Mathematics for Machine Learning

Vazgen Mikayelyan

August 25, 2020



Math for ML

× /	A 411		
v	- Milli	kavi	elvar

Definition

For a function $f : \mathbb{R}^n \to \mathbb{R}$ (of *n* variables x_1, \ldots, x_n)

V.	Mikay	/el\	/ar

Definition

For a function $f : \mathbb{R}^n \to \mathbb{R}$ (of *n* variables x_1, \ldots, x_n) we define the **partial derivatives** as

$$f'_{x_1}(\mathbf{x}) = \frac{\partial f}{\partial x_1}(\mathbf{x}) = \lim_{h \to 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(x_1, x_2, \dots, x_n)}{h}$$

$$\vdots$$

$$f'_{x_n}(\mathbf{x}) = \frac{\partial f}{\partial x_n}(\mathbf{x}) = \lim_{h \to 0} \frac{f(x_1, x_2, \dots, x_n + h) - f(x_1, x_2, \dots, x_n)}{h}.$$

< □ > < 同 >

Definition

For a function $f : \mathbb{R}^n \to \mathbb{R}$ (of *n* variables x_1, \ldots, x_n) we define the partial derivatives as

$$f'_{x_1}(\mathbf{x}) = \frac{\partial f}{\partial x_1}(\mathbf{x}) = \lim_{h \to 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(x_1, x_2, \dots, x_n)}{h}$$

$$\vdots$$

$$f'_{x_n}(\mathbf{x}) = \frac{\partial f}{\partial x_n}(\mathbf{x}) = \lim_{h \to 0} \frac{f(x_1, x_2, \dots, x_n + h) - f(x_1, x_2, \dots, x_n)}{h}.$$

The row vector

$$\nabla f = \frac{df}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} & \frac{\partial f(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{1 \times n},$$

is called the gradient of f or the Jacobian.

V.	Mikay	/el\	/ar
		_	

< □ > < 凸

Differentiation Rules

If the functions f and g have partial derivatives, then Sum Rule: $\frac{\partial}{\partial x_i}(f(\mathbf{x}) + g(\mathbf{x})) = \frac{\partial f(\mathbf{x})}{\partial x_i} + \frac{\partial g(\mathbf{x})}{\partial x_i}$ Product Rule: $\frac{\partial}{\partial x_i}(f(\mathbf{x}) \cdot g(\mathbf{x})) = \frac{\partial f(\mathbf{x})}{\partial x_i}g(\mathbf{x}) + f(\mathbf{x})\frac{\partial g(\mathbf{x})}{\partial x_i}$

Differentiation Rules

If the functions f and g have partial derivatives, then Sum Rule: $\frac{\partial}{\partial x_i}(f(\mathbf{x}) + g(\mathbf{x})) = \frac{\partial f(\mathbf{x})}{\partial x_i} + \frac{\partial g(\mathbf{x})}{\partial x_i}$ Product Rule: $\frac{\partial}{\partial x_i}(f(\mathbf{x}) \cdot g(\mathbf{x})) = \frac{\partial f(\mathbf{x})}{\partial x_i}g(\mathbf{x}) + f(\mathbf{x})\frac{\partial g(\mathbf{x})}{\partial x_i}$

Example

Find the gradient of the following functions:

a)
$$f : \mathbb{R}^2 \to \mathbb{R}$$
 defined by $f(x_1, x_2) = (2x_1 + 3x_2)^3$

b)
$$f: \mathbb{R}^3 \to \mathbb{R}$$
 defined by $f(x, y, z) = e^{2x} + y^2 z^3$

Differentiation Rules

If the functions f and g have partial derivatives, then Sum Rule: $\frac{\partial}{\partial x_i}(f(\mathbf{x}) + g(\mathbf{x})) = \frac{\partial f(\mathbf{x})}{\partial x_i} + \frac{\partial g(\mathbf{x})}{\partial x_i}$ Product Rule: $\frac{\partial}{\partial x_i}(f(\mathbf{x}) \cdot g(\mathbf{x})) = \frac{\partial f(\mathbf{x})}{\partial x_i}g(\mathbf{x}) + f(\mathbf{x})\frac{\partial g(\mathbf{x})}{\partial x_i}$

Example

Find the gradient of the following functions:

a)
$$f: \mathbb{R}^2 \to \mathbb{R}$$
 defined by $f(x_1, x_2) = (2x_1 + 3x_2)^3$

b)
$$f: \mathbb{R}^3 \to \mathbb{R}$$
 defined by $f(x, y, z) = e^{2x} + y^2 z^3$

Example

Given
$$z(x, y) = x^2 + y^2$$
 where $x(r, t) = r \cos(t)$ and $y(r, t) = r + t$,
determine the value of $\frac{\partial z}{\partial t}$ and $\frac{\partial z}{\partial r}$.

