

## 矩阵求导术（上）



长驱鬼侠  
数学爱好者

取消关注

编辑推荐

王赞 Maigo 等 4,894 人赞同了该文章

矩阵求导的技术，在统计学、控制论、机器学习等领域有广泛的应用。鉴于我看过的一些资料或言之不详、或繁乱无绪，本文来做个科普，分作两篇，上篇讲标量对矩阵的求导术，下篇讲矩阵对矩阵的求导术。本文使用小写字母 $x$ 表示标量，粗体小写字母 $\mathbf{x}$ 表示（列）向量，大写字母 $X$ 表示矩阵。

首先来琢磨一下定义，标量 $f$ 对矩阵 $X$ 的导数，定义为 $\frac{\partial f}{\partial X} = \left[ \frac{\partial f}{\partial X_{ij}} \right]$ ，即 $f$ 对 $X$ 逐元素求导排成与 $X$ 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**整体性**。试想，为何要将 $f$ 看做矩阵 $X$ 而不是各元素 $X_{ij}$ 的函数呢？答案是用矩阵运算更整洁。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）与微分有联系： $df = f'(x)dx$ ；多元微积分中的梯度（标量对向量的导数）也与微分有联系： $df = \sum_{i=1}^n \frac{\partial f}{\partial x_i} dx_i = \frac{\partial f^T}{\partial \mathbf{x}} d\mathbf{x}$ ，这里第一个等号是全微分公式，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：

$df = \sum_{i=1}^m \sum_{j=1}^n \frac{\partial f}{\partial X_{ij}} dX_{ij} = \text{tr} \left( \frac{\partial f^T}{\partial X} dX \right)$ 。其中 $\text{tr}$ 代表迹(trace)是方阵对角线元素之和，满足性

质：对尺寸相同的矩阵 $A, B$ ， $\text{tr}(A^T B) = \sum_{i,j} A_{ij} B_{ij}$ ，即 $\text{tr}(A^T B)$ 是矩阵 $A, B$ 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 $\frac{\partial f}{\partial X}$ （ $m \times n$ ）与微分矩阵 $dX$ （ $m \times n$ ）的内积。

然后来建立运算法则。回想遇到较复杂的一元函数如 $f = \log(2 + \sin x)e^{\sqrt{x}}$ ，我们是如何求导的呢？通常不是从定义开始求极限，而是先建立了初等函数求导和四则运算、复合等法则，再来运用这些法则。故而，我们来创立常用的矩阵微分的运算法则：

1. 加减法： $d(X \pm Y) = dX \pm dY$ ；

▲ 赞同 4894 ▼

● 373 条评论

➤ 分享

★ 收藏

...



$d(X^T) = (dX)^T$ ；迹： $d\text{tr}(X) = \text{tr}(dX)$ 。

2. 逆： $dX^{-1} = -X^{-1}dXX^{-1}$ 。此式可在  $XX^{-1} = I$  两侧求微分来证明。

3. 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，详见张贤达《矩阵分析与应用》第279页。

4. 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。

5. 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  $df = \text{tr}\left(\frac{\partial f^T}{\partial X} dX\right)$ ，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(trace trick)：

1. 标量套上迹： $a = \text{tr}(a)$

2. 转置： $\text{tr}(A^T) = \text{tr}(A)$ 。

3. 线性： $\text{tr}(A \pm B) = \text{tr}(A) \pm \text{tr}(B)$ 。

4. 矩阵乘法交换： $\text{tr}(AB) = \text{tr}(BA)$ ，其中  $A$  与  $B^T$  尺寸相同。两侧都等于  $\sum_{i,j} A_{ij} B_{ji}$ 。

5. 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  $\sum_{i,j} A_{ij} B_{ij} C_{ij}$ 。

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  $dX$  左侧，对照导数与微分的联系  $df = \text{tr}\left(\frac{\partial f^T}{\partial X} dX\right)$ ，即能得到导数。

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

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。我们直接从微分入手建立复合法则：先写出  $df = \text{tr}\left(\frac{\partial f^T}{\partial Y} dY\right)$ ，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，即可以得到

▲ 赞同 4894 ▼

● 373 条评论

➤ 分享

★ 收藏

..



