

**LAPORAN AKHIR
PENELITIAN MANDIRI**



**Combination of Genetic Algorithm and Brill Tagger Algorithm for Part of
Speech Tagging Bahasa Madura**

TIM PENGUSUL

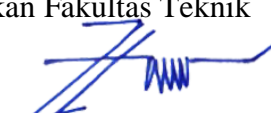
1. Nindian Puspa Dewi M.Kom (Ketua)
NIDN. 710413344
2. Ubaidi M.Kom (Anggota)
NIDN. 0707118603

**UNIVERSITAS MADURA
2020**

HALAMAN PENGESAHAN

- 1. Judul Penelitian** : EFEKTIFITAS PENGGUNAAN
TEKNOLOGI PADA PESANTREN
MODERNDALAM MENGHADAPI
REVOLUSI INDUSTRI 4.0
- 2. Ketua Tim Pengusul**
- a. Nama Lengkap : Nindian Puspa Dewi M.Kom
 - b. NIDN : 710413344
 - c. Jabatan/Golongan : Asisten Ahli/
 - d. Program Studi : Informatika
 - e. Perguruan Tinggi : Universitas Madura
 - f. Bidang Keahlian : Informatika
 - g. Alamat Kantor/Telp/surel : Jl. Panglegur Km 3,5 Pamekasan
- 3. Jangka waktu Pelaksanaan** : 1 tahun
- 4. Biaya Total** : Rp. 3500000,-

Mengetahui,
Dekan Fakultas Teknik


Ir. H. Moch. Hazin Mukti, MT, MM.
NIP.195906051987031002

Pamekasan, 30 Maret 2020
Ketua Peneliti,

Nindian Puspa Dewi M.Kom
NIDN. 710413344

Menyetujui,
Ketua LPPM Universitas Madura


MOH. ZALI, S.Pt, M.Agr
NIDN. 0706088401

DAFTAR ISI

Halaman Sampul	i
Halaman Pengesahan	ii
Daftar Isi	iii
Ringkasan.....	iv
BAB 1 PENDAHULUAN	
1.1. Latar Belakang Penelitian.....	1
1.2. Tujuan Penelitian.....	2
BAB 2 Tinjauan Pustaka	
2.1. Solusi yang Ditawarkan.....	3
BAB 3 METODE PELAKSANAAN	
3.1. Metode Pendekatan Masalah.....	4
BAB 4 HASIL DAN PEMBAHASAN	
4.1 KBM Lembaga Formal dan Non Formal.....	6
4.2 Kegiatan Kebutuhan Santri.....	8
4.3 Kegiatan Pembayaran Kewajiban Santri.....	8
BAB 5 PENUTUP	
5.1. Kesimpulan.....	9
5.2. Saran	9
DAFTAR PUSTAKA	10
LAMPIRAN - LAMPIRAN	

SUMMARY

Part of speech (POS) is commonly known as word types in a sentence such as verbs, adjectives, nouns, and so on. Part of Speech (POS) Tagging is a process of marking the word class or part of speech in every word in a sentence. Part of Speech Tagging has an important role to be used as a basis for research in Natural Language Processing. That is why research on Part of Speech Tagging for Bahasa Madura as an effort to preserve and develop the use of regional languages. In this research, POS Tagging is done using the Brill Tagger Algorithm which is combined with the Genetic Algorithm. Brill Tagger is a POS Tagging Algorithm that has the best level of accuracy when implemented in other languages. Genetic Algorithms used in the contextual learner process with consideration in previous studies can increase the speed of the training process so that it is more efficient. The results of this study are then compared with the results of the previous study so that we can find out suitable algorithms used for the development of text processing in Bahasa Madura. From a series of experiments, the average accuracy obtained by using Brill Tagger is 86.4% with the highest accuracy of 86.7%, while using GA Brill Tagger shows an average accuracy of 86.5% with the highest accuracy of 86.6%. Testing by observing OOV (Out of Vocabulary) achieves an average accuracy of 67.7% for Brill Taggers and 64.6% for GA Brill Taggers. Testing by considering multiple POS with Brill Tagger produces an average accuracy of 73.3% while testing using GA Brill Tagger produces an average accuracy of 90.9%. This shows that the accuracy with GA Brill Tagger is better than Brill Tagger, especially if considering multiple POS. This is because GA Brill Tagger can generate rules for handling the existence of multiple POS more than pure Brill Tagger.

