目录

1	课程总结2	
	1.1 课程内容 1.2 课程形式 1.3 课程收获	
2	说明3 -	
3 :	实验一:实现 IBM Model 1 3 -	
	3.1 实验要求3.2 IBM Model 1 简介3.3 具体实现3.4 结果展示	
4	实验二:基于短语的模型	
附	录 A: 实验 1 部分代码	
附	录 B: 实验 2 部分代码 10 -	

1 课程总结

1.1 课程内容

课程主要讲解了以下内容:

- 1、 机器翻译的几个层次;
- 2、 基于单词的模型;
- 3、 基于短语的模型;
- 4、解码过程。

1.2 课程形式

课程采取"阅读教材+老师提问"的形式。个人认为这种课程形式对于小班授课来说,真真是极好的。概括的来说,有两个优点:

- 1、 课程的节奏非常慢,能够使学生对知识的掌握非常细致。同时,老师可以很好的根据学生的掌握,调整上课的速度。
- 2、 这种课程形式其实极大地激发了学生的主动性,有助于提高学生的自 学能力。

我个人是非常喜欢这种课程形式的。大部分的课程主要采取的是,单纯的"教师讲授"这种方法。我每每遇到这种课程,往往昏昏欲睡,没有上课的劲头;待到考试将近,才着急地开始"真正的"学习。

1.3 课程收获

学习本课程,我有以下的收获:

- 1、 EM **算法**。很多算法虽学习过,但如不结合具体的例子,我还是懵懵懂懂、一脑袋浆糊。此番花了很多时间在 EM 算法上,我不仅对它的原理和实施过程有了进一步的了解,同时我相信以后我学习这类算法也会比以前更加轻松。
- 2、 阅读英文书籍。以前,阅读英文书籍给我的印象是,阴森恐怖、困难 重重。如今读了这一遭,我发现英语书籍也不是那么的令人读不下去。 语言障碍并没有想象的大。只是需要花些时间和耐心。
- 3、 **锻炼代码能力**。说实话,自从大一大二以后,代码写的非常少;而且非常多的一部分是在调 API。实验虽未完全实现,但是我的代码能力还是有所提高。
- 4、 **入门机器翻译**。虽不打算从事机器翻译这个方向,但是作为一个自然语言方向的学生,我觉得还是有必要对这个领域有所了解的。

2 说明

2.1 符号说明

为了以下内容的严谨性,我认为有必要对一些符号进行说明。

表 1 符号说明

符号	意义
\mathbf{e},\mathbf{f}	句子
e, f	单词
$rac{e,f}{\overline{e},\overline{f}}$	短语
a, A	对齐
$t(e \mid f), \phi(\overline{e} \mid \overline{f})$	翻译概率

2.2 其他说明

详细代码见我的 github 仓库。1

3 实验一: 实现 IBM Model 1

3.1 实验要求

实现 IBM Model 1,并在给出的双语语料库 fbis.en.10k 和fbis.zh.10k 上进行实验,估计单词的翻译概率。

3.2 IBM Model 1 简介

给定一系列双语语句对,IBM Model 1 将句子之间的对应分解为单词之间的对应。对于具体的句对,单词之间的对应通过对齐 a 进行表示。如对于双语句对 $< e_1e_2e_3, f_1f_2f_3>$,一个可行的对齐为

$$a: \{1 \to 1, 2 \to 2, 3 \to 3\}.$$
 (1)

而对于所有双语句对,单词之间的对应则体现在英语单词e 和外语单词 f 的翻译概率 $t(e \mid f)$ 上。

对齐a是针对具体的某个句子而言的,而单词间翻译概率则是全局的。如果已知二者中的一个,则可以对另一个进行估计。

¹ https://github.com/1140310118/2018-course-Machine-Translation/

3.2.1 由对齐到翻译概率

假设有大量的双语句对及它们的对齐 a_i ,我们可以使用计数的方法对 $t(e\mid f)$ 进行估计。如给定如下双语句对及对齐

表 2 双语句对及它们的对齐

双语句对 <e,f></e,f>	对齐 a
<a pen,="" 一根="" 钢笔="">	$a1: \{1 \to 1, 2 \to 2\}$
<a book,="" 一本="" 书="">	$a2: \{1 \to 1, 2 \to 2\}$

