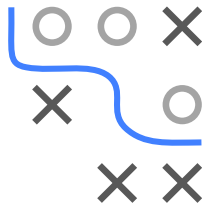


Optimization in Machine Learning

Mathematical Concepts

Matrix Calculus



δ

Learning goals

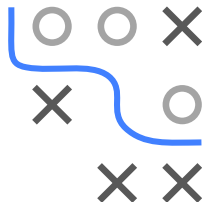
- Rules of matrix calculus
- Connection of gradient, Jacobian and Hessian

SCOPE

- \mathcal{X}/\mathcal{Y} denote space of **independent/dependent** variables
- Identify dependent variable y with a **function**
 $f : \mathcal{X} \rightarrow \mathcal{Y}, x \mapsto f(x)$
- Assume y sufficiently smooth
- In matrix calculus, x and y can be **scalars**, **vectors**, or **matrices**
- We denote vectors/matrices in **bold** lowercase/uppercase letters

Type	scalar x	vector \mathbf{x}	matrix \mathbf{X}
scalar y	dy/dx	$dy/d\mathbf{x}$	$dy/d\mathbf{X}$
vector \mathbf{y}	$d\mathbf{y}/dx$	$d\mathbf{y}/d\mathbf{x}$	—
matrix \mathbf{Y}	$d\mathbf{Y}/dx$	—	—

- This notation is also referred to as *Leibniz notation*



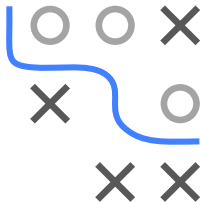
LEIBNIZ NOTATION CONVENTION

- Instead of writing $f(x)$ everywhere, we replace the function f with the variable y .
- This helps clarify relationships when multiple functions or variables are involved, especially in contexts like partial derivatives or matrix calculus.
- Also applicable to partial derivatives: For $y = f : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{x} \mapsto f(\mathbf{x}) = f(x_1, x_2, \dots, x_n)$, the partial derivative of y w.r.t. x_i is dy/dx_i .

- **Examples:**

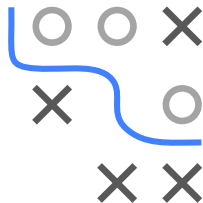
- $y = x^3 + 5x \implies \frac{dy}{dx} = 3x^2 + 5$

- $y = x_1^2 + 3x_2 \implies \frac{dy}{dx_1} = 2x_1, \quad \frac{dy}{dx_2} = 3$



DERIVATIVES OF SCALAR-VALUED FUNCTIONS

Type	scalar x	vector \mathbf{x}	matrix \mathbf{X}
scalar y	dy/dx	$dy/d\mathbf{x}$	$dy/d\mathbf{X}$
vector \mathbf{y}	$d\mathbf{y}/dx$	$d\mathbf{y}/d\mathbf{x}$	—
matrix \mathbf{Y}	$d\mathbf{Y}/dx$	—	—

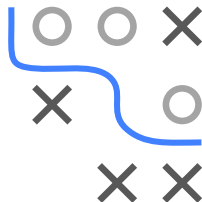


- $dy/d\mathbf{x}$ is the gradient from the previous slide deck
- When the input is a matrix the concept remains the same, i.e. for $y = f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$, $\mathbf{X} \mapsto f(\mathbf{X})$

$$\frac{dy}{d\mathbf{X}} = \begin{pmatrix} \frac{\partial f}{\partial x_{11}} & \cdots & \frac{\partial f}{\partial x_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_{m1}} & \cdots & \frac{\partial f}{\partial x_{mn}} \end{pmatrix} \in \mathbb{R}^{m \times n}$$

DERIVATIVES: UNIVARIATE AND JACOBIAN

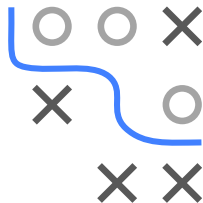
Type	scalar x	vector \mathbf{x}	matrix \mathbf{X}
scalar y	dy/dx	$dy/d\mathbf{x}$	$dy/d\mathbf{X}$
vector \mathbf{y}	$d\mathbf{y}/dx$	$d\mathbf{y}/d\mathbf{x}$	—
matrix \mathbf{Y}	$d\mathbf{Y}/dx$	—	—



- dy/dx is the univariate derivative y'
- $d\mathbf{y}/d\mathbf{x}$ is the Jacobian from the previous slide deck

DERIV. OF FUNCTIONS WITH SCALAR INPUTS

Type	scalar x	vector \mathbf{x}	matrix \mathbf{X}
scalar y	dy/dx	$dy/d\mathbf{x}$	$dy/d\mathbf{X}$
vector \mathbf{y}	$d\mathbf{y}/dx$	$d\mathbf{y}/d\mathbf{x}$	—
matrix \mathbf{Y}	$d\mathbf{Y}/dx$	—	—



