
自然语言处理及应用

实验报告



名称 机器翻译 machine translation

姓名 罗福杰

班级 硕0078班

学号 3120305208

Email 1626027173@qq.com

日期 2020-10-31

1. 实验目的

This assignment focuses on machine translation. The goal is to implement two translation models, IBM model 1 and IBM model 2, and apply these models to predict English/Chinese word alignments.

Available resources:

- *English/Chinese parallel corpus: <http://www.datatang.com/data/14779/>
- (using “cn.txt” and “en.txt” as parallel corpus or make parallel corpus by yourself)
- GIZA++: <https://code.google.com/p/giza-pp/>

- 1) 理解并掌握统计机器翻译中的 IBM module 1 和 IBM module 2 中的概念和原理。
- 2) 了解EM算法，并理解该算法在统计机器学习中的应用。

2. 实验环境

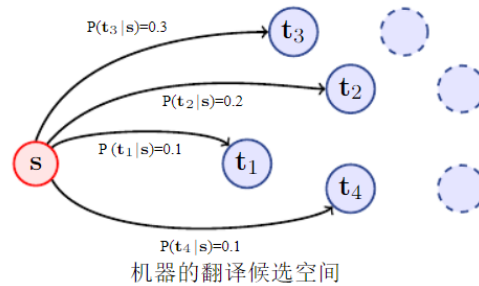
- Win 10 操作系统 64位；
- Spyder编辑器，python 3.6版本；
- 注：因为实验过程中要读取数据，因此，如果运行源代码，要修改文件的相应的路径。该实验的数据在处理过程中较大，可能会出现内存溢出的情况（memory out）。

3. 实验方法

3.1 噪声信道模型

噪声信道模型（Noise Channel Model）是Shannon在上世纪40年代提出的，该模型在语音识别领域得到了广泛的应用，后来被应用在统计机器翻译中。IBM模型也是基于噪声信道模型。

在统计机器翻译中，对于源语言句子 s ，所有目标语词串 t 都是可能的译文，只是译文的可能性大小不同。即每对 (s, t) 都有一个概率值 $P(t|s)$ 来描述 s 翻译成 t 的好与坏。下图展示了机器的翻译候选空间。



在噪声信道模型中，源语言 s （信宿）被看作是由目标语言句子 t （信源）经过一个由噪声的信道得到的。如果知道了 s 和信道的性质，可以通过 $P(t|s)$ 得到信源的信息。对于汉译英的翻译任务，汉语句 s 可以被看作是英语句子 t 加入噪声通过信道后得到的结果。于是需要根据观察到汉语的特征，通过概率 $P(t|s)$ 猜测最为可能的英语句子。这个找到最可能的目标语句（信源）的过程被称为解码。这个过程可以表述为：给定输入 s ，找到最可能的输出 t ，使得 $P(t|s)$ 达到最大：

$$\hat{t} = \underset{t}{\operatorname{argmax}} P(t|s)$$

在IBM 模型中，可以使用贝叶斯准则对 $P(t|s)$ 进行变换，得到的公式如下：

$$P(t|s) = \frac{P(s, t)}{P(s)} = \frac{P(s|t)P(t)}{P(s)}$$

在该式中：

- $P(s|t)$ 是由译文 t 得到源语言句子 s 的翻译模型；
- $P(t)$ 是语言模型，表示目标语言句子 t 出现的可能性；
- $P(s)$ 是源语言句子 s 出现的可能性，因为 s 是输入的不变量，而且 $P(s) > 0$ ，所以省略该项（分母），不影响最大值的求解。

于是翻译模型的目标被重新定义为：给定源语言句子 s ，寻找这样的译文 t ，它使得翻译模型 $P(t|s)$ 和语言模型 $P(t)$ 的乘积最大。

$$\text{Target} = \operatorname{argmax} \{P(s|t)P(t)\}$$

由此，解决上述问题需要三个步骤：

- 建模：如何建立翻译模型 $P(t|s)$ 和语言模型 $P(t)$ 。
- 训练：如何获得 $P(t|s)$ 和 $P(t)$ 所需的参数。
- 解码：如何完成搜索最优解的过程，即完成 argmax 的过程。

3.2 基于词对齐的翻译模型

词对齐描述了源语言句子和目标语言句子之间单词级别的对应。具体是：给定源语言句子 $S = s_1^m \equiv s_1 s_2 \dots s_m$ 和目标语言译文 $T = t_1 \equiv t_1 t_2 \dots t_l$ 。IBM模型假设词对齐通常满足：一个源语言单词只能对应一个目标语单词；源单词可以翻译为空。由于句子之间的对应关系可以由单词之间的对应进行表示，于是，句子翻译的概率可以被转化为词对齐生成的概率：

$$P(s|t) = \sum_a P(s, a|t) \quad \text{式A}$$

其中， $P(s, a|t)$ 可以表示为：

$$P(s, a|t) = P(m|t) \prod_{j=1}^m P(a_j | a_1^{j-1}, s_1^{j-1}, m, t) P(s_j | a_1^{j-1}, s_1^{j-1}, m, t) \quad \text{式B}$$