最常见的情形是  $Y = AXB$ ，此时

$df = \text{tr} \left( \frac{\partial f}{\partial Y}^T dY \right) = \text{tr} \left( \frac{\partial f}{\partial Y}^T AdXB \right) = \text{tr} \left( B \frac{\partial f}{\partial Y}^T AdX \right) = \text{tr} \left( (A^T \frac{\partial f}{\partial Y} B^T)^T dX \right)$ ，可得到  $\frac{\partial f}{\partial X} = A^T \frac{\partial f}{\partial Y} B^T$ 。注意这里  $dY = (dA)XB + AdXB + AXdB = AdXB$ ，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  $\frac{\partial f}{\partial Y}^T AdX$  与  $B$ 。

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：  $f = a^T Xb$ ，求  $\frac{\partial f}{\partial X}$ 。其中  $a$  是  $m \times 1$  列向量， $X$  是  $m \times n$  矩阵， $b$  是  $n \times 1$  列向量， $f$  是标量。

解：先使用矩阵乘法法则求微分， $df = da^T Xb + a^T dXb + a^T Xdb = a^T dXb$ ，注意这里的  $a, b$  是常量， $da = 0, db = 0$ 。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： $df = \text{tr}(a^T dXb) = \text{tr}(ba^T dX) = \text{tr}((ab^T)^T dX)$ ，注意这里我们根据  $\text{tr}(AB) = \text{tr}(BA)$  交换了  $a^T dX$  与  $b$ 。对照导数与微分的联系  $df = \text{tr} \left( \frac{\partial f}{\partial X}^T dX \right)$ ，得到  $\frac{\partial f}{\partial X} = ab^T$ 。

注意：这里不能用  $\frac{\partial f}{\partial X} = a^T \frac{\partial X}{\partial X} b = ?$ ，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，这是逻辑上解释不通的。

例2：  $f = a^T \exp(Xb)$ ，求  $\frac{\partial f}{\partial X}$ 。其中  $a$  是  $m \times 1$  列向量， $X$  是  $m \times n$  矩阵， $b$  是  $n \times 1$  列向量， $\exp$  表示逐元素求指数， $f$  是标量。

解：先使用矩阵乘法、逐元素函数法则求微分： $df = a^T (\exp(Xb) \odot (dXb))$ ，再套上迹并做交换： $df = \text{tr}(a^T (\exp(Xb) \odot (dXb))) = \text{tr}((a \odot \exp(Xb))^T dXb) = \text{tr}(b(a \odot \exp(Xb))^T dX) = \text{tr}(((a \odot \exp(Xb))b^T)^T dX)$ ，注意这里我们先根据  $\text{tr}(A^T(B \odot C)) = \text{tr}((A \odot B)^T C)$  交换了  $a \odot \exp(Xb)$  与  $dXb$ ，再根据  $\text{tr}(AB) = \text{tr}(BA)$  交换了  $(a \odot \exp(Xb))^T dX$  与  $b$ 。对照导数与微分的联系  $df = \text{tr} \left( \frac{\partial f}{\partial X}^T dX \right)$ ，得到  $\frac{\partial f}{\partial X} = (a \odot \exp(Xb))b^T$ 。

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



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

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

$df = \text{tr}\left(\frac{\partial f}{\partial Y} dY\right)$ ，再将dY用dX表示出来代入，并使用矩阵乘法/逐元素乘法交换：

$df = \text{tr}\left(\frac{\partial f}{\partial Y} (\sigma'(WX) \odot (WdX))\right) = \text{tr}\left(\left(\frac{\partial f}{\partial Y} \odot \sigma'(WX)\right)^T WdX\right)$ ，对照导数与微分的联系，得到

