# Explanatory Notes for 6.390

Shaunticlair Ruiz (Current TA)

Fall 2022

# X. Matrix Derivatives

In general, we want to be able 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

$$w^\mathsf{T} x = \begin{bmatrix} w_1 & w_2 & \cdots & w_m \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} = \sum_{i=1}^{m} x_i w_i \tag{1}$$

  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

$$\frac{\partial f}{\partial x_1} = 10 \qquad \frac{\partial f}{\partial x_2} = -5 \tag{2}$$

  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**.

- This "data" is 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 $b$.

$$\overbrace{(\mathbf{a}\times\mathbf{b})}^{X}\;\overbrace{(\mathbf{b}\times\mathbf{c})}^{Y} \tag{3}$$

We can't do anything that would **violate** this rule: otherwise, our **equations** don't make sense, and we get stuck. This means we need to build our math carefully.

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

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

## X.1   Review: 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 used when there may be **multiple variables** in our functions.

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

---

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

$$\frac{\partial f}{\partial x} = 2(2x)y \tag{4}$$

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.

If this is too complicated, we can change those 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 $\Delta x_i$ 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 $\{x_1, x_2, \ldots x_k\}$ we can generalize this as:

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

---

## X.2   Thinking about derivatives

The typical definition of derivatives

$$\lim_{h \to 0} \frac{f(x + h) - f(x)}{h} \tag{5}$$

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".

$$A = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} \qquad B = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} \tag{6}$$

**Example:** Suppose we want a derivative $\partial B / \partial A$: $\Delta a_1, \Delta a_2$, and $\Delta a_3$ could each, separately, have an effect on $\Delta b_1$ and/or $\Delta b_2$. That requires 6 different derivatives, $\partial b_i / \partial a_j$. _____ | 3 dimensions of A times 2 dimensions of B: 6 combinations.

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

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?

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

## X.3  Derivatives: Approximation

To answer this, we need to ask ourselves *why* we care about these derivatives: their **structure** will be based on what we need them for.

- We care about the **direction of greatest decrease**, the gradient. For example, we might want to adjust weight vector $w$ to reduce $\mathcal{L}$.

- We also want other derivatives that have the **same** behavior, so we can combine them using the **chain rule**.

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 $\mathcal{L}$.

$$\frac{\partial \mathcal{L}}{\partial w} \approx \frac{\text{Change in } \mathcal{L}}{\text{Change in } w} = \frac{\Delta\mathcal{L}}{\Delta w} \tag{7}$$

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

> All we do is multiply both sides by $\Delta w$.

$$\Delta\mathcal{L} \approx \frac{\partial \mathcal{L}}{\partial w}\Delta w \tag{8}$$

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 $\Delta w$, and returns an **approximation** of the loss.

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

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

## X.4    The Gradient: a vector input, scalar output

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

First, we consider:

$$\frac{\partial (\text{Scalar})}{\partial (\text{Vector})} = \frac{\partial s}{\partial v} \tag{9}$$

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

$$\Delta v \longrightarrow \boxed{f} \longrightarrow \Delta s \tag{10}$$

How do we make sense of this? Well, let's write $\Delta v_i$ explicitly:

$$\overbrace{\begin{bmatrix} \Delta v_1 \\ \Delta v_2 \\ \vdots \\ \Delta v_m \end{bmatrix}}^{\Delta v} \longrightarrow \Delta s \tag{11}$$

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 $v_i$.

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

## X.5    Finding the scalar/vector derivative

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

$$\Delta s \approx \frac{\partial s}{\partial v} \star \Delta v \qquad \text{or} \qquad \Delta s \approx \frac{\partial s}{\partial v} \star \overbrace{\begin{bmatrix} \Delta v_1 \\ \Delta v_2 \\ \vdots \\ \Delta v_m \end{bmatrix}}^{\Delta v} \qquad (12)$$

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

> It's usually possible to change each $v_i$, so we have to look at every one of them.

One problem, though: how can we treat each **derivative** separately? Each $\Delta v_i$ will move our position, which can change a different derivative $v_k$: they can **affect** each other.

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

## X.6    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.

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

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

- If we go from $x = 1 \to 2$, then our derivative goes from $f'(x) = 2 \to 4$.

- Let's **shrink** our step. We go from $x = 1 \to 1.01$, our derivative goes from $f'(x) = 2 \to 2.02$.

  – Our derivative is almost the same!

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

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

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

---

**Concept 6**

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

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

---

**Clarification**

This section 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. $\Delta f'(x)$ is **small**.

Let's say it can't change more then $\delta$.

If you want

- $\Delta f'(x)$ to be very small ($|\Delta f'(x)| < \delta$)

- It has been proven that...