其中：

- s_j 和 a_j 分别表示第j个源语言单词及第j个源语言单词对应到目标位置；
- s_1^{j-1} 表示前（j-1）个源语言单词， a_1^{j-1} 表示前(j-1)个源语言的词对齐；
- m 表示源语句子的长度。

3.3 IBM module 1

IBM module 1对式B进行了化简：

- 假设 $P(m|t)$ 为常数 ε ，即源语言的长度的生成概率服从均匀分布： $P(m|t)=\varepsilon$ ；
- 对齐概率 $P(a_j | a_1^{j-1}, s_1^{j-1}, m, t)$ 仅依赖译文长度l，即每个词对齐连接的概率也服从均匀分布，即任何源语言位置j对齐到目标语言任何位置都是等概率的。即 $P(a_j | a_1^{j-1}, s_1^{j-1}, m, t) = \frac{1}{l+1}$ 。
- 源单词 s_j 的生成概率 $P(s_j | a_1^{j-1}, s_1^{j-1}, m, t)$ 仅依赖与其对齐的译文单词 t_{a_j} ，即词汇翻译概率 $f(s_j | t_{a_j})$ 。

将上述三个假设代入公式A中，得到 $P(s|t)$ 的表达式C：

$$\begin{aligned} P(s|t) &= \sum_a P(s, a|t) = \sum_a P(m|t) \prod_{j=1}^m P(a_j | a_1^{j-1}, s_1^{j-1}, m, t) P(s_j | a_1^{j-1}, s_1^{j-1}, m, t) \\ &= \sum_a \varepsilon \prod_{j=1}^m \frac{1}{l+1} f(s_j | t_{a_j}) = \sum_a \frac{\varepsilon}{(l+1)^m} \prod_{j=1}^m f(s_j | t_{a_j}) \end{aligned} \quad \text{式C}$$

在该式C中，需要遍历所有的词对齐，即 \sum_a 。可以进一步表示为：

$$P(s|t) = \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \frac{\varepsilon}{(l+1)^m} \prod_{j=1}^m f(s_j|t_{a_j}) \quad \text{式D}$$

式D分为两个部分：

- 1) 遍历所有的对齐a。其中a由 $\{a_1 \dots a_m\}$ 组成，每个 a_j 从译文开始位置（0）循环到截至位置（l）；
- 2) 对于每个a累加对齐的概率 $P(s, a|t) = \frac{\varepsilon}{(l+1)^m} \prod_{j=1}^m f(s_j|t_{a_j})$ 。

这样就得到了IBM模型1中句子翻译概率的计算方式。可以看出IBM module 1 的假设把翻译模型化简成了很简单形式，对于给定的s, a, t，只要知道 ε 和 $f(s_j|t_{a_j})$ 就可以计算 $P(s|t)$ ，从而求出 $P(t|s)$ 。

3.4 IBM module 2

IBM module 1 虽然很好地化简了问题，但是由于使用了过强的假设，导致模型和实际情况有较大的差异，其中一个假设是词对齐的生成概率模型服从均匀分布，这会使模型忽略目标语言单词的位置信息。因此，当单词翻译相同但顺序不同时，翻译的概率一样。同时，由于源语言单词是错误位置的目标语单词生成的，不合理的对齐也会导致不合理的词汇翻译概率。

因此，IBM module 2抛弃了对 $P(a_j|a_1^{j-1}, s_1^{j-1}, m, t)$ 服从均匀分布的假设，IBM module 2 认为词对齐具有倾向性的，它要与源语言单词的位置和目标语言单词的位置有关。具体是：对齐位置信息 a_j 的生成概率与位置j，源语句子长度m和译文长度l有关，表示为：

$$P(a_j|a_1^{j-1}, s_1^{j-1}, m, t) = a(a_j|j, m, l)$$

IBM module 2的其他假设均与module 1相同，因此IBM module 2的数学公式为：

$$P(s|t) = \sum_a P(s, a|t) = \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \varepsilon \prod_{j=1}^m a(a_j|j, m, l) f(s_j|t_{a_j}) \quad \text{式E}$$

类似module 1，module 2可以也可以看成两部分：

- 1) 遍历所有的a；
- 2) 对于每个a累加对齐的概率 $P(s, a|t)$ ，即计算对齐概率 $a(a_j|j, m, l)$ 和词汇翻译概率 $f(s_j|t_{a_j})$ 对于所有源语言位置的乘积。

