

The copyright © of this thesis belongs to its rightful author and/or other copyright owner. Copies can be accessed and downloaded for non-commercial or learning purposes without any charge and permission. The thesis cannot be reproduced or quoted as a whole without the permission from its rightful owner. No alteration or changes in format is allowed without permission from its rightful owner.



**AN IMPROVED LEVENSHTTEIN ALGORITHM FOR SPELLING
CORRECTION WORD CANDIDATE LIST GENERATION**

HANAN NAJM ABDULKHUDHUR



COLLEGE OF ARTS AND SCIENCE

UNIVERSITI UTARA MALAYSIA

2016

Permission to Use

In presenting this dissertation in fulfillment of the requirements for a postgraduate degree from Universiti Utara Malaysia, I agree that the Universiti Library may make it freely available for inspection. I further agree that permission for the copying of this dissertation in any manner, in whole or in part, for the scholarly purpose may be granted by my supervisor(s) or, in their absence, by the Dean of Awang Had Salleh Graduate School of Arts and Sciences. It is understood that any copying or publication or use of this dissertation or parts thereof for financial gain shall not be allowed without my written permission. It is also understood that due recognition shall be given to me and to Universiti Utara Malaysia for any scholarly use which may be made of any material from my dissertation.

Requests for permission to copy or to make other use of materials in this dissertation, in whole or in part should be addressed to:



Dean of Awang Had Salleh Graduate School of Arts and Sciences

UUM College of Arts and Sciences

Universiti Utara Malaysia

06010 UUM Sintok a Malaysia

Abstrak

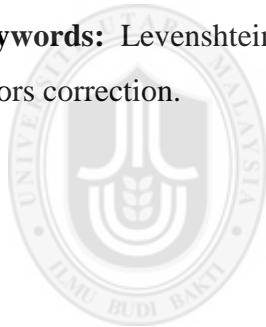
Senarai calon terhasil dalam pembetulan ejaan adalah satu proses untuk mencari kata-kata dari leksikon yang hampir sama dengan perkataan yang tidak tepat. Algoritma paling banyak digunakan untuk menjana senarai calon untuk kata-kata yang tidak tepat adalah berdasarkan jarak Levenshtein. Walau bagaimanapun, algoritma ini mengambil masa yang terlalu lama apabila terdapat bilangan besar kesilapan ejaan. Sebabnya ialah bahawa pengiraan algoritma Levenshtein termasuk operasi yang menghasilkan jajaran dan pengisian sel-sel jajaran dengan membandingkan huruf-huruf perkataan yang tidak betul dengan huruf-huruf perkataan dari leksikon. Oleh kerana kebanyakan leksikon mengandungi berjuta-juta perkataan, maka operasi ini akan diulang berjuta-juta kali bagi setiap perkataan tidak tepat untuk menjana senarai calonnya. Kajian ini menambahbaikkan algoritma Levenshtein dengan merekabentuk teknik operasi yang telah dimasukkan dalam algoritma ini. Teknik operasi yang dicadangkan meningkatkan algoritma Levenshtein dari segi masa pemprosesan perlaksanaannya tanpa menjejaskan ketepatanannya. Ia mengurangkan langkah operasi yang diperlukan untuk mengukur nilai sel-sel dalam baris dan lajur pertama, baris dan lajur kedua serta baris dan lajur ketiga dalam jajaran Levenshtein. Algoritma Levenshtein yang telah tingkatan telah dibandingkan dengan algoritma asal. Hasil kajian menunjukkan bahawa prestasi algoritma yang dicadangkan melebihi prestasi algoritma Levenshtein asal dari segi masa pemprosesan, iaitu sebanyak 36.45% manakala ketepatan kedua-dua algoritma adalah masih sama.

Kata Kunci: Algoritma Levenshtein, Masa pemprosesan, Penghasilan senarai calon, Pmbetulan Kesilapan.

