

# 机器学习导论 习题四

221300066, 季千<sup>①</sup>, qkjia@smail.nju.edu.cn

2024 年 5 月 28 日

## 作业提交注意事项

1. 作业所需的 LaTeX 及 Python 环境配置要求请参考: [Link];
2. 请在 LaTeX 模板中第一页填写个人的学号、姓名、邮箱;
3. 本次作业需提交的文件为:
  - (a) 作答后的 LaTeX 代码 — HW4.tex;
  - (b) 由 (a) 编译得到的 PDF 文件 — HW4.pdf;
  - (c) 题目 2.(2) 的求解代码文件 — svm\_qp\_dual.py
  - (d) 题目 3 的求解代码文件 — svm\_kernel\_solution.py

其他文件 (如其他代码、图片等) 无需提交. 请将以上文件**打包为** 学号\_姓名.zip (例如 221300001\_张三.zip) 后提交;
3. 若多次提交作业, 则在命名 .zip 文件时加上版本号, 例如 221300001\_张三\_v1.zip” (批改时以版本号最高的文件为准);
4. 本次作业提交截止时间为 **5 月 28 日 23:59:59**. 未按照要求提交作业, 提交作业格式不正确, **作业命名不规范**, 将会被扣除部分作业分数; 除特殊原因 (如因病缓交, 需出示医院假条) 逾期未交作业, 本次作业记 0 分; **如发现抄袭, 抄袭和被抄袭双方成绩全部取消**;
5. 学习过程中, 允许参考 ChatGPT 等生成式语言模型的生成结果, 但必须在可信的信息源处核实信息的真实性; **不允许直接使用模型的生成结果作为作业的回答内容**, 否则将视为作业非本人完成并取消成绩;
6. 本次作业提交地址为 [Link], 请大家预留时间提前上交, 以防在临近截止日期时, 因网络等原因无法按时提交作业.

## 1 [35pts] Soft Margin

考虑软间隔 SVM 问题, 其原问题形式如下:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi_i} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^m \xi_i^p \\ \text{s.t.} \quad & y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0, i \in [m]. \end{aligned} \quad (1.1)$$

其中, 松弛变量  $\boldsymbol{\xi} = \{\xi_i\}_{i=1}^m, \xi_i > 0$  表示样本  $\mathbf{x}_i$  对应的间隔约束不满足的程度, 在优化问题中加入惩罚  $C \sum_{i=1}^m \xi_i^p, C > 0, p \geq 1$  使得不满足约束的程度尽量小 ( $\xi_i \rightarrow 0$ ). 课本式 (6.35) 即为  $p = 1$  时对应的情况, 此时, 所有违反约束的样本都会受到相同比例的惩罚, 而不考虑它们违反约束的程度. 这可能导致模型对较大偏差的样本不够敏感, 不足以强调更严重的违规情况. 下面将考虑一些该问题的变式:

- (1) [2+7pts] 我们首先考虑  $p = 2$  的情况, 它对于违反约束程度较大的样本提供了更大的惩罚.
  - (a) 如课本式 (6.34)-(6.35) 所述,  $p = 1$  的情况对应了 hinge 损失  $\ell_{hinge} : x \rightarrow \max(0, 1 - x)$ . 请直接写出  $p = 2$  的情况下对应的损失函数.
  - (b) 请推导  $p = 2$  情况下软间隔 SVM 的对偶问题.
- (2) [14pts]  $p = 1$  的情况下, 相当于对向量  $\boldsymbol{\xi}$  使用  $L_1$  范数惩罚:  $\|\boldsymbol{\xi}\|_1 = \sum_i |\xi_i|$ . 现在, 我们考虑使用  $L_\infty$  范数惩罚:  $\|\boldsymbol{\xi}\|_\infty = \max_i \xi_i$ , 这会使得模型着重控制最大的违背约束的程度, 从而促使模型在最坏情况下的表现尽可能好. 请推导使用  $L_\infty$  范数惩罚的原问题和对偶问题.
- (3) [4+8pts] 在(1.1)中, 正例和负例在目标函数中分类错误的“惩罚”是相同的. 然而在实际场景中, 很多时候正例和负例错分的“惩罚”代价是不同的 (参考教材 2.3.4 节). 比如考虑癌症诊断问题, 将一个确实患有癌症的人误分类为健康人, 以及将健康人误分类为患有癌症, 产生的错误影响以及代价不应该认为是等同的. 所以我们考虑对负例违反间隔约束的样本施加  $k > 0$  倍于正例中违反间隔约束的样本的“惩罚”.
  - (a) 令(1.1)中  $p = 1$ , 并令所有正例样本的集合为  $D_+$ , 负例样本的集合为  $D_-$ . 请给出相应的 SVM 优化问题.
  - (b) 请给出相应的对偶问题.

