

矩阵求导术（下）



长驱鬼侠
数学爱好者

612 人赞同了该文章

本文承接上篇 zhuanlan.zhihu.com/p/24...，来讲矩阵对矩阵的求导术。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。矩阵对矩阵的求导采用了向量化的思路，常应用于二阶方法求解优化问题。

首先来琢磨一下定义。矩阵对矩阵的导数，需要什么样的定义？第一，矩阵 $F(p \times q)$ 对矩阵 $X(m \times n)$ 的导数应包含所有 $mnpq$ 个偏导数 $\frac{\partial F_{kl}}{\partial X_{ij}}$ ，从而不损失信息；第二，导数与微分有简明的联系，因为在计算导数和应用中需要这个联系；第三，导数有简明的从整体出发的算法。我们先

定义向量 $\mathbf{f}(p \times 1)$ 对向量 $\mathbf{x}(m \times 1)$ 的导数 $\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_2}{\partial x_1} & \cdots & \frac{\partial f_p}{\partial x_1} \\ \frac{\partial f_1}{\partial x_2} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_p}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_1}{\partial x_m} & \frac{\partial f_2}{\partial x_m} & \cdots & \frac{\partial f_p}{\partial x_m} \end{bmatrix} (m \times p)$ ，有 $d\mathbf{f} = \frac{\partial \mathbf{f}}{\partial \mathbf{x}} d\mathbf{x}$

；再定义矩阵的（按列优先）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

$\text{vec}(dF) = \frac{\partial F}{\partial X}^T \text{vec}(dX)$ 。几点说明如下：

1. 按此定义，标量 f 对矩阵 $X(m \times n)$ 的导数 $\frac{\partial f}{\partial X}$ 是 $mn \times 1$ 向量，与上篇的定义不兼容，不过二者容易相互转换。为避免混淆，用记号 $\nabla_X f$ 表示上篇定义的 $m \times n$ 矩阵，则有 $\frac{\partial f}{\partial X} = \text{vec}(\nabla_X f)$ 。虽然本篇的技术可以用于标量对矩阵求导这种特殊情况，但使用上篇中的技术更方便。读者可以通过上篇中的算例试验两种方法的等价转换。
2. 标量对矩阵的二阶导数，又称Hessian矩阵，定义为 $\nabla_X^2 f = \frac{\partial^2 f}{\partial X^2} = \frac{\partial \nabla_X f}{\partial X}$ ($mn \times mn$)，是对称矩阵。对向量 $\frac{\partial f}{\partial X}$ 或矩阵 $\nabla_X f$ 求导都可以得到Hessian矩阵，但从矩阵 $\nabla_X f$ 出发更方便。
3. $\frac{\partial F}{\partial X} = \frac{\partial \text{vec}(F)}{\partial X} = \frac{\partial F}{\partial \text{vec}(X)} = \frac{\partial \text{vec}(F)}{\partial \text{vec}(X)}$ ，求导时矩阵被向量化，弊端是这在一定程度破坏了矩阵的结构，会导致结果变得形式复杂；好处是多元微积分中关于梯度、Hessian矩阵的结论可以沿用过来，只需将矩阵向量化。例如优化问题中，牛顿法的更新 ΔX ，满足 $\text{vec}(\Delta X) = -(\nabla_X^2 f)^{-1} \text{vec}(\nabla_X f)$ 。
4. 在资料中，矩阵对矩阵的导数还有其它定义，比如 $\frac{\partial F}{\partial X} = \left[\frac{\partial F_{kl}}{\partial X} \right]$ ($mp \times nq$)，或是 $\frac{\partial F}{\partial X} = \left[\frac{\partial F}{\partial X_{ij}} \right]$ ($mp \times nq$)，它能兼容上篇中的标量对矩阵导数的定义，但微分与导数的联系（ dF 等于 $\frac{\partial F}{\partial X}$ 中逐个 $m \times n$ 子块分别与 dX 做内积）不够简明，不便于计算和应用。文献[5]综述了以上定义，并批判它们是“坏”定义，能配合微分运算的才是“好”定义（the procedure based on differentials is superior to other methods of differentiation）。

