

Derivative wrt Scales: $\frac{\partial \vec{x}}{\partial \alpha}$

Kronecker product with $\frac{\partial \vec{x}}{\partial \alpha}$

Derivative wrt vector \vec{x}

Kronecker product with $\frac{\partial \vec{x}}{\partial \alpha}$: $J = \begin{bmatrix} \frac{\partial g_i}{\partial x_1} & \dots & \frac{\partial g_i}{\partial x_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_n}{\partial x_1} & \dots & \frac{\partial g_n}{\partial x_m} \end{bmatrix}$, $\vec{y} = \begin{bmatrix} \alpha_1^m g_1(\alpha) & \dots & \alpha_1^m g_n(\alpha) \\ \vdots & \ddots & \vdots \\ \alpha_m^m g_1(\alpha) & \dots & \alpha_m^m g_n(\alpha) \end{bmatrix}$

gradient vector: $\frac{\partial f(\vec{x})}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \dots & \frac{\partial f_n}{\partial x_m} \end{bmatrix}$

Jacobian matrix: $J = \frac{\partial \vec{y}}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial g_1}{\partial x_1} & \dots & \frac{\partial g_1}{\partial x_m} & \frac{\partial g_1}{\partial x_1} & \dots & \frac{\partial g_1}{\partial x_m} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \frac{\partial g_n}{\partial x_1} & \dots & \frac{\partial g_n}{\partial x_m} & \frac{\partial g_n}{\partial x_1} & \dots & \frac{\partial g_n}{\partial x_m} \end{bmatrix}$

Derivative wrt matrix X : $\frac{\partial \vec{y}}{\partial X} = \begin{bmatrix} \frac{\partial g_1}{\partial X} & \dots & \frac{\partial g_1}{\partial X} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_n}{\partial X} & \dots & \frac{\partial g_n}{\partial X} \end{bmatrix}$ long matrix

Matrix Differentiation Utilities

$\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \vec{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, A = \begin{bmatrix} a_{ij} \dots a_{in} \\ \vdots \\ a_{nj} \dots a_{nn} \end{bmatrix}, B = \begin{bmatrix} b_{ij} \dots b_{in} \\ \vdots \\ b_{nj} \dots b_{nn} \end{bmatrix}$

Variable operations: x_1^m

(a) Let $E = x_1^m$
 $\frac{\partial E}{\partial x_i} = x_1^{m-1} \times x_2^m + x_1^m \times x_2^{m-1} + \dots + x_1^m \times x_n^m = 2x_i^m - 2x_i^m = 0$
 $\Rightarrow \frac{\partial E}{\partial x_i} = 0$

(b) Let $E = x_1^m x_2^m$
 $\frac{\partial E}{\partial x_i} = x_1 x_2^{m-1} + x_1^{m-1} x_2 = x_1^{m-1} x_2 + x_1 x_2^{m-1}$
 $\frac{\partial E}{\partial x_i} = \frac{\partial}{\partial x_i} [x_1 x_2^m] \oplus [x_1^{m-1} x_2]$
 $= [x_1, x_2, \dots, x_n]$

$\text{Hence } \frac{\partial E}{\partial x_i} = j_i$

DERIVATIVES OF SCALAR, VECTOR, TENSOR

scalar	vector	tensor	
scalar	scalar	vector	matrix
vector	gradient tensor	Jacobian	+ big
tensor	tensor	tensor	+ big

dependent → scalar vector matrix

AUTOMATIC DIFFERENTIATION

$f(x) = \ln(x^2 + 1)$

Forward mode:

U	Value	Derivative
$x^2 + 1$	$x^2 + 1$	$2x$
$b(x)$	$\ln(x^2 + 1)$	$\frac{1}{2x + 1} = \frac{1}{x^2 + 1}$

Total derivative: $\frac{\partial y}{\partial x} = \frac{2x}{x^2 + 1}$

Reverse mode:

$$w_j = \sum_{j \in \text{backward}(i)} w_j \cdot \frac{\partial y}{\partial x_i}$$

$f(y) = y \cdot \sin(b(x))$. Find $\frac{\partial f}{\partial x}$, $\frac{\partial f}{\partial y}$ using reverse mode.

