

第2章 统计学习基础

宗成庆

中国科学院自动化研究所

cqzong@nlpr.ia.ac.cn

本章内容



1. 概率论略览
2. 齐夫定律
3. 信息论基础
4. 统计学习概念
5. 应用举例
6. 习题
7. 附录



1. 概率论略览

- 概率 (probability)
- 最大似然估计 (maximum likelihood estimation)
- 条件概率 (conditional probability)
- 全概率公式 (full probability)
- 贝叶斯决策理论 (Bayesian decision theory)
- 贝叶斯法则 (Bayes' theorem)
- 二项式分布 (binomial distribution)
- 期望 (expectation)
- 方差 (variance)
- 随机过程 (stochastic process)

“语言是稳态的可遍历性随机过程”



1. 概率论略览

- 随机过程 (stochastic process)

假设 $\{\xi_t, t \in T\}$ 是一族随机变量, T 是一个实数集合, 如果对于任意实数 $t \in T$, ξ_t 是一个随机变量, 则称 $\{\xi_t, t \in T\}$ 为随机过程。

任意一个样本点 t , ξ_t 都对应一个实数, 而 ξ_t 是随着试验结果不同而变化的一个变量, 则称 ξ_t 为随机变量。

(有时将 ξ_t 写为: $\xi(t)$, 或者用大写 X 等表示。)

1. 概率论略览

随机过程的平稳性(stationary): 在数学中平稳过程又称严格平稳过程 或者 强平稳过程, 是一种特殊的随机过程, 在任一时间段($\Delta t = t_2 - t_1$)或空间里的联合概率分布, 与将这段时间任意平移后的新时间段里的联合概率分布相等, 即:

$$f_X(x_1, \dots, x_n, t_1 + \Delta t, \dots, t_n + \Delta t) = f_X(x_1, \dots, x_n, t_1, \dots, t_n)$$

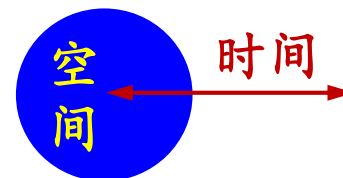


1. 概率论略览

随机过程的遍历性(ergodic): 对于一个平稳的随机变量 X , 如果它的所有样本函数在某一固定时刻的统计特性和单一样本函数在长时间内的统计特性一致, 我们则称 X 为各态遍历, 即随机变量单一样本函数随时间变化的过程可以包括该变量所有样本函数的取值经历。



隐藏含义是: 如果一个随机变量是遍历性的, 那么该随机变量的时间统计特性等于其空间统计信息。





本章内容

1. 概率论略览
- ➡ 2. 齐夫定律
3. 信息论基础
4. 统计学习概念
5. 应用举例
6. 习题
7. 附录



2. 齐夫定律

齐夫定律 (Zipf's law) 是1949年哈佛大学语言学家乔治·金斯利·齐夫 (George Kingsley Zipf) 发表的**实验定律**。

在自然语言的大规模文本数据上统计，一个单词出现的频率与它在频率表中的名次（序号）成反比。

例如，在100万单词的Brown语料库中，the出现的频率最高，出现了69971次，占比大约7%，名列第一。单词of 出现了36411次，占比约3.5%，名列第二。and 出现了28852次，约占2.9%，名列第三。粗略的规律是：按词频顺序，第 r (r 为自然数)个词汇出现的频率为第1个词出现频率的 $1/r$ 倍。如英语中的 the, of, and 这三个词的词频和名次约为：6: 3: 2。

在Brown语料库中，前135个词在整个语料库中约占一半。

2. 齐夫定律

有人在2000万字（14829个词汇）的汉语语料库上进行了统计，频次最高的前10个字（单字词）为：

名次	字（词）	频率	占比(%)	累计占比(%)
1	的	744863	7.7946	7.7946
2	了	130191	1.3624	9.1570
3	在	118823	1.2434	10.4004
4	是	118527	1.2403	11.6407
5	和	83958	0.8786	12.5193
6	一	81119	0.8489	13.3682
7	这	65146	0.6817	14.0499
8	有	53556	0.5604	14.6103
9	他	52912	0.5537	15.1640
10	我	52728	0.5518	15.7158

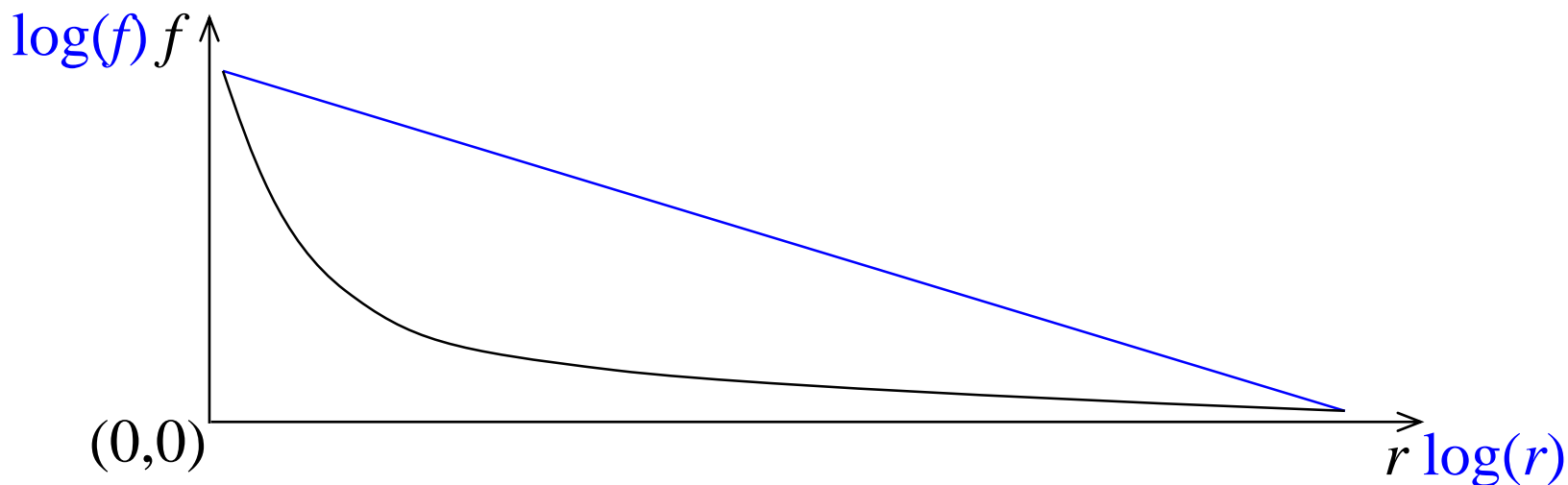
摘自网站：<https://zhuanlan.zhihu.com/p/44646312>

说明：具体结果取决于选取的统计语料类型和规模。

2. 齐夫定律

一般而言，假设词汇 w 出现的频率为 f ，在排序列表中处于 r 号位置上，那么， f 和 r 的乘积趋于一个常数，即 $f \times r = C$ ， C 为常数。

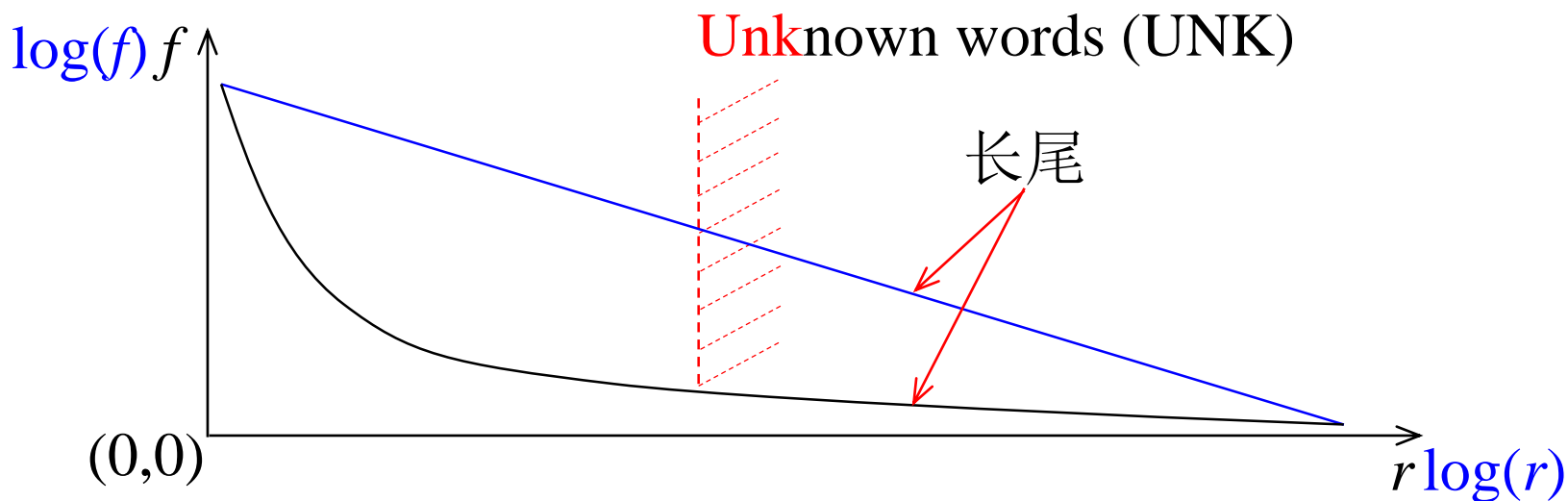
如果横轴表示 r 的对数值 $\log(r)$ ，纵轴表示 f 的对数值 $\log(f)$ ， $\log(r)$ 与 $\log(f)$ 的取值关系近似为一条直线，如下图中的蓝线。



2. 齐夫定律

从统计结果看，少数高频词占了整个语料规模的大部分比例，而大部分词汇属于低频词。这种现象通常被称为长尾效应(long tail effect)，相应的词汇称为长尾词(long tail word)。

