

**On a simple matrix based neural network
model
NAME**

Summary

This short paper will be examining a small scale feed-forward neural network, comprising of 2 layers of neurons, comprising of a layer with 800 and one of 10.

Contents

1	The behaviour of a neural network	4
1.1	The input and output of a neural network	4
1.2	Traits of the inner layers of the NN	4
1.3	The shape of the input vector	5
1.4	The shape of the weight and bias terms in a single neuron . .	5
1.5	The shape of the weight and bias terms in a layer	5
2	The feedforward algorithm	5
2.1	Intuition	5
2.2	The variables in each layer of a NN	5
2.3	The feedforward function	6
2.4	What is loss?	6
2.5	The shape of the weight and bias terms in a layer(2)	6
3	The back propogation algorithm	7
3.1	Intuition	7
3.2	Demonstration of the multivariate chain rule	7
3.3	Minimizing loss	7
3.4	Applying the chain rule to the loss function	7
3.5	Calculating the partial derivatives	8
3.6	The algorithm itself	9
4	Efficiency and computation	10
4.1	Context	10
4.2	Analysis of time complexity of inner functions	10
4.3	Time complexity of the feed forward algorithm	10
4.4	Time complexity of the backpropogation algorithm	10
4.5	Consolidating both into one complexity	11
5	Computation in batches	11
5.1	Intuition	11
5.2	Redefining the shape of the input and output	12

6	Formulae and notation	13
7	Appendix	13
	A	13
	B	13
	C	13
	D	14
	E	15
	F	15

Introduction

Neural networks (NNs) are an extremely powerful tool. The standard model features a number of layers, with a feed-forward model and a back-propagation algorithm in order to "train" each layer in the model. The NN being discussed will be a classification NN, with a purpose of recognizing a handwritten digit. The (nth) section will be examining the behaviour of a NN, the (nth) section will be examining the feed-forward aspect of the NN, and the (nth) sections will be examining the back-propagation algorithm. Finally, in the (nth) section, we will be discussing the application of batches to NN's, as well as how the gradient changes when dealing with many images.

1 The behaviour of a neural network

This section will be discussing the behaviour of a neural at a high level.

1.1 The input and output of a neural network

A NN, in principle, is a function. The NN takes in some form of an input, and returns the NN's certainty that it was a certain digit. A network of this form could thus be represented as a map:

$$f : \mathbb{R}^{1 \times m} \rightarrow \mathbb{R}^{1 \times n} \quad (1)$$

where m represents the dimension. The output, \mathbf{Z} , represents the confidence that the number is the index position. That is,

$$\%certainty_j = \mathbf{Z}_j$$

1.2 Traits of the inner layers of the NN

Each layer has a number of neurons. Each neuron has a number of weights, as well as a bias term. Intuitively, the weights and bias modify the input such that the output can be normalized to find the certainty of each digit. The NN being examined will have 2 layers, with 800 neurons in the first layer, while having 10 neurons in the second layer.

1.3 The shape of the input vector

As seen before in equation (1), the function f is defined to take $\mathbb{R}^{1 \times m}$ to $\mathbb{R}^{1 \times n}$. This choice of variables was intentional. Our input variable, \mathbf{x} , will be of form $\mathbb{R}^{1 \times m}$, where $m = 768$, in order to fit the input data.

1.4 The shape of the weight and bias terms in a single neuron

For now, let's define the weight term to be of form $\mathbb{R}^{768 \times 1}$, and the bias term to be of form \mathbb{R} . The reasoning for these dimension choices will be justified in the next section.

1.5 The shape of the weight and bias terms in a layer

Since a layer comprises multiple neurons, the weight and bias term can be expressed as a vector or vectors, or a matrix. We will therefore choose the weight and bias of an arbitrary layer j to be $\mathbb{R}^{n_{j-1} \times n_j}$ and $\mathbb{R}^{1 \times n_j}$ respectively, where n_j denotes the number of neurons in layer j .

2 The feedforward algorithm

This section will be demonstrating the feedforward algorithm. Matrix and vector addition in this section will be defined according to appendix (B).

2.1 Intuition

The feedforward algorithm can intuitively be thought of as the following: the goal of the network is to take an input, and output how certain the network is that it is a specific object. Since we have separate layers, it is immediately clear that each layer should modify the input, with more layers being able to perform more sensitive operations on the input. This leads to the feedforward algorithm, which takes an input, modifies it on a layer, and then passes the output on to the next layer.