- Here, for univariate $f_{ij} : \mathbb{R} \rightarrow \mathbb{R}$, \mathbf{y} ($n = 1$) or \mathbf{Y} ($n > 1$) are equal to a function $f : \mathbb{R} \rightarrow \mathbb{R}^{m \times n}$, $x \mapsto (f_{ij}(x))_{i=1,\dots,m; j=1,\dots,n}$ and the derivatives are, respectively given by

$$\frac{d\mathbf{y}}{dx} = \begin{pmatrix} \frac{\partial f_1}{\partial x} \\ \vdots \\ \frac{\partial f_m}{\partial x} \end{pmatrix} \in \mathbb{R}^m; \quad \frac{d\mathbf{Y}}{dx} = \begin{pmatrix} \frac{\partial f_{11}}{\partial x} & \cdots & \frac{\partial f_{1n}}{\partial x} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_{m1}}{\partial x} & \cdots & \frac{\partial f_{mn}}{\partial x} \end{pmatrix} \in \mathbb{R}^{m \times n}$$

MULTIVARIATE DIFFERENTIATION RULES

- Basic rules from single-variable calculus still apply.
- But, for $\mathbf{x} \in \mathbb{R}^n$: gradients are vectors/matrices (order matters).

Key Rules:

- **Sum:**

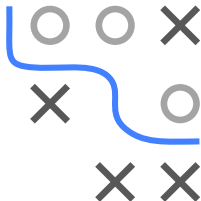
$$\frac{d}{d\mathbf{x}}(f + g) = \frac{df}{d\mathbf{x}} + \frac{dg}{d\mathbf{x}}$$

- **Product:**

$$\frac{d}{d\mathbf{x}}(fg) = \frac{df}{d\mathbf{x}} g + f \frac{dg}{d\mathbf{x}}$$

- **Chain:**

$$\frac{d}{d\mathbf{x}}((f \circ g)(\mathbf{x})) = \frac{d}{d\mathbf{x}}(f(g(\mathbf{x}))) = \frac{df}{dg} \frac{dg}{d\mathbf{x}}$$



DETAILS ON THE CHAIN RULE I

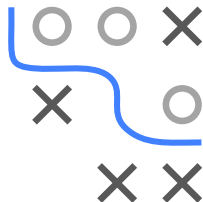
- Suppose
 - we have functions $\mathbf{g} : S \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $\mathbf{f} : T \subseteq \mathbb{R}^m \rightarrow \mathbb{R}^\ell$
 - $\mathbf{a} \in S$ is a point such that $\mathbf{g}(\mathbf{a}) \in T \Rightarrow \mathbf{f} \circ \mathbf{g}(\mathbf{x}) = \mathbf{f}(\mathbf{g}(\mathbf{x}))$ is well-defined for all \mathbf{x} close to \mathbf{a}
- Then, if \mathbf{g} is differentiable at \mathbf{a} and \mathbf{f} is differentiable at $\mathbf{g}(\mathbf{a})$
 $\Rightarrow \mathbf{f} \circ \mathbf{g}$ is differentiable at \mathbf{a} , and the derivative $\frac{d\mathbf{f}}{d\mathbf{g}} \frac{d\mathbf{g}}{d\mathbf{x}}$, is equal to

$$\nabla_{\mathbf{a}} \mathbf{f} \circ \mathbf{g} \hat{=} \mathbf{J}_{\mathbf{f} \circ \mathbf{g}}(\mathbf{a}) = \mathbf{J}_{\mathbf{f}}(\mathbf{g}(\mathbf{a})) \mathbf{J}_{\mathbf{g}}(\mathbf{a}) \hat{=} \nabla_{\mathbf{g}(\mathbf{a})} \mathbf{f} \nabla_{\mathbf{a}} \mathbf{g} \in \mathbb{R}^{\ell \times n}$$

(See Chapter 2.3 of the UofT course *MAT237 - Multivariable Calculus* for proof)

- We can also write \mathbf{f} as a function of $\mathbf{y} = (y_1, \dots, y_m) \in \mathbb{R}^m$, and \mathbf{g} as a function of $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$. Then for each $k = 1, \dots, \ell$ and $j = 1, \dots, n$

$$[\mathbf{J}_{\mathbf{f} \circ \mathbf{g}}(\mathbf{a})]_{kj} = \frac{\partial}{\partial x_j} (f_k \circ \mathbf{g})(\mathbf{a}) = \sum_{i=1}^m \frac{\partial f_k}{\partial y_i}(\mathbf{g}(\mathbf{a})) \frac{\partial g_i}{\partial x_j}(\mathbf{a})$$

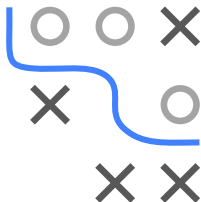


HELPFUL CALCULATION RULES

Let \mathbf{a} , \mathbf{b} denote vectors, \mathbf{X} , \mathbf{A} matrices, and $f(\mathbf{X})^{-1}$ the inverse of $f(\mathbf{X})$ if it exists.

