

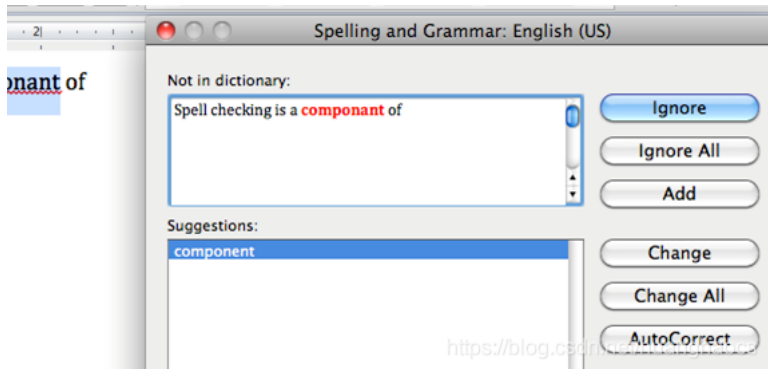
# 错 (Spelling Correction)

1:51:41 692 收藏 版权

标签: NLP query rewrite query 改写 拼写纠错

拼纠错的一节，感觉讲的比较好，虽然课程很老，但是讲的知识，在目前的query改写、拼写纠错还是很

，又称拼写检查（Spelling Checker），往往被用于字处理软件、输入法和搜索引擎中，如下所示：



u mean: 自然语言处理

任务：  
错误类型不同，分为Non-word Errors和Real-word Errors。前者指那些拼写错误后的词本身就不合法，如  
者指那些拼写错误后的词仍然是合法的情况，如将“there”错误拼写为“three”（形近），将“peace”错误拼  
写为“too”（同音）。  
纠错，如把“hte”自动校正为“the”，或者给出一个最可能的拼写建议，甚至一个拼写

点赞Mark关注该博主, 随时了解TA的最新博文

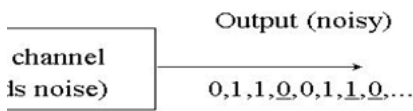
点赞 评论 分享 收藏 举报 关注 一键三连

被词典所包含的word均被当作spelling error，识别准确率依赖词典的规模和质量。  
词典中与error最近似的word，常见的方法有Shortest weighted edit distance和Highest noisy channel

rd都作为spelling error candidate。  
和拼写等角度，查找与word最近似的words集合作为拼写建议，常见的方法有Highest noisy channel

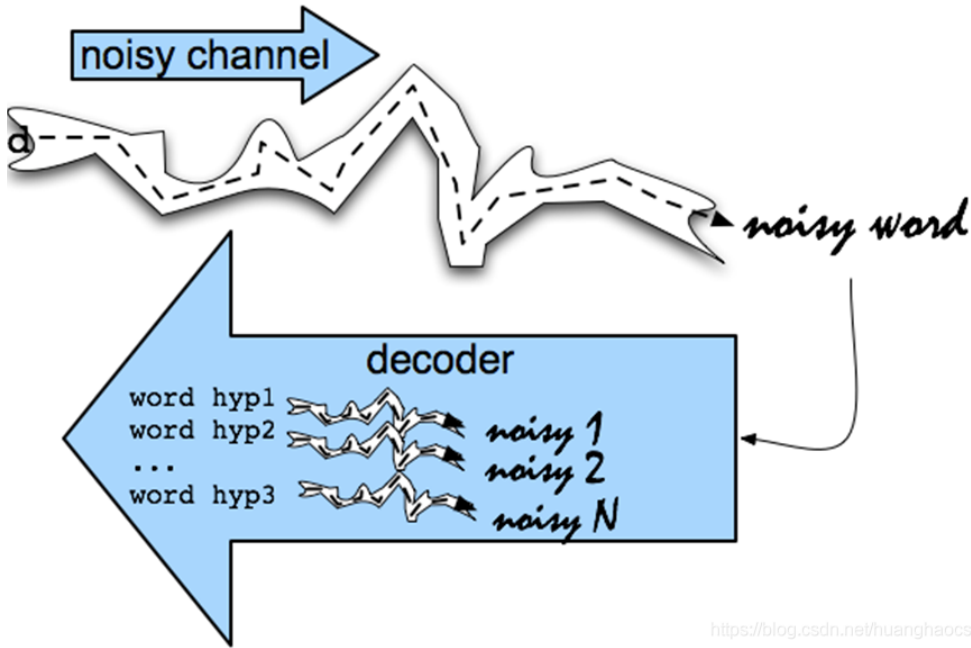
odel的拼写纠错

模型，或称信源信道模型，这是一个普适性的模型，被用于语音识别、拼写纠错、机器翻译、中文分  
互用领域。其形式很简单，如下图所示：



信号恢复输入信号，形式化定义为：

$$\frac{P(O|I)P(I)}{P(O)} = \arg \max_I P(O|I)P(I)$$



<https://blog.csdn.net/huanghaocs>

写作Original word通过noisy channel转换得到，现在已知noisy word（用x表示）如何求得最大可能的  
如下：

$$P(w|x)$$

$$\frac{P(x|w)P(w)}{P(x)}$$

贝叶斯公式

$$P(x|w)P(w)$$

对于特定的x，P(x)不变

<https://blog.csdn.net/huanghaocs>

点赞Mark关注该博主, 随时了解TA的最新博文

率，二者可以基于训练语料库建立语言模型和转移

点赞

评论

分享

收藏

举报

关注

一键三连

error correction的例子加以解释：  
词典匹配容易确定为“Non-word spelling error”；然后通过计算最小编辑距离获取最相似的candidate  
这里的最小编辑距离涉及四种操作：

