

Pattern Recognition*

Homework II

* Teacher: Jianxin Wu. TA: ...

1st 张逸凯 171840708 计算机科学与技术系 本科生

Department of Computer Science and Technology

Nanjing University

zykhelloha@gmail.com

张逸凯 171840708 计算机科学与技术系 本科生

03_Framework 2. K-means
04_Error 2. Linear regression
04_Error 5. PR curve
04_Error 6. bias-variance decomposition
05_PCA 5. Jacobi eigenvalue algorithm

03_Framework 2. K-means

解：

1. 对于 $\forall x_i$ ，需要找到 x_i 的类别 $c_i = \arg \min_j \|x_i - \mu_j\|$ ，等价于：

$$J(x_i) = \begin{cases} \|x_i - \mu_j\| & (j = c_i) \\ 0 & (j \neq c_i) \end{cases}$$

对于 c_i, μ_j 的计算，令 C_i 为属于第 i 类样本的集合，详细地说为：

$$C_i = \{x_n : \|x_i - \mu_k\| \leq \text{all } \|x_i - \mu_j\|\}$$
$$\mu_k = \frac{1}{|C_i|} \sum_{x_i \in C_i} x_i \quad (1)$$

此时可以将 $J(x_i)$ 对 i 的每一项同乘 γ_{ij} ，当 $j = c_i$ 时 $\gamma_{ij} = 1$ ，反之 $\gamma_{ij} = 0$ 。这样就得到了K-means需要优化的目标：

$$\arg \min_{\gamma_{ij}, \mu_i} \sum_{i=1}^K \sum_{j=1}^M \gamma_{ij} \|x_j - \mu_i\|^2 \quad (2)$$

2. update γ_{ij} ：

$$\gamma_{ij} = \begin{cases} 1 & (\arg \min_j \|x_i - \mu_j\| = j) \\ 0 & (\text{others}) \end{cases}$$

update μ_i ：

$$\mu_i = \frac{1}{\sum_j \gamma_{ij}} \sum_j \gamma_{ij} x_j \quad (3)$$

3. 由上述K-means goal，可以定义loss函数如下：

$$f(\gamma, \mu) = \sum_{i=1}^K \sum_{j=1}^M \gamma_{ij} \|\mathbf{x}_j - \mu_i\|^2 \quad (4)$$

首先证明以下几个引理从而说明迭代步中**loss**函数是减小的。

Lemma 1: $\mathbb{E}\|X - z\|^2 = \mathbb{E}\|X - \mathbb{E}X\|^2 + \|z - \mathbb{E}X\|^2$, 其中 z 是其他数据点。
证明:

$$\begin{aligned} & \mathbb{E}\|X - \mathbb{E}X\|^2 + \|z - \mathbb{E}X\|^2 \\ &= \mathbb{E}[\|X\|^2 + \|\mathbb{E}X\|^2 - 2X \cdot \mathbb{E}X] + [\|z\|^2 + \|\mathbb{E}X\|^2 - 2z \cdot \mathbb{E}X] \\ &= \mathbb{E}\|X\|^2 + \|\mathbb{E}X\|^2 - 2\mathbb{E}X \cdot \mathbb{E}X + \|z\|^2 + \|\mathbb{E}X\|^2 - 2z \cdot \mathbb{E}X \\ &= \mathbb{E}\|X - z\|^2 \end{aligned}$$

Lemma 2: 以上代价在**k-means**的迭代(**b**题中的两个步骤)中都是单调减的。

证明: 令 $\gamma^{(t)}$ 是第 t 次迭代的结果, 对于算法的两个步骤:

- **update** γ_{ij} 之后 $f(\gamma^{(t+1)}, \mu) \leq f(\gamma^{(t)}, \mu)$ (易见因为 x 被分配给了最近的中心, 显然比之前更远的中心欧氏距离更小)。
- **update** μ_i 之后 $f(\gamma, \mu^{t+1}) \leq f(\gamma, \mu^t)$ (因为**Lemma 1**)。

