

## 矩阵求导的理解（重要！）



codingli...

公众号：L的算法成长之路

关注她

99 人赞同了该文章

《矩阵求导术》重点笔记

### 首先是标量对矩阵的求导

一元微积分（标量对标量）中的导数与微分的关系： $df = f'(x)dx$

多元微积分（标量对向量）中的梯度与微分的关系： $df = \sum_{i=1}^n \frac{\partial f}{\partial x_i} dx_i = \frac{\partial f^T}{\partial x} dx$

（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  $\frac{\partial f}{\partial x}$  (nx1)与微分向量  $dx$  (nx1)的内积）

矩阵导数与微分建立联系： $df = \sum_{i=1}^m \sum_{j=1}^n \frac{\partial f}{\partial X_{ij}} dX_{ij} = \text{tr}(\frac{\partial f^T}{\partial X} dX)$

其中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。

这里表示，全微分  $df$  是导数  $\frac{\partial f}{\partial X}$  (m×n)与微分矩阵  $dX$  (m×n)的内积。

矩阵微分的运算法则：

1.加减法： $d(X \pm Y) = dX \pm dY$

矩阵乘法： $d(XY) = (dX)Y + XdY$

转置： $d(X^T) = (dX)^T$

迹： $d\text{tr}(X) = \text{tr}(dX)$

2.逆： $dX^{-1} = -X^{-1}dXX^{-1}$

3.行列式： $d|X| = \text{tr}(X^* dX)$ ， $X^*$  表示X的伴随矩阵

如果X可逆，上式可写成  $d|X| = |X|\text{tr}(X^{-1}dX)$

4.逐元素乘法： $d(X \odot Y) = dX \odot Y + X \odot dY$ ， $\odot$  代表尺寸相同的矩阵逐元素相乘

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。

一些迹技巧：

1.标量套上迹： $a = \text{tr}(a)$

2.转置： $\text{tr}(A^T) = \text{tr}(A)$



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  $\sum_{i,j} A_{ij}B_{ij}$

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，即能得到导数。

关于复合：

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  $\frac{\partial f}{\partial x} = \frac{\partial f}{\partial y} \frac{\partial y}{\partial x}$ ，但这里我们不能沿用链式法则，因为矩阵对矩阵的导数  $\frac{\partial Y}{\partial X}$  截止目前仍是未定义的。我们直接从微分入手建立复合法则：先写出  $df = \text{tr}(\frac{\partial f}{\partial Y} dY)$ ，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  $\frac{\partial f}{\partial X}$ 。

来看几个例子：

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

解：先使用矩阵乘法法则求微分，这里的  $\mathbf{a}, \mathbf{b}$  是常量， $d\mathbf{a} = \mathbf{0}, d\mathbf{b} = \mathbf{0}$ ，得到： $df = \mathbf{a}^T d\mathbf{X} \mathbf{b}$ ，再套上迹并做矩阵乘法交换： $df = \text{tr}(\mathbf{a}^T d\mathbf{X} \mathbf{b}) = \text{tr}(\mathbf{b} \mathbf{a}^T d\mathbf{X})$ ，注意这里我们根据  $\text{tr}(AB) = \text{tr}(BA)$  交换了  $\mathbf{a}^T d\mathbf{X}$  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  $\frac{\partial f}{\partial \mathbf{X}} = (\mathbf{b} \mathbf{a}^T)^T = \mathbf{a} \mathbf{b}^T$ 。

