



SRI-LM toolkit

2003/12/16
Louis Tsai
Speech Lab, CSIE, NTNU
louis@csie.ntnu.edu.tw



SRILM

- 由SRI International所開發出來的tool，用在Statistical Language Models
- <http://www.speech.sri.com/projects/srilm>
- 目前版本 Ver 1.3.3 2003/03



Create language model

$$\tilde{W} = \arg \max_W P(W | X) \cong \arg \max_W P(X | W)P(W)$$

- SRI-LM toolkit :
 - Step 1. Calculate *ngram count* from *training data*
 - Step 2. Create *LM* from *ngram count*



training data (ssp.train)

已斷好詞的文件

關心孩子學業
可是又苦無時間為孩子挑選參考書的家長們
從今天起
多了一個輕鬆的選擇
只要趁工作之餘逛到巷口的
填個國小參考書預購單
就能以八五折的價格在限定時間內取貨了
這項服務明年還將擴展到國中參考書
由於國小教科書版本不同
參考書也相當多
各書店進書狀況和數量也不相同



Vocabulary

(Lexicon2003-72k.txt)

巴
八
扒
叭
...
巴巴
叭叭
八寶
八方
八德
八堵
吧女
...

共有 71695 個 words

ngram count file (ssp.n3.count)

在整個training data中出現的次數

Unigram

Bigram

Trigram

想像得到	4
想像得到的	3
想像得到的事	2
想像得到的 </s>	1
想像得到球員	1
想像得到球員可以	1
鳳凰	123
鳳凰 駱	4
鳳凰 駱 </s>	4
鳳凰 花	5
鳳凰 花開	1
鳳凰 花 </s>	4
鳳凰 </s>	20
鳳凰 自行車	1
鳳凰 自行車 </s>	1
鳳凰 城	20
鳳凰 城 </s>	8
鳳凰 城 大學	1

鳳凰 城 的	1
鳳凰 城 警方	1
鳳凰 城 太陽	5
鳳凰 城 電	3
鳳凰 城 水星	1
鳳凰 機場	1
鳳凰 機場 供	1
鳳凰 一樣	1
鳳凰 一樣 快速	1
鳳凰 中隊	1
鳳凰 中隊 </s>	1
鳳凰 自然	1
鳳凰 自然 來	1
鳳凰 般	3
鳳凰 般 重	1
鳳凰 般 再度	2
鳳凰 科技	6
.....	

LM Format (ssp.n3.kn.lm)

以\data\開頭

```
\data\  
ngram 1=71697  
ngram 2=2795822  
ngram 3=1700527
```

各種ngram的個數

```
\1-grams:
```

\1-grams:是
unigram的開頭

```
-1.804081 </s>  
-99 <s> -0.5467538  
-2.548589 - -0.6464233  
-3.989615 -- -0.2506604  
-6.685257 --恐怖 -0.1400223  
-6.685257 --恐怖攻擊 -0.1128442  
-6.075799 -丁點 -0.1827688  
-4.105376 -九九 -1.257224  
-5.808214 -了百了 -0.308177  
-5.921716 -刀兩斷 -0.3545529  
-3.940887 -下 -0.367599
```

Log probability
(Base 10)

Log of backoff
weight (Base 10)

最高的ngram沒
有backoff weight

```
...  
-0.4833062 糊裏糊塗  
-0.6593975 餐廳裏往  
-0.6593975 館裏大家  
-0.2190889 鎮裏不同  
-0.4833062 鐵窗裏竄  
-1.223669 <s> 恒大  
-1.223669 <s> 恒生  
-0.3077702 <s> 恒春  
-0.3159989 小型 恒指  
-0.4833062 屏東 恒春  
-0.2190889 范 恒山  
-0.4833062 從 恒春  
-0.4833062 邊 恒雄  
-0.6593975 上粧 </s>  
-0.2759993 化粧品
```

```
\end\  