**Solution.** 此处用于写解答 (中英文均可)

(1) a. 对于  $p = 2$  的情况, 对应的损失函数为:

$$\ell_{hinge}(x) = \begin{cases} 0 & \text{if } x \geq 1 \\ (1 - x)^2 & \text{otherwise} \end{cases}$$

b.  $p = 2$  情况下软间隔 SVM 的对偶问题的推导:

$$L(\mathbf{w}, \mathbf{b}, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\mu}) = \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^m \xi_i^2 + \sum_{i=1}^m \alpha_i (1 - \xi_i - y_i (\mathbf{w}^\top \mathbf{x}_i + b)) - \sum_{i=1}^m \mu_i \xi_i$$

其中  $\alpha_i \geq 0, \mu_i \geq 0$  是拉格朗日乘子.

对  $\mathbf{w}, \mathbf{b}, \xi$  求导可得:

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i = 0 \Rightarrow \hat{\mathbf{w}} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial L}{\partial b} = -\sum_{i=1}^m \alpha_i y_i = 0 \Rightarrow \sum_{i=1}^m \alpha_i y_i = 0$$

$$\frac{\partial L}{\partial \xi_i} = 2C\xi_i - \alpha_i - \mu_i = 0 \Rightarrow \alpha_i + \mu_i = 2C\xi_i$$

带入可得:

$$L(\boldsymbol{\alpha}, \boldsymbol{\mu}) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j - \frac{1}{4C} \sum_{i=1}^n (\alpha_i + \mu_i)^2$$

$$\begin{aligned} \max_{\boldsymbol{\alpha}, \boldsymbol{\mu}} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j - \frac{1}{4C} \sum_{i=1}^n (\alpha_i + \mu_i)^2 \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \\ & 0 \leq \alpha_i, i = 1, 2, \dots, m, \quad 0 \leq \mu_i, i = 1, 2, \dots, m. \end{aligned} \quad (1.2)$$

(2) 原问题的优化目标为:

$$\min_{\mathbf{w}, b, \xi, t} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + Ct$$

约束条件为:

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, i \in [m]$$

$$t \geq \xi_i \geq 0, i \in [m]$$

$$L(\mathbf{w}, \mathbf{b}, \xi, \alpha, \mu, \beta, t) = \frac{1}{2} \|\mathbf{w}\|_2^2 + Ct + \sum_{i=1}^m \alpha_i (1 - \xi_i - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) - \sum_{i=1}^m \mu_i \xi_i + \sum_{i=1}^m \beta_i (\xi_i - t)$$

其中  $\alpha_i \geq 0, \mu_i \geq 0, \beta_i \geq 0$  是拉格朗日乘子.

对  $\mathbf{w}, \mathbf{b}, \xi, t$  求导可得:

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i = 0 \Rightarrow \mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial L}{\partial b} = -\sum_{i=1}^m \alpha_i y_i = 0 \Rightarrow \sum_{i=1}^m \alpha_i y_i = 0$$

$$\frac{\partial L}{\partial \xi_i} = \beta_i - \alpha_i - \mu_i = 0 \Rightarrow \alpha_i + \mu_i = \beta_i$$

$$\frac{\partial L}{\partial t} = C - \sum_{i=1}^m \beta_i = 0 \Rightarrow \sum_{i=1}^m \beta_i = C$$

带入可得对偶问题为:

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \\ & 0 \leq \alpha_i, i = 1, 2, \dots, m, \quad \sum_{i=1}^m \alpha_i \leq C. \end{aligned} \quad (1.3)$$

(3) a. 对于  $p = 1$  的情况, 正例和负例违反间隔约束的“惩罚”代价不同, 可以引入权重系数  $k > 0$ 。此时, SVM 优化问题为:

$$\min_{\mathbf{w}, b, \xi} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i \in D_+} k \xi_i + kC \sum_{i \in D_-} k \xi_i$$

约束条件为：

$$y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i,$$

$$\xi_i \geq 0, i \in [m]$$

b. 对应的对偶问题为：

$$\begin{aligned} \max_{\boldsymbol{\alpha}, \boldsymbol{\mu}} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C, i \in D_+, \quad 0 \leq \alpha_i \leq kC, i \in D_-, . \end{aligned} \tag{1.4}$$

