# Appendix D

# Matrix calculus

*From too much study, and from extreme passion, cometh madnesse.*

*—*Isaac Newton [86, §5]

## D.1 Directional derivative, Taylor series

### D.1.1 Gradients

*Gradient* of a differentiable real function $f(x) : \mathbb{R}^K \to \mathbb{R}$ with respect to its vector domain is defined

$$\nabla f(x) \triangleq \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_K} \end{bmatrix} \in \mathbb{R}^K \tag{1354}$$

while the second-order gradient of the twice differentiable real function with respect to its vector domain is traditionally called the *Hessian*;

$$\nabla^2 f(x) \triangleq \begin{bmatrix} \frac{\partial^2 f(x)}{\partial^2 x_1} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_K} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial^2 x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_K} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_K \partial x_1} & \frac{\partial^2 f(x)}{\partial x_K \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial^2 x_K} \end{bmatrix} \in \mathbb{S}^K \tag{1355}$$

The gradient of vector-valued function $v(x) : \mathbb{R} \to \mathbb{R}^N$ on real domain is a row-vector

$$\nabla v(x) \triangleq \left[ \begin{array}{cccc} \frac{\partial v_1(x)}{\partial x} & \frac{\partial v_2(x)}{\partial x} & \cdots & \frac{\partial v_N(x)}{\partial x} \end{array} \right] \in \mathbb{R}^N \qquad (1356)$$

while the second-order gradient is

$$\nabla^2 v(x) \triangleq \left[ \begin{array}{cccc} \frac{\partial^2 v_1(x)}{\partial x^2} & \frac{\partial^2 v_2(x)}{\partial x^2} & \cdots & \frac{\partial^2 v_N(x)}{\partial x^2} \end{array} \right] \in \mathbb{R}^N \qquad (1357)$$

Gradient of vector-valued function $h(x) : \mathbb{R}^K \to \mathbb{R}^N$ on vector domain is

$$\nabla h(x) \triangleq \left[ \begin{array}{cccc} \frac{\partial h_1(x)}{\partial x_1} & \frac{\partial h_2(x)}{\partial x_1} & \cdots & \frac{\partial h_N(x)}{\partial x_1} \\ \frac{\partial h_1(x)}{\partial x_2} & \frac{\partial h_2(x)}{\partial x_2} & \cdots & \frac{\partial h_N(x)}{\partial x_2} \\ \vdots & \vdots & & \vdots \\ \frac{\partial h_1(x)}{\partial x_K} & \frac{\partial h_2(x)}{\partial x_K} & \cdots & \frac{\partial h_N(x)}{\partial x_K} \end{array} \right] \qquad (1358)$$

$$= \left[ \nabla h_1(x) \ \ \nabla h_2(x) \ \cdots \ \nabla h_N(x) \right] \in \mathbb{R}^{K \times N}$$

while the second-order gradient has a three-dimensional representation dubbed *cubix* ;[D.1]

$$\nabla^2 h(x) \triangleq \left[ \begin{array}{cccc} \nabla \frac{\partial h_1(x)}{\partial x_1} & \nabla \frac{\partial h_2(x)}{\partial x_1} & \cdots & \nabla \frac{\partial h_N(x)}{\partial x_1} \\ \nabla \frac{\partial h_1(x)}{\partial x_2} & \nabla \frac{\partial h_2(x)}{\partial x_2} & \cdots & \nabla \frac{\partial h_N(x)}{\partial x_2} \\ \vdots & \vdots & & \vdots \\ \nabla \frac{\partial h_1(x)}{\partial x_K} & \nabla \frac{\partial h_2(x)}{\partial x_K} & \cdots & \nabla \frac{\partial h_N(x)}{\partial x_K} \end{array} \right] \qquad (1359)$$

$$= \left[ \nabla^2 h_1(x) \ \ \nabla^2 h_2(x) \ \cdots \ \nabla^2 h_N(x) \right] \in \mathbb{R}^{K \times N \times K}$$

where the gradient of each real entry is with respect to vector $x$ as in (1354).

---

[D.1]The word *matrix* comes from the Latin for *womb*; related to the prefix *matri-* derived from *mater* meaning *mother*.

The gradient of real function $g(X) : \mathbb{R}^{K \times L} \to \mathbb{R}$ on matrix domain is

$$
\nabla g(X) \triangleq
\begin{bmatrix}
\frac{\partial g(X)}{\partial X_{11}} & \frac{\partial g(X)}{\partial X_{12}} & \cdots & \frac{\partial g(X)}{\partial X_{1L}} \\
\frac{\partial g(X)}{\partial X_{21}} & \frac{\partial g(X)}{\partial X_{22}} & \cdots & \frac{\partial g(X)}{\partial X_{2L}} \\
\vdots & \vdots & & \vdots \\
\frac{\partial g(X)}{\partial X_{K1}} & \frac{\partial g(X)}{\partial X_{K2}} & \cdots & \frac{\partial g(X)}{\partial X_{KL}}
\end{bmatrix}
\in \mathbb{R}^{K \times L}
$$

$$
=
\begin{bmatrix}
\nabla_{X(:,1)}\, g(X) & & & \\
& \nabla_{X(:,2)}\, g(X) & & \\
& & \ddots & \\
& & & \nabla_{X(:,L)}\, g(X)
\end{bmatrix}
\in \mathbb{R}^{K \times 1 \times L}
$$

$$(1360)$$

where the gradient $\nabla_{X(:,i)}$ is with respect to the $i^{\text{th}}$ column of $X$. The strange appearance of (1360) in $\mathbb{R}^{K \times 1 \times L}$ is meant to suggest a third dimension perpendicular to the page (not a diagonal matrix). The second-order gradient has representation

$$
\nabla^2 g(X) \triangleq
\begin{bmatrix}
\nabla \frac{\partial g(X)}{\partial X_{11}} & \nabla \frac{\partial g(X)}{\partial X_{12}} & \cdots & \nabla \frac{\partial g(X)}{\partial X_{1L}} \\
\nabla \frac{\partial g(X)}{\partial X_{21}} & \nabla \frac{\partial g(X)}{\partial X_{22}} & \cdots & \nabla \frac{\partial g(X)}{\partial X_{2L}} \\
\vdots & \vdots & & \vdots \\
\nabla \frac{\partial g(X)}{\partial X_{K1}} & \nabla \frac{\partial g(X)}{\partial X_{K2}} & \cdots & \nabla \frac{\partial g(X)}{\partial X_{KL}}
\end{bmatrix}
\in \mathbb{R}^{K \times L \times K \times L}
$$

$$
=
\begin{bmatrix}
\nabla \nabla_{X(:,1)}\, g(X) & & & \\
& \nabla \nabla_{X(:,2)}\, g(X) & & \\
& & \ddots & \\
& & & \nabla \nabla_{X(:,L)}\, g(X)
\end{bmatrix}
\in \mathbb{R}^{K \times 1 \times L \times K \times L}
$$

$$(1361)$$

where the gradient $\nabla$ is with respect to matrix $X$.

Gradient of vector-valued function $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}^N$ on matrix domain is a cubix

$$\nabla g(X) \triangleq \begin{array}{c} \left[ \nabla_{X(:,1)}\, g_1(X) \;\; \nabla_{X(:,1)}\, g_2(X) \;\; \cdots \;\; \nabla_{X(:,1)}\, g_N(X) \right. \\ \nabla_{X(:,2)}\, g_1(X) \;\; \nabla_{X(:,2)}\, g_2(X) \;\; \cdots \;\; \nabla_{X(:,2)}\, g_N(X) \\ \ddots \qquad\qquad \ddots \qquad\qquad \ddots \\ \left. \nabla_{X(:,L)}\, g_1(X) \;\; \nabla_{X(:,L)}\, g_2(X) \;\; \cdots \;\; \nabla_{X(:,L)}\, g_N(X) \right] \end{array}$$

$$= \left[ \nabla g_1(X) \;\; \nabla g_2(X) \;\; \cdots \;\; \nabla g_N(X) \right] \in \mathbb{R}^{K \times N \times L} \qquad (1362)$$

while the second-order gradient has a five-dimensional representation;

$$\nabla^2 g(X) \triangleq \begin{array}{c} \left[ \nabla\nabla_{X(:,1)}\, g_1(X) \;\; \nabla\nabla_{X(:,1)}\, g_2(X) \;\; \cdots \;\; \nabla\nabla_{X(:,1)}\, g_N(X) \right. \\ \nabla\nabla_{X(:,2)}\, g_1(X) \;\; \nabla\nabla_{X(:,2)}\, g_2(X) \;\; \cdots \;\; \nabla\nabla_{X(:,2)}\, g_N(X) \\ \ddots \qquad\qquad \ddots \qquad\qquad \ddots \\ \left. \nabla\nabla_{X(:,L)}\, g_1(X) \;\; \nabla\nabla_{X(:,L)}\, g_2(X) \;\; \cdots \;\; \nabla\nabla_{X(:,L)}\, g_N(X) \right] \end{array}$$

$$= \left[ \nabla^2 g_1(X) \;\; \nabla^2 g_2(X) \;\; \cdots \;\; \nabla^2 g_N(X) \right] \in \mathbb{R}^{K \times N \times L \times K \times L} \qquad (1363)$$

The gradient of matrix-valued function $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}^{M \times N}$ on matrix domain has a four-dimensional representation called *quartix*

$$\nabla g(X) \triangleq \begin{bmatrix} \nabla g_{11}(X) & \nabla g_{12}(X) & \cdots & \nabla g_{1N}(X) \\ \nabla g_{21}(X) & \nabla g_{22}(X) & \cdots & \nabla g_{2N}(X) \\ \vdots & \vdots & & \vdots \\ \nabla g_{M1}(X) & \nabla g_{M2}(X) & \cdots & \nabla g_{MN}(X) \end{bmatrix} \in \mathbb{R}^{M \times N \times K \times L} \quad (1364)$$

while the second-order gradient has six-dimensional representation

$$\nabla^2 g(X) \triangleq \begin{bmatrix} \nabla^2 g_{11}(X) & \nabla^2 g_{12}(X) & \cdots & \nabla^2 g_{1N}(X) \\ \nabla^2 g_{21}(X) & \nabla^2 g_{22}(X) & \cdots & \nabla^2 g_{2N}(X) \\ \vdots & \vdots & & \vdots \\ \nabla^2 g_{M1}(X) & \nabla^2 g_{M2}(X) & \cdots & \nabla^2 g_{MN}(X) \end{bmatrix} \in \mathbb{R}^{M \times N \times K \times L \times K \times L}$$

$$(1365)$$

and so on.

