

从“词”开始：基于单词的翻译模型

肖桐 朱靖波

xiaotong@mail.neu.edu.cn

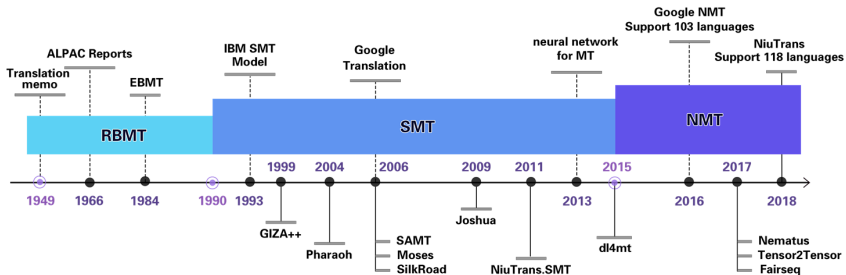
zhujingbo@mail.neu.edu.cn

东北大学 自然语言处理实验室

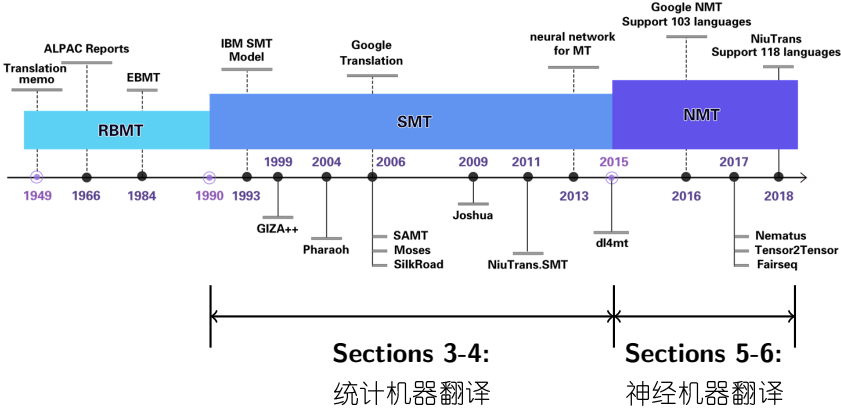
<http://www.nlplab.com>



Landscape



Landscape

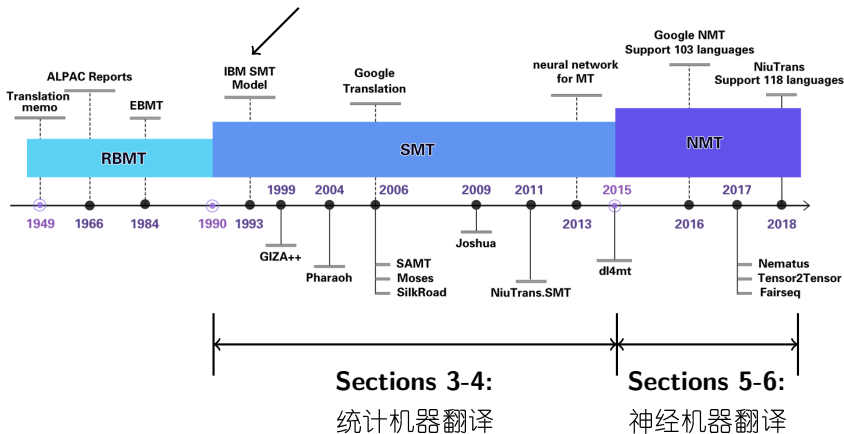


Landscape

本章内容 (Section 3) :

统计机器翻译开山之作 - IBM模型

*The Mathematics of Statistical Machine Translation:
Parameter Estimation* by Peter E. Brown et al., 1993.



按单词进行翻译

基于单词的(统计)机器翻译思想

两个互译句子之间存在一种**单词**间的对应，而整句的翻译是由**单词**的翻译”组成”

我 对 你 感到 满意



I am satisfied with you

按单词进行翻译

基于单词的(统计)机器翻译思想

两个互译句子之间存在一种**单词**间的对应，而整句的翻译是由**单词**的翻译”组成”



按单词进行翻译

基于单词的(统计)机器翻译思想

两个互译句子之间存在一种**单词**间的对应，而整句的翻译是由**单词**的翻译”组成”



- 传统观点下的(基于词)的翻译过程 - 三个步骤

① 分析：将输入句子进行分词



按单词进行翻译

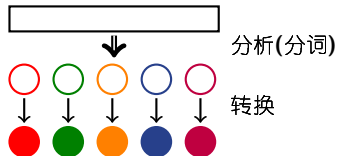
基于单词的(统计)机器翻译思想

两个互译句子之间存在一种**单词**间的对应，而整句的翻译是由**单词**的翻译”组成”



• 传统观点下的(基于词)的翻译过程 - 三个步骤

- ① 分析：将输入句子进行分词
- ② 转换：获得每个源语单词的翻译



按单词进行翻译

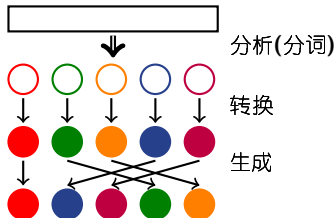
基于单词的(统计)机器翻译思想

两个互译句子之间存在一种**单词**间的对应，而整句的翻译是由**单词**的翻译“组成”



• 传统观点下的(基于词)的翻译过程 - 三个步骤

- 1 分析：将输入句子进行分词
- 2 转换：获得每个源语单词的翻译
- 3 生成：将单词的翻译组成通顺的整句译文



按单词进行翻译

基于单词的(统计)机器翻译思想

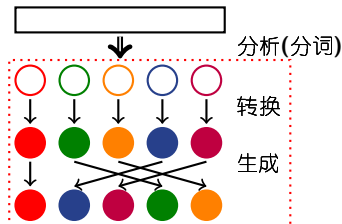
两个互译句子之间存在一种**单词**间的对应，而整句的翻译是由**单词**的翻译”组成”



• 传统观点下的(基于词)的翻译过程 - 三个步骤

- 1 分析：将输入句子进行分词
- 2 转换：获得每个源语单词的翻译
- 3 生成：将单词的翻译组成通顺的整句译文

本章的重点内容为**转换**和**生成**



1. 一个简单的翻译实例

2. IBM模型

- 建模
- 解码
- 模型训练

我们人是如何进行翻译的？

待翻译句子(已经分词):

我 对 你 表示 满意

我们人是如何进行翻译的？

待翻译句子(已经分词):

我



1 I

1 me

1 I'm

对



2 to

2 with

2 for

你



3 you

表示



4 ϕ

4 show

满意



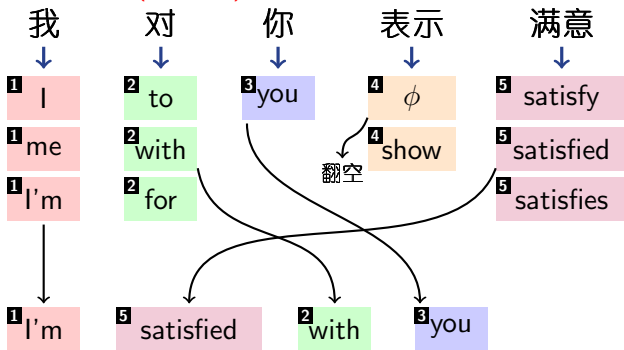
5 satisfy

5 satisfied

5 satisfies

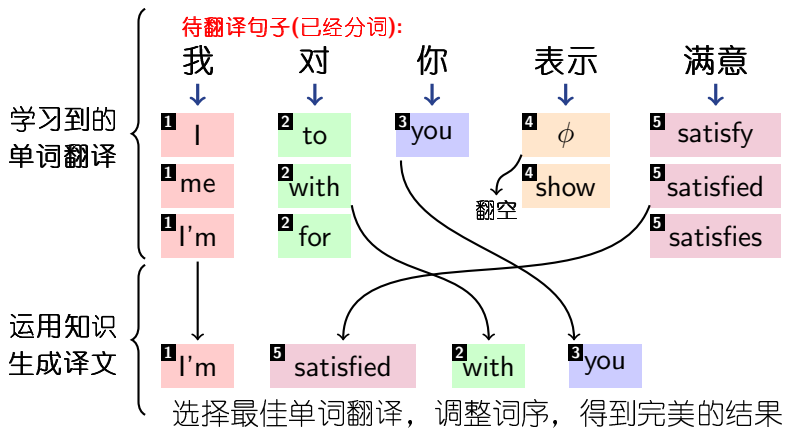
我们人是如何进行翻译的？

待翻译句子(已经分词):



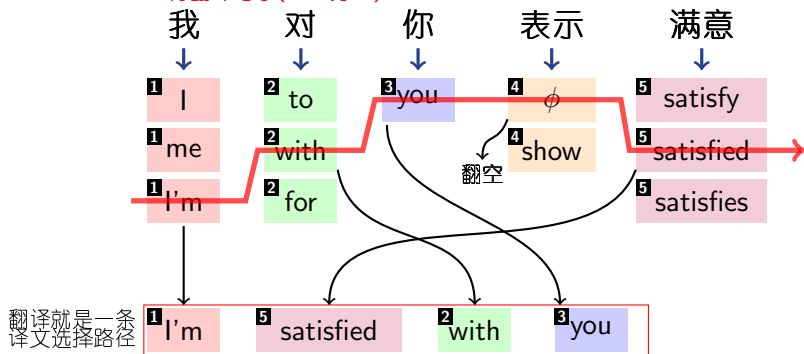
选择最佳单词翻译，调整词序，得到完美的结果

我们人是如何进行翻译的？



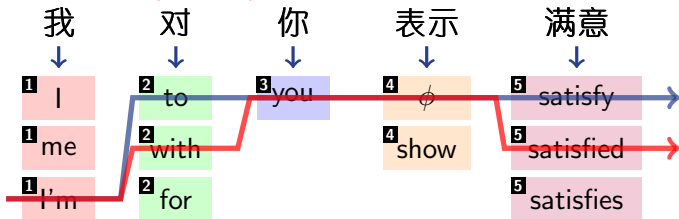
机器翻译系统是如何进行翻译的？

待翻译句子(已经分词):



机器翻译系统是如何进行翻译的？

待翻译句子(已经分词):



翻译就是一条
译文选择路径

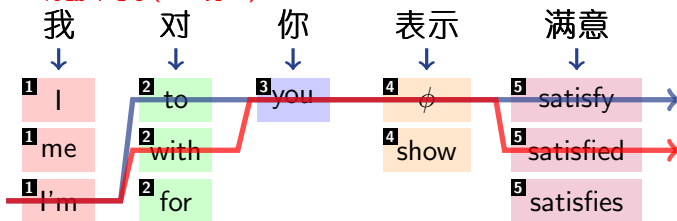
1 | I'm 5 | satisfied 2 | with 3 | you

不同的译文对
应不同的路径

1 | I'm 5 | satisfy 2 | to 3 | you

机器翻译系统是如何进行翻译的？

待翻译句子(已经分词):



翻译就是一条
译文选择路径

1 I'm 5 satisfied 2 with 3 you

不同的译文对
应不同的路径

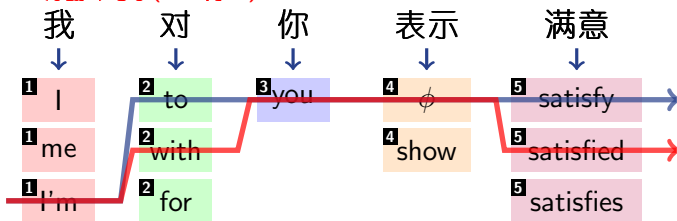
1 I'm 5 satisfy 2 to 3 you

单词翻译的词
序也可能不同

1 I'm 5 satisfy 3 you 2 to

机器翻译系统是如何进行翻译的？

待翻译句子(已经分词):



翻译就是一条
译文选择路径

1 I'm 5 satisfied 2 with 3 you

不同的译文对
应不同的路径

1 I'm 5 satisfy 2 to 3 you

单词翻译的词
序也可能不同

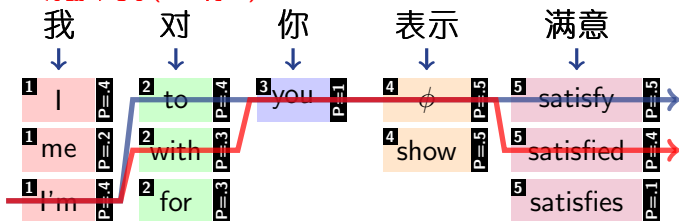
1 I'm 5 satisfy 3 you 2 to

可能的翻译路
径非常多

...

机器翻译系统是如何进行翻译的？

待翻译句子(已经分词):



翻译就是一条
译文选择路径



不同的译文对
应不同的路径



单词翻译的词
序也可能不同



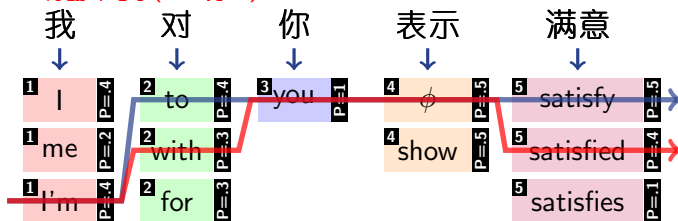
可能的翻译路
径非常多

...

所有翻译单元都是概率化的 P=概率

机器翻译系统是如何进行翻译的？

待翻译句子(已经分词):



系统给每个译文
都赋予一个模型得分

翻译就是一条
译文选择路径



$P=0.042$

不同的译文对
应不同的路径



$P=0.006$

单词翻译的词
序也可能不同



$P=0.003$

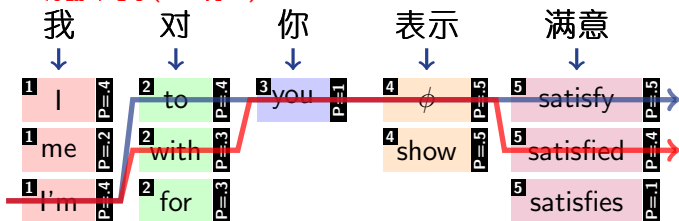
可能的翻译路
径非常多

...

所有翻译单元都是概率化的 $P=\text{概率}$

机器翻译系统是如何进行翻译的？

待翻译句子(已经分词):



系统给每个译文
都赋予一个模型得分

翻译就是一条
译文选择路径



$P=0.042$ ok

不同的译文对
应不同的路径



$P=0.006$

单词翻译的词
序也可能不同



$P=0.003$

可能的翻译路
径非常多

...

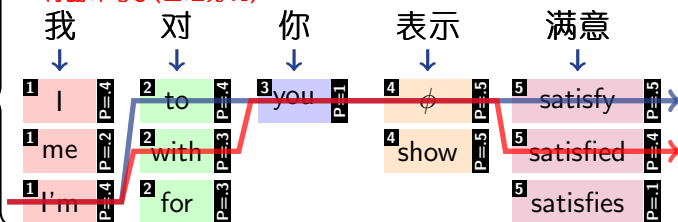
所有翻译单元都是概率化的 $P=\text{概率}$

选择得分
最高的译文

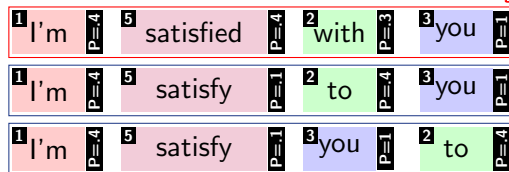
机器翻译系统是如何进行翻译的？

从双语数据中自动学习词典 (训练)

待翻译句子(已经分词):



利用概率化的词典进行翻译 (解码)



系统给每个译文都赋予一个模型得分

$P=0.042$ ok

$P=0.006$

$P=0.003$

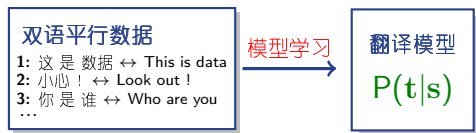
选择得分最高的译文

... 所有翻译单元都是概率化的 $P=\text{概率}$

实现一个简单的机器翻译系统：学习单词翻译概率

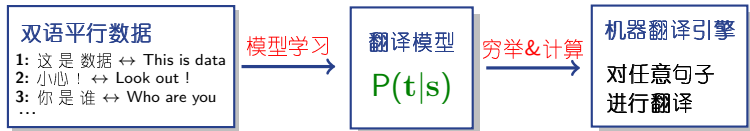
- 构建一个简单的机器翻译系统

- ① 模型学习 - 从双语平行数据中学习单词翻译知识



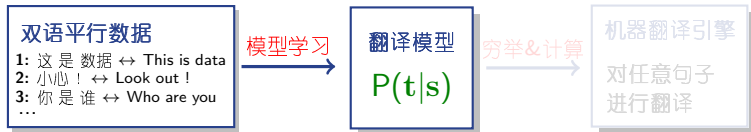
实现一个简单的机器翻译系统：学习单词翻译概率

- 构建一个简单的机器翻译系统
 - 模型学习 - 从双语平行数据中学习单词翻译知识
 - 解码 - 利用单词翻译知识，对新的句子进行翻译



实现一个简单的机器翻译系统：学习单词翻译概率

- 构建一个简单的机器翻译系统
 - 模型学习 - 从双语平行数据中学习单词翻译知识
 - 解码 - 利用单词翻译知识，对新的句子进行翻译



- 步骤1：构建单词翻译表 - 翻译词典**

对于任意的源语言单词 x ，要获得它所有可能的译文 Y 。给定一个互译句对 (s, t) ，对于 $y \in Y$ ，定义 $P(x \leftrightarrow y; s, t)$ 表示 x 和 y 在 (x, y) 中互译的概率，我们用 x 和 y 的联合概率表示：

$$\begin{aligned} P(x \leftrightarrow y; s, t) &\equiv P(x, y; s, t) \\ &= \frac{c(x, y; s, t)}{\sum_{x', y'} c(x', y'; s, t)} \end{aligned}$$

$c(x, y; s, t)$ 表示 (x, y) 在 (s, t) 中共现的次数；

$\sum_{x', y'} c(x', y'; s, t)$ 表示 (s, t) 中任意源/译文单词共现的总次数

实现一个简单的机器翻译系统：学习单词翻译概率(2)

$$P(x, y; \mathbf{s}, \mathbf{t}) = \frac{c(x, y; \mathbf{s}, \mathbf{t})}{\sum_{x', y'} c(x', y'; \mathbf{s}, \mathbf{t})}$$

\mathbf{s} = 机器 翻译 就是 用 计算机 来 进行 翻译

\mathbf{t} = machine translation is just translation by computer

- $c(\text{'翻译'}, \text{'translation'}; \mathbf{s}, \mathbf{t}) = ?$

实现一个简单的机器翻译系统：学习单词翻译概率(2)

$$P(x, y; \mathbf{s}, \mathbf{t}) = \frac{c(x, y; \mathbf{s}, \mathbf{t})}{\sum_{x', y'} c(x', y'; \mathbf{s}, \mathbf{t})}$$

\mathbf{s} = 机器 翻译 就是用计算机来进行翻译

\mathbf{t} = machine translation is just translation by computer

- $c(\text{'翻译'}, \text{'translation'}; \mathbf{s}, \mathbf{t}) = 1$

实现一个简单的机器翻译系统：学习单词翻译概率(2)

$$P(x, y; \mathbf{s}, \mathbf{t}) = \frac{c(x, y; \mathbf{s}, \mathbf{t})}{\sum_{x', y'} c(x', y'; \mathbf{s}, \mathbf{t})}$$

\mathbf{s} = 机器 翻译 就是用计算机来进行翻译

\mathbf{t} = machine translation is just translation by computer

- $c(\text{'翻译'}, \text{'translation'}; \mathbf{s}, \mathbf{t}) = 1 + 1$

实现一个简单的机器翻译系统：学习单词翻译概率(2)

$$P(x, y; \mathbf{s}, \mathbf{t}) = \frac{c(x, y; \mathbf{s}, \mathbf{t})}{\sum_{x', y'} c(x', y'; \mathbf{s}, \mathbf{t})}$$