\therefore 原命题得证。

以上证明了**loss function**是减小的, 下面证明收敛性: 我们知道最多只有 $\left\{ \begin{smallmatrix} n \\ k \end{smallmatrix} \right\}$ (第二类 **Stirling**数)种方法将 N 个数据划分成 k 类, 这是有限的。注意到每次迭代, 有且仅有以下两种影响:

1. 上述 f 值在两次迭代之间不变, 则接下来也不会变, 收敛。
2. 由上**Lemma2**的证明我们知道, 如果两次迭代后 f 值不同, 则迭代后 f 减小。又划分数有限, 所以收敛。

综上所述, 收敛性得证。

(ps: **k-means**确实是很有趣的东西, 偶然看了shen教授的Yinyang K-Means: A Drop-In Replacement of the Classic K-Means with Consistent Speedup, 优化方法很妙)

04_Error 2. Linear regression

解:

$$1. \quad cost = \frac{1}{N} \sum_{i=1}^n (\hat{y}_i - (x_i^T \beta + \epsilon))^2 \quad (5)$$

其中 \hat{y}_i 是Ground Truth。可以对 β, ϵ 求偏导来得到最优解。

2. 令 $\mathbf{E} \in \mathbb{R}^d$, $\mathbf{E} = (\epsilon, \dots, \epsilon)^T$ 。则有

$$\begin{aligned} \mathbf{y} &= \mathbf{X}^T \beta + \mathbf{E} \\ cost &= \frac{1}{N} (\hat{\mathbf{y}} - \mathbf{X}^T \beta - \mathbf{E})^T (\hat{\mathbf{y}} - \mathbf{X}^T \beta - \mathbf{E}) \end{aligned} \quad (6)$$

3. 不妨令 $\mathbf{E} = N \text{ loss}$, 因为同时对 β, ϵ 进行优化, 可以做一个吸收: $\beta := (\beta; \epsilon)$, 并把数据 \mathbf{X} 表示成一个 $n \times (d+1)$ 维度的矩阵。即:

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1d} & 1 \\ x_{21} & x_{22} & \dots & x_{2d} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nd} & 1 \end{pmatrix} = \begin{pmatrix} x_1^T & 1 \\ x_2^T & 1 \\ \vdots & \vdots \\ x_n^T & 1 \end{pmatrix} \quad (7)$$

$$\begin{aligned}
E &= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \boldsymbol{\beta} - \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y} + \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} \\
\frac{\partial E}{\partial \boldsymbol{\beta}} &= \frac{\partial \mathbf{y}^T \mathbf{y}}{\partial \boldsymbol{\beta}} - \frac{\partial \mathbf{y}^T \mathbf{X} \boldsymbol{\beta}}{\partial \boldsymbol{\beta}} - \frac{\partial \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y}}{\partial \boldsymbol{\beta}} + \frac{\partial \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\beta}}{\partial \boldsymbol{\beta}} + 0 \\
&= 0 - \mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{y} + (\mathbf{X}^T \mathbf{X} + \mathbf{X}^T \mathbf{X}) \boldsymbol{\beta} \\
&= 2\mathbf{X}^T (\mathbf{X} \boldsymbol{\beta} - \mathbf{y})
\end{aligned} \tag{8}$$

上式求偏导后为0，等价于 $\boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ ，其最后一列即为吸收的 ϵ 的值。

4. $\because \mathbf{X} \in \mathbb{R}^{n \times d}$. $\therefore \text{rank}(\mathbf{X}) = \text{rank}(\mathbf{X}^T \mathbf{X}) \leq n < d$

$\mathbf{X}^T \mathbf{X}$ 不满秩，所以不可逆。

5. ridge regression所加的L2可以防止过拟合，因为大多数的拟合是因为参数过大，使新的数据与原数据有一点差异就带来输出的巨大变化，加入正则化项后限制了参数 $\boldsymbol{\beta}$ 的大小，从另一个角度也减小了模型的复杂度，有效防止过拟合。

6. ridge regression的loss function可表达为：

$$\text{cost} = \frac{1}{N} (\hat{\mathbf{y}} - \mathbf{X}^T \boldsymbol{\beta} - \mathbf{E})^T (\hat{\mathbf{y}} - \mathbf{X}^T \boldsymbol{\beta} - \mathbf{E}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta} \tag{9}$$

