

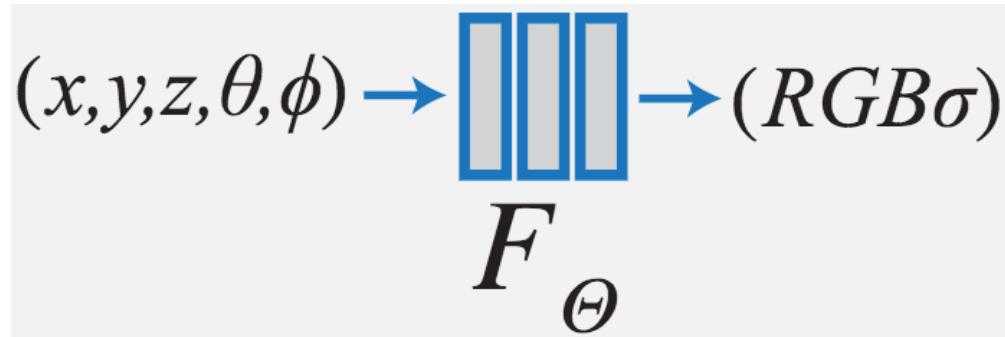
Vector Calculus

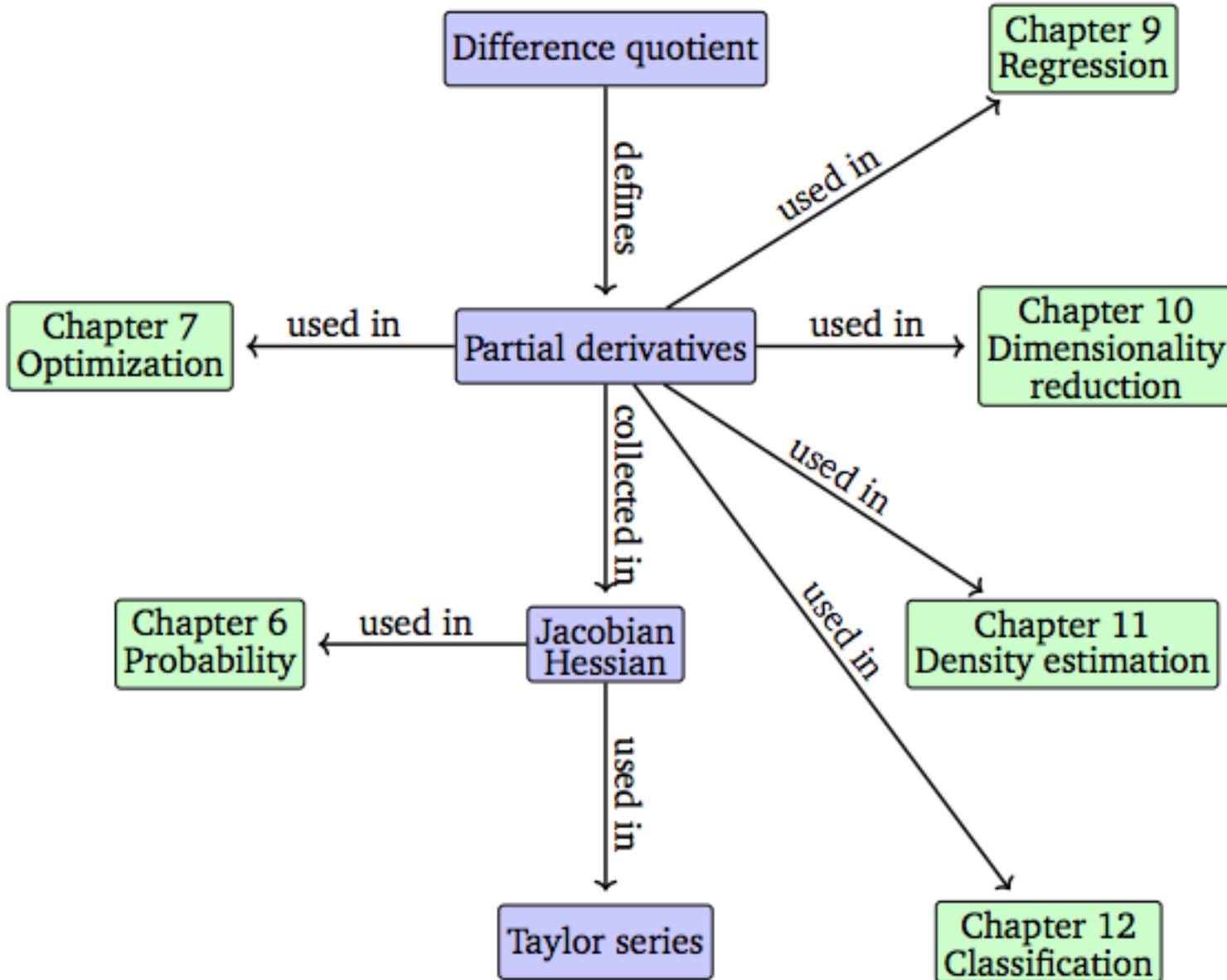
Liang Zheng

Australian National University

liang.zheng@anu.edu.au

NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. Mildenhall et al., ECCV 2020





5 Vector Calculus

- We discuss functions

$$\begin{aligned} f : \mathbb{R}^D &\rightarrow \mathbb{R} \\ x &\mapsto f(x) \end{aligned}$$

where \mathbb{R}^D is the **domain** of f , and the function values $f(x)$ are the **image/codomain** of f .

- Example (dot product)
- Previously, we write dot product as

$$f(x) = x^T x, \quad x \in \mathbb{R}^2$$

- In this chapter, we write it as

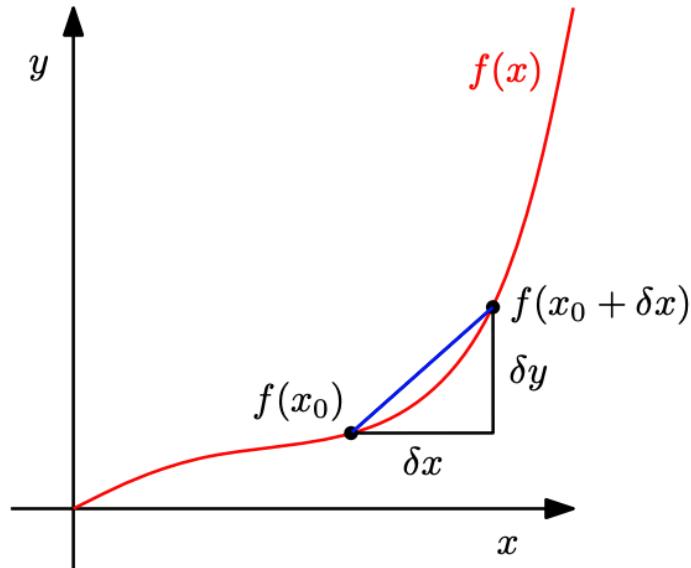
$$\begin{aligned} f : \mathbb{R}^2 &\rightarrow \mathbb{R} \\ x &\mapsto x_1^2 + x_2^2 \end{aligned}$$

5.1 Differentiation of Univariate Functions

- Given $y = f(x)$, the **difference quotient** is defined as

$$\frac{\delta y}{\delta x} := \frac{f(x + \delta x) - f(x)}{\delta x}$$

- It computes the slope of the secant line through two points on the graph of f . In this figure, these are the points with x -coordinates x_0 and $x_0 + \delta x_0$.
- In the limit for $\delta x \rightarrow 0$, we obtain the tangent of f at x (if f is differentiable). The tangent is then the derivative of f at x .



5.1 Differentiation of Univariate Functions

- For $h > 0$, the **derivative** of f at x is defined as the limit

$$\frac{df}{dx} := \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h}$$

- The derivative of f points in the direction of steepest ascent of f .

- Example - Derivative of a Polynomial

- Compute the derivative of $f(x) = x^n$, $n \in \mathbb{N}$. (From our high school knowledge, the derivative is nx^{n-1} .)

$$\begin{aligned}\frac{df}{dx} &:= \lim_{h \rightarrow 0} \frac{f(x + h) - f(x)}{h} \\ &= \lim_{h \rightarrow 0} \frac{(x + h)^n - x^n}{h} \\ &= \lim_{h \rightarrow 0} \frac{\sum_{i=0}^n \binom{n}{i} x^{n-i} h^i - x^n}{h}\end{aligned}$$

we see that $x^n = \binom{n}{0} x^{n-0} h^0$. By starting the sum at 1, the x^n cancels.

