

**ALGORITMA JARO WINKLER UNTUK MENGUKUR SIMILARITY
BERITA ONLINE**

SKRIPSI



disusun oleh

Teguh Efriyanto

18.11.2388

**PROGRAM SARJANA
PROGRAM STUDI INFORMATIKA
FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA
2022**

**ALGORITMA JARO WINKLER UNTUK MENGUKUR SIMILARITY
BERITA ONLINE**

SKRIPSI

untuk memenuhi sebagian persyaratan
mencapai gelar Sarjana
pada Program Studi Informatika



disusun oleh

Teguh Efriyanto

18.11.2388

**PROGRAM SARJANA
PROGRAM STUDI INFORMATIKA
FAKULTAS ILMU KOMPUTER
UNIVERSITAS AMIKOM YOGYAKARTA
YOGYAKARTA
2022**

PERSETUJUAN

SKRIPSI

**ALGORITMA JARO WINKLER UNTUK MENGUKUR SIMILARITY
BERITA ONLINE**

yang dipersiapkan dan disusun oleh

Teguh Efriyanto

18.11.2388

telah disetujui oleh Dosen Pembimbing Skripsi
pada tanggal 08 Oktober 2021

Dosen Pembimbing,

Mardhiya Hayaty, S.T., M.Kom.

NIK. 190302108

PENGESAHAN

SKRIPSI

**ALGORITMA JARO WINKLER UNTUK MENGUKUR SIMILARITY
BERITA ONLINE**

yang dipersiapkan dan disusun oleh

Teguh Efriyanto

18.11.2388

telah dipertahankan di depan Dewan Penguji
pada tanggal 22 Maret 2022

Susunan Dewan Penguji

Nama Penguji

Heri Sismoro, M.Kom

NIK. 190302057

Yuli Astuti, M.Kom

NIK. 190302146

Mardhiya Hayaty, S.T., M.Kom.

NIK. 190302108

Tanda Tangan

Skripsi ini telah diterima sebagai salah satu persyaratan
untuk memperoleh gelar Sarjana Komputer
Tanggal ... Maret 2022

DEKAN FAKULTAS ILMU KOMPUTER

Hanif Al Fatta, M.Kom

NIK. 190302096

PERNYATAAN

Saya yang bertandatangan dibawah ini menyatakan bahwa, skripsi ini merupakan karya saya sendiri (ASLI), dan isi dalam skripsi ini tidak terdapat karya yang pernah diajukan oleh orang lain untuk memperoleh gelar akademis di suatu institusi pendidikan tinggi manapun, dan sepanjang pengetahuan saya juga tidak terdapat karya atau pendapat yang pernah ditulis dan/atau diterbitkan oleh orang lain, kecuali yang secara tertulis diacu dalam naskah ini dan disebutkan dalam daftar pustaka.

Segala sesuatu yang terkait dengan naskah dan karya yang telah dibuat adalah menjadi tanggungjawab saya pribadi.

Yogyakarta, 22 Maret 2022



Handwritten signature of Teguh I. Priyanto in black ink, overlaid on a purple and yellow geometric watermark.

Teguh I. Priyanto

NIM. 18.11.2388

MOTTO

“Tidak ada kesuksesan tanpa kerja keras. Tidak ada keberhasilan tanpa kebersamaan. Tidak ada kemudahan tanpa doa.”

(Ridwan Kamil)

“Sukses adalah guru yang buruk. Sukses menggoda orang yang tekun berpikir bahwa mereka tidak bisa gagal.”

(Bill Gates)

“Angin tidak berhembus untuk menggoyangkan pepohonan, melainkan menguji kekuatan akarnya.”

(Ali bin Abi Thalib)



PERSEMBAHAN

Puji syukur kehadirat Allah SWT yang telah melimpahkan rahmat, taufik dan hidayah-Nya sehingga penulis dapat menyelesaikan skripsi ini dengan baik dan lancar. Selama pengerjaan skripsi ini banyak sekali bantuan, dukungan, dan do'a dari berbagai pihak, sehingga penulis sampaikan rasa terimakasih sedalam-dalamnya kepada :

1. Kedua orang tua dan saudara-saudara yang telah memberikan do'a dan dukungan kepada penulis secara moril maupun materil sehingga skripsi ini dapat selesai dengan baik.
2. Dosen pembimbing, Ibu Mardhiya Hayaty, S.T., M.Kom. yang telah membimbing dalam penyusunan skripsi ini hingga selesai.
3. Calon istri, Rifa Indriyani yang telah senantiasa memberikan motivasi dan membantu dalam segala hal terutama telah bersedia menyediakan alat untuk saya bisa menyelesaikan skripsi ini.
4. Teman-teman Informatika angkatan 2018 khususnya kelas IF-09 yang senantiasa memberi motivasi dan berjuang bersama selama menjadi mahasiswa.
5. Semua pihak yang telah banyak membantu dalam penyusunan skripsi ini yang tidak bisa penulis sebutkan satu-persatu.

