

矩阵微分

花半秒钟就看透事物本质的人,和花一辈子都看不清事物本质的人,注定是截然不同的命运。——马里奥·普佐

在一般的线性代数中，主要是用代数的方法来研究矩阵，没有涉及到极限和微积分的运算规律。为了书写简便，我们通常把单个函数对多个变量或者多元函数对单个变量的偏导数**写成向量和矩阵的形式，使其可以被当做一个整体处理**。矩阵微积分是多元微积分的**一种表达方式**，即使用矩阵和向量来表示因变量每个成分关于自变量每个成分的偏导数。

矩阵分析理论和数学分析一样，是建立在极限的概念之上。

向量序列的极限

对于向量序列 $\{X^{(k)}\}$ ，其中 $X^{(k)} = (x_1^{(k)}, x_2^{(k)}, \dots, x_n^{(k)})$ ，若其每一个分量 $x_i^{(k)}$ 当 $k \rightarrow \infty$ ，都有极限 x_i ，即 $\lim_{k \rightarrow \infty} x_i^{(k)} = x_i, i = 1, 2, \dots, n$ ，则称向量序列 $\{X^{(k)}\}$ 有极限 $X = (x_1, x_2, \dots, x_n)$ ，或称 $\{X^{(k)}\}$ 收敛于 X ，记为 $\lim_{k \rightarrow \infty} X^{(k)} = X$ 。

由于实数序列的极限是唯一的，所以向量序列的极限必定是唯一的。

矩阵序列的极限

矩阵序列的收敛性和向量序列类似。

设有同阶矩阵序列 $\{A_k\} = \{[a_{ij}^k]_{m \times n}\}$ ，若存在矩阵 $A = [a_{ij}]_{m \times n}$ ，当 $k \rightarrow \infty$ 时， a_{ij}^k 收敛于 a_{ij} 。则称矩阵序列 $\{A_k\}$ 收敛于 A 。即 $\lim_{k \rightarrow \infty} A_k = A$ 。矩阵 A 称为矩阵序列 $\{A_k\}$ 的极限。

收敛的矩阵序列与收敛的标量数列有类似的性质。

函数矩阵

设矩阵

$$A(x) = \begin{bmatrix} a_{11}(x) & a_{12}(x) & \cdots & a_{1n}(x) \\ a_{21}(x) & a_{22}(x) & \cdots & a_{2n}(x) \\ \cdots & \cdots & \cdots & \cdots \\ a_{m1}(x) & a_{m2}(x) & \cdots & a_{mn}(x) \end{bmatrix}$$

其中每个元素都是实变量 x 的函数，则称矩阵 $A(x)$ 为函数矩阵。

若对每一个元素 $a_{ij}(x)$ ，当 $x \rightarrow x_0$ 时，都存在着极限 a_{ij} （常数），则称矩阵 $A(x)$ 在 x_0 处有极限，且极限值为矩阵 $A = [a_{ij}]_{m \times n}$ 。

函数矩阵的极限，具有函数极限的类似性质。

函数矩阵的微分

本节所研究的内容：是用**矩阵**来描述微积分中的若干结果，在工程实际中常见的三个问题，即是**函数矩阵关于自变量的微分和积分**；**标量函数关于矩阵的微分**；**向量函数关于向量的微分**。

标量函数关于向量的偏导数

对于向量 \boldsymbol{x} 和函数 $y = f(\boldsymbol{x})$ ，若 y 对于 \boldsymbol{x} 各元素均可微分，则 y 关于 \boldsymbol{x} 的偏导数为 $\frac{\partial y}{\partial \boldsymbol{x}} = [\frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_M}]^T$ 。

至于为什么写成转置的形式，是因为这里才用了**分母布局**，即列向量的形式（分子布局：行向量的形式），以下同理。

向量函数关于标量的偏导数

对于标量 x 和函数 $\boldsymbol{y} = f(x)$ ，若 \boldsymbol{y} 中元素对于 x 均可微分，则 \boldsymbol{y} 关于 x 的偏导数为 $[\frac{\partial y_1}{\partial x}, \dots, \frac{\partial y_N}{\partial x}]^T$ 。

向量函数关于向量的偏导数

不难想象，向量函数关于向量的偏导数理所应当是一个矩阵：

对于向量 \boldsymbol{x} 和函数 $\boldsymbol{y} = f(\boldsymbol{x})$ ，则 $f(\boldsymbol{x})$ 关于 \boldsymbol{x} 的偏导数为

$$\frac{\partial f(\boldsymbol{x})}{\partial \boldsymbol{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_N}{\partial x_1} \\ \dots & \dots & \dots \\ \frac{\partial y_1}{\partial x_M} & \dots & \frac{\partial y_N}{\partial x_M} \end{bmatrix}$$

