

机器学习-第二讲课后作业

姓名：周宝航 学号：2120190442 专业：计算机科学与技术

一、问题描述

要求：使用 Python 语言编程实现 Find-S 算法和候选消除算法（参见《Machine Learning》一书相关内容），并使用这两种算法归纳给定训练样本的目标概念。最后，对归纳结果进行分析讨论。

难点：深入理解 Find-S 算法和候选消除算法，并使用 Python 语言编程实现它们。

二、基本思路

根据讲义内容和参考书籍，我们先理解 Find-S 算法和候选消除法的流程思路。针对不同的算法，我们对训练样本进行不同的处理方式。比如：Find-S 算法只考虑正例样本，而候选消除算法是分别对正例与负例样本进行处理。所以，在对算法有了深入理解后再去编程实现。最后根据两种算法的运行结果进行深入分析讨论。

三、解题步骤

1. 样本数据分析

在本次实验中，我们考虑进行归纳的样本一共有 14 个，其中关于概念 PlayTennis 为正例的样本有 9 个，而负例一共有 5 个。我们考虑到样本涉及的特征有：Outlook、Temperature、Humidity、Wind。如

下表 1，我们对每个特征涉及的取值进行了统计，以为后面的实验提供基础。

特征	Outlook	Temperature	Humidity	Wind
取值	Rain	Hot	High	True
	Sunny	Cool	Normal	False
	Overcast	Mild		

表 1 数据特征统计

2. 算法描述

2.1 Find-S 算法

输入： 训练样本、假设 h

- 1) 将假设 h 初始化为假设空间 H 的最特殊假设
- 2) 对每个训练样本正例 x

对假设 h 中的每个属性约束 $a_i = v$

如果 x 不满足 h 中的约束 $a_i = v$

则将 $a_i = v$ 替换为 x 满足的更一般约束

输出： 假设 h

2.2 候选消除算法

输入： 训练样本、边界集合 G 、 S

- 1) 将 G 集合初始化为 H 中极大一般假设
- 2) 将 S 集合初始化为 H 中极大特殊假设
- 3) 对每个训练样本 d ，进行下一步操作：
 - a. 如果 d 样本为正例：
 - a) 从 G 中移去所有与 d 不一致的假设

b) 对 S 中每个与 d 不一致的假设 s

从 S 中移去 s

把 s 的所有的极小一般化式 h 加入到 S 中, 其中 h 满足: h
与 d 一致, 而且 G 的某个成员比 h 更一般

从 S 中移去这样的假设: 它比 S 中另一假设更一般

b. 如果 d 样本为负例:

a) 从 S 中移去所有与 d 不一致的假设

b) 对 G 中每个与 d 不一致的假设 g

从 G 中移除 g

把 g 的所有的极小特殊化式 h 加入到 G 中, 其中 h 满足:
h 与 d 一致, 而且 S 的某个成员比 h 更特殊

从 G 中移去所有这样的假设: 它比 G 中另一假设更特殊

输出: 边界集合 G、S

3. 算法实现

3.1 Find-S 算法

```
1. def FindS(data):
2.     # 目标概念
3.     h = None
4.     for i, d in enumerate(data):
5.         # 若为正例样本
6.         if d[-1] == 'Yes':
7.             if h is None:
8.                 h = d[:-1]
9.             else:
10.                for j, v in enumerate(d[:-1]):
11.                    # 若不符合约束, 替换为更宽泛的约束
12.                    if h[j] != v:
13.                        h[j] = '?'
14.     return h
```

在上面我们实现的 Find-S 算法中，我们循环遍历样本。当遇到正例样本时，我们将其与假设 h 进行比较，对于不满足约束的属性，我们替换为更一般的约束，并用 '?' 符号代表一般约束。最后返回所求的假设 h 。

3.2 候选消除算法

我们将该算法封装为一个类，而各个成员函数负责算法的不同步骤。我们将其分开一一列出：

```
1. def collect_attr(self, data):
2.     """
3.     统计样本的所有属性
4.     """
5.     attr = []
6.     data = np.array(data)
7.     n_attr = data.shape[1] - 1
8.     for i in range(n_attr):
9.         attr.append(np.unique(data[:,i]).tolist())
10.    return attr
```

