

Explanatory Notes for 6.390

Shauntclair Ruiz (Current TA)

Spring 2024

Contents

A Matrix Derivatives	3
A.1 Introduction and Review	3
A.1.1 Partial Derivatives	4
A.1.2 Thinking about derivatives	6
A.1.3 Derivatives: Approximation	7
A.2 Derivative: Scalar/Vector (Gradient)	9
A.2.1 Finding the scalar/vector derivative	9
A.2.2 Review: Planar Approximation	10
A.2.3 Our completed scalar/vector derivative	13
A.3 Derivative: Vector/Scalar	16
A.3.1 Working with the vector derivative	16
A.4 Derivative: Vector/Vector	19
A.4.1 The vector/vector derivative	19
A.5 General derivative (Vector/Vector)	21
A.5.1 More about the vector/vector derivative	24
A.6 Derivative: matrix/scalar	26
A.7 Derivative: scalar/matrix	28
A.8 Tensors	30
A.8.1 Other Derivatives	30
A.8.2 Dimensions (Optional)	32
A.8.3 Dealing with Tensors	34
A.9 Chapter 7 Derivatives	37
A.9.1 The loss derivative	37
A.9.2 The weight derivative	37
A.9.3 Linking Layers $\ell - 1$ and ℓ	42
A.9.4 Activation Function	44

A.9.5 Element-wise multiplication	46
A.10 Terms	48

APPENDIX A

Matrix Derivatives

A.1 Introduction and Review

Our goal here, is to combine the powers of matrices and calculus:

- **Matrices:** the ability to store lots of **data**, and do fast linear operations on all that data at the **same time**.

Example: Consider

$$\mathbf{w}^T \mathbf{x} = \begin{bmatrix} w_1 & w_2 & \cdots & w_m \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} = \sum_{i=1}^m x_i w_i \quad (\text{A.1})$$

In this case, we're able to do m different **multiplications** at the same time! This is what we like about matrices.

In this case, we're thinking about vectors as $(m \times 1)$ matrices.

- **Calculus:** analyzing the way different variables are **related**: how does changing x affect y ?

Example: Suppose we have

$$\frac{\partial f}{\partial x_1} = 10 \quad \frac{\partial f}{\partial x_2} = -5 \quad (\text{A.2})$$

Now we know that, if we increase x_1 , we increase f . This **understanding** of variables

is what we like about derivatives.

Concept 1

Matrix derivatives allow us to find **relationships** between large volumes of **data**.

- These "relationships" are **derivatives**: consider dy/dx . How does y change if we modify x ? Currently, we only have **scalar derivatives**.
- These "data" are stored as **matrices**: blocks of data, that we can do linear operations (matrix multiplication) on.

Our goal is to work with many scalar derivatives at the **same time**.

In order to do that, we can apply some **derivative** rules, but we have to do it in a way that **agrees** with **matrix** math.

Our work is a careful balancing act between getting the **derivatives** we want, without violating the **rules** of matrices (and losing what makes them useful!)

Example: When we multiply two matrices, their inner shape has to match: in the below case, they need to share a dimension b .

$$\underbrace{(a \times \color{red}{b})}_X \underbrace{(\color{red}{b} \times c)}_Y \quad (A.3)$$

We can't do anything that would **violate** matrix rules like these: otherwise, we're not really doing "matrix **calculus**" anymore. This means we need to build our math carefully.

First, we'll look at the **properties** of derivatives. Then, we figure out how to usefully apply them to **vectors**, and finally, to **matrices**.

A.1.1 Partial Derivatives

One more comment, though - we may have many different variables floating around. This means we **have** to use the multivariable **partial derivative**.

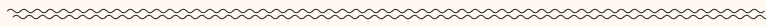
Definition 2

The **partial derivative**

$$\frac{\partial B}{\partial A}$$

Is useful when there may be **multiple variables** in our functions.

The rule of the partial derivative is that we keep every variable **other** than A and B **fixed**.



- This is different from the "total" derivative $\frac{dB}{dA}$, where we consider how B affects A, without keeping other variables fixed.
- The total derivative is also used in multi-variable calculus, but serves a different role: we'll come back to this later.

Example: Consider $f(x, y) = 2x^2y$.

$$\frac{\partial f}{\partial x} = 2(2x)y \quad (\text{A.4})$$

Here, we kept y *fixed* - we treat it as if it were an unchanging **constant**.

Using the partial derivative lets us keep our work tidy:

- If **many** variables were allowed to **change** at the same time, it could get very confusing. _____
- Since this is too complicated, we'll instead change these variables *one at a time*. We get a partial derivative for each of them, holding the others **constant**.

Imagine keeping track of k different variables x_i with k different changes Δx_i at the same time! That's a headache.

Our **total** derivative is the result of all of those different variables, **added** together. This is how we get the **multi-variable chain rule**.

Definition 3

The **multi-variable chain rule** in 3-D $\{(x, y, z)\}$ is given as

$$\frac{df}{ds} = \overbrace{\frac{\partial f}{\partial x} \frac{\partial x}{\partial s}}^{\text{only modify } x} + \overbrace{\frac{\partial f}{\partial y} \frac{\partial y}{\partial s}}^{\text{only modify } y} + \overbrace{\frac{\partial f}{\partial z} \frac{\partial z}{\partial s}}^{\text{only modify } z}$$

If we have k variables $\{x_1, x_2, \dots, x_k\}$ we can generalize this as:

$$\frac{df}{ds} = \sum_{i=1}^k \overbrace{\frac{\partial f}{\partial x_i} \frac{\partial x_i}{\partial s}}^{x_i \text{ component}}$$

- ~~~~~
- We justify a lot of aspects of the multivariable chain rule, and the gradient, in our Gradient Descent chapter. Feel free to review there.

A.1.2 Thinking about derivatives

The typical definition of derivatives

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad (\text{A.5})$$

Gives an *idea* of what sort of things we're looking for. It reminds us of one piece of information we need:

- Our derivative **depends** on the **current position** x we are taking the derivative at.

We need this because derivative are **local**: the relationship between our variables might change if we move to a different **position**.

But, the problem with vectors is that each component can act **separately**: if we have a vector, we can change in many different "directions".

$$A = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} \quad B = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \quad (\text{A.6})$$

Example: Suppose we want a derivative $\partial B / \partial A$: $\Delta a_1, \Delta a_2$, and Δa_3 could each, separately, have an effect on Δb_1 and/or Δb_2 . That requires 6 different derivatives, $\partial b_i / \partial a_j$.

- Every component of the input A can potentially modify **every** component of the output B .

3 dimensions of A times
2 dimensions of B : 6
combinations.

One solution we could try is to just collect all of these derivatives into a **vector** or **matrix**.

Concept 4

For the **derivative** between two objects (scalars, vectors, matrices) **A** and **B**

$$\frac{\partial B}{\partial A}$$

We need to get the **derivatives**

$$\frac{\partial b_j}{\partial a_i}$$

between every **pair** of elements **a_i**, **b_j**: each pair of elements could have a **relationship**.

The total number of elements (or "size") is...

$$\text{Size}\left(\frac{\partial B}{\partial A}\right) = \text{Size}(B) * \text{Size}(A)$$

Collecting these values into a **matrix** will gives us all the information we need.

But, how do we gather them? What should the **shape** look like? Should we **transpose** our matrix or not?

A.1.3 Derivatives: Approximation

To answer this, we need to ask ourselves *why* we care about these derivatives: we'll define them based on what we **need** them for.

- We care about the **direction of greatest decrease**: the negative of the **gradient**.
 - One application: we might want to adjust weight vector **w** to reduce loss \mathcal{L} .
- To find the gradient, we'll need some other **matrix** derivatives that will be useful in a **chain rule**.
 - The chain rule creates lots of intermediate derivatives, that we need to work with.

Let's focus on the first point: we want to **minimize** \mathcal{L} . Our focus is the **change** in \mathcal{L} , $\Delta\mathcal{L}$.

We want to take steps that reduce our loss \mathcal{L} .

$$\frac{\partial \mathcal{L}}{\partial w} \approx \frac{\text{Change in } \mathcal{L}}{\text{Change in } w} = \frac{\Delta \mathcal{L}}{\Delta w} \quad (\text{A.7})$$

Thus, we **solve** for $\Delta\mathcal{L}$:

All we do is multiply both sides by Δw .

$$\Delta\mathcal{L} \approx \frac{\partial\mathcal{L}}{\partial w} \Delta w \quad (\text{A.8})$$

Since this derivation was gotten using scalars, we might need a **different** type of multiplication for our **vector** and **matrix** derivatives.

Concept 5

We can use derivatives to **approximate** the change in our output based on our input:

$$\Delta\mathcal{L} \approx \frac{\partial\mathcal{L}}{\partial w} \star \Delta w$$

Where the \star symbol represents some type of **multiplication**.

We can think of this as a **function** that takes in change in Δw , and returns an **approximation** of the loss.

We already understand **scalar** derivatives, so let's move on to the **gradient**.

~~~~~

## A.2 Derivative: Scalar/Vector (Gradient)

Our plan is to look at every derivative combination of scalars, vectors, and matrices we can.