注意：这里不能用  $\frac{\partial f}{\partial \mathbf{X}} = \mathbf{a}^T \frac{\partial \mathbf{X}}{\partial \mathbf{X}} \mathbf{b} = ?$ ，导数与乘常数矩阵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： $f = \mathbf{a}^T \exp(\mathbf{X} \mathbf{b})$ ，求  $\frac{\partial f}{\partial \mathbf{X}}$ 。其中  $\mathbf{a}$  是  $m \times 1$  列向量， $\mathbf{X}$  是  $m \times n$  矩阵， $\mathbf{b}$  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： $df = \mathbf{a}^T (\exp(\mathbf{X} \mathbf{b}) \odot (d\mathbf{X} \mathbf{b}))$ ，再套上迹并做  $df = \text{tr}(\mathbf{a}^T (\exp(\mathbf{X} \mathbf{b}) \odot (d\mathbf{X} \mathbf{b}))) = \text{tr}((\mathbf{a} \odot \exp(\mathbf{X} \mathbf{b}))^T d\mathbf{X} \mathbf{b}) = \text{tr}(\mathbf{b} (\mathbf{a} \odot \exp(\mathbf{X} \mathbf{b}))^T d\mathbf{X})$ ，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  $\mathbf{a} \odot \exp(\mathbf{X} \mathbf{b})$  与  $d\mathbf{X} \mathbf{b}$ ，再根据  $\text{tr}(AB) = \text{tr}(BA)$  交换了  $(\mathbf{a} \odot \exp(\mathbf{X} \mathbf{b}))^T d\mathbf{X}$  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  $\frac{\partial f}{\partial \mathbf{X}} = (\mathbf{b} (\mathbf{a} \odot \exp(\mathbf{X} \mathbf{b}))^T)^T = (\mathbf{a} \odot \exp(\mathbf{X} \mathbf{b})) \mathbf{b}^T$ 。

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  $l \times m$  列向量， $\mathbf{X}$  是  $m \times n$  矩阵， $\mathbf{Y}$  是  $l \times n$  矩阵， $\mathbf{M}$  是  $l \times l$  对称矩阵， $\sigma$  是逐元素函数， $f$  是标量。

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

$df = \text{tr}((d\mathbf{Y})^T \mathbf{M} \mathbf{Y}) + \text{tr}(\mathbf{Y}^T \mathbf{M} d\mathbf{Y}) = 2\text{tr}(\mathbf{Y}^T \mathbf{M} d\mathbf{Y})$ ，对照导数与微分的联系，得到  $\frac{\partial f}{\partial \mathbf{Y}} = 2\mathbf{M} \mathbf{Y}$ 。为求  $\frac{\partial f}{\partial \mathbf{X}}$ ，写出  $df = \text{tr}(\frac{\partial f}{\partial \mathbf{Y}} d\mathbf{Y})$ ，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：严格来说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：

$dl = (Xdw)^T (Xw - y) + (Xw - y)^T (Xdw) = 2(Xw - y)^T Xdw$ 。对照导数与微分的联系  $dl = \frac{\partial l}{\partial w} dw$ ，得到  $\frac{\partial l}{\partial w} = (2(Xw - y)^T X)^T = 2X^T (Xw - y)$ 。 $\frac{\partial l}{\partial w}$  的零点即  $w$  的最小二乘估计为  $w = (X^T X)^{-1} X^T y$ 。

例5【方差的极大似然估计】：样本  $x_1, \dots, x_n \sim N(\mu, \Sigma)$ ，求方差  $\Sigma$  的极大似然估计。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。

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交

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

$\Sigma^{-1} (x_i - \bar{x})$  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

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

例6【多元logistic回归】： $l = -y^T \log \text{softmax}(Wx)$ ，求  $\frac{\partial l}{\partial W}$ 。其中  $y$  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， $x$  是  $n \times 1$  列向量， $l$  是标量；

$\text{softmax}(a) = \frac{\exp(a)}{1^T \exp(a)}$ ，其中  $\exp(a)$  表示逐元素求指数， $1$  代表全1向量。

解：首先将softmax函数代入并写成

$l = -y^T (\log(\exp(Wx)) - 1 \log(1^T \exp(Wx))) = -y^T Wx + \log(1^T \exp(Wx))$ ，这里要注意逐元素log满足等式  $\log(u/c) = \log(u) - 1 \log(c)$ ，以及  $y$  满足  $y^T 1 = 1$ 。求微分，使用矩阵乘法、逐

元素函数等法则： $dl = -y^T dWx + \frac{1^T (\exp(Wx) \odot (dWx))}{1^T \exp(Wx)}$ 。再套上述并做交换，注意可化简