## D.1.2 Product rules for matrix-functions

Given dimensionally compatible matrix-valued functions of matrix variable $f(X)$ and $g(X)$

$$\nabla_X \big(f(X)^T g(X)\big) = \nabla_X(f)\, g \,+\, \nabla_X(g)\, f \tag{1366}$$

while [35, §8.3] [205]

$$\nabla_X \text{tr}\big(f(X)^T g(X)\big) = \nabla_X \Big(\text{tr}\big(f(X)^T g(Z)\big) \,+\, \text{tr}\big(g(X)\, f(Z)^T\big)\Big)\Big|_{Z \leftarrow X} \tag{1367}$$

These expressions implicitly apply as well to scalar-, vector-, or matrix-valued functions of scalar, vector, or matrix arguments.

**D.1.2.0.1 Example.** *Cubix.*
Suppose $f(X) : \mathbb{R}^{2\times 2} \to \mathbb{R}^2 = X^T a$ and $g(X) : \mathbb{R}^{2\times 2} \to \mathbb{R}^2 = Xb$. We wish to find

$$\nabla_X\big(f(X)^T g(X)\big) = \nabla_X\, a^T X^2 b \tag{1368}$$

using the product rule. Formula (1366) calls for

$$\nabla_X\, a^T X^2 b = \nabla_X(X^T a)\, Xb \,+\, \nabla_X(Xb)\, X^T a \tag{1369}$$

Consider the first of the two terms:

$$\begin{aligned}\nabla_X(f)\, g \,&=\, \nabla_X(X^T a)\, Xb \\ &=\, \big[\, \nabla(X^T a)_1 \quad \nabla(X^T a)_2 \,\big]\, Xb \end{aligned} \tag{1370}$$

The gradient of $X^T a$ forms a cubix in $\mathbb{R}^{2\times 2\times 2}$.

$$\nabla_X(X^T a)\, Xb = \begin{bmatrix} \dfrac{\partial(X^T a)_1}{\partial X_{11}} & \cdots\cdots & \dfrac{\partial(X^T a)_2}{\partial X_{11}} \\ & \dfrac{\partial(X^T a)_1}{\partial X_{12}} & \cdots\cdots & \dfrac{\partial(X^T a)_2}{\partial X_{12}} \\ \dfrac{\partial(X^T a)_1}{\partial X_{21}} & & \dfrac{\partial(X^T a)_2}{\partial X_{21}} \\ & \dfrac{\partial(X^T a)_1}{\partial X_{22}} & \cdots\cdots & \dfrac{\partial(X^T a)_2}{\partial X_{22}} \end{bmatrix} \begin{bmatrix} (Xb)_1 \\ \\ (Xb)_2 \end{bmatrix} \in \mathbb{R}^{2\times 1\times 2} \tag{1371}$$

Because gradient of the product (1368) requires total change with respect to change in each entry of matrix $X$, the $Xb$ vector must make an inner product with each vector in the second dimension of the cubix (indicated by dotted line segments);

$$
\nabla_X(X^T a)\, Xb = \begin{bmatrix} a_1 & & 0 & \\ & 0 & & a_1 \\ a_2 & & 0 & \\ & 0 & & a_2 \end{bmatrix} \begin{bmatrix} b_1 X_{11} + b_2 X_{12} \\ b_1 X_{21} + b_2 X_{22} \end{bmatrix} \in \mathbb{R}^{\mathbf{2 \times 1 \times 2}}
$$

$$
= \begin{bmatrix} a_1(b_1 X_{11} + b_2 X_{12}) & a_1(b_1 X_{21} + b_2 X_{22}) \\ a_2(b_1 X_{11} + b_2 X_{12}) & a_2(b_1 X_{21} + b_2 X_{22}) \end{bmatrix} \in \mathbb{R}^{\mathbf{2 \times 2}}
$$

$$
= ab^T X^T
$$

(1372)

where the cubix appears as a complete $2 \times 2 \times 2$ matrix. In like manner for the second term $\nabla_X(g)\, f$

$$
\nabla_X(Xb)\, X^T a = \begin{bmatrix} b_1 & & 0 & \\ & b_2 & & 0 \\ 0 & & b_1 & \\ & 0 & & b_2 \end{bmatrix} \begin{bmatrix} X_{11} a_1 + X_{21} a_2 \\ X_{12} a_1 + X_{22} a_2 \end{bmatrix} \in \mathbb{R}^{\mathbf{2 \times 1 \times 2}}
$$

$$
= X^T ab^T \in \mathbb{R}^{\mathbf{2 \times 2}}
$$

(1373)

The solution

$$
\nabla_X\, a^T X^2 b = ab^T X^T + X^T ab^T
$$

(1374)

can be found from Table **D.2.1** or verified using (1367).  □

### D.1.2.1   Kronecker product

A partial remedy for venturing into *hyperdimensional* representations, such as the cubix or quartix, is to first vectorize matrices as in (29). This device gives rise to the Kronecker product of matrices $\otimes$ ; `a.k.a`, *direct product* or *tensor product*. Although it sees reversal in the literature, [211, §2.1] we adopt the definition: for $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{p \times q}$

$$
B \otimes A \triangleq \begin{bmatrix} B_{11}A & B_{12}A & \cdots & B_{1q}A \\ B_{21}A & B_{22}A & \cdots & B_{2q}A \\ \vdots & \vdots & & \vdots \\ B_{p1}A & B_{p2}A & \cdots & B_{pq}A \end{bmatrix} \in \mathbb{R}^{pm \times qn}
$$

(1375)

One advantage to vectorization is existence of a traditional two-dimensional matrix representation for the second-order gradient of a real function with respect to a vectorized matrix. For example, from §A.1.1 no.22 (§D.2.1) for square $A\,,B\in\mathbb{R}^{n\times n}$ [96, §5.2] [10, §3]

$$\nabla^2_{\text{vec}\,X}\,\text{tr}(AXBX^T) \;=\; \nabla^2_{\text{vec}\,X}\,\text{vec}(X)^T(B^T\otimes A)\,\text{vec}\,X \;=\; B\otimes A^T + B^T\otimes A \in \mathbb{R}^{n^2\times n^2} \tag{1376}$$

To disadvantage is a large new but known set of algebraic rules and the fact that its mere use does not generally guarantee two-dimensional matrix representation of gradients.

## D.1.3 Chain rules for composite matrix-functions

Given dimensionally compatible matrix-valued functions of matrix variable $f(X)$ and $g(X)$ [137, §15.7]

$$\nabla_X\,g\big(f(X)^T\big) = \nabla_X f^T\,\nabla_f\,g \tag{1377}$$

$$\nabla^2_X\,g\big(f(X)^T\big) \;=\; \nabla_X\big(\nabla_X f^T\,\nabla_f\,g\big) \;=\; \nabla^2_X f\,\nabla_f\,g \;+\; \nabla_X f^T\,\nabla^2_f g\,\nabla_X f \tag{1378}$$

### D.1.3.1 Two arguments

$$\nabla_X\,g\big(f(X)^T,\,h(X)^T\big) = \nabla_X f^T\,\nabla_f\,g \;+\; \nabla_X h^T\,\nabla_h\,g \tag{1379}$$

**D.1.3.1.1 Example.** *Chain rule for two arguments.* [28, §1.1]

$$g\big(f(x)^T,\,h(x)^T\big) = (f(x) + h(x))^T A\,(f(x) + h(x)) \tag{1380}$$

$$f(x) = \begin{bmatrix} x_1 \\ \varepsilon x_2 \end{bmatrix}, \qquad h(x) = \begin{bmatrix} \varepsilon x_1 \\ x_2 \end{bmatrix} \tag{1381}$$

$$\nabla_x\,g\big(f(x)^T,\,h(x)^T\big) = \begin{bmatrix} 1 & 0 \\ 0 & \varepsilon \end{bmatrix}(A + A^T)(f + h) + \begin{bmatrix} \varepsilon & 0 \\ 0 & 1 \end{bmatrix}(A + A^T)(f + h) \tag{1382}$$

$$\nabla_x\,g\big(f(x)^T,\,h(x)^T\big) = \begin{bmatrix} 1+\varepsilon & 0 \\ 0 & 1+\varepsilon \end{bmatrix}(A + A^T)\left(\begin{bmatrix} x_1 \\ \varepsilon x_2 \end{bmatrix} + \begin{bmatrix} \varepsilon x_1 \\ x_2 \end{bmatrix}\right) \tag{1383}$$

$$\lim_{\varepsilon \to 0} \nabla_x g\big(f(x)^T, \, h(x)^T\big) = (A + A^T)x \qquad (1384)$$

from Table **D.2.1**.                                                                    □

These formulae remain correct when the gradients produce hyperdimensional representations:

## D.1.4  First directional derivative

Assume that a differentiable function $g(X) : \mathbb{R}^{K \times L} \to \mathbb{R}^{M \times N}$ has continuous first- and second-order gradients $\nabla g$ and $\nabla^2 g$ over $\mathrm{dom}\, g$ which is an open set. We seek simple expressions for the first and second directional derivatives in direction $Y \in \mathbb{R}^{K \times L}$, $\overset{\to Y}{dg} \in \mathbb{R}^{M \times N}$ and $\overset{\to Y}{dg^2} \in \mathbb{R}^{M \times N}$ respectively.

