

Practical Guide to Matrix Calculus for Deep Learning

Andrew DeLong

andrew.delong@gmail.com

Abstract

Several learning algorithms require computing the gradient of a training objective. This document is a guide to expressing such gradients in *vectorized* form, *i.e.* where inputs, parameters, and intermediate values are all *matrices*. A vectorized gradient expression can be directly implemented in Matlab/Numpy, making use of highly-optimized numerical libraries.

1 A Simple Example

Before reviewing matrix calculus, we give a simple example of what the guide is all about.

Assume we are given t training examples where the n -dimensional inputs are in matrix $\mathbf{X} \in \mathbb{R}^{t \times n}$ and the m -dimensional outputs in matrix $\mathbf{Y} \in \mathbb{R}^{t \times m}$. We can feed *all* the input examples \mathbf{X} through a neural network in matrix form:

$$\text{output} = f(\mathbf{XW} + \mathbf{b}). \quad (1)$$

This network is parameterized by a weight matrix $\mathbf{W} \in \mathbb{R}^{n \times m}$, a bias vector $\mathbf{b} \in \mathbb{R}^{1 \times m}$, and an *activation function* $f(\cdot)$ that is applied element-wise to its input. (Here “+ \mathbf{b} ” is understood to broadcast row-wise.) Row i of the $t \times m$ output matrix corresponds to example i from input \mathbf{X} . Vectorized Matlab code for sending \mathbf{X} through this network might look like:

```
function Z = eval_nnet(X,W,b)
    Z = tanh(bsxfun(@plus,X*W,b));    % f(X*W + b) where f = tanh
end
```

```
>> X = rand(20,2);    Y = rand(20,3);    % t = 20, n = 2, m = 3
>> W = rand(2,3);    b = rand(1,3);
>> Z = eval_nnet(X,W,b);
>> size(Z)
ans =
    20     3                % matrix of 3-dimensional outputs
```

We can train the model by minimizing a standard training objective J such as

$$J(\mathbf{W}, \mathbf{b}) = \frac{1}{2} \|\mathbf{f}(\mathbf{XW} + \mathbf{b}) - \mathbf{Y}\|^2. \quad (2)$$

Here $\|\cdot\|^2$ is understood to be the sum of squares of the matrix elements. This cost function can of course be evaluated in Matlab with code like:

```
function J = eval_cost(X,Y,W,b)
    Z = eval_nnet(X,W,b);
    E = (Z-Y).^2;           % squared errors
    J = 0.5 * sum(E(:));    % sum of squared errors
end
```

To optimize this cost function with gradient descent, we need an expression for the gradient—preferably in vectorized form so that we can program it directly in Matlab.

By straight-forward application of the matrix calculus rules in this guide, the vectorized gradient of J with respect to parameters \mathbf{W} and \mathbf{b} turns out to be

$$\frac{\partial J}{\partial \mathbf{W}} = \mathbf{X}^T((\mathbf{Z} - \mathbf{Y}) \odot f'(\mathbf{XW} + \mathbf{b})) \quad (3)$$

$$\frac{\partial J}{\partial \mathbf{b}} = \mathbf{1}^T((\mathbf{Z} - \mathbf{Y}) \odot f'(\mathbf{XW} + \mathbf{b})) \quad (4)$$

where matrix derivatives $\frac{\partial J}{\partial \mathbf{W}} \in \mathbb{R}^{n \times m}$, $\frac{\partial J}{\partial \mathbf{b}} \in \mathbb{R}^{1 \times m}$, and operator \odot denotes the element-wise product (Hadamard product). The gradient $\nabla J = (\frac{\partial J}{\partial \mathbf{W}}, \frac{\partial J}{\partial \mathbf{b}})$ can be easily computed in Matlab.

```
function [dW,db] = eval_grad(X,Y,W,b)
    Z = eval_nnet(X,W,b); % forward pass; Z = tanh(XW + b)
    D = (Z-Y) .* (1-Z.^2); % backward pass; note that tanh'(x) = 1 - tanh(x)^2
    dW = X'*D;           % X^T D
    db = sum(D);         % 1^T D
end
```

