

Explanatory Notes for 6.390

Shauntclair Ruiz (Current TA)

Fall 2022

7.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:

$$\frac{\partial(\text{Scalar})}{\partial(\text{Vector})} = \frac{\partial \mathbf{s}}{\partial \mathbf{v}} \quad (1)$$

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**.

$$\Delta \mathbf{v} \longrightarrow \boxed{f} \longrightarrow \Delta \mathbf{s} \quad (2)$$

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

$$\overbrace{\begin{bmatrix} \Delta \mathbf{v}_1 \\ \Delta \mathbf{v}_2 \\ \vdots \\ \Delta \mathbf{v}_m \end{bmatrix}}^{\Delta \mathbf{v}} \longrightarrow \Delta \mathbf{s} \quad (3)$$

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 .

7.X.5 Finding the scalar/vector derivative

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

$$\Delta \mathbf{s} \approx \frac{\partial \mathbf{s}}{\partial \mathbf{v}} \star \Delta \mathbf{v} \quad \text{or} \quad \Delta \mathbf{s} \approx \frac{\partial \mathbf{s}}{\partial \mathbf{v}} \star \overbrace{\begin{bmatrix} \Delta \mathbf{v}_1 \\ \Delta \mathbf{v}_2 \\ \vdots \\ \Delta \mathbf{v}_m \end{bmatrix}}^{\Delta \mathbf{v}} \quad (4)$$

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**.

One problem, though: how can we treat each **derivative** separately? Each $\Delta \mathbf{v}_i$ will move our position, which can change a different derivative \mathbf{v}_k : they can **affect** each other.

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

7.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.

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 Δv_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 1

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 than δ .

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 2

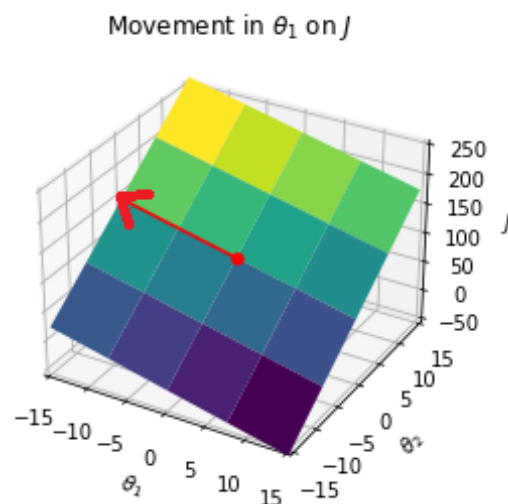
Why is this **true**? Because a **hyperplane** can be represented using our **linear** function

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

If we take a derivative:

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

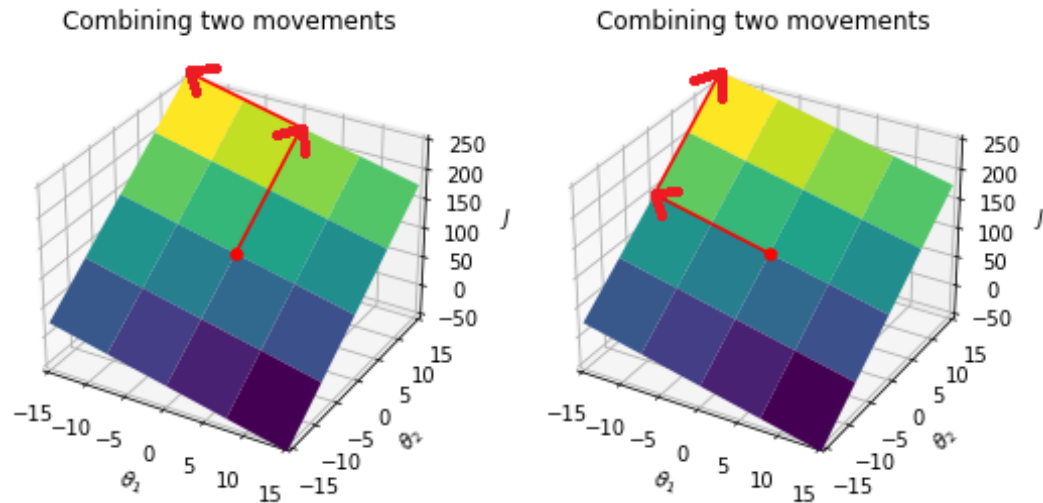
That derivative is a **constant**! It's doesn't change based on **position**.



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

Why does this help? If our derivative doesn't **change**, we can combine multiple steps. 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:

$$\Delta s = \Delta s_{\text{from } v_1} + \Delta s_{\text{from } v_2} + \cdots + \Delta s_{\text{from } v_m} \quad (5)$$

7.X.7 Our scalar/vector derivative

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

Definition 3

The **multivariable chain rule approximation** looks similar to the multivariable chain rule, but for finite changes Δx_i .

In 3-D, we get

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

In general, we have

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

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

This gives us what we're looking for:

$$\Delta \mathbf{s} \approx \sum_{i=1}^m \frac{\partial \mathbf{s}}{\partial \mathbf{v}_i} \Delta \mathbf{v}_i \quad (6)$$

If we circle back around to our original approximation:

$$\sum_{i=1}^m \frac{\partial \mathbf{s}}{\partial \mathbf{v}_i} \Delta \mathbf{v}_i = \frac{\partial \mathbf{s}}{\partial \mathbf{v}} \star \overbrace{\begin{bmatrix} \Delta \mathbf{v}_1 \\ \Delta \mathbf{v}_2 \\ \vdots \\ \Delta \mathbf{v}_m \end{bmatrix}}^{\Delta \mathbf{v}} \quad (7)$$

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

$$\sum_{i=1}^m \frac{\partial \mathbf{s}}{\partial \mathbf{v}_i} \Delta \mathbf{v}_i = \overbrace{\begin{bmatrix} \partial \mathbf{s} / \partial \mathbf{v}_1 \\ \partial \mathbf{s} / \partial \mathbf{v}_2 \\ \vdots \\ \partial \mathbf{s} / \partial \mathbf{v}_m \end{bmatrix}}^{\partial \mathbf{s} / \partial \mathbf{v}} \cdot \overbrace{\begin{bmatrix} \Delta \mathbf{v}_1 \\ \Delta \mathbf{v}_2 \\ \vdots \\ \Delta \mathbf{v}_m \end{bmatrix}}^{\Delta \mathbf{v}} \quad (8)$$

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

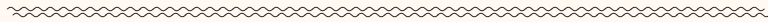
Definition 4

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:

$$\overbrace{\Delta s}^{(1 \times 1)} = \overbrace{\left(\frac{\partial s}{\partial \mathbf{v}} \right)^T}^{(1 \times m)} \overbrace{\Delta \mathbf{v}}^{(m \times 1)} \quad (9)$$