KATA PENGANTAR

Puji dan syukur kehadiran Allah SWT, karena rahmat dan hidayah-Nya sehingga skripsi yang berjudul “Algoritma Jaro Winkler Untuk Mengukur Similarity Berita Online” ini dapat selesai dengan baik sebagai salah syarat untuk dapat menempuh ujian sarjana pada Fakultas Ilmu Komputer (FIK) program studi Informatika di Universitas Amikom Yogyakarta.

Pengerjaan skripsi ini tidak lepas dari bantuan beberapa pihak. Oleh karena itu, penulis ingin menyampaikan rasa hormat dan terimakasih kepada :

1. Prof. Dr. M. Suyanto, MM. selaku Rektor Universitas Amikom Yogyakarta.
2. Ibu Windha Mega Pradnya Dhuhita, M.Kom. selaku Ketua Program Studi S1 Informatika Universitas Amikom Yogyakarta.
3. Ibu Mardhiya Hayaty, S.T., M.Kom. selaku dosen pembimbing yang telah membantu dan membimbing penulis dengan saran dan waktunya.

Penulis menyadari bahwa dalam penyusunan skripsi ini masih banyak kekurangan sehingga saran dan kritik yang membangun sangat penulis harapkan dan penulis berharap semoga skripsi ini bisa memberikan manfaat kepada para pembaca khususnya bagi penulis secara pribadi.

Yogyakarta,

Teguh Efriyanto

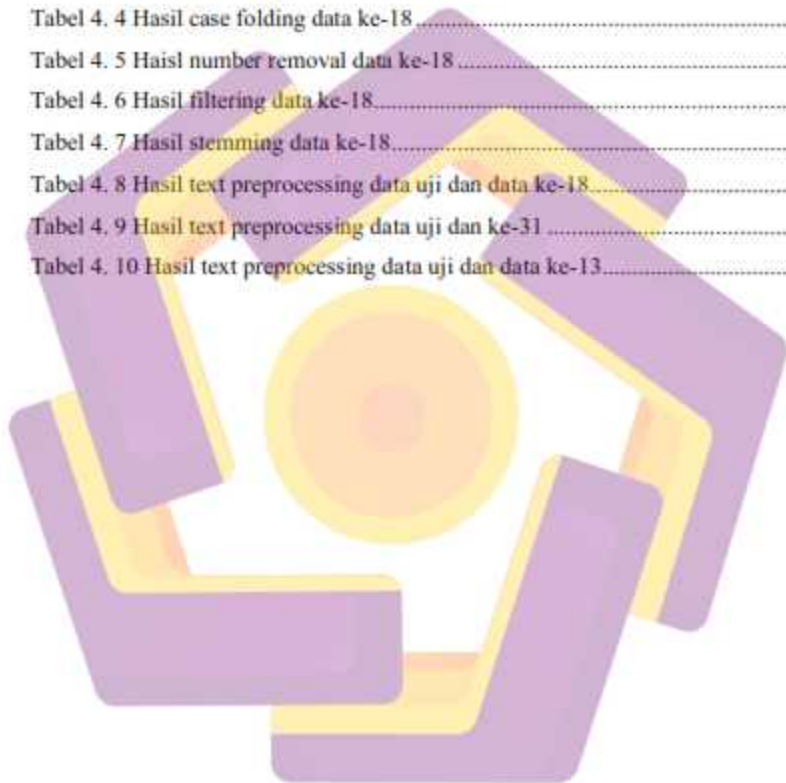
DAFTAR ISI

JUDUL.....	I
PERSETUJUAN.....	ii
PENGESAHAN.....	iii
PERNYATAAN.....	iv
MOTTO.....	v
PERSEMBAHAN.....	vi
KATA PENGANTAR.....	vii
DAFTAR ISI.....	viii
DAFTAR TABEL.....	x
DAFTAR GAMBAR.....	xi
INTISARI.....	xii
<i>ABSTRACT</i>	xiii
BAB I PENDAHULUAN.....	1
1.1 LATAR BELAKANG.....	1
1.2 RUMUSAN MASALAH.....	2
1.3 BATASAN MASALAH.....	2
1.4 MAKSUD DAN TUJUAN PENELITIAN.....	3
1.5 MANFAAT PENELITIAN.....	3
1.6 METODE PENELITIAN.....	3
1.7 SISTEMATIKA PENULISAN.....	4
BAB II LANDASAN TEORI.....	6
2.1 KAJIAN PUSTAKA.....	6
2.2 DASAR TEORI.....	8
2.2.1 <i>PLAGIARISME</i>	8
2.2.2 <i>TEXT PREPROCESSING</i>	9
2.2.3 <i>ALGORITMA JARO WINKLER</i>	10
BAB III METODE PENELITIAN.....	13
3.1 ALUR PENELITIAN.....	13