Assuming that the limit exists, we may state the partial derivative of the $mn^{\text{th}}$ entry of $g$ with respect to the $kl^{\text{th}}$ entry of $X$ ;

$$\frac{\partial g_{mn}(X)}{\partial X_{kl}} = \lim_{\Delta t \to 0} \frac{g_{mn}(X + \Delta t \, e_k e_l^T) - g_{mn}(X)}{\Delta t} \in \mathbb{R} \qquad (1385)$$

where $e_k$ is the $k^{\text{th}}$ standard basis vector in $\mathbb{R}^K$ while $e_l$ is the $l^{\text{th}}$ standard basis vector in $\mathbb{R}^L$. The total number of partial derivatives equals $KLMN$ while the gradient is defined in their terms; the $mn^{\text{th}}$ entry of the gradient is

$$\nabla g_{mn}(X) = \begin{bmatrix} \frac{\partial g_{mn}(X)}{\partial X_{11}} & \frac{\partial g_{mn}(X)}{\partial X_{12}} & \cdots & \frac{\partial g_{mn}(X)}{\partial X_{1L}} \\ \frac{\partial g_{mn}(X)}{\partial X_{21}} & \frac{\partial g_{mn}(X)}{\partial X_{22}} & \cdots & \frac{\partial g_{mn}(X)}{\partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial g_{mn}(X)}{\partial X_{K1}} & \frac{\partial g_{mn}(X)}{\partial X_{K2}} & \cdots & \frac{\partial g_{mn}(X)}{\partial X_{KL}} \end{bmatrix} \in \mathbb{R}^{K \times L} \qquad (1386)$$

while the gradient is a quartix

$$\nabla g(X) = \begin{bmatrix} \nabla g_{11}(X) & \nabla g_{12}(X) & \cdots & \nabla g_{1N}(X) \\ \nabla g_{21}(X) & \nabla g_{22}(X) & \cdots & \nabla g_{2N}(X) \\ \vdots & \vdots & & \vdots \\ \nabla g_{M1}(X) & \nabla g_{M2}(X) & \cdots & \nabla g_{MN}(X) \end{bmatrix} \in \mathbb{R}^{M \times N \times K \times L} \qquad (1387)$$

By simply rotating our perspective of the four-dimensional representation of the gradient matrix, we find one of three useful transpositions of this quartix (connoted $T_1$):

$$\nabla g(X)^{T_1} = \begin{bmatrix} \frac{\partial g(X)}{\partial X_{11}} & \frac{\partial g(X)}{\partial X_{12}} & \cdots & \frac{\partial g(X)}{\partial X_{1L}} \\ \frac{\partial g(X)}{\partial X_{21}} & \frac{\partial g(X)}{\partial X_{22}} & \cdots & \frac{\partial g(X)}{\partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial g(X)}{\partial X_{K1}} & \frac{\partial g(X)}{\partial X_{K2}} & \cdots & \frac{\partial g(X)}{\partial X_{KL}} \end{bmatrix} \in \mathbb{R}^{K \times L \times M \times N} \qquad (1388)$$

When the limit for $\Delta t \in \mathbb{R}$ exists, it is easy to show by substitution of variables in (1385)

$$\frac{\partial g_{mn}(X)}{\partial X_{kl}} Y_{kl} = \lim_{\Delta t \to 0} \frac{g_{mn}(X + \Delta t\, Y_{kl}\, e_k e_l^T) - g_{mn}(X)}{\Delta t} \in \mathbb{R} \qquad (1389)$$

which may be interpreted as the change in $g_{mn}$ at $X$ when the change in $X_{kl}$ is equal to $Y_{kl}$, the $kl^{\text{th}}$ entry of any $Y \in \mathbb{R}^{K \times L}$. Because the total change in $g_{mn}(X)$ due to $Y$ is the sum of change with respect to each and every $X_{kl}$, the $mn^{\text{th}}$ entry of the directional derivative is the corresponding total differential [137, §15.8]

$$dg_{mn}(X)|_{dX \to Y} = \sum_{k,l} \frac{\partial g_{mn}(X)}{\partial X_{kl}} Y_{kl} = \text{tr}\left(\nabla g_{mn}(X)^T Y\right) \qquad (1390)$$

$$= \sum_{k,l} \lim_{\Delta t \to 0} \frac{g_{mn}(X + \Delta t\, Y_{kl}\, e_k e_l^T) - g_{mn}(X)}{\Delta t} \qquad (1391)$$

$$= \lim_{\Delta t \to 0} \frac{g_{mn}(X + \Delta t\, Y) - g_{mn}(X)}{\Delta t} \qquad (1392)$$

$$= \frac{d}{dt}\bigg|_{t=0} g_{mn}(X + t\, Y) \qquad (1393)$$

where $t \in \mathbb{R}$. Assuming finite $Y$, equation (1392) is called the *Gâteaux differential* [27, App.A.5] [125, §D.2.1] [234, §5.28] whose existence is implied by the existence of the *Fréchet differential*, the sum in (1390). [157, §7.2] Each may be understood as the change in $g_{mn}$ at $X$ when the change in $X$ is equal

in magnitude and direction to $Y$ .[D.2]  Hence the directional derivative,

$$
\overset{\rightarrow Y}{dg}(X) \triangleq
\left.\begin{bmatrix}
dg_{11}(X) & dg_{12}(X) & \cdots & dg_{1N}(X) \\
dg_{21}(X) & dg_{22}(X) & \cdots & dg_{2N}(X) \\
\vdots & \vdots & & \vdots \\
dg_{M1}(X) & dg_{M2}(X) & \cdots & dg_{MN}(X)
\end{bmatrix}\right|_{dX \to Y}
\in \mathbb{R}^{M \times N}
$$

$$
=
\begin{bmatrix}
\mathrm{tr}\big(\nabla g_{11}(X)^T Y\big) & \mathrm{tr}\big(\nabla g_{12}(X)^T Y\big) & \cdots & \mathrm{tr}\big(\nabla g_{1N}(X)^T Y\big) \\
\mathrm{tr}\big(\nabla g_{21}(X)^T Y\big) & \mathrm{tr}\big(\nabla g_{22}(X)^T Y\big) & \cdots & \mathrm{tr}\big(\nabla g_{2N}(X)^T Y\big) \\
\vdots & \vdots & & \vdots \\
\mathrm{tr}\big(\nabla g_{M1}(X)^T Y\big) & \mathrm{tr}\big(\nabla g_{M2}(X)^T Y\big) & \cdots & \mathrm{tr}\big(\nabla g_{MN}(X)^T Y\big)
\end{bmatrix}
$$

$$
=
\begin{bmatrix}
\sum_{k,l} \frac{\partial g_{11}(X)}{\partial X_{kl}} Y_{kl} & \sum_{k,l} \frac{\partial g_{12}(X)}{\partial X_{kl}} Y_{kl} & \cdots & \sum_{k,l} \frac{\partial g_{1N}(X)}{\partial X_{kl}} Y_{kl} \\
\sum_{k,l} \frac{\partial g_{21}(X)}{\partial X_{kl}} Y_{kl} & \sum_{k,l} \frac{\partial g_{22}(X)}{\partial X_{kl}} Y_{kl} & \cdots & \sum_{k,l} \frac{\partial g_{2N}(X)}{\partial X_{kl}} Y_{kl} \\
\vdots & \vdots & & \vdots \\
\sum_{k,l} \frac{\partial g_{M1}(X)}{\partial X_{kl}} Y_{kl} & \sum_{k,l} \frac{\partial g_{M2}(X)}{\partial X_{kl}} Y_{kl} & \cdots & \sum_{k,l} \frac{\partial g_{MN}(X)}{\partial X_{kl}} Y_{kl}
\end{bmatrix}
\tag{1394}
$$

from which it follows

$$
\overset{\rightarrow Y}{dg}(X) = \sum_{k,l} \frac{\partial g(X)}{\partial X_{kl}} Y_{kl}
\tag{1395}
$$

Yet for all $X \in \mathrm{dom}\, g$, any $Y \in \mathbb{R}^{K \times L}$, and some open interval of $t \in \mathbb{R}$

$$
g(X + t\,Y) = g(X) + t\,\overset{\rightarrow Y}{dg}(X) + o(t^2)
\tag{1396}
$$

which is the first-order Taylor series expansion about $X$. [137, §18.4] [85, §2.3.4] Differentiation with respect to $t$ and subsequent $t$-zeroing isolates the second term of the expansion. Thus differentiating and zeroing $g(X+t\,Y)$ in $t$ is an operation equivalent to individually differentiating and zeroing every entry $g_{mn}(X + t\,Y)$ as in (1393). So the directional derivative of $g(X)$ in any direction $Y \in \mathbb{R}^{K \times L}$ evaluated at $X \in \mathrm{dom}\, g$ becomes

$$
\overset{\rightarrow Y}{dg}(X) = \left.\frac{d}{dt}\right|_{t=0} g(X + t\,Y) \in \mathbb{R}^{M \times N}
\tag{1397}
$$

---

[D.2]Although $Y$ is a matrix, we may regard it as a vector in $\mathbb{R}^{KL}$.

