



Department of Informatics, Systems and Communication

---

# An HMM-based approach for misspelling correction

---

PROBABILISTIC MODELES FOR DECISION

Giorgia Adorni  
806787

Elia Cereda  
807539

Nassim Habbash  
808292

**Academic Year 2018-2019**

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Problem Formulation</b>	<b>5</b>
2.1	Purpose . . . . .	5
2.2	Design Choices . . . . .	5
2.3	Software . . . . .	6
<b>3</b>	<b>Dataset</b>	<b>7</b>
3.1	Dataset Acquisition . . . . .	7
3.1.1	Error Model . . . . .	7
3.1.2	Language Model . . . . .	8
3.1.3	Transition Model . . . . .	8
3.1.4	Perturbated Dataset . . . . .	8
<b>4</b>	<b>Models</b>	<b>10</b>
4.1	Noisy Channel Model . . . . .	10
4.2	Most Likely State Sequence . . . . .	12
4.3	Hidden Markov Model . . . . .	13
<b>5</b>	<b>Experimental Results</b>	<b>14</b>
<b>6</b>	<b>Conclusions</b>	<b>16</b>
6.1	Future Works . . . . .	16
	<b>References</b>	<b>17</b>

# List of Figures

4.1	Noisy channel . . . . .	11
4.2	Trellis example . . . . .	13

# List of Tables

5.1	Typos performance evaluation . . . . .	14
5.2	Sentences performance evaluation . . . . .	15

# Chapter 1

## Introduction

In recent decades, technology has had a strong impact on everyone's life. It plays an important role in the communication process, simplifying different activities for both individuals and businesses.

Nowadays we have advanced communication technology tools available, such as smartphones, tablets and computers that have simplified the way humans communicate.

Companies can write an e-mail and deliver it to all their consumers in a few minutes. People can message their friends at every moment and share an interest with new friends from different countries.

This advancement in communication technology has made it necessary to equip our technological tools with a series of programs and software that control and correct automatically the misspelt words typed.

In this project, we propose and evaluate an automatic spelling correction algorithm, modelling the typing process as an *Hidden Markov Model* (HMM).

**FIXME:** keyboard?

# Chapter 2

## Problem Formulation

**FIXME:** Given an observation sequence, typically a phrase, and the model parameters, we are interested in detect and correct errors estimating the optimal state sequence.

In our HMM, the hidden states represent the intended words and the observations are the typed words.

The initial state probabilities  $\pi$  are given by the word frequencies and the state transitions  $A_{ij}$ , that is the probability of one word given its predecessor obtained from ??**FIXME**

The emission probabilities,  $B_{ij}$ , represent the probability of a typed word given the intended one, and depend on the confusion probabilities.

**FIXME:** We also detect the non-word error?

### 2.1 Purpose

### 2.2 Design Choices

We have chosen the English language for different reasons. First of all, the great majority of material in literature deals with the problem in question in the English language. Moreover it is a simple language, both from a grammatical and a lexical point of view: it lacks in certain symbols, like accents and apostrophes, and genres. Furthermore all punctuation and special character symbols were not considered, but only letters and sometimes numbers.

We considered that a typed word only depends on the previous one, being in the framework of Markov chains. If we know the probability of a word given its predecessor, the frequency of each word, and the probability to type word  $x$  when

word  $y$  is intended, we have all the necessary ingredients to use Hidden Markov Models.

## 2.3 Software

**FIXME** We have developed the project in **Python**.  
The interface is a native macOS application written in **Swift**.

# Chapter 3

## Dataset

### 3.1 Dataset Acquisition

#### 3.1.1 Error Model

This dataset was collected from the following resources <sup>1 2 3 4</sup>:

- BIRKBECK: contains 36 133 misspellings of 6136 words, taken from the native-speaker section (British and American) of the Birkbeck spelling error corpus.
- HOLBROOK: contains 1791 misspellings of 1200 words, taken from the book "English for the Rejected" by David Holbrook (Cambridge University Press - 1964).
- ASPELL: contains 531 misspellings of 450 words, taken from one assembled by Atkinson for testing the GNU Aspell spellchecker.
- WIKIPEDIA: contains 2455 misspellings of 1922 words, taken from the misspellings made by Wikipedia editors.
- URBAN-DICTIONARY-VARIANTS: contains 716 variant spellings, taken from the text scraped from Urban Dictionary (in UK English).
- SPELL-SET: contains 670 typos.
- TWEET-TYPO: contains 39 172 typos, taken from Twitter.

