

**PEMILIHAN KOMBINASI FITUR UNTUK *AUTHOR*
VERIFICATION DENGAN METODE *CLUSTERING K-MEANS*
PADA TEKS BERBAHASA INDONESIA**

Skripsi



Diajukan oleh:

THOMAS WIDIARYA BUDIMAN

71160018

PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNOLOGI INFORMASI
UNIVERSITAS KRISTEN DUTA WACANA
YOGYAKARTA

2020

**PEMILIHAN KOMBINASI FITUR UNTUK *AUTHOR*
VERIFICATION DENGAN METODE *CLUSTERING K-MEANS*
PADA TEKS BERBAHASA INDONESIA**

Skripsi



Diajukan kepada Fakultas Teknologi Informasi Program Studi Informatika
Universitas Kristen Duta Wacana
Sebagai salah satu syarat dalam memperoleh gelar Sarjana Komputer

Diajukan oleh:

THOMAS WIDIARYA BUDIMAN

71160018

PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNOLOGI INFORMASI
UNIVERSITAS KRISTEN DUTA WACANA
YOGYAKARTA

2020

HALAMAN PERNYATAAN PERSETUJUAN PUBLIKASI
SKRIPSI/TESIS/DISERTASI UNTUK KEPENTINGAN AKADEMIS

Sebagai sivitas akademika Universitas Kristen Duta Wacana, saya yang bertanda tangan di bawah ini:

Nama : Thomas Widiarya Budiman
NIM : 71160018
Program studi : Informatika
Fakultas : Teknologi Informasi
Jenis Karya : Skripsi

demi pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada Universitas Kristen Duta Wacana **Hak Bebas Royalti Noneksklusif** (*None-exclusive Royalty Free Right*) atas karya ilmiah saya yang berjudul:

**“PEMILIHAN KOMBINASI FITUR UNTUK AUTHOR VERIFICATION
DENGAN METODE CLUSTERING K-MEANS PADA TEKS BERBAHASA
INDONESIA”**

beserta perangkat yang ada (jika diperlukan). Dengan Hak Bebas Royalti/Noneksklusif ini Universitas Kristen Duta Wacana berhak menyimpan, mengalih media/formatkan, mengelola dalam bentuk pangkalan data (*database*), merawat dan mempublikasikan tugas akhir saya selama tetap mencantumkan nama kami sebagai penulis/pencipta dan sebagai pemilik Hak Cipta.

Demikian pernyataan ini saya buat dengan sebenarnya.

Dibuat di : Yogyakarta
Pada Tanggal : 12 April 2020

Yang menyatakan



Thomas Widiarya Budiman
NIM.71160018

PERNYATAAN KEASLIAN SKRIPSI

Saya menyatakan dengan sesungguhnya bahwa skripsi dengan judul:

PEMILIHAN KOMBINASI FITUR UNTUK AUTHOR VERIFICATION DENGAN METODE CLUSTERING K-MEANS PADA TEKS BERBAHASA INDONESIA

yang saya kerjakan untuk melengkapi sebagian persyaratan menjadi Sarjana Komputer pada pendidikan Sarjana Program Studi Informatika Fakultas Teknologi Informasi Universitas Kristen Duta Wacana, bukan merupakan tiruan atau duplikasi dari skripsi kesarjanaan di lingkungan Universitas Kristen Duta Wacana maupun di Perguruan Tinggi atau instansi manapun, kecuali bagian yang sumber informasinya dicantumkan sebagaimana mestinya.

Jika dikemudian hari didapati bahwa hasil skripsi ini adalah hasil plagiasi atau tiruan dari skripsi lain, saya bersedia dikenai sanksi yakni pencabutan gelar kesarjanaan saya.

Yogyakarta, 10 Juli 2020



THOMAS WIDIARYA BUDIMAN
71160018

HALAMAN PERSETUJUAN

Judul : Pemilihan Kombinasi Fitur untuk *Author Verification*
dengan Metode *Clustering K-Means* pada Teks
Berbahasa Indonesia

Nama : Thomas Widiarya Budiman

NIM : 71160018

Mata Kuliah : Skripsi

Kode : TI0366

Semester : Genap

Tahun akademik : 2019/2020

Telah diperiksa dan disetujui
Di Yogyakarta,
Pada Tanggal 7 Juli 2020

Dosen Pembimbing I



Dr. Phil. Lucia Dwi K., SS., M.A.

Dosen Pembimbing II



Laurentius Kuncoro P. S., S.T., M.Eng

HALAMAN PENGESAHAN

PEMILIHAN KOMBINASI FITUR UNTUK AUTHOR VERIFICATION DENGAN METODE CLUSTERING K-MEANS PADA TEKS BERBAHASA INDONESIA

Oleh: THOMAS WIDIARYA BUDIMAN / 71160018

Dipertahankan di depan Dewan Penguji Skripsi
Program Studi Informatika Fakultas Teknologi Informasi
Universitas Kristen Duta Wacana - Yogyakarta
Dan dinyatakan diterima untuk memenuhi salah satu syarat memperoleh gelar
Sarjana Komputer
pada tanggal 23 Juli 2020

Yogyakarta, 28 Juli 2020
Mengesahkan,

Dewan Penguji:

1. Lucia Dwi Krisnawati, Dr. Phil.
2. Laurentius Kuncoro Probo Saputra, S.T.,
M.Eng.
3. Willy Sudiarto Raharjo, S.Kom., M.Cs.
4. Aditya Wikan Mahastama, S.Kom., M.Cs.



Dekan

(Restyandito, S.Kom., MSIS., Ph.D.)

Ketua Program Studi

(Gloria Virginia, Ph.D.)