Abstract

Candidates' list generation in spelling correction is a process of finding words from a lexicon that should be close to the incorrect word. The most widely used algorithm for generating candidates' list for incorrect words is based on Levenshtein distance. However, this algorithm takes too much time when there is a large number of spelling errors. The reason is that calculating Levenshtein algorithm includes operations that create an array and fill the cells of this array by comparing the characters of an incorrect word with the characters of a word from a lexicon. Since most lexicons contain millions of words, then these operations will be repeated millions of times for each incorrect word to generate its candidates list. This dissertation improved Levenshtein algorithm by designing an operational technique that has been included in this algorithm. The proposed operational technique enhances Levenshtein algorithm in terms of the processing time of its executing without affecting its accuracy. It reduces the operations required to measure cells' values in the first row, first column, second row, second column, third row, and third column in Levenshtein array. The improved Levenshtein algorithm was evaluated against the original algorithm. Experimental results show that the proposed algorithm outperforms Levenshtein algorithm in terms of the processing time by 36.45% while the accuracy of both algorithms is still the same.

Keywords: Levenshtein Algorithm, Processing time, Candidates' list generation, Errors correction.



UUM
Universiti Utara Malaysia

Acknowledgement

Each part of this dissertation is guided, inspired, and supported by many people. Firstly, I would like to thank all the members of my family especially my husband and my parents for their unconditional support. My goal would not be achieved without them. The most important support and guidance were from my research supervisors Dr. Shahrul Azmi Mohd Yusof and Prof. Madya Dr. Yuhanis Yusof. Thank you very much for your great help and support. It is an honor for me to do a research under your supervisions. I would like to thank all the academic and technical staff in School of Computing of Utara Universiti Malaysia for their help in the study process and providing all the excellent facilities. Finally, I would like to thank all my friends for their support.



Table of Contents

Permission to Use.....	i
Abstrak	ii
Abstract	iii
Acknowledgement.....	iv
Table of Contents	v
List of Tables.....	viii
List of Figures	ix
List of Appendices	x
List of Abbreviations.....	xi
CHAPTER ONE INTRODUCTION	1
1.1 Background	1
1.2 Problem Statement	4
1.3 Research Questions	6
1.4 Research Objectives	6
1.5 Significance of the Dissertation	6
1.6 Scope of the Dissertation.....	7
1.7 Organization of the Dissertation	8
CHAPTER TWO LITERATURE REVIEW	9
2.1 Introduction.....	9
2.2 Spelling Correction Stages.....	9
2.2.1 Error Detection Stage.....	9
2.2.2 Candidates List Generation Stage	10
2.2.3 Error Correction Stage	11
2.3 Techniques of Candidates List Generation Based on Edit Distance.....	11
2.3.1 Levenshtein Distance	12
2.3.2 Hamming Distance.....	16
2.3.3 Damerau-Levenshtein Distance	17
2.3.4 Longest Common Subsequence Distance	17
2.5.5 Threshold Levenshtein Distance	18
2.5.6 Levenshtein Automata Distance	18
2.4 Other Techniques of Candidates List Generation	20
2.4.1 N-gram Distance	20

2.4.2 Bag Distance	21
2.4.3 N-gram language Model	21
2.4.4 Topic Model	22
2.5 Evaluation Measurement	23
2.6 Summary	25
CHAPTER THREE RESEARCH METHODOLOGY	26
3.1 Introduction	26
3.2 Research Phases	26
3.3 Theoretical Study	27
3.4 Design Phase	28
3.5 Development Phase	30
3.6 Evaluation	30
3.6.1 Experimental Design	30
3.6.1.1 How to Measure Processing Time	31
3.6.1.2 How to Measure Accuracy	33
3.6.1.3 Statistical Test	35
3.6.2 Testing Dataset	36
3.7 Summary	38
CHAPTER FOUR PROPOSED DIFFERENTIATION TECHNIQUE	39
4.1 Introduction	39
4.2 ILA-OT Steps	39
4.2.1 First Step	40
4.2.2 Second Step	41
4.2.3 Third Step	45
4.3 Comparison between operations of LA and ILA-OT	49
4.4 ILA-OT Algorithm	55
4.5 Summary	60
CHAPTER FIVE EXPERIMENTAL RESULTS AND DISCUSSION	62
5.1 Introduction	62
5.2 Processing Time	62
5.3 Accuracy	66
5.4 Statistical Test	67
5.5 Results Discussion	68