## 2 [20pts] Primal and Dual Problem

给定一个包含  $m$  个样本的数据集  $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ , 其中每个样本的特征维度为  $d$ , 即  $\mathbf{x}_i \in \mathbb{R}^d$ . 软间隔 SVM 的原问题和对偶问题可以表示为:

原问题:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, \forall i \in [m] \\ & \xi_i \geq 0, \forall i \in [m] \end{aligned}$$

对偶问题:

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^\top Q \alpha - \mathbf{1}_m^\top \alpha \\ \text{s.t.} \quad & \mathbf{y}^\top \alpha = 0 \\ & 0 \leq \alpha_i \leq C, \forall i \in [m] \end{aligned}$$

其中, 对于任意  $i, j \in [m]$  有  $Q_{ij} \equiv y_i y_j \mathbf{x}_i^\top \mathbf{x}_j$ .

上述的原问题和对偶问题都是二次规划 (Quadratic Programming) 问题, 都可以使用相关软件包求解. 本题目中我们将通过实践来学习凸优化软件包的使用, 并以软间隔 SVM 为例了解原问题、对偶问题在优化方面的特性.

- (1) [2pts] 请直接写出原问题和对偶问题的参数量 (注意参数只包含分类器所保存的参数, 不包含中间变量).
- (2) [10pts] 请参考 lab2/svm\_qp.py 中对于原问题的求解代码, 编写对偶问题的求解代码 lab2/svm\_qp\_dual.py. (这里使用了 CVXPY 求解 QP 问题.) 请将代码提交至下方的解答处.
- (3) [8pts] 设特征维度和样例数量的比值  $r = \frac{d}{m}$ , 请绘制原问题和对偶问题的求解速度随着这个比值变化的曲线图. 并简述: 何时适合求解原问题, 何时适合求解对偶问题?

**Solution.** 此处用于写解答 (中英文均可)

- (1) 原问题的参数量是  $d + 1 + m$ , 对偶问题的参数量是  $m$
- (2) 对偶问题的求解代码为:

```
1 def solve_dual(X, y, C):
2     '''
3     :参数 X: ndarray, 形状为(m, d), 样例矩阵
4     :参数 y: ndarray, 形状为(m), 样例标签向量
5     :参数 C: 标量, 含义与教材式(6.35)中C相同
6     :返回: alpha, SVM的对偶变量
7     '''
8     m, d = X.shape
9     alpha = cp.Variable(m)
10    y = y.astype(float).reshape(-1, 1)
11    c_ = np.full((m,1),C)
12
13    K = X @ X.T
14    q_1 = np.outer(y, y) * K
15
16    prob = cp.Problem(
17        cp.Minimize(0.5 * cp.quad_form(alpha, q_1) - cp.sum(alpha)), # 目标函数
18        [alpha >= 0, alpha <= c_, y.T @ alpha == 0] # 约束
```

```

19     )
20     prob.solve()
21     return alpha.value

```

(3) 曲线图为:

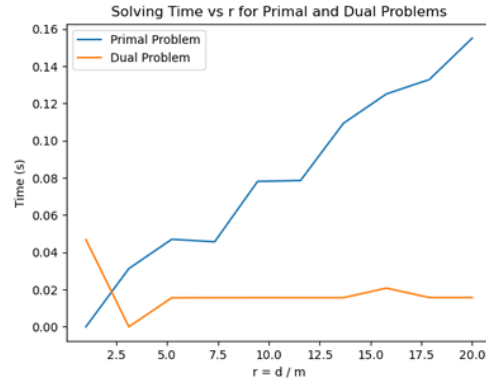


图 1: 原问题和对偶问题的求解速度随比值变化的曲线图

简述题: 当特征维度和样例数量的比值  $r = \frac{d}{m}$  较大时, 适合求解对偶问题; 当  $r$  较小时, 适合求解原问题。这是因为在高维空间中, 数据点之间的距离较远, 对偶问题的参数量相对较少, 使用对偶问题可以更好地处理这种情况。同理在低维空间中, 数据点之间的距离较近, 使用原问题可以更好地处理这种情况。

### 3 [15pts] Kernel Function in Practice

lab3/svm\_kernel.py 中构造了异或 (XOR) 问题, 如图 2 所示. 该问题是线性不可分的.

