

report

by Haonan Li

Submission date: 06-Apr-2018 01:26PM (UTC+1000)

Submission ID: 94 1928805

File name: report.pdf (947.66K)

Word count: 2330

Character count: 11731

Spelling Correction Methods Evaluation

1 Introduction

Spelling correction is an old NLP task. From traditional non-neural methods (Hall and Dowling, 1980a; Gadd, 1988; Brill and Moore, 2000), to neural methods presents in recent years (Han and Baldwin, 2011; Eger et al., 2016; Silfverberg et al., 2016). The performance of spelling correction methods have been better and better and relatively mature.

In this paper, we investigate some old but really effective and impactful algorithms' performance on spelling correction task. Use a newly UrbanDictionary¹ dataset. The algorithm been evaluated include Soundex, N-gram, Edit-distance and Editex.

2 Methods

In this section, we simply introduce the algorithms we evaluated in our paper.

Soundex was developed by Odell and Russell, and patented in 1918 (Hall and Dowling, 1980b). Uses codes based on the sound of each letter to translate a word into an at most 4 character's string. We call it Soundex code in this paper. Figure 1 shows the letter translation table and Algorithm 1 describes the 4 steps processing.

Algorithm 1 Soundex

Input: One string

Output: Soundex code of the string

- 1: Except for first character, translate string characters according to table.
- 2: Eliminate any adjacent repetitions of codes.
- 3: Eliminate all occurrences of code 0.
- 4: Return the first four character.

N-Gram methods are string distance methods based on n-gram counts, where a n-gram

¹A crowd-compiled dictionary of informal words and slang with over 7 million entries. <http://urbandictionary.com>

aehiouwy	-->	0 (vowels)
bpfv	-->	1 (labials)
cgjkqsz	-->	2 (misc: velars, fricatives, etc.)
dt	-->	3 (dentals)
l	-->	4 (lateral)
mn	-->	5 (nasals)
r	-->	6 (rhotic)

Figure 1: Soundex letter translation table

of string s is any substring of s of some fixed length. This algorithm defines similarity of words by calculating n-gram distance, which was proposed by (Ukkonen, 1992) defines as:

$$|N_s| + |N_t| - 2|N_s N_t| \quad (1)$$

where N_s is the set of n-gram in string s .

Edit distance presented by (Levenshtein, 1965), also know as Levenshtein distance, which is defined as the minimum number of elementary edit operations needed to transform one string into another. Algorithm 2 gives the main part of one kind of edit distance (global edit distance) with particular distance calculate strategy.

Algorithm 2 Global edit distance

Input: Two strings s and t

Output: Edit distance between s and t

- 1: $ls = \text{length}(s)$, $lt = \text{length}(t)$
- 2: Initialize $\text{edit}[0][0] = 0$
- 3: Initialize $\text{edit}[i][0] = i$, $\text{edit}[0][j] = j$
- 4: **for** $i = 0$ to ls **do**
- 5: **for** $j = 0$ to lt **do**
- 6: $\text{edit}[i][j] = \min(\text{edit}[i-1][j] + 1,$
- 7: $\text{edit}[i][j-1] + 1,$
- 8: $\text{edit}[i-1][j-1] + \text{equal}(s[i-1], t[j-1]))$
- 9: **end for**
- 10: **end for**
- 11: **return** $\text{edit}[ls][lt]$

aeiouy --> 1 dt --> 4 fpv--> 7
 bp --> 2 lr --> 5 sxz--> 8
 ckq--> 3 gj--> 6 csz --> 9

Figure 2: Editex letter groups

Item	#Word
Testset size	716
Dictionary size	393954
Misspelled in Dictionary	175
Correct not in Dictionary	122

Table 1: Dataset statistics (# stands for the number of words)

Editex presented by (Zobel and Dart, 1996), is a method combines phonetic matching and edit distance. In this method, alphabets also be grouped similar with Soundex, shows in Figure 2. Then, edit distance will be calculated with the grouping info, the distance between two letters in the same group defines 1 while the different groups defines 2. This scheme makes distance between similar phonetic words closer.²


3 Experiments

3.1 Data


For our experiments we use a particular dataset: a sub-sample of actual data posted to UrbanDictionary, in which a number of head-words taken from have been automatically identified as being misspelled (Saphra, 2016). This dataset simulates the spelling error data in the real world as much as possible so that is suitable for our task. Table 1 shows the statistics of our dataset.

3.2 Setting

Soundex Calculate the Soundex code for every word and then matched with global edit distance.