5.6 Summary	69
CHAPTER SIX CONCLUSION AND FUTURE WORK	70
6.1 Introduction	70
6.2 Research Summary.....	70
6.3 Research Contributions	71
6.4 Recommendation for Future Work	73
6.5 Summary	74
REFERENCES	75
APPENDICES	80



UUM
Universiti Utara Malaysia

List of Tables

Table 3.1 Main variables resulted from t-test.	36
Table 3.2 Characteristics of the testing dataset.	37
Table 4.1 Theoretical comparison between operations of LA and ILA-OT.	52
Table 4.2 Practical comparison between LA and ILA-OT.	54
Table 5.1 A sample of the processing time values of both algorithms	62
Table 5.2 Processing times of both algorithms LA and ILA-OT.	64
Table 5.3 A sample of the distances of both algorithms for several words.	66
Table 5.4 T-test results between LA and ILA-OT.	67



List of Figures

Figure 1.1. The scope of this research.	7
Figure 2.1. Levenshtein distance example	13
Figure 2.2. Bigram distance example.....	20
Figure 3.1. Research phases.	27
Figure 4.1. Levenshtein array before and after applying the first step.	40
Figure 4.2. The second step examples for the first row.	42
Figure 4.3. The second step examples for the first column.	44
Figure 4.4. Third step examples for the second row	46
Figure 4.5. The third step examples for the second column.	48
Figure 4.6. Cells affected by three steps proposed for this dissertation.....	54
Figure 5.1. Line graph for the processing time values listed in Table 5.2.	65



UUM
Universiti Utara Malaysia

List of Appendices

Appendix A Sample of Output Results of Word Length 3.....	80
Appendix B Sample of Output Results of Word Length 4.....	81
Appendix C Sample of Output Results of Word Length 5.....	82
Appendix D Sample of Output Results of Word Length 6.....	83
Appendix E Sample of Output Results of Word Length 7.....	84
Appendix F Sample of Output Results of Word Length 8.....	85
Appendix G Sample of Output Results of Word Length 9.....	86
Appendix H Sample of Output Results of Word Length 10.....	87
Appendix I Sample of Output Results of Word Length 11.....	88
Appendix J Sample of Output Results of Word Length 12.....	89



UUM
Universiti Utara Malaysia

List of Abbreviations

LA	Levenshtein algorithm
HD	Hamming distance
DD	Damerau distance
LCS	Longest common subsequence
OT	Operational technique
ILA-OT	Improved Levenshtein algorithm by using the proposed OT
MVFRFC	Measure values of the first row and first column
MVSRSC	Measure values of the second row and second column
PT	Processing time
PD	Percentage decrease



UUM
Universiti Utara Malaysia

CHAPTER ONE

INTRODUCTION

1.1 Background

Spelling correction is the process of detecting and repairing spelling errors in a text. Research in spelling correction is not new; it started in the mid of 1960, and many algorithms for spelling correction have been suggested since then (Mahdi, 2012). Spelling correction can be either manual or automatic. The first type allows intervention of humans in the correction process. The second type, a system will decide the correction to replace an incorrect word by choosing the best candidate word without human's intervention (Bassil & Alwani, 2012b).

Most methods of automatic spelling correction have three functions: error detection, generation of candidates, and error correction (Naseem & Hussain, 2007). The first function is to find incorrect words in the output text. The second function is to generate candidate words from a lexicon for each of the incorrect words. Candidate list generation is a process of finding words from a lexicon that should be close to the incorrect word. For example, the candidates' list generated from a lexicon for the incorrect word "czp" are "cup", "cap", and "cop". The last function is to correct all incorrect words by selecting the best candidate to replace with each incorrect word.