5.1 Differentiation of Univariate Functions

$$\begin{aligned}\frac{df}{dx} &= \lim_{h \rightarrow 0} \frac{\sum_{i=0}^n \binom{n}{i} x^{n-i} h^i - x^n}{h} \\&= \lim_{h \rightarrow 0} \frac{\sum_{i=1}^n \binom{n}{i} x^{n-i} h^i}{h} \\&= \lim_{h \rightarrow 0} \sum_{i=1}^n \binom{n}{i} x^{n-i} h^{i-1} \\&= \lim_{h \rightarrow 0} \left\{ \binom{n}{1} x^{n-1} + \underbrace{\sum_{i=2}^n \binom{n}{i} x^{n-i} h^{i-1}}_{\rightarrow 0 \text{ as } h \rightarrow 0} \right\} \\&= nx^{n-1}\end{aligned}$$

5.1.2 Differentiation Rules

- Product rule

$$(f(x)g(x))' = f'(x)g(x) + f(x)g'(x)$$

- Quotient rule:

$$\left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2}$$

- Sum rule:

$$(f(x) + g(x))' = f'(x) + g'(x)$$

- Chain rule:

$$(g(f(x)))' = (g \circ f)'(x) = g'(f(x))f'(x)$$

Here, $g \circ f$ denotes function composition $g(f(x))$

Example -- Chain rule

- Compute the derivative of the function $h(x) = (2x + 1)^4$
- We write

$$\begin{aligned} h(x) &= (2x + 1)^4 = g(f(x)) \\ f(x) &= 2x + 1 \\ g(f) &= f^4 \end{aligned}$$

- We obtain the derivatives of f and g as,

$$\begin{aligned} f'(x) &= 2 \\ g'(f) &= 4f^3 \end{aligned}$$

- The derivative of h is given as

$$h'(x) = g'(f) f'(x) = (4f^3) \cdot 2 = 4(2x + 1)^3 \cdot 2 = 8(2x + 1)^3$$

5.2 Partial Differentiation and Gradients

- Instead of considering $x \in \mathbb{R}$, we consider $\mathbf{x} \in \mathbb{R}^n$, e.g., $f(\mathbf{x}) = f(x_1, x_2)$
- The generalization of the derivative to functions of several variables is the **gradient**.
- We find the gradient of the function f with respect to \mathbf{x} by
 - varying one variable at a time and keeping the others constant.
 - The gradient is the collection of these **partial derivatives**.
- For a function $f: \mathbb{R}^n \rightarrow \mathbb{R}, \mathbf{x} \mapsto f(\mathbf{x}), \mathbf{x} \in \mathbb{R}^n$ of n variables x_1, \dots, x_n , we define the **partial derivatives** as

$$\frac{\partial f}{\partial x_1} := \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(x)}{h}$$
$$\vdots$$
$$\frac{\partial f}{\partial x_n} := \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_{n-1}, x_n + h) - f(x)}{h}$$

and collect them in the row vector

$$\nabla_{\mathbf{x}} f = \text{grad } f = \frac{df}{d\mathbf{x}} = \left[\frac{\partial f(\mathbf{x})}{\partial x_1} \quad \frac{\partial f(\mathbf{x})}{\partial x_n} \quad \dots \quad \frac{\partial f(\mathbf{x})}{\partial x_n} \right] \in \mathbb{R}^{1 \times n}$$

5.2 Partial Differentiation and Gradients

- $\nabla_x f = \text{grad } f = \frac{df}{dx} = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} & \frac{\partial f(x)}{\partial x_n} & \dots & \frac{\partial f(x)}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{1 \times n}$
- n is the number of variables and 1 is the dimension of the image/range/codomain of f
- The row vector $\nabla_x f \in \mathbb{R}^{1 \times n}$ is called the **gradient** of f or the **Jacobian**.

- Example - Partial Derivatives Using the Chain Rule

- For $f(x, y) = (x + 2y^3)^2$, we obtain the partial derivatives

$$\frac{\partial f(x, y)}{\partial x} = 2(x + 2y^3) \frac{\partial}{\partial x} (x + 2y^3) = 2(x + 2y^3)$$

$$\frac{\partial f(x, y)}{\partial y} = 2(x + 2y^3) \frac{\partial}{\partial y} (x + 2y^3) = 12(x + 2y^3)y^2$$

5.2 Partial Differentiation and Gradients

- For $f(x_1, x_2) = x_1^2 x_2 + x_1 x_2^3 \in \mathbb{R}$, the partial derivatives (i.e., the derivatives of f with respect to x_1 and x_2) are

$$\frac{\partial f(x_1, x_2)}{\partial x_1} = 2x_1 x_2 + x_2^3$$
$$\frac{\partial f(x_1, x_2)}{\partial x_2} = x_1^2 + 3x_1 x_2^2$$