CHAPTER 1

INTRODUCTION

1. Background

Today, technological advances have learned about human language. Many studies have been conducted to process natural language into a computational model. This allows interaction between humans and computers to occur using human language (natural language). Research in this field became known as Natural Language Processing. One study in Natural Language Processing is Part-of-Speech Tagging. Part-of-Speech (POS) is known as word types in a sentence [1] such as verbs, adjectives, nouns, etc. Part-of-Speech (POS) tagging is a process of marking the word class for each word in a sentence. POS Tagging is a basis of research in Natural Language Processing, such as in Word Sense Disambiguation, Stemming in Information Retrieval, and Question and Answering [2]. Research on Part-of-Speech Tagging in Indonesia has been carried out using various methods including POS Indonesian Tagging with Hidden Markov Model and Rule Based [3], Probabilistic Part of Speech Tagging for Indonesian [4] using 37 tag set, Brill Tagger Implementation to provide POS Tagging on Indonesian Language Documents [5], Toward a Standardized and More Accurate Indonesian Part-of-Speech Tagging [6], and On Part of Speech Tagger for Indonesian Language [7]. From several studies that have been done, the highest accuracy value is by using Brill Tagger [8]. Brill Tagger was first introduced by [9]. The Tagger process is a transformation or rules of learning outcomes from detecting error values [10]. From several studies on POS Tagging, the highest accuracy value is to use the Brill Tagger method. Brill Tagger has also applied in many languages, such as English, Kadazan, and Bahasa Indonesia. POS Tagging research using genetic algorithms such as Part-of-Speech Tagging using Genetic Algorithms [11], A New Approach to the POS Tagging Problem Using Evolutionary Computation [12], and Genetic Algorithm (GA) Implementation for Feature Selection in POS Tagging Manipuri [13]. Research that combines Brill Tagger and Genetic Algorithm, was carried out by Wilson who included GA in Brill Tagger to improve time efficiency compared to using Brill Tagger alone [14]. Another study, Genetic Algorithms in the Brill Tagger written by Johannes Bjerva, explained that Brill GA-Tagger performed much better than standard Brill tagger in all 9 target languages [15]

Bahasa Madura is a regional language used by ethnic Madurese, both living on Madura Island and outside the island, as a means of daily communication. The area of Bahasa Madura usage is not only limited to Madura Island but also extends to other places outside the island such as Sapudi, Raas, Goat, Kangean, and other surrounding islands because the majority of the islands are inhabited by Bahasa Madura. Bahasa Madura as a regional language needs to be fostered and developed, especially as a means of developing regional culture and national culture [16]. In previous studies, we have conducted POS Tagging research in Bahasa Madura using the Brill Tagger Algorithm [17], [18].

2. Objective of Research

In this study, we used Brill Tagger combined with genetic algorithms (GA Brill Tagger). The difference with previous research, besides using GA Brill Tagger for POS Tagging in Bahasa Madura, this research also conducted experiments using words that have multiple POS. Multiple POS means words that have more than one class of words or tagset, such as the word "bisa" in Indonesian that can have tagset modals (MD) and tagset Noun (NN). The results of this study are then compared with the results of the previous study so that we can find out suitable algorithms used for the development of text processing in Bahasa Madura.

CHAPTER 2

Literature Review

Brill Tagger introduced by Eric Brill in 1992. Generally, Brill Tagger is also called Transformation-based Error-driven Learning (TEL). Brill Taggers are the basis of transformation or rules and learn from detecting error values [9]. Brill Tagger can give the right word class to a word by using lexical and contextual rules. Lexical rules are the result of lexical learners. Lexical rules are rules used to label words based on word affixes. Contextual rules are rules that pay attention to the existence of tags around the word being checked or searched for labels [19]. Contextual rules are the result of contextual learners.

Genetic algorithm is a search method based on the natural evolutionary process [20], namely the formation of a random initial population consisting of individuals with traits that depend on genes on their chromosomes. Individuals carry out reproductive processes to give birth to offspring. Offspring formed from a combination of the properties of the two parents. Like natural processes that inspire computational processes, populations in Genetic Algorithms also consists of many individuals called chromosomes. If in natural processes chromosomes contain unique individual characteristics, then in the Genetic Algorithm, chromosomes are representations of problem solving that are still symbolic. As with the natural selection process, only fit individuals survive in the population. Each generation, chromosomes will undergo an evaluation process using the fitness function. The fitness value of a chromosome shows the quality of a chromosome in the population. The higher the fitness value of a chromosome, the higher the possibility to be maintained in the next population. The initial chromosomes formed randomly and referred to as the parent. The chromosomes created from the parent chromosome pair are called child (offspring). The process of making a child from its parent is called a crossover operator. This process allows the child to inherit the properties of both parents [21]. In genetic algorithms, there is also a mutation operator (mutations). It is a process that can change genes in a chromosome.