We can check the above gradient by comparing it to the numerical gradient:

```
>> [dW,db] = eval_grad(X,Y,W,b);
>> [nW,nb] = eval_grad_numeric(X,Y,W,b);
>> relerr = @(A,B) max(abs(A(:)-B(:))./(abs(A(:))+abs(B(:))));
>> relerr(dW,nW), relerr(db,nb)
ans =
    1.7391e-010    % very small difference between symbolic gradient dW,db
ans =
    1.5546e-010    % and numeric gradient nW, nb
```

Where the numerical gradient is computed by something like:

```
function [dW,db] = eval_grad_numeric(X,Y,W,b)
    dW = zeros(size(W));
    for i=1:numel(W)
        step = zeros(size(W)); step(i) = 1e-5;
        dW(i) = (eval_cost(X,Y,W+step,b) - eval_cost(X,Y,W-step,b)) / 2e-5;
    end

    db = zeros(size(b));
    for i=1:numel(b)
        step = zeros(size(b)); step(i) = 1e-5;
        db(i) = (eval_cost(X,Y,W,b+step) - eval_cost(X,Y,W,b-step)) / 2e-5;
    end
end
```

2 Matrix Calculus for Learning

Given a feed-forward training objective $J(\cdot)$ in vectorized form, the rules here help to derive the gradient $\nabla J(\cdot)$ in vectorized form. Here \mathbf{A} is a matrix, \mathbf{A}_i the i^{th} row, and \mathbf{A}_{ij} the $(i, j)^{\text{th}}$ element.

Frobenius product. The scalar (dot) product of $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ is the sum $\mathbf{A} \cdot \mathbf{B} = \sum_{ij} \mathbf{A}_{ij} \mathbf{B}_{ij}$.

Hadamard product. The element-wise product of $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ is matrix $(\mathbf{A} \odot \mathbf{B})_{ij} = \mathbf{A}_{ij} \mathbf{B}_{ij}$.

Row-wise product. The row-wise product of $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ is the vector $(\mathbf{A} \ominus \mathbf{B})_i = \mathbf{A}_i \cdot \mathbf{B}_i$.

As usual \mathbf{AB} denotes matrix product. Some simple matrix identities will be useful.

$$\mathbf{A} \cdot (\mathbf{B} \odot \mathbf{C}) = (\mathbf{A} \odot \mathbf{B}) \cdot \mathbf{C} \quad (5)$$

$$\mathbf{A} \cdot (\mathbf{BC}) = (\mathbf{B}^T \mathbf{A}) \cdot \mathbf{C} = (\mathbf{AC}^T) \cdot \mathbf{B} \quad (6)$$

$$\mathbf{1}^T (\mathbf{A} \ominus \mathbf{B}) = \mathbf{A} \cdot \mathbf{B} \quad (7)$$

Some standard differential identities will be useful as well [1]. Here \mathbf{X}, \mathbf{Y} are matrix-valued functions.

$$d\mathbf{C} = \mathbf{0} \quad (\text{if } \mathbf{C} \text{ is constant}) \quad (8)$$

$$d(\alpha \mathbf{X}) = \alpha(d\mathbf{X}) \quad (9)$$

$$d(\mathbf{X}^T) = (d\mathbf{X})^T \quad (10)$$

$$d(\mathbf{X} \pm \mathbf{Y}) = d\mathbf{X} \pm d\mathbf{Y} \quad (11)$$

$$d(\mathbf{XY}) = (d\mathbf{X})\mathbf{Y} + \mathbf{X}(d\mathbf{Y}) \quad (12)$$

$$d(\mathbf{X} \cdot \mathbf{Y}) = (d\mathbf{X}) \cdot \mathbf{Y} + \mathbf{X} \cdot (d\mathbf{Y}) \quad (13)$$

$$d(\mathbf{X} \odot \mathbf{Y}) = (d\mathbf{X}) \odot \mathbf{Y} + \mathbf{X} \odot (d\mathbf{Y}) \quad (14)$$

