

白马负金羁

数据挖掘 | 统计分析 | 图像处理 | 程序设计

目录视图

摘要视图

RSS 订阅

个人资料



白马负金羁

关注

发私信

访问：2117601次

积分：26143

等级：BLOG > 7

排名：第252名

原创：326篇

转载：15篇

译文：0篇

评论：3985条

算法之美



算法之美

:隐匿在数据结构背后的原理

(C++版)

源码获取, 读者答疑, 请加算法学习群

群容量有限, 未购书者勿扰

(495573865)

图像处理

图灵赠书——程序员11月书单

【思考】Python这么厉害的原因竟然是！

感恩节赠书：《深度学习》等异步社区优秀图书和作译者评选启动！

每周荐书：京东架构、Linux内核、Python全栈

自然语言处理中的N-Gram模型详解

标签：NLP N-Gram 自然语言处理 模糊匹配 编辑距离

2016-04-29 21:32

44413人阅读

评论(2)

收藏

举报

分类：

自然语言处理与信息检索 (16)

版权声明：本文为博主原创文章，未经博主允许不得转载。

目录(?)

[+]

N-Gram（有时也称为N元模型）是自然语言处理中一个非常重要的概念，通常在NLP中，人们基于一定的语料库，可以利用N-Gram来预计或者评估一个句子是否合理。另外一方面，N-Gram的另外一个作用是用来评估两个字符串之间的差异程度。这是模糊匹配中一段。本文将从此开始，进而向读者展示N-Gram在自然语言处理中的各种power

基于N-Gram模型定义的字符串距离

利用N-Gram模型评估语句是否合理

使用N-Gram模型时的数据平滑算法

欢迎关注白马负金羁的博客 <http://blog.csdn.net/baimafujinji>，为保证公式、图表得以正确显示，强烈建议你从该地址上查看原版博文。本博客主要关注方向包括：数字图像处理、算法设计与分析、数据结构、机器学习、数据挖掘、统计分析方法、自然语言处理。

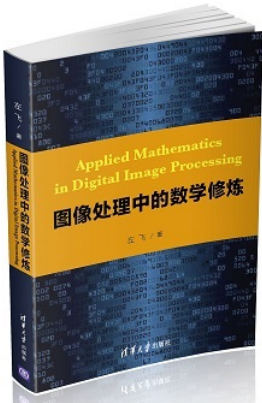
基于N-Gram模型定义的字符串距离

在自然语言处理时，最常用也最基础的一个操作是就是“模式匹配”，或者称为“字符串查找”。而模式匹配（字符串查找）又分为精确匹配和模糊匹配两种。

所谓精确匹配，大家应该并不陌生，比如我们要统计一篇文章中关键词“information”出现的次数，这时所使用的方法就是精确的模式匹配。这方面的算法也比较多，而且应该是计算机相关专业必修的基础课中都会涉及到的内容，例如KMP算法、BM算法和BMH算法等等。

另外一种匹配就是所谓的模糊匹配，它的应用也随处可见。例如，一般的文字处理软件（例如，Microsoft Word等）都会提供拼写检查功能。当你输入一个错误的单词，例如“informtaion”时，系统会提示你是否要输入的词其实是“information”。将一个可能错拼单词映射到一个推荐的正确拼写上所采用的技术就是模糊匹配。

模糊匹配的关键在于如何衡量两个长得很像的单词（或字符串）之间的“差异”。这种差异通常又称为“距离”。这方面的具体算法有很多，例如基于编辑距离的概念，人们设计出了



《图像处理中的数学修炼》

图像处理书籍读者群
(155911675)

还未购入本书者，切勿尝试加群，本群
谢绝吃瓜群众以及毫无诚信者围观，妄
图浑水摸鱼的行为最终都只能自取其辱

联系方式

- 1. 在博客文章下留言，[博客私信一律不回](#)。
- 2. 邮件fzuo#foxmail.com，将#换成@。

文章分类

- 编程语言与程序设计 (26)
- 图像与信号处理 (26)
- 数据结构与算法 (24)
- 其他杂文 (15)
- 应用技巧 (25)
- 经济研究 (15)
- 机器学习 (24)
- 数据挖掘十大算法 (15)
- 自然语言处理与信息检索 (17)
- 图像处理中的数学 (38)
- 线性代数 (20)
- 多核编程与并行计算 (15)
- 废言集 (28)
- 文学与诗歌 (10)
- 学习方法与方法论 (14)
- 已出版图书的相关资源 (15)
- 深度学习与TensorFlow (12)
- 有关LLVM的一切 (2)