First, we consider:

$$\frac{\partial(\text{Scalar})}{\partial(\text{Vector})} = \frac{\partial \mathbf{s}}{\partial \mathbf{v}} \quad (\text{A.9})$$

We'll take  $\mathbf{s}$  to be our scalar, and  $\mathbf{v}$  to be our vector. So, our input is a **vector**, and our output is a **scalar**.

$$\Delta \mathbf{v} \longrightarrow \boxed{\mathbf{f}} \longrightarrow \Delta \mathbf{s} \quad (\text{A.10})$$

How do we make sense of this? Well, let's write  $\Delta \mathbf{v}_i$  explicitly:

$$\overbrace{\begin{bmatrix} \Delta \mathbf{v}_1 \\ \Delta \mathbf{v}_2 \\ \vdots \\ \Delta \mathbf{v}_m \end{bmatrix}}^{\Delta \mathbf{v}} \longrightarrow \Delta \mathbf{s} \quad (\text{A.11})$$

We can see that we have  $m$  different **inputs** we can change in order to change our **one** output.

So, our derivative needs to have  $m$  different **elements**: one for each element  $\mathbf{v}_i$ .

### A.2.1 Finding the scalar/vector derivative

But how do we shape our matrix? Let's look at our **rule**.

$$\Delta \mathbf{s} \approx \frac{\partial \mathbf{s}}{\partial \mathbf{v}} \star \Delta \mathbf{v} \quad (\text{A.12})$$

We can transform using our shapes:

$$\Delta \mathbf{s} \approx \frac{\partial \mathbf{s}}{\partial \mathbf{v}} \star \overbrace{\begin{bmatrix} \Delta \mathbf{v}_1 \\ \Delta \mathbf{v}_2 \\ \vdots \\ \Delta \mathbf{v}_m \end{bmatrix}}^{\Delta \mathbf{v}} \quad (\text{A.13})$$

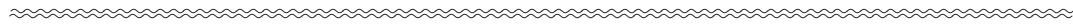
How do we get  $\Delta \mathbf{s}$ ? We have so many variables. Let's focus on them one at a time: breaking  $\Delta \mathbf{v}$  into  $\Delta \mathbf{v}_i$ , so we'll try to consider each  $\mathbf{v}_i$  **separately**.

It's usually possible to change each  $\mathbf{v}_i$ , so we have to look at every one of them.

One problem, though: how can we treat each **derivative** separately?

- Suppose we only start by considering  $\Delta v_1$ : it'll move us along one axis.
- But our derivative is based on our position! If we've just moved, all of our derivatives have changed.
- So, choosing  $\Delta v_1$  first, had an effect on the other  $\Delta v_i$  terms: they're not really "separate".

So, what do we do?



## A.2.2 Review: Planar Approximation

We'll resolve this the same way we did in chapter 3, **gradient** descent: by taking advantage of the "planar approximation".

The solution is this: assume your function is **smooth**. The **smaller** a step you take, the **less** your derivative has a chance to change.

**Example:** Take  $f(x) = x^2$ .

- If we go from  $x = 1 \rightarrow 2$ , then our derivative goes from  $f'(x) = 2 \rightarrow 4$ .
- Let's **shrink** our step. We go from  $x = 1 \rightarrow 1.01$ , our derivative goes from  $f'(x) = 2 \rightarrow 2.02$ .
  - Our derivative is almost the same!

This isn't true for big steps, but eventually, if your step is small enough, then the derivative will barely change.

If we take a small enough step  $\Delta v_i$ , then, if our function is **smooth**, then the derivative will hardly change!

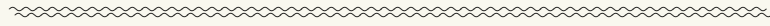
So, if we zoom in enough (shrink the scale of change), then we can **pretend** the derivative is **constant**.

You could imagine repeatedly shrinking the size of our step, until the change in the derivatives is basically unnoticeable.

**Concept 6**

If you have a **smooth function**, then...

If you take sufficiently **small steps**, then you can treat the derivatives as **constant**.



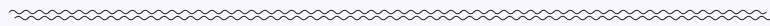
This is called the **planar approximation** because that's how a smooth surface looks, at small scales:

- The derivatives of a flat plane are the same at every position: they're **constant** (we show this below).
- This is, of course, an **approximation**: we can't use it at every position, just "very nearby positions".

**Clarification**

This clarification is **optional**.

We can describe "sufficiently small steps" in a more mathematical way:



Our goal is for  $f'(x)$  to be **basically constant**: it doesn't change much. In other words,  $\Delta f'(x)$  is **small**.

Let's say we don't want it to change more than  $\delta$ : this is our definition of "very small".

If you want

- $\Delta f'(x)$  to be very small ( $|\Delta f'(x)| < \delta$ )
- It has been proven that...
  - it's possible to take a small enough step  $|\Delta x| < \epsilon$ , and to get that result.

One way to describe this is to say that our function is (locally) **flat**: it looks like some kind of plane/hyperplane.

The word "locally" represents the small step size: we stay in the "local area".

**Clarification 7**

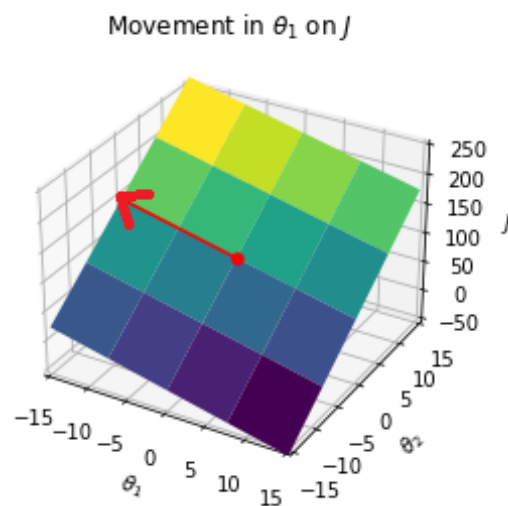
Why is this **true**? Because a **hyperplane** can be represented using our **linear** function

$$f(\mathbf{x}) \approx \boldsymbol{\theta}^T \mathbf{x} + \theta_0 = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_m x_m$$

If we take a derivative:

$$\frac{\partial f}{\partial x_i} = \theta_i$$

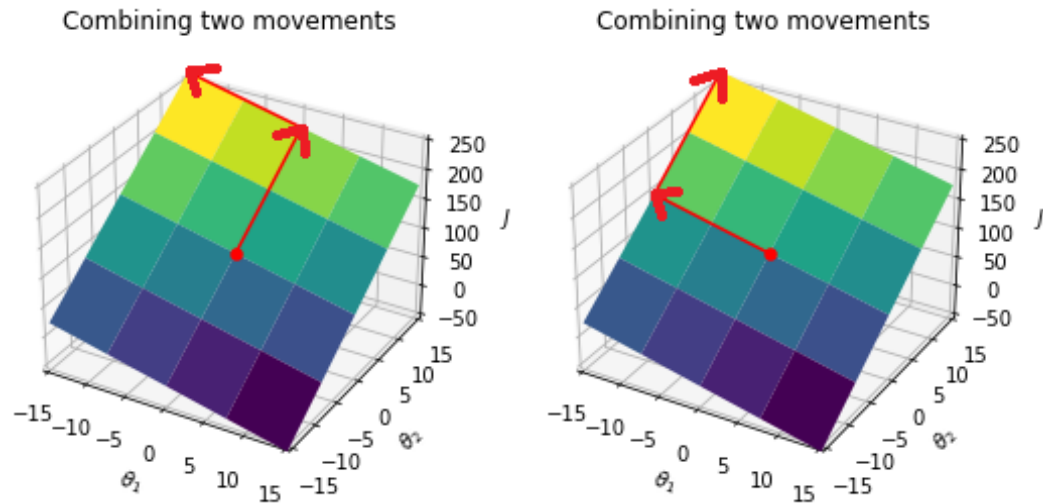
That derivative is a **constant**! It's doesn't change based on **position**.



If we take very small steps, we can approximate our function as **flat**.

Why does this help? If our derivative doesn't **change**, you can take multiple steps  $\Delta \mathbf{v}_i$  and the order doesn't matter.

So, you can combine your steps or separate them easily.



We can break up our big step into two smaller steps that are truly independent: order doesn't matter.

With that, we can add up all of our changes:

$$\Delta s = \Delta s_{\text{from } v_1} + \Delta s_{\text{from } v_2} + \cdots + \Delta s_{\text{from } v_m} \quad (\text{A.14})$$

### A.2.3 Our completed scalar/vector derivative

From this, we can get an **approximated** version of the MV chain rule.

#### Definition 8

The **multivariable chain rule approximation** looks similar to the multivariable chain rule, but for finite changes  $\Delta x_i$ .

In 3-D, we get

$$\Delta f \approx \overbrace{\frac{\partial f}{\partial x} \Delta x}^{\text{x component}} + \overbrace{\frac{\partial f}{\partial y} \Delta y}^{\text{y component}} + \overbrace{\frac{\partial f}{\partial z} \Delta z}^{\text{z component}}$$

In general, we have

$$\Delta f \approx \sum_{i=1}^m \overbrace{\frac{\partial f}{\partial x_i} \Delta x_i}^{\text{x}_i \text{ component}}$$

This function lets us add up the effect each component has on our output, using **derivatives**.

