

矩阵求导术 (上)



长躯鬼侠
数学爱好者

编辑推荐

6,687 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：

$$df = \sum_{i=1}^n \frac{\partial f}{\partial x_i} dx_i = \frac{\partial f}{\partial \mathbf{x}}^T d\mathbf{x}, \text{ 这里第一个等号是全微分公式，第二个等号表达了梯度与微}$$

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

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

(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**。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

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。

▲ 赞同 6687

● 461 条评论

➤ 分享

♥ 喜欢

★ 收藏

📄 申请转载

...

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： $d(X \pm Y) = dX \pm dY$ ；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： $dX^{-1} = -X^{-1}dXX^{-1}$ 。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 $df = \text{tr} \left(\frac{\partial f}{\partial X}^T dX \right)$ ，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 $df = \text{tr} \left(\frac{\partial f}{\partial X}^T dX \right)$ ，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 $df = \text{tr} \left(\frac{\partial f}{\partial Y}^T dY \right)$ ，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，即可得到 $\frac{\partial f}{\partial X}$ 。

最常见的情形是 $Y = AXB$ ，此时

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

，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： $f = \mathbf{a}^T X \mathbf{b}$ ，求 $\frac{\partial f}{\partial X}$ 。其中 \mathbf{a} 是 $m \times 1$ 列向量， 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， $df = d\mathbf{a}^T X \mathbf{b} + \mathbf{a}^T dX \mathbf{b} + \mathbf{a}^T X d\mathbf{b} = \mathbf{a}^T dX \mathbf{b}$ ，注意这里的 \mathbf{a}, \mathbf{b} 是常量， $d\mathbf{a} = 0, d\mathbf{b} = 0$ 。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： $df = \text{tr}(\mathbf{a}^T dX \mathbf{b}) = \text{tr}(\mathbf{b} \mathbf{a}^T dX) = \text{tr}((\mathbf{a} \mathbf{b}^T)^T dX)$ ，注意这里我们根据 $\text{tr}(AB) = \text{tr}(BA)$ 交换了 $\mathbf{a}^T dX$ 与 \mathbf{b} 。对照导数与微分的联系

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

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: $df = \mathbf{a}^T (\exp(\mathbf{X}\mathbf{b}) \odot (d\mathbf{X}\mathbf{b}))$, 再套上迹并做交换: $df = \text{tr}(\mathbf{a}^T (\exp(\mathbf{X}\mathbf{b}) \odot (d\mathbf{X}\mathbf{b}))) = \text{tr}((\mathbf{a} \odot \exp(\mathbf{X}\mathbf{b}))^T d\mathbf{X}\mathbf{b})$
 $= \text{tr}(\mathbf{b}(\mathbf{a} \odot \exp(\mathbf{X}\mathbf{b}))^T d\mathbf{X}) = \text{tr}(((\mathbf{a} \odot \exp(\mathbf{X}\mathbf{b}))\mathbf{b}^T)^T d\mathbf{X})$, 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} 、 $\exp(\mathbf{X}\mathbf{b})$ 与 $d\mathbf{X}\mathbf{b}$, 再根据 $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$ 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
 $df = \text{tr} \left(\frac{\partial f}{\partial \mathbf{X}}^T d\mathbf{X} \right)$, 得到 $\frac{\partial f}{\partial \mathbf{X}} = (\mathbf{a} \odot \exp(\mathbf{X}\mathbf{b}))\mathbf{b}^T$ 。

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 $l \times m$ 矩阵, \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, σ 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:

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

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

照导数与微分的联系, 得到

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

$m \times 1$ 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

$dl = \frac{\partial l}{\partial w}^T dw$, 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 Σ 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, Σ 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

$\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^T d\Sigma^{-1} (x_i - \bar{x}) = -\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^T \Sigma^{-1} d\Sigma \Sigma^{-1} (x_i - \bar{x})$ 。再

给第二项套上迹做交换: $\text{tr} \left(\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^T \Sigma^{-1} d\Sigma \Sigma^{-1} (x_i - \bar{x}) \right)$

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

$= \frac{1}{N} \sum_{i=1}^N \text{tr}(\Sigma^{-1} (x_i - \bar{x})(x_i - \bar{x})^T \Sigma^{-1} d\Sigma) = \text{tr}(\Sigma^{-1} S \Sigma^{-1} d\Sigma)$, 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