\mathbf{s} = 机器 翻译 就是用 计算机 来进行 翻译

\mathbf{t} = machine translation is just translation by computer

- $c(\text{'翻译'}, \text{'translation'}; \mathbf{s}, \mathbf{t}) = 1 + 1 + 1$

实现一个简单的机器翻译系统：学习单词翻译概率(2)

$$P(x, y; \mathbf{s}, \mathbf{t}) = \frac{c(x, y; \mathbf{s}, \mathbf{t})}{\sum_{x', y'} c(x', y'; \mathbf{s}, \mathbf{t})}$$

\mathbf{s} = 机器 翻译 就是 用 计算机 来 进行 翻译

\mathbf{t} = machine translation is just translation by computer

- $c(\text{'翻译'}, \text{'translation'}; \mathbf{s}, \mathbf{t}) = 1 + 1 + 1 + 1 = 4$

实现一个简单的机器翻译系统：学习单词翻译概率(2)

$$P(x, y; \mathbf{s}, \mathbf{t}) = \frac{c(x, y; \mathbf{s}, \mathbf{t})}{\sum_{x', y'} c(x', y'; \mathbf{s}, \mathbf{t})}$$

\mathbf{s} = 机器 翻译 就是 用 计算机 来 进行 翻译

\mathbf{t} = machine translation is just translation by computer

- $c(\text{'翻译'}, \text{'translation'}; \mathbf{s}, \mathbf{t}) = 1 + 1 + 1 + 1 = 4$
- $\sum_{x', y'} c(x', y'; \mathbf{s}, \mathbf{t}) = \text{使劲数...} = 63$

实现一个简单的机器翻译系统：学习单词翻译概率(2)

$$P(x, y; \mathbf{s}, \mathbf{t}) = \frac{c(x, y; \mathbf{s}, \mathbf{t})}{\sum_{x', y'} c(x', y'; \mathbf{s}, \mathbf{t})}$$

\mathbf{s} = 机器 翻译 就是 用 计算机 来 进行 翻译

\mathbf{t} = machine translation is just translation by computer

- $c(\text{'翻译'}, \text{'translation'}; \mathbf{s}, \mathbf{t}) = 1 + 1 + 1 + 1 = 4$
- $\sum_{x', y'} c(x', y'; \mathbf{s}, \mathbf{t}) = \text{使劲数...} = 63 = 9 \times 7 = |\mathbf{s}| \times |\mathbf{t}|$
 - ▶ $|\cdot|$ 表示句子长度

实现一个简单的机器翻译系统：学习单词翻译概率(2)

$$P(x, y; \mathbf{s}, \mathbf{t}) = \frac{c(x, y; \mathbf{s}, \mathbf{t})}{\sum_{x', y'} c(x', y'; \mathbf{s}, \mathbf{t})}$$

\mathbf{s} = 机器 翻译 就是用 计算机 来进行 翻译

\mathbf{t} = machine translation is just translation by computer

- $c(\text{'翻译'}, \text{'translation'}; \mathbf{s}, \mathbf{t}) = 1 + 1 + 1 + 1 = 4$
- $\sum_{x', y'} c(x', y'; \mathbf{s}, \mathbf{t}) = \text{使劲数} \dots = 63 = 9 \times 7 = |\mathbf{s}| \times |\mathbf{t}|$
 - ▶ $|\cdot|$ 表示句子长度
- '翻译'和'translation'的互译概率为

$$P(\text{'翻译'}, \text{'translation'}; \mathbf{s}, \mathbf{t}) = 4/63$$

类似的

$$P(\text{'机器'}, \text{'translation'}; \mathbf{s}, \mathbf{t}) = 2/63$$

$$P(\text{'机器'}, \text{'look'}; \mathbf{s}, \mathbf{t}) = 0/63$$

实现一个简单的机器翻译系统：学习单词翻译概率(3)

- 很多时候，我们有多组互译句对 $(\mathbf{s}^{[1]}, \mathbf{t}^{[1]}), \dots, (\mathbf{s}^{[n]}, \mathbf{t}^{[n]})$ ，称之为**双语平行数据(语料)**。翻译概率可以被定义为

$$P(x, y) = \frac{\sum_{i=1}^n c(x, y; \mathbf{s}^{[i]}, \mathbf{t}^{[i]})}{\sum_{i=1}^n \sum_{x', y'} c(x', y'; \mathbf{s}^{[i]}, \mathbf{t}^{[i]})}$$

实现一个简单的机器翻译系统：学习单词翻译概率(3)

- 很多时候，我们有多组互译句对 $(\mathbf{s}^{[1]}, \mathbf{t}^{[1]}), \dots, (\mathbf{s}^{[n]}, \mathbf{t}^{[n]})$ ，称之为**双语平行数据(语料)**。翻译概率可以被定义为

$$P(x, y) = \frac{\sum_{i=1}^n c(x, y; \mathbf{s}^{[i]}, \mathbf{t}^{[i]})}{\sum_{i=1}^n \sum_{x', y'} c(x', y'; \mathbf{s}^{[i]}, \mathbf{t}^{[i]})}$$

- 说白了就是计算 (x, y) 的频次时，在每个句子上累加
 \mathbf{s}^1 = 机器 翻译 就是 用 计算机 进行 翻译
 \mathbf{t}^1 = Machine translation is just translation by computer
 \mathbf{s}^2 = 那 人工 翻译 呢？
 \mathbf{t}^2 = So , what is human translation ?

$$\begin{aligned} & P(\text{'翻译'}, \text{'translation'}) \\ &= \frac{c(\text{'翻译'}, \text{'translation'}; \mathbf{s}^{[1]}, \mathbf{t}^{[1]}) + c(\text{'翻译'}, \text{'translation'}; \mathbf{s}^{[2]}, \mathbf{t}^{[2]})}{\sum_{x', y'} c(x', y'; \mathbf{s}^{[1]}, \mathbf{t}^{[1]}) + \sum_{x', y'} c(x', y'; \mathbf{s}^{[2]}, \mathbf{t}^{[2]})} \end{aligned}$$

实现一个简单的机器翻译系统：学习单词翻译概率(3)

- 很多时候，我们有多组互译句对 $(\mathbf{s}^{[1]}, \mathbf{t}^{[1]}), \dots, (\mathbf{s}^{[n]}, \mathbf{t}^{[n]})$ ，称之为**双语平行数据(语料)**。翻译概率可以被定义为

$$P(x, y) = \frac{\sum_{i=1}^n c(x, y; \mathbf{s}^{[i]}, \mathbf{t}^{[i]})}{\sum_{i=1}^n \sum_{x', y'} c(x', y'; \mathbf{s}^{[i]}, \mathbf{t}^{[i]})}$$

- 说白了就是计算 (x, y) 的频次时，在每个句子上累加
 \mathbf{s}^1 = 机器 翻译 就是 用 计算机 进行 翻译
 \mathbf{t}^1 = Machine translation is just translation by computer
 \mathbf{s}^2 = 那 人工 翻译 呢？
 \mathbf{t}^2 = So , what is human translation ?

$$\begin{aligned} & P(\text{'翻译'}, \text{'translation'}) \\ &= \frac{c(\text{'翻译'}, \text{'translation'}; \mathbf{s}^{[1]}, \mathbf{t}^{[1]}) + c(\text{'翻译'}, \text{'translation'}; \mathbf{s}^{[2]}, \mathbf{t}^{[2]})}{\sum_{x', y'} c(x', y'; \mathbf{s}^{[1]}, \mathbf{t}^{[1]}) + \sum_{x', y'} c(x', y'; \mathbf{s}^{[2]}, \mathbf{t}^{[2]})} \\ &= \frac{4 + 1}{|\mathbf{s}^{[1]}| \times |\mathbf{t}^{[1]}| + |\mathbf{s}^{[2]}| \times |\mathbf{t}^{[2]}|} = \frac{4 + 1}{9 \times 7 + 5 \times 7} = \frac{5}{98} \end{aligned}$$

实现一个简单的机器翻译系统：句子级翻译模型



- **步骤2:** 对任意的句对(s, t)计算句子级翻译概率 $P(t|s)$

实现一个简单的机器翻译系统：句子级翻译模型



- **步骤2:** 对任意的句对(s, t)计算句子级翻译概率 $P(t|s)$
用一种比较简单的思路：定义(s, t)上的一种分数 $g(s, t)$
 - ▶ $g(s, t)$ 的值越大翻译质量越好
 - ▶ $g(s, t)$ 的值越小翻译质量越差

实现一个简单的机器翻译系统：句子级翻译模型



- **步骤2:** 对任意的句对 (s, t) 计算句子级翻译概率 $P(t|s)$
用一种比较简单的思路：定义 (s, t) 上的一种分数 $g(s, t)$
 - ▶ $g(s, t)$ 的值越大翻译质量越好
 - ▶ $g(s, t)$ 的值越小翻译质量越差

于是，我们进一步定义

$$P(t|s) = \frac{g(s, t)}{\sum_{t'} g(s, t')}$$

实际上就是对 $g(s, t)$ 在所有可能的译文集合上做归一化，使其具有概率意义

实现一个简单的机器翻译系统：句子级翻译模型(2)

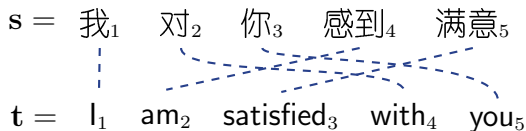
- 两个问题
 - ① 如何计算 $g(\mathbf{s}, \mathbf{t})$
 - ② 如何计算 $\sum_{\mathbf{t}'} g(\mathbf{s}, \mathbf{t}')$

实现一个简单的机器翻译系统：句子级翻译模型(2)

- 两个问题
 - ① 如何计算 $g(\mathbf{s}, \mathbf{t})$ - 最关键的建模问题，马上开始
 - ② 如何计算 $\sum_{\mathbf{t}'} g(\mathbf{s}, \mathbf{t}')$ - 实际上不用计算，后面再说

实现一个简单的机器翻译系统：句子级翻译模型(2)

- 两个问题
 - ① 如何计算 $g(s, t)$ - 最关键的建模问题，马上开始
 - ② 如何计算 $\sum_{t'} g(s, t')$ - 实际上不用计算，后面再说
- 对 $g(s, t)$ 建模：根据本章第一页的假设， s 与 t 之间存在一种单词间的对应，我们称之为词对齐关系



每根虚线代表一个对齐链接，记为 (j, i) ，意思是源语言的第 j 个单词对应目标语言第 i 个单词，即 s_j 与 t_i 对应。

- 所有对齐链接构成集合 A ，对于上例： $A = \{(1, 1), (2, 4), (3, 5), (4, 2), (5, 3)\}$


实现一个简单的机器翻译系统：句子级翻译模型(3)

- 给定一个句对(s, t)，及它们之间的(最优)词对齐 \hat{A} ，可以定义模型得分为：

$$g(s, t) \equiv \prod_{(j,i) \in \hat{A}} P(s_j, t_i)$$

显然每个单词翻译概率都高，那么整句的模型得分也高

s = 我₁ 对₂ 你₃ 感到₄ 满意₅



t = I₁ am₂ satisfied₃ with₄ you₅

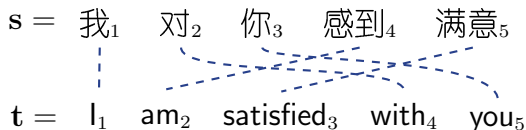
$$g(s, t) = P(\text{'我', 'I'}) \times P(\text{'对', 'with'}) \times P(\text{'你', 'you'}) \times P(\text{'感到', 'am'}) \times P(\text{'满意', 'satisfied'})$$

实现一个简单的机器翻译系统：句子级翻译模型(3)

- 给定一个句对(s, t)，及它们之间的(最优)词对齐 \hat{A} ，可以定义模型得分为：

$$g(s, t) \equiv \prod_{(j,i) \in \hat{A}} P(s_j, t_i)$$

显然每个单词翻译概率都高，那么整句的模型得分也高

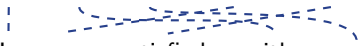



$$g(s, t) = P(\text{'我'}, \text{'I'}) \times P(\text{'对'}, \text{'with'}) \times P(\text{'你'}, \text{'you'}) \times P(\text{'感到'}, \text{'am'}) \times P(\text{'满意'}, \text{'satisfied'})$$

- 其它问题：**如果有单词翻译为空怎么办？
 - 简单的解决方案是平滑，给予这种翻译一个很小的概率值
 - 或者在建模中考虑，见后面的IBM模型

实现一个简单的机器翻译系统：句子级翻译模型(4)

- 但是，这样设计的 $g(s, t)$ 没有考虑词序的信息。相同译词出现在不同的位置，得分相同 - 无法选择流畅的译文

		$\prod_{(j,i) \in \hat{A}} P(s_j, t_i)$
$s =$ 我 ₁ 对 ₂ 你 ₃ 感到 ₄ 满意 ₅		0.0023
$t' =$ I ₁ am ₂ satisfied ₃ with ₄ you ₅		
$s =$ 我 ₁ 对 ₂ 你 ₃ 感到 ₄ 满意 ₅		0.0023
$t'' =$ I ₁ with ₂ you ₃ am ₄ satisfied ₅		

实现一个简单的机器翻译系统：句子级翻译模型(4)

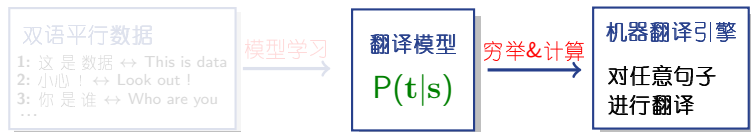
- 但是，这样设计的 $g(s, t)$ 没有考虑词序的信息。相同译词出现在不同的位置，得分相同 - 无法选择流畅的译文

	$\prod_{(j,i) \in \hat{A}} P(s_j, t_i) \times P_{lm}(t)$
$s =$ 我 ₁ 对 ₂ 你 ₃ 感到 ₄ 满意 ₅ $t' =$ I ₁ am ₂ satisfied ₃ with ₄ you ₅	0.0023×0.0107
$s =$ 我 ₁ 对 ₂ 你 ₃ 感到 ₄ 满意 ₅ $t'' =$ I ₁ with ₂ you ₃ am ₄ satisfied ₅	0.0023×0.0009

- 解决方案：引入语言模型 $P_{lm}(t)$ 来度量译文的流畅度
 $P_{2\text{-gram}}(w_1 \dots w_m) =$
 $P(w_1) \times P(w_2|w_1) \times P(w_3|w_2) \dots \times P(w_m|w_{m-1})$
- 最终

$$g(s, t) = \prod_{(j,i) \in \hat{A}} P(s_j, t_i) \times P_{lm}(t)$$

实现一个简单的机器翻译系统：解码

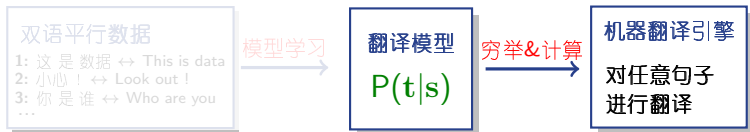


- **步骤3：解码** - 对任意的 s ，找到翻译概率最大的译文 \hat{t}

$$\hat{t} = \arg \max_t P(t|s)$$

这里 $\arg \max_a f(a)$ 表示找到使 $f(a)$ 达到最大的 a 输出

实现一个简单的机器翻译系统：解码



- **步骤3：解码** - 对任意的 s ，找到翻译概率最大的译文 \hat{t}

$$\hat{t} = \arg \max_t P(t|s)$$

这里 $\arg \max_a f(a)$ 表示找到使 $f(a)$ 达到最大的 a 输出

- 现在我们可以对任意的 (s, t) 计算 $P(t|s) = \frac{g(s, t)}{\sum_{t'} g(s, t')}$
 - ▶ 给定 s ， $\sum_{t'} g(s, t')$ 是个常数(因为 $\sum_{t'} g(s, t')$ 的变量只有 s)
 - ▶ 这样，我们得到解码步骤的形式化描述为

$$\begin{aligned}\hat{t} &= \arg \max_t \frac{g(s, t)}{\sum_{t'} g(s, t')} \\ &= \arg \max_t g(s, t)\end{aligned}$$

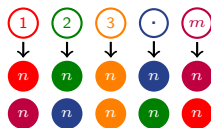
实现一个简单的机器翻译系统：解码(2)

- 解码的核心问题是在所有可能的翻译结果中找到使 $g(s, t)$ 达到最大的译文
 - 若 s 有 m 个词，每个词有 n 个翻译候选 - 共有 n^m 种组合



实现一个简单的机器翻译系统：解码(2)

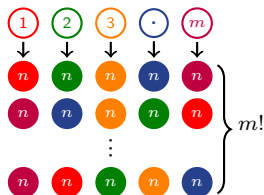
- 解码的核心问题是在所有可能的翻译结果中找到使 $g(s, t)$ 达到最大的译文
 - ▶ 若 s 有 m 个词，每个词有 n 个翻译候选 - 共有 n^m 种组合
 - ▶ 词的翻译候选可以任意调序



实现一个简单的机器翻译系统：解码(2)

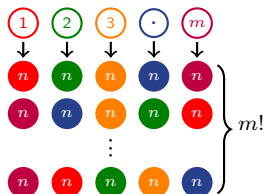
- 解码的核心问题是在所有可能的翻译结果中找到使 $g(s, t)$ 达到最大的译文

- ▶ 若 s 有 m 个词，每个词有 n 个翻译候选 - 共有 n^m 种组合
- ▶ 词的翻译候选可以任意调序
- ▶ s 对应可能的译文至少有 $n^m \cdot m!$



实现一个简单的机器翻译系统：解码(2)

- 解码的核心问题是在所有可能的翻译结果中找到使 $g(s, t)$ 达到最大的译文
 - 若 s 有 m 个词，每个词有 n 个翻译候选 - 共有 n^m 种组合
 - 词的翻译候选可以任意调序
 - s 对应可能的译文至少有 $n^m \cdot m!$
- 潜在翻译结果数量巨大，无法穷尽，因此我们需要高效的搜索算法找到理想的解



m	n	$n^m \cdot m!$
1	1	1
1	10	10
2	10	200
10	10	3628800000000000000
20	10	$2.43290200817664 \times 10^{38}$
20	30	$8.48300477127188 \times 10^{47}$

实现一个简单的机器翻译系统：解码(3)

我们这里使用一种贪婪的解码算法：每轮选择一个源语言单词，利用它的翻译候选扩展译文，保留“最好”的结果，反复进行直到所有源语言单词都被翻译

输入：源句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{true}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

实现一个简单的机器翻译系统：解码(3)

我们这里使用一种**贪婪**的解码算法：每轮选择一个源语言单词，利用它的翻译候选扩展译文，保留“最好”的结果，反复进行直到所有源语言单词都被翻译

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{true}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

获取每个单词
的翻译候选



实现一个简单的机器翻译系统：解码(3)

我们这里使用一种**贪婪**的解码算法：每轮选择一个源语言单词，利用它的翻译候选扩展译文，保留“最好”的结果，反复进行直到所有源语言单词都被翻译

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{true}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

获取每个单词
的翻译候选



$best$ 用于保存当前最好的翻译结果

h 用于保存每步生成的所有译文候选

实现一个简单的机器翻译系统：解码(3)

我们这里使用一种**贪婪**的解码算法：每轮选择一个源语言单词，利用它的翻译候选扩展译文，保留“最好”的结果，反复进行直到所有源语言单词都被翻译

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{true}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

获取每个单词
的翻译候选



$best$ 用于保存当前最好的翻译结果

h 用于保存每步生成的所有译文候选

$\text{JOIN}(a, b)$ 返回 $\begin{bmatrix} a1 \\ a2 \end{bmatrix} \times \begin{bmatrix} b1 \\ b2 \end{bmatrix} = \begin{bmatrix} a1b1 & a1b2 \\ a2b1 & a2b2 \end{bmatrix}$
 a 和 b 的所有组合

实现一个简单的机器翻译系统：解码(3)

我们这里使用一种**贪婪**的解码算法：每轮选择一个源语言单词，利用它的翻译候选扩展译文，保留“最好”的结果，反复进行直到所有源语言单词都被翻译

输入：源语句句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{true}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

获取每个单词
的翻译候选



$best$ 用于保存当前最好的翻译结果

h 用于保存每步生成的所有译文候选

$\text{JOIN}(a, b)$ 返回 $\begin{bmatrix} a1 \\ a2 \end{bmatrix} \times \begin{bmatrix} b1 \\ b2 \end{bmatrix} = \begin{bmatrix} a1b1 & a1b2 \\ a2b1 & a2b2 \end{bmatrix}$
 a 和 b 的所有组合

PRUNEFORTOP1

保留得分最高的结果

0.234	← top1
0.197	
0.083	

实现一个简单的机器翻译系统：解码(3)

我们这里使用一种**贪婪**的解码算法：每轮选择一个源语言单词，利用它的翻译候选扩展译文，保留“最好”的结果，反复进行直到所有源语言单词都被翻译

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{true}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

获取每个单词
的翻译候选



$best$ 用于保存当前最好的翻译结果

h 用于保存每步生成的所有译文候选

JOIN(a, b) 返回
 a 和 b 的所有组合

$$\begin{bmatrix} a1 \\ a2 \end{bmatrix} \times \begin{bmatrix} b1 \\ b2 \end{bmatrix} = \begin{bmatrix} a1b1 & a1b2 \\ a2b1 & a2b2 \end{bmatrix}$$

PRUNEFORTOP1

保留得分最高的结果

0.234	← top1
0.197	
0.083	

记录已经翻译过
的源语单词



实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANSOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

输入：待翻译句子(已经分词)

我 对 你 表示 满意

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

输入：待翻译句子(已经分词)

我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2
am .1	to .1		means .1	it .1
...