This gives us what we're looking for:

$$\Delta \mathbf{s} \approx \sum_{i=1}^m \frac{\partial \mathbf{s}}{\partial \mathbf{v}_i} \Delta \mathbf{v}_i \quad (\text{A.15})$$

If we circle back around to our original approximation:

$$\sum_{i=1}^m \frac{\partial \mathbf{s}}{\partial \mathbf{v}_i} \Delta \mathbf{v}_i = \frac{\partial \mathbf{s}}{\partial \mathbf{v}} \star \overbrace{\begin{bmatrix} \Delta \mathbf{v}_1 \\ \Delta \mathbf{v}_2 \\ \vdots \\ \Delta \mathbf{v}_m \end{bmatrix}}^{\Delta \mathbf{v}} \quad (\text{A.16})$$

When we look at the left side, we're multiplying pairs of components, and then adding them. That sounds similar to a **dot product**.

$$\sum_{i=1}^m \frac{\partial \mathbf{s}}{\partial \mathbf{v}_i} \Delta \mathbf{v}_i = \overbrace{\begin{bmatrix} \partial \mathbf{s} / \partial \mathbf{v}_1 \\ \partial \mathbf{s} / \partial \mathbf{v}_2 \\ \vdots \\ \partial \mathbf{s} / \partial \mathbf{v}_m \end{bmatrix}}^{\partial \mathbf{s} / \partial \mathbf{v}} \cdot \overbrace{\begin{bmatrix} \Delta \mathbf{v}_1 \\ \Delta \mathbf{v}_2 \\ \vdots \\ \Delta \mathbf{v}_m \end{bmatrix}}^{\Delta \mathbf{v}} \quad (\text{A.17})$$

This gives us our derivative: it contains all of the **element-wise** derivatives we need, and in a **useful** form!

**Definition 9**

If  $s$  is a **scalar** and  $\mathbf{v}$  is an  $(m \times 1)$  **vector**, then we define the derivative or **gradient**  $\partial s / \partial \mathbf{v}$  as fulfilling:

$$\Delta s = \frac{\partial s}{\partial \mathbf{v}} \cdot \Delta \mathbf{v}$$

Or, equivalently,

$$\Delta s = \left( \frac{\partial s}{\partial \mathbf{v}} \right)^T \Delta \mathbf{v}$$

Thus, our derivative must be an  $(m \times 1)$  vector

$$\frac{\partial s}{\partial \mathbf{v}} = \begin{bmatrix} \partial s / \partial v_1 \\ \partial s / \partial v_2 \\ \vdots \\ \partial s / \partial v_m \end{bmatrix} = \begin{bmatrix} \frac{\partial s}{\partial v_1} \\ \frac{\partial s}{\partial v_2} \\ \vdots \\ \frac{\partial s}{\partial v_m} \end{bmatrix}$$

We can see the shapes work out in our matrix multiplication:

$$\overbrace{\Delta s}^{(1 \times 1)} = \overbrace{\left( \frac{\partial s}{\partial \mathbf{v}} \right)^T}^{(1 \times m)} \overbrace{\Delta \mathbf{v}}^{(m \times 1)} \quad (\text{A.18})$$



## A.3 Derivative: Vector/Scalar

Now, we want to try the flipped version: we swap our vector and our scalar.

$$\frac{\partial(\text{Vector})}{\partial(\text{Scalar})} = \frac{\partial \mathbf{w}}{\partial s} \quad (\text{A.19})$$

We'll take  $s$  to be our scalar, and  $\mathbf{w}$  to be our vector. So, our input is a **scalar**, and our output is a **vector**.

$$\Delta s \longrightarrow \boxed{f} \longrightarrow \Delta \mathbf{w} \quad (\text{A.20})$$

Written explicitly, like before:

$$\Delta s \longrightarrow \overbrace{\begin{bmatrix} \Delta w_1 \\ \Delta w_2 \\ \vdots \\ \Delta w_n \end{bmatrix}}^{\Delta \mathbf{w}} \quad (\text{A.21})$$

We have 1 **input**, that can affect  $n$  different **outputs**. So, our derivative needs to have  $n$  elements.

Again, let's look at our **approximation** rule:

$$\Delta \mathbf{w} \approx \frac{\partial \mathbf{w}}{\partial s} \star \Delta s \quad \text{or} \quad \overbrace{\begin{bmatrix} \Delta w_1 \\ \Delta w_2 \\ \vdots \\ \Delta w_n \end{bmatrix}}^{\Delta \mathbf{w}} \approx \frac{\partial \mathbf{w}}{\partial s} \star \Delta s \quad (\text{A.22})$$

Here, we can't do a **dot product**: we're multiplying our derivative by a **scalar**.

### A.3.1 Working with the vector derivative

How do we get each of our terms  $\Delta w_i$ ?

Well, each term is **separately** affected by  $\Delta s$ : we have our terms  $\partial w_i / \partial s$ .

So, if we take one of these terms **individually**, treating it as a scalar derivative, we get:

$$\Delta w_i = \frac{\partial w_i}{\partial s} \Delta s \quad (\text{A.23})$$

Since we only have **one** input, we don't have to worry about **planar** approximations: we only take one step, in the  $s$  direction.

Note that we're using vector  $\mathbf{w}$  instead of  $\mathbf{v}$  this time: this will be helpful for our vector/vector derivative: we'll need two different symbols for "a vector".

If you're ever confused with matrix math, thinking about individual elements is often a good way to figure it out!

In our matrix, we get:

$$\mathbf{w} = \begin{bmatrix} \Delta w_1 \\ \Delta w_2 \\ \vdots \\ \Delta w_n \end{bmatrix} = \begin{bmatrix} \Delta s (\partial w_1 / \partial s) \\ \Delta s (\partial w_2 / \partial s) \\ \vdots \\ \Delta s (\partial w_n / \partial s) \end{bmatrix} \quad (\text{A.24})$$

This works out for our equation above!

It could be tempting to think of our derivative  $\partial \mathbf{w} / \partial s$  as a **column vector**: we just take  $w$  and just differentiate each element. Easy!

- In fact, this *is* a valid convention. However, this conflicts with our previous derivative: they're both column vectors!
- Not only is it **confusing**, but it also will make it harder to do our **vector/vector** derivative.

So, what do we do? We refer back to the equation we used last time:

$$\Delta \mathbf{w} = \left( \frac{\partial \mathbf{w}}{\partial s} \right)^T \Delta s \quad (\text{A.25})$$

We take the **transpose**! That way, one derivative is a column vector, and the other is a row vector. And, we know that this equation works out from the work we just did.

$$\Delta \mathbf{w} = \left[ \frac{\partial w_1}{\partial s}, \frac{\partial w_2}{\partial s}, \dots, \frac{\partial w_n}{\partial s} \right]^T \Delta s \quad (\text{A.26})$$

#### Clarification 10

We mentioned that it is a valid **convention** to have that **vector derivative** be a **column vector**, and have our **gradient** be a **row vector**.

This is **not** the convention we will use in this class - you will be confused if we try!

That means, for whatever **notation** we use here, you might see the **transposed** version elsewhere. They mean exactly the **same** thing!

$$\overbrace{\Delta \mathbf{w}}^{(n \times 1)} = \overbrace{\left( \frac{\partial \mathbf{w}}{\partial s} \right)^T}^{(n \times 1)} \overbrace{\Delta s}^{(1 \times 1)} \quad (\text{A.27})$$

As we can see, the dimensions check out.

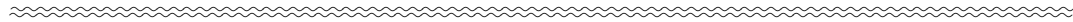
**Definition 11**

If  $s$  is a **scalar** and  $w$  is an  $(n \times 1)$  **vector**, then we define the **vector derivative**  $\partial w / \partial s$  as fulfilling:

$$\Delta w = \left( \frac{\partial w}{\partial s} \right)^T \Delta s$$

Thus, our derivative must be a  $(1 \times n)$  vector

$$\frac{\partial w}{\partial s} = \left[ \frac{\partial w_1}{\partial s}, \frac{\partial w_2}{\partial s}, \dots, \frac{\partial w_n}{\partial s} \right]$$



## A.4 Derivative: Vector/Vector

We'll be combining our two previous derivatives:

$$\frac{\partial(\text{Vector})}{\partial(\text{Vector})} = \frac{\partial \mathbf{w}}{\partial \mathbf{v}} \quad (\text{A.28})$$

$\mathbf{v}$  and  $\mathbf{w}$  are both **vectors**: thus, input and output are both **vectors**.

$$\Delta \mathbf{v} \longrightarrow \boxed{f} \longrightarrow \Delta \mathbf{w} \quad (\text{A.29})$$

Written out, we get:

$$\begin{array}{c} \overbrace{\begin{bmatrix} \Delta v_1 \\ \Delta v_2 \\ \vdots \\ \Delta v_m \end{bmatrix}}^{\Delta \mathbf{v}} \longrightarrow \begin{bmatrix} \Delta w_1 \\ \Delta w_2 \\ \vdots \\ \Delta w_n \end{bmatrix}^{\Delta \mathbf{w}} \end{array} \quad (\text{A.30})$$

Something pretty complicated! We have  $m$  inputs and  $n$  outputs. Every input can interact with every output.

So, our derivative needs to have  $mn$  different elements. That's a lot!

### A.4.1 The vector/vector derivative

We return to our rule from before. We'll skip the star notation, and jump right to the equation we've gotten for both of our two previous derivatives:

$$\Delta \mathbf{w} = \left( \frac{\partial \mathbf{w}}{\partial \mathbf{v}} \right)^T \Delta \mathbf{v} \quad (\text{A.31})$$

Hopefully, since we're combining two different derivatives, we should be able to use the same rule here.