---

<sup>1</sup><https://www.dcs.bbk.ac.uk/~ROGER/corpora.html>

<sup>2</sup><https://www.kaggle.com/rtatman/spelling-variation-on-urban-dictionary>

<sup>3</sup><https://www.kaggle.com/bittlingmayer/spelling>

<sup>4</sup><http://luululu.com/tweet>



All the datasets are joined and cleaned. After that, it contains 79 677 row, in each of which there is the typo and the correct word. The dataset obtained is divided into two corpora: 80% is used as a train set (63 679 row) and 20% is used as a test set (15 998 row).

### 3.1.2 Language Model

The words frequency dataset includes wordlists derived from the [Google's ngram corpora](#). In particular we use the dataset `frequency-alpha-gcide.txt`, a smaller version of the original dataset, cleaned up and limited to only the top 65 538 words.

Each row of the corpus contains the ranking of the word, the word itself, the count of the occurrences of the word, a percentage of how often each word was being used and a cumulative percentage.

With the following dataset we found some problems due to the lack of proper names, city names, countries, brands etc.

The language model dataset is a concatenation of public domain book excerpts from [Project Gutenberg](#) and lists of most frequent words from [Wiktionary](#) and the [British National Corpus](#).

This dataset contains about a million word. It is use to estimate the probability of a word by counting the number of times each word appears in this dataset. **FIXME:** spiegare per cosa usiamo questo dataset To do that, the text was breaks into words, then a variable holds a counter of how often each word appears. Than the probability of each word, based on this Counter was estimated

### 3.1.3 Transition Model

Abbiamo cambiato con il signore degli anelli

### 3.1.4 Perturbated Dataset

In order to evaluate our algorithm on whole sentences, we create a new perturbed dataset starting from the dataset `big` described in the section 3.1.3.

The disturbance introduced presents an error dependent on the error model previously presented, with the difference that it is created starting from the typos belonging to the test dataset.

Three different texts, of approximately 29 000 sentences each, have been generated, each of which has a percentage of errors in the text of 10-15-20% respectively.

**FIXME:** Estimates for the frequency of spelling errors in human-typed text vary from 1-2% for carefully retyping already printed text to 10-15% for web queries.

We implemented a perturbation algorithm, which for each line of our input file generates a new perturbed string.

The input text is perturbed accordingly the following steps:

1. The probability that a word has an edit is computed by multiplying the value of  $p$ , coming from the error model, by the percentage of errors desired (10-15-20%).
2. For each word of length  $n$ , the number of edits to be introduced  $x$  is calculated according to the relation  $x \sim \text{Bin}(n, p)$  **FIXME**: why
3. The characters to be changed within each word are chosen randomly.
4. **FIXME**: The typo of edit to be introduced within each word is chosen according to the probability of each type of error. The disturbance goes to alter the letters designated not in a random manner. In fact, we use four different probabilities to define whether a letter will be deleted from the index in question, if a new letter will be inserted after the actual character, or if the current character will be replaced with one of the possible letters according to the error model probability, or if the current character will be swapped with the next or the previous one.

Swap errors are only introduced if there are no further changes in the word. Cases of elimination of a whole word are excluded, as these would heavily influence the evaluation metrics as they are inconsistent with our model.

# Chapter 4

## Models

We will consider two optimality criteria. The first one chooses the states that are individually most likely and maximizes the expected number of correct individual states. The second criterion estimates the most likely state sequence, or *trellis path*. The algorithm used to implement these criteria is the **Viterbi** algorithm.

### 4.1 Noisy Channel Model

**Non-word errors** are detected by looking for any word not found in a dictionary. To correct them we first generate **candidates**, according to a distance given as a parameter to the model (`edit_distance`), that are real words with a similar letter sequence to the error.

The intuition of the noisy channel model is to treat the misspelled word as if a correctly spelled word had been “distorted” by being passed through a noisy communication channel. This channel introduces “noise” in the form of substitutions or other changes to the letters, making it hard to recognise the “true” word.

We see an observation  $x$  (a misspelled word) and we want to find the word  $w$  that generated this misspelled word (the intended word). Out of all possible words in the vocabulary  $V$  we want to find the word  $\hat{w}$  such that  $P(\hat{w}|x)$  is highest among all the candidates  $w$

$$\hat{w} = \arg \max_{w \in V} P(w|x). \quad (4.1)$$

This noisy channel model is, therefore, a kind of Bayesian inference, in which it is possible to transform the equation 4.1 into a set of other probabilities

$$\hat{w} = \arg \max_{w \in V} \frac{P(x|w)P(w)}{P(x)}. \quad (4.2)$$

Since  $P(x)$  doesn't change for each word because we are always asking about the most likely word for the same observed error  $x$ , we can conveniently simplify the equation 4.2 by dropping the denominator

$$\hat{w} = \arg \max_{w \in V} P(x|w)P(w). \quad (4.3)$$

**FIXME:** We apply the noisy channel approach to correcting non-word spelling errors by taking any word not in our spell dictionary, for example *adventhre*, generating a list of candidate words like *adventure*, *adventurer*, *adventured*, ranking them according to 4.3, and picking the highest-ranked one, *adventure*. **FIXME:** We choose the most likely candidate, the one with the highest probability.

