

# 特征学习与稀疏学习

# 1. 子集搜索与评价

相关特征：对当前的学习任务有用的属性称为“相关特征”

特征选择：从给定特征集合中选择出相关特征子集的过程

冗余特征：在机器学习问题中，某些特征提供的信息与其他特征已经提供的信息高度相似或重复，对模型的性能提升很小或者没有贡献。这样的特征不会为模型带来额外的信息，却可能增加模型的复杂性，降低模型的可解释性，甚至在某些情况下可能会导致过拟合。

**假定：**本章数据中不涉及冗余特征，并且假定初始的特征集合包含了所有重要信息

思路：产生一个**候选子集**

**如何 根据评价结果获取下一个候选特征子集？**

**如何评价候选特征子集的好坏？**

# 1. 子集搜索与评价

- 子集搜索 (subset search)

- 前向搜索

给定特征集合  $\{a_1, a_2, \dots, a_d\}$ , 将每个特征看作一个候选子集, 对这  $d$  个候选单特征子集进行评价, 假定  $\{a_2\}$  最优, 于是将  $\{a_2\}$  作为第一轮选定集; 然后, 在上一轮的选定集中加入一个特征, 构成包含两个特征的候选子集, 假定在这  $d - 1$  个候选两特征子集中  $\{a_2, a_4\}$  最优, 且优于  $\{a_2\}$ , 于是将  $\{a_2, a_4\}$  作为本轮的选定集; ……假定在第  $k + 1$  轮时, 最优的候选  $(k + 1)$  特征子集不如上一轮的选定集, 则停止生成候选子集, 并将上一轮选定的  $k$  特征集合作为特征选择结果.

- 后向搜索

若从完整的特征集合开始, 每次尝试去掉一个无关特征

- 双向搜索

每一轮逐渐增加选定相关特征 (这些特征在后续轮中不会被去除)、同时减少无关特征

# 1. 子集搜索与评价

- 子集评价(subset evaluation)

给定数据集  $D$ , 假定  $D$  中第  $i$  类样本所占的比例为  $p_i$  ( $i = 1, 2, \dots, |\mathcal{Y}|$ ).

对属性子集  $A$ , 假定根据其取值将  $D$  分成了  $V$  个子集  $\{D^1, D^2, \dots, D^V\}$ ,

每个子集中的样本在  $A$  上取值相同, 于是我们可计算属性子集  $A$  的信息增益

其中信息熵定义为

$$\text{Gain}(A) = \text{Ent}(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v),$$

$$\text{Ent}(D) = - \sum_{i=1}^{|\mathcal{Y}|} p_k \log_2 p_k$$

- 信息增益Gain越大, 意味着特征子集A包含的有助于分类的信息越多

常见的特征选择方法大致分为三类:

- 过滤式filter
- 包裹式wrapper
- 嵌入式embedding

# 1. 子集搜索与评价

假设有一个简单的数据集如下:

天气	玩耍
晴天	是
晴天	是
晴天	否
晴天	是
晴天	否
阴天	否
阴天	是
雨天	否

1. 计算数据集的熵  $H(D) = -(\frac{4}{8} \log_2 \frac{4}{8} + \frac{4}{8} \log_2 \frac{4}{8}) = 1$

2. 计算 D 的条件熵  $H(D|\text{天气}) = \sum_{V \in \{\text{晴, 阴, 雨}\}} \frac{|D_V|}{|D|} H(D_V)$

①  $|D_{\text{晴}}| = 3$  且  $H(D_{\text{晴}}) = -(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3}) \approx 0.918$

②  $|D_{\text{阴}}| = 3$  且  $H(D_{\text{阴}}) = -(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3}) \approx 0.918$

③  $|D_{\text{雨}}| = 2$  且  $H(D_{\text{雨}}) = -(0 \log_2 0 + 1 \log_2 1) = 0$

$$H(D|\text{天气}) = (\frac{3}{8} \times 0.918 + \frac{3}{8} \times 0.918 + \frac{2}{8} \times 0) \approx 0.688$$

3. 计算信息增益  $\text{Gain}(D, \text{天气}) = H(D) - H(D|\text{天气}) = 1 - 0.688 = 0.312$

## 2. 过滤式选择

- 过滤式方法先对数据集进行特征选择，然后再训练学习器，**特征选择过程与后续学习器无关**。这相当于先用特征选择过程对初始特征进行“过滤”，再用过滤后的特征来训练模型。换言之，过滤式特征选择是一种在特征选择过程中独立于具体机器学习模型的方法。
- Relief：设计了一个“相关统计量”来度量特征的重要性

