



Bayesian Replacement for Good-Turing

*Introduction to
MacKay (1994) “Hierarchical Dirichlet Language Model”*

Daichi Mochihashi
daichi.mochihashi@atr.jp

ATR SLT Dept.2 Regular Meeting

May 19, 2004

Modified for Über SVM 2004

Aug 3, 2004

Overview

- ⑥ N -gram smoothing is a crucial machinery in speech recognition and machine translation.

- ⑥ But N -gram parameters are so numerous, and there are not much data (e.g. MAD)



We can't make an exact prediction from such data.

But..

- ⑥ Taking our uncertainty about the parameters into the model, we **can** make a stable prediction. (This is called a Bayesian Method.)

- ⑥ We get a theoretically sound smoothing formula.

Introduction

By restricting ourselves to bigram for simplicity,

- ⑥ Empirical (Maximum Likelihood) estimate

$$\hat{p}_{i|j} \equiv \hat{p}(w_i|w_j) = \frac{f_{i|j}}{f_j} \quad (1)$$

($f_{i|j}, f_j$: frequency of $\langle w_j \rightarrow w_i \rangle$ and w_j)

- △ Probability 0 for unseen words after w_j
 - e.g. $p(\text{an}|\text{quite}) = 0$ simply if “quite an” accidentally did not appear in the training data.
- △ Some smoothing is required.

Existent smoothing

- ⑥ “Adding some” method
 - △ Adding some count to every N-gram
 - △ Interpreted as an interpolation between \hat{p} and uniform probability
 - △ Laplace smoothing, Lidstone’s law, Jeffreys-Perks law, ...
- ⑥ Good-Turing smoothing
 - △ uses “Bins of N-gram” (number of freq. 1 N-gram, ..)
 - △ only applicable when $f_{i|j} < \theta$.
 - △ shares several flaws also (next slide)

Problem of Existent smoothing

- ⑥ Uniform probability to unseen words
 - △ $p(\text{well}|\text{quite}) = p(\text{epistemological}|\text{quite})?$
- ⑥ ad hoc threshold (Good-Turing)
- ⑥ frequency of context is ignored.
 - △ probability $0.5 = 50/100 = 2/4?$
 - △ the more frequent the context is, the more stable \hat{p} should be (requires less smoothing)
 - △ But this information is discarded in the ordinary approach.

Example of Context Frequency

$$\begin{cases} \text{he} & \rightarrow 1000 \text{ times} \\ \text{he does} & \rightarrow 200 \text{ times} \end{cases} \therefore p(\text{does}|\text{he}) = \frac{200}{1000} = 0.2.$$

This estimate is very reliable.

$$\begin{cases} \text{alice} & \rightarrow 5 \text{ times} \\ \text{alice wandered} & \rightarrow 1 \text{ time} \end{cases} \therefore p(\text{wandered}|\text{alice}) = \frac{1}{5} = 0.2$$

$p(\text{does}|\text{he}) = p(\text{wandered}|\text{alice})$?

The latter may have been 0.3 or 0.1

⇓

Context frequency (1000 and 5) should be considered.

Bayesian Hierarchical model

- ⑥ Bigrams are governed by a probability table

$$q_{i|j} = p(w_i|w_j).$$

- ⑥ But we are not confident exactly what q is



Consider (infinite) possible q 's, and average them.

- ⑥ In fact,
 - △ Introducing “probability of probability table q ” and taking expectation of the prediction from each q
 - △ What governs above “probability of q ” is a hyperparameter α of the Dirichlet distribution.

Result of Bayesian Hierarchical model

- Resultant smoothing is a linear interpolation using empirical probability $\hat{p}_{i|j}$ and hyperparameter α

$$E[p(w_i|w_j)] = \frac{f_{i|j} + \alpha_i}{\sum_i (f_{i|j} + \alpha_i)} \quad (2)$$

$$= \frac{f_j}{f_j + \alpha_0} \cdot \hat{p}_{i|j} + \frac{\alpha_0}{f_j + \alpha_0} \cdot \bar{\alpha}_i \quad (3)$$

$$\text{where } \alpha_0 = \sum_k \alpha_k \text{ and } \bar{\alpha}_i = \frac{\alpha_i}{\alpha_0}$$

- also depends on the frequency f_j of context w_j
 - non-uniform interpolation like back-off (α_i)
- $\bar{\alpha}_i = p(w_i)$? (unigram) \rightarrow No.

