

信息论基础

熵: $H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$

约定 $0 \log 0 = 0$

$p(x) \downarrow \rightarrow H(X) \uparrow \rightarrow$ 不确定性 \uparrow

[信源: 每个符号所提供的平均信息量]

[在自然语言处理中的应用:

使熵 $H(X)$ 值最大的模型用来推断某种语言现象存在的可能性]

联合熵: $H(X, Y) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(x, y)$

描述一对随机变量平均所需要的信息量

条件熵: $H(Y|X) = \sum_{x \in X} p(x) H(Y|X=x)$

$$= \sum_{x \in X} p(x) \left[- \sum_{y \in Y} p(y|x) \log_2 p(y|x) \right]$$

$$= - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(y|x)$$

链式法则 $H(X, Y) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log [p(x) p(y|x)]$

$$= - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(x) - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(y|x)$$

$$= - \sum_{x \in X} p(x) \log p(x) - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(y|x)$$

$$= H(X) + H(Y|X) \quad \text{同理 } H(Y) + H(X|Y)$$

例: $H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$

$$= - \left(\frac{1}{4} \times \log_2 \frac{1}{4} + \frac{1}{4} \times \log_2 \frac{1}{4} + \frac{1}{8} \times \log_2 \frac{1}{8} + \frac{1}{8} \times \log_2 \frac{1}{8} \right)$$

$$= - \left(-\frac{1}{2} + (-\frac{1}{2}) + (-\frac{1}{8}) \right) = \frac{7}{8} \text{ (bits)}$$

~~$H(X|Y) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(x|y)$~~

~~$= - \sum_{y \in Y} p(y) \log p(y)$~~

$$= \sum_{y \in Y} p(y=1) H(X|Y=1)$$

$$= \frac{1}{4} \cdot \left(- \sum_{x \in X} p(x=1|y=1) \log p(x=1|y=1) \right) + \dots$$

$$= \frac{1}{4} \cdot \left[- \left(\frac{1}{2} \times \log_2 \frac{1}{2} \right) + \dots \right] + \dots$$

结论: $H(Y|X) \neq H(X|Y)$



假设尼西语 $\frac{1}{8} \frac{1}{8} \frac{1}{8} \frac{1}{8} \frac{1}{8} \frac{1}{8}$

$p: \frac{1}{8} \quad t: \frac{1}{8} \quad k: \frac{1}{8} \quad a: \frac{1}{8} \quad i: \frac{1}{8} \quad u: \frac{1}{8}$

$$H(p) = 2.5 \text{ bits}$$

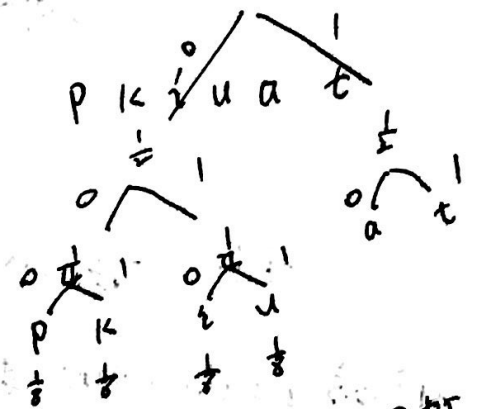
100 00 101 01 110 111

元 $V = \{a, i, u\}$

辅 $C = \{p, t, k\}$

元-辅

假设所有单词都由 CV (consonant-vowel) 音节构成，~~每个~~联合概率与边缘概率分布如下：



带权路径长度最小的二叉树

$C \backslash V$	p	t	k	$P(C, V)$
a	$\frac{1}{16}$	$\frac{3}{8}$	$\frac{1}{16}$	$\frac{1}{4}$
i	$\frac{1}{16}$	$\frac{3}{16}$	0	$\frac{1}{8}$
u	0	$\frac{3}{16}$	$\frac{1}{16}$	$\frac{1}{8}$
$P(C, \cdot)$	$\frac{1}{8}$	$\frac{3}{4}$	$\frac{1}{8}$	

求联合熵：

$p: \frac{1}{16} \quad t: \frac{3}{8} \quad k: \frac{1}{16} \quad a: \frac{1}{16} \quad i: \frac{1}{8} \quad u: \frac{1}{8}$

$$H(C, V) = H(C) + H(V|C)$$

$$H(C) = 1.061 \text{ bits}$$

$$H(V|C) =$$

熵率：对于一个长度为 n 的信息，每一个子符或子的熵为

$$H_{\text{rate}} = \frac{1}{n} H(X_1, \dots, X_n) = -\frac{1}{n} \sum_{i=1}^n p(x_i) \log p(x_i)$$