3.5 解码和计算优化

对计算IBM module 1 和 2的公式进行化简，得到一下公式：

$$\text{IBM module 1:} \quad P(s|t) = \frac{\varepsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i).$$

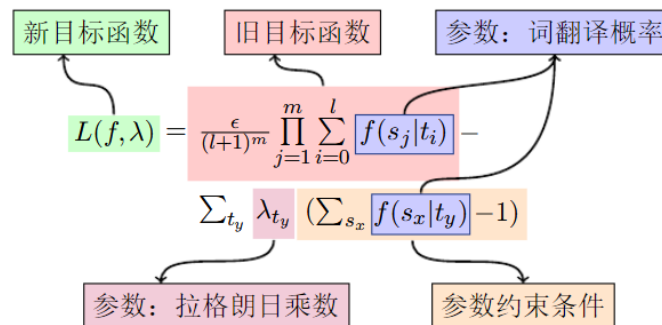
$$\text{IBM module 2:} \quad P(s|t) = \varepsilon \prod_{j=1}^m \sum_{i=0}^l a(i|j, m, l) f(s_j|t_i)$$

有了上面的模型和化简后的结果，对其参数进行估计，然后将问题看成优化问题。其优化的目标函数是 $P(s|t)$ ，以IBM module 1 为例，最终带约束的优化目标可以表示为：

$$\max\left\{\frac{\varepsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)\right\} \quad \text{s.t. 任意单词 } t_y: \sum_{s_x} f(s_x|t_y) = 1$$

其中： $\sum_{s_x} f(s_x|t_y) = 1$ 是优化的约束条件，以保证翻译概率归一化的要求。

通过上式，已将IBM module转化为了带约束的目标函数优化的问题，由于目标函数是可微分函数，解决这类问题的一种方法是把拉格朗日乘数法将带约束问题转化为不带约束问题，它的思想是把含有 n 个变量和 m 个约束条件的优化问题转化为含有 $(n+m)$ 个变量的无约束优化问题。下图以IBM module 1 为例，表示了引入拉格朗日乘数的各项的意义：（来源：机器翻译统计建模与深度学习方法）



然后对该式进行求导，令导数为0，并化简等一系列运算，可以得到下式：

$$f(s_u|t_v) = \lambda_{t_v}^{-1} \frac{\varepsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i) \sum_{j=1}^m \delta(s_j, s_u) \sum_{i=0}^l \delta(t_i, t_v) \frac{f(s_u|t_v)}{\sum_{i=0}^l f(s_u|t_i)}$$

可以看出，这不是计算 $f(s_u|t_v)$ 的解析式，因为等式右边包含了 $f(s_u|t_v)$ ，不过可以用期望最大法（**expectation maximization, EM**）的方法进行求解。使用EM方法可以利用上式迭代计算 $f(s_u|t_v)$ ，使其最终收敛到最优解。

EM方法的思想是：用当前的参数，求似然函数的期望，之后最大化这个期望同时得到新的一组参数的值。而在IBM 模型中，为了方便 $f(s_u|t_v)$ 的计算，用某词语翻译出现的总次数代替一部分的累加和乘积运算。具体过程如下图所示：

$$f(s_u|t_v) = \lambda_{t_v}^{-1} \frac{\epsilon}{(t+1)^m} \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)$$

||
P(s|t)

$$\frac{\sum_{j=1}^m \delta(s_j, s_u) \sum_{i=0}^l \delta(t_i, t_v)}{\sum_{i=0}^l f(s_u|t_i)}$$

“ t_v 翻译为 s_u ” 这个事件
出现次数的期望的估计
称之为期望频次 (Expected Count)

最终得到IBM module 的迭代求解公式：

$$f(s_u|t_v) = \frac{\sum_{i=1}^N c_E(s_u|t_v; s^{[i]}, t^{[i]})}{\sum_{s_u} \sum_{i=1}^N c_E(s_u|t_v; s^{[i]}, t^{[i]})}$$

利用这个公式计算新的 $f(s_u|t_v)$ 值 反复执行 用当前的 $f(s_u|t_v)$ 计算期望频次 $c_E(\cdot)$

以IBM 模型 1为例，其伪代码如下所示：

IBM 模型 1 的训练 (EM 算法)

输入: 平行语料 $(s^{[1]}, t^{[1]}), \dots, (s^{[N]}, t^{[N]})$

输出: 参数 $f(\cdot|\cdot)$ 的最优值

1: **Function** EM($\{(s^{[1]}, t^{[1]}), \dots, (s^{[N]}, t^{[N]})\}$)