And now some differential identities involving scalar-to-scalar function $y = f(x)$ and vector-to-scalar function $y = g(\mathbf{x})$. It should therefore be understood that $\mathbf{Y} = f(\mathbf{X})$ is applied element-wise, and that $f'(\mathbf{X})$ is the derivative of f applied element-wise. Likewise $\mathbf{y} = g(\mathbf{X})$ results in a column vector with i^{th} element $y_i = g(\mathbf{X}_i)$, and $\nabla g(\mathbf{X})$ is a matrix with i^{th} row $\nabla g(\mathbf{X}_i)$.

$$\mathbf{Y} = f(\mathbf{X}) \quad \Rightarrow \quad d\mathbf{Y} = f'(\mathbf{X}) \odot d\mathbf{X} \quad (15)$$

$$\mathbf{y} = g(\mathbf{X}) \quad \Rightarrow \quad d\mathbf{y} = \nabla g(\mathbf{X}) \odot d\mathbf{X} \quad (16)$$

Finally, let $y = h(\mathbf{X}^1, \dots, \mathbf{X}^K)$ be a scalar-valued function of several matrices. If dy can then be manipulated into the form below, then each matrix \mathbf{A}^k is the partial derivative with respect to \mathbf{X}^k .

$$dy = \sum_k \mathbf{A}^k \cdot d\mathbf{X}^k \quad \Leftrightarrow \quad \nabla y = \left(\frac{\partial y}{\partial \mathbf{X}^1}, \dots, \frac{\partial y}{\partial \mathbf{X}^K} \right) = (\mathbf{A}^1, \dots, \mathbf{A}^K) \quad (17)$$

3 Feed-Forward Neural Networks

Let integers n_0, \dots, n_k denote the number of units in each layer. A feed-forward network is defined by weight matrices¹ $\mathbf{W} = (\mathbf{W}_1, \dots, \mathbf{W}_k)$ with $\mathbf{W}_j \in \mathbb{R}^{n_{j-1} \times n_j}$, bias vectors $\mathbf{b} = (\mathbf{b}_1, \dots, \mathbf{b}_k)$ with $\mathbf{b}_j \in \mathbb{R}^{1 \times n_j}$, and activation functions f_1, \dots, f_k . We want to evaluate t input examples of dimension n_0 stacked in a matrix $\mathbf{X} \in \mathbb{R}^{t \times n_0}$. Layer j of the network computes a matrix $\mathbf{Z}_j \in \mathbb{R}^{t \times n_j}$ as

$$\mathbf{Z}_j = \begin{cases} \mathbf{X} & j = 0 \\ f_j(\mathbf{A}_j) & j > 0 \end{cases} \quad \text{where } \mathbf{A}_j = \mathbf{Z}_{j-1} \mathbf{W}_j + \mathbf{b}_j. \quad (18)$$

Choose some vector-to-scalar loss function $\ell(\cdot)$ and apply it row-wise to \mathbf{Z}_k in training objective

$$J(\mathbf{W}, \mathbf{b}) = \frac{1}{t} \mathbf{1}^T \ell(\mathbf{Z}_k) + p(\mathbf{W}, \mathbf{b}) \quad (19)$$

where p is a penalty on the model parameters, *e.g.* $\|\cdot\|_1$ or $\|\cdot\|_2$ weight decay.

