

HW8

PB21111653 李宇哲

T1

Question:

试证明对于不含冲突数据集（即特征向量完全相同但标记不同）的训练集，必存在与训练集一致（即训练误差为 0）的决策树。

Answer:

用递归的方式构建一棵与训练集一致的决策树

首先，选择一个划分属性，将数据集根据该属性划分为若干个子集，然后再每个子集中递归构建决策树，直到所有子集中的数据都属于同一类别或者无法再次选择划分属性为止

由于数据集不含冲突数据，每个子集的数据都属于同一类别，递归构建决策树时，每个叶子结点都属于一个数据的类别，且每个叶子结点的数据集都只包含属于该类别的数据

每个叶子节点对应于一个数据的类别，而且每个非叶子节点对应于一个划分属性，将数据集划分为若干个子集。由于数据集不含冲突数据，因此在每个子集中，所有的数据都属于同一类别。这样的决策树可以将每个特征向量映射到其正确的类别，训练误差为 0。

T2

Question:

2.最小二乘学习方法在求解 $\min_w (Xw - y)^2$ 问题后得到闭式解 $w^* = (X^T X)^{-1} X^T y$ （为简化问题，我们忽略偏差项 b ）。如果我们知道数据中部分特征有较大的误差，在不修改损失函数的情况下，引入规范化项 $\lambda w^T D w$ ，其中 D 为对角矩阵，由我们取值。相应的最小二乘分类学习问题转换为以下形式的优化问题：

$$\min_w (Xw - y)^2 + \lambda w^T D w$$

(1).请说明选择规范化项 $w^T D w$ 而非 L_2 规范化项 $w^T w$ 的理由是什么。 D 的对角线元素 D_{ii} 有何意义，它的取值越大意味着什么？

(2).请对以上问题进行求解。

Answer:

(1)

选择规范化项 $w^T D w$ 的原因是，可以对不同特征的权重进行不同程度的约束，从而避免某些特征对预测结果的贡献过大

取值越大表示该特征的权重受到的约束越强，即该特征对预测结果的贡献更小

(2)

$\min_w (|Xw - y|^2 + \lambda w^T D w)$ 为了求解上述优化问题，我们可以先对 w 进行展开：

$w = [w_1 w_2 \cdots w_d]$ 其中， d 为特征的数量。则有：

$$\begin{aligned}\lambda w^T D w &= \lambda \sum_{i=1}^d w_i^2 D_{ii} \\ &= w^T \begin{bmatrix} \lambda D_{11} & 0 & \cdots & 0 & 0 & \lambda D_{22} & \cdots & 0 & \vdots & \vdots & \ddots & \vdots & 0 & 0 & \cdots & \lambda D_{dd} \end{bmatrix} w \\ &= w^T \Lambda w\end{aligned}$$

其中， Λ 即为第二步中的对角矩阵，其对角线元素为 λD_{ii} 。

此时问题转化为求解：

$$\min_w (|Xw - y|^2 + w^T \Lambda w)$$

对 w 求导，令导数为0，可得： $X^T X w - X^T y + \Lambda w = 0$ $w^* = (X^T X + \Lambda)^{-1} X^T y$

T3

Question:

3. 假设有 n 个数据点 x_1, \dots, x_n 以及一个映射 $\varphi: x \rightarrow \varphi(x)$,以此定义核函数 $K(x, x') = \varphi(x) \cdot \varphi(x')$ 。试证明由该核函数所决定的核矩阵 $K: K_{i,j} = K(x_i, x_j)$ 有以下性质:

(1). K 是一个对称矩阵;

(2). K 是一个半正定矩阵, 即 $\forall z \in R^n, z^T K z \geq 0$ 。

Answer:

(1)

证明:

对于 $\forall i, j \in [n]$, 对于矩阵 K 我们有:

$$\begin{aligned}K_{i,j} &= \varphi(x_i) \cdot \varphi(x_j) \\ &= \varphi(x_j) \cdot \varphi(x_i) \\ &= K_{j,i}\end{aligned}$$

因此 K 是对称矩阵

(2)