    - can 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

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

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**.
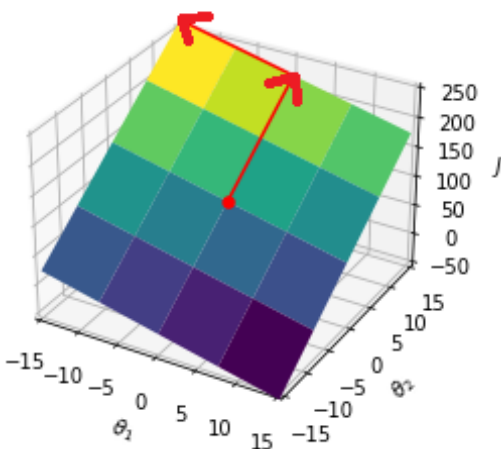
Movement in $\theta_1$ on $J$



If we take very small steps, we can approximate our function as **flat**.
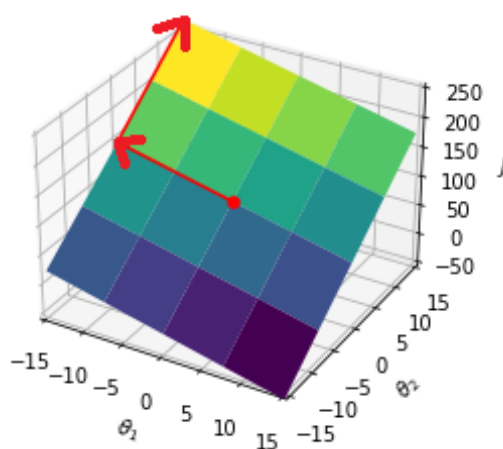
Why does this help? If our derivative doesn't **change**, we can combine multiple steps You can take multiple steps $\Delta v_i$ and the order doesn't matter. _____

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

Combining two movements                    Combining two movements



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:

$$\Delta s = \Delta s_{\text{from } v_1} + \Delta s_{\text{from } v_2} + \cdots + \Delta s_{\text{from } v_m} \tag{13}$$

## X.7 Our 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

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

In general, we have

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

---

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

This gives us what we're looking for:

$$\Delta s \approx \sum_{i=1}^{m} \frac{\partial s}{\partial v_i}\Delta v_i \tag{14}$$

If we circle back around to our original approximation:

$$\sum_{i=1}^{m} \frac{\partial s}{\partial v_i}\Delta v_i = \frac{\partial s}{\partial v} \star \overbrace{\begin{bmatrix} \Delta v_1 \\ \Delta v_2 \\ \vdots \\ \Delta v_m \end{bmatrix}}^{\Delta v} \tag{15}$$

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

$$\sum_{i=1}^{m} \frac{\partial s}{\partial v_i} \Delta v_i \;=\; \overbrace{\begin{bmatrix} \partial s/\partial v_1 \\[6pt] \partial s/\partial v_2 \\[6pt] \vdots \\[6pt] \partial s/\partial v_m \end{bmatrix}}^{\partial s/\partial v} \cdot \overbrace{\begin{bmatrix} \Delta v_1 \\[6pt] \Delta v_2 \\[6pt] \vdots \\[6pt] \Delta v_m \end{bmatrix}}^{\Delta v} \tag{16}$$

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 $v$ is an $(m \times 1)$ **vector**, then we define the **derivative** or **gradient** $\partial s/\partial v$ as fulfilling:

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

Or, equivalently,

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

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

$$\frac{\partial s}{\partial v} = \begin{bmatrix} \partial s/\partial v_1 \\[6pt] \partial s/\partial v_2 \\[6pt] \vdots \\[6pt] \partial s/\partial v_m \end{bmatrix} = \begin{bmatrix} \dfrac{\partial s}{\partial v_1} \\[10pt] \dfrac{\partial s}{\partial v_2} \\[10pt] \vdots \\[10pt] \dfrac{\partial s}{\partial v_m} \end{bmatrix}$$

---

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

$$\overbrace{\Delta s}^{(1\times 1)} \;=\; \overbrace{\left( \frac{\partial s}{\partial v} \right)^{\mathsf{T}}}^{(1\times m)} \overbrace{\Delta v}^{(m\times 1)} \tag{17}$$

## X.8    Vector derivative: a scalar input, vector output

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

$$\frac{\partial(\text{Vector})}{\partial(\text{Scalar})} = \frac{\partial w}{\partial s} \tag{18}$$

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

> Note that we're using vector $w$ instead of $v$ this time: this will be helpful for our vector/vector derivative: we can use both.

$$\Delta s \longrightarrow \boxed{f} \longrightarrow \Delta w \tag{19}$$

Written explicitly, like before:

$$\Delta s \longrightarrow \overbrace{\begin{bmatrix} \Delta w_1 \\ \Delta w_2 \\ \vdots \\ \Delta w_n \end{bmatrix}}^{\Delta w} \tag{20}$$

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:

$$\Delta w \approx \frac{\partial w}{\partial s} \star \Delta s \qquad \text{or} \qquad \overbrace{\begin{bmatrix} \Delta w_1 \\ \Delta w_2 \\ \vdots \\ \Delta w_n \end{bmatrix}}^{\Delta w} \approx \frac{\partial w}{\partial s} \star \Delta s \tag{21}$$

Here, we can't do a **dot product**: we're multiplying our derivative by a **scalar**. Plus, we'd get the **same shape** as before: we might **mix up** our derivatives.

〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜

## X.9    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 these terms **individually**, treating it as a scalar derivative, we get:

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

$$\Delta w_i = \frac{\partial w_i}{\partial s} \Delta s \tag{22}$$

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

In our matrix, we get:

$$w = \begin{bmatrix} \Delta w_1 \\ \Delta w_2 \\ \vdots \\ \Delta w_n \end{bmatrix} = \begin{bmatrix} \Delta s(\partial w_1/\partial s) \\ \Delta s(\partial w_2/\partial s) \\ \vdots \\ \Delta s(\partial w_n/\partial s) \end{bmatrix} \tag{23}$$

This works out for our equation above!

It could be tempting to think of our derivative $\partial w/\partial s$ 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:

$$\Delta w = \left( \frac{\partial w}{\partial s} \right)^\mathsf{T} \Delta s \tag{24}$$

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.

$$\Delta w = \begin{bmatrix} \dfrac{\partial w_1}{\partial s}, & \dfrac{\partial w_2}{\partial s}, & \cdots & \dfrac{\partial w_n}{\partial s} \end{bmatrix}^\mathsf{T} \Delta s \tag{25}$$

> **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!

$$\overbrace{\Delta w}^{(n \times 1)} = \overbrace{\left( \frac{\partial w}{\partial s} \right)^\mathsf{T}}^{(n \times 1)} \overbrace{\Delta s}^{(1 \times 1)} \tag{26}$$

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:

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

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

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

---

## X.10    Vectors and vectors: vector input, vector output

We'll be combining our two previous derivatives:

$$\frac{\partial (\text{Vector})}{\partial (\text{Vector})} = \frac{\partial w}{\partial v} \tag{27}$$

$v$ and $w$ are both **vectors**: thus, input and output are both **vectors**.

$$\Delta v \longrightarrow \boxed{f} \longrightarrow \Delta w \tag{28}$$

Written out, we get:

$$\overbrace{\begin{bmatrix} \Delta v_1 \\ \Delta v_2 \\ \vdots \\ \Delta v_m \end{bmatrix}}^{\Delta v} \longrightarrow \overbrace{\begin{bmatrix} \Delta w_1 \\ \Delta w_2 \\ \vdots \\ \Delta w_n \end{bmatrix}}^{\Delta w} \tag{29}$$

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!

---

## X.11    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:

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

$$\Delta w = \left( \frac{\partial w}{\partial v} \right)^{\mathsf{T}} \Delta v \tag{30}$$

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

$$\begin{bmatrix} \Delta w_1 \\ \Delta w_2 \\ \vdots \\ \Delta w_n \end{bmatrix} = \left( \frac{\partial w}{\partial v} \right)^{\mathsf{T}} \begin{bmatrix} \Delta v_1 \\ \Delta v_2 \\ \vdots \\ \Delta v_m \end{bmatrix} \tag{31}$$

**One input**

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

$$\overbrace{\begin{bmatrix} \Delta w_1 \\ \Delta w_2 \\ \vdots \\ \Delta w_n \end{bmatrix}}^{\Delta w \text{ from } v_i} = \left( \frac{\partial w}{\partial v_i} \right)^{\mathsf{T}} \Delta v_i \tag{32}$$

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

$$\frac{\partial w}{\partial v_i} = \begin{bmatrix} \frac{\partial w_1}{\partial v_i}, & \frac{\partial w_2}{\partial v_i}, & \cdots & \frac{\partial w_n}{\partial v_i} \end{bmatrix} \tag{33}$$

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

**One output**

$$\Delta w_j = \left( \frac{\partial w_j}{\partial v} \right)^{\mathsf{T}} \begin{bmatrix} \Delta v_1 \\ \Delta v_2 \\ \vdots \\ \Delta v_m \end{bmatrix} \tag{34}$$

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

$$\frac{\partial w_j}{\partial v} = \begin{bmatrix} \partial w_j/\partial v_1 \\\\ \partial w_j/\partial v_2 \\\\ \vdots \\\\ \partial w_j/\partial v_m \end{bmatrix} \tag{35}$$

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.

〰〰〰〰〰〰〰〰〰〰〰〰〰〰〰〰〰〰〰〰〰〰〰〰〰〰

## X.12    General derivative

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

To prove it to ourselves, we can **combine** these concepts. We'll handle solve as if we have one vector, and then **substitute** in the second one.

---

**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 **scalar**, has to match how the **vector** would behave.

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

---

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 $\Delta v$ as a scalar so we can get the simplified derivative.

$$\Delta w = \left( \frac{\partial w}{\partial v} \right)^{\mathsf{T}} \Delta v \tag{36}$$

We'll only expand **one** of our vectors, since we know how to manage **one** of them.

$$\begin{bmatrix} \Delta w_1 \\ \Delta w_2 \\ \vdots \\ \Delta w_n \end{bmatrix} = \left(\frac{\partial w}{\partial v}\right)^{\mathsf{T}} \Delta v \qquad (37)$$

This time, notice that we **didn't** simplify $v$ to $v_i$. We didn't **remove** the other elements - we still have a full **vector**. But, let's treat it as if it *were* a scalar.