$1^T (\exp(Wx) \odot (dWx)) = \exp(Wx)^T dWx$ ，这是根据等式  $1^T (u \odot v) = u^T v$ ，故

$dl = \text{tr} \left( -y^T dWx + \frac{\exp(Wx)^T dWx}{1^T \exp(Wx)} \right) = \text{tr}(x(\text{softmax}(Wx) - y)^T dW)$ 。对照导数与微分的联

系，得到  $\frac{\partial l}{\partial W} = (\text{softmax}(Wx) - y)x^T$ 。

另解：定义  $a = Wx$ ，则  $l = -y^T \log \text{softmax}(a)$ ，先如上求出  $\frac{\partial l}{\partial a} = \text{softmax}(a) - y$ ，再利用复

合法则： $dl = \text{tr} \left( \frac{\partial l}{\partial a} da \right) = \text{tr} \left( \frac{\partial l}{\partial a} dWx \right) = \text{tr} \left( x \frac{\partial l}{\partial a}^T dW \right)$ ，得到  $\frac{\partial l}{\partial W} = \frac{\partial l}{\partial a} x^T$ 。

先定义向量  $f$  ( $p \times 1$ ) 对向量  $x$  ( $m \times 1$ ) 的导数:

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_2}{\partial x_1} & \cdots & \frac{\partial f_p}{\partial x_1} \\ \frac{\partial f_1}{\partial x_2} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_p}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_1}{\partial x_m} & \frac{\partial f_2}{\partial x_m} & \cdots & \frac{\partial f_p}{\partial x_m} \end{bmatrix} \quad (m \times p)$$

有  $df = \frac{\partial f^T}{\partial x} dx$  ;

再定义矩阵的 (按列优化) 向量化

$\text{vec}(X) = [X_{11}, \dots, X_{m1}, X_{12}, \dots, X_{m2}, \dots, X_{1n}, \dots, X_{mn}]^T$  ( $mn \times 1$ ), 并定义矩阵  $F$  对矩阵  $X$  的导数  $\frac{\partial F}{\partial X} = \frac{\partial \text{vec}(F)}{\partial \text{vec}(X)}$  ( $mn \times pq$ )。导数与微分有联系:

$$\text{vec}(dF) = \frac{\partial F}{\partial X}^T \text{vec}(dX)$$

向量化的技巧:

1. 线性:  $\text{vec}(A + B) = \text{vec}(A) + \text{vec}(B)$

2. 矩阵乘法:  $\text{vec}(AXB) = (B^T \otimes A) \text{vec}(X)$ , 其中  $\otimes$  代表 Kronecker 积,  $A(m \times n)$  与  $B(p \times q)$  的 Kronecker 积是  $A \otimes B = [A_{ij}B]$  ( $mp \times nq$ )。

3. 转置:  $\text{vec}(A^T) = K_{mn} \text{vec}(A)$ ,  $A$  是  $m \times n$  矩阵, 其中  $K_{mn}$  ( $mn \times mn$ ) 是交换矩阵 (commutation matrix)。

4. 逐元素乘法:  $\text{vec}(A \odot X) = \text{diag}(\text{vec}(A)) \text{vec}(X)$ , 其中  $\text{diag}(A)$  ( $mn \times mn$ ) 是用  $A$  的元素 (按列优先) 排成的对角阵。

观察一下可以断言, 若矩阵函数  $F$  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, 再做向量化并使用技巧将其它项交换至  $\text{vec}(dX)$  左侧, 即能得到导数。

再谈一谈复合: 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? 从导数与微分的联系入手,

$$\text{vec}(dF) = \frac{\partial F}{\partial Y}^T \text{vec}(dY) = \frac{\partial F}{\partial Y}^T \frac{\partial Y}{\partial X}^T \text{vec}(dX), \text{ 可以推出链式法则 } \frac{\partial F}{\partial X} = \frac{\partial Y}{\partial X} \frac{\partial F}{\partial Y}.$$