$$v \triangleq \begin{bmatrix} \nabla_x \hat{f}(\alpha) \\ \\ \xrightarrow{\nabla_x \hat{f}(\alpha)} \frac{1}{2} d\hat{f}(\alpha) \end{bmatrix}$$
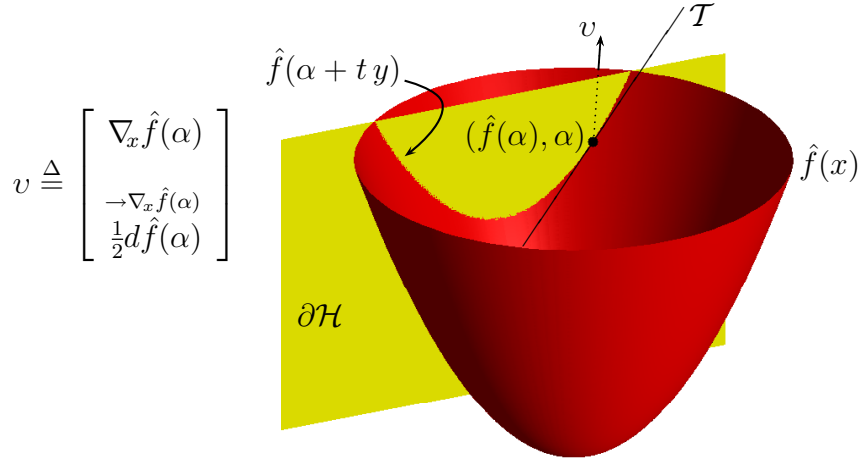
Figure 97: Drawn is a convex quadratic bowl in $\mathbb{R}^\mathbf{2} \times \mathbb{R}$; $\hat{f}(x) = x^T x : \mathbb{R}^\mathbf{2} \to \mathbb{R}$ *versus* $x$ on some open disc in $\mathbb{R}^\mathbf{2}$. Plane slice $\partial\mathcal{H}$ is perpendicular to function domain. Slice intersection with domain connotes bidirectional vector $y$. Tangent line $\mathcal{T}$ slope at point $(\alpha, \hat{f}(\alpha))$ is directional derivative value $\nabla_x \hat{f}(\alpha)^T y$ (1424) at $\alpha$ in slice direction $y$. Recall, negative gradient $-\nabla_x \hat{f}(x) \in \mathbb{R}^\mathbf{2}$ is always steepest descent direction [248]. [137, §15.6] When vector $v \in \mathbb{R}^\mathbf{3}$ entry $v_3$ is half directional derivative in gradient direction at $\alpha$ and when $\begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \nabla_x \hat{f}(\alpha)$, then $-v$ points directly toward bowl bottom.

[177, §2.1, §5.4.5] [25, §6.3.1] which is simplest. The derivative with respect to $t$ makes the directional derivative (1397) resemble ordinary calculus (§D.2); *e.g.*, when $g(X)$ is linear, $\overset{\to Y}{dg}(X) = g(Y)$. [157, §7.2]

### D.1.4.1 Interpretation directional derivative

In the case of any differentiable real function $\hat{f}(X) : \mathbb{R}^{K \times L} \to \mathbb{R}$, the directional derivative of $\hat{f}(X)$ at $X$ in any direction $Y$ yields the slope of $\hat{f}$ along the line $X + tY$ through its domain (parametrized by $t \in \mathbb{R}$) evaluated at $t = 0$. For higher-dimensional functions, by (1394), this slope interpretation can be applied to each entry of the directional derivative. Unlike the gradient, directional derivative does not expand dimension; *e.g.*, directional derivative in (1397) retains the dimensions of $g$.

Figure **97**, for example, shows a plane slice of a real convex bowl-shaped function $\hat{f}(x)$ along a line $\alpha + t\,y$ through its domain. The slice reveals a one-dimensional real function of $t$; $\hat{f}(\alpha + t\,y)$. The directional derivative at $x = \alpha$ in direction $y$ is the slope of $\hat{f}(\alpha + t\,y)$ with respect to $t$ at $t = 0$. In the case of a real function having vector argument $h(X) : \mathbb{R}^K \to \mathbb{R}$, its directional derivative in the normalized direction of its gradient is the gradient magnitude. (1424) For a real function of real variable, the directional derivative evaluated at any point in the function domain is just the slope of that function there scaled by the real direction. (*confer* §3.1.1.4)

**D.1.4.1.1   Theorem.**   *Directional derivative condition for optimization.*
[157, §7.4]  Suppose $\hat{f}(X) : \mathbb{R}^{K \times L} \to \mathbb{R}$ is minimized on convex set $\mathcal{C} \subseteq \mathbb{R}^{p \times k}$ by $X^\star$, and the directional derivative of $\hat{f}$ exists there. Then for all $X \in \mathcal{C}$

$$d\hat{f}(X) \overset{\to X - X^\star}{} \geq 0 \qquad (1398)$$

$\diamond$

**D.1.4.1.2   Example.**   *Simple bowl.*
Bowl function (Figure **97**)

$$\hat{f}(x) : \mathbb{R}^K \to \mathbb{R} \overset{\Delta}{=} (x - a)^T (x - a) - b \qquad (1399)$$

has function offset $-b \in \mathbb{R}$, axis of revolution at $x = a$, and positive definite Hessian (1355) everywhere in its domain (an open *hyperdisc* in $\mathbb{R}^K$); *id est*, strictly convex quadratic $\hat{f}(x)$ has unique global minimum equal to $-b$ at $x = a$. A vector $-\upsilon$ based anywhere in $\operatorname{dom}\hat{f} \times \mathbb{R}$ pointing toward the unique bowl-bottom is specified:

$$\upsilon \propto \begin{bmatrix} x - a \\ \hat{f}(x) + b \end{bmatrix} \in \mathbb{R}^K \times \mathbb{R} \qquad (1400)$$

Such a vector is

$$\upsilon = \begin{bmatrix} \nabla_x \hat{f}(x) \\[6pt] \frac{1}{2} d\hat{f}(x) \overset{\to \nabla_x \hat{f}(x)}{} \end{bmatrix} \qquad (1401)$$

since the gradient is

$$\nabla_x \hat{f}(x) = 2(x - a) \qquad (1402)$$

and the directional derivative in the direction of the gradient is (1424)

$$
d\hat{f}(x) \stackrel{\rightarrow \nabla_x \hat{f}(x)}{=} \nabla_x \hat{f}(x)^T \nabla_x \hat{f}(x) = 4(x-a)^T(x-a) = 4\left(\hat{f}(x)+b\right) \quad (1403)
$$

$\square$

## D.1.5  Second directional derivative

By similar argument, it so happens: the second directional derivative is equally simple. Given $g(X) : \mathbb{R}^{K \times L} \to \mathbb{R}^{M \times N}$ on open domain,

$$
\nabla \frac{\partial g_{mn}(X)}{\partial X_{kl}} = \frac{\partial \nabla g_{mn}(X)}{\partial X_{kl}} = \begin{bmatrix} \frac{\partial^2 g_{mn}(X)}{\partial X_{kl}\partial X_{11}} & \frac{\partial^2 g_{mn}(X)}{\partial X_{kl}\partial X_{12}} & \cdots & \frac{\partial^2 g_{mn}(X)}{\partial X_{kl}\partial X_{1L}} \\ \frac{\partial^2 g_{mn}(X)}{\partial X_{kl}\partial X_{21}} & \frac{\partial^2 g_{mn}(X)}{\partial X_{kl}\partial X_{22}} & \cdots & \frac{\partial^2 g_{mn}(X)}{\partial X_{kl}\partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial^2 g_{mn}(X)}{\partial X_{kl}\partial X_{K1}} & \frac{\partial^2 g_{mn}(X)}{\partial X_{kl}\partial X_{K2}} & \cdots & \frac{\partial^2 g_{mn}(X)}{\partial X_{kl}\partial X_{KL}} \end{bmatrix} \in \mathbb{R}^{K \times L} \quad (1404)
$$

$$
\nabla^2 g_{mn}(X) = \begin{bmatrix} \nabla \frac{\partial g_{mn}(X)}{\partial X_{11}} & \nabla \frac{\partial g_{mn}(X)}{\partial X_{12}} & \cdots & \nabla \frac{\partial g_{mn}(X)}{\partial X_{1L}} \\ \nabla \frac{\partial g_{mn}(X)}{\partial X_{21}} & \nabla \frac{\partial g_{mn}(X)}{\partial X_{22}} & \cdots & \nabla \frac{\partial g_{mn}(X)}{\partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \nabla \frac{\partial g_{mn}(X)}{\partial X_{K1}} & \nabla \frac{\partial g_{mn}(X)}{\partial X_{K2}} & \cdots & \nabla \frac{\partial g_{mn}(X)}{\partial X_{KL}} \end{bmatrix} \in \mathbb{R}^{K \times L \times K \times L}
$$

$$
(1405)
$$

$$
= \begin{bmatrix} \frac{\partial \nabla g_{mn}(X)}{\partial X_{11}} & \frac{\partial \nabla g_{mn}(X)}{\partial X_{12}} & \cdots & \frac{\partial \nabla g_{mn}(X)}{\partial X_{1L}} \\ \frac{\partial \nabla g_{mn}(X)}{\partial X_{21}} & \frac{\partial \nabla g_{mn}(X)}{\partial X_{22}} & \cdots & \frac{\partial \nabla g_{mn}(X)}{\partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial \nabla g_{mn}(X)}{\partial X_{K1}} & \frac{\partial \nabla g_{mn}(X)}{\partial X_{K2}} & \cdots & \frac{\partial \nabla g_{mn}(X)}{\partial X_{KL}} \end{bmatrix}
$$

Rotating our perspective, we get several views of the second-order gradient:

$$
\nabla^2 g(X) = \begin{bmatrix} \nabla^2 g_{11}(X) & \nabla^2 g_{12}(X) & \cdots & \nabla^2 g_{1N}(X) \\ \nabla^2 g_{21}(X) & \nabla^2 g_{22}(X) & \cdots & \nabla^2 g_{2N}(X) \\ \vdots & \vdots & & \vdots \\ \nabla^2 g_{M1}(X) & \nabla^2 g_{M2}(X) & \cdots & \nabla^2 g_{MN}(X) \end{bmatrix} \in \mathbb{R}^{M \times N \times K \times L \times K \times L} \quad (1406)
$$

$$\nabla^2 g(X)^{T_1} = \begin{bmatrix} \nabla\frac{\partial g(X)}{\partial X_{11}} & \nabla\frac{\partial g(X)}{\partial X_{12}} & \cdots & \nabla\frac{\partial g(X)}{\partial X_{1L}} \\ \nabla\frac{\partial g(X)}{\partial X_{21}} & \nabla\frac{\partial g(X)}{\partial X_{22}} & \cdots & \nabla\frac{\partial g(X)}{\partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \nabla\frac{\partial g(X)}{\partial X_{K1}} & \nabla\frac{\partial g(X)}{\partial X_{K2}} & \cdots & \nabla\frac{\partial g(X)}{\partial X_{KL}} \end{bmatrix} \in \mathbb{R}^{K \times L \times M \times N \times K \times L} \quad (1407)$$

$$\nabla^2 g(X)^{T_2} = \begin{bmatrix} \frac{\partial \nabla g(X)}{\partial X_{11}} & \frac{\partial \nabla g(X)}{\partial X_{12}} & \cdots & \frac{\partial \nabla g(X)}{\partial X_{1L}} \\ \frac{\partial \nabla g(X)}{\partial X_{21}} & \frac{\partial \nabla g(X)}{\partial X_{22}} & \cdots & \frac{\partial \nabla g(X)}{\partial X_{2L}} \\ \vdots & \vdots & & \vdots \\ \frac{\partial \nabla g(X)}{\partial X_{K1}} & \frac{\partial \nabla g(X)}{\partial X_{K2}} & \cdots & \frac{\partial \nabla g(X)}{\partial X_{KL}} \end{bmatrix} \in \mathbb{R}^{K \times L \times K \times L \times M \times N} \quad (1408)$$