$u = x + y$, $v = \sin u$, $f = yv$.

Here, for all $i \neq j$, derivative becomes 0 for f .

$$\frac{\partial f}{\partial x} = \text{diag}\left[\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n}\right] = \text{diag}\left[\frac{\partial f}{\partial x_1} \otimes g(x)\right], \dots, \text{diag}\left[\frac{\partial f}{\partial x_n} \otimes g(x)\right]$$

Say you consider vector addition:
 $f_1(\vec{w}) + g(\vec{z}) = u + v$
 $\frac{\partial y}{\partial \vec{w}} = \text{diag}\left[\frac{\partial}{\partial w_1} (u + v), \dots, \frac{\partial}{\partial w_n} (u + v)\right] = I$

Operator	Derivative
+	I
-	I
\odot	$\text{diag}\left[\frac{\partial}{\partial u_1} \dots \frac{\partial}{\partial u_n}\right] = \vec{x}$
\oslash	$\text{diag}\left[\frac{\partial}{\partial u_1} \dots \frac{\partial}{\partial u_n}\right] = \text{diag}\left[\frac{1}{u_1} \dots \frac{1}{u_n}\right]$
\oslash	$\text{diag}\left[\frac{\partial}{\partial u_1} \dots \frac{\partial}{\partial u_n}\right] = \text{diag}\left[-\frac{1}{u_1^2} \dots -\frac{1}{u_n^2}\right]$

Fraction of scalar-vector multiplication

You can do the same as above by considering the scalar as a vector.

(1) $\vec{f} = \vec{x} + \vec{z} = \frac{f(\vec{x}) + g(\vec{z})}{\vec{x}}$
 $J(\vec{f}) = \text{diag}\left[\dots \frac{\partial}{\partial x_1} (\vec{x} + \vec{z}) \dots\right] = \text{diag}[I_n] = I$ same if you diff wrt z

(2) $\vec{f} = \vec{x} \vec{z} = f(\vec{x}) \odot g(\vec{z})$
 $J(\vec{f}) = \text{diag}\left[\dots \frac{\partial}{\partial x_i} [\vec{x} \vec{z}] \dots\right] = \text{diag}[\dots z_i \dots] = IZ$
 $\frac{\partial}{\partial z} (\vec{x} \vec{z}) = \text{diag}\left[\dots \frac{\partial}{\partial z_i} [\vec{x} \vec{z}] \dots\right] = \text{diag}[\dots x_i \dots] = \vec{x}$

Scalar-vector operations: Jacobian

(1) $x + \vec{y} = \vec{x} + \vec{y}$ scalar vector
 $J_y = \text{diag}\left(\frac{\partial}{\partial y_1}, \dots, \frac{\partial}{\partial y_n}\right) = \text{diag}(I) = I_n$ (scalar matrix)

(2) $x \vec{y} = \vec{x} \vec{y}$
 $J(\vec{y}) = \text{diag}\left(\frac{\partial}{\partial y_1}, \dots, \frac{\partial}{\partial y_n}\right) = \text{diag}(y) = I_n$ (scalar matrix)
 $J(x) = \text{diag}\left(\frac{\partial}{\partial x}, \dots, \frac{\partial}{\partial x}\right) = \text{diag}(y_i) = \vec{y}$

Dataflow: $\vec{x} \rightarrow \vec{g}(\vec{x}) \xrightarrow{f} \vec{f}(\vec{g}(\vec{x}))$

