

**TUGAS AKHIR**

**PENERAPAN METODE *WORD2VEC* UNTUK  
MENDETEKSI KEMIRIPAN DOKUMEN**



**ADE RIYANI**

**15102002**

**PROGRAM STUDI S1 INFORMATIKA  
FAKULTAS TEKNOLOGI INDUSTRI DAN INFORMATIKA  
INSTITUT TEKNOLOGI TELKOM PURWOKERTO  
2019**

**TUGAS AKHIR**

**PENERAPAN METODE *WORD2VEC* UNTUK  
MENDETEKSI KEMIRIPAN DOKUMEN**

***IMPLEMENTATION WORD2VEC METHOD TO  
DETECT DOCUMENT SIMILARITY***

Disusun Sebagai Salah Satu Syarat untuk Memperoleh Gelar Sarjana Komputer



**ADE RIYANI  
15102002**

**PROGRAM STUDI TEKNIK INFORMATIKA  
FAKULTAS TEKNOLOGI INDUSTRI DAN INFORMATIKA  
INSTITUT TEKNOLOGI TELKOM PURWOKERTO  
2019**

Lembar Pengesahan Pembimbing

**PENERAPAN METODE *WORD2VEC* UNTUK  
MENDETEKSI KEMIRIPAN DOKUMEN**

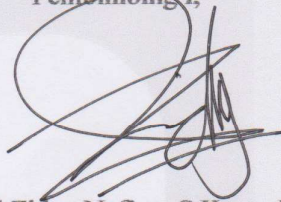
***IMPLEMENTATION WORD2VEC METHOD TO  
DETECT DOCUMENT SIMILARITY***

Dipersiapkan dan Disusun Oleh

**ADE RIYANI  
15102002**

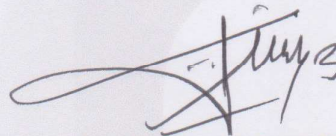
Telah Diujikan dan Dipertahankan dalam sidang Ujian Tugas Akhir  
Pada hari kamis, 7 Februari 2019

Pembimbing I,



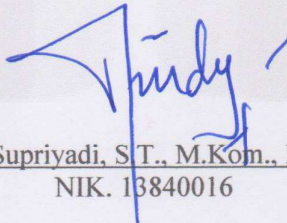
Muhammad Zidny Naf'an, S.Kom., M.Kom  
NIDN. 0626128801

Pembimbing II,



Auliya Burhanuddin, S.Si., M.Kom  
NIDN. 0630058202

Tugas Akhir ini diterima sebagai salah satu persyaratan  
untuk memperoleh gelar Sarjana Komputer  
Tanggal 28 Desember 2019  
Dekan



Didi Supriyadi, S.T., M.Kom., ITILE  
NIK. 13840016



Lembar Penetapan Penguji

**PENERAPAN METODE *WORD2VEC* UNTUK  
MENDETEKSI KEMIRIPAN DOKUMEN**

**PENERAPAN METODE *WORD2VEC* UNTUK  
MENDETEKSI KEMIRIPAN DOKUMEN**

Dipersiapkan dan Disusun Oleh

**ADE RIYANI  
15102002**

**Tugas Akhir Telah diuji dan Dinilai Panitia Penguji Program Studi Teknik  
Informatika Fakultas Teknologi Industri dan Informatika  
Institut Teknologi Telkom Purwokerto  
Pada Tanggal : 7 Februari 2019**

**Ketua  
Penguji**

**(Tri Ginanjar Laksana, S.Kom., M.C.S., M.Kom.)  
NIDN. 0407088502**

**Anggota  
Penguji I**

**(Agi Prasetiadi, S.T., M. Eng.)  
NIDN. 0617098802**

**Anggota  
Penguji II**

**(Muhammad Fajar Sidiq, S.T., M.T.)  
NIDN. 0619029102**



## HALAMAN PERNYATAAN KEASLIAN TUGAS AKHIR

Yang bertandatangan di bawah ini,

Nama Mahasiswa : Ade Riyani  
NIM : 15102002  
Program Studi : Teknik Informatika

Menyatakan bahwa Tugas Akhir dengan judul berikut:  
**PENERAPAN METODE *WORD2VEC* UNTUK MENDETEKSI  
KEMIRIPAN DOKUMEN**