N-Gram We evaluate the N-Gram algorithm for n in range 1 to 9.  a particular n , we first pad $(n-1)$ “#” in the front and end of every word. This guarantee the building of n -gram set. For example, 5-gram set for word “he” defined as:

$$\{\####h, \###he, \##he\#, \#he\#\#, he\###, e#####\} \quad (2)$$

Edit Distance There are two kinds of edit distance algorithm, **local edit distance (LED)**  **global edit distance (GED)**. In this paper, we

²Because of the space limitation, we can not introduce detailed algorithm, it is illustrated in our code.

N	Predicted	Right	Precision	Recall
1	7150	183	2.56	25.56
2	1484	151	10.19	21.09
3	1429	149	10.43	20.81
4	1426	148	10.38	20.67

Table 2: N-gram algorithm results

Scheme	Predicted	Right	Precision	Recall
GED-1	5528	253	4.57	35.34
GED-2	2497	204	8.16	28.49
LED	727774	133	0.02	18.58

Table 3: Edit distance algorithm results

implement both of them and evaluate them. For global edit distance, we have two distance calculate scheme: 1) (+1) for indel and mismatch and nothing to do for match; 2) (+1) for indel and mismatch and (-1) for match. The different between them will discuss in Section 4. For local distance algorithm, we use (-1) for indel and mismatch and (+1) for match, and always assign 0 if 0 is better.

Editex We calculate the Editex follows (Zobel and Dart, 1996) setting.

4 Results & Evaluation

4.1 Evaluation Method

For all experiments in this paper. We keep all parallel best results as matched corrections for a misspelling. For example, the Soundex code for word “adn” is “a35”, while for “attain”, “attainabilities”, “adamas”, “atom” etc. are the same. No matter how many words in dictionary have the same code with “adn”, all of them seem equally. Although this may result in a bad performance on precision when the predicted set is large.³

4.2 N-Gram Algorithm Evaluation

We employ the original results to evaluate since the purpose of this paper is to illustrates the characteristics of different string matching algorithms through a spelling correction task, rather than find a best algorithm with parameters to correct the misspelled datasets.⁴ We use recall

³We also do more experiments on parameter optimization. For example, randomly selecting three best matches as predicted words, which may increase precision for Soundex and local edit distance and impact recall for all algorithm. Due to space limitation, we could not list it here.

⁴The dataset we used contains only 716 words, even find best algorithm with parameters, it might be some kind of overfitting of the particular datasets.

Method	Predicted	Right	Precision	Recall
Soundex	495146	436	0.09	60.89
Local Edit Distance	727774	133	0.02	18.58
Global Edit Distance	2497	204	8.16	28.49
N-Gram (N=2)	1484	151	10.18	21.09
Editex	2830	230	8.13	32.12

Table 4: Full Scale results

and precision as measured of algorithm performance follows (Raghavan et al., 1989).

The results of n-gram algorithm shows in Table 2, in which we see that when $N = 1$, the number of predicted words is large and recall is the highest. This illustrates that around 25.56% of the misspelled words are character shifted without deletions and insertion. the best matches for $N = 1$ means two words have the same consist of characters but maybe different order. However, the precision of $N = 1$ is really low which makes it not a good parameter for the algorithm. $N = 2$ and $N = 3$ gives more balanced and significant results. For $N \geq 4$, the results are all the same. The reason might be that padding $n - 1$ “#” strategy for n-gram result in the related similarity tend to be stable when n is large. Needs more Analysis

4.3 Edit Distance Algorithm Evaluation

Table 3 shows the results of edit distance algorithm. Which indicates that although the first two lines are both global edit distance algorithm, different distance calculate strategy leads to divergent results. Both strategies with (+1) for indels and mismatches and the difference is that first strategy do nothing for matches while the second (-1). This difference always makes the first strategy matches more result. For example, calculate the edit distance between “against” with “against” and “agist”. The first strategy’s result is the all equals 1 because “against” add one character and “agist” delete one character compare with “against”. But the second strategy calculate -5 for “against” and -4 for “agist”, which means “against” have more similar elements (characters) with “against”. We could not say which one is always better in practice because it depends on the application scenarios.

