

## 补充习题答案

在作业要求掌握的内容之外，还必须掌握的：

**习题 1** 能运用LU分解、Cholesky分解求解适定方程组

利用作业中习题4.1的结果，解线性方程组

$$Ax = b$$

其中  $b = \begin{pmatrix} 24 \\ 24 \\ 12 \end{pmatrix}$

**解** 先解  $Ly = b$ ，得

$$y = \begin{pmatrix} 24 \\ 12 \\ -4 \end{pmatrix}$$

再解  $Ux = y$ ，得

$$x = \begin{pmatrix} 6 \\ 6 \\ 6 \end{pmatrix}$$

**习题 2** 能运用法方程组、Cholesky分解、QR分解求解方程组的最小二乘解

利用作业中习题4.3的结果，解线性方程组的最小二乘解

$$Ax = b$$

其中  $b = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$

$$\text{解 } Q = \begin{pmatrix} \frac{3}{\sqrt{11}} & \frac{-1}{\sqrt{66}} & \frac{-1}{\sqrt{6}} \\ \frac{1}{\sqrt{11}} & \frac{7}{\sqrt{66}} & \frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{11}} & -\frac{4}{\sqrt{66}} & \frac{2}{\sqrt{6}} \end{pmatrix}, R = \begin{pmatrix} \sqrt{11} & \frac{8}{\sqrt{11}} \\ 0 & \frac{12}{\sqrt{66}} \\ 0 & 0 \end{pmatrix}, Q^T b = \begin{pmatrix} \frac{3}{\sqrt{11}} \\ \frac{-1}{\sqrt{66}} \\ \frac{-1}{\sqrt{6}} \end{pmatrix} \text{ 解}$$

$$\begin{pmatrix} \sqrt{11} & \frac{8}{\sqrt{11}} \\ 0 & \frac{12}{\sqrt{66}} \end{pmatrix} x = \begin{pmatrix} \frac{3}{\sqrt{11}} \\ \frac{-1}{\sqrt{66}} \end{pmatrix}$$

得最小二乘解:

$$x = \begin{pmatrix} 1/3 \\ -1/12 \end{pmatrix}$$

**习题 3** 熟练运用含有迹、行列式、矩阵的逆及其复合微分的计算方法

- $f(X) = \ln |X|$ ,  $X$  可逆, 求  $\frac{\partial f}{\partial X}$
- $f(X) = \|AX^{-1}\|_F^2$ ,  $X$  可逆, 求  $\frac{\partial f}{\partial X}$

**解** (1)  $f(X) = \ln |X|$

$$\begin{aligned} \frac{\partial f}{\partial X} &= \frac{\partial f}{\partial |X|} \frac{\partial |X|}{\partial X} \\ &= \frac{1}{|X|} |X| X^{-T} \\ &= X^{-T} \end{aligned}$$

因此

$$\frac{\partial f}{\partial X} = X^{-T}$$

(2)

$$\begin{aligned} f(X) &= \|X^{-1}\|_F^2 = \text{Tr}(X^{-T} A^T A X^{-1}) \\ df(X) &= \text{Tr}[d(X^{-T} A^T A X^{-1})] \\ &= \text{Tr}[(dX^{-T}) A^T A X^{-1} + X^{-T} A^T A dX^{-1}] \\ &= \text{Tr}[(X^{-1} dX X^{-1})^T A^T A X^{-1} + X^{-T} A^T A X^{-1} dX X^{-1}] \\ &= -\text{Tr}[X^{-T} dX^T X^{-T} A^T A X^{-1} + X^{-T} A^T A X^{-1} dX X^{-1}] \\ &= -\text{Tr}[(X^{-T} A^T A X^{-1} X^{-T})^T dX + (X^{-T} A^T A X^{-1} X^{-T})^T dX] \end{aligned}$$

所以

$$\frac{\partial f}{\partial \mathbf{X}} = -2\mathbf{X}^{-T} \mathbf{A}^T \mathbf{A} \mathbf{X}^{-1} \mathbf{X}^{-T}$$

#### 习题 4 理解MLE、MAP的思想

给定 $N$ 个独立同分布样本,  $\mathbf{x}_t, t = 1, 2, \dots, N$ , 服从多元正态分布

$$G(\mathbf{x}_t | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} \det(\boldsymbol{\Sigma})^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x}_t - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_t - \boldsymbol{\mu})\right\}$$

其中 $\mathbf{x}_t, \boldsymbol{\mu} \in \mathbb{R}^d$ ,  $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$ 是可逆对称矩阵。利用极大似然估计(MLE)估计参数 $\boldsymbol{\mu}, \boldsymbol{\Sigma}$ 。

解

$$p(\mathbf{X} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{t=1}^N G(\mathbf{x}_t | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

对上式取对数得到对数似然

$$L = \ln p(\mathbf{X} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{t=1}^N \ln G(\mathbf{x}_t | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