```

以\end\結尾

ngram-count

- 計算training data中每個出現的ngram的出現次數

3 表示 trigram

要計算count的Training data

詞典檔

```
ngram-count \  
-order 3 \  
-vocab Lexicon2003-72k.txt \  
-text ssp.train \  
-unk \  
-write ssp.n3.count
```

有-unk：把OOV當作一個<unk>符號處理
無-unk：遇到OOV就移除

輸出的 ngram count 檔名

Good-Turing Discounting and Katz Backoff

預設的
discounting跟
backoff方法

```
smooth=""  
ngram-count \  
  -order 3 \  
  $smooth \  
  -read ssp.n3.count \  
  -vocab Lexicon2003-72k.txt \  
  -lm ssp.n3.gt.lm
```

前一步 ngram count 所產生的檔

要輸出的LM檔

Compute perplexity

```
ngram \  
-lm ssp.n3.gt.lm \  
-pp1 ssp.test
```

LM檔

要算 perplexity 的 test data

輸出範
例：

```
file ssp.test: 89 sentences, 6979 words, 0 OOVs  
0 zeroprobs, logprob= -21066.5 pp1= 956.182 pp11= 1043.64
```

$$10^{-\frac{\text{logprob}}{\text{句數}+\text{字數}}}$$

$$10^{-\frac{\text{logprob}}{\text{字數}}}$$

Good-Turing Discounting Parameters

$$P_{Katz}(w_i|w_{i-1}) = \begin{cases} C(w_{i-1}w_i)/C(w_{i-1}) & r > k \\ d_r C(w_{i-1}w_i)/C(w_{i-1}) & k \geq r > 0 \\ \alpha(w_{i-1})P(w_i) & r = 0 \end{cases}$$

maxcount

mincount

$$\text{where } d_r = \frac{\frac{r^*}{r} - \frac{(k+1)n_{k+1}}{n_1}}{1 - \frac{(k+1)n_{k+1}}{n_1}} \text{ and } \alpha(w_{i-1}) = \frac{1 - \sum_{w_i:r>0} P_{Katz}(w_i|w_{i-1})}{1 - \sum_{w_i:r>0} P_{Katz}(w_i)}$$

自行指定mincount跟maxcount：

```
smooth=""  
  
ngram-count \  
  -order 3 \  
  $smooth \  
  -gt1min 3 -gt1max 7 -gt2min 3 -gt2max 7 -gt3min 3 -gt3max 7 \  
  -read ssp.n3.count \  
  -vocab Lexicon2003-72k.txt \  
  -lm ssp.n3.gt.lm
```



Modified Kneser-Ney Discounting

```
smooth="-kndiscount1 -kndiscount2 -kndiscount3"

ngram-count \
  -order 3 \
  $smooth \
  -read ssp.n3.count \
  -vocab Lexicon2003-72k.txt \
  -lm ssp.n3.kn.lm
```

- Take the **conditional probabilities** of bigrams (m -gram, $m=2$) for example:

$$P_{KN}(w_i|w_{i-1}) = \begin{cases} \frac{\max\{C[w_{i-1}, w_i] - D, 0\}}{C[w_{i-1}]} & \text{if } C[w_{i-1}, w_i] > 0 \\ \alpha(w_{i-1})P_{KN}(w_i) & \text{otherwise} \end{cases} \quad 0 \leq D \leq 1$$

$$1. P_{KN}(w_i) = C[\bullet w_i] / \sum_{w_j} C[\bullet w_j]$$

$C[\bullet w_i]$ is the unique words preceding w_i

2. **normalizing constant**

$$\alpha(w_{i-1}) = \frac{1 - \sum_{w_i: C[w_{i-1}w_i] > 0} \frac{\max\{C[w_{i-1}w_i] - D, 0\}}{C[w_{i-1}]}}{\sum_{w_i: C[w_{i-1}w_i] = 0} P_{KN}(w_i)}$$



Witten-Bell Discounting

```
smooth="-wbdisc1 -wbdisc2 -wbdisc3"  
  
ngram-count \  
  -order 3 \  
  $smooth \  
  -read ssp.n3.count \  
  -vocab Lexicon2003-72k.txt \  
  -lm ssp.n3.wb.lm
```