Dosen Pembimbing Utama : Muhammad Zidny Naf'an, S.Kom.,M.Kom

Dosen Pembimbing Pendamping : Auliya Burhanuddin, S.SI., M.Kom

1. Karya tulis ini adalah benar-benar ASLI dan BELUM PERNAH diajukan untuk mendapatkan gelar akademik, baik di Institut Teknologi Telkom Purwokerto maupun di Perguruan Tinggi lainnya.
2. Karya tulis ini merupakan gagasan, rumusan, dan penelitian Saya Sendiri, tanpa bantuan pihak lain kecuali arahan dari Tim Dosen Pembimbing.
3. Dalam Karya tulis ini tidak terdapat karya atau pendapat orang lain, kecuali secara tertulis dengan jelas dicantumkan sebagai acuan dalam naskah dengan disebutkan nama pengarang dan disebutkan dalam Daftar Pustaka pada karya tulis ini.
4. Perangkat lunak yang digunakan dalam penelitian ini sepenuhnya menjadi tanggungjawab Saya, bukan tanggungjawab Institut Teknologi Telkom Purwokerto.
5. Pernyataan ini Saya buat dengan sesungguhnya, apabila dikemudian hari terdapat penyimpangan dan ketidakbenaran dalam pernyataan ini, maka Saya bersedia menerima Sanksi Akademik dengan pencabutan gelar yang sudah diperoleh serta sanksi lainnya sesuai dengan norma yang berlaku di Perguruan Tinggi.

Purwokerto, 28 Januari 2019,

Yang Menyatakan,

  
(Ade Riyani, 

## KATA PENGANTAR

Puji syukur penulis panjatkan kehadirat Allah SWT, karena atas berkat rahmat, taufik, hidayah dan karunia-Nya penulis dapat menyelesaikan penyusunan skripsi ini di Institut Teknologi Telkom Purwokerto tempat penulis belajar.

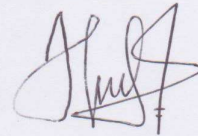
Adapun judul skripsi penelitian ini yaitu **PENERAPAN METODE *WORD2VEC* UNTUK MENDETEKSI KEMIRIPAN DOKUMEN**. Dalam menyelesaikan penulisan laporan skripsi ini tentunya tidak terlepas dari bimbingan, doa serta bantuan dari berbagai pihak. Oleh karena itu dalam kesempatan ini penulis mengucapkan terimakasih kepada:

1. Bapak Dr. Ali Rokhman, M.Si selaku Rektor Institut Teknologi Telkom Purwokerto.
2. Bapak Didi Supriyadi, S.T., M.Kom selaku Dekan Fakultas Teknik Industri dan Informatika Institut Teknologi Telkom Purokerto.
3. Muhammad Zidny Naf'an, S.Kom., M.Kom selaku ketua program studi S1 Informatika dan selaku Dosen Pembimbing pertama yang telah memberikan bimbingan dan pengarahan pada saat penyusunan Tugas Akhir.
4. Auliya Burhanuddin, S.SI., M.Kom. selaku Dosen Pembimbing kedua yang telah memberikan bimbingan dan pengarahan pada saat penyusunan Tugas Akhir.
5. Orang tua yang telah memberikan doa dan dukungan untuk penulis pada saat mengerjakan Tugas Akhir.
6. Asyhar Nurrochman yang telah menemani dan memberikan semangat untuk menyelesaikan Tugas Akhir.
7. Cuplis kucing kesayangan yang selalu menemani penulis bermain diwaktu senggang.
8. Alhamda Adisoka Bimantara sebagai konsultan pada saat program saya gagal atau eror.
9. Teman-teman mahasiswa Institut Teknologi Telkom Purwokerto yang tidak bisa disebutkan satu persatu.

Dalam penyusunan penelitian Tugas Akhir ini, penulis menyadari masih banyak kekurangan. Untuk itu, saran dan kritik pembaca untuk kesempurnaan laporan Tugas Akhir.

Akhirnya, penulis berharap semoga tugas akhir ini dapat bermanfaat dan menambah wawasan bagi pembaca.