As for results of local edit distance algorithm. The predicted words is explosive large but both precision and recall are not satisfactory. Which indicates that this algorithm is not suitable for the task. The reason might be that misspelled

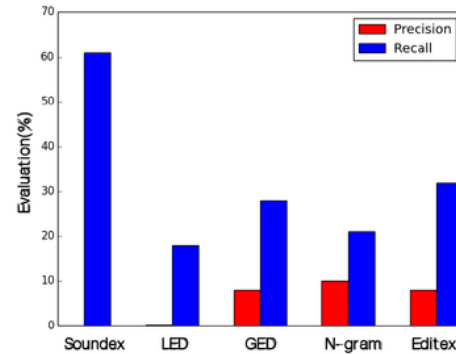


Figure 3: Evaluation visualization

words are not tend to be a part of the corresponding corrections or reverse. Good Crit. Analysis It does not means local edit distance is a bad algorithm. It is very powerful in other tasks such as gene alignments (Altschul and Erickson, 1986).

4.4 Full Scale Evaluation and Demonstration

Table 4 shows a full scale evaluation and Figure 3 is the visualization of it. Table 5 demonstrates output of the best match of two words for each algorithm.

From Table 4 we knows that both Soundex and local edit distance predicted plenty of best matches. The reason for Soundex is that for words with a fixed first letter, the number of different Soundex code is only 1000 (3 digitals from 0-9). This will lead to many words have the same Soundex code. The first two lines of Table 5 demonstrates the Soundex’s results, from which we find the pronounce of the matched set is somewhat similar but the spelling of them various, which will leads to a number of meaningless predicts.

Although the precision of Soundex is extremely low, it seems much better than local edit distance for both precision and recall, which also demonstrates in Table 5. Obviously, “actually” and any string with a substring “actually” are treated the same. This is the trait

Method	Misspelled	Correct	Matched set
Soundex	accually ahain	actually again	akal, axile, azalea, asylees, auxiliar, ... awin, annoy, aani, aoyama, anne, ayme, anay, ...
Local Edit Distance	accually ahain	actually again	actually, tactually, unactually, contactually, ... disenchain, rechain, toolchain, toolchains, ...
Global Edit Distance	accually ahain	actually again	actually chain, amain, arain, again, ghain, alain, hain
N-Gram (N=2)	accually ahain	actually again	actually, ally ain
Editex	accually ahain	actually again	usually, actually, annually, casually, chally, ... amain, attain, arain, again, alain, hain

Table 5: Output demonstration for all algorithms.

of local edit distance, part matching.

As for global edit distance matching, it is such a pretty good algorithm for this spelling correction task that get a acceptable precision and recall. Based on the analysis of Soundex's high recall and global edit distance's high precision and recall, it is not hard to think of computing edit distance based on the Soundex results. Editex algorithm make use of Soundex information and combine it with edit distance algorithm.

We find that Editex gets almost the same precision with global edit distance but higher recall, which might because Editex combines the advantage of Soundex with edit distance. Form the last two rows of Table 5, we find that Editex matched words with small edit distance and also consider pronunciation.

5 Conclusion

We have investigated several classical approximate string matching algorithms with a particular spelling correction task and evaluated by recall and prision. We find that Editex algorithm achieves better performance then others. Besides, we also illustrates the characteristics of each method include n-gram, local edit distance, global edit distance with different parameters, Soundex and Editex in the paper.

References

- S. F. Altschul and B. W. Erickson. 1986. Locally optimal subalignments using nonlinear similarity functions. *Bulletin of Mathematical Biology*, 48(5-6):633–660.
- Eric Brill and Robert C Moore. 2000. An improved error model for noisy channel spelling correction. In *Proc. Meeting of the Association for Computational Linguistics*, pages 286–293.
- Steffen Eger, Tim Vor Der Brck, and Alexander Mehler. 2016. A comparison of

four character-level string-to-string translation models for (ocr) spelling error correction. *Prague Bulletin of Mathematical Linguistics*, 105(1):77–99.

- T. N. Gadd. 1988. fishing fore werds: phonetic retrieval of written text in information systems. *Program Electronic Library & Information Systems*, 22(3):222–237.

Patrick A. V. Hall and Geoff R. Dowling. 1980a. Approximate string matching. *ACM Comput. Surv.*

Patrick A. V. Hall and Geoff R. Dowling. 1980b. Approximate string matching. *ACM Comput. Surv.*, 12:381–402.

Bo Han and Timothy Baldwin. 2011. Lexical normalisation of short text messages: Maku sens a #twitter. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies - Volume 1*, pages 368–378.

V. I Levenshtein. 1965. Binary codes capable of correcting deletions, insertions and reversals. *Dokl.akad.nauk Sssr*, 10(1).