This comes out to:

$$\overbrace{\frac{\partial w}{\partial v} = \begin{bmatrix} \dfrac{\partial w_1}{\partial v}, & \dfrac{\partial w_2}{\partial v}, & \cdots & \dfrac{\partial w_n}{\partial v} \end{bmatrix}}^{\text{Column j matches } w_j} \qquad (38)$$

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

Now that we've taken care of $\partial w_j$ (one for each column), we can expand our derivatives in terms of $\partial v_i$.

First, for $w_1$:

$$\frac{\partial w}{\partial v} = \overbrace{\left.\begin{bmatrix} \begin{bmatrix} \dfrac{\partial w_1}{\partial v_1} \\[2mm] \dfrac{\partial w_1}{\partial v_2} \\[2mm] \vdots \\[2mm] \dfrac{\partial w_1}{\partial v_m} \end{bmatrix}, & \dfrac{\partial w_2}{\partial v}, & \cdots & \dfrac{\partial w_n}{\partial v} \end{bmatrix}\right\}}^{\text{Column j matches } w_j} \text{Row i matches } v_i \qquad (39)$$

And again, for $w_2$:

$$\frac{\partial w}{\partial v} = \overbrace{\left.\begin{bmatrix} \begin{bmatrix} \dfrac{\partial w_1}{\partial v_1} \\[2mm] \dfrac{\partial w_1}{\partial v_2} \\[2mm] \vdots \\[2mm] \dfrac{\partial w_1}{\partial v_m} \end{bmatrix}, & \begin{bmatrix} \dfrac{\partial w_2}{\partial v_1} \\[2mm] \dfrac{\partial w_2}{\partial v_2} \\[2mm] \vdots \\[2mm] \dfrac{\partial w_2}{\partial v_m} \end{bmatrix}, & \cdots & \dfrac{\partial w_n}{\partial v} \end{bmatrix}\right\}}^{\text{Column j matches } w_j} \text{Row i matches } v_i \qquad (40)$$

And again, for $w_n$:

$$\frac{\partial w}{\partial v} = \underbrace{\left[ \begin{bmatrix} \dfrac{\partial w_1}{\partial v_1} \\[2ex] \dfrac{\partial w_1}{\partial v_2} \\[2ex] \vdots \\[2ex] \dfrac{\partial w_1}{\partial v_m} \end{bmatrix}, \begin{bmatrix} \dfrac{\partial w_2}{\partial v_1} \\[2ex] \dfrac{\partial w_2}{\partial v_2} \\[2ex] \vdots \\[2ex] \dfrac{\partial w_2}{\partial v_m} \end{bmatrix}, \cdots \begin{bmatrix} \dfrac{\partial w_n}{\partial v_1} \\[2ex] \dfrac{\partial w_n}{\partial v_2} \\[2ex] \vdots \\[2ex] \dfrac{\partial w_n}{\partial v_m} \end{bmatrix} \right]}_{\text{Column } j \text{ matches } w_j} \left.\vphantom{\begin{bmatrix} a \\ b \\ c \end{bmatrix}}\right\} \text{Row } i \text{ matches } v_i \qquad (41)$$

We have column vectors in our row vector... let's just combine them into a **matrix**.

---

**Definition 13**

If

- $v$ is an $(m \times 1)$ **vector**

- $w$ is an $(n \times 1)$ **vector**

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

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

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

$$\frac{\partial w}{\partial v} = \underbrace{\begin{bmatrix} \dfrac{\partial w_1}{\partial v_1} & \dfrac{\partial w_2}{\partial v_1} & \cdots & \dfrac{\partial w_n}{\partial v_1} \\[2ex] \dfrac{\partial w_1}{\partial v_2} & \dfrac{\partial w_2}{\partial v_2} & \cdots & \dfrac{\partial w_n}{\partial v_2} \\[2ex] \vdots & \vdots & \ddots & \vdots \\[2ex] \dfrac{\partial w_1}{\partial v_m} & \dfrac{\partial w_2}{\partial v_m} & \cdots & \dfrac{\partial w_n}{\partial v_m} \end{bmatrix}}_{\text{Column } j \text{ matches } w_j} \left.\vphantom{\begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix}}\right\} \text{Row } i \text{ matches } v_i$$

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.

〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜

## X.13  More about the vector/vector derivative

Let's show a specific example: $w$ is $(3 \times 1)$, $v$ is $(2 \times 1)$.

$$
\frac{\partial w}{\partial v} =
\begin{bmatrix}
\overbrace{\dfrac{\partial w_1}{\partial v_1}}^{w_1} & \overbrace{\dfrac{\partial w_2}{\partial v_1}}^{w_2} & \overbrace{\dfrac{\partial w_3}{\partial v_1}}^{w_3} \\[4mm]
\dfrac{\partial w_1}{\partial v_2} & \dfrac{\partial w_2}{\partial v_2} & \dfrac{\partial w_3}{\partial v_2}
\end{bmatrix}
\begin{array}{l} \left.\rule{0pt}{6mm}\right\} v_1 \\[4mm] \left.\rule{0pt}{6mm}\right\} v_2 \end{array}
\tag{42}
$$

Another way to describe the general case:

> **Notation 14**
>
> Our matrix $\partial w/\partial v$ is entirely filled with **scalar derivatives**
>
> $$\frac{\partial w_j}{\partial v_i} \tag{43}$$
>
> 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 $\partial w/\partial v$ can be written as

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

or

$$
\frac{\partial w}{\partial v} = \left. \begin{bmatrix} \dfrac{\partial w}{\partial v_1} \\[2mm] \dfrac{\partial w}{\partial v_2} \\[2mm] \vdots \\[2mm] \dfrac{\partial w}{\partial v_m} \end{bmatrix} \right\} \text{Row i matches } v_i
$$

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

≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈≈

## X.14    Derivative: matrix/scalar

Now, we have our general form for creating derivatives.