然后来建立运算法则。仍然要利用导数与微分的联系 $\text{vec}(dF) = \frac{\partial F}{\partial X}^T \text{vec}(dX)$ ，求微分的方法与上篇相同，而从微分得到导数需要一些向量化的技巧：

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 表示Kronecker积， $A(m \times n)$ 与 $B(p \times q)$ 的Kronecker积是 $A \otimes B = [A_{ij}B]$ ($mp \times nq$)。此式证明见张贤达《矩阵分析与应用》第107-108页。
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： $\text{vec}(A \odot X) = \text{diag}(A)\text{vec}(X)$ ，其中 $\text{diag}(A)$ (mn×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，对照导数与微分的联系 $\text{vec}(dF) = \frac{\partial F^T}{\partial X} \text{vec}(dX)$ ，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 $d\mathbf{f} = \frac{\partial \mathbf{f}^T}{\partial \mathbf{x}} d\mathbf{x}$ ，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， $\text{vec}(dF) = \frac{\partial F^T}{\partial Y} \text{vec}(dY) = \frac{\partial F^T}{\partial Y} \frac{\partial Y^T}{\partial X} \text{vec}(dX)$ ，可以推出链式法则 $\frac{\partial F}{\partial X} = \frac{\partial Y}{\partial X} \frac{\partial F}{\partial Y}$ 。

和标量对矩阵的导数相比，矩阵对矩阵的导数形式更加复杂，从不同角度出发常会得到形式不同的结果。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. $\text{vec}(\mathbf{a}\mathbf{b}^T) = \mathbf{b} \otimes \mathbf{a}$ 。
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，直接求导得到 $\frac{\partial F}{\partial X} = (AC) \otimes (BD)$ ；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. $K_{pm}(A \otimes B)K_{nq} = B \otimes A$ ，A是m×n矩阵，B是p×q矩阵。可以对 AXB^T 做向量化来证明，一方面， $\text{vec}(AXB^T) = (B \otimes A)\text{vec}(X)$ ；另一方面， $\text{vec}(AXB^T) = K_{pm}\text{vec}(BX^T A^T) = K_{pm}(A \otimes B)\text{vec}(X^T) = K_{pm}(A \otimes B)K_{nq}\text{vec}(X)$ 。

接下来演示一些算例。

例1： $F = AX$ ，X是m×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： $\text{vec}(dF) = \text{vec}(AdX) = (I_n \otimes A)\text{vec}(dX)$ ，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 $\mathbf{f} = A\mathbf{x}$ ，则根据向量的导数与微分的关系 $d\mathbf{f} = \frac{\partial \mathbf{f}^T}{\partial \mathbf{x}} d\mathbf{x}$ ，得到 $\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = A^T$ 。

例2： $f = \log|X|$ ，X是n×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 $\text{vec}(d\nabla_X f) = -K_{nn}\text{vec}(X^{-1}dXX^{-1}) = -K_{nn}(X^{-1T} \otimes X^{-1})\text{vec}(dX)$ ，对照导数与微分的联系，得到 $\nabla_X^2 f = -K_{nn}(X^{-1T} \otimes X^{-1})$ ，注意它是对称矩阵。在X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×m矩阵，X是m×n矩阵，B是n×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：

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

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

$\text{vec}(dF) = (I_p \otimes A)\text{diag}(\exp(XB))(B^T \otimes I_m)\text{vec}(dX)$ ，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】： $l = -\mathbf{y}\mathbf{x}^T\mathbf{w} + \log(1 + \exp(\mathbf{x}^T\mathbf{w}))$ ，求 $\nabla_{\mathbf{w}}l$ 和 $\nabla_{\mathbf{w}}^2l$ 。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 $\nabla_{\mathbf{w}}^2l$ ，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 $\nabla_{\mathbf{w}}^2l = \mathbf{x}\sigma'(\mathbf{x}^T\mathbf{w})\mathbf{x}^T$ 。