Vijay Raghavan, Gwang S. Jung, and Peter Bollmann. 1989. A critical investigation of recall and precision as measures of retrieval system performance. 7:205–229, 07.

Naomi Saphra. 2016. Evaluating informal-domain word representations with urbandictionary. In *RepEval@ACL*.

Miikka Silfverberg, Pekka Kauppinen, and Kristin Lindn. 2016. Data-driven spelling correction using weighted finite-state methods. In *Sigfsm Workshop on Statistical Nlp and Weighted Automata*, pages 51–59.

Esko Ukkonen. 1992. Approximate string matching with q-grams and maximal matches. *Theor. Comput. Sci.*, 92:191–211.

Justin Zobel and Philip W. Dart. 1996. Phonetic string matching: Lessons from information retrieval. In *SIGIR*.

report

GRADEMARK REPORT

FINAL GRADE

GENERAL COMMENTS

Instructor

11 /14

PAGE 1



Comment 1

not needed



Comment 2

not needed



Comment 3

You waste your first page. We assume the reader is familiar with all this.

PAGE 2



Comment 4

Where are the results from 5 to 9?



Comment 5

too much!



Comment 6

why?

PAGE 3



Comment 7

why? It is 25.56% with 1 [modification | indel]

QM

Good Crit. Analysis

This is a good use of Critical Analysis.



Needs more Analysis

This would have been a good opportunity for more Critical Analysis; perhaps you could have examined more tokens from the data, or looked into what the system was producing in more detail.



Comment 8

not a surprise

PAGE 4



Good Crit. Analysis

This is a good use of Critical Analysis.

METHOD (20%)

8 / 10

10 (10)	<ul style="list-style-type: none"> • All technical tasks accomplished and discussed in the report • Parameters identified and contrasted where necessary • Evaluation is systematic and logical • Methods are clearly illustrative to the knowledge gained
8 OR 9 (9)	<ul style="list-style-type: none"> • All technical tasks accomplished and discussed in the report • Parameters identified but perhaps only weakly contrasted • Evaluation is systematic and logical • Methods are adequate for the knowledge gained
8 OR 9 (8)	<ul style="list-style-type: none"> • All technical tasks accomplished and discussed in the report • Parameters identified but perhaps only weakly contrasted • Evaluation is systematic and logical • Methods are adequate for the knowledge gained
7 (7)	<ul style="list-style-type: none"> • Technical tasks mostly accomplished, with perhaps some shortcomings in logic or presentation • Parameters not identified or not contrasted • Evaluation is logical but not systematic • Methods limit the prospective of gaining knowledge
5 OR 6 (6)	<ul style="list-style-type: none"> • Some major issues with the logic of a task or tasks • Evaluation is illogical or absent • Methods are inadequate and prevent meaningful analysis
5 OR 6 (5)	<ul style="list-style-type: none"> • Some major issues with the logic of a task or tasks • Evaluation is illogical or absent • Methods are inadequate and prevent meaningful analysis
0 TO 4 (4)	<ul style="list-style-type: none"> • Tasks are essentially incomplete or not attempted
0 TO 4 (3)	<ul style="list-style-type: none"> • Tasks are essentially incomplete or not attempted
0 TO 4 (2)	<ul style="list-style-type: none"> • Tasks are essentially incomplete or not attempted
0 TO 4 (1)	<ul style="list-style-type: none"> • Tasks are essentially incomplete or not attempted
0 TO 4 (0)	<ul style="list-style-type: none"> • Tasks are essentially incomplete or not attempted

CR. ANALYSIS (50%)

7 / 10

10 (10)	<ul style="list-style-type: none"> • Clearly identifies the knowledge gained about the task • Argumentation is logical and incontrovertibly supported by evidence • Theoretical properties of methods are well understood and linked to practical behaviour • Demonstrates a very high level of abstract thought • Admirably situated with respect to the academic community • Publishable with perhaps minor changes
8 OR 9 (9)	<ul style="list-style-type: none"> • Clearly identifies the knowledge gained about the task • Argumentation is logical and thoroughly supported by evidence • Theoretical properties of methods are well-understood and linked to practical behaviour • Demonstrates a moderate level of

abstract thought • Attempts to situate with respect to the academic community, but perhaps not clearly

8 OR 9
(8)

• Clearly identifies the knowledge gained about the task • Argumentation is logical and thoroughly supported by evidence • Theoretical properties of methods are well-understood and linked to practical behaviour • Demonstrates a moderate level of abstract thought • Attempts to situate with respect to the academic community, but perhaps not clearly