We'll get our derivative of the form

$$
\frac{\partial(\text{Matrix})}{\partial(\text{Scalar})} = \frac{\partial M}{\partial s} \tag{44}
$$

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

$$
M = \begin{bmatrix} m_{11} & m_{12} & \cdots & m_{1r} \\ m_{21} & m_{22} & \cdots & m_{2r} \\ \vdots & \vdots & \ddots & \vdots \\ m_{k1} & m_{k2} & \cdots & m_{kr} \end{bmatrix} \tag{45}
$$

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/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**.

$$M = \begin{bmatrix} M_1 \\ M_2 \\ \vdots \\ M_k \end{bmatrix} \tag{46}$$

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

$$\frac{\partial M}{\partial s} = \begin{bmatrix} \dfrac{\partial M_1}{\partial s} & \dfrac{\partial M_2}{\partial s} & \cdots & \dfrac{\partial M_r}{\partial s} \end{bmatrix} \tag{47}$$

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

$$\frac{\partial M}{\partial s} = \overbrace{\begin{bmatrix} \begin{bmatrix} \dfrac{\partial m_{11}}{\partial s} \\ \dfrac{\partial m_{12}}{\partial s} \\ \vdots \\ \dfrac{\partial m_{1r}}{\partial s} \end{bmatrix}, & \begin{bmatrix} \dfrac{\partial m_{21}}{\partial s} \\ \dfrac{\partial m_{22}}{\partial s} \\ \vdots \\ \dfrac{\partial m_{2r}}{\partial s} \end{bmatrix}, & \cdots & \begin{bmatrix} \dfrac{\partial m_{k1}}{\partial s} \\ \dfrac{\partial m_{k2}}{\partial s} \\ \vdots \\ \dfrac{\partial m_{kr}}{\partial s} \end{bmatrix} \end{bmatrix}}^{\text{Column j matches } m_{j?}} \left.\vphantom{\begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix}}\right\} \text{Row i matches } m_{?i} \tag{48}$$

**Definition 17**

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

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

$$
\frac{\partial M}{\partial s} =
\overbrace{
\begin{bmatrix}
\dfrac{\partial m_{11}}{\partial s} & \dfrac{\partial m_{21}}{\partial s} & \cdots & \dfrac{\partial m_{k1}}{\partial s} \\[2ex]
\dfrac{\partial m_{12}}{\partial s} & \dfrac{\partial m_{22}}{\partial s} & \cdots & \dfrac{\partial m_{k2}}{\partial s} \\[2ex]
\vdots & \vdots & \ddots & \vdots \\[2ex]
\dfrac{\partial m_{1r}}{\partial s} & \dfrac{\partial m_{2r}}{\partial s} & \cdots & \dfrac{\partial m_{kr}}{\partial s}
\end{bmatrix}
}^{\text{Column j matches } m_{j?}}
\left.\vphantom{\begin{bmatrix}1\\1\\1\\1\end{bmatrix}}\right\} \text{Row i matches } m_{?i}
$$

This matrix has the transpose of the shape of $M$.

## X.15    Derivative: scalar/matrix

We'll get our derivative of the form

$$
\frac{\partial(\text{Scalar})}{\partial(\text{Matrix})} = \frac{\partial s}{\partial M} \tag{49}
$$

We have a matrix $M$ 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**.

$$
M = \begin{bmatrix} M_1 \\ M_2 \\ \vdots \\ M_k \end{bmatrix} \tag{50}
$$

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

$$\frac{\partial s}{\partial M} = \begin{bmatrix} \dfrac{\partial s}{\partial M_1} \\[2ex] \dfrac{\partial s}{\partial M_2} \\[2ex] \vdots \\[2ex] \dfrac{\partial s}{\partial M_k} \end{bmatrix} \tag{51}$$

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

$$\frac{\partial s}{\partial M} = \begin{bmatrix} \left[ \dfrac{\partial s}{\partial m_{11}} \quad \dfrac{\partial s}{\partial m_{12}} \quad \cdots \quad \dfrac{\partial s}{\partial m_{1r}} \right] \\[3ex] \left[ \dfrac{\partial s}{\partial m_{21}} \quad \dfrac{\partial s}{\partial m_{22}} \quad \cdots \quad \dfrac{\partial s}{\partial m_{2r}} \right] \\[3ex] \vdots \\[3ex] \left[ \dfrac{\partial s}{\partial m_{k1}} \quad \dfrac{\partial s}{\partial m_{k2}} \quad \cdots \quad \dfrac{\partial s}{\partial m_{kr}} \right] \end{bmatrix} \tag{52}$$