上面的成员函数负责统计样本数据中每个特征所有的可能取值。

```
1. def consistent(self, h, d):
2.     """
3.     判断样本 d 是否符合假设 h
4.     """
5.     flag = True
6.     # 遍历样本和假设的属性
7.     for x, y in zip(h, d):
8.         # 若假设 h 为最特殊假设
9.         if x == 'o':
10.            flag = False
11.            break
12.        # 若样本 d 不满足假设 h
13.        if x != '?' and x != y:
14.            flag = False
15.            break
16.    return flag
```

该成员函数负责判断：一个样本是否符合一个假设。

```

1. def normal(self, g, h):
2.     ....
3.         判断假设 g 是否比 h 更一般
4.     ...
5.     flag = True
6.     # 考虑假设 g 中的一般性
7.     n_g = len([1 for i in g if i == '?'])
8.     # 考虑假设 g 中的最特殊性
9.     o_g = len([1 for i in g if i == 'o'])
10.    # 考虑假设 h 中的一般性
11.    n_h = len([1 for i in h if i == '?'])
12.    # 考虑假设 h 中的最特殊性
13.    o_h = len([1 for i in h if i == 'o'])
14.    # 若 g 更特殊
15.    if o_g > 0 and o_h == 0:
16.        flag = False
17.    # 若 h 更特殊
18.    elif o_g == 0 and o_h > 0:
19.        flag = True
20.    # 若 h 与 g 都为特殊
21.    elif o_g > 0 and o_g == o_h:
22.        flag = True
23.    # 若 g 更一般
24.    elif n_g > n_h:
25.        flag = True
26.    # 若 h 更一般
27.    elif n_g < n_h:
28.        flag = False
29.    elif n_g == n_h:
30.        flag = True
31.        if g == h:
32.            flag = False
33.    return flag

```

该成员函数主要负责判断：一个假设是否比另一个假设更一般。

```

1. def findMinSpecial(self, d):
2.     ....
3.         寻找极小特殊化式加入 G
4.     ...
5.     tmp = []
6.     for g in self.G:
7.         # 若 d 符合假设 g
8.         if not self.consistent(g, d):
9.             tmp.append(g)

```

```

10.         else:
11.             # 若 d 不符合假设 g
12.                 i = 0
13.                 for x, y in zip(d, g):
14.                     if y == '?' or x == y:
15.                         # 根据属性构造极小特殊化式 h
16.                         for a in self.attr[i]:
17.                             if a != x:
18.                                 h = g[:]
19.                                 h[i] = a
20.                                 # 若 h 与 d 一致
21.                                 if not self.consistent(h, d):
22.                                     flag = False
23.                                     # S 的某个成员比 h 更特殊
24.                                     for s in self.S:
25.                                         flag = self.normal(h, s)
26.                                         if flag:
27.                                             break
28.                                     if flag:
29.                                         tmp.append(h)
30.                                     i += 1
31.         G = []
32.         # 遍历添加极小化特殊式到 G
33.         for i, g1 in enumerate(tmp):
34.             flag = True
35.             for j, g2 in enumerate(tmp):
36.                 if i != j:
37.                     # 判断 g 是否比 G 中另一假设更特殊
38.                     flag = self.normal(g1, g2)
39.                     if not flag:
40.                         break
41.             # 若 g 满足条件则不被移去
42.             if flag:
43.                 G.append(g1)
44.         self.G = G

```

该函数主要负责寻找最小特殊化式并加入集合 G。

```

1. def findMinNormal(self, d):
2.     ....
3.     寻找极小一般化式加入 S
4.     ...
5.     tmp = []
6.     for s in self.S:
7.         # 若 d 符合假设 s