$\frac{\partial f}{\partial X} = W^T \left(\frac{\partial f}{\partial Y} \odot \sigma'(WX)\right) = W^T ((2M\sigma(WX)) \odot \sigma'(WX))$ 。

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

解：这是标量对向量的导数，不过可以把向量看做矩阵的特例。先将向量模平方改写成向量与自身的内积： $l = (Xw - y)^T (Xw - y)$ ，求微分，使用矩阵乘法、转置等法则：

$dl = (Xdw)^T (Xw - y) + (Xw - y)^T (Xdw) = 2(Xw - y)^T Xdw$ ，注意这里  $Xdw$  和  $Xw - y$  是向量，两个向量的内积满足  $u^T v = v^T u$ 。对照导数与微分的联系  $dl = \frac{\partial l}{\partial w} dw$ ，得到

$\frac{\partial l}{\partial w} = 2X^T (Xw - y)$ 。  $\frac{\partial l}{\partial w} = 0$  即  $X^T Xw = X^T y$ ，得到  $w$  的最小二乘估计为  $w = (X^T X)^{-1} X^T y$ 。

例5【方差的最大似然估计】：样本  $x_1, \dots, x_N \sim \mathcal{N}(\mu, \Sigma)$ ，求方差  $\Sigma$  的最大似然估计。写成数学式是： $l = \log |\Sigma| + \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^T \Sigma^{-1} (x_i - \bar{x})$ ，求  $\frac{\partial l}{\partial \Sigma}$  的零点。其中  $x_i$  是  $m \times 1$  列向量，

$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$  是样本均值， $\Sigma$  是  $m \times m$  对称正定矩阵， $l$  是标量， $\log$  表示自然对数。

解：首先求微分，使用矩阵乘法、行列式、逆等运算法则，第一项是

$d \log |\Sigma| = |\Sigma|^{-1} d|\Sigma| = \text{tr}(\Sigma^{-1} d\Sigma)$ ，第二项是

$\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^T d\Sigma^{-1} (x_i - \bar{x}) = -\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^T d\Sigma \Sigma^{-1} (x_i - \bar{x})$ 。再给第二项套上述做法

换： $\text{tr}\left(\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^T \Sigma^{-1} d\Sigma \Sigma^{-1} (x_i - \bar{x})\right) = \frac{1}{N} \sum_{i=1}^N \text{tr}((x_i - \bar{x})^T \Sigma^{-1} d\Sigma \Sigma^{-1} (x_i - \bar{x}))$

$= \frac{1}{N} \sum_{i=1}^N \text{tr}(\Sigma^{-1} (x_i - \bar{x})(x_i - \bar{x})^T \Sigma^{-1} d\Sigma) = \text{tr}(\Sigma^{-1} S \Sigma^{-1} d\Sigma)$ ，其中先交换迹与求和，然后将

$\Sigma^{-1} (x_i - \bar{x})(x_i - \bar{x})^T$  交换到左边，最后再交换迹与求和，并定义  $S = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$  为样本方



差矩阵。得到  $dl = \text{tr}((\Sigma^{-1} - \Sigma^{-1}S\Sigma^{-1})d\Sigma)$ 。对照导数与微分的联系，有

$$\frac{\partial l}{\partial \Sigma} = (\Sigma^{-1} - \Sigma^{-1}S\Sigma^{-1})^T, \text{ 其零点即 } \Sigma \text{ 的最大似然估计为 } \Sigma = S。$$

例6【多元logistic回归】： $l = -\mathbf{y}^T \log \text{softmax}(W\mathbf{x})$ ，求  $\frac{\partial l}{\partial W}$ 。其中  $\mathbf{y}$  是除一个元素为1外其它元素为0的  $m \times 1$  列向量， $W$  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， $\text{softmax}(\mathbf{a}) = \frac{\exp(\mathbf{a})}{\mathbf{1}^T \exp(\mathbf{a})}$ ，其中  $\exp(\mathbf{a})$  表示逐元素求指数， $\mathbf{1}$  代表全1向量。