同上题使用lagrange乘子法求得闭式解：

$$\boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \tag{10}$$

7. $(\mathbf{X} \mathbf{x})^T (\mathbf{X} \mathbf{x}) = \|\mathbf{X} \mathbf{x}\|^2 \geq 0$ (11)

所以 $\mathbf{X}^T \mathbf{X}$ 是半正定的， $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ 是正定的，可逆。这样 $\boldsymbol{\beta}$ 的闭式解存在。

8. $\circ \lambda = 0$ ，则ridge regression退化为一般的linear regression.

$\circ \lambda \rightarrow +\infty$ ，则loss function中正则化项的影响很大，优化时 $\boldsymbol{\beta} \rightarrow 0$.

9. λ 是不能作为普通的参数被模型学得的。

\circ 从梯度下降优化角度： $\lambda \boldsymbol{\beta}^T \boldsymbol{\beta}$ 是线性的，所以沿着负梯度方向， λ 会一直减小，这在 $\lambda \geq 0$ 下是无意义的。

\circ 从 λ 的作用角度： λ 是可以调整正则化项对模型复杂程度(或者说优化速率)的影响，使SGD作用在其他参数上时在真正的目标上获得了良好的性能，而不是通过直接优化 λ 来最小化loss function.

04_Error 5. PR curve

解：

	index	label	score	precision	recall	AUC-PR	AP
	0			1.0000	0.0000	—	—
	1	1	1.0	1.0000	0.2000	0.2000	0.2000
	2	2	0.9	0.5000	0.2000	0.0000	0.0000
	3	1	0.8	0.6667	0.4000	0.1167	0.1333
	4	1	0.7	0.7500	0.6000	0.1417	0.1500
1.	5	2	0.6	0.6000	0.6000	0.0000	0.0000
	6	1	0.5	0.6667	0.8000	0.1267	0.1333
	7	2	0.4	0.5714	0.8000	0.0000	0.0000
	8	2	0.3	0.5000	0.8000	0.0000	0.0000
	9	1	0.2	0.5556	1.0000	0.1056	0.1111
	10	2	0.1	0.5000	1.0000	0.0000	0.0000
						0.6906	0.7278

(12)

2. AUC-PR和AP是相似的，因为他们只是在计算PR曲线下面积时 运用了不同的插值计算近似方法。

3. 调换数据位置后的表为：

index	label	score	precision	recall	AUC-PR	AP
0			1.0000	0.0000	—	—
Same	as	the	above...
9	2	0.2	0.4444	0.8000	0.0000	0.0000
10	1	0.1	0.5000	1.0000	0.0944	0.1000
					0.6794	0.7167

(13)

4. 代码如下:

```

1 import csv
2 import numpy as np
3
4 def readData():
5     x = [
6         [1, 1, 1.0],
7         [2, 2, 0.9],
8         [3, 1, 0.8],
9         [4, 1, 0.7],
10        [5, 2, 0.6],
11        [6, 1, 0.5],
12        [7, 2, 0.4],
13        [8, 2, 0.3],
14        [9, 1, 0.2],
15        [10, 2, 0.1]
16    ]
17    return np.array(x).astype('float')
18
19 # 计算当前阈值下的TP, FP, FN, TN
20 def calculateTPFPFNTN(data, threshold):
21     TP, FP, FN, TN = 0, 0, 0, 0
22     for i in data:
23         if i[2] >= threshold:
24             if i[1] == 1:
25                 TP += 1
26             else:
27                 FP += 1
28         else:
29             if i[1] == 1:
30                 FN += 1
31             else:
32                 TN += 1
33     return TP, FP, FN, TN
34
35 def calculateAUCPRandAP(data):
36     befRec, befPre = 0, 1 # 上一步的rec和pre.
37     sumAUCPR, sumAP = 0, 0
38     for i in data:
39         TP, FP, FN, TN = calculateTPFPFNTN(data, i[2])
40         curPre, curRec = TP / (TP + FP), TP / (TP + FN)
41         AUCPR = (curRec - befRec) * (curPre + befPre) / 2 # 当前插值梯形的面
42         AP = curPre * (curRec - befRec) # 当前矩形的面积.
43         befRec, befPre = curRec, curPre
44         sumAUCPR += AUCPR
45         sumAP += AP
46         print("%.4f\t%.4f\t%.4f\t%.4f" % (curPre, curRec, AUCPR, AP))
47     print("#.##\t#.##\t%.4f\t%.4f" % (sumAUCPR, sumAP))
48
49 def main():
50     data = readData()
51     calculateAUCPRandAP(data)
52
53 if __name__ == '__main__':
54     main()