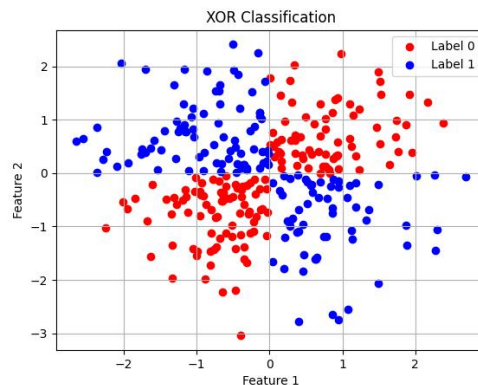


图 2: 异或 (XOR) 问题

本题中我们将通过实验了解核函数的选择对于 SVM 解决非线性问题的影响. 请使用 sklearn 包中的 SVM 分类器完成下述实验:

- (1) [6pts] 请分别训练线性核 SVM 分类器和核 (RBF 核) SVM 分类器, 并绘制出各自的决策边界.
- (2) [6pts] sklearn 还允许自定义核函数, 参考 lab3/svm\_kernel\_custom.py 的用法, 编写核函数  $\kappa(\mathbf{x}, \mathbf{x}') = \frac{1}{1 + \|\mathbf{x} - \mathbf{x}'\|_2^2}$ , 训练该核函数的 SVM 分类器, 并绘制出决策边界.

具体的实验要求可以参考 lab3/svm\_kernel.py 的 main 部分. 请将 lab3/svm\_kernel\_solution.py 中的代码和三个核函数分别对应的决策边界图提交至下方的解答处.

最后, 请直接回答 [3pts]: 三个核函数, 各自能够解决异或 (XOR) 分类问题吗?

**Solution.** 此处用于写解答 (中英文均可)

求解代码为:

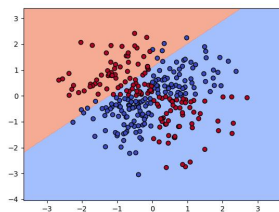
```
1 from sklearn.svm import SVC
2 import numpy as np
3
4 def svm_kernel_linear(X, Y):
5     '''
6     :参数 X: ndarray, 形状(m, d), 样例矩阵
7     :参数 Y: ndarray, 形状(m), 样例标签向量
8     :返回: clf_linear, 训练好的分类器
9     '''
10    clf_linear = SVC(kernel='linear')
11    clf_linear.fit(X, Y)
12    return clf_linear
13
14 def svm_kernel_rbf(X, Y):
15     '''
16     :参数 X: ndarray, 形状(m, d), 样例矩阵
17     :参数 Y: ndarray, 形状(m), 样例标签向量
```

```

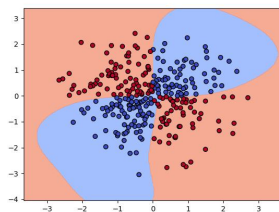
18     :返回: clf_rbf, 训练好的分类器
19     '''
20     clf_rbf = SVC(kernel='rbf')
21     clf_rbf.fit(X, Y)
22     return clf_rbf
23
24 def custom_kernel(X1, X2):
25     '''
26     :参数 X1: ndarray, 形状(m, d)
27     :参数 X2: ndarray, 形状(n, d)
28     :返回: 形状为(m, n)的Gram矩阵, 第(i,j)个元素为X1[i]和X2[j]之间的核函数值
29     '''
30     gram_matrix = np.sum(X1**2, axis=1).reshape(-1, 1) + np.sum(X2**2, axis=1) - 2 * np.
        dot(X1, X2.T)
31     return 1 / (1 + gram_matrix)
32
33 def svm_kernel_custom(X, Y):
34     '''
35     :参数 X: ndarray, 形状(m, d), 样例矩阵
36     :参数 Y: ndarray, 形状(m), 样例标签向量
37     :返回: clf_custom, 训练好的分类器
38     '''
39     # hint: 需要使用 custom_kernel 函数
40     clf_custom = SVC(kernel=custom_kernel)
41     clf_custom.fit(X, Y)
42     return clf_custom

```

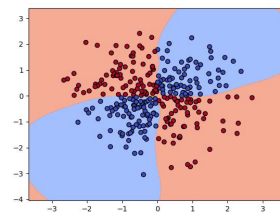
决策边界为:



(a) clf\_linear



(b) clf\_rbf



(c) clf\_custom

能否解决异或 (XOR) 分类问题: 线性核函数不能解决, 高斯核函数和自定义核函数可以解决.