Assuming the limits exist, we may state the partial derivative of the $mn^{\text{th}}$ entry of $g$ with respect to the $kl^{\text{th}}$ and $ij^{\text{th}}$ entries of $X$;

$$\frac{\partial^2 g_{mn}(X)}{\partial X_{kl}\,\partial X_{ij}} = \lim_{\Delta\tau,\Delta t \to 0} \frac{g_{mn}(X+\Delta t\, e_k e_l^T + \Delta\tau\, e_i e_j^T) - g_{mn}(X+\Delta t\, e_k e_l^T) - \left(g_{mn}(X+\Delta\tau\, e_i e_j^T) - g_{mn}(X)\right)}{\Delta\tau\,\Delta t}$$

$$(1409)$$

Differentiating (1389) and then scaling by $Y_{ij}$

$$\frac{\partial^2 g_{mn}(X)}{\partial X_{kl}\,\partial X_{ij}} Y_{kl} Y_{ij} = \lim_{\Delta t \to 0} \frac{\partial g_{mn}(X+\Delta t\, Y_{kl}\, e_k e_l^T) - \partial g_{mn}(X)}{\partial X_{ij}\,\Delta t} Y_{ij} \quad (1410)$$

$$= \lim_{\Delta\tau,\Delta t \to 0} \frac{g_{mn}(X+\Delta t\, Y_{kl}\, e_k e_l^T + \Delta\tau\, Y_{ij}\, e_i e_j^T) - g_{mn}(X+\Delta t\, Y_{kl}\, e_k e_l^T) - \left(g_{mn}(X+\Delta\tau\, Y_{ij}\, e_i e_j^T) - g_{mn}(X)\right)}{\Delta\tau\,\Delta t}$$

which can be proved by substitution of variables in (1409).   The $mn^{\text{th}}$ second-order total differential due to any $Y \in \mathbb{R}^{K \times L}$ is

$$d^2 g_{mn}(X)|_{dX \to Y} = \sum_{i,j} \sum_{k,l} \frac{\partial^2 g_{mn}(X)}{\partial X_{kl}\,\partial X_{ij}} Y_{kl} Y_{ij} = \operatorname{tr}\left(\nabla_X \operatorname{tr}\left(\nabla g_{mn}(X)^T Y\right)^T Y\right) \quad (1411)$$

$$= \sum_{i,j} \lim_{\Delta t \to 0} \frac{\partial g_{mn}(X+\Delta t\, Y) - \partial g_{mn}(X)}{\partial X_{ij}\,\Delta t} Y_{ij} \quad (1412)$$

$$= \lim_{\Delta t \to 0} \frac{g_{mn}(X+2\Delta t\, Y) - 2g_{mn}(X+\Delta t\, Y) + g_{mn}(X)}{\Delta t^2} \quad (1413)$$

$$= \frac{d^2}{dt^2}\bigg|_{t=0} g_{mn}(X+t\, Y) \quad (1414)$$

Hence the second directional derivative,

$$\overset{\rightarrow Y}{dg^2}(X) \triangleq \left.\begin{bmatrix} d^2g_{11}(X) & d^2g_{12}(X) & \cdots & d^2g_{1N}(X) \\ d^2g_{21}(X) & d^2g_{22}(X) & \cdots & d^2g_{2N}(X) \\ \vdots & \vdots & & \vdots \\ d^2g_{M1}(X) & d^2g_{M2}(X) & \cdots & d^2g_{MN}(X) \end{bmatrix}\right|_{dX \to Y} \in \mathbb{R}^{M \times N}$$

$$= \begin{bmatrix} \operatorname{tr}\left(\nabla\operatorname{tr}\left(\nabla g_{11}(X)^T Y\right)^T Y\right) & \operatorname{tr}\left(\nabla\operatorname{tr}\left(\nabla g_{12}(X)^T Y\right)^T Y\right) & \cdots & \operatorname{tr}\left(\nabla\operatorname{tr}\left(\nabla g_{1N}(X)^T Y\right)^T Y\right) \\ \operatorname{tr}\left(\nabla\operatorname{tr}\left(\nabla g_{21}(X)^T Y\right)^T Y\right) & \operatorname{tr}\left(\nabla\operatorname{tr}\left(\nabla g_{22}(X)^T Y\right)^T Y\right) & \cdots & \operatorname{tr}\left(\nabla\operatorname{tr}\left(\nabla g_{2N}(X)^T Y\right)^T Y\right) \\ \vdots & \vdots & & \vdots \\ \operatorname{tr}\left(\nabla\operatorname{tr}\left(\nabla g_{M1}(X)^T Y\right)^T Y\right) & \operatorname{tr}\left(\nabla\operatorname{tr}\left(\nabla g_{M2}(X)^T Y\right)^T Y\right) & \cdots & \operatorname{tr}\left(\nabla\operatorname{tr}\left(\nabla g_{MN}(X)^T Y\right)^T Y\right) \end{bmatrix}$$

$$= \begin{bmatrix} \sum_{i,j}\sum_{k,l}\frac{\partial^2 g_{11}(X)}{\partial X_{kl}\partial X_{ij}}Y_{kl}Y_{ij} & \sum_{i,j}\sum_{k,l}\frac{\partial^2 g_{12}(X)}{\partial X_{kl}\partial X_{ij}}Y_{kl}Y_{ij} & \cdots & \sum_{i,j}\sum_{k,l}\frac{\partial^2 g_{1N}(X)}{\partial X_{kl}\partial X_{ij}}Y_{kl}Y_{ij} \\ \sum_{i,j}\sum_{k,l}\frac{\partial^2 g_{21}(X)}{\partial X_{kl}\partial X_{ij}}Y_{kl}Y_{ij} & \sum_{i,j}\sum_{k,l}\frac{\partial^2 g_{22}(X)}{\partial X_{kl}\partial X_{ij}}Y_{kl}Y_{ij} & \cdots & \sum_{i,j}\sum_{k,l}\frac{\partial^2 g_{2N}(X)}{\partial X_{kl}\partial X_{ij}}Y_{kl}Y_{ij} \\ \vdots & \vdots & & \vdots \\ \sum_{i,j}\sum_{k,l}\frac{\partial^2 g_{M1}(X)}{\partial X_{kl}\partial X_{ij}}Y_{kl}Y_{ij} & \sum_{i,j}\sum_{k,l}\frac{\partial^2 g_{M2}(X)}{\partial X_{kl}\partial X_{ij}}Y_{kl}Y_{ij} & \cdots & \sum_{i,j}\sum_{k,l}\frac{\partial^2 g_{MN}(X)}{\partial X_{kl}\partial X_{ij}}Y_{kl}Y_{ij} \end{bmatrix}$$

$$(1415)$$

from which it follows

$$\overset{\rightarrow Y}{dg^2}(X) = \sum_{i,j}\sum_{k,l}\frac{\partial^2 g(X)}{\partial X_{kl}\partial X_{ij}}Y_{kl}Y_{ij} = \sum_{i,j}\frac{\partial}{\partial X_{ij}}\overset{\rightarrow Y}{dg}(X)\,Y_{ij} \qquad (1416)$$

Yet for all $X \in \operatorname{dom} g$, any $Y \in \mathbb{R}^{K \times L}$, and some open interval of $t \in \mathbb{R}$

$$g(X + tY) = g(X) + t\,\overset{\rightarrow Y}{dg}(X) + \frac{1}{2!}t^2\,\overset{\rightarrow Y}{dg^2}(X) + o(t^3) \qquad (1417)$$

which is the second-order Taylor series expansion about $X$. [137, §18.4] [85, §2.3.4] Differentiating twice with respect to $t$ and subsequent $t$-zeroing isolates the third term of the expansion. Thus differentiating and zeroing $g(X + tY)$ in $t$ is an operation equivalent to individually differentiating and zeroing every entry $g_{mn}(X + tY)$ as in (1414). So the second directional derivative becomes

$$\overset{\rightarrow Y}{dg^2}(X) = \left.\frac{d^2}{dt^2}\right|_{t=0} g(X + tY) \in \mathbb{R}^{M \times N} \qquad (1418)$$

[177, §2.1, §5.4.5] [25, §6.3.1] which is again simplest. (*confer* (1397))

### D.1.6    Taylor series

Series expansions of the differentiable matrix-valued function $g(X)$, of matrix argument, were given earlier in (1396) and (1417). Assuming $g(X)$ has continuous first-, second-, and third-order gradients over the open set $\operatorname{dom} g$, then for $X \in \operatorname{dom} g$ and any $Y \in \mathbb{R}^{K \times L}$ the complete Taylor series on some open interval of $\mu \in \mathbb{R}$ is expressed

$$g(X+\mu Y) = g(X) + \mu \overset{\to Y}{dg}(X) + \frac{1}{2!}\mu^2 \overset{\to Y}{dg^2}(X) + \frac{1}{3!}\mu^3 \overset{\to Y}{dg^3}(X) + o(\mu^4) \quad (1419)$$

or on some open interval of $\|Y\|$

$$g(Y) = g(X) \ + \ \overset{\to Y - X}{dg(X)} \ + \ \frac{1}{2!}\overset{\to Y - X}{dg^2(X)} \ + \ \frac{1}{3!}\overset{\to Y - X}{dg^3(X)} \ + \ o(\|Y\|^4) \qquad (1420)$$

which are third-order expansions about $X$. The *mean value theorem* from calculus is what insures the finite order of the series. [28, §1.1] [27, App.A.5] [125, §0.4] [137]