```

04_Error 6. bias-variance decomposition

解：(注：本题公式由 `mathpix` 软件识别)

1.

$$\begin{aligned} \mathbb{E} &\Leftrightarrow \mathbb{E}_D \\ \mathbb{E} [(y - f(x; D))^2] \\ &= \mathbb{E} [(F(x) - f(x; D) + \epsilon)^2] \\ &= \mathbb{E} [(F(x) - f(x; D))^2] + \mathbb{E} [\epsilon^2] + \mathbb{E}[2(F(x) - f(x; D))\epsilon] \\ &= (\mathbb{E}[F(x) - f(x; D)])^2 + \text{Var}[F(x) - f(x; D)] + (\mathbb{E}[\epsilon])^2 + \text{Var}(\epsilon) \\ &= (\mathbb{E}[F(x) - f(x; D)])^2 + \text{Var}(f(x; 0)) + \sigma^2 \\ &= (F(x) - \mathbb{E}[f(x; D)])^2 + \mathbb{E} [(f(x; D) - \mathbb{E}[f(x; D)])^2] + \sigma^2 \end{aligned}$$

其中第一项是偏差，第二项是方差，第三项是噪声。

2.

$$\begin{aligned} \because E[F] &= F \\ \mathbb{E}[f] \\ &= \mathbb{E} \left[\frac{1}{k} \sum_{i=1}^k y_{nn(i)} \right] \\ &= \frac{1}{k} \mathbb{E} \left[\sum_{i=1}^k F(x_{nn(i)}) + \epsilon \right] \\ &= \frac{1}{k} \mathbb{E} \left[\sum_{i=1}^k F \right] + 0 \\ &= \frac{\sum_{i=1}^k F(x_{nn(i)})}{k} \end{aligned}$$

3.

$$\begin{aligned} (\mathbb{E}[F - f])^2 &= (F - \mathbb{E}[f])^2 \\ &= \left(F(x) - \frac{\sum_{i=1}^k F(x_{nn(i)})}{k} \right)^2 \end{aligned}$$

$\because \epsilon_i$ is i.i.d.

$$\begin{aligned} \text{Var} \left(\frac{1}{k} \sum_{i=1}^k y_{nn(i)} \right) \\ &= \frac{1}{k^2} \text{Var} \left(\sum_{i=1}^k F - \epsilon_i \right) \\ &= \frac{1}{k^2} \text{Var} \left(\sum_{i=1}^k \epsilon_i \right) \\ &= \frac{\sigma^2}{k} \end{aligned}$$

(a)'s result :

$$\begin{aligned} \mathbb{E} [(y - f)^2] &= \mathbb{E} [(F - f)^2] + \sigma^2 \\ &= (\mathbb{E}[F - f])^2 + \text{Var}[F - f] + \sigma^2 \\ &= (\mathbb{E}[F - f])^2 + \text{Var}(f) + \sigma^2 \\ &= \left(F(x) - \frac{\sum_k F(x_{nn(i)})}{k} \right)^2 + \frac{\sigma^2}{k} + \sigma^2 \end{aligned}$$

4. 由(c)题知, variance term 为: $\frac{\sigma^2}{k}$, 它与 k 成反比. 随着 k 的增大而减小.

5. 由(c)题知, squared bias term 为 $\left(F(x) - \frac{\sum_k F(x_{nn(i)})}{k}\right)^2$, 考虑提示中 $k = n$ 时, squared bias term 将变大, 这是可解释的, 因为 k 较大时考虑的最近邻的点更多了, 也就是他们的平均更有可能远离这个点; 反之如果 k 较小, 所考虑的是更小最近邻集合, 那么他们的平均值更有可能接近这个点, 所以 bias 会更小.
- 即 squared bias term 随着 k 的增大而增大.