解1：首先将softmax函数代入并写成

$l = -\mathbf{y}^T (\log(\exp(W\mathbf{x})) - \mathbf{1} \log(\mathbf{1}^T \exp(W\mathbf{x}))) = -\mathbf{y}^T W\mathbf{x} + \log(\mathbf{1}^T \exp(W\mathbf{x}))$ ，这里要注意逐元素  $\log$  满足等式  $\log(\mathbf{u}/c) = \log(\mathbf{u}) - \mathbf{1} \log(c)$ ，以及  $\mathbf{y}$  满足  $\mathbf{y}^T \mathbf{1} = 1$ 。求微分，使用矩阵乘法、逐元素函数等法则： $dl = -\mathbf{y}^T dW\mathbf{x} + \frac{\mathbf{1}^T (\exp(W\mathbf{x}) \odot (dW\mathbf{x}))}{\mathbf{1}^T \exp(W\mathbf{x})}$ 。再套上迹并做交换，注意可化简

$\mathbf{1}^T (\exp(W\mathbf{x}) \odot (dW\mathbf{x})) = \exp(W\mathbf{x})^T dW\mathbf{x}$ ，这是根据等式  $\mathbf{1}^T (\mathbf{u} \odot \mathbf{v}) = \mathbf{u}^T \mathbf{v}$ ，故

$$dl = \text{tr} \left( -\mathbf{y}^T dW\mathbf{x} + \frac{\exp(W\mathbf{x})^T dW\mathbf{x}}{\mathbf{1}^T \exp(W\mathbf{x})} \right) = \text{tr}(-\mathbf{y}^T dW\mathbf{x} + \text{softmax}(W\mathbf{x})^T dW\mathbf{x}) = \text{tr}(\mathbf{x}(\text{softmax}(W\mathbf{x}) - \mathbf{y})^T dW)$$

。对照导数与微分的联系，得到  $\frac{\partial l}{\partial W} = (\text{softmax}(W\mathbf{x}) - \mathbf{y})\mathbf{x}^T$ 。

解2：定义  $\mathbf{a} = W\mathbf{x}$ ，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： $dl = \text{tr} \left( \frac{\partial l}{\partial \mathbf{a}}^T d\mathbf{a} \right) = \text{tr} \left( \frac{\partial l}{\partial \mathbf{a}}^T dW\mathbf{x} \right) = \text{tr} \left( \mathbf{x} \frac{\partial l}{\partial \mathbf{a}}^T dW \right)$ ，得到  $\frac{\partial l}{\partial W} = \frac{\partial l}{\partial \mathbf{a}} \mathbf{x}^T$ 。

最后一例留给经典的神经网络。神经网络的求导术是学术史上的重要成果，还有个专门的名字叫做BP算法，我相信如今很多人在初次推导BP算法时也会颇费一番脑筋，事实上使用矩阵求导术来推导并不复杂。为简化起见，我们推导二层神经网络的BP算法。

例7【二层神经网络】： $l = -\mathbf{y}^T \log \text{softmax}(W_2 \sigma(W_1 \mathbf{x}))$ ，求  $\frac{\partial l}{\partial W_1}$  和  $\frac{\partial l}{\partial W_2}$ 。其中  $\mathbf{y}$  是除一个元素为1外其它元素为0的  $m \times 1$  列向量， $W_2$  是  $m \times p$  矩阵， $W_1$  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， $\text{softmax}(\mathbf{a}) = \frac{\exp(\mathbf{a})}{\mathbf{1}^T \exp(\mathbf{a})}$  同上， $\sigma$  是逐元素sigmoid函数

