

语言不是由随机抽取的一些单词组成的序列，词与词之间是有联系的。这种规律用概率模型刻画。我们用语言模型来称呼单词序列的统计模型。

例1 一个人所说的话中每100个句子有一句是“OK”，则可认为 $P(\text{OK}) \approx 1\%$

例2. _____ 和阿里、腾讯 并称为中国互联网BAT三巨头。

$P(\text{百度和阿里、腾讯并称为...}) \approx 1$

例2揭示了一件事：上下文信息越强，词语可预测性越高。

设词语序列 $S = w_1 w_2 \dots w_L$ ，根据链式法则。

$$P(S) = P(w_1 w_2 \dots w_L) = P(w_1 | \langle \text{BOS} \rangle) P(w_2 | w_1) P(w_3 | w_1 w_2) \dots P(w_L | w_1 \dots w_{L-1})$$

即 $P(S) = \prod_{i=1}^L P(w_i | w_1 \dots w_{i-1})$

每个条件概率都是模型的一个参数，训练模型本质上是

① 确定模型参数集合 ② 确定各参数的值

在 n -gram 语言模型中，我们用最大似然估计来算条件概率。

$$P(w_i | w_1 w_2 \dots w_{i-1}) = \frac{C(w_1 w_2 \dots w_{i-1} w_i)}{\sum_{w_i} C(w_1 w_2 \dots w_{i-1} w_i)}$$

C 是 Count 的缩写，分子表示 $(w_1 \dots w_i)$ 在语料库中出现次数
分母表示所有可能的下一个词 w_i 与前面序列组合出现总次数

如 $w_1 \dots w_{i-1}$ 表示“我喜欢”，分子表示 $C(\text{我喜欢吃})$ (w_i 固定为“吃”) 而分母表示 $C(\text{我喜欢吃}) + C(\text{我喜欢看}) + C(\text{我喜欢你}) + \dots$
即“我喜欢 + 任意词”的次数总和

存在问题：参数空间过大， $P(w_i | w_1 \dots w_{i-1})$ 在训练语料中趋近于 0。
数据严重稀疏

n元语法

马尔可夫假设：一个词的出现仅依赖它前一个或几个词，受离它近的词影响大

$$P(\text{什么} | \text{你今天} \text{在干}) \approx P(\text{什么} | \text{干})$$

二元语法模型：一个词的出现仅依赖其前一个词

$$P(w_1 \dots w_L) = \prod_{i=1}^L P(w_i | w_{i-1})$$

$$P(\text{我爱你}) = P(\text{我} | \langle \text{BOS} \rangle) \times P(\text{爱} | \text{我}) \times P(\text{你} | \text{爱}) \times P(\langle \text{EOS} \rangle | \text{你})$$

注：
<BOS> beginning of sequence 序列起始符
<EOS> End of sequence 序列终止符
<PAD> padding 填充

例

二元语法模型参数 $P(w_i | w_{i-1})$ 的最大似然估计

例

序号	w_{i-1}	w_i	$C(w_{i-1}, w_i)$	$P(w_i w_{i-1})$
1	父亲	写	50	50/100 = 0.5
2	父亲	看	40	40/100 = 0.4
3	父亲	吃	10	10/100 = 0.1
总计			100	1

$$P(w_i | w_{i-1}) = \frac{P(w_{i-1} w_i)}{P(w_{i-1})} = \frac{C(w_{i-1} w_i) / N}{C(w_{i-1}) / N} = \frac{C(w_{i-1} w_i)}{C(w_{i-1})}$$

$C(w_a w_b) \rightarrow w_a w_b$ 在给定文本中出现的次数

$N \rightarrow$ 训练语料中词的个数

例、训练语料：中国万岁 中国中国 万岁中国 万岁万岁

一元语法： $P(\text{中国}) = \frac{4}{8}$ $P(\text{万岁}) = \frac{4}{8}$

二元语法： $P(\text{中国} | \langle \text{BOS} \rangle) = \frac{2}{4}$ $P(\text{万岁} | \langle \text{BOS} \rangle) = \frac{2}{4}$