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 $\nabla_{\mathbf{w}}^2l$ 。

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

$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$ ，将 l 写成矩阵形式 $l = -\mathbf{y}^T X\mathbf{w} + \mathbf{1}^T \log(\mathbf{1} + \exp(X\mathbf{w}))$ ，进而可以使用上篇中的技

术求得 $\nabla_{\mathbf{w}}l = X^T(\sigma(X\mathbf{w}) - \mathbf{y})$ 。为求 $\nabla_{\mathbf{w}}^2l$ ，先求微分，再用逐元素乘法的技巧：

$d\nabla_{\mathbf{w}}l = X^T(\sigma'(X\mathbf{w}) \odot (X d\mathbf{w})) = X^T \text{diag}(\sigma'(X\mathbf{w}))X d\mathbf{w}$ ，对照导数与微分的联系，得到

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

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 $\nabla_W l$ 和 $\nabla_W^2 l$ 。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 $\mathbf{a} = W\mathbf{x}$ ，

$$d\nabla_W l = \left(\frac{\exp(\mathbf{a}) \odot d\mathbf{a}}{\mathbf{1}^T \exp(\mathbf{a})} - \frac{\exp(\mathbf{a})(\mathbf{1}^T (\exp(\mathbf{a}) \odot d\mathbf{a}))}{(\mathbf{1}^T \exp(\mathbf{a}))^2} \right) \mathbf{x}^T = \left(\frac{\text{diag}(\exp(\mathbf{a}))}{\mathbf{1}^T \exp(\mathbf{a})} - \frac{\exp(\mathbf{a}) \exp(\mathbf{a})^T}{(\mathbf{1}^T \exp(\mathbf{a}))^2} \right) d\mathbf{a} \mathbf{x}^T$$

$= (\text{diag}(\text{softmax}(\mathbf{a})) - \text{softmax}(\mathbf{a})\text{softmax}(\mathbf{a})^T) d\mathbf{a} \mathbf{x}^T$ ，注意这里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 $\mathbf{1}^T (\exp(\mathbf{a}) \odot d\mathbf{a}) = \exp(\mathbf{a})^T d\mathbf{a}$ 。定义矩阵

$D(\mathbf{a}) = \text{diag}(\text{softmax}(\mathbf{a})) - \text{softmax}(\mathbf{a})\text{softmax}(\mathbf{a})^T$ ， $d\nabla_W l = D(\mathbf{a})d\mathbf{a} \mathbf{x}^T = D(W\mathbf{x})dW\mathbf{x} \mathbf{x}^T$ ，

做向量化并使用矩阵乘法的技巧，得到 $\nabla_W^2 l = (\mathbf{x} \mathbf{x}^T) \otimes D(W\mathbf{x})$ 。

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 $d\mathbf{f} = \text{tr}(\nabla_X^T \mathbf{f} dX)$ ，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

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

1. 张贤达. 矩阵分析与应用. 清华大学出版社有限公司, 2004.
2. Fackler, Paul L. "Notes on matrix calculus." *North Carolina State University*(2005).
3. Petersen, Kaare Brandt, and Michael Syskind Pedersen. "The matrix cookbook." *Technical University of Denmark* 7 (2008): 15.
4. HU, Pili. "Matrix Calculus: Derivation and Simple Application." (2012).

5. Magnus, Jan R., and Heinz Neudecker. "Matrix differential calculus with applications to simple, Hadamard, and Kronecker products." *Journal of Mathematical Psychology* 29.4 (1985): 474-492.

编辑于昨天 03:04

- 矩阵分析
- 机器学习
- 优化

118 条评论

⇌ 切换为时间排序

写下你的评论...

