



deepshare.net

深度之眼

# 《Deep Learning》 ——数学基础

导师: Johnson

---

# 数学基础

Machine Learning Basics

---



deepshare.net

深度之眼

1. 矩阵对角化, SVD分解以及应用
2. 逆矩阵, 伪逆矩阵
3. PCA原理与推导
4. 极大似然估计, 误差的高斯分布与最小二乘估计的等价性
5. 最优化, 无约束, 有约束, 拉格朗日乘子的意义, KKT条件

# 1. 矩阵对角化, SVD分解以及应用



实用性质:

$$A(B+C) = AB + AC \quad A(B+C) = AB + AC \quad (\text{分配率})$$

$$A(BC) = (AB)C \quad A(BC) = (AB)C \quad (\text{结合律})$$

$$AB \neq BA \quad AB \neq BA \quad (\text{一般不满足交换律})$$

$$(AB)^T = B^T A^T \quad (AB)^T = B^T A^T \quad (\text{转置})$$

$$x^T y = (x^T y)^T = y^T x \quad (\text{转置})$$

## 单位矩阵：

任意向量或矩阵和单位矩阵相乘，都不会改变，记为  $I$ 。所有沿主对角线的元素都是1，而所有其他位置的元素都是 0。

$$I_1 = [1], I_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, I_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \dots, I_n = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}$$

**矩阵逆：** 矩阵（方阵）的逆满足如下条件

$$A^{-1}A = AA^{-1} = I_n$$

矩阵B（方阵）的对角化  $P^{-1}AP = B$

其中A为对角矩阵，P为单位正交矩阵

一般的矩阵不一定能对角化，但是对称矩阵一定可以对角化(特别是对称正定矩阵，得到的 $\lambda_i$ 都是正数)

曾经一道面试题矩阵的压缩表示

最小 $n+1$

$$P^T P = P P^T = I$$

$$P^T = P^{-1}$$

所以:  $B = P^T A P$

设  $P^T = (u_1, u_2, \dots, u_n), u_i \in \mathbb{R}^n$

$$A = \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \lambda_3 \end{pmatrix}$$

$$\text{则 } B = (u_1, u_2, \dots, u_n) \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \lambda_3 \end{pmatrix} \begin{pmatrix} u_1^T \\ \vdots \\ u_n^T \end{pmatrix}$$

$$= \lambda_1 u_1 u_1^T + \lambda_2 u_2 u_2^T + \lambda_3 u_3 u_3^T + \dots + \lambda_n u_n u_n^T$$


$$\begin{pmatrix} & & \\ & & \\ & & \end{pmatrix}_{n \times 1} \begin{pmatrix} & & \\ & & \\ & & \end{pmatrix}_{1 \times n}$$

一般矩阵的svd分解:  $A^T A$  为对称正定矩阵  
 $(A^T A)^T = A^T (A^T)^T = A^T A \Rightarrow$  对称性  
 $x^T A^T A x = (x^T A^T)(Ax) = (Ax)^T (Ax) \geq 0 \Rightarrow$  正定性

$$A_{m \times n}$$

$$(A^T A)_{n \times n} = U^T D U$$

$$(A A^T)_{m \times m} = V^T D V$$

$$\Rightarrow A_{m \times n} = V_{m \times m}^T \begin{pmatrix} \lambda_1^{\frac{1}{2}} & & \\ & \lambda_2^{\frac{1}{2}} & \\ & & \ddots \end{pmatrix} U_{n \times n}$$

其中  $\lambda_i$  为  $D$  的对角元素

$$\text{令 } V_{m \times m}^T = (v_1, v_2, \dots, v_m)$$

$$U_{n \times n}^T = (u_1, u_2, \dots, u_n)$$

$$\Rightarrow A_{m \times n} = \lambda_1^{\frac{1}{2}} v_1 u_1^T + \lambda_2^{\frac{1}{2}} v_2 u_2^T + \dots$$

