



華東師範大學
EAST CHINA NORMAL UNIVERSITY

数学基础

微积分：向量和矩阵求导

考虑 $\frac{dA}{dB}$, A 和 B 都可能是标量、向量或矩阵, 共有9种不同的导数

自变量\因变量	标量 ϕ	向量 f	矩阵 F
标量 ξ	$\frac{d\phi}{d\xi}$	$\frac{df}{d\xi}$	$\frac{dF}{d\xi}$
向量 x	$\frac{d\phi}{dx}$	$\frac{df}{dx}$	$\frac{dF}{dx}$
矩阵 X	$\frac{d\phi}{dX}$	$\frac{df}{dX}$	$\frac{dF}{dX}$

微积分：向量和矩阵求导

一、向量 f 对标量 ξ 的导数 $\frac{df}{d\xi}$

定义1 对于 n 维向量函数(vector function)

$$f(\xi) = (f_1(\xi), f_2(\xi), \dots, f_n(\xi))^T$$

定义它对 ξ 的导数为：

$$\frac{df(\xi)}{d\xi} \triangleq \left(\frac{df_1(\xi)}{d\xi}, \frac{df_2(\xi)}{d\xi}, \dots, \frac{df_n(\xi)}{d\xi} \right)^T$$

也可以对行向量求导：

$$\frac{df(\xi)^T}{d\xi} \triangleq \left(\frac{df_1(\xi)}{d\xi}, \frac{df_2(\xi)}{d\xi}, \dots, \frac{df_n(\xi)}{d\xi} \right)$$

微积分：向量和矩阵求导

二、矩阵 F 对标量 ξ 的导数 $\frac{dF}{d\xi}$

定义2. 对于 $n \times m$ 维矩阵函数(matrix function)

$$F(\xi) = \begin{pmatrix} F_{11}(\xi), & F_{12}(\xi), & \cdots & F_{1m}(\xi) \\ \vdots & \vdots & \ddots & \vdots \\ F_{n1}(\xi), & F_{n2}(\xi), & \cdots & F_{nm}(\xi) \end{pmatrix} = (F_{ij}(\xi))_{nm}$$

定义它对 ξ 的导数为：

$$\begin{aligned} \frac{dF(\xi)}{d\xi} &\triangleq \begin{pmatrix} \frac{dF_{11}(\xi)}{d\xi}, & \frac{dF_{12}(\xi)}{d\xi}, & \cdots & \frac{dF_{1m}(\xi)}{d\xi} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{dF_{n1}(\xi)}{d\xi}, & \frac{dF_{n2}(\xi)}{d\xi}, & \cdots & \frac{dF_{nm}(\xi)}{d\xi} \end{pmatrix} \\ &= \left(\frac{dF_{ij}(\xi)}{d\xi} \right)_{nm} \end{aligned}$$

上述两个定义是
统一的！

微积分：向量和矩阵求导

运算性质：

(1) 加法运算公式

$$\frac{d(F \pm G)}{d\xi} = \frac{dF}{d\xi} \pm \frac{dG}{d\xi}$$

(2) 数乘运算公式 ($\lambda(\xi)$ 是标量函数)

$$\frac{d(\lambda(\xi)F)}{d\xi} = \frac{d\lambda(\xi)}{d\xi}F + \lambda(\xi)\frac{dF}{d\xi}$$

(3) 乘法运算公式

$$\frac{d(FG)}{d\xi} = \frac{dF}{d\xi}G + F\frac{dG}{d\xi}$$

例1：求 $x^T A x$ 对 ξ 的导数，其中 x 是 ξ 的 n 维向量函数， A 是 $n \times n$ 对称常数矩阵。

微积分：向量和矩阵求导

三、标量 ϕ 对向量 x 的导数 $\frac{d\phi}{dx}$

设函数 $\phi(x) = \phi(x_1, \dots, x_n)$ 是以 x 为自变量的数量函数
定义3. 标量函数 ϕ 对列向量 x 的导数为：

$$\frac{d\phi}{dx} \triangleq \left(\frac{\partial \phi}{\partial x_1}, \frac{\partial \phi}{\partial x_2}, \dots, \frac{\partial \phi}{\partial x_n} \right)^T$$

也称为函数 ϕ 的梯度，记为 $\text{grad } \phi$ 或 $\nabla \phi$

也可以对行向量求导：

$$\frac{d\phi}{dx^T} \triangleq \left(\frac{\partial \phi}{\partial x_1}, \frac{\partial \phi}{\partial x_2}, \dots, \frac{\partial \phi}{\partial x_n} \right)$$