给定训练集  $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$

- 猜中近邻 (near-hit)：在 $X_i$ 的同类样本中寻找其最近  $\mathbf{x}_{i,nh}$
- 猜错近邻 (near-miss)：从异类样本中寻找其近邻  $\mathbf{x}_{i,nm}$
- 相关统计量对应于属性 $j$ 的分量为

$$\delta^j = \sum_i -\text{diff}(x_i^j, x_{i,nh}^j)^2 + \text{diff}(x_i^j, x_{i,nm}^j)^2$$

## 2. 过滤式选择

$$\delta^j = \sum_i -\text{diff}(x_i^j, x_{i,\text{nh}}^j)^2 + \text{diff}(x_i^j, x_{i,\text{nm}}^j)^2$$

$x_a^j$  表示样本  $\mathbf{x}_a$  在属性  $j$  上的取值

$\text{diff}(x_a^j, x_b^j)$  取决于属性  $j$  的类型: 若属性  $j$  为离散型, 则  $x_a^j = x_b^j$  时  $\text{diff}(x_a^j, x_b^j) = 0$ , 否则为 1;  
若属性  $j$  为连续型, 则  $\text{diff}(x_a^j, x_b^j) = |x_a^j - x_b^j|$

- 若  $\mathbf{x}_i$  与其猜中近邻  $\mathbf{x}_{i,\text{nh}}$  在属性  $j$  上的距离小于  $\mathbf{x}_i$  与其猜错近邻  $\mathbf{x}_{i,\text{nm}}$  的距离, 则说明属性  $j$  对区分同类与异类样本是有益的, 于是增大属性  $j$  所对应的统计量分量;
- 若  $\mathbf{x}_i$  与其猜中近邻  $\mathbf{x}_{i,\text{nh}}$  在属性  $j$  上的距离大于  $\mathbf{x}_i$  与其猜错近邻  $\mathbf{x}_{i,\text{nm}}$  的距离, 则说明属性  $j$  起负面作用, 于是减小属性  $j$  所对应的统计量分量.

# 3. 包裹式选择

- 包裹式特征选择**直接把最终将要使用的学习器的性能作为特征子集的评价准则**。换言之，包裹式特征选择的目的是为给定学习器选择最有利于其性能，“量身定做”的特征子集
- LVW (Las Vegas Wrapper) .在拉斯维加斯方法的框架下使用随机策略来进行子集搜索

```
输入: 数据集  $D$ ;  
      特征集  $A$ ;  
      学习算法  $\mathcal{L}$ ;  
      停止条件控制参数  $T$ .  
  
过程:  
1:  $E = \infty$ ;  
2:  $d = |A|$ ;  
3:  $A^* = A$ ;  
4:  $t = 0$ ;  
5: while  $t < T$  do  
6:   随机产生特征子集  $A'$ ;  
7:    $d' = |A'|$ ;  
8:    $E' = \text{CrossValidation}(\mathcal{L}(D^{A'}))$ ;  
9:   if  $(E' < E) \vee ((E' = E) \wedge (d' < d))$  then  
10:     $t = 0$ ;  
11:     $E = E'$ ;  
12:     $d = d'$ ;  
13:     $A^* = A'$   
14:   else  
15:     $t = t + 1$   
16:   end if  
17: end while  
  
输出: 特征子集  $A^*$ 
```

在数据集  $D$  上, 使用交叉验证法来估计学习器  $\mathcal{L}$  的误差.



## 4. 嵌入式选择与L1正则化

- 嵌入式特征选择是**将特征选择过程与学习器训练过程融为一体**，即在学习器训练过程中自动进行了特征选择

给定数据集  $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ , 其中  $\mathbf{x} \in \mathbb{R}^d, y \in \mathbb{R}$ .

- 考虑最简单的线性回归模型，以平方误差为损失函数，则优化目标为

$$\min_{\mathbf{w}} \sum_{i=1}^m (y_i - \mathbf{w}^T \mathbf{x}_i)^2 .$$

- 为了缓解过拟合问题，引入正则化项目，若使用L1范数正则化，则

$$\min_{\mathbf{w}} \sum_{i=1}^m (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_1 .$$

## 4. 嵌入式选择与L1正则化

- 采用L1范数, 则

$$\min_{\mathbf{w}} \sum_{i=1}^m (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_1 .$$