有一些 Kronecker 积和交换矩阵相关的恒等式, 可用来做等价变形:

$$1. (A \otimes B)^T = A^T \otimes B^T.$$

$$2. \text{vec}(ab^T) = b \otimes a.$$

3.  $(A \otimes B)(C \otimes D) = (AC) \otimes (BD)$ 。可以对  $F = D^T B^T X A C$  求导来证明, 一方面, 直接求得  $\frac{\partial F}{\partial X} = (AC) \otimes (BD)$ ; 另一方面, 引入  $Y = B^T X A$ , 有  $\frac{\partial F}{\partial Y} = C \otimes D$ ,  $\frac{\partial Y}{\partial X} = A \otimes B$ , 用链式法则得到  $\frac{\partial F}{\partial X} = (A \otimes B)(C \otimes D)$ 。

$$4. K_{mn} = K_{nm}^T, K_{mn} K_{nm} = I.$$

5.  $K_{mn}(A \otimes B)K_{nm} = B \otimes A$ 。  $A$  是  $m \times n$  矩阵,  $B$  是  $n \times m$  矩阵。可以对  $A Y B^T$  做向量化来证明, 一

$$\text{vec}(AXB^T) = K_{pm} \text{vec}(BX^T A^T) = K_{pm}(A \otimes B) \text{vec}(X^T) = K_{pm}(A \otimes B) K_{nq} \text{vec}(X)。$$



例子：

例1:  $F = AX$ ,  $X$ 是 $m \times n$ 矩阵, 求  $\frac{\partial F}{\partial X}$ 。

解: 先求微分:  $dF = AdX$ , 再做向量化, 使用矩阵乘法的技巧, 注意在 $dX$ 右侧添加单位阵:

$$\text{vec}(dF) = \text{vec}(AdX) = (I_n \otimes A) \text{vec}(dX), \text{ 对照导数与微分的联系得到 } \frac{\partial F}{\partial X} = I_n \otimes A^T。$$

特例: 如果 $X$ 退化为向量, 即  $f = Ax$ , 则根据向量的导数与微分的关系  $df = \frac{\partial f^T}{\partial x} dx$ , 得到

$$\frac{\partial f}{\partial x} = A^T。$$

例2:  $f = \log|X|$ ,  $X$ 是 $n \times n$ 矩阵, 求  $\nabla_X f$  和  $\nabla_X^2 f$ 。

解: 使用上篇中的技术可求得  $\nabla_X f = X^{-1T}$ 。为求  $\nabla_X^2 f$ , 先求微分:  $d\nabla_X f = -(X^{-1}dXX^{-1})^T$ , 再做向量化, 使用转置和矩阵乘法的技巧

$$\text{vec}(d\nabla_X f) = -K_{nn} \text{vec}(X^{-1}dXX^{-1}) = -K_{nn}(X^{-1T} \otimes X^{-1}) \text{vec}(dX), \text{ 对照导数与微分的联系, 得到 } \nabla_X^2 f = -K_{nn}(X^{-1T} \otimes X^{-1}), \text{ 注意它是对称矩阵。在 } X \text{ 是对称矩阵时, 可简化为 } \nabla_X^2 f = -X^{-1} \otimes X^{-1}。$$

例3:  $F = A \exp(XB)$ ,  $A$ 是 $l \times m$ 矩阵,  $X$ 是 $m \times n$ 矩阵,  $B$ 是 $n \times p$ 矩阵,  $\exp$ 为逐元素函数, 求  $\frac{\partial F}{\partial X}$ 。

解: 先求微分:  $dF = A(\exp(XB) \odot (dXB))$ , 再做向量化, 使用矩阵乘法的技巧:

$$\text{vec}(dF) = (I_p \otimes A) \text{vec}(\exp(XB) \odot (dXB)), \text{ 再用逐元素乘法的技巧:}$$

$$\text{vec}(dF) = (I_p \otimes A) \text{diag}(\exp(XB)) \text{vec}(dXB), \text{ 再用矩阵乘法的技巧:}$$

$$\text{vec}(dF) = (I_p \otimes A) \text{diag}(\exp(XB))(B^T \otimes I_m) \text{vec}(dX), \text{ 对照导数与微分的联系得到}$$

$$\frac{\partial F}{\partial X} = (B \otimes I_m) \text{diag}(\exp(XB))(I_p \otimes A^T)。$$