UCAPAN TERIMA KASIH

Puji syukur senantiasa peneliti panjatkan kepada Tuhan Yang Maha Esa atas berkat, penyertaan, dan rahmat-Nya, sehingga peneliti bisa menyelesaikan skripsi. Peneliti juga berterimakasih atas dukungan, saran dan bimbingan selama proses pengerjaan skripsi ini kepada:

1. Bapak Restyandito, S.Kom., MSIS, Ph.D. selaku Fakultas Teknologi Informasi UKDW.
2. Ibu Gloria Virginia, S.Kom., MAI, Ph.D. selaku Ketua Program Studi Informatika Fakultas Teknologi Informasi UKDW.
3. Ibu Dr. Phil. Lucia Dwi K.,SS., M.A. selaku dosen pembimbing I yang bersedia membantu dan memberikan arahan serta dukungan dari awal perencanaan penelitian hingga skripsi ini selesai. Saya sangat berterimakasih karena sudah menyempatkan waktu untuk konsultasi secara luring walaupun sedang terjadi pandemic Covid-19.
4. Bapak Laurentius Kuncoro Probo Saputro, S.T., M.Eng. selaku dosen pembimbing II yang bersedia membantu dan memberikan arahan serta dukungan dari awal perencanaan penelitian hingga skripsi ini selesai.
5. Seluruh dosen di Fakultas Teknologi Informasi, khususnya dosen Program Studi Informatika yang telah memberikan bekal pengetahuan selama penulis menempuh pendidikan di Universitas Kristen Duta Wacana.

Peneliti berharap semoga pengorbanan dan segala sesuatunya yang dengan tulus dan ikhlas telah diberikan akan selalu mendapat limpahan rahmat-Nya.

INTISARI

Pemilihan Kombinasi Fitur untuk *Author Verification* dengan Metode *Clustering K-Means* pada Teks Berbahasa Indonesia

Penipuan dengan menggunakan identitas orang lain umumnya dilakukan melalui media tertulis, karena pelaku tidak perlu memperlihatkan fisik maupun suaranya. Maka dari itu diperlukan sistem yang bisa membedakan identitas tulisan seseorang. Identitas tulisan seseorang bisa dianalisa menggunakan fitur *stylometry*, fitur ini merupakan salah satu faktor sistem verifikasi penulis bisa berjalan dengan baik. Penelitian ini mengasumsikan jika kombinasi fitur *stylometry* berhasil merepresentasikan gaya penulisan seseorang maka proses klasterisasi atau klasifikasi dokumen berdasarkan gaya penulisan seseorang akan menghasilkan nilai evaluasi klasterisasi (*purity*) dan klasifikasi (akurasi, presisi, sensitivitas, *FScore*) yang memuaskan. Hasil penelitian dengan dokumen berbahasa Indonesia pada *klasifier* MKNN, KNN, dan SVM menunjukkan bahwa nilai *purity* bisa menggambarkan nilai evaluasi klasifikasi walaupun tidak terlalu tepat, kombinasi fitur mendapatkan nilai *purity* kurang dari 0,5 memiliki kemungkinan besar untuk mendapatkan nilai evaluasi klasifikasi kurang dari 0,5 juga, hal ini bisa menghemat proses jika kombinasi fitur dengan nilai *purity* kurang dari 0,5 tidak digunakan untuk proses klasifikasi. Kombinasi fitur yang mengandung fitur frekuensi relatif tanda baca (Fitur Sintaksis), frekuensi relatif stopword (Fitur Sintaksis) cenderung memiliki nilai evaluasi klasifikasi maupun klasterisasi yang lebih baik dari pada kombinasi fitur yang tidak mengandung kedua fitur tersebut dan ketika ditambah dengan fitur rata-rata panjang paragraf (Fitur Struktural) kombinasi fitur ini menjadi kombinasi fitur terbaik pada semua *klasifier*.

Kata kunci: identitas tulisan seseorang, verifikasi penulis, fitur *stylometry*

ABSTRACT

Selection of Feature Combinations for Author Verification with K-Means Clustering Method in Indonesian Language Text

Fraud by using other people's identities is generally done through written media, because the perpetrators do not need to show physical or voice. Therefore a system that can distinguish someone's writing identity is needed. The identity of a person's writing can be analyzed using the stylometry feature, this feature is one of the factors the author's verification system can work well. This research assumes that if a combination of stylometry features successfully represents a person's writing style, the clustering or classification of documents based on one's writing style will result in a satisfactory classification metric (purity) and classification metric (accuracy, precision, sensitivity, FScore). The results of research with Indonesian language documents in the MKNN, KNN, and SVM classifiers show that the value of purity can describe the value of classification evaluation even though it is not very precise, a combination of features getting a purity value of less than 0.5 has a high likelihood to get a classification evaluation value of less than 0.5 also, this can save the process if a combination of features with a purity value of less than 0.5 is not used for the classification process. Combinations of features that contain punctuation relative frequency features (Syntactic Features), stopword relative frequencies (Syntax Syntactic) tend to have better classification and classification evaluation values than feature combinations that do not contain both features and when added to the average length feature paragraph (Structural Features) this combination of features becomes the best combination of features in all classifiers.

keywrods: author writing identity, author verification, stylometry feature

DAFTAR ISI

PERNYATAAN KEASLIAN SKRIPSI.....	iii
HALAMAN PERSETUJUAN.....	iv
HALAMAN PENGESAHAN.....	v
UCAPAN TERIMA KASIH.....	vi
INTISARI.....	vii
<i>ABSTRACT</i>	viii
DAFTAR ISI.....	ix
DAFTAR GAMBAR.....	xiii
DAFTAR TABEL.....	xvii
1. BAB I.....	1
1.1 Latar Belakang.....	1
1.2 Rumusan Masalah.....	2
1.3 Batasan Masalah.....	2
1.4 Tujuan Penelitian.....	2
1.5 Manfaat Penelitian.....	3
1.6 Metodologi Penelitian.....	3
1.7 Sistematika Penulisan.....	4
2. BAB II.....	6
2.1 Tinjauan Pustaka.....	6
2.2 Landasan Teori.....	7
2.2.1 Verifikasi Penulis.....	7
2.2.2 Fitur Gaya Penulisan (<i>Stylometry</i>).....	7
2.2.3 Ekstraksi Fitur.....	8
2.2.4 <i>Stopword</i>	10
2.2.5 Normalisasi Skala Nilai.....	10
2.2.6 K-Means.....	10
2.2.7 <i>Euclidean Distance</i>	11
2.2.8 Evaluasi <i>Purity Klasterisasi</i>	11
2.2.9 <i>Modified K-Nearest Neighbor (MKNN)</i>	12