Example of Bayesian model (2)

We assume $\alpha(\text{does}) = 1.5$, $\alpha(\text{wandered}) = 0.01$, $\alpha_0 = 10$.
Then because $f_{\text{he}} = 1000$ and $f_{\text{alice}} = 5$,

$$p(\text{does}|\text{he}) = \frac{1000}{1000 + 10} \cdot 0.2 + \frac{10}{1000 + 10} \cdot \frac{1.5}{10} = 0.1995.$$

$$p(\text{wandered}|\text{alice}) = \frac{5}{5 + 10} \cdot 0.2 + \frac{10}{5 + 10} \cdot \frac{0.01}{10} = 0.0673.$$

Very intuitive and different from any conventional methods that give equal probability 0.2 to both cases!

How to derive α ?

- Only what remains is a hyperparameter α .
- Most reasonable point estimate is the α which maximizes the probability of observed counts $F = \{f_{i|j}\}$ (called “evidence” in Bayesian statistics)

$$\begin{aligned} p(F|\boldsymbol{\alpha}) &= \int p(F|\mathbf{q})p(\mathbf{q}|\boldsymbol{\alpha})d\mathbf{q} \\ &= \prod_{j=1}^L \int_0^1 \cdots \int_0^1 \prod_{i=1}^L q_i^{f_{i|j}} \cdot \frac{\Gamma(\alpha_0)}{\prod_i \Gamma(\alpha_i)} \prod_{i=1}^L q_i^{\alpha_i-1} dq_1 \cdots dq_L \\ &= \prod_{j=1}^L \left[\frac{\Gamma(\alpha_0)}{\prod_i \Gamma(\alpha_i)} \cdot \frac{\prod_i \Gamma(f_{i|j} + \alpha_i)}{\Gamma(f_j + \alpha_0)} \right] \end{aligned} \quad (4)$$

How to derive α ? (2)

$$p(F|\boldsymbol{\alpha}) = \prod_{j=1}^L \left[\frac{\Gamma(\alpha_0)}{\prod_i \Gamma(\alpha_i)} \cdot \frac{\prod_i \Gamma(f_{i|j} + \alpha_i)}{\Gamma(f_j + \alpha_0)} \right] \quad (5)$$

- ⑥ This evidence (likelihood) is convex in α , and has a global maximum
- ⑥ Maximum of α can be obtained by an iterative optimization (MacKay 1994, Minka 2003)
 - △ 77 lines of MATLAB code last week
 - △ Taking a few hours to calculate (for small data).

Minka's Exact Method

- ⑥ Minka (2003) "Estimating a Dirichlet distribution"

$$\alpha_i^{(t+1)} = \alpha_i^{(t)} \cdot \frac{\sum_j \Psi(f_{i|j} + \alpha_j) - \Psi(\alpha_j)}{\sum_j \Psi(f_j + \sum_k \alpha_k) - \Psi(\sum_k \alpha_k)}. \quad (6)$$

where $\Psi(x) = \frac{d}{dx} \log \Gamma(x)$

- ⑥ For bigrams, it takes about 30 minutes on P4 2GHz (dependent on data)
- ⑥ MATLAB code available on request.

MacKay's Approximation

- MacKay (1994) approximates $\Psi(x)$ by expansion:

$$K(\alpha) = \sum_{j=1}^L \log \frac{f_j + \alpha}{\alpha} + \frac{1}{2} \sum_{j=1}^L \frac{f_j}{\alpha(f_j + \alpha)} \quad (7)$$

$V(i)$ = number of contexts before word i

$G(i), H(i)$ = sufficient statistics from the n-gram table

Then, (no proof is given!)

$$\alpha'_i = 2V(i) / \left[K(\alpha_i) - G(i) + \sqrt{(K(\alpha_i) - G(i))^2 + 4H(i)V(i)} \right] \quad (8)$$

- Consistent to the exact answer (while difference of performance needs to be examined.)

Is it perfect?

- ⑥ Yes, almost perfect.
- ⑥ But, in general history h , the formula is:

$$E[p(w_i|h)] = \frac{f_h}{f_h + \alpha_0} \cdot \hat{p}_{i|h} + \frac{\alpha_0}{f_h + \alpha_0} \cdot \bar{\alpha}_i \quad (9)$$

- ⑥ It uses a MLE (no smoothing) for history frequency f_h !
- ⑥ We must estimate f_h *recursively* also by a Bayesian method. (current work)
 - △ Due to the point estimate of hyperparameter and the assumption of uniform hyperprior.

Gamma function

- ⑥ Gamma function $\Gamma(x)$ is a continuous analogue of the factorial
- ⑥ $\Gamma(x) = (x - 1)!$ if x is an integer
- ⑥ $\Gamma(x)$ is defined by: $\Gamma(x) = \int_0^{\infty} \exp(-\theta)\theta^{x-1}d\theta.$