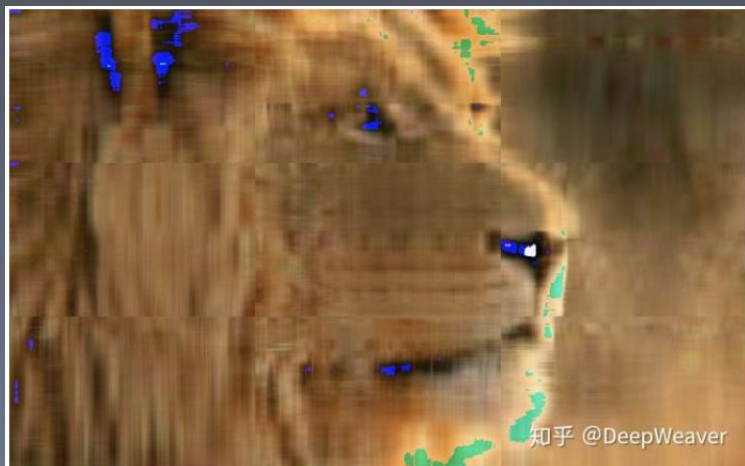
图像的压缩存储最小需要  $m+n+1$

当压缩存储量为  $(m+n+1)*K$  时, 误差为

$$error = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^{\min(m,n)} \lambda_i}$$



$$A_{m \times n} = \lambda_1^{\frac{1}{2}} v_1 u_1^T + \lambda_2^{\frac{1}{2}} v_2 u_2^T + \dots$$



保留了前10个(压缩率122)



前30个(压缩率31)



前50个(压缩率17)

传统网络图片传输与现代传输的原理



$$A_{m \times n} = \lambda_1^{\frac{1}{2}} v_1 u_1^T + \lambda_2^{\frac{1}{2}} v_2 u_2^T + \dots$$

$$A_{200 \times 100} = \lambda_1^{\frac{1}{2}} \begin{pmatrix} v_1 \\ \vdots \\ v_{10} \end{pmatrix}_{200 \times 1} \begin{pmatrix} u_1 & \dots & u_{10} \end{pmatrix}_{1 \times 100} + \lambda_2^{\frac{1}{2}} \begin{pmatrix} v_2 \\ \vdots \\ v_{10} \end{pmatrix}_{200 \times 1} \begin{pmatrix} u_2 & \dots & u_{10} \end{pmatrix}_{1 \times 100} + \dots$$

$$= \begin{pmatrix} v_1 & v_2 & \dots & v_{10} \end{pmatrix}_{200 \times 10} \begin{pmatrix} \lambda_1^{\frac{1}{2}} & & & \\ & \lambda_2^{\frac{1}{2}} & & \\ & & \ddots & \\ & & & \lambda_{10}^{\frac{1}{2}} \end{pmatrix}_{10 \times 10} \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_{10} \end{pmatrix}_{10 \times 100}$$

$$\approx M_{200 \times 10} N_{10 \times 100}$$

$$A_{200 \times 10} X_{100 \times 1} \text{ 算 } 200 \times 100 \text{ 次乘法}$$

$$M_{200 \times 10} N_{10 \times 100} X_{100 \times 1} \text{ 算 } 10 \times 100 + 10 \times 200 = 3000 \text{ 次}$$

能够极大的减小算法复杂度，  
在深度神经网络中有着广泛的应用



## 2. 逆矩阵, 伪逆矩阵, 最小二乘解, 最小范数解

$$x_1, x_2, \dots, x_N, x_i \in \mathbb{R}^n$$

$$y_1, y_2, \dots, y_N, y_i \in \mathbb{R}^1$$

$$y_1 = x_{11}a_1 + x_{12}a_2 + \dots + x_{1n}a_n$$

$$y_2 = x_{21}a_1 + x_{22}a_2 + \dots + x_{2n}a_n$$

$\vdots$

$$y_N = x_{N1}a_1 + x_{N2}a_2 + \dots + x_{Nn}a_n$$

$$\begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{Nn} \end{pmatrix} \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

$$X_{N \times n} a_{n \times 1} = Y_{N \times 1}$$

当  $N = n$  且  $X_{N \times n}$  可逆时:

$$a = X^{-1}Y_S$$

一般情况:  $N \neq n$

$$\min \|xa - Y\|^2 = J \quad \frac{\partial J}{\partial a} = x^T(xa - Y) = 0$$

$$x^T xa = x^T Y \quad x^T x \text{ 是否可逆?}$$

1.  $N > n$

如  $N = 5, n = 3$   $(x^T x)_{3 \times 3}$  一般是可逆的

$$a = (x^T x)^{-1} x^T Y$$

2.  $N < n$

如  $N = 3, n = 5$   $(x^T x)_{5 \times 5}$

$$R(x^T x) \leq R(x) \leq 3$$

故  $x^T x$  不可逆

此刻  $J = \|xa - Y\|^2 + \lambda \|a\|^2$

$$\frac{\partial J}{\partial a} = x^T x a - x^T Y + \lambda a = 0$$

$$(x^T x + \lambda I) a = x^T Y$$

$$a = (x^T x + \lambda I)^{-1} x^T Y$$

$x^T x$  为对称矩阵可对角化

$$x^T x = p^{-1} \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} p$$

$$|x^T x| = \lambda_1 \lambda_2 \cdots \lambda_n$$

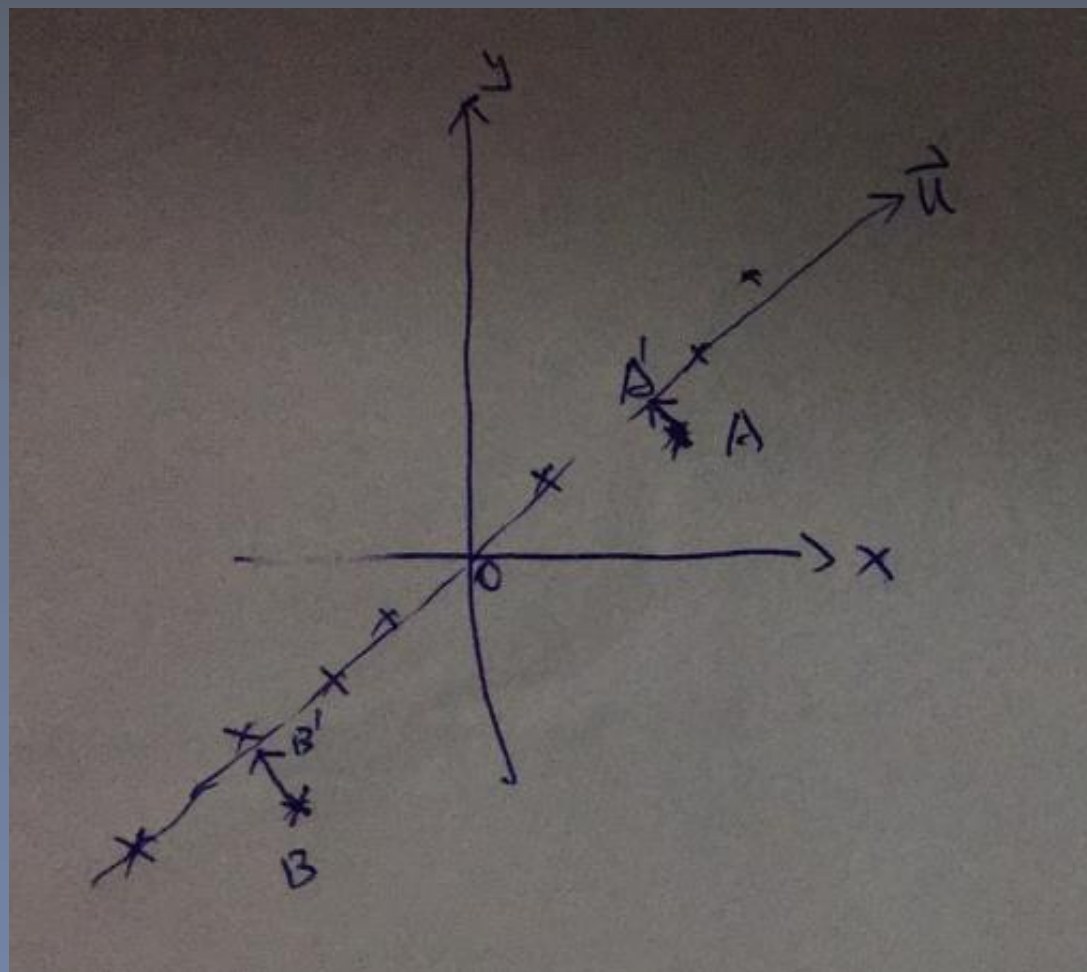
1.  $a^T (x^T x) a = (xa)^T (xa) \geq 0 \rightarrow \lambda_i \geq 0$