2: Initialize $f(\cdot|\cdot)$ ▷ 比如给 $f(\cdot|\cdot)$ 一个均匀分布

3: Loop until $f(\cdot|\cdot)$ converges

4: **foreach** $k = 1$ to N **do**

5: $c_E(s_u|t_v; s^{[k]}, t^{[k]}) = \sum_{j=1}^{|s^{[k]}|} \delta(s_j, s_u) \sum_{i=0}^{|t^{[k]}|} \delta(t_i, t_v) \cdot \frac{f(s_u|t_v)}{\sum_{i=0}^l f(s_u|t_i)}$

6: **foreach** t_v appears at least one of $\{t^{[1]}, \dots, t^{[N]}\}$ **do**

7: $\lambda'_{t_v} = \sum_{s_u} \sum_{k=1}^N c_E(s_u|t_v; s^{[k]}, t^{[k]})$

8: **foreach** s_u appears at least one of $\{s^{[1]}, \dots, s^{[N]}\}$ **do**

9: $f(s_u|t_v) = \sum_{k=1}^N c_E(s_u|t_v; s^{[k]}, t^{[k]}) \cdot (\lambda'_{t_v})^{-1}$

10: **return** $f(\cdot|\cdot)$

4.实验结果

4.1 实验过程

1) 数据预处理:

本实验需要从两个文件里提取数据，并对数据进行预处理。步骤：先从.txt文件中分别获取信息，然后通过正则表达式去除一些标点符号和非法的字符，如：“÷”，“a”这种没有意义的字符。处理完成后，对语句进行切分，将得到的内容存入两个列表中。得到的结果如下：

```
an, extra, [at, this, moment, she, noticed,
point, the, speaker, was, interrupted], [at, th
[at, this, point, we, are, ready, to, proceed
able, to, afford, a, holiday], [at, this, rat
holiday], [at, this, school, only, ten, people
that, may, seem, a, small, number, but, it,
this, stage, i, will, take, anything], [at, t
curie, joined, with, her, in, search, of, the
another, fundamental, change, occurred], [at, thi
their, fields], [at, this, time, of, the, yea
down, late], [at, this, time, of, the, year,
suddenly, much, colder], [at, this, time, of,
the, morning, and, goes, down, at, about, eig
time, the, watch, and, camera, industries, were
yesterday], [at, this, very, moment, there, are
fallow], [at, this, they, all, rushed, out], [
is, always, you, who, put, your, finger, on,
[at, times, he, does, his, school, work, very
very, poorly], [at, times, my, heart, cries,
things], [at, times, solemn, speeches, sober, d
wonder, why, there, is, no, complete, system,
perfectly], [at, twenty, two, he, s, just, be
clock, the, fighters, fell, in], [at, two, th
disagreement], [at, what, speed, are, you, driv
[at, what, time, did, he, get, in], [at, wh
station], [at, what, time, shall, we, arrive,
[at, what, time], [at, whose, door, does, the
to, sell, the, house, so, as, to, pay, off
```

```
胆汁, 尤其, 用作, 鱼, 和, 蔬菜, 的, 配料, [丰富
汤, 一种, 由, 肉, 鱼, 或, 贝壳, 制成, 的, 含,
但, 经常, 是, 很, 奢侈, 的, 人], [糖果, 一种,
果, 通常, 用, 水果, 或, 坚果, 调味, 或与, 之,
的, 浓郁, 的, 甜, 葡萄酒], [回旋, 轮, 在, 游乐场
形, 轨迹, 运动, 的, 快速, 突, 转, 的, 小车], [
N, N, 地, 飞奔而来], [茎沟, 与, 沟, 之间, 的,
鞭, 痕], [条痕, 因, 鞭打, 或, 猛击, 或, 有时,
起, 或, 肿块], [冰, 丘, 冰原, 上, 的, 冰脊, 或
件], [后座力, 很大的, 步枪], [直角, 弯, 肘, 导管,
灰, 泥板, 灰, 胶纸, 柏, 板, 由, 石膏, 灰泥, 作
板, 制成, 的, 模, 板, 在, 建筑, 中, 用于, 代
用, 木模, 或, 金属模, 做成, 的, 坚硬, 的, 构架,
钱, 和, 坊, 的, 边缘], [武装部队, 国威, 一, 圈
或, 一, 圈, 绳子, 用来, 升起, 桅杆, 或, 防止,
或, 酒, 醉的, 节日, 狂欢], [水, 面上, 的, 等,
上升, 或, 升高], [上涨, 增加, 价格, 或, 招生, 人
或, 焚, 洗礼, 着, 吹气, 的, 一种, 宗教仪式], [一
新, 信用, 购物, 方案], [在, 山谷, 中, 蜿蜒, 的,
流, 用船, 有, 龙骨, 但无, 帆], [车行道, 汽车, 或
客, 共产主义, 的, 道路], [道路, 为, 从, 一地, 到
或, 高速公路], [烤肉, 叉, 可以, 转动, 的, 烤肉,
违法, 的, 匪帮], [抢劫, 或, 偷窃], [手术, 服, 在
为, 防止, 感染, 或, 传染, 而, 穿的, 长袍, 或,
高昂, 的, 年轻人], [水晶石, 石英, 断片, 直径, 4,
尤, 指, 天然, 圆, 形的, 石英, 颗粒], [树枝, 石
岩石, 或, 矿物], [钻井, 机, 一种, 在, 采矿, 中
乱, 的, 工具], [火箭, 铁路, 带回, 我们, 用, 其他
```