The process of generating candidates list can be achieved by using a specific algorithm. An algorithm is a set of operations that will be performed on some data to solve a specific problem. In general, algorithms can be classified according to their optimal solution into two categories: exact and approximate. In execution, exact algorithms will reach an optimal solution while approximation algorithms can be

The contents of
the thesis is for
internal user
only

REFERENCES

- Adhitama, P., Kim, S. H., & Na, I. S. (2014). Lexicon-Driven Word Recognition Based on Levenshtein Distance. *International Journal of Software Engineering and Its Applications*, 8(2), 11-20.
- Ahmed, B. (2015). *Lexical Normalisation of Twitter Data*. Paper presented at the Science and Information Conference (SAI), 2015.
- Al-Bakry, A. M., & Al-Rikaby, M. K. (2015). *Enhanced Levenshtein Edit Distance Method functioning as a String-to-String Similarity Measure*. Paper presented at the Proceedings of the Networks Security and Distributed Systems (NSDS'2015).
- Al-Masoudi, A. F. R., & Al-Obeidi, H. S. R. (2015). Smoothing Techniques Evaluation of N-gram Language Model for Arabic OCR Post-processing. *Journal of Theoretical and Applied Information Technology*, 82(3), 432-439.
- Al-Zaydi, Z. Q., & Salam, H. (2015). Multiple Outputs Techniques Evaluation for Arabic Character Recognition. *International Journal of Computer Techniques (IJCT)*, 2(5), 1-7.
- Alobaedy, M. M. T. (2015). *Hybrid Ant Colony System Algorithm For Static And Dynamic Job Scheduling In Grid Computing*. (PhD thesis, Universiti Utara Malaysia, Kedah, Malaysia). Retrieved from <http://etd.uum.edu.my/id/eprint/5382>
- Andoni, A., & Krauthgamer, R. (2012). The smoothed complexity of edit distance. *ACM Transactions on Algorithms (TALG)*, 8(4), 44.
- Andoni, A., & Onak, K. (2012). Approximating edit distance in near-linear time. *SIAM Journal on Computing*, 41(6), 1635-1648.
- Attia, M., Pecina, P., Samih, Y., Shaalan, K. F., & van Genabith, J. (2012). *Improved Spelling Error Detection and Correction for Arabic*. Paper presented at the International Conference on Computational Linguistics (COLING), Mumbai, India.
- Bard, G. V. (2007). *Spelling-error tolerant, order-independent pass-phrases via the Damerau-Levenshtein string-edit distance metric*. Paper presented at the Proceedings of the fifth Australasian symposium on ACSW frontiers-Volume 68.
- Bassil, Y., & Alwani, M. (2012a). Context-sensitive Spelling Correction Using Google Web 1T 5-Gram Information. *Computer and Information Science*, 5(3), 37-48.

- Bassil, Y., & Alwani, M. (2012b). Ocr post-processing error correction algorithm using google online spelling suggestion. *Journal of Emerging Trends in Computing and Information Sciences*, 3(1), 90-99.
- Batawi, Y., & Abulnaja, O. (2012). Accuracy Evaluation of Arabic Optical Character Recognition Voting Technique: Experimental Study. *IJECS: International Journal of Electrical & Computer Sciences*, 12(1), 29-33.
- Bergroth, L., Hakonen, H., & Raita, T. (2000). *A survey of longest common subsequence algorithms*. Paper presented at the Proceedings of the Seventh International Symposium on String Processing Information Retrieval (SPIRE'00).
- Burkhardt, S., & Kärkkäinen, J. (2002). *One-gapped q-gram filters for Levenshtein distance*. Paper presented at the Proceedings of the Combinatorial Pattern Matching.
- Cooper, L., & Cooper, M. W. (1981). *Introduction to dynamic programming*: Pergamon Press New York.
- Crumrine, K. T., Ritschel, J. D., & White, E. (2014). Earned Schedule 10 Years Later Analyzing Military Programs: DTIC Document.
- Daðason, J. F. (2012). *Post-Correction of Icelandic OCR Text*. (Master's thesis), University of Iceland, Reykjavik, Iceland.
- Federico, M., & Cettolo, M. (2007). *Efficient handling of n-gram language models for statistical machine translation*. Paper presented at the Proceedings of the Second Workshop on Statistical Machine Translation.
- Formiga Fanals, L., & Rodríguez Fonollosa, J. A. (2012). *Dealing with input noise in statistical machine translation*. Paper presented at the International Conference on Computational Linguistics (COLING), Mumbai, India.
- Gomaa, W. H., & Fahmy, A. A. (2013). A survey of text similarity approaches. *International Journal of Computer Applications*, 68(13).
- Google. (2015). Google Ngram Retrieved september 05, 2015, from <http://storage.googleapis.com/books/ngrams/books/datasetsv2.html>
- Grzebala, P. B. (2016). *Private Record Linkage: A Comparison of Selected Techniques for Name Matching*. Master thesis, Wright State University, Ohio, United States.
- Gupta, B., Bhatt, G., & Mittal, A. (2016). Language Identification and Disambiguation in Indian Mixed-Script. *Distributed Computing and Internet Technology* (pp. 113-121): Springer.

