

2024 秋季高级机器学习 习题二参考答案

2025.1.17

一. (30 points) 特征选择与稀疏学习

1. (20 points) 教材中提到, 为了缓解过拟合问题, 可对损失函数引入正则化项。给定包含 m 个样例的数据集 $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$, 其中 $y_i \in \mathbb{R}$ 为 \mathbf{x}_i 的实数标记, $\mathbf{x}_i = (x_{i1}; x_{i2}; \dots; x_{id}) \in \mathbb{R}^d$ 。针对数据集 D 中的 m 个示例, 以平方误差为损失函数, 使用 $\sum_j |w_j|^q$ 作为正则项, 可以得到带正则化的误差项

$$\sum_{i=1}^m (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda \sum_{j=1}^d |w_j|^q, \quad (1)$$

其中 \mathbf{w} 是待学习参数, $\lambda > 0$ 是正则化系数。

- (1) (10 points) 试说明最小化以上不带约束的问题与最小化下面带约束的问题等价。(提示: 可以利用拉格朗日乘子)

$$\begin{aligned} & \text{minimize}_{\mathbf{w}} \quad \sum_{i=1}^m (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 \\ & \text{subject to} \quad \sum_{j=1}^d |w_j|^q \leq \eta, \end{aligned} \quad (2)$$

- (2) (10 points) 在 (1) 的基础上, 请讨论 η 和 λ 之间的联系。(提示: 可以考虑 KKT 条件)

2. (10 points) 字典学习与压缩感知都有对稀疏性的利用, 请你分析两者对稀疏性利用的异同点。

解:

1. 注: 本题题干中“ $\lambda > 0$ ”存在争议。 $\lambda > 0$ 方便了第 (2) 小问的讨论, 但是当 η 约束对原始最小二乘问题最优解没有影响时, 优化问题 (1) 中 λ 取 0 才可等价, 因此严格来说 $\lambda \geq 0$ 时两个优化问题等价。(例如: 考虑由 $y = x$ 产生的数据集 $D = \{(1, 1), (2, 2)\}$, 取 $q = 1$, 显然无约束的最小二乘问题最优解 $w = 1$, 若此时 $\eta \geq 1$, 则优化问题 (2) 中的约束项不影响最优解, 若要使优化问题 (1) 中最优解也为 $w = 1$, 则必须要使得正则化项对优化没有影响, 即有 $\lambda = 0$ 。)

(1) 记优化问题 (1) 的最优解为 \mathbf{w}_1^* , 优化问题 (2) 的最优解为 \mathbf{w}_2^* , 要证等价, 即证

$$\forall \lambda \geq 0, \exists \eta \geq 0, \mathbf{w}_1^* = \mathbf{w}_2^* \quad (\dagger)$$

$$\forall \eta \geq 0, \exists \lambda \geq 0, \mathbf{w}_1^* = \mathbf{w}_2^* \quad (\dagger\dagger)$$

记

$$L_1 = \sum_{i=1}^m (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda \sum_{j=1}^d |w_j|^q,$$

最优解 \mathbf{w}_1^* 应满足

$$\frac{\partial L_1}{\partial \mathbf{w}} = -2 \sum_{i=1}^m \mathbf{x}_i (y_i - \mathbf{w}^\top \mathbf{x}_i) + \lambda q \mathbf{w}^{q-1} \cdot \text{sign}(\mathbf{w}) = 0 \quad (\dagger)$$

其中

$$\mathbf{w}^{q-1} = \begin{bmatrix} |w_1|^{q-1} \\ |w_2|^{q-1} \\ \vdots \\ |w_d|^{q-1} \end{bmatrix},$$

$\text{sign}(\cdot)$ 表示符号函数。为优化问题 (2) 引入拉格朗日函数

$$L_2 = L(\mathbf{w}, \nu) = \sum_{i=1}^m (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \nu \left(\sum_{j=1}^d |w_j|^q - \eta \right)$$

最优解 \mathbf{w}_2^* 应满足的 KKT 条件为

$$\begin{cases} \nu \geq 0 \\ \nu \left(\sum_{j=1}^d |w_j|^q - \eta \right) = 0 \\ \sum_{j=1}^d |w_j|^q - \eta \leq 0 \\ \frac{\partial L_2}{\partial \mathbf{w}} = -2 \sum_{i=1}^m \mathbf{x}_i (y_i - \mathbf{w}^\top \mathbf{x}_i) + \nu q \mathbf{w}^{q-1} \cdot \text{sign}(\mathbf{w}) = 0 \end{cases} \quad (\ddagger)$$