2) 组成语料库和词典:

由上面的分词结果，得到语料库corpus，以及这个语料库中出现的中文和英文的词典，其中corpus的结构中，每一句的英文内容放在同一个list中，共有100000条数据，其中每个list包含如下内容：

中文	[科幻小说, '不能, '简单, '地, '看成是, '供, '消遣, '的, '而, '实际上, '它, '给, '读者, '展示, '更, '深刻, '的, '内容]
英文	['a, 'science, 'fiction, 'cannot, 'not, 'be, 'regarded, 'as, 'a, 'mere, 'entertainment, 'but, 'in, 'fact, 'it, 'tells, 'the, 'reader, 'much, 'more']

实现结果如下：（下图依次是语料库corpus，英文词典english_words，中文词典foreign_words的部分结果示意图）

['tgilbert', 'a', 'aaron', 'ababa', 'aback', 'abacus', 'abaft', 'abalone', 'abandon', 'abandoned', 'abandoning', 'abandonment', 'abasement', 'abash', 'abashed', 'abate', 'abated', 'abatement', 'abbey', 'abbot', 'abbott', 'abbreviated', 'abbreviation', 'abbreviations', 'abby', 'abc', 'abdicte', 'abdication', 'abdomen', 'abdominal', 'abduct', 'abducted', 'abduction', 'abe', 'abel', 'aberrant', 'aberration', 'abeyance', 'abhorrence', 'abhorrent', 'abidance', 'abide', 'abiding', 'abilities', 'ability', 'abject', 'able', 'abloom', 'ably', 'abnormal', 'abnormality', 'abnormally', 'aboard', 'abode', 'abolish', 'abolished', 'abolition', 'abominable', 'aboriginal', 'aborigines', 'abort', 'aborted', 'abortion', 'abortions', 'abound', 'abounding', 'abounds', 'about', 'above', 'abraham', 'abrasion', 'abrasive', 'abrasives', 'abreast', 'abridge', 'abroad', 'abrupt', 'abruptly', 'absciss', 'absconded', 'absence', 'absences', 'absent', 'absentee', 'absentees', 'absentia', 'absolute', 'absolutely', 'absolved', 'absorb',

世，二个，'二个人'，'二个月'，'二义性'，'二价'，'二任'，'二份'，'二位'，'二元'，'二元论'，'二分之一'，'二分钟'，'二分音
好'，'二副'，'二十'，'二十一'，'二十七'，'二十万'，'二十三'，'二十九'，'二十二'，'二十五'，'二十五日'，'二十八'，'二十六'
六'，'二十九'，'二十分'，'二十四'，'二十四日'，'二十年'，'二十日'，'二千'，'二千二百万'，'二千人'，'二又'，'二垒'，'二'
来'，'二头肌'，'二局'，'二层'，'二年'，'二年级'，'二度'，'二战'，'二手'，'二手货'，'二手车'，'二无'，'二晚'，'二月'，'二'
杯'，'二极管'，'二档'，'二次大战'，'二次方程'，'二氧化氮'，'二氧化碳'，'二点钟'，'二环'，'二班'，'二百'，'二百万'，
二百三十'，'二百人'，'二百元'，'二种'，'二笔'，'二等'，'二等品'，'二等奖'，'二级'，'二维'，'二老'，'二者'，'二者之间'，'二'
善必居其一'，'二话没说'，'二轮'，'二进制'，'二进制位'，'二进制数'，'二进制码'，'二郎腿'，'二重唱'，'二重奏'，'二页'，'二'
首'，'于'，'于事无补'，'于归'，'于我'，'于斯'，'于是'，'于是乎'，'于此'，'亏'，'亏了'，'亏待'，'亏心'，'亏损'，'亏本'，'亏'
空'，'亏缺'，'云'，'云块'，'云天'，'云层'，'云彩'，'云影'，'云散'，'云朵'，'云杉'，'云梯'，'云母'，'云海'，'云雀'，'云集'，
云雾'，'云霞'，'互'，'互不'，'互为条件'，'互交'，'互利'，'互助'，'互向'，'互惠'，'互换'，'互敦互学'，'互敬'，'互斥'，'互'
无'，'互相'，'互相冲突'，'互相制约'，'互相合作'，'互相帮助'，'互相攻击'，'互相矛盾'，'互相联系'，'互相配合'，'互联网'，'互'
见'，'互让'，'互访'，'互谅互让'，'互赠'，'互连'，'互通式'，'互通有无'，'五'，'五一节'，'五万'，'五万元'，'五个人'，'五个'
月'，'五份'，'五体投地'，'五元'，'五六'，'五分钟'，'五分线'，'五加'，'五十'，'五十七'，'五十万'，'五十二'，'五十五'，'五十'
十'，'五十分'，'五千五百'，'五千五百元'，'五千米'，'五叶'，'五国'，'五折'，'五天'，'五十六'，'五十八'，'五十九'，'五十'