$$\sigma(a) = \frac{1}{1 + \exp(-a)}。$$

解：定义  $\mathbf{a}_1 = W_1 \mathbf{x}$ ， $\mathbf{h}_1 = \sigma(\mathbf{a}_1)$ ， $\mathbf{a}_2 = W_2 \mathbf{h}_1$ ，则  $l = -\mathbf{y}^T \log \text{softmax}(\mathbf{a}_2)$ 。在前例中已求出  $\frac{\partial l}{\partial \mathbf{a}_2} = \text{softmax}(\mathbf{a}_2) - \mathbf{y}$ 。使用复合法则，

$$dl = \text{tr} \left( \frac{\partial l}{\partial \mathbf{a}_2}^T d\mathbf{a}_2 \right) = \text{tr} \left( \frac{\partial l}{\partial \mathbf{a}_2}^T dW_2 \mathbf{h}_1 \right) + \text{tr} \left( \frac{\partial l}{\partial \mathbf{a}_2}^T W_2 d\mathbf{h}_1 \right)，使用矩阵乘法交换的迹技巧从$$

第一项得到  $\frac{\partial l}{\partial W_2} = \frac{\partial l}{\partial \mathbf{a}_2} \mathbf{h}_1^T$ ，从第二项得到  $\frac{\partial l}{\partial \mathbf{h}_1} = W_2^T \frac{\partial l}{\partial \mathbf{a}_2}$ 。接下来对第二项继续使用复合求导  $\frac{\partial l}{\partial \mathbf{a}_1}$ ，并利用矩阵乘法和逐元素乘法交换的迹技巧：

$$d\mathbf{l}_2 = \text{tr} \left( \frac{\partial l}{\partial \mathbf{h}_1}^T d\mathbf{h}_1 \right) = \text{tr} \left( \frac{\partial l}{\partial \mathbf{h}_1}^T (\sigma'(\mathbf{a}_1) \odot d\mathbf{a}_1) \right) = \text{tr} \left( \left( \frac{\partial l}{\partial \mathbf{h}_1} \odot \sigma'(\mathbf{a}_1) \right)^T d\mathbf{a}_1 \right), \text{ 得到}$$

$\frac{\partial l}{\partial \mathbf{a}_1} = \frac{\partial l}{\partial \mathbf{h}_1} \odot \sigma'(\mathbf{a}_1)$ 。为求  $\frac{\partial l}{\partial W_1}$ ，再用一次复合法则：

$$d\mathbf{l}_2 = \text{tr} \left( \frac{\partial l}{\partial \mathbf{a}_1}^T d\mathbf{a}_1 \right) = \text{tr} \left( \frac{\partial l}{\partial \mathbf{a}_1}^T dW_1 \mathbf{x} \right) = \text{tr} \left( \mathbf{x} \frac{\partial l}{\partial \mathbf{a}_1}^T dW_1 \right), \text{ 得到 } \frac{\partial l}{\partial W_1} = \frac{\partial l}{\partial \mathbf{a}_1} \mathbf{x}^T.$$

推广：样本  $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)$ ， $l = -\sum_{i=1}^N \mathbf{y}_i^T \log \text{softmax}(W_2 \sigma(W_1 \mathbf{x}_i + \mathbf{b}_1) + \mathbf{b}_2)$ ，其中  $\mathbf{b}_1$  是  $p \times 1$  列向量， $\mathbf{b}_2$  是  $m \times 1$  列向量，其余定义同上。

解1：定义  $\mathbf{a}_{1,i} = W_1 \mathbf{x}_i + \mathbf{b}_1$ ， $\mathbf{h}_{1,i} = \sigma(\mathbf{a}_{1,i})$ ， $\mathbf{a}_{2,i} = W_2 \mathbf{h}_{1,i} + \mathbf{b}_2$ ，则