3.2	METODE PERANCANGAN SISTEM	14
3.3	DATASET	15
3.4	TEXT PREPROCESSING	16
3.5	IMPLEMENTASI ALGORITMA JARO WINKLER	17
3.6	EVALUASI	17
BAB IV HASIL DAN PEMBAHASAN		18
4.1	HASIL	18
4.2	PEMBAHASAN	21
4.2.1	DATASET	21
4.2.2	<i>TEXT PREPROCESSING</i>	23
4.2.2.1	<i>REMOVE HTML TAG</i>	23
4.2.2.2	<i>REMOVE SPECIAL CHARACTER</i>	24
4.2.2.3	<i>CASE FOLDING</i>	25
4.2.2.4	<i>NUMBER REMOVAL</i>	26
4.2.2.5	<i>FILTERING</i>	27
4.2.2.6	<i>STEMMING</i>	28
4.2.3	IMPLEMENTASI ALGORITMA JARO WINKLER	28
4.2.4	EVALUASI	31
4.2.5	TAMPILAN SISTEM	35
BAB V PENUTUP		37
5.1	KESIMPULAN	37
5.2	SARAN	37
DAFTAR PUSTAKA		38
LAMPIRAN		42

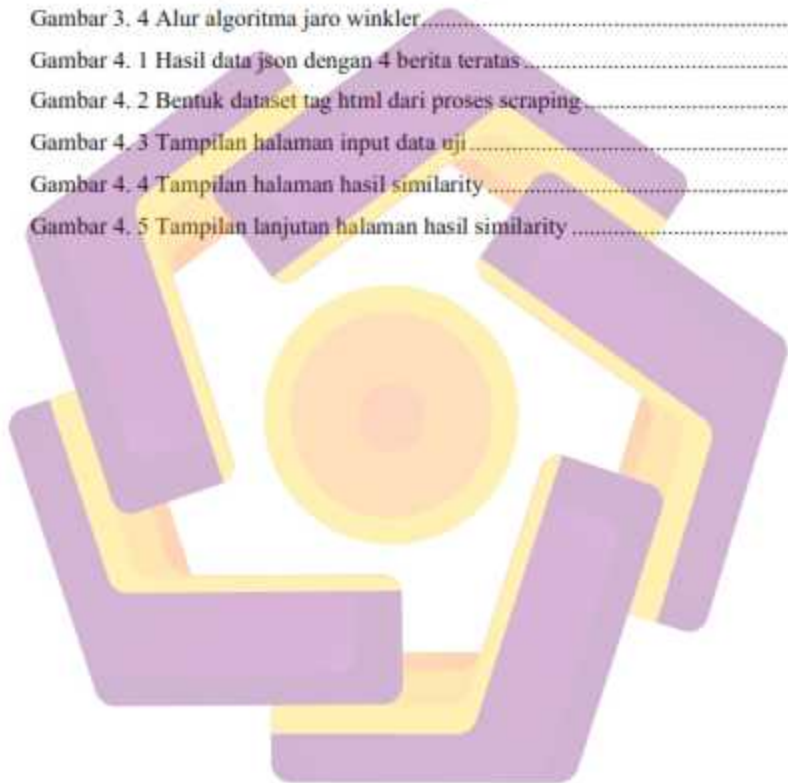
DAFTAR TABEL

Tabel 4. 1 Hasil similarity pada dataset 20 portal berita.....	18
Tabel 4. 2 Hasil remove html tag data ke-18	24
Tabel 4. 3 Hasil remove special character data ke-18.....	25
Tabel 4. 4 Hasil case folding data ke-18	25
Tabel 4. 5 Hasil number removal data ke-18	26
Tabel 4. 6 Hasil filtering data ke-18.....	27
Tabel 4. 7 Hasil stemming data ke-18.....	28
Tabel 4. 8 Hasil text preprocessing data uji dan data ke-18.....	29
Tabel 4. 9 Hasil text preprocessing data uji dan ke-31	32
Tabel 4. 10 Hasil text preprocessing data uji dan data ke-13.....	33