$|x^T x|$  仍可能为0, 不一定可逆

2.  $a^T (x^T x + \lambda I) a = (xa)^T (xa) + \lambda a^T a > 0 \rightarrow \lambda_i > 0$

$|x^T x + \lambda I| > 0$  恒成立, 一定可逆

# 3. PCA原理与推导



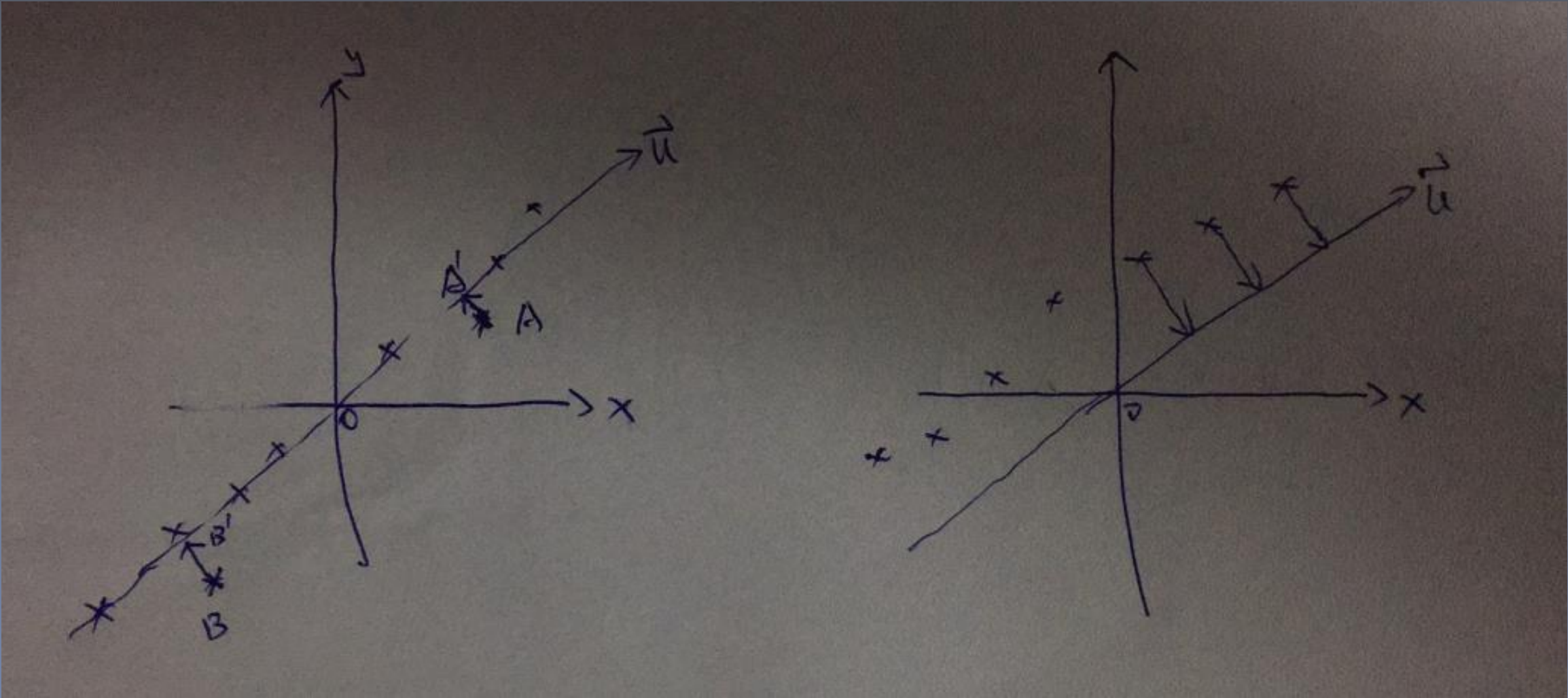
PCA仍然是一种数据压缩的算法

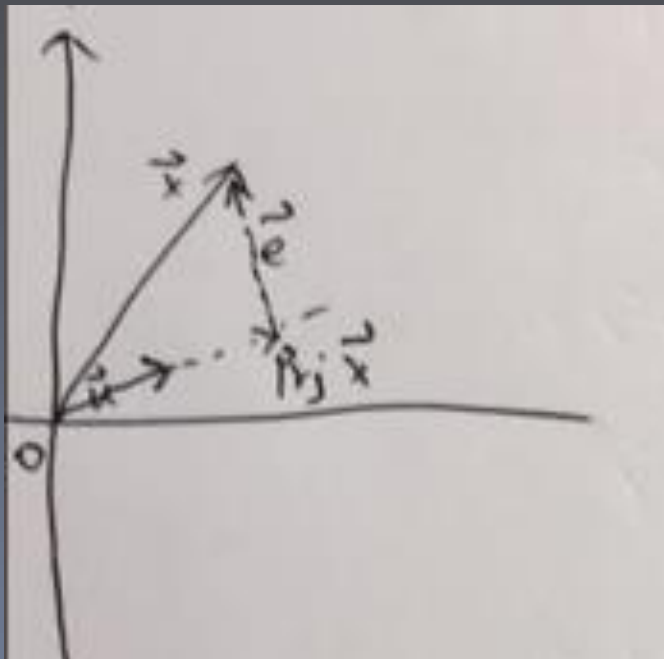
A点需要 $x, y$ 两个坐标来表示, 假设A在向量 $u$ 上面的投影点为 $A'$ , 则 $A'$  仅仅需要一个参数就能表示, 就是 $OA'$  的长度(即 $A'$  在 $u$ 上的坐标), 我们就想着用 $A'$  来替换A, 这样 $N$ 个点(原来要 $2*N$ 个参数), 现在只需要 $(N+2)$ 个参数( $u$ 也需要2个参数)

但是此时就带来了误差, 如 $AA'$  和 $BB'$ , 所以我们要能够找到这样一个方向 $u$ , 使得所有原始点与投影点之间的误差最小

最小重构误差

样本点中心化





$$\vec{e} = \vec{x} - p_{rj} \vec{x}$$

$$= \vec{x} - \langle \vec{x}, \vec{u} \rangle \vec{u}$$

$$= x - (x^T u)u; x, u \in \mathbb{R}^n \quad \text{且} \|u\| = 1, u^T u = 1$$

$$J = \|\vec{e}\|^2 = e^T e = [x - (x^T u)u]^T [x - (x^T u)u]$$

$$= [x^T - (x^T u)u^T][x - (x^T u)u]$$