V. Mikayelyan

V.	Mika	vel	van
			,

2

メロト メポト メヨト メヨト

Let z be a function of two variables, x, y and each of these variables x, y be in turn functions of two variables, t, s.

Let z be a function of two variables, x,y and each of these variables x,y be in turn functions of two variables, $t,s.{\rm Then}$

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

Let z be a function of two variables, x, y and each of these variables x, y be in turn functions of two variables, t, s. Then

$$\frac{\partial z}{\partial t} = \frac{\partial z}{\partial x}\frac{\partial x}{\partial t} + \frac{\partial z}{\partial y}\frac{\partial y}{\partial t} = \begin{bmatrix} \frac{\partial z}{\partial x} & \frac{\partial z}{\partial y} \end{bmatrix} \begin{bmatrix} \frac{\partial x}{\partial t} \\ \frac{\partial y}{\partial t} \end{bmatrix} = \nabla z\frac{\partial \mathbf{x}}{\partial t}$$

Let z be a function of two variables, x,y and each of these variables x,y be in turn functions of two variables, $t,s.{\rm Then}$

$$\frac{\partial z}{\partial t} = \frac{\partial z}{\partial x}\frac{\partial x}{\partial t} + \frac{\partial z}{\partial y}\frac{\partial y}{\partial t} = \begin{bmatrix} \frac{\partial z}{\partial x} & \frac{\partial z}{\partial y} \end{bmatrix} \begin{bmatrix} \frac{\partial x}{\partial t} \\ \frac{\partial y}{\partial t} \end{bmatrix} = \nabla z\frac{\partial \mathbf{x}}{\partial t}$$

$$\frac{\partial z}{\partial s} = \frac{\partial z}{\partial x}\frac{\partial x}{\partial s} + \frac{\partial z}{\partial y}\frac{\partial y}{\partial s} = \begin{bmatrix} \frac{\partial z}{\partial x} & \frac{\partial z}{\partial y} \end{bmatrix} \begin{bmatrix} \frac{\partial x}{\partial s} \\ \frac{\partial y}{\partial s} \end{bmatrix} = \nabla z\frac{\partial \mathbf{x}}{\partial s}.$$

Let z be a function of two variables, x,y and each of these variables x,y be in turn functions of two variables, $t,s.{\rm Then}$

$$\frac{\partial z}{\partial t} = \frac{\partial z}{\partial x}\frac{\partial x}{\partial t} + \frac{\partial z}{\partial y}\frac{\partial y}{\partial t} = \begin{bmatrix} \frac{\partial z}{\partial x} & \frac{\partial z}{\partial y} \end{bmatrix} \begin{bmatrix} \frac{\partial x}{\partial t} \\ \frac{\partial y}{\partial t} \end{bmatrix} = \nabla z\frac{\partial \mathbf{x}}{\partial t}$$

$$\frac{\partial z}{\partial s} = \frac{\partial z}{\partial x}\frac{\partial x}{\partial s} + \frac{\partial z}{\partial y}\frac{\partial y}{\partial s} = \begin{bmatrix} \frac{\partial z}{\partial x} & \frac{\partial z}{\partial y} \end{bmatrix} \begin{bmatrix} \frac{\partial x}{\partial s} \\ \frac{\partial y}{\partial s} \end{bmatrix} = \nabla z\frac{\partial \mathbf{x}}{\partial s}.$$



Given $z(x,y) = x^2 + y^2$ where $x(r,t) = r\cos(t)$ and $y(r,t) = r\sin(t)$, determine the value of $\frac{\partial z}{\partial t}$ and $\frac{\partial z}{\partial r}$ using the chain rule. Verify the results by expressing z as a function of r, t and computing the partial derivatives directly.

Given $z(x, y) = x^2 + y^2$ where $x(r, t) = r \cos(t)$ and $y(r, t) = r \sin(t)$, determine the value of $\frac{\partial z}{\partial t}$ and $\frac{\partial z}{\partial r}$ using the chain rule. Verify the results by expressing z as a function of r, t and computing the partial derivatives directly.

