

# 机器学习常用数学公式汇总（一）

李新春 lxcnju@163.com

2017 年 5 月 7 日

本文主要根据周志华老师的《机器学习》、李航老师的《统计学习方法》、《模式识别》、《矩阵论》、《数理统计》、《优化方法》以及网上博客对机器学习中常用的数学公式进行总结。主要分为矩阵、优化、概率统计三个部分。本次以矩阵为主进行总结其中的公式，然后列出优化里面常用的牛顿法和拟牛顿法公式。概率统计主要涉及常见分布族及其数值特征和各种检验方法等，概率统计将单独总结为一篇文章。

## 一、矩阵

### 1、矩阵秩

$$R(\mathbf{A}) = R(\mathbf{A}^T) = R(\mathbf{A}^T \mathbf{A}) = R(\mathbf{A} \mathbf{A}^T) \quad (0.1)$$

$$\max\{R(\mathbf{A}), R(\mathbf{B})\} \leq R(\mathbf{A}, \mathbf{B}) \leq R(\mathbf{A}) + R(\mathbf{B}) \quad (0.2)$$

$$R(\mathbf{A} \pm \mathbf{B}) \leq R(\mathbf{A}) + R(\mathbf{B}) \quad (0.3)$$

$$R(\mathbf{A} \mathbf{B}) \leq \min\{R(\mathbf{A}), R(\mathbf{B})\} \quad (0.4)$$

### 2、矩阵迹（ $\mathbf{a} \mathbf{b}^T$ 为列向量）

$$\text{tr}(\mathbf{A} \mathbf{B}) = \text{tr}(\mathbf{B} \mathbf{A}) \quad (0.5)$$

$$\text{tr}(\mathbf{A} \mathbf{B} \mathbf{C}) = \text{tr}(\mathbf{B} \mathbf{C} \mathbf{A}) = \text{tr}(\mathbf{C} \mathbf{A} \mathbf{B}) \quad (0.6)$$

$$\frac{\partial \text{tr}(\mathbf{A} \mathbf{B})}{\partial \mathbf{A}} = \mathbf{B}^T, \quad \frac{\partial \text{tr}(\mathbf{A} \mathbf{B})}{\partial \mathbf{B}} = \mathbf{A}^T \quad (0.7)$$

$$\frac{\partial \text{tr}(\mathbf{A}^T \mathbf{B})}{\partial \mathbf{A}} = \mathbf{B}, \quad \frac{\partial \text{tr}(\mathbf{A}^T \mathbf{B})}{\partial \mathbf{B}} = \mathbf{A} \quad (0.8)$$

$$\frac{\partial \text{tr}(\mathbf{A})}{\partial \mathbf{A}} = \mathbf{I} \quad (0.9)$$

$$\text{tr}(\mathbf{a} \mathbf{b}^T) = \mathbf{a}^T \mathbf{b} \quad (0.10)$$

$$\text{tr}(\mathbf{x} \mathbf{x}^T \mathbf{A}) = \mathbf{x}^T \mathbf{A} \mathbf{x} \quad (0.11)$$

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{A}} = \mathbf{x} \mathbf{x}^T \quad (0.12)$$

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A}^T + \mathbf{A}) \mathbf{x} \quad (0.13)$$

$$\frac{\partial \text{tr}(\mathbf{A} \mathbf{B} \mathbf{A}^T)}{\partial \mathbf{A}} = \mathbf{A}(\mathbf{B} + \mathbf{B}^T) \quad (0.14)$$

$$\frac{\partial \text{tr}(\mathbf{A} \mathbf{B} \mathbf{A}^T \mathbf{C})}{\partial \mathbf{A}} = \mathbf{C} \mathbf{A} \mathbf{B} + \mathbf{C}^T \mathbf{A} \mathbf{B}^T \quad (0.15)$$