$l = -\sum_{i=1}^N \mathbf{y}_i^T \log \text{softmax}(\mathbf{a}_{2,i})$ 。先同上可求出  $\frac{\partial l}{\partial \mathbf{a}_{2,i}} = \text{softmax}(\mathbf{a}_{2,i}) - \mathbf{y}_i$ 。使用复合法则，

$$d\mathbf{l} = \text{tr} \left( \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{2,i}}^T d\mathbf{a}_{2,i} \right) = \text{tr} \left( \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{2,i}}^T dW_2 \mathbf{h}_{1,i} \right) + \underbrace{\text{tr} \left( \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{2,i}}^T W_2 d\mathbf{h}_{1,i} \right)}_{d_2} + \text{tr} \left( \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{2,i}}^T d\mathbf{b}_2 \right)$$

，从第一项得到  $\frac{\partial l}{\partial W_2} = \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{2,i}} \mathbf{h}_{1,i}^T$ ，从第二项得到  $\frac{\partial l}{\partial \mathbf{h}_{1,i}} = W_2^T \frac{\partial l}{\partial \mathbf{a}_{2,i}}$ ，从第三项得到

$\frac{\partial l}{\partial \mathbf{b}_2} = \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{2,i}}$ 。接下来对第二项继续使用复合法则，得到  $\frac{\partial l}{\partial \mathbf{a}_{1,i}} = \frac{\partial l}{\partial \mathbf{h}_{1,i}} \odot \sigma'(\mathbf{a}_{1,i})$ 。为求

$\frac{\partial l}{\partial W_1}$ ， $\frac{\partial l}{\partial \mathbf{b}_1}$ ，再用一次复合法则：

$$d\mathbf{l}_2 = \text{tr} \left( \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{1,i}}^T d\mathbf{a}_{1,i} \right) = \text{tr} \left( \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{1,i}}^T dW_1 \mathbf{x}_i \right) + \text{tr} \left( \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{1,i}}^T d\mathbf{b}_1 \right), \text{ 得到}$$

$$\frac{\partial l}{\partial W_1} = \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{1,i}} \mathbf{x}_i^T, \quad \frac{\partial l}{\partial \mathbf{b}_1} = \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{1,i}}.$$

解2：可以用矩阵来表示N个样本，以简化形式。定义  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ ，

$$\mathbf{A}_1 = [\mathbf{a}_{1,1}, \dots, \mathbf{a}_{1,N}] = W_1 \mathbf{X} + \mathbf{b}_1 \mathbf{1}^T, \quad \mathbf{H}_1 = [\mathbf{h}_{1,1}, \dots, \mathbf{h}_{1,N}] = \sigma(\mathbf{A}_1),$$

$$\mathbf{A}_2 = [\mathbf{a}_{2,1}, \dots, \mathbf{a}_{2,N}] = W_2 \mathbf{H}_1 + \mathbf{b}_2 \mathbf{1}^T, \text{ 注意这里使用全1向量来扩展维度。先同上求出}$$

$\frac{\partial l}{\partial \mathbf{A}_2} = [\text{softmax}(\mathbf{a}_{2,1}) - \mathbf{y}_1, \dots, \text{softmax}(\mathbf{a}_{2,N}) - \mathbf{y}_N]$ 。使用复合法则，

$$d\mathbf{l} = \text{tr} \left( \frac{\partial l}{\partial \mathbf{A}_2}^T d\mathbf{A}_2 \right) = \text{tr} \left( \frac{\partial l}{\partial \mathbf{A}_2}^T dW_2 \mathbf{H}_1 \right) + \underbrace{\text{tr} \left( \frac{\partial l}{\partial \mathbf{A}_2}^T W_2 d\mathbf{H}_1 \right)}_{d_2} + \text{tr} \left( \frac{\partial l}{\partial \mathbf{A}_2}^T d\mathbf{b}_2 \mathbf{1}^T \right), \text{ 从第一}$$

项得到  $\frac{\partial l}{\partial W_2} = \frac{\partial l}{\partial \mathbf{A}_2} \mathbf{H}_1^T$ ，从第二项得到  $\frac{\partial l}{\partial \mathbf{H}_1} = W_2^T \frac{\partial l}{\partial \mathbf{A}_2}$ ，从第三项得到  $\frac{\partial l}{\partial \mathbf{b}_2} = \frac{\partial l}{\partial \mathbf{A}_2} \mathbf{1}$ 。接下

来对第二项继续使用复合法则，得到

▲ 赞同 4894 ▼

● 373 条评论

➤ 分享

★ 收藏

..

