

HMM-命名实体识别 算法说明

学号：20373594 姓名：魏少杭

HMM-命名实体识别 算法说明

- 、整体思路和步骤
- 一、学习过程
- 二、推理过程
 - 1.维特比算法说明
- 三、评估
 - 评价结果：
- 四、其他说明

○、整体思路和步骤

第一，学习过程。

在学习过程中，通过./train.json中的所有text中他们的label种类的统计，训练得到各隐状态之间的概率转移矩阵A、每一个label状态下得到各观测值的概率矩阵B、初始概率矩阵（即第一个字符的label类型的分布）pi。

第二，推理过程。

编写Viterbi算法对输入的dev.json观测序列进行解码，得到预测状态序列。

第三，评估。

在步骤二种，得到的状态序列与dev.json的真实标签进行对比，计算出精确率Precision、召回率Recall和F1值。

一、学习过程

定义隐状态概率转移矩阵A，表示i到j状态转移的概率；定义每一个状态下的观测值矩阵B，表示从某状态观测到某字符的概率；定义pi列表为初始的概率分布。为了统计出以上三者，利用似然估计方法，通过频数的比值代替概率值。

```
1 A = {}
2 B = {}
3 pi = {}
4 # 为了能够统计出pi，我们定义count_state，分别对每一个状态进行统计出现次数
5 count_state = {}
6 # 统计各个状态下，出现各观测值的个数
7 count_s_o = {}
8 # 统计从i状态到j状态的个数
9 count_sij = {}
10 # 统计所有的状态
11 count_1_s = []
```

然后对./train.json文件按行解析出text和label两个部分的字典。第一，遍历所有label标签，并用hidden_states 存储一行text每一个字符的隐状态。

```
1 # 下面扫描所有的label
```

```

2         # k: e.g. 'address'
3         for k in label_now.keys():
4             index1 = labelDict[k]      # e.g. address-> 1
5             index2 = 2 * index1-1      # e.g. B-ADDRESS 的index就是address的2
倍-1, 2*1-1=1
6             # obs:观测字典, 如{"汉堡": [[3, 4]], "汉诺威": [[16, 18]]}
7             obs = label_now[k]
8
9             #locations 如: [[3,4], [6,7], [10,11]]
10            for locations in obs.values():
11                list_location_ob_word = locations      #
list_location_ob_word是定位的列表, 可能有多次出现的位置, 如[[3,4],[7,8]]
12                for item in list_location_ob_word:
13                    hidden_states[item[0]] = index2      # B-xxx
14
15                    for i1 in range(item[0]+1, item[1]+1):
16                        hidden_states[i1] = index2 + 1    # I-xxx

```

扫描整个text的 hidden_states, 获得各个频数的更新:

```

1         # 下面进行counter更新
2         for obs_index in range(len(text_now)):
3             # if hidden_states[obs_index] == 0:
4             count_state[hidden_states[obs_index]] += 1      # 各状态出现次数计
数
5
6             # 如果之前没有统计过这个 状态-观测值的出现次数, 现在就令其为 1
7             if text_now[obs_index] not in
count_s_o[hidden_states[obs_index]].keys():
8                 # 状态-观测值概率
9                 count_s_o[hidden_states[obs_index]][text_now[obs_index]] = 1
10            else:      # 否则, 出现次数需要加1
11                count_s_o[hidden_states[obs_index]][text_now[obs_index]] +=
1
12            for obs_index in range(1, len(text_now)):
13                # 状态i到状态j的转移频数
14                count_sij[hidden_states[obs_index-1]][hidden_states[obs_index]]
+= 1

```

用似然估计方法对A、B、pi进行估计。理论依据是:

$$A[i][j] = \frac{P(S_i S_j)}{P(S_i)} = \frac{N(S_i S_j)}{N(S_i)}, \quad B[i][O] = \frac{P(S_i O)}{P(S_i)} = \frac{N(S_i O)}{N(S_i)}, \quad \pi_i = P(S_i | t = 1) = \frac{N(S_i)}{\sum_j N(S_j)}$$

```

1     sum_1_s = sum(count_1_s)
2     for i in range(n):
3         for j in range(n):
4             # 状态转移概率矩阵估计 (最大似然估计)
5             A[i][j] = count_sij[i][j] / count_state[i]
6         for ob in count_s_o[i].keys():
7             # 状态-观测值 估计
8             B[i][ob] = count_s_o[i][ob] / count_state[i]
9         # 初始状态的概率分布 估计
10        pi[i] = count_1_s[i] / sum_1_s