- and the gradient is then

$$\frac{df}{dx} = \begin{bmatrix} \frac{\partial f(x_1, x_2)}{\partial x_1} & \frac{\partial f(x_1, x_2)}{\partial x_2} \end{bmatrix} = [2x_1 x_2 + x_2^3 \quad x_1^2 + 3x_1 x_2^2] \in \mathbb{R}^{1 \times 2}$$

5.2.1 Basic Rules of Partial Differentiation

- Product rule:

$$\frac{\partial}{\partial \mathbf{x}}(f(\mathbf{x})g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}}g(\mathbf{x}) + f(\mathbf{x})\frac{\partial g}{\partial \mathbf{x}}$$

- Sum rule:

$$\frac{\partial}{\partial \mathbf{x}}(f(\mathbf{x}) + g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}} + \frac{\partial g}{\partial \mathbf{x}}$$

- Chain rule:

$$\frac{\partial}{\partial \mathbf{x}}(g \circ f)(\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}}\left(g(f(\mathbf{x}))\right) = \frac{\partial g}{\partial f}\frac{\partial f}{\partial \mathbf{x}}$$

5.2.2 Chain Rule

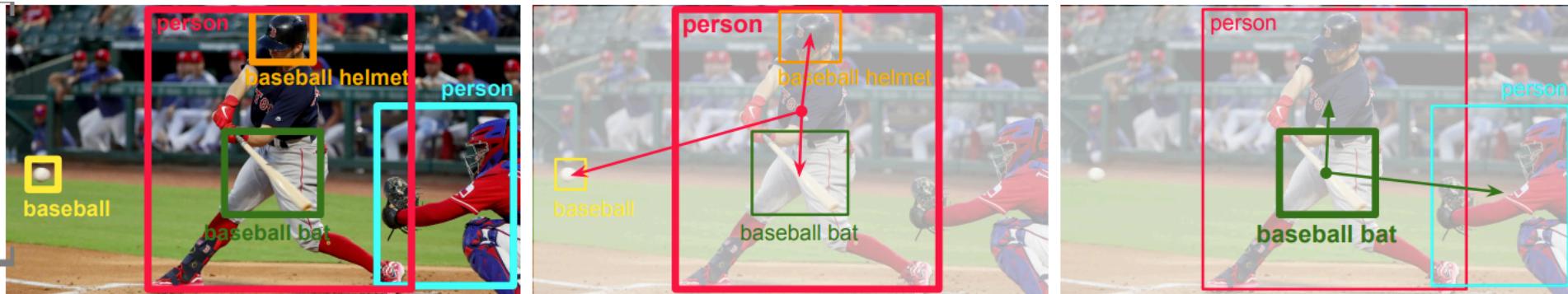
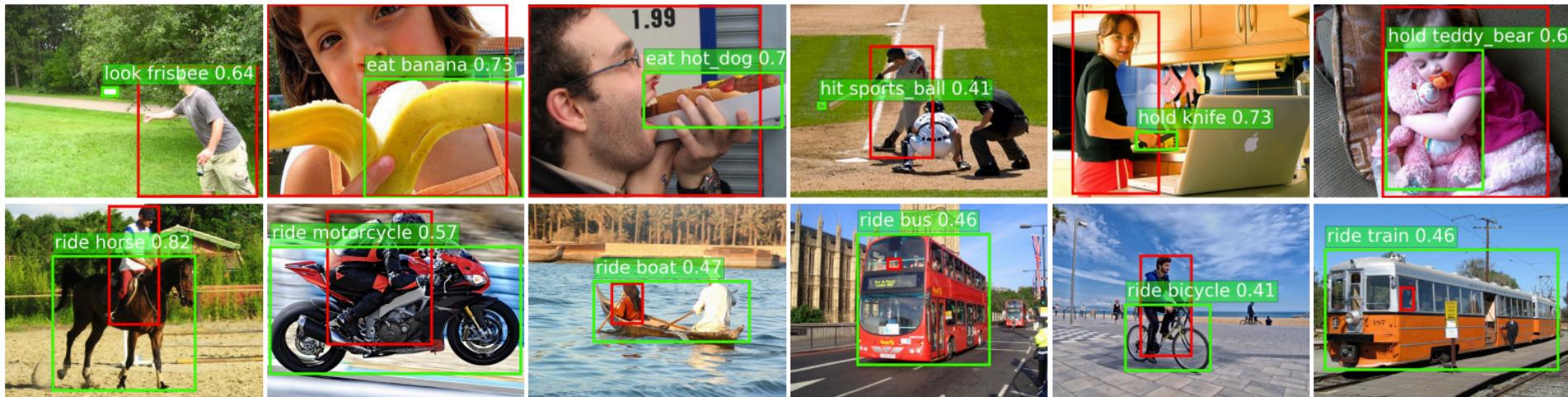
- Consider a function $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ of two variables x_1 and x_2 .
- $x_1(t)$ and $x_2(t)$ are themselves functions of t .
- To compute the gradient of f with respect to t , we apply the chain rule:

$$\frac{df}{dt} = \frac{\partial f}{\partial x} \frac{\partial x}{\partial t} = \left[\frac{\partial f}{\partial x_1} \frac{\partial f}{\partial x_2} \right] \begin{bmatrix} \frac{\partial x_1(t)}{\partial t} \\ \frac{\partial x_2(t)}{\partial t} \end{bmatrix} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t}$$

Where ∂ denotes the gradient and $\frac{\partial}{\partial}$ partial derivates.

- Example
- Consider $f(x_1, x_2) = x_1^2 + 2x_2$, where $x_1 = \sin t$ and $x_2 = \cos t$, then
$$\begin{aligned}\frac{df}{dt} &= \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t} \\ &= 2 \sin t \frac{\partial \sin t}{\partial t} + 2 \frac{\partial \cos t}{\partial t} \\ &= 2 \sin t \cos t - 2 \sin t = 2 \sin t(\cos t - 1)\end{aligned}$$
- The above is the corresponding derivative of f with respect to t .

DRG: Dual Relation Graph for Human-Object Interaction Detection. Gao et al., ECCV 2020



(a) Object detection

(b) Human-centric

(c) Object-centric

5.2.2 Chain Rule

- If $f(x_1, x_2)$ is a function of x_1 and x_2 , where $f: \mathbb{R}^2 \rightarrow \mathbb{R}$, $x_1(s, t)$ and $x_2(s, t)$ are themselves functions of two variables s and t , the chain rule yields the partial derivatives

$$\frac{\partial f}{\partial s} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial s} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial s}$$

$$\frac{\partial f}{\partial t} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t}$$

- The gradient can be obtained by the matrix multiplication

$$\frac{df}{d(s, t)} = \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial (s, t)} = \underbrace{\begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix}}_{= \frac{\partial f}{\partial \mathbf{x}}} \underbrace{\begin{bmatrix} \frac{\partial x_1}{\partial s} & \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial s} & \frac{\partial x_2}{\partial t} \end{bmatrix}}_{= \frac{\partial \mathbf{x}}{\partial (s, t)}}$$

5.3 Gradients of Vector-Valued Functions

- We discussed partial derivatives and gradients of function $f: \mathbb{R}^n \rightarrow \mathbb{R}$
- We will generalize the concept of the gradient to vector-valued functions (vector fields) $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$, where $n \geq 1$ and $m > 1$.
- For a function $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ and a vector $\mathbf{x} = [x_1, \dots, x_n]^T \in \mathbb{R}^n$, the corresponding vector of function values is given as

$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^m$$

- Writing the vector-valued function in this way allows us to view a vector valued function $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ as a vector of functions $[f_1, \dots, f_m]^T$, $f_i: \mathbb{R}^n \rightarrow \mathbb{R}$ that map onto \mathbb{R} .
- The differentiation rules for every f_i are exactly the ones we discussed before.

5.3 Gradients of Vector-Valued Functions

- The partial derivative of a vector-valued function $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ with respect to $x_i \in \mathbb{R}$, $i = 1, \dots, n$, is given as the vector

$$\frac{\partial \mathbf{f}}{\partial x_i} = \begin{bmatrix} \frac{\partial f_1}{\partial x_i} \\ \vdots \\ \frac{\partial f_m}{\partial x_i} \end{bmatrix} = \begin{bmatrix} \lim_{h \rightarrow 0} \frac{f_1(x_1, \dots, x_{i-1}, x_i + h, x_{i+1}, \dots, x_n) - f_1(\mathbf{x})}{h} \\ \vdots \\ \lim_{h \rightarrow 0} \frac{f_m(x_1, \dots, x_{i-1}, x_i + h, x_{i+1}, \dots, x_n) - f_m(\mathbf{x})}{h} \end{bmatrix} \in \mathbb{R}^m$$

