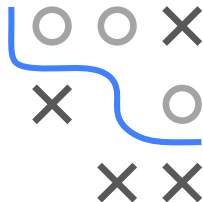


Optimization in Machine Learning

Mathematical Concepts

Matrix Calculus



Learning goals

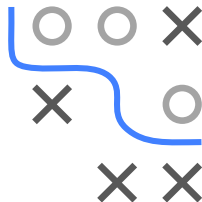
- Rules of matrix calculus
- Connection of gradient, Jacobian and Hessian

SCOPE

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

Type	scalar x	vector \mathbf{x}	matrix \mathbf{X}
scalar y	$\partial y / \partial x$	$\partial y / \partial \mathbf{x}$	$\partial y / \partial \mathbf{X}$
vector \mathbf{y}	$\partial \mathbf{y} / \partial x$	$\partial \mathbf{y} / \partial \mathbf{x}$	–
matrix \mathbf{Y}	$\partial \mathbf{Y} / \partial x$	–	–

- We denote vectors/matrices in **bold** lowercase/uppercase letters



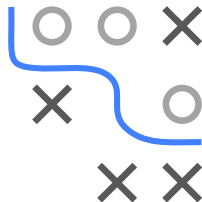
NUMERATOR LAYOUT

- **Matrix calculus:** collect derivative of each component of dependent variable w.r.t. each component of independent variable
- We use so-called **numerator layout** convention:

$$\frac{\partial y}{\partial \mathbf{x}} = \left(\frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_d} \right) = \nabla y^T \in \mathbb{R}^{1 \times d}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \left(\frac{\partial y_1}{\partial \mathbf{x}}, \dots, \frac{\partial y_m}{\partial \mathbf{x}} \right)^T \in \mathbb{R}^m$$

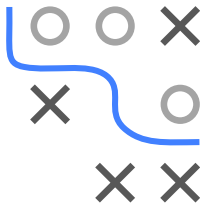
$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial \mathbf{x}} \\ \vdots \\ \frac{\partial y_m}{\partial \mathbf{x}} \end{pmatrix} = \left(\frac{\partial \mathbf{y}}{\partial x_1} \dots \frac{\partial \mathbf{y}}{\partial x_d} \right) = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_d} \end{pmatrix} = \mathbf{J}_y \in \mathbb{R}^{m \times d}$$



SCALAR-BY-VECTOR

Let $\mathbf{x} \in \mathbb{R}^d$, $y, z : \mathbb{R}^d \rightarrow \mathbb{R}$ and \mathbf{A} be a matrix.

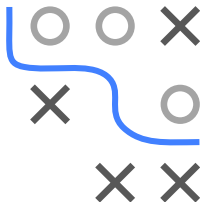
- If y is a **constant** function: $\frac{\partial y}{\partial \mathbf{x}} = \mathbf{0}^T \in \mathbb{R}^{1 \times d}$
- **Linearity**: $\frac{\partial(a \cdot y + z)}{\partial \mathbf{x}} = a \frac{\partial y}{\partial \mathbf{x}} + \frac{\partial z}{\partial \mathbf{x}}$ (a constant)
- **Product** rule: $\frac{\partial(y \cdot z)}{\partial \mathbf{x}} = y \frac{\partial z}{\partial \mathbf{x}} + \frac{\partial y}{\partial \mathbf{x}} z$
- **Chain** rule: $\frac{\partial g(y)}{\partial \mathbf{x}} = \frac{\partial g(y)}{\partial y} \frac{\partial y}{\partial \mathbf{x}}$ (g scalar-valued function)
- **Second** derivative: $\frac{\partial^2 y}{\partial \mathbf{x} \partial \mathbf{x}^T} = \nabla^2 y^T (= \nabla^2 y$ if $y \in \mathcal{C}^2)$ (Hessian)
- $\frac{\partial(\mathbf{x}^T \mathbf{A} \mathbf{x})}{\partial \mathbf{x}} = \mathbf{x}^T (\mathbf{A} + \mathbf{A}^T)$
- $\frac{\partial(\mathbf{y}^T \mathbf{A} \mathbf{z})}{\partial \mathbf{x}} = \mathbf{y}^T \mathbf{A} \frac{\partial \mathbf{z}}{\partial \mathbf{x}} + \mathbf{z}^T \mathbf{A}^T \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ (\mathbf{y}, \mathbf{z} vector-valued functions of \mathbf{x})



VECTOR-BY-SCALAR

Let $x \in \mathbb{R}$ and $\mathbf{y}, \mathbf{z} : \mathbb{R} \rightarrow \mathbb{R}^m$.

- If \mathbf{y} is a **constant** function: $\frac{\partial \mathbf{y}}{\partial x} = \mathbf{0} \in \mathbb{R}^m$
- **Linearity:** $\frac{\partial (a \cdot \mathbf{y} + \mathbf{z})}{\partial x} = a \frac{\partial \mathbf{y}}{\partial x} + \frac{\partial \mathbf{z}}{\partial x}$ (a constant)
- **Chain rule:** $\frac{\partial \mathbf{g}(\mathbf{y})}{\partial x} = \frac{\partial \mathbf{g}(\mathbf{y})}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial x}$ (\mathbf{g} vector-valued function)
- $\frac{\partial (\mathbf{A}\mathbf{y})}{\partial x} = \mathbf{A} \frac{\partial \mathbf{y}}{\partial x}$ (\mathbf{A} matrix)



VECTOR-BY-VECTOR

Let $\mathbf{x} \in \mathbb{R}^d$ and $\mathbf{y}, \mathbf{z} : \mathbb{R}^d \rightarrow \mathbb{R}^m$.

- If \mathbf{y} is a **constant** function: $\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{0} \in \mathbb{R}^{m \times d}$
- $\frac{\partial \mathbf{x}}{\partial \mathbf{x}} = \mathbf{I} \in \mathbb{R}^{d \times d}$
- **Linearity:** $\frac{\partial (a \cdot \mathbf{y} + \mathbf{z})}{\partial \mathbf{x}} = a \frac{\partial \mathbf{y}}{\partial \mathbf{x}} + \frac{\partial \mathbf{z}}{\partial \mathbf{x}}$ (a constant)
- **Chain rule:** $\frac{\partial \mathbf{g}(\mathbf{y})}{\partial \mathbf{x}} = \frac{\partial \mathbf{g}(\mathbf{y})}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ (\mathbf{g} vector-valued function)
- $\frac{\partial (\mathbf{A}\mathbf{x})}{\partial \mathbf{x}} = \mathbf{A}$, $\frac{\partial (\mathbf{x}^T \mathbf{B})}{\partial \mathbf{x}} = \mathbf{B}^T$ (\mathbf{A}, \mathbf{B} matrices)