## 4 [30pts] Maximum Likelihood Estimation

给定由  $m$  个样本组成的训练集  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ , 其中  $\mathbf{x}_i \in \mathbb{R}^d$  是第  $i$  个示例,  $y_i \in \mathbb{R}$  是对应的实值标记. 令  $\mathbf{X} \in \mathbb{R}^{m \times d}$  表示整个训练集中所有样本特征构成的矩阵, 并令  $\mathbf{y} \in \mathbb{R}^m$  表示训练集中所有样本标记构成的向量. 线性回归的目标是寻找一个参数向量  $\mathbf{w} \in \mathbb{R}^d$ , 使得在训练集上模型预测的结果和真实标记之间的差距最小. 对于一个样本  $\mathbf{x}$ , 线性回归给出的预测为  $\hat{y} = \mathbf{w}^\top \mathbf{x}$ ,<sup>1</sup> 它与真实标记  $y$  之间的差距可以用平方损失  $(\hat{y} - y)^2$  来描述. 因此, 在整个训练集上最小化损失函数的过程可以写作如下的优化问题:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 \quad (4.1)$$

- (1) [8pts] 考虑这样一种概率观点: 样本  $\mathbf{x}$  的标记  $y$  是从一个高斯分布  $\mathcal{N}(\mathbf{w}^\top \mathbf{x}, \sigma^2)$  中采样得到的. 这个高斯分布的均值由样本特征  $\mathbf{x}$  和模型参数  $\mathbf{w}$  共同决定, 而方差是一个额外的参数  $\sigma^2$ . 基于这种概率观点, 我们可以基于观测数据对高斯分布中的参数  $\mathbf{w}$  做极大似然估计. 请证明:  $\mathbf{w}$  的极大似然估计结果  $\mathbf{w}_{\text{MLE}}$  与式 (4.1) 中的  $\mathbf{w}^*$  相等;
- (2) [9pts] 极大似然估计容易过拟合, 一种常见的解决办法是采用最大后验估计: 沿着上一小问的思路, 现在我们在概率建模下对参数  $\mathbf{w}$  做最大后验估计. 为此, 引入参数  $\mathbf{w}$  上的先验  $p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \lambda \mathbf{I})$ . 其中, 均值  $\mathbf{0}$  是  $d$  维的全 0 向量,  $\mathbf{I}$  是  $d$  维单位矩阵,  $\lambda > 0$  是一个控制方差的超参数. 现在, 请推导对  $\mathbf{w}$  做最大后验估计的目标函数, 并讨论一下该结果与“带有  $L_2$  范数正则项的线性回归”之间的关系;
- (3) [9pts] 沿着上一小问的思路, 我们尝试给参数  $\mathbf{w}$  施加一个拉普拉斯先验. 简便起见, 我们假设参数  $\mathbf{w}$  的  $d$  个维度之间是独立的, 且每一维都服从 0 均值的一元拉普拉斯分布, 即:

$$p(\mathbf{w}) = \prod_{j=1}^d p(w_j), \quad (4.2)$$

$$p(w_j) = \text{Lap}(w_j | 0, \lambda), \quad j = 1, 2, \dots, d.$$

请推导对  $\mathbf{w}$  做最大后验估计的目标函数, 并讨论一下该结果与“带有  $L_1$  范数正则项的线性回归”之间的关系;

Note: 由参数  $\mu, \lambda$  确定的一元拉普拉斯分布的概率密度函数为:

$$\text{Lap}(w | \mu, \lambda) = \frac{1}{2\lambda} \exp\left(-\frac{|w - \mu|}{\lambda}\right). \quad (4.3)$$

- (4) [4pts] 基于 (2) 和 (3) 的结果, 从概率角度讨论为什么  $L_1$  范数能使模型参数更稀疏.

**Solution.** 此处用于写解答 (中英文均可)

- (1) 对于高斯分布, 其对数似然函数为:

$$\begin{aligned} \log p(\mathbf{y}|\mathbf{w}) &= \log \prod_{i=1}^m \mathcal{N}(y_i | \mathbf{w}^\top \mathbf{x}_i, \sigma^2) \\ n &= \sum_{i=1}^m \log \mathcal{N}(y_i | \mathbf{w}^\top \mathbf{x}_i, \sigma^2) \\ n &= \frac{m}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 \end{aligned}$$

<sup>1</sup> 本题不考虑偏移  $b$ , 可参考教材第 3 章将偏移  $b$  吸收进  $\mathbf{w}$ .