2.2 The variables in each layer of a NN

Each layer in a NN will have 5 variables. The 5 variables are the following: the input, the net, the output, the weight, and the bias. The input, weight and bias combine to form the net, which is then run through an activation

function to obtain out. From the Formula and notation section, r will be the activation function for layer 1, with s for layer 2.

2.3 The feedforward function

Given a layer, define net to be

$$\mathbf{net} = \mathbf{input} \cdot \mathbf{weights} + \mathbf{bias} \quad (2)$$

Let us denote the activation function of layer j to be $a_j(\mathbf{z})$. Then, in accordance with section (2.2), the out will be defined to be

$$\mathbf{out} = a(\mathbf{net}) \quad (3)$$

Combining equations (2) and (3) yields

$$\mathbf{out} = a(\mathbf{input} \cdot \mathbf{weights} + \mathbf{bias}) \quad (4)$$

Using equation (4), we can find the output of the second layer to be

$$\mathbf{out}_2 = s(a(\mathbf{input}_1 \cdot \mathbf{weights}_1 + \mathbf{bias}_1) \cdot \mathbf{weights}_2 + \mathbf{bias}_2) \quad (5)$$

$$\implies \mathbf{out}_2 = s(\mathbf{net}_2) \quad (6)$$

2.4 What is loss?

Loss is defined to be a value that represents how far off \mathbf{out}_2 is from the correct input. This function is denoted $L_{\mathbf{y}}(\mathbf{out}_2)$ where \mathbf{y} is an one hot encoded vector (see appendix A) and the definition can be found in the Formula and notation section.

2.5 The shape of the weight and bias terms in a layer(2)

In section (1.4), we defined the shape of the weight term to be $\mathbb{R}^{n_{j-1} \times n_j}$. Now that we have introduced the loss calculation, we can justify this shape. Let the weight matrix in layer 1 be $\mathbb{R}^{a \times b}$. Substituting the dimensions of the variables in layer 1 into equation (2) yields

$$\mathbf{out} = \mathbb{R}^{1 \times 768} \cdot \mathbb{R}^{a \times b} + \mathbb{R} \quad (7)$$

by the definition of matrix multiplication, we see that a must be 768.

\mathbf{out} in this case, is of shape $1 \times b$, meaning that the weight matrix in layer 2 must have shape $b \times c$, where c is some arbitrary constant. In that case, \mathbf{out}_2 is of shape $1 \times c$. For convenience purposes, we will consider c to be 10.

3 The back propogation algorithm

In this section, we will be examining the back propogation algorithm

3.1 Intuition

As discussed in section 2, each weight plays a significant part in the final output. Seeing as how we want the final output to be as close as possible, we will need to measure how much each variable in the equation needs to change. This is the essence of the back propogation algorithm.

3.2 Demonstration of the multivariate chain rule

Given some function $f(g(\mathbf{t}))$, the derivative with respect to \mathbf{t} is

$$\frac{\partial f}{\partial \mathbf{t}} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial \mathbf{t}}$$

This is very important for us as it allows us to find the derivative with respect to each paramater relatively easily.

3.3 Minimizing loss

As mentioned in section (2.4), loss is the measure of how different the output was compared to the correct output. Informally, this can be thought of as how incorrect the output is. Naturally, the goal would be to reduce the loss, in order to make the order more correct. This can be done by finding the so called gradient, and subtracting that from every matrix. This is represented with the general formula:

let \mathbf{a}^n denote the nth iteration of \mathbf{a} .

given some $f(\mathbf{a})$, the **local** minimum can be found with the following iterative process:

$$\mathbf{a}^{n+1} = \mathbf{a}^n - \eta \nabla_{\mathbf{a}^n} f$$

where $\eta \in (0, 1)$ is the learning rate of the network and $\nabla_{\mathbf{a}^n}$ is the gradient with respect to \mathbf{a}^n (see appendix D)

3.4 Applying the chain rule to the loss function

In this subsection, I will be demonstrating how to use the chain rule to find the partial derivative with respect to each non-input variable, i.e weights and

biases.