知乎

首发于
深度学习于图网络

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：

$$dl = -\mathbf{y}^T dW\mathbf{x} + \frac{\mathbf{1}^T (\exp(W\mathbf{x}) \odot (dW\mathbf{x}))}{\mathbf{1}^T \exp(W\mathbf{x})}。再套上述并做交换，注意可化简$$

$\mathbf{1}^T (\exp(W\mathbf{x}) \odot (dW\mathbf{x})) = \exp(W\mathbf{x})^T dW\mathbf{x}$ ，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

σ 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,

$$dl = \text{tr} \left(\frac{\partial l}{\partial \mathbf{a}_2}^T d\mathbf{a}_2 \right) = \text{tr} \left(\frac{\partial l}{\partial \mathbf{a}_2}^T dW_2 \mathbf{h}_1 \right) + \underbrace{\text{tr} \left(\frac{\partial l}{\partial \mathbf{a}_2}^T W_2 d\mathbf{h}_1 \right)}_{dl_2}, \text{ 使用矩阵乘}$$

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法则来求 $\frac{\partial l}{\partial \mathbf{a}_1}$, 并利用矩阵乘法和逐元素乘法交换的迹技巧:

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

, 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:

$$dl_2 = \text{tr} \left(\frac{\partial l}{\partial \mathbf{a}_1}^T d\mathbf{a}_1 \right) = \text{tr} \left(\frac{\partial l}{\partial \mathbf{a}_1}^T dW_1 \mathbf{x} \right) = \text{tr} \left(\mathbf{x} \frac{\partial l}{\partial \mathbf{a}_1}^T dW_1 \right), \text{ 得到}$$

$$\frac{\partial l}{\partial W_1} = \frac{\partial l}{\partial \mathbf{a}_1} \mathbf{x}^T.$$

推广: 样本 $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)$,

$$l = - \sum_{i=1}^N \mathbf{y}_i^T \log \text{softmax}(W_2 \sigma(W_1 \mathbf{x}_i + \mathbf{b}_1) + \mathbf{b}_2), \text{ 其中 } \mathbf{b}_1 \text{ 是 } p \times 1 \text{ 列向量, } \mathbf{b}_2$$

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

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

复合法则,

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

三项得到到 $\frac{\partial l}{\partial \mathbf{b}_2} = \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{2,i}}$ 。接下来对第二项继续使用复合法则，得到

$$\frac{\partial l}{\partial \mathbf{a}_{1,i}} = \frac{\partial l}{\partial \mathbf{h}_{1,i}} \odot \sigma'(\mathbf{a}_{1,i})。为求 \frac{\partial l}{\partial \mathbf{W}_1}, \frac{\partial l}{\partial \mathbf{b}_1}, 再用一次复合法则:$$

$$dl_2 = \text{tr} \left(\sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{1,i}}^T d\mathbf{a}_{1,i} \right) = \text{tr} \left(\sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{1,i}}^T d\mathbf{W}_1 \mathbf{x}_i \right) + \text{tr} \left(\sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{1,i}}^T d\mathbf{b}_1 \right)$$

, 得到 $\frac{\partial l}{\partial \mathbf{W}_1} = \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{1,i}} \mathbf{x}_i^T, \frac{\partial l}{\partial \mathbf{b}_1} = \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{1,i}}$ 。

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

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

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

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

$$dl = \text{tr} \left(\frac{\partial l}{\partial \mathbf{A}_2}^T d\mathbf{A}_2 \right) = \text{tr} \left(\frac{\partial l}{\partial \mathbf{A}_2}^T d\mathbf{W}_2 \mathbf{H}_1 \right) + \underbrace{\text{tr} \left(\frac{\partial l}{\partial \mathbf{A}_2}^T \mathbf{W}_2 d\mathbf{H}_1 \right)}_{dl_2} + \text{tr} \left(\frac{\partial l}{\partial \mathbf{A}_2}^T d\mathbf{b}_2 \mathbf{1}^T \right)$$