```

```

8.         if self.consistent(s, d):
9.             tmp.append(s)
10.        else:
11.            # 若 d 不符合假设 s
12.            if s[0] == 'o':
13.                # 若 S 为最特殊假设
14.                tmp.append(d)
15.            else:
16.                # 根据样本 d 构造极小一般化式 h
17.                h = d[:]
18.                for i, v in enumerate(s):
19.                    # 不符合的属性替换为一般约束'?'
20.                    if v != h[i]:
21.                        h[i] = '?'
22.                # 若 h 与 d 一致
23.                if self.consistent(h, d):
24.                    flag = False
25.                    for g in self.G:
26.                        # 若 G 的某个成员比 h 更一般
27.                        flag = self.normal(g, h)
28.                        if flag:
29.                            break
30.                    if flag:
31.                        tmp.append(h)
32.        S = []
33.        # 遍历添加极小一般化式到 S
34.        for i, s1 in enumerate(tmp):
35.            flag = False
36.            for j, s2 in enumerate(tmp):
37.                if i != j:
38.                    # 判断 s 是否比 S 中另一假设更一般
39.                    flag = self.normal(s1, s2)
40.                    if flag:
41.                        break
42.            # 若 s 满足条件则不被移去
43.            if not flag:
44.                S.append(s2)
45.        self.S = S

```

该函数主要负责寻找最小一般化式并加入集合 S。

```

1. class CandidateElimination(object):
2.     def __init__(self, data):
3.         self.data = data
4.         self.attr = self.collect_attr(data)

```

```

5.         self.S = [['o'] * len(self.attr)]
6.         self.G = [['?'] * len(self.attr)]
7.     def __call__(self):
8.         ....
9.         算法主入口
10.        ...
11.        for d in self.data:
12.            # 特征
13.            x = d[:-1]
14.            # 标签
15.            y = d[-1]
16.            if y == 'Yes':
17.                # 从 G 中移去与 d 不一致的假设
18.                tmp = []
19.                for g in self.G:
20.                    if self.consistent(g, x):
21.                        tmp.append(g)
22.                self.G = tmp
23.                # 对 S 进行一般化
24.                self.findMinNormal(x)
25.            else:
26.                # 从 S 中移去与 d 不一致的假设
27.                tmp = []
28.                for s in self.S:
29.                    if not self.consistent(s, d):
30.                        tmp.append(s)
31.                self.S = tmp
32.                # 对 G 进行特殊化
33.                self.findMinSpecial(x)
34.            print(f"Data: {d}")
35.            print(f"S:{self.S}")
36.            print(f"G:{self.G}")
37.            print(".....")

```

上面是算法类的定义以及算法流程主函数。

四、结果与分析

1. 实验内容与步骤

我们采用 Find-S 算法和候选消除算法对给定训练样本进行训练，

得到对目标概念的学习结果。

2. 实验结果

2.1 Find-S 算法结果

```
Data: ['Overcast', 'Hot', 'High', 'False', 'Yes']
h: ['Overcast', 'Hot', 'High', 'False']
.....
Data: ['Rain', 'Mild', 'High', 'False', 'Yes']
h: ['?', '?', 'High', 'False']
.....
Data: ['Rain', 'Cool', 'Normal', 'False', 'Yes']
h: ['?', '?', '?', 'False']
.....
Data: ['Overcast', 'Cool', 'Normal', 'True', 'Yes']
h: ['?', '?', '?', '?']
.....
Data: ['Sunny', 'Cool', 'Normal', 'False', 'Yes']
h: ['?', '?', '?', '?']
.....
Data: ['Rain', 'Mild', 'Normal', 'False', 'Yes']
h: ['?', '?', '?', '?']
.....
Data: ['Sunny', 'Mild', 'Normal', 'True', 'Yes']
h: ['?', '?', '?', '?']
.....
Data: ['Overcast', 'Mild', 'High', 'True', 'Yes']
h: ['?', '?', '?', '?']
.....
Data: ['Overcast', 'Hot', 'Normal', 'False', 'Yes']
h: ['?', '?', '?', '?']
.....
Find-S算法结果
['?', '?', '?', '?']
```