3.1 Regression Networks

Regression problems typically use *squared-error* loss $\ell(\mathbf{z}) = \frac{1}{2} \|\mathbf{z} - \mathbf{y}\|^2$ where \mathbf{y} is a target. The matrix of loss gradients is simply $\nabla \ell(\mathbf{Z}) = \mathbf{Z} - \mathbf{Y}$ in that case. The final activation $f_k(\cdot)$ of a regression network is typically a scalar-to-scalar element-wise function (*e.g.* linear, sigmoid, ReLU). For simplicity, assume no penalty. Applying the differential operator d to $J(\cdot)$ we have

$$\begin{aligned} dJ &= \frac{1}{t} \mathbf{1}^T d\ell(\mathbf{Z}_k) && \text{by (8),(12)} \\ &= \frac{1}{t} \mathbf{1}^T (\nabla \ell(\mathbf{Z}_k) \odot d\mathbf{Z}_k) && \text{by (16)} \\ &= \frac{1}{t} \nabla \ell(\mathbf{Z}_k) \cdot d\mathbf{Z}_k && \text{by (7)} \\ &= \frac{1}{t} \nabla \ell(\mathbf{Z}_k) \cdot (f'_k(\mathbf{A}_k) \odot d\mathbf{A}_k) && \text{by (15) since } f_k \text{ element-wise} \\ &= \left(\frac{1}{t} \nabla \ell(\mathbf{Z}_k) \odot f'_k(\mathbf{A}_k) \right) \cdot d\mathbf{A}_k && \text{by (5)} \\ &= \Delta_k \cdot d\mathbf{A}_k && (20) \\ &= \Delta_k \cdot (\mathbf{Z}_{k-1} d\mathbf{W}_k + d\mathbf{b}_k + d\mathbf{Z}_{k-1} \mathbf{W}_k) && \text{by (12)} \\ &= \mathbf{Z}_{k-1}^T \Delta_k \cdot d\mathbf{W}_k + \mathbf{1}^T \Delta_k \cdot d\mathbf{b}_k + \Delta_k \mathbf{W}_k^T \cdot d\mathbf{Z}_{k-1} && \text{by (6) gives } \frac{\partial J}{\partial \mathbf{W}_k} \text{ and } \frac{\partial J}{\partial \mathbf{b}_k} \\ &= \dots + \Delta_k \mathbf{W}_k^T \cdot (f'_{k-1}(\mathbf{A}_{k-1}) \odot d\mathbf{A}_{k-1}) && \text{by (15)} \\ &= \dots + (\Delta_k \mathbf{W}_k^T \odot f'_{k-1}(\mathbf{A}_{k-1})) \cdot d\mathbf{A}_{k-1} && \text{by (5)} \\ &= \dots + \Delta_{k-1} \cdot d\mathbf{A}_{k-1} && (21) \\ &= \sum_{j=1}^k \mathbf{Z}_{j-1}^T \Delta_j \cdot d\mathbf{W}_j + \mathbf{1}^T \Delta_j \cdot d\mathbf{b}_j && \text{unroll as (20-21) for } j = k..1 \quad (22) \end{aligned}$$

By (17) we can get complete gradient² $\nabla J = (\frac{\partial J}{\partial \mathbf{W}_1}, \frac{\partial J}{\partial \mathbf{b}_1}, \dots, \frac{\partial J}{\partial \mathbf{W}_k}, \frac{\partial J}{\partial \mathbf{b}_k})$. For completeness, we also add the possibility of a penalty $p(\cdot)$ on the individual parameter matrices.

$$\frac{\partial J}{\partial \mathbf{W}_j} = \mathbf{Z}_{j-1}^T \Delta_j + \frac{\partial p}{\partial \mathbf{W}_j} \quad \text{where } \Delta_j = f'_j(\mathbf{A}_j) \odot \begin{cases} \frac{1}{t} (\mathbf{Z}_k - \mathbf{Y}) & j = k \\ \Delta_{j+1} \mathbf{W}_{j+1}^T & j < k \end{cases} \quad (23)$$

¹Now we use subscript \mathbf{A}_j to denote the j^{th} matrix in a sequence $\mathbf{A}_1, \dots, \mathbf{A}_k$.

²Note that it is slightly more efficient to compute $\mathbf{1}^T \Delta$ by directly summing the rows of Δ .