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{w}_2} &= \frac{\partial L}{\partial \mathbf{out}_2} \frac{\partial \mathbf{out}_2}{\partial \mathbf{net}_2} \frac{\partial \mathbf{net}_2}{\partial \mathbf{w}_2} \\ \frac{\partial L}{\partial \mathbf{b}_2} &= \frac{\partial L}{\partial \mathbf{out}_2} \frac{\partial \mathbf{out}_2}{\partial \mathbf{net}_2} \frac{\partial \mathbf{net}_2}{\partial \mathbf{b}_2}\end{aligned}$$

As it is clear, $\frac{\partial L}{\partial \mathbf{out}_2} \frac{\partial \mathbf{out}_2}{\partial \mathbf{net}_2}$ appears quite often. Denote this quantity δ_2 , and the equations simplify to

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{w}_2} &= \delta_2 \frac{\partial \mathbf{net}_2}{\partial \mathbf{w}_2} \\ \frac{\partial L}{\partial \mathbf{b}_2} &= \delta_2 \frac{\partial \mathbf{net}_2}{\partial \mathbf{b}_2}\end{aligned}$$

Now, let's consider the derivatives with respect to the values in the first layer

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{w}_1} &= \frac{\partial L}{\partial \mathbf{out}_2} \frac{\partial \mathbf{out}_2}{\partial \mathbf{net}_2} \frac{\partial \mathbf{net}_2}{\partial \mathbf{out}_1} \frac{\partial \mathbf{out}_1}{\partial \mathbf{net}_1} \frac{\partial \mathbf{net}_1}{\partial \mathbf{w}_1} \\ \frac{\partial L}{\partial \mathbf{b}_1} &= \frac{\partial L}{\partial \mathbf{out}_2} \frac{\partial \mathbf{out}_2}{\partial \mathbf{net}_2} \frac{\partial \mathbf{net}_2}{\partial \mathbf{out}_1} \frac{\partial \mathbf{out}_1}{\partial \mathbf{net}_1} \frac{\partial \mathbf{net}_1}{\partial \mathbf{b}_1}\end{aligned}$$

Again, let us simplify this expression.

$$\delta_1 \equiv \frac{\partial L}{\partial \mathbf{out}_2} \frac{\partial \mathbf{out}_2}{\partial \mathbf{net}_2} \frac{\partial \mathbf{net}_2}{\partial \mathbf{out}_1} \frac{\partial \mathbf{out}_1}{\partial \mathbf{net}_1}$$

As is clear, this can be defined as the following.

$$\delta_1 = \delta_2 \frac{\partial \mathbf{net}_2}{\partial \mathbf{out}_1} \frac{\partial \mathbf{out}_1}{\partial \mathbf{net}_1}$$

In our network, this is the extent of what we need. However, a general form is provided in appendix E.

3.5 Calculating the partial derivatives

Due to the nature of the L function and the s function, it is difficult to obtain each derivative in the δ_{22} term. Instead of each step of the partial differential, we will skip directly to $\frac{\partial L}{\partial \mathbf{net}_2}$. A derivation of this can be found in appendix A.F. This value is equal to $(\mathbf{out}_2 - \mathbf{y})$. This implies that

$$\delta_2 = (\mathbf{out}_2 - \mathbf{y})$$

Now that we have this information, let us calculate δ_1 .

$$\begin{aligned}\delta_1 &= \delta_2 \frac{\partial \mathbf{net}_2}{\partial \mathbf{out}_1} \frac{\partial \mathbf{out}_1}{\partial \mathbf{net}_1} \\ \frac{\partial \mathbf{net}_2}{\partial \mathbf{out}_1} &= \mathbf{w}_2^T \\ \implies \delta_1 &= ((\mathbf{out}_2 - \mathbf{y}) \cdot \mathbf{w}_2^T) \otimes \frac{\partial \mathbf{out}_1}{\partial \mathbf{net}_1}\end{aligned}$$

with this information, the rest of the required derivatives are easy to calculate;

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{w}_2} &= \delta_2 \frac{\partial \mathbf{net}_2}{\partial \mathbf{w}_2} \\ \frac{\partial L}{\partial \mathbf{b}_2} &= \delta_2 \frac{\partial \mathbf{net}_2}{\partial \mathbf{b}_2} \\ \frac{\partial L}{\partial \mathbf{w}_1} &= \delta_1 \frac{\partial \mathbf{net}_1}{\partial \mathbf{w}_1} \\ \frac{\partial L}{\partial \mathbf{b}_1} &= \delta_1 \frac{\partial \mathbf{net}_1}{\partial \mathbf{b}_1}\end{aligned}$$

3.6 The algorithm itself

Let us consider a 2 layer network with input $\mathbf{x} \in \mathbb{R}^{1 \times 768}$ and output $C \equiv L(\mathbf{out}_2) \in \mathbb{R}$. The goal is to minimize C . In order to do so, we can use the general formula described in section (3.3), $\mathbf{a}^{n+1} = \mathbf{a}^n - \eta \nabla_{\mathbf{a}^n} f$. We will go about doing so as follows.

