

§1 fundamentals

$$x \in \mathbb{R}^n \xrightarrow{f} y = f(x) \in \mathbb{R}^m$$

commonly, we write

$$f: \underbrace{\mathbb{R}^n}_{\text{domain / input space}} \rightarrow \underbrace{\mathbb{R}^m}_{\text{codomain / output or target space}}$$

def for some function $y = f(x)$, $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ let

$$\frac{\partial y}{\partial x} := \left[\frac{\partial y_i}{\partial x_j} \right]_{ij}$$

matrix whose (i,j) entry is $\partial y_i / \partial x_j$: the partial derivative of y_i w.r.t. x_j

sometimes, we also write $\frac{\partial f}{\partial x}$

$$= \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

notice if we think of this matrix as a linear function it has the same "shape" as f : takes input in \mathbb{R}^n and gives output in \mathbb{R}^m

def for $y = f(x)$, we say the jacobian of f at some point x_0 is the derivative above evaluated at $x = x_0$

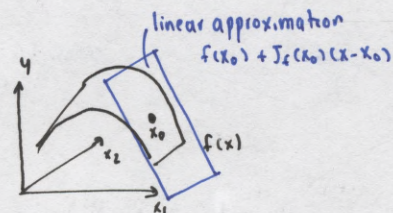
$$J_f(x_0) := \left. \frac{\partial y}{\partial x} \right|_{x=x_0}$$

the derivatives change at different points, this notation just allows us specify which point we are looking at

prop. let $y = f(x)$, $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$. for some point $x_0 \in \mathbb{R}^n$

$$f(x_0 + \Delta x) \approx f(x_0) + J_f(x_0) \Delta x$$

or $f(x) \approx f(x_0) + J_f(x_0)(x - x_0)$



intuition.

$$f(x_0 + \Delta x) - f(x_0) \approx J_f(x_0) \Delta x$$

how much does y_i change when I move x a bit?

$$= \begin{bmatrix} \frac{\partial y_i}{\partial x_1} & \frac{\partial y_i}{\partial x_2} & \dots & \frac{\partial y_i}{\partial x_n} \end{bmatrix} \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \vdots \\ \Delta x_n \end{bmatrix}$$

add up all changes

$$\Delta y_i = \underbrace{\frac{\partial y_i}{\partial x_1} \Delta x_1}_{\text{how much } y_i \text{ changes based off } x_1} + \underbrace{\frac{\partial y_i}{\partial x_2} \Delta x_2}_{\text{how much } y_i \text{ changes based off } x_2} + \dots + \frac{\partial y_i}{\partial x_n} \Delta x_n$$

example consider $f: \mathbb{R}^n \rightarrow \mathbb{R}$ (f is a scalar function)
vs a vector map

$$y = f(x) = \sum_{i=1}^n x_i^2 = \|x\|_2^2$$

$$J_f(x) = \left[\frac{\partial y}{\partial x_1} \quad \frac{\partial y}{\partial x_2} \quad \dots \quad \frac{\partial y}{\partial x_n} \right] \quad \text{important: this is a row vector } \mathbb{R}^{1 \times n}$$

$$= \begin{bmatrix} 2x_1 & 2x_2 & \dots & 2x_n \end{bmatrix}$$

$$= 2x^T$$

since $\frac{\partial y}{\partial x_k} = \frac{\partial}{\partial k} \sum_{i=1}^n x_i^2$

$$= \frac{\partial}{\partial k} x_k^2$$

$$= 2x_k$$

only 1 term actually depends on x_k

Example consider $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$

$y = f(x) = -\log(x)$ common abuse of notation: write regular functions, to mean componentwise application

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \xrightarrow{f} \begin{bmatrix} -\log x_1 \\ -\log x_2 \\ \vdots \\ -\log x_n \end{bmatrix}$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} -1/x_1 & 0 & \dots & 0 \\ 0 & -1/x_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & -1/x_n \end{bmatrix}$$

notice only diagonal entries are non-zero!