将 K 做矩阵展开，并把二次型写为和式

$$\begin{aligned}
z^T K z &= \sum_{i=1}^n \sum_{j=1}^n z_i K_{i,j} z_j \\
&= \sum_{i=1}^n \sum_{j=1}^n z_i \varphi(x_i)^T \varphi(x_j) z_j \\
&= \sum_{i=1}^n z_i \varphi(x_i)^T \sum_{j=1}^n z_j \varphi(x_j) \\
&= \left\| \sum_{i=1}^n z_i \varphi(x_i) \right\|^2 \\
&\geq 0
\end{aligned}$$

T4

Question:

4. 已知正例点 $x_1 = (1, 2)^T, x_2 = (2, 3)^T, x_3 = (3, 3)^T$ ，负例点 $x_4 = (2, 1)^T, x_5 = (3, 2)^T$ ，试求Hard Margin SVM 的最大间隔分离超平面和分类决策函数，并在图上画出分离超平面、间隔边界及支持向量。

Answer:

最大间隔分离超平面也就是以下问题的条件极值问题：

$\min_{w,b} \frac{1}{2} \|w\|^2$ subject to $\forall i, y_i(w^T x_i + b) \geq 1$ 将数据点代入约束条件，得到：

$$y_1(w^T x_1 + b) \geq 1 \quad y_2(w^T x_2 + b) \geq 1 \quad y_3(w^T x_3 + b) \geq 1 \quad y_4(w^T x_4 + b) \geq 1 \quad y_5(w^T x_5 + b) \geq 1$$

显然这个问题是一个二维问题，因此可以将 w 表示为 $w = (w_1, w_2)^T$ 。再将上面的约束转

化为等式约束，令拉格朗日乘子 $\alpha_i \geq 0$ 可以得到：

$L = \frac{1}{2} (w_1^2 + w_2^2) + \sum_{i=1}^5 \alpha_i (1 - y_i(w^T x_i + b))$ 将上述等式代入原优化问题，得到：

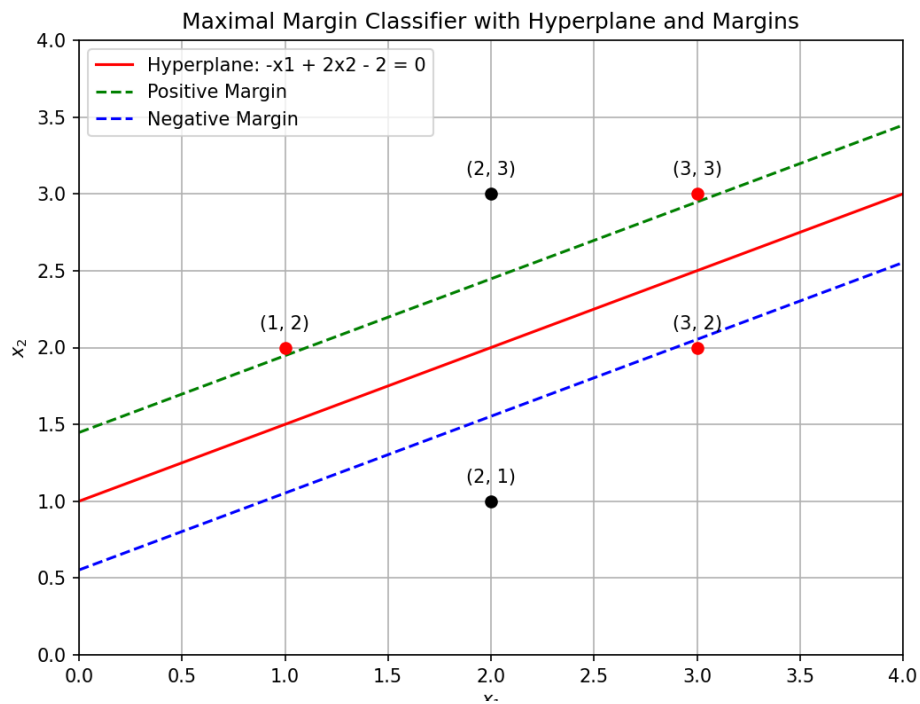
$$\begin{aligned}
\max_{\alpha} \quad & \sum_{i=1}^5 \alpha_i = 1 \\
\text{s.t.} \quad & \sum_{i=1}^5 \alpha_i y_i = 0
\end{aligned}$$

$$\alpha_i \geq 0, \quad i = 1, 2, \dots, 5$$

经过求解，得到 $w_1 = -1, w_2 = 2, b = -2$ ，因此最大间隔分离超平面的解析式为：

$x^{(1)} - 2x^{(2)} + 2 = 0$ 分类决策函数为： $f(x) = \text{sign}(x^{(1)} - 2x^{(2)} + 2)$ 支持向量为：