---

**Definition 18**

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

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

Column j matches $m_{?j}$

$$\frac{\partial s}{\partial M} = \begin{bmatrix} \dfrac{\partial s}{\partial m_{11}} & \dfrac{\partial s}{\partial m_{12}} & \cdots & \dfrac{\partial s}{\partial m_{1r}} \\[3ex] \dfrac{\partial s}{\partial m_{21}} & \dfrac{\partial s}{\partial m_{22}} & \cdots & \dfrac{\partial s}{\partial m_{2r}} \\[3ex] \vdots & \vdots & \ddots & \vdots \\[3ex] \dfrac{\partial s}{\partial m_{k1}} & \dfrac{\partial s}{\partial m_{k2}} & \cdots & \dfrac{\partial s}{\partial m_{kr}} \end{bmatrix} \quad \text{Row i matches } m_{i?}$$

This matrix has the same shape as $M$.

---

## X.16    Other Derivatives

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

$$\frac{\partial v}{\partial M}? \qquad \frac{\partial M}{\partial v}? \qquad \frac{\partial M}{\partial M^2}? \qquad\qquad (53)$$

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): $v_1, v_2, v_3$...

$$\left.\begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_m \end{bmatrix}\right\} \text{Axis 1}$$

  〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜〜

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

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

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? Jump to X.17 for a clarification.

$$
\overbrace{\phantom{\frac{\partial w_1}{\partial v_1}}}^{\text{Column j: vertical axis of } w}
$$

$$
\frac{\partial w}{\partial v} = \left.\begin{bmatrix} \dfrac{\partial w_1}{\partial v_1} & \dfrac{\partial w_2}{\partial v_1} & \cdots & \dfrac{\partial w_n}{\partial v_1} \\[2ex] \dfrac{\partial w_1}{\partial v_2} & \dfrac{\partial w_2}{\partial v_2} & \cdots & \dfrac{\partial w_n}{\partial v_2} \\[2ex] \vdots & \vdots & \ddots & \vdots \\[2ex] \dfrac{\partial w_1}{\partial v_m} & \dfrac{\partial w_2}{\partial v_m} & \cdots & \dfrac{\partial w_n}{\partial v_m} \end{bmatrix}\right\} \text{Row i: vertical axis of } v
$$

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:

- $\partial s / \partial v$: we need one axis to represent each term $v_i$.

    - 0 axis + 1 axis $\rightarrow$ 1 axis: the output is a (column) **vector**.

- $\partial v / \partial s$: we need one axis to represent each term $w_j$.

    - 1 axis + 0 axis $\rightarrow$ 1 axis: the output is a (row) **vector**.

- $\partial w / \partial v$: 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**.

- $\partial M / \partial s$: we need one axis to represent the rows of $M$, and another to represent the columns of $M$.

    - 2 axis + 0 axis $\rightarrow$ 2 axes: the output is a **matrix**.

- $\partial s / \partial M$: we need one axis to represent the rows of $M$, and another to represent the columns of $M$.

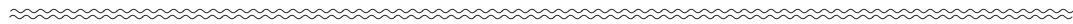    - 0 axis + 2 axis $\rightarrow$ 2 axes: the output is a **matrix**.

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:
>
> $$\text{Axes}\left(\frac{\partial y}{\partial x}\right) = \text{Axes}(y) + \text{Axes}(x) \tag{54}$$

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.

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

## X.17    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**.

$$\left.\begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_k \end{bmatrix}\right\} \text{We can adjust each of our } k \text{ scalars.}$$

- 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.

If you jumped here from X.16, feel free to follow this link back. Otherwise, continue on.

∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿∿

## X.18   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**.

Note that what we call a tensor is **not** a mathematical (or physics) tensor: we do not often use the "tensor product", or other tensor properties.

Our tensor can be better thought of as a "**generalized matrix**".

---

**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.

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

> These examples aren't especially important, but you'll see different variations in different softwares!

- If a tensor has a pattern of zeroes, we might be able to flatten it into a matrix.

    - For example, if we wanted to flatten a matrix into a vector (which we sometimes do!), we could do

$$
\begin{bmatrix} 3 & 0 & 0 \\ 0 & 9 & 0 \\ 0 & 0 & 4 \end{bmatrix} \rightarrow \begin{bmatrix} 3 \\ 9 \\ 4 \end{bmatrix} \tag{55}
$$

- We can also flatten it into a matrix or vector by placing the layers next to each other.

- We cleverly do regular matrix multiplication in a way that's compatible with our tensors.

    - Note that tensors do not have a matrix multiplication-like multiplication by default: several have been designed, however.

- We ignore the structure of the tensor, and just look at the individual elements: we take the scalar chain rule for each of them, without respecting the overall tensor.

> **Clarification 24**
>
> 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.
>
> These techniques are part of why: there are **different** approaches for how to approach these objects.

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

However, because they don't have a natural matrix multiplication, we'll try to convert it into a matrix in most cases.

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

### X.19    The loss derivative

Finally, we apply this to our common derivatives in section 7.5.

$$\overbrace{\frac{\partial \mathcal{L}}{\partial A^L}}^{(n^L \times 1)} \tag{56}$$

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 $A^L$.