在自然语言处理中，考虑到计算量、存储量和运算效率等因素，通常只考虑出现词频高于某个阈值的词，而将低于该阈值的长尾词当作生词(unknown word)处理。





本章内容

1. 概率论略览
2. 齐夫定律
- ➡ 3. 信息论基础
4. 统计学习概念
5. 应用举例
6. 习题
7. 附录

3. 信息论基础

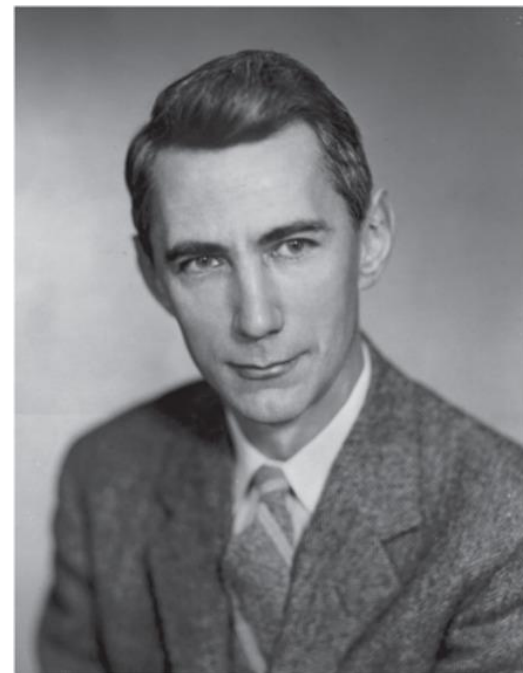
◆熵(entropy)

克劳德·艾尔伍德·香农(Claude Elwood Shannon)于1940年获得MIT数学博士学位和电子工程硕士学位后, 1941年加入贝尔实验室数学部, 工作到1972年。1956年他成为MIT客座教授, 并于1958年成为终生教授, 1978年成为名誉教授。

1948年6月和10月由贝尔实验室出版的《贝尔系统技术》杂志连载了香农博士的文章《通讯的数学原理》, 该文奠定了香农信息论的基础。

熵是信息论中重要的基本概念。

克劳德·艾尔伍德·香农
(1916.4.30~2001.2.24)





3. 信息论基础

如果 X 是一个离散型随机变量，其概率分布为：

$p(x) = P(X = x)$, $x \in X$ 。 X 的熵 $H(X)$ 为：

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x) \quad (1)$$

其中，约定 $0 \log 0 = 0$ 。

$H(X)$ 也可以写为 $H(p)$ 。通常熵的单位为二进制位比特 (bit)。

熵表示随机变量 X 每个具体取值(如信源每发射一个信号)所提供的平均信息量。熵也可以被视为描述一个随机变量的不确定性的数量。一个随机变量的熵越大，它的不确定性越大。那么，正确估计其值的可能性就越小。越不确定的随机变量越需要大的信息量用以确定其值。



3. 信息论基础

例2-1: 计算下列两种情况下英文(26个字母和1个空格, 共27个字符)信息源的熵: (1)假设27个字符等概率出现; (2)假设英文字母的概率分布如下:

字母	空格	E	T	O	A	N	I	R	S
概率	0.1956	0.105	0.072	0.0654	0.063	0.059	0.055	0.054	0.052

字母	H	D	L	C	F	U	M	P	Y
概率	0.047	0.035	0.029	0.023	0.0225	0.0225	0.021	0.0175	0.012

字母	W	G	B	V	K	X	J	Q	Z
概率	0.012	0.011	0.0105	0.008	0.003	0.002	0.001	0.001	0.001



3. 信息论基础

解：(1) 等概率出现情况：

$$\begin{aligned} H(X) &= - \sum_{x \in X} p(x) \log_2 p(x) \\ &= 27 \times \left\{ -\frac{1}{27} \log_2 \frac{1}{27} \right\} = \log_2 27 = 4.75 \quad (\text{bits/letter}) \end{aligned}$$

(2) 按表中概率计算：

$$H(X) = - \sum_{i=1}^{27} p(x_i) \log_2 p(x_i) = 4.02 \quad (\text{bits/letter})$$

说明： 考虑了英文字母和空格实际出现的概率后，英文信源的平均不确定性，比把字母和空格看作等概率出现时英文信源的平均不确定性要小。

3. 信息论基础

法语、意大利语、西班牙语、英语、俄语字母及汉字的熵：

语言	熵 (bits)
法语	3.98
意大利语	4.00
西班牙语	4.01
英语	4.03 (英语单词的熵约为10。)
俄语	4.35
汉字	9.71 (汉语词的熵约为11.46。)

[冯志伟, 1989]

说明：西文字母的熵与汉语字的熵之间没有对比意义。

3. 信息论基础

北京、香港、台北三地汉语词的熵[Tsou,2003]

北京5年		台北5年		香港5年		京、港、台5年	
A1	A2	B1	B2	C1	C2	D1	D2
11.45	11.11	11.69	11.36	11.96	11.64	11.96	11.60

其中，A1, B1, C1 分别是从小港城市大学建立的语料库(LIVAC)中北京、台北、香港三地5年各约1000万字文本中所提取的数据；A2, B2, C2 为三地文本剔除专用名词之后的数据。D1, D2分别为三地文本合并之后剔除专用名词前后的数据。

专用名词又称命名实体(named entity)，主要指：人名、地名、组织机构名、时间、数字及货币等。



3. 信息论基础

◆ 联合熵(joint entropy)

如果 X, Y 是一对离散型随机变量 $X, Y \sim p(x, y)$, X, Y 的联合熵 $H(X, Y)$ 为:

$$H(X, Y) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(x, y) \quad (2)$$

联合熵实际上就是描述一对随机变量平均所需要的信息量。



3. 信息论基础

◆ 条件熵(conditional entropy)

给定随机变量 X 的情况下，随机变量 Y 的条件熵定义为：

$$\begin{aligned} H(Y | X) &= \sum_{x \in X} p(x) H(Y | X = x) \\ &= \sum_{x \in X} p(x) \left[- \sum_{y \in Y} p(y | x) \log_2 p(y | x) \right] \\ &= - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(y | x) \end{aligned} \quad (3)$$

$$p(x) \cdot p(y|x) = p(x, y)$$

3. 信息论基础

将 (2) 式: $H(X, Y) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(x, y)$ 中的

$\log_2 p(x, y)$ 根据概率公式展开:

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

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

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

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

$$= H(X) + H(Y | X) \quad (4) \quad (\text{连锁规则})$$

3. 信息论基础

例2-2: 假设 (X, Y) 服从如下联合概率分布:

$Y \backslash X$	1	2	3	4
1	$1/8$	$1/16$	$1/32$	$1/32$
2	$1/16$	$1/8$	$1/32$	$1/32$
3	$1/16$	$1/16$	$1/16$	$1/16$
4	$1/4$	0	0	0

请计算 $H(X)$ 、 $H(Y)$ 、 $H(X|Y)$ 、 $H(Y|X)$ 和 $H(X, Y)$ 各是多少?

3. 信息论基础

例2-2: 假设 (X, Y) 服从如下联合概率分布: $\rightarrow (X, Y)$

$Y \backslash X$	1	2	3	4
1	1/8	1/16	1/32	1/32
2	1/16	1/8	1/32	1/32
3	1/16	1/16	1/16	1/16
4	1/4	0	0	0
$p(X, \bullet)$	1/2	1/4	1/8	1/8

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

$$= - \left(\frac{1}{2} \times \log_2 \left(\frac{1}{2} \right) + \frac{1}{4} \times \log_2 \left(\frac{1}{4} \right) + \frac{1}{8} \times \log_2 \left(\frac{1}{8} \right) + \frac{1}{8} \times \log_2 \left(\frac{1}{8} \right) \right)$$

$$= \frac{7}{4}$$

3. 信息论基础

类似地，可以计算 $H(Y)$ 。

$Y \backslash X$	1	2	3	4	$p(\bullet, Y)$
1	1/8	1/16	1/32	1/32	1/4
2	1/16	1/8	1/32	1/32	1/4
3	1/16	1/16	1/16	1/16	1/4
4	1/4	0	0	0	1/4

$$H(Y) = -\sum_{y \in Y} p(y) \log_2 p(y) = 2 \text{ (bits)}$$

3. 信息论基础

$Y \backslash X$	1	2	3	4	$p(\bullet, Y)$
1	1/8	1/16	1/32	1/32	1/4
2	1/16	1/8	1/32	1/32	1/4
3	1/16	1/16	1/16	1/16	1/4
4	1/4	0	0	0	1/4
$p(X, \bullet)$	1/2	1/4	1/8	1/8	

$$p(x_1 | y_1) = \frac{p(x_1, y_1)}{p(y_1)} = \frac{1}{8} \times \frac{4}{1} = \frac{1}{2}$$

$$p(x_2 | y_1) = \frac{p(x_2, y_1)}{p(y_1)} = \frac{1}{16} \times \frac{4}{1} = \frac{1}{4}$$

$$p(x_3 | y_1) = \frac{p(x_3, y_1)}{p(y_1)} = \frac{1}{32} \times \frac{4}{1} = \frac{1}{8}$$

$$p(x_4 | y_1) = \frac{p(x_4, y_1)}{p(y_1)} = \frac{1}{32} \times \frac{4}{1} = \frac{1}{8}$$

.....

3. 信息论基础

$$\begin{aligned}
 H(X | Y) &= \sum_{i=1}^4 p(y=i) H(X | Y=i) \\
 &= \frac{1}{4} H\left(\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8}\right) + \frac{1}{4} H\left(\frac{1}{4}, \frac{1}{2}, \frac{1}{8}, \frac{1}{8}\right) \\
 &\quad + \frac{1}{4} H\left(\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}\right) + \frac{1}{4} H(1,0,0,0) \\
 &= \frac{1}{4} \times \frac{7}{4} + \frac{1}{4} \times \frac{7}{4} + \frac{1}{4} \times 2 + \frac{1}{4} \times 0 = \frac{11}{8} \text{ (bits)}
 \end{aligned}$$

$-\sum_{i=1}^4 p(x_i | y_1) \log p(x_i | y_1)$
 $-\sum_{i=1}^4 p(x_i | y_2) \log p(x_i | y_2)$

类似地， $H(Y|X)=13/8$ (bits), $H(X, Y)=H(X)+H(Y|X)=27/8$ (bits) 。

$H(Y|X) \neq H(X|Y)$? **No!**



3. 信息论基础

例2-3: 简单的波利尼西亚语(Polynesian)是一些随机的字符序列, 其中部分字符出现的概率为:

p: 1/8, t: 1/4, k: 1/8, a: 1/4, i: 1/8, u: 1/8

那么, 每个字符的熵为:

$$\begin{aligned} H(P) &= - \sum_{i \in \{p, t, k, a, i, u\}} P(i) \log P(i) \\ &= - \left[4 \times \frac{1}{8} \log \frac{1}{8} + 2 \times \frac{1}{4} \log \frac{1}{4} \right] = 2 \frac{1}{2} \text{ (bits)} \end{aligned}$$