😊

 James Liu2 年前


先赞为敬

👍 1

 方不觉2 年前


赞

👍 赞

 王赞 Maigo2 年前


矩阵对矩阵求导的结果中有好多Kronecker积啊.....不过也没办法，把一个本来是四维的东西压成两维已经不容易了。

👍 8

 知乎用户2 年前


我只是想吐槽知乎这个狗屎编辑器到现在都不支持 MathJax，真是废物。怕是坦桑尼亚的民族风情网站都支持 MathJax 了。可能知乎的前端一没高分屏，二没大学文凭。

👍 13

 Shinji Fujiwara2 年前

赞

👍 赞

 排骨郎2 年前

例3中，与B的转置做 kronecker product 的单位矩阵应该是 I_m 不是 I_n



陌烛

1 年前

你好，请问下，原文中有句话：“若矩阵函数F是矩阵X经加减乘法、行列式、逆、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vec(dX)左侧，即能得到导数”，那么如果F是由X卷积操作得到的，那么，对于这个卷积的运算法则是什么呢？👀👀👀👀👀👀👀

👍 赞



长躯鬼侠 (作者) 回复 陌烛

1 年前

对于卷积，你可以自己推导一下，运算法则也可以用卷积来表示，对full、valid模式在细节上有些差异。

👍 1



陌烛

1 年前

你好，我还想知道下，克罗内克积和矩阵乘积，哪个的优先级大啊？😓

👍 赞



长躯鬼侠 (作者) 回复 陌烛

1 年前

我没有指定Kronecker积和矩阵乘积哪个优先级高，所以都加括号了啊。

👍 赞



陌烛 回复 长躯鬼侠 (作者)

1 年前

蟹蟹，我果然不适合推导数学公式😭

👍 赞

[查看全部 9 条回复](#)



孙培钦

1 年前

我想问一下，假设我有个等式 $S = WX$ ， S 是 $m \times 1$ 向量， X 是 $n \times 1$ 向量， W 是 $m \times n$ 矩阵，使用上述的求导术去求 S 向量对 X 向量的导数，我得到的是一个 $n \times n$ 的单位阵与 W 转置的kronecker积啊，这个积的尺寸应该是 $nn \times mn$ ，明显不对啊。

👍 赞



长躯鬼侠 (作者) 回复 孙培钦

1 年前

$ds = Wdx$ ， x 的右面是 1×1 的单位阵，不是 $n \times n$ 的。

👍 1



暮日落流年

1 年前

收益匪浅，赞一个！

👍 赞



知乎用户

1 年前

能区分一下行列求导就更完美了，同时使用线性变换解释一下导数与微分的关系，可能更加恰当，和利于直觉。

👍 1



neilfvhv

1 年前

问个问题， $\frac{\partial F}{\partial X} = \frac{\partial \text{vec} F}{\partial \text{vec} X}$ 这一步是怎么得到的？

👍 赞



长躯鬼侠 (作者) 回复 neilfvhv

1 年前

这是定义

👍 赞



张云

1 年前

你好，不怎么懂那个例四推广里的log的和是怎么转化为矩阵形式的，前一半可以理解，后一半就有点懵了

👍 赞

 长躯鬼侠 (作者) 回复 张云

1 年前

你是说 $\sum_i \log(1 + \exp(\vec{x}_i^T \vec{w})) = \vec{1}^T \log(\vec{1} + \exp(X\vec{w}))$ 这个嘛？因为 $\log(1 + \exp(\vec{x}_i^T \vec{w}))$ 是向量 $\log(\vec{1} + \exp(X\vec{w}))$ 的第i个元素，左边的求和就等于右边的 $\vec{1}^T$ 乘以向量。

👍 赞

 张云 回复 长躯鬼侠 (作者)

1 年前

懂了。那个1的向量的转置是为了把矩阵里面的每个元素都求和。

哇，这些算是本科学的吗？感觉有点酷炫诶😂

👍 赞

展开其他 1 条回复