则：
$$d\mathbf{l}_2 = \text{tr} \left( \frac{\partial \mathbf{l}}{\partial \mathbf{A}_1}{}^T d\mathbf{A}_1 \right) = \text{tr} \left( \frac{\partial \mathbf{l}}{\partial \mathbf{A}_1}{}^T d\mathbf{W}_1 \mathbf{X} \right) + \text{tr} \left( \frac{\partial \mathbf{l}}{\partial \mathbf{A}_1}{}^T d\mathbf{b}_1 \mathbf{1}^T \right), \text{ 得到 } \frac{\partial \mathbf{l}}{\partial \mathbf{W}_1} = \frac{\partial \mathbf{l}}{\partial \mathbf{A}_1} \mathbf{X}^T$$

$$\frac{\partial \mathbf{l}}{\partial \mathbf{b}_1} = \frac{\partial \mathbf{l}}{\partial \mathbf{A}_1} \mathbf{1}.$$

下篇见[zhuanlan.zhihu.com/p/24...](https://zhuanlan.zhihu.com/p/24709748)

编辑于 03-06

机器学习    矩阵分析    优化

推荐阅读

让向量、矩阵和张量的求导更简洁些吧

本文是我在阅读Erik Learned-Miller的《Vector, Matrix, and Tensor Derivatives》时的记录，点此下载。本文的主要内容是帮助你学习如何进行向量、矩阵以及高阶张量（三维及以上的数组）的...  
李是Lyapunov的李

矩阵求导术（下）

本文承接上篇  
<https://zhuanlan.zhihu.com/p/24709748>  
来讲矩阵对矩阵的求导术。使用小写字母 $\mathbf{x}$ 表示标量，粗体小写字母 $\mathbf{x}$ 表示列向量，大写字母 $\mathbf{X}$ 表示矩阵。矩阵对矩阵的求...  
长躯鬼侠

干货 | 万数!

梁勇

373 条评论

切换为时间排序

写下你的评论...

henryWang

2017-01-21

之前看吴恩达的cs229就对这个矩阵求导和迹的联系十分懵，多谢作者

1

你说啥研究我不懂

2017-01-22

赞同 4894

373 条评论

分享

收藏

刚开始搞CV时候，发现典型的CVPR paper总有这么几个矩阵求导公式，与那些只知道刷baseline的文章相比，顿时逼格高大上啊！唉，作者你要是早几年写这个专栏一定能解救一批苦逼的PhD啊！

3

花花 2017-02-05

写得真好！

赞

郑建东 2017-02-12

感谢大神指导，有点茅塞顿开的感觉。

赞

小成 2017-02-14

小白求教：

请问在例3【logistic回归】的解答中第一行：

首先将softmax函数代入并写成  $l = \dots$ ，这个公式中右侧第二项为什么不用和第一项一样乘上y的转置？

1

长躯鬼侠 (作者) 回复 小成 2017-02-18

这里我跳了一步，y是除一个元素为1外其它元素为0的向量，log()是个标量，化简一下就是这样。

赞

小成 回复 长躯鬼侠 (作者) 2017-02-19

明白啦，感谢大神回复！

赞

查看全部 6 条回复

郭洋洋 2017-02-15

十分好用 NBNBNNBNBNBNBNB

赞

Martins3 2017-02-18

很强，很酷

赞

shiryaev 2017-02-27

[cal.cs.illinois.edu/~jo...](http://cal.cs.illinois.edu/~joc...)