阅读排行

- 在Eclipse中进行C/C++开发的... (46747)
- 自然语言处理中的N-Gram模... (44390)
- 暗通道优先的图像去雾算法 (... (36738)
- 如何学好图像处理——从小白... (29400)
- 在R中使用支持向量机 (SVM... (28411)
- 自己动手用C++写的图像处理... (27656)
- 暗通道优先的图像去雾算法 (... (26183)
- 从泊松方程到泊松融合 (Poiss... (24473)
- 图像的泊松(Poisson)编辑、泊... (23796)

Smith-Waterman 算法和Needleman-Wunsch 算法，其中后者还是历史上最早的应用动态规划思想设计的算法之一。现在Smith-Waterman 算法和Needleman-Wunsch 算法在生物信息学领域也有重要应用，研究人员常常用它们来计算两个DNA序列片段之间的“差异”（或称“距离”）。甚至于在LeetCode上也有一道 [“No.72 Edit Distance”](#)，其本质就是在考察上述两种算法的实现。可见相关问题离我们并不遥远。

N-Gram在模糊匹配中的应用

事实上，笔者在新出版的[《算法之美——隐匿在数据结构背后的原理》](#)一书中已经详细介绍了包括Needleman-Wunsch算法、Smith-Waterman算法、N-Gram算法、Soundex算法、Phonix算法等在内的多种距离定义算法（或模糊匹配算法）。而今天为了引出N-Gram模型在NLP中的其他应用，我们首先来介绍一下如何利用N-Gram来定义字符串之间的距离。

我们除了可以定义两个字符串之间的编辑距离（通常利用Needleman-Wunsch算法或Smith-Waterman算法）之外，还可以定义它们之间的N-Gram距离。N-Gram（有时也称为N元模型）是自然语言处理中一个非常重要的概念。假设有一个字符串 s ，那么该字符串的N-Gram就表示按长度 N 切分原词得到的词段，也就是 s 中所有长度为 N 的子字符串。设想如果有两个字符串，然后分别求它们的N-Gram，那么就可以从它们的共有子串的数量这个角度去定义两个字符串间的N-Gram距离。但是仅仅是简单地对共有子串进行计数显然也存在不足，这种方案显然忽略了两个字符串长度差异可能导致的问题。比如字符串 `girl` 和 `girlfriend`，二者所拥有的公共子串数量显然与 `girl` 和其自身所拥有的公共子串数量相等，但是我们并不能据此认为 `girl` 和 `girlfriend` 是两个等同的匹配。

为了解决该问题，有学者便提出以非重复的N-Gram分词为基础来定义 N-Gram距离这一概念，可以用下面的公式来表述：

$$|G_N(s)| + |G_N(t)| - 2 \times |G_N(s) \cap G_N(t)|$$

此处， $|G_N(s)|$ 是字符串 s 的 N-Gram集合， N 值一般取2或者3。以 $N = 2$ 为例对字符串 `Gorbachev`和`Gorbechyov`进行分段，可得如下结果（我们用下划线标出了其中的公共子串）。

Go, or, rb, ba, ac, ch, he, ev
Go, or, rb, be, ec, ch, hy, yo, ov

结合上面的公式，即可算得两个字符串之间的距离是 $8 + 9 - 2 \times 4 = 9$ 。显然，字符串之间的距离越小，它们就越接近。当两个字符串完全相等的时候，它们之间的距离就是0。

利用N-Gram计算字符串间距离的Java实例

在[《算法之美——隐匿在数据结构背后的原理》](#)一书中，我们给出了在C++下实现的计算两个字符串间N-Gram距离的函数，鉴于全书代码已经在本博客中发布，这里不再重复列出。事实上，很多语言的函数库或者工具箱中都已经提供了封装好的计算 N-Gram 距离的函数，下面这个例子演示了在Java中使用N-Gram 距离的方法。

针对这个例子，这里需要说明的是：

基于直方图的图像增强算法 (... (23630)

最新评论

K-means算法原理与R语言实例qq_35925656 : @Awon_Lee:好的 谢谢哦

图像处理之让手心长出眼睛，其实嘴也可...白马负金羁 : @Ingrid_W:这是伪代码

TF-IDF算法解析与Python实现剑与星辰 : def make_count(text): #将原始的text文件生成预处理后的count文件...