In general, assume z is a function of n variables, x_1, \ldots, x_n and each of these variables are in turn functions of m variables, t_1, t_2, \ldots, t_m . Then for any variable $t_i, i = 1, 2, \ldots, m$ we have the following,

$$\frac{\partial z}{\partial t_i} = \frac{\partial z}{\partial x_1} \frac{\partial x_1}{\partial t_i} + \frac{\partial z}{\partial x_2} \frac{\partial x_2}{\partial t_i} + \dots + \frac{\partial z}{\partial x_n} \frac{\partial x_n}{\partial t_i}$$

Given $z(x, y) = x^2 + y^2$ where $x(r, t) = r \cos(t)$ and $y(r, t) = r \sin(t)$, determine the value of $\frac{\partial z}{\partial t}$ and $\frac{\partial z}{\partial r}$ using the chain rule. Verify the results by expressing z as a function of r, t and computing the partial derivatives directly.

In general, assume z is a function of n variables, x_1, \ldots, x_n and each of these variables are in turn functions of m variables, t_1, t_2, \ldots, t_m . Then for any variable $t_i, i = 1, 2, \ldots, m$ we have the following,

$$\frac{\partial z}{\partial t_i} = \frac{\partial z}{\partial x_1} \frac{\partial x_1}{\partial t_i} + \frac{\partial z}{\partial x_2} \frac{\partial x_2}{\partial t_i} + \dots + \frac{\partial z}{\partial x_n} \frac{\partial x_n}{\partial t_i}$$

Example

Find the partial derivatives $\frac{\partial z}{\partial t_i}$, i = 1, 2, 3 of the function z(x, y) where $x = t_1 + 2t_2 + 4t_3$ and $y = t_1 - 3t_2 + 5t_3$.

Let $\mathbf{f}: \mathbb{R}^n \to \mathbb{R}^m$ be a vector valued function.

Let $\mathbf{f}: \mathbb{R}^n \to \mathbb{R}^m$ be a vector valued function. Then for a vector $\mathbf{x} = \begin{bmatrix} x_1 & \dots & x_n \end{bmatrix}^T \in \mathbb{R}^n$, the value of the function \mathbf{f} at point \mathbf{x} is a *vector* given as

$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^m$$

Let $\mathbf{f}: \mathbb{R}^n \to \mathbb{R}^m$ be a vector valued function. Then for a vector $\mathbf{x} = \begin{bmatrix} x_1 & \dots & x_n \end{bmatrix}^T \in \mathbb{R}^n$, the value of the function \mathbf{f} at point \mathbf{x} is a *vector* given as

$$\mathbf{f}(\mathbf{x}) = egin{bmatrix} f_1(\mathbf{x}) \ f_2(\mathbf{x}) \ dots \ f_m(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^m$$

Here the functions $f_i : \mathbb{R}^n \to \mathbb{R}$ are real-valued functions.

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

$$\frac{\partial \mathbf{f}}{\partial x_i} = \begin{bmatrix} \frac{\partial f_1}{\partial x_i} \\ \frac{\partial f_2}{\partial x_i} \\ \vdots \\ \frac{\partial f_m}{\partial x_i} \end{bmatrix} \in \mathbb{R}^m.$$

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

$$\frac{\partial \mathbf{f}}{\partial x_i} = \begin{bmatrix} \frac{\partial f_1}{\partial x_i} \\ \frac{\partial f_2}{\partial x_i} \\ \vdots \\ \frac{\partial f_m}{\partial x_i} \end{bmatrix} \in \mathbb{R}^m.$$

Hence the gradient of the vector-valued function $\mathbf{f}:\mathbb{R}^n \to \mathbb{R}^m$ is

$$\nabla \mathbf{f} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \dots & \frac{\partial \mathbf{f}}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

Note that gradient of the vector-valued function can also be represented as follows

$$\nabla \mathbf{f} = \begin{bmatrix} \nabla f_1 \\ \vdots \\ \nabla f_m \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

Image: A matrix and a matrix

Note that gradient of the vector-valued function can also be represented as follows

$$\nabla \mathbf{f} = \begin{bmatrix} \nabla f_1 \\ \vdots \\ \nabla f_m \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