比对式 (†) 与式 (‡), 对于 (†), 取 $\nu = \lambda, \eta = \sum_{j=1}^d |w_{1j}^*|^q$ 即可使得 \mathbf{w}_1^* 满足优化问题 (2); 对于 (‡), 取 $\lambda = \nu$ 即可使得 \mathbf{w}_2^* 满足优化问题 (1), 因此等价。

(2) 当 $\lambda > 0$ 时, 正则化约束存在, 因此优化问题 (1) 中最优解可表示成为关于 λ 的函数 $\mathbf{w}_1^*(\lambda)$, 由第 (1) 小问中等价讨论可知, 取 $\nu = \lambda > 0$, 因此有 $\eta = \sum_{j=1}^d |w_{1j}^*(\lambda)|^q$ 。

当 $\lambda = 0$ 时, 正则化约束消失, 优化问题 (1) 退化成一个无约束的最小二乘问题, 由第 (1) 小问中讨论并结合 (‡), 取 $\nu = \lambda = 0, \eta \geq \sum_{j=1}^d |w_{1j}^*(\lambda)|^q$ 仍可使两问题等价。

2. 相同点: 都基于信号或数据在某个基底或字典下具有稀疏表示的假设, 并利用了稀疏表示在提取信号本质特征方面的优势。

不同点: 字典学习旨在从给定的训练数据中学习一个字典, 使得原始数据可以在该字典下获得最稀疏的表示, 而压缩感知关注的是如何利用信号本身所具有的稀疏性, 从部分观测样本中恢复原信号。字典学习需要优化字典, 而压缩感知中通常测量矩阵和稀疏基 (类似于字典) 已知。

二. (40 points) 半监督学习

生成式方法 (generative methods) 是直接基于生成式模型的方法。此类方法假设所有数据 (无论是否有标记) 都是由同一个潜在的模型 “生成” 的。这个假设使得我们能通过潜在模型的参数将未标记数据与学习目标联系起来, 而未标记数据的标记则可看作模型的缺失参数, 通常可基于 EM 算法进行极大似然估计求解。我们接下来探究高斯混合模型的参数估计过程。

给定有标记样本集 $D_l = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_l, y_l)\}$ 和未标记样本集 $D_u = \{\mathbf{x}_{l+1}, \mathbf{x}_{l+2}, \dots, \mathbf{x}_{l+u}\}$, $l \ll u$, $l + u = m$ 。假设所有样本独立同分布, 且都是由同一个高斯混合模型生成的。用极大似然法来估计高斯混合模型的参数 $\{(\alpha_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \mid 1 \leq i \leq N\}$, $D_l \cup D_u$ 的对数似然是:

$$\begin{aligned} LL(D_l \cup D_u) = & \sum_{(\mathbf{x}_j, y_j) \in D_l} \ln \left(\sum_{i=1}^N \alpha_i \cdot p(\mathbf{x}_j \mid \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \cdot p(y_j \mid \Theta = i, \mathbf{x}_j) \right) \\ & + \sum_{\mathbf{x}_j \in D_u} \ln \left(\sum_{i=1}^N \alpha_i \cdot p(\mathbf{x}_j \mid \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \right) \end{aligned} \quad (3)$$

上式由两项组成: 基于有标记数据 D_l 的有监督项和基于未标记数据 D_u 的无监督项。我们将用 EM 算法求解高斯混合模型参数。

1. (10 points) **E 步更新公式**: 根据当前模型参数计算未标记样本 \mathbf{x}_j 属于各高斯混合成分的概率为:

$$\gamma_{ji} = \frac{\alpha_i \cdot p(\mathbf{x}_j \mid \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)}{\sum_{i=1}^N \alpha_i \cdot p(\mathbf{x}_j \mid \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)} \quad (4)$$