$P(\text{中国} | \text{中国}) = \frac{1}{4}$ $P(\text{万岁} | \text{中国}) = \frac{1}{4}$

$P(\text{中国} | \text{万岁}) = \frac{1}{4}$ $P(\text{万岁} | \text{万岁}) = \frac{1}{4}$

$$P(\langle \text{EOS} \rangle | \text{中国}) = \frac{1}{4} \quad P(\langle \text{EOS} \rangle | \text{万岁}) = \frac{1}{4}$$

$$\begin{aligned} \therefore P(\text{中国万岁}) &= P(\text{中国} | \langle \text{EOS} \rangle) \times P(\text{万岁} | \text{中国}) \times \\ &P(\langle \text{EOS} \rangle | \text{万岁}) = \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} = \frac{1}{16} \end{aligned}$$

推广: n 元语法模型: 一个词的出现依赖其前 $n-1$ 个词

$$P(w_1, \dots, w_n) = \prod_{i=1}^n P(w_i | w_{i-n+1}, \dots, w_{i-1})$$

n 的选择 $\left\{ \begin{array}{l} \text{大} \quad \text{约束的, 辨别强} \\ \text{小} \quad \text{训练集数据多, 更可靠} \end{array} \right.$

理论上 n 越大越好, 实际上二元和三元用得较多

Laplace 平滑

最大似然估计的问题: 训练语料不可能囊括所有合理的词语搭配, 一旦出现未在训练语料中囊括的词语搭配, 那么其 P 就为 0, 模型完全无法预测

$$\text{Laplace 平滑: } P_{\text{MLE}}(w_i | w_{i-1}) = \frac{C(w_{i-1} w_i)}{\sum w_i C(w_{i-1} w_i)} = \frac{C(w_{i-1} w_i)}{C(w_{i-1})}$$



$$P_{\text{Lap}}(w_i | w_{i-1}) = \frac{C(w_{i-1} w_i) + 1}{\sum w_i [C(w_{i-1} w_i) + 1]} = \frac{C(w_{i-1} w_i) + 1}{C(w_{i-1}) + |V|}$$

$|V|$ 表示词汇表中单词个数

例、

• $p(\text{Father read a book})$

Father read Holy Bible
Mother read a text book
He read a book by Grandpa

$V = \{\text{Father}, \text{read}, \dots, \text{Grandpa} \langle \text{BOS} \rangle, \langle \text{EOS} \rangle\}$
 $|V| = 13$

$$p(\text{Father read a book}) = p(\text{Father} \langle \text{BOS} \rangle) \times p(\text{read} | \text{Father}) \times p(a | \text{read}) \times p(\text{book} | a) \times p(\langle \text{EOS} \rangle | \text{book})$$

MLE:

LAP:

$$= \frac{1}{3} \times 1 \times \frac{2}{3} \times \frac{1}{2} \times \frac{1}{2} = 0.06$$

$$= \frac{2}{16} \times \frac{2}{14} \times \frac{3}{16} \times \frac{2}{15} \times \frac{2}{15} \approx 0.00006$$

• $p(\text{Grandpa read a book})$

Father read Holy Bible
Mother read a text book
He read a book by Grandpa

$V = \{\text{Father}, \text{read}, \dots, \text{Grandpa} \langle \text{BOS} \rangle, \langle \text{EOS} \rangle\}$
 $|V| = 13$

$$p(\text{Grandpa read a book}) = p(\text{Grandpa} \langle \text{BOS} \rangle) \times p(\text{read} | \text{Grandpa}) \times p(a | \text{read}) \times p(\text{book} | a) \times p(\langle \text{EOS} \rangle | \text{book})$$

$$= \frac{1}{16} \times \frac{1}{14} \times \frac{3}{16} \times \frac{2}{15} \times \frac{2}{15} \approx 0.000015$$