X_n 表示 $\dots, (X_1, \dots, X_n)$ 有时将 X_n 写作 X_1^n

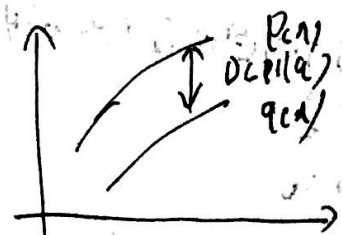
相对熵：或称 KL 散度 Kullback-Leibler divergence

relative entropy

$$D(p||q) = \sum p(x) \log \frac{p(x)}{q(x)}$$

两个概率分布 $p(x)$ 和 $q(x)$ 的相对熵定义为

$$\text{其中 } 0 \log 0/q = 0 \quad p \log q/q = \infty$$



交叉熵:

$X \sim p(x)$ $q(x)$ 为用于近似 $p(x)$ 的概率分布, 那么, 随机变量 X 的模型 q 之间的交叉熵定义为:

$$H(X, q) = H(X) + D(p||q) \\ = -\sum_x p(x) \log q(x)$$

推导: $-\sum_x p(x) \log_2 p(x) + \sum_x p(x) \log_2 \frac{p(x)}{q(x)} = -\sum_x p(x) \log q(x)$

对于语言 $L = (x_i) \sim p(x_i)$ 与其模型 q 的交叉熵定义为:

$$H(L, q) = -\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n p(x_i) \log q(x_i)$$

$x_i^1 = x_1, \dots, x_n$ 为语言 L 的固定列(样本)

$p(x_i^1)$ 为 x_i^1 的概率

$q(x_i^1)$ 为估计值

定理: 假定语言 L 是独立同分布随机过程, x_i^1 为 L 的样本, 那么, 有

$$H(L, q) = -\lim_{n \rightarrow \infty} \frac{1}{n} \log q(x_i^1)$$

困惑度:

在设计语言模型时, 常用困惑度来代替交叉熵衡量拟合度

语言 L 的样本: $u^1 = u_1, \dots, u_n$

困惑度 PP_q 定义为 $PP_q = 2^{H(L, q)} \approx 2^{-\frac{1}{n} \log q(u^1)}$

$$= \log_2 [q(u^1)]^{-\frac{1}{n}}$$

互信息: $(X, Y) \sim p(x, y)$ X, Y 之间的互信息 $I(X; Y)$ 定义为:

$$I(X; Y) = H(X) - H(X|Y)$$

$$= -\sum_{x,y} p(x,y) \log p(x,y) + \sum_{x,y} p(x,y) \log p(x|y)$$

$$= \sum_{x,y} \sum_{y'} p(x,y) (\log p(x|y) - \log p(x|y'))$$

$$= \sum_{x,y} \sum_{y'} p(x,y) (\log \frac{p(x,y)}{p(x,y')})$$

$$= \sum_{x,y} \sum_{y'} p(x,y) \log \frac{p(x,y)}{p(x) p(y)}$$

表示知道 Y 的值以后 X 的不确定性的减少量

定性的减少量

$$H(X|X) = 0$$

$$H(X) = H(X,Y) - H(X|Y)$$

$$= I(X; Y)$$



可利用互信息值估计两个汉字结合的程度。

$$I(x; y) = \log_2 \frac{p(x, y)}{p(x)p(y)} = \log_2 \frac{p(y|x)}{p(y)}$$

互信息值越大，表示两个汉字之间的结合越紧密

双字耦合度

不考虑两个字不连续出现的情况

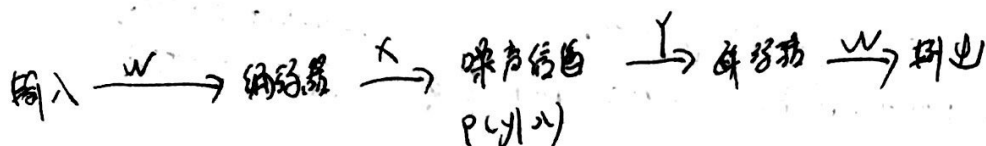
$$\text{Couple}(c_i, c_{i+1}) = \frac{N(c_i c_{i+1})}{N(c_i c_{i+1}) + N(c_i \dots c_{i+1} \dots)}$$

统计样本中 (c_i, c_{i+1}) 连续出现在一个词中的次数 / 连续出现的总次数

噪声信道模型

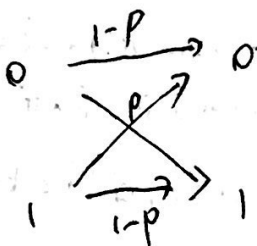
基本假设：一个信道的输出以一定的概率依赖于输入

目标：优化噪声信道中信号传输的吞吐量与准确性



二进制对称信道：

p ：差错概率



信道容量：

$$C = \max_{p(x)} I(x; Y)$$

输入与输出的互信息达到最大值

词义消歧问题：word sense disambiguation, wsd)