$$\frac{\partial}{\partial x_i} \underbrace{-\log x_i}_y = \begin{cases} -1/x_i & : i=j \\ 0 & : i \neq j \end{cases}$$

(if $i \neq j$, it's like we are differentiating a constant)

def. for a scalar function $f: \mathbb{R}^n \rightarrow \mathbb{R}$, we write the gradient as

$$\nabla_x f(x) = \left(\frac{\partial f}{\partial x} \right)^T$$

$\in \mathbb{R}^n$ ~ always has same dimension/shape as input

recall this was a row vector — transposing makes it a column

$$\underset{\mathbb{R}^n}{x^{(t+1)}} = \underset{\mathbb{R}^n}{x^{(t)}} - \alpha \underset{\mathbb{R}^n}{\nabla_x f(x^{(t)})}$$

(gradient update rule)

→ gradients are convenient because they share the shape of the input

example $f(x) = \|x\|_2^2$

$$\nabla_x f(x) = \left(\frac{\partial f}{\partial x} \right)^T = (2x^T)^T = 2x \quad \checkmark$$

see above

82 derivative rules

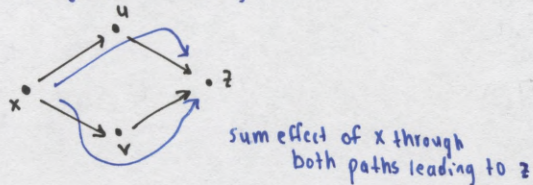
prop (chain rule, multivariate).

$$\begin{aligned} \text{let } u &= f(x) \\ v &= g(x) \\ z &= h(u, v) \\ &= h(f(x), g(x)) \end{aligned}$$

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial u} \frac{\partial u}{\partial x} + \frac{\partial z}{\partial v} \frac{\partial v}{\partial x}$$

how x affects z through u how x affects z through v

"Math 53 version"



prop (chain rule, multivariate) *

$$x \in \mathbb{R}^n \xrightarrow{f} y \in \mathbb{R}^m \xrightarrow{g} z \in \mathbb{R}^k$$

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

really easy to remember!
exactly the same as the regular chain rule...

remember, these are now jacobians

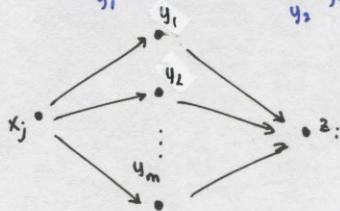
order matters now, this is a matrix multiplication

$$\begin{bmatrix} \frac{\partial z}{\partial x} \end{bmatrix}_{1 \times n} = \begin{bmatrix} \frac{\partial z}{\partial y_1} & \frac{\partial z}{\partial y_2} & \dots & \frac{\partial z}{\partial y_m} \end{bmatrix}_{1 \times m} \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_m}{\partial x} \end{bmatrix}_{m \times n}$$

so this tells us:

$$\frac{\partial z_i}{\partial x_j} = \frac{\partial z_i}{\partial y_1} \frac{\partial y_1}{\partial x_j} + \frac{\partial z_i}{\partial y_2} \frac{\partial y_2}{\partial x_j} + \dots + \frac{\partial z_i}{\partial y_m} \frac{\partial y_m}{\partial x_j}$$

effect through y_1 effect through y_2 ...



exact same idea:
sum effect through all paths to z_i

transposing both sides, we get:

$$\left(\frac{\partial z}{\partial x} \right)^T = \left(\frac{\partial y}{\partial x} \right)^T \left(\frac{\partial z}{\partial y} \right)^T$$

remember $(AB)^T = B^T A^T$

useful to write out some rules for gradients!