In the case of a real function $g(X) : \mathbb{R}^{K \times L} \to \mathbb{R}$, all the directional derivatives are in $\mathbb{R}$:

$$\overset{\to Y}{dg}(X) = \operatorname{tr}\big(\nabla g(X)^T Y\big) \tag{1421}$$

$$\overset{\to Y}{dg^2}(X) = \operatorname{tr}\Big(\nabla_X \operatorname{tr}\big(\nabla g(X)^T Y\big)^T Y\Big) = \operatorname{tr}\Big(\nabla_X \overset{\to Y}{dg}(X)^T Y\Big) \tag{1422}$$

$$\overset{\to Y}{dg^3}(X) = \operatorname{tr}\Big(\nabla_X \operatorname{tr}\big(\nabla_X \operatorname{tr}\big(\nabla g(X)^T Y\big)^T Y\big)^T Y\Big) = \operatorname{tr}\Big(\nabla_X \overset{\to Y}{dg^2}(X)^T Y\Big) \tag{1423}$$

In the case $g(X) : \mathbb{R}^K \to \mathbb{R}$ has vector argument, they further simplify:

$$\overset{\to Y}{dg}(X) = \nabla g(X)^T Y \tag{1424}$$

$$\overset{\to Y}{dg^2}(X) = Y^T \nabla^2 g(X) Y \tag{1425}$$

$$\overset{\to Y}{dg^3}(X) = \nabla_X \big(Y^T \nabla^2 g(X) Y\big)^T Y \tag{1426}$$

and so on.

**D.1.6.0.1    Exercise.**    *log det.*                            (*confer* [39, p.644])
Find the first two terms of the Taylor series expansion (1420) for $\log \det X$.
▼

### D.1.7 Correspondence of gradient to derivative

From the foregoing expressions for directional derivative, we derive a relationship between the gradient with respect to matrix $X$ and the derivative with respect to real variable $t$ :

#### D.1.7.1 first-order

Removing from (1397) the evaluation at $t = 0$ ,[D.3] we find an expression for the directional derivative of $g(X)$ in direction $Y$ evaluated anywhere along a line $X + tY$ (parametrized by $t$) intersecting $\operatorname{dom} g$

$$\overset{\rightarrow Y}{dg}(X + tY) = \frac{d}{dt}g(X + tY) \tag{1427}$$

In the general case $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}^{M \times N}$, from (1390) and (1393) we find

$$\operatorname{tr}\big(\nabla_X g_{mn}(X + tY)^T Y\big) = \frac{d}{dt}g_{mn}(X + tY) \tag{1428}$$

which is valid at $t = 0$, of course, when $X \in \operatorname{dom} g$. In the important case of a real function $g(X) : \mathbb{R}^{K \times L} \rightarrow \mathbb{R}$, from (1421) we have simply

$$\operatorname{tr}\big(\nabla_X g(X + tY)^T Y\big) = \frac{d}{dt}g(X + tY) \tag{1429}$$

When, additionally, $g(X) : \mathbb{R}^K \rightarrow \mathbb{R}$ has vector argument,

$$\nabla_X g(X + tY)^T Y = \frac{d}{dt}g(X + tY) \tag{1430}$$

---

[D.3]Justified by replacing $X$ with $X + tY$ in (1390)-(1392); beginning,

$$dg_{mn}(X + tY)|_{dX \rightarrow Y} = \sum_{k,l} \frac{\partial g_{mn}(X + tY)}{\partial X_{kl}} Y_{kl}$$

**D.1.7.1.1   Example.**   *Gradient.*
$g(X) = w^T X^T X w$, $\ X \in \mathbb{R}^{K \times L}$, $\ w \in \mathbb{R}^L$.  Using the tables in §D.2,

$$\operatorname{tr}\left(\nabla_X g(X + t\,Y)^T Y\right) \ = \ \operatorname{tr}\left(2 w w^T (X^T + t\,Y^T) Y\right) \tag{1431}$$
$$= \ 2 w^T (X^T Y + t\,Y^T Y) w \tag{1432}$$

Applying the equivalence (1429),

$$\frac{d}{dt} g(X + t\,Y) \ = \ \frac{d}{dt} w^T (X + t\,Y)^T (X + t\,Y) w \tag{1433}$$
$$= \ w^T \left(X^T Y + Y^T X + 2t\,Y^T Y\right) w \tag{1434}$$
$$= \ 2 w^T (X^T Y + t\,Y^T Y) w \tag{1435}$$

which is the same as (1432); hence, equivalence is demonstrated.
    It is easy to extract $\nabla g(X)$ from (1435) knowing only (1429):

$$\begin{aligned}
\operatorname{tr}\left(\nabla_X g(X + t\,Y)^T Y\right) \ &= \ 2 w^T (X^T Y + t\,Y^T Y) w \\
&= \ 2 \operatorname{tr}\left(w w^T (X^T + t\,Y^T) Y\right) \\
\operatorname{tr}\left(\nabla_X g(X)^T Y\right) \ &= \ 2 \operatorname{tr}\left(w w^T X^T Y\right) \\
&\Leftrightarrow \\
\nabla_X g(X) \ &= \ 2 X w w^T
\end{aligned} \tag{1436}$$

$\square$

**D.1.7.2   second-order**

Likewise removing the evaluation at $t = 0$ from (1418),

$$\overset{\rightarrow Y}{dg^2}(X + t\,Y) = \frac{d^2}{dt^2} g(X + t\,Y) \tag{1437}$$

we can find a similar relationship between the second-order gradient and the second derivative: In the general case $g(X) : \mathbb{R}^{K \times L} \to \mathbb{R}^{M \times N}$ from (1411) and (1414),

$$\operatorname{tr}\left(\nabla_X \operatorname{tr}\left(\nabla_X g_{mn}(X + t\,Y)^T Y\right)^T Y\right) = \frac{d^2}{dt^2} g_{mn}(X + t\,Y) \tag{1438}$$

In the case of a real function $g(X) : \mathbb{R}^{K \times L} \to \mathbb{R}$ we have, of course,

$$\operatorname{tr}\left(\nabla_X \operatorname{tr}\left(\nabla_X g(X + t\,Y)^T Y\right)^T Y\right) = \frac{d^2}{dt^2} g(X + t\,Y) \tag{1439}$$

From (1425), the simpler case, where the real function $g(X) : \mathbb{R}^K \to \mathbb{R}$ has vector argument,

$$Y^T \nabla_X^2 g(X + t\,Y)Y = \frac{d^2}{dt^2} g(X + t\,Y) \tag{1440}$$

**D.1.7.2.1 Example.** *Second-order gradient.*
Given real function $g(X) = \log \det X$ having domain $\operatorname{int} \mathbb{S}_+^K$, we want to find $\nabla^2 g(X) \in \mathbb{R}^{K \times K \times K \times K}$. From the tables in §D.2,

$$h(X) \triangleq \nabla g(X) = X^{-1} \in \operatorname{int} \mathbb{S}_+^K \tag{1441}$$

so $\nabla^2 g(X) = \nabla h(X)$. By (1428) and (1396), for $Y \in \mathbb{S}^K$

$$\operatorname{tr}\!\left(\nabla h_{mn}(X)^T Y\right) \;=\; \frac{d}{dt}\bigg|_{t=0} h_{mn}(X + t\,Y) \tag{1442}$$

$$= \left(\frac{d}{dt}\bigg|_{t=0} h(X + t\,Y)\right)_{mn} \tag{1443}$$

$$= \left(\frac{d}{dt}\bigg|_{t=0} (X + t\,Y)^{-1}\right)_{mn} \tag{1444}$$

$$= -\left(X^{-1}Y X^{-1}\right)_{mn} \tag{1445}$$

Setting $Y$ to a member of the standard basis $E_{kl} = e_k e_l^T$, for $k, l \in \{1 \ldots K\}$, and employing a property of the trace function (31) we find

$$\nabla^2 g(X)_{mnkl} \;=\; \operatorname{tr}\!\left(\nabla h_{mn}(X)^T E_{kl}\right) \;=\; \nabla h_{mn}(X)_{kl} \;=\; -\left(X^{-1} E_{kl} X^{-1}\right)_{mn} \tag{1446}$$

$$\nabla^2 g(X)_{kl} \;=\; \nabla h(X)_{kl} \;=\; -\left(X^{-1} E_{kl} X^{-1}\right) \in \mathbb{R}^{K \times K} \tag{1447}$$

$\square$

From all these first- and second-order expressions, we may generate new ones by evaluating both sides at arbitrary $t$ (in some open interval) but only after the differentiation.

## D.2   Tables of gradients and derivatives

[96] [43]

- When proving results for symmetric matrices algebraically, it is critical to take gradients ignoring symmetry and to then substitute symmetric entries afterward.

- $a\,,b\in\mathbb{R}^n,\quad x\,,y\in\mathbb{R}^k,\quad A\,,B\in\mathbb{R}^{m\times n},\quad X,Y\in\mathbb{R}^{K\times L},\quad t\,,\mu\in\mathbb{R}\,,$
  $i,j,k,\ell,K,L\,,m\,,n\,,M\,,N$ are integers, unless otherwise noted.

- $x^\mu$ means $\delta\big(\delta(x)^\mu\big)$ for $\mu\in\mathbb{R}$ ; *id est*, entrywise vector exponentiation. $\delta$ is the main-diagonal linear operator (1036). $x^0\overset{\Delta}{=}\mathbf{1},\ \ X^0\overset{\Delta}{=}I$ if square.