证明: 这里仅给出最后一个式子的证明:  $\mathbf{A} \in \mathbf{R}^{m \times n}, \mathbf{B} \in \mathbf{R}^{n \times n}, \mathbf{C} \in \mathbf{R}^{m \times m}$ , 记 $\mathbf{A}$ 按列分块结果为 $(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n)$ , 记 $\mathbf{B}_{ij}$ 为 $\mathbf{B}$ 的第 $i$ 行第 $j$ 列元素, 其余符号含义类推。

$$\begin{aligned}
tr(\mathbf{A}^T \mathbf{C} \mathbf{A} \mathbf{B}) &= tr \left( \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_n^T \end{bmatrix} \mathbf{C} \begin{bmatrix} \mathbf{a}_1 & \mathbf{a}_2 & \dots & \mathbf{a}_n \end{bmatrix} \mathbf{B} \right) \\
&= tr \left( \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_n^T \end{bmatrix} \begin{bmatrix} \mathbf{C} \mathbf{a}_1 & \mathbf{C} \mathbf{a}_2 & \dots & \mathbf{C} \mathbf{a}_n \end{bmatrix} \mathbf{B} \right) \\
&= tr \left( \begin{bmatrix} \mathbf{a}_1^T \mathbf{C} \mathbf{a}_1 & \mathbf{a}_1^T \mathbf{C} \mathbf{a}_2 & \dots & \mathbf{a}_1^T \mathbf{C} \mathbf{a}_n \\ \mathbf{a}_2^T \mathbf{C} \mathbf{a}_1 & \mathbf{a}_2^T \mathbf{C} \mathbf{a}_2 & \dots & \mathbf{a}_2^T \mathbf{C} \mathbf{a}_n \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{a}_n^T \mathbf{C} \mathbf{a}_1 & \mathbf{a}_n^T \mathbf{C} \mathbf{a}_2 & \dots & \mathbf{a}_n^T \mathbf{C} \mathbf{a}_n \end{bmatrix} \mathbf{B} \right) \\
&= \sum_{i,j} \mathbf{a}_i^T \mathbf{C} \mathbf{a}_j \mathbf{B}_{ij} \\
\frac{\partial \sum_{i,j} \mathbf{a}_i^T \mathbf{C} \mathbf{a}_j \mathbf{B}_{ij}}{\partial \mathbf{a}_i} &= \sum_j \mathbf{C} \mathbf{a}_j \mathbf{B}_{ij} + \sum_j \mathbf{C}^T \mathbf{a}_j \mathbf{B}_{ji} \\
&\triangleq \mathbf{p} \mathbf{a}_i \\
\frac{\partial tr(\mathbf{A} \mathbf{B} \mathbf{A}^T \mathbf{C})}{\partial \mathbf{A}} &= \frac{\partial tr(\mathbf{A}^T \mathbf{C} \mathbf{A} \mathbf{B})}{\partial \mathbf{A}} \\
&= \frac{\partial \sum_{i,j} \mathbf{a}_i^T \mathbf{C} \mathbf{a}_j \mathbf{B}_{ij}}{\partial \mathbf{A}} \\
&= \begin{bmatrix} \mathbf{p} \mathbf{a}_1 & \mathbf{p} \mathbf{a}_2 & \dots & \mathbf{p} \mathbf{a}_n \end{bmatrix} \\
&= \mathbf{C} \mathbf{A} \mathbf{B} + \mathbf{C}^T \mathbf{A} \mathbf{B}^T
\end{aligned}$$

### 3、矩阵导数

$$\left( \frac{\partial \mathbf{a}}{\partial x} \right)_i = \frac{\partial \mathbf{a}_i}{\partial x} \quad (\mathbf{a} \in \mathbf{R}^n, x \in \mathbf{R}) \quad (0.16)$$

$$\left( \frac{\partial \mathbf{A}}{\partial x} \right)_{ij} = \frac{\partial \mathbf{A}_{ij}}{\partial x} \quad (\mathbf{A} \in \mathbf{R}^{m \times n}, x \in \mathbf{R}) \quad (0.17)$$

$$\left( \frac{\partial x}{\partial \mathbf{a}} \right)_i = \frac{\partial x}{\partial \mathbf{a}_i} \quad (\mathbf{a} \in \mathbf{R}^n, x \in \mathbf{R}) \quad (0.18)$$

$$\left( \frac{\partial x}{\partial \mathbf{A}} \right)_{ij} = \frac{\partial x}{\partial \mathbf{A}_{ij}} \quad (\mathbf{A} \in \mathbf{R}^{m \times n}, x \in \mathbf{R}) \quad (0.19)$$