① 基于上下文分类的消歧方法

1.1 基于贝叶斯分类器

数学描述：假设某个词 w 所处的上下文语境为 C ，如果 w 的多个词义记作 $s_i (i \geq 1, 2)$ ，那么，可通过计算 $\arg \max_{s_i} p(s_i | C)$ 确定 w 的词义

根据贝叶斯公式：
$$p(s_i | C) = \frac{p(C | s_i) p(s_i)}{p(C)} = \frac{p(C | s_i) p(s_i)}{p(C)}$$

考虑词义的不变性，并应用如下独立性假设

$$p(C | s_i) = \prod_{v \in C} p(v | s_i)$$

因此
$$\hat{s}_i = \arg \max_{s_i} [p(s_i) \prod_{v \in C} p(v | s_i)]$$



概率 $P(V_2|s_1)$ 和 $P(s_1)$ 都可用最大似然估计求得

$$P(V_2|s_1) = \frac{n(V_2, s_1)}{n(s_1)} \quad P(s_1) = \frac{n(s_1)}{n(W)}$$

$n(s_1)$: 训练数据中词 w 用于语义 s_1 时的次数

$n(V_2, s_1)$: w 用于语义 s_1 时词 V_2 出现在 w 的上下文的次数

$n(w)$: 多义词 w 在训练数据中出现的总次数

② 基于最大熵的消歧方法

估计在条件 $b \in B$ 下 (已知信息), 发生某个事件 (未知信息) 的概率 $P(a|b)$, 似概率使熵 $H(P(a|b))$ 最大

$$P^*(a, b) = \frac{1}{Z(b)} \exp \left(\sum_{j=1}^J \lambda_j \cdot f_j(a, b) \right)$$

λ_j : 特征权重

其中 $Z(b) = \sum_a \exp \left(\sum_{j=1}^J \lambda_j \cdot f_j(a, b) \right)$

$f_j(a, b)$: 特征函数

$Z(b)$ 为保证对所有 b , 使得 $\sum_a P(a|b) = 1$ 的归一常量

J 代表有 J 个特征

确定特征函数:

对于语义消歧而言, 设 A 为某一多义词所有义项的集合, B 为所有上下文的集合
定义 $\{0, 1\}$ 域上的二值函数 $f(a, b)$ 来表示上下文条件与义项之间的关系

$$f(a, b) = \begin{cases} 1 & \text{若 } (a, b) \in (A, B), \text{ 且满足某种条件} \\ 0 & \end{cases}$$

上下文条件 b 包含: 同形信息, 同义信息, 同形+同义

表示方法: ① 顺序无关

词袋模型 (值都用向量表示)

$\{5, 1, 1, 2, 0, \dots\}$

取二维窗口范围内的同形

② 位置有关

模板表示

$$f(a, b) = \begin{cases} 1 & \text{if } a = s_1 \text{ and } b = \langle (L, V), (V, R) \rangle \\ 0 & \text{otherwise} \end{cases}$$



- 特征选择
- ① 从候选特征集中选择那些在训练数据中出现频数超过一定阈值的
 - ② 利用互信息作为评价指标从候选特征集中选择满足一定互信息要求的特征
 - ③ 利用增量式特征选择方法

选取 n 个, 则 $L = n+1$

入的确定:

GIS 算法:

$(a, b) \in A \times B$ 则特征函数之和为一度量 L , 即:

$$\sum_{j=1}^n f_j(a, b) = L$$

若条件不满足, 则根据训练集求: $L = \max_{a \in A, b \in B} \sum_{j=1}^n f_j(a, b)$

并增加一修正特征 f_L : $f_L(a, b) = L - \sum_{j=1}^n f_j(a, b)$

$f_L(a, b)$ 的取值范围为 $0 \sim L$

2. 初始化

2.1 计算每个 f_j 的期望 $E(f_j)$

3. 迭代

4. 终止

限定迭代次数

对数似然 $L(p)$ 的值极小

$$|L_M - L| < \epsilon$$

$$L(p) = \sum_{a,b} \tilde{p}(a, b) \log p(a|b)$$

修正公式:

$$\lambda^{(n+1)} = \lambda^{(n)} + \frac{1}{\epsilon} \ln \left(\frac{E \tilde{p}(f_j)}{E_{p^{(n)}}(f_j)} \right)$$