请尝试推导上式。

2. (30 points) **M 步更新公式**: 基于 γ_{ji} 更新模型参数, 其中 l_i 表示第 i 类的有标记样本数目为:

$$\boldsymbol{\mu}_i = \frac{1}{\sum_{\mathbf{x}_j \in D_u} \gamma_{ji} + l_i} \left(\sum_{\mathbf{x}_j \in D_u} \gamma_{ji} \mathbf{x}_j + \sum_{(\mathbf{x}_j, y_j) \in D_l \wedge y_j = i} \mathbf{x}_j \right), \quad (5)$$

$$\begin{aligned} \boldsymbol{\Sigma}_i = & \frac{1}{\sum_{\mathbf{x}_j \in D_u} \gamma_{ji} + l_i} \left(\sum_{\mathbf{x}_j \in D_u} \gamma_{ji} (\mathbf{x}_j - \boldsymbol{\mu}_i)(\mathbf{x}_j - \boldsymbol{\mu}_i)^\top \right. \\ & \left. + \sum_{(\mathbf{x}_j, y_j) \in D_l \wedge y_j = i} (\mathbf{x}_j - \boldsymbol{\mu}_i)(\mathbf{x}_j - \boldsymbol{\mu}_i)^\top \right), \end{aligned} \quad (6)$$

$$\alpha_i = \frac{1}{m} \left(\sum_{\mathbf{x}_j \in D_u} \gamma_{ji} + l_i \right). \quad (7)$$

请根据先前给出的对数似然函数, 计算推导出以上 3 个参数的更新公式。

解:

1. 未标记样本 \mathbf{x}_j 属于第 i 个高斯混合成分的概率为

$$\begin{aligned} \gamma_{ji} &= p(\Theta = i \mid \mathbf{x}_j) \\ &= \frac{p(\Theta = i, \mathbf{x}_j)}{p(\mathbf{x}_j)} \\ &= \frac{p(\Theta = i) p(\mathbf{x}_j \mid \Theta = i)}{\sum_{k=1}^N p(\Theta = k) p(\mathbf{x}_j \mid \Theta = k)} \\ &= \frac{\alpha_i \cdot p(\mathbf{x}_j \mid \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)}{\sum_{k=1}^N \alpha_k \cdot p(\mathbf{x}_j \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \end{aligned}$$

2. 注：严格来说求解第三个更新公式时还应考虑混合系数的约束 $\alpha_i \geq 0$ ，部分同学遗漏

$$\begin{aligned}
 \frac{\partial LL(D_l \cup D_u)}{\partial \mu_i} &= \sum_{(\mathbf{x}_j, y_j) \in D_l \wedge y_j = i} \frac{\partial \ln(\alpha_i \cdot p(\mathbf{x}_j | \mu_i, \Sigma_i))}{\partial p(\mathbf{x}_j | \mu_i, \Sigma_i)} \cdot \frac{\partial p(\mathbf{x}_j | \mu_i, \Sigma_i)}{\partial \mu_i} \\
 &\quad + \frac{\partial}{\partial p(\mathbf{x}_j | \mu_i, \Sigma_i)} \left(\sum_{\mathbf{x}_j \in D_u} \ln \left(\sum_{k=1}^N \alpha_k \cdot p(\mathbf{x}_j | \mu_k, \Sigma_k) \right) \right) \cdot \frac{\partial p(\mathbf{x}_j | \mu_i, \Sigma_i)}{\partial \mu_i} \\
 &= \sum_{(\mathbf{x}_j, y_j) \in D_l \wedge y_j = i} \frac{\alpha_i \cdot p(\mathbf{x}_j | \mu_i, \Sigma_i)}{\alpha_i \cdot p(\mathbf{x}_j | \mu_i, \Sigma_i)} \Sigma_i^{-1} (\mathbf{x}_j - \mu_i) \\
 &\quad + \sum_{\mathbf{x}_j \in D_u} \frac{\alpha_i \cdot p(\mathbf{x}_j | \mu_i, \Sigma_i)}{\sum_{k=1}^N \alpha_k \cdot p(\mathbf{x}_j | \mu_k, \Sigma_k)} \Sigma_i^{-1} (\mathbf{x}_j - \mu_i) \\
 &= \sum_{(\mathbf{x}_j, y_j) \in D_l \wedge y_j = i} [\Sigma_i^{-1} (\mathbf{x}_j - \mu_i)] + \sum_{\mathbf{x}_j \in D_u} \gamma_{ji} [\Sigma_i^{-1} (\mathbf{x}_j - \mu_i)] \\
 &= \Sigma_i^{-1} \left[-\mu_i \left(l_i + \sum_{\mathbf{x}_j \in D_u} \gamma_{ji} \right) + \left(\sum_{(\mathbf{x}_j, y_j) \in D_l \wedge y_j = i} \mathbf{x}_j + \sum_{\mathbf{x}_j \in D_u} \gamma_{ji} \mathbf{x}_j \right) \right]
 \end{aligned}$$