那么可以得到翻译概率:

表 3 翻译概率 t(elf)

е	f	t(e f)
a	一根	0.5
a	一本	0.5
pen	钢笔	1.0
book	书	1.0

3.2.2 由翻译概率到对齐

由翻译概率到对齐则复杂些。对于具体的某个句对,假设已知单词之间的翻译概率 $t(e \mid f)$,为了选择一个最佳的对齐,IBM Model 1 将对齐出现的概率定义为

$$p(a \mid \mathbf{e}, \mathbf{f}) = \prod_{j} \frac{t(e_j \mid f_{a(j)})}{\sum_{i} t(e_j \mid f_i)}.$$
 (2)

公式(2)的推导过程在这里不再赘述。此公式表明, 句子 \mathbf{e} 中单词 \mathbf{e}_i 应该与使得翻译概率 $t(\mathbf{e}_i \mid f_i)$ 最大的单词 f_i 对齐。

3.2.3 IBM Model 1 的训练

如果已知对齐a和翻译概率 $t(e \mid f)$ 二者中的一个,我们就能推得另一个。但不幸的是,二者都是未知的。

为了解决这个问题, IBM Model 1 采取了 EM 算法:

- 1. 使用均匀分布,对 $t(e \mid f)$ 进行初始化;
- 2. 使用 $t(e \mid f)$ 去估计对齐 a;
- 3. 使用对齐 a 更新 $t(e \mid f)$;
- 4. 迭代步骤 2、3 直至收敛。

3.3 具体实现

3.3.1 算法伪代码

```
Input: set of sentence pairs (e, f)
                                                                    // collect counts
Output: translation prob. t(e|f)
                                                        15:
                                                                   for all words e in e do
                                                                    for all words f in f do
 1: initialize t(e|f) uniformly
                                                                         \begin{array}{lll} \operatorname{count}(e \,|\, f) & += & \frac{t(e|f)}{s - \operatorname{total}(e)} \\ \operatorname{total}(f) & += & \frac{t(e|f)}{s - \operatorname{total}(e)} \end{array}
2: while not converged do
                                                        17:
                                                        18:
      // initialize
     count(e|f) = 0 for all e, f 19:

total(f) = 0 for all f 20:

for all sentence pairs (e,f) do 21:
                                                                       end for
                                                                    end for
                                                                 end for
           // compute normalization
                                                        22:
                                                                 // estimate probabilities
8:
          for all words e in e do
                                                         23:
                                                                 for all foreign words f do
                                                                  for all English words e do
t(e|f) = \frac{\text{count}(e|f)}{\text{total}(f)}
9:
             s-total(e) = 0
                                                         24:
10:
             for all words f in f do
                                                         25:
11:
                s-total(e) += t(e|f)
                                                                    end for
12:
           end for
                                                         27:
                                                                 end for
        end for
                                                         28: end while
13:
```

图 1 使用 EM 算法训练 IBM Model 1 的伪代码

3.3.3 实现细节

(1) 使用数字代替单词²

为了减少程序的内存占用,采取的策略是建立单词与数字之间的一一映射,训练时直接使用数字代替单词。

(2) 初始化

假设单词 e 可能的翻译选项为 $\{f_1, f_2, \dots, f_n\}$, 那么 $t(e \mid f_i) = 1 / n$.

我的另一想法是,不进行平均初始化,而是根据 f_i 与 e 的共现频率以及 f_i 本身出现的频率对 $t(e \mid f_i)$ 进行初始化。即采取类似 TF_IDF 的策略进行初始化。这种初始化策略有可能加快后续的收敛。但实际代码中未采取此策略。