With  $mn$  different elements, this could get messy very fast. Let's see if we can focus on only **part** of our problem:

$$\begin{bmatrix} \Delta w_1 \\ \Delta w_2 \\ \vdots \\ \Delta w_n \end{bmatrix} = \left( \frac{\partial \mathbf{w}}{\partial \mathbf{v}} \right)^T \begin{bmatrix} \Delta v_1 \\ \Delta v_2 \\ \vdots \\ \Delta v_m \end{bmatrix} \quad (\text{A.32})$$

#### One input

We could try focusing on just a single **input** or a single **output**, to simplify things. Let's start with a single  $v_i$ .

$$\overbrace{\begin{bmatrix} \Delta w_1 \\ \Delta w_2 \\ \vdots \\ \Delta w_n \end{bmatrix}}^{\Delta w \text{ from } v_i} = \left( \frac{\partial \mathbf{w}}{\partial v_i} \right)^T \Delta v_i \quad (\text{A.33})$$

We now have a simpler case:  $\partial \text{Vector} / \partial \text{Scalar}$ . We're familiar with this case!

$$\frac{\partial \mathbf{w}}{\partial v_i} = \begin{bmatrix} \frac{\partial w_1}{\partial v_i}, & \frac{\partial w_2}{\partial v_i}, & \dots & \frac{\partial w_n}{\partial v_i} \end{bmatrix} \quad (\text{A.34})$$

We get a vector. What if the **output** is a scalar instead?

**One output**

$$\Delta w_j = \left( \frac{\partial w_j}{\partial \mathbf{v}} \right)^T \begin{bmatrix} \Delta v_1 \\ \Delta v_2 \\ \vdots \\ \Delta v_m \end{bmatrix} \quad (\text{A.35})$$

We have  $\partial \text{Scalar} / \partial \text{Vector}$ :

$$\frac{\partial w_j}{\partial \mathbf{v}} = \begin{bmatrix} \partial w_j / \partial v_1 \\ \partial w_j / \partial v_2 \\ \vdots \\ \partial w_j / \partial v_m \end{bmatrix} \quad (\text{A.36})$$

So, our vector-vector derivative is a **generalization** of the two derivatives we did before!

- It seems that extending along the **vertical** axis changes our  $v_i$  value, while moving along the **horizontal** axis changes our  $w_j$  value.

~~~~~

A.5 General derivative (Vector/Vector)

You might have a hint of what we get: one derivative stretches us along **one** axis, the other along the **second**.

But now, let's prove it to ourselves.

~~~~~

Our biggest problem is that we developed these tools for working with **vectors**, and now we have a **matrix**.

Rather than giving up this perspective, we'll instead take it further:

- We can think of a matrix as a "stack of vectors".

$$M = \begin{bmatrix} a_1 & b_1 & \cdots & z_1 \\ a_2 & b_2 & \cdots & z_2 \\ \vdots & \vdots & \ddots & \vdots \\ a_n & b_n & \cdots & z_n \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} \cdots \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{bmatrix} \quad (\text{A.37})$$

- Or, a **row vector of column vectors**.

$$M = \begin{bmatrix} \vec{a} & \vec{b} & \cdots & \vec{z} \end{bmatrix} \quad (\text{A.38})$$

So, that's our plan: we first view our matrix as a row vector, hiding the column vectors inside.

- This requires treating our column vectors as if they were "scalars".

Then, we'll expand those column vectors, and see what we get.

### Concept 12

One way to **simplify** our work is to treat **vectors** as **scalars**, and then convert them back into **vectors** after applying some math.