- In above, every partial derivative $\frac{\partial \mathbf{f}}{\partial x_i}$ is a column vector
- Recall that the gradient of \mathbf{f} with respect to a vector is the row vector of the partial derivatives
- Therefore, we obtain the gradient of $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ with respect to $\mathbf{x} \in \mathbb{R}^n$, by collecting these partial derivatives:

$$\frac{d\mathbf{f}(\mathbf{x})}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

5.3 Gradients of Vector-Valued Functions

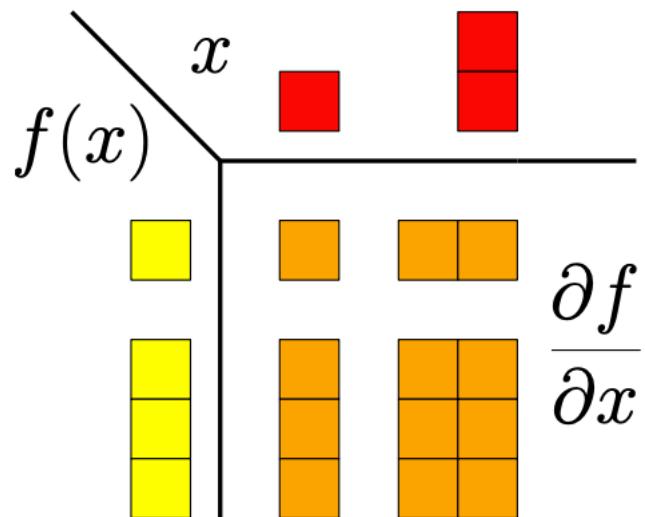
- The collection of all first-order partial derivatives of a vector-valued function $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is called the **Jacobian**. The Jacobian \mathbf{J} is an $m \times n$ matrix, which we define and arrange as follows:

$$\begin{aligned}\mathbf{J} = \nabla_{\mathbf{x}} \mathbf{f} &= \frac{d\mathbf{f}(\mathbf{x})}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_n} \end{bmatrix} \\ &= \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{bmatrix} \\ \mathbf{x} &= \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad J(i, j) = \frac{\partial f_i}{\partial x_j}\end{aligned}$$

- The elements of \mathbf{f} define the rows and the elements of \mathbf{x} define the columns of the corresponding Jacobian
- Special case: for a function $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^1$ which maps a vector $\mathbf{x} \in \mathbb{R}^n$ onto a scalar, i.e., $m = 1$, the Jacobian is a row vector of dimension $1 \times n$.

5.3 Gradients of Vector-Valued Functions

- If $f: \mathbb{R} \rightarrow \mathbb{R}$, the gradient is a scalar
- If $f: \mathbb{R}^D \rightarrow \mathbb{R}$, the gradient is a $1 \times D$ row vector
- If $f: \mathbb{R} \rightarrow \mathbb{R}^E$, the gradient is a $E \times 1$ column vector
- If $f: \mathbb{R}^D \rightarrow \mathbb{R}^E$, the gradient is an $E \times D$ matrix



Example - Gradient of a Vector-Valued Function

- We are given $\mathbf{f}(\mathbf{x}) = \mathbf{Ax}$, $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^M$, $\mathbf{A} \in \mathbb{R}^{M \times N}$, $\mathbf{x} \in \mathbb{R}^N$.
- To compute the gradient $d\mathbf{f}/d\mathbf{x}$ we first determine the dimension of $d\mathbf{f}/d\mathbf{x}$: Since $\mathbf{f}: \mathbb{R}^N \rightarrow \mathbb{R}^M$, it follows that $d\mathbf{f}/d\mathbf{x} \in \mathbb{R}^{M \times N}$.
- Then, we determine the partial derivatives of \mathbf{f} with respect to every x_j :

$$f_i(\mathbf{x}) = \sum_{j=1}^N A_{ij}x_j \Rightarrow \frac{\partial f_i}{\partial x_j} = A_{ij}$$

- We collect the partial derivatives in the Jacobian and obtain the gradient

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \dots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} = \begin{bmatrix} A_{11} & \dots & A_{1N} \\ \vdots & & \vdots \\ A_{M1} & \dots & A_{MN} \end{bmatrix} = \mathbf{A} \in \mathbb{R}^{M \times N}$$