图像处理之让手心长出眼睛，其实嘴也可...Ingrid_W : 你好，请问adjacent ()和pixel () 两个函数是自己写的两个功能函数嘛？

暗通道优先的图像去雾算法 (上)夕拾li : @qq_40668683:你好，能问你要一下老师给的代码吗？邮箱1552216813@qq.com

算法之美隆重上市欢迎关注 (另附勘误表...luckstudent : 第99页，注释中指出preCure指向要删除的节点，而在代码中写成了curPre。另外delete c...

蒙特卡洛采样之拒绝采样 (Reject Sampling...liuxiaoyang_222 : 楼主，为啥我没有找到重要性抽样的文章？是没有吗？还是我没找到？

蒙特卡洛采样之拒绝采样 (Reject Sampling...liuxiaoyang_222 : 楼主厉害，正好用的上，讲的很明白

TF-IDF算法解析与Python实现白马负金羁 : @nailunhan7929:参考 http://scikit-learn.org/stable/m...

TF-IDF算法解析与Python实现nailunhan7929 : 你好，请问from sklearn.feature_extraction.text import T...

- 调用函数需要引用lucene的JAR包，我所使用的是lucene-suggest-5.0.0.jar
- 前面我们所给出的算法计算所得为一个绝对性的距离分值。而Java中所给出的函数在此基础上进行了归一化，也就是说所得之结果是一个介于0~1之间的浮点数，即0的时候表示两个字符串完全不同，而1则表示两个字符串完全相同。

```
1 import org.apache.lucene.search.spell.*;
2
3 public class NGram_distance {
4
5     public static void main(String[] args) {
6
7         NGramDistance ng = new NGramDistance();
8         float score1 = ng.getDistance("Gorbachev", "Gorbechyov");
9         System.out.println(score1);
10        float score2 = ng.getDistance("girl", "girlfriend");
11        System.out.println(score2);
12    }
13 }
```

有兴趣的读者可以在引用相关JAR包之后在Eclipse中执行上述Java程序，你会发现，和我们预期的一样，字符串Gorbachev和Gorbechyov所得之距离评分较高（=0.7），说明二者很接近；而girl和girlfriend所得之距离评分并不高（=0.3999），说明二者并不很接近。

利用N-Gram模型评估语句是否合理

从现在开始，我们所讨论的N-Gram模型跟前面讲过N-Gram模型从外在来看已经大不相同，但是请注意它们内在的联系（或者说本质上它们仍然是统一的概念）。

为了引入N-Gram的这个应用，我们从几个例子开始。
首先，从统计的角度来看，自然语言中的一个句子 s 可以由任何词串构成，不过概率 $P(s)$ 有大有小。例如：

- s_1 = 我刚吃过晚饭
- s_2 = 刚我吃过晚饭吃

显然，对于中文而言 s_1 是一个通顺而有意义的句子，而 s_2 则不是，所以对于中文来说， $P(s_1) > P(s_2)$ 。但不同语言来说，这两个概率值的大小可能会反转。

其次，另外一个例子是，如果我们给出了某个句子的一个节选，我们其实可以能够猜测后续的词应该是什么，例如

- the large green __ . Possible answer may be "mountain" or "tree" ?
- Kate swallowed the large green __ . Possible answer may be "pill" or "broccoli" ?

显然，如果我们知道这个句子片段更多前面的内容的前提下，我们会得到一个更加准确的答案。这就告诉我们，前面的（历史）信息越多，对后面未知信息的约束就越强。

如果我们有一个由 m 个词组成的序列（或者说一个句子），我们希望算得概率 $P(w_1, w_2, \dots, w_m)$ ，根据链式规则，可得

$$P(w_1, w_2, \dots, w_m) = P(w_1)P(w_2|w_1)P(w_3|w_1, w_2) \cdots P(w_m|w_1, \dots, w_{m-1})$$