2.2.10	<i>Confusion Matrix</i>	13
2.2.11	<i>KFold Cross Validasi</i>	13
2.2.12	Evaluasi Klasifikasi.....	15
3.	BAB III	17
3.1	Perancangan Kebutuhan Sistem	17
3.1.1	Spesifikasi Perangkat Keras yang Digunakan.....	17
3.1.2	Kebutuhan Perangkat Lunak.....	17
3.2	Perancangan Sistem.....	18
3.2.1	Perancangan Pengumpulan Data.....	18
3.2.2	Perancangan Penyimpanan Data	19
3.2.3	Perancangan Prapemrosesan <i>Case Folding</i>	20
3.2.4	Perancangan Ekstraksi Fitur	20
3.2.5	Perancangan Kombinasi Fitur.....	20
3.2.6	Perancangan Prapemrosesan Normalisasi Nilai.....	21
3.2.7	Perancangan Klasterisasi.....	22
3.2.8	Perancangan Evaluasi Klasterisasi.....	22
3.2.9	Perancangan Pembagian Data Latih dan Data Uji	22
3.2.10	Perancangan Klasifikasi	22
3.2.11	Perancangan Evaluasi Klasifikasi	22
3.3	Desain Antarmuka	23
3.4.1	Blok Diagram Sistem User	23
3.4.2	Blok Diagram Sistem Admin.....	24
3.4.3	Mode User.....	26
3.4.1	Mode Admin	28
3.4	Perancangan Pengujian.....	29
3.4.1	Skenario Pengujian pada Percobaan Pertama	30
3.4.2	Skenario Pengujian pada Percobaan Kedua	32
3.4.3	Skenario Pengujian pada Percobaan Ketiga.....	32
4.	BAB IV	35
4.1	Implementasi Sistem	35
4.1.1	Implementasi Pengumpulan Data	35

4.1.2	Implementasi Penyimpanan Data.....	37
4.1.3	Implementasi Prapemrosesan <i>Case Folding</i>	38
4.1.4	Implementasi Ekstraksi Fitur	38
4.1.5	Implementasi Kombinasi Fitur.....	45
4.1.6	Implementasi Pembagian Dataset	46
4.1.7	Implementasi Prapemrosesan Normalisasi Nilai Kombinasi Fitur	48
4.1.8	Implementasi Klasterisasi	49
4.1.9	Implementasi Evaluasi Klasterisasi.....	50
4.1.10	Implementasi Klasifikasi.....	50
4.1.11	Implementasi Evaluasi Klasifikasi.....	52
4.2	Implementasi Antarmuka Sistem	53
4.2.1	Mode User.....	53
4.2.2	Mode Admin	56
4.3	Hasil dan Analisis Pengujian.....	64
4.3.1	Hasil dan Analisis Pengujian pada Percobaan Pertama	64
4.3.2	Hasil dan Analisis Pengujian pada Percobaan Kedua.....	69
4.3.3	Hasil dan Analisis Pengujian pada Percobaan Ketiga.....	70
4.3.4	Hasil dan Analisis Pengujian pada Percobaan Keempat.....	74
4.3.5	Analisis dari Empat Percobaan	81
5.	BAB V.....	88
5.1	Kesimpulan.....	88
5.2	Saran.....	89
	Daftar Pustaka	90
	LAMPIRAN A.....	1
A.1.	Tabel Evaluasi pada Pengujian Percobaan Pertama.....	1
A.2.	Tabel Evaluasi pada Pengujian Percobaan Ketiga	23
A.3.	Tabel Evaluasi pada Pengujian Percobaan Keempat dengan <i>Klasifier</i> KNN.....	28
A.4.	Tabel Evaluasi pada Pengujian Percobaan Keempat dengan <i>Klasifier</i> SVM.....	50
	LAMPIRAN B	1

B.1.	MKNN.py	1
B.2.	urls.py.....	2
B.3.	views.py.....	3
B.4.	utils.py.....	13
B.5.	feature_extraction.py	14
B.6.	author_verification.py	15
B.7.	base.html.....	19
B.8.	admin.html.....	20
B.9.	dokumen.html.....	26
B.10.	percobaan.html	29
B.11.	user.html	45
B.12.	index.css	46
B.13.	admin.js	51
B.14.	dokumen.js	56
B.15.	latih_dokumen.js	60
B.16.	user.js.....	62
LAMPIRAN C		1
C.1.	Kartu Konsultasi dengan Bu Lucia.....	1
C.2.	Kartu Konsultasi dengan Pak Kuncoro	2
C.3.	Formulir Perbaikan (Revisi) Skripsi.....	3