- $\frac{d}{dx}\overset{\Delta}{=}\begin{bmatrix}\frac{d}{dx_1}\\ \vdots\\ \frac{d}{dx_k}\end{bmatrix},\quad \overset{\rightarrow y}{dg}(x)\,,\quad \overset{\rightarrow y}{dg^2}(x)$ (directional derivatives §D.1), $\ \log x\,,$
  $\operatorname{sgn}x\,,\ \sin x\,,\ x/y$ (Hadamard quotient), $\sqrt{x}$ (entrywise square root), *etcetera*, are maps $f:\mathbb{R}^k\to\mathbb{R}^k$ that maintain dimension; *e.g.*, (§A.1.1)

$$\frac{d}{dx}x^{-1}\ \overset{\Delta}{=}\ \nabla_x\,\mathbf{1}^T\delta(x)^{-1}\mathbf{1} \tag{1448}$$

- The standard basis: $\big\{E_{kl}=e_k e_\ell^T\in\mathbb{R}^{K\times K}\mid k\,,\ell\in\{1\ldots K\}\big\}$

- For $A$ a scalar or matrix, we have the Taylor series [45, §3.6]

$$e^A=\sum_{k=0}^{\infty}\frac{1}{k!}A^k \tag{1449}$$

Further, [215, §5.4]
$$e^A\succ0\qquad\forall\,A\in\mathbb{S}^m \tag{1450}$$

- For all square $A$ and integer $k$

$$\det{}^k A=\det A^k \tag{1451}$$

- Table entries with notation $X\in\mathbb{R}^{\mathbf{2\times2}}$ have been algebraically verified in that dimension but may hold more broadly.

## D.2.1  Algebraic

$$\nabla_x\, x = \nabla_x\, x^T = I \in \mathbb{R}^{k \times k}$$

$$\nabla_X X = \nabla_X X^T \triangleq I \in \mathbb{R}^{K \times L \times K \times L} \qquad \text{(identity)}$$

$$\nabla_x (Ax - b) = A^T$$

$$\nabla_x \left(x^T A - b^T\right) = A$$

$$\nabla_x (Ax - b)^T (Ax - b) = 2A^T (Ax - b)$$

$$\nabla_x^2 (Ax - b)^T (Ax - b) = 2A^T A$$

$$\nabla_x \left(x^T A x + 2x^T B y + y^T C y\right) = (A + A^T)x + 2By$$

$$\nabla_x^2 \left(x^T A x + 2x^T B y + y^T C y\right) = A + A^T$$

$$\nabla_X\, a^T X b = \nabla_X\, b^T X^T a = a b^T$$

$$\nabla_X\, a^T X^2 b = X^T a b^T + a b^T X^T$$

$$\nabla_X\, a^T X^{-1} b = -X^{-T} a b^T X^{-T}$$

$$\nabla_X (X^{-1})_{kl} = \frac{\partial X^{-1}}{\partial X_{kl}} = -X^{-1} E_{kl}\, X^{-1}, \;\; confer\,(1388)(1447)$$

$$\nabla_x\, a^T x^T x b = 2x a^T b$$

$$\nabla_X\, a^T X^T X b = X(a b^T + b a^T)$$

$$\nabla_x\, a^T x x^T b = (a b^T + b a^T)x$$

$$\nabla_X\, a^T X X^T b = (a b^T + b a^T)X$$

$$\nabla_x\, a^T x^T x a = 2x a^T a$$

$$\nabla_X\, a^T X^T X a = 2X a a^T$$

$$\nabla_x\, a^T x x^T a = 2a a^T x$$

$$\nabla_X\, a^T X X^T a = 2a a^T X$$

$$\nabla_x\, a^T y x^T b = b a^T y$$

$$\nabla_X\, a^T Y X^T b = b a^T Y$$

$$\nabla_x\, a^T y^T x b = y b^T a$$

$$\nabla_X\, a^T Y^T X b = Y a b^T$$

$$\nabla_x\, a^T x y^T b = a b^T y$$

$$\nabla_X\, a^T X Y^T b = a b^T Y$$

$$\nabla_x\, a^T x^T y b = y a^T b$$

$$\nabla_X\, a^T X^T Y b = Y b a^T$$

**Algebraic** continued

$$\frac{d}{dt}(X + tY) = Y$$

$$\frac{d}{dt}B^T(X + tY)^{-1}A = -B^T(X + tY)^{-1}Y(X + tY)^{-1}A$$

$$\frac{d}{dt}B^T(X + tY)^{-T}A = -B^T(X + tY)^{-T}Y^T(X + tY)^{-T}A$$

$$\frac{d}{dt}B^T(X + tY)^{\mu}A = \dots, \quad -1 \leq \mu \leq 1, \quad X, Y \in \mathbb{S}_+^M$$

$$\frac{d^2}{dt^2}B^T(X + tY)^{-1}A = 2B^T(X + tY)^{-1}Y(X + tY)^{-1}Y(X + tY)^{-1}A$$

$$\frac{d}{dt}\big((X + tY)^T A(X + tY)\big) = Y^T A X + X^T A Y + 2t Y^T A Y$$

$$\frac{d^2}{dt^2}\big((X + tY)^T A(X + tY)\big) = 2Y^T A Y$$

$$\frac{d}{dt}\big((X + tY)A(X + tY)\big) = YAX + XAY + 2t YAY$$

$$\frac{d^2}{dt^2}\big((X + tY)A(X + tY)\big) = 2YAY$$

## D.2.2   Trace Kronecker

$$\nabla_{\mathrm{vec}\,X}\mathrm{tr}(AXBX^T) = \nabla_{\mathrm{vec}\,X}\mathrm{vec}(X)^T(B^T \otimes A)\,\mathrm{vec}\,X = (B \otimes A^T + B^T \otimes A)\,\mathrm{vec}\,X$$

$$\nabla^2_{\mathrm{vec}\,X}\mathrm{tr}(AXBX^T) = \nabla^2_{\mathrm{vec}\,X}\mathrm{vec}(X)^T(B^T \otimes A)\,\mathrm{vec}\,X = B \otimes A^T + B^T \otimes A$$

## D.2.3   Trace

$\nabla_x \, \mu \, x = \mu I$

$\nabla_X \operatorname{tr} \mu X = \nabla_X \mu \operatorname{tr} X = \mu I$

$\nabla_x \mathbf{1}^T \delta(x)^{-1} \mathbf{1} = \frac{d}{dx} x^{-1} = -x^{-2}$

$\nabla_X \operatorname{tr} X^{-1} = -X^{-2T}$

$\nabla_x \mathbf{1}^T \delta(x)^{-1} y = -\delta(x)^{-2} y$

$\nabla_X \operatorname{tr}(X^{-1} Y) = \nabla_X \operatorname{tr}(Y X^{-1}) = -X^{-T} Y^T X^{-T}$

$\frac{d}{dx} x^{\mu} = \mu x^{\mu - 1}$

$\nabla_X \operatorname{tr} X^{\mu} = \mu X^{(\mu-1)T}, \qquad\qquad\qquad X \in \mathbb{R}^{\mathbf{2 \times 2}}$

$\nabla_X \operatorname{tr} X^j = j X^{(j-1)T}$

$\nabla_x (b - a^T x)^{-1} = (b - a^T x)^{-2} a$

$\nabla_X \operatorname{tr}\big((B - AX)^{-1}\big) = \big((B - AX)^{-2} A\big)^T$

$\nabla_x (b - a^T x)^{\mu} = -\mu (b - a^T x)^{\mu - 1} a$

$\nabla_x \, x^T y = \nabla_x \, y^T x = y$

$\nabla_X \operatorname{tr}(X^T Y) = \nabla_X \operatorname{tr}(Y X^T) = \nabla_X \operatorname{tr}(Y^T X) = \nabla_X \operatorname{tr}(X Y^T) = Y$

$\nabla_X \operatorname{tr}(A X B X^T) = \nabla_X \operatorname{tr}(X B X^T A) = A^T X B^T \ + \ A X B$

$\nabla_X \operatorname{tr}(A X B X) \ \ = \nabla_X \operatorname{tr}(X B X A) \ \ = A^T X^T B^T + B^T X^T A^T$

$\nabla_X \operatorname{tr}(A X A X A X) = \nabla_X \operatorname{tr}(X A X A X A) = 3 (A X A X A)^T$

$\nabla_X \operatorname{tr}(Y X^k) = \nabla_X \operatorname{tr}(X^k Y) = \sum_{i=0}^{k-1} \big(X^i Y X^{k-1-i}\big)^T$

$\nabla_X \operatorname{tr}(Y^T X X^T Y) = \nabla_X \operatorname{tr}(X^T Y Y^T X) = 2 Y Y^T X$

$\nabla_X \operatorname{tr}(Y^T X^T X Y) \ = \nabla_X \operatorname{tr}(X Y Y^T X^T) = 2 X Y Y^T$

$\nabla_X \operatorname{tr}\big((X + Y)^T (X + Y)\big) = 2(X + Y)$

$\nabla_X \operatorname{tr}\big((X + Y)(X + Y)\big) \ \ = 2(X + Y)^T$

$\nabla_X \operatorname{tr}(A^T X B) \ \ \ \ = \nabla_X \operatorname{tr}(X^T A B^T) \ \ = \ \ \ \ \ \ \ \ A B^T$

$\nabla_X \operatorname{tr}(A^T X^{-1} B) = \nabla_X \operatorname{tr}(X^{-T} A B^T) = -X^{-T} A B^T X^{-T}$

$\nabla_X \, a^T X b \ \ = \nabla_X \operatorname{tr}(b a^T X) \ \ = \nabla_X \operatorname{tr}(X b a^T) \ \ = a b^T$

$\nabla_X \, b^T X^T a \ = \nabla_X \operatorname{tr}(X^T a b^T) = \nabla_X \operatorname{tr}(a b^T X^T) = a b^T$

$\nabla_X \, a^T X^{-1} b \ = \nabla_X \operatorname{tr}(X^{-T} a b^T) = -X^{-T} a b^T X^{-T}$