Corollary. $x \in \mathbb{R}^n \xrightarrow{f} y \in \mathbb{R} \xrightarrow{g} z \in \mathbb{R}$

$$f: \mathbb{R}^n \rightarrow \mathbb{R} \quad g \circ f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\begin{aligned} \nabla_x z &= \left(\frac{\partial z}{\partial x} \right)^T \\ &= \left(\frac{\partial y}{\partial x} \right)^T \left(\frac{\partial z}{\partial y} \right) \quad \text{both scalars!} \\ &\quad \text{gradient!} \in \mathbb{R} \\ &= \frac{dz}{dy} \nabla_x y \end{aligned}$$

$$\boxed{\nabla_x z = \frac{dz}{dy} \nabla_x y}$$

Example. compute the gradient of $\|x\|_2^4$ (really messy to do component-wise...)

$$f(x) = \|x\|_2^2$$

$$g(y) = y^2$$

$$\|x\|_2^4 = g(f(x))$$

$$\begin{aligned} \therefore \nabla_x \|x\|_2^4 &= \frac{dy^2}{dy} \cdot \nabla_x \|x\|_2^2 \quad \text{here } y = f(x) = \|x\|_2^2 \\ &= 2y \cdot 2x \\ &= 4\|x\|_2^2 x \quad \checkmark \end{aligned}$$

Corollary $x \in \mathbb{R}^n \xrightarrow{f} y \in \mathbb{R}^m \xrightarrow{g} z \in \mathbb{R}$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m \quad g \circ f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\begin{aligned} \nabla_x z &= \left(\frac{\partial z}{\partial x} \right)^T \\ &= \left(\frac{\partial y}{\partial x} \right)^T \left(\frac{\partial z}{\partial y} \right)^T \\ &= \left(\frac{\partial y}{\partial x} \right)^T \nabla_y z \end{aligned}$$

example recall the ℓ_2^2 loss for linear regression, given weights w is defined as:

$$L(w) = \|Xw - y\|_2^2$$

now let $f(w) = Xw - y$
 $g(v) = \|v\|_2^2$

$$\begin{aligned} \therefore \nabla_w L(w) &= \left(\frac{\partial f}{\partial w} \right)^T \nabla_v g(v) \\ &= \boxed{X^T} \cdot 2v \\ &= 2X^T(Xw - b) \end{aligned}$$

we can show $\frac{\partial}{\partial w}(Xw - y) = X$
 but this should make sense, since a linear approximation of a linear function should be linear

term pops out when you compose w/ a linear function

here $v = f(w) = Xw - b$

derivative rules

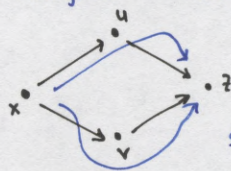
prop (chain rule, multivariate).

$$\begin{aligned} \text{let } u &= f(x) \\ v &= g(x) \\ z &= h(u, v) \\ &= h(f(x), g(x)) \end{aligned}$$

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial u} \frac{\partial u}{\partial x} + \frac{\partial z}{\partial v} \frac{\partial v}{\partial x}$$

how x affects z through u
how x affects z through v

"Math 53 version"



sum effect of x through both paths leading to z

prop (chain rule, multivariate) *

$$x \in \mathbb{R}^n \xrightarrow{f} y \in \mathbb{R}^m \xrightarrow{g} z \in \mathbb{R}^k$$

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

remember, these are now jacobians

order matters now, this is a matrix multiplication

really easy to remember! exactly the same as the regular chain rule...

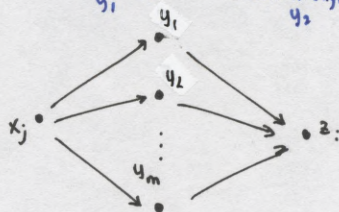
$$\begin{bmatrix} \frac{\partial z_i}{\partial x_1} \\ \vdots \\ \frac{\partial z_i}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial z_i}{\partial y_1} & \frac{\partial z_i}{\partial y_2} & \dots & \frac{\partial z_i}{\partial y_m} \end{bmatrix} \begin{bmatrix} \frac{\partial y_1}{\partial x_i} \\ \frac{\partial y_2}{\partial x_i} \\ \vdots \\ \frac{\partial y_m}{\partial x_i} \end{bmatrix}$$

