

# CS276 Programming Assignment 2

*Raj Bandyopadhyay, Rafael Guerrero (rajb2, rplatero)*

*May 1st, 2013*

## Table of Contents

- [Introduction](#)
- [Model Generation](#)
  - [Language Model](#)
  - [Noisy-Channel Model](#)
- [Candidate Generation](#)
  - [Optimizations](#)
  - [Handling spaces](#)
- [Candidate Scoring](#)
  - [1. Uniform Noisy Channel](#)
  - [2. Empirical Noisy Channel](#)
- [Query Scoring](#)
  - [Combining candidates for multiple words](#)
- [Extra Credit \(Stupid Backoff\)](#)
- [Parameter Tuning](#)
  - [1\) lambda for interpolation](#)
  - [2\) uniform edit cost](#)
  - [3\) Jaccard coefficient cutoff](#)

## Introduction

In this assignment we have implemented a spelling corrector in Python, using a probabilistic approach to model the noisy channel. This report describes the different design options and trade-offs we have considered during the implementation of every module of the system i.e. model generation, candidate generation, filtering and scoring; as well as the results obtained.

## Model Generation

### Language Model

We scan the corpus files and generate these dictionaries:

- word\_language\_model: how probable is each word of the corpus
- biword\_language\_model: how probable is each biword of the corpus
- word\_counter: how frequent is each word of the corpus
- biword\_counter: how frequent is each biword of the corpus

We also generate word indexes, bi-character and tri-character indexes that will be used later on candidate generation.

### Noisy-Channel Model

For the noisy-channel model we have followed the approach described by *Kernighan, Church and*

The model is represented by four confusion matrices (deletion, substitution, transposition and insertion), one unichar counter and one bichar counter. The entries of the confusion matrices are defined as:

- $\text{del}[(x,y)] = \text{count}(xy \text{ typed as } x)$
- $\text{sub}[(x,y)] = \text{count}(y \text{ typed as } x)$
- $\text{tra}[(x,y)] = \text{count}(xy \text{ typed as } yx)$
- $\text{ins}[(x,y)] = \text{count}(x \text{ typed as } xy)$

## Candidate Generation

The approach we have followed for correcting a potential misspelled query phrase is as follows:

1. *Split the query phrase into bi-words and identify bi-words with one potential misspelled word* : A bi-word is considered as potentially misspelled when the bi-word language model shows no hits for it. Likewise a word is considered potentially misspelled when the word index shows no hits for it
2. *Generate candidates for the potential misspelled words using bigram and trigram character indices*: We are using an n-gram index approach; this guarantees that all candidates being considered for a given word, are also part of the dictionary, which reduces the size of the candidate set. Specifically:
  - When the word is 10 characters or less we use a bigram index.
  - When the word is greater than 10 characters, we use a trigram index.
  - If the word occurs in the dictionary, it is always included as a candidate

## Optimizations

While the candidates are been selected using the bi-char or tri-char index we applied two more heuristics to prune the final set of candidates:

1. Assume that the first character of each query word is typed correctly
2. Jaccard Coefficient with cutoff threshold: If the Jaccard Coefficient between the word and the candidate is smaller or equal than the cutoff threshold (0.2 by default) then the candidate is discarded. If the potential misspelled word has more than 10 characters then we force the cutoff threshold to be at least 0.5. This is because for smaller words, a single edit can influence Jaccard distance drastically compared to longer words.
3. Edit Distance: we discard any candidate at edit distance greater than 2.

## Handling spaces

We assume that every word has a potential missing space in it, and handle it as follows:

- Generate each possible split of the word.
- For each split
  - If one of the two parts occurs in the dictionary, add it as a candidate and don't generate any further candidates for that part. Generate candidates as before for the other part, but restrict them to edit distance 1 since the space counts as one edit.
  - If neither of the two parts occurs in the dictionary, skip candidate generation for this split and don't add either of them as candidates.
  - If both parts of the split occur in the dictionary, just add them as candidates and don't generate any more candidates.

## Candidate Scoring

After generating a reduced list of potential candidates for each word, we rank them using a probabilistic approach for the Noisy Channel model. When we rank the bi-gram (R,Q) where Q is a candidate for R, we do the following calculation

$$\text{rank}(Q,R) = P(Q|R) = P(R|Q) * P(Q)$$

$P(Q)$  comes from our language model, while  $P(R|Q)$  is modeled using one of the two options used to model the noisy channel  $P(w_2|w_1)$ :

### 1. Uniform Noisy Channel

We model the cost of the edit operations as:

$$\begin{aligned} P(R|Q) &= (\text{cost})^d && \text{when } Q \neq R \\ &= 0.9 \text{ (default)} && \text{when } Q == R \end{aligned}$$

Where  $d$  is the edit distance between  $R$  and  $Q$ , and  $\text{cost}$  is the uniform edit cost (set to 0.01, which was discovered through experimentation).