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

输入：待翻译句子(已经分词)

我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2
am .1	to .1		means .1	it .1
...

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

输入：待翻译句子(已经分词)

只保留前 m 个结果

我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANSOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

输入：待翻译句子(已经分词)

我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

$best.translation = \phi$

$best.j = -1$

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANSOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

输入：待翻译句子(已经分词)

我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm	to	you	ϕ	satisfy
I	with	your	show	satisfied
me	for		shows	satisfies
.4	.3	.7	.4	.3
.3	.3	.3	.2	.2
.1	.2		.1	.2

$best.translation = \phi$

$best.j = -1$

h存放临时翻译结果

$g(\mathbf{s}, \mathbf{t})$	翻译结果

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

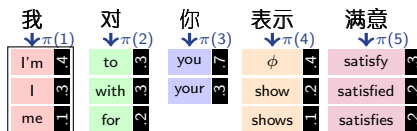
9: $used[best.j] = \text{true}$

10: **return** $best.translation$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)



$best.translation = \phi$

$best.j = -1$

$i = 1, j = 1$

$g(s, t)$ | 翻译结果

.1744
.1674
.1041

I'm
I
me

h存放临时翻译结果

JOIN

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translation$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)

我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm	to	you	ϕ	satisfy
I	with	your	show	satisfied
me	for		shows	satisfies
.4	.3	.7	.4	.3
.3	.3	.3	.2	.2
.1	.2		.1	.2

$best.translation = \phi$
 $best.j = -1$

$i = 1, j = 2$

$g(s, t)$ | 翻译结果

h存放临时翻译结果

.1744
.1674
.1041
.0561
.0013
.0013

I'm
I
me
to
with
for

JOIN

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANSOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

输入：待翻译句子(已经分词)

我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

$best.translation = \phi$

$best.j = -1$

$i = 1, j = 5$

$g(\mathbf{s}, \mathbf{t})$ | 翻译结果

.1744	I'm ← top1
.1674	I
.1041	me
.0561	to
.0013	with
.0013	for
...	...
.1452	satisfies

h存放临时翻译结果

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANSOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

输入：待翻译句子(已经分词)

我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm	to	you	ϕ	satisfy
I	with	your	show	satisfied
me	for		shows	satisfies
.4	.3	.7	.4	.3
.3	.3	.3	.2	.2
.1	.2		.1	.2

$best.translation =$ I'm

$best.j = 1$

$i = 1, j = 5$

$g(\mathbf{s}, \mathbf{t})$ | 翻译结果

.1744	I'm ← top1
.1674	I
.1041	me
.0561	to
.0013	with
.0013	for
...	...
.1452	satisfies

h存放临时翻译结果

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

输入：待翻译句子(已经分词)

used				
我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

$best.translation =$ I'm

$best.j = 1$

$i = 1, j = 5$

$g(s, t)$ | 翻译结果

h存放临时翻译结果

.1744	I'm ← top1
.1674	I
.1041	me
.0561	to
.0013	with
.0013	for
...	...
.1452	satisfies

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

输入：待翻译句子(已经分词)

used				
我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm	to	you	ϕ	satisfy
I	with	your	show	satisfied
me	for		shows	satisfies
.4	.3	.7	.4	.3
.3	.3	.3	.2	.2
.1	.2		.1	.2

$best.translation =$ I'm

$best.j = 1$

$i = 1, j = 5$

$g(s, t)$ 翻译结果

h存放临时翻译结果

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

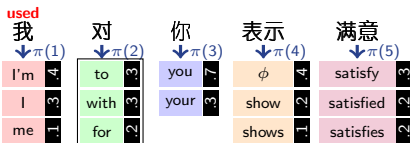
9: $used[best.j] = \text{true}$

10: **return** $best.translation$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)



$best.translation =$ I'm
 $best.j = 1$

$i = 2, j = 2$

$g(s, t)$ | 翻译结果

.0099	I'm	to
.0300	I'm	with
.0231	I'm	for

h存放临时翻译结果



实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translation$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)

used				
我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm	to	you	ϕ	satisfy
I	with	your	show	satisfied
me	for		shows	satisfies
.4	.3	.7	.4	.3
.3	.3	.3	.2	.2
.1	.2		.1	.2

$best.translation =$

I'm

$best.j = 1$

$i = 2, j = 5$

$g(s, t)$ 翻译结果

h存放临时翻译结果

.0099	I'm	to
.0300	I'm	with
.0231	I'm	for
...	...	
.0121	I'm	satisfy
.0314	I'm	satisfied
.0081	I'm	satisfies



实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)

used				
我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

$best.translation =$ I'm

$best.j = 1$

$i = 2, j = 5$

$g(s, t)$ | 翻译结果

h存放临时翻译结果	.0099	I'm to	
	.0300	I'm with	
	.0231	I'm for	
	
	.0121	I'm satisfy	
	.0314	I'm satisfied	← top1
	.0081	I'm satisfies	

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)

used				
我	对	你	表示	满意
↓ $\pi(1)$	↓ $\pi(2)$	↓ $\pi(3)$	↓ $\pi(4)$	↓ $\pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

$best.translation =$ I'm satisfied

$best.j = 5$

$i = 2, j = 5$

$g(s, t)$ | 翻译结果

h存放临时翻译结果	.0099	I'm to	
	.0300	I'm with	
	.0231	I'm for	
	
	.0121	I'm satisfy	
	.0314	I'm satisfied	← top1
	.0081	I'm satisfies	

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANSOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

输入：待翻译句子(已经分词)

used				used
我	对	你	表示	满意
↓ $\pi(1)$	↓ $\pi(2)$	↓ $\pi(3)$	↓ $\pi(4)$	↓ $\pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

$best.translation =$ I'm satisfied

$best.j = 5$

$i = 2, j = 5$

$g(\mathbf{s}, \mathbf{t})$ | 翻译结果

h存放临时翻译结果	.0099	I'm to	
	.0300	I'm with	
	.0231	I'm for	
	
	.0121	I'm satisfy	
	.0314	I'm satisfied	← top1
	.0081	I'm satisfies	

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

输入：待翻译句子(已经分词)

^{used} 我 ↓ $\pi(1)$	对 ↓ $\pi(2)$	你 ↓ $\pi(3)$	表示 ↓ $\pi(4)$	^{used} 满意 ↓ $\pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

$best.translation =$ I'm satisfied

$best.j = 5$

$i = 2, j = 5$

$g(s, t)$ | 翻译结果

h存放临时翻译结果

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

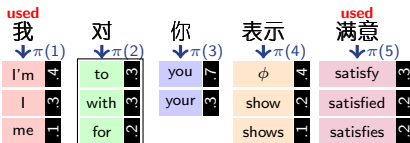
9: $used[best.j] = \text{true}$

10: **return** $best.translation$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)



$best.translation =$ I'm satisfied
 $best.j = 5$

$i = 3, j = 2$



h存放临时翻译结果

$g(s, t)$ | 翻译结果

.0082	I'm satisfied	to
.0103	I'm satisfied	with
.0073	I'm satisfied	for

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

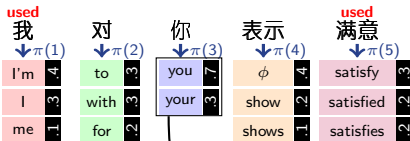
9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)



$best.translation =$

$best.j = 5$

$i = 3, j = 3$

$g(s, t)$ 翻译结果

h存放临时翻译结果

.0082	I'm satisfied	to
.0103	I'm satisfied	with
.0073	I'm satisfied	for
.0033	I'm satisfied	you
.0039	I'm satisfied	your



实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANSOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

输入：待翻译句子(已经分词)

used				used
我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

$best.translation =$

I'm satisfied

$best.j = 5$

$i = 3, j = 4$



h存放临时翻译结果

$g(\mathbf{s}, \mathbf{t})$ 翻译结果

.0082	I'm satisfied	to
.0103	I'm satisfied	with
.0073	I'm satisfied	for
.0033	I'm satisfied	you
.0039	I'm satisfied	your
.0106	I'm satisfied	ϕ
.0035	I'm satisfied	show
.0028	I'm satisfied	shows

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANSOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)

used				used
我	对	你	表示	满意
↓ $\pi(1)$	↓ $\pi(2)$	↓ $\pi(3)$	↓ $\pi(4)$	↓ $\pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

$best.translation =$ I'm satisfied

$best.j = 5$

$i = 3, j = 4$

$g(\mathbf{s}, \mathbf{t})$ | 翻译结果

h存放临时翻译结果

.0082	I'm satisfied	to
.0103	I'm satisfied	with
.0073	I'm satisfied	for
.0033	I'm satisfied	you
.0039	I'm satisfied	your
.0106	I'm satisfied	ϕ ← top1
.0035	I'm satisfied	show
.0028	I'm satisfied	shows

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)

^{used} 我 ↓ $\pi(1)$	对 ↓ $\pi(2)$	你 ↓ $\pi(3)$	表示 ↓ $\pi(4)$	^{used} 满意 ↓ $\pi(5)$
I'm 0.4	to 0.3	you 0.7	ϕ 0.4	satisfy 0.3
I 0.3	with 0.3	your 0.3	show 0.2	satisfied 0.2
me 0.1	for 0.2		shows 0.1	satisfies 0.2

$best.translation =$ I'm satisfied

$best.j = 4$

$i = 3, j = 4$

$g(s, t)$ | 翻译结果

h存放临时翻译结果

.0082	I'm satisfied	to
.0103	I'm satisfied	with
.0073	I'm satisfied	for
.0033	I'm satisfied	you
.0039	I'm satisfied	your
.0106	I'm satisfied	ϕ ← top1
.0035	I'm satisfied	show
.0028	I'm satisfied	shows

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)

used				used	used
我	对	你	表示	满意	
↓ $\pi(1)$	↓ $\pi(2)$	↓ $\pi(3)$	↓ $\pi(4)$	↓ $\pi(5)$	
I'm	to	you	ϕ	satisfy	
I	with	your	show	satisfied	
me	for		shows	satisfies	

$best.translation =$ I'm satisfied

$best.j = 4$

$i = 3, j = 4$

$g(s, t)$ | 翻译结果

h存放临时翻译结果

.0082	I'm satisfied	to
.0103	I'm satisfied	with
.0073	I'm satisfied	for
.0033	I'm satisfied	you
.0039	I'm satisfied	your
.0106	I'm satisfied	ϕ ← top1
.0035	I'm satisfied	show
.0028	I'm satisfied	shows

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANSOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

输入：待翻译句子(已经分词)

used				used	used
我	对	你	表示	满意	
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$	
I'm	to	you	ϕ	satisfy	
I	with	your	show	satisfied	
me	for		shows	satisfies	
.4	.3	.7	.4	.3	
.3	.3	.3	.2	.2	
.1	.2		.1	.2	

$best.translation =$ I'm satisfied

$best.j = 4$

$i = 3, j = 4$

$g(\mathbf{s}, \mathbf{t})$ 翻译结果

h存放临时翻译结果

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

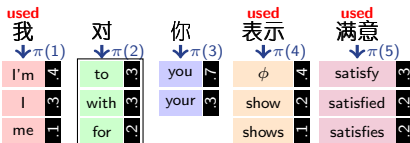
9: $used[best.j] = \text{true}$

10: **return** $best.translation$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)



$best.translation =$ I'm satisfied

$best.j = 4$

$i = 4, j = 2$

$g(s, t)$ 翻译结果

h存放临时翻译结果

.0032	I'm satisfied	to
.0040	I'm satisfied	with
.0028	I'm satisfied	for



实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

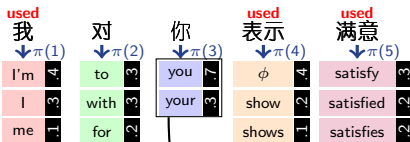
9: $used[best.j] = \text{true}$

10: **return** $best.translation$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)



$best.translation =$

I'm satisfied

$best.j = 4$

$i = 4, j = 3$

$g(\mathbf{s}, \mathbf{t})$ | 翻译结果

h存放临时翻译结果

.0032	I'm satisfied	to
.0040	I'm satisfied	with
.0028	I'm satisfied	for
.0012	I'm satisfied	you
.0014	I'm satisfied	your



实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)

used			used	used
我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

$best.translation =$ I'm satisfied

$best.j = 4$

$i = 4, j = 3$

$g(s, t)$ 翻译结果

h存放临时翻译结果

.0032	I'm satisfied	to	
.0040	I'm satisfied	with	← top1
.0028	I'm satisfied	for	
.0012	I'm satisfied	you	
.0014	I'm satisfied	your	

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANSOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)

used			used	used
我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

$best.translation =$ I'm satisfied with

$best.j = 2$

$i = 4, j = 3$

$g(\mathbf{s}, \mathbf{t})$ 翻译结果

h存放临时翻译结果

.0032	I'm satisfied	to	
.0040	I'm satisfied	with	← top1
.0028	I'm satisfied	for	
.0012	I'm satisfied	you	
.0014	I'm satisfied	your	

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)

used	used		used	used
我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

$best.translation =$ I'm satisfied with

$best.j = 2$

$i = 4, j = 3$

$g(\mathbf{s}, \mathbf{t})$ 翻译结果

h存放临时翻译结果

.0032	I'm satisfied	to	
.0040	I'm satisfied	with	← top1
.0028	I'm satisfied	for	
.0012	I'm satisfied	you	
.0014	I'm satisfied	your	

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANSOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

输入：待翻译句子(已经分词)

used	used		used	used
我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm	to	you	ϕ	satisfy
I	with	your	show	satisfied
me	for		shows	satisfies
.4	.3	.7	.4	.3
.3	.3	.3	.2	.2
.1	.2		.1	.2

$best.translation =$ I'm satisfied with

$best.j = 2$

$i = 4, j = 3$

$g(\mathbf{s}, \mathbf{t})$ 翻译结果

h存放临时翻译结果

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANSOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

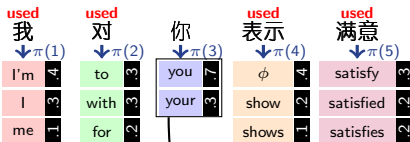
9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)



$best.translation = \begin{bmatrix} \text{'I'm} & \text{satisfied} & \text{with} \end{bmatrix}$
 $best.j = 2$

$i = 5, j = 3$

$g(s, t)$ | 翻译结果

.0011 | $\begin{bmatrix} \text{'I'm} & \text{satisfied} & \text{with} & \text{you} \end{bmatrix}$
.0010 | $\begin{bmatrix} \text{'I'm} & \text{satisfied} & \text{with} & \text{your} \end{bmatrix}$

h存放临时翻译结果



实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度： $O(m^2 \cdot n)$

空间复杂度： $O(m \cdot n)$

输入：待翻译句子(已经分词)

used	used		used	used
我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

$best.translation =$ I'm satisfied with
 $best.j = 2$

$i = 5, j = 3$

$g(s, t)$ | 翻译结果

.0011	I'm satisfied with	you	$\leftarrow \text{top1}$
.0010	I'm satisfied with	your	

h存放临时翻译结果

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $\mathbf{s} = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(\mathbf{s})

1: $\pi = \text{GETTRANSOPTIONS}(\mathbf{s})$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

输入：待翻译句子(已经分词)

used	used		used	used
我	对	你	表示	满意
$\downarrow \pi(1)$	$\downarrow \pi(2)$	$\downarrow \pi(3)$	$\downarrow \pi(4)$	$\downarrow \pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

$best.translation =$ I'm satisfied with you
 $best.j = 3$

$i = 5, j = 3$

$g(\mathbf{s}, \mathbf{t})$ | 翻译结果

.0011	I'm satisfied with	you	← top1
.0010	I'm satisfied with	your	

h存放临时翻译结果

实现一个简单的机器翻译系统：解码(4)

来看一个解码过程的运行实例

输入：源语句子 $s = s_1 \dots s_m$

输出：找的最佳译文

Function WORDDECODING(s)

1: $\pi = \text{GETTRANOPTIONS}(s)$

2: $best = \phi$

3: **for** i in $[1, m]$ **do**

4: $h = \phi$

5: **foreach** j in $[1, m]$ **do**

6: **if** $used[j] = \text{false}$ **then**

7: $h = h \cup \text{JOIN}(best, \pi[j])$

8: $best = \text{PRUNEFORTOP1}(h)$

9: $used[best.j] = \text{true}$

10: **return** $best.translatoin$

时间复杂度: $O(m^2 \cdot n)$

空间复杂度: $O(m \cdot n)$

输入：待翻译句子(已经分词)

used 我 ↓ $\pi(1)$	used 对 ↓ $\pi(2)$	used 你 ↓ $\pi(3)$	used 表示 ↓ $\pi(4)$	used 满意 ↓ $\pi(5)$
I'm .4	to .3	you .7	ϕ .4	satisfy .3
I .3	with .3	your .3	show .2	satisfied .2
me .1	for .2		shows .1	satisfies .2

$best.translation =$ I'm satisfied with you
 $best.j = 3$

$i = 5, j = 3$

$g(s, t)$ | 翻译结果

.0011	I'm satisfied with you ← top1
.0010	I'm satisfied with your

h存放临时翻译结果

即将开始正式的内容

- 基于词的统计机器翻译很简单 - 建议实现一下前面所描述的方法，确实不难
 - ▶ 是的，思想很简单
 - ▶ 不对，前面的例子主要面向实现，还有很多问题没有回答
 - 建模：还需要更严密的数学模型来描述翻译过程
 - 训练：需要有一个科学的优化方法
 - 解码：是否还有其它更系统的解码方法？A*, Bottom-up decoding?
 - 其它：空翻译问题、调序建模问题等等

即将开始正式的内容

- 基于词的统计机器翻译很简单 - 建议实现一下前面所描述的方法，确实不难
 - ▶ 是的，思想很简单
 - ▶ 不对，前面的例子主要面向实现，还有很多问题没有回答
 - 建模：还需要更严密的数学模型来描述翻译过程
 - 训练：需要有一个科学的优化方法
 - 解码：是否还有其它更系统的解码方法？ A^* ，Bottom-up decoding？
 - 其它：空翻译问题、调序建模问题等等
- 下面，继续对一些问题进行进一步介绍及分析
 - ① 统计机器翻译基础模型和IBM模型
 - 统计机器翻译模型基础(生成模型)
 - IBM Models 1-2的建模及训练
 - IBM Models 3-5的建模
 - ② 统计词对齐
 - 基于HMMs的词对齐
 - 基于判别式模型的词对齐