Chain rule is extended for vectors using Jacobians

Vector Chain Rule

Let $\vec{x}, \vec{f}, \vec{g}$ such that $\vec{x} \rightarrow \vec{f} \rightarrow \vec{g}$ and $|x| \cdot n \cdot (f \rightarrow m \cdot (g \rightarrow k))$

$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$, $\frac{\partial g}{\partial f} = \begin{bmatrix} \frac{\partial g_1}{\partial f_1} & \dots & \frac{\partial g_1}{\partial f_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_k}{\partial f_1} & \dots & \frac{\partial g_k}{\partial f_m} \end{bmatrix}$

$\frac{\partial g}{\partial x} = \frac{\partial g}{\partial f} \cdot \frac{\partial f}{\partial x}$

$x \in \mathbb{R}^n$, $\vec{f} \in \mathbb{R}^k$, $f_i : \mathbb{R}^k \rightarrow \mathbb{R}^{k_i}$, $f : \mathbb{R}^k \rightarrow \mathbb{R}^m$, $f_i : \mathbb{R}^m \rightarrow \mathbb{R}^{m_i}$

(1) $A \vec{x} = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix}$

(2) Let scalar $\alpha \in \mathbb{R}$, $a = \alpha A \vec{x}$, where $\alpha \in \mathbb{R}$, $\vec{x} \in \mathbb{R}^n$, $a \in \mathbb{R}^m$.
 $\frac{\partial a}{\partial x_i} = a_i \vec{x}_i \Rightarrow \frac{\partial a}{\partial x_i} = \alpha \vec{x}_i \Rightarrow \alpha \vec{x}_i = \frac{\partial a}{\partial x_i}$

(3) $\frac{\partial a}{\partial x} = A \vec{x}$

(4) $y = A \vec{x}$, $x \in \mathbb{R}^n$, $A \in \mathbb{R}^{m \times n}$.
 $\frac{\partial y}{\partial x} = A$

(5) $y = A \vec{x}$, $x \in \mathbb{R}^n$, $A \in \mathbb{R}^{m \times n}$.
 $\frac{\partial y}{\partial x} = A$

Matrix Differentiation Utilities

$\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \vec{g} = \begin{bmatrix} g_1 \\ \vdots \\ g_n \end{bmatrix}, A = \begin{bmatrix} a_{ij} \dots a_{in} \\ \vdots \\ a_{nj} \dots a_{nn} \end{bmatrix}, B = \begin{bmatrix} b_{ij} \dots b_{in} \\ \vdots \\ b_{nj} \dots b_{nn} \end{bmatrix}$

(1) Let $E = x_1^m$
 $\frac{\partial E}{\partial x_i} = x_1^{m-1} \times x_2^m + x_1^m \times x_2^{m-1} + \dots + x_1^m \times x_n^m = 2x_i^m - 2x_i^m = 0$
 $\Rightarrow \frac{\partial E}{\partial x_i} = 0$

(2) Let $E = x_1^m x_2^m$
 $\frac{\partial E}{\partial x_i} = x_1 x_2^{m-1} + x_1^{m-1} x_2 = x_1^{m-1} x_2 + x_1 x_2^{m-1}$
 $\frac{\partial E}{\partial x_i} = \frac{\partial}{\partial x_i} [x_1 x_2^m] \oplus [x_1^{m-1} x_2]$
 $= [x_1, x_2, \dots, x_n]$

(3) $\frac{\partial E}{\partial x_i} = 2x_i^m$

(4) $\frac{\partial E}{\partial x_i} = 2(x_1^m) \vec{x}$
 $\frac{\partial E}{\partial x_i} = 2x_1^m \frac{\partial \vec{x}}{\partial x_i}$

MATRIX DIFFERENTIATION UTILITIES

(1) $y = mx_1$, $x \in \mathbb{R}^n$, independent

(2) $y = \vec{x}^T A \vec{x}$, $\vec{x} \in \mathbb{R}^n$

(3) Trace of scalar α
 $\alpha = [x_1]$
 $\frac{\partial \alpha}{\partial x_i} = \delta_{i1}$