此为函数 $f(\boldsymbol{x})$ 的**jacobi**矩阵的转置。

标量函数关于矩阵的偏导数

设函数 $f(\boldsymbol{x})$ 是以矩阵 $\boldsymbol{X} = [x_{ij}]_{m \times n}$ 中的 $m \times n$ 元素 x_{ij} 为自变量的**可微标量**函数，即

$$f(\boldsymbol{X}) = f(x_{11}, x_{12}, \dots, x_{1n}; x_{21}, x_{22}, \dots, x_{2n}; \dots; x_{m1}, x_{m2}, \dots, x_{mn})$$

则 $f(\boldsymbol{x})$ 对矩阵 \boldsymbol{X} 的导数有如下定义：

设矩阵 $\boldsymbol{X} = [x_{ij}]_{m \times n}$ ，若标量函数 $f(\boldsymbol{x})$ 对自变量 x_{ij} 可微，则定义 f 对矩阵 \boldsymbol{X} 的导数如下：

$$\frac{d}{d\boldsymbol{X}} = \begin{bmatrix} \frac{\partial f}{\partial x_{11}} & \frac{\partial f}{\partial x_{12}} & \dots & \frac{\partial f}{\partial x_{1n}} \\ \frac{\partial f}{\partial x_{21}} & \frac{\partial f}{\partial x_{22}} & \dots & \frac{\partial f}{\partial x_{2n}} \\ \dots & \dots & \dots & \dots \\ \frac{\partial f}{\partial x_{m1}} & \frac{\partial f}{\partial x_{m2}} & \dots & \frac{\partial f}{\partial x_{mn}} \end{bmatrix}$$

或 $\frac{\partial f}{\partial \boldsymbol{X}} = [\frac{\partial f}{\partial x_{ij}}]_{m \times n}$ 。即**标量函数 f 对矩阵 \boldsymbol{X} 逐元素求导数而已**。特别地，以 $\boldsymbol{x} = (x_1, x_2, \dots, x_n)^T$ 为自变量

的函数 $f(\boldsymbol{x})$ 的导数 $\frac{df}{d\boldsymbol{x}} = (\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n})$ 称为**数量函数对向量变量的导数**。

设矩阵

$$\boldsymbol{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \dots & \dots & \dots & \dots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}$$

和标量函数

$$f(X) = x_{11}^2 + x_{12}^2 + \cdots + x_{1n}^2 + x_{21}^2 + x_{22}^2 + \cdots + x_{2n}^2 + \cdots + x_{m1}^2 + x_{m2}^2 + \cdots + x_{mn}^2$$

求 $\frac{df}{dX}$ 。

解: 因为 $\frac{\partial f}{\partial x_{i,j}} = 2x_{i,j} (i = 1, 2, \cdots, m; j = 1, 2, \cdots, n)$, 所以

$$\frac{d}{dX} f(x) = 2 \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ x_{m1} & x_{m2} & \cdots & x_{mn} \end{bmatrix} = 2X。$$

又因为 $f(X) = \text{tr}(XX^T)$, 由此可以得出一些关于迹的重要公式:

- $\frac{d}{dX} \text{tr}(XX^T) = 2X$
- $\frac{d}{dX} \text{tr}(BX) = \frac{d}{dX} \text{tr}(X^T B^T) = B^T$
- $\frac{d}{dX} \text{tr}(X^T AX) = (A + A^T)X$

其中 $X = [x_{ij}]_{m \times n}$; $B = [b_{ij}]_{n \times m}$; $A = [a_{ij}]_{m \times m}$ 。

此外, 还可以验证对于标量函数 $f(X)$ 、 $g(X)$ 有如下运算法则:

$$\begin{aligned} \frac{d}{dX} (f + g) &= \frac{df}{dX} + \frac{dg}{dX} \\ \frac{d}{dX} fg &= \frac{df}{dX} g + f \frac{dg}{dX} \end{aligned}$$

公式中的乘积可以交换。

例 1 设 $a = (a_1, a_2, \cdots, a_n)^T$ 为给定的向量, $x = (x_1, x_2, \cdots, x_n)^T$ 是向量变量, 且

$$f(x) = a^T x = x^T a,$$

求 $\frac{df}{dx}$ 。

解 由 $f(x) = \sum_{i=1}^n a_i x_i$ 得

$$\frac{\partial f}{\partial x_i} = a_i, i = 1, 2, \cdots, n,$$

所以

$$\frac{df}{dx} = (a_1, a_2, \cdots, a_n)^T = a。$$

例 2 设 $A = (a_{ij})_{n \times n}$ 为给定的矩阵, $x = (x_1, \dots, x_n)^T$ 是向量变量, $f(x) = x^T A x$, 求 $\frac{df}{dx}$.

