

第5章 语言模型

基本概念

基于大规模语料库和统计方法, 可以

- 发现语言使用的普遍规律
- 进行机器学习, 自动获取语言知识
- 对未知语言现象进行预测

计算一段文字(如句子)的概率

(以词为单位) $p(w_1) \times p(w_2) \times \dots \times p(w_n)$

语句 $s = w_1 w_2 \dots w_m$ 的先验概率:

$$p(s) = p(w_1) \times p(w_2 | w_1) \times p(w_3 | w_1 w_2) \times \dots \times p(w_m | w_1 \dots w_{m-1})$$

$$= \prod_{i=1}^m p(w_i | w_1 \dots w_{i-1})$$

当 $i=1$ 时, $p(w_1 | w_0) = p(w_1)$

注: 1) w_i 可以是字、词、短语或词类符, 称为统计单元

2) 由 w_1, \dots, w_{i-1} 构成的一个序列, 称为 w_i 的历史

问题: 自由参数的个数可达 L^m

解决方法: 减少历史单元的个数, 即将 w_1, w_2, \dots, w_{i-1} 映射到等价类

$sc(w_1, w_2, \dots, w_{i-1})$, 使等价类的数目远小于原来不同历史单元的数目

$$p(w_i | w_1, \dots, w_{i-1}) = p(w_i | sc(w_1, \dots, w_{i-1}))$$

等价类的划分: 假设只与前 $n-1$ 个单元有关

即: n 元文法模型 (n -gram)

$n=1$ 一元文法 独立于历史 unigram monogram
 $n=2$ 二元文法 独立于历史 bigram monogram
 $n=3$ 三元文法 独立于历史 trigram monogram

为保证条件概率在 $i=1$ 时有意义, 同时为了保证句内所有字符串的概率和为 1, 我们:

$$\langle BOS \rangle w_1 w_2 \dots w_m \langle EOS \rangle$$

不失一般性, 对于 $n \geq 2$ 的 n -gram, $p(s)$ 可以写为:

$$p(s) = \prod_{i=1}^{m+n-1} p(w_i | w_{i-n+1}^{i-1})$$

其中, w_i^j 表示词序列 $w_1 \dots w_j$, w_{i-n+1}^{i-1} 从 w_0 开始, w_0 为 $\langle BOS \rangle$, w_{m+1} 为 $\langle EOS \rangle$

例: $\langle BOS \rangle$ John read a book $\langle EOS \rangle$

基于二元文法的概率为:

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

$$p(\text{read} | \text{John}) \times p(\text{a} | \text{read}) \times$$

$$p(\text{book} | \text{a}) \times p(\langle EOS \rangle | \text{book})$$



应用: 1. 音节转写问题

$$\begin{aligned}\hat{c}_{string} &= \arg \max_{c_{string}} P(c_{string} | p_{nyn}) \\ &= \arg \max_{c_{string}} \frac{P(p_{nyn} | c_{string}) \times P(c_{string})}{P(p_{nyn})}\end{aligned}$$

$$= \arg \max_{c_{string}} P(c_{string})$$

c_{string} 的概率即可使用 2-gram 计算
一元语法样本空间为 \mathcal{N} , 二元为 $\mathcal{N}^2 \dots$

汉字: 四元模型

应用实例: 破译断与输入法.

2. 汉字转写问题

5.2 参数估计

训练语料 training data:

最大似然估计 (maximum likelihood evaluation, MLE): 相对频率计算概率

对于 n -gram, 参数 $P(w_i | w_{i-n+1}^{i-1})$ 可由最大似然估计求得

$$P(w_i | w_{i-n+1}^{i-1}) = f(w_i | w_{i-n+1}^{i-1}) = \frac{C(w_i | w_{i-n+1}^{i-1})}{\sum_{w_i} C(w_i | w_{i-n+1}^{i-1})}$$

例: training data: "John read Moby Dick"

...

$$P(\text{John} | \langle \text{BOS} \rangle) = \dots$$

数据平滑: 解决数据匮乏 (稀疏) (sparse data) 引起的零概率问题
data smoothing

基本目标: 测试样本的语言模型困惑度越小越好

基本约束: $\sum_{w_i} P(w_i | w_1, w_2, \dots, w_{i-1}) = 1$

对于平滑的 n -gram, 其概率为 $P(w_i | w_{i-n+1}^{i-1})$, 可以计算句子的概率:

$$P(S) = \prod_{i=1}^{n+1} P(w_i | w_{i-n+1}^{i-1})$$

假设测试语料 T 由 L_T 个句子构成 (t_1, \dots, t_{L_T}) , 则整个测试集的概率为

$$P(T) = \prod_{i=1}^{L_T} P(t_i)$$



模型 $p(u_2/u_{1:n+1})$ 对于测试语料的交叉熵:

$$H_p(T) = -\frac{1}{w_T} \log_2 p(T)$$

其中, w_T 是训练文本 T 的词数

模型 p 的困惑度 $PP_p(T)$ 定义为: $PP_p(T) = 2^{H_p(T)}$

ngram 对于英语文本的困惑度,

50 ~ 1000

交叉熵: 6 ~ 10 bits/word

数据平滑方法:

1. 加1法: $\frac{1}{6} \quad 0 \quad \frac{5}{6}$
Additive smoothing $\frac{1}{6} \quad \frac{1}{6} \quad \frac{3}{6}$

对于 2-gram: $p(u_2/u_{1:n}) = \frac{(1 + c(u_1, u_2))}{\sum u_2 (1 + c(u_1, u_2))}$

2. 成值法/折扣法 (Discounting)

修改 training data 中的实际计数, 使每个的 ~~概率~~ 概率 < 1 . 效果分配

① Good-Turning 估计

N n_r (样本中正好出现 r 次的事件的数目)

$$N = \sum_{r=1}^{\infty} n_r r$$

$$= \sum_{r=0}^{\infty} n_r r^* = \sum_{r=0}^{\infty} (r+1) n_{r+1} \quad c r^* \approx c r \text{ 较小}$$

$$\text{即 } r^* = (r+1) \frac{n_{r+1}}{n_r}$$

那么, Good-Turning 估计在样本中出现 r 次的事件的概率为: $p_r = \frac{r^*}{N}$

适用于: 大词汇量产生的符号的统计分布的大量观测数据

归一化 $\hat{p}_r = \frac{p_r}{\sum_r p_r}$

② Back-off 估计/后退/方法 Katz 后退法

对非 0 事件按 Good-Turning 法计算成值, 若值出来的概率按估计分布归一化

0 概率事件

$$p_{\text{Katz}}(u_2/u_{1:n}) = \begin{cases} dr \frac{c(u_1, u_2)}{c(u_1)} & \text{if } (u_1, u_2) = r > 0 \\ d(u_1) p_{\text{ML}}(u_2) & \text{if } (u_1, u_2) = 0 \end{cases}$$

$dr (0 < dr < 1)$ 折扣率
 $\approx r^*/r$

$p_{\text{ML}}(u_2) = \lfloor u_2 \rfloor$ 的最大似然估计



$\alpha(w_{i-1})$ 的确定:

$$\sum_{u_i} p_{katz}(u_i | w_{i-1}) = 1$$

$$\sum_{u_i: r=0} \alpha(w_{i-1}) p_{ML}(u_i) + \sum_{u_i: r>0} p_{katz}(u_i | w_{i-1}) = 1$$

$$\alpha(w_{i-1}) = \frac{1 - \sum_{u_i: r>0} p_{katz}(u_i | w_{i-1})}{\sum_{u_i: r=0} p_{ML}(u_i)}$$

适用于 2 元/2 元上模型

③ 绝对成值法 (Absolute discounting)

从每个计数 r 中减去同样的量, 使得后未见事件

$$p_r = \begin{cases} \frac{r-b}{n} & r > 0 \\ \frac{b(R-n_0)}{n n_0} & r = 0 \end{cases}$$

n_0 : 未出现事件的数目 b : b 叫减去的常量

$\frac{b(R-n_0)}{n}$: 由于成值而产生的校正概率

④ 线性成值法.

$$p_r = \begin{cases} \frac{c(r)r}{n} & r > 0 \\ \frac{c}{n_0} & r = 0 \end{cases}$$

$c = \frac{n_1}{n}$ 或从每个计数 r 中减去与 r 成比例的量

绝对成值优于线性成值

3. 线性插值法.

用低阶语法估计高阶语法

3-gram \rightarrow 4-gram

2-gram \rightarrow 3-gram

1-gram \rightarrow 2-gram

$$p(w_3 | w_{1:2}) = \lambda_1 p'(w_3 | w_1 w_2) + \lambda_2 p'(w_3 | w_2) + \lambda_3 p'(w_3)$$

$$\lambda_1 + \lambda_2 + \lambda_3 = 1$$

$\lambda_1, \lambda_2, \lambda_3$ 的确定: 将训练语料分为 2 部分, 即从原始中删一部分作为验证数据, 另一部分用于



估计 $p'(w_3|w_1w_2)$, $p'(w_3|w_2)$ 及 $p'(w_3)$
 需一部分用于计算 $\lambda_1, \lambda_2, \lambda_3$, 使用越左越小

5.4 语言模型的自适应

问题: ①. 综合性语料难以反映不同领域之间在语言使用规律上的差异, 而语言模型恰恰对训练文本的类型、主题、风格等都十分敏感

②. n -元语言模型的独立性假设在很多时候都明显不成立
 自适应方法:

① 基于缓存的语言模型 (cache-based LM)

② 基于混合方法的语言模型

③ 基于最大熵的语言模型

①. 在文本中刚出现的词在后边句子再次出现的可能性往往较大
 线性插值: $p(w_i|w_1^{i-1}) = \lambda p_{\text{cache}}(w_i|w_1^{i-1}) + (1-\lambda) p_{\text{gram}}(w_i|w_1^{i-n+1})$

λ 可通过 EM 算法求得.

cache: 1000 ~ 3000

$$\hat{p}_{\text{cache}}(w_i|w_1^{i-1}) = \frac{1}{L} \sum_{j=i-L}^{i-1} I\{w_j = w_i\}$$