L<sub>1</sub> 正则化问题的求解可使用近端梯度下降 (Proximal Gradient Descent, 简称 PGD) [Boyd and Vandenberghe, 2004]. 具体来说, 令  $\nabla$  表示微分算子, 对优化目标

$$\min_{\mathbf{x}} f(\mathbf{x}) + \lambda \|\mathbf{x}\|_1 , \quad (11.8)$$

若  $f(\mathbf{x})$  可导, 且  $\nabla f$  满足  $L$ -Lipschitz 条件, 即存在常数  $L > 0$  使得

$$\|\nabla f(\mathbf{x}') - \nabla f(\mathbf{x})\|_2^2 \leq L \|\mathbf{x}' - \mathbf{x}\|_2^2 \quad (\forall \mathbf{x}, \mathbf{x}') , \quad (11.9)$$

## 4. 嵌入式选择与L1正则化

$$\begin{aligned} & \min_x f(x) + \lambda \|x\|_1 \\ \Rightarrow \text{泰勒: } & f(x) = f(a) + f'(a)(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \dots \\ \text{即: } & f(x) \approx f(x_k) + \nabla f(x_k)(x-x_k) + \frac{L}{2} \|x-x_k\|^2 \\ & = \frac{L}{2} \left\| x - \left( x_k - \frac{1}{L} \nabla f(x_k) \right) \right\|_2^2 + \text{const.} \end{aligned}$$

则在  $x_k$  附近可将  $f(x)$  通过二阶泰勒展式近似为

$$\begin{aligned} \hat{f}(x) & \simeq f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{L}{2} \|x - x_k\|^2 \\ & = \frac{L}{2} \left\| x - \left( x_k - \frac{1}{L} \nabla f(x_k) \right) \right\|_2^2 + \text{const}, \quad (11.10) \end{aligned}$$

于是, 若通过梯度下降法对  $f(x)$  进行最小化, 则每一步梯度下降迭代实际上等价于最小化二次函数  $\hat{f}(x)$ . 将这个思想推广到式(11.8), 则能类似地得到其每一步迭代应为

$$x_{k+1} = \arg \min_x \frac{L}{2} \left\| x - \left( x_k - \frac{1}{L} \nabla f(x_k) \right) \right\|_2^2 + \lambda \|x\|_1, \quad (11.12)$$

## 4. 嵌入式选择与L1正则化

即在每一步对  $f(\mathbf{x})$  进行梯度下降迭代的同时考虑  $L_1$  范数最小化.

对于式(11.12), 可先计算  $\mathbf{z} = \mathbf{x}_k - \frac{1}{L}\nabla f(\mathbf{x}_k)$ , 然后求解

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \frac{L}{2} \|\mathbf{x} - \mathbf{z}\|_2^2 + \lambda \|\mathbf{x}\|_1. \quad (11.13)$$

$$x_{k+1}^i = \arg \min_{x^i} \frac{L}{2} (x^i - z^i)^2 + \lambda |x^i|$$

$$\frac{L}{2} x^2 - Lz x + \frac{L}{2} z^2 + \lambda |x|$$

$$x \geq 0: \frac{L}{2} x^2 + (\lambda - Lz)x + \frac{L}{2} z^2 \quad x^* = \frac{zL - \lambda}{L}$$

$$x < 0: \frac{L}{2} x^2 + (L\lambda - Lz)x + \frac{L}{2} z^2 \quad x^* = \frac{zL + \lambda}{L}$$

令  $x^i$  表示  $\mathbf{x}$  的第  $i$  个分量, 将式(11.13)按分量展开可看出, 其中不存在  $x^i x^j$  ( $i \neq j$ ) 这样的项, 即  $\mathbf{x}$  的各分量互不影响, 于是式(11.13)有闭式解

$$x_{k+1}^i = \begin{cases} z^i - \lambda/L, & \lambda/L < z^i; \\ 0, & |z^i| \leq \lambda/L; \\ z^i + \lambda/L, & z^i < -\lambda/L, \end{cases} \quad (11.14)$$

# 5. 稀疏表示与字典学习

- 考虑另外一种稀疏性：数据集D所对应的矩阵中存在很多零元素，但这些零元素并不是以整列整行的形式存在
- 字典学习（稀疏编码）：为普通稠密表达的样本找到合适的字典，将样本转化为合适的稀疏表达形式，从而使学习任务得以简化，模型的复杂度得以降低。
- 本质
  - 对庞大数据集的降维表示
  - 尝试学习蕴藏在样本背后最质朴的特征