7
(7)

• Attempts to identify the knowledge gained about the task, but vague or unclear • Argumentation is logical, but evidence is lacking in some areas • Theoretical properties of methods are understood, but not clearly linked to practical behaviour • Demonstrates abstract thought, but extended analysis not always clear or successful • Little connection to the academic community

5 OR 6
(6)

• Knowledge gained about the task is fundamentally flawed or lacking • Argumentation is illogical in places, and evidence is inadequate or contradictory • Theoretical properties of methods are not in evidence • No signs of abstract thought and/or analysis • No connection to the academic community

5 OR 6
(5)

• Knowledge gained about the task is fundamentally flawed or lacking • Argumentation is illogical in places, and evidence is inadequate or contradictory • Theoretical properties of methods are not in evidence • No signs of abstract thought and/or analysis • No connection to the academic community

0 TO 4
(4)

• No indication of knowledge gained about the task • Argumentation is generally absent • Mostly data without corresponding analysis • Theoretical properties of methods are not in evidence • No connection to the academic community

0 TO 4
(3)

• No indication of knowledge gained about the task • Argumentation is generally absent • Mostly data without corresponding analysis • Theoretical properties of methods are not in evidence • No connection to the academic community

0 TO 4
(2)

• No indication of knowledge gained about the task • Argumentation is generally absent • Mostly data without corresponding analysis • Theoretical properties of methods are not in evidence • No connection to the academic community

0 TO 4
(1)

• No indication of knowledge gained about the task • Argumentation is generally absent • Mostly data without corresponding analysis • Theoretical properties of methods are not in evidence • No connection to the academic community

0 TO 4
(0)

• No indication of knowledge gained about the task • Argumentation is generally absent • Mostly data without corresponding analysis • Theoretical properties of methods are not in evidence • No connection to the academic community

REP. QUALITY (30%)

8 / 10

10
(10)

• Ideas and arguments are cohesive, where the components of the report clearly indicate how they relate to the whole • Report structure is logical and formal, in line with typical standards in academic writing • Generally clear and easy-to-follow • References are suitably synthesised and chosen discriminately with respect to the

given problem • Adequately concise and meets word limits

8 OR 9 (9)	• Ideas and arguments are coherent, and generally the work fits together as a unit • Report structure is logical and formal, with small divergences from typical academic standards • Generally clear, with small disruptions in flow • References are suitably synthesised, but are too few or chosen indiscriminately • Adequately concise and meets word limits
8 OR 9 (8)	• Ideas and arguments are coherent, and generally the work fits together as a unit • Report structure is logical and formal, with small divergences from typical academic standards • Generally clear, with small disruptions in flow • References are suitably synthesised, but are too few or chosen indiscriminately • Adequately concise and meets word limits
7 (7)	• Ideas and arguments are mostly coherent, but do not come together in a unified way • Report structure is logical, but possibly informal or out-of-line with academic standards • Some unclear sections that do not detract from the overall work • References are present, but terse or disconnected from the problem at hand • Perhaps small divergences from the word limits
5 OR 6 (6)	• Ideas and arguments are notably incoherent • Report structure is flawed • Some unclear sections which detract from the overall work • References are disconnected or absent • Possibly way off the word limits
5 OR 6 (5)	• Ideas and arguments are notably incoherent • Report structure is flawed • Some unclear sections which detract from the overall work • References are disconnected or absent • Possibly way off the word limits
0 TO 4 (4)	• Ideas and arguments are missing or impossible to follow • Report has no structure or references • Not a formal report, even at a stretch
0 TO 4 (3)	• Ideas and arguments are missing or impossible to follow • Report has no structure or references • Not a formal report, even at a stretch
0 TO 4 (2)	• Ideas and arguments are missing or impossible to follow • Report has no structure or references • Not a formal report, even at a stretch
0 TO 4 (1)	• Ideas and arguments are missing or impossible to follow • Report has no structure or references • Not a formal report, even at a stretch
0 TO 4 (0)	• Ideas and arguments are missing or impossible to follow • Report has no structure or references • Not a formal report, even at a stretch

CRITERION 4 (0%)

0 / 10

10
(10)

8 OR 9
(9)

8 OR 9

(8)

7

(7)

5 OR 6

(6)

5 OR 6

(5)

0 TO 4

(4)

0 TO 4

(3)

0 TO 4

(2)

0 TO 4

(1)

0 TO 4

(0)