}_{\vec{b}} C$  and  $D_{\vec{w}} C$  at each layer. Finally, we adjust the bias vector  $\vec{b}$  and weight matrix  $W$  at each layer by a small amount in the direction opposite those derivatives. This "small amount" is computed by scaling  $\vec{d}_{\vec{b}} C$  and  $D_{\vec{w}} C$  by a small factor  $\eta$ , often called the "learning rate" (since it determines how much biases and weights change).

Note:  $\sigma'(x) = \sigma(x)(1 - \sigma(x))$   
We use this substitution in the algorithm below.

Here are the formal algorithms for backpropagation and training.

Backpropagate ( $y, A$ )

- Given weight matrices  $W^2, \dots, W^L$  and bias vectors  $\vec{b}^2, \dots, \vec{b}^L$  for each layer, and learning rate  $\eta$ ,
- Declare vars:
  - $\vec{d}_{\vec{a}} C$ , the derivative of cost w.r.t current layer's activations
  - $\vec{d}_{\vec{z}} C$ , same w.r.t current layer's  $\vec{z}$  values
  - $\Delta \vec{b}$ , change in  $\vec{b}$  at current layer
  - $\Delta W$ , change in  $W$  at current layer
  - $\vec{a}$ , current layer's activations
  - $\vec{a}'$ , previous layer's activations
  - $\vec{y}$ , vector corresponding to classification  $y$

- $\vec{y} = \vec{e}_y$
- $\vec{a}$  = last element of  $A$
- $\vec{d}_{\vec{a}} C = \vec{a} - \vec{y}$
- $\vec{d}_{\vec{z}} C = \vec{d}_{\vec{a}} C \odot (\vec{a} (\vec{I} - \vec{a}'))$
- For each layer  $l$  from  $L$  to  $2$ :
  - $\vec{a} = l^{\text{th}}$  elem. of  $A$
  - $\vec{a}' = (l-1)^{\text{th}}$  elem. of  $A$
  - $\Delta \vec{b} = -\eta \vec{d}_{\vec{z}} C$
  - $\Delta W = -\eta \vec{d}_{\vec{z}} C \otimes \vec{a}'$
  - $\vec{d}_{\vec{a}} C = (W^l)^t \vec{d}_{\vec{z}} C$
  - $\vec{d}_{\vec{z}} C = \vec{d}_{\vec{a}} C \odot (\vec{a}' (\vec{I} - \vec{a}'))$
  - $\vec{b} += \Delta \vec{b}$
  - $W += \Delta W$

Train ( $\vec{x}_1, \dots, \vec{x}_n, y_1, \dots, y_n$ )

- Declare vars:
  - $A$ , the list of activations returned by Feedforward.
- For each  $\vec{x}, y$ :
  - $A = \text{Feedforward}(\vec{x})$
  - $\text{Backpropagate}(y, A)$

# Derivative Notation

/ Appendix A /

1 in, 1 out  $f(x)$

•  $df_x =$  deriv. of  $f$   
wrt.  $x$   
 $\frac{df}{dx}$

$n$  in, 1 out  $f(\vec{x})$

•  $\partial_{x_i} f =$  partial deriv. of  
 $f$  wrt.  $x_i$   
 $\frac{\partial f}{\partial x_i}$

•  $\vec{d}_{\vec{x}} f =$  "gradient" of  $f$ .  
Derivs. wrt. each  
elem. of  $\vec{x}$   
 $\begin{bmatrix} \partial_{x_1} f \\ \vdots \\ \partial_{x_n} f \end{bmatrix}$

1 in,  $m$  out  $\vec{f}(x)$

•  $d_x f_i =$  deriv. of  $f_i$  wrt.  $x$

•  $\vec{d}_x \vec{f} =$  derivs of each elem. of  $\vec{f}$   
wrt.  $x$   
 $\begin{bmatrix} d_x f_1 \\ \vdots \\ d_x f_m \end{bmatrix}$

$n$  in,  $m$  out  $\vec{f}(\vec{x}) = [f_j(x_1, \dots, x_n)]$   $1 \leq j \leq m$

•  $\partial_{x_i} f_j =$  partial deriv. of  $f_j$   
wrt.  $x_i$

•  $\vec{d}_{\vec{x}} f_j =$  "gradient" of  $f_j$ .  
Derivs wrt. each  
elem. of  $\vec{x}$   
 $\begin{bmatrix} \partial_{x_1} f_j \\ \vdots \\ \partial_{x_n} f_j \end{bmatrix}$

•  $\vec{\partial}_{x_i} \vec{f} =$  partial derivs. of each  
elem. of  $\vec{f}$  wrt.  $x_i$   
 $\begin{bmatrix} \partial_{x_i} f_1 \\ \vdots \\ \partial_{x_i} f_m \end{bmatrix}$

•  $D_{\vec{x}} \vec{f} =$  matrix of partial derivs:  
each elem. of  $\vec{f}$  wrt. each  
elem. of  $\vec{x}$

$\begin{matrix} x_1 \\ \vdots \\ x_n \end{matrix} \downarrow \begin{bmatrix} \partial_{x_1} f_1 & \dots & \partial_{x_1} f_m \\ \vdots & \ddots & \vdots \\ \partial_{x_n} f_1 & \dots & \partial_{x_n} f_m \end{bmatrix}$   
 $\xrightarrow{f_1 \quad \quad \quad f_m}$

Some extensions of this

$f(A) = f \left( \begin{matrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{matrix} \right)$  f'n of matrix  $A$

$\rightarrow D_A f =$  derivs of  $f$  wrt. each elem. of  $A$   
 $\begin{bmatrix} \partial_{a_{11}} f & \dots & \partial_{a_{1n}} f \\ \vdots & \ddots & \vdots \\ \partial_{a_{m1}} f & \dots & \partial_{a_{mn}} f \end{bmatrix}$

$\vec{f}(\vec{x}(\vec{u}))$

$\rightarrow D_{\vec{u}} \vec{f} = (D_{\vec{u}} \vec{x}) (D_{\vec{x}} \vec{f})$  by chain rule