Penulis

A handwritten signature in black ink, appearing to be 'Ade Riyani', written in a cursive style.

Ade Riyani

## DAFTAR ISI

HALAMAN SAMBUNG.....	ii
HALAMAN PENGESAHAN.....	iii
HALAMAN PERNYATAAN KEASLIAN TUGAS AKHIR.....	v
KATA PENGANTAR.....	vi
ABSTRAK.....	viii
ABSTRACT.....	ix
DAFTAR ISI.....	x
DAFTAR TABEL.....	xii
DAFTAR GAMBAR.....	xiii
DAFTAR LAMPIRAN.....	xiv
BAB I PENDAHULUAN.....	1
1.1 Latar Belakang.....	1
1.2 Rumusan Masalah.....	3
1.3 Tujuan Penelitian.....	3
1.4 Batasan Masalah.....	4
BAB II TINJAUAN PUSTAKA.....	5
2.1 Penelitian Sebelumnya.....	5
2.2 Landasan Teori.....	10
2.2.1 Plagiarisme.....	10
2.2.2 Dokumen.....	11
2.2.3 Text Mining.....	13
2.2.4 Text Processing.....	14
2.2.5 <i>Word2vec</i> .....	17
2.2.6 <i>Term Frequency Inverse Document Frequency (TF-IDF)</i> .....	18
2.2.7 <i>Cosine Similarity</i> .....	19
BAB III METODE PENELITIAN.....	20
3.1 Tahapan Penelitian.....	20
3.1.1 Studi Literatur.....	21
3.1.2 Pengumpulan Data.....	21
3.1.3 <i>Preprocessing</i> .....	23
3.1.4 Penerapan Metode.....	23
3.1.5 Hasil Kemiripan.....	25



3.2	Gambaran Umum Sistem .....	25
BAB IV HASIL PENELITIAN DAN ANALISIS.....		26
4.1	Pengumpulan Data.....	26
4.2	<i>Preprocessing</i> .....	26
4.3	Penerapan Metode .....	31
4.3.1	Pembobotan TF-IDF dan <i>Cosine Similarity</i> .....	31
4.3.2	Metode <i>Word2vec</i> dan <i>Cosine Similarity</i> .....	34
4.4	Analisis .....	38
BAB V PENUTUP .....		42
5.1	Kesimpulan .....	42
5.2	Saran .....	42
Daftar Pustaka .....		43
LAMPIRAN .....		46

## DAFTAR TABEL

Tabel 2. 1 Penelitian Terdahulu.....	8
Tabel 3. 1 Abstrak Skripsi.....	22
Tabel 4. 1 Sebelum dan Sesudah <i>Preprocessing</i> .....	26
Tabel 4. 1 Sebelum dan Sesudah Preprocessing .....	26
Tabel 4. 2 Dokumen pembandingan sebelum dan sesudah preprocessing.....	28
Tabel 4. 3 Nilai Kemiripan Dokumen Uji dengan Pembandingan dengan TF-IDF ....	33
Tabel 4. 4 Nilai Kemiripan Dokumen Uji dengan Pembandingan dengan <i>Word2vec</i> .	36
Tabel 4. 5 Perbandingan Nilai Cosine Similarity tanpa Stemming dan Stemming...	38
Tabel 4. 6 Tabel Perbandingan.....	40

## DAFTAR GAMBAR

Gambar 2. 1 Tahapan proses <i>case folding</i> . .....	14
Gambar 2. 2 Tahapan proses <i>tokenizing</i> . .....	15
Gambar 2. 3 Tahapan proses stopword removal atau filtering.....	16
Gambar 2. 4 Tahapan proses <i>stemming</i> .....	17
Gambar 2. 5 CBOW Skip-gram .....	18
Gambar 3. 1 Tahapan penelitian.....	20
Gambar 3. 2 <i>Website repository</i> perpustakaan.....	21
Gambar 3. 3 Gambaran pemodelan <i>Word2vec</i> .....	24
Gambar 3. 4 Gambaran Umum Sistem.....	25
Gambar 4. 1 Hasil Plagiarisme Detector.....	39