显然：
$$\frac{d\phi(x)}{dx} = \left(\frac{d\phi(x)}{dx^T} \right)^T$$

微积分：向量和矩阵求导

运算性质：

(1) 加法运算公式

$$\frac{d(\varphi \pm \phi)}{dx} = \frac{d\varphi}{dx} \pm \frac{d\phi}{dx}$$

(2) 乘法运算公式

$$\frac{d(\varphi\phi)}{dx} = \frac{d\varphi}{dx}\phi + \varphi\frac{d\phi}{dx}$$

例2： 求函数 $\phi(x) = x^T x = x_1^2 + x_2^2 + \cdots + x_n^2$ 对 x 的导数。

微积分：向量和矩阵求导

四、向量 f 对向量 x 的导数 $\frac{df}{dx}$

设函数 $f(x) = (f_1(x), \dots, f_m(x))^T$ 是 x 的 m 维列向量函数

定义4. $n \times m$ 阶矩阵函数：

$$\begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_1} \\ \vdots & & \vdots \\ \frac{\partial f_1}{\partial x_n} & \dots & \frac{\partial f_m}{\partial x_n} \end{pmatrix} \triangleq \left(\frac{\partial f_j}{\partial x_i} \right)_{n \times m}$$

称为 m 维行向量函数 $f(x)^T$ 对 n 维列向量 x 的导数，

记为： $\frac{df(x)^T}{dx}$

微积分：向量和矩阵求导

定义4续. $m \times n$ 阶矩阵函数:

$$\begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix} \triangleq \left(\frac{\partial f_i}{\partial x_j} \right)_{m \times n}$$

称为 m 维列向量函数 $f(x)$ 对 n 维行向量 x^T 的导数。

记为: $\frac{df(x)}{dx^T}$

显然: $\frac{df(x)^T}{dx} = \left(\frac{df(x)}{dx^T} \right)^T$

微积分：向量和矩阵求导

运算性质：

(1) 加法运算公式

$$\frac{d(f^T \pm g^T)}{dx} = \frac{df^T}{dx} \pm \frac{dg^T}{dx}$$

(2) 数乘运算公式

$$\frac{d(\lambda(x)f^T)}{dx} = \frac{d\lambda(x)}{dx}f^T + \lambda(x)\frac{df^T}{dx}$$

(3) 乘法运算公式

$$\frac{d(f^T g)}{dx} = \frac{df^T}{dx}g + \frac{dg^T}{dx}f$$

注意：我们有 $\frac{dx^T}{dx} = I$ 和 $\frac{dx}{dx^T} = I$

微积分：向量和矩阵求导

例3： (1) 求行向量 $x^T A$ 对 x 的导数

(2) 求列向量 Bx 对 x^T 的导数，其中 A 、 B 是常数矩阵，但不一定是方阵。

例 4：求二次型 $x^T A x$ 对 x 的导数

例 5：求数量函数 $p^T A x$ 对 x 的导数，其中 p^T 是 $1 \times n$ 行向量， A 是 $n \times n$ 矩阵，都是常量。 x 是 $n \times 1$ 列向量

微积分：向量和矩阵求导

复合函数微分法

1. 数量函数的公式

公式1. 设 $\phi = \phi(x)$, $x = x(\xi)$, 则

$$\frac{d\phi}{d\xi} = \frac{d\phi}{dx^T} \frac{dx}{d\xi} = \frac{dx^T}{d\xi} \frac{d\phi}{dx}$$

公式2. 设 $\phi = \phi(y)$, $y = y(x)$, 则

$$\frac{d\phi}{dx} = \frac{dy^T}{dx} \frac{d\phi}{dy}$$

$$\frac{d\phi}{dx^T} = \frac{d\phi}{dy^T} \frac{dy}{dx^T}$$

微积分：向量和矩阵求导

公式3. 设 $\phi = \phi(x, y)$, $y = y(x)$, 则

$$\frac{d\phi}{dx} = \frac{dy^T}{dx} \frac{\partial \phi}{\partial y} + \frac{\partial \phi}{\partial x}$$

$$\frac{d\phi}{dx^T} = \frac{\partial \phi}{\partial y^T} \frac{dy}{dx^T} + \frac{\partial \phi}{\partial x^T}$$