解 由 $f(x) = x^T A x = \sum_{s=1}^n \sum_{k=1}^n x_s a_{sk} x_k$,

$$\frac{\partial f}{\partial x_j} = \sum_{s=1}^n a_{sj} x_s + \sum_{k=1}^n a_{jk} x_k,$$

所以

$$\frac{df}{dx} = A^T x + A x = (A^T + A) x.$$

特别地, 当 A 是对称矩阵时, 有

$$\frac{df}{dx} = 2Ax.$$

函数矩阵关于标量的偏导数

设函数矩阵 $A(x) = [a_{ij}(x)]_{m \times n}$ ($i = 1, 2, \dots, m; j = 1, 2, \dots, n$), 如果对于所有元素 $a_{ij}(x)$ 在 $x = x_0$ 点或某一区间上是可微的, 则称该函数矩阵 $A(x)$ 在 $x = x_0$ 点或该区间上是可微的。并定义

$$\frac{d}{dx} A(x) = A'(x) = \begin{bmatrix} a'_{11}(x) & a'_{12}(x) & \cdots & a'_{1n}(x) \\ a'_{21}(x) & a'_{22}(x) & \cdots & a'_{2n}(x) \\ \cdots & \cdots & \cdots & \cdots \\ a'_{m1}(x) & a'_{m2}(x) & \cdots & a'_{mn}(x) \end{bmatrix}$$

为函数矩阵 $A(x)$ 对 x 的导数, 即是函数矩阵 A 逐元素对标量 x 求导数而已。

由函数矩阵的导数定义, 可以推算出关于函数矩阵的导数运算法则。

设函数矩阵 $A(x)$ 与 $B(x)$ 是可微的, 则有:

- **矩阵加法:** $A(x) + B(x)$ 也是可微的, 且满足 $\frac{d}{dx} [A(x) + B(x)] = \frac{d}{dx} A(x) + \frac{d}{dx} B(x)$
- **标量乘法:** $\frac{d}{dx} [kA(x)] = k \frac{d}{dx} A(x)$, 其中 k 是任意常数。
- **标量函数乘法:** $\frac{d(aB)}{dx} = (\frac{da}{dx})B + a(\frac{dB}{dx})$ 。其中 a 是 x 的标量函数, 且对 x 可微分。
- **矩阵乘法:** 在 $A(x)B(x)$ 有定义的情况下, $A(x)B(x)$ 也是可微的, 且有 $\frac{d}{dx} [A(x)B(x)] = [\frac{d}{dx} A(x)]B(x) + A(x)\frac{d}{dx} B(x)$ 。公式中的乘积次序是不能交换的。
- **复合:** 设矩阵 A 是 u 的函数, $u = f(x)$, 则有 $\frac{d}{dx} [A[f(x)]] = \frac{d}{du} A(u) \cdot f'(x)$ 或者 $\frac{d}{dx} [A[f(x)]] = \frac{d}{du} f'(x) \cdot A(u)$ 。
- **逆:** 若矩阵 $A(x)$ 为正则矩阵, 则其逆矩阵 $A^{-1}(x)$ 也可微分, $\frac{d}{dx} A^{-1}(x) = -A^{-1}(x) \frac{dA(x)}{dx} A^{-1}(x)$ 。

常见向量函数的导数

$$\frac{\partial \boldsymbol{x}}{\partial \boldsymbol{x}} = \boldsymbol{I}$$

$$\frac{\partial \boldsymbol{x}^T \boldsymbol{a}}{\partial \boldsymbol{x}} = \frac{\partial \boldsymbol{a}^T \boldsymbol{x}}{\partial \boldsymbol{x}} = \boldsymbol{a}$$

$$\frac{\partial \boldsymbol{A} \boldsymbol{x}}{\partial \boldsymbol{x}} = \boldsymbol{A}^T$$

$$\frac{\partial \boldsymbol{x}^T \boldsymbol{A}}{\partial \boldsymbol{x}} = \boldsymbol{A}$$

矩阵求导

本文使用小写字母 x 表示标量，粗体小写字母 \mathbf{x} 表向量，大写粗体字母 \mathbf{X} 表示矩阵。

首先来琢磨一下定义，**标量 f 对矩阵 \mathbf{X} 的导数**，定义为 $\frac{\partial f}{\partial \mathbf{X}} = [\frac{\partial f}{\partial X_{ij}}]$ ，即 f 对 \mathbf{X} 逐元素求导排成与 \mathbf{X} 尺寸相同的矩阵。