The two components of the equation are, respectively,  $P(x|w)$  the **channel model** and  $P(w)$  the prior probability of a hidden word (a candidate). The prior probability of each correction is the language model probability of the word  $w$  in context, which is computed using a **FIXME unigram** language model. The likelihood is estimated just using the number of times that the a letter  $i$  was substituted for the letter  $j$  in the large corpus of errors **FIXME QUALE=?** To compute the probability for each edit we used a confusion matrix that contains counts of errors. **FIXME** How?

Each candidate is scored by  $P(\text{candidate})P(\text{typo}|\text{candidate})$  and then normalised by the sum of the scores for all proposed candidates.

**FIXME:** (How is estimated the prior? How are computed the conditional probabilities?)

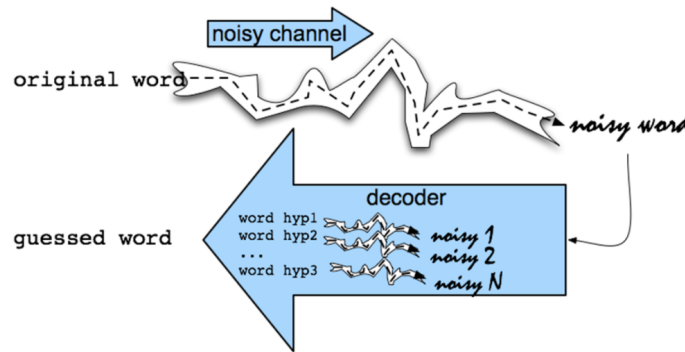


Figure 4.1: Noisy channel

**FIXME:** When for a input typo we do not have any candidate ??

**FIXME:** Real word spelling error detection is a much more difficult task, since any word in the input text could be an error. So we don't...

To find this list of candidates we use a minimum edit distance algorithm in the extended version with transposition called **Damerau-Levenshtein** edit distance ??.

The following type of errors were considered in our model:

- DIGITS DELETIONS
- DIGITS INSERTIONS
- SUBSTITUTIONS OF DIGITS: this is the probability to type a character  $i$  when the character  $j$  was intended ( $P(i|j)$ ). This probability was determined experimentally. **FIXME**: how? from where?
- TRANSPOSITIONS OF ADJACENT DIGITS: two letters are swapped.

## 4.2 Most Likely State Sequence

The Viterbi algorithm calculates the most probable sequence of hidden states, the words intended. The Viterbi algorithm is a probabilistic extension of minimum edit distance. Instead of computing the “minimum edit distance” between two strings, Viterbi computes the “maximum probability alignment” of one string with another.

The initial probability of being in a state  $i$ ,  $\pi_i$ , in our case the probability of intend a word  $i$ , and the transition probabilities  $A_{ij}$ , or the transition from the word  $i$  to the next word  $j$ , are given. Since we have observed the output  $y_1, y_2, \dots, y_t$ , that is the sentence written with typos, it is possible to compute the most likely state sequence  $x_1, x_2, \dots, x_t$ , the sentence intended, starting from the following expression:

$$\begin{aligned}
 V_{1,t+1} &= P(x_1, \dots, x_t, x_{t+1}, y_1, \dots, y_t, y_{t+1}) = \\
 &= \arg \max_{x_{1:t}} p(x_1, \dots, x_t | y_1, \dots, y_t) = \\
 &= \alpha \cdot p(y_{t+1} | x_{t+1}) \cdot \max_{x_t} \left( p(x_{t+1} | x_t) \max p(x_1, \dots, x_t | y_1, \dots, y_t) \right)
 \end{aligned} \tag{4.4}$$

The initial state probabilities  $\pi$  are actually the word frequencies (? We don’t estimate it in a proper way), the state transition probabilities are given by the probability of a word given its predecessor, and the emission probabilities are the probabilities to type word  $i$  when word  $j$  was intended.

In our implementation, we construct the *trellis* choosing the locally/global best state. As we can see in the picture below, we display an example of the trellis drawing only the best predecessor of each state.