Outline

1. 一个简单的翻译实例

2. IBM模型

- 建模
- 解码
- 模型训练

统计机器翻译

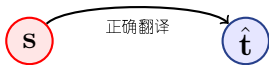
- 统计机器翻译源于1941年的Warren Weaver的思想，而正式开始于
 - ▶ **A Program for Aligning Sentences in Bilingual Corpora**
William A. Gale; Kenneth W. Church. 1993.
 - ▶ **The Mathematics of Statistical Machine Translation: Parameter Estimation**
Peter E. Brown; Stephen A. Della Pietra; Vincent J. Della Pietra; Robert L. Mercer. 1993.
- 特别是IBM Waston研究中心Brown等人的基于词的统计机器翻译模型，成为机器翻译领域至今**经典中的经典**

Brown等人的工作(IBM Waston) - on Wikipedia

The first ideas of statistical machine translation were introduced by Warren Weaver in 1949, including the ideas of applying Claude Shannon's information theory. Statistical machine translation was re-introduced in 1991 by researchers at IBM's Thomas J. Watson Research Center and has contributed to the significant resurgence in interest in machine translation in recent years. Nowadays it is by far the most widely-studied machine translation method

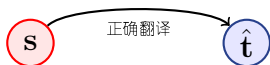
机器翻译的统计建模

- 一个人在做翻译时：对于给定的源语言句子 s ，可以了翻译为一个（或者若干个）正确的译文 \hat{t}
 - ▶ 也就是说除了正确的译文，其它的翻译都是不正确的

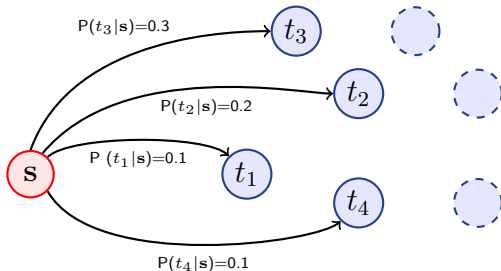


机器翻译的统计建模

- 一个人在做翻译时：对于给定的源语言句子 s ，可以了翻译为一个（或者若干个）正确的译文 \hat{t}
 - 也就是说除了正确的译文，其它的翻译都是不正确的

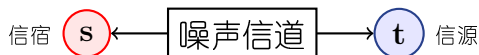


- 统计机器翻译的思想是：对于 s ，所有可能的目标语词串 t 都是可能的译文。每一对 (s, t) 都有一个概率值 $P(t|s)$ 来描述 s 翻译为 t 的好与坏



噪声信道模型

- **噪声信道模型**：源语言句子 s (信宿)是由目标语句子 t (信源)经过一个有噪声的信道得到的。如果知道了 s 和信道的性质，我们可以通过 $P(t|s)$ 得到可能的信源的概率。

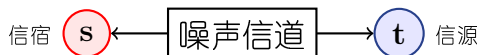


而通过上述过程找到最可能的信源的过程被称之为**解码**

$$\hat{t} = \arg \max_t P(t|s)$$

噪声信道模型

- **噪声信道模型**：源语言句子 s (信宿)是由目标语句子 t (信源)经过一个有噪声的信道得到的。如果知道了 s 和信道的性质，我们可以通过 $P(t|s)$ 得到可能的信源的概率。



而通过上述过程找到最可能的信源的过程被称之为**解码**

$$\hat{t} = \arg \max_t P(t|s)$$

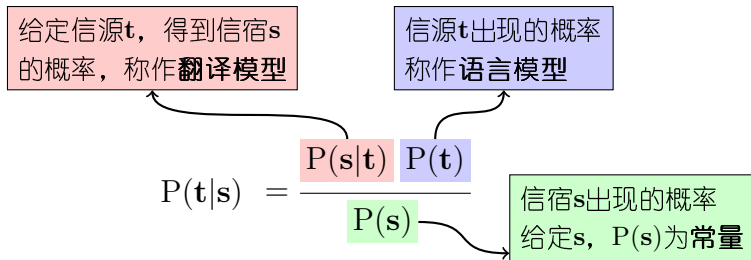
- 贝叶斯变换

$$\begin{aligned} P(t|s) &= \frac{P(s, t)}{P(s)} \\ &= \frac{P(s|t)P(t)}{P(s)} \end{aligned}$$

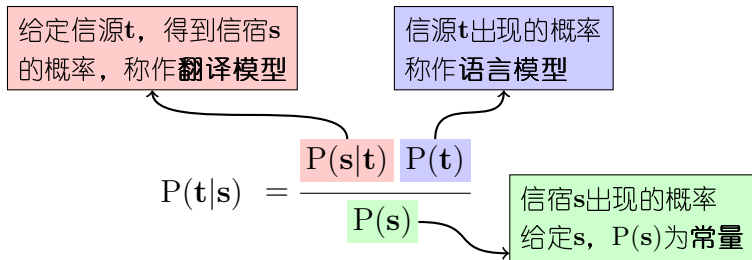
噪声信道模型(2)

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

噪声信道模型(2)



噪声信道模型(2)



- 因此，翻译问题可以被重新描述为

$$\begin{aligned}\hat{t} &= \arg \max_t \frac{P(s|t)P(t)}{P(s)} \\ &= \arg \max_t P(s|t)P(t)\end{aligned}$$

即，在所有可能的译文中找到使翻译模型 $P(s|t)$ 和语言模型 $P(t)$ 乘积最大的译文

基本问题

$$\hat{\mathbf{t}} = \arg \max_{\mathbf{t}} P(\mathbf{s}|\mathbf{t})P(\mathbf{t})$$

- 三个基本问题

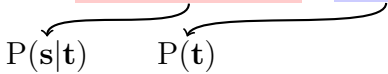
- ① 建模：如何描述计算 $P(\mathbf{s}|\mathbf{t})$ 和 $P(\mathbf{t})$ 的计算方式
- ② 训练：如何获得计算 $P(\mathbf{s}|\mathbf{t})$ 和 $P(\mathbf{t})$ 所需的参数
- ③ 解码：如何完成搜索最优解的过程 *argmax*

基本问题

$$\hat{\mathbf{t}} = \arg \max_{\mathbf{t}} P(\mathbf{s}|\mathbf{t})P(\mathbf{t})$$

• 三个基本问题

- ① 建模：如何描述计算 $P(\mathbf{s}|\mathbf{t})$ 和 $P(\mathbf{t})$ 的计算方式
 - ② 训练：如何获得计算 $P(\mathbf{s}|\mathbf{t})$ 和 $P(\mathbf{t})$ 所需的参数
 - ③ 解码：如何完成搜索最优解的过程 $argmax$
- 回忆一下本章开始的实例，是不是有似曾相识的感觉？

$$g(\mathbf{s}, \mathbf{t}) = \prod_{(j,i) \in \hat{A}} P(s_j, t_i) \times P_{lm}(\mathbf{t})$$


$P(\mathbf{s}|\mathbf{t})$ $P(\mathbf{t})$
翻译模型 语言模型

- ▶ 词汇翻译概率和ngram概率使用相对频率估计(第7页)
- ▶ $argmax$ 使用第16页所描述的解码算法

IBM模型中的假设 - 词对齐

- $P(t)$ 和解码在前面的内容中有介绍，下面重点求解 $P(s|t)$ ，即：
 - ▶ 翻译模型建模 - $P(s|t)$ 的计算方法
 - ▶ 翻译模型参数估计 - 计算 $P(s|t)$ 所需的参数

IBM模型中的假设 - 词对齐

- $P(t)$ 和解码在前面的内容中有介绍，下面重点求解 $P(s|t)$ ，即：
 - ▶ 翻译模型建模 - $P(s|t)$ 的计算方法
 - ▶ 翻译模型参数估计 - 计算 $P(s|t)$ 所需的参数
- IBM模型的假设： $s = s_1...s_m$ 和 $t = t_1...t_n$ 之间有单词一级的对应，称作**单词对齐**或者**词对齐**。此外：
 - ▶ 约束：一个源语言单词只能对应一个目标语单词

谢谢 你
| |
thank you

Yes

谢谢 你
| |
thank you

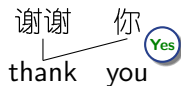
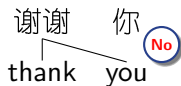
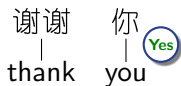
No

谢谢 你
| |
thank you

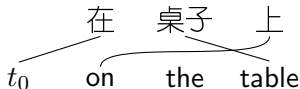
Yes

IBM模型中的假设 - 词对齐

- $P(t)$ 和解码在前面的内容中有介绍，下面重点求解 $P(s|t)$ ，即：
 - ▶ 翻译模型建模 - $P(s|t)$ 的计算方法
 - ▶ 翻译模型参数估计 - 计算 $P(s|t)$ 所需的参数
- IBM模型的假设： $s = s_1...s_m$ 和 $t = t_1...t_n$ 之间有单词一级的对应，称作**单词对齐**或者**词对齐**。此外：
 - ▶ 约束：一个源语言单词只能对应一个目标语单词

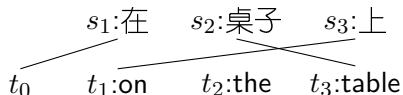


- ▶ **空翻译(删词)**：源语言单词可以翻空，这时它对应到一个虚拟的目标语单词 t_0



建模 - $P(s|t)$

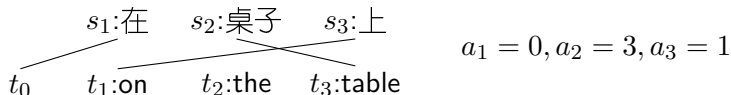
- 给定 s 和 t ，它们之间的词对齐被记为 $\mathbf{a} = a_1 \dots a_m$
 - a_j 表示第 j 个源语单词 s_j 对应的目标语单词的位置



$$a_1 = 0, a_2 = 3, a_3 = 1$$

建模 - $P(s|t)$

- 给定 s 和 t ，它们之间的词对齐被记为 $\mathbf{a} = a_1 \dots a_m$
 - a_j 表示第 j 个源语单词 s_j 对应的目标语单词的位置

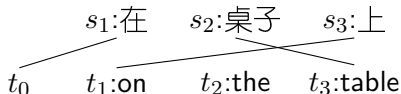


- 建模!!:** $P(s|t)$ 被表示为所有可能的词对齐的生成概率

$$P(s|t) = \sum_{\mathbf{a}} P(s, \mathbf{a}|t)$$

建模 - $P(s|t)$

- 给定 s 和 t ，它们之间的词对齐被记为 $\mathbf{a} = a_1 \dots a_m$
 - a_j 表示第 j 个源语单词 s_j 对应的目标语单词的位置



$$a_1 = 0, a_2 = 3, a_3 = 1$$

- 建模!!:** $P(s|t)$ 被表示为所有可能的词对齐的生成概率

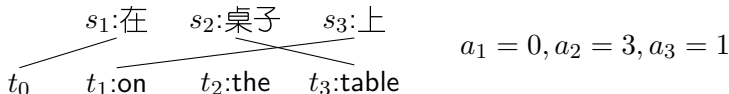
$$P(s|t) = \sum_{\mathbf{a}} P(s, \mathbf{a}|t)$$

每一种 \mathbf{a} 对应一个 $P(s, \mathbf{a}|t)$



建模 - $P(s|t)$

- 给定 s 和 t ，它们之间的词对齐被记为 $\mathbf{a} = a_1 \dots a_m$
 - a_j 表示第 j 个源语单词 s_j 对应的目标语单词的位置



- 建模!!:** $P(s|t)$ 被表示为所有可能的词对齐的生成概率

$$P(s|t) = \sum_{\mathbf{a}} P(s, \mathbf{a}|t)$$

每一种 \mathbf{a} 对应一个 $P(s, \mathbf{a}|t)$

$$P\left(\begin{array}{ccc} \text{谢谢} & \text{你} & \\ \swarrow & \nearrow & \\ t_0 & \text{thank} & \text{you} \end{array}\right) + P\left(\begin{array}{ccc} \text{谢谢} & \text{你} & \\ \swarrow & \nearrow & \\ t_0 & \text{thank} & \text{you} \end{array}\right) + P\left(\begin{array}{ccc} \text{谢谢} & \text{你} & \\ \swarrow & \nearrow & \\ t_0 & \text{thank} & \text{you} \end{array}\right) +$$

$$P\left(\begin{array}{ccc} \text{谢谢} & \text{你} & \\ \swarrow & \nearrow & \\ t_0 & \text{thank} & \text{you} \end{array}\right) + P\left(\begin{array}{ccc} \text{谢谢} & \text{你} & \\ \swarrow & \nearrow & \\ t_0 & \text{thank} & \text{you} \end{array}\right) + P\left(\begin{array}{ccc} \text{谢谢} & \text{你} & \\ \swarrow & \nearrow & \\ t_0 & \text{thank} & \text{you} \end{array}\right) +$$

$$P\left(\begin{array}{ccc} \text{谢谢} & \text{你} & \\ \swarrow & \nearrow & \\ t_0 & \text{thank} & \text{you} \end{array}\right) + P\left(\begin{array}{ccc} \text{谢谢} & \text{你} & \\ \swarrow & \nearrow & \\ t_0 & \text{thank} & \text{you} \end{array}\right) + P\left(\begin{array}{ccc} \text{谢谢} & \text{你} & \\ \swarrow & \nearrow & \\ t_0 & \text{thank} & \text{you} \end{array}\right) = P(s|t)$$

建模 - $P(\mathbf{s}, \mathbf{a}|\mathbf{t})$

- **进一步建模!!**: 对于源语句子 $\mathbf{s} = s_1 \dots s_m$ (m 个词)、目标语译文 $\mathbf{t} = t_0 \dots t_n$ (n 个词)和词对齐 $\mathbf{a} = a_1 \dots a_m$, 按如下方式计算 $P(\mathbf{s}, \mathbf{a}|\mathbf{t})$

► 符号定义: $s_x^y = s_x \dots s_y$, $a_x^y = a_x \dots a_y$

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

建模 - $P(s, a|t)$

- **进一步建模!!**: 对于源语句子 $s = s_1 \dots s_m$ (m 个词)、目标语译文 $t = t_0 \dots t_n$ (n 个词) 和词对齐 $a = a_1 \dots a_m$, 按如下方式计算 $P(s, a|t)$

► 符号定义: $s_x^y = s_x \dots s_y$, $a_x^y = a_x \dots a_y$

$$P(s, a|t) = \underbrace{P(m|t)}_{\textcircled{1}} \prod_{j=1}^m P(a_j | a_1^{j-1}, s_1^{j-1}, m, t) P(s_j | a_1^j, s_1^{j-1}, m, t)$$

► **生成模型**: 给定译文 t 生成源文 s 和对齐 a

① 根据译文 t 选择源文的长度 m

建模 - $P(s, a|t)$

- **进一步建模!!**: 对于源语句子 $s = s_1 \dots s_m$ (m 个词)、目标语译文 $t = t_0 \dots t_n$ (n 个词) 和词对齐 $a = a_1 \dots a_m$, 按如下方式计算 $P(s, a|t)$

► 符号定义: $s_x^y = s_x \dots s_y$, $a_x^y = a_x \dots a_y$

$$P(s, a|t) = \underbrace{P(m|t)}_{\textcircled{1}} \underbrace{\prod_{j=1}^m}_{\textcircled{2}} P(a_j | a_1^{j-1}, s_1^{j-1}, m, t) P(s_j | a_1^j, s_1^{j-1}, m, t)$$

► **生成模型**: 给定译文 t 生成源文 s 和对齐 a

- ① 根据译文 t 选择源文的长度 m
- ② 循环源文的每个位置 j

建模 - $P(s, a|t)$

- **进一步建模!!**: 对于源语句子 $s = s_1 \dots s_m$ (m 个词)、目标语译文 $t = t_0 \dots t_n$ (n 个词) 和词对齐 $a = a_1 \dots a_m$, 按如下方式计算 $P(s, a|t)$

► 符号定义: $s_x^y = s_x \dots s_y$, $a_x^y = a_x \dots a_y$

$$P(s, a|t) = \underbrace{P(m|t)}_{\textcircled{1}} \underbrace{\prod_{j=1}^m}_{\textcircled{2}} \underbrace{P(a_j|a_1^{j-1}, s_1^{j-1}, m, t)}_{\textcircled{3}} P(s_j|a_1^j, s_1^{j-1}, m, t)$$

► **生成模型**: 给定译文 t 生成源文 s 和对齐 a

- ① 根据译文 t 选择源文的长度 m
- ② 循环源文的每个位置 j
- ③ 根据译文 t 、源文长度 m 、已经生成的源语单词 s_1^{j-1} 和对齐 a_1^{j-1} , 生成第 j 个位置的对齐结果 a_j

建模 - $P(s, a|t)$

- **进一步建模!!**: 对于源语句子 $s = s_1 \dots s_m$ (m 个词)、目标语译文 $t = t_0 \dots t_n$ (n 个词) 和词对齐 $a = a_1 \dots a_m$, 按如下方式计算 $P(s, a|t)$

► 符号定义: $s_x^y = s_x \dots s_y$, $a_x^y = a_x \dots a_y$

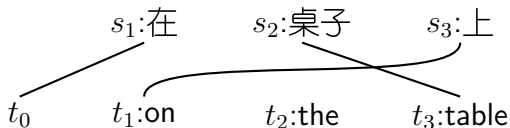
$$P(s, a|t) = \underbrace{P(m|t)}_{\textcircled{1}} \underbrace{\prod_{j=1}^m}_{\textcircled{2}} \underbrace{P(a_j|a_1^{j-1}, s_1^{j-1}, m, t)}_{\textcircled{3}} \underbrace{P(s_j|a_1^j, s_1^{j-1}, m, t)}_{\textcircled{4}}$$

► **生成模型**: 给定译文 t 生成源文 s 和对齐 a

- ① 根据译文 t 选择源文的长度 m
- ② 循环源文的每个位置 j
- ③ 根据译文 t 、源文长度 m 、已经生成的源语单词 s_1^{j-1} 和对齐 a_1^{j-1} , 生成第 j 个位置的对齐结果 a_j
- ④ 根据译文 t 、源文长度 m 、已经生成的源语单词 s_1^{j-1} 和对齐 a_1^j , 生成第 j 个位置的源语言单词 s_j (注意: 这时 a_j 已经生成了)

实例 - $P(\mathbf{s}, \mathbf{a}|\mathbf{t})$

\mathbf{s} = 在 桌子 上 $\mathbf{t} = t_0$ on the table $\mathbf{a} = \{1-0, 2-3, 3-1\}$



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

实例 - $P(\mathbf{s}, \mathbf{a}|\mathbf{t})$

\mathbf{s} = 在 桌子 上 $\mathbf{t} = t_0$ on the table $\mathbf{a} = \{1-0, 2-3, 3-1\}$

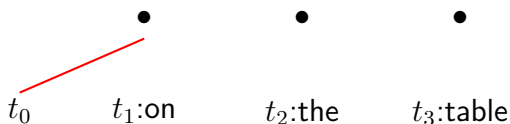
• • •

t_0 t_1 :on t_2 :the t_3 :table

$$\begin{aligned} P(\mathbf{s}, \mathbf{a}|\mathbf{t}) &= P(m|\mathbf{t}) \prod_{j=1}^m P(a_j|a_1^{j-1}, s_1^{j-1}, m, \mathbf{t}) P(s_j|a_1^j, s_1^{j-1}, m, \mathbf{t}) \\ &= P(m = 3 \mid 't_0 \text{ on the table}')$$

实例 - $P(s, a|t)$

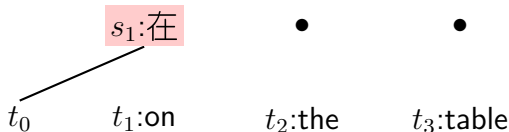
s = 在 桌子 上 $t = t_0$ on the table $a = \{1-0, 2-3, 3-1\}$



$$\begin{aligned} P(s, a|t) &= P(m|t) \prod_{j=1}^m P(a_j | a_1^{j-1}, s_1^{j-1}, m, t) P(s_j | a_1^j, s_1^{j-1}, m, t) \\ &= P(m = 3 \mid 't_0 \text{ on the table} ') \times \\ &\quad P(a_1 = 0 \mid \phi, \phi, 3, 't_0 \text{ on the table} ') \end{aligned}$$