然而，这个定义在计算中并不好用，实用上的原因是对函数较复杂的情形**难以逐元素求导**；哲理上的原因是逐元素求导破坏了矩阵**整体性**。试想，为何我们要改变常见想法，将 f 看做矩阵 \mathbf{X} 而不是各元素 X_{ij} 的函数呢？答案是用**矩阵运算更简便整洁**，**软件工具对矩阵也有相关优化**。所以在求导时**不宜拆开矩阵，而是要找一个从整体出发的算法**。

为此，我们来回顾，一元微积分中的导数（**标量对标量的导数**）与微分的联系： $df = f'(x)dx$ 。

多元微积分中的梯度（**标量对向量的导数**）也与微分有联系：

$$df = \sum_{i=1}^n \frac{\partial f}{\partial x_i} dx_i = \left(\frac{\partial f}{\partial \mathbf{x}} \right) d\mathbf{x}$$

值得注意的是： $\sum ab$ 往往可以转化成向量内积的形式，对应程序中的实现就是 for 循环变为向量运算。而 $\sum \sum ab$ 往往可以转化成矩阵内积的形式，对应程序中的实现就是双层 for 循环变为矩阵运算。

- 第一个等号是全微分公式（对每个分量 x_i 分别求微分，然后相加）
- 第二个等号表达了**梯度与微分**的联系：全微分 df 是梯度向量 $\frac{\partial f}{\partial \mathbf{x}} (n \times 1)$ ，与微分向量 $d\mathbf{x} (n \times 1)$ 的**内积**；受此启发，我们将矩阵导数与微分建立联系：

$$df = \sum_{i=1}^m \sum_{j=1}^n \frac{\partial f}{\partial X_{ij}} dX_{ij} = tr\left(\frac{\partial f}{\partial \mathbf{X}} d\mathbf{X}\right)$$

其中 tr 代表迹：方阵对角线元素之和。满足性质：对尺寸相同的矩阵 \mathbf{A}, \mathbf{B} ， $tr(\mathbf{A}^T \mathbf{B}) = \sum_{i,j} \mathbf{A}_{i,j} \mathbf{B}_{i,j}$ ，即

$tr(\mathbf{A}^T \mathbf{B})$ 是矩阵 \mathbf{A}, \mathbf{B} 的**内积**。与梯度相似，这里第一个等号是全微分公式，第二个等号表达了**矩阵导数与微分**的联系：全微分 df 是导数 $\frac{\partial f}{\partial \mathbf{X}} (m \times n)$ 与微分矩阵 $d\mathbf{X} (m \times n)$ 的内积。

- 补充：
 - 设矩阵 $\mathbf{A} = (a_{ij})$ ，把矩阵 \mathbf{A} 的元素按行的顺序排列成一个列向量， $vec \mathbf{A} = (a_{11}, a_{12}, \dots, a_{1n}, a_{21}, a_{22}, \dots, a_{2n}, \dots, a_{m1}, a_{m2}, \dots, a_{mn})^T$ ，则称向量 $vec \mathbf{A}$ 为矩阵 \mathbf{A} 按行拉直的列向量。
 - 设矩阵 \mathbf{A}, \mathbf{B} ，称 $\mathbf{A} \cdot \mathbf{B} = \langle \mathbf{A}, \mathbf{B} \rangle = tr(\mathbf{A}^T \mathbf{B}) = \sum_{i=1}^n \sum_{j=1}^n a_{ij} b_{ij} = (vec \mathbf{A}^T vec \mathbf{B})$ 为矩阵 \mathbf{A}, \mathbf{B} 的**内积**。

运算法则

回想遇到较复杂的一元函数如 $f(x) = \log(2 + \sin x)e^{\sqrt{x}}$ ，我们是如何求导的呢？通常不是从定义开始由极限来求出结果，而是先求出了**初等函数的导数**，然后建立**四则运算、复合函数等法则**，根据两者求出复杂函数的导数。故而，我们来建立常用的矩阵微分的运算法则。

- **加减法**： $d(\mathbf{X} \pm \mathbf{Y}) = d\mathbf{X} \pm d\mathbf{Y}$ ；
- **矩阵乘法**： $d(\mathbf{XY}) = (d\mathbf{X})\mathbf{Y} + \mathbf{X}d(\mathbf{Y})$ ；**注意乘法因数不能改变位置**。
- **转置**： $d(\mathbf{X}^T) = (d\mathbf{X})^T$ ；
- **迹**： $d(tr(\mathbf{X})) = tr(d\mathbf{X})$ 。
- **逆运算**： $d\mathbf{X}^{-1} = -\mathbf{X}^{-1}(d\mathbf{X})\mathbf{X}^{-1}$ 。此式可在 $\mathbf{XX}^{-1} = \mathbf{I}$ 两侧求微分来证明。