Let the superscript denote the n 'th training data point passed to the network. Then, we apply the formula to each variable. That is,

$$\begin{aligned}\mathbf{w}_j^{n+1} &= \mathbf{w}_j^n - \eta \frac{\partial L}{\partial \mathbf{w}_j^n} \\ \mathbf{b}_j^{n+1} &= \mathbf{b}_j^n - \eta \frac{\partial L}{\partial \mathbf{b}_j^n}\end{aligned}$$

, where $\eta \in (0, 1)$.

In theory, after all training data in the given dataset is processed, \mathbf{w}_j^{final} and \mathbf{b}_j^{final} should be fine tuned enough to the point where it can accurately match any image in the dataset. There will inherently be issues dealing with elements outside of the dataset, but the fine tuning should converge to a point where it is very accurate ($> 95\%$)

4 Efficiency and computation

This section will examine the dataset in detail, as well as talk about the efficiency and the time complexity for the NN

4.1 Context

The dataset being used is the standard MNIST database, with 60000 training points, consisting of digits identified by humans. The program itself is run on Java, using no pre-existing framework. Since the implementation was all done naively, the time complexity is enormous.

4.2 Analysis of time complexity of inner functions

In the neural network, the major functions used are matrix multiplication, elementwise matrix multiplication, as well as the R, L and s functions. The implementation of matrix multiplication is done in such a way that it is $O(n^3)$, elementwise matrix multiplication is $O(n^2)$, and the rest of the mentioned functions are $O(n)$. By inspection, it is clear that the time complexity is expected to blow up. However, there are hard limits, such as $s(\mathbf{z})$. In theory, this function is $O(n)$. However, we know that \mathbf{z} is capped at 10 elements, so the complexity collapses to $O(1)$. In subsequent analysis we will be considering the original complexities.

4.3 Time complexity of the feed forward algorithm

We know that given some functions f, g with time complexities $O(a), O(b)$, then $O(f \circ g) = O(a \cdot cb)$ where $c \in \mathbb{R}$. For the sake of brevity, we will just denote this quantity as $O(ab)$, since the time complexities of the used functions are all polynomial.

In the feedForward network there are two instances of matrix multiplication, as well as 2 activation functions. This works out to be $O(n^8)$.

4.4 Time complexity of the backpropagation algorithm

Thankfully, since the partial derivatives were already derived analytically, we can easily show the time complexity. First, let us show the time complexity

of δ_1 and δ_2 , keeping in mind that the values are already calculated.

$$\begin{aligned}\delta_2 &= (\mathbf{out}_2 - \mathbf{y}) \\ &= O(1) \\ \delta_1 &= ((\mathbf{out}_2 - \mathbf{y}) \cdot \mathbf{w}_2^T) \otimes \frac{\partial \mathbf{out}_1}{\mathbf{net}_1} \\ &= O(1 \cdot n^3 \cdot n^2) \\ &= O(n^5)\end{aligned}$$

4.5 Consolidating both into one complexity

Now that we know that the feed forward algorithm is $O(n^8)$, and the back propagation algorithm $O(n^5)$, we can find the combined time complexity. Knowing that the algorithms are run in sequence, the total time complexity $O(n^8 + n^5)$, and since we take the highest order term, this collapses into $O(n^8)$, giving our final result of the runtime of the 2 algorithms being $O(n^8)$

Important note

This time complexity is exclusively for the two algorithms on one piece of training data. It does not say anything about the total runtime of the NN.

5 Computation in batches

This is the final section and will be discussing performing computation in batches instead of processing one input at a time.

5.1 Intuition

Instead of doing 60000 $O(n^8)$ operations, we can instead perform significantly less by using "batches". The concept is to pick a sufficiently large size and construct a column vector of row vectors. The functions are to be redefined to work with matrices.