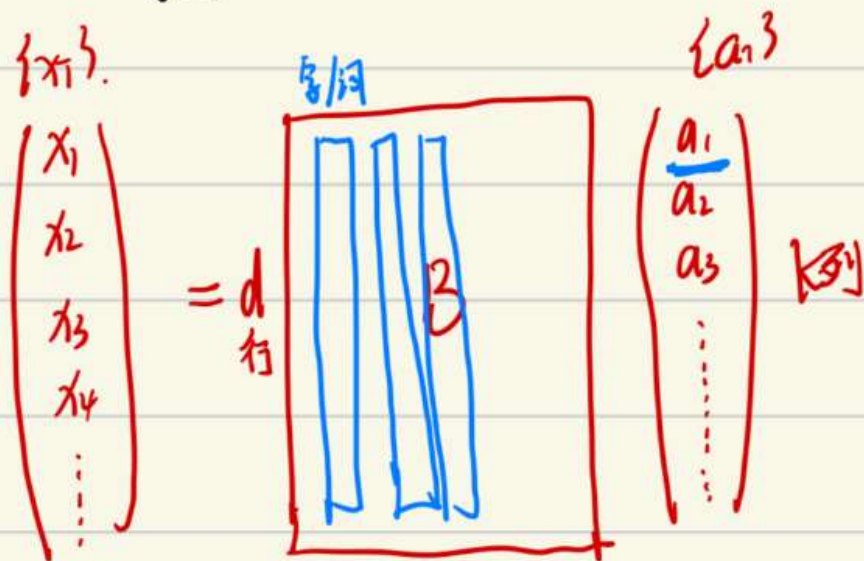
给定数据集  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ ，字典学习最简单的形式为

$$\min_{\mathbf{B}, \alpha_i} \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{B}\alpha_i\|_2^2 + \lambda \sum_{i=1}^m \|\alpha_i\|_1, \quad (11.15)$$

# 5. 稀疏表示与字典学习

给定数据集  $\{x_1, x_2, \dots, x_m\}$ . 字典学习最简形式

$$\min_{B, a_i} \sum_{i=1}^m \|x_i - Ba_i\|_2^2 + \lambda \sum_{i=1}^m \|a_i\|_1$$



$d$ 维  $\xrightarrow{\quad}$   $k$ 维. (字典  $k$ 维空间中每点每个部分都可能非0, 0的值可能更多)

每一列  $\times a_i$  (权重)

## 5. 稀疏表示与字典学习

- 采用变量交替优化的策略来求解（字典B和稀疏矩阵 $\alpha_i$ ）
- 第一步：固定住字典B

$$\min_{\alpha_i} \|\mathbf{x}_i - \mathbf{B}\alpha_i\|_2^2 + \lambda \|\alpha_i\|_1. \quad (11.16)$$

- 第二步：固定住 $\alpha_i$ 来更新字典B

$$\min_{\mathbf{B}} \|\mathbf{X} - \mathbf{B}\mathbf{A}\|_F^2, \quad (11.17)$$

- 基于逐列更新策略的KSVD 令  $\mathbf{b}_i$  表示字典矩阵  $\mathbf{B}$  的第  $i$  列,  $\alpha^i$  表示稀疏矩阵  $\mathbf{A}$  的第  $i$  行, 式(11.17)可重写为

$$\begin{aligned} \min_{\mathbf{B}} \|\mathbf{X} - \mathbf{B}\mathbf{A}\|_F^2 &= \min_{\mathbf{b}_i} \left\| \mathbf{X} - \sum_{j=1}^k \mathbf{b}_j \alpha^j \right\|_F^2 \\ &= \min_{\mathbf{b}_i} \left\| \left( \mathbf{X} - \sum_{j \neq i} \mathbf{b}_j \alpha^j \right) - \mathbf{b}_i \alpha^i \right\|_F^2 \\ &= \min_{\mathbf{b}_i} \|\mathbf{E}_i - \mathbf{b}_i \alpha^i\|_F^2. \end{aligned} \quad (11.18)$$