实例 - $P(s, a|t)$

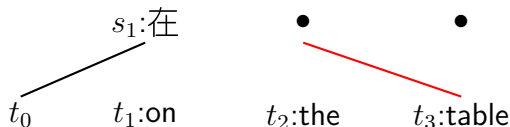
s = 在 桌子 上 $t = t_0$ on the table $a = \{1-0, 2-3, 3-1\}$



$$\begin{aligned} P(s, a|t) &= P(m|t) \prod_{j=1}^m P(a_j | a_1^{j-1}, s_1^{j-1}, m, t) P(s_j | a_1^j, s_1^{j-1}, m, t) \\ &= P(m = 3 \mid 't_0 \text{ on the table} ') \times \\ &\quad P(a_1 = 0 \mid \phi, \phi, 3, 't_0 \text{ on the table} ') \times \\ &\quad P(f_1 = \text{在} \mid \{1-0\}, \phi, 3, 't_0 \text{ on the table} ') \end{aligned}$$

实例 - $P(\mathbf{s}, \mathbf{a}|\mathbf{t})$

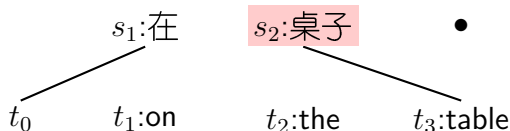
\mathbf{s} = 在 桌子 上 $\mathbf{t} = t_0$ on the table $\mathbf{a} = \{1-0, 2-3, 3-1\}$



$$\begin{aligned} P(\mathbf{s}, \mathbf{a}|\mathbf{t}) &= P(m|\mathbf{t}) \prod_{j=1}^m P(a_j|a_1^{j-1}, s_1^{j-1}, m, \mathbf{t}) P(s_j|a_1^j, s_1^{j-1}, m, \mathbf{t}) \\ &= P(m = 3 \mid 't_0 \text{ on the table} ') \times \\ &\quad P(a_1 = 0 \mid \phi, \phi, 3, 't_0 \text{ on the table} ') \times \\ &\quad P(f_1 = \text{在} \mid \{1-0\}, \phi, 3, 't_0 \text{ on the table} ') \times \\ &\quad P(a_2 = 3 \mid \{1-0\}, \text{在}, 3, 't_0 \text{ on the table} ') \end{aligned}$$

实例 - $P(s, a|t)$

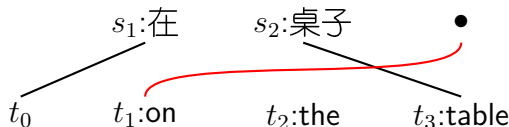
s = 在 桌子 上 $t = t_0$ on the table $a = \{1-0, 2-3, 3-1\}$



$$\begin{aligned} P(s, a|t) &= P(m|t) \prod_{j=1}^m P(a_j | a_1^{j-1}, s_1^{j-1}, m, t) P(s_j | a_1^j, s_1^{j-1}, m, t) \\ &= P(m = 3 \mid 't_0 \text{ on the table} ') \times \\ &\quad P(a_1 = 0 \mid \phi, \phi, 3, 't_0 \text{ on the table} ') \times \\ &\quad P(f_1 = \text{在} \mid \{1-0\}, \phi, 3, 't_0 \text{ on the table} ') \times \\ &\quad P(a_2 = 3 \mid \{1-0\}, \text{'在'}, 3, 't_0 \text{ on the table} ') \times \\ &\quad P(f_2 = \text{桌子} \mid \{1-0, 2-3\}, \text{'在'}, 3, 't_0 \text{ on the table} ') \end{aligned}$$

实例 - $P(s, a|t)$

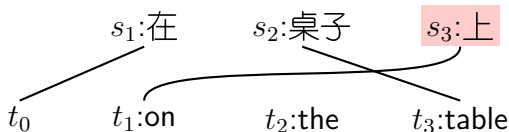
s = 在 桌子 上 $t = t_0$ on the table $a = \{1-0, 2-3, 3-1\}$



$$\begin{aligned}
 P(s, a|t) &= P(m|t) \prod_{j=1}^m P(a_j | a_1^{j-1}, s_1^{j-1}, m, t) P(s_j | a_1^j, s_1^{j-1}, m, t) \\
 &= P(m = 3 \mid 't_0 \text{ on the table}') \times \\
 &\quad P(a_1 = 0 \mid \phi, \phi, 3, 't_0 \text{ on the table}') \times \\
 &\quad P(f_1 = \text{在} \mid \{1-0\}, \phi, 3, 't_0 \text{ on the table}') \times \\
 &\quad P(a_2 = 3 \mid \{1-0\}, \text{'在'}, 3, 't_0 \text{ on the table}') \times \\
 &\quad P(f_2 = \text{桌子} \mid \{1-0, 2-3\}, \text{'在'}, 3, 't_0 \text{ on the table}') \times \\
 &\quad P(a_3 = 1 \mid \{1-0, 2-3\}, \text{'在 桌子'}, 3, 't_0 \text{ on the table}')
 \end{aligned}$$

实例 - $P(s, a|t)$

s = 在 桌子 上 $t = t_0$ on the table $a = \{1-0, 2-3, 3-1\}$



$$\begin{aligned}
 P(s, a|t) &= P(m|t) \prod_{j=1}^m P(a_j | a_1^{j-1}, s_1^{j-1}, m, t) P(s_j | a_1^j, s_1^{j-1}, m, t) \\
 &= P(m = 3 \mid 't_0 \text{ on the table} ') \times \\
 &\quad P(a_1 = 0 \mid \phi, \phi, 3, 't_0 \text{ on the table} ') \times \\
 &\quad P(f_1 = \text{在} \mid \{1-0\}, \phi, 3, 't_0 \text{ on the table} ') \times \\
 &\quad P(a_2 = 3 \mid \{1-0\}, \text{'在'}, 3, 't_0 \text{ on the table} ') \times \\
 &\quad P(f_2 = \text{桌子} \mid \{1-0, 2-3\}, \text{'在'}, 3, 't_0 \text{ on the table} ') \times \\
 &\quad P(a_3 = 1 \mid \{1-0, 2-3\}, \text{'在 桌子'}, 3, 't_0 \text{ on the table} ') \times \\
 &\quad P(f_3 = \text{上} \mid \{1-0, 2-3, 3-1\}, \text{'在 桌子'}, 3, 't_0 \text{ on the table} ')
 \end{aligned}$$

问题来了

- **重点**：重新回顾一下这个模型 - 两个公式

$$P(\mathbf{s}|\mathbf{t}) = \sum_{\mathbf{a}} P(\mathbf{s}, \mathbf{a}|\mathbf{t})$$

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

- **两个严重问题**

- ① 第一个公式：如何遍历所有的对齐 \mathbf{a}
- ② 第二个公式：如何计算 $P(m|\mathbf{t})$ 、 $P(a_j|a_1^{j-1}, s_1^{j-1}, m, \mathbf{t})$ 和 $P(s_j|a_1^j, s_1^{j-1}, m, \mathbf{t})$

问题来了

- **重点**：重新回顾一下这个模型 - 两个公式

$$P(\mathbf{s}|\mathbf{t}) = \sum_{\mathbf{a}} P(\mathbf{s}, \mathbf{a}|\mathbf{t})$$

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

- **两个严重问题**

- ① 第一个公式：如何遍历所有的对齐 \mathbf{a}
- ② 第二个公式：如何计算 $P(m|\mathbf{t})$ 、 $P(a_j|a_1^{j-1}, s_1^{j-1}, m, \mathbf{t})$ 和 $P(s_j|a_1^j, s_1^{j-1}, m, \mathbf{t})$

- Brown等人(1993)的解决方法：对问题进行化简

- ▶ IBM Models 1-2: 化简公式2的参数，精确高效求解公式1
→本教程内容
- ▶ IBM Models 3-5: 简单化简公式2的参数，近似求解公式1
→自学内容

IBM模型1的假设

- IBM模型1假设

- ① 源语长度概率为常数 ϵ

$$P(m|\mathbf{t}) \equiv \epsilon$$

- ② 对齐概率 $P(a_j|a_1^{j-1}, s_1^{j-1}, m, \mathbf{t})$ 仅依赖于译文长度 $l + 1$ (均匀分布)

$$P(a_j|a_1^{j-1}, s_1^{j-1}, m, \mathbf{t}) \equiv \frac{1}{l + 1}$$

- ③ 源语单词 s_j 生成概率 $P(s_j|a_1^j, s_1^{j-1}, m, \mathbf{t})$ 仅依赖与其对齐的译文单词 t_{a_j} ，即词汇翻译概率 $f(s_j|t_{a_j})$ ($\sum_{s_j} f(s_j|t_{a_j}) = 1$)

$$P(s_j|a_1^j, s_1^{j-1}, m, \mathbf{t}) \equiv f(s_j|t_{a_j})$$

IBM模型1的假设

- IBM模型1假设

- ① 源语长度概率为常数 ϵ

$$P(m|\mathbf{t}) \equiv \epsilon$$

- ② 对齐概率 $P(a_j|a_1^{j-1}, s_1^{j-1}, m, \mathbf{t})$ 仅依赖于译文长度 $l+1$ (均匀分布)

$$P(a_j|a_1^{j-1}, s_1^{j-1}, m, \mathbf{t}) \equiv \frac{1}{l+1}$$

- ③ 源语单词 s_j 生成概率 $P(s_j|a_1^j, s_1^{j-1}, m, \mathbf{t})$ 仅依赖与其对齐的译文单词 t_{a_j} ，即词汇翻译概率 $f(s_j|t_{a_j})$ ($\sum_{s_j} f(s_j|t_{a_j}) = 1$)

$$P(s_j|a_1^j, s_1^{j-1}, m, \mathbf{t}) \equiv f(s_j|t_{a_j})$$

- 核心思想是把复杂参数化简为简单参数

- ▶ 比如： $P(a_j|a_1^{j-1}, s_1^{j-1}, m, \mathbf{t}) \equiv \frac{1}{l+1}$ 把参数空间 $(a_1^j, s_1^{j-1}, m, \mathbf{t})$ 化简为 l
- ▶ **优点**：模型大大化简；**缺点**：化简导致模型不准确

IBM模型1 - 完整模型

- 代入假设

$$\begin{aligned} P(\mathbf{s}, \mathbf{a} | \mathbf{t}) &= P(m | \mathbf{t}) \prod_{j=1}^m P(a_j | a_1^{j-1}, s_1^{j-1}, m, \mathbf{t}) P(s_j | a_1^j, s_1^{j-1}, m, \mathbf{t}) \\ &\quad \begin{array}{ccc} ||| & & ||| \\ \epsilon & & \frac{1}{l+1} \\ & & f(s_j | t_{a_j}) \end{array} \\ &= \epsilon \prod_{j=1}^m \frac{1}{l+1} f(s_j | t_{a_j}) \\ &= \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m f(s_j | t_{a_j}) \end{aligned}$$

IBM模型1 - 完整模型

- 代入假设

$$\begin{aligned} P(\mathbf{s}, \mathbf{a} | \mathbf{t}) &= P(m | \mathbf{t}) \prod_{j=1}^m P(a_j | a_1^{j-1}, s_1^{j-1}, m, \mathbf{t}) P(s_j | a_1^j, s_1^{j-1}, m, \mathbf{t}) \\ &\quad \begin{array}{ccc} ||| & & ||| \\ \epsilon & & \frac{1}{l+1} \end{array} \quad \begin{array}{c} ||| \\ f(s_j | t_{a_j}) \end{array} \\ &= \epsilon \prod_{j=1}^m \frac{1}{l+1} f(s_j | t_{a_j}) \\ &= \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m f(s_j | t_{a_j}) \end{aligned}$$

- 将上式代入 $P(\mathbf{s} | \mathbf{t}) = \sum_{\mathbf{a}} P(\mathbf{s}, \mathbf{a} | \mathbf{t})$

$$\begin{aligned} P(\mathbf{s} | \mathbf{t}) &= \sum_{\mathbf{a}} P(\mathbf{s}, \mathbf{a} | \mathbf{t}) \\ &= \sum_{\mathbf{a}} \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m f(s_j | t_{a_j}) \end{aligned}$$

IBM模型1 - 完整模型(2)

- $P(\mathbf{s}|\mathbf{t}) = \sum_{\mathbf{a}} \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m f(s_j|t_{a_j})$ 中需要对遍历所有的对齐, 即 $\sum_{\mathbf{a}}$ 。这个过程可以被重新表示为

$$P(\mathbf{s}|\mathbf{t}) = \sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m f(s_j|t_{a_j})$$

IBM模型1 - 完整模型(2)

- $P(\mathbf{s}|\mathbf{t}) = \sum_{\mathbf{a}} \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m f(s_j|t_{a_j})$ 中需要对遍历所有的对齐, 即 $\sum_{\mathbf{a}}$ 。这个过程可以被重新表示为

$$P(\mathbf{s}|\mathbf{t}) = \underbrace{\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l}_{\textcircled{1}} \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m f(s_j|t_{a_j})$$

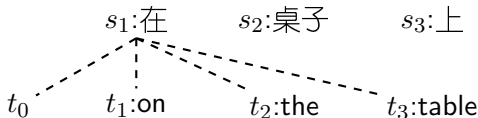
- ① 遍历所有的的对齐 \mathbf{a} 。 \mathbf{a} 由 $\{a_1, \dots, a_m\}$ 组成, 每个 $a_j \in \{a_1, \dots, a_m\}$ 从第译文开始位置(0)循环到截止位置(l)

IBM模型1 - 完整模型(2)

- $P(\mathbf{s}|\mathbf{t}) = \sum_{\mathbf{a}} \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m f(s_j|t_{a_j})$ 中需要对遍历所有的对齐, 即 $\sum_{\mathbf{a}}$ 。这个过程可以被重新表示为

$$P(\mathbf{s}|\mathbf{t}) = \underbrace{\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l}_{\textcircled{1}} \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m f(s_j|t_{a_j})$$

- ① 遍历所有的的对齐 \mathbf{a} 。 \mathbf{a} 由 $\{a_1, \dots, a_m\}$ 组成, 每个 $a_j \in \{a_1, \dots, a_m\}$ 从第译文开始位置(0)循环到截止位置(l)

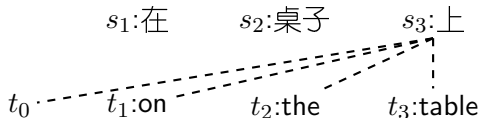


IBM模型1 - 完整模型(2)

- $P(\mathbf{s}|\mathbf{t}) = \sum_{\mathbf{a}} \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m f(s_j|t_{a_j})$ 中需要对遍历所有的对齐, 即 $\sum_{\mathbf{a}}$ 。这个过程可以被重新表示为

$$P(\mathbf{s}|\mathbf{t}) = \underbrace{\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l}_{\textcircled{1}} \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m f(s_j|t_{a_j})$$

- ① 遍历所有的的对齐 \mathbf{a} 。 \mathbf{a} 由 $\{a_1, \dots, a_m\}$ 组成, 每个 $a_j \in \{a_1, \dots, a_m\}$ 从第译文开始位置(0)循环到截止位置(l)

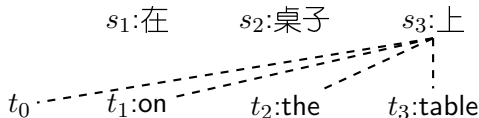


IBM模型1 - 完整模型(2)

- $P(\mathbf{s}|\mathbf{t}) = \sum_{\mathbf{a}} \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m f(s_j|t_{a_j})$ 中需要对遍历所有的对齐, 即 $\sum_{\mathbf{a}}$ 。这个过程可以被重新表示为

$$P(\mathbf{s}|\mathbf{t}) = \underbrace{\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l}_{\textcircled{1}} \underbrace{\frac{\epsilon}{(l+1)^m} \prod_{j=1}^m f(s_j|t_{a_j})}_{\textcircled{2}}$$

- ① 遍历所有的的对齐 \mathbf{a} 。 \mathbf{a} 由 $\{a_1, \dots, a_m\}$ 组成, 每个 $a_j \in \{a_1, \dots, a_m\}$ 从第译文开始位置(0)循环到截止位置(l)



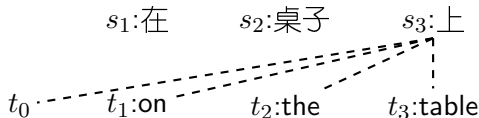
- ② 对于每个 \mathbf{a} 累加对齐概率 $P(\mathbf{s}, \mathbf{a}|\mathbf{t})$

IBM模型1 - 完整模型(2)

- $P(\mathbf{s}|\mathbf{t}) = \sum_{\mathbf{a}} \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m f(s_j|t_{a_j})$ 中需要对遍历所有的对齐, 即 $\sum_{\mathbf{a}}$ 。这个过程可以被重新表示为

$$\text{IBM模型1: } P(\mathbf{s}|\mathbf{t}) = \sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m f(s_j|t_{a_j})$$

- 1 遍历所有的的对齐 \mathbf{a} 。 \mathbf{a} 由 $\{a_1, \dots, a_m\}$ 组成, 每个 $a_j \in \{a_1, \dots, a_m\}$ 从第译文开始位置(0)循环到截止位置(l)



- 2 对于每个 \mathbf{a} 累加对齐概率 $P(\mathbf{s}, \mathbf{a}|\mathbf{t})$

IBM模型2的假设

- IBM模型1虽然形式非常简单，但是问题是忽略了翻译的调序问题。对于同一个源语言句子，单词翻译相同，但是译文顺序不同，模型得分（翻译概率相同）

s = 在 桌子 上 s = 在 桌子 上
 ↓ ↓
t = on the table t' = table on the

IBM模型1: $P(s|t)$ = $P(s|t')$

理想: $P(s|t)$ > $P(s|t')$

IBM模型2 - 完整模型

- 代入假设

$$\begin{aligned} P(\mathbf{s}, \mathbf{a}|\mathbf{t}) &= P(m|\mathbf{t}) \prod_{j=1}^m P(a_j|a_1^{j-1}, s_1^{j-1}, m, \mathbf{t}) P(s_j|a_1^j, s_1^{j-1}, m, \mathbf{t}) \\ &\quad \begin{array}{ccc} ||| & & ||| \\ \in & & \end{array} \\ &\quad \quad \quad a(a_j|j, m, l) \quad \quad \quad f(s_j|t_{a_j}) \\ &= \epsilon \prod_{j=1}^m a(a_j|j, m, l) f(s_j|t_{a_j}) \end{aligned}$$

IBM模型2 - 完整模型

- 代入假设

$$\begin{aligned} P(\mathbf{s}, \mathbf{a} | \mathbf{t}) &= P(m | \mathbf{t}) \prod_{j=1}^m P(a_j | a_1^{j-1}, s_1^{j-1}, m, \mathbf{t}) P(s_j | a_1^j, s_1^{j-1}, m, \mathbf{t}) \\ &\quad \begin{array}{ccc} \text{|||} & & \text{|||} \\ \epsilon & & f(s_j | t_{a_j}) \end{array} \\ &\quad a(a_j | j, m, l) \\ &= \epsilon \prod_{j=1}^m a(a_j | j, m, l) f(s_j | t_{a_j}) \end{aligned}$$

- 将上式代入 $P(\mathbf{s} | \mathbf{t}) = \sum_{\mathbf{a}} P(\mathbf{s}, \mathbf{a} | \mathbf{t})$

$$P(\mathbf{s} | \mathbf{t}) = \underbrace{\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l}_{\textcircled{1}} \underbrace{\epsilon \prod_{j=1}^m a(a_j | j, m, l) f(s_j | t_{a_j})}_{\textcircled{2}}$$