t 是用来保存不同外文单词翻译成不同英文单词的概率。其核心代码如下:

t 是用来保存不同外文单词翻译成不同英文单词的概率。其核心代码如下:

```
def create_t(corpus,init_val):
    t = { }
    for sp in corpus:
        for f_w in sp[0]:
            for e_w in sp[1]:
                if f_w not in t:
                    t[f_w] = { }
                t[f_w][e_w] = init_val
    return t
```

通过打印输出，得到t共有4261308个。

序列对的总个数有: 4261308

在参数初始化部分，主要是初始化训练的轮数epoch_num的值，然后开始进行训练，本次实验为了对比效果，用不同的轮数训练了多次。其中一次的过程图如下所示：

```

epoch : 0 ALL: 100000 NOW processing : 98000
epoch : 0 ALL: 100000 NOW processing : 99000
-----epoch 2-----
epoch : 1 ALL: 100000 NOW processing : 0
epoch : 1 ALL: 100000 NOW processing : 1000
epoch : 1 ALL: 100000 NOW processing : 2000
epoch : 1 ALL: 100000 NOW processing : 3000
epoch : 1 ALL: 100000 NOW processing : 4000
epoch : 1 ALL: 100000 NOW processing : 5000
epoch : 1 ALL: 100000 NOW processing : 6000
epoch : 1 ALL: 100000 NOW processing : 7000
epoch : 1 ALL: 100000 NOW processing : 8000
epoch : 1 ALL: 100000 NOW processing : 9000
epoch : 1 ALL: 100000 NOW processing : 10000
epoch : 1 ALL: 100000 NOW processing : 11000
epoch : 1 ALL: 100000 NOW processing : 12000
epoch : 1 ALL: 100000 NOW processing : 13000
epoch : 1 ALL: 100000 NOW processing : 14000
epoch : 1 ALL: 100000 NOW processing : 15000
epoch : 1 ALL: 100000 NOW processing : 16000
epoch : 1 ALL: 100000 NOW processing : 17000
epoch : 1 ALL: 100000 NOW processing : 18000
epoch : 1 ALL: 100000 NOW processing : 19000
epoch : 1 ALL: 100000 NOW processing : 20000
epoch : 1 ALL: 100000 NOW processing : 21000
epoch : 1 ALL: 100000 NOW processing : 22000
epoch : 1 ALL: 100000 NOW processing : 23000

```

4.2 实验结果

在实验过程中，为了更好地对比实验结果，在代码中将对词对齐概率的结果进行了降序排列，并将结果写入到.csv文件中，最终得到的结果如下图所示。

	A	B	C	D
1	汉语	英语	概率	
2	受领	accreditee	1	
3	西班牙文	adios	1	
4	字母数字的	alphanumeric	1	
5	采写	assignment	1	
6	伯利兹	belize	1	
7	领头羊	bellwether	1	
8	出生地点	birthplace	1	
9	电脑化	cyberspace	1	
10	搬运费	cartage	1	
11	小箱	caske	1	
37	有销路	popular	1	
38	打印输出	printout	1	
39	招股	prospectus	1	
40	翻两番	quadruple	1	
41	险别	risk	1	
42	共享软件	shareware	1	
43	花絮	sidebar	1	
44	附加费	surcharge	1	
45	制表机	tabulator	1	
46	皮重	tare	1	
47	三个	three	0.929443	
48	土壤	soil	0.918559	
49	水泥	cement	0.912625	
50	英语	english	0.9073	
51	现代	modern	0.896114	
52	两	two	0.891141	
53	乔治	george	0.889408	
54	玻璃	glass	0.888017	
55	新	new	0.887645	
56	小姐	miss	0.886876	
57	杰克	jack	0.885751	
58	她那	her	0.885629	
59	旧	old	0.88544	
60	黑色	black	0.882724	
61	或者	or	0.88229	
62	比尔	bill	0.87792	
63	服务器	server	0.87725	
64	能量	energy	0.876527	
65	头发	hair	0.876416	
66	棉花	cotton	0.874191	
67	她的	her	0.873476	
68	红色	red	0.873398	
69	白色	white	0.872832	
70	红	red	0.871632	
71	牛肉	beef	0.870173	
72	三	three	0.86883	
73	由...组成	composed of	0.86722	