DAFTAR GAMBAR

Gambar 2.1 Contoh Kasus Purity.....	11
Gambar 2.2 Contoh Cara Kerja <i>KFold Cross</i> Validasi dengan k adalah 5	14
Gambar 2.3 Contoh Cara Kerja <i>StratifiedKFold Cross</i> Validasi dengan k adalah 5	15
Gambar 3.1 Struktur Penyimpanan Data	19
Gambar 3.2 Alur Kerja Sistem Berdasarkan Skema Penggunaan User.....	24
Gambar 3.3 Alur Penginputan Data Website dengan Skema Penggunaan Admin25	
Gambar 3.4 Alur Pemilihan Kombinasi Fitur Serta Pelatihan Model Klasifikasi	26
Gambar 3.5 Desain Website pada Mode User (Tampilan Awal)	27
Gambar 3.6 Desain Website pada Mode User (Keluaran Sistem Benar)	27
Gambar 3.7 Desain Website pada Mode User (Keluaran Sistem Salah)	28
Gambar 3.8 Desain Website pada Mode Admin (Tampilan Awal)	28
Gambar 3.9 Desain Website pada Mode Admin (Input Data).....	29
Gambar 3.10 Desain Website pada Mode Admin (Output Sistem dan Pemilihan Kombinasi Fitur dan K pada MKNN).....	29
Gambar 3.11 Alur Pengujian Percobaan Pertama	30
Gambar 3.12 Alur Pengujian Percobaan Kedua	32
Gambar 3.13 Alur Pengujian Percobaan Ketiga	33
Gambar 4.1 Penggunaan Pandas untuk Penyimpanan Sementara	37
Gambar 4.2 Contoh Data Sebelum <i>Case Folding</i>	38
Gambar 4.3 Contoh Data sesudah <i>Case Folding</i>	38
Gambar 4.4 Contoh Data Tulisan.....	39
Gambar 4.5 String <i>punctuation</i>	39
Gambar 4.6 Contoh Kasus Perhitungan Frekuensi Relatif Tanda Baca	40
Gambar 4.7 Contoh Cara Menghitung Frekuensi Relatif Tanda Baca	40
Gambar 4.8 Contoh Hasil Ekstraksi Fitur Frekuensi Relatif Tanda Baca	40
Gambar 4.9 Stopword dari Sastrawi	41
Gambar 4.10 Contoh Kata Stopword dari Data Tulisan	41
Gambar 4.11 Contoh Perhitungan Fitur Frekuensi Relatif Stopword.....	42

Gambar 4.12 Contoh Hasil Ekstraksi Fitur Frekuensi Relatif Stopword.....	42
Gambar 4.13 Langkah Pertama Ekstraksi Fitur Rata-Rata Panjang Kalimat	43
Gambar 4.14 Langkah Kedua Ekstraksi Fitur Rata-Rata Panjang Kalimat	43
Gambar 4.15 Langkah Pertama Ekstraksi Fitur Rata-Rata Panjang Paragraf.....	43
Gambar 4.16 Langkah Kedua Ekstraksi Fitur Rata-Rata Panjang Paragraf	44
Gambar 4.17 Contoh Token	44
Gambar 4.18 Contoh Type.....	44
Gambar 4.19 Contoh <i>StratifiedKfold</i> dengan persebaran kelas merata	46
Gambar 4.20 Contoh Hasil <i>StratifiedKfold</i> dengan persebaran kelas merata.....	46
Gambar 4.21 Contoh <i>StratifiedKfold</i> dengan persebaran kelas tidak merata	46
Gambar 4.22 Contoh Hasil <i>StratifiedKfold</i> dengan persebaran kelas tidak merata	46
Gambar 4.23 Dataframe Data Latih untuk Percobaan Kedua dan Ketiga	47
Gambar 4.24 Dataframe Data Latih untuk Percobaan Kedua dan Ketiga	48
Gambar 4.25 Normaliasi nilai dengan <i>MinMaxScaler</i> untuk klusterisasi	48
Gambar 4.26 Normaliasi nilai dengan <i>MinMaxScaler</i> untuk klasifikasi	49
Gambar 4.27 Kode Program untuk Membuat Objek KMeans.....	49
Gambar 4.28 Kode Program untuk Melakukan Proses Klusterisasi	49
Gambar 4.29 Algoritma Fungsi Perhitungan Purity	50
Gambar 4.30 <i>Pseudocode MKNN</i>	50
Gambar 4.31 Kode Program untuk Membuat Objek MKNN	51
Gambar 4.32 Kode Program untuk Melakukan Proses Pelatihan MKNN.....	51
Gambar 4.33 Kode Program untuk Melakukan Klasifikasi pada MKNN	51
Gambar 4.34 Kode Program untuk Menyimpan Objek MKNN ke dalam File	52
Gambar 4.35 Kode Program untuk Evaluasi Akurasi Klasifikasi pada MKNN...	52
Gambar 4.36 Kode Program untuk Evaluasi Presisi Klasifikasi pada MKNN.....	52
Gambar 4.37 Kode Program untuk Evaluasi Sensitivitas Klasifikasi pada MKNN	52
Gambar 4.38 Kode Program untuk Evaluasi Fscore Klasifikasi pada MKNN.....	52
Gambar 4.39 Halaman Awal Mode User	53
Gambar 4.40 Tampilan Error Halaman Awal Mode User 1	54

Gambar 4.41 Tampilan Error Halaman Awal Mode User 2	54
Gambar 4.42 Tampilan Error Halaman Awal Mode User 3	55
Gambar 4.43 Tampilan Ketika Output dan Input Penulis Dokumen Cocok.....	55
Gambar 4.44 Tampilan Ketika Output dan Input Penulis Dokumen Tidak Cocok	56
Gambar 4.45 Halaman Awal Mode Admin (Model Belum Dilatih)	57
Gambar 4.46 Halaman Awal Mode Admin (Model Sudah Dilatih).....	57
Gambar 4.47 Antarmuka Tambah Data Mode Admin.....	58
Gambar 4.48 Antarmuka Hapus Penulis Mode Admin	58
Gambar 4.49 Antarmuka Edit Penulis Mode Admin	59
Gambar 4.50 Halaman RUD Dokumen Mode Admin.....	60
Gambar 4.51 Aksi dari Tombol Show More Mode Admin 1	61
Gambar 4.52 Aksi dari Tombol Show More Mode Admin 2	61
Gambar 4.53 Antarmuka Edit Dokumen Mode Admin 1	62
Gambar 4.54 Antarmuka Edit Dokumen Mode Admin 2.....	62
Gambar 4.55 Antarmuka Hapus Dokumen Mode Admin	63
Gambar 4.56 Antarmuka Tampilan Kombinasi Fitur	63
Gambar 4.57 Antarmuka Tampilan Input Pelatihan Model.....	64
Gambar 4.58 Visualisasi Nilai Evaluasi Akurasi dari Pengujian Percobaan Pertama.....	65
Gambar 4.59 Visualisasi Nilai Evaluasi Presisi dari Pengujian Percobaan Pertama	65
Gambar 4.60 Visualisasi Nilai Evaluasi Sensitivitas dari Pengujian Percobaan Pertama.....	66
Gambar 4.61 Visualisasi Nilai Evaluasi <i>FScore</i> dari Pengujian Percobaan Pertama	66
Gambar 4.62 Hasil Visualisasi pada Masing-Masing Nilai Evaluasi di Percobaan Kedua	69
Gambar 4.63 Visualisasi Nilai Evaluasi Akurasi dari Pengujian Percobaan Ketiga	70

