

矩阵求导

3种标准导数(梯度)公式

1) 自变量是一个标量(Scalar)时:

$$Df(x) = \lim_{t \rightarrow 0} \frac{f(x+t) - f(x)}{t}$$

2) 自变量是一个向量(Vector)时:

$$D_w f(\mathbf{x}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{w}) - f(\mathbf{x})}{t}$$

(w的维数和x一致) 这个导数的含义是, 在n维空间中f(x)所定义的(超)平面上的某个坐标点x相对于w的斜率。

3) 自变量是一个矩阵(Matrix)时:

$$D_W f(\mathbf{X}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{X} + t\mathbf{W}) - f(\mathbf{X})}{t}$$

含义和2)类似。(已经无法想象了)

标量f对矩阵X的导数-Trace derivative

标量f对矩阵X的导数, 定义为:

$$\frac{\partial f}{\partial \mathbf{X}} = \left[\frac{\partial f}{\partial X_{ij}} \right]$$

矩阵导数与微分建立联系:

$$df = \sum_{i=1}^m \sum_{j=1}^n \frac{\partial f}{\partial X_{ij}} dX_{ij} = \text{tr} \left(\frac{\partial f}{\partial \mathbf{X}}^T d\mathbf{X} \right)$$

矩阵运算法则

- 加减法: $d(X \pm Y) = dX \pm dY$; 矩阵乘法: $d(XY) = (dX)Y + X(dY)$; 转置: $d(X^T) = (dX)^T$; 迹: $d\text{tr}(X) = \text{tr}(dX)$ 。
- 逆: $dX^{-1} = -X^{-1}(dX)X^{-1}$ 。此式可在 $XX^{-1} = I$ 两侧求微分来证明。
- 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。
- 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。
- 逐元素函数: $d\sigma(X) = \sigma'(X) \odot dX$, $\sigma(X) = [\sigma(X_{ij})]$ 是逐元素运算的标量函数。

Trace trick

- 标量套上迹: $a = \text{tr}(a)$
- 转置: $\text{tr}(A^T) = \text{tr}(A)$ 。
- 线性: $\text{tr}(A \pm B) = \text{tr}(A) \pm \text{tr}(B)$ 。
- 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 $\sum_{i,j} A_{ij}B_{ji}$ 。
- 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 $\sum_{i,j} A_{ij}B_{ij}C_{ij}$ 。

Example:

【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: 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: $l = (X\mathbf{w} - \mathbf{y})^T (X\mathbf{w} - \mathbf{y})$, 求微分, 使用矩阵乘法、转置等法则 $d(X^T Y) = (dX)^T Y + X^T (dY)$:

$$dl = (Xd\mathbf{w})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + (\mathbf{X}\mathbf{w} - \mathbf{y})^T (Xd\mathbf{w}) = 2(\mathbf{X}\mathbf{w} - \mathbf{y})^T X d\mathbf{w}$$

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

$$\mathbf{w} = (X^T X)^{-1} X^T \mathbf{y}$$

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

解：首先将softmax函数代入并写成

$$l = -\mathbf{y}^T \left(\log(\exp(W\mathbf{x})) - \mathbf{1} \log(\mathbf{1}^T \exp(W\mathbf{x})) \right) = -\mathbf{y}^T W\mathbf{x} + \log(\mathbf{1}^T \exp(W\mathbf{x}))$$

，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

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

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

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

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

，得到

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

【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； $\text{softmax}(\mathbf{a}) = \frac{\exp(\mathbf{a})}{\mathbf{1}^T \exp(\mathbf{a})}$ 同例2， $\sigma(\cdot)$ 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 $l = -\mathbf{y}^T \log \text{softmax}(\mathbf{a}_2)$ 。在例2中已求出 $\frac{\partial l}{\partial \mathbf{a}_2} = \text{softmax}(\mathbf{a}_2) - \mathbf{y}$ 。使用复合法则，注意此处 \mathbf{h}_1, W_2 都是变量：

$$dl = \text{tr} \left(\frac{\partial l}{\partial \mathbf{a}_2} d\mathbf{a}_2 \right) = \text{tr} \left(\frac{\partial l}{\partial \mathbf{a}_2} dW_2 \mathbf{h}_1 \right) + \text{tr} \left(\frac{\partial l}{\partial \mathbf{a}_2} W_2 d\mathbf{h}_1 \right)$$

，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 $\frac{\partial l}{\partial \mathbf{h}_1} = W_2^T \frac{\partial l}{\partial \mathbf{a}_2}$ 。接下来求 $\frac{\partial l}{\partial \mathbf{a}_1}$ ，继续使用复合法则，并利用矩阵乘法和逐元素乘法交换的迹技巧：

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

，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：

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

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

矩阵 $F(p \times q)$ 对矩阵 $X(m \times n)$ 的导数