微积分：向量和矩阵求导

2. 向量函数的公式

公式4. 设 $z = z(y)$, $y = y(\xi)$, 则

$$\frac{dz}{d\xi} = \frac{dz}{dy^T} \frac{dy}{d\xi}$$

公式5. 设 $z = z(y)$, $y = y(x)$, 则

$$\frac{dz^T}{dx} = \frac{dy^T}{dx} \frac{dz^T}{dy}$$

$$\frac{dz}{dx^T} = \frac{dz}{dy^T} \frac{dy}{dx^T}$$

微积分：向量和矩阵求导

公式6. 设 $z = z(x, y)$, $y = y(x)$, 则

$$\frac{dz^T}{dx} = \frac{dy^T}{dx} \frac{\partial z^T}{\partial y} + \frac{\partial z^T}{\partial x}$$

$$\frac{dz}{dx^T} = \frac{\partial z}{\partial y^T} \frac{dy}{dx^T} + \frac{\partial z}{\partial x^T}$$

微积分：向量和矩阵求导

五、标量 ϕ 对矩阵 X 的导数 $\frac{d\phi}{dX}$

设函数 $\phi = \phi(X)$ 是以 $p \times m$ 矩阵 X 的 $p \times m$ 个元 X_{ij} 为自变量的数量函数，简称以矩阵 X 为自变量的数量函数。

$$\begin{aligned} \text{函数 } f &= X_{11}^3 + (1 + X_{12})X_{11}^2 + (X_{21} + X_{22} + X_{23})X_{11} + X_{21} + X_{22} \\ &= (X_{11} \quad 1) \begin{pmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \end{pmatrix} \begin{pmatrix} X_{11} \\ 1 \end{pmatrix} = f(X) \end{aligned}$$

就是以

$$X = \begin{pmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \end{pmatrix}$$

为自变量的函数。

微积分：向量和矩阵求导

定义5. $p \times m$ 矩阵函数:

$$\begin{pmatrix} \frac{\partial \phi}{\partial X_{11}} & \cdots & \frac{\partial \phi}{\partial X_{1m}} \\ \vdots & & \vdots \\ \frac{\partial \phi}{\partial X_{p1}} & \cdots & \frac{\partial \phi}{\partial X_{pm}} \end{pmatrix} = \left(\frac{\partial \phi}{\partial X_{ij}} \right)_{p \times m}$$

称为数量矩阵 ϕ 对矩阵 X 的导数。

记为: $\frac{d\phi}{dX}$

例 6 求 $\phi = x^T A x$ 对矩阵 A 的导数, 其中向量 x 是定常的。

微积分：向量和矩阵求导

运算性质：

(1) 加法运算公式

$$\frac{d(\phi \pm \varphi)}{dX} = \frac{d\phi}{dX} \pm \frac{d\varphi}{dX}$$

(2) 乘法运算公式

$$\frac{d(\phi^T \varphi)}{dX} = \frac{d\varphi}{dX} \phi + \varphi \frac{d\phi}{dX}$$

微积分：向量和矩阵求导

五、标量 ϕ 对矩阵 X 的导数 $\frac{d\phi}{dX}$

设函数 $\phi = \phi(X)$ 是以 $p \times m$ 矩阵 X 的 $p \times m$ 个元 X_{ij} 为自变量的数量函数，简称以矩阵 X 为自变量的数量函数。

$$\begin{aligned} \text{函数 } f &= a^2 X_{11} + (X_{21} + X_{12})a + X_{22} \\ &= (a \quad 1) \begin{pmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \end{pmatrix} \begin{pmatrix} a \\ 1 \end{pmatrix} = f(X) \end{aligned}$$

就是以

$$X = \begin{pmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \end{pmatrix}$$

为自变量的函数。

微积分：向量和矩阵求导

定义5. $p \times m$ 矩阵函数:

$$\begin{pmatrix} \frac{\partial \phi}{\partial X_{11}} & \cdots & \frac{\partial \phi}{\partial X_{1m}} \\ \vdots & & \vdots \\ \frac{\partial \phi}{\partial X_{p1}} & \cdots & \frac{\partial \phi}{\partial X_{pm}} \end{pmatrix} = \left(\frac{\partial \phi}{\partial X_{ij}} \right)_{p \times m}$$