> **Notation 25**
>
> Our derivative
>
> $$\frac{\partial \mathcal{L}}{\partial A^L} \tag{57}$$
>
> Is a scalar/vector derivative, and thus the shape $(n^L \times 1)$.

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

### X.20    The weight derivative

$$\overbrace{\frac{\partial Z^\ell}{\partial W^\ell}}^{(m^\ell \times 1)?} \tag{58}$$

This derivative is difficult - it's a derivative in the form vector/matrix. With **three** axes, we might imagine representing as a 3-tensor.

In fact, this can be manipulated into multiple different interesting **shapes** based on your **interpretation**: as we mentioned, there's no consistent rule for these variables.

But, our goal is to use this for the **chain rule**: so, we need to make the shapes **match**. This is why we do that strange transposing for our complete derivative.

$$\frac{\partial \mathcal{L}}{\partial W^\ell} = \overbrace{\frac{\partial Z^\ell}{\partial W^\ell}}^{\text{Weight link}} \cdot \overbrace{\left(\frac{\partial \mathcal{L}}{\partial Z^\ell}\right)^{\mathsf{T}}}^{\text{Other layers}} \tag{59}$$

Our problem is we have **too many axes**: the easiest way to resolve this to **break up** our matrix. So, for now, we focus on only **one neuron** at a time: it has a column vector $W_i$.

$$W = \begin{bmatrix} W_1 & W_2 & \cdots & W_n \end{bmatrix} \tag{60}$$

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

Notice that, this time, we broke it into **column vectors**, rather than row vectors: each neuron's **weights** are represented by a column vector.

We'll ignore everything except $W_i$.

$$W_i = \begin{bmatrix} w_{1i} \\ w_{2i} \\ \vdots \\ w_{mi} \end{bmatrix} \tag{61}$$

Finally, we get into our equation: notice that a **single** neuron has only **one** pre-activation $z_i$, so we don't need the whole vector.

$$z_i = W_i^{\mathsf{T}} A \tag{62}$$

Wait: there's something to notice, right off the bat. $z_i$ is **only** a function of $W_i$: that means the derivative for every other term $\partial/\partial W_k$ is **zero**!

> For example, changing $W_2$ would have **no** effect on $z_1$.

---

**Concept 26**

The $i^{\text{th}}$ neuron's **weights**, $W_i$, have **no effect** on a different neuron's **pre-activation** $z_j$.
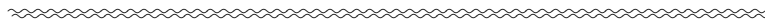
So, if the **neurons** don't match, then our derivative is zero:

- $i$ is the neuron for pre-activation $z_i$

- $j$ is the $j^{\text{th}}$ **weight** in a neuron.

- $k$ is the neuron for weight vector $W_k$

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

So, our only nonzero derivatives are

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

---

With that done, let's substitute in our values:

$$z_i = \begin{bmatrix} w_{1i} & w_{2i} & \cdots & w_{mi} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix} \tag{63}$$

And we'll do our **matrix multiplication**:

$$z_i = \sum_{j=1}^{n} W_{ji} a_j \tag{64}$$

Finally, we can get our derivatives:
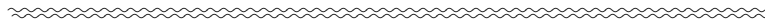
$$\frac{\partial z_i}{\partial W_{ji}} = a_j \tag{65}$$

So, if we combine that into a vector, we get:

$$\frac{\partial z_i}{\partial W_i} = \begin{bmatrix} \frac{\partial z_i}{\partial W_{1i}} \\[2ex] \frac{\partial z_i}{\partial W_{2i}} \\[2ex] \vdots \\[2ex] \frac{\partial z_i}{\partial W_{mi}} \end{bmatrix} \tag{66}$$

We can use our equation:

$$\frac{\partial z_i}{\partial W_i} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix} = A \tag{67}$$

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 don't affect different neurons.

$$\frac{\partial z_i}{\partial W_{jk}} = 0 \qquad \text{if } i \neq k \tag{68}$$

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

$$\frac{\partial z_i}{\partial W_k} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \vec{0} \qquad \text{if } i \neq k \tag{69}$$

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

$$\frac{\partial z_i}{\partial W_k} = \begin{cases} A & \text{if } i = k \\ \vec{0} & \text{if } i \neq k \end{cases} \tag{70}$$

These are all the elements of our matrix $\partial z_i / \partial W_k$: so, we can get our result.

$$\frac{\partial Z}{\partial W} = \begin{bmatrix} A & \vec{0} & \cdots & \vec{0} \\ \vec{0} & A & \cdots & \vec{0} \\ \vdots & \vdots & \ddots & \vec{0} \\ \vec{0} & \vec{0} & \vec{0} & A \end{bmatrix} \tag{71}$$

We have our result: it turns out, despite being stored in a **matrix**-like format, this is actually a **3-tensor**! Each entry of our **matrix** is a **vector**: 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 27**

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.

$$\begin{bmatrix} A & \vec{0} & \cdots & \vec{0} \\ \vec{0} & A & \cdots & \vec{0} \\ \vdots & \vdots & \ddots & \vec{0} \\ \vec{0} & \vec{0} & \vec{0} & A \end{bmatrix} \longrightarrow \begin{bmatrix} A \\ A \\ \vdots \\ A \end{bmatrix} \tag{72}$$

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

$$\begin{bmatrix} A \\ A \\ \vdots \\ A \end{bmatrix} \longrightarrow A \tag{73}$$

We have our big result!