- $\frac{d\mathbf{x}^\top \mathbf{a}}{d\mathbf{x}} = \mathbf{a}^\top$, $\frac{d\mathbf{a}^\top \mathbf{x}}{d\mathbf{x}} = \mathbf{a}^\top$
- $\frac{d\mathbf{X}\mathbf{a}}{d\mathbf{a}} = \mathbf{X}$, $\frac{d\mathbf{a}^\top \mathbf{X}}{d\mathbf{a}} = \mathbf{X}^\top$
- $\frac{d\mathbf{a}^\top \mathbf{X}\mathbf{b}}{d\mathbf{X}} = \mathbf{a}\mathbf{b}^\top$
- $\frac{d\mathbf{x}^\top \mathbf{A}\mathbf{x}}{d\mathbf{x}} = \mathbf{x}^\top (\mathbf{A} + \mathbf{A}^\top)$
- For a symmetric matrix \mathbf{W} ,
$$\frac{d}{ds}(\mathbf{x} - \mathbf{A}\mathbf{s})^\top \mathbf{W}(\mathbf{x} - \mathbf{A}\mathbf{s}) = -2(\mathbf{x} - \mathbf{A}\mathbf{s})^\top \mathbf{W}\mathbf{A}$$
- $\frac{d}{d\mathbf{X}} \mathbf{f}(\mathbf{X})^\top = \left(\frac{d\mathbf{f}(\mathbf{X})}{d\mathbf{X}} \right)^\top$
- $\frac{d}{d\mathbf{X}} \mathbf{f}(\mathbf{X})^{-1} = -\mathbf{f}(\mathbf{X})^{-1} \frac{d\mathbf{f}(\mathbf{X})}{d\mathbf{X}} \mathbf{f}(\mathbf{X})^{-1}$
- $\frac{d\mathbf{a}^\top \mathbf{X}^{-1} \mathbf{b}}{d\mathbf{X}} = -(\mathbf{X}^{-1})^\top \mathbf{a}\mathbf{b}^\top (\mathbf{X}^{-1})^\top$

Note: to compute gradients of matrices with respect to vectors (or other matrices) we need *tensors*, see chapter 5.4 of [Deisenroth](#) for more.



EXAMPLE: LOGISTIC REGRESSION I

- Let's say, for data in $\mathbb{R}^{n \times m}$ we're trying to minimize the risk in logistic regression by finding the gradient for negative log loss:

$$-\ell(\theta) = \sum_{i=1}^n -y^{(i)} \log \left(\pi \left(\mathbf{x}^{(i)} \mid \theta \right) \right) - \left(1 - y^{(i)} \right) \log \left(1 - \pi \left(\mathbf{x}^{(i)} \mid \theta \right) \right)$$

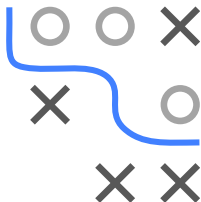
where $\pi(\mathbf{x} \mid \theta) = s(f(\theta, \mathbf{x}))$

with $f(\theta, \mathbf{x}) = \theta^\top \mathbf{x}$ and $s(x) = \frac{1}{1 + \exp(-x)}$.

⇒ We want to find

$$\nabla_{\theta} - \ell(\theta) = -\nabla_{\theta} \ell(\theta) = - \sum_{i=1}^n \underbrace{y^{(i)} \log(s(f(\theta, \mathbf{x}^{(i)}))) + (1 - y^{(i)}) \log(1 - s(f(\theta, \mathbf{x}^{(i)})))}_{=: y_i \log(s_i) + (1 - y_i) \log(1 - s_i)}$$

- $s_i = s(f_i)$; $h_i := -[y_i \log(s_i) + (1 - y_i) \log(1 - s_i)]$; $f_i := \theta^\top \mathbf{x}^{(i)}$



EX: LOG. REGR. – ALTERNATIVE NOTATION

Alternatively, we could write this example out as follows:

$$-[y \log(s(f(\theta, \mathbf{x}))) + (1 - y) \log(1 - s(f(\theta, \mathbf{x})))]$$

$$\nabla_{\theta} h(y, \cdot) \circ \mathbf{s} \circ f(\cdot, \mathbf{x}).$$

$$\nabla_{\theta} h(y, \cdot) \circ s \circ f(\cdot, \mathbf{x}) = \frac{dh}{ds} \frac{ds}{df} \frac{df}{d\theta} \in \mathbb{R}^{1 \times p}.$$