3. 信息论基础

这个结果表明，我们可以设计一种编码，传输一个字符平均只需要2.5个比特：

p	t	k	a	i	u
100	00	101	01	110	111

这种语言的字符分布并不是随机变量，但是，我们可以近似地将其看作随机变量。如果将字符按元音和辅音分成两类，元音随机变量 $V=\{a, i, u\}$ ，辅音随机变量 $C=\{p, t, k\}$ 。

3. 信息论基础

假定所有的单词都由CV(consonant-vowel)音节序列组成，其联合概率分布 $P(C, V)$ 、边缘分布 $P(C, \bullet)$ 和 $P(\bullet, V)$ 如下表所示：

$V \backslash C$	p	t	k	$P(\bullet, V)$
a	1/16	3/8	1/16	1/2
i	1/16	3/16	0	1/4
u	0	3/16	1/16	1/4
$P(C, \bullet)$	1/8	3/4	1/8	



3. 信息论基础

注意，这里的边缘概率是基于每个音节的，其值是基于每个字符的概率的两倍，因此，每个字符的概率值应该为相应边缘概率的 $1/2$ ，即：

$$p: 1/16 \quad t: 3/8 \quad k: 1/16 \quad a: 1/4 \quad i: 1/8 \quad u: 1/8$$

现在我们来求联合熵为多少？



3. 信息论基础

利用连锁规则求联合熵：

$$\begin{aligned} H(C) &= - \sum_{c=p,t,k} p(c) \log p(c) = -2 \times \frac{1}{8} \times \log \frac{1}{8} - \frac{3}{4} \times \log \frac{3}{4} \\ &= \frac{9}{4} - \frac{3}{4} \log 3 \approx 1.061 \text{ (bits)} \end{aligned}$$

$$\begin{aligned} H(V | C) &= \sum_{c=p,t,k} p(C=c) H(V | C=c) \\ &= \frac{1}{8} H\left(\frac{1}{2}, \frac{1}{2}, 0\right) + \frac{3}{4} H\left(\frac{1}{2}, \frac{1}{4}, \frac{1}{4}\right) + \frac{1}{8} H\left(\frac{1}{2}, 0, \frac{1}{2}\right) = \frac{11}{8} = 1.375 \text{ (bits)} \end{aligned}$$

因此，

$$H(C, V) = H(C) + H(V | C) = \frac{9}{4} - \frac{3}{4} \log 3 + \frac{11}{8} \approx 2.44 \text{ (bits)}$$



3. 信息论基础

- ◆ 相对熵(relative entropy, 或称 Kullback-Leibler divergence, K-L 距离, 或K-L散度)

两个概率分布 $p(x)$ 和 $q(x)$ 的相对熵定义为:

$$D(p \parallel q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)} \quad (5)$$

该定义中约定 $0 \log (0/q) = 0$, $p \log (p/0) = \infty$ 。

3. 信息论基础

相对熵常被用以衡量两个随机分布的差距。当两个随机分布相同时,其相对熵为0。两个随机分布的差别增加时,其相对熵也增加。

$$D(p \parallel q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)}$$

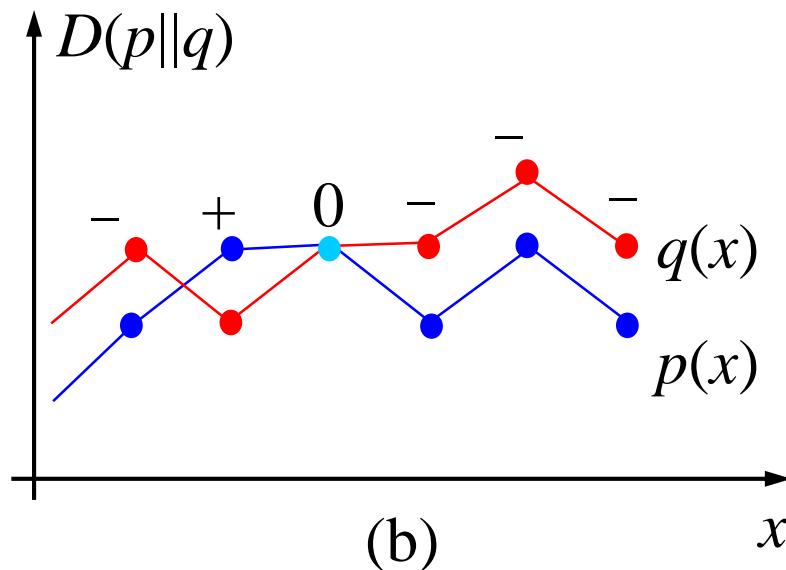
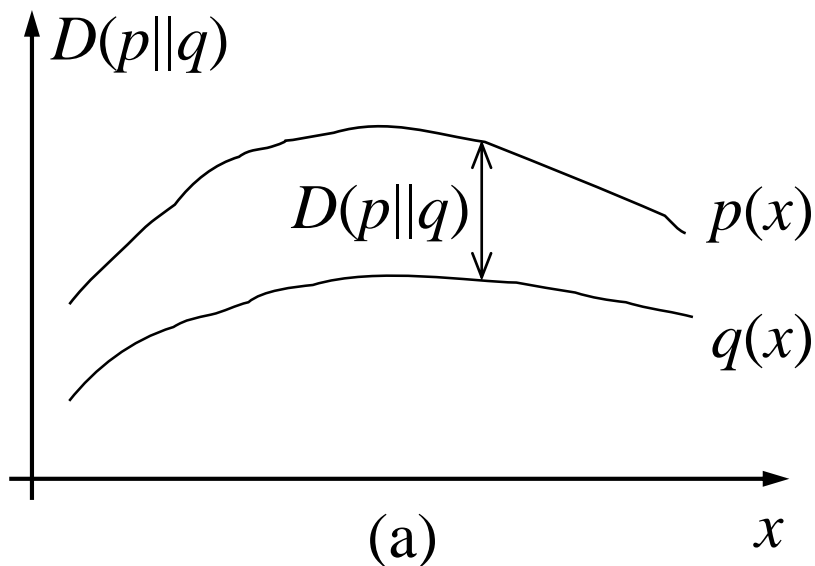


图2-1. 相对熵示意图



3. 信息论基础

◆交叉熵(cross entropy)

如果一个随机变量 $X \sim p(x)$, 理论模型 $q(x)$ 为用于近似 $p(x)$ 的概率分布, 那么, 统计分布 p 和模型 q 之间的交叉熵定义为:

$$\begin{aligned} H(X, q) &= H(X) + D(p \parallel q) \\ &= -\sum_{x \in X} p(x) \log p(x) + \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)} \\ &= -\sum_{x \in X} p(x) \log q(x) \end{aligned} \quad (6)$$

交叉熵衡量的也是两个模型分布之间的差异。

3. 信息论基础

对于语言 $L = (X) \sim p(x)$ 与其理论模型 q 的交叉熵定义为：

$$H(L, q) = -\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{x_1^n} p(x_1^n) \log q(x_1^n) \quad (7)$$

其中， $x_1^n = x_1 \dots x_n$ 为语言 L 的词序列（已知样本）；

$p(x_1^n)$ 为 x_1^n 的概率（在已知样本上的统计分布）；

$q(x_1^n)$ 为模型 q 对 x_1^n 的概率估计值（理论模型）。



3. 信息论基础

定理：假定语言 L 是稳态(stationary)的可遍历性(ergodic)随机过程, x_1^n 为 L 的样本, 那么, 有:

$$H(L, q) = -\lim_{n \rightarrow \infty} \frac{1}{n} \log q(x_1^n) \quad (8)$$

证明见本章讲义**附录1**。

由此, 我们可以根据模型 q 和一个含有大量数据的 L 的样本来计算交叉熵。在设计模型 q 时, 我们的目的是使交叉熵最小, 从而使模型最接近真实的概率分布 $p(x)$ 。



3. 信息论基础

问题： 相对熵和交叉熵的作用有什么区别呢？

根据定义，相对熵：
$$D(p \parallel q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)}$$

交叉熵：
$$H(p, q) = - \sum_{x \in X} p(x) \log q(x) = \sum_{x \in X} p(x) \log \frac{1}{q(x)}$$

$H(p, q) = H(p) + D(p \parallel q)$ ， 即： **交叉熵 = 熵 + 相对熵**。



3. 信息论基础

说明：在机器学习中经常用 $p(x)$ 表示真实数据的概率分布，由于真实数据的概率分布往往无法获得，所以一般通过大量的训练数据来近似。假设我们通过某个模型得到了训练数据的概率分布 $q(x)$ ，**由于真实数据的概率分布 $p(x)$ 往往是不变的，因此最小化交叉熵 $H(p, q)$ 等效于最小化相对熵 $D(p||q)$ 。**习惯上机器学习算法中通常采用交叉熵计算损失函数。

例如，在某机器学习任务中定义损失函数为**交叉熵**： $\text{Loss}=H(p, q)$ ，假设我们训练到得到一个非常好的模型，即 $p(x)\approx q(x)$ ，此时Loss不会降低为0，而是一个很小的值，如 $\text{Loss}=2$ ，它表示真实数据自身的熵为 $H(p)=2$ 。如果选择**相对熵**作为损失函数，即 $\text{Loss}=D(p||q)$ ，同样假设我们训练得到一个非常好的模型，即 $p(x)\approx q(x)$ ，此时， $\text{Loss}=0$ ，意味着两个概率分布几乎一样。

实际上，上述两种方法所得到的Loss仅仅是数值上的区别，训练得到的模型是完全一样的，即两个概念的作用一样。



3. 信息论基础

◆ 困惑度(perplexity)

在设计语言模型时，我们通常用困惑度来代替交叉熵衡量语言模型的好坏。给定语言 L 的样本 $x_1^n = x_1 \dots x_n$ ， L 的困惑度 PP_q 定义为：

$$PP_q = 2^{H(L,q)} \approx 2^{-\frac{1}{n} \log q(l_1^n)} = [q(x_1^n)]^{-\frac{1}{n}} \quad (9)$$

语言模型设计的任务就是寻找困惑度最小的模型，使其最接近真实的语言分布。



3. 信息论基础

◆ 互信息(mutual information)

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

$$I(X; Y) = H(X) - H(X | Y) \quad (10)$$

根据 $H(X)$ 和 $H(X|Y)$ 的定义:

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

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



3. 信息论基础

$$\begin{aligned} I(X; Y) &= H(X) - H(X | Y) \\ &= -\sum_{x \in X} p(x) \log_2 p(x) + \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(x | y) \\ &= \sum_{x \in X} \sum_{y \in Y} p(x, y) (\log_2 p(x | y) - \log_2 p(x)) \\ &= \sum_{x \in X} \sum_{y \in Y} p(x, y) \left(\log_2 \frac{p(x | y)}{p(x)} \right) \\ I(X; Y) &= \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 \frac{p(x, y)}{p(x)p(y)} \end{aligned} \quad (11)$$

互信息 $I(X; Y)$ 是在知道了 Y 的值以后 X 的不确定性的减少量，即 Y 的值透露了多少关于 X 的信息量。

3. 信息论基础

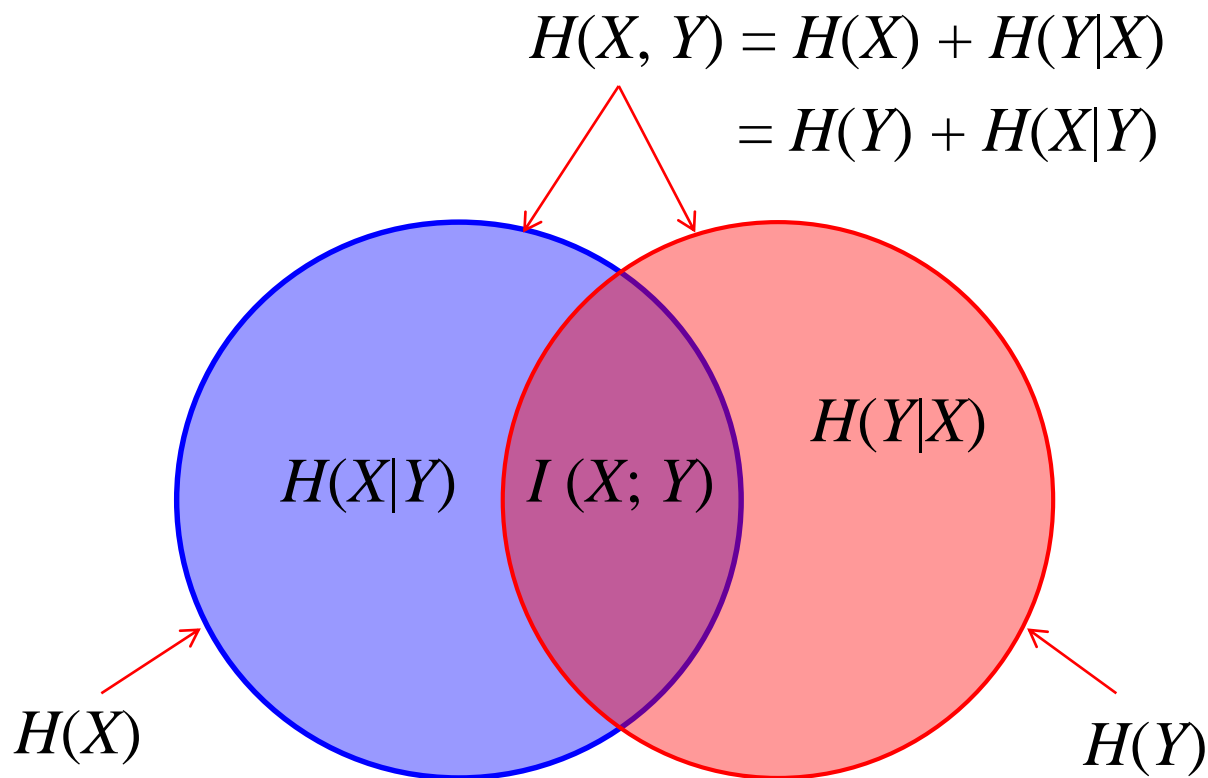


图 2-2. 互信息、条件熵与联合熵



3. 信息论基础

由于 $H(X|X) = 0$, 所以,

$$H(X) = H(X) - H(X|X) = I(X; X) \quad (12)$$

这就是为什么熵又称为自信息(self-information)。这也意味着两个完全相互依赖的变量之间的互信息并不是一个常量，而是取决于它们的熵。



3. 信息论基础

例如：汉语分词问题

为人 民服 务。

利用互信息值估计两个汉字结合的程度：

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

理论上，互信息值越大，表示两个汉字之间的结合越紧密，越可能成词。反之，断开的可能性越大。当两个汉字 x 和 y 关联度较强时，其互信息值 $I(x, y) > 0$ ； x 与 y 关系弱时， $I(x, y) \approx 0$ ；而当 $I(x, y) < 0$ 时， x 与 y 称为“互补分布”。



3. 信息论基础

说明：两个单个离散事件 (x_i, y_j) 之间的互信息 $I(x_i, y_j)$ 通常称为点式互信息(point-wise mutual information)，点式互信息可能为负值。两个随机变量 (X, Y) 之间的互信息 $I(X, Y)$ 称为平均互信息，**平均互信息不可能为负值。**

关于两个随机变量之间平均互信息为非负值的证明见本课件**附录2**。

3. 信息论基础

在汉语分词研究的实践证明中，用互信息作为衡量两个汉字之间是否作为切分边界点的判别依据，效果并不好。

例如：**教务**以连续字符串形式在统计样本中共出现了16次，而**教**字出现了14945次，**务**字出现了6015次。**教**和**务**之间的互信息只有 -0.5119 。如果用互信息来判断的话，这两个字应该被切开。但实际上，**教**和**务**这两个字在文本集中出现的16次全部都是**教务**、**教务长**、**教务处**等词汇。也就是说，这两个字一旦连续同现，一定成词。因此，用两个邻近的字在训练样本中构成词的比率作为切分依据，比互信息效果更好。



本章内容

1. 概率论略览
2. 齐夫定律
3. 信息论基础
- ➡ 4. 统计学习概念
5. 应用举例
6. 习题
7. 附录



4. 统计学习概念

◆ 统计学习方法 (1980s~)

- ✧ 语音识别，如 SPHINX 语音识别系统(CMU)
- ✧ 训练控制系统用于驾驶车辆，如 ALVINN 系统
- ✧ 在各种大规模数据库中发现隐藏的一般规律，如美国国家航空和航天局(NASA)使用决策树进行天体分类
- ✧ 世界级水平的西洋双陆棋博弈
- 学习过程：通过经验提高模型的性能
- 理论基础：统计学、信息论、计算复杂性理论、人工智能、神经生物学



4. 统计学习概念

◆ 统计学习方法

- ✧ 数据驱动
- ✧ 对数据进行预测与分析
- ✧ 以方法为中心，构建模型

◆ 统计学习类型

- 监督学习(supervised learning)
- 非监督学习(unsupervised learning)
- 半监督学习(semi-supervised learning)
- 强化学习(reinforcement learning)



4. 统计学习概念

● 监督学习(supervised learning)

- 给定有限的、人工标注好的大量数据，假设这些数据是独立同分布产生的(训练集, training data)
- 假设要学习的模型属于某个函数的集合，即假设空间(hypothesis space)
- 应用某（些）个评价准则(evaluation criterion)，从假设空间中选取最优的模型，使其对已知的训练数据和未知的测试数据(test data)在给定的评价准则下有最优的预测



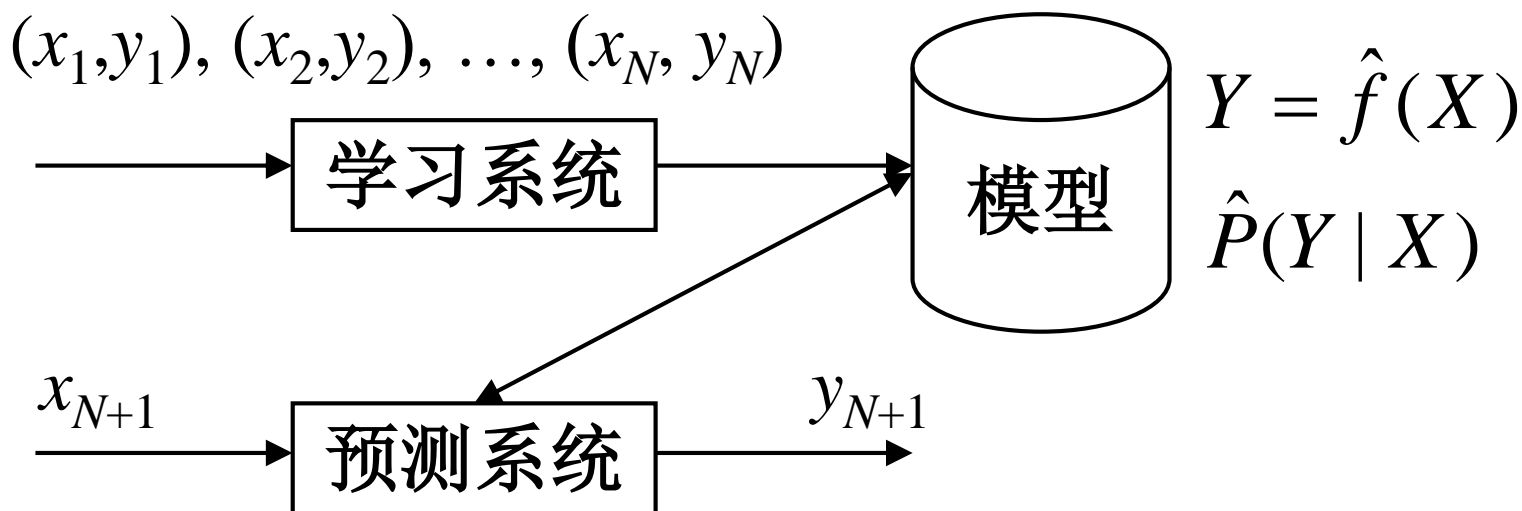
4. 统计学习概念

一般步骤:

- ① 获得一个有限的训练数据集合
- ② 确定包含所有可能的模型的假设空间，即学习模型的集合
- ③ 确定模型选择的准则，即学习的策略
- ④ 通过学习方法选择最优模型
- ⑤ 利用学习到的最优模型对新数据进行预测或分析

4. 统计学习概念

问题的形式化:



给定一个训练数据集: $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$, 其中, (x_i, y_i) , $i=1, 2, \dots, N$, 称为样本。 x_i 是输入的观测值, 也称输入或实例;
 y_i 是输出的观测值, 也称输出。