$$= x^T x - (x^T u)(x^T u) - (x^T u)(x^T u) + (x^T u)^2 u^T u$$

$$= \|x\|^2 - (x^T u)^2 - (x^T u)^2 + (x^T u)^2$$

$$= \|x\|^2 - (x^T u)^2$$



$$\max (x^T u)^2$$

$$\Leftrightarrow \max (x^T u)(x^T u) \Leftrightarrow \max (u^T x)(x^T u)$$

$$\Leftrightarrow \max u^T (x^T x) u$$

共有  $N$  个样本

$$\max \sum_{i=1}^N u^T (x_i^T x_i) u = u^T \left( \sum_{i=1}^N x_i^T x_i \right) u, \text{ 且 } \|u\| = 1$$

$$x_i^T x_i = X$$

$$\max u^T (X) u, \text{ st: } \|u\| = 1$$

$$L(u, \lambda) = u^T X u + \lambda(1 - u^T u)$$

$$\frac{\partial L}{\partial u} = 0 \Rightarrow Xu - \lambda u = 0$$

$$Xu = \lambda u$$

$$\frac{\partial L}{\partial \lambda} = 0 \Rightarrow u^T u = 1$$



**deepshare.net**

深度之眼

联系我们：

电话：18001992849

邮箱：[service@deepshare.net](mailto:service@deepshare.net)

Q Q：2677693114



公众号



客服微信