CHAPTER 3

RESEARCH METHOD

The experiment was carried out using a tagset consisting of 34 tagset as in [17]. The compilation of datasets was carried out by collecting articles of Bahasa Madura totaling 10,535 words [18] and manually tagged using a tagset. The results of this labeling are then referred to as Manually Tagged copus (Goal Corpus). The structure of Bahasa Madura is almost the same as Bahasa Indonesia, so the determination of the word class is also not much different. It's just that there are a number of word classes broken down as if in Bahasa Indonesia [23], verbs are simply given a verb word class (VB), then in this study, it is divided into transitive verbs (VBT) and intransitive verbs (VBI).

CHAPTER 4

RESULT AND DISCUSSION

We conducted experiments using computers with the specifications of Intel Corei5 1.7 GHz, 8 GB, and Windows 10 64bit. The POS Tagging application created using the C# programming language. The training process is carried out by changing the threshold to see its effect on the acquisition of rules, both on the lexical learner and contextual learner (GA). For testing, a trial is conducted to find out the accuracy of the training that has been done. For the calculation of the accuracy value, three types of calculations are used, namely the standard calculation without regard to OOV (Out of Vocabulary), the calculation by taking into account multiple POS and the calculation of accuracy by paying attention to the existence of OOV using equation.

From the experimental results of the lexical learner for the 10 threshold produces 48 rules, the threshold of 20 to 40 has decreased the number of rules that is only 32 rules. Likewise for the 50 threshold produces the same rule as many as 13 rules. This shows that the smaller the threshold value, the more rules are produced. The greater the threshold value, the fewer rule will be produced. The same thing happens in contextual learners with Brill Tagger, using threshold 2 produces 48 rules, threshold 3 produces 33 rules and threshold 4 produces 24 rules.

After conducting several contextual learner experiments with Brill Tagger and GA Brill Tagger by making threshold. changes, the number of contextual rules is quite varied depending on the results of randomization. But the smaller the threshold, the more rules are obtained and the greater the threshold, the fewer rules are obtained. Accuracy has increased from lexical results, from an accuracy of 85.81% to 86.61%. Besides that, it is shown that the more rules that are produced (the smaller the accuracy), the better the accuracy tends to be. But in certain cases certain rules can justify the tag of a word and also give the wrong tag for other words. This stage depends on the rules generated in the genetic process where the resulting rules depend on the randomization process. Seen in table IV, the results with GA Brill Tagger for the word *nèp-krennèp* (glittering decoration) get the correct tag because of the rule "NN Prev1 / VBT" which means change the tag to NN if 1 tag was previously VBT, and for Brill Tagger results, the word *mennang* (win) gets the correct rule because the rule "NN JJ PREVWD sè" means that if the initial tag is NN and is located after the word sè then change the tag to JJ. The experiment also be conducted by taking into account the existence of multiple POS. For example for the word *dháddi* which can have the tag as VBT in the sentence *èpateppa 'dháddi bhágus* (to be good) and as the SC in the *dháddi manabi sampèyan songkan entar ka dokter* (so

if you are sick go to the doctor). From 2405 words and symbols in the corpus, there are several words have more than one POS

CHAPTER 5

CONCLUSSION

Testing using Brill Tagger produces an average accuracy of 86.4% with the highest accuracy of 86.7% while testing using GA Brill Tagger produces an average accuracy of 86.5% with the highest accuracy of 86.6%. Testing by considering multiple POS with Brill Tagger produces an average accuracy of 73.4% while testing using GA Brill Tagger produces an average accuracy of 90.9%. Testing with OOV produces an average accuracy of 67.2% with Brill Tagger and an accuracy of 64.6% with GA Brill Tagger. This shows that the accuracy with GA Brill Tagger is better than Brill Tagger, especially if considering multiple POS. This is because GA Brill Tagger can generate rules for handling the existence of multiple POS more than pure Brill Tagger. For future work, the results of this study can be used to conduct other research on Bahasa Madura in the field of Natural Language Preprocessing such as Stemming, Question and Answering. This research can also be utilized for E-learning Bahasa Madura, and this very good because now Bahasa Madura has been abandoned by many Madurese people, especially among young people.

DAFTAR PUSTAKA

- C. D. Manning, S. Hinrich, "Foundation of Statistical Natural Language Processing," Cambridge: MIT Press Textbook on statistical and probabilistic methods in NLP, 1999.
- [2] E. R. Setyaningsih, "Penetapan Tagset dan Modifikasi Brill Tagger untuk Part-of Speech Bahasa Indonesia," *Dinamika Teknologi*, vol. 9 no.2, pp.37-42, 2017.
- [3] K. Widhiyanti, A. Harjoko, "POS Tagging Bahasa Indonesia dengan HMM dan Rule Based," *Jurnal Informatika*, vol. 8 no.2, pp.151-167, 2012.
- [4] F. Pisceldo, M. Adriani, R. Manurung, "Probabilistic Part Of Speech Tagging for Bahasa Indonesia," *Third International MALINDO*

Workshop, 2009.