```

二、推理过程

针对./dev.json下的文件中的每一个数据进行处理，通过观测各行字符串，根据之前求取的A, B和pi, 推测出各行字符串的各字符的隐状态label。

1.维特比算法说明

维特比算法首先定义了Viterbi变量 $\delta_t(i)$

定义： Viterbi 变量 $\delta_t(i)$ 是在时间 t 时，模型沿着某一条路径到达 S_i ，输出观察序列 $O=O_1O_2 \dots O_t$ 的最大概率为：

$$\delta_t(i) = \max_{q_1, q_2, \dots, q_{t-1}} p(q_1, q_2, \dots, q_t = S_i, O_1 O_2 \dots O_t | \mu) \quad \dots (6.22)$$

我们只要求得每一个时间t下使得变量 $\delta_t(i)$ 最大的那个状态*i*，就可以得到最可能生成当前text的一组隐藏序列。

递归计算： $\delta_{t+1}(i) = \max_j [\delta_t(j) \cdot a_{ji}] \cdot b_i(O_{t+1}) \quad \dots (6.23)$

● 算法6.3： Viterbi 算法描述

(1) 初始化： $\delta_1(i) = \pi_i b_i(O_1), \quad 1 \leq i \leq N$

概率最大的路径变量： $\psi_1(i) = 0$

(2) 递推计算：

$$\delta_t(j) = \max_{1 \leq i \leq N} [\delta_{t-1}(i) \cdot a_{ij}] \cdot b_j(O_t), \quad 2 \leq t \leq T, \quad 1 \leq j \leq N$$

$$\psi_t(j) = \operatorname{argmax}_{1 \leq i \leq N} [\delta_{t-1}(i) \cdot a_{ij}] \cdot b_j(O_t), \quad 2 \leq t \leq T, \quad 1 \leq j \leq N$$

(3) 结束：

$$\hat{Q}_T = \underset{1 \leq i \leq N}{\operatorname{argmax}} [\delta_T(i)], \quad \hat{p}(\hat{Q}_T) = \max_{1 \leq i \leq N} \delta_T(i)$$

(4) 通过回溯得到路径（状态序列）：

$$\hat{q}_t = \psi_{t+1}(\hat{q}_{t+1}), \quad t = T-1, T-2, \dots, 1$$

算法的时间复杂度： $O(N^2T)$

而需要主要到，在HMM模型中使用动态规划思想求解该问题的前提必须是两个基本假设：

- 1.马尔科夫的，则当前状态的转移概率只与前一步的状态有关
- 2.观测值是独立产生的，也就是只与当前的状态*i*有关，与其他均无关

以下是初始化、动态规划迭代部分：

```
1      # 初始化derta
2      derta = {}
3      phi = {}
4      for t in range(length):
5          derta[t] = {}
6          phi[t] = {}
7          if t == 0:
8              for i in range(n):
9                  # 对于t时刻下的观测值，其是由i状态产生的概率是这样的
10                 if seq[t] in B[i].keys():
11                     derta[t][i] = pi[i] * B[i][seq[t]]
12                 else: # 如果之前从没有在i状态下观测到过B
13                     # derta[t][i] = 0
14                     # if i == 0:
15                     derta[t][i] = pi[i] * (1 / count_state[i]) # 尝试用(1
16                     / count_state[j])代替B[j][seq[t]]
17                     # else:
18                     # derta[t][i] = 0
19                 # 这里实际上直接应用到了HMM的两个基本假设
20                 # 1.马尔科夫的，则当前状态的转移概率只与前一步的状态有关
21                 # 2.观测值是独立产生的，也就是只与当前的状态i有关，与其他均无关
22                 phi[t][i] = 0
23
24     # 更新迭代
25     for t in range(1, length):
26         for j in range(n):
27             max_derta_1 = -100000
28             max_phi = 0
29             for i in range(n): # 遍历上一个时刻，状态i
```

```
29         # Attention! temp没有乘以B[j][O_t]
30         temp = derta[t-1][i] * A[i][j]
31         if max_derta_1 < temp:
32             max_derta_1 = temp
33             max_phi = i
34         if seq[t] in B[j].keys():
35             derta[t][j] = max_derta_1 * B[j][seq[t]]
36         else: # 如果在j状态下从没观察到过当前的观测值
37             # derta[t][j] = 0
38             # if j == 0: # TODO:当完全得到了模型评价之后, 再看如果没出
现过当前的观测值应该怎么办
39             derta[t][j] = max_derta_1 * (1 / count_state[j]) # 尝试用
(1 / count_state[j])代替B[j][seq[t]]
40             # else:
41             #     # derta[t][j] = 0
42             # 更新当前的phi值
43             # derta[t][0] = 1
44             phi[t][j] = max_phi
```