, 从第一项得到 $\frac{\partial l}{\partial \mathbf{W}_2} = \frac{\partial l}{\partial \mathbf{A}_2} \mathbf{H}_1^T$, 从第二项得到 $\frac{\partial l}{\partial \mathbf{H}_1} = \mathbf{W}_2^T \frac{\partial l}{\partial \mathbf{A}_2}$, 从第三项得到到

$$\frac{\partial l}{\partial \mathbf{b}_2} = \frac{\partial l}{\partial \mathbf{A}_2} \mathbf{1}。接下来对第二项继续使用复合法则, 得到 \frac{\partial l}{\partial \mathbf{A}_1} = \frac{\partial l}{\partial \mathbf{H}_1} \odot \sigma'(\mathbf{A}_1)。为$$

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

$$dl_2 = \text{tr} \left(\frac{\partial l}{\partial \mathbf{A}_1}^T d\mathbf{A}_1 \right) = \text{tr} \left(\frac{\partial l}{\partial \mathbf{A}_1}^T d\mathbf{W}_1 \mathbf{X} \right) + \text{tr} \left(\frac{\partial l}{\partial \mathbf{A}_1}^T d\mathbf{b}_1 \mathbf{1}^T \right), \text{ 得到}$$

$$\frac{\partial l}{\partial \mathbf{W}_1} = \frac{\partial l}{\partial \mathbf{A}_1} \mathbf{X}^T, \quad \frac{\partial l}{\partial \mathbf{b}_1} = \frac{\partial l}{\partial \mathbf{A}_1} \mathbf{1}。$$

下篇见 zhuanlan.zhihu.com/p/24709748。

编辑于 2020-03-06

文章被以下专栏收录



深度学习于图网络
不定时更新图网络学习笔记



数学

推荐阅读

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 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

矩阵求导浅析（一）

本文主要关注标量函数对矩阵的求导，并提供一种简明直观易操作的矩阵求导方法。推荐矩阵求导的专栏文章：矩阵求导术（上）矩阵求导术（下）机器学习中的向量求导 1.内积 向量 ...

倚楼

461 条评论

⇌ 切换为时间排序

写下你的评论...



resurrectcore

2017-01-07

网上有个pdf叫做 The Matrix Cookbook.

👍 74



知乎用户 回复 resurrectcore

2017-01-07

我看过，相信写这篇文章的人也看过。感觉这里的方法比硬背公式简单，而且对于搞机

matrix COOKBOOK是当之无愧的。

👍 2

查看全部 9 条回复



王赞 Maigo

2017-01-07

赞啊！终于找到系统的方法了！

👍 59



方不觉

2017-01-07

nbnb，一直对矩阵求导感觉无从下手，这几条法则比背公式好记多了！

👍 30



知乎用户

2017-01-09

参考文献大概是：Matrix Calculus: Derivation and Simple Application HU Pili

👍 15



知乎用户

2019-07-27

matrixcalculus.org 曾经打算记过 后来不怎么用就放弃了 要用的时候直接查公式（逃

👍 10



黑与白 回复 知乎用户

2019-10-15

大佬，请问这个网站是什么网站呀

👍 赞



DreamYun 回复 知乎用户

2020-04-28

这个牛逼了，不得了！！

👍 赞

展开其他 3 条回复



知乎用户

2017-01-23

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！

👍 6

算

👍 6



程许 回复 maja

2020-08-08

我知道你的意思是，假如 a, b, x 都是 $n+1$ 维矩阵，那么从 $a^T x = b^T x$ 并不能必然推出 $a = b$ ，因为存在恰好是位于零空间的可能。我想这里的gap可能可以类似于多项式插值定理来弥补：如果 $f(x) = g(x)$ 在 $n+2$ 个不同点上都成立，而且 f, g 的次数都不超过 n 次，那么 f 和 g 恒等。不过作者没有在本文证明或举出这种类型的定理，所以确实可以认为文章的推导只是intuitive，并不严格。

👍 赞



长躯鬼侠 (作者) 回复 程许

2020-08-08

给定 $a^T x = b^T x$ 对任意 x 成立，是可以推出 $a = b$ 的。

👍 3



知乎用户

2017-01-07

写的真好，醍醐灌顶

👍 6



lang的飞起

2017-11-07

刚开始学，学习了。那个线性回归那里，怎么一步就推出2倍了，没看明白

👍 5



长躯鬼侠 (作者) 回复 lang的飞起

2017-11-07

两个向量的内积， $a^T b = b^T a$ ，加起来就是2倍了

👍 2



虚有

2020-09-30

两年前看:这是什么玩意儿?