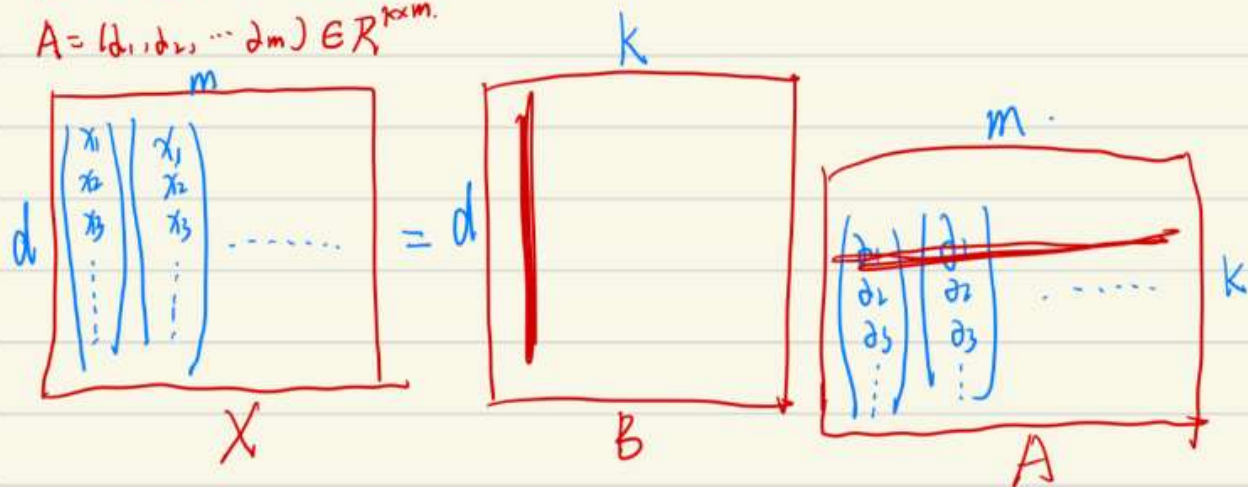
# 5. 稀疏表示与字典学习

求: 1. 给定  $B$ .  $\min_{d_i} \|x_i - B d_i\|_2^2 + \lambda \|d_i\|_1$   
 $d_i \times k$ .  
 向量 = 范数平方和

2. 给定  $X$ .  $\min_B \|X - BA\|_F^2$   
 范数 = 每一列的平方和加起来

$$X = (x_1, x_2, \dots, x_m) \in \mathbb{R}^{d \times m}$$

$$A = (d_1, d_2, \dots, d_m) \in \mathbb{R}^{k \times m}$$



要使得所有元素的平方和最小. 采用逐列更新策略

$B$  的一列为一个词. 每次更新一列

$b_i$  表示  $B$  中的一列.  $d_i$  表示稀疏矩阵  $A$  的一行

$$\text{即: } \min_{b_i} \|X - BA\|_F^2 = \min_{b_i} \left\| X - \sum_{j=1}^k b_j d_j^i \right\|_F^2$$

$d_i \times m$  且秩为 1 的矩阵.  
 $k$  个矩阵相加

$$= \min_{b_i} \left\| (X - \sum_{j \neq i} b_j d_j^i) - b_i d_i^i \right\|_F^2$$

将第  $i$  列单独列出来

$$= \min_{b_i} \left\| E_i - b_i d_i^i \right\|_F^2$$

要更新的一列, 让范数最小.

对  $E_i$  进行奇异值分解. 以取行最大奇异值所对应的正交向量.



# 课后习题

11.1 试编程实现 Relief 算法, 并考察其在西瓜数据集 3.0 上的运行结果.

11.2 试写出 Relief-F 的算法描述.

Relief-F 算法对  $\delta^j$  的计算方式进行了扩展.

Relief-F 算法步骤如下:

- 假定数据集 D 中样本共有  $|Y|$  个类别,
- 其中样本  $x_i$  属于第 k 类: |
- (1) 先在第 k 类样本中寻找最近邻  $x_{i,nh}$ , 称为“猜中近邻”;
- (2) 再在第 k 类样本外寻找最近邻  $x_{i,nm}$ , 称为“猜错近邻”;
- (3) 计算相关统计量对应于属性的分量:

$$\delta^j = \sum_i -\text{diff}(x_i^j, x_{i,nh}^j)^2 + p_l \cdot \text{diff}(x_i^j, x_{i,nm}^j)^2$$

CSDN @将月藏进诗尾

其中,  $p_l$  为第 l 类样本在数据集 D 中所占比例.

- (4) 将基于不同样本得到的估计结果取均值.

11.2.

输入: 样本数据 X, 类别 Y.

过程: 计算样本数 m, 邻居数 n, 类别数  $N = |Y|$ , 各类样本所占比例.

初始  $\delta = \text{zero}(L, N)$

for  $i = 1, 2, \dots, m$

    所属类别 k

    在  $x_i$  中的同类样本找最近邻同类样本  $x_{i,nh}$

    for  $j = 1, 2, \dots, n$

$\delta_j = \delta_j - \text{diff}(x_i^j - x_{i,nh}^j)^2$

    for  $l = 1, 2, 3, \dots, k-1, k+1, \dots, N$

        在属于 l 类别的样本子集中找到与  $x_i$  的最近邻样本  $x_{i,nm}$

        for  $j = 1, 2, \dots, n$ :

$\delta_j = \delta_j + p_l \times \text{diff}(x_i^j - x_{i,nm}^j)^2$

输出: 相关统计量  $\delta$ .