要找到一个  $\mathbf{w}$  使得  $\ln L(\hat{\mathbf{w}}|\mathbf{X}, \mathbf{y})$  最大, 即  $\hat{\mathbf{w}}_{MLE} = \arg \max_{\mathbf{w}} \ln L(\hat{\mathbf{w}}|\mathbf{X}, \mathbf{y})$  前一项显然与  $\mathbf{w}$  无关, 后一项和式 (4.1) 成反比. 最大化对数似然函数等价于最小化平方损失函数. 因此,  $\mathbf{w}_{MLE}$  与式 (4.1) 中的  $\mathbf{w}^*$  显然相等.

(2) 由题意知, 此分布的概率密度函数为:

$$p(\mathbf{w}) = \frac{1}{(2\pi\lambda)^{d/2}} \exp\left(-\frac{\|\mathbf{w}\|^2}{2\lambda}\right)$$

我们需要计算参数  $\mathbf{w}$  的后验分布  $p(\mathbf{w}|\mathbf{y}, \mathbf{X})$ , 根据贝叶斯公式:

$$p(\mathbf{w}|\mathbf{y}, \mathbf{X}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{p(\mathbf{y}|\mathbf{X})}$$

$$\log p(\mathbf{w}|\mathbf{y}, \mathbf{X}) = \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) + \log p(\mathbf{w}) - \log p(\mathbf{y}|\mathbf{X}) = \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) + \log p(\mathbf{w})$$

$$\log p(\mathbf{w}|\mathbf{y}, \mathbf{X}) = -\frac{m}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^m (y_i - \mathbf{w}^T \mathbf{x}_i)^2 - \frac{d}{2} \log(2\pi\lambda) - \frac{1}{2\lambda} \|\mathbf{w}\|^2$$

$$\mathbf{w}_{MAP} = \arg \min_{\mathbf{w}} \left( \frac{1}{2\sigma^2} \sum_{i=1}^m (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \frac{1}{2\lambda} \|\mathbf{w}\|^2 \right)$$

$$\mathbf{w}_{MAP} = \arg \min_{\mathbf{w}} (\|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \sigma^2 \lambda \|\mathbf{w}\|^2)$$

这个目标函数正是带有  $L_2$  正则项的线性回归, 其正则化系数为  $\frac{\sigma^2}{\lambda}$ .

(3) 推导过程:

$$p(\mathbf{w}) = \prod_{j=1}^d p(\mathbf{w}_j), \quad p(\mathbf{w}_j) = \text{Lap}(\mathbf{w}_j|0, \lambda) = \frac{1}{2\lambda} \exp\left(-\frac{|\mathbf{w}_j|}{\lambda}\right)$$

$$\log p(\mathbf{w}|\mathbf{y}, \mathbf{X}) = \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) + \log p(\mathbf{w}) - \log p(\mathbf{y}|\mathbf{X}) = \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) + \log p(\mathbf{w})$$

$$\log p(\mathbf{w}|\mathbf{y}, \mathbf{X}) = -\frac{m}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^m (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \sum_{j=1}^d \log \frac{1}{2\lambda} \exp(-|\mathbf{w}_j|/\lambda)$$

$$\log p(\mathbf{w}|\mathbf{y}, \mathbf{X}) = -\frac{m}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^m (y_i - \mathbf{w}^T \mathbf{x}_i)^2 - \sum_{j=1}^d \log(2\lambda) + |\mathbf{w}_j|/\lambda$$

$$\mathbf{w}_{MAP} = \arg \min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \frac{2\sigma^2}{\lambda} \|\mathbf{w}\|_1$$

这个目标函数正是带有  $L_1$  正则项的线性回归, 其正则化系数为  $\frac{2\sigma^2}{\lambda}$ .

(4) 从概率角度来说,  $L_1$  范数对应于拉普拉斯分布, 这种分布鼓励模型参数中的许多值接近零, 因为拉普拉斯分布在 0 处的陡峭, 可以使模型参数更稀疏. 它对参数的绝对值进行惩罚, 使得参数更倾向于取较小的值. 这样, 在最大后验估计中, 参数的先验分布会偏向于较小的值, 从而使得模型参数更加稀疏.

## Acknowledgments

允许与其他同样未完成作业的同学讨论作业的内容, 但需在此注明并加以致谢; 如在作业过程中, 参考了互联网上的资料 (包括生成式模型的结果), 且对完成作业有帮助的, 亦需注明并致谢.