> **Notation 28**
>
> Our derivative
>
> $$\frac{\overbrace{\partial Z^\ell}^{(m^\ell \times 1)}}{\partial W^\ell} = A^{\ell-1}$$
>
> 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 derivative order, etc.)

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

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

## X.21    Linking Layers

$$\frac{\partial Z^\ell}{\partial A^{\ell-1}} \tag{74}$$

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 $\ell$, as before.

$$Z = W^\mathsf{T} A \tag{75}$$

We'll use the same approach we did last time: $W$ is a vector, and we'll focus on $W_i$. This will allow us to break it up **element-wise**, and get all of our **derivatives**.

We could treat $W$ as a whole matrix, but this will give us our results without as much clutter: the only **difference** is that we would have to depict every $W_i$ at **once**.

$$W = \begin{bmatrix} W_1 & W_2 & \cdots & W_n \end{bmatrix} \qquad W_i = \begin{bmatrix} w_{1i} \\ w_{2i} \\ \vdots \\ w_{mi} \end{bmatrix} \tag{76}$$

Here's our equation:

$$z_i = \begin{bmatrix} w_{1i} & w_{2i} & \cdots & w_{mi} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix} \tag{77}$$

We matrix multiply:

$$z_i = \sum_{j=1}^{n} W_{ji} a_j \tag{78}$$

The derivative can be gotten from here -

$$\frac{\partial z_i}{\partial a_j} = W_{ji} \tag{79}$$

We look at our whole matrix derivative:

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

$$\frac{\partial Z}{\partial A} = \overbrace{\begin{bmatrix} \ddots & \vdots & \cdot^{\cdot^{\cdot}} \\ \cdots & \dfrac{\partial z_i}{\partial a_j} & \cdots \\ \cdot^{\cdot^{\cdot}} & \vdots & \ddots \end{bmatrix}}^{\text{Column j matches } z_i} \left.\vphantom{\begin{bmatrix} \ddots \\ \cdots \\ \cdot \end{bmatrix}}\right\} \text{Row i matches } a_j \tag{80}$$

Wait.

- The derivative $\partial z_i / \partial a_j$ is in the $j^{\text{th}}$ row, $i^{\text{th}}$ column.

- $W_{ji}$ represents the element in the $j^{\text{th}}$ row, $i^{\text{th}}$ column.

They're the same matrix!

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

We get our final result:

---

**Notation 29**

Our derivative

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

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

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

---

## X.22    Activation Function

$$\frac{\partial A^{\ell}}{\partial Z^{\ell}} \tag{81}$$

The last derivative is less unusual than it looks.

$$A^{\ell} = f(Z^{\ell}) \longrightarrow \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = f\left( \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{bmatrix} \right) \tag{82}$$

We can apply our function element-wise:

$$\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = \begin{bmatrix} f(z_1) \\ f(z_2) \\ \vdots \\ f(z_n) \end{bmatrix} \tag{83}$$

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

---

**Concept 30**

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

$$\frac{\partial a_j}{\partial z_i}$$

---

As for our remaining term, we'll describe any row of the above vectors:

$$a_i = f(z_i) \tag{84}$$

Our derivative is:

$$\frac{\partial a_i}{\partial z_i} = f'(z_i) \tag{85}$$

In general, including the non-diagonals:

$$\frac{\partial a_i}{\partial z_i} = \begin{cases} f'(z_i) & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases} \tag{86}$$

This gives us our result:

---

**Notation 31**

Our derivative

$$\overset{(n^\ell \times n^\ell)}{\frac{\partial A^\ell}{\partial Z^\ell}} = \overset{\text{Column } j \text{ matches } a_j}{\begin{bmatrix} f'(z_1^\ell) & 0 & 0 & \cdots & 0 \\ 0 & f'(z_2^\ell) & 0 & \cdots & 0 \\ 0 & 0 & f'(z_3^\ell) & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & 0 \\ 0 & 0 & 0 & 0 & f'(z_n^\ell) \end{bmatrix}} \left.\vphantom{\begin{bmatrix} f'(z_1^\ell) \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}}\right\} \text{Row } i \text{ matches } z_i \tag{87}$$

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

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

---

## X.23    Element-wise multiplication

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

$$\begin{bmatrix} f'(z_1^\ell) \\ f'(z_2^\ell) \\ \vdots \\ f'(z_n^\ell) \end{bmatrix} \tag{88}$$

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 32**

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

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

Then we can write their product each of these ways:

$$
RQ = \overbrace{R * Q}^{\text{Element-wise multiplication}} = \begin{bmatrix} r_1 q_1 \\ r_2 q_2 \\ r_3 q_3 \\ \vdots \\ r_n q_n \end{bmatrix} \tag{89}
$$

So, we can substitute the chain rule this way.