- **行列式**: $d|X| = \text{tr}(X^* dX)$, 其中 X^* 表示 X 的伴随矩阵, 在 X 可逆时又可以写作 $d|X| = |X| \text{tr}(X^{-1} dX)$ 。此式可用 *Laplace* 展开来证明, 见张贤达《矩阵分析与应用》第279页。
- **逐元素乘法**: $d(X \odot Y) = dX \odot Y + X \odot dY$, \odot 表示尺寸相同的矩阵 X, Y 逐元素相乘, 即 *Hadamard* 积。
- **逐元素函数**: $d\sigma(X) = \sigma'(X) \odot dX$, $\sigma(X) = [\sigma(X_{ij})]$ 是逐元素标量函数运算, $\sigma'(X) = [\sigma'(X_{ij})]$ 是逐元素求导数。例如

$$X = \begin{bmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \end{bmatrix}, d\sin(X) = \begin{bmatrix} \cos X_{11} dX_{11} & \cos X_{12} dX_{12} \\ \cos X_{21} dX_{21} & \cos X_{22} dX_{22} \end{bmatrix} = \cos(X) \odot dX$$

我们试图利用矩阵导数与微分的联系 $df = \text{tr}(\frac{\partial f^T}{\partial X} dX)$, 在求出左侧的微分 df 后, 该如何写成右侧的形式并得到导数呢? 这需要一些迹技巧:

- **标量的迹**: $a = \text{tr}(a)$
- **转置**: $\text{tr}(A^T) = \text{tr}(A)$
- **线性**: $\text{tr}(A \pm B) = \text{tr}(A) \pm \text{tr}(B)$
- **矩阵乘法交换**: $\text{tr}(AB) = \text{tr}(BA) = \sum_{i,j} A_{ij} B_{ji}$, 其中 A 与 B^T 尺寸大小相同。
- **矩阵乘法/逐元素乘法交换**: $\text{tr}(A^T(B \odot C)) = \text{tr}((A \odot B)^T C)$, 其中 A, B, C 尺寸相同。两侧都等于 $\sum_{i,j} A_{ij} B_{ij} C_{ij}$

方法思想:

观察一下可以断言, 若标量函数 f 是矩阵 X 经加减乘法、逆、行列式、逐元素函数等运算构成, 则**使用相应的运算法则对 f 求微分**, 再使用迹技巧给 df 套上迹并将其它项交换至等式右侧的 dX 左侧, 对照导数与微分的联系

$df = \text{tr}(\frac{\partial f^T}{\partial X} dX)$, 就能得到导数。

特别地, 若矩阵退化为向量, 对照导数与微分的联系 $df = \frac{\partial f^T}{\partial x} dx$, 就能得到导数。

在建立法则的最后, 来谈一谈**复合**: 假设已求得 $\frac{\partial f}{\partial Y}$, 而 Y 是 X 的函数, 如何求 $\frac{\partial f}{\partial X}$ 呢? 在微积分中有标量求导的链式法则 $\frac{\partial f}{\partial x} = \frac{\partial f}{\partial y} \frac{\partial y}{\partial x}$, 但这里我们**不能随意沿用标量的链式法则**, 因为矩阵对矩阵的导数 $\frac{\partial Y}{\partial X}$ 截至目前仍是未定义的。于是我们继续追本溯源, 链式法则是从何而来? 源头仍然是微分。

故我们直接从微分入手建立复合法则: 先写出 $df = \text{tr}(\frac{\partial f^T}{\partial Y} dY)$, 再将 dY 用 dX 表示出来代入, 并使用迹技巧将其他项交换至 dX 左侧, 即可得到 $\frac{\partial f}{\partial X}$ 。

最常见的情形是 $Y = AXB$, 此时

$$df = \text{tr}(\frac{\partial f^T}{\partial Y} dY) = \text{tr}(\frac{\partial f^T}{\partial Y} A dX B) = \text{tr}(B \frac{\partial f^T}{\partial Y} A dX) = \text{tr}((A^T \frac{\partial f}{\partial Y} B^T)^T dX)$$

可得到 $\frac{\partial f}{\partial X} = A^T \frac{\partial f}{\partial Y} B^T$ 。注意推导过程中的 $dY = (dA)XB + AdXB = AdXB$, 由于 A, B 是常量, $dA = 0, dB = 0$, 以及我们使用矩阵乘法交换的迹技巧交换了 $\frac{\partial f^T}{\partial Y} AdX$ 与 B 。

演示例子

接下来演示一些算例。特别提醒要依据已经建立的运算法则来计算，不能随意套用微积分中标量导数的结论，比如认为 AX 对 X 的导数为 A ，这是没有根据、意义不明的。