## DAFTAR LAMPIRAN

Lampiran 1 Nilai Kemiripan Menggunakan TF-IDF .....	46
Lampiran 2 Nilai Kemiripan Menggunakan <i>Word2vec</i> .....	51
Lampiran 3 Program TF-IDF .....	56
Lampiran 4 Program <i>Word2vec</i> .....	59
Lampiran 5 Dokumen Abstrak .....	65
Lampiran 6 Plagiarisme Detector .....	79

# **BAB I**

## **PENDAHULUAN**

### **1.1 Latar Belakang**

Pemanfaatan teknologi informasi memiliki dampak pada banyak sektor, termasuk sektor Pendidikan. Dalam dunia Pendidikan, sangatlah diperlukan informasi untuk mendukung pembelajaran. Komputer dan internet sangat bermanfaat untuk pencarian informasi. Selain adanya kemudahan dalam pencarian informasi, hadirnya teknologi informasi ternyata juga memiliki dampak negatif. Salah satu dampak negatif di dunia Pendidikan adalah plagiarisme oleh civitas akademika. Tindakan plagiarisme tidak hanya melibatkan siswa maupun mahasiswa, namun plagiarisme juga dapat melibatkan dosen bahkan mahasiswa yang bergelar doktor. Diungkapkan pada Kumparan, 2018 yang melaporkan bahwa Ombudsman RI menemukan plagiarisme dalam tiga karya ilmiah Rektor Universitas Halu Oleo (UHO) Kendari, Sulawesi Tenggara yang dilakukan oleh Muhammad Zamrun Firihu. Berdasarkan hasil analisis yang 30 guru besar UHO, zamrun terbukti bersalah melakukan tindakan plagiarisme dalam tiga jurnal internasional. Tingkat kesamaanya bahkan mencapai 78% bahkan lebih. Atas kasus tersebut, Ombudsman menilai ada pelanggaran berupa maladministrasi yang dilakukan Kementerian Riset, Teknologi, dan Pendidikan tinggi (Kemendiknas) atas pelantikan Zamrun sebagai rektor. Dari kasus tersebut, plagiat dapat merusak moral civitas akademika dalam dunia Pendidikan terutama di Indonesia.

Plagiarisme merupakan salah satu perbuatan yang melanggar kode etik dalam dunia penulisan, karena tindakan tersebut telah mengambil karya orang lain dan mengakuinya sebagai karya sendiri [1]. Selain melanggar kode etik, plagiarisme merusak moral dan mempengaruhi kualitas akan suatu dokumen sehingga keasliannya tidak dapat dipercaya. Pada tahun 2010, Direktorat Jenderal Pendidikan Tinggi telah mengeluarkan peraturan tentang cara pencegahan dan penanggulangan plagiarisme termasuk sanksi untuk dosen, mahasiswa, dan calon guru besar sekalipun. Praktek ini kerap dihubungkan dalam dunia pendidikan khususnya mahasiswa dalam melaksanakan tugasnya. Merujuk pada Permendiknas No. 17 Tahun 2010 tentang

Pencegahan dan Penanggulangan Plagiat di Perguruan Tinggi disebutkan bahwa dalam melaksanakan otonomi keilmuan dan kebebasan akademik mahasiswa, dosen, peneliti, tenaga kependidikan wajib menjunjung tinggi kejujuran dan etika akademik, terutama larangan untuk melakukan plagiat dalam menghasilkan karya ilmiah, sehingga kreativitas dalam bidang akademik dapat tumbuh dan berkembang [2]. Namun jika pengambilan suatu karangan disertai dengan mencantumkan asalnya atau nama pengarang serta judul karangan yang diambil, maka tindakan pengambilan karangan tersebut bukan merupakan plagiat. Oleh karena itu, perlu dilakukan pemeriksaan kemiripan antar dokumen, dalam hal ini adalah dokumen teks, sebagai langkah validasi keterkaitan dan hubungan antar dokumen tersebut. Sepasang kata dinyatakan mirip apabila memiliki kesamaan dari sisi makna atau konsep. Perhitungan deteksi kemiripan dilatar belakangi oleh suatu masalah yaitu mesin belum dapat menyamakan persepsi manusia dengan baik, untuk itu kemiripan semantik digunakan untuk membantu mesin memahami bahasa manusia.