- 遍历所有的的对齐 \mathbf{a}
- 对于每个 \mathbf{a} 累加对齐概率 $P(\mathbf{s}, \mathbf{a} | \mathbf{t})$ ，即计

$$\text{算} \epsilon \prod_{j=1}^m a(a_j | j, m, l) f(s_j | t_{a_j})$$

IBM模型2 - 完整模型

- 代入假设

$$\begin{aligned} P(\mathbf{s}, \mathbf{a} | \mathbf{t}) &= P(m | \mathbf{t}) \prod_{j=1}^m P(a_j | a_1^{j-1}, s_1^{j-1}, m, \mathbf{t}) P(s_j | a_1^j, s_1^{j-1}, m, \mathbf{t}) \\ &\quad \begin{array}{ccc} ||| & & ||| \\ \in & & \end{array} \\ &\quad \quad \quad a(a_j | j, m, l) \quad \quad \quad f(s_j | t_{a_j}) \\ &= \epsilon \prod_{j=1}^m a(a_j | j, m, l) f(s_j | t_{a_j}) \end{aligned}$$

- 将上式代入 $P(\mathbf{s} | \mathbf{t}) = \sum_{\mathbf{a}} P(\mathbf{s}, \mathbf{a} | \mathbf{t})$

IBM模型2: $P(\mathbf{s} | \mathbf{t}) = \sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \epsilon \prod_{j=1}^m a(a_j | j, m, l) f(s_j | t_{a_j})$

- 遍历所有的的对齐 \mathbf{a}
- 对于每个 \mathbf{a} 累加对齐概率 $P(\mathbf{s}, \mathbf{a} | \mathbf{t})$, 即计算 $\epsilon \prod_{j=1}^m a(a_j | j, m, l) f(s_j | t_{a_j})$

停顿一下，重新看看这两个模型

$$\text{IBM模型1 : } P(\mathbf{s}|\mathbf{t}) = \frac{\epsilon}{(l+1)^m} \sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m f(s_j|t_{a_j})$$

$$\text{IBM模型2 : } P(\mathbf{s}|\mathbf{t}) = \epsilon \sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m a(a_j|j, m, l) f(s_j|t_{a_j})$$

- 参数
 - ▶ ϵ 为常量

停顿一下，重新看看这两个模型

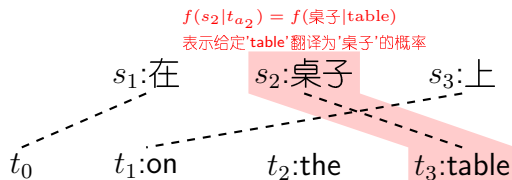
$$\text{IBM模型1: } P(s|t) = \frac{\epsilon}{(l+1)^m} \sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m f(s_j|t_{a_j})$$

$$\text{IBM模型2: } P(s|t) = \epsilon \sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m a(a_j|j, m, l) f(s_j|t_{a_j})$$

- 参数

- ▶ ϵ 为常量

- ▶ $f(s_j|t_{a_j})$: 给定单词 t_{a_j} 翻译为 s_j 的概率



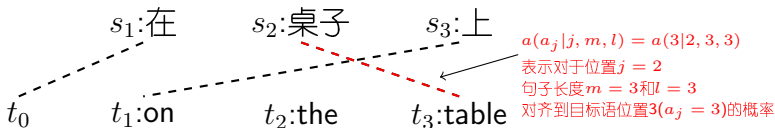
停顿一下，重新看看这两个模型

$$\text{IBM模型1: } P(s|t) = \frac{\epsilon}{(l+1)^m} \sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m f(s_j|t_{a_j})$$

$$\text{IBM模型2: } P(s|t) = \epsilon \sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m a(a_j|j, m, l) f(s_j|t_{a_j})$$

- 参数

- ▶ ϵ 为常量
- ▶ $f(s_j|t_{a_j})$: 给定单词 t_{a_j} 翻译为 s_j 的概率
- ▶ $a(a_j|j, m, l)$: 给定位置 j 和句对长度 m 和 l ，对齐 a_j 的概率



停顿一下，重新看看这两个模型(2)

$$\text{IBM模型1: } P(\mathbf{s}|\mathbf{t}) = \frac{\epsilon}{(l+1)^m} \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \prod_{j=1}^m f(s_j|t_{a_j})$$

$$\text{IBM模型2: } P(\mathbf{s}|\mathbf{t}) = \epsilon \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \prod_{j=1}^m a(a_j|j, m, l) f(s_j|t_{a_j})$$

- 对于翻译模型 $P(\mathbf{s}|\mathbf{t})$ ，再来回顾一下统计机器翻译的三个基本问题
 - ① 建模：如何描述 $P(\mathbf{s}|\mathbf{t})$
 - ② 解码：给定模型参数 ϵ 、 $a(a_j|j, m, l)$ 和 $f(s_j|t_{a_j})$ ，如何利用上面的公式计算 $P(\mathbf{s}|\mathbf{t})$ (语言模型计算暂不讨论)，并找到最佳译文 $\hat{\mathbf{t}}$
 - ③ 训练：如何从数据中自动学习模型参数 ϵ 、 $a(a_j|j, m, l)$ 和 $f(s_j|t_{a_j})$

停顿一下，重新看看这两个模型(2)

$$\text{IBM模型1: } P(\mathbf{s}|\mathbf{t}) = \frac{\epsilon}{(l+1)^m} \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \prod_{j=1}^m f(s_j|t_{a_j})$$

$$\text{IBM模型2: } P(\mathbf{s}|\mathbf{t}) = \epsilon \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \prod_{j=1}^m a(a_j|j, m, l) f(s_j|t_{a_j})$$

- 对于翻译模型 $P(\mathbf{s}|\mathbf{t})$ ，再来回顾一下统计机器翻译的三个基本问题
 - ① 建模：如何描述 $P(\mathbf{s}|\mathbf{t})$ ← 已解！见上面两个公式
 - ② 解码：给定模型参数 ϵ 、 $a(a_j|j, m, l)$ 和 $f(s_j|t_{a_j})$ ，如何利用上面的公式计算 $P(\mathbf{s}|\mathbf{t})$ （语言模型计算暂不讨论），并找到最佳译文 $\hat{\mathbf{t}}$ ← 下面讨论
 - ③ 训练：如何从数据中自动学习模型参数 ϵ 、 $a(a_j|j, m, l)$ 和 $f(s_j|t_{a_j})$ ← 下面讨论

Outline

1. 一个简单的翻译实例

2. IBM模型

- 建模
- 解码
- 模型训练

解码 = 搜索 + 模型得分计算

- 何为模型：统计描述 + 参数

问题的统计描述: $P(\mathbf{s}|\mathbf{t}) = \epsilon \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \prod_{j=1}^m a(a_j|j, m, l) f(s_j|t_{a_j})$

模型的参数: $\epsilon = ?; \quad \forall a_j, j, m, l : a(a_j|j, m, l) = ?, f(s_j|t_{a_j}) = ?$

解码 = 搜索 + 模型得分计算

- 何为模型：统计描述 + 参数

问题的统计描述: $P(\mathbf{s}|\mathbf{t}) = \epsilon \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \prod_{j=1}^m a(a_j|j, m, l) f(s_j|t_{a_j})$

模型的参数: $\epsilon = ?; \forall a_j, j, m, l : a(a_j|j, m, l) = ?, f(s_j|t_{a_j}) = ?$

- 何为解码：模型得分计算 + 找到模型得分最高的译文

模型得分计算：对任意的 \mathbf{s} 和 \mathbf{t} ，(高效地)计算 $P(\mathbf{s}|\mathbf{t})$ (同时计算 $P(\mathbf{t})$)

搜索：对所有可能的 \mathbf{t} ，找到模型得分($P(\mathbf{s}|\mathbf{t})P(\mathbf{t})$)最高的译文输出

- ▶ 搜索(解码)问题前面的实例已经描述了一种解法（见第16页）：自左向右添加译文单词 + 剪枝技术。这里不再讨论，可以自行学习
- ▶ 剩下的问题是：对于任意的 \mathbf{s} 和 \mathbf{t} ，如何高效地计算 $P(\mathbf{s}|\mathbf{t})$

模型得分的计算

- $O((l+1)^m \cdot m)$ - IBM模型得分的直接计算几乎不可能!

$$P(\mathbf{s}|\mathbf{t}) = \frac{\epsilon}{(l+1)^m} \underbrace{\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l}_{(l+1)^m \text{次循环}} \underbrace{\prod_{j=1}^m f(s_j|t_{a_j})}_{m \text{次循环}}$$

模型得分的计算

- $O((l+1)^m \cdot m)$ - IBM模型得分的直接计算几乎不可能!

$$P(\mathbf{s}|\mathbf{t}) = \frac{\epsilon}{(l+1)^m} \underbrace{\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l}_{(l+1)^m \text{次循环}} \underbrace{\prod_{j=1}^m f(s_j|t_{a_j})}_{m \text{次循环}}$$

- $O(l \cdot m)$ - 实际上我们可以做的更好

$$\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m f(s_j|t_{a_j}) = \underbrace{\prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)}_{m \cdot l \text{次循环}}$$

模型得分的计算

- $O((l+1)^m \cdot m)$ - IBM模型得分的直接计算几乎不可能!

$$P(\mathbf{s}|\mathbf{t}) = \frac{\epsilon}{(l+1)^m} \underbrace{\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l}_{(l+1)^m \text{次循环}} \underbrace{\prod_{j=1}^m f(s_j|t_{a_j})}_{m \text{次循环}}$$

- $O(l \cdot m)$ - 实际上我们可以做的更好

$$\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m f(s_j|t_{a_j}) = \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)$$

IBM模型1: $P(\mathbf{s}|\mathbf{t}) = \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)$

模型得分的计算

- $O((l+1)^m \cdot m)$ - IBM模型得分的直接计算几乎不可能!

$$P(\mathbf{s}|\mathbf{t}) = \frac{\epsilon}{(l+1)^m} \underbrace{\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l}_{(l+1)^m \text{次循环}} \underbrace{\prod_{j=1}^m f(s_j|t_{a_j})}_{m \text{次循环}}$$

- $O(l \cdot m)$ - 实际上我们可以做的更好

$$\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m f(s_j|t_{a_j}) = \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)$$

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

$$\text{类似的, IBM模型2: } P(\mathbf{s}|\mathbf{t}) = \epsilon \prod_{j=1}^m \sum_{i=0}^l a(i|j, m, l) f(s_j|t_i)$$

$$\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m f(s_j|t_{a_j}) = \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)$$

- 定义 $\alpha(x, y)$ 为一个函数， x 和 y 是两个变量

$$\begin{aligned} &\alpha(1, 0)\alpha(2, 0) + \alpha(1, 0)\alpha(2, 1) + \alpha(1, 0)\alpha(2, 2) + \\ &\alpha(1, 1)\alpha(2, 0) + \alpha(1, 1)\alpha(2, 1) + \alpha(1, 1)\alpha(2, 2) + \\ &\alpha(1, 2)\alpha(2, 0) + \alpha(1, 2)\alpha(2, 1) + \alpha(1, 2)\alpha(2, 2) \end{aligned}$$

$$\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m f(s_j|t_{a_j}) = \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)$$

- 定义 $\alpha(x, y)$ 为一个函数, x 和 y 是两个变量

$$\begin{aligned} &\alpha(1, 0)\alpha(2, 0) + \alpha(1, 0)\alpha(2, 1) + \alpha(1, 0)\alpha(2, 2) + \\ &\alpha(1, 1)\alpha(2, 0) + \alpha(1, 1)\alpha(2, 1) + \alpha(1, 1)\alpha(2, 2) + \\ &\alpha(1, 2)\alpha(2, 0) + \alpha(1, 2)\alpha(2, 1) + \alpha(1, 2)\alpha(2, 2) \end{aligned}$$

$$\begin{aligned} &\sum_{y_1=0}^2 \sum_{y_2=0}^2 \alpha(1, y_1)\alpha(2, y_2) \\ &= \sum_{y_1=0}^2 \sum_{y_2=0}^2 \prod_{x=1}^2 \alpha(x, y_x) \end{aligned}$$

$$\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m f(s_j|t_{a_j}) = \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)$$

- 定义 $\alpha(x, y)$ 为一个函数, x 和 y 是两个变量

$$\begin{aligned} &\alpha(1, 0)\alpha(2, 0) + \alpha(1, 0)\alpha(2, 1) + \alpha(1, 0)\alpha(2, 2) + \\ &\alpha(1, 1)\alpha(2, 0) + \alpha(1, 1)\alpha(2, 1) + \alpha(1, 1)\alpha(2, 2) + \\ &\alpha(1, 2)\alpha(2, 0) + \alpha(1, 2)\alpha(2, 1) + \alpha(1, 2)\alpha(2, 2) \end{aligned}$$

$$\begin{aligned} &\sum_{y_1=0}^2 \sum_{y_2=0}^2 \alpha(1, y_1)\alpha(2, y_2) \\ &= \sum_{y_1=0}^2 \sum_{y_2=0}^2 \prod_{x=1}^2 \alpha(x, y_x) \end{aligned}$$

$$\begin{aligned} &(\alpha(1, 0) + \alpha(1, 1) + \alpha(1, 2)) \cdot \\ &(\alpha(2, 0) + \alpha(2, 1) + \alpha(2, 2)) \\ &= \prod_{x=1}^2 \sum_{y=0}^2 \alpha(x, y) \end{aligned}$$

$$\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m f(s_j|t_{a_j}) = \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)$$

- 定义 $\alpha(x, y)$ 为一个函数, x 和 y 是两个变量

$$\begin{aligned} &\alpha(1, 0)\alpha(2, 0) + \alpha(1, 0)\alpha(2, 1) + \alpha(1, 0)\alpha(2, 2) + \\ &\alpha(1, 1)\alpha(2, 0) + \alpha(1, 1)\alpha(2, 1) + \alpha(1, 1)\alpha(2, 2) + \\ &\alpha(1, 2)\alpha(2, 0) + \alpha(1, 2)\alpha(2, 1) + \alpha(1, 2)\alpha(2, 2) \end{aligned}$$

$$\begin{aligned} &\sum_{y_1=0}^2 \sum_{y_2=0}^2 \alpha(1, y_1)\alpha(2, y_2) \\ &= \sum_{y_1=0}^2 \sum_{y_2=0}^2 \prod_{x=1}^2 \alpha(x, y_x) \end{aligned}$$

=

$$\begin{aligned} &(\alpha(1, 0) + \alpha(1, 1) + \alpha(1, 2)) \cdot \\ &(\alpha(2, 0) + \alpha(2, 1) + \alpha(2, 2)) \\ &= \prod_{x=1}^2 \sum_{y=0}^2 \alpha(x, y) \end{aligned}$$

$$\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m f(s_j|t_{a_j}) = \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)$$

- 定义 $\alpha(x, y)$ 为一个函数, x 和 y 是两个变量

$$\begin{aligned} &\alpha(1, 0)\alpha(2, 0) + \alpha(1, 0)\alpha(2, 1) + \alpha(1, 0)\alpha(2, 2) + \\ &\alpha(1, 1)\alpha(2, 0) + \alpha(1, 1)\alpha(2, 1) + \alpha(1, 1)\alpha(2, 2) + \\ &\alpha(1, 2)\alpha(2, 0) + \alpha(1, 2)\alpha(2, 1) + \alpha(1, 2)\alpha(2, 2) \end{aligned}$$

$$\begin{aligned} &\sum_{y_1=0}^2 \sum_{y_2=0}^2 \alpha(1, y_1)\alpha(2, y_2) \\ &= \sum_{y_1=0}^2 \sum_{y_2=0}^2 \prod_{x=1}^2 \alpha(x, y_x) \end{aligned}$$

=

$$\begin{aligned} &(\alpha(1, 0) + \alpha(1, 1) + \alpha(1, 2)) \cdot \\ &(\alpha(2, 0) + \alpha(2, 1) + \alpha(2, 2)) \\ &= \prod_{x=1}^2 \sum_{y=0}^2 \alpha(x, y) \end{aligned}$$

$$\sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m f(s_j|t_{a_j}) = \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)$$

Outline

1. 一个简单的翻译实例

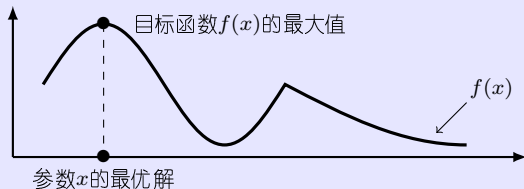
2. IBM模型

- 建模
- 解码
- 模型训练

训练 - 目标函数

训练 = 优化

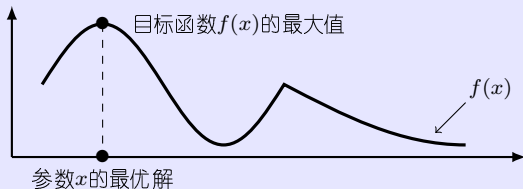
在训练集上调整参数使得目标函数的值达到最大（最小），此时的参数称作该模型在该目标函数下的最优解



训练 - 目标函数

训练 = 优化

在训练集上调整参数使得目标函数的值达到最大（最小），此时的参数称作该模型在该目标函数下的最优解



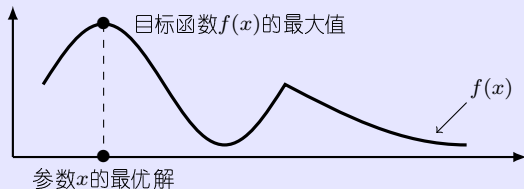
- IBM模型的训练：对于给定的句对 (s, t) ，最大化翻译概率 $P(s|t)$ 。这里用符号 $P_{\theta}(s|t)$ 表示概率由参数 θ 决定

$$\hat{\theta} = \arg \max_{\theta} P_{\theta}(s|t)$$

训练 - 目标函数

训练 = 优化

在训练集上调整参数使得目标函数的值达到最大（最小），此时的参数称作该模型在该目标函数下的最优解



- IBM模型的训练：对于给定的句对 (s, t) ，最大化翻译概率 $P(s|t)$ 。这里用符号 $P_{\theta}(s|t)$ 表示概率由参数 θ 决定

$$\hat{\theta} = \underset{\theta}{\operatorname{arg\,max}} P_{\theta}(s|t)$$

求最优参数

目标函数

极大似然估计

- $P(s|t)$ 可以被看做是 (s, t) 上的似然函数($L(s, t; \theta)$)。所谓极大似然估计，就是要找到使 $L(s, t; \theta)$ 达到最大的 θ ：

$$\{\hat{\theta}\} \subseteq \{\arg \max_{\theta \in \Theta} L(s, t; \theta)\}$$

$L(s, t; \theta)$ 表示 $L(\cdot)$ 依赖模型参数 θ （注意分号）， $\{\hat{\theta}\}$ 表示可能有多组结果， Θ 表示参数空间

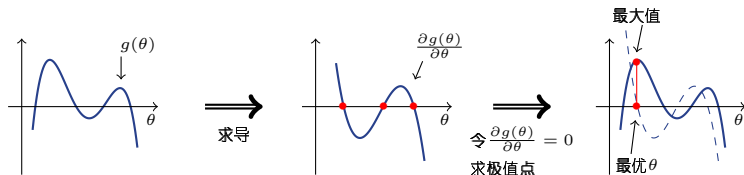
极大似然估计

- $P(s|t)$ 可以被看做是 (s, t) 上的似然函数($L(s, t; \theta)$)。所谓极大似然估计，就是要找到使 $L(s, t; \theta)$ 达到最大的 θ ：

$$\{\hat{\theta}\} \subseteq \{\arg \max_{\theta \in \Theta} L(s, t; \theta)\}$$

$L(s, t; \theta)$ 表示 $L(\cdot)$ 依赖模型参数 θ （注意分号）， $\{\hat{\theta}\}$ 表示可能有多组结果， Θ 表示参数空间

- 先不用考虑上面的公式。我们还是回归到原始问题：如何找到一组 θ 使 $P_{\theta}(s|t)$ 达到最大？
 - ▶ 求函数最大值问题。比如，我们可以对 $P_{\theta}(s|t)$ 求导，令导数为零，得到极值点