称为数量矩阵 ϕ 对矩阵 X 的导数。

记为: $\frac{d\phi}{dX}$

例 6 求 $\phi = x^T A x$ 对矩阵 A 的导数, 其中向量 x 是定常的。

微积分：向量和矩阵求导

运算性质：

(1) 加法运算公式

$$\frac{d(\phi \pm \varphi)}{dX} = \frac{d\phi}{dX} \pm \frac{d\varphi}{dX}$$

(2) 乘法运算公式

$$\frac{d(\phi\varphi)}{dX} = \frac{d\varphi}{dX}\phi + \varphi \frac{d\phi}{dX}$$

微积分：向量和矩阵求导

矩阵微分法

标量函数对矩阵的导数，即 $\frac{d\phi}{dX}$

微分法：

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

微分性质

(1) 加法运算公式：

$$d(X \pm Y) = dX \pm dY$$

$$d(XY) = d(X)Y + Xd(Y)$$

$$d(X^T) = (dX)^T$$

$$d \text{tr}(X) = \text{tr}(dX)$$

(2) 逆： $dX^{-1} = -X^{-1}dX X^{-1}$

微积分：向量和矩阵求导

- (3) 行列式: $d|X| = \text{tr}(X^*dX)$
- (4) 逐元素乘法: $d(X \odot Y) = dX \odot Y + X \odot dY$
- (5) 逐元素函数: $d\sigma(X) = \sigma'(X) \odot dX$

关于迹的运算

- (1) 标量的迹: $\xi = \text{tr}(\xi)$
- (2) 转置: $\text{tr}(X) = \text{tr}(X^T)$
- (3) 线性: $\text{tr}(X \pm Y) = \text{tr}(X) \pm \text{tr}(Y)$
- (4) 矩阵乘法: $\text{tr}(XY) = \text{tr}(YX)$
- (5) 矩阵乘法/逐元素乘法: $\text{tr}(X^T(Y \odot Z)) = \text{tr}((X \odot Y)^T Z)$
其中 X, Y, Z 具有相同的尺寸。

微积分：向量和矩阵求导

例 1. $\phi = a^T X b$, 计算 $\frac{d\phi}{dX}$. 其中 a 是 m 维列向量, X 是 $m \times n$ 矩阵, b 是 n 维列向量.

例 2. $\phi = a^T \exp(Xb)$, 计算 $\frac{d\phi}{dX}$. 其中 a 是 m 维列向量, X 是 $m \times n$ 矩阵, b 是 n 维列向量.

例 3. $\phi = \text{tr}(Y^T M Y)$, $Y = \sigma(WX)$, 计算 $\frac{d\phi}{dX}$. 其中 W 是 $l \times m$ 矩阵, X 是 $m \times n$ 矩阵, Y 是 $l \times n$ 矩阵, M 是 $l \times l$ 对称矩阵, σ 是逐元素函数

例 4. [线性回归]: $l = \|Xw - y\|^2$, 求 w 的最小二乘估计, 即求 $\frac{dl}{dw}$ 的零点。其中 y 是 m 维列向量, X 是 $m \times n$ 矩阵, w 是 n 维列向量。

信息论

信息是对不确定性的消除。

自信息: $I(a_i) = -\log(a_i)$

性质: 设 a_1, a_2 为两个随机事件

(1) 若 $P(a_1) \geq P(a_2)$, 则 $I(a_1) \leq I(a_2)$

(2) 若 $P(a_1) = 1$, 则 $I(a_1) = 0$

(3) 若 $P(a_1) = 0$, 则 $I(a_1) = \infty$

(4) 若 a_1, a_2 为独立事件, 则 $I(a_1, a_2) = I(a_1) + I(a_2)$

信息论

信息熵是指随机系统的总体信息量，用所有随机事件自信息的统计平均来表示。

$$H(X) = f(p_1, p_2, \dots, p_N) = - \sum_{x \in X} p(x) \log p(x)$$

性质：

- 1、对称性： $f(p_1, p_2, \dots, p_N) = f(p_{k(1)}, p_{k(2)}, \dots, p_{k(N)})$
- 2、非负性： $H(X) \geq 0$
- 3、可加性： $H(X, Y) = H(X) + H(Y|X)$
- 4、条件减少熵： $H(X|Y) \leq H(X)$
- 5、最大离散熵定理： $f(p_1, p_2, \dots, p_N) \leq f(\frac{1}{N}, \frac{1}{N}, \dots, \frac{1}{N}) = \log N = \log |X|$