## f two adjacent letters

andidate orrection	Correct Letter	Error Letter	Type
ctress	t	–	deletion
ress	–	a	insertion
aress	ca	ac	transposition
ccess	c	r	substitution
cross	o	e	substitution
cres	–	s	insertion
cres	–	s	insertion

为1，几乎所有的拼写错误编辑距离小于等于2，基于此，可以减少大量不必要的计算。  
建议候选集（candidate w），此时，我们希望选择概率最大的w作为最终的拼写建议，基于噪声信道模型  
P(w)。  
可以很容易建立语言模型，即可得到P(w)，如下表所示，计算Unigram Prior Probability（word总数：

cy of word	P(word)
9, 321	. 0000230573
220	. 0000005442
686	. 0000016969
37, 038	. 0000916207
120, 844	. 0002989314
12, 874	. 0000318463

and word x=x1 x2 x3 ... xm, correct word w=w1 w2 w3 ... wn>pair计算del、ins、sub和trans四种转移矩

$\text{count}(xy \text{ typed as } x)$   
 $\text{count}(x \text{ typed as } xy)$   
 $\text{count}(x \text{ typed as } y)$   
 $:$   $\text{count}(xy \text{ typed as } yx)$

$\frac{\text{del}[w_{i-1}, w_i]}{\text{count}[w_{i-1} w_i]}$ , if deletion  
 $\frac{\text{ins}[w_{i-1}, x_i]}{\text{count}[w_{i-1}]}$ , if insertion  
 $\frac{\text{sub}[x_i, w_i]}{\text{count}[w_i]}$ , if substitution  
 $\frac{\text{trans}[w_i, w_{i+1}]}{\text{count}[w_i w_{i+1}]}$ , if transposition

Correct Letter	Error Letter	x   w	P(x   word)
	-	c   ct	.000117
	a	a   #	.00000144
a	ac	ac   ca	.00000164
	r	r   c	.000000209
	e	e   o	.00000093
	s	es   e	.0000321
	s	ss   s	.0000342

点赞Mark关注该博主, 随时了解TA的最新博文

点赞

评论

分享

收藏

举报

关注

一键三连

Error Letter	x   w	P(x   word)	P(word)	10 <sup>9</sup> * P(x   w)P(w)
–	c   ct	.000117	.0000231	2.7
a	a   #	.00000144	.000000544	.00078
ac	ac   ca	.00000164	.00000170	.0028
r	r   c	.000000209	.0000916	.019
e	e   o	.0000093	.000299	2.8
s	es   e	.0000321	.0000318	1.0
s	ss   s	.0000342	.0000318	1.0

注更大。  
am，也可以推广到bigram，甚至更高阶，以较好的融入上下文信息。  
ess whose combination of sass and glamour...”，计算bigram为：  
/hose|actress) = .0010  
/hose|across) = .000006