(3) 加速训练

每次迭代结束,淘汰掉那些翻译概率接近0的单词对< e, f>。实验发现,采取此策略,每次迭代的时间将逐渐减少。

(4) 终止条件

最初程序的终止条件是设定固定的迭代次数。但这是不科学的。

改进后的终止条件如下。由于每次迭代结束后会淘汰一部分单词对,因此单词对的总数将慢慢减少。这里设定的终止条件是: 当单词对减少的比例小于一个阈值,则迭代终止。

² 此策略是从课程中的某个学长那学习到的。

3.4 结果展示

序号	е	f	t(e f)
1	introduction	引言	1.000000
2	other	其他	0.997382
3	advanced	先进	0.987018
4	1997	1997	0.979666
5	wuhan	武汉	0.977492
6	or	或	0.977167
7	1998	1998	0.972462
8	military	军事	0.962291
9	authorities	当局	0.952336
10	1999	1999	0.948960

图 2 实验一的部分结果

4 实验二: 基于短语的模型

4.1 实验要求

实现基于短语的模型,输出的短语对及其概率。

4.2 基于短语的模型简介

IBM Model 属于基于单词的模型。与之相比,基于短语的模型有如下优点:

- 1. 由于单词往往存在着一对多的情况,因此单词不是最佳的翻译单元。
- 2. 基于短语进行翻译,有利于解决二义性的问题。
- 3. 如果语料充足,我们可以学习到非常长的短语,甚至可以记住整个句子的翻译。
- 4. 从概念上来看,基于短语的模型更简单。

训练基于短语的模型的前提,我们通过 IBM Model 1 或者其他模型得到了单词间的翻译概率。接下来,通过翻译概率获得双语句对间的对齐A。然后,通过对齐A发现对应的短语对。接着,简单地通过计数的方式计算短语对的翻译概率。

4.2.1 由对齐 a 到对齐 A

对齐实际上就是句对中单词间的映射。在 IBM Model 1 中,一个句对中的对齐 a 是一个单射,而在基于短语的模型中对齐 A 则不是。对齐 A 允许 $i \to j_1, i \to j_2$ 或 $i_1 \to j, i_2 \to j$ 同时在 A 中出现;但不允许 $i_1 \to j_1, i_1 \to j_2, j_2 \to j_1$ 同时在 A 中出现。下图是一个对齐 A 的一个例子。

我们通过调换双语句对的顺序,运行 IBM Model 1 两次,获得两

个对齐 a_{e2f} 和 a_{f2e} ,然后通过算法 GROW-DIAG 获得 A 。具体算法 这里不进行赘述。

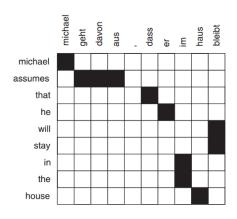


图 3 对齐 A 举例

4.2.2 从对齐 A 中学习短语对

对于对齐 A 而言,如果短语对 $< \overline{e}, \overline{f} >$ 满足如下条件,那么短语 对 $< \overline{e}, \overline{f} >$ 是一致的:

- 1. $\forall ei \in \overline{e} : (ei, fj) \in A \Rightarrow fj \in \overline{f} ;$
- 2. $\forall f_j \in \overline{f} : (e_i, f_j) \in A \Rightarrow e_i \in \overline{e}$;
- 3. $\exists e_i \in \overline{e}, f_j \in \overline{f} : (e_i, f_j) \in A$.

4.2.3 计算短语间的翻译概率

使用列表 P 保存双语句对中所有一致的短语对, $\operatorname{count}(\overline{e},\overline{f})$ 则 表示短语对 $<\overline{e},\overline{f}>$ 在列表 P 中出现的次数。于是,短语间的翻译概率为

$$\phi(\overline{f} \mid \overline{e}) = \frac{count(\overline{e}, \overline{f})}{\sum_{\overline{f}_i} count(\overline{e}, \overline{f}_i)}$$
(3)