信息论

联合熵：一对离散随机变量 (X, Y) 的联合熵定义为：

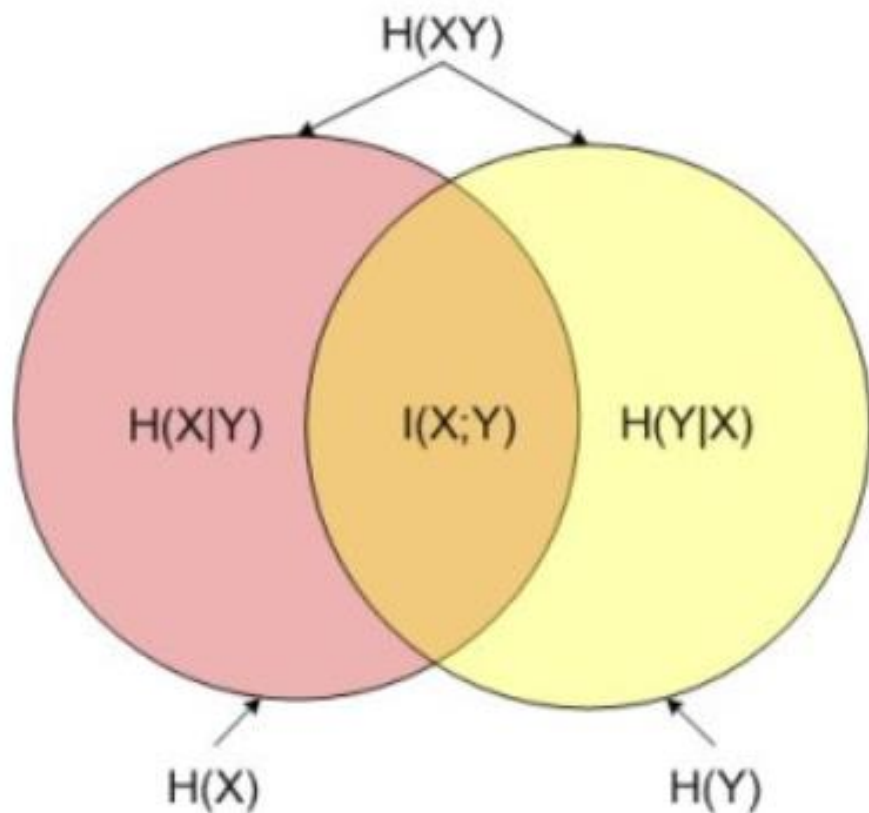
$$H(X, Y) = - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x, y)$$

条件熵： $H(Y|X) = - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(y|x)$

$$H(Y|X) = \sum_{x \in \mathcal{X}} p(x) H(Y|X = x)$$

信息论

互信息: $I(X;Y) = H(X) - H(X|Y)$ 或 $I(X;Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$



$$I(X;Y) = H(X) - H(X|Y)$$

$$I(X;Y) = H(Y) - H(Y|X)$$

$$I(X;Y) = H(X) + H(Y) - H(XY)$$

信息论

相对熵（KL散度）：

设 $p(x)$ 、 $q(x)$ 是关于随机变量 X 的两个概率分布，则 p 相对于 q 的相对熵是：

$$D_{KL}(p\|q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

性质：1、如果 $p(x)$ 和 $q(x)$ 两个分布相同，那么相对熵等于0

2、 $D_{KL}(p\|q) \neq D_{KL}(q\|p)$

3、 $D_{KL}(p\|q) \geq 0$

信息论

交叉熵 (Cross entropy) :

设 $p(x)$ 、 $q(x)$ 是关于随机变量 X 的两个概率分布, 使用分布 $q(x)$ 表示目标分布 $p(x)$ 的困难程度:

$$H(p, q) = - \sum_i p(x_i) \log q(x_i)$$

性质: 1、 $D_{KL}(p, q) = H(p, q) - H(p)$

信息论

在机器学习中，

- (1) 希望学到的模型的分布和真实分布一致: $P(model) \simeq P(real)$
- (2) 但是真实分布不可知，假设训练数据时从真实数据中独立同分布采样的: $P(train) \simeq P(real)$
- (3) 因此，我们希望学到的模型分布至少和训练数据的分布一致: $P(train) \simeq P(model)$

交叉熵可以用来计算学习模型分布与训练分布之间的差异

$$L(Y, P(Y|\mathbf{X})) = -\log P(Y|\mathbf{X}) = -\frac{1}{N} \sum_{i=1}^N (1 - y_i) \log (1 - p(\mathbf{x}_i)) - y_i \log p(\mathbf{x}_i)$$



THE END