这个概率显然并不好算，不妨利用马尔科夫链的假设，即当前这个词仅仅跟前面几个有限的词相关，因此也就不必追溯到最开始的那个词，这样便可以大幅缩减计算式的长度。即

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

特别地，对于 n 取得较小值的情况

当 $n = 1$, 一个一元模型 (unigram model) 即为

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i)$$

当 $n = 2$, 一个二元模型 (bigram model) 即为

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i | w_{i-1})$$

当 $n = 3$, 一个三元模型 (trigram model) 即为

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i | w_{i-2} w_{i-1})$$

接下来的思路就比较明确了，可以利用最大似然法来求出一组参数，使得训练样本的概率取得最大值。

- 对于unigram model而言，其中 $c(w_1, \dots, w_n)$ 表示 n -gram w_1, \dots, w_n 在训练语料中出现的次数， M 是语料库中的总字数（例如对于 yes no no no yes 而言， $M = 5$ ）

$$P(w_i) = \frac{C(w_i)}{M}$$

- 对于bigram model而言，

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

- 对于 n -gram model 而言，

$$P(w_i | w_{i-n+1}, \dots, w_{i-1}) = \frac{C(w_{i-n+1}, \dots, w_i)}{C(w_{i-n+1}, \dots, w_{i-1})}$$

来看一个具体的例子，假设我们现在有一个语料库如下，其中 $< s1 > < s2 >$ 是句首标记， $< /s2 > < /s1 >$ 是句尾标记：

$< s1 > < s2 > \text{yes no no no no yes} < /s2 > < /s1 >$
 $< s1 > < s2 > \text{no no no yes yes yes no} < /s2 > < /s1 >$

下面我们的任务是来评估如下这个句子的概率：

$$< s1 > < s2 > yes \quad no \quad no \quad yes < /s2 > < /s1 >$$

我们来演示利用trigram模型来计算概率的结果

$$P(yes | < s1 > < s2 >) = \frac{1}{2}, \quad P(no | < s2 > yes) = 1$$
$$P(no | yes \quad no) = \frac{1}{2}, \quad P(yes | no \quad no) = \frac{2}{5}$$
$$P(< /s2 > | no \quad yes) = \frac{1}{2}, \quad P(< /s1 > | yes < /s2 >) = 1$$

所以我们要求的概率就等于：

$$\frac{1}{2} \times 1 \times \frac{1}{2} \times \frac{2}{5} \times \frac{1}{2} \times 1 = 0.05$$

再举一个来自文献[1]的例子，假设现在有一个语料库，我们统计了下面一些词出现的数量

i	want	to	eat	chinese	food	lunch	spend
2533	927	2417	746	158	1093	341	278

下面这个概率作为其他一些已知条件给出：

$$P(i | < s >) = 0.25 \qquad P(english | want) = 0.0011$$
$$P(food | english) = 0.5 \qquad P(< /s > | food) = 0.68$$

下面这个表给出的是基于Bigram模型进行计数之结果

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

例如，其中第一行，第二列 表示给定前一个词是 “i” 时，当前词为 “want” 的情况一共出现了827次。据此，我们便可以算得相应的频率分布表如下。

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

因为我们从表1中知道 “i” 一共出现了2533次，而其后出现 “want” 的情况一共有827次，所以 $P(want|i) = 827/2533 \approx 0.33$

现在设 $s_1 = \langle s \rangle i \text{ want english food } \langle /s \rangle$ ，则可以算得

$$\begin{aligned} P(s_1) &= P(i | \langle s \rangle) P(want|i) P(english|want) P(food|english) P(\langle /s \rangle | foo) \\ &= 0.25 \times 0.33 \times 0.0011 \times 0.5 \times 0.68 = 0.000031 \end{aligned}$$

使用N-Gram模型时的数据平滑算法

有研究人员用150万词的训练语料来训练 trigram 模型，然后用同样来源的测试语料来做验证，结果发现23%的 trigram 没有在训练语料中出现过。这其实就意味着上一节我们所计算的那些概率有空为 0，这就导致了数据稀疏的可能性，我们的表3中也确实有些为0的情况。对语言而言，由于数据稀疏的存在，极大似然法不是一种很好的参数估计办法。

这时的解决办法，我们称之为“平滑技术”（Smoothing）或者“减值”（Discounting）。其主要策略是把在训练样本中出现过的事件的概率适当减小，然后把减小得到的概率密度分配给训练语料中没有出现过的事件。实际中平滑算法有很多种，例如：

- Laplacian (add-one) smoothing
- Add-k smoothing
- Jelinek-Mercer interpolation
- Katz backoff
- Absolute discounting
- Kneser-Ney