Tugas deteksi kemiripan yaitu untuk menghitung nilai kedekatan antara dua buah kata. Karena mesin hanya dapat membaca angka, maka kata-kata yang ada harus diterjemahkan dalam bentuk angka, maka kata-kata yang ada harus diterjemahkan dalam bentuk angka terlebih dahulu. Untuk mencari nilai kedekatan antara dua buah kata dapat dilakukan dengan cara menghitung nilai *vector* kata tersebut. Supaya mendapatkan nilai kemiripan dengan *vector* yaitu dengan cara menghitung perbedaan sudut antara dua buah *vector* dengan rumus *Cosine Similarity*.

Beberapa penelitian untuk mengetahui tingkat kesamaan dokumen telah dilakukan diantaranya menggunakan metode *Cosine Similarity* [3], Jaccard [4], TF-IDF [5], *Support Vector Regression* [6], *Single pass clustering* [7], *Levenshtein Distance* [8], dan *Latent Semantic Analysis* [9]. Berdasarkan referensi penelitian yang telah dilakukan Metode *Word2vec* belum banyak digunakan untuk kemiripan dokumen karena metode *Word2vec* baru dikembangkan pada tahun 2013. *Word2vec* merupakan salah satu metode berbasis *vector* untuk mencari kemiripan semantik sehingga penulis menggunakan metode tersebut untuk membuat sistem dan menggabungkan dengan algoritma *Cosine Similarity*. Pendekatan *Word2vec* dipilih pada penelitian tugas akhir ini sebab *Word2vec* dapat mempelajari hubungan antar



kata-kata secara otomatis dan dapat mempelajari berbagai macam *text* sebagai data latih [10]. *Word2vec* merupakan representasi kata dalam bentuk vektor yang digunakan untuk menghasilkan *word embeddings* [11]. Sedangkan *Cosine Similarity* digunakan untuk menghitung nilai kemiripan antar kalimat dan menjadi salah satu teknik untuk mengukur kemiripan teks yang populer. Kelebihan dari algoritma *Cosine Similarity* adalah tidak terpengaruh pada panjang pendeknya suatu dokumen. Hasil akhir yang akan diberikan dalam sistem ini adalah presentase nilai *similarity* dokumen yang diuji dengan dokumen pembanding sebanyak 116 abstrak yang didapatkan dari *repository* Institut Teknologi Telkom Purwokerto.

## 1.2 Rumusan Masalah

Rumusan masalah dari penelitian ini adalah metode *Word2vec* banyak digunakan pada teks Bahasa Inggris dan belum banyak digunakan pada teks Bahasa Indonesia terutama tentang kemiripan dokumen. Berdasarkan rumusan masalah tersebut maka munculah pertanyaan penelitian sebagai berikut:

1. Apakah metode *Word2vec* dapat digunakan untuk menghitung tingkat kemiripan pada dokumen dalam teks Bahasa Indonesia?
2. Berapa nilai kemiripan pada metode *Word2vec* dibandingkan dengan sistem plagiarism detector?

## 1.3 Tujuan Penelitian

Berdasarkan dekomposisi masalah dalam rumusan masalah tujuan pada penelitian ini adalah sebagai berikut.

1. Menerapkan metode *Word2vec* kedalam teks Bahasa Indonesia.
2. Mengetahui nilai kemiripan pada dokumen menggunakan metode *Word2vec*.

## 1.4 Manfaat

Berdasarkan tujuan penelitian yang sudah dijelaskan, adapun manfaat dari penelitian ini adalah sebagai berikut:

1. Mempercepat proses pengoreksian dengan jumlah besar dengan cepat dan efektif.
2. Mengetahui presentase kemiripan sehingga dapat digunakan sebagai bahan pertimbangan untuk mendeteksi adanya tindakan meniru.

### **1.5 Batasan Masalah**

Batasan masalah pada penelitian ini adalah sebagai berikut:

1. Dokumen yang digunakan dalam penelitian ini adalah dokumen Abstrak Tugas Akhir dan laporan Praktik Kerja Lapangan.
2. Tipe dokumen yang digunakan dalam penelitian ini adalah .pdf.
3. Tidak memperhatikan kesalahan penulisan dalam dokumen.