Example - Chain Rule

- Consider the function $h: \mathbb{R} \rightarrow \mathbb{R}$, $h(t) = (f \circ g)(t)$ with

$$f: \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$g: \mathbb{R} \rightarrow \mathbb{R}^2$$

$$f(\mathbf{x}) = \exp(x_1 x_2^2)$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = g(t) = \begin{bmatrix} t \cos t \\ t \sin t \end{bmatrix}$$

- We compute the gradient of h with respect to t . Since $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ and $g: \mathbb{R} \rightarrow \mathbb{R}^2$ we note that

$$\frac{\partial f}{\partial \mathbf{x}} \in \mathbb{R}^{1 \times 2}, \quad \frac{\partial g}{\partial t} \in \mathbb{R}^{2 \times 1}$$

- The desired gradient is computed by applying the chain rule:

$$\begin{aligned} \frac{dh}{dt} &= \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial t} = \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \right] \begin{bmatrix} \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial t} \end{bmatrix} \\ &= [\exp(x_1 x_2^2) x_2^2 \quad 2\exp(x_1 x_2^2) x_1 x_2] \begin{bmatrix} \cos t - t \sin t \\ \sin t + t \cos t \end{bmatrix} \\ &= \exp(x_1 x_2^2) \left(x_2^2 (\cos t - t \sin t) + 2x_1 x_2 (\sin t + t \cos t) \right) \end{aligned}$$

where $x_1 = t \cos t$ and $x_2 = t \sin t$

Example - Gradient of a Least-Squares Loss in a Linear Model

- Let us consider the linear model

$$\mathbf{y} = \Phi\boldsymbol{\theta}$$

where $\boldsymbol{\theta} \in \mathbb{R}^D$ is a parameter vector, $\Phi \in \mathbb{R}^{N \times D}$ are input features and $\mathbf{y} \in \mathbb{R}^N$ are the corresponding observations. We define the functions

$$\begin{aligned} L(\mathbf{e}) &:= \| \mathbf{e} \|^2, \\ \mathbf{e}(\boldsymbol{\theta}) &:= \mathbf{y} - \Phi\boldsymbol{\theta} \end{aligned}$$

- We seek $\frac{\partial L}{\partial \boldsymbol{\theta}}$, and we will use the chain rule for this purpose. L is called a least-squares loss function.
- First, we determine the dimensionality of the gradient as

$$\frac{\partial L}{\partial \boldsymbol{\theta}} \in \mathbb{R}^{1 \times D}$$

- The chain rule allows us to compute the gradient as

$$\frac{\partial L}{\partial \boldsymbol{\theta}} = \frac{\partial L}{\partial \mathbf{e}} \frac{\partial \mathbf{e}}{\partial \boldsymbol{\theta}}$$

Example - Gradient of a Least-Squares Loss in a Linear Model

- We know that $\|e\|^2 = e^T e$ and determine

$$\frac{\partial L}{\partial e} = 2e^T \in \mathbb{R}^{1 \times N}$$

- Further, we obtain

$$\frac{\partial e}{\partial \theta} = -\Phi \in \mathbb{R}^{N \times D}$$

- Our desired derivative is

$$\frac{\partial L}{\partial \theta} = -2e^T \Phi = - \underbrace{2(y^T - \theta^T \Phi^T)}_{1 \times N} \underbrace{\Phi}_{N \times D} \in \mathbb{R}^{1 \times D}$$

5.4 Gradients of Matrices

- Consider the following example

$$\mathbf{f} = \mathbf{A}\mathbf{x}, \quad \mathbf{f} \in \mathbb{R}^M, \quad \mathbf{A} \in \mathbb{R}^{M \times N}, \quad \mathbf{x} \in \mathbb{R}^N$$

- We seek the gradient $\frac{d\mathbf{f}}{d\mathbf{A}}$
- First, we determine the dimension of the gradient

$$\frac{d\mathbf{f}}{d\mathbf{A}} \in \mathbb{R}^{M \times (M \times N)}$$

- By definition, the gradient is the collection of the partial derivatives:

$$\frac{d\mathbf{f}}{d\mathbf{A}} = \begin{bmatrix} \frac{\partial f_1}{\partial \mathbf{A}} \\ \vdots \\ \frac{\partial f_M}{\partial \mathbf{A}} \end{bmatrix}, \quad \frac{\partial f_i}{\partial \mathbf{A}} \in \mathbb{R}^{1 \times (M \times N)}$$