- Habeeb, I. Q., Yusof, S. A., & Ahmad, F. B. (2014). Two Bigrams Based Language Model for Auto Correction of Arabic OCR Errors. *International Journal of Digital Content Technology and its Applications*, 8(1), 72 - 80.
- Haldar, R., & Mukhopadhyay, D. (2011). Levenshtein distance technique in dictionary lookup methods: An improved approach. *arXiv preprint arXiv:1101.1232*.
- Hassan, A., Noeman, S., & Hassan, H. (2008). *Language Independent Text Correction using Finite State Automata*. Paper presented at the Third International Joint Conference on Natural Language Processing (IJCNLP), Hyderabad, India.
- Huldén, M. (2009). Fast approximate string matching with finite automata. *Procesamiento del lenguaje natural*, 43, 57-64.
- Islam, A., & Inkpen, D. (2009). *Real-word spelling correction using Google Web IT n-gram with backoff*. Paper presented at the International Conference on Natural Language Processing and Knowledge Engineering.
- Jurafsky, D., & Martin, J. H. (2009). *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition* (2nd ed.): Pearson Education India.
- Levenshtein, V. I. (1966). *Binary codes capable of correcting deletions, insertions and reversals*. Paper presented at the Soviet physics doklady.
- Li, M., Zhang, Y., Zhu, M., & Zhou, M. (2006). *Exploring distributional similarity based models for query spelling correction*. Paper presented at the Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics.
- Lounis, O., Guermeche, B., Eddine, S., Saoudi, L., & Benaicha, S. E. (2014). *A new algorithm for detecting SQL injection attack in Web application*. Paper presented at the Science and Information Conference (SAI), 2014.
- Lu, W., Du, X., Hadjieleftheriou, M., & Ooi, B. C. (2014). Efficiently Supporting Edit Distance based String Similarity Search Using B-trees. *IEEE Transactions on Knowledge and Data Engineering*, 26(12), 2983-2996.
- Lund, W. B. (2014). *Ensemble Methods for Historical Machine-Printed Document Recognition*. (PhD dissertation, Brigham Young University, Utah, United States).
- Ma, D., & Agam, G. (2013). *A super resolution framework for low resolution document image OCR*. Paper presented at the IS&T/SPIE Electronic Imaging.
- Maarif, H., Akmeliawati, R., Htike, Z., & Gunawan, T. S. (2014). *Complexity Algorithm Analysis for Edit Distance*. Paper presented at the International Conference on Computer and Communication Engineering (ICCCE), 2014