### 2. Empirical Noisy Channel

The cost of each edit operation (deletion, substitution, transposition and insertion) has been calculated following the approach described by *Kernighan et al.*

$$\begin{aligned} P(R|Q) &= P_{\text{oper}}(R|Q) && \text{when } Q \neq R \\ &= 0.9 \text{ (default)} && \text{when } Q == R \end{aligned}$$

where:

- $\text{edit\_distance}(R|Q) = 1$
- $\text{one\_edit\_distance\_operation}(R|Q) = \text{'oper'}$
- $P_{\text{oper}}(R|Q) = \text{cost of changing } Q \text{ into } R \text{ as defined by the Noisy Channel Model}$

There may be candidates  $Q$  that are at edit distance 2 from  $R$ . This means that the noisy channel transformation from  $Q$  into  $R$  is:  $Q \rightarrow \text{edit\_operation\_1} \rightarrow R' \rightarrow \text{edit\_operation\_2} \rightarrow R$ . We have modelled this 2 edit-transformation as:

$$P(R|Q) = P_{\text{oper}}(R|Q) * \text{uniform\_edit\_operation\_cost}$$

where:

- $P_{\text{oper}}(R|Q)$  models the cost of the first edit operation, using the empirical noisy model
- $\text{uniform\_edit\_operation\_cost}$  models the cost of the second edit operation

Using a constant for the second edit operation allowed us to avoid the computational burden of calculating the exact edit operation at the expenses of a less precise estimated cost.

## Query Scoring

When best candidates have been selected we run a cross-product for the different bigrams involved,

generating a set of potential query candidates. We score each query phrase using the pre-computed Language Model e.g. if  $(w_1, w_2, w_3)$  is a query phrase candidate we score it using:

$$P(w_1, w_2, w_3) = P(w_1) * P(w_2 | w_1) * P(w_3 | w_2)$$

When we hit a biword that is not part of our corpus we use n-gram interpolation with  $\lambda = 0.2$ .

## Combining candidates for multiple words

We limit the total number of candidates for a query phrase to  $M=500$ . and the number of candidates for a single word to  $M/N$ , where  $N$  is the number of words in the phrase. The candidates for each word in the query, as well as the entire query at the end are ranked using the scoring technique described above and the designated number of candidates is returned.

## Extra Credit (Stupid Backoff)

In stupid backoff, we use just the weighted prior probability of a single word instead of the bi-word probability for bi-words that don't occur in the dictionary. i.e.

$$P(w_2 | w_1) = \text{count}(w_1, w_2) / \text{count}(w_1) \quad \text{if } (w_1, w_2) \text{ exists in the language model} \\ = \alpha * \text{count}(w_2) / T \quad \text{otherwise}$$

As suggested in literature,  $\alpha$  is set to 0.4.

## Parameter Tuning

We attempted to manually tune many parameters: e.g. the  $\lambda$  for interpolation, the uniform edit cost, the prior probability  $P(R|R)$  and so on. The following results show the numbers we obtained for parameters that made a significant difference through tuning.

### 1) $\lambda$ for interpolation

An overly small value of  $\lambda$  drastically reduces accuracy, but beyond 0.2, it makes little difference. Hence, we stuck with 0.2.

<b><math>\lambda</math> for interpolation</b>	<b>Accuracy (% correct on dev set)</b>
0.1	50.8
0.2	78.9
0.4	78.6

## 2) uniform edit cost

0.01 was the optimal value we found for this cost.

uniform edit cost	Accuracy (% correct on dev set)
0.003	77.3
0.01	78.9
0.03	77.8
0.1	75.8

## 3) Jaccard coefficient cutoff

There is a tradeoff between accuracy and runtime. For small Jaccard cutoffs, the accuracy is high, but runtime increases due to large numbers of candidates. For large cutoffs, too few candidates are generated. We decided on 0.2 as a default cutoff.

jaccard cutoff	Accuracy (% correct on dev set)	Runtime (secs)
0.1	78.9	350
0.2	78.9	257
0.3	78	222
0.4	74.7	206