最后, 遍历最终得到的 δ 数组中的最大值, 以及最大值对应的状态q:

```
1 # 结束遍历
2 #首先定义最优序列q
3 q = np.zeros(length)
4 p_t_max = max(derta[length-1].values())
5 for i in range(n):
6     # 这里需要思考, 如果出现了derta相同的情况, 选择最大值如何选取, 我是直接选取
的从前到后出现的第一个最大的
7     if derta[length-1][i] == p_t_max:
8         q[length-1] = i
9         break # 这里我就直接让从前到后遍历到的第一个最大值为q
10
11 # 回溯以得到路径
12 for t in range(length-2, -1, -1):
13     q[t] = phi[t+1][q[t+1]]
```

最终得到了一行字符串中的最有可能的隐藏状态序列 (在q列表中)

三、评估

思路: 首先, 遍历统计所有的预测标签和真实标签。分别用TP FP FN TN存储在每一个标签维度下的预测/真实|正确/错误的个数列表,index就是各类标签

对每一个label进行考察, 设置Positive表示“当前label”, Negative表示“其他label”, 所以可以得到如下的表

混淆矩阵		真实值	
		P (是当前label)	N (不是当前label)
预测值	P' (是当前label)	TP	NP
	N' (不是当前label)	FN	TN

```

1 # 定义和初始化评估参数
2 TP = []
3 FP = []
4 # TN = []
5 FN = []
6
7 for i in range(21):
8     TP.append(0)
9     FP.append(0)
10    # TN.append(0)
11    FN.append(0)

```

针对每一行的数据：

```

1 # 在dev.json每一行根据当前的句子情况，更新TP，FP，FN
2 for i in range(length):
3     hid_state = hidden_states_dev[i]
4     pred_state = pred_states_dev[i]
5     if hid_state == pred_state:
6         TP[hid_state] += 1
7     else:
8         FP[int(pred_state)] += 1
9         FN[hid_state] += 1
10    # 不必计算TN

```

有了TP,FP,FN后计算precision, recall和F1 score

```

1 # 下面进行一些precision、recall和f1-score的计算
2 TP = np.array(TP)
3 FP = np.array(FP)
4 FN = np.array(FN)
5 precision = np.zeros(n)
6 recall = np.zeros(n)
7 F1 = np.zeros(n)
8
9 precision = TP / (TP + FP)
10 recall = TP / (TP + FN)
11 F1 = 2 * precision * recall / (precision + recall)

```

评价结果：

见 `output版本1.txt`，结果如下

Label	Precision	Recall	F1-score
O	0.961471	0.870602	0.913783
B-ADDRESS	0.510836	0.442359	0.474138
I-ADDRESS	0.513346	0.622272	0.562585
B-BOOK	0.505952	0.551948	0.527950
I-BOOK	0.468317	0.539339	0.501325
B-COM	0.611702	0.608466	0.610080
I-COM	0.559494	0.672243	0.610708
B-GAME	0.626087	0.732203	0.675000
I-GAME	0.624623	0.759912	0.685658
B-GOV	0.433702	0.635628	0.515599
I-GOV	0.469476	0.813670	0.595409