最大化 $P(s|t)$

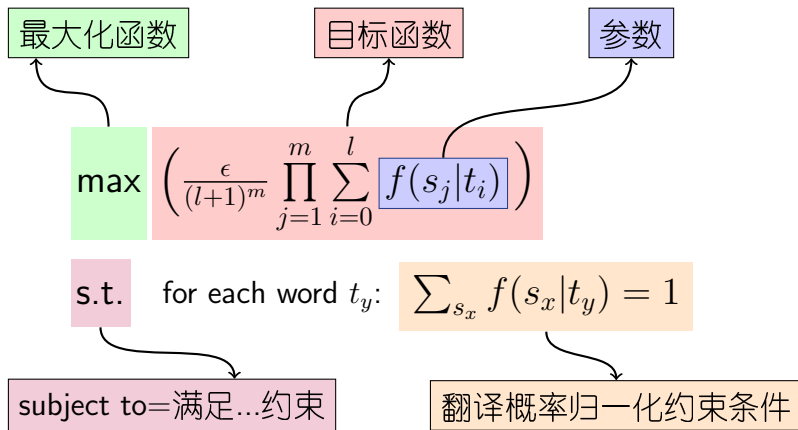
- 以IBM模型1为例，它所对应的(基于极大似然估计)的模型训练问题可以被描述为：

$$\max \left(\frac{\epsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i) \right)$$

$$\text{s.t.} \quad \text{for each word } t_y: \quad \sum_{s_x} f(s_x|t_y) = 1$$

最大化 $P(s|t)$

- 以IBM模型1为例，它所对应的(基于极大似然估计)的模型训练问题可以被描述为：



- 注意： $\{f(s_x|t_y)\}$ 对应很多参数，每个源语言单词和每个目标语单词的组合都对应一个 $f(s_x|t_y)$

拉格朗日乘数法 - Lagrange multiplier

- 含有约束的优化问题: 不好解

目标: $\max(P_{\theta}(s|t))$ + 约束: $\forall t_y : \sum_{s_x} P(s_x|t_y) = 1$

拉格朗日乘数法 - Lagrange multiplier

- 含有约束的优化问题: 不好解

目标: $\max(P_\theta(s|t))$ + 约束: $\forall t_y : \sum_{s_x} P(s_x|t_y) = 1$

- 解决方法: 含有约束优化 \Rightarrow 不含约束优化

拉格朗日乘数法 - 百度百科

拉格朗日乘数法是一种寻找变量受一个或多个条件所限制的多元函数的极值的方法。这种方法将一个有 n 个变量与 k 个约束条件的最优化问题转换为一个有 $n + k$ 个变量的方程组的极值问题，其变量不受任何约束。这种方法引入了一种新的标量未知数，即拉格朗日乘数。

- ▶ 对于每个 t_y 所对应的约束，引入一个拉格朗日乘数 λ_{t_y}

$$L(f, \lambda) = \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l \prod_{j=1}^m f(s_j|t_i) - \sum_{t_y} \lambda_{t_y} \left(\sum_{s_x} f(s_x|t_y) - 1 \right)$$

拉格朗日乘数法 - Lagrange multiplier (2)

原始问题	转化后的问题
$\max(P(\mathbf{s} \mathbf{t}))$ s.t. $\forall t_y : \sum_{s_x} f(s_x t_y) = 1$	$\max(L(f, \lambda))$

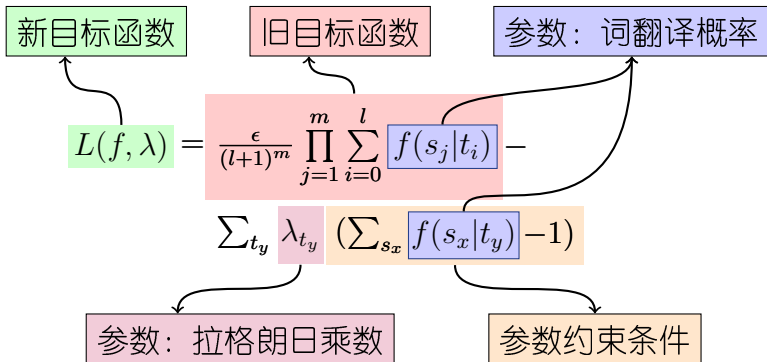
新目标函数

旧目标函数

$$L(f, \lambda) = \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i) - \sum_{t_y} \lambda_{t_y} (\sum_{s_x} f(s_x|t_y) - 1)$$

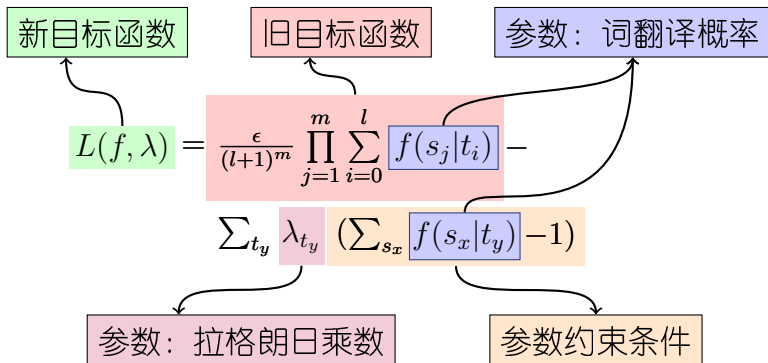
拉格朗日乘数法 - Lagrange multiplier (2)

原始问题	转化后的问题
$\max(P(s t))$ s.t. $\forall t_y : \sum_{s_x} f(s_x t_y) = 1$	$\max(L(f, \lambda))$



拉格朗日乘数法 - Lagrange multiplier (2)

原始问题	转化后的问题
$\max(P(s t))$ s.t. $\forall t_y : \sum_{s_x} f(s_x t_y) = 1$	$\max(L(f, \lambda))$



- 剩下的问题：计算 $\frac{\partial L(f, \lambda)}{\partial f(s_u|t_v)}$ ，并求解 $\frac{\partial L(f, \lambda)}{\partial f(s_u|t_v)} = 0$

对 L 求导

- 令 s_u 和 t_v 表示任意一个源语单词和一个目标语单词

$$\begin{aligned}\frac{\partial L(f, \lambda)}{\partial f(s_u | t_v)} &= \frac{\partial \left[\frac{\epsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l f(s_j | t_i) \right]}{\partial f(s_u | t_v)} - \\ &\quad \frac{\partial \left[\sum_{t_y} \lambda_{t_y} (\sum_{s_x} f(s_x | t_y) - 1) \right]}{\partial f(s_u | t_v)} \\ &= \frac{\epsilon}{(l+1)^m} \cdot \frac{\partial \left[\prod_{j=1}^m \sum_{i=0}^l f(s_j | t_i) \right]}{\partial f(s_u | t_v)} - \lambda_{t_v}\end{aligned}$$

对 L 求导

- 令 s_u 和 t_v 表示任意一个源语单词和一个目标语单词

$$\begin{aligned}\frac{\partial L(f, \lambda)}{\partial f(s_u | t_v)} &= \frac{\partial \left[\frac{\epsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l f(s_j | t_i) \right]}{\partial f(s_u | t_v)} - \\ &\quad \frac{\partial \left[\sum_{t_y} \lambda_{t_y} (\sum_{s_x} f(s_x | t_y) - 1) \right]}{\partial f(s_u | t_v)} \\ &= \frac{\epsilon}{(l+1)^m} \cdot \frac{\partial \left[\prod_{j=1}^m \sum_{i=0}^l f(s_j | t_i) \right]}{\partial f(s_u | t_v)} - \lambda_{t_v}\end{aligned}$$

- 为了求 $\frac{\partial \left[\prod_{j=1}^m \sum_{i=0}^l f(s_j | t_i) \right]}{\partial f(s_u | t_v)}$ ，先看一个简单的例子：
 $g(z) = \alpha z^\beta$ 为一个变量 z 的多项式函数，显然

$$\frac{\partial g(z)}{\partial z} = \alpha \beta z^{\beta-1} = \frac{\beta}{z} \alpha z^\beta = \frac{\beta}{z} g(z)$$

对 L 求导(2)

- 这里可以把 $\prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)$ 看做 $g(z) = \alpha z^\beta$ 的实例
 - ▶ 令 $z = \sum_{i=0}^l f(s_u|t_i)$, 注意 s_u 为给定的源语单词
 - ▶ 那么 β 为 $\sum_{i=0}^l f(s_u|t_i)$ 在 $\prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)$ 中出现的次数, 即源语句子中与 s_u 相同的单词的个数

$$\beta = \sum_{j=1}^m \delta(s_j, s_u)$$

$\delta(x, y)$: 当 $x = y$ 时为1, 否则为0

对 L 求导(2)

- 这里可以把 $\prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)$ 看做 $g(z) = \alpha z^\beta$ 的实例
 - 令 $z = \sum_{i=0}^l f(s_u|t_i)$, 注意 s_u 为给定的源语单词
 - 那么 β 为 $\sum_{i=0}^l f(s_u|t_i)$ 在 $\prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)$ 中出现的次数, 即源语句子中与 s_u 相同的单词的个数

$$\beta = \sum_{j=1}^m \delta(s_j, s_u)$$

$\delta(x, y)$: 当 $x = y$ 时为1, 否则为0

- 根据 $\frac{\partial g(z)}{\partial z} = \frac{\beta}{z} g(z)$

$$\frac{\partial g(z)}{\partial z} = \frac{\partial \left[\prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i) \right]}{\partial \left[\sum_{i=0}^l f(s_u|t_i) \right]} = \frac{\sum_{j=1}^m \delta(s_j, s_u)}{\sum_{i=0}^l f(s_u|t_i)} \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)$$

对 L 求导(3)

- 如果把 z 看成 f 的函数, 有 $\frac{\partial g(z)}{\partial f} = \frac{\partial g(z)}{\partial z} \cdot \frac{\partial z}{\partial f}$
 - ▶ 因为 $z = \sum_{i=0}^l f(s_u|t_i)$, 很容易得到

$$\frac{\partial z}{\partial f} = \frac{\partial [\sum_{i=0}^l f(s_u|t_i)]}{\partial f(s_u|t_v)} = \sum_{i=0}^l \delta(t_i, t_v)$$

即目标语译文单词中与 t_v 相同的个数

对 L 求导(3)

- 如果把 z 看成 f 的函数, 有 $\frac{\partial g(z)}{\partial f} = \frac{\partial g(z)}{\partial z} \cdot \frac{\partial z}{\partial f}$
 - 因为 $z = \sum_{i=0}^l f(s_u|t_i)$, 很容易得到

$$\frac{\partial z}{\partial f} = \frac{\partial [\sum_{i=0}^l f(s_u|t_i)]}{\partial f(s_u|t_v)} = \sum_{i=0}^l \delta(t_i, t_v)$$

即目标语译文单词中与 t_v 相同的个数

- 根据 $\frac{\partial g(z)}{\partial z}$ 和 $\frac{\partial z}{\partial f}$ 计算的结果, 可以得到

$$\begin{aligned} \frac{\partial \left[\prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i) \right]}{\partial f(s_u|t_v)} &= \frac{\partial \left[\prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i) \right]}{\partial \left[\sum_{i=0}^l f(s_u|t_i) \right]} \cdot \frac{\partial \left[\sum_{i=0}^l f(s_u|t_i) \right]}{\partial f(s_u|t_v)} \\ &= \frac{\sum_{j=1}^m \delta(s_j, s_u)}{\sum_{i=0}^l f(s_u|t_i)} \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i) \cdot \sum_{i=0}^l \delta(t_i, t_v) \end{aligned}$$

对 L 求导(4)

- 将 $\frac{\partial [\prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)]}{\partial f(s_u|t_v)}$ 进一步代入 $\frac{\partial L(f, \lambda)}{\partial f(s_u|t_v)}$

$$\frac{\partial L(f, \lambda)}{\partial f(s_u|t_v)}$$

$$\begin{aligned} &= \frac{\epsilon}{(l+1)^m} \cdot \frac{\partial [\prod_{j=1}^m \sum_{i=0}^l f(s_j|t_{a_j})]}{\partial f(s_u|t_v)} - \lambda_{t_v} \\ &= \frac{\epsilon}{(l+1)^m} \cdot \frac{\sum_{j=1}^m \delta(s_j, s_u) \cdot \sum_{i=0}^l \delta(t_i, t_v)}{\sum_{i=0}^l f(s_u|t_i)} \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i) - \lambda_{t_v} \end{aligned}$$

对 L 求导(4)

- 将 $\frac{\partial [\prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)]}{\partial f(s_u|t_v)}$ 进一步代入 $\frac{\partial L(f, \lambda)}{\partial f(s_u|t_v)}$

$$\frac{\partial L(f, \lambda)}{\partial f(s_u|t_v)}$$

$$= \frac{\epsilon}{(l+1)^m} \cdot \frac{\partial [\prod_{j=1}^m \sum_{i=0}^l f(s_j|t_{a_j})]}{\partial f(s_u|t_v)} - \lambda_{t_v}$$

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

- 令 $\frac{\partial L(f, \lambda)}{\partial f(s_u|t_v)} = 0$, 有

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

求 $f(s_u|t_v)$ 的最优解 - 这不是一个解析式!!!

- 注意: 这不是一个计算 $f(s_u|t_v)$ 的解析式, 因为等式右端仍含有 $f(s_u|t_v)$

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

求 $f(s_u|t_v)$ 的最优解 - 这不是一个解析式!!!

- 注意: 这不是一个计算 $f(s_u|t_v)$ 的解析式, 因为等式右端仍含有 $f(s_u|t_v)$

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

- 怎么办? - 我们仍可以利用上式迭代地计算 $f(s_u|t_v)$, 使其最终收敛到最优值
 - 这是一个非常经典的方法, 称作期望最大化(Expectation Maximization), 简称EM方法(或算法)
 - 思想: 用当前的参数, 求一个似然函数的期望, 之后最大化这个期望值同时得到新的一组参数的值
 - 对于IBM模型, 很简单: 反复使用上式即可!

新的参数值

旧的参数值

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

求 $f(s_u|t_v)$ 的最优解 - 这不是一个解析式!!!

- 注意: 这不是一个计算 $f(s_u|t_v)$ 的解析式, 因为等式右端仍含有 $f(s_u|t_v)$

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

- 怎么办? - 我们仍可以利用上式迭代地计算 $f(s_u|t_v)$, 使其最终收敛到最优值
 - 这是一个非常经典的方法, 称作期望最大化(Expectation Maximization), 简称EM方法(或算法)
 - 思想: 用当前的参数, 求一个似然函数的期望, 之后最大化这个期望值同时得到新的一组参数的值
 - 对于IBM模型, 很简单: 反复使用上式即可!

新的参数值 问题: $\lambda_{t_v} = ?$ 旧的参数值

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

对 $f(s_u|t_v)$ 的计算进行化简

- 为了计算 $f(s_u|t_v)$ ，重新组织一下上面的公式

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

对 $f(s_u|t_v)$ 的计算进行化简

- 为了计算 $f(s_u|t_v)$ ，重新组织一下上面的公式

翻译概率 $P(\mathbf{s}|\mathbf{t})$

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

||
 $P(\mathbf{s}|\mathbf{t})$

对 $f(s_u|t_v)$ 的计算进行化简

- 为了计算 $f(s_u|t_v)$ ，重新组织一下上面的公式

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

||
 $P(\mathbf{s}|\mathbf{t})$

(s_u, t_v) 在句对 (\mathbf{s}, \mathbf{t}) 中
配对的总次数

对 $f(s_u|t_v)$ 的计算进行化简

- 为了计算 $f(s_u|t_v)$ ，重新组织一下上面的公式

	翻译概率 $P(\mathbf{s} \mathbf{t})$	(s_u, t_v) 在句对 (\mathbf{s}, \mathbf{t}) 中配对的总次数	$f(s_u t_v)$ 对于所有的 t_i 的相对值
$f(s_u t_v) = \lambda_{t_v}^{-1} \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l f(s_j t_i)$		$\sum_{j=1}^m \delta(s_j, s_u) \sum_{i=0}^l \delta(t_i, t_v)$	$\frac{f(s_u t_v)}{\sum_{i=0}^l f(s_u t_i)}$
$P(\mathbf{s} \mathbf{t})$			

对 $f(s_u|t_v)$ 的计算进行化简

- 为了计算 $f(s_u|t_v)$ ，重新组织一下上面的公式

	翻译概率 $P(\mathbf{s} \mathbf{t})$	(s_u, t_v) 在句对 (\mathbf{s}, \mathbf{t}) 中配对的总次数	$f(s_u t_v)$ 对于所有的 t_i 的相对值
$f(s_u t_v) = \lambda_{t_v}^{-1} \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l f(s_j t_i)$		$\sum_{j=1}^m \delta(s_j, s_u) \sum_{i=0}^l \delta(t_i, t_v)$	$\frac{f(s_u t_v)}{\sum_{i=0}^l f(s_u t_i)}$
\parallel	$P(\mathbf{s} \mathbf{t})$	$\underbrace{\hspace{10em}}$ ' 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	

对 $f(s_u|t_v)$ 的计算进行化简

- 为了计算 $f(s_u|t_v)$ ，重新组织一下上面的公式

$$f(s_u|t_v) = \lambda_{t_v}^{-1} \underbrace{\frac{\epsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)}_{\text{翻译概率 } P(\mathbf{s}|\mathbf{t})} \underbrace{\sum_{j=1}^m \delta(s_j, s_u) \sum_{i=0}^l \delta(t_i, t_v)}_{\substack{\text{'t_v 翻译为 s_u' 这个事件} \\ \text{出现次数的期望的估计} \\ \text{称之为期望频次 expected count}}} \underbrace{\frac{f(s_u|t_v)}{\sum_{i=0}^l f(s_u|t_i)}}_{\substack{(s_u, t_v) \text{ 在句对 } (\mathbf{s}, \mathbf{t}) \text{ 中} \\ \text{配对的总次数}} \quad f(s_u|t_v) \text{ 对于所有} \\ \text{的 } t_i \text{ 的相对值}}$$

\parallel
 $P(\mathbf{s}|\mathbf{t})$

- 啥是事件的期望频次？ - 事件在其分布下出现的次数的期望 $c_{\mathbb{E}}(X) = \sum_i c(x_i) \cdot P(x_i)$

x_i	$c(x_i)$	x_i	$c(x_i)$	$P(x_i)$	$c(x_i) \cdot P(x_i)$
x_1	2	x_1	2	0.1	0.2
x_2	1	x_2	1	0.3	0.3
x_3	5	x_3	5	0.2	1.0

对 $f(s_u|t_v)$ 的计算进行化简

- 为了计算 $f(s_u|t_v)$ ，重新组织一下上面的公式

$$f(s_u|t_v) = \lambda_{t_v}^{-1} \underbrace{\frac{\epsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)}_{\text{翻译概率 } P(\mathbf{s}|\mathbf{t})} \underbrace{\sum_{j=1}^m \delta(s_j, s_u) \sum_{i=0}^l \delta(t_i, t_v)}_{\substack{(s_u, t_v) \text{ 在句对 } (\mathbf{s}, \mathbf{t}) \text{ 中} \\ \text{配对的总次数}}} \underbrace{\frac{f(s_u|t_v)}{\sum_{i=0}^l f(s_u|t_i)}}_{f(s_u|t_v) \text{ 对于所有} t_i \text{ 的相对值}$$

||

$P(\mathbf{s}|\mathbf{t})$

‘ 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**

- 啥是事件的期望频次？？？ - 事件在其分布下出现的次数的期望 $c_{\mathbb{E}}(X) = \sum_i c(x_i) \cdot P(x_i)$

x_i	$c(x_i)$	x_i	$c(x_i)$	$P(x_i)$	$c(x_i) \cdot P(x_i)$
x_1	2	x_1	2	0.1	0.2
x_2	1	x_2	1	0.3	0.3
x_3	5	x_3	5	0.2	1.0
$c(X) = 8$					