Gambar 4.64 Visualisasi Nilai Evaluasi Presisi dari Pengujian Percobaan Ketiga	71
Gambar 4.65 Visualisasi Nilai Evaluasi Sensitivitas dari Pengujian Percobaan Ketiga	71
Gambar 4.66 Visualisasi Nilai Evaluasi <i>Fscore</i> dari Pengujian Percobaan Ketiga	72
Gambar 4.67 Visualisasi Nilai Evaluasi Akurasi dari Pengujian Percobaan Keempat dengan <i>Klasifier</i> KNN	75
Gambar 4.68 Visualisasi Nilai Evaluasi Presisi dari Pengujian Percobaan Keempat dengan <i>Klasifier</i> KNN	75
Gambar 4.69 Visualisasi Nilai Evaluasi Sensitivitas dari Pengujian Percobaan Keempat dengan <i>Klasifier</i> KNN	76
Gambar 4.70 Visualisasi Nilai Evaluasi <i>FScore</i> dari Pengujian Percobaan Keempat dengan <i>Klasifier</i> KNN	76
Gambar 4.71 Visualisasi Nilai Evaluasi Akurasi dari Pengujian Percobaan Keempat dengan <i>Klasifier</i> SVM	79
Gambar 4.72 Visualisasi Nilai Evaluasi Presisi dari Pengujian Percobaan Keempat dengan <i>Klasifier</i> SVM	79
Gambar 4.73 Visualisasi Nilai Evaluasi Presisi dari Pengujian Percobaan Keempat dengan <i>Klasifier</i> SVM	80
Gambar 4.74 Visualisasi Nilai Evaluasi <i>FScore</i> dari Pengujian Percobaan Keempat dengan <i>Klasifier</i> SVM	80
Gambar 4.75 Contoh Data Pada Masalah MKNN	85
Gambar 4.76 Visualisasi Contoh Data Pada Masalah MKNN	85
Gambar 4.77 Contoh Nilai Validasi MKNN pada H adalah 1 pada Penjelasan Masalah Performa MKNN	86
Gambar 4.78 Contoh Nilai Validasi MKNN pada H adalah 3 pada Penjelasan Masalah Performa MKNN	86

DAFTAR TABEL

Tabel 2.1 Contoh Fitur Gaya Penulisan Perkategori.....	8
Tabel 2.2 Kode/Notasi Fitur yang Terdapat Penelitian ini	8
Tabel 2.3 Contoh <i>Confusion Matrix</i>	13
Tabel 3.1 Spesifikasi Perangkat Keras yang Digunakan	17
Tabel 3.2 Kebutuhan Perangkat Lunak.....	17
Tabel 3.3 Contoh Data yang di Representasikan ke Ascii	18
Tabel 3.4 Format Penyimpanan Data Excel.....	19
Tabel 3.5 Contoh Tahap Case Folding.....	20
Tabel 3.6 Kombinasi Fitur <i>Stylometry</i> pada Penelitian Ini	21
Tabel 3.7 Tabel Evaluasi untuk Percobaan Pertama.....	31
Tabel 3.8 Tabel Voting untuk Percobaan.....	31
Tabel 3.9 Tabel Rangkuman Hasil Voting.....	31
Tabel 3.10 Tabel Evaluasi untuk Percobaan Kedua.....	32
Tabel 3.11 Tabel Evaluasi untuk Percobaan Ketiga	34
Tabel 4.1 Infomari dari Hasil Pengumpulan Data	35
Tabel 4.2 Persebaran Kategori di Seluruh Data	35
Tabel 4.3 Persebaran Kategori di Tiap-Tiap Penulis	36
Tabel 4.4 Informasi Bentuk, Ukuran, dan Isi dari Kombinasi Fitur <i>stylometry</i>	45
Tabel 4.5 Hasil Tabel Voting pada Percobaan Pertama.....	67
Tabel 4.6 Hasil Pengurutan Tabel Voting pada Percobaan Pertama	68
Tabel 4.7 Hasil Tabel Voting pada Percobaan Ketiga	72
Tabel 4.8 Hasil Rangkuman Tabel Voting pada Percobaan Ketiga.....	73
Tabel 4.9 Hasil Tabel Voting pada Percobaan Keempat dengan <i>Klasifier</i> KNN .	77
Tabel 4.10 Hasil Rangkuman Tabel Voting pada Percobaan dengan <i>Klasifier</i> KNN	78
Tabel 4.11 Hasil Pengurutan Kombinasi Fitur terbaik dari Pengujian Percobaan dengan <i>Klasifier</i> SVM	81
Tabel 4.12 Probabilitas Hasil Evaluasi Klasifikasi pada Nilai Purity diatas dan dibawah 0,5 pada percobaan Pertama	83

Tabel 4.13 Probabilitas Hasil Evaluasi Klasifikasi pada Nilai Purity diatas dan dibawah 0,5 pada Percobaan Ketiga	83
Tabel 4.14 Probabilitas Hasil Evaluasi Klasifikasi pada Nilai Purity diatas dan dibawah 0,5 pada Percobaan Keempat dengan <i>Klasifier</i> KNN	83
Tabel 4.15 Probabilitas Hasil Evaluasi Klasifikasi pada Nilai Purity diatas dan dibawah 0,5 pada Percobaan Keempat dengan <i>Klasifier</i> SVM	84