- 例1. $f = \mathbf{a}^T \mathbf{X} \mathbf{b}$, 求 $\frac{\partial f}{\partial \mathbf{X}}$ 。其中 \mathbf{a} 是 $m \times 1$ 的列向量， \mathbf{X} 是 $m \times n$ 的矩阵， \mathbf{b} 是 $n \times 1$ 列向量， f 是标量。

- 先使用矩阵乘法法则求微分：

$df = (d\mathbf{a}^T) \mathbf{X} \mathbf{b} + \mathbf{a}^T (d\mathbf{X}) \mathbf{b} + \mathbf{a}^T \mathbf{X} (d\mathbf{b}) = \mathbf{a}^T (d\mathbf{X}) \mathbf{b}$, 注意这里的 \mathbf{a}, \mathbf{b} 是常量，所以 $d\mathbf{a}^T = 0, d\mathbf{b} = 0$ 。由于 df 是标量，它的迹等于自身， $df = \text{tr}(df)$ 。

- 套上迹并做矩阵乘法交换：

$df = \text{tr}(\mathbf{a}^T d\mathbf{X} \mathbf{b}) = \text{tr}(\mathbf{b} \mathbf{a}^T d\mathbf{X}) = \text{tr}((\mathbf{a} \mathbf{b}^T)^T d\mathbf{X})$ 。对照导数与微分的联系

$df = \text{tr}(\frac{\partial f}{\partial \mathbf{X}} d\mathbf{X})$, 得到 $\frac{\partial f}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^T$ 。

- 注意：这里不能用 $\frac{\partial f}{\partial \mathbf{X}} = \mathbf{a}^T \frac{\partial \mathbf{X}}{\partial \mathbf{X}} \mathbf{b} = ?$ ，导数与矩阵乘法的交换是不合法则的运算（而微分是合法的）。有些资料在计算矩阵导数时，会略过求微分这一步，这在逻辑上是解释不通的。
-

- 例2. $f = \mathbf{a}^T e^{\mathbf{X} \mathbf{b}}$, 求 $\frac{\partial f}{\partial \mathbf{X}}$ 。其中 \mathbf{a} 是 $m \times 1$ 的列向量， \mathbf{X} 是 $m \times n$ 的矩阵， \mathbf{b} 是 $n \times 1$ 列向量， $e^{\mathbf{X} \mathbf{b}}$ 是逐元素求指数， f 是标量。

- 先使用矩阵乘法、逐元素函数法则求微分：

$df = \mathbf{a}^T (e^{\mathbf{X} \mathbf{b}} \odot (d\mathbf{X} \mathbf{b}))$ 。

- 再套上迹并做交换：

$df = \text{tr}(\mathbf{a}^T e^{\mathbf{X} \mathbf{b}} \odot (d\mathbf{X} \mathbf{b})) = \text{tr}((\mathbf{a} \odot e^{\mathbf{X} \mathbf{b}}) d\mathbf{X} \mathbf{b}) = \text{tr}(\mathbf{b} (\mathbf{a} \odot e^{\mathbf{X} \mathbf{b}})^T d\mathbf{X}) = \text{tr}(((\mathbf{a} \odot e^{\mathbf{X} \mathbf{b}}) \mathbf{b}^T)^T d\mathbf{X})$, 注意这里我们先根据 $\text{tr}(\mathbf{A}^T (\mathbf{B} \odot \mathbf{C})) = \text{tr}((\mathbf{A} \odot \mathbf{B})^T \mathbf{C})$ 交换了 $\mathbf{a}, e^{\mathbf{X} \mathbf{b}}, d\mathbf{X} \mathbf{b}$ ，再根据

$\text{tr}(\mathbf{A} \mathbf{B}) = \text{tr}(\mathbf{B} \mathbf{A})$ 交换了 $(\mathbf{a} \odot e^{\mathbf{X} \mathbf{b}})^T d\mathbf{X}$ 与 \mathbf{b} 。对照导数与微分的联系， $df = \text{tr}(\frac{\partial f}{\partial \mathbf{X}} d\mathbf{X})$ ，得到 $\frac{\partial f}{\partial \mathbf{X}} = (\mathbf{a} \odot e^{\mathbf{X} \mathbf{b}}) \mathbf{b}^T$ 。

- 例3. $f = \text{tr}(\mathbf{Y}^T \mathbf{M} \mathbf{Y}), \mathbf{Y} = \sigma(\mathbf{W} \mathbf{X})$, 求 $\frac{\partial f}{\partial \mathbf{X}}$ 。其中 \mathbf{W} 是 $l \times m$ 矩阵， \mathbf{X} 是 $m \times n$ 矩阵， \mathbf{Y} 是 $l \times n$ 矩阵， \mathbf{M} 是 $l \times l$ 是对称矩阵， σ 是逐元素函数， f 是标量。