(4) Frobenius norm
 $\text{Frob}(AB) = \sqrt{\text{tr}(A^T B)}$

(5) Cyclic permutation
 $\frac{\partial \text{tr}(AB)}{\partial x} = \vec{x}^T [D(A)] + (A^T) [D(\vec{x})]$

(6) Trace of scalar product
 $\langle x, y \rangle = \vec{x}^T \vec{y} = \sum_{i=1}^n x_i y_i$

Trace → sum of diagonal elements $\sum_{i=1}^n y_i$:

(1) $\text{tr}(A+B) = \text{tr}(A) + \text{tr}(B)$

(2) $\text{tr}(CA) = C \text{tr}(A)$

(3) $\text{tr}(AB) = \text{tr}(BA)$

(4) $\text{tr}(A) = \text{tr}(A^T)$

(5) $\text{tr}(x^T Ax) = \sum_{i=1}^n x_i^2$

(6) Frobenius norm
 $\|\vec{x}\|_F^2 = \text{tr}(x^T x) = \sum_{i=1}^n x_i^2$

(7) $d(D(V^T V)) = V^T D(V) + V^T D(V)$

(8) $\frac{\partial D(x)}{\partial x} = \frac{\partial x}{\partial x} = I$

(9) $D(U^T V) = U^T D(V) + V^T D(U)$

(10) $D(A^T A) = A^T D(A) + D(A)A$

(11) $D(x^T A) = x^T D(A) + A^T D(x)$

(12) $D(A^T A) = \text{tr}(A^T A) = \text{tr}(A)^2$

(13) $\text{tr}(A^2) = \text{tr}(AA)$

(14) $\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(ACB)$

(15) $\text{tr}(A^T B^T) = \text{tr}(B^T A^T) = \text{tr}(B^T A)$

(16) $\text{tr}(AB) = \text{tr}(BA)$

(17) $d(D(x)) = \vec{x}^T \frac{\partial D(x)}{\partial x} = \vec{x}^T \frac{\partial x}{\partial x} = \vec{x}^T I = \vec{x}^T$

(18) $\text{tr}(A^T A) = \text{tr}(A^T A)$

(19) $\text{tr}(A^T A) = \text{tr}(A^T A)$

(7) non-singular matrix whose elements = $f(x)$
 $\frac{\partial A^{-1}}{\partial x} = -A^{-1} \frac{\partial A}{\partial x} A^{-1}$

Proof:
 $AA^{-1} = I$ invertible
 $A^T \frac{\partial A}{\partial x} + \frac{\partial A^T}{\partial x} A = 0$
 $\frac{\partial A^T}{\partial x} A = -A^T \frac{\partial A}{\partial x}$
 $\text{Multiply } A^T \text{ on both sides}$
 $\frac{\partial A^T}{\partial x} = -A^T \frac{\partial A}{\partial x} A$

Matrix Differentiation Utilities

(1) Define scalar $\alpha \in \mathbb{R}$, $a = \alpha x$, $\vec{x} \in \mathbb{R}^n$, $a \in \mathbb{R}^m$.
 $\frac{\partial a}{\partial x} = \alpha \vec{x}$

(2) Define scalar $\alpha \in \mathbb{R}$, $a = \alpha x$, $\vec{x} \in \mathbb{R}^n$, $a \in \mathbb{R}^m$.
 $\frac{\partial a}{\partial x} = \alpha \vec{x}$

(3) $\frac{\partial \alpha}{\partial x} = \vec{V} \alpha$

(4) $\frac{\partial \alpha}{\partial x} = \alpha \vec{V}$

(5) $\frac{\partial \alpha}{\partial x} = \vec{V} \alpha$

(6) $\frac{\partial \alpha}{\partial x} = \vec{V} \alpha$

(7) $\frac{\partial \alpha}{\partial x} = \vec{V} \alpha$