Definition

The Jacobian matrix is the matrix of all first-order partial derivatives of a vector-valued function $\mathbf{f} : \mathbb{R}^n \to \mathbb{R}^m$. The Jacobian matrix \mathbf{J} of \mathbf{f} is an $m \times n$ matrix, usually defined and arranged as follows:

$$\mathbf{J} = \nabla \mathbf{f} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \dots & \frac{\partial \mathbf{f}}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}, \quad \mathbf{J}_{ij} = \frac{\partial f_i}{\partial x_j}.$$



Figure: The dimension of the Jacobian ${\bf J_f}$



Figure: The dimension of the Jacobian ${\bf J_f}$

In particular, the Jacobian of a function $f : \mathbb{R}^n \to \mathbb{R}^1$, which maps a vector $\mathbf{x} \in \mathbb{R}^n$ onto a scalar, is a row vector (matrix of dimension $1 \times n$).



Figure: The dimension of the Jacobian ${\bf J_f}$

In particular, the Jacobian of a function $f : \mathbb{R}^n \to \mathbb{R}^1$, which maps a vector $\mathbf{x} \in \mathbb{R}^n$ onto a scalar, is a row vector (matrix of dimension $1 \times n$).

Example

b

Find the Jacobian of the following functions:

a)
$$f : \mathbb{R}^2 \to \mathbb{R}$$
, given by $f(\mathbf{x}) = x_1 + x_2^3$

)
$$\mathbf{f}: \mathbb{R}^2 \to \mathbb{R}^3$$
, given by $f(\mathbf{x}) = [2x_1, x_1x_2, x_1 + 3x_2]^T$

Consider the vector valued function $\mathbf{f}:\mathbb{R}^2\to\mathbb{R}^2$ given by

$$\mathbf{f}(\mathbf{x}) = \mathbf{f}(x, y) = \begin{bmatrix} 3x \\ -2y \end{bmatrix}$$

•

Consider the vector valued function $\mathbf{f}:\mathbb{R}^2\to\mathbb{R}^2$ given by

$$\mathbf{f}(\mathbf{x}) = \mathbf{f}(x, y) = \begin{bmatrix} 3x \\ -2y \end{bmatrix}$$



V.	Mika	velvar	ı

f

Consider the vector valued function $\mathbf{f}:\mathbb{R}^2\to\mathbb{R}^2$ given by

$$\mathbf{f}(\mathbf{x}) = \mathbf{f}(x, y) = \begin{bmatrix} 3x \\ -2y \end{bmatrix}$$

Note that the image of the $[0,1]^2$ (the blue square) is the rectangle $[0,3] \times [-2,0]$ (depicted in red). The quotient of the areas of the rectangle and the square is 6 and it is equal to $|\det \mathbf{J_f}| = \left|\det \begin{bmatrix} 3 & 0\\ 0 & -2 \end{bmatrix}\right|_{0}$

10 / 22





A nonlinear map $f: \mathbb{R}^2 \to \mathbb{R}^2$ sends a small square (left, in red) to a distorted parallelogram (right, in red). The Jacobian at a point gives the best linear approximation of the distorted parallelogram near that point (right, in translucent white), and the Jacobian determinant gives the ratio of the area of the approximating parallelogram to that of the original square.

Find the gradient of the vector-valued function $\mathbf{f} : \mathbb{R}^2 \to \mathbb{R}^2$, given by $\mathbf{f}(\mathbf{x}) = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$.

Find the gradient of the vector-valued function $\mathbf{f} : \mathbb{R}^2 \to \mathbb{R}^2$, given by $\mathbf{f}(\mathbf{x}) = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$.

Example

Prove that the Jacobian of the vector-valued function $\mathbf{f} : \mathbb{R}^n \to \mathbb{R}^m$, given by $\mathbf{f}(\mathbf{x}) = A\mathbf{x}$, where $A \in \mathbb{R}^{m \times n}$ and $\mathbf{x} \in \mathbb{R}^n$ is

Find the gradient of the vector-valued function $\mathbf{f} : \mathbb{R}^2 \to \mathbb{R}^2$, given by $\mathbf{f}(\mathbf{x}) = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$.

Example

Prove that the Jacobian of the vector-valued function $\mathbf{f} : \mathbb{R}^n \to \mathbb{R}^m$, given by $\mathbf{f}(\mathbf{x}) = A\mathbf{x}$, where $A \in \mathbb{R}^{m \times n}$ and $\mathbf{x} \in \mathbb{R}^n$ is the matrix A, i.e.