56749	大员	newsreels	0.045648
56750	头戴	newsreels	0.045648
56751	令人讨厌	dreadful	0.045178
56752	千亿	tripling	0.044905
56753	增加三倍	tripling	0.044905
56754	愿闻其详	tripling	0.044905
56755	猖狂	ills	0.04469
56756	施压	hurd	0.04298
56757	检方	inconsister	0.042566
56758	十尺	bangs	0.035562
56759	侧门	bangs	0.035562
56760	未了	bangs	0.035562
56761	后墙	bangs	0.035562
56762	报亭	bangs	0.035562
56763			
56764			

有结果可知，共有（56763-1）个汉语词语，并将其对应的英语概率最大的单词的挑选出来了，以及输出了相应的概率。此结果可以看出，概率排名在前面单词翻译的效果不错，基本满足要求。最后的汉语词语翻译的结果还有待改进。

然后，提高轮数，观察对于相同的单词翻译的是否相同，如果相同，它的概率会不会随着轮数的增加而上升。

中文	英文	Epoch = 5	Epoch = 15
小箱	casket	1	1
三个	three	0.929443	0.983434
现代	modern	0.896114	0.969343
那么些	such	0.163423	0.163773

除此之外，5轮迭代的实验，准确率在0.8以上的词队有：191组，15轮迭代的实验，准确率在0.8以上的词队有：344组。如下图所示（左图迭代5轮，右图迭代15轮）。

180	黄色	yellow	0.803589
181	下午	afternoon	0.803511
182	政府	governme	0.80311
183	网络	network	0.802837
184	学校	school	0.802076
185	和平	peace	0.801533
186	生意	business	0.801036
187	常常	often	0.800951
188	其他	other	0.8007
189	哦	oh	0.800687
190	简	jane	0.800323
191	冬季	winter	0.800244
192	北京	beijing	0.800019
193	母亲	mother	0.799307
194	战争	war	0.799195
195	面包	bread	0.799126
196	会计部	accounting	0.798944
197	更多		0.79886
336	球队	team	0.804047
337	光	light	0.804001
338	性虐待	abuse	0.803709
339	右	right	0.803387
340	理论	theory	0.803314
341	句子	sentence	0.802548
342	摄像机	camera	0.802005
343	罐装	canned	0.801845
344	世界	world	0.800242
345	按揭	mortgage	0.800052
346	事故	accident	0.799895
347	早上	morning	0.799695
348	通常	usually	0.799625
349	水果	fruit	0.799474
350	野战	field	0.798908
351	棕色	brown	0.798661
352	之前	before	0.79859
353	明天	tomorrow	0.798125
354	总是	always	0.798098

由此可见，迭代的轮数对实验的准确率有着重要的影响，轮数的增加可以提高

汉语翻译为英语的概率的提高。

经过第三章中的原理，IBM module 2抛弃了module 1对 $P(a_j|a_1^{j-1}, s_1^{j-1}, m, t)$ 服从均匀分布的假设，IBM module 2认为词对齐具有倾向性的，它要与源语言单词的位置和目标语言单词的位置有关。具体是：对齐位置信息 a_j 的生成概率与位置 j ，源句子长度 m 和译文长度 l 有关，表示为：

$$P(a_j|a_1^{j-1}, s_1^{j-1}, m, t) = a(a_j|j, m, l)$$

IBM module 2在实验过程上和IBM module 1的过程相同，训练模型，并将实验结果保存在csv文件中，如上一样，对其结果进行分析：得到在训练相同轮数的情况下，使用module 2训练得到的概率大于0.8的词对数目多于使用module 1得到的结果。如下图所示：（作图为module 1的结果，右图为module 2的结果）。

335	妹妹	sister	0.80453	353	出厂	vintage	0.802553
336	球队	team	0.804047	354	球队	team	0.802421
337	光	light	0.804001	355	算术	arithmetic	0.802416
338	性虐待	abuse	0.803709	356	聚会	party	0.802399
339	右	right	0.803387	357	尤其	especially	0.801618
340	理论	theory	0.803314	358	进去	in	0.801604
341	句子	sentence	0.802548	359	字典	dictionary	0.801529
342	摄像机	camera	0.802005	360	这座	the	0.801457
343	罐装	canned	0.801845	361	暴力	violence	0.80056
344	世界	world	0.800242	362	李斯特	liszt	0.800546
345	按揭	mortgage	0.800052	363	王子	prince	0.799829
346	事故	accident	0.799895	364	请	please	0.799318
347	早上	morning	0.799695	365	宣布独立	independe	0.79905
348	通常	usually	0.799625	366	婴儿	baby	0.798972
349	水果	fruit	0.799474	367	世界	world	0.798898
350	野战	field	0.798908	368	宾夕法尼亚	pennsylvan	0.798702
351	棕色	brown	0.798661	369	之前	before	0.798422
352	之前	before	0.79859	370	明天	tomorrow	0.798391
353	明天	tomorrow	0.798125	371	按揭	mortgage	0.798243
354	总是	always	0.798098				