3.2 Classification Networks

In classification networks, the final activation function $f_k(\cdot)$ is typically a *softmax*:

$$\mathbf{z} = \text{softmax}(\mathbf{a}) = \frac{e^{\mathbf{a}}}{\mathbf{1}^T e^{\mathbf{a}}}. \quad (24)$$

For softmax it makes sense to use the *negative-log likelihood* (NLL) loss $\ell(\mathbf{z}) = -\mathbf{y} \cdot \ln \mathbf{z}$ where \mathbf{y} has a single non-zero entry indicating the target class. Note that $\nabla \ell(\mathbf{z}) = -\mathbf{y} \odot (\mathbf{z}^{-1})$ in this case, and the differential operator applied to \mathbf{z} is

$$d\mathbf{z} = \frac{e^{\mathbf{a}} \odot d\mathbf{a}}{\mathbf{1}^T e^{\mathbf{a}}} - \frac{e^{\mathbf{a}}(e^{\mathbf{a}} \cdot d\mathbf{a})}{(\mathbf{1}^T e^{\mathbf{a}})^2} = \mathbf{z} \odot (d\mathbf{a} - \mathbf{z} \cdot d\mathbf{a}) \quad (25)$$

Applying the differential operator d to $J(\cdot)$ we start as before

$$\begin{aligned} dJ &= \frac{1}{t} \mathbf{1}^T d\ell(\mathbf{Z}_k) && \text{by (8),(12)} \\ &= \frac{1}{t} \mathbf{1}^T (\nabla \ell(\mathbf{Z}_k) \odot d\mathbf{Z}_k) && \text{by (16)} \\ &= \frac{1}{t} \nabla \ell(\mathbf{Z}_k) \cdot d\mathbf{Z}_k && \text{by (7)} \end{aligned}$$

Let us carefully expand the contribution $\nabla \ell(\mathbf{z}) \cdot d\mathbf{z}$ of each row above,

$$\begin{aligned} \nabla \ell(\mathbf{z}) \cdot d\mathbf{z} &= -(\mathbf{y} \odot (\mathbf{z}^{-1})) \cdot (\mathbf{z} \odot (d\mathbf{a} - \mathbf{z} \cdot d\mathbf{a})) && \text{by (25)} \\ &= -\mathbf{y} \cdot (d\mathbf{a} - \mathbf{z} \cdot d\mathbf{a}) \\ &= (\mathbf{z} - \mathbf{y}) \cdot d\mathbf{a} \end{aligned} \quad (26)$$

We can now continue where we left computing dJ . Let \mathbf{Y} be the stack of target vectors, then

$$\begin{aligned} &= \frac{1}{t} (\mathbf{Z}_k - \mathbf{Y}) \cdot d\mathbf{A}_k && \text{by (26)} \\ &= \mathbf{\Delta}_k \cdot d\mathbf{A}_k \end{aligned}$$

The rest of the evaluation is the same as for the regression case. So, with respect to computing the gradient, the only change is to set $\mathbf{\Delta}_k = \frac{1}{t} (\mathbf{Z}_k - \mathbf{Y})$.

3.3 Auto-Encoder with Tied Weights

An auto-encoder is a regression network, and we can assume layer sizes n_0, \dots, n_k satisfy $n_j = n_{k-j}$ and there are an odd number of layers (k is even). Tied weights imply $\mathbf{W}_{k-j+1} = \mathbf{W}_j^T$, so there are $\frac{k}{2}$ unique weight matrices but k (untied) bias vectors $\mathbf{b} = (\mathbf{b}_1, \dots, \mathbf{b}_k)$. To simplify indexing, it helps to define $\tilde{\mathbf{W}} = (\tilde{\mathbf{W}}_1, \dots, \tilde{\mathbf{W}}_k)$ as ‘logical’ matrices, each of which identifies to an ‘actual’ weight matrix from among $\mathbf{W} = (\mathbf{W}_1, \dots, \mathbf{W}_{\frac{k}{2}})$. Specifically, let

$$\tilde{\mathbf{W}}_j \equiv \begin{cases} \mathbf{W}_j & 1 \leq j \leq \frac{k}{2} \\ \mathbf{W}_{k-j+1}^T & \frac{k}{2} < j \leq k. \end{cases}$$