4. 统计学习概念

◆ 模型的类别

- 生成式方法(generative model)
- 区分式方法/判别式方法(discriminative model)



4. 统计学习概念

● 生成式方法 (generative model)

假设 O 是观察值， Q 是模型。生成式模型首先在训练数据集上建立概率模型 Q ，然后，利用模型 Q 对输入实例进行预测，即通过概率 $P(O'|Q)$ 进行预测。它建立在统计学和Bayes理论的基础之上，要求标注的训练样本无穷多或者足够多。

代表性模型：

- n 元文法模型(n -gram)
- 隐马尔可夫模型(Hidden Markov Model, HMM)



4. 统计学习概念

- 区分式/判别式方法 (Discriminative Model)

假设 O 是观察值， Q 是模型，区分式方法对 $P(Q|O)$ 进行建模。其基本思路是：在有限样本条件下建立判别函数，不考虑样本的产生模型，直接研究预测模型，寻找不同类别之间的最优分类面，反映的是不同类别数据之间的差异性。

代表性模型：

➤ 各种分类器



4. 统计学习概念

◆ 相关概念

- 语料(corpus)/ 语料库(corpus base): 语言资源 (数据集)
- 训练集(training data/set): 用于模型参数训练
- 开发集(development data/set): 用于模拟测试, 优化参数
- 测试集(test data/set): 测试模型或模型实际处理的数据
- 过拟合(overfitting): 模型只在训练集上性能优良
- 欠拟合(under-fitting): 模型在训练集上性能远未达到最优
- 鲁棒性(robustness): 测试集的差异对模型性能影响不大



4. 统计学习概念

◆常用的统计模型

● 生成式模型

- 语言模型 (language model)
- 隐马尔可夫模型(hidden Markov model, HMM)

● 区分式模型

- 朴素贝叶斯法(naïve Bayes): 多类分类问题
- 最大熵(maximum entropy): 多类分类问题
- 条件随机场(conditional random field, CRF): 序列标注
- k -近邻法(k -nearest neighbor, k -NN): 多类分类问题
- 决策树(decision tree): 多类分类问题
- 感知机(perceptron): 二类分类
- 支持向量机(support vector machine, SVM): 二类分类



本章内容

1. 概率论略览
2. 齐夫定律
3. 信息论基础
4. 统计学习概念
-  5. 应用举例
6. 习题
7. 附录



5. 应用举例

例2-4: 词汇歧义消解

●问题提出

任何一种语言中，一词多义（歧义）现象是普遍存在的。如何区分不同上下文中的词汇语义，就是词汇歧义消解问题，或称词义消歧(word sense disambiguation, WSD)。

词义消歧是自然语言处理中的基本问题之一。



5. 应用举例

以“打”字为例：

- | | |
|-------------------|-----------------|
| (1) 他会打鼓。 | (13) 给他打个电话吧。 |
| (2) 他把碗打破了。 | (14) 他把款打过去了。 |
| (3) 他在学校打架了。 | (15) 你别打杈。 |
| (4) 他想打官司。 | (16) 你打两瓶水去。 |
| (5) 他用土打了一堵墙。 | (17) 他想打车票回家。 |
| (6) 他会用木头打家具。 | (18) 他以打鱼为生。 |
| (7) 她用面打浆糊贴对联。 | (19) 他放学后去打猪草了。 |
| (8) 他打铺盖卷儿走人了。 | (20) 你打个草稿再写。 |
| (9) 她会用毛线打毛衣。 | (21) 八路军会打游击。 |
| (10) 他用尺子在纸上打了格子。 | (22) 我们一起打扑克吧。 |
| (11) 他打开了井盖子。 | (23) 他给她打了个手势。 |
| (12) 这种人打着灯笼也难找。 | (24) 你别打官腔/马虎眼。 |



5. 应用举例

● 解决思路

每个词表达不同的含意时其上下文（语境）往往不同，也就是说，不同的词义对应不同的上下文，因此，如果能够将多义词的上下文区别开，其词义自然就明确了。**词义消歧变成一个上下文分类任务。**

他/P 打/V 鼓/N 很/D 在行/A 。/PU
-1 0 +1 +2

基本的上下文信息：词、词性、位置。



5. 应用举例

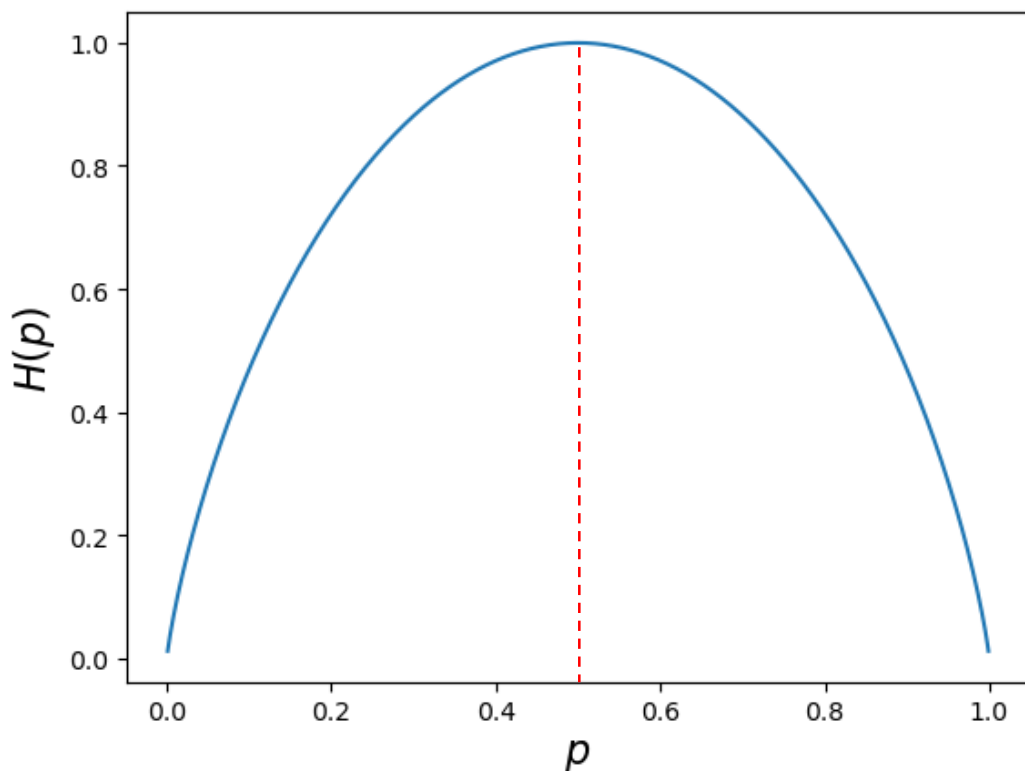
● 基于最大熵的消歧方法

➤ 基本原理:

在只掌握关于未知分布的部分知识的情况下，符合已知知识的概率分布可能有多个，使熵值最大的概率分布能够最真实地反映事件的分布情况，由于熵定义了随机变量的不确定性，当熵最大时，随机变量最不确定。也就是说，在已知部分知识的前提下，关于未知分布最合理的推断应该是符合已知知识最不确定或最大随机的推断。

5. 应用举例

假设变量 x 有两种可能的取值: a 和 b , 那么, $p(a)+p(b)=1$, 将 $p(a)$ 简记为 p , p 的取值与其熵 $H(p)$ 的关系如下图所示:



在没有先验知识的情况下, 均匀分布时熵最大, 反映的理念是对所有可能发生的情况一视同仁, 赋予同等的概率, 这是最公平、最合理的处理策略。



5. 应用举例

如果已知样本向我们提供了关于该随机变量另一些取值的概率分布，如 $p(a') + p(b') = \alpha$ ，那么，我们在拥有这个先验知识的前提下，应该如何处理 $p(a)$ 和 $p(b)$ 的概率分布呢？

$$\begin{cases} p(a') + p(b') = \alpha \\ p(a') + p(b') + p(a) + p(b) = 1 \end{cases}$$

让熵最大是对未知分布推断最合理的准则。换句话说，推断未知分布合理的做法是将已知的先验知识作为约束，让未知分布的熵最大。



5. 应用举例

➤ 模型定义

已知训练数据集 $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$, x_i 表示输入条件, y_i 表示预测值。训练集中每一种情况的概率 $\tilde{p}(x, y)$ 可以通过简单的统计计算出来:

$$\tilde{p}(x, y) \equiv \underline{\text{样本中含有}(x, y)\text{的数量}} / N$$

对于训练集 T 中的所有样本可通过特征函数(feature function)描述 $x \in X$ 和 $y \in Y$ 之间基于某种条件的关系:

$$f(x, y) = \begin{cases} 1, & \text{当} x \text{和} y \text{之间满足某种条件时} \\ 0, & \text{否则} \end{cases}$$

$f(\bullet)$ 实际上是克罗内克(Kronecker)函数。

5. 应用举例

那么， $f(x, y)$ 在训练集上关于经验分布 $\tilde{p}(x, y)$ 的期望值可通过下面的式子计算出来：

$$E_{\tilde{p}}(f) = \sum_{x,y} \tilde{p}(x, y) f(x, y) \quad (17)$$

而理论值为：

$$E_p(f) = \sum_{x,y} \underline{p(x, y)} f(x, y) \quad (18)$$

由于 $p(x, y) = p(x)p(y|x)$ ，而且所建立的理论模型应该符合训练集中的概率分布（近似相等），即 $p(x) = \tilde{p}(x)$ ，因此，理论期望值计算公式(18)式可以写为：

$$E_p(f) = \sum_{x,y} \tilde{p}(x) p(y|x) f(x, y) \quad (19)$$

5. 应用举例

约束: $E_p(f) = E_{\tilde{p}}(f)$, 或者

$$\sum_{x,y} \tilde{p}(x) \cdot p(y|x) \cdot f(x,y) = \sum_{x,y} \tilde{p}(x,y) f(x,y) \quad (20)$$