$$\frac{\partial f}{\partial \mathbf{A}} = \begin{bmatrix} \frac{\partial f}{\partial \mathbf{A}_{11}} & \frac{\partial f}{\partial \mathbf{A}_{12}} & \dots & \frac{\partial f}{\partial \mathbf{A}_{1n}} \\ \frac{\partial f}{\partial \mathbf{A}_{21}} & \frac{\partial f}{\partial \mathbf{A}_{22}} & \dots & \frac{\partial f}{\partial \mathbf{A}_{2n}} \\ \vdots & \ddots & \ddots & \vdots \\ \frac{\partial f}{\partial \mathbf{A}_{m1}} & \frac{\partial f}{\partial \mathbf{A}_{m2}} & \dots & \frac{\partial f}{\partial \mathbf{A}_{mn}} \end{bmatrix} \quad (\mathbf{X} \in \mathbf{R}^{m \times n}) \quad (0.20)$$

$$\begin{aligned}
F(\mathbf{X}) &= \begin{bmatrix} f_{11}(\mathbf{X}) & \cdots & f_{1s}(\mathbf{X}) \\ \vdots & \ddots & \vdots \\ f_{r1}(\mathbf{X}) & \cdots & f_{rs}(\mathbf{X}) \end{bmatrix} \quad (\mathbf{X} \in \mathbf{R}^{m \times n}) \\
\frac{\partial F(\mathbf{X})}{\partial \mathbf{X}} &= \begin{bmatrix} \frac{\partial F}{\partial \mathbf{X}_{11}} & \cdots & \frac{\partial F}{\partial \mathbf{X}_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial F}{\partial \mathbf{X}_{m1}} & \cdots & \frac{\partial F}{\partial \mathbf{X}_{mn}} \end{bmatrix} \\
\frac{\partial F}{\partial \mathbf{X}_{ij}} &= \begin{bmatrix} \frac{\partial f_{11}}{\partial \mathbf{X}_{ij}} & \cdots & \frac{\partial f_{1n}}{\partial \mathbf{X}_{ij}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_{m1}}{\partial \mathbf{X}_{ij}} & \cdots & \frac{\partial f_{mn}}{\partial \mathbf{X}_{ij}} \end{bmatrix} \tag{0.21}
\end{aligned}$$

$$\begin{aligned}
\mathbf{f}(\mathbf{x}) &= \begin{pmatrix} f_1(\mathbf{x}) & \cdots & f_m(\mathbf{x}) \end{pmatrix}^T \quad \mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_n) \\
J(\mathbf{x}) &= \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial \mathbf{x}_1} & \cdots & \frac{\partial f_1}{\partial \mathbf{x}_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial \mathbf{x}_1} & \cdots & \frac{\partial f_m}{\partial \mathbf{x}_n} \end{bmatrix} \tag{0.22}
\end{aligned}$$

$$H(\mathbf{x}) = \frac{\partial \nabla f(\mathbf{x})}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial^2 f}{\partial \mathbf{x}_1 \partial \mathbf{x}_1} & \frac{\partial^2 f}{\partial \mathbf{x}_1 \partial \mathbf{x}_2} & \cdots & \frac{\partial^2 f}{\partial \mathbf{x}_1 \partial \mathbf{x}_n} \\ \frac{\partial^2 f}{\partial \mathbf{x}_2 \partial \mathbf{x}_1} & \frac{\partial^2 f}{\partial \mathbf{x}_2 \partial \mathbf{x}_2} & \cdots & \frac{\partial^2 f}{\partial \mathbf{x}_2 \partial \mathbf{x}_n} \\ \vdots & \ddots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial \mathbf{x}_n \partial \mathbf{x}_1} & \frac{\partial^2 f}{\partial \mathbf{x}_n \partial \mathbf{x}_2} & \cdots & \frac{\partial^2 f}{\partial \mathbf{x}_n \partial \mathbf{x}_n} \end{bmatrix} \tag{0.23}$$

$$\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a} \tag{0.24}$$

$$\frac{\partial \mathbf{A} \mathbf{B}}{\partial \mathbf{x}} = \frac{\partial \mathbf{A}}{\partial \mathbf{x}} \mathbf{B} + \mathbf{A} \frac{\partial \mathbf{B}}{\partial \mathbf{x}} \tag{0.25}$$

