

# 矩阵求导

## • 标量对矩阵求导

- 一元微分学中倒数与微分的联系

$$df = f'(x) dx$$

- 将其推广到向量

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

即，在向量中全微分 $df$ 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

- 在矩阵中微分与矩阵倒数的关系可以写成

$$df = \sum_{i=1}^m \sum_{j=1}^n \frac{df}{dx_{i,j}} dx_{i,j} \text{ 变成乘积模式为 } df = \text{tr} \left( \frac{df}{d\mathbf{X}}^T * d\mathbf{X} \right)$$

即，在矩阵中全微分 $df$ 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

- 所以根据标量对向量与矩阵的求导形式可变形替换得到结果

- 例4【线性回归】： $l = \|X\mathbf{w} - \mathbf{y}\|^2$ ，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 $m \times 1$ 列向量， $X$ 是 $m \times n$ 矩阵， $\mathbf{w}$ 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： $l = (X\mathbf{w} - \mathbf{y})^T (X\mathbf{w} - \mathbf{y})$ ，求微分，使用矩阵乘法、转置等法则：

$dl = (X d\mathbf{w})^T (X\mathbf{w} - \mathbf{y}) + (X\mathbf{w} - \mathbf{y})^T (X d\mathbf{w}) = 2(X\mathbf{w} - \mathbf{y})^T X d\mathbf{w}$ ，注意这里 $X d\mathbf{w}$ 和 $X\mathbf{w} - \mathbf{y}$ 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

$$dl = \frac{\partial l}{\partial \mathbf{w}}^T d\mathbf{w}, \text{ 得到 } \frac{\partial l}{\partial \mathbf{w}} = 2X^T (X\mathbf{w} - \mathbf{y}). \frac{\partial l}{\partial \mathbf{w}} = \mathbf{0} \text{ 即 } X^T X\mathbf{w} = X^T \mathbf{y},$$

得到 $\mathbf{w}$ 的最小二乘估计为 $\mathbf{w} = (X^T X)^{-1} X^T \mathbf{y}$ 。

例5【方差的最大似然估计】：样本 $\mathbf{x}_1, \dots, \mathbf{x}_N \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ ，求方差 $\Sigma$ 的最大似然估计。

写成数学式是： $l = \log |\Sigma| + \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^T \Sigma^{-1} (\mathbf{x}_i - \bar{\mathbf{x}})$ ，求 $\frac{\partial l}{\partial \Sigma}$ 的零点。其

中 $\mathbf{x}_i$ 是 $m \times 1$ 列向量， $\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$ 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

$$\frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^T d\Sigma^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}) = -\frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^T \Sigma^{-1} d\Sigma \Sigma^{-1} (\mathbf{x}_i - \bar{\mathbf{x}})$$

$$\begin{aligned} & \text{。再给第二项套上迹做交换：} \text{tr} \left( \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^T \Sigma^{-1} d\Sigma \Sigma^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}) \right) \\ &= \frac{1}{N} \sum_{i=1}^N \text{tr}((\mathbf{x}_i - \bar{\mathbf{x}})^T \Sigma^{-1} d\Sigma \Sigma^{-1} (\mathbf{x}_i - \bar{\mathbf{x}})) \\ &= \frac{1}{N} \sum_{i=1}^N \text{tr}(\Sigma^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \Sigma^{-1} d\Sigma) = \text{tr}(\Sigma^{-1} S \Sigma^{-1} d\Sigma), \text{ 其中先} \end{aligned}$$

交换迹与求和，然后将  $\Sigma^{-1} (\mathbf{x}_i - \bar{\mathbf{x}})$  交换到左边，最后再交换迹与求和，并定义

$$S = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \text{ 为样本方差矩阵。得到}$$

$dl = \text{tr}((\Sigma^{-1} - \Sigma^{-1} S \Sigma^{-1}) d\Sigma)$ 。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：

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

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

$$\mathbf{1}^T (\mathbf{u} \odot \mathbf{v}) = \mathbf{u}^T \mathbf{v}, \text{ 故}$$

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

。对照导数与微分的联系，得到结果

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:

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

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

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

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

, 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

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

使用复合法则,

$$dl = \text{tr} \left( \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{2,i}}^T d\mathbf{a}_{2,i} \right) = \text{tr} \left( \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{2,i}}^T dW_2 \mathbf{h}_{1,i} \right) + \underbrace{\text{tr} \left( \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{2,i}}^T W_2 d\mathbf{h}_{1,i} \right)}_{dl_2} + \text{tr} \left( \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{2,i}}^T d\mathbf{b}_2 \right)$$

, 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:

$$dl_2 = \text{tr} \left( \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{1,i}}^T d\mathbf{a}_{1,i} \right) = \text{tr} \left( \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{1,i}}^T dW_1 \mathbf{x}_i \right) + \text{tr} \left( \sum_{i=1}^N \frac{\partial l}{\partial \mathbf{a}_{1,i}}^T d\mathbf{b}_1 \right)$$

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

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

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

,  $\mathbf{A}_2 = [\mathbf{a}_{2,1}, \dots, \mathbf{a}_{2,N}] = W_2 \mathbf{H}_1 + \mathbf{b}_2 \mathbf{1}^T$ , 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,

$$dl = \text{tr} \left( \frac{\partial l}{\partial \mathbf{A}_2}^T d\mathbf{A}_2 \right) = \text{tr} \left( \frac{\partial l}{\partial \mathbf{A}_2}^T dW_2 \mathbf{H}_1 \right) + \underbrace{\text{tr} \left( \frac{\partial l}{\partial \mathbf{A}_2}^T W_2 d\mathbf{H}_1 \right)}_{dl_2} + \text{tr} \left( \frac{\partial l}{\partial \mathbf{A}_2}^T d\mathbf{b}_2 \mathbf{1}^T \right)$$

, 从第一项得到  $\frac{\partial l}{\partial W_2} = \frac{\partial l}{\partial \mathbf{A}_2} \mathbf{H}_1^T$ , 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  $\frac{\partial l}{\partial \mathbf{A}_1} = \frac{\partial l}{\partial \mathbf{H}_1} \odot \sigma'(\mathbf{A}_1)$ 。

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

$$dl_2 = \text{tr} \left( \frac{\partial l}{\partial \mathbf{A}_1}^T d\mathbf{A}_1 \right) = \text{tr} \left( \frac{\partial l}{\partial \mathbf{A}_1}^T dW_1 \mathbf{X} \right) + \text{tr} \left( \frac{\partial l}{\partial \mathbf{A}_1}^T d\mathbf{b}_1 \mathbf{1}^T \right), \text{ 得}$$

到  $\frac{\partial l}{\partial W_1} = \frac{\partial l}{\partial \mathbf{A}_1} \mathbf{X}^T$ ,  $\frac{\partial l}{\partial \mathbf{b}_1} = \frac{\partial l}{\partial \mathbf{A}_1} \mathbf{1}$ 。

## · 矩阵对矩阵求导

- 向量 $m$ 对向量 $n$ 的求导得到一个 $m \times n$ 的矩阵
- 对于矩阵，可以进行行/列优先的向量转化
- 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 $\frac{\partial F}{\partial X} = \frac{\partial \text{vec}(F)}{\partial \text{vec}(X)}$  ( $mn \times pq$ )。
- 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： $\text{vec}(dF) = \text{vec}(AdX) = (I_n \otimes A)\text{vec}(dX)$ 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 $\frac{\partial f}{\partial x} = A^T$ 。