4.3 具体实现

图 4 探索法对齐的伪代码

```
\textbf{Input:} \ \text{word alignment A for sentence pair } (\textbf{e},\textbf{f})
Output: set of phrase pairs BP
 1: for e<sub>start</sub> = 1 ... length(e) do

2: for e<sub>end</sub> = e<sub>start</sub> ... length(e) do

3: // find the minimally matching foreign phrase
                (f_{\text{start}}, f_{\text{end}}) = (\text{length}(\mathbf{f}), 0)
for all (e, f) \in A do
                   if e_{\text{start}} \le e \le e_{\text{end}} then

f_{\text{start}} = \min(f, f_{\text{start}})

f_{\text{end}} = \max(f, f_{\text{end}})

end if
10:
                  end for
11:
                 add \operatorname{extract}(f_{\operatorname{start}}, f_{\operatorname{end}}, e_{\operatorname{start}}, e_{\operatorname{end}}) to set BP
           end for
13: end for
function extract(f_{start}, f_{end}, e_{start}, e_{end})

1: return {} if f_{end} == 0 // check if at least one alignment point

2: // check if alignment points violate consistency

3: for all (e,f) \in A do

4: return {} if e < e_{start} or e > e_{end}
  5: end for
 6: // add pharse pairs (incl. additional unaligned f) 7: E = \{\}
9: repeat

10: f_e = f_{end}

11: repeat
12:
              add phrase pair (e_{	ext{start}} .. e_{	ext{end}}, f_s .. f_e) to set E
13:
15: f_S - -
16: until f_S aligned
17: return E
```

图 5 一致短语对抽取的伪代码

4.4 结果展示

可以看出,实现了一致短语对的抽取,但翻译概率的估计仍不准确。可能的原因有 3 个: (a) 数据量太少; (b) 翻译概率的计算方法不恰当; (c) 代码中存在错误。

表 4 实验 2 的部分结果

序号	\overline{e}	\overline{f}	$\phi(\overline{f} \mid \overline{e})$
1	zhongxing 189 dual-band mobile	自主 知识 产权 的 中兴	0.1
	telephone characterized by	一八九双频 手机	
2	zhongxing 189 dual-band mobile	自主 知识 产权 的 中兴	0.1
	telephone characterized by	一八九双频	
3	zhongxing 189 dual-band mobile	知识 产权 的 中兴 一八	0.1
	telephone characterized by	九双频 手机 仔细	
4	zhongxing 189 dual-band mobile	知识 产权 的 中兴 一八	0.1
	telephone characterized by	九双频 手机	

附录 A: 实验 1 部分代码

1、单词转数字 word2num (sentences)

```
1
    def word2num(sentences):
 2
         sentences new = []
 3
         n2w = {} {}
         w2n = \{\}
 4
 5
         n = 0
         for sentence in sentences:
 6
 7
              sent new = []
 8
              for w in sentence.split():
 9
                   if w in _w2n:
10
                        sent_new.append(_w2n[w])
11
                   else:
12
                        sent_new.append(n)
13
                        _w2n[w] = n
14
                        n2w[n] = w
                        n += 1
15
16
              sentences new.append(sent new)
17
         return sentences_new,n2w
```

2、初始化 init t(E,F)

```
1
    def init_t(E,F):
2
         option = defaultdict(set)
3
         t = defaultdict(int)
 4
         for p in range(len(E)):
 5
              for e in E:
                   for f in F:
 6
 7
                       t[(e,f)] += 1
8
                       option[f].add(e)
 9
         for f in option:
10
              option of f = len(option[f])
11
12
              for e in option[f]:
13
                  t[(e,f)] = 1/option_of_f
14
         return t
```

3、一次迭代 one_loop_for_train(E,F,t)

```
def one_loop_for_train(E,F,t,zero=0.001):
 1
 2
         count = defaultdict(int)
 3
         total = defaultdict(int)
         s_total = defaultdict(int)
 4
 5
         for p in range(len(E)):
 6
              s_total.clear()
 7
              for e in E[p]:
 8
                   for f in F[p]:
 9
                        if (e,f) in t:
10
                             s_total[e] += t[(e,f)]
11
12
              for e in E[p]:
                   for f in F[p]:
1.3
                        if (e,f) in t:
                             count[(e,f)] += t[(e,f)]/s_total[e]
1.5
                             total[f] += t[(e,f)]/s_total[e]
16
17
18
         non_zero_num = 0
19
         total num = 0
20
         for e,f in list(t):
21
              t_ef = count[(e,f)]/total[f]
22
              if t_ef > zero:
23
                   t[(e,f)] = t ef
24
                   non_zero_num += 1
25
              else:
26
                   t.pop((e,f))
2.7
              total_num += 1
         return non zero num, total num
```