DAFTAR GAMBAR

Gambar 3. 1 Alur penelitian.....	13
Gambar 3. 2 Prototype Halaman Cek Similarity	15
Gambar 3. 3 Prototype Halaman Hasil Similarity Berita.....	15
Gambar 3. 4 Alur algoritma jaro winkler.....	17
Gambar 4. 1 Hasil data json dengan 4 berita teratas.....	22
Gambar 4. 2 Bentuk dataset tag html dari proses scraping.....	23
Gambar 4. 3 Tampilan halaman input data uji.....	35
Gambar 4. 4 Tampilan halaman hasil similarity.....	36
Gambar 4. 5 Tampilan lanjutan halaman hasil similarity.....	36



INTISARI

Perkembangan zaman di era industri 4.0 telah membuat media berita online menjadi sangat penting di kehidupan sehari-hari masyarakat untuk mencari sumber informasi. Hal ini berdampak pada wartawan media berita online untuk dapat mencari informasi berita yang cepat dan akurat, tidak menutup kemungkinan wartawan yang bekerja secara bersama-sama di lapangan melakukan tindakan plagiarisme ke wartawan lain atau mengambil bahan berita dari situs media berita lain dan menggunakannya untuk dimuat dimediannya tanpa mencantumkan sumbernya. sehingga untuk mengatasi hal tersebut dibutuhkan sebuah algoritma yang dapat mengukur similarity berita online untuk mengetahui tingkat plagiarisme antar berita online.

Salah satu algoritma yang dapat digunakan untuk mengetahui tingkat plagiarisme dengan menghitung nilai similarity adalah algoritma jaro winkler. Pada algoritma jaro winkler semakin tinggi nilai jaro winkler pada dua string, maka semakin mirip string tersebut. Nilai 1 menandakan kesamaan antar string dan nilai 0 menandakan ketidaksamaan antar string. Data yang digunakan adalah data isi berita dari 20 situs media berita online daerah Kalimantan Tengah, yang diperoleh dengan proses scraping yang disorting dengan Google Custom Search Engine dengan memanfaatkan Custom Search JSON API dan menggunakan keyword untuk mendapatkan berita dengan topik yang sama dan dilakukan proses text preprocessing. Dan pada metode perancangan sistem pada penelitian ini menggunakan metode prototype.

Pada penelitian ini menghasilkan nilai rata-rata similarity antar berita online sebesar 74,49% dengan data sebanyak 55 data berita, dimana diperoleh 43 data berita dengan tingkat plagiarisme berat dan 12 data berita dengan tingkat plagiarisme sedang dan terdapat kelemahan pada algoritma jaro winkler dalam menghitung nilai similarity pada data yang diperoleh, yang mana terdapat beberapa data yang tidak terdeteksi yang seharusnya tingkat plagiarismenya berat namun tidak berat dan sebaliknya.

Kata Kunci: Jaro Winkler, Similarity, Plagiarisme, Text Preprocessing, Berita Online

ABSTRACT

The development of the era in the industrial era 4.0 has made online news media very important in people's daily lives to find sources of information. This has an impact on online news media journalists to be able to find news information quickly and accurately, it is possible for journalists who work together in the field to commit plagiarism to other journalists or take news material from other news media sites and use it to be published in their media. without citing the source. so that to overcome this we need an algorithm that can measure the similarity of online news to determine the level of plagiarism between online news.

One of the algorithms that can be used to determine the level of plagiarism by calculating the similarity value is the Jaro Winkler algorithm. In the Jaro Winkler algorithm, the higher the Jaro Winkler value on two strings, the more similar the strings are. The value 1 indicates the similarity between the strings and the value 0 indicates the inequality between the strings. The data used is news content data from 20 online news media sites in the Central Kalimantan area, which is obtained by a scraping process which is sorted by Google Custom Search Engine by utilizing the Custom Search JSON API and using keywords to get news with the same topic and text preprocessing is carried out. . And the system design method in this study uses the prototype method.

In this study, the average value of similarity between online news was 74.49% with 55 news data, which obtained 43 news data with severe plagiarism levels and 12 news data with moderate plagiarism levels and there were weaknesses in the Jaro Winkler algorithm in calculating the similarity value in the data obtained, where there are some undetected data which should have a heavy plagiarism level but not heavy and vice versa.

Keyword: *Jaro Winkler, Similarity, Plagiarism, Text Preprocessing, Online News*