例4【一元logistic回归】:  $l = -y\mathbf{x}^T \mathbf{w} + \log(1 + \exp(\mathbf{x}^T \mathbf{w}))$ , 求  $\nabla_{\mathbf{w}} l$  和  $\nabla_{\mathbf{w}}^2 l$ 。其中  $y$  是取值0或1的标量,  $\mathbf{x}, \mathbf{w}$  是  $n \times 1$  列向量。

解: 使用上篇中的技术可求得  $\nabla_{\mathbf{w}} l = \mathbf{x}(\sigma(\mathbf{x}^T \mathbf{w}) - y)$ , 其中  $\sigma(a) = \frac{\exp(a)}{1 + \exp(a)}$  为sigmoid函数。

为求  $\nabla_{\mathbf{w}}^2 l$ , 先求微分:  $d\nabla_{\mathbf{w}} l = \mathbf{x}\sigma'(\mathbf{x}^T \mathbf{w})\mathbf{x}^T d\mathbf{w}$ , 其中  $\sigma'(a) = \frac{\exp(a)}{(1 + \exp(a))^2}$  为sigmoid函数的导数, 对照导数与微分的联系, 得到  $\nabla_{\mathbf{w}}^2 l = \mathbf{x}\sigma'(\mathbf{x}^T \mathbf{w})\mathbf{x}^T$ 。

推广: 样本  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ ,  $l = \sum_{i=1}^N (-y_i \mathbf{x}_i^T \mathbf{w} + \log(1 + \exp(\mathbf{x}_i^T \mathbf{w})))$ , 求  $\nabla_{\mathbf{w}} l$  和  $\nabla_{\mathbf{w}}^2 l$ 。有

两种方法, 方法一: 先对每个样本求导, 然后相加; 方法二: 定义矩阵  $X = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}$ , 向量

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \text{ 将 } l \text{ 写成矩阵形式 } l = -\mathbf{y}^T X \mathbf{w} + \mathbf{1}^T \log(1 + \exp(X \mathbf{w})), \text{ 进而可以求得}$$

$$\nabla_{\mathbf{w}} l = X^T(\sigma(X \mathbf{w}) - \mathbf{y}), \quad \nabla_{\mathbf{w}}^2 l = X^T \text{diag}(\sigma'(X \mathbf{w})) X。$$

例5【多元logistic回归】:  $l = -\mathbf{y}^T \log \text{softmax}(W \mathbf{x}) = -\mathbf{y}^T W \mathbf{x} + \log(\mathbf{1}^T \exp(W \mathbf{x}))$ , 求  $\nabla_{\mathbf{w}} l$  和  $\nabla_{\mathbf{w}}^2 l$ 。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。



解：上篇中已求得  $\nabla_W l = (\text{softmax}(W\mathbf{x}) - \mathbf{y})\mathbf{x}^T$ 。为求  $\nabla_W^2 l$ ，先求微分：定义  $\mathbf{a} = W\mathbf{x}$ ，  

$$d\text{softmax}(\mathbf{a}) = \frac{\exp(\mathbf{a}) \odot d\mathbf{a}}{\mathbf{1}^T \exp(\mathbf{a})} - \frac{\exp(\mathbf{a})(\mathbf{1}^T (\exp(\mathbf{a}) \odot d\mathbf{a}))}{(\mathbf{1}^T \exp(\mathbf{a}))^2}$$
，这里需要化简去掉逐元素乘法，第一项中  $\exp(\mathbf{a}) \odot d\mathbf{a} = \text{diag}(\exp(\mathbf{a}))d\mathbf{a}$ ，第二项中  $\mathbf{1}^T (\exp(\mathbf{a}) \odot d\mathbf{a}) = \exp(\mathbf{a})^T d\mathbf{a}$ ，故有  