4、主循环 mainloop (E,F,t)

```
def mainloop(E,F,t,threshold=0.05):
2
        print ("第1次迭代")
        _,last_total_num = one_loop_for_train(E,F,t)
3
 4
        i = 1
        while True:
 5
 6
             i += 1
 7
             print ("第%d次迭代: t中减少了"%i,end=' ')
             _,total_num = one_loop_for_train(E,F,t)
              _rate = (last_total_num-total_num)/last_total_num
 9
10
             if _rate<threshold:
11
                  break
12
             last_total_num = total_num
```

附录 B: 实验 2 部分代码

1、探索法对齐 GROW_DIAG(e2f,f2e,ee_len,ff_len)

```
def GROW DIAG(e2f, f2e,ee len,ff len):
2
         alignment = e2f & f2e
         union_ali = e2f | f2e
3
 4
         for i in range(ee_len):
 5
              for j in range(ff_len):
 6
                   for ni,nj in ((-1,0),(0,-1),(1,0),(0,1),
                            (-1,-1), (-1,1), (1,-1), (1,1)):
 7
                       inew = max(min(0,i+ni),ee len-1)
 8
                       jnew = max(min(0,j+nj),ff_len-1)
 9
                       if (inew, jnew) in union ali:
10
                            if noaligned(inew,alignment,0) \
                               or noaligned(jnew,alignment,1):
11
                                 alignment.add((inew,jnew))
         return alignment
```

2、一致短语对抽取 extract(f_start,f_end,e_start,e_end,

```
A, ff len)
```

```
1
    def extract(f start, f end, e start, e end, A, ff len):
2
         if f_{end} == -1:
3
               return set()
 4
         for e,f in A:
 5
              if f_start<=f<=f_end:
 6
                    if not e_start<=e<=e_end:
 7
                         return set()
 8
         E = set()
         f_s = f_start
9
10
         while f s>=0:
11
               f_e = f_end
12
               while f e<ff len:
                    if 0 < f e - f s <= 6 and 0 < e end - e start <= 6:
13
                         E.add((e_start,e_end,f_s,f_e))
14
15
16
               f s = 1
         return E
```

3、计算短语间的翻译概率 cal_translation_probability_ for pharse(pharses dic)

```
1  def cal_translation_probability_for_pharse(pharses_dic):
2      count = defaultdict(lambda : defaultdict(int))
3      t_pha = {}
4      for pe,pf in pharses_dic:
5      count[pf][pe] += 1
```

```
6
         for pf in count:
 7
             sum = len(count[pf])
 8
             if sum_<10:
 9
                 continue
10
             for pe in count[pf]:
11
                  p = count[pf][pe]/sum_
12
                  t pha[(pe,pf)] = p
13
         return t_pha
```

4、主循环

```
zh=txt_to_lst('data/fbis.zh.10k')
 1
    en=txt to lst('data/fbis.en.10k')
    t1 = read_translation_probability_for_word('data/pharse_\
 3
         base_model_p_result1.txt')
    t1 = read_translation_probability_for_word('data/pharse_\
 4
         base_model_p_result1.txt')
 5
    pharses_dic = defaultdict(int)
 6
    for i in range(len(zh)):
 7
         ee = en[i].lower().split()
         ff = zh[i].lower().split()
 8
 9
         ee_len,ff_len = len(ee),len(ff)
10
11
         e2f = get_aligment_for_sentence(t1,ee,ff,ee_len,\
               ff_len)
12
         f2e = get_aligment_for_sentence(t2,ff,ee,ff_len,\
               ee len)
         A = GROW_DIAG(e2f, f2e, ee_len, ff_len)
13
14
         for e_start in range(ee_len):
1.5
              for e_end in range(e_start,ee_len):
16
17
                  f_start,f_end = ff_len-1,-1
18
                  for e,f in A:
19
                       if e_start<=e<=e_end:</pre>
2.0
                            f_start = min(f, f_start)
21
                            f end = max(f, f end)
22
                  E = extract(f_start, f_end, e_start, e_end, A, \
23
                      ff len)
24
                   read_pharse(E,ee,ff,pharses_dic)
25
26 t_pha=cal_translation_probability_for_pharse(pharses_dic)
```