- Mahdi, A. A. M. (2012). *Spell checking and correction for Arabic text recognition*. (Master thesis), King Fahd university of petroleum and minerals, Saudi Arabia.
- Mainsah, B. O., Morton, K. D., Collins, L. M., Sellers, E. W., & Throckmorton, C. S. (2015). Moving away from error-related potentials to achieve spelling correction in P300 spellers. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 23(5), 737-743.
- Mihov, S., & Schulz, K. U. (2004). Fast approximate search in large dictionaries. *Computational Linguistics*, 30(4), 451-477.
- Musa, S. M., Opara, E. U., Shayib, M. A., & Oliver, J. (2016). Measurement and Test Performance for Integrated Digital Loop Carrier for White Noise Impairment Using Interleaved Mode. *Communications of the IIMA*, 14(3), 4.
- Naseem, T. (2004). *A Hybrid Approach for Urdu Spell Checking*. (Master thesis), National University of Computer & Emerging Sciences, Islamabad, Pakistan.
- Naseem, T., & Hussain, S. (2007). A novel approach for ranking spelling error corrections for Urdu. *Language Resources and Evaluation*, 41(2), 117-128. doi: 10.1007/s10579-007-9028-6
- Navarro, G. (2001). A Guided tour to approximate string matching. *ACM Computing Surveys (CSUR)*, 33(1), 31-88.
- Navarro, G., Grabowski, S., Mäkinen, V., & Deorowicz, S. (2005). Improved time and space complexities for transposition invariant string matching *Technical Report TR/DCC-2005-4, Department of Computer Science: University of Chile*.
- Pal, S., & Rajasekaran, S. (2015). *Improved algorithms for finding edit distance based motifs*. Paper presented at the IEEE International Conference on Bioinformatics and Biomedicine (BIBM), 2015.
- Phillips, C. R. (2015). *Employing an Efficient and Scalable Implementation of the Cost Sensitive Alternating Decision Tree Algorithm to Efficiently Link Person Records*. Master thesis, Texas State University, San Marcos, United States.
- Polyanovsky, V., Roytberg, M. A., & Tumanyan, V. G. (2011). Comparative analysis of the quality of a global algorithm and a local algorithm for alignment of two sequences. *Algorithms for Molecular Biology*, 6, 25.
- Popescu, O., & Vo, N. P. A. (2014). *Fast and Accurate Misspelling Correction in Large Corpora*. Paper presented at the Conference on Empirical Methods in Natural Language Processing (EMNLP).
- Pradhan, N., Gyanchandani, M., & Wadhvani, R. (2015). A Review on Text Similarity Technique used in IR and its Application. *International Journal of Computer Applications*, 120(9).

- Rardin, R. L., & Uzsoy, R. (2001). Experimental evaluation of heuristic optimization algorithms: A tutorial. *Journal of Heuristics*, 7(3), 261-304.
- Rieck, K., & Wressnegger, C. (2016). Harry: A Tool for Measuring String Similarity. *Journal of Machine Learning Research*, 17(9), 1-5.
- Ruoro, S. W. (2009). *A Parallel Corpus Based Translation Using Sentence Similarity*. PhD dissertation, University of Nairobi, Nairobi, Kenya.
- Sheng, C., Tao, Y., & Li, J. (2012). Exact and approximate algorithms for the most connected vertex problem. *ACM Transactions on Database Systems (TODS)*, 37(2), 12.
- Siklósi, B., Novák, A., & Prószéky, G. (2016). Context-aware correction of spelling errors in Hungarian medical documents. *Computer Speech & Language*, 35, 219-233.
- Singh, S. P., Kumar, A., Darbari, H., Chauhan, S., Srivastava, N., & Singh, P. (2015). Evaluation of Similarity metrics for translation retrieval in the Hindi-English Translation Memory. *International Journal of Advanced Research in Computer and Communication Engineering*, 4(8).
- Ukkonen, E. (1985). Finding approximate patterns in strings. *Journal of algorithms*, 6(1), 132-137.
- Vargas, S. G. J. (2008). *A Knowledge-Based information Extraction Prototype for Data-Rich Documents in the Information Technology Domain*. (Master dissertation), National University of Colombia, Bogotá, Colombia.
- Wick, M. L., Ross, M. G., & Learned-Miller, E. G. (2007). *Context-sensitive error correction: Using topic models to improve OCR*. Paper presented at the Ninth International Conference on Document Analysis and Recognition (ICDAR).