13	B-MOVIE	0.507937	0.635762	0.564706
14	I-MOVIE	0.471240	0.762332	0.582441
15	B-NAME	0.685466	0.679570	0.682505
16	I-NAME	0.526074	0.671890	0.590108
17	B-ORG	0.561680	0.583106	0.572193
18	I-ORG	0.500000	0.528979	0.514081
19	B-POS	0.538899	0.655889	0.591667
20	I-POS	0.538618	0.690104	0.605023
21	B-SCENE	0.360153	0.449761	0.400000
22	I-SCENE	0.432702	0.601108	0.503188
23				

四、其他说明

在“三、评估”部分所给出的结果是我认为较为合理的其中一种方案。之所以称之为较为合理，是因为我观察到训练数据的时候，我们的状态-观测值矩阵中，可能没有包含dev.json等除了训练数据之外的可能性。所以，在进行动态规划算法实现的时候需要合理地处理这一部分新的状态-观测值的矩阵。

在以上“二、推理过程”用维特比算法进行初始化和迭代的代码部分

```

1      # 初始化derta
2      derta = {}
3      phi = {}
4      for t in range(length):
5          derta[t] = {}
6          phi[t] = {}
7          if t == 0:
8              for i in range(n):
9                  # 对于t时刻下的观测值，其是由i状态产生的概率是这样的
10                 if seq[t] in B[i].keys():
11                     derta[t][i] = pi[i] * B[i][seq[t]]
12                 else: # 如果之前从没有在i状态下观测到过B
13                     # derta[t][i] = 0
14                     # if i == 0:
15                     derta[t][i] = pi[i] * (1 / count_state[i]) # 尝试用(1
/ count_state[j])代替B[j][seq[t]]
16                     # else:
17                     # derta[t][i] = 0
18                 # 这里实际上直接应用到了HMM的两个基本假设
19                 # 1.马尔科夫的，则当前状态的转移概率只与前一步的状态有关
20                 # 2.观测值是独立产生的，也就是只与当前的状态i有关，与其他均无关
21                 phi[t][i] = 0
22
23     # 更新迭代
24     for t in range(1, length):
25         for j in range(n):
26             max_derta_1 = -100000
27             max_phi = 0
28             for i in range(n): # 遍历上一个时刻，状态i
29                 # Attention! temp没有乘以B[j][O_t]
30                 temp = derta[t-1][i] * A[i][j]
31                 if max_derta_1 < temp:
32                     max_derta_1 = temp
33                     max_phi = i
34                 if seq[t] in B[j].keys():
35                     derta[t][j] = max_derta_1 * B[j][seq[t]]
36                 else: # 如果在j状态下从没观察到过当前的观测值

```

```

37         # derta[t][j] = 0
38         # if j == 0:      # TODO:当完全得到了模型评价之后, 再看如果没出
    现过当前的观测值应该怎么处理
39         derta[t][j] = max_derta_1 * (1 / count_state[j]) # 尝试用
    (1 / count_state[j])代替B[j][seq[t]]
40         # else:
41         #     derta[t][j] = 0
42         # 更新当前的phi值
43         # derta[t][0] = 1
44         phi[t][j] = max_phi

```

这段代码中10-15行、34-39行代码就是在针对从未在某个状态下出现过当前观测值的情况, 对数据进行合理处理的部分。尤其是在第15、39行代码, 在这种特殊情况下计算 $\delta(S_i, O_t)$ 的时候, 由于 $B(S_i, O_t) = 0$, 正常情况下 $\delta(S_i, O_t)$ 也应当为0, 但是这可能会导致由于训练数据的不完备性导致该状态 S_i 本身合理但被直接忽略掉, 进而导致验证/应用效果降低。因此, 我用 $\frac{1}{N(\text{隐藏状态 } S_i \text{ 在训练数据中的总数})}$ 来代替 $B(S_i, O_t)$, 这个比值能够约等于 $\frac{N(S_i \text{ 导致了当前新出现的训练数据})}{N(\text{隐藏状态 } S_i \text{ 在训练数据中的总数})}$, 所以这个处理方法也具有数值上的可行性。