- We have to be **careful** - any operation we apply to the pretend "**scalar**", has to match how the **vector** would behave.

~~~~~

This is **equivalent** to when just focused on one scalar inside our vector, and then stacked all those scalars back into the vector.

Here's how we apply this to our situation.

- First, we treat each **column** as a **scalar**, and the whole object as a **vector**.
 - This will use one familiar derivative.

- Then, we'll expand each **column** into a whole **vector**. Then, we have a **matrix**.
 - This will us our other derivative.

This isn't just a cute trick: it relies on an understanding that, at its **basic** level, we're treating **scalars** and **vectors** and **matrices** as the same type of object: a structured array of numbers.

We'll get into "arrays" later.

As always, our goal is to **simplify** our work, so we can handle each piece of it.

- We treat Δv as a **scalar** so we can get the **simplified** derivative.

$$\Delta w = \left(\frac{\partial w}{\partial v} \right)^T \Delta v \quad (\text{A.39})$$

We'll only expand Δw , because that's a simpler case we know how to manage.

$$\begin{bmatrix} \Delta w_1 \\ \Delta w_2 \\ \vdots \\ \Delta w_n \end{bmatrix} = \left(\frac{\partial w}{\partial v} \right)^T \Delta v \quad (\text{A.40})$$

- Notice that we **didn't** simplify v to v_i . We aren't using only one element of v : we're pretending as if the whole vector v (including all elements) is a scalar.

We compute our derivative, based on earlier work:

$$\frac{\partial w}{\partial v} = \overbrace{\left[\frac{\partial w_1}{\partial v}, \frac{\partial w_2}{\partial v}, \dots, \frac{\partial w_n}{\partial v} \right]}^{\text{Column } j \text{ matches } w_j} \quad (\text{A.41})$$

- Our "answer" is a row vector. But, each of those derivatives is a **column** vector!

Now that we've taken care of ∂w_j (one for each column), we can **expand** our derivatives in terms of ∂v_i : each is a **column vector**!

First, for w_1 :

$$\frac{\partial w}{\partial v} = \left[\begin{array}{c} \frac{\partial w_1}{\partial v_1} \\ \frac{\partial w_1}{\partial v_2} \\ \vdots \\ \frac{\partial w_1}{\partial v_m} \end{array} \right], \frac{\partial w_2}{\partial v}, \dots, \frac{\partial w_n}{\partial v} \quad \left. \vphantom{\begin{array}{c} \frac{\partial w_1}{\partial v_1} \\ \frac{\partial w_1}{\partial v_2} \\ \vdots \\ \frac{\partial w_1}{\partial v_m} \end{array}} \right\} \begin{array}{l} \text{Column } j \text{ matches } w_j \\ \text{Row } i \text{ matches } v_i \end{array} \quad (\text{A.42})$$

And again, for w_2 :

$$\frac{\partial w}{\partial v} = \left[\begin{array}{c} \frac{\partial w_1}{\partial v_1} \\ \frac{\partial w_1}{\partial v_2} \\ \vdots \\ \frac{\partial w_1}{\partial v_m} \end{array} \right], \left[\begin{array}{c} \frac{\partial w_2}{\partial v_1} \\ \frac{\partial w_2}{\partial v_2} \\ \vdots \\ \frac{\partial w_2}{\partial v_m} \end{array} \right], \dots, \frac{\partial w_n}{\partial v} \quad \left. \vphantom{\begin{array}{c} \frac{\partial w_1}{\partial v_1} \\ \frac{\partial w_1}{\partial v_2} \\ \vdots \\ \frac{\partial w_1}{\partial v_m} \end{array}} \right\} \begin{array}{l} \text{Column } j \text{ matches } w_j \\ \text{Row } i \text{ matches } v_i \end{array} \quad (\text{A.43})$$

And again, for w_n :

$$\frac{\partial w}{\partial v} = \left[\begin{array}{c} \frac{\partial w_1}{\partial v_1} \\ \frac{\partial w_1}{\partial v_2} \\ \vdots \\ \frac{\partial w_1}{\partial v_m} \end{array} \right], \left[\begin{array}{c} \frac{\partial w_2}{\partial v_1} \\ \frac{\partial w_2}{\partial v_2} \\ \vdots \\ \frac{\partial w_2}{\partial v_m} \end{array} \right], \dots, \left[\begin{array}{c} \frac{\partial w_n}{\partial v_1} \\ \frac{\partial w_n}{\partial v_2} \\ \vdots \\ \frac{\partial w_n}{\partial v_m} \end{array} \right] \quad \left. \vphantom{\begin{array}{c} \frac{\partial w_1}{\partial v_1} \\ \frac{\partial w_1}{\partial v_2} \\ \vdots \\ \frac{\partial w_1}{\partial v_m} \end{array}} \right\} \begin{array}{l} \text{Column } j \text{ matches } w_j \\ \text{Row } i \text{ matches } v_i \end{array} \quad (\text{A.44})$$

We have column vectors in our row vector... based on our new perspective, this is basically the same as a **matrix**.

Definition 13

If

- \mathbf{v} is an $(m \times 1)$ **vector**
- \mathbf{w} is an $(n \times 1)$ **vector**

Then we define the **vector derivative** $\partial \mathbf{w} / \partial \mathbf{v}$ as fulfilling:

$$\Delta \mathbf{w} = \left(\frac{\partial \mathbf{w}}{\partial \mathbf{v}} \right)^T \Delta \mathbf{v}$$

Thus, our derivative must be a $(m \times n)$ vector

$$\frac{\partial \mathbf{w}}{\partial \mathbf{v}} = \begin{matrix} & \text{Column } j \text{ matches } \mathbf{w}_j \\ \left[\begin{array}{cccc} \frac{\partial \mathbf{w}_1}{\partial v_1} & \frac{\partial \mathbf{w}_2}{\partial v_1} & \cdots & \frac{\partial \mathbf{w}_n}{\partial v_1} \\ \frac{\partial \mathbf{w}_1}{\partial v_2} & \frac{\partial \mathbf{w}_2}{\partial v_2} & \cdots & \frac{\partial \mathbf{w}_n}{\partial v_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \mathbf{w}_1}{\partial v_m} & \frac{\partial \mathbf{w}_2}{\partial v_m} & \cdots & \frac{\partial \mathbf{w}_n}{\partial v_m} \end{array} \right] & \left. \vphantom{\begin{array}{c} \frac{\partial \mathbf{w}_1}{\partial v_1} \\ \frac{\partial \mathbf{w}_1}{\partial v_2} \\ \vdots \\ \frac{\partial \mathbf{w}_1}{\partial v_m} \end{array}} \right\} \text{Row } i \text{ matches } v_i \end{matrix}$$

This general form can be used for **any** of our matrix derivatives.

So, our matrix can represent any **combination** of two elements! We just assign each **row** to a v_i component, and each **column** with a w_j component.

A.5.1 More about the vector/vector derivativeLet's show a specific example: \mathbf{w} is (3×1) , \mathbf{v} is (2×1) .

$$\frac{\partial \mathbf{w}}{\partial \mathbf{v}} = \left[\begin{array}{ccc} \overbrace{\frac{\partial \mathbf{w}_1}{\partial v_1}}^{w_1} & \overbrace{\frac{\partial \mathbf{w}_2}{\partial v_1}}^{w_2} & \overbrace{\frac{\partial \mathbf{w}_3}{\partial v_1}}^{w_3} \\ \frac{\partial \mathbf{w}_1}{\partial v_2} & \frac{\partial \mathbf{w}_2}{\partial v_2} & \frac{\partial \mathbf{w}_3}{\partial v_2} \end{array} \right] \left\{ \begin{array}{l} v_1 \\ v_2 \end{array} \right. \quad (\text{A.45})$$

Another way to describe the general case:

Notation 14

Our matrix $\partial \mathbf{w} / \partial \mathbf{v}$ is entirely filled with **scalar derivatives**

$$\frac{\partial w_j}{\partial v_i}$$

Where any one **derivative** is stored in

- Row i
 - m rows total
- Column j
 - n columns total

We can also compress it along either axis (just like how we did to derive this result):

Notation 15

Our matrix $\partial \mathbf{w} / \partial \mathbf{v}$ can be written as

$$\frac{\partial \mathbf{w}}{\partial \mathbf{v}} = \overbrace{\left[\frac{\partial w_1}{\partial \mathbf{v}}, \frac{\partial w_2}{\partial \mathbf{v}}, \dots, \frac{\partial w_n}{\partial \mathbf{v}} \right]}^{\text{Column } j \text{ matches } w_j}$$

or

$$\frac{\partial \mathbf{w}}{\partial \mathbf{v}} = \left[\begin{array}{c} \frac{\partial w}{\partial v_1} \\ \frac{\partial w}{\partial v_2} \\ \vdots \\ \frac{\partial w}{\partial v_m} \end{array} \right] \overbrace{\left. \right\}}^{\text{Row } i \text{ matches } v_i}$$

These compressed forms will be useful for deriving our new and final derivatives, **matrix-scalar** pairs.

A.6 Derivative: matrix/scalar

Now, we have our general form for creating derivatives.

We'll get our derivative of the form

$$\frac{\partial(\text{Matrix})}{\partial(\text{Scalar})} = \frac{\partial \mathbf{M}}{\partial s} \quad (\text{A.46})$$

We have a matrix \mathbf{M} in the shape $(r \times k)$ and a scalar s . Our **input** is a **scalar**, and our **output** is a **matrix**.

$$\mathbf{M} = \begin{bmatrix} m_{11} & m_{12} & \cdots & m_{1k} \\ m_{21} & m_{22} & \cdots & m_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ m_{r1} & m_{r2} & \cdots & m_{rk} \end{bmatrix} \quad (\text{A.47})$$

This may seem concerning: before, we divided **inputs** across **rows**, and **outputs** across **columns**. But in this case, we have **no** input axes, and **two** output axes.

Well, let's try to make this work anyway.

What did we do before, when we didn't know how to handle a **new** derivative?

- We compared it to **old** versions: we built our vector/vector case using the vector/s-scalar case and the scalar/vector case.
- We did this by **compressing** one of our *vectors* into a *scalar* temporarily: this works, because we want to treat each of these objects the **same way**.

We don't know how to work with **Matrix/Scalar**, but what's the **closest** thing we do know? **Vector/Scalar**.

How do we accomplish that? As we saw above, a matrix is a **vector** of **vectors**. We could turn it into a **vector** of **scalars**.

Concept 16

A **matrix** can be thought of as a **column vector** of **row vectors** (or vice versa).

So, we can use our earlier technique and convert the **row vectors** into **scalars**.

We'll replace the **row vectors** in our matrix with **scalars**.

$$\mathbf{M} = \begin{bmatrix} M_1 \\ M_2 \\ \vdots \\ M_r \end{bmatrix} \quad (\text{A.48})$$

Now, we can pretend our matrix is a vector! We've got a derivative for that:

$$\frac{\partial \mathbf{M}}{\partial \mathbf{s}} = \begin{bmatrix} \frac{\partial \mathbf{M}_1}{\partial \mathbf{s}} & \frac{\partial \mathbf{M}_2}{\partial \mathbf{s}} & \cdots & \frac{\partial \mathbf{M}_r}{\partial \mathbf{s}} \end{bmatrix} \quad (\text{A.49})$$

Aha - we have the same form that we did for our vector/scalar derivative! Each derivative is a column vector. Let's expand it out:

$$\frac{\partial \mathbf{M}}{\partial \mathbf{s}} = \left[\begin{array}{c} \left[\begin{array}{c} \frac{\partial \mathbf{m}_{11}}{\partial \mathbf{s}} \\ \frac{\partial \mathbf{m}_{12}}{\partial \mathbf{s}} \\ \vdots \\ \frac{\partial \mathbf{m}_{1k}}{\partial \mathbf{s}} \end{array} \right], \left[\begin{array}{c} \frac{\partial \mathbf{m}_{21}}{\partial \mathbf{s}} \\ \frac{\partial \mathbf{m}_{22}}{\partial \mathbf{s}} \\ \vdots \\ \frac{\partial \mathbf{m}_{2k}}{\partial \mathbf{s}} \end{array} \right], \cdots, \left[\begin{array}{c} \frac{\partial \mathbf{m}_{r1}}{\partial \mathbf{s}} \\ \frac{\partial \mathbf{m}_{r2}}{\partial \mathbf{s}} \\ \vdots \\ \frac{\partial \mathbf{m}_{rk}}{\partial \mathbf{s}} \end{array} \right] \end{array} \right] \quad \left. \begin{array}{l} \text{Column } j \text{ matches } \mathbf{m}_{j?} \\ \text{Row } i \text{ matches } \mathbf{m}_{?i} \end{array} \right\} \quad (\text{A.50})$$

Definition 17

If \mathbf{M} is a matrix in the shape $(r \times k)$ and \mathbf{s} is a scalar,

Then we define the **matrix derivative** $\partial \mathbf{M} / \partial \mathbf{s}$ as the $(k \times r)$ matrix:

$$\frac{\partial \mathbf{M}}{\partial \mathbf{s}} = \left[\begin{array}{cccc} \frac{\partial \mathbf{m}_{11}}{\partial \mathbf{s}} & \frac{\partial \mathbf{m}_{21}}{\partial \mathbf{s}} & \cdots & \frac{\partial \mathbf{m}_{r1}}{\partial \mathbf{s}} \\ \frac{\partial \mathbf{m}_{12}}{\partial \mathbf{s}} & \frac{\partial \mathbf{m}_{22}}{\partial \mathbf{s}} & \cdots & \frac{\partial \mathbf{m}_{r2}}{\partial \mathbf{s}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \mathbf{m}_{1k}}{\partial \mathbf{s}} & \frac{\partial \mathbf{m}_{2k}}{\partial \mathbf{s}} & \cdots & \frac{\partial \mathbf{m}_{rk}}{\partial \mathbf{s}} \end{array} \right] \quad \left. \begin{array}{l} \text{Column } j \text{ matches } \mathbf{m}_{j?} \\ \text{Row } i \text{ matches } \mathbf{m}_{?i} \end{array} \right\}$$

- This matrix has the shape of \mathbf{M}^T .

A.7 Derivative: scalar/matrix

We'll get our derivative of the form

$$\frac{\partial(\text{Scalar})}{\partial(\text{Matrix})} = \frac{\partial \mathbf{s}}{\partial \mathbf{M}} \quad (\text{A.51})$$

We have a matrix \mathbf{M} in the shape $(r \times k)$ and a scalar s . Our **input** is a **matrix**, and our **output** is a **scalar**.

Let's do what we did last time: break it into **row vectors**.

$$\mathbf{M} = \begin{bmatrix} \mathbf{M}_1 \\ \mathbf{M}_2 \\ \vdots \\ \mathbf{M}_r \end{bmatrix} \quad (\text{A.52})$$

The gradient for this "vector" gives us a **column vector**:

$$\frac{\partial \mathbf{s}}{\partial \mathbf{M}} = \begin{bmatrix} \frac{\partial s}{\partial \mathbf{M}_1} \\ \frac{\partial s}{\partial \mathbf{M}_2} \\ \vdots \\ \frac{\partial s}{\partial \mathbf{M}_r} \end{bmatrix} \quad (\text{A.53})$$

This time, each derivative is a **row vector**. Let's **expand**:

$$\frac{\partial \mathbf{s}}{\partial \mathbf{M}} = \begin{bmatrix} \left[\frac{\partial s}{\partial m_{11}} & \frac{\partial s}{\partial m_{12}} & \cdots & \frac{\partial s}{\partial m_{1k}} \right] \\ \left[\frac{\partial s}{\partial m_{21}} & \frac{\partial s}{\partial m_{22}} & \cdots & \frac{\partial s}{\partial m_{2k}} \right] \\ \vdots \\ \left[\frac{\partial s}{\partial m_{r1}} & \frac{\partial s}{\partial m_{r2}} & \cdots & \frac{\partial s}{\partial m_{rk}} \right] \end{bmatrix} \quad (\text{A.54})$$

Definition 18

If \mathbf{M} is a matrix in the shape $(r \times k)$ and s is a scalar,

Then we define the **matrix derivative** $\partial s / \partial \mathbf{M}$ as the $(r \times k)$ matrix:

$$\frac{\partial s}{\partial \mathbf{M}} = \begin{matrix} & \text{Column } j \text{ matches } \mathbf{m}_{\cdot j} \\ \left[\begin{array}{cccc} \frac{\partial s}{\partial m_{11}} & \frac{\partial s}{\partial m_{12}} & \cdots & \frac{\partial s}{\partial m_{1k}} \\ \frac{\partial s}{\partial m_{21}} & \frac{\partial s}{\partial m_{22}} & \cdots & \frac{\partial s}{\partial m_{2k}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial s}{\partial m_{r1}} & \frac{\partial s}{\partial m_{r2}} & \cdots & \frac{\partial s}{\partial m_{rk}} \end{array} \right] & \left. \vphantom{\begin{array}{c} \frac{\partial s}{\partial m_{11}} \\ \frac{\partial s}{\partial m_{21}} \\ \vdots \\ \frac{\partial s}{\partial m_{r1}} \end{array}} \right\} \text{Row } i \text{ matches } \mathbf{m}_i \end{matrix}$$

This matrix has the same shape as \mathbf{M} .

A.8 Tensors

A.8.1 Other Derivatives

After these, you might ask yourself, what about other derivative combinations?

$$\frac{\partial \mathbf{v}}{\partial \mathbf{M}}? \quad \frac{\partial \mathbf{M}}{\partial \mathbf{v}}? \quad \frac{\partial \mathbf{M}}{\partial \mathbf{M}^2}? \quad (\text{A.55})$$

There's a problem with all of these: the total number of axes is **too large**.

What do we mean by an **axis**?

Definition 19

An **axis** is one of the **indices** we can adjust to get a different scalar in our array: each index is a "direction" we can move along our object to **store** numbers.

- A **scalar** has **0 axes**: we only have one scalar, so we have no indices to adjust.