©UKDW

BAB I PENDAHULUAN

1.1 Latar Belakang

Penggunaan identitas palsu atau identitas orang lain merupakan sebuah tindakan yang sudah tidak asing lagi, hal ini dikarenakan kemudahan seseorang untuk mendapatkan identitas orang lain melalui dunia maya. Karena kemudahan ini banyak tindakan kriminal yang bisa saja terjadi di dunia maya dengan menggunakan identitas orang lain. Contoh kasus tindakan kriminal yang dapat dilakukan dengan menggunakan identitas orang lain adalah kasus penipuan melalui aplikasi sosial media dengan menggunakan identitas palsu (Amiruddin, 2019)

Penipuan dengan menggunakan identitas orang lain umumnya dilakukan melalui media tertulis. Media tertulis dianggap aman karena pelaku tidak perlu memperlihatkan fisik maupun suaranya. Namun faktanya gaya penulisan juga dapat dijadikan sebagai penanda identitas tulisan seseorang. J. K. Rowling yang menggunakan nama samaran Robert Galbraith ketika menulis novel *The Cuckoo's Calling* terungkap identitas tulisannya ketika gaya penulisannya dianalisis dengan menggunakan sistem verifikasi penulis (*Author Verification*) oleh Juloa (2013). Sistem verifikasi penulis merupakan sistem yang memiliki tugas untuk menjawab pertanyaan ‘apakah penulis x yang menulis dokumen ini?’ (Luyckx & Daelemans, 2008). Salah satu faktor yang membuat sistem verifikasi penulis ini bekerja dengan baik adalah fitur gaya penulisan *stylometry* yang digunakan, fitur *stylometry* yang bisa merepresentasikan gaya penulisan seseorang akan membuat sistem verifikasi penulis bekerja lebih baik.

Berangkat dari kasus J. K. Rowling, penelitian ini berfokus membangun sistem verifikasi penulis pada teks berbahasa Indonesia dengan melakukan uji coba berbagai kombinasi fitur gaya penulisan (*stylometry*) dan ingin melihat hubungan antara nilai purity (hasil klasterisasi dengan KMeans) dengan nilai akurasi, presisi, sensitivitas dan *Fscore* (hasil klasifikasi dengan MKNN) pada masing-masing kombinasi fitur yang di teliti. Penelitian ini mengasumsikan jika kombinasi fitur *stylometry* yang dipilih merupakan kombinasi fitur yang merepresentasikan gaya

penulisan seseorang maka proses klasterisasi maupun klasifikasi dokumen berdasarkan penulis akan menghasilkan nilai evaluasi yang memuaskan.

1.2 Rumusan Masalah

Melihat latar belakang yang ada, berikut adalah rumusan masalah dalam penelitian ini:

- 1.2.1 Kombinasi fitur gaya penulisan (*stylometry*) apa saja yang dapat digunakan untuk memverifikasi penulis dalam dokumen berbahasa Indonesia, sehingga ketika dilakukan klasterisasi dengan K-means, dokumen berbahasa Indonesia ini akan terklaster dengan tepat berdasarkan gaya penulisan seseorang ?
- 1.2.2 Apakah kombinasi fitur gaya penulisan (*stylometry*) memiliki nilai purity yang berbanding lurus dengan nilai metode evaluasi klasifikasi pada metode MKNN ?

1.3 Batasan Masalah

Dalam penelitian ini terdapat batasan-batasan masalah seperti berikut:

- a Sumber dokumen berbentuk opini dan diambil dari website Seward.com yang di publikasikan pada tahun 2019 sampai 2020.
- b Dokumen yang digunakan merupakan dokumen berbahasa Indonesia.
- c Setiap dokumen hanya ditulis oleh 1 penulis.
- d Jumlah penulis 11 orang dengan 25 dokumen masing-masing penulis.

1.4 Tujuan Penelitian

Penelitian ini bertujuan untuk menemukan kombinasi fitur yang dapat merepresentasikan gaya penulisan seseorang pada dokumen berbahasa Indonesia dan untuk melihat hubungan nilai purity dengan nilai f1score yang dimiliki oleh masing-masing kombinasi fitur gaya penulisan.

1.5 Manfaat Penelitian

Manfaat yang ingin didapat dari penelitian ini adalah sebagai berikut:

- 1.5.1 Penelitian ini diharapkan mampu mendorong terbentuknya sistem verifikasi penulis yang dapat membantu pemerintah Indonesia untuk mengatasi kejahatan di dunia maya.
- 1.5.2 Diharapkan sistem yang dibangun ini dapat membantu penelitian di bidang teks forensik selanjutnya untuk mendapatkan kombinasi fitur yang dapat merepresentasikan gaya penulisan seseorang sehingga proses pengerjaan penelitian selanjutnya dapat dipercepat dengan adanya penelitian ini.

1.6 Metodologi Penelitian

Berikut merupakan metode-metode yang peneliti lakukan untuk memecahkan permasalahan yang peneliti tulis di latar belakang.

1.6.1 Studi Pustaka

Tahap ini adalah merupakan tahap awal dari penelitian yang dilakukan oleh peneliti, dimana peneliti membaca jurnal dan situs-situs yang mengangkat tentang verifikasi penulis.

1.6.2 Pengumpulan Data

Tahap selanjutnya adalah tahap pengumpulan data, pengumpulan data dilakukan secara semi otomatis, yaitu dengan bantuan Bahasa pemrograman python. Mula-mula peneliti mengambil kunci API yang dimiliki oleh website sword.com, selanjutnya peneliti mengambil beberapa *username* dari beberapa penulis sword.com, lalu peneliti membuat script untuk menembak API yang disediakan oleh sword.com untuk mengambil data tulisan dari masing-masing penulis berdasarkan *username* yang dipilih oleh peneliti. Data yang diambil merupakan data yang berbentuk opini, berbahasa Indonesia, bisa saja memiliki topik yang berbeda dan ditulis oleh penulis tunggal. Data yang diambil adalah data memiliki genre yang sama yaitu opini karena akan lebih susah untuk melakukan verifikasi penulis pada banyak genre (Kestemont, Luyckx, Daelemans, & Crombez, 2012).