$$d\text{softmax}(\mathbf{a}) = D\text{softmax}(\mathbf{a})d\mathbf{a}$$
，其中  $D\text{softmax}(\mathbf{a}) = \frac{\text{diag}(\exp(\mathbf{a}))}{\mathbf{1}^T \exp(\mathbf{a})} - \frac{\exp(\mathbf{a})\exp(\mathbf{a})^T}{(\mathbf{1}^T \exp(\mathbf{a}))^2}$ ，代入有  

$$d\nabla_W l = D\text{softmax}(\mathbf{a})d\mathbf{a}\mathbf{x}^T = D\text{softmax}(W\mathbf{x})dW\mathbf{x}\mathbf{x}^T$$
，做向量化并使用矩阵乘法的技巧，得到  

$$\nabla_W^2 l = (\mathbf{x}\mathbf{x}^T) \otimes D\text{softmax}(W\mathbf{x})$$
。

最后做个总结。我们发展了从**整体**出发的矩阵求导的技术，**导数与微分的联系是计算的枢纽**，标量对矩阵的导数与微分的联系是  $d\mathbf{f} = \text{tr}(\nabla_X^T \mathbf{f} d\mathbf{X})$ ，先对f求微分，再使用迹技巧可求得导数，特别地，标量对向量的导数与微分的联系是  $d\mathbf{f} = \nabla_{\mathbf{x}}^T \mathbf{f} d\mathbf{x}$ ；矩阵对矩阵的导数与微分的联系是  $\text{vec}(d\mathbf{F}) = \frac{\partial \mathbf{F}^T}{\partial \mathbf{X}} \text{vec}(d\mathbf{X})$ ，先对F求微分，再使用向量化的技巧可求得导数，特别地，向量对向量的导数与微分的联系是  $d\mathbf{f} = \frac{\partial \mathbf{f}^T}{\partial \mathbf{x}} d\mathbf{x}$ 。

参考资料：

[Matrix calculus - Wikipedia](#)

[通过一个例子快速上手矩阵求导 - NoGeek - CSDN博客](#)

[矩阵求导术（上）](#)

[矩阵求导术（下）](#)

发布于 2018-11-19

[理解](#)

## 文章被以下专栏收录



### 算法修炼之路

机器学习和深度学习的相关算法等

## 推荐阅读

### 矩阵求导与矩阵微分

矩阵求导与矩阵微分符号定义 使用大写的粗体字母表示矩阵  $\mathbf{A}$ 、 $\mathbf{F}$  使用小写的粗体字母表示向量  $\mathbf{x}$ 、 $\mathbf{f}$ ，这里默认为列向量 使用小写的正体字母表示标量  $x$ 、 $f$ ...

说谎的傻子

发表于科学与技术

### 微分的四种理解

一 微分是无穷小？物理人喜欢把微分看做是一个很小的量，这在计算时总是很方便的，但是给人一种不严谨的感觉。实际上，它确实不严谨，第二次数学危机就是因此产生的。严谨性与明晰性是互补...

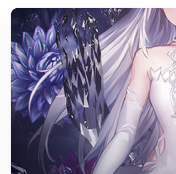
像牛一样的猫

### 向量恒等式

这次的内容相当平凡。最近学磁流体力学，发现对向量恒等式还是不熟悉，以前只是记一下  $\mathbf{A} \times (\mathbf{B} \times \mathbf{C}) = (\mathbf{A} \cdot \mathbf{C})\mathbf{B} - (\mathbf{A} \cdot \mathbf{B})\mathbf{C}$ ...

魂魄妖妖梦

发表于A Tri...







### Stokes定理八分 分外乘积(楔积)

来自虚空的...





写下你的评论...




 冷雨04-19

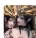
这句话是啥意思？因为矩阵对矩阵的导数截至目前是未定义的。为啥？

 赞


 codingling (作者) 回复 冷雨04-20


因为不知道y对x的导数，没法用链式法则

 赞

 阳光陈靖文01-22

刚看完长驱鬼侠的矩阵求导术，就读到你的笔记哈哈哈，总结的不错



 赞

 雁过无声2020-06-06

666

 赞

 han2019-08-27

清晰，易懂

 赞