~~~~~

- A **vector** has **1 axis**: we can get different scalars by moving **vertically** (for column vectors):  $v_1, v_2, v_3 \dots$

$$\left[ \begin{array}{c} v_1 \\ v_2 \\ \vdots \\ v_m \end{array} \right] \left. \vphantom{\begin{array}{c} v_1 \\ v_2 \\ \vdots \\ v_m \end{array}} \right\} \text{Axis 1}$$

~~~~~

- A **matrix** has **2 axes**: we can move **horizontally** or **vertically**.

$$\overbrace{\left[\begin{array}{cccc} m_{11} & m_{12} & \cdots & m_{1r} \\ m_{21} & m_{22} & \cdots & m_{2r} \\ \vdots & \vdots & \ddots & \vdots \\ m_{k1} & m_{k2} & \cdots & m_{kr} \end{array} \right]}^{\text{Axis 2: Columns}} \left. \vphantom{\left[\begin{array}{cccc} m_{11} & m_{12} & \cdots & m_{1r} \\ m_{21} & m_{22} & \cdots & m_{2r} \\ \vdots & \vdots & \ddots & \vdots \\ m_{k1} & m_{k2} & \cdots & m_{kr} \end{array} \right]} \right\} \text{Axis 1: Rows}$$

These can also be called **dimensions**.

Why does the number of **axes** matter? Remember that, so far, for our derivatives, each axis of the output represented an axis of the **input** or **output**.

Note that last bit: we're saying a vector has one dimension. Can't a vector have **multiple** dimensions? We'll clarify this in an optional following section.

$$\frac{\partial \mathbf{w}}{\partial \mathbf{v}} = \left[\begin{array}{cccc} \frac{\partial w_1}{\partial v_1} & \frac{\partial w_2}{\partial v_1} & \dots & \frac{\partial w_n}{\partial v_1} \\ \frac{\partial w_1}{\partial v_2} & \frac{\partial w_2}{\partial v_2} & \dots & \frac{\partial w_n}{\partial v_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial w_1}{\partial v_m} & \frac{\partial w_2}{\partial v_m} & \dots & \frac{\partial w_n}{\partial v_m} \end{array} \right]$$

Column j: vertical axis of \mathbf{w}

Row i: vertical axis of \mathbf{v}

The way we currently build derivatives, we try to get **every pair** of input-output variables: we use **one** axis for each **axis** of either the **input** or **output**.

Take some examples:

- $\partial \mathbf{s} / \partial \mathbf{v}$: we need one axis to represent each term v_i .
 - 0 axis + 1 axis \rightarrow 1 axis: the output is a (column) **vector**.
- $\partial \mathbf{v} / \partial \mathbf{s}$: we need one axis to represent each term w_j .
 - 1 axis + 0 axis \rightarrow 1 axis: the output is a (row) **vector**.
- $\partial \mathbf{w} / \partial \mathbf{v}$: we need one axis to represent each term v_i , and another to represent each term w_j .
 - 1 axis + 1 axis \rightarrow 2 axes: the output is a **matrix**.
- $\partial \mathbf{M} / \partial \mathbf{s}$: we need one axis to represent the rows of \mathbf{M} , and another to represent the columns of \mathbf{M} .
 - 2 axis + 0 axis \rightarrow 2 axes: the output is a **matrix**.
- $\partial \mathbf{s} / \partial \mathbf{M}$: we need one axis to represent the rows of \mathbf{M} , and another to represent the columns of \mathbf{M} .
 - 0 axis + 2 axis \rightarrow 2 axes: the output is a **matrix**.

Notice the pattern!

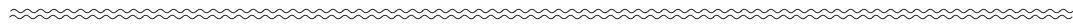
Concept 20

A **matrix derivative** needs to be able to account for each type/**index** of variable in the input **and** the output.

So, if the **input** x has m axes, and the **output** y has n axes, then the derivative needs to have the same **total** number:

$$\text{Axes}\left(\frac{\partial y}{\partial x}\right) = \text{Axes}(y) + \text{Axes}(x)$$

This is where our problem comes in: if we have a vector and a matrix, we need **3 axes**! That's more than a matrix.



A.8.2 Dimensions (Optional)

Here's a quick aside to clear up possible confusion from the last section: our definition of axes and "dimensions".

We said a vector has 1 axis, or "dimension" of movement. But, can't a vector have **multiple** dimensions?

Clarification 21

We have two competing definition of **dimension**: this explains why we can say seemingly conflicting things about derivatives.

So far, by "**dimension**", we mean, "a separate **value** we can **adjust**".

- Under this definition, a $(k \times 1)$ column **vector** has **k** dimensions: it contains **k** different scalars we can **adjust**.

$$\left[\begin{array}{c} v_1 \\ v_2 \\ \vdots \\ v_k \end{array} \right] \left. \vphantom{\begin{array}{c} v_1 \\ v_2 \\ \vdots \\ v_k \end{array}} \right\} \text{We can adjust each of our } k \text{ scalars.}$$

- You might say a $(k \times r)$ **matrix** has **k** dimensions, too: based on the **dimensionality** of its column vectors.
 - Since we prioritize the size of the vectors, we could say this is a very "vector-centric" definition.

In this section, by "dimension", we mean, "an **index** we can **adjust** (move along) to find another scalar.

- Under this definition, a $(k \times 1)$ column **vector** has **1** dimension: we only have **1** axis of **movement**.
- You might say a $(k \times r)$ **matrix** has **2** dimensions: a **horizontal** one, and a **vertical** one.
 - This **definition** is the kind we use in the following sections.

In other words:

- Vector-style dimensionality: number of separate scalars in your vector.
- Matrix(tensor)-style dimensionality: number of separate indices we can use to find scalars.

If you jumped here from X.16, feel free to follow this [link](#) back. Otherwise, continue on.

A.8.3 Dealing with Tensors

If a vector looks like a "**line**" of numbers, and a matrix looks like a "**rectangle**" of numbers, then a **3-axis** version would look like a "**box**" of numbers. How do we make sense of this?

First, what is this kind of object we've been working with? Vectors, matrices, etc. This collection of numbers, organized neatly, is an **array**.

Definition 22

An **array** of objects is an **ordered sequence** of them, stored together.

The most typical example is a **vector**: an ordered sequence of **scalars**.

A **matrix** can be thought of as a **vector** of **vectors**. For example: it could be a row vector, where every column is a column vector.

- So, we think of a matrix as a "two-dimensional array".

We can extend this to any number of dimensions. We call this kind of generalization a **tensor**.

Definition 23

In machine learning, we think of a **tensor** as a "**multidimensional array**" of numbers.

- Each "dimension" is what we have been calling an "**axis**".
- A tensor with c axes is called a **c-Tensor**.

Example: The 3-D box we are talking about above is called a 3-Tensor. We can simply think of it as a stack of matrices.

- We can imagine stacking the following two matrices in the third dimension, with the leftmost in front, and the rightmost in the back, in a $(2 \times 3 \times 2)$ block:

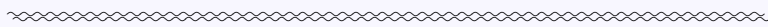
$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \quad \begin{bmatrix} 5 & 7 & 11 \\ 13 & 17 & 19 \end{bmatrix} \quad (\text{A.56})$$

Clarification 24

Note that what we call a tensor is **not** a mathematical (or physics) tensor.

- We don't make use of most of the crucial properties of a mathematical tensor.

Our tensor can be better thought of as a "**generalized matrix**", or a "multidimensional array".



This is important, because a "mathematical" tensor has properties that can be confusing from an ML perspective:

- The "tensor product" doesn't behave at all like a matrix product.
- Rather, the generalized version of the **matrix product** is something called **tensor contraction**.

Tensor contraction examples, like the einsum function in numpy ("einstein summation notation"), can get very complex for higher-dimensional tensors.

If tensors don't really behave quite like matrices, and their math is complex, how do we handle **tensors**?

Simply, we convert them into regular **matrices** in some way, and then do our usual math on them:

- If a tensor has a pattern of **zeroes**, we might be able to flatten it into a matrix.
 - For example, in this matrix, we make a vector out of the diagonal:

$$\begin{bmatrix} 3 & 0 & 0 \\ 0 & 9 & 0 \\ 0 & 0 & 4 \end{bmatrix} \rightarrow \begin{bmatrix} 3 \\ 9 \\ 4 \end{bmatrix} \quad (\text{A.57})$$

These examples aren't especially important, but you can see different variations used for different problems!

- We can also flatten it into a matrix or vector by **stacking** layers next to each other in the same dimension.
 - For example, we could stack the two columns of a matrix:

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 3 \\ 2 \\ 4 \end{bmatrix} \quad (\text{A.58})$$

This version is used when we are willing to give up our n-D structure: we lose a lot of information.

- We cleverly "**hide**" dimensions of a matrix, as we have before: treat a scalar as a vector, etc.

This works best for high-level, conceptual stuff. We'll use it below.

- We **systematically** multiply and add our elements in a way that gives us the derivatives we want, replicating matrix multiplication.
 - This is **tensor contraction**.
 - This method is often used by softwares to implement the chain rule.

Clarification 25

If you look into **derivatives** that would result in a **3-tensor** or higher, you'll find that there's no consistent **notation** for what these derivatives look like: shape, structure, etc.

These techniques are part of why: there are **different** approaches for how to approach these objects.

The solution tends to be directly computing the chain rule, rather than figuring out the structure of the abstract object.

As we will see in the next chapter, tensors are **very** important to machine learning.

However, because they're so troublesome to work with, we'll convert to matrices in most cases.