Laplace 平滑一般形式:

$$P_{\text{lap}}(w_i | w_{i-n+1}^{i-1}) = \frac{C(w_{i-n+1}^i) + 1}{\sum_{w_i} C(w_{i-n+1}^i) + |V|}$$
$$= \frac{C(w_{i-n+1}^i) + 1}{C(w_{i-n+1}^{i-1}) + |V|}$$

其中 $w_i^j = w_i \dots w_j$

在前一个例子中, 我们发现 Laplace 平滑后句子出现概率非常小, 这是由于直接加 1 会将概率稀释得很厉害, 所以我们有:

Lidstone 法则:

$$P_{\text{lid}}(w_i | w_{i-(n-1)}^{i-1}) = \frac{C(w_{i-(n-1)}^{i-1}) + \delta}{C(w_{i-(n-1)}^{i-1}) + \delta |V|}$$

其中 δ 是一个很小的正数, 如 0.1, 0.01 等.

比 Laplace 更温和、实用

- 元先验平滑: 给所有同一个“保底概率”, 无论其是否在训练集中出现过.

$$P(w_i) = \frac{C(w_{i-1}w_i) + m P(w_i)}{C(w_{i-1}) + m}$$

Good-Turning 估计

对于任何一个出现 r 次的词 (这记作 n 元语法) 都假设其发生 r^* 次

$$r^* = (r+1) \frac{n_{r+1}}{n_r} \quad n_r \text{ 表示训练语料中出现 } r \text{ 次词的数量}$$

r	0	1	2	3	4	...
n_r	n_0	n_1	n_2	n_3	n_4	...
r^*	$\frac{n_1}{n_0}$	$2 \frac{n_2}{n_1}$	$3 \frac{n_3}{n_2}$	$4 \frac{n_4}{n_3}$	$5 \frac{n_5}{n_4}$...

这样便平滑了低频次
的频率, 避免了“出现 0
次的词概率为 0”

总词数验证: $\sum r \cdot n_r = n_1 + 2n_2 + 3n_3 + \dots = N$

修正后: $\sum r^* n_r = n_1 + 2n_2 + 3n_3 + \dots = N$

∴ 总词数无变化

$$P_{GT}(w) = \frac{r^*}{N} = (r+1) \frac{n_{r+1}}{n_r \cdot N}$$

缺陷: 在实际统计中, 会发现一个现象: 低频次出现次数远高于高频词, 且高频词数量波动大, 这会导致公式低估低频词而严重失真

改进: Katz 平滑

只对 $r \leq k$ 的词修正 $r > k$ 的词直接用原次数
(k 一般取 5)

$$\text{即 } \begin{cases} r=0 & r^* = \frac{n_1}{n_0} \\ 1 \leq r \leq k & r^* = \frac{(r+1) \frac{n_{r+1}}{n_r} - r \frac{(k+1)n_{k+1}}{n_1}}{1 - \frac{(k+1)n_{k+1}}{n_1}} \\ r > k & r^* = r \end{cases}$$

$$P_{\text{katz}}(w) = \begin{cases} \frac{r^*}{N} & r \leq k \\ \frac{r}{N} & r > k \end{cases}$$

绝对折扣

$$P_{\text{abs}}(w) = \begin{cases} \frac{r - \delta}{N} & r > 0 \\ \frac{(\sum_{r=1}^{\infty} n_r) \delta}{n_0 N} & r = 0 \end{cases}$$

$r > 0$: 这个词出现次数减去一个固定值 δ (通常很小)
再除以总词数 N

$r = 0$: 未登录词的概率 = 所有词扣掉的总折扣除以未登录词的数量与总词数之积, 即将扣掉的概率平均分给未登录词

r (次数)	n_r (词数)	$P_{\text{MLE}}(w)$	$P_{\text{abs}}(w)$ ($\delta=0.01$)
0	16	0	$(40-16) \times 0.01 \div 16 \div 50 = 0.0003$
1	10	0.02	$0.99 / 50 = 0.0198$
2	6	0.04	$1.99 / 50 = 0.0398$
3	5	0.06	$2.99 / 50 = 0.0598$
4	2	0.08	$3.99 / 50 = 0.0798$
5	1	0.10	$4.99 / 50 = 0.0998$
总计	$V=40$ $N=50$		