假设训练集中有 $n \in N$ 个特征函数 $f_j(x, y)$, 它们在建模过程中都对输出结果有影响, 也就是说有 n 个约束条件, 而理论上能够满足这些约束的模型有很多, 它们构成一个集合:

$$P = \{p \mid E_p(f_j) = E_{\tilde{p}}(f_j), j \in \{1, 2, \dots, n\}\}$$

5. 应用举例

在所有满足约束的模型中，使条件熵最大的模型就是最大熵模型，也就是我们要寻找的是最合理的模型。最大熵模型的学习等价于求解条件约束的优化问题：

$$\begin{aligned} p^*(y|x) &= \arg \max_{p \in P} H(p) \\ &= \arg \max_{p \in P} \left\{ - \sum_{x,y} \tilde{p}(x) p(y|x) \log p(y|x) \right\} \end{aligned} \quad (21)$$

$$\text{s. t. } E_p(f_j) = E_{\tilde{p}}(f_j), \quad j = 1, 2, \dots, n$$

$$\sum_y p(y|x) = 1$$

5. 应用举例

经推导(见本章附录3)，有：

$$p^*(y | x) = \frac{1}{Z(x)} \exp\left(\sum_{j=1}^l \lambda_j \cdot f_j(x, y)\right) \quad (22)$$

$$\text{其中, } Z(x) = \sum_y \exp\left(\sum_{j=1}^l \lambda_j \cdot f_j(x, y)\right) \quad (23)$$

$Z(x)$ 为归一化常量，确保对于所有的 x 使 $\sum_y p(y | x) = 1$ 。

5. 应用举例

➤ 确定特征函数

对于词义消歧问题, 设 X 为上下文条件, Y 为某个多义词所有义项的集合。可定义 $\{0, 1\}$ 域上的二值函数 $f(x,y)$ 表示上下文条件与义项之间的关系:

$$f(x, y) = \begin{cases} 1 & \text{若}(x, y) \in (X, Y), \text{ 且满足某种条件} \\ 0 & \text{否则} \end{cases}$$

如: “打” 字的动词义项集合: $Y = \{s_1, s_2, s_3, \dots, s_{24}\}$,
 $X = \{\text{“打” 字出现的上下文}\}$ 。

5. 应用举例

上下文条件(X)表示有:

(1) 词形信息:

他 很 会 与 人 打 交道 。
..... ↑

(2) 词性信息:

他/PN 很/D 会/V 与/C 人/N 打 交道/N 。 /Pu
..... ↑

(3) 词形+词性信息:

他/PN 很/D 会/V 与/C 人/N 打 交道/N 。 /Pu
..... ↑

5. 应用举例

上下文有两种表示方法：

①位置无关：目标词周围的词形、词性或其组合构成的集合，
如取 ± 2 窗口范围内的词形：

{与，人，交道，。}
={交道，与，。，人}

词袋模型
(通常用向量表示)



他/PN 很/D 会/V 与/C 人/N 打 交道/N 。/Pu

②位置有关：词形(± 2)： $\langle \text{与}_{-2}, \text{人}_{-1}, \text{交道}_{+1}, \text{。}_{+2} \rangle$

$\neq \langle \text{与}_{-2}, \text{。}_{+2}, \text{交道}_{+1}, \text{人}_{-1} \rangle$

模板表示



5. 应用举例

他/P 很/D 会/V 与/C 人/N 打/V s_3 交道/N 。/Pu
↑

假设以词形和词性等表示条件，可以构造特征函数：

$$f_1(x, y) = \begin{cases} 1 & \text{If } x = \{\text{与, 人, 交道, 。}\} \text{ and } y = s_3 \\ 0 & \text{Otherwise} \end{cases}$$

$$f_2(x, y) = \begin{cases} 1 & \text{If } x = \{C, N, N, Pu\} \text{ and } y = s_3 \\ 0 & \text{Otherwise} \end{cases}$$

...



5. 应用举例

他/P 很/D 会/V 与/C 人/N 打/ V_{s_3} 交道/N 。/Pu
跟/C 美国/N 人/N 打/ V_{s_3} 交道/N 他/P 有/V 经验/N 。/Pu
杯子/N 今天/N 被/J 他/P 打/ V_{s_2} 破/A 了/Aux 。/Pu
她/P 打/ V_{s_1} 起/X 腰鼓/N 来/Aux 很/D 好看/A 。/Pu
打/ V_{s_1} 锣/N 敲/V 鼓/N 是/V 他/P 的/Aux 强项/N 。/Pu
..... (假设已有大规模标注样本)

取 ± 2 个词为语境窗口，词形、词性为特征，那么，

$f(\text{她}, s_1)=1; f(\text{起}, s_1)=1; f(\text{腰鼓}, s_1)=1; f(\text{锣}, s_1)=1; f(\text{敲}, s_1)=1; f(\text{P}, s_1)=1;$
 $f(\text{X}, s_1)=1; f(\text{N}, s_1)=1; f(\text{V}, s_1)=1; f(\text{他}, s_2)=1; f(\text{被}, s_2)=1; f(\text{破}, s_2)=1;$
 $f(\text{了}, s_2)=1; f(\text{P}, s_2)=1; f(\text{J}, s_2)=1; f(\text{A}, s_2)=1; f(\text{Aux}, s_2)=1; f(\text{人}, s_3)=1;$
 $f(\text{与}, s_3)=1; f(\text{交道}, s_3)=1; f(\text{。}, s_3)=1; f(\text{美国}, s_3)=1; \dots\dots$



5. 应用举例

如果上下文条件由如下三类信息表示：

- (1)特征的类型：词形、词性、词形+词性，3种情况；
- (2)上下文窗口大小：当前词的左右2个词，1种情况；
- (3)是否考虑位置：是或否，2种情况。

上述3种情况组合，可得到如下 n 种特征模板：

$$n=3 \times 1 \times 2=6$$

考虑到词形、词性、位置等又可以组合出很多种可能，可以构造出若干特征函数，因此需要对特征进行筛选。

5. 应用举例

特征选择一般有三种方法：

- ① 从候选特征集中选择那些在训练数据中出现频次超过一定阈值的特征；
- ② 利用互信息作为评价尺度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(Della Pietra *et al.*)从候选特征集中选择特征。(比较复杂)

最终选定 k ($k > 0$) 个特征，对应 k 个特征函数 f 。在以下叙述中不再区分特征和特征函数。



5. 应用举例

➤ 获取 λ 参数

— 利用GIS(generalized iterative scaling) 算法

GIS 迭代过程要求对于训练集中每个实例的任意 $(x, y) \in X \times Y$, k 个特征函数之和为一常量 C , 即:

$$\sum_{j=1}^k f_j(x, y) = C$$

若该条件不满足, 则根据训练集取: $C = \max_{x \in X, y \in Y} \sum_{j=1}^k f_j(x, y)$

并增加一个修正特征 f_l : $f_l(x, y) = C - \sum_{j=1}^k f_j(x, y)$

$l=k+1$ 。 $f_l(x, y)$ 与其它特征函数不一样, 其取值范围为: $0 \sim C$ 。

5. 应用举例

➤ GIS算法描述

(a) 初始化: $\lambda[1..l]=0$;

(b) 计算每个特征函数 f_j 的训练样本期望值 $E_{\tilde{p}}(f_j)$;

(c) 迭代计算特征函数的模型期望值 $E_p(f_j)$:

①利用公式(23)、(22)计算概率 p^* ;

$$Z(x) = \sum_y \exp\left(\sum_{j=1}^l \lambda_j \cdot f_j(x, y)\right)$$
$$p^*(y|x) = \frac{1}{Z(x)} \exp\left(\sum_{j=1}^l \lambda_j \cdot f_j(x, y)\right)$$

λ 修正方法:

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

②若满足终止条件, 则结束迭代;

否则, **修正 λ** , 继续下轮迭代。

(d) 算法结束, 确定 λ , 算出每个 p^* 。



5. 应用举例

迭代终止条件:

- (a) 限定迭代次数;
- (b) 对数似然($L(p)$)的变化小到可以忽略:

$$|L_{i+1} - L_i| < \varepsilon$$

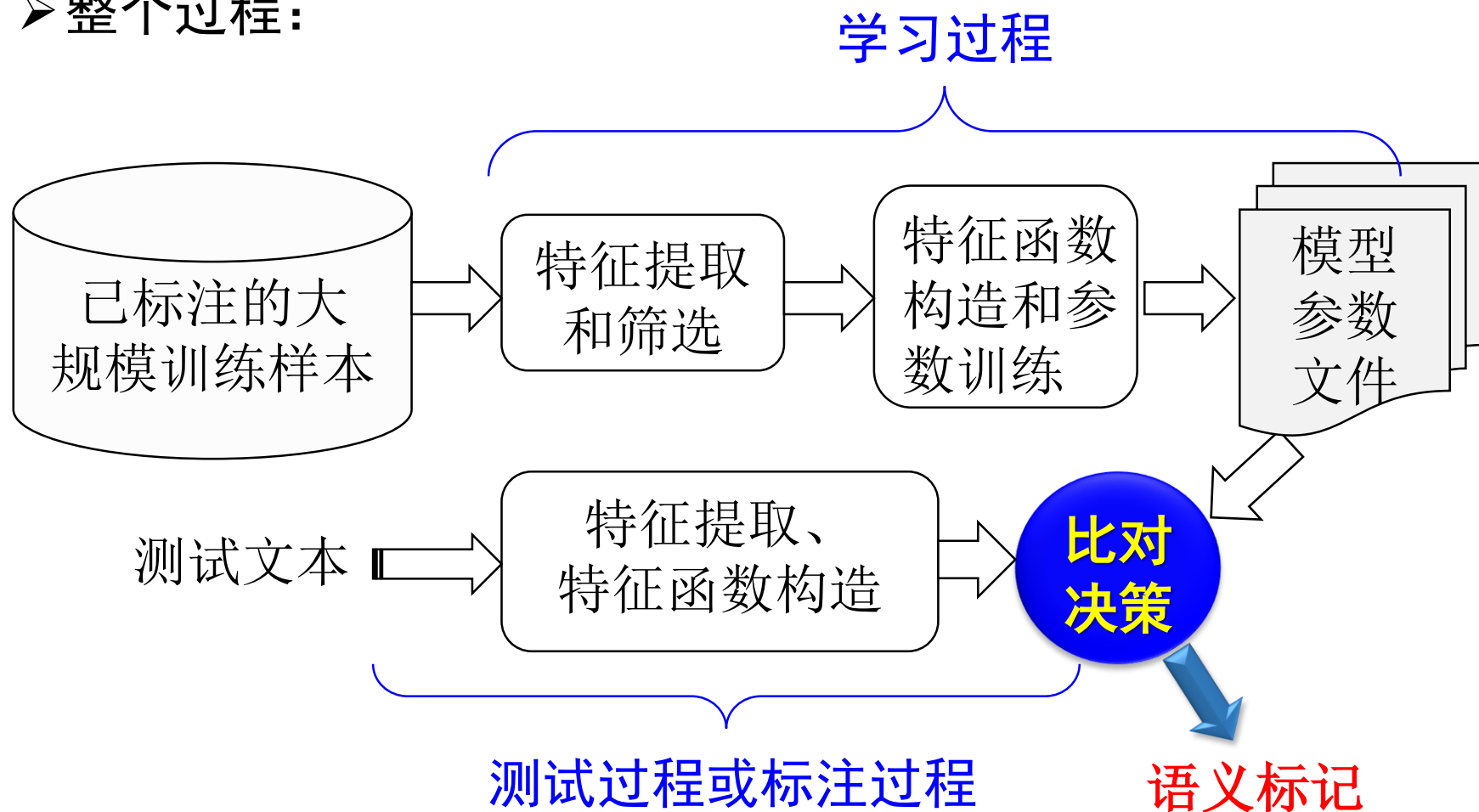
$$L(p) = \sum_{x,y} \tilde{p}(x, y) \log p(y | x)$$