$\nabla_X \, a^T X^{\mu} b = \dots$

**Trace** continued

$$\tfrac{d}{dt}\operatorname{tr} g(X+t\,Y) = \operatorname{tr}\tfrac{d}{dt}\,g(X+t\,Y)$$

$$\tfrac{d}{dt}\operatorname{tr}(X+t\,Y) = \operatorname{tr} Y$$

$$\tfrac{d}{dt}\operatorname{tr}^j(X+t\,Y) = j\operatorname{tr}^{j-1}(X+t\,Y)\operatorname{tr} Y$$

$$\tfrac{d}{dt}\operatorname{tr}(X+t\,Y)^j = j\operatorname{tr}((X+t\,Y)^{j-1}\,Y) \qquad\qquad (\forall\, j)$$

$$\tfrac{d}{dt}\operatorname{tr}((X+t\,Y)Y) = \operatorname{tr} Y^2$$

$$\tfrac{d}{dt}\operatorname{tr}\big((X+t\,Y)^k\,Y\big) = \tfrac{d}{dt}\operatorname{tr}(Y(X+t\,Y)^k) = k\operatorname{tr}\big((X+t\,Y)^{k-1}Y^2\big)\,, \quad k\in\{0,1,2\}$$

$$\tfrac{d}{dt}\operatorname{tr}\big((X+t\,Y)^k\,Y\big) = \tfrac{d}{dt}\operatorname{tr}(Y(X+t\,Y)^k) = \operatorname{tr}\sum_{i=0}^{k-1}(X+t\,Y)^i\,Y(X+t\,Y)^{k-1-i}\,Y$$

$$\tfrac{d}{dt}\operatorname{tr}((X+t\,Y)^{-1}Y) \quad\;\; = -\operatorname{tr}((X+t\,Y)^{-1}Y(X+t\,Y)^{-1}Y)$$
$$\tfrac{d}{dt}\operatorname{tr}\big(B^T(X+t\,Y)^{-1}A\big) = -\operatorname{tr}\big(B^T(X+t\,Y)^{-1}Y(X+t\,Y)^{-1}A\big)$$
$$\tfrac{d}{dt}\operatorname{tr}\big(B^T(X+t\,Y)^{-T}A\big) = -\operatorname{tr}\big(B^T(X+t\,Y)^{-T}Y^T(X+t\,Y)^{-T}A\big)$$
$$\tfrac{d}{dt}\operatorname{tr}\big(B^T(X+t\,Y)^{-k}A\big) = ...\,, \quad k>0$$
$$\tfrac{d}{dt}\operatorname{tr}\big(B^T(X+t\,Y)^{\mu}A\big) \;\; = ...\,, \quad -1\le\mu\le 1,\;\; X,Y\in\mathbb{S}_+^M$$

$$\tfrac{d^2}{dt^2}\operatorname{tr}\big(B^T(X+t\,Y)^{-1}A\big) = 2\operatorname{tr}\big(B^T(X+t\,Y)^{-1}Y(X+t\,Y)^{-1}Y(X+t\,Y)^{-1}A\big)$$

$$\tfrac{d}{dt}\operatorname{tr}\big((X+t\,Y)^T A(X+t\,Y)\big) = \operatorname{tr}\big(Y^T AX + X^T AY + 2t\,Y^T AY\big)$$

$$\tfrac{d^2}{dt^2}\operatorname{tr}\big((X+t\,Y)^T A(X+t\,Y)\big) = 2\operatorname{tr}\big(Y^T AY\big)$$

$$\tfrac{d}{dt}\operatorname{tr}\big((X+t\,Y)A(X+t\,Y)\big) = \operatorname{tr}(YAX + XAY + 2t\,YAY)$$

$$\tfrac{d^2}{dt^2}\operatorname{tr}\big((X+t\,Y)A(X+t\,Y)\big) = 2\operatorname{tr}(YAY)$$

## D.2.4  Log determinant

$x \succ 0$, $\det X > 0$ on some neighborhood of $X$, and $\det(X + tY) > 0$ on some open interval of $t$; otherwise, $\log(\,)$ would be discontinuous.

| | |
|---|---|
| $\frac{d}{dx} \log x = x^{-1}$ | $\nabla_X \log \det X = X^{-T}$ <br><br> $\nabla_X^2 \log \det(X)_{kl} = \dfrac{\partial X^{-T}}{\partial X_{kl}} = -(X^{-1} E_{kl} X^{-1})^T$, *confer* (1405)(1447) |
| $\frac{d}{dx} \log x^{-1} = -x^{-1}$ | $\nabla_X \log \det X^{-1} = -X^{-T}$ |
| $\frac{d}{dx} \log x^\mu = \mu x^{-1}$ | $\nabla_X \log \det^\mu X = \mu X^{-T}$ |
| | $\nabla_X \log \det X^\mu = \mu X^{-T}$, $\qquad\qquad\qquad\qquad X \in \mathbb{R}^{\mathbf{2 \times 2}}$ |
| | $\nabla_X \log \det X^k = \nabla_X \log \det^k X = k X^{-T}$ |
| | $\nabla_X \log \det^\mu(X + tY) = \mu(X + tY)^{-T}$ |
| $\nabla_x \log(a^T x + b) = a \frac{1}{a^T x + b}$ | $\nabla_X \log \det(AX + B) = A^T(AX + B)^{-T}$ |
| | $\nabla_X \log \det(I \pm A^T X A) = \ldots$ |
| | $\nabla_X \log \det(X + tY)^k = \nabla_X \log \det^k(X + tY) = k(X + tY)^{-T}$ |
| | $\frac{d}{dt} \log \det(X + tY) = \mathrm{tr}\left((X + tY)^{-1} Y\right)$ |
| | $\frac{d^2}{dt^2} \log \det(X + tY) = -\mathrm{tr}\left((X + tY)^{-1} Y (X + tY)^{-1} Y\right)$ |
| | $\frac{d}{dt} \log \det(X + tY)^{-1} = -\mathrm{tr}\left((X + tY)^{-1} Y\right)$ |
| | $\frac{d^2}{dt^2} \log \det(X + tY)^{-1} = \mathrm{tr}\left((X + tY)^{-1} Y (X + tY)^{-1} Y\right)$ |
| | $\frac{d}{dt} \log \det\left(\delta(A(x + ty) + a)^2 + \mu I\right)$ <br> $\quad = \mathrm{tr}\left(\left(\delta(A(x + ty) + a)^2 + \mu I\right)^{-1} 2\delta(A(x + ty) + a)\delta(Ay)\right)$ |

### D.2.5    Determinant

$\nabla_X \det X = \nabla_X \det X^T = \det(X) X^{-T}$

$\nabla_X \det X^{-1} = -\det(X^{-1}) X^{-T} = -\det(X)^{-1} X^{-T}$

$\nabla_X \det^\mu X = \mu \det^\mu(X) X^{-T}$

$\nabla_X \det X^\mu = \mu \det(X^\mu) X^{-T} \,,$                                        $X \in \mathbb{R}^{\mathbf{2 \times 2}}$

$\nabla_X \det X^k = k \det^{k-1}(X) \big( \operatorname{tr}(X) I - X^T \big) \,,$                      $X \in \mathbb{R}^{\mathbf{2 \times 2}}$

$\nabla_X \det X^k = \nabla_X \det^k X = k \det(X^k) X^{-T} = k \det^k(X) X^{-T}$

$\nabla_X \det^\mu(X + t\, Y) = \mu \det^\mu(X + t\, Y)(X + t\, Y)^{-T}$

$\nabla_X \det(X + t\, Y)^k = \nabla_X \det^k(X + t\, Y) = k \det^k(X + t\, Y)(X + t\, Y)^{-T}$

$\frac{d}{dt} \det(X + t\, Y) = \det(X + t\, Y) \operatorname{tr}((X + t\, Y)^{-1} Y)$

$\frac{d^2}{dt^2} \det(X + t\, Y) = \det(X + t\, Y) \big( \operatorname{tr}^2((X + t\, Y)^{-1} Y) - \operatorname{tr}((X + t\, Y)^{-1} Y (X + t\, Y)^{-1} Y) \big)$

$\frac{d}{dt} \det(X + t\, Y)^{-1} = -\det(X + t\, Y)^{-1} \operatorname{tr}((X + t\, Y)^{-1} Y)$

$\frac{d^2}{dt^2} \det(X + t\, Y)^{-1} = \det(X + t\, Y)^{-1} \big( \operatorname{tr}^2((X + t\, Y)^{-1} Y) + \operatorname{tr}((X + t\, Y)^{-1} Y (X + t\, Y)^{-1} Y) \big)$

$\frac{d}{dt} \det^\mu(X + t\, Y) = \ldots$

## D.2.6   Logarithmic

$$\frac{d}{dt}\log(X+tY)^{\mu} = \dots, \quad -1 \le \mu \le 1, \quad X, Y \in \mathbb{S}_+^M \quad [128, \S 6.6, \text{prob.20}]$$

## D.2.7   Exponential

[45, §3.6, §4.5] [215, §5.4]

$$\nabla_X e^{\text{tr}(Y^T X)} = \nabla_X \det e^{Y^T X} = e^{\text{tr}(Y^T X)} Y \qquad\qquad (\forall\, X, Y)$$

$$\nabla_X \text{tr}\, e^{YX} = e^{Y^T X^T} Y^T = Y^T e^{X^T Y^T}$$

log-sum-exp & geometric mean [39, p.74]...

$$\frac{d^j}{dt^j} e^{\text{tr}(X+tY)} = e^{\text{tr}(X+tY)} \text{tr}^j(Y)$$

$$\frac{d}{dt} e^{tY} = e^{tY} Y = Y e^{tY}$$

$$\frac{d}{dt} e^{X+tY} = e^{X+tY} Y = Y e^{X+tY}, \qquad\qquad XY = YX$$

$$\frac{d^2}{dt^2} e^{X+tY} = e^{X+tY} Y^2 = Y e^{X+tY} Y = Y^2 e^{X+tY}, \quad XY = YX$$

$e^X$ for symmetric $X$ of dimension less than 3 [39, pg.110]...