- 先求 $\frac{\partial f}{\partial \mathbf{Y}}$ 求微分，使用矩阵乘法、转置法则：

$df = \text{tr}((d\mathbf{Y})^T \mathbf{M} \mathbf{Y}) + \text{tr}(\mathbf{Y}^T \mathbf{M} d\mathbf{Y}) = \text{tr}(\mathbf{Y}^T \mathbf{M}^T d\mathbf{Y}) + \text{tr}(\mathbf{Y}^T \mathbf{M} d\mathbf{Y}) = \text{tr}(\mathbf{Y}^T (\mathbf{M} + \mathbf{M}^T) d\mathbf{Y})$, 对照导数与微分的联系，得到 $\frac{\partial f}{\partial \mathbf{Y}} = (\mathbf{M} + \mathbf{M}^T) \mathbf{Y} = 2\mathbf{M} \mathbf{Y}$ ，注意这里 \mathbf{M} 是对称矩阵。

- 为求 $\frac{\partial f}{\partial \mathbf{X}}$ ，写出 $df = \text{tr}(\frac{\partial f}{\partial \mathbf{Y}} d\mathbf{Y})$ ，再将 $d\mathbf{Y}$ 用 $d\mathbf{X}$ 表示代入，并使用矩阵乘法/逐元素乘法交换：

$df = \text{tr}\left(\frac{\partial f}{\partial \mathbf{Y}} (\sigma'(\mathbf{W} \mathbf{X}) \odot (\mathbf{W} d\mathbf{X}))\right) = \text{tr}\left(\left(\frac{\partial f}{\partial \mathbf{Y}} \odot \sigma'(\mathbf{W} \mathbf{X})\right)^T \mathbf{W} d\mathbf{X}\right)$ ，对照导数与

微分的联系，得到 $\frac{\partial f}{\partial \mathbf{X}} = \mathbf{W}^T \left(\frac{\partial f}{\partial \mathbf{Y}} \odot \sigma(\mathbf{W} \mathbf{X})\right) = \mathbf{W}^T ((2\mathbf{M} \sigma(\mathbf{W} \mathbf{X})) \odot \sigma'(\mathbf{W} \mathbf{X}))$ 。

- 例4. 线性回归: $l = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$, 求 \mathbf{w} 的最小二乘估计, 即求 $\frac{\partial l}{\partial \mathbf{w}}$ 的零点。其中 \mathbf{y} 是 $m \times 1$ 的列向量, \mathbf{X} 是 $m \times n$ 矩阵, \mathbf{w} 是 $n \times 1$ 列向量, l 是标量。

这里是标量对向量的导数, 不过可以把向量看做矩阵的特例。先将向量的模(2范数)改写成向量与自身的内积: $l = (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y})$ 。

- 求微分, 使用矩阵乘法、转置等法则:

$$dl = (\mathbf{X}d\mathbf{w})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}d\mathbf{w}) = 2(\mathbf{X}\mathbf{w} - \mathbf{y})^T\mathbf{X}d\mathbf{w}, \text{ 注意这里}\mathbf{X}d\mathbf{w}\text{和}\mathbf{X}\mathbf{w} - \mathbf{y}\text{是向量, 两个向量的内积满足}\mathbf{u}^T\mathbf{v} = \mathbf{v}^T\mathbf{u}.$$

- 对照导数与微分的联系:

$$dl = \frac{\partial l}{\partial \mathbf{w}} d\mathbf{w}, \text{ 得到 } \frac{\partial l}{\partial \mathbf{w}} = 2\mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{y}). \text{ 令其为0, 即 } \mathbf{X}^T\mathbf{X}\mathbf{w} = \mathbf{X}^T\mathbf{y}, \text{ 得到}\mathbf{w}\text{的最小二乘估计为}\mathbf{w} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

- 例5. 多元logistic回归: $l = -\mathbf{y}^T \log \text{softmax}(\mathbf{W}\mathbf{x})$, 求 $\frac{\partial l}{\partial \mathbf{W}}$ 。其中 \mathbf{y} 是除一个元素为1外其他元素为0的 $m \times 1$ 列向量, \mathbf{W} 是 $m \times n$ 矩阵, \mathbf{x} 是 $n \times 1$ 列向量, l 是标量; \log 表示自然对数, $\text{softmax}(\mathbf{a}) = \frac{e^{\mathbf{a}}}{\mathbf{E}^T e^{\mathbf{a}}}$, 其中 $e^{\mathbf{a}}$ 表示逐元素求指数, \mathbf{E} 代表全1向量。