- To compute the partial derivatives, we explicitly write out the matrix vector multiplication

$$f_i = \sum_{j=1}^N A_{ij}x_j, \quad i = 1, \dots, M,$$

$$f_i = \sum_{j=1}^N A_{ij}x_j, \quad i = 1, \dots, M,$$

- The partial derivatives are then given as

$$\frac{\partial f_i}{\partial A_{iq}} = x_q$$

- Partial derivatives of f_i with respect to a row of \mathbf{A} are given as

$$\frac{\partial f_i}{\partial A_{i,:}} = \mathbf{x}^T \in \mathbb{R}^{1 \times 1 \times N}, \quad \frac{\partial f_i}{\partial A_{k \neq i,:}} = \mathbf{0}^T \in \mathbb{R}^{1 \times 1 \times N}$$

- Since f_i maps onto \mathbb{R} and each row of \mathbf{A} is of size $1 \times N$, we obtain a $1 \times 1 \times N$ sized tensor as the partial derivative of f_i with respect to a row of \mathbf{A} .
- We stack the partial derivatives and get the desired gradient

$$\frac{\partial f_i}{\partial \mathbf{A}} = \begin{bmatrix} \mathbf{0}^T \\ \vdots \\ \mathbf{0}^T \\ \mathbf{x}^T \\ \mathbf{0}^T \\ \vdots \\ \mathbf{0}^T \end{bmatrix} \in \mathbb{R}^{1 \times (M \times N)}$$

Example - Gradient of Matrices with Respect to Matrices

- Consider a matrix $\mathbf{R} \in \mathbb{R}^{M \times N}$ and $f: \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{N \times N}$ with
$$f(\mathbf{R}) = \mathbf{R}^T \mathbf{R} =: \mathbf{K} \in \mathbb{R}^{N \times N}$$

- We seek the gradient $\frac{d\mathbf{K}}{d\mathbf{R}}$

- First, the dimension of the gradient is given as

$$\frac{d\mathbf{K}}{d\mathbf{R}} \in \mathbb{R}^{(N \times N) \times (M \times N)}$$

$$\frac{dK_{pq}}{d\mathbf{R}} \in \mathbb{R}^{1 \times M \times N}$$

for $p, q = 1, \dots, N$, where K_{pq} is the pq th entry of $\mathbf{K} = f(\mathbf{R})$.

- Denoting the i th column of \mathbf{R} by \mathbf{r}_i , every entry of \mathbf{K} is given by the dot product of two columns of \mathbf{R} , i.e.,

$$K_{pq} = \mathbf{r}_p^T \mathbf{r}_q = \sum_{m=1}^M R_{mp} R_{mq}$$

Example - Gradient of Matrices with Respect to Matrices

- Denoting the i th column of \mathbf{R} by \mathbf{r}_i , every entry of \mathbf{K} is given by the dot product of two columns of \mathbf{R} , i.e.,

$$K_{pq} = \mathbf{r}_p^T \mathbf{r}_q = \sum_{m=1}^M R_{mp} R_{mq}$$

- We now compute the partial derivative $\frac{\partial K_{pq}}{\partial R_{ij}}$, we obtain

$$\frac{\partial K_{pq}}{\partial R_{ij}} = \sum_{m=1}^M \frac{\partial}{\partial R_{ij}} R_{mp} R_{mq} = \partial_{pqij}$$

$$\partial_{pqij} = \begin{cases} R_{iq} & \text{if } j = p, p \neq q \\ R_{ip} & \text{if } j = q, p \neq q \\ 2R_{iq} & \text{if } j = p, p = q \\ 0 & \text{otherwise} \end{cases}$$

- The desired gradient has the dimension $(N \times N) \times (M \times N)$, and every single entry of this tensor is given by ∂_{pqij} , where $p, q, j = 1, \dots, N$ and $i = 1, \dots, M$

5.5 Useful Identities for Computing Gradients

- Some useful gradients that are frequently required in machine learning
- $\text{tr}(\cdot)$: trace $\det(\cdot)$: determinant $f(X)^{-1}$: the inverse of $f(X)$

$$\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \mathbf{a}^T$$

$$\frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}^T$$

$$\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^T$$

$$\frac{\partial \mathbf{x}^T \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = \mathbf{x}^T (\mathbf{B} + \mathbf{B}^T)$$

$$\frac{\partial}{\partial \mathbf{s}} (\mathbf{x} - \mathbf{A}\mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{A}\mathbf{s}) = -2(\mathbf{x} - \mathbf{A}\mathbf{s})^T \mathbf{W} \mathbf{A} \quad \text{for symmetric } \mathbf{W}$$

You should be able to calculate these gradients