赞同 4894 373 条评论 分享 收藏



这篇更基础点，希望能给更多的小伙伴帮忙



4



相思作坊半世離殤 回复 shiryaev

2017-03-20

感谢分享资料。看了这篇pdf和楼主写的文章，明白了许多

1



bb just 回复 相思作坊半世離殤

2019-06-04

链接打不开了，请问你还有这个质量么。

赞

查看全部 7 条回复



黄角兽

2017-03-08

功德无量

赞



排骨郎

2017-03-16

写的很好

赞



我不是那样的烟火

2017-03-27

虽然我也是理解一点的，可是，你这篇文章简直不要太赞好吧

赞



Winston

2017-05-12

很实用，感谢！

赞



赵小呆

2017-05-24

之前看cs229,上面写要拆开运算，结果每次写cs229的习题都要运算好久，还容易出错。楼主伟大！但是这里有一个问题，矩阵运算和常量运算有时运算法则一致，有时不一致，楼主有推荐资料解决这个问题吗？

赞



长躯鬼侠 (作者) 回复 赵小呆

2017-05-25

我不确定你具体想问哪种情况，或许你可以举个例子来讨论。简略地讲，求微分的运算法则，对矩阵和标量基本一致；从微分获得导数，对矩阵需要特别的技术，比如迹技巧。

赞



亦然

赞同 4894

373 条评论

分享

收藏

..

想问一下复数矩阵求导且对非方阵求导应做哪些改变和处理？

👍 赞



穆文

2017-07-30

文章讲的很赞，非常感谢分享。说两点我的理解，跟作者探讨下

1. 我的理解就是，其实矩阵不过是另一套『数学语言』，把原来标量世界的各种运算法则，转换到矩阵世界，用矩阵、矩阵的乘积、矩阵秩、矩阵迹、行列式等，来刻画新的运算法则，不知这个理解对不对？
2. 对原文『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』这句话持保留态度，这里X是矩阵的时候，没有明确定义，但如果X是向量，结论是A没错吧？因为AX刻画了原始空间X到新的空间AX的映射，即向量对向量求导，是有明确定义

👍 赞



长躯鬼侠 (作者) 回复 穆文

2017-07-31

我在这里讲的是标量对矩阵的导数。而向量对向量的导数，是我的下篇里矩阵对矩阵导数的一种特殊情况，按下篇中的定义有A对x的导数是A<sup>T</sup>。

👍 赞



陌烛

2017-08-26

你好，请问下，上篇中，例三中的log都是当作ln使用的吗？写作log用作ln？

👍 赞



长躯鬼侠 (作者) 回复 陌烛

2017-08-28

是的，log当做ln用了。

👍 1



陌烛 回复 长躯鬼侠 (作者)

2017-09-07

知道了，蟹蟹

👍 赞



X'plore

2017-09-14

$$df = \sum_i \frac{\partial f}{\partial x_i} dx_i = \frac{\partial f}{\partial \boldsymbol{x}}^T d\boldsymbol{x}$$
 请问右面这一步如何理解，为什么等于转置的偏微分？

👍 赞



长躯鬼侠 (作者) 回复 X'plore

2017-09-14

根据定义

👍 赞

▲ 赞同 4894 ▼

● 373 条评论

➦ 分享

★ 收藏

..



武宽Quin Wu

2017-09-1



博主您好，关于您提到的一个标量对矩阵X的求导，我有一点不解的是，您这里没有提到wiki中所谓的 分子布局 跟分母布局，如果是分子布局的话。 $\frac{\partial f}{\partial X} = \frac{\partial f}{\partial X_{ij}}$ 吧，而不应该是 $X_{ij}$ 。在分母布局下， $\frac{\partial f}{\partial X}$ 就是您所描述的  $\frac{\partial f}{\partial X_{ij}}$ 。是这样嘛？

👍 1



长驱鬼侠 (作者) 回复 武宽Quin Wu

2017-09-16

我使用的定义第二个。

👍 1



徐敬

2017-09-21

厉害了

👍 赞



侯佐

2017-10-03

很厉害，解决了很长时间的疑问

👍 赞