Applying the differential operator to $J(\tilde{\mathbf{W}}, \mathbf{b})$ lets us start from expression (22):

$$\begin{aligned}
dJ &= \sum_{j=1}^k (\mathbf{Z}_{j-1}^T \Delta_j) \cdot d\tilde{\mathbf{W}}_j + (\mathbf{1}^T \Delta_j) \cdot d\mathbf{b}_j \\
&= \sum_{j=1}^{\frac{k}{2}} (\mathbf{Z}_{j-1}^T \Delta_j) \cdot d\mathbf{W}_j + \sum_{j=\frac{k}{2}+1}^k (\mathbf{Z}_{j-1}^T \Delta_j) \cdot d\mathbf{W}_{k-j+1}^T + \sum_{j=1}^k (\mathbf{1}^T \Delta_j) \cdot d\mathbf{b}_j \\
&= \sum_{j=1}^{\frac{k}{2}} (\mathbf{Z}_{j-1}^T \Delta_j) \cdot d\mathbf{W}_j + \sum_{j=\frac{k}{2}+1}^k (\Delta_j^T \mathbf{Z}_{j-1}) \cdot d\mathbf{W}_{k-j+1} + \sum_{j=1}^k (\mathbf{1}^T \Delta_j) \cdot d\mathbf{b}_j \\
&= \sum_{j=1}^{\frac{k}{2}} (\mathbf{Z}_{j-1}^T \Delta_j + \Delta_{k-j+1}^T \mathbf{Z}_{k-j}) \cdot d\mathbf{W}_j + \sum_{j=1}^k (\mathbf{1}^T \Delta_j) \cdot d\mathbf{b}_j
\end{aligned}$$

By (17) we again find components of the gradient ∇J are

$$\begin{aligned}
\frac{\partial J}{\partial \tilde{\mathbf{W}}_j} &= \mathbf{Z}_{j-1}^T \Delta_j + \Delta_{k-j+1}^T \mathbf{Z}_{k-j} + \frac{\partial p}{\partial \tilde{\mathbf{W}}_j} \\
\frac{\partial J}{\partial \mathbf{b}_j} &= \mathbf{1}^T \Delta_j + \frac{\partial p}{\partial \mathbf{b}_j}
\end{aligned}
\quad \text{where } \Delta_j = f'_j(\mathbf{A}_j) \odot \begin{cases} \frac{1}{t}(\mathbf{Z}_k - \mathbf{Y}) & j = k \\ \Delta_{j+1} \mathbf{W}_{k-j} & \frac{k}{2} \leq j < k \\ \Delta_{j+1} \mathbf{W}_{j+1}^T & 1 \leq j < \frac{k}{2} \end{cases} \quad (27)$$

3.4 Auto-Encoder with Tied Weights and Scaled Activations

The idea here is to tie the weights, but allow them to effectively have different scales. We add $\frac{k}{2} + 1$ scalar parameters $\boldsymbol{\alpha} = (1, \dots, 1, \alpha_{\frac{k}{2}}, \dots, \alpha_k) > 0$ and let $\mathbf{Z}_j = \alpha_j f_j(\mathbf{A}_j)$. The training objective is

$$J(\mathbf{W}, \mathbf{b}, \boldsymbol{\alpha}) = \frac{1}{t} \mathbf{1}^T \ell(\mathbf{Z}^k) + p(\mathbf{W}, \mathbf{b}, \boldsymbol{\alpha})$$