图 1 Find-S 算法求解过程图

上图展示了 Find-S 算法根据正例样本求解假设 h 的过程。我们可以看到，随着正例样本的数量增多，假设 h 的属性会逐渐变得更加一般。这证明 Find-S 算法是在对目标概念进行学习。但在本次实验中，我们的目标概念 h 最后变成最一般假设。这种情况说明，根据这些训练样本，我们对目标概念的学习是不成功的。

2.1 候选消除算法结果

```
Data: ['Sunny', 'Hot', 'High', 'False', 'No']
S: [['o', 'o', 'o', 'o']]
G: [['Overcast', '?', '?', '?'], ['Rain', '?', '?', '?'], ['?', 'Cool', '?', '?'], ['?', 'Mild', '?', '?'], ['?', '?', 'Normal', '?'], ['?', '?', '?', 'True']]
.....
Data: ['Sunny', 'Hot', 'High', 'True', 'No']
S: [['o', 'o', 'o', 'o']]
G: [['Overcast', '?', '?', '?'], ['Rain', '?', '?', '?'], ['?', 'Cool', '?', '?'], ['?', 'Mild', '?', '?'], ['?', '?', 'Normal', '?'], ['?', '?', '?', 'False']]
.....
Data: ['Overcast', 'Hot', 'High', 'False', 'Yes']
S: [['Overcast', 'Hot', 'High', 'False']]
G: [['Overcast', '?', '?', '?'], ['?', '?', '?', 'False']]
.....
Data: ['Rain', 'Mild', 'High', 'False', 'Yes']
S: [['?', '?', 'High', 'False']]
G: [['?', '?', '?', 'False']]
.....
Data: ['Rain', 'Cool', 'Normal', 'False', 'Yes']
S: []
G: [['?', '?', '?', 'False']]
.....
Data: ['Rain', 'Cool', 'Normal', 'True', 'No']
S: []
G: [['?', '?', '?', 'False']]
.....
Data: ['Overcast', 'Cool', 'Normal', 'True', 'Yes']
S: []
G: []
.....
Data: ['Sunny', 'Mild', 'High', 'False', 'No']
S: []
G: []
.....
Data: ['Sunny', 'Cool', 'Normal', 'False', 'Yes']
S: []
G: []
.....
Data: ['Rain', 'Mild', 'Normal', 'False', 'Yes']
S: []
G: []
.....
Data: ['Sunny', 'Mild', 'Normal', 'True', 'Yes']
S: []
G: []
.....
Data: ['Overcast', 'Mild', 'High', 'True', 'Yes']
S: []
G: []
.....
Data: ['Overcast', 'Hot', 'Normal', 'False', 'Yes']
S: []
G: []
.....
Data: ['Rain', 'Mild', 'High', 'True', 'No']
S: []
G: []
.....
候选消除结果
特殊边界: []
一般边界: []
```

图 2 候选消除算法求解过程图

上图展示了候选消除算法根据训练样本求解边界集合的过程。我们可以看到，起始阶段的负例样本使得一般边界 G 集合变得更加特殊。而遇到正例样本，特殊集合 S 将变得更加一般，且不符合样本的 G 集合中的假设被移去。如此一来， G 集合与 S 集合将相互靠近。最后，两个集合变为空集，说明对目标概念的学习是不成功的。

3. 结果分析

Find-S 算法的结果为: $[?, ?, ?, ?]$ ，而候选消除算法的结果为: G 、 S 集合均为空。Find-S 算法的结果表明: 若要能够满足所有的正例样本，假设 h 的所有属性需要全部变为最一般化。可是这个假设并不能拒绝任一负例样本，所以该算法对目标概念的学习是不成功的。候选消除

算法的结果表明：若要满足所有的训练样本，特殊边界 S 最终无法一般化至某个假设集合，而一般边界 G 最终无法特殊化至某个假设集合。该算法本应将边界集合 G 、 S 分别变化至两个假设集合，且两边界之间的假设可以满足所有训练样本。可是，该算法作用在训练数据上，最终结果没有收敛，说明对目标概念的学习也是不成功的。