 $\mathbf{J_f} = A.$

Chain Rule (Matrix form)

Let $\mathbf{g}:\mathbb{R}^n\to\mathbb{R}^m$ and $\mathbf{f}:\mathbb{R}^m\to\mathbb{R}^k$ are differentiable functions, then

$$\mathbf{J}_{\mathbf{f} \circ \mathbf{g}}(\mathbf{a}) = \mathbf{J}_{\mathbf{f}}(\mathbf{g}(\mathbf{a})) \mathbf{J}_{\mathbf{g}}(\mathbf{a}), \quad \mathbf{a} \in \mathbb{R}^{n}.$$

Chain Rule (Matrix form)

Let $\mathbf{g}:\mathbb{R}^n\to\mathbb{R}^m$ and $\mathbf{f}:\mathbb{R}^m\to\mathbb{R}^k$ are differentiable functions, then

 $\mathbf{J}_{\mathbf{f} \circ \mathbf{g}}(\mathbf{a}) = \mathbf{J}_{\mathbf{f}}(\mathbf{g}(\mathbf{a})) \mathbf{J}_{\mathbf{g}}(\mathbf{a}), \quad \mathbf{a} \in \mathbb{R}^n.$

Note that $\mathbf{J}_{\mathbf{f} \circ \mathbf{g}}(\mathbf{a}) \in \mathbb{R}^{k \times n}, \ \mathbf{J}_{\mathbf{f}}(\mathbf{g}(\mathbf{a})) \in \mathbb{R}^{k \times m}, \text{ and } \mathbf{J}_{\mathbf{g}}(\mathbf{a}) \in \mathbb{R}^{m \times n}.$

Chain Rule (Matrix form)

Let $\mathbf{g}:\mathbb{R}^n\to\mathbb{R}^m$ and $\mathbf{f}:\mathbb{R}^m\to\mathbb{R}^k$ are differentiable functions, then

 $\mathbf{J}_{\mathbf{f} \circ \mathbf{g}}(\mathbf{a}) = \mathbf{J}_{\mathbf{f}}(\mathbf{g}(\mathbf{a})) \mathbf{J}_{\mathbf{g}}(\mathbf{a}), \quad \mathbf{a} \in \mathbb{R}^n.$

Note that $\mathbf{J}_{\mathbf{f} \circ \mathbf{g}}(\mathbf{a}) \in \mathbb{R}^{k \times n}$, $\mathbf{J}_{\mathbf{f}}(\mathbf{g}(\mathbf{a})) \in \mathbb{R}^{k \times m}$, and $\mathbf{J}_{\mathbf{g}}(\mathbf{a}) \in \mathbb{R}^{m \times n}$. Equivalently, if $\mathbf{z} = \mathbf{f}(\mathbf{y})$ and $\mathbf{y} = \mathbf{g}(\mathbf{x})$ then

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

Chain Rule (Matrix form)

Let $\mathbf{g}:\mathbb{R}^n\to\mathbb{R}^m$ and $\mathbf{f}:\mathbb{R}^m\to\mathbb{R}^k$ are differentiable functions, then

$$\mathbf{J}_{\mathbf{f} \circ \mathbf{g}}(\mathbf{a}) = \mathbf{J}_{\mathbf{f}}(\mathbf{g}(\mathbf{a})) \mathbf{J}_{\mathbf{g}}(\mathbf{a}), \quad \mathbf{a} \in \mathbb{R}^{n}.$$

Note that $\mathbf{J}_{\mathbf{f} \circ \mathbf{g}}(\mathbf{a}) \in \mathbb{R}^{k \times n}$, $\mathbf{J}_{\mathbf{f}}(\mathbf{g}(\mathbf{a})) \in \mathbb{R}^{k \times m}$, and $\mathbf{J}_{\mathbf{g}}(\mathbf{a}) \in \mathbb{R}^{m \times n}$. Equivalently, if $\mathbf{z} = \mathbf{f}(\mathbf{y})$ and $\mathbf{y} = \mathbf{g}(\mathbf{x})$ then

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

Example

Find the Jacobians of the functions $f \circ \mathbf{g}$ and $\mathbf{g} \circ f$, where $f : \mathbb{R}^2 \to \mathbb{R}$, given by $f(\mathbf{x}) = x_1 + x_2^2$ and $\mathbf{g} : \mathbb{R} \to \mathbb{R}^2$ given by $\mathbf{g}(t) = \begin{bmatrix} e^t \\ t^2 \end{bmatrix}$.