Applying the differential operator to $J(\tilde{\mathbf{W}}, \mathbf{b}, \boldsymbol{\alpha})$ we get

$$\begin{aligned}
dJ &= \nabla \ell(\mathbf{Z}_k) \cdot d\mathbf{Z}_k \\
&= \nabla \ell(\mathbf{Z}_k) \cdot (f_k(\mathbf{A}_k) d\alpha_k + \alpha_k f'_k(\mathbf{A}_k) \odot d\mathbf{A}_k) \\
&= \left(\frac{1}{\alpha_k} \nabla \ell(\mathbf{Z}_k) \cdot \mathbf{Z}_k \right) \cdot d\alpha_k + \left(\alpha_k \nabla \ell(\mathbf{Z}_k) \odot f'_k(\mathbf{A}_k) \right) \cdot d\mathbf{A}_k \\
&= \delta_k \cdot d\alpha_k + \Delta_k \cdot d\mathbf{A}_k \\
&= \dots + \Delta_k \cdot (\mathbf{Z}_{k-1} d\tilde{\mathbf{W}}_k + \mathbf{1} d\mathbf{b}_k + d\mathbf{Z}_{k-1} \tilde{\mathbf{W}}_k) \\
&= \dots + \mathbf{Z}_{k-1}^T \Delta_k \cdot d\tilde{\mathbf{W}}_k + \mathbf{1}^T \Delta_k \cdot d\mathbf{b}_k + \Delta_k \tilde{\mathbf{W}}_k^T \cdot d\mathbf{Z}_{k-1} \\
&= \dots + \dots + \Delta_k \tilde{\mathbf{W}}_k^T \cdot (f_{k-1}(\mathbf{A}_{k-1}) d\alpha_{k-1} + \alpha_{k-1} f'_{k-1}(\mathbf{A}_{k-1}) \odot d\mathbf{A}_{k-1}) \\
&= \dots + \dots + \left(\frac{1}{\alpha_{k-1}} \Delta_k \tilde{\mathbf{W}}_k^T \cdot \mathbf{Z}_{k-1} \right) \cdot d\alpha_{k-1} + \left(\alpha_{k-1} \Delta_k \tilde{\mathbf{W}}_k^T \odot f'_{k-1}(\mathbf{A}_{k-1}) \right) \cdot d\mathbf{A}_{k-1} \\
&= \dots + \dots + \delta_{k-1} \cdot d\alpha_{k-1} + \Delta_{k-1} \cdot d\mathbf{A}_{k-1} \\
&= \sum_{j=1}^k \mathbf{Z}_{j-1}^T \Delta_j \cdot d\tilde{\mathbf{W}}_j + \mathbf{1}^T \Delta_j \cdot d\mathbf{b}_j + \delta_j \cdot d\alpha_j \quad \text{unroll as (28–29) for } j = k..1 \quad (30)
\end{aligned} \quad (28)$$

So we get $\frac{\partial J}{\partial \tilde{\mathbf{W}}_j}$ and $\frac{\partial J}{\partial \mathbf{b}_j}$ the same, $\frac{\partial J}{\partial \alpha} = (0, \dots, 0, \delta_{\frac{k}{2}} + \frac{\partial p}{\partial \alpha_{\frac{k}{2}}}, \dots, \delta_k + \frac{\partial p}{\partial \alpha_k})$, and intermediate values:

$$\Delta_j = \alpha_j f'_j(\mathbf{A}_j) \odot \begin{cases} \frac{1}{t}(\mathbf{Z}_k - \mathbf{Y}) & j = k \\ \Delta_{j+1} \mathbf{W}_{k-j} & \frac{k}{2} \leq j < k \\ \Delta_{j+1} \mathbf{W}_{j+1}^T & 1 \leq j < \frac{k}{2} \end{cases} \quad \delta_j = \frac{1}{\alpha_j} \mathbf{Z}_j \cdot \begin{cases} \frac{1}{t}(\mathbf{Z}_k - \mathbf{Y}) & j = k \\ \Delta_{j+1} \mathbf{W}_{k-j} & \frac{k}{2} \leq j < k \\ 0 & 1 \leq j < \frac{k}{2} \end{cases} \quad (31)$$

References

- [1] Domke, J. (2011) Lecture Notes on Statistical Machine Learning. *Rochester Institute of Technology*, [<http://people.rit.edu/jcdicsa/courses/SML/01background.pdf>].