5.2 Redefining the shape of the input and output

originally, the input was in the shape of $\mathbb{R}^{1 \times 768}$. Following the definition

outlined in section (5.1), define $\mathbf{x} = \begin{bmatrix} \mathbb{R}^{1 \times 768} \\ \vdots \\ ntimes \\ \vdots \\ \mathbb{R}^{1 \times 768} \end{bmatrix} \in \mathbb{R}^{n \times 768}$. Thus, \mathbf{out}_1 will be

of the shape ,

6 Formulae and notation

$\mathbb{R}^{a \times b}$ denotes matrices of dimension $a \times b$ over the field \mathbb{R}

All figures marked in **bold** denote row vectors

All figured marked in ***bold italics*** denote matrices;

Given some function $f: \mathbb{R}^{1 \times n} \rightarrow \mathbb{R}^{1 \times m}$, $\frac{\partial f}{\partial \mathbf{z}} = \nabla_{\mathbf{z}} f$

\mathbf{Z}_j represents the value at the j -th index of \mathbf{Z}

Activation function of layer 1: $r(\mathbf{z}) : \mathbb{R}^a \rightarrow \mathbb{R}^a = [\max(0, \mathbf{z}_j)]_j$

Activation function of layer 2: $s(\mathbf{z}) : \mathbb{R}^a \rightarrow \mathbb{R}^a = [\frac{e^{\mathbf{z}_j}}{\sum_{n=0}^a e^{\mathbf{z}_n}}]_j$

Loss function: $L_{\mathbf{y}}(\mathbf{z}) : \mathbb{R}^a \rightarrow \mathbb{R} = -\sum_{n=0}^a \mathbf{y}_n \ln(s(\mathbf{z})_n)$

Kronecker delta: $\delta_{ij} = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases}$

7 Appendix

A

An one hot encoded vector is a variable $\mathbf{y} \in \mathbb{Z}^{n+}$, where \mathbb{Z}^{n+} denotes the integers greater than, or equal to 0, such that the sum of the elements of $\mathbf{y} = 1$

B

Initially, an addition between $\mathbb{R}^{m \times 1}$ and \mathbb{R} seems ill defined. we will then define it as:

let $a \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^{1 \times n}$

define $a + b$ to be $\begin{pmatrix} a_1 + b \\ a_2 + b \\ \vdots \\ a_m + b \end{pmatrix}$, where a_j denotes the j 'th row of a and $a_j + b$

is standard vector addition

C

Since differentiation with respect to a matrix is ill-defined under many formalisms of calculus, we will set some ground rules:

$$\begin{aligned}\frac{\partial \mathbf{x} \mathbf{W}}{\partial \mathbf{x}} &= \mathbf{W}^T \\ \frac{\partial \mathbf{x} \mathbf{W}}{\partial \mathbf{W}} &= \mathbf{x}^T\end{aligned}$$

Another rule we will set:

$$\begin{aligned}\text{let } f : \mathbb{R}^j &\rightarrow \mathbb{R} \\ \text{let } \mathbf{W} &\in \mathbb{R}^{m \times n} \\ \implies \frac{\partial f}{\partial \mathbf{W}} &\in \mathbb{R}^{m \times n}\end{aligned}$$

Since the chain rule raises many vector valued arguments, we will need to do products thereof elementwise, in order to maintain the shape of the arguments. That is, an operation \otimes will be used when multiplying differentials. An example being some function $f(g(x))$:

$$\frac{\partial f}{\partial \mathbf{t}} = \frac{\partial f}{\partial g} \otimes \frac{\partial g}{\partial \mathbf{t}}$$

This will be shortened to

$$\frac{\partial f}{\partial g} \frac{\partial g}{\partial \mathbf{t}}$$

for the sake of brevity.

An important thing to note is that the hadamard product behaves differently under certain conditions.

Let us consider $\frac{\partial L}{\partial a} \frac{\partial a}{\partial \mathbf{W}}$. Contrary to the claim made above, this is equivalent to $\mathbf{W}^T \cdot \frac{\partial L}{\partial a}$

D

The gradient of a function $f(\mathbf{z}) : \mathbb{R}^n \rightarrow \mathbb{R}$ with respect to some vector \mathbf{a} is

notated $\nabla_{\mathbf{a}}$ and is defined
$$\begin{bmatrix} [h] \frac{\partial f}{\partial \mathbf{z}_1}(\mathbf{a}_1) \\ \frac{\partial f}{\partial \mathbf{z}_2}(\mathbf{a}_2) \\ \vdots \\ \frac{\partial f}{\partial \mathbf{z}_n}(\mathbf{a}_n) \end{bmatrix}$$