Example
Using the formula
$$\frac{\partial \mathbf{x}^T B \mathbf{x}}{\partial \mathbf{x}} = \mathbf{x}^T (B + B^T)$$
 deduce that $\frac{\partial \mathbf{x}^T \mathbf{x}}{\partial \mathbf{x}} = 2\mathbf{x}^T$.

	Scalar y	Vector \mathbf{y} (size m)
	Notation Type	Notation Type
Scalar x	$\frac{\partial y}{\partial x}$ scalar	$\frac{\partial \mathbf{y}}{\partial x}$ size- <i>m</i> col. vector
Vector \mathbf{x} (size n)	$rac{\partial y}{\partial \mathbf{x}}$ size- n row vector	$rac{\partial \mathbf{y}}{\partial \mathbf{x}} \ m imes n$ matrix
Matrix $\mathbf X$ (size $p imes q$)	$rac{\partial y}{\partial \mathbf{X}} \; p imes q$ matrix	$rac{\partial \mathbf{y}}{\partial \mathbf{X}} \ m imes (p imes q)$ tensor

æ

	Scalar y	Vector \mathbf{y} (size m)
	Notation Type	Notation Type
Scalar x	$rac{\partial y}{\partial x}$ scalar	$\frac{\partial \mathbf{y}}{\partial x}$ size- m col. vector
Vector \mathbf{x} (size n)	$rac{\partial y}{\partial \mathbf{x}}$ size- n row vector	$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \ m imes n$ matrix
Matrix $\mathbf X$ (size $p imes q$)	$rac{\partial y}{\partial \mathbf{X}} \; p imes q$ matrix	$rac{\partial \mathbf{y}}{\partial \mathbf{X}} \ m imes (p imes q)$ tensor

(Gradient of Scalars with respect to Matrices) Let

$$y = y(\mathbf{X}) = \operatorname{tr}(\mathbf{X}), \,$$
 where $\mathbf{X} \in \mathbb{R}^{p imes p}.$

Find the gradient $\frac{\partial y}{\partial \mathbf{X}}$.

Image: A matrix

э

(Gradient of Vectors with respect to Matrices) Let $\mathbf{v} \in \mathbb{R}^q$ be a fixed vector and $\mathbf{f} : \mathbb{R}^{p \times q} \to \mathbb{R}^p$ be a function given by

$$\mathbf{f}(\mathbf{X}) = \mathbf{X}\mathbf{v}, \ \textit{where} \ \mathbf{X} \in \mathbb{R}^{p imes q}.$$

Find the gradient $\frac{\partial \mathbf{y}}{\partial \mathbf{X}}$ of the function $\mathbf{y} = \mathbf{f}(\mathbf{X})$.

	Matrix (size $m imes k$)
	Notation Type
Scalar x	$rac{\partial \mathbf{Y}}{\partial x} \; m imes k \;$ matrix
Vector \mathbf{x} (size n)	$\frac{\partial \mathbf{Y}}{\partial \mathbf{x}} (m \times k) \times n$ tensor
Matrix $\mathbf X$ (size $p imes q$)	$\frac{\partial \mathbf{Y}}{\partial \mathbf{X}} \ (m \times k) \times (p \times q)$ tensor

2

Matrix (size $m \times k$)	
	Notation Type
Scalar x	$rac{\partial \mathbf{Y}}{\partial x} \; m imes k \;$ matrix
Vector \mathbf{x} (size n)	$\frac{\partial \mathbf{Y}}{\partial \mathbf{x}} (m \times k) \times n$ tensor
Matrix $\mathbf X$ (size $p imes q$)	$\frac{\partial \mathbf{Y}}{\partial \mathbf{X}} \ (m \times k) \times (p \times q)$ tensor

(Gradient of Matrices with respect to Matrices) Let $\mathbf{f} : \mathbb{R}^{p \times q} \to \mathbb{R}^{q \times q}$ be a function given by

$$\mathbf{f}(\mathbf{X}) = \mathbf{X}^T \mathbf{X}, \text{ where } \mathbf{X} \in \mathbb{R}^{p imes q}.$$

Find the gradient
$$\frac{\partial}{\partial \mathbf{X}}$$
 of the function = $\mathbf{f}(\mathbf{X})$.