得

$$L = -\frac{Nd}{2} \ln(2\pi) - \frac{N}{2} \ln|\boldsymbol{\Sigma}| - \frac{1}{2} \sum_t (\mathbf{x}_t - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_t - \boldsymbol{\mu})$$

对数似然在 $\frac{\partial L}{\partial \boldsymbol{\mu}} = \mathbf{0}$ ,  $\frac{\partial L}{\partial \boldsymbol{\Sigma}} = \mathbf{0}$ 时取得极大值, 因此可以利用MLE估计参数。

##### 1. 估计 $\boldsymbol{\mu}$

$$\frac{\partial L}{\partial \boldsymbol{\mu}} = \sum_t \boldsymbol{\Sigma}^{-1} (\mathbf{x}_t - \boldsymbol{\mu}) \quad \text{令} \quad \frac{\partial L}{\partial \boldsymbol{\mu}} = \mathbf{0}$$

$$\boldsymbol{\mu} = \frac{1}{N} \sum_t \mathbf{x}_t$$

##### 2. 估计 $\boldsymbol{\Sigma}$

$$dL = d\left[\frac{N}{2} \ln|\boldsymbol{\Sigma}|\right] - d\left[\frac{1}{2} \sum_t (\mathbf{x}_t - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_t - \boldsymbol{\mu})\right]$$

第一项为

$$d\left[\frac{N}{2} \ln|\boldsymbol{\Sigma}|\right] = -\frac{N}{2} d[\ln|\boldsymbol{\Sigma}|] = -\frac{N}{2} \text{Tr}[\boldsymbol{\Sigma}^{-1} d\boldsymbol{\Sigma}]$$

第二项为

$$\begin{aligned}
 d\left[\frac{1}{2} \sum_t (\mathbf{x}_t - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_t - \boldsymbol{\mu})\right] &= -\frac{1}{2} d\text{Tr}\left[\sum_t (\mathbf{x}_t - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_t - \boldsymbol{\mu})\right] \\
 &= -\frac{1}{2} d\text{Tr}\left[\sum_t (\mathbf{x}_t - \boldsymbol{\mu})(\mathbf{x}_t - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}\right] \\
 &= -\frac{1}{2} \text{Tr}\left[\sum_t (\mathbf{x}_t - \boldsymbol{\mu})(\mathbf{x}_t - \boldsymbol{\mu})^T (-\boldsymbol{\Sigma}^{-1}(d\boldsymbol{\Sigma})\boldsymbol{\Sigma}^{-1})\right] \\
 &= \frac{1}{2} \text{Tr}[\boldsymbol{\Sigma}^{-1} \sum_t (\mathbf{x}_t - \boldsymbol{\mu})(\mathbf{x}_t - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} d\boldsymbol{\Sigma}]
 \end{aligned}$$

得到

$$\frac{\partial L}{\partial \boldsymbol{\Sigma}} = -\frac{N}{2} \boldsymbol{\Sigma}^{-1} + \frac{1}{2} \boldsymbol{\Sigma}^{-1} \sum_t (\mathbf{x}_t - \boldsymbol{\mu})(\mathbf{x}_t - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}$$

令  $\frac{\partial L}{\partial \boldsymbol{\Sigma}} = 0$ , 易得  $\boldsymbol{\Sigma} = \frac{1}{N} \sum_t (\mathbf{x}_t - \boldsymbol{\mu})(\mathbf{x}_t - \boldsymbol{\mu})^T$

**习题 5** 理解熵、互信息、交叉熵、KL散度的含义并掌握计算方法

证明，在多分类问题中，利用交叉熵函数作为损失函数和用KL散度作为损失函数是等价的。

**解** 设第 $i$ 个样本 $\mathbf{x}_i$ 属于 $y_i$ 类，则其真实标签分布为 $\mathbf{p}_i$ ， $\mathbf{p}_i$ 的第 $y_i$ 个分量为1，其余分量为0的向量。设预测器 $f$ 预测样本 $\mathbf{x}_i$ 标签分布为 $\mathbf{q}_i = f(\mathbf{x}_i; \boldsymbol{\theta})$ ， $\mathbf{q}_i$ 的每个分量为预测器预测的 $\mathbf{x}_i$ 属于每个分类的概率， $\boldsymbol{\theta}$ 是预测器中要学习的参数。

$$\text{KL散度} = (\mathbf{p}_i^T \log \mathbf{p}_i - \mathbf{p}_i^T \log \mathbf{q}_i)$$

$$\text{交叉熵} = (-\mathbf{p}_i^T \log \mathbf{q}_i)$$

由于真实标签分布是真实存在、固定不变的，因此

$$\arg \min_{\boldsymbol{\theta}} \text{KL散度} = \arg \min_{\boldsymbol{\theta}} \text{交叉熵}$$