线性折扣

$$P_{\text{LD}}(w) = \begin{cases} \frac{(1-\alpha)r}{N} & r > 0 \\ \frac{\alpha}{n_0} & r = 0 \end{cases}$$

α 是折扣比例 (通常很小)

$r > 0$: 每个词按比例 $(1-\alpha)$ 打折

$r=0$: 将 α 平均分给未登录词

组合估计

假定我们要在一堆语料上构建三元语法模型

其中 $C(\text{he send the}) = 0$ $C(\text{he send thou}) = 0$

无论是 Lap, GT, abs 都有:

$$P(\text{the} | \text{he send}) = P(\text{thou} | \text{he send})$$

但是直觉上 $P(\text{the} | \text{he send})$ 更大一些

这是因为在二元语法模型中

$$P(\text{the} | \text{send}) > P(\text{thou} | \text{send})$$

低阶模型可以给高阶模型提供有效信息。

常用组合估计模型 { 线性插值 (混合一、二、三元)
回退 (三元优先, 二元次之, 一元最后)

线性插值

三元:

$$P_{\text{interp}}(w_i | w_{i-2} w_{i-1}) = \lambda_3 P_{ML}(w_i | w_{i-2} w_{i-1}) + \lambda_2 P_{ML}(w_i | w_{i-1}) + \lambda_1 P_{ML}(w_i)$$

$$0 \leq \lambda_i \leq 1 \text{ 且 } \sum \lambda_i = 1$$

假设 $C(\text{he send the}) = C(\text{he send thou}) = 0$

$$P_{\text{interp}}(\text{the} | \text{he send}) = \lambda_3 \frac{C(\text{he send the})}{C(\text{he send})} + \lambda_2 \frac{C(\text{send the})}{C(\text{send})} + \lambda_1 \frac{C(\text{the})}{N}$$

$$P_{\text{interp}}(\text{thou} | \text{he send}) = \lambda_3 \frac{C(\text{he send thou})}{C(\text{he send})} + \lambda_2 \frac{C(\text{send thou})}{C(\text{send})} + \lambda_1 \frac{C(\text{thou})}{N}$$

二元

$$P_{\text{interp}}(w_i | w_{i-1}) = \lambda P_{\text{ML}}(w_i | w_{i-1}) + (1-\lambda) P_{\text{ML}}(w_i)$$

语言模型性能评价

- 外部评估：将语言模型用于其它任务，以任务的指标衡量模型（问题：不够客观）
- 内部评估：困惑度（问题：不够精确）

困惑度 (perplexity)：测试集概率的倒数，直观来说，就是模型预测某句话的下一个词，猜得准，困惑度低

$$\begin{aligned} PP(W) &= P(w_1 \dots w_n)^{-\frac{1}{n}} = \sqrt[n]{\frac{1}{P(w_1 \dots w_n)}} \\ &= \sqrt[n]{\prod_{i=1}^n \frac{1}{P(w_i | w_1 \dots w_{i-1})}} \quad (W = w_1 \dots w_n) \end{aligned}$$

对于二元语法：
$$PP(W) = \sqrt[n]{\prod_{i=1}^n \frac{1}{P(w_i | w_{i-1})}}$$

$PP(W)$ 是对 W 这个句子的困惑度，而模型的困惑度会综合许多 $PP(W)$ 来计算（如取均值）

例：句子由完全随机的数字 (0~9) 构成，则困惑度是多少

即：每个位置出现数字 i 的概率为 $\frac{1}{10}$

$$PP(W) = \sqrt[n]{\left(\frac{1}{10}\right)^n} = 10$$