关于GIS算法及改进, 请参阅如下文献:

- [1] J. N. Darroch, D. Ratcliff. Generalized Iterative Scaling for Log-linear Models. *Annals of Math. Statistics*, 1972, 43: 1470-1480
- [2] A. L. Berger. The improved iterative scaling algorithm: A gentle introduction, *Technical report*, Carnegie Mellon University, 1997
- [3] D. Pietra et al. Inducing Features of Random Fields, *IEEE Trans. on PAMI*, 1997, 19(4): 380-393

5. 应用举例

➤ 整个过程：





5. 应用举例

➤ 实验结果：

- 训练数据：用2000年1月1～28日28天的《人民日报》标注文本作为训练数据（全部进行了词义标注）；
- 测试数据：2000年1月29～31日三天的文本作为测试数据，利用所建立的最大熵模型算法对其进行义项标注实验，多义词有4931个；
- 特征模板：特征类型=词形，窗口大小=全句，不考虑位置特征；
- 标注结果：正确率为 94.34%。

张仰森：面向语言资源建设的汉语词义消歧与标注方法研究，北京大学博士后出站报告，2006年12月



5. 应用举例

➤ 关于最大熵方法在NLP中的应用，请参阅：

- [1] Adam L. Berger, Stephen A. Della Pietra, Vincent J. Della Pietra. A Maximum Entry Approach to Natural Language Processing. *Computational Linguistics*, Vol. 22, No. 1, 1996. Pages 39-71 **[ToT Paper Award, ACL-IJCNLP'2021]**
- [2] A. Ratnaparkhi. A Simple Introduction to Maximum Entropy Models for Natural Language Processing. *Technical Report IRCS-97-08*, Dept. of Computer Science, UPenn., 1997
- [3] A. Ratnaparkhi. Maximum Entropy Models for Natural Language Ambiguity Resolution, PhD Dissertation, UPenn., 1998

➤ 最大熵开源工具：

- ✧ OpenNLP: <https://opennlp.apache.org/docs/2.0.0/manual/opennlp.html#opennlp.ml.maxent>
- ✧ Malouf: <http://tadm.sourceforge.net/>



本章小结

◆ 概率论基础

- 基本概念和方法（回顾）

◆ 齐夫定律

◆ 信息论基础

- 熵
- 互信息
- 交叉熵
- 联合熵
- 相对熵
- 困惑度

◆ 统计学习的基本概念

◆ 应用举例

- 词义消歧—— 掌握最大熵分类器的使用方法



本章内容

1. 概率论略览
2. 齐夫定律
3. 信息论基础
4. 统计学习概念
5. 应用举例

 **6. 习题**

7. 附录



5. 习题

1. 分别收集尽量多的英语和汉语文本，编写程序计算这些文本中英语字母和汉字的熵，对比本章课件第17页上表中给出的结果。然后逐步扩大文本规模，如每次增加2M，重新计算文本规模扩大之后的熵，分析多次增加之后熵的变化情况。

要求：

作业

- ① 利用爬虫工具从互联网上收集样本，并对样本进行处理，如清洗乱码等；
- ② 设计算法并编程实现在收集样本上字母/汉字的概率和熵的计算；
- ③ 当改变样本规模时重新计算字母/汉字的概率和熵，并对比计算结果；
- ④ 完成一份技术报告，在报告中写明利用什么爬虫工具从哪些网站上收集的样本，如何进行的样本清洗，清洗后样本的规模，在不同样本规模下计算的结果等。实验分析有较大的伸缩空间。



5. 习题

2. 利用上一个题目收集的文本，分别计算任意两个汉字之间或任意两个英语字母之间的点式互信息。
3. 设 $X \sim p(x)$ ， $q(x)$ 为用于近似 $p(x)$ 的一个概率分布，则 $p(x)$ 与 $q(x)$ 的交叉熵定义为： $H(p, q) = H(p) + D(p \parallel q)$ 。请证明：
$$H(p, q) = -\sum_x p(x) \log q(x)$$
4. 根据本章内容设计并实现一种算法，计算任意两个句子之间的相似性。
5. 举例说明自然语言中两个事件（如“字”“词”或“短语”等）之间的相关性、因果关系、组合成更大单位的可能性与它们之间点式互信息值大小的关系。



本章内容

1. 概率论略览
2. 齐夫定律
3. 信息论基础
4. 统计学习概念
5. 应用举例
6. 习题

 7. 附录

7. 附录

1. 证明前面P31公式(8): $H(L, q) = -\lim_{n \rightarrow \infty} \frac{1}{n} \log q(x_1^n)$

(又见《统计自然语言处理》P28, 公式(2-39))

2. 证明: 两个随机变量之间的平均互信息为非负值。

3. 前面P57上概率 $p^*(a|b)$ 的推导说明

7. 附录

1. 证明: $H(L, q) = -\lim_{n \rightarrow \infty} \frac{1}{n} \log q(x_1^n)$

布莱曼渐近均分性(Breiman's AEP)定理: 如果 X 是稳态的遍历性随机过程, 那么

$$H_{\text{rate}}(X) = -\lim_{n \rightarrow \infty} \frac{1}{n} \log p(x_1, x_2, \dots, x_n) = -\lim_{n \rightarrow \infty} \frac{1}{n} \log p(x_1^n)$$

该定理的证明:

假设 (x_1, x_2, \dots, x_n) 符合独立同分布。根据熵率的定义, 左边:

$$\begin{aligned} H_{\text{rate}}(X) &= \lim_{n \rightarrow \infty} \frac{1}{n} H(x_1^n) \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} H(x_1, x_2, \dots, x_n) \\ &= \lim_{n \rightarrow \infty} \left\{ -\frac{1}{n} \sum_{x_1 x_2 \dots x_n} p(x_1, x_2, \dots, x_n) \log p(x_1, x_2, \dots, x_n) \right\} \end{aligned}$$

7. 附录

$$\begin{aligned}
 &= \lim_{n \rightarrow \infty} \left\{ -\frac{1}{n} \sum_{x_1 x_2 \dots x_n} p(x_1, x_2, \dots, x_n) \left(\sum_{x_i} \log p(x_i) \right) \right\} \\
 &= \lim_{n \rightarrow \infty} \left\{ -\frac{1}{n} \left[\sum_{x_1 x_2 \dots x_n} p(x_1, x_2, \dots, x_n) \log p(x_1) + \sum_{x_1 x_2 \dots x_n} p(x_1, x_2, \dots, x_n) \log p(x_2) + \dots \right. \right. \\
 &\quad \left. \left. + \sum_{x_1 x_2 \dots x_n} p(x_1, x_2, \dots, x_n) \log p(x_n) \right] \right\} \\
 &= \lim_{n \rightarrow \infty} \left\{ -\frac{1}{n} \left[\sum_{x_1 x_2 \dots x_n} p(x_1) p(x_2) \dots p(x_n) \log p(x_1) + \sum_{x_1 x_2 \dots x_n} p(x_1) p(x_2) \dots p(x_n) \log p(x_2) \right. \right. \\
 &\quad \left. \left. + \dots + \sum_{x_1 x_2 \dots x_n} p(x_1) p(x_2) \dots p(x_n) \log p(x_n) \right] \right\}
 \end{aligned}$$

7. 附录

$$\begin{aligned}
 &= \lim_{n \rightarrow \infty} \left\{ -\frac{1}{n} \left[\sum_{x_1} p(x_1) \log p(x_1) \left(\sum_{x_2 x_3 \dots x_n} p(x_2) p(x_3) \dots p(x_n) \right) \right. \right. \\
 &\quad + \sum_{x_2} p(x_2) \log p(x_2) \left(\sum_{x_1, x_3 x_4 \dots x_n} p(x_1) p(x_3) p(x_4) \dots p(x_n) \right) \\
 &\quad \left. \left. + \dots + \sum_{x_n} p(x_n) \log p(x_n) \left(\sum_{x_1 x_2 \dots x_{n-1}} p(x_1) p(x_2) \dots p(x_{n-1}) \right) \right] \right\}
 \end{aligned}$$

由于 (x_1, x_2, \dots, x_n) 符合概率同分布，所以红线部分可以被看作是联合概率分布对所有的可能取值的求和，其值为1。

$$\begin{aligned}
 &= \lim_{n \rightarrow \infty} \left\{ -\frac{1}{n} \left[\sum_{x_1} p(x_1) \log p(x_1) + \sum_{x_2} p(x_2) \log p(x_2) + \dots + \sum_{x_n} p(x_n) \log p(x_n) \right] \right\} \\
 &= \lim_{n \rightarrow \infty} \left\{ -\frac{1}{n} \left[\sum_{x_1} p(x_1) \log p(x_1) + \sum_{x_2} p(x_2) \log p(x_2) + \dots + \sum_{x_n} p(x_n) \log p(x_n) \right] \right\} \\
 &= \lim_{n \rightarrow \infty} \left\{ \frac{1}{n} [H(x_1) + H(x_2) + \dots + H(x_n)] \right\} \\
 &= \lim_{n \rightarrow \infty} \left\{ \frac{1}{n} \times n \times H(x_i) \right\} \\
 &= H(x_i)
 \end{aligned}$$

7. 附录

而定理右边:

$$\begin{aligned} -\lim_{n \rightarrow \infty} \frac{1}{n} \log p(x_1^n) &= -\lim_{n \rightarrow \infty} \frac{1}{n} \log p(x_1, x_2, \dots, x_n) \\ &= -\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{x_i} \log p(x_i) \end{aligned} \quad (\text{基于概率同分布})$$