I : 指示函数, 如果 z 表示的情况出现, 则 $I_z = 1$, 否则 $I_z = 0$

缺陷: 缓存里的词对当前词的影响是一样的

改进: $\hat{p}_{\text{cache}}(w_i|w_1^{i-1}) = \beta \sum_{j=1}^{i-1} I_{\{w_j = w_i\}} e^{-\alpha(i-j)}$

α : 衰变率 β : 归一化系数 $\sum_{w \in V} \hat{p}_{\text{cache}}(w_i|w_1^{i-1}) = 1$

2). 来将不同, 组成同族 (homogeneous)

划分为不同 n 个子模型 m_1, \dots, m_n (聚类)

$$\hat{p}(w_i|w_1^{i-1}) = \sum_{j=1}^n \lambda_j p_{m_j}(w_i|w_1^{i-1})$$

其中, $0 \leq \lambda_j \leq 1$, $\sum_{j=1}^n \lambda_j = 1$

λ 值通过 EM 算法计算



EM 迭代计算维数参数：

a) 对于 n 个类，随机初始化维数参数

b) 计算新的维数参数期望

c) $\lambda_{ij}^t = \frac{\lambda_{ij}^{t-1} p_{ij}(w|h)}{\sum_{i=1}^n \lambda_{ij}^{t-1} p_{ij}(w|h)}$ h 为历史

d) 迭代直至收敛

理论收敛，实际一般为伪

⑤. 两个语言模型 M_1 和 M_2

$$p_{M_1}(w_i | w_{1:i-1}) = f(w_i, w_{1:i-1})$$

M_2 是距离为 2 的 2-gram 模型 $\hat{p}_{M_2}(w_i | w_{1:i-1}) = g(w_i, w_{i-2})$

M_1 是标准 2-gram 模型

线性插值取这两个概率估计的平均值，使用 GIS 算法选择使熵最大的模型。

5.5 语言模型应用举例

1. 汉语分词 $\hat{w} \approx \arg \max_w p(w)$

其中 $s = s_1 s_2 \dots s_m$, $w = w_1 w_2 \dots w_k$ 是一种可能的切分。

具体实现时，可将汉语词汇分成如下几类：

1) 词典中

2) 词典规则派生出来的词

3) 与数字相关的术语

4) 专用名词

把一个可能的词序列 w 转换成词类序列 $C = c_1 c_2 \dots c_n$

① c_i 为 PV , c_i 为 LV , 机构为 OV

② 日期 dat , 名词 nm , 形容词 pr , 动词 mv

③. 词类组合词 mv , 词类词 LV , 每个词单独一类



将同义列写成类别列

$$C = \arg \max p(C|S) \\ = \arg \max p(C) \times p(S|C)$$

↑ ↑
语言模型 生成模型

$$p(C) = p(c_1) \times p(c_2|c_1) \times \prod_{i=3}^N p(c_i|c_{i-2})$$

极大似然估计

$$p(S|C) \approx \prod_{i=1}^N p(s_i|c_i)$$

近似假设

例如：如果“教授”是同表中的词，那么 $p(s_2 = \text{教授} | c_1 = \text{kw}) = 1$ ，否则 $p(s_2 | c_1) = 0$

生成模型则为判断同义类别

2. 分词与同义标注一体化系统

句子： $w = w_1 w_2 \dots w_n$

单词 $w_i (1 \leq i \leq n)$ 的同义标注为 t_i ， $T = t_1 t_2 \dots t_n$

任务：在 S 所对应的各种标注和标注形式中，寻找最优的 $p(w, T)$

(1) 基于同义的三元统计模型

$$p(w, T) = p(w|T) \times p(T) \approx \prod_{i=1}^n p(w_i|t_i) \times p(t_i|t_{i-1}, t_{i-2}) \dots$$

$p(w|T)$ 为生成模型， $p(T)$ 为基于同义的统计模型。

(2) 基于单词的三元统计模型

$$p(w, T) = p(T|w) \times p(w) \approx \prod_{i=1}^n p(t_i|w_i) \times p(w_i|w_{i-1}, w_{i-2})$$

(3) 分词与同义标注一体化模型

$$p^*(w, T) = \alpha \prod_{i=1}^n p(w_i|t_i) \times p(t_i|t_{i-1}, t_{i-2}) + \beta \prod_{i=1}^n p(t_i|w_i) \times p(w_i|w_{i-1}, w_{i-2})$$

指导思想：通过调整参数 α 和 β 的值来确定两个子模型在整个分词与同义标注中所发挥作用的比重。



分析公式 (2)

$p(t_i | w)$ 对当前词无帮助, 且作为词确证后对词性标注会增添偏见
因此删掉, 并令 $\alpha = 1$.

$$p(w, T) = \prod_{i=1}^n p(w_i | t_i) \times p(t_i | t_{i-1}, t_{i+1}) + \beta \prod_{i=1}^n p(w_i | w_{i-1}, w_{i+1})$$

确定 β 的值: $\beta = \frac{\text{词典中词 } w \text{ 的个数}}{\text{词性 } t \text{ 的种类数}}$