Image: A matrix and a matrix

э

Useful Identities for Computing Gradients

$$\frac{\partial}{\partial \mathbf{X}} f(\mathbf{X})^{\top} = \left(\frac{\partial f(\mathbf{X})}{\partial \mathbf{X}}\right)^{\top}$$
$$\frac{\partial}{\partial \mathbf{X}} \operatorname{tr}(f(\mathbf{X})) = \operatorname{tr}\left(\frac{\partial f(\mathbf{X})}{\partial \mathbf{X}}\right)$$
$$\frac{\partial}{\partial \mathbf{X}} \det(f(\mathbf{X})) = \det(f(\mathbf{X}))\operatorname{tr}\left(f^{-1}(\mathbf{X})\frac{\partial f(\mathbf{X})}{\partial \mathbf{X}}\right)$$
$$\frac{\partial}{\partial \mathbf{X}} f^{-1}(\mathbf{X}) = -f^{-1}(\mathbf{X})\frac{\partial f(\mathbf{X})}{\partial \mathbf{X}} f^{-1}(\mathbf{X})$$
$$\frac{\partial a^{\top} \mathbf{X}^{-1} \mathbf{b}}{\partial \mathbf{X}} = -(\mathbf{X}^{-1})^{\top} \mathbf{a} \mathbf{b}^{\top} (\mathbf{X}^{-1})^{\top}$$
$$\frac{\partial a^{\mathbf{x}} \mathbf{a}}{\partial \mathbf{x}} = \mathbf{a}^{\top}$$
$$\frac{\partial a^{\mathbf{x}} \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^{\top}$$
$$\frac{\partial a^{\mathbf{x}} \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^{\top}$$
$$\frac{\partial \mathbf{a}^{\mathbf{x}} \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = \mathbf{x}^{\top} (\mathbf{B} + \mathbf{B}^{\top})$$
$$\frac{\partial}{\partial s} (\mathbf{x} - \mathbf{A}s)^{\top} \mathbf{W} (\mathbf{x} - \mathbf{A}s) = -2(\mathbf{x} - \mathbf{A}s)^{\top} \mathbf{W} \mathbf{A} \text{ for symmetric } \mathbf{W}$$

< A[™]

Limits of Multivariable Functions

Let $a \in \mathbb{R}^n$ and $\varepsilon > 0$. Denote $B(a, \varepsilon) = \{x \in \mathbb{R}^n : ||x - a|| < \varepsilon\}.$

Image: A matrix and a matrix

Let $a \in \mathbb{R}^n$ and $\varepsilon > 0$. Denote $B(a, \varepsilon) = \{x \in \mathbb{R}^n : ||x - a|| < \varepsilon\}$.

Definition

Let $f: X \to \mathbb{R}^m$, $X \subset \mathbb{R}^n$, $a \in \mathbb{R}^n$ and $A \in \mathbb{R}^m$. We will say that $\lim_{x \to a} f(x) = A$ if for all $\varepsilon > 0$ there exists $\delta > 0$, such that from $0 < ||x - a||_n < \delta, x \in X$, follows that $||f(x) - A||_m < \varepsilon$. Let $a \in \mathbb{R}^n$ and $\varepsilon > 0$. Denote $B(a, \varepsilon) = \{x \in \mathbb{R}^n : ||x - a|| < \varepsilon\}$.

Definition

Let $f: X \to \mathbb{R}^m$, $X \subset \mathbb{R}^n$, $a \in \mathbb{R}^n$ and $A \in \mathbb{R}^m$. We will say that $\lim_{x \to a} f(x) = A$ if for all $\varepsilon > 0$ there exists $\delta > 0$, such that from $0 < \|x - a\|_n < \delta, x \in X$, follows that $\|f(x) - A\|_m < \varepsilon$.

Definition

Let $f: X \times Y \to \mathbb{R}$, $X, Y \subset \mathbb{R}$, $(x_0, y_0) \in \mathbb{R}^2$ and $A \in \mathbb{R}$. We will say that $\lim_{\substack{x \to x_0 \\ y \to y_0}} f(x) = A$ if for all $\varepsilon > 0$ there exists $\delta > 0$, such that from $|x - x_0| < \delta, |y - y_0| < \delta, (x_0, y_0) \neq (0, 0), x \in X, y \in Y$, follows that $|f(x) - A| < \varepsilon$.