05_PCA 5. Jacobi eigenvalue algorithm

解: (注: 此题是在上完 Indian Institute of Technology Roorkee 的 Professor Dr. Sanjeev Kumar 的 Numerical Methods 课程 Lecture 13 后作答的)

1. $\because G^T G = I, \therefore \|x\| = \sqrt{x^T x} = \sqrt{(Gx)^T Gx} = \|G^T x\|$

同理: $\|x\| = \sqrt{(G^T x)^T G^T x} = \|G^T x\|$

2.

$$\text{tr}(AB) = \text{tr}(BA) \quad (14)$$

$$\begin{aligned} \|G^T X G\|_F &= \sqrt{\text{tr}(G^T X G G^T X^T G)} \\ &= \sqrt{\text{tr}(G^T X X^T G)} \\ &= \sqrt{\text{tr}(X^T G G^T X)} \\ &= \sqrt{\text{tr}(X^T X)} = \sqrt{\text{tr}(X X^T)} \\ &= \|X\|_F \end{aligned}$$

3. 首先证明一些引理:

Lemma: if $G^T G = I$, $\|GX\|_F^2 = \|X\|_F^2$.

证明:

$$\|GX\|_F^2 = \text{tr}(GX X^T G^T) = \text{tr}(X^T G^T G X) = \text{tr}(X^T X) = \|X\|_F^2 \quad (15)$$

Lemma: 若 A, B 相似, A 和 B 有相同特征值.

证明:

$$\begin{aligned} A &= GBG^{-1} \\ A\xi_i &= \lambda_i \xi_i \Rightarrow BG^{-1}\xi_i = \lambda_i G^{-1}\xi_i \end{aligned} \quad (16)$$

$\because J^T J = I, \therefore J^T X J, X$ 两矩阵相似, 相似矩阵有相同特征值.

对于任意的 J , 我们有 $\text{off}(J^T X J) < \text{off}(X)$, 也即在 J 连续作用之后(迭代步):

$\text{off}((J^T)^N X J^N)$, 即 X 除了对角线的元素平方和在不断减小, 最后只剩下对角线的元素, 又因为相似矩阵特征值相同, 所以最后对角线留下的是特征值, 在最后一小题将证明迭代步是收敛的.

4. (由于前几题使用 G , 本小题中 G, J 等价.) 因为 J 是正交的, 并且等价于旋转操作, 不失一般性考虑 Givens rotation matrix J :

$$J(i, j, \theta) = \begin{bmatrix} 1 & \cdots & 0 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & & \vdots & & \vdots \\ 0 & \cdots & \cos(\theta) & \cdots & -\sin(\theta) & \cdots & 0 \\ \vdots & & \vdots & \ddots & \vdots & & \vdots \\ 0 & \cdots & \sin(\theta) & \cdots & \cos(\theta) & \cdots & 0 \\ \vdots & & \vdots & & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 0 & \cdots & 1 \end{bmatrix} \quad \begin{matrix} \\ \\ (idx : i) \\ \\ (idx : j) \\ \\ \end{matrix}$$

其中 $J_{ii} = \cos(\theta)$, $J_{ij} = \sin(\theta)$, 依次类推 $J_{ji} = -\sin(\theta)$, $J_{jj} = \cos(\theta)$.

以上的Jacobi旋转矩阵 J 作用在 X 上时, 等价于做了一次相似(旋转)变换: J 改变了第 i, j 行和第 i, j 列.