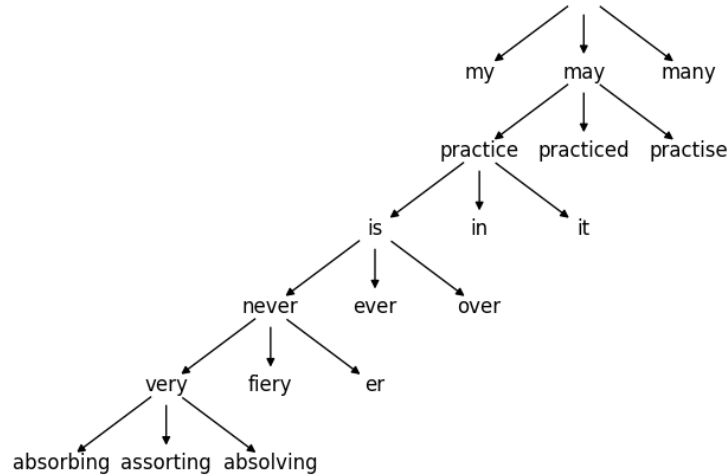


Figure 4.2: Trellis example

In this example, we observed the sentence *my practice iu never iery absorcing* and our algorithm return the most likely state sequence *may practice is never very absorbing*. In this case, the algorithm partially fails, because the intended sentence was *my practice is never very absorbing*.

**FIXME:** (pi la probabilità iniziale del prima stato non lo abbiamo, non lo facciamo perchè dal nostro dataset non abbiamo sempre la prima parola. . . )

We decide to implement the **Viterbi** based algorithm instead the Forward-Backward algorithm relying on the experiments carried out in the literature. The HMM-Based Error Correction Mechanism for Five-Key Chording Keyboards article [1] explains that the Forward-Backward algorithm estimates the most likely state for each observation, but the resulting state sequence may not be a valid succession of words in natural language (or a very unlikely word sequence) and produce inferior results.

### 4.3 Hidden Markov Model

**FIXME:** Hmm approach **FIXME:** description of hmm

**FIXME:** our application ...

# Chapter 5

## Experimental Results

### **FIXME** comparison **FIXME** performance

The fact that words are generally more frequent than their misspellings can be used in candidate suggestion, by building a set of words and spelling variations that have similar contexts, sorting by frequency, treating the most frequent variant as the source, and learning an error model from the difference

**FIXME** Mostrare il nostro modello con p e la tabella delle probabilità We detected some problems with our dataset, in particular it lacks of plural forms and other things

In order to avoid this problem, we decided to try a new approach that use lemmatisation (stemmisationj .....

Analysis of spelling error data has shown that the majority of spelling errors consist of a single-letter change and so we often make the simplifying assumption that these candidates have an edit distance of 1 from the error word.

Evaluating spell correction algorithms is generally done by holding out a training, development and test set from lists of errors like those on the Norvig and Mitton sites mentioned above.

	Time (sec)	Accuracy Top1	Accuracy Top3	Accuracy Top5
Train	20 985.70	41,38%	57,28 %	61,60 %
Test	5270.06	41,54%	57,18 %	61,15 %

Table 5.1: Typos performance evaluation

---

	Time (sec)	Accuracy	Initial Error	Accuracy Top5
Test	2695.01	47,99%	15,01%	61,15 %

---

Table 5.2: Sentences performance evaluation



# Chapter 6

## Conclusions

it is important to use larger language models than unigrams

For this reason modern systems often use much larger dictionaries automatically derived from very large lists of unigrams like the Google N-gram corpus

### 6.1 Future Works

- update the various dataset: in particular we can use a new frequency dataset consistent with the language model one.
- Compute the initial probabilities  $\pi$  from the language model
- Consider additional type of error in the model
- improving the noisy channel error model to consider neighbouring characters - Use other algorithms, such as the forward-backward (smoothing) instead Viterbi
- Optimise the code, in particular the lemmatisation/stemming
- Use recurrent neural network instead the hidden markov models.
- We can also improve the performance of the noisy channel model by changing how the prior and the likelihood are combined. In the standard model they are just multiplied together.

# Bibliography

- [1] Adrian Tarniceriu, Bixio Rimoldi, and Pierre Dillenbourg. “HMM-based error correction mechanism for five-key chording keyboards”. In: *2015 International Symposium on Signals, Circuits and Systems (ISSCS)*. IEEE. 2015, pp. 1–4 (cit. on p. 13).
- [2] Yanen Li, Huizhong Duan, and ChengXiang Zhai. “A generalized hidden markov model with discriminative training for query spelling correction”. In: *Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval*. ACM. 2012, pp. 611–620.
- [3] Grzegorz Szymanski and Zygmunt Ciota. “Hidden Markov models suitable for text generation”. In: *WSEAS International Conference on Signal, Speech and Image Processing (WSEAS ICOSIP 2002)*, pp. 3081–3084.
- [4] James H Martin and Daniel Jurafsky. *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition*. Pearson/Prentice Hall Upper Saddle River, 2009.