对于这些算法的详细介绍，我们将在后续的文章中结合一些实例再来进行讨论。

A Final Word

如果你能从前面那些繁冗、复杂的概念和公式中挺过来，恭喜你，你对N-Gram模型已经有所认识了。尽管，我们还没来得及探讨平滑算法（但它即将出现在我的下一篇博文里，如果你觉得还未过瘾的话），但是其实你已经掌握了一个相对powerful的工具。你可以可能会问，在实践中N-Gram模型有哪些具体应用，作为本文的结束，主页君便在此补充几个你曾见过的或者曾经好奇它是如何实现例子。

Eg.1

搜索引擎（Google或者Baidu）、或者输入法的猜想或者提示。你在用百度时，输入一个或几个词，搜索框通常会以下拉菜单的形式给出几个像下图一样的备选，这些备选其实是在猜想你想要搜索的那个词串。再者，当你用输入法输入一个汉字的时候，输入法通常可以联系出一个完整的词，例如我输入一个“刘”字，通常输入法会提示我是否要输入的是“刘备”。通过上面的介绍，你应该能够很敏锐的发觉，这其实是以N-Gram模型为基础来实现的，如果你能有这种觉悟或者想法，那我不得不恭喜你，都学会抢答了！

lots of |

lots of love

lots of fish

lots of discharge

lots of lollies

Press Enter to search.

Eg.2

某某作家或者语料库风格的文本自动生成。这是一个相当有趣的话题。来看下面这段话（该例子取材自文献【1】）：

“You are uniformly charming!” cried he, with a smile of associating and now and then I bowed and they perceived a chaise and four to wish for.

你应该还没有感觉到它有什么异样吧。但事实上这并不是由人类写出的句子，而是计算机根据 Jane Austen 的语料库利用 trigram 模型自动生成的文段。（Jane Austen 是英国著名女作家，代表作有《傲慢与偏见》等）

再来看两个例子，你是否能看出它们是按照哪位文豪（或者语料库）的风格生成的吗？

- This shall forbid it should be branded, if renown made it empty.
- They also point to ninety nine point six billion dollars from two hundred four oh three percent of the rates of interest stores as Mexico and Brazil on market conditions.

答案是第一个是莎士比亚，第二个是华尔街日报。最后一个问题留给读者思考，你觉得上面两个文段所运用的 n-gram 模型中，n 应该等于多少？

推荐阅读和参考文献：

[1] Speech and Language Processing. Daniel Jurafsky & James H. Martin, 3rd. Chapter 4

[2] 本文中的一些例子和描述来自 北京大学 常宝宝 以及 The University of Melbourne “Web Search and Text Analysis” 课程的幻灯片素材

顶 18
踩 1

- 上一篇 从小蝌蚪找妈妈谈“机器学习VS数据挖掘”
- 下一篇 机器学习中的隐马尔科夫模型（HMM）详解

相关文章推荐

- TensorFlow 聊天机器人
 - MySQL在微信支付下的高可用运营--莫晓东
 - cnn、rnn相结合进行文本分类
 - 容器技术在58同城的实践--姚远
 - 自然语言处理中的N-Gram模型详解
 - SDCC 2017之容器技术实战线上峰会
 - 自然语言处理中的N-Gram模型详解
 - SDCC 2017之数据库技术实战线上峰会
- 自然语言处理3-N-gram模型
 - 腾讯云容器服务架构实现介绍--董晓杰
 - 统计自然语言处理——n元语法（马尔可夫模型）...
 - 微博热点事件背后的数据库运维心得--张冬洪
 - 统计自然语言处理（统计推理：稀疏数据集上的n元..
 - 基于统计模型的自然语言处理NLP技术应用简介
 - 隐马尔可夫模型及其在自然语言处理中的应用
 - 自然语言处理的形式模型 冯志伟

查看评论



迷雾forest

2楼 2017-06-29 16:20发表

很好的文章，谢谢分享！！



独上高楼望天涯

1楼 2016-10-25 11:21发表

最后一个问题：我觉得对于莎士比亚的N=2~3，对于华尔街日报的话N=3~5，可能这样比较合适吧。总之应该是后面的N应该要比前面的大一些。

您还没有登录,请[\[登录\]](#)或[\[注册\]](#)

* 以上用户言论只代表其个人观点，不代表CSDN网站的观点或立场