对 $f(s_u|t_v)$ 的计算进行化简

- 为了计算 $f(s_u|t_v)$ ，重新组织一下上面的公式

$$f(s_u|t_v) = \lambda_{t_v}^{-1} \underbrace{\frac{\epsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l f(s_j|t_i)}_{\text{翻译概率 } P(\mathbf{s}|\mathbf{t})} \underbrace{\sum_{j=1}^m \delta(s_j, s_u) \sum_{i=0}^l \delta(t_i, t_v)}_{\substack{(s_u, t_v) \text{ 在句对 } (\mathbf{s}, \mathbf{t}) \text{ 中} \\ \text{配对的总次数}}} \underbrace{\frac{f(s_u|t_v)}{\sum_{i=0}^l f(s_u|t_i)}}_{\substack{f(s_u|t_v) \text{ 对于所有} \\ \text{的 } t_i \text{ 的相对值}}}$$

||

$P(\mathbf{s}|\mathbf{t})$

'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**

- 啥是事件的期望频次？ - 事件在其分布下出现的次数的期望 $c_{\mathbb{E}}(X) = \sum_i c(x_i) \cdot P(x_i)$

x_i	$c(x_i)$
x_1	2
x_2	1
x_3	5
$c(X) = 8$	

x_i	$c(x_i)$	$P(x_i)$	$c(x_i) \cdot P(x_i)$
x_1	2	0.1	0.2
x_2	1	0.3	0.3
x_3	5	0.2	1.0
$c_{\mathbb{E}}(X) = 0.2 + 0.3 + 1.0 = 1.5$			

期望频次

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

期望频次

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

' t_v 翻译为 s_u '在所有
对齐中出现的次数

$c('t_v$ 翻译为 $s_u')$

$$\frac{f(s_u|t_v)}{\sum_{i=0}^l f(s_u|t_i)}$$

期望频次

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

' t_v ' 翻译为 ' s_u ' 在所有
对齐中出现的次数

$c('t_v$ 翻译为 ' s_u ')

$$\frac{f(s_u|t_v)}{\sum_{i=0}^l f(s_u|t_i)}$$

' t_v ' 翻译为 ' s_u ' 在所有
对齐中出现的相对概率

$P('t_v$ 翻译为 ' s_u ')

期望频次

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

' t_v 翻译为 s_u '在所有
对齐中出现的次数

$$\frac{f(s_u|t_v)}{\sum_{i=0}^l f(s_u|t_i)}$$

' t_v 翻译为 s_u '在所有
对齐中出现的相对概率

$$c('t_v \text{ 翻译为 } s_u') \times P('t_v \text{ 翻译为 } s_u')$$
$$= c_{\mathbb{E}}('t_v \text{ 翻译为 } s_u')$$

期望频次

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

' t_v 翻译为 s_u '在所有
对齐中出现的次数

$$\frac{f(s_u|t_v)}{\sum_{i=0}^l f(s_u|t_i)}$$

' t_v 翻译为 s_u '在所有
对齐中出现的相对概率

$$c('t_v \text{ 翻译为 } s_u') \times P('t_v \text{ 翻译为 } s_u') \\ = c_{\mathbb{E}}('t_v \text{ 翻译为 } s_u')$$

- 定义：在 $P(\mathbf{s}|\mathbf{t})$ 中， t_v 翻译(连接)到 s_u 的期望频次为

$$c_{\mathbb{E}}(s_u|t_v; \mathbf{s}, \mathbf{t}) \equiv \sum_{j=1}^m \delta(s_j, s_u) \sum_{i=0}^l \delta(t_i, t_v) \cdot \frac{f(s_u|t_v)}{\sum_{i=0}^l f(s_u|t_i)}$$

期望频次

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

' t_v 翻译为 s_u '在所有
对齐中出现的次数

$$\frac{f(s_u|t_v)}{\sum_{i=0}^l f(s_u|t_i)}$$

' t_v 翻译为 s_u '在所有
对齐中出现的相对概率

$$c('t_v \text{ 翻译为 } s_u') \times P('t_v \text{ 翻译为 } s_u') \\ = c_{\mathbb{E}}('t_v \text{ 翻译为 } s_u')$$

- 定义：在 $P(\mathbf{s}|\mathbf{t})$ 中， t_v 翻译(连接)到 s_u 的期望频次为

$$c_{\mathbb{E}}(s_u|t_v; \mathbf{s}, \mathbf{t}) \equiv \sum_{j=1}^m \delta(s_j, s_u) \sum_{i=0}^l \delta(t_i, t_v) \cdot \frac{f(s_u|t_v)}{\sum_{i=0}^l f(s_u|t_i)}$$

- 重写 $f(s_u|t_v)$!!!

$$f(s_u|t_v) = \lambda_{t_v}^{-1} \cdot P(\mathbf{s}|\mathbf{t}) \cdot c_{\mathbb{E}}(s_u|t_v; \mathbf{s}, \mathbf{t})$$

通过期望频次计算 $f(s_u|t_v)$

- 一个小**trick**: 令 $\lambda'_{t_v} = \frac{\lambda_{t_v}}{P(\mathbf{s}|\mathbf{t})}$

$$\begin{aligned} f(s_u|t_v) &= \lambda_{t_v}^{-1} \cdot P(\mathbf{s}|\mathbf{t}) \cdot c_{\mathbb{E}}(s_u|t_v; \mathbf{s}, \mathbf{t}) \\ &= (\lambda'_{t_v})^{-1} \cdot c_{\mathbb{E}}(s_u|t_v; \mathbf{s}, \mathbf{t}) \end{aligned}$$

通过期望频次计算 $f(s_u|t_v)$

- 一个小**trick**: 令 $\lambda'_{t_v} = \frac{\lambda_{t_v}}{P(\mathbf{s}|\mathbf{t})}$

$$\begin{aligned} f(s_u|t_v) &= \lambda_{t_v}^{-1} \cdot P(\mathbf{s}|\mathbf{t}) \cdot c_{\mathbb{E}}(s_u|t_v; \mathbf{s}, \mathbf{t}) \\ &= (\lambda'_{t_v})^{-1} \cdot c_{\mathbb{E}}(s_u|t_v; \mathbf{s}, \mathbf{t}) \end{aligned}$$

- λ'_{t_v} 究竟是什么? - 回忆一下IBM模型对 $f(\cdot|\cdot)$ 的约束

$$\forall t_y : \sum_{s_x} f(s_x|t_y) = 1$$

为了满足 $f(\cdot|\cdot)$ 的概率归一化约束, 显然

$$\lambda'_{t_v} = \sum_{s_u} c_{\mathbb{E}}(s_u|t_v; \mathbf{s}, \mathbf{t})$$

通过期望频次计算 $f(s_u|t_v)$

- 一个小**trick**: 令 $\lambda'_{t_v} = \frac{\lambda_{t_v}}{P(\mathbf{s}|\mathbf{t})}$

$$\begin{aligned} f(s_u|t_v) &= \lambda_{t_v}^{-1} \cdot P(\mathbf{s}|\mathbf{t}) \cdot c_{\mathbb{E}}(s_u|t_v; \mathbf{s}, \mathbf{t}) \\ &= (\lambda'_{t_v})^{-1} \cdot c_{\mathbb{E}}(s_u|t_v; \mathbf{s}, \mathbf{t}) \end{aligned}$$

- λ'_{t_v} 究竟是什么? - 回忆一下IBM模型对 $f(\cdot|\cdot)$ 的约束

$$\forall t_y : \sum_{s_x} f(s_x|t_y) = 1$$

为了满足 $f(\cdot|\cdot)$ 的概率归一化约束, 显然

$$\lambda'_{t_v} = \sum_{s_u} c_{\mathbb{E}}(s_u|t_v; \mathbf{s}, \mathbf{t})$$

- $f(s_u|t_v)$ 的计算式

$$f(s_u|t_v) = \frac{c_{\mathbb{E}}(s_u|t_v; \mathbf{s}, \mathbf{t})}{\sum_{s_u} c_{\mathbb{E}}(s_u|t_v; \mathbf{s}, \mathbf{t})}$$

在整个数据集上计算

- 更真实的情况：我们拥有一系列互译的句对（称作平行语料），记为 $\{(\mathbf{s}^{[1]}, \mathbf{t}^{[1]}), (\mathbf{s}^{[2]}, \mathbf{t}^{[2]}), \dots, (\mathbf{s}^{[N]}, \mathbf{t}^{[N]})\}$ 。对于这 N 个训练用句对，定义 $f(s_u|t_v)$ 的期望频次为

$$c_{\mathbb{E}}(s_u|t_v) = \sum_{i=1}^N c_{\mathbb{E}}(s_u|t_v; \mathbf{s}^{[i]}, \mathbf{t}^{[i]})$$

在整个数据集上计算

- 更真实的情况：我们拥有一系列互译的句对（称作平行语料），记为 $\{(\mathbf{s}^{[1]}, \mathbf{t}^{[1]}), (\mathbf{s}^{[2]}, \mathbf{t}^{[2]}), \dots, (\mathbf{s}^{[N]}, \mathbf{t}^{[N]})\}$ 。对于这 N 个训练用句对，定义 $f(s_u|t_v)$ 的期望频次为

$$c_{\mathbb{E}}(s_u|t_v) = \sum_{i=1}^N c_{\mathbb{E}}(s_u|t_v; \mathbf{s}^{[i]}, \mathbf{t}^{[i]})$$

- 于是

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

在整个数据集上计算

- 更真实的情况：我们拥有一系列互译的句对（称作平行语料），记为 $\{(\mathbf{s}^{[1]}, \mathbf{t}^{[1]}), (\mathbf{s}^{[2]}, \mathbf{t}^{[2]}), \dots, (\mathbf{s}^{[N]}, \mathbf{t}^{[N]})\}$ 。对于这 N 个训练用句对，定义 $f(s_u|t_v)$ 的期望频次为

$$c_{\mathbb{E}}(s_u|t_v) = \sum_{i=1}^N c_{\mathbb{E}}(s_u|t_v; \mathbf{s}^{[i]}, \mathbf{t}^{[i]})$$

- 于是

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

用当前的 $f(s_u|t_v)$
计算期望频次 $c_{\mathbb{E}}(\cdot)$

在整个数据集上计算

- 更真实的情况：我们拥有一系列互译的句对（称作**平行语料**），记为 $\{(\mathbf{s}^{[1]}, \mathbf{t}^{[1]}), (\mathbf{s}^{[2]}, \mathbf{t}^{[2]}), \dots, (\mathbf{s}^{[N]}, \mathbf{t}^{[N]})\}$ 。对于这 N 个训练用句对，定义 $f(s_u|t_v)$ 的期望频次为

$$c_{\mathbb{E}}(s_u|t_v) = \sum_{i=1}^N c_{\mathbb{E}}(s_u|t_v; \mathbf{s}^{[i]}, \mathbf{t}^{[i]})$$

- 于是

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

利用这个公式计算新的 $f(s_u|t_v)$ 值

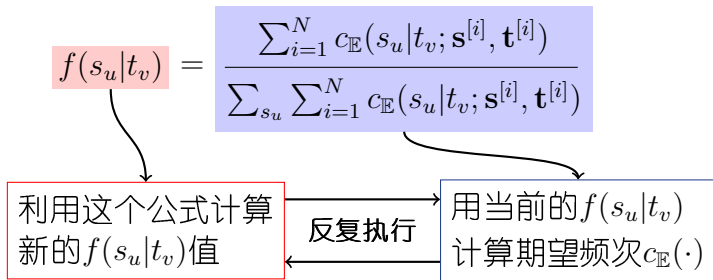
用当前的 $f(s_u|t_v)$ 计算期望频次 $c_{\mathbb{E}}(\cdot)$

在整个数据集上计算

- 更真实的情况：我们拥有一系列互译的句对（称作**平行语料**），记为 $\{(\mathbf{s}^{[1]}, \mathbf{t}^{[1]}), (\mathbf{s}^{[2]}, \mathbf{t}^{[2]}), \dots, (\mathbf{s}^{[N]}, \mathbf{t}^{[N]})\}$ 。对于这 N 个训练用句对，定义 $f(s_u|t_v)$ 的期望频次为

$$c_{\mathbb{E}}(s_u|t_v) = \sum_{i=1}^N c_{\mathbb{E}}(s_u|t_v; \mathbf{s}^{[i]}, \mathbf{t}^{[i]})$$

- 于是



完整的模型训练方法 - EM算法

IBM模型1的训练 (EM算法)

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

输出: 参数 $f(\cdot|\cdot)$ 的最优值

- 1: **Function** TRAINITWITHEM($\{(\mathbf{s}^{[1]}, \mathbf{t}^{[1]}), \dots, (\mathbf{s}^{[N]}, \mathbf{t}^{[N]})\}$)
- 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:
$$c_{\mathbb{E}}(s_u|t_v; \mathbf{s}^{[k]}, \mathbf{t}^{[k]}) = \sum_{j=1}^{|\mathbf{s}^{[k]}|} \delta(s_j, s_u) \sum_{i=0}^{|\mathbf{t}^{[k]}|} \delta(t_i, t_v) \cdot \frac{f(s_u|t_v)}{\sum_{i=0}^{|\mathbf{t}^{[k]}|} f(s_u|t_i)}$$
- 6: **foreach** t_v appears at least one of $\{\mathbf{t}^{[1]}, \dots, \mathbf{t}^{[N]}\}$ **do**
- 7: $\lambda'_{t_v} = \sum_{s_u} \sum_{k=1}^N c_{\mathbb{E}}(s_u|t_v; \mathbf{s}^{[k]}, \mathbf{t}^{[k]})$
- 8: **foreach** s_u appears at least one of $\{\mathbf{s}^{[1]}, \dots, \mathbf{s}^{[N]}\}$ **do**
- 9: $f(s_u|t_v) = \sum_{k=1}^N c_{\mathbb{E}}(s_u|t_v; \mathbf{s}^{[k]}, \mathbf{t}^{[k]}) \cdot (\lambda'_{t_v})^{-1}$
- 10: **return** $f(\cdot|\cdot)$

完整的模型训练方法 - EM算法

IBM模型1的训练 (EM算法)

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

输出: 参数 $f(\cdot|\cdot)$ 的最优值

- 1: **Function** TRAINITWITHEM($\{(\mathbf{s}^{[1]}, \mathbf{t}^{[1]}), \dots, (\mathbf{s}^{[N]}, \mathbf{t}^{[N]})\}$)
- 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:
$$c_{\mathbb{E}}(s_u|t_v; \mathbf{s}^{[k]}, \mathbf{t}^{[k]}) = \sum_{j=1}^{|\mathbf{s}^{[k]}|} \delta(s_j, s_u) \sum_{i=0}^{|\mathbf{t}^{[k]}|} \delta(t_i, t_v) \cdot \frac{f(s_u|t_v)}{\sum_{i=0}^t f(s_u|t_i)}$$
- 6: **foreach** t_v appears at least one of $\{\mathbf{t}^{[1]}, \dots, \mathbf{t}^{[N]}\}$ **do**
- 7:
$$\lambda'_{t_v} = \sum_{s_u} \sum_{k=1}^N c_{\mathbb{E}}(s_u|t_v; \mathbf{s}^{[k]}, \mathbf{t}^{[k]})$$
- 8: **foreach** s_u appears at least one of $\{\mathbf{s}^{[1]}, \dots, \mathbf{s}^{[N]}\}$ **do**
- 9:
$$f(s_u|t_v) = \sum_{k=1}^N c_{\mathbb{E}}(s_u|t_v; \mathbf{s}^{[k]}, \mathbf{t}^{[k]}) \cdot (\lambda'_{t_v})^{-1}$$
- 10: **return** $f(\cdot|\cdot)$

- **E-Step:** 行4-5, 计算 $c_{\mathbb{E}}(\cdot)$; **M-Step:** 行6-9, 计算 $f(\cdot)$

EM算法与IBM模型2

- EM算法可以直接用于训练IBM模型2(这里省去推导)

- ▶ 参数1: $f(s_j|t_i)$ - 与IBM模型1一样
- ▶ 参数2: $a(i|j, m, l)$ - 调序概率

1 E-Step (对于句对(s, t), $m = |s|, l = |t|$)

$$c_{\mathbb{E}}(s_u|t_v; \mathbf{s}, \mathbf{t}) = \frac{\sum_{j=1}^m \sum_{i=0}^l \frac{f(s_u|t_v)a(i|j, m, l)\delta(s_j, s_u)\delta(t_i, t_v)}{\sum_{k=0}^l f(s_u|t_k)a(k|j, m, l)}}$$

$$c_{\mathbb{E}}(i|j, m, l; \mathbf{s}, \mathbf{t}) = \frac{f(s_j|t_i)a(i|j, m, l)}{\sum_{k=0}^l f(s_j|t_k)a(k|j, m, l)}$$

2 M-Step

$$f(s_u|t_v) = \frac{\sum_{k=0}^K c_{\mathbb{E}}(s_u|t_v; \mathbf{s}^{[k]}, \mathbf{t}^{[k]})}{\sum_{s_u} \sum_{k=0}^K c_{\mathbb{E}}(s_u|t_v; \mathbf{s}^{[k]}, \mathbf{t}^{[k]})}$$

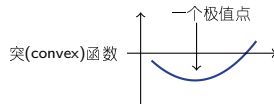
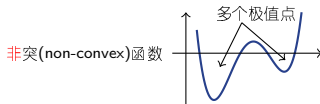
$$a(i|j, m, l) = \frac{\sum_{k=0}^K c_{\mathbb{E}}(i|j; \mathbf{s}^{[k]}, \mathbf{t}^{[k]})}{\sum_i \sum_{k=0}^K c_{\mathbb{E}}(i|j; \mathbf{s}^{[k]}, \mathbf{t}^{[k]})}$$

IBM模型1和2基本介绍完了

- 再来回顾一下，IBM模型1和2包括以下三方面内容：
 - ① 建模 - 第25页~第36页
 - ② 解码 - 第38页
 - ③ 训练 - 第56页
- 还有很多很多问题没有讨论
 - ▶ 解码：第16页所描述的解码算法很粗糙
 - 需要考虑更多的翻译假设 - beam search - 随后章节会介绍

IBM模型1和2基本介绍完了

- 再来回顾一下，IBM模型1和2包括以下三方面内容：
 - 1 建模 - 第25页~第36页
 - 2 解码 - 第38页
 - 3 训练 - 第56页
- 还有很多很多问题没有讨论
 - ▶ 解码：第16页所描述的解码算法很粗糙
 - 需要考虑更多的翻译假设 - beam search - 随后章节会介绍
 - ▶ 训练：EM算法的问题是：
 - 在目标函数为非突(凹)集函数时，容易陷入局部最优



比较幸运的是IBM模型1有单个极值点，EM可以得到全局最优。而且通常会把IBM模型1的参数作为高阶模型(如模型2)的输入，保证后面的训练从较好的初始值开始

- 会有degenerate analysis的问题
- Google一下variational EM和bayesian inference

- IBM模型3-5不是本教程的重点
 - ▶ 建模相对复杂一些，当然模型也更强大
 - ▶ 训练和解码不太简单
 - ▶ 推荐一些阅读材料
 - **A Program for Aligning Sentences in Bilingual Corpora.** William A. Gale and Kenneth W. Church. 1993.
 - **A systematic comparison of various statistical alignment models.** Franz Och and Hermann Ney. 2003.
 - **基于IBM模型的统计机器翻译技术的研究.** 肖桐. 2008. (硕士论文)

GIZA++

- IBM模型3-5不是本教程的重点
 - ▶ 建模相对复杂一些，当然模型也更强大
 - ▶ 训练和解码不太简单
 - ▶ 推荐一些阅读材料
 - **A Program for Aligning Sentences in Bilingual Corpora.** William A. Gale and Kenneth W. Church. 1993.
 - **A systematic comparison of various statistical alignment models.** Franz Och and Hermann Ney. 2003.
 - **基于IBM模型的统计机器翻译技术的研究.** 肖桐. 2008. (硕士论文)
- 不得不提的**GIZA++**：实现了IBM模型+HMM模型，现在广泛用于双语平行语料的自动词对齐
 - ▶ 完全开源
<http://code.google.com/p/giza-pp/>
 - ▶ NiuTrans中也有使用方式的介绍
<http://www.nlplab.com/NiuPlan/NiuTrans.YourData.ch.html>

这章内容先写到这里

这章内容写了很多，但是还没弄完 :(
后面(可能过两年)补上单词自动对齐的内容 :)

谢谢	大家
	\
thank	you all