两年后:卧槽，这不就是矩阵理论里没讲的章节吗? np。

给大佬点赞，讲的太好太详细了，看了这么多本矩阵理论的教材，不是缺这个就是缺那个。还是您讲的好。

👍 4



知乎用户

2017-02-27



相思作坊半世離殤 回复 知乎用户

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](#)

Zzzzzzz

2017-01-07

行列式微分的那个可以把行列式放进tr里变成伴随矩阵。行列式对其中元素的偏导显然就是其代数余子式，应该不用行列式不为0。这公式叫jacobi's formula

👍 4



长躯鬼侠 (作者) 回复 Zzzzzzz

2017-01-07

嗯，你说得对。

👍 赞



知乎用户 回复 Zzzzzzz

2017-01-10

他的意思是当行列式不为零时，可以把伴随矩阵写成行列式乘以逆吧

👍 赞



maja

2017-01-08

没有回应，可能你没看懂我的意思：这样吧我们从梯度下降流来考虑：

1) 全微分的出发点是没有问题的，但用迹是没有解决问题的，问题的本质是如何定义符号“ df / dX ”， f 必然是 n^2 元函数，根据标量对 n 元函数有 $df = \sum (\partial f / \partial X_{ij}) dX_{ij}$ 如果定义了 dX 为矩阵形式就决定了 df / dX 。如果用矩阵内积运算 则 就确定了 df / dX 的矩阵形式 衡量的就是 df 对于 dX 变化率。

回到基础问题— $f = a^T \cdot X \cdot b$ 现在我们要根据 $df = df/dX \cdot dX$ 来求解. 所以要确定 df 的形式 矩阵乘法始终是一个记号！和行列式通过伴随矩阵建立联系. 通俗的运算就是加 乘 = 定义记号-累乘- $S(i, a) = \sum_{j=1}^n a_j$ 则 $f = S(i, a) S(k, x_{ik} \cdot b_k)$ 然后两边对 X 诸元素全微分, $df = S(<i, k>, dX) \{S(p, a) \cdot S(q, b)\}$. 提取公因式, $df = a^T \cdot b \cdot dX \cdot 1$

这个是我们可以用矩阵描述梯度下降流的朴素原理，现在考虑第二个问题：

这就是为什么呢答主用迹会出错，因为就不对。

👍 3



知乎用户 回复 maja

2020-09-05

楼主这一套方法，来源应该就是张贤达的《矩阵分析与应用》里面的，方法从计算层面上应该是没问题的，但是一个潜在的问题就是，这一套方法和用到jacobi或者符合函数的时候，总是会差一个转置，使用 df / dx 行向量的约定，jacobi和符合求导能够更简洁，也和流形上的微积分的符号一致。而像这里的全部都用列向量的情况，和流形上微积分的公式就不一致，差转置

👍 2



平常心 回复 知乎用户

2020-12-19

对

👍 赞

展开其他 2 条回复



紫杉

2017-01-07

期待下集。。我都忘了很多这部分内容了

👍 3



阳光陈靖文 回复 紫杉

01-22

去主页把链接贴过来了[长躯鬼侠：矩阵求导术（下）](#)不用谢



👍 2



谭含章

2017-01-07

写的很清楚，很实用！

👍 3



白夜行

2020-07-04



2



柴鹰

2017-12-28

虽然写了很多，但是看得不是很懂。看来真的是要靠背公式活下去了。

2



黑与白 回复 柴鹰

2019-10-15

我能感觉到答主写的很有用，但是我也确实没有理解，而现在又着急使用，请问您是背哪些公式？可以分享一下吗

赞



老大 回复 黑与白

03-29

我猜你俩是没看懂全微分公式，导致后面的推导没法convince你自己，其实楼主确实全微分公式那里省略了一点，不过可以自己写写看

赞



知乎用户

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

2



甄景贤

2017-01-08

这个超有用，功德无量 :)

2



CoolMan

2020-05-18

牛逼牛逼，有些比张贤达的书写得还清楚

1

1 2 3 4 ... 15 下一页