00021\*.0010 = 210 x10<sup>-10</sup>  
00021\*.000006 = 1 x10<sup>-10</sup>

%的拼写错误都属于Real-word类型，与Non-word类型相比，纠错难度更大，因为句子中的每个word都被  
去分两步：

## rd in sentence

andidate set

itself

letter edits that are English words

at are homophones

candidates

annel model

cific classifier

<https://blog.csdn.net/huanghaocs>

wn，为每个wi生成candidate set，如下：

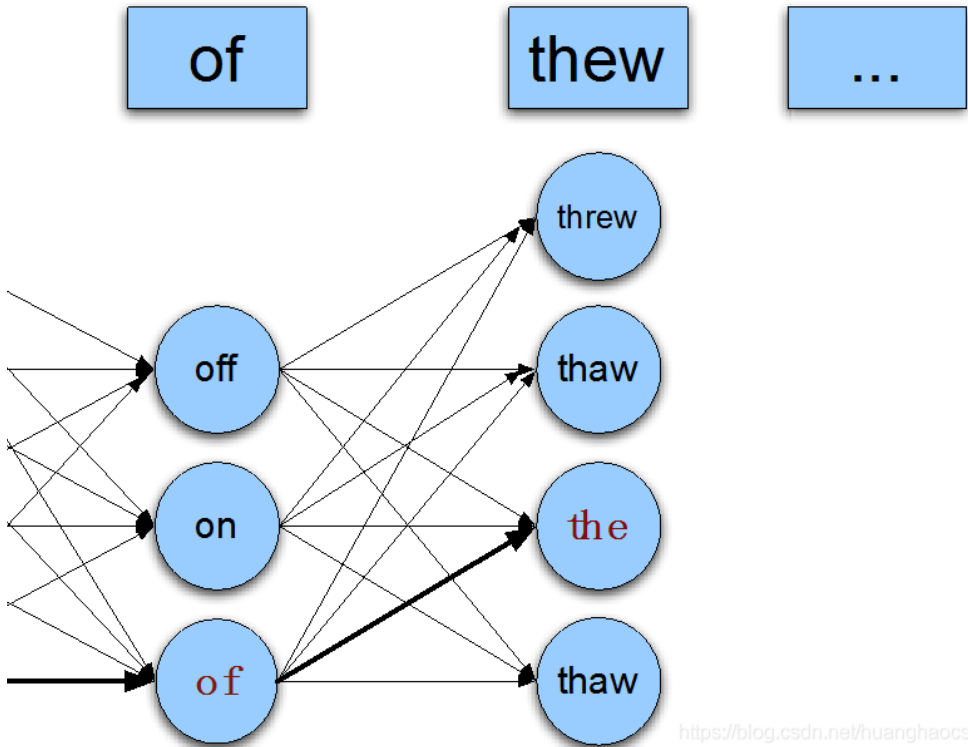
w""1 ,...}  
w""2 ,...}

w""n ,...}

自动纠错后的句子，与中文分词、音字转换等应用相同，可以表示成词网格形式，转化为HMM的解码过

点赞Mark关注该博主, 随时了解TA的最新博文

点赞 评论 分享 收藏 举报 关注 一键三连



中最多有一个word存在spelling error（事实上，所出现的情况也的确如此）。

下HCI（Human Computer Interface）准则：

confident in correction

orrect

nfident

he best correction

nfident

a correction list

ident

ag as an error

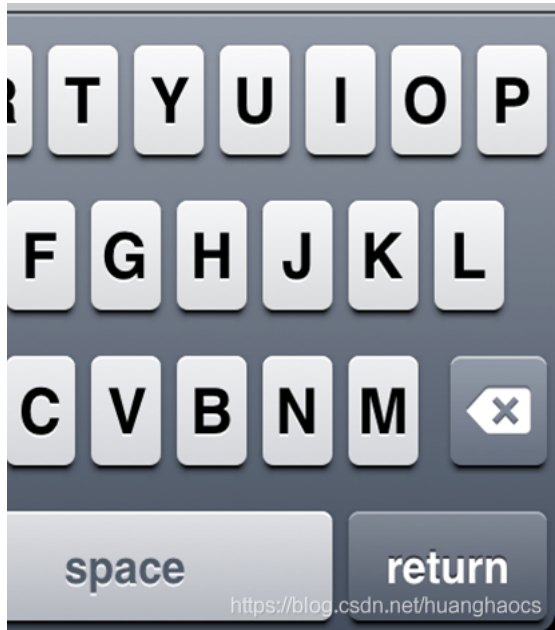
sitivity），需要对语言模型进行特别的处理，如：

$x | w)P(w)^\lambda$  Learn  $\lambda$  from a development test set

点赞Mark关注该博主, 随时了解TA的最新博文

形同音导致，所以，有些系统将“error model”转化为“l 点赞 评论 分享 收藏 举报 关注 一键三连

引入spelling error pair，据此可以对转移矩阵进行加权。



与分类问题，通过构建训练语料库，抽取features，训练分类模型，预测新实例等一系列过程解决，如

ifier for a specific pair like:

weather

” within +/- 10 words

VERB

not

<https://blog.csdn.net/huanghaocs>

haimizhao的白板

2073

y-Document(摘要)看成源语言和目标语言，用翻译模型计算二者的短语与短语之间的对齐关系，扩展Query的同时起...

拼写纠错 (Spelling Correction) ”

overstack的专栏

3168

我在Coursera启动了在线自然语言处理课程，由NLP领域大牛Dan Jurafsky 和 Chirs Manning教授授课： <https://class...>

高权重

抢沙发

评论

weixin\_33835103的博客

851

题还是在于用户query和商品描述之间存在GAP，特别是中长尾query。把问题分成以下几种类型： 多种描述：划痕...

mingo220的博客

450

纠错技术和架构\_chengyong...

12-25

成的query纠错相对更加精准,每个词存在上下文关系约束,整个query的意图更加明确。通过对query分词,查找每个词...

(NMT)的语法纠错算法\_怎...

1-4

rection,下同)是NLP的一个子领域,在搜索query纠错、语音纠错、舆情文本纠错等诸多直接与普通大众交互的场景得...

纠错技术和架构

点赞Mark关注该博主,随时了解TA的最新博文

技术和架构 1 背景 如今，搜索引擎是人们的获取信息最重要

点赞 评论 分享 收藏 举报 关注 一键三连