$\mathbb{R}^{k \times n}$ $\mathbb{R}^{k \times m}$ $\mathbb{R}^{m \times n}$

so this tells us:

$$\frac{\partial z_i}{\partial x_j} = \frac{\partial z_i}{\partial y_1} \frac{\partial y_1}{\partial x_j} + \frac{\partial z_i}{\partial y_2} \frac{\partial y_2}{\partial x_j} + \dots + \frac{\partial z_i}{\partial y_m} \frac{\partial y_m}{\partial x_j}$$

effect through y_1
effect through y_2 ...



exact same idea: sum effect through all paths to z_i

transposing both sides, we get:

$$\left(\frac{\partial z}{\partial x} \right)^T = \left(\frac{\partial y}{\partial x} \right)^T \left(\frac{\partial z}{\partial y} \right)^T$$

remember $(AB)^T = B^T A^T$

useful to write out some rules for gradients!

prop (addition rule).

$$\frac{\partial}{\partial x} \sum_{i=1}^m f_i(x) = \sum_{i=1}^m \frac{\partial}{\partial x} f_i(x)$$

$$\nabla_x \sum_{i=1}^m f_i(x) = \sum_{i=1}^m \nabla_x f_i(x)$$

taking transposes

"the derivative is linear"

abuse of notation:

here $\{f_1, f_2, \dots, f_m\}$ are different functions,
not the coordinate functions of some map f

proof idea.

let $f: \mathbb{R}^n \rightarrow \mathbb{R}$

$$\nabla_x \sum_{i=1}^m f_i(x) = \begin{bmatrix} \frac{\partial}{\partial x_1} \sum_{i=1}^m f_i(x) \\ \vdots \\ \frac{\partial}{\partial x_n} \sum_{i=1}^m f_i(x) \end{bmatrix} = \sum_{i=1}^m \begin{bmatrix} \frac{\partial}{\partial x_1} f_i(x) \\ \vdots \\ \frac{\partial}{\partial x_n} f_i(x) \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial}{\partial x_1} f_1(x) + \dots + \frac{\partial}{\partial x_1} f_m(x) \\ \vdots \\ \frac{\partial}{\partial x_n} f_1(x) + \dots + \frac{\partial}{\partial x_n} f_m(x) \end{bmatrix}$$

$$= \sum_{i=1}^m \nabla_x f_i(x)$$

prop (scalar multiplication rule)

$$\frac{\partial}{\partial x} c f(x) = c \frac{\partial}{\partial x} f(x)$$

$$\nabla_x c f(x) = c \nabla_x f(x)$$

$$\left[\frac{\partial}{\partial x} c f(x) \right]_{ij} = \frac{\partial}{\partial x_i} c f_i(x)$$

refers to i th coordinate of the output $f(x)$

$$= c \cdot \frac{\partial}{\partial x_i} f_i(x)$$

$$= c \left[\frac{\partial}{\partial x} f(x) \right]_{ij}$$

prop (gradient dot product rule).

$$x \in \mathbb{R}^n \xrightarrow{f} u \in \mathbb{R}^m$$

$$x \in \mathbb{R}^n \xrightarrow{g} v \in \mathbb{R}^m$$

$$u^T v \in \mathbb{R}$$

$$u \in \mathbb{R}^m$$

$$v \in \mathbb{R}^m$$

$$\nabla_x u^T v = \nabla_x f(x)^T g(x)$$

$$= \underbrace{\left(\frac{\partial u}{\partial x} \right)^T}_{\mathbb{R}^{n \times m}} \underbrace{v}_{\mathbb{R}^m} + \underbrace{\left(\frac{\partial v}{\partial x} \right)^T}_{\mathbb{R}^{n \times m}} u$$

looks similar to regular product rule!

$$\frac{\partial}{\partial x_i} u^T v = \frac{\partial}{\partial x_i} \cdot \sum_{j=1}^m u_j v_j$$

apply scalar product rule 2
split into two summations

$$= \sum_{j=1}^m \frac{\partial u_j}{\partial x_i} v_j + \sum_{j=1}^m u_j \frac{\partial v_j}{\partial x_i}$$

looks a bit like matrix multiplication:

$$\begin{bmatrix} \vdots \\ \frac{\partial u}{\partial x} \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} v_1 \\ \vdots \\ v_m \end{bmatrix}$$