$x_1 = (1, 2)^T, x_3 = (3, 3)^T, x_5 = (3, 2)^T$ 图像如下



T5

Question:

5. 计算 $\frac{\partial}{\partial w_j} L_{CE}(w, b)$ ，其中

$$L_{CE}(w, b) = -[y \log \sigma(w \cdot x + b) + (1 - y) \log(1 - \sigma(w \cdot x + b))]$$

为Logistic Regression的Loss Function。

已知

$$\begin{aligned} \frac{\partial}{\partial z} \sigma(z) &= \frac{\partial}{\partial z} \frac{1}{1+e^{-z}} = -\left(\frac{1}{1+e^{-z}}\right)^2 \times (-e^{-z}) \\ &= \sigma^2(z) \left(\frac{1-\sigma(z)}{\sigma(z)}\right) = \sigma(z)(1-\sigma(z)) \end{aligned}$$

Answer:

根据链式法则，可以得到： $\frac{\partial}{\partial w_j} L_{CE}(w, b) = \frac{\partial}{\partial \sigma} L_{CE}(w, b) \cdot \frac{\partial \sigma}{\partial z} \cdot \frac{\partial z}{\partial w_j}$ 首先计算第一项： $\frac{\partial L_{CE}(w, b)}{\partial \sigma} = -\left(\frac{y}{\sigma(w \cdot x + b)} - \frac{1-y}{1-\sigma(w \cdot x + b)}\right)$ 然后计算第二项： $\frac{\partial \sigma}{\partial z} = \sigma(z)(1-\sigma(z))$ 然后计算第三项： $\frac{\partial z}{\partial w_j} = \frac{\partial}{\partial w_j}(w \cdot x + b) = x_j$ 代入原式可以得到： $\frac{\partial}{\partial w_j} L_{CE}(w, b) = -\left(\frac{y}{\sigma(w \cdot x + b)} - \frac{1-y}{1-\sigma(w \cdot x + b)}\right) \cdot \sigma(w \cdot x + b)(1-\sigma(w \cdot x + b)) \cdot x_j$ 最后化简一下得到： $\frac{\partial}{\partial w_j} L_{CE}(w, b) = (\sigma(w \cdot x + b) - y) \cdot x_j$

T6

Question:

6. $K - means$ 算法是否一定会收敛? 如果是，给出证明过程; 如果不是，给出说明

Answer:

该算法最后会收敛，证明需要针对算法的两个步骤作证明：

- 算法第一步会找到将所有的点定位到距离其最近的中心点，处理后会让其距离中心点的距离更小，这样更新的距离一定会比这个点上次的距离要短，因此所有点和中心点的距离在这一步必定减小。
- 算法第二步重新确定了每个聚类的中心点，我们需要证明：对于任意聚类 A 和其中的点 a_1, \dots, a_n ，定义其新中心点 $C(A) = \frac{1}{A} \sum_{i=1}^n a_i$ ，则对于任意的点 x ， $\sum_{i=1}^n |a_i - C(A)|^2 \leq \sum_{i=1}^n |a_i - x|^2$ 。

下面来对右边的的式子进行变换：

$$\begin{aligned} \sum_i |a_i - x|^2 &= \sum_i |a_i - C(A) + C(A) - x|^2 \\ &= \sum_i |a_i - C(A)|^2 + \sum_i (a_i - C(A)) \cdot (C(A) - x) + \sum_i |C(A) - x|^2 \\ &= \sum_i |a_i - C(A)|^2 + (C(A) - x) \cdot \sum_i (a_i - C(A)) + |A| |C(A) - x|^2 \\ &= \sum_i |a_i - C(A)|^2 + |A| |C(A) - x|^2 \\ &\geq \sum_i |a_i - C(A)|^2 \end{aligned}$$

倒数第二步用到了 $\sum_i (a_i - C(A)) = 0$ ，定理得证