~~~~~

## A.9 Chapter 7 Derivatives

### A.9.1 The loss derivative

Finally, we apply all of our new knowledge to the common derivatives in section 7.5.

$$\overbrace{\frac{\partial \mathcal{L}}{\partial \mathbf{A}^L}}^{(n^L \times 1)} \quad (\text{A.59})$$

Loss is not given, so we can't compute it. But, we can get the shape: we have a scalar/vector derivative, so the shape matches  $\mathbf{A}^L$ .

#### Notation 26

Our derivative

$$\frac{\partial \mathcal{L}}{\partial \mathbf{A}^L}$$

Is a scalar/vector derivative, and thus the shape  $(n^L \times 1)$ .

### A.9.2 The weight derivative

$$\overbrace{\frac{\partial \mathbf{Z}^\ell}{\partial \mathbf{W}^\ell}}^{(m^\ell \times 1)?} \quad (\text{A.60})$$

This derivative is difficult - it's a derivative in the form vector/matrix. With **three** axes, a 3-tensor seems suitable.

But, our goal is to use this for the **chain rule**: so, we need to make matrix shapes that **match**.

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^\ell} = \overbrace{\frac{\partial \mathbf{Z}^\ell}{\partial \mathbf{W}^\ell}}^{\text{Weight link}} \cdot \overbrace{\left( \frac{\partial \mathcal{L}}{\partial \mathbf{Z}^\ell} \right)^T}^{\text{Other layers}} \quad (\text{A.61})$$

Remember that we had to do weird reordering and transposing in chapter 7 to flip the order of the chain rule? This is why: shape consistency.

Our problem is we have **too many axes**: the easiest way to resolve this to **break up** our matrix.

$$\mathbf{W} = \begin{bmatrix} \mathbf{W}_1 & \mathbf{W}_2 & \cdots & \mathbf{W}_n \end{bmatrix} \quad (\text{A.62})$$

Notice that, this time, we broke  $\mathbf{W}$  into **column vectors**  $\mathbf{W}_i$ , rather than row vectors. Each **neuron's** weights are represented by one column vector.

So, for now, we focus on only **one neuron** at a time: it has a column vector  $W_i$ .  
We'll ignore everything except  $W_i$ .

For simplicity, we're gonna ignore the  $\ell$  notation: just be careful, because  $Z$  and  $A$  are from two different layers!

$$W_i = \begin{bmatrix} w_{1i} \\ w_{2i} \\ \vdots \\ w_{mi} \end{bmatrix} \quad (\text{A.63})$$

### Concept 27

When you have a derivative that has **too many dimensions**, it's easier to only focus on **one** of the elements/dimensions, and find the derivative of that first.

- In this example, rather than finding  $\frac{\partial Z}{\partial W}$ , we're finding  $\frac{\partial Z}{\partial W_i}$ .
- We're simplifying from **matrix**  $W$  to **vector**  $W_i$ .

~~~~~  
So, we're doing $\frac{\partial Z}{\partial W_i}$: we need equations relating these variables.

- W_i represents the weights of one neuron, while Z represents the pre-activation of all the neurons.

Let's focus on one pre-activation neuron at a time.

$$z_j = W_j^T A \quad (\text{A.64})$$

Notably, this equation reminds us that pre-activation z_j is only affected by weights from the same neuron, W_j .

- Weights from a different neuron have no effect on z_j .
- So, the derivative between them will be 0: we can check this using the equation above.

Concept 28

The i^{th} neuron's **weights**, W_i , have **no effect** on a different neuron's **pre-activation** z_j .

So, if the neurons don't match, then our derivative is **zero**:

- i is the neuron for **pre-activation** z_i
- j is the j^{th} weight **in neuron** k .
- k is the **neuron** for weight vector W_k

$$\frac{\partial z_i}{\partial W_{jk}} = 0 \quad \text{if } i \neq k$$

So, our only nonzero derivatives are

$$\frac{\partial z_i}{\partial W_{ji}}$$

~~~~~  
With that out of the way, let's actually **expand** our expression from above, to compute  $\frac{\partial z_i}{\partial W_{ji}}$ .

$$z_i = W_i^T A \tag{A.65}$$

$$z_i = \begin{bmatrix} w_{1i} & w_{2i} & \cdots & w_{mi} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix} \tag{A.66}$$

This matrix multiplication can be written as an easier-to-use **sum**:

$$z_i = \sum_{j=1}^n W_{ji} a_j \tag{A.67}$$

Finally, we can get our derivatives, for the **non-zero** terms:

$$\frac{\partial z_i}{\partial W_{ji}} = a_j \tag{A.68}$$

Now, we can combine them into a vector.

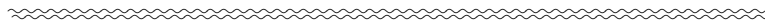


$$\frac{\partial z_i}{\partial \mathbf{W}_i} = \begin{bmatrix} \frac{\partial z_i}{\partial W_{1i}} \\ \frac{\partial z_i}{\partial W_{2i}} \\ \vdots \\ \frac{\partial z_i}{\partial W_{mi}} \end{bmatrix} \quad (\text{A.69})$$

We plug in our new values:

$$\frac{\partial z_i}{\partial \mathbf{W}_i} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix} = \mathbf{A} \quad (\text{A.70})$$

We get a result!



What if the pre-activation  $z_i$  and weights  $W_k$  don't match? We've already seen: the derivative is 0: weights from one neuron don't affect different neurons.

$$\frac{\partial z_i}{\partial W_{jk}} = 0 \quad \text{if } i \neq k \quad (\text{A.71})$$

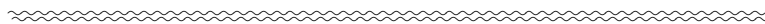
We can combine these into a **zero vector**:

$$\frac{\partial z_i}{\partial W_k} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \vec{0} \quad \text{if } i \neq k \quad (\text{A.72})$$

So, now, we can describe all of our vector components:

$$\frac{\partial z_i}{\partial W_k} = \begin{cases} a & \text{if } i = k \\ \vec{0} & \text{if } i \neq k \end{cases} \quad (\text{A.73})$$

Again, we see: "same neuron? they're related. Different neurons? They're not."



This derivative, despite being a vector, can be represented with a **single** symbol ( $\vec{0}$  or  $A$ ), for each element  $\frac{\partial z_i}{\partial W_k}$ .

- For the sake of trying to wrangle our 3-tensor, we'll "hide" one dimension using this fact.
- If this is the case, then we have to pretend that  $W_k$  is a "scalar": so,  $W$  is now a "vector".