$\frac{\partial u}{\partial x} \in \mathbb{R}^{m \times n}$

$$\nabla_x u^T v = \begin{bmatrix} \frac{\partial u_1}{\partial x_1} & \dots & \frac{\partial u_m}{\partial x_1} \end{bmatrix} \begin{bmatrix} v_1 \\ \vdots \\ v_m \end{bmatrix} + \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} v_1 \\ \vdots \\ v_m \end{bmatrix}$$

\mathbb{R}^n $\begin{pmatrix} \frac{\partial u}{\partial x} \end{pmatrix}^T v$ $\begin{pmatrix} \frac{\partial v}{\partial x} \end{pmatrix}^T u$
 $\mathbb{R}^{n \times m} \cdot \mathbb{R}^{m \times 1}$ $\mathbb{R}^{n \times m} \mathbb{R}^{m \times 1}$

$$\therefore \nabla_x u^T v = \left(\frac{\partial u}{\partial x} \right)^T v + \left(\frac{\partial v}{\partial x} \right)^T u$$

example let $x \in \mathbb{R}^n$
 $A \in \mathbb{R}^{n \times n}$ be some matrix, and define

$$f(x) = \underbrace{x^T}_u \underbrace{Ax}_v$$

set $u = x$
 $v = Ax$

$$\begin{aligned} \therefore \nabla_x f(x) &= \nabla_x u^T v \\ &= \left(\frac{\partial u}{\partial x} \right)^T v + \left(\frac{\partial v}{\partial x} \right)^T u \\ &= (I)^T (Ax) + (A)^T x \\ &= (A + A^T) x \end{aligned}$$

$$\left[\frac{\partial x}{\partial x} \right]_{ij} = \frac{\partial x_i}{\partial x_j}$$

clearly 1 if $i=j$,
 0 otherwise

$\frac{\partial}{\partial x} Ax = A$, appealing once again to the notion of a derivative as a linear approximation

§3 common derivatives

(linear functions)

$$\nabla_x c^T x = c$$

$$\frac{\partial}{\partial x} c^T x = c^T$$

$$\nabla_x \overset{\text{constant}}{c} = 0$$

$$\frac{\partial}{\partial x} \overset{\text{constant}}{c} = 0$$

(quadratic functions)

$$\begin{aligned} \nabla_x x^T A x &= (A + A^T) x \\ &= 2Ax \quad (\text{when } A \text{ is symmetric}) \end{aligned}$$

$$\begin{aligned} \nabla_x x^T A x + b^T x + c &= (A + A^T) x + b \\ &= 2Ax + b \quad (\text{when } A \text{ is symmetric}) \end{aligned}$$

$$\begin{aligned} \nabla_x \|Ax + b\|_2^2 &= \nabla_x (Ax + b)^T (Ax + b) \\ &= \nabla_x [x^T \underbrace{A^T A}_{\text{symmetric}} x + 2b^T A x + b^T b] \\ &= 2A^T A x + A^T b \end{aligned}$$

$$\nabla_x \|Ax - b\|_2^2 = 2A^T A x - A^T b$$

$$\nabla_x \|x\|_2^2 = 2x$$

$$\nabla_x (Ax + b)^T C (Dx + e) = D^T C^T (Ax + b) + A^T C (Dx + e)$$

$$\begin{aligned} \nabla_x (Ax - b)^T W (Ax - b) &= A^T W^T (Ax - b) + A^T W (Ax - b) \\ &= 2A^T W (Ax - b) \quad \text{when } W \text{ is symmetric} \end{aligned}$$

$$\nabla_x a^T x x^T b = (ab^T + ba^T) x$$

(other functions)

$$\nabla_x \|x - a\|_2 = \frac{x - a}{\|x - a\|_2}$$