另外一种处理方式, 较为简单粗暴, 即, 如果当前状态 S_i 下对当前观测值没有发射矩阵的值与之对应, 则先判断 S_i 是否为其他类型 O , 如果是, 则直接 $\frac{1}{N(\text{隐藏状态 } S_i \text{ 在训练数据中的总数})}$ 来代替 $B(S_i, O_t)$; 如果不是 O 状态, 则令该状态下的 $\delta(S_i, O_t)$ 为0, 在遍历过程中, 后续的从该时间 t 以状态 S_i 开始的概率变为0。

具体实现方式

```

1     # 初始化derta
2     derta = {}
3     phi = {}
4     for t in range(length):
5         derta[t] = {}
6         phi[t] = {}
7         if t == 0:
8             for i in range(n):
9                 # 对于t时刻下的观测值, 其是由i状态产生的概率是这样的
10                if seq[t] in B[i].keys():
11                    derta[t][i] = pi[i] * B[i][seq[t]]
12                else: # 如果之前从没有在i状态下观测到过B
13                    # derta[t][i] = 0
14                    if i == 0:
15                        derta[t][i] = pi[i] * (1 / count_state[i]) # 尝
    试用(1 / count_state[j])代替B[j][seq[t]]
16                    else:
17                        derta[t][i] = 0
18                # 这里实际上直接应用到了HMM的两个基本假设
19                # 1. 马尔科夫的, 则当前状态的转移概率只与前一步的状态有关
20                # 2. 观测值是独立产生的, 也就是只与当前的状态i有关, 与其他均无关
21                phi[t][i] = 0
22
23     # 更新迭代
24     for t in range(1, length):
25         for j in range(n):
26             max_derta_1 = -100000
27             max_phi = 0
28             for i in range(n): # 遍历上一个时刻, 状态i
29                 # Attention! temp没有乘以B[j][O_t]
30                 temp = derta[t-1][i] * A[i][j]

```



```

31         if max_derta_1 < temp:
32             max_derta_1 = temp
33             max_phi = i
34         if seq[t] in B[j].keys():
35             derta[t][j] = max_derta_1 * B[j][seq[t]]
36         else:    # 如果在j状态下从没观察到过当前的观测值
37                 # derta[t][j] = 0
38                 if j == 0:    # TODO:当完全得到了模型评价之后，再来看如果没出现
过当前的观测值应该怎么处理
39                     derta[t][j] = max_derta_1 * (1 / count_state[j]) #
尝试用(1 / count_state[j])代替B[j][seq[t]]
40                 else:
41                     derta[t][j] = 0
42                 # 更新当前的phi值
43                 # derta[t][0] = 1
44                 phi[t][j] = max_phi

```

在这个情况下，可以得到模型评价结果如下，同时可见 output版本2.txt 文件

1	Label	Precision	Recall	F1-score
2	O	0.943051	0.917054	0.929871
3	B-ADDRESS	0.511905	0.461126	0.485190
4	I-ADDRESS	0.580622	0.604214	0.592183
5	B-BOOK	0.618644	0.474026	0.536765
6	I-BOOK	0.622977	0.438997	0.515050
7	B-COM	0.629333	0.624339	0.626826
8	I-COM	0.574553	0.659316	0.614023
9	B-GAME	0.649860	0.786441	0.711656
10	I-GAME	0.678295	0.770925	0.721649
11	B-GOV	0.492163	0.635628	0.554770
12	I-GOV	0.567531	0.779026	0.656669
13	B-MOVIE	0.582822	0.629139	0.605096
14	I-MOVIE	0.637014	0.698430	0.666310
15	B-NAME	0.684989	0.696774	0.690832
16	I-NAME	0.546851	0.697356	0.613000
17	B-ORG	0.611260	0.621253	0.616216
18	I-ORG	0.558681	0.529899	0.543909
19	B-POS	0.595833	0.660508	0.626506
20	I-POS	0.611628	0.684896	0.646192
21	B-SCENE	0.480000	0.459330	0.469438
22	I-SCENE	0.602826	0.531856	0.565121

这个细节的两个处理方法评价结果差异并不大（对比output版本1.txt（即第三部分）和output版本2.txt（第二种处理方式）），原因可能是O状态占了大多数情况，而其他状态均较少。当我们认为在验证数据中的某个状态下出现新的某个字符，将这个字符出现的原因只归因为“其他标签”也是十分合理的。这与将该字符出现归因到所有标签状态下均分概率，差别并不大。