Witten-Bell Discounting

- **Key Concept—Things Seen Once:** Use the count of things you've seen once to help estimate the count of things you've never seen
- So we estimate the *total* probability mass of all the zero N -grams with the number of types divided by the number of tokens plus observed types:

$$\sum_{i:c_i=0} p_i^* = \frac{T}{N + T}$$

N : the number of tokens

T : observed types



Witten-Bell Discounting

- $T/(N+T)$ gives the total “probability of unseen N -grams”, we need to divide this up among all the zero N -grams
- We could just choose to divide it **equally**

$$Z = \sum_{i:c_i=0} 1$$

Z is the total number of N -grams with count zero

$$p_i^* = \frac{T}{Z(N+T)}$$



Witten-Bell Discounting

$$p_i^* = \frac{c_i}{N+T} \quad \text{if } (c_i > 0)$$

Alternatively, we can represent the smoothed counts directly as:

$$c_i^* = \begin{cases} \frac{T}{Z} \frac{N}{N+T}, & \text{if } c_i = 0 \\ c_i \frac{N}{N+T}, & \text{if } c_i > 0 \end{cases}$$



Witten-Bell Discounting

$$\begin{aligned} & \frac{T}{Z(N+T)} \cdot Z + \sum_{i:c_i>0} \frac{c_i}{N+T} \\ &= \frac{T}{N+T} + \frac{N}{N+T} \\ &= 1 \end{aligned}$$

$$\begin{aligned} & \frac{T}{Z} \frac{N}{N+T} \cdot Z + \sum_{i:c_i>0} c_i \frac{N}{N+T} \\ &= \frac{NT}{N+T} + \frac{N^2}{N+T} \\ &= \frac{N(N+T)}{N+T} \\ &= N \end{aligned}$$



Witten-Bell Discounting

- For bigram

$$\sum_{i:c(w_x w_i)=0} p^*(w_i | w_x) = \frac{T(w_x)}{N(w_x) + T(w_x)}$$

T : the number of bigram types, N : the number of bigram token

$$Z(w_x) = \sum_{i:c(w_x w_i)=0} 1$$

$$p^*(w_i | w_{i-1}) = \frac{T(w_{i-1})}{Z(w_{i-1})(N(w_{i-1}) + T(w_{i-1}))} \quad \text{if}(c_{w_{i-1}w_i} = 0)$$

$$\sum_{i:c(w_x w_i)>0} p^*(w_i | w_x) = \frac{c(w_x w_i)}{c(w_x) + T(w_x)}$$



Absolute Discounting

```
smooth="-cdiscount1 0.5 -cdiscount2 0.5 -cdiscount3 0.5"  
  
ngram-count \  
  -order 3 \  
  $smooth \  
  -read ssp.n3.count \  
  -vocab Lexicon2003-72k.txt \  
  -lm ssp.n3.abs.lm
```

Smoothing

absolute discounting

- Absolute discounting subtracting a fixed discount $D <= 1$ from each nonzero count

$$P_{abs}(w_i | w_{i-n+1} \dots w_{i-1}) \\ = \frac{C(w_{i-n+1} \dots w_{i-1}) - D}{\sum_{w_i} C(w_{i-n+1} \dots w_i)} + (1 - \lambda_{w_{i-n+1} \dots w_{i-1}}) P_{abs}(w_i | w_{i-n+2} \dots w_{i-1})$$



Ristad's Natural Discounting

```
smooth="-ndiscount1 -ndiscount2 -ndiscount3"

ngram-count \
  -order 3 \
  $smooth \
  -read ssp.n3.count \
  -vocab Lexicon2003-72k.txt \
  -lm ssp.n3.nd.lm
```

Eric Ristad's Natural Law of Succession -- The discount factor d is identical for all counts,

$$d = (n(n+1) + q(1- q)) / (n^2 + n + 2q)$$

where n is the total number events tokens, q is the number of observed event types. If q equals the vocabulary size no discounting is performed.