**习题 6** 运用拉格朗日对偶函数解最优化问题

已知矩阵 $\mathbf{A} \in \mathbb{R}^{p \times q}$ ， $\mathbf{B} \in \mathbb{R}^{p \times r}$ 。 $\text{rank}(\mathbf{A}) = \min(p, q)$ 。未知矩阵 $\mathbf{X} \in \mathbb{R}^{q \times r}$ ，列出以下优化问题并求解。

- 若 $p < q$ ，求Frobenius范数最小的矩阵 $\mathbf{X}$ ，使得 $\mathbf{A}\mathbf{X} = \mathbf{B}$ 。

优化问题为

$$\min f(\mathbf{X}) = \frac{1}{2} \|\mathbf{X}\|_F^2$$

$$\text{s.t. } \mathbf{A}\mathbf{X} = \mathbf{B}$$

- 若  $p > q$ , 求矩阵  $\mathbf{X}$ , 使得矩阵  $\mathbf{A}\mathbf{X} - \mathbf{B}$  的Frobenius范数最小。  
优化问题为

$$\min f(\mathbf{X}) = \|\mathbf{A}\mathbf{X} - \mathbf{B}\|_F^2$$

**解** (1) 优化问题为

$$\min f(\mathbf{X}) = \frac{1}{2} \|\mathbf{X}\|_F^2$$

$$\text{s.t. } \mathbf{A}\mathbf{X} = \mathbf{B}$$

Lagrange函数为:

$$L(\mathbf{X}, \mathbf{\Lambda}) = \frac{1}{2} \mathbf{X}^T \mathbf{X} - \text{Tr}(\mathbf{\Lambda}^T (\mathbf{A}\mathbf{X} - \mathbf{B}))$$

$$\frac{\partial L}{\partial \mathbf{X}} = \mathbf{X} - \mathbf{A}^T \mathbf{\Lambda}$$

令  $\frac{\partial L}{\partial \mathbf{X}} = 0$ , 有:

$$\mathbf{X} = \mathbf{A}^T \mathbf{\Lambda}$$

$$g(\mathbf{\Lambda}) = -\frac{1}{2} \mathbf{\Lambda}^T \mathbf{A} \mathbf{A}^T \mathbf{\Lambda} + \mathbf{\Lambda}^T \mathbf{B}$$

令  $\frac{\partial g}{\partial \mathbf{\Lambda}} = 0$ :

$$-\mathbf{A} \mathbf{A}^T \mathbf{\Lambda} + \mathbf{B} = 0$$

由  $\mathbf{A} \in \mathbb{R}^{p \times q}$ ,  $\text{rank}(\mathbf{A}) = p$  得  $\mathbf{A} \mathbf{A}^T$  可逆, 因此

$$\mathbf{\Lambda} = (\mathbf{A} \mathbf{A}^T)^{-1} \mathbf{B}$$

因此,  $\mathbf{X}$  满足  $\mathbf{A}\mathbf{X} = \mathbf{B}$  的最小二范数解:

$$\mathbf{X} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T)^{-1} \mathbf{B}$$

因此,

$$\mathbf{X} = \mathbf{A}^T (\mathbf{A} \mathbf{A}^T)^{-1} \mathbf{B}$$

是满足  $\mathbf{A}\mathbf{X} = \mathbf{B}$  的Frobenius范数最小的矩阵。

(2) 若  $p > q$ , 优化问题

$$\min f(\mathbf{X}) = \|\mathbf{A}\mathbf{X} - \mathbf{B}\|_F^2$$

$f(\mathbf{X}) = \text{Tr}((\mathbf{A}\mathbf{X} - \mathbf{B})^T(\mathbf{A}\mathbf{X} - \mathbf{B}))$ , 当 $f$ 关于 $\mathbf{X}$ 梯度为零时 $f$ 最小。解方程

$$\frac{\partial f}{\partial \mathbf{X}} = \mathbf{A}^T(\mathbf{A}\mathbf{X} - \mathbf{B}) = 0$$

12

由 $p > q, \text{rank } \mathbf{A} = q$ 可得 $\mathbf{A}^T \mathbf{A}$ 可逆。得

$$\mathbf{X} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{B}$$

因此使得 $\mathbf{A}\mathbf{X} - \mathbf{B}$ 的Frobenius范数最小的矩阵 $\mathbf{X} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{B}$

14

**习题 7** 运用梯度下降法和牛顿法迭代求解单变量及多变量的最优化问题

第二次作业习题15, 优化问题改为 $\min_x x^3 - ax$ , 其中 $a > 0$ , 定义域 $\{x | x > 0\}$ 。

解

$$f(x) = x^3 - ax$$

$f'(x) = 3x^2 - a$ ,  $f''(x) = 6x$  所以迭代格式

$$x_{n+1} = x_n - \frac{3x_n^2 - a}{6x_n}$$