注：(0.16) 是函数向量的导数，(0.17) 是函数矩阵的导数，(0.18) 是变量对向量的导数，(0.19) 是变量对矩阵的导数，(0.20) 是函数对矩阵的导数，(0.21) 是向量值函数对矩阵的导数，(0.22) 是Jacobi矩阵（可由0.21推导出来），(0.23) 是Hessian矩阵。

4、矩阵逆、广义逆（用 $\mathbf{A}^+$ 表示）

$$\frac{\partial \mathbf{A}^{-1}}{\partial \mathbf{x}} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial \mathbf{x}} \mathbf{A}^{-1} \tag{0.26}$$

$$\text{if } \mathbf{A} \mathbf{G} \mathbf{A} = \mathbf{A}, \mathbf{G} \mathbf{A} \mathbf{G} = \mathbf{G}, (\mathbf{A} \mathbf{G})^T = \mathbf{A} \mathbf{G}, (\mathbf{G} \mathbf{A})^T = \mathbf{G} \mathbf{A}, \text{ then } \mathbf{A}^+ = \mathbf{G} \tag{0.27}$$

$$(\mathbf{A}^+)^T = (\mathbf{A}^T)^+, (\mathbf{A} \mathbf{A}^T)^+ = (\mathbf{A}^T)^+ \mathbf{A}^+ \tag{0.28}$$

$$\text{if } \mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T, \text{ then } \mathbf{A}^+ = \mathbf{V} \mathbf{\Sigma}^{-1} \mathbf{U}^T \tag{0.29}$$

5、矩阵Frobenius范数

$$\|\mathbf{A}\|_F = (\text{tr}(\mathbf{A}^T \mathbf{A}))^{1/2} \tag{0.30}$$

$$\frac{\partial \|\mathbf{A}\|_F^2}{\partial \mathbf{A}} = \frac{\partial \mathbf{A}^T \mathbf{A}}{\partial \mathbf{A}} = 2\mathbf{A} \tag{0.31}$$

注：关于矩阵广义逆和矩阵范数还有很多知识点，以后会有更详细的总结，这里先省去。（矩阵广义逆有很多种，其和最小二乘法以及矩阵投影之间有着密切的联系；矩阵范数更是重要的一点，矩阵的收敛性、条件数、特征值等与其密切相关。）

## 二、优化

1、牛顿法（0.32是由泰勒展开式得到；假设 $\mathbf{G}_k$ 正定，利用一阶最优性条件得到0.33）

$$f(\mathbf{x}_k + \mathbf{s}) = f(\mathbf{x}_k) + \mathbf{s}^T \mathbf{g}_k + \frac{1}{2} \mathbf{s}^T \mathbf{G}_k \mathbf{s} \quad (0.32)$$

$$\mathbf{G}_k \mathbf{s} = -\mathbf{g}_k \Rightarrow \mathbf{s}_k = -\mathbf{G}_k^{-1} \mathbf{g}_k \quad (0.33)$$

2、拟牛顿法（公式0.34是通过梯度函数 $\mathbf{g}(\mathbf{x})$ 在 $\mathbf{x}_{k+1}$ 点的近似得到；式子0.36是通过利用 $\mathbf{B}_{k+1}$ 和 $\mathbf{H}_{k+1}$ 来分别近似 $\mathbf{G}_{k+1}$ 和 $\mathbf{G}_{k+1}^{-1}$ 得到，称之为拟牛顿方式；下面公式是通过各种校正方法进行求解。）

$$g(\mathbf{x}_k) \approx \mathbf{g}_{k+1} + \mathbf{G}_{k+1}(\mathbf{x}_k - \mathbf{x}_{k+1}) \quad (0.34)$$

$$\text{let } \mathbf{y}_k = \mathbf{g}_{k+1} - \mathbf{g}_k, \quad \mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k \quad \text{then } \mathbf{y}_k \approx \mathbf{G}_{k+1} \mathbf{s}_k \quad (0.35)$$

$$\mathbf{y}_k = \mathbf{B}_{k+1} \mathbf{s}_k, \quad \mathbf{s}_k = \mathbf{H}_{k+1} \mathbf{y}_k \quad (0.36)$$