< □ > < □ > < □ > < □ > < □ > < □ >

Theorem

If
$$\lim_{\substack{x \to x_0 \\ y \to y_0}} f(x, y) = A$$
 and $\lim_{x \to x_0} f(x, y) = \varphi(y)$ for all $y \in Y$, $y \neq y_0$, then
$$\lim_{y \to y_0} \varphi(y) = \lim_{y \to y_0} \lim_{x \to x_0} f(x, y) = A.$$

Theorem

 $\underset{y \rightarrow y_{0}}{ If} \lim_{x \rightarrow x_{0}} f\left(x,y\right) = A \text{ and } \lim_{x \rightarrow x_{0}} f\left(x,y\right) = \varphi\left(y\right) \text{ for all } y \in Y, \, y \neq y_{0} \text{, then } x \in Y_{0} \text{, } y \neq y_{0} \text{, then } y \in Y_{0} \text{, } y \neq y_{0} \text{, then } y \in Y_{0} \text{, } y \neq y_{0} \text{, }$

$$\lim_{y \to y_0} \varphi\left(y\right) = \lim_{y \to y_0} \lim_{x \to x_0} f\left(x, y\right) = A.$$

Example

•
$$f(x,y) = x \sin \frac{1}{y}, (x_0,y_0) = (0,0),$$

V.	Mikay	/el	∕an

(日) (四) (日) (日) (日)

3

Theorem

$$\underset{y \to y_{0}}{lf} \lim_{x \to x_{0}} f(x, y) = A \text{ and } \lim_{x \to x_{0}} f(x, y) = \varphi(y) \text{ for all } y \in Y, \ y \neq y_{0}, \text{ then } y \in Y_{0}, \ y \neq y_{0}, \text{ then } y \in Y_{0}, \ y \neq y_{0}, \text{ then } y \in Y_{0}, \ y \neq y_{0}, \text{ then } y \in Y_{0}, \ y \neq y_{0},$$

$$\lim_{y \to y_0} \varphi\left(y\right) = \lim_{y \to y_0} \lim_{x \to x_0} f\left(x, y\right) = A.$$

Example

•
$$f(x,y) = x \sin \frac{1}{y}, (x_0, y_0) = (0,0),$$

• $f(x,y) = \begin{cases} 0, & \text{if } x \neq y, \\ 1, & \text{if } x = y \end{cases}, (x_0, y_0) = (0,0).$

э

Let $f: X \to \mathbb{R}$, $X \subset \mathbb{R}^2$ and (x_0, y_0) is an interior point of X.

Image: A matrix

Let $f: X \to \mathbb{R}$, $X \subset \mathbb{R}^2$ and (x_0, y_0) is an interior point of X.

Definition

f is called differentiable at the point (x_0,y_0) if there exists $A,B\in\mathbb{R}$ such that

$$f(x_{0} + \Delta x, y_{0} + \Delta y) = f(x_{0}, y_{0}) + A\Delta x + B\Delta y + o(\rho), \rho \to 0,$$

where $\rho = \sqrt{\Delta x^2 + \Delta y^2}$.

Let $f: X \to \mathbb{R}$, $X \subset \mathbb{R}^2$ and (x_0, y_0) is an interior point of X.

Definition

f is called differentiable at the point (x_0,y_0) if there exists $A,B\in\mathbb{R}$ such that

$$f(x_{0} + \Delta x, y_{0} + \Delta y) = f(x_{0}, y_{0}) + A\Delta x + B\Delta y + o(\rho), \rho \to 0,$$

where $\rho = \sqrt{\Delta x^2 + \Delta y^2}$.

Theorem

If partial derivatives of the first degree of f are continuous at (x_0, y_0) then it is differentiable at (x_0, y_0) . The inverse is not true.

$$f(x,y) = \begin{cases} x^2 \sin \frac{1}{x}, & \text{if } (x,y) \neq (0,0), \\ 0, & \text{if } (x,y) = (0,0) \end{cases}$$

• • • • • • • •

æ

$$f(x,y) = \begin{cases} x^2 \sin \frac{1}{x}, & \text{if } (x,y) \neq (0,0), \\ 0, & \text{if } (x,y) = (0,0) \end{cases}$$

Definition

$$df(x_0, y_0) = \frac{\partial f}{\partial x}(x_0, y_0) \Delta x + \frac{\partial f}{\partial y}(x_0, y_0) \Delta y \text{ is called differential of } f.$$

イロト イヨト イヨト

æ