令 $X' = J^T X J$, 注意到变换作用后:

$$\begin{aligned} X'_{ii} &= \cos(\theta)^2 X_{ii} - 2\sin(\theta)\cos(\theta)X_{ij} + \sin(\theta)^2 X_{jj} \\ X'_{jj} &= \sin(\theta)^2 X_{ii} + 2\sin(\theta)\cos(\theta)X_{ij} + \cos(\theta)^2 X_{jj} \\ X'_{ij} &= X'_{ji} = (\cos(\theta)^2 - \sin(\theta)^2) X_{ij} + \sin(\theta)\cos(\theta) (X_{ii} - X_{jj}) \end{aligned} \quad (17)$$

要使 $X'_{ij} = X'_{ji} = 0$, 等价于

$$\theta = \frac{1}{2} \arctan \left(\frac{2X_{ij}}{X_{jj} - X_{ii}} \right) \quad (18)$$

即 θ 取上述值即可满足题意.

5. 不妨令 $X' = J^T X J$, 我们有 $\|X'\|_F = \|X\|_F$, 注意到当 $m, n \neq i, j$ 时, $X'_{mn} = X_{mn}$, 即不在第 i, j 行或列上的元素不变. 从而我们可以证明 $\|X'\|_F - \sum_{i=1}^n (X'_i)^2 < \|X\|_F - \sum_{i=1}^n (X_i)^2$ 即可. 也就是对角线元素平方和(不妨令为 X_{diag})变大.

这里 p, q 就是上一题中的 i, j , X' 对应元素变换:

$$\begin{aligned} X'_{pp} &= \cos(\theta)^2 X_{pp} - 2\sin(\theta)\cos(\theta)X_{pq} + \sin(\theta)^2 X_{qq} \\ X'_{qq} &= \sin(\theta)^2 X_{pp} + 2\sin(\theta)\cos(\theta)X_{pq} + \cos(\theta)^2 X_{qq} \\ X'_{pq} &= X'_{qp} = (\cos(\theta)^2 - \sin(\theta)^2) X_{pq} + \sin(\theta)\cos(\theta) (X_{pp} - X_{qq}) \end{aligned} \quad (19)$$

注意 $|X_{pq}| = \max_{i \neq j} |X_{ij}|$ 在证明过程的最后一个 $>$ 中被使用到.

$$\begin{aligned} (X')^2_{diag} &= \sum_{i=1}^n (X'_i)^2 \\ &= \sum_{i \neq p, q} (X'_{ii})^2 + ((X'_{pp})^2 + (X'_{qq})^2) \\ &= \sum_{i \neq p, q} X_{ii}^2 + (X_{pp}^2 + X_{qq}^2 + 2X_{pq}^2) \\ &= \sum_{i=1}^n X_{ii}^2 + 2X_{pq}^2 \\ &> (X)^2_{diag} \end{aligned} \quad (20)$$

上面的证明过程化简是极麻烦的(二次项的平方), 所以我直接从 off 来证明:

$$\begin{aligned} k &\neq p, q \\ X'_{pk} &= \cos X_{pk} - \sin X_{qk} \\ X'_{qk} &= \sin X_{pk} + \cos X_{qk} \\ X' \text{ 的上半角部分 (因为是对称阵)} \\ &= \sum_k X'^2_{pk} + X'^2_{qk} + X'^2_{pq} \\ &= \sum_k (\sin(\theta)X_{pk} + \cos(\theta)X_{qk})^2 + (\cos(\theta)X_{pk} - \sin(\theta)X_{qk})^2 \\ &= X_{pk}^2 + X_{qk}^2 \\ &< X_{pk}^2 + X_{qk}^2 + X_{pq}^2 \end{aligned}$$

6. 延用上题中 $|X_{pq}| = \max_{i \neq j} |X_{ij}|$, 不妨令其等于 m , 又 $X \in \mathbb{R}^{n \times n}$, 所以除了对角线以外的元素有 $n(n-1)$ 个.

$$\begin{aligned} \text{off}(X)^2 &\leq n(n-1)X_{pq}^2 \\ \Rightarrow \text{off}(X')^2 &\leq \left(1 - \frac{2}{n(n-1)}\right) \text{off}(X)^2 \end{aligned} \quad (21)$$

上式已经可以说明收敛性了，不妨令 $X'^N = (J)^{N^T} X J^N$ ，对 $\forall \epsilon > 0, \exists N > N_0$ ，使

$$\left(1 - \frac{2}{n(n-1)}\right)^N < \epsilon.$$

收敛性得证。