该式中, $\frac{1}{n} \sum_{x_i} \log p(x_i)$ 可以看作是 $\log p(x_i)$ 的均值, 而 $E(\log p(x_i))$

为其期望值(相当于下面式子中的 μ)。根据**辛钦大数定律**(Wiener-Khinchin Law of Large Numbers): **样本均值依概率收敛于期望值 μ** , 即

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| < \varepsilon) = 1$$

7. 附录

因此,
$$\lim_{n \rightarrow \infty} P \left(\left| \frac{1}{n} \sum_{x_i} \log p(x_i) - E(\log p(x_i)) \right| < \varepsilon \right) = 1$$

即
$$-\frac{1}{n} \sum_{x_i} \log p(x_i) \rightarrow -E(\log p(x_i)) = H(x_i)$$

(依概率)

所以, 左边等于等式右边。

参阅如下文献(英文第2版P58, 中文第2版P33):

Thomas M.Cover and Joy A.Thomas. Elements of Information Theory,
2nd edition, John Wiley & Sons. 2006

7. 附录

类似地，可以证明在 (x_1, x_2, \dots, x_n) 不满足独立同分布的条件下，该定理同样成立。请参阅下面的文献。

根据布莱曼渐近均分性定理(Breiman's AEP)，可以推广到：

$$\lim_{n \rightarrow \infty} \frac{1}{n} E(\log p(x_1^n)) = \lim_{n \rightarrow \infty} \frac{1}{n} \log p(x_1^n)$$

Paul H.Algoet and Thomas M.Cover. A Sanwich Proof of The Shannon-McMillan-Breiman Theorem. In The Annals of Probability 1998, Vol. 16, No.2 899-909



7. 附录

对于本章讲义公式(8), 《统计自然语言处理》 P28, 公式(2-39):

$$\begin{aligned} H(L, q) &= -\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{x_1^n} p(x_1^n) \log q(x_1^n) \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{x_1^n} p(x_1^n) \log q(x_1^n)^{-1} \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} E(\log q(x_1^n)^{-1}) \\ &\quad \text{(利用布莱曼渐进均分性定理的推广)} \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \log q(x_1^n)^{-1} \\ &= -\lim_{n \rightarrow \infty} \frac{1}{n} \log q(x_1^n) \end{aligned}$$

证毕。



7. 附录

1. 证明前面P31公式(8): $H(L, q) = -\lim_{n \rightarrow \infty} \frac{1}{n} \log q(x_1^n)$

(又见《统计自然语言处理》P28, 公式(2-39))

2. 证明: 两个随机变量之间的平均互信息为非负值。

3. 前面P57上概率 $p^*(a|b)$ 的推导说明

7. 附录

2. 证明：两个随机变量之间的平均互信息为非负值。

方法一：

证明：根据琴生不等式(Jensen inequality)的积分形式：

$$\frac{\int_a^b f(g(x))p(x)dx}{\int_a^b p(x)dx} \geq f\left(\frac{\int_a^b g(x)p(x)dx}{\int_a^b p(x)dx}\right)$$

其中， $f(x)$ 是凸函数， $g(x)$ 为任意函数。那么，

$$\begin{aligned} I(X,Y) &= \sum_{x \in X} \sum_{y \in Y} p(x,y) \log \left(\frac{p(x,y)}{p(x)p(y)} \right) = \int \int p(x,y) \log \left(\frac{p(x,y)}{p(x)p(y)} \right) dx dy \\ &\geq \int \int p(x,y) \left[-\log \left(\frac{p(x)p(y)}{p(x,y)} \right) \right] dx dy \\ &= -\log \left[\int \int p(x,y) \frac{p(x)p(y)}{p(x,y)} dx dy \right] = 0 \quad \text{证毕。} \end{aligned}$$

7. 附录

方法二：

证明：根据互信息的定义：
$$I(X;Y) = \sum_{x_i \in X} \sum_{y_j \in Y} p(x_i, y_j) \log_2 \frac{p(x_i, y_j)}{p(x_i)p(y_j)}$$

那么，

$$-I(X;Y) = \sum_{x_i \in X} \sum_{y_j \in Y} p(x_i, y_j) \log_2 \frac{p(x_i)p(y_j)}{p(x_i y_j)}$$

利用不等式： $\ln z \leq z-1$ ，且 $\log_2 z = \ln z \cdot \log_2 e$

所以， $\log_2 z \leq (z-1) \cdot \log_2 e$ ， $\log_2 \frac{p(x_i)p(y_j)}{p(x_i y_j)} \leq \left[\frac{p(x_i)p(y_j)}{p(x_i y_j)} - 1 \right] \cdot \log_2 e$

7. 附录

$$\begin{aligned} -I(X; Y) &= \sum_{x_i \in X} \sum_{y_j \in Y} p(x_i, y_j) \log_2 \frac{p(x_i)p(y_j)}{p(x_i, y_j)} \\ &\leq \sum_{x_i \in X} \sum_{y_j \in Y} p(x_i, y_j) \left[\frac{p(x_i)p(y_j)}{p(x_i, y_j)} - 1 \right] \log_2 e \\ &= \left[\sum_{x_i \in X} \sum_{y_j \in Y} p(x_i)p(y_j) - \sum_{x_i \in X} \sum_{y_j \in Y} p(x_i, y_j) \right] \log_2 e \\ &= \left[\sum_{x_i \in X} p(x_i) \sum_{y_j \in Y} p(y_j) - \sum_{x_i \in X} \sum_{y_j \in Y} p(x_i, y_j) \right] \log_2 e = 0 \end{aligned}$$

$$I(X; Y) \geq 0$$

证毕。

根据自然对数的性质： $\ln z \leq z-1, z > 0$, 当且仅当 $z=1$ 时取等号, 因此, 当且仅当 $\frac{p(x_i)p(y_j)}{p(x_i, y_j)}=1$ 时, 即 $p(x_i, y_j) = p(x_i)p(y_j)$ 时, $I(X; Y) = 0$ 。



7. 附录

1. 证明前面P31公式(8): $H(L, q) = -\lim_{n \rightarrow \infty} \frac{1}{n} \log q(x_1^n)$

(又见《统计自然语言处理》P28, 公式(2-39))

2. 证明: 两个随机变量之间的平均互信息为非负值。

3. 关于最大熵模型的推导说明

3. 概率 $p^*(y|x)$ 的推导说明

根据最大熵方法的基本思路，估计概率 $p(y|x)$ 时应满足如下两个基本约束：

$$p^* = \arg \max_{p \in P} H(p)$$

假设存在 k 个特征 $f_j (j = 1, 2, \dots, k)$ ，它们都在建模过程中对输出有影响，我们所建立的模型应满足所有这些特征，即所建立的模型 p 应该属于这 k 个特征约束下所产生的所有模型的集合 P ：

$$P = \{p \mid E_p(f_j) = E_{\tilde{p}}(f_j), j \in \{1, 2, \dots, k\}\}$$

7. 附录

根据条件熵的定义(理论值):

$$\begin{aligned} H(p) &= H(Y | X) = \sum_{x \in X} p(x) H(Y | X = x) \\ &= - \sum_{x, y} p(x) p(y | x) \log p(y | x) \end{aligned}$$

由于所建模型的概率分布 $p(x)$ 应符合已知样本中的概率分布 $\tilde{p}(x)$, 即: $p(x) = \tilde{p}(x)$, 因此,

$$H(p) = - \sum_{x, y} \tilde{p}(x) p(y | x) \log p(y | x) \quad (24)$$

即求解使 $H(p)$ 值最大的条件概率 $p^*(y|x)$:

$$p^*(y | x) = \arg \max_{p \in P} H(p) = \arg \max_{p \in P} \left(- \sum_{x, y} \tilde{p}(x) p(y | x) \log p(y | x) \right)$$

目标函数

7. 附录

如果有特征函数 $f_j(x, y)$ ，它在已知样本中的经验概率分布 $\tilde{p}(x, y)$ 可由下式计算得出：

$$\tilde{p}(x, y) \approx \frac{Count(x, y)}{\sum_{X, Y} Count(x, y)}$$

其中， $Count(x, y)$ 为 (x, y) 在训练语料中出现的次数。

f_j 在训练样本中关于经验概率分布的数学期望为：

$$E_{\tilde{p}}(f_j) = \sum_{x, y} \tilde{p}(x, y) f_j(x, y) \quad (25)$$

7. 附录

假设所建模型的理论分布为 $p(x, y)$ ，则特征 f_j 关于 $p(x, y)$ 的数学期望(理论值)为：

$$E_p(f_j) = \sum_{x,y} p(x, y) f_j(x, y) \quad (26)$$

由于 $p(x, y) = p(x)p(y|x)$ ，且 $p(x) = \tilde{p}(x)$ ，由此，(26)式变为：

$$E_p(f_j) = \sum_{x,y} \tilde{p}(x) p(y|x) f_j(x, y) \quad (27)$$

如果特征 f_j 对所建的模型是有用的，那么，所建模型中特征 f_j 的数学期望与它在已知样本中的数学期望应该相同，即：

$$E_p(f_j) = E_{\tilde{p}}(f_j) \quad (28)$$

该式称为该问题建模的约束方程，简称约束。

7. 附录

归纳上述各点：

$$p^* = \arg \max_{p \in P} H(p)$$

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

$$E_{\tilde{p}}(f_j) = \sum_{x,y} \tilde{p}(x, y) f_j(x, y)$$

$$E_p(f_j) = \sum_{x,y} \tilde{p}(x) p(y | x) f_j(x, y)$$

$$P = \{ p \mid E_p(f_j) = E_{\tilde{p}}(f_j), j \in \{1, 2, \dots, k\} \}$$

$$\sum_y p(y | x) = 1$$



7. 附录

这样，问题就变成了在满足一组约束的条件下求最优解的问题，可用拉格朗日乘子法解决此问题，从而证明，满足(28)式约束条件的解具有如下形式：

$$p^*(y|x) = \frac{1}{Z(x)} \exp\left(\sum_{j=1}^l \lambda_j \cdot f_j(x, y)\right)$$

$$\text{其中, } Z(x) = \sum_y \exp\left(\sum_{j=1}^l \lambda_j \cdot f_j(x, y)\right)$$

$Z(x)$ 为保证对所有 x ，使得 $\sum_y p(y|x) = 1$ 的归一化常量。

更详细的推导，可参阅：李航著，《统计学习方法》第2版，清华大学出版社，2019

谢谢!

Thanks!