SR1 :

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \mathbf{v}_k \mathbf{v}_k^T \quad (0.37)$$

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{(\mathbf{s}_k - \mathbf{H}_k \mathbf{y}_k)(\mathbf{s}_k - \mathbf{H}_k \mathbf{y}_k)^T}{(\mathbf{s}_k - \mathbf{H}_k \mathbf{y}_k)^T \mathbf{y}_k} \quad (0.38)$$

$$\mathbf{B}_{k+1} = \mathbf{B}_k + \frac{(\mathbf{y}_k - \mathbf{B}_k \mathbf{s}_k)(\mathbf{y}_k - \mathbf{B}_k \mathbf{s}_k)^T}{(\mathbf{y}_k - \mathbf{B}_k \mathbf{s}_k)^T \mathbf{s}_k} \quad (0.39)$$

DFP :

$$\mathbf{H}_{k+1} = \mathbf{H}_k + a \mathbf{v}_k \mathbf{v}_k^T + b \mathbf{u}_k \mathbf{u}_k^T \quad (0.40)$$

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{\mathbf{s}_k \mathbf{s}_k^T}{\mathbf{s}_k^T \mathbf{y}_k} - \frac{\mathbf{H}_k \mathbf{y}_k \mathbf{y}_k^T \mathbf{H}_k}{\mathbf{y}_k^T \mathbf{H}_k \mathbf{y}_k} \quad (0.41)$$

$$\mathbf{B}_{k+1} = \left( \mathbf{I} - \frac{\mathbf{y}_k \mathbf{s}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} \right) \mathbf{B}_k \left( \mathbf{I} - \frac{\mathbf{s}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} \right) + \frac{\mathbf{y}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} \quad (0.42)$$

BFGS :

$$\mathbf{B}_{k+1} = \mathbf{B}_k + a \mathbf{v}_k \mathbf{v}_k^T + b \mathbf{u}_k \mathbf{u}_k^T \quad (0.43)$$

$$\mathbf{B}_{k+1} = \mathbf{B}_k + \frac{\mathbf{y}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} - \frac{\mathbf{B}_k \mathbf{s}_k \mathbf{s}_k^T \mathbf{B}_k}{\mathbf{s}_k^T \mathbf{B}_k \mathbf{s}_k} \quad (0.44)$$

$$\mathbf{H}_{k+1} = \left( \mathbf{I} - \frac{\mathbf{s}_k \mathbf{y}_k^T}{\mathbf{s}_k^T \mathbf{y}_k} \right) \mathbf{H}_k \left( \mathbf{I} - \frac{\mathbf{y}_k \mathbf{s}_k^T}{\mathbf{s}_k^T \mathbf{y}_k} \right) + \frac{\mathbf{s}_k \mathbf{s}_k^T}{\mathbf{s}_k^T \mathbf{y}_k} \quad (0.45)$$

Broyden :

$$\mathbf{H}_{k+1}^\alpha = (1 - \alpha) \mathbf{H}_k^{DFP} + \alpha \mathbf{H}_k^{BFGS} \quad (0.46)$$

Morrison :

$$(\mathbf{A} + \mathbf{u} \mathbf{v}^T)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1} \mathbf{u} \mathbf{v}^T \mathbf{A}^{-1}}{1 + \mathbf{v}^T \mathbf{A}^{-1} \mathbf{u}} \quad (0.47)$$

注：关于最优化理论，还有很多优化算法。共轭梯度法、约束优化、二次规划、最小二乘、罚函数等可能会在以后专题详细介绍。因为其中涉及东西太多，这里不全面给出。同时，拟牛顿法只是给出了主要公式，其全局收敛性证明、超线性收敛性证明及条件等异常繁杂，请读者参考最优化理论的书籍。

参考资料（仅列出博客）：

1. 矩阵的迹 求导等公式 <http://blog.sina.com.cn/s/blog74b69f7101016tg5.html>