E

Assuming a m layer network, δ_n has a recursive formula: $\delta_n = \begin{cases} \frac{\partial L}{\partial \text{out}_n} \frac{\partial \text{out}_n}{\partial \text{net}_n} & n = m \\ \delta_{n+1} \frac{\partial \text{net}_{n+1}}{\partial \text{out}_n} \frac{\partial \text{out}_n}{\partial \text{net}_n} & n \neq m \end{cases}$

F

Denote the activation function of the second layer $s(\mathbf{z})$. Since this is a vector valued function, returning a vector, we expect the derivative to be the Jacobian matrix, a matrix denoted $J_{\mathbf{z}}(s)$. The following derivation will show that. Denote s_k to be the k th element of the output of s .

$$J_{\mathbf{z}}(s) = \begin{bmatrix} \frac{\partial s_1}{\partial \mathbf{z}_1} & \dots & \frac{\partial s_1}{\partial \mathbf{z}_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial s_n}{\partial \mathbf{z}_1} & \dots & \frac{\partial s_n}{\partial \mathbf{z}_n} \end{bmatrix}$$

Note that since $s_i = \frac{e^{\mathbf{z}_i}}{\sum_{n=0}^{\dim(z)} e^{\mathbf{z}_n}}$, it is easier to use logarithmic differentiation

$$\begin{aligned} \frac{\partial \ln(s_i)}{\partial \mathbf{z}_j} &= \frac{1}{s_i} \frac{\partial s_i}{\partial \mathbf{z}_j} \\ \implies \frac{\partial s_i}{\partial \mathbf{z}_j} &= \frac{\partial \ln(s_i)}{\partial \mathbf{z}_j} s_i \\ &= s_i \frac{\partial}{\partial \mathbf{z}_j} \ln \left(\frac{e^{\mathbf{z}_i}}{\sum_{n=0}^{\dim(z)} e^{\mathbf{z}_n}} \right) \\ &= s_i \frac{\partial}{\partial \mathbf{z}_j} (\ln(e^{\mathbf{z}_i}) - \ln(\sum_{n=0}^{\dim(z)} e^{\mathbf{z}_n})) \\ &= s_i \frac{\partial}{\partial \mathbf{z}_j} (\mathbf{z}_i - \ln(\sum_{n=0}^{\dim(z)} e^{\mathbf{z}_n})) \\ &= s_i (\delta_{ij} - \frac{\partial}{\partial \mathbf{z}_j} \ln(\sum_{n=0}^{\dim(z)} e^{\mathbf{z}_n})) \\ &= s_i (\delta_{ij} - \frac{1}{\sum_{n=0}^{\dim(z)} e^{\mathbf{z}_n}} e^{\mathbf{z}_j}) \\ &= s_i (\delta_{ij} - s_j) \end{aligned}$$

Denote the loss function with respect to some one hot encoded vector \mathbf{y} $L_{\mathbf{y}}(\mathbf{z})$

In the following derivation, we will be differentiating L with respect to each component of \mathbf{z} and building it back into a vector at the end

$$\begin{aligned}
\frac{\partial L_{\mathbf{y}}}{\partial \mathbf{z}_j} &= -\frac{\partial}{\partial \mathbf{z}_j} \sum_{n=0}^a \mathbf{y}_n \ln(s(\mathbf{z})_n) \\
&= -\sum_{n=0}^a \mathbf{y}_n \frac{\partial}{\partial \mathbf{z}_j} \ln(s(\mathbf{z})_n) \\
&= -\sum_{n=0}^a \mathbf{y}_n \frac{\partial}{\partial \mathbf{z}_j} \ln(s(\mathbf{z})_n) \\
&= -\sum_{n=0}^a \mathbf{y}_n \frac{\partial}{\partial \mathbf{z}_j} \ln(s_n) \\
&= -\sum_{n=0}^a \mathbf{y}_n \frac{1}{s_n} \frac{\partial s_n}{\partial \mathbf{z}_j} \\
&= -\sum_{n=0}^a \mathbf{y}_n \frac{1}{s_n} s_n (\delta_{nj} - s_j) \\
&= -\sum_{n=0}^a \mathbf{y}_n (\delta_{nj} - s_j) \\
&= -\mathbf{y}_j - \sum_{n=0}^a -s_j \mathbf{y}_n \\
&= -\mathbf{y}_j + s_j \sum_{n=0}^a \mathbf{y}_n \\
&= s_j - \mathbf{y}_j
\end{aligned}$$