令 $\frac{\partial LL(D_l \cup D_u)}{\partial \mu_i} = \mathbf{0}$ 即得第一个更新公式。

$$\begin{aligned}
 \frac{\partial LL(D_l \cup D_u)}{\partial \Sigma_i} &= \sum_{(\mathbf{x}_j, y_j) \in D_l \wedge y_j = i} \frac{\partial \ln(\alpha_i \cdot p(\mathbf{x}_j | \mu_i, \Sigma_i))}{\partial p(\mathbf{x}_j | \mu_i, \Sigma_i)} \cdot \frac{\partial p(\mathbf{x}_j | \mu_i, \Sigma_i)}{\partial \Sigma_i} \\
 &\quad + \frac{\partial}{\partial p(\mathbf{x}_j | \mu_i, \Sigma_i)} \left(\sum_{\mathbf{x}_j \in D_u} \ln \left(\sum_{k=1}^N \alpha_k \cdot p(\mathbf{x}_j | \mu_k, \Sigma_k) \right) \right) \cdot \frac{\partial p(\mathbf{x}_j | \mu_i, \Sigma_i)}{\partial \Sigma_i} \\
 &= \sum_{(\mathbf{x}_j, y_j) \in D_l \wedge y_j = i} \frac{\alpha_i \cdot p(\mathbf{x}_j | \mu_i, \Sigma_i)}{\alpha_i \cdot p(\mathbf{x}_j | \mu_i, \Sigma_i)} \cdot \frac{1}{2} \left[-\Sigma_i^{-1} + \Sigma_i^{-1} (\mathbf{x}_j - \mu_i) (\mathbf{x}_j - \mu_i)^\top \Sigma_i^{-1} \right] \\
 &\quad + \sum_{\mathbf{x}_j \in D_u} \frac{\alpha_i \cdot p(\mathbf{x}_j | \mu_i, \Sigma_i)}{\sum_{k=1}^N \alpha_k \cdot p(\mathbf{x}_j | \mu_k, \Sigma_k)} \cdot \frac{1}{2} \left[-\Sigma_i^{-1} + \Sigma_i^{-1} (\mathbf{x}_j - \mu_i) (\mathbf{x}_j - \mu_i)^\top \Sigma_i^{-1} \right] \\
 &= \sum_{(\mathbf{x}_j, y_j) \in D_l \wedge y_j = i} \frac{1}{2} \left[-\Sigma_i^{-1} + \Sigma_i^{-1} (\mathbf{x}_j - \mu_i) (\mathbf{x}_j - \mu_i)^\top \Sigma_i^{-1} \right] \\
 &\quad + \sum_{\mathbf{x}_j \in D_u} \frac{1}{2} \gamma_{ji} \left[-\Sigma_i^{-1} + \Sigma_i^{-1} (\mathbf{x}_j - \mu_i) (\mathbf{x}_j - \mu_i)^\top \Sigma_i^{-1} \right] \\
 &= \frac{1}{2} \Sigma_i^{-1} \left[-\Sigma_i \left(l_i + \sum_{\mathbf{x}_j \in D_u} \gamma_{ji} \right) + \sum_{(\mathbf{x}_j, y_j) \in D_l \wedge y_j = i} (\mathbf{x}_j - \mu_i) (\mathbf{x}_j - \mu_i)^\top \right. \\
 &\quad \left. + \sum_{\mathbf{x}_j \in D_u} \gamma_{ji} (\mathbf{x}_j - \mu_i) (\mathbf{x}_j - \mu_i)^\top \right] \Sigma_i^{-1}
 \end{aligned}$$

令 $\frac{\partial LL(D_l \cup D_u)}{\partial \Sigma_i} = \mathbf{0}$ 即得第二个更新公式。

考虑到混合系数的约束 $\alpha_i \geq 0$, $\sum_{i=1}^N \alpha_i = 1$, 记拉格朗日函数

$$L = L(\alpha, \nu, \lambda) = LL(D_l \cup D_u) + \sum_{i=1}^N \lambda_i (-\alpha_i) + \nu \left(\sum_{i=1}^N \alpha_i - 1 \right)$$