We compute  $\frac{\partial Z}{\partial W}$  as a "vector/vector" derivative: this is a **matrix**!

$$\frac{\partial Z}{\partial W} = \begin{bmatrix} A & \vec{0} & \dots & \vec{0} \\ \vec{0} & A & \dots & \vec{0} \\ \vdots & \vdots & \ddots & \vec{0} \\ \vec{0} & \vec{0} & \vec{0} & A \end{bmatrix} \quad (\text{A.74})$$

We have our result: it turns out, despite being stored in a **matrix**-like format, this is actually a **3-tensor**! Each "scalar" of our **matrix** is a **vector**: we really have 3 axes.

~~~~~

But, we don't really... *want* a tensor. It doesn't have the right shape, and we can't do matrix multiplication.

We'll solve this by **simplifying**, without losing key information.

Concept 29

For many of our "tensors" resulting from matrix derivatives, they contain **empty** rows or **redundant** information.

Based on this, we can **simplify** our tensor into a fewer-dimensional (fewer axes) object.

We can see two types of **redundancy** above:

- Every element **off** the diagonal is 0.
- Every element **on** the diagonal is the same.

Let's fix the first one: we'll go from a diagonal matrix to a column vector.

$$\begin{bmatrix} A & \vec{0} & \dots & \vec{0} \\ \vec{0} & A & \dots & \vec{0} \\ \vdots & \vdots & \ddots & \vec{0} \\ \vec{0} & \vec{0} & \vec{0} & A \end{bmatrix} \rightarrow \begin{bmatrix} A \\ A \\ \vdots \\ A \end{bmatrix} \quad (\text{A.75})$$

Then, we'll combine all of our redundant A values.

$$\begin{bmatrix} \mathbf{A} \\ \mathbf{A} \\ \vdots \\ \mathbf{A} \end{bmatrix} \rightarrow \mathbf{A} \quad (\text{A.76})$$

We have our final answer!

Notation 30

Our derivative

$$\overbrace{\frac{\partial \mathbf{Z}^\ell}{\partial \mathbf{W}^\ell}}^{(m^\ell \times 1)} = \mathbf{A}^{\ell-1}$$

Is a vector/matrix derivative, and thus should be a 3-tensor.

But, we have turned it into the shape $(m^\ell \times 1)$.

This is as **condensed** as we can get our information: if we compress to a scalar, we lose some of our elements.

- Even with this derivative, we still have to do some clever **reshaping** to get the result we need (transposing, changing multiplication order of the chain rule, etc.)

However, at the end, we get the right shape for our chain rule!

This is the weird order change we mentioned in chapter 7, where we reverse the order of elements.

A.9.3 Linking Layers $\ell - 1$ and ℓ

$$\frac{\partial \mathbf{Z}^\ell}{\partial \mathbf{A}^{\ell-1}} \quad (\text{A.77})$$

This derivative is much more manageable: it's just the derivative between a vector and a vector. Let's look at our equation again:

Ignoring superscripts ℓ , as before.

$$\mathbf{Z} = \mathbf{W}^T \mathbf{A} \quad (\text{A.78})$$

We'll use the same approach we did last section: \mathbf{W} is a "vector", and we'll focus on \mathbf{W}_i . This will allow us to break it up **element-wise**.

- We could treat \mathbf{W} as a whole matrix, but our approach will be cleaner: otherwise, we'd have to depict every \mathbf{W}_i at **once**.

$$W = \begin{bmatrix} W_1 & W_2 & \cdots & W_n \end{bmatrix} \quad W_i = \begin{bmatrix} w_{1i} \\ w_{2i} \\ \vdots \\ w_{mi} \end{bmatrix} \quad (\text{A.79})$$

Here's our equation:

$$z_i = \begin{bmatrix} w_{1i} & w_{2i} & \cdots & w_{mi} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix} \quad (\text{A.80})$$

We matrix multiply:

$$z_i = \sum_{j=1}^n W_{ji} a_j \quad (\text{A.81})$$

The derivative can be gotten from here -

$$\frac{\partial z_i}{\partial a_j} = W_{ji} \quad (\text{A.82})$$

We look at our whole matrix derivative:

$$\frac{\partial Z}{\partial A} = \left[\begin{array}{ccc} \ddots & \vdots & \ddots \\ \cdots & \frac{\partial z_i}{\partial a_j} & \cdots \\ \ddots & \vdots & \ddots \end{array} \right] \quad \left. \begin{array}{l} \text{Column } i \text{ matches } z_i \\ \text{Row } j \text{ matches } a_j \end{array} \right\} \quad (\text{A.83})$$

This notation looks a bit weird, but it's just a way to represent that all of our elements follow the same pattern.

Wait.

- The derivative $\partial z_i / \partial a_j$ is in the j^{th} row, i^{th} column of $\frac{\partial Z}{\partial A}$.
- W_{ji} is the element of W in the j^{th} row, i^{th} column.

W_{ji} goes in the same place in both matrices! They're the same matrix: $W = \frac{\partial Z}{\partial A}$

We get our final result:

If two matrices have exactly the same shape and elements, they're the same matrix.

Notation 31

Our derivative

$$\overbrace{\frac{\partial \mathbf{Z}^\ell}{\partial \mathbf{A}^{\ell-1}}}^{(m^\ell \times n^\ell)} = \mathbf{W}^\ell$$

Is a vector/vector derivative, and thus a matrix.

It takes the shape $(m^\ell \times n^\ell)$.

- ~~~~~
- This matches the intuition we have from the simplified, 1d version, where we have $z = wa$, instead of $Z = W^T A$.

A.9.4 Activation Function

$$\frac{\partial \mathbf{A}^\ell}{\partial \mathbf{Z}^\ell} \tag{A.84}$$

The last derivative is less strange than its solution looks.

$$\mathbf{A}^\ell = f(\mathbf{Z}^\ell) \longrightarrow \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = f \left(\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{bmatrix} \right) \tag{A.85}$$

We can apply our function element-wise:

$$\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} f(z_1) \\ f(z_2) \\ \vdots \\ f(z_n) \end{bmatrix} \tag{A.86}$$

As we can see, each activation is a function of only **one** pre-activation.

Concept 32

Each **activation** is only affected by the **pre-activation** in the **same neuron**.

So, if the **neurons** don't match, then our derivative is zero:

- i is the neuron for pre-activation z_i
- j is the neuron for activation a_j

$$\frac{\partial a_j}{\partial z_i} = 0 \quad \text{if } i \neq j$$

So, our only nonzero derivatives are

$$\frac{\partial a_i}{\partial z_i}$$

As for our **non-zero** terms, they all rely on the equation:

$$a_i = f(z_i) \tag{A.87}$$

Our derivative is:

$$\frac{\partial a_i}{\partial z_i} = f'(z_i) \tag{A.88}$$

In general, including the non-diagonals:

$$\frac{\partial a_i}{\partial z_j} = \begin{cases} f'(z_i) & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases} \tag{A.89}$$

This gives us our result:

Notation 33

Our derivative

$$\underbrace{\frac{\partial \mathbf{A}^\ell}{\partial \mathbf{Z}^\ell}}_{(n^\ell \times n^\ell)} = \left[\begin{array}{ccccc} \text{Column } j \text{ matches } \mathbf{a}_j & & & & \\ f'(z_1^\ell) & 0 & 0 & \cdots & 0 \\ 0 & f'(z_2^\ell) & 0 & \cdots & 0 \\ 0 & 0 & f'(z_3^\ell) & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & f'(z_n^\ell) \end{array} \right] \left. \vphantom{\begin{array}{c} f'(z_1^\ell) \\ f'(z_2^\ell) \\ f'(z_3^\ell) \\ \vdots \\ f'(z_n^\ell) \end{array}} \right\} \text{Row } i \text{ matches } z_i$$

Is a vector/vector derivative, and thus a matrix.

It takes the shape $(n^\ell \times n^\ell)$.**A.9.5 Element-wise multiplication**

Notice that, in the previous section, we could've compressed this matrix down to remove the unnecessary 0's:

$$\begin{bmatrix} f'(z_1^\ell) \\ f'(z_2^\ell) \\ \vdots \\ f'(z_n^\ell) \end{bmatrix} \quad (\text{A.90})$$

This is a valid way to interpret this matrix! The only thing we need to be careful of: if we were to use this in a chain rule, we couldn't do normal matrix multiplication.

However, because of how this matrix works, you can just do **element-wise** multiplication instead!

You can check it for yourself: each index is separately scaled.

Concept 34

When multiplying two vectors R and Q , if they take the form

$$R = \begin{bmatrix} r_1 & 0 & 0 & \cdots & 0 \\ 0 & r_2 & 0 & \cdots & 0 \\ 0 & 0 & r_3 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & r_n \end{bmatrix} \quad Q = \begin{bmatrix} q_1 \\ q_2 \\ q_3 \\ \vdots \\ q_n \end{bmatrix}$$

Then we can write their product each of these ways:

$$RQ = \overbrace{R * Q}^{\text{Element-wise multiplication}} = \begin{bmatrix} r_1 q_1 \\ r_2 q_2 \\ r_3 q_3 \\ \vdots \\ r_n q_n \end{bmatrix}$$

So, we can substitute the chain rule this way.

A.10 Terms

- Matrix/scalar derivative
- Scalar/Matrix derivative
- Axis
- Dimension (vector)
- Dimension (array)
- Array
- "Tensor" (Generalized matrix)
- c-Tensor (2-tensor, 3-tensor, etc.)