定义 $\mathbf{a} = \mathbf{W}\mathbf{x}$, 则 $l = -\mathbf{y}^T \log \text{softmax}(\mathbf{a})$, 先同上求出 $\frac{\partial l}{\partial \mathbf{a}} = \text{softmax}(\mathbf{a}) - \mathbf{y}$, 再利用复合法则: $dl = \text{tr}\left(\frac{\partial l}{\partial \mathbf{a}} d\mathbf{a}\right) = \text{tr}\left(\frac{\partial l}{\partial \mathbf{a}} d\mathbf{W}\mathbf{x}\right) = \text{tr}\left(\mathbf{x} \frac{\partial l}{\partial \mathbf{a}} d\mathbf{W}\right)$, 得到 $\frac{\partial l}{\partial \mathbf{W}} = \frac{\partial l}{\partial \mathbf{a}} \mathbf{x}^T$ 。

- 例6. 二层神经网络: $l = \mathbf{y}^T \log \text{softmax}(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}))$, 求 $\frac{\partial l}{\partial \mathbf{W}_1}$ 和 $\frac{\partial l}{\partial \mathbf{W}_2}$ 。其中 \mathbf{y} 是除一个元素为1外其他元素为0的 $m \times 1$ 列向量, \mathbf{W}_2 是 $m \times p$ 矩阵, \mathbf{W}_1 是 $p \times n$ 矩阵, \mathbf{x} 是 $n \times 1$ 列向量, l 是标量; \log 表示自然对数, $\text{softmax}(\mathbf{a}) = \frac{e^{\mathbf{a}}}{\mathbf{E}^T e^{\mathbf{a}}}$, 同上, $\sigma(\mathbf{a})$ 是逐元素sigmoid函数 $\sigma(a) = \frac{1}{1 + e^{-a}}$ 。

定义 $\mathbf{a}_1 = \mathbf{W}_1 \mathbf{x}$, $\mathbf{h}_1 = \sigma(\mathbf{a}_1)$, $\mathbf{a}_2 = \mathbf{W}_2 \mathbf{h}_1$, 则 $l = -\mathbf{y}^T \log \text{softmax}(\mathbf{a}_2)$ 。在前例中已经求出 $\frac{\partial l}{\partial \mathbf{a}_2} = \text{softmax}(\mathbf{a}_2) - \mathbf{y}$ 。使用复合法则,

$$dl = \text{tr}\left(\frac{\partial l}{\partial \mathbf{a}_2} d\mathbf{a}_2\right) = \text{tr}\left(\frac{\partial l}{\partial \mathbf{a}_2} d\mathbf{W}_2 \mathbf{h}_1\right) + \text{tr}\left(\frac{\partial l}{\partial \mathbf{a}_2} \mathbf{W}_2 d\mathbf{h}_1\right), \text{ 使用矩阵乘法交换的迹技巧从第一项得到 } \frac{\partial l}{\partial \mathbf{W}_2} = \frac{\partial l}{\partial \mathbf{a}_2} \mathbf{h}_1^T, \text{ 从第二项得到 } \frac{\partial l}{\partial \mathbf{h}_1} = \mathbf{W}_2^T \frac{\partial l}{\partial \mathbf{a}_2}. \text{ 接下来对第二项继续使用复合法则来求 } \frac{\partial l}{\partial \mathbf{a}_1}, \text{ 并利用矩阵乘法和逐元素乘法交换的迹技巧:}$$

$$dl_2 = \text{tr}\left(\frac{\partial l}{\partial \mathbf{h}_1} d\mathbf{h}_1\right) = \text{tr}\left(\frac{\partial l}{\partial \mathbf{h}_1} (\sigma'((\mathbf{a}_1) \odot \mathbf{a}_1))\right) = \text{tr}\left(\left(\frac{\partial l}{\partial \mathbf{h}_1} \odot \sigma'(\mathbf{a}_1)\right)^T d\mathbf{a}_1\right), \text{ 得到}$$

$$\frac{\partial l}{\partial \mathbf{a}_1} = \frac{\partial l}{\partial \mathbf{h}_1} \odot \sigma'(\mathbf{a}_1). \text{ 为求 } \frac{\partial l}{\partial \mathbf{W}_1}, \text{ 再用一次复合法则:}$$

$$dl_2 = \text{tr}\left(\frac{\partial l}{\partial \mathbf{a}_1} d\mathbf{a}_1\right) = \text{tr}\left(\frac{\partial l}{\partial \mathbf{a}_1} d\mathbf{W}_1 \mathbf{x}\right) = \text{tr}\left(\mathbf{x} \frac{\partial l}{\partial \mathbf{a}_1} d\mathbf{W}_1\right), \text{ 得到 } \frac{\partial l}{\partial \mathbf{W}_1} = \frac{\partial l}{\partial \mathbf{a}_1} \mathbf{x}^T.$$