则有

$$\begin{aligned} \frac{\partial L}{\partial \alpha_i} &= \sum_{(\mathbf{x}_j, y_j) \in D_l \wedge y_j = i} \frac{\partial \ln(\alpha_i \cdot p(\mathbf{x}_j | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i))}{\partial \alpha_i} + \frac{\partial}{\partial \alpha_i} \left(\sum_{\mathbf{x}_j \in D_u} \ln \left(\sum_{k=1}^N \alpha_k \cdot p(\mathbf{x}_j | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) \right) - \lambda_i + \nu \\ &= \sum_{(\mathbf{x}_j, y_j) \in D_l \wedge y_j = i} \frac{1}{\alpha_i} + \sum_{\mathbf{x}_j \in D_u} \frac{p(\mathbf{x}_j | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)}{\sum_{k=1}^N \alpha_k \cdot p(\mathbf{x}_j | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} - \lambda_i + \nu \\ &= \frac{l_i}{\alpha_i} + \sum_{\mathbf{x}_j \in D_u} \frac{\gamma_{ji}}{\alpha_i} - \lambda_i + \nu \end{aligned}$$

考虑到 KKT 条件, 可得以下式子

$$\begin{cases} \frac{\partial L}{\partial \alpha_i} = \frac{l_i}{\alpha_i} + \sum_{\mathbf{x}_j \in D_u} \frac{\gamma_{ji}}{\alpha_i} - \lambda_i + \nu = 0 \\ \sum_{i=1}^N \alpha_i - 1 = 0 \\ \lambda_k \alpha_k = 0 \quad (1 \leq k \leq N, k \in \mathbb{Z}) \end{cases} \quad (*)$$

结合 (*) 中第一和第三个式子, 有

$$l_i + \sum_{\mathbf{x}_j \in D_u} \gamma_{ji} = \alpha_i (\lambda_i - \nu) = -\alpha_i \nu \quad (**)$$

考虑所有的 N 个成分, 并结合 (*) 中第二个式子有

$$-\sum_{i=1}^N \alpha_i \nu = -\nu = \sum_{i=1}^N \left(l_i + \sum_{\mathbf{x}_j \in D_u} \gamma_{ji} \right) = l + \sum_{\mathbf{x}_j \in D_u} \sum_{i=1}^N \frac{\alpha_i \cdot p(\mathbf{x}_j | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)}{\sum_{k=1}^N \alpha_k \cdot p(\mathbf{x}_j | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} = l + u = m$$

将 $\nu = -m$ 代入式 (**) 即得第三个更新公式。

三. (30 points) 方法讨论

- (10 points) LoRA (Low-Rank Adaptation) 是当前常见的模型微调技术之一, 它通过在预训练模型的基础上引入低秩矩阵来调整模型参数, 从而实现模型的微调。请先对 LoRA 方法进行描述, 并讨论 LoRA 有作用的原因 (可以结合教材第 11 章内容进行讨论)。
- (20 points) 给定 $D_l = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_l, y_l)\}$ 和 $D_u = \{\mathbf{x}_{l+1}, \mathbf{x}_{l+2}, \dots, \mathbf{x}_{l+u}\}$, $l \ll u$, 且 $l + u = m$ 。我们可将其映射为一个图, 数据集中每个样本对应于图中一个结点, 若两个样本之间的相似度很高 (或相关性很强), 则对应的结点之间存在一条边, 边的“强度” (strength) 正比于样本之间的相似度 (或相关性)。我们先基于 $D_l \cup D_u$ 构建一个图 $G = (V, E)$, 其中结点集 $V = \{\mathbf{x}_1, \dots, \mathbf{x}_l, \mathbf{x}_{l+1}, \dots, \mathbf{x}_{l+u}\}$, 边集 E 可表示为一个亲和矩阵 (affinity matrix), 常基于高斯函

数定义为

$$(W)_{ij} = \begin{cases} \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right), & \text{if } i \neq j; \\ 0, & \text{otherwise,} \end{cases} \quad (8)$$