[5] V. Christanti, J. Pragantha, E. Purnamasari, "Implementasi Brill Tagger untuk memberikan POS-Tagging pada Dokumen Bahasa Indonesia," *Jurnal Teknik dan Ilmu Komputer*, 1(3), pp. 301–315, 2007.

[6] K. Kurniawan, A. F. Aji, "Toward a Standardized and More Accurate Indonesian Part-of-Speech Tagging," *International Conference on Asian Language Processing (IALP)*, 10 September 2018, pp. 303–307, 2018.

[7] R. S. Yuwana, A. R. Yuliani, H. F. Pardede, "On Part of Speech Tagger for Indonesian Language," *International conferences on Information Technology, Information Systems and Electrical Engineering (ICITISEE)*, 1-2 Nopember 2017, pp. 369-372, 2017.

[8] F. M. Hasan, N. Uzzaman, M. Khan, "Comparison of different POS Tagging Techniques (N-Gram, HMM and Brill's tagger) for Bangla," *Advances and Innovations in Systems, Computing Sciences and Software Engineering*, pp.121-126, 2007.

[9] E. Brill, "A simple rule-based part of speech tagger," *Proc. third Conf. Appl. Nat. Lang. Process*, pp. 152, 1992.

[10] N. P. M. Sriyati, "Part-Of-Speech Tagging Untuk Dokumen Bahasa Bali Menggunakan Algoritma Brill Tagger," *Fakultas Matematika dan Ilmu Pengetahuan Alam. Tugas Akhir. Universitas Udayana*, 2016.

[11] K. Singh, "Part-of-Speech Tagging using Genetic Algorithms," *International Journal of Simulation - Systems, Science & Technology*, 16(6), pp.111-117, 2015.

[12] A. P. Silva, A. Silva, I. Rodrigues, "A New Approach to the POS Tagging Problem Using Evolutionary Computation," *Proceedings of Recent Advances in Natural Language Processing*, 7-13 September 2013, pp. 619-625, 2013.

[13] K. Nongmeikapam, S. Bandyopadhyay, "Genetic Algorithm (GA) Implementation for Feature Selection in Manipuri POS Tagging," *Proceedings of the 13th International Conference on Natural Language Processing*, Desember 2016, pp. 267–274, 2016.

[14] Wilson, Garnett, Malcolm Heywood, "Use of a Genetic Algorithm in

Brill's Transformation-Based Part-of-Speech Tagger," Faculty of Computer Science Dalhousie University, 2005.

[15] J. Bjerva, "Genetic Algorithms in the Brill Tagger-Moving towards language independence," Department of Linguistics. Thesis. Stockholm University, 2013.

[16] A. Halim, "Politik Bahasa Nasional 1 dan 2," Jakarta: Aneka Ilmu, 1976.

[17] N. P. Dewi, Ubaidi, "Lexical Rule dan Pengaruh Penggunaan Lexicon Pada Pos Tagging Bahasa Madura," Jurnal Matrik, 18(1) pp.69-70, 2018.

[18] N. P. Dewi, Ubaidi, "POS Tagging Bahasa Madura dengan Menggunakan Algoritma Brill Tagger," Jurnal Teknologi Informasi dan Ilmu Komputer, unpublished.

[19] A. G. Ayana, "Improving Brill's Tagger Lexical and Transformation Rule for Afaan Oromo Language," PeerJ PrePrints, pp.1-11, 2015.

[20] M. Mitchell, "An Introduction to Genetic Algorithms," Cambridge, MA: MIT Press, 1996.

[21] A. E. Eiben, "Genetic algorithms with multi-parent recombination," PPSN III: Proceedings of the International Conference on Evolutionary Computation. The Third Conference on Parallel Problem Solving from Nature: 78-87, 1994.

[22] D. Goldberg, "Genetic Algorithms in Search, Optimization and Machine Learning," Reading, MA: Addison-Wesley Professional, 1989.

[23] A. Dinakaramani, F. Rashel, A. Luthfi, R. Manurung, "Designing an Indonesian Part of speech Tagset and Manually Tagged Indonesian Corpus," International Conference on Asian Language Processing (IALP), 20-22 Oktober 2014, pp. 66-69, 2014.

[24] Wicaksono, A. Farizki, A. Purwarianti, "HMM Based Part-of-Speech Tagger for Bahasa Indonesia," On Proceedings of 4th International MALINDO (Malay and Indonesian Language) Workshop, 2nd August 2010