5.遇到问题及解决思路

1) 问题：迭代了多轮，但参数更新不动；

- 问题描述：在测试代码的正确性时遇到了问题：参数迭代了好几轮，但是概率值一直没有更新；如下面所示：

```

-----epoch 1-----
我的 my 0.5
书 my 0.5
一本 a 0.5
杂志 a 0.5
-----epoch 2-----
我的 my 0.5
书 my 0.5
一本 a 0.5
杂志 a 0.5
-----epoch 3-----
我的 my 0.5
书 my 0.5
一本 a 0.5
杂志 a 0.5
-----epoch 4-----
我的 my 0.5
书 my 0.5
一本 a 0.5
杂志 a 0.5
-----epoch 5-----
我的 my 0.5
书 my 0.5
一本 a 0.5
杂志 a 0.5

```

- **原因：**其实代码的主题功能没问题，而是语料库的问题，此时的语料库是：

```
corpus = [[['我的','书'], ['my','book']], [['一本','杂志'], ['a','magazine']]]
```

当语料库中的句子完全不相关时，就没办法就行参数的迭代和更新了。

- **解决方法：**在预料库中添加相关的语句，如以这个为例，把语料库扩充为：

```
corpus = [[['我的','书'], ['my','book']], [['一本','杂志'], ['a','magazine']], [['我的','杂志'], ['my','magazine']]]
```

当迭代五轮后，得到的结果如下：可见取得了较好的结果。

```

-----epoch 1-----
我的 my 0.5
书 my 0.5
一本 a 0.5
杂志 magazine 0.5
-----epoch 2-----
我的 my 0.6363636363636364
书 book 0.5714285714285715
一本 a 0.5714285714285715
杂志 magazine 0.6363636363636364
-----epoch 3-----
我的 my 0.7478974515333993
书 book 0.6533864541832669
一本 a 0.653386454183267
杂志 magazine 0.7478974515333995
-----epoch 4-----
我的 my 0.8344395715124338
书 book 0.7244788210478862
一本 a 0.7244788210478862
杂志 magazine 0.8344395715124338
-----epoch 5-----
我的 my 0.8960831177730377
书 book 0.7817398718164896
一本 a 0.7817398718164896
杂志 magazine 0.8960831177730378

```

2) 中文乱码

- **问题描述：**将结果写入到csv文件中时，中文字符出现了乱码的情况；
- **原因：**编码参数出错了，encoding = “utf-8-sig”
- **解决方案：**encoding = “utf-8-sig”，其中，sig是signature，即带有签名的utf-8，可以解决中文编码的错误。结果如下：

	A	B	C
1	汉语	英语	概率
2	杂志	magazine	0.896083
3	我的	my	0.896083
4	书	book	0.78174
5	一本	a	0.78174

3) 提升概率的方法

- 问题描述：有些单词的概率较低和单词翻译出错；
- 原因：部分原因是数据预处理的问题，将一些没有实际意义的非法字符处理为一个单词了，虽然这部分字符比较少，但非法字符的存在仍会影响词对齐概率的结果。
- 解决方案：在数据的预处理部分，将非法字符：如“o÷”，“a”这种字符去掉。

4) Module 2 中a的理解

- 问题描述：module 2 中将均匀分布的概率换成了一个变量a，但对a的理解不好；
- 解决方法：通过查找资料，查看了module 2 的伪代码，对a的理解进一步加深，其中，module 2 的伪代码如下：

Input: A training corpus (f^k, e^k) and $1 \leq k \leq S$
Initialization: Initialize t and a parameters (e.g. to random values)
Algorithm:
Do{
 -set all counts to 0;
 -traverse all of the sentence pair
 *for j = 1...m
 ~for i=0...l
 $c(f_j|e_i; f, e) += \alpha(k, l, j)$
 $c(e_i) += \alpha(k, l, j)$
 $c(i|j, l, m) += \alpha(k, l, j)$
 $c(j, l, m) += \alpha(k, l, j)$
 where $\alpha(k, l, j) = \frac{t(f_j|e_i)a(i|j, m, l)}{\sum_{k=0}^l t(f_j|e_k)a(k|j, m, l)}$
 -set $t(f|e) = c(f|e; f, e)/c(e)$ $a(i|j, m, l) = c(i|j, m, l)/c(j, m, l)$
}until t & a converge
Output: parameters t(f|e) a(i|j, m, l)