其中 $i, j \in \{1, 2, \dots, m\}$, $\sigma > 0$ 是用户指定的高斯函数带宽参数。

在上述情景中, 我们可将有标记样本所对应的结点想象为染过色, 而未标记样本所对应的结点尚未染色, 于是, 半监督学习就对应于“颜色”在图上扩散或传播的过程。该算法亦被称为标记传播方法 (label propagation)。我们接下来仅考虑二分类场景, 希望从图 $G = (V, E)$ 学得一个实值函数 $f: V \rightarrow \mathbb{R}$, 其对应的分类规则为: $y_i = \text{sign}(f(\mathbf{x}_i))$, $y_i \in \{-1, +1\}$, 并定义关于 f 的“能量函数”(energy function):

$$E(f) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (W)_{ij} (f(\mathbf{x}_i) - f(\mathbf{x}_j))^2 \quad (9)$$

请尝试利用上述的条件, 推导出未标记节点的函数值 f_u 的预测公式。你的答案可以写为矩阵乘法的形式。

解:

1. LoRA 中, 预训练的网络权重 $W_0 \in \mathbb{R}^{d \times k}$ 会固定不动, 选择引入可更新参数的矩阵 $B \in \mathbb{R}^{d \times r}$, $A \in \mathbb{R}^{r \times k}$, 其中 $r \ll \min(d, k)$ 。微调时网络更新方式为

$$W_0 + \Delta W = W_0 + BA$$

有效的原因: (1) LoRA 的低秩分解大幅降低待优化参数量, 加速训练。(2) 预训练模型已经拥有了一定的基础能力, 为了更好适应下游任务实际需要更新的参数空间可能是低维的, LoRA 的低秩分解能够选择捕捉重要的特征。与全参数训练相比, 在下游任务上微调性能也不会显著降低。

2. 注: 此处求解过程与教材 13.4 中内容相同, 部分同学采用课外文献求解得到结果, 也可以

$$\begin{aligned} E(f) &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (\mathbf{W})_{ij} (f(\mathbf{x}_i) - f(\mathbf{x}_j))^2 \\ &= \frac{1}{2} \left(\sum_{i=1}^m d_i f^2(\mathbf{x}_i) + \sum_{j=1}^m d_j f^2(\mathbf{x}_j) - 2 \sum_{i=1}^m \sum_{j=1}^m (\mathbf{W})_{ij} f(\mathbf{x}_i) f(\mathbf{x}_j) \right) \\ &= \sum_{i=1}^m d_i f^2(\mathbf{x}_i) - \sum_{i=1}^m \sum_{j=1}^m (\mathbf{W})_{ij} f(\mathbf{x}_i) f(\mathbf{x}_j) \\ &= \mathbf{f}^T (\mathbf{D} - \mathbf{W}) \mathbf{f} \\ &= \begin{pmatrix} \mathbf{f}_l^T & \mathbf{f}_u^T \end{pmatrix} \left(\begin{bmatrix} \mathbf{D}_{ll} & \mathbf{0}_{lu} \\ \mathbf{0}_{ul} & \mathbf{D}_{uu} \end{bmatrix} - \begin{bmatrix} \mathbf{W}_{ll} & \mathbf{W}_{lu} \\ \mathbf{W}_{ul} & \mathbf{W}_{uu} \end{bmatrix} \right) \begin{bmatrix} \mathbf{f}_l \\ \mathbf{f}_u \end{bmatrix} \\ &= \mathbf{f}_l^T (\mathbf{D}_{ll} - \mathbf{W}_{ll}) \mathbf{f}_l - 2 \mathbf{f}_u^T \mathbf{W}_{ul} \mathbf{f}_l + \mathbf{f}_u^T (\mathbf{D}_{uu} - \mathbf{W}_{uu}) \mathbf{f}_u \end{aligned}$$

其中

$$\begin{cases} \mathbf{f} = \left(\mathbf{f}_l^\top \mathbf{f}_u^\top \right)^\top \\ \mathbf{f}_l = (f(\mathbf{x}_1); f(\mathbf{x}_2); \cdots; f(\mathbf{x}_l)) \\ \mathbf{f}_u = (f(\mathbf{x}_{l+1}); f(\mathbf{x}_{l+2}); \cdots; f(\mathbf{x}_{l+u})) \\ \mathbf{D} = \text{diag}(d_1, d_2, \cdots, d_{l+u}) \\ d_i = \sum_{j=1}^{l+u} (\mathbf{W})_{ij} \end{cases}$$

令 $\frac{\partial E(f)}{\partial \mathbf{f}_u} = \mathbf{0}$, 可得

$$\mathbf{f}_u = (\mathbf{D}_{uu} - \mathbf{W}_{uu})^{-1} \mathbf{W}_{ul} \mathbf{f}_l$$