1.6.3 Pengembangan Sistem

Setelah data terkumpul tahap selanjutnya adalah membangun sistem, tahap ini meliputi pembagian data, normalisasi data, prapemrosesan data, ekstraksi fitur, kombinasi fitur, klusterisasi dan klasifikasi. Secara garis besar tahap ini merupakan tahap pembangunan sistem verifikasi.

1.6.4 Evaluasi dan Validasi Sistem

Setelah sistem verifikasi terbentuk tahap selanjutnya adalah melakukan verifikasi dan validasi terhadap sistem tersebut. Peneliti melakukan beberapa skenario percobaan untuk menilai kinerja dari sistem yang sudah dibangun ini. Tahap validasi pada penelitian ini menggunakan KFold yang *ter-stratified*. Sedangkan tahap evaluasi klusterisasi menggunakan purity dan klasifikasi menggunakan akurasi, presisi, sensitivitas, dan *FScore*.

1.6.5 Analisa Hasil Evaluasi

Tahap ini merupakan tahap terakhir dalam penelitian ini, tahap ini merupakan tahap dimana peneliti melakukan analisa berdasarkan hasil evaluasi yang didapatkan pada pengujian yang dilakukan pada penelitian ini.

1.7 Sistematika Penulisan

BAB I PENDAHULUAN membahas tentang gambaran umum yang akan dilakukan pada penelitian ini. Gambaran umum pada bab ini meliputi latar belakang, rumusan masalah, batasan masalah, tujuan penelitian, metode penelitian dan sistematika penulisan. Pada sub bab 1.1 menjelaskan tentang kasus nyata mengenai verifikasi penulis dan rencana penelitian secara umum. Sub bab 1.2 menjelaskan permasalahan apa yang menjadi fokus dalam penelitian ini, pada sub bab 1.3 berisi tentang ruang lingkup dari penelitian ini. Tujuan penelitian tertulis pada sub bab 1.4 sedangkan manfaat dari penelitian ini ditulis pada sub bab 1.5. Pada sub bab 1.6 dijelaskan metode apa saja yang digunakan dalam penelitian ini, sedangkan pada sub bab 1.7 berfokus pada menjelaskan seluruh isi bab yang ada dalam penelitian ini secara singkat.

BAB II TINJAUAN PUSTAKA DAN LANDASAN TEORI mencakup 2 bagian utama yaitu tinjauan pustaka dan landasan teori. Sub bab 2.1 berisi tentang tinjauan dari berbagai penelitian sebagai pendukung penelitian yang akan dilakukan dan perbedaan penelitian ini dengan penelitian yang sudah dilakukan sebelumnya. Sub bab 2.2 berisi tentang landasan dasar teori yang menjadi prinsip utama ataupun konsep dilakukannya penelitian ini.

BAB III PERANCANGAN SISTEM membahas tentang perancangan sistem verifikasi penulis yang dibangun selama proses penelitian. pada sub bab 3.1 akan menjelaskan tentang kebutuhan perangkat keras maupun perangkat lunak dalam membangun sistem. Sub bab selanjutnya yaitu sub bab 3.2 menjelaskan tentang perancangan sistem verifikasi penulis. Pada sub bab 3.3 akan dijelaskan mengenai alur sistem beserta desain antarmukanya. Pada sub bab 3.4 akan dijelaskan rancangan pengujian sistem.

BAB IV IMPLEMENTASI DAN ANALISIS DATA mencakup tentang detail implementasi dari perancangan sistem penelitian. Sub bab 4.1 akan memaparkan hasil implementasi sistem dari penelitian ini. Sub bab 4.2 memaparkan hasil implementasi desain antarmuka yang dibangun. Pada sub bab berikutnya yaitu sub bab 4.3 akan dijelaskan hasil dan analisis pengujian dari penelitian.

BAB V KESIMPULAN DAN SARAN membahas tentang penarikan kesimpulan dari hasil penelitian ini yang akan ditulis pada sub bab 5.1 dan sub bab 5.2 berisi saran-saran yang dapat digunakan pada penelitian yang berkaitan di masa yang akan datang.

BAB V

KESIMPULAN DAN SARAN

5.1 Kesimpulan

Berdasarkan hasil pengujian dalam penelitian ini beserta analisa yang peneliti buat, peneliti mendapatkan beberapa kesimpulan sebagai berikut

- a. Kombinasi fitur *stylometry* dari fitur frekuensi relatif tanda baca (Fitur Sintaksis), frekuensi relatif *stopword* (Fitur Sintaksis), rata-rata panjang paragraf (Fitur Struktural) merupakan kombinasi fitur *stylometry* terbaik yang ditemukan pada penelitian ini, pernyataan ini dibuktikan dengan kombinasi fitur *stylometry* ini mendapatkan nilai voting terbanyak pada percobaan pengujian pertama, ketiga dan keempat. Selain itu kombinasi fitur ini mendapatkan nilai *FScore* terbaik pada pengujian pertama, keempat dengan *klasifier* KNN, dan keempat dengan *klasifier* SVM dengan nilai 0,86, 0,91, dan 0,88.
- b. Kombinasi fitur *stylometry* yang mengandung fitur frekuensi relatif tanda baca (Fitur Sintaksis), frekuensi relatif *stopword* (Fitur Sintaksis) cenderung memiliki nilai evaluasi klasifikasi maupun klasterisasi yang lebih baik dari pada kombinasi fitur *stylometry* yang tidak mengandung kedua fitur tersebut. Hal ini ditunjukkan pada kombinasi fitur *stylometry* dengan kode adalah 'KF1', 'KF2', 'KF3', dan 'KF4' hampir selalu menempati peringkat 1 sampai 4 pada proses voting di percobaan pertama, percobaan ketiga, dan percobaan keempat. Hal ini disebabkan karena keempat kombinasi fitur ini berhasil menghasilkan nilai *FScore* dan nilai *purity* yang tinggi.
- c. Nilai *purity* memberikan gambaran mengenai hasil evaluasi klasifikasi walaupun tidak terlalu tepat hal ini dapat dilihat ketika nilai *purity* pada suatu kombinasi fitur *stylometry* lebih besar dari kombinasi fitur *stylometry* lain, tidak menjamin kombinasi fitur *stylometry* tersebut memiliki nilai evaluasi klasifikasi yang lebih besar juga.
- d. Kombinasi fitur *stylometry* yang memiliki nilai *purity* kurang dari 0,5 memiliki kemungkinan yang besar akan menghasilkan nilai evaluasi