## BAB II

### TINJAUAN PUSTAKA

#### 2.1 Penelitian Sebelumnya

Pada penelitian ini penulis menggunakan studi literatur sebagai sarana untuk kelengkapan data sekaligus untuk mempertajam masalah yang penulis kaji. Penulis telah mengkaji delapan jurnal penelitian terdahulu. Kedelapan jurnal tersebut penulis pilih berdasarkan topik dan tema yang sesuai untuk penelitian ini. Berikut merupakan penjelasan lebih lanjut.

Agung Wardhana Z. Nasution, dkk [6]. Melakukan Analisa dan implementasi perhitungan nilai semantic kemiripan kata pada ayat Al-Quran terjemahan Bahasa Inggris menggunakan pendekatan *Word alignment* dan *vector semantic Word2vec* [6]. Sebagai evaluasi dari penelitian ini digunakan model *Support Vector Regression*. Hasil dari penelitian ini adalah pengujian Sistem yang menggunakan data 2017 sebagai data test dengan inputan 350 ayat terjemahan Bahasa Inggris Ibnu Katsir dan 50 ayat indeks sematik memperoleh nilai korelasi *alignment* 0,81958, nilai *Word2vec* 0,64395 dan nilai korelasi SVR yang diperoleh 0,81221. Nilai SVR ini, memiliki nilai lebih kecil 0,00737 dibandingkan *alignment* TF-IDF dikarenakan kompleksitas pada masing-masing data dan pengaruh fitur *Word2vec*. hal tersebut tidak terjadi jika data *test* menggunakan data 2016 dengan fitur tanpa *stemming* pada *Word2vec* menghasilkan nilai SVR 0,9106 lebih besar 0,0016 dari hasil *alignment* [6]. Kesamaan penelitian ini dengan penelitian penulis adalah metode yang digunakan yaitu *Word2vec*. Kekurangan dalam penelitian ini adalah kompleksitas pasangan ayat yang tidak dapat diidentifikasi oleh sistem. Sistem yang dibangun tidak sepenuhnya dapat mengidentifikasi kata-kata serapan yang berhubungan dengan AlQuran kedalam bentuk *vector*. Namun penelitian pada [11] mencantumkan alasan tersebut pada bagian analisis sehingga transparansi dalam penelitian tersebut dapat dilihat dengan jelas dan kejujuran menjadi nilai lebih dalam melaksanakan penelitiannya.

Kemal Ade Sekarwati, dkk [12]. Membangun sistem pengukuran kemiripan dokumen berbahasa Indonesia dengan menggunakan metode *Latent Semantic Analysis* (LSA) sedangkan *tools* yang digunakan untuk penelitian ini adalah *Gensim*. *Gensim* merupakan *open-source* model ruang vektor dan *toolkit* topik modeling [12].



Hasil dari penelitian ini yaitu Pengujian dokumen dilakukan terhadap dokumen yang diduga mempunyai kemiripan dengan dokumen lain. Dokumen 1 merupakan dokumen yang diuji, sedangkan dokumen 2 merupakan dokumen asli. Menghasilkan nilai persentase kemiripan 100% bagi dokumen yang sama [12].

Sugiyamta [13] membuat sistem deteksi kemiripan dokumen (*document similarity*) menggunakan Algoritma *Cosine Similarity* dan teknik mengelompokkan dokumen dengan Algoritma *Single Pass Clustering*. Hasil dari penelitian adalah pengukuran kemiripan dengan *Cosine Similarity* untuk dokumen abstrak memiliki tingkat akurasi 99%. *Single Pass Clustering* dapat membantu menemukan dokumen yang ada dalam satu *Cluster* dengan *query* yang dimasukkan oleh pengguna [13].

Na'fairul Hasna [8] melakukan percobaan mendeteksi kemiripan dokumen teks menggunakan algoritma *Levenshtein Distance* sehingga dapat digunakan untuk membantu menentukan plagiarisme. Tipe dokumen yang diuji adalah .pdf .docx dan .txt. Hasil dari penelitian ini adalah pengukuran kemiripan dengan algoritma *Levenshtein Distance* menghasilkan nilai *similarity* yang tinggi yaitu 85% untuk dokumen yang tingkat kemiripannya tinggi. Sedangkan untuk dokumen yang tingkat kemiripannya rendah menghasilkan nilai *