klasifikasi kurang dari 0,5 juga. Tetapi Kombinasi fitur *stylometry* yang memiliki nilai purity lebih dari sama dengan 0,5 memiliki kemungkinan yang besar akan menghasilkan nilai evaluasi klasifikasi lebih dari sama dengan 0,5 juga.

- e. Kombinasi fitur yang memiliki nilai purity tertinggi hampir selalu mendapatkan voting terbanyak dan memiliki nilai evaluasi klasifikasi yang tertinggi juga.
- f. Kombinasi fitur *stylometry* yang dijelaskan pada penelitian Halvani (2015) dan Luyckx dan Daelemans (2008) bisa digunakan untuk teks bahasa Indonesia.
- g. Fitur sintaksis juga bisa bekerja lebih baik daripada fitur pada level token lainnya pada teks bahasa Indonesia.

5.2 Saran

Saran dari peneliti yang peneliti dapatkan setelah melakukan penelitian ini adalah sebagai berikut.

- a. Mencoba fitur *Part-of-speech* dan frekuensi relatif *Part-of-speech*.
- b. Mencoba fitur pola kemunculan *stopword*.
- c. Mencoba kombinasi fitur *stylometry* yang isinya keseluruhan merupakan fitur dengan jenis sintaksis.
- d. Mencoba model klasifikasi lain apakah hasil dari klasifikasi juga berbanding lurus dengan hasil dari klasterisasi.
- e. Mencoba verifikasi penulis dengan paradigma berbasis profil.
- f. Mengatur nilai H pada MKNN untuk meningkatkan hasil klasifikasi.

Daftar Pustaka

- Amiruddin, H. (2019, Februari 21). *Polisi Tangkap 3 Penipu Online Jaringan Internasional*. Diambil kembali dari www.news.okezone.com:
<https://news.okezone.com/read/2019/02/21/609/2021204/polisi-tangkap-3-penipu-online-jaringan-internasional>
- Atmoko, L. (2018). *Sistem Analisis Gaya Penulisan Untuk Pengelompokan Dokumen Dengan Metode K-means*. Diambil kembali dari
<http://sinta.ukdw.ac.id>.
- Diez, P. (2018). *Smart Wheelchairs and Brain-Computer Interfaces*. Argentina: Elsevier.
- Gonçalves, D. N., Gonçalves, C. d., Assis, T. F., & Silva, M. A. (2014). Analysis of the difference between the Euclidean distance and the actual road distance in Brazil. *Transportation Research Procedia* 3, 876 – 885.
- Halvani, O. (2015). Register & Genre Seminar: Towards Intrinsic Plagiarism Detection.
- Indraloka, D. S., & Santosa, B. (2017). Penerapan Text Mining untuk Melakukan Clustering Data Tweet Shopee Indonesia. *JURNAL SAINS DAN SENI ITS*, 2337-3520.
- Juola, P. (2013, Agustus 20). *How a Computer Program Helped Show J.K. Rowling write A Cuckoo's Calling*. Diambil kembali dari www.scientificamerican.com:
<https://www.scientificamerican.com/article/how-a-computer-program-helped-show-jk-rowling-write-a-cuckoos-calling/>
- Kestemont, M., Luyckx, K., Daelemans, W., & Crombez, T. (2012). Cross-Genre Authorship Verification Using Unmasking. *English Studies*, 340–356.
- Li, Y., & Wu, H. (2012). A Clustering Method Based on K-Means Algorithm. *Physics Procedia* 25, 1104 – 1109.
- Luyckx, K., & Daelemans, W. (2008). Authorship Attribution and Verification with Many Authors and Limited. *Proceedings of the 22nd International Conference on Computational Linguistics*, 513-520.

- Parvin, H., Alizadeh, H., & Minaei-Bidgoli, B. (2008). MKNN: Modified K-Nearest Neighbor. *Proceedings of the World Congress on Engineering and Computer Science*, 22 - 24.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., . . . Duchesnay, E. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830.
- Potha, N., & Stamatatos, E. (2014). A Profile-Based Method for Authorship Verification. *Springer International Publishing Switzerland*, 313–326.
- Sokolova, M., & Lapalme, G. (2009). A systematic analysis of performance measures for classification tasks. *Information Processing and Management* 45, 427–437.
- Sripada, S. C., & Rao, D. M. (2011). COMPARISON OF PURITY AND ENTROPY OF K-MEANS CLUSTERING AND FUZZY C MEANS CLUSTERING. *Indian Journal of Computer Science and Engineering (IJCSSE)*, 343-346.
- Stamatatos, E., Tschuggnall, M., Verhoeven, B., Daelemans, W., Specht, G., Stein, B., & Potthast, M. (2016). Clustering by Authorship Within and Across Documents. *CLEF*, 691-715.
- Vartapetian, A., & Gillam, L. (2016). A Big Increase in Known Unknowns: from Author Verification to Author Clustering. *CLEF*, 1008-1013.
- Wafiyah, F., Hidayat, N., & Perdana, R. S. (2017). Implementasi Algoritma Modified K-Nearest Neighbor (MKNN) untuk Klasifikasi Penyakit Demam. *Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer*, 1210-1219.
- Zainuddin, N., & Selamat, A. (2014). Sentiment Analysis Using Support Vector Machine. *International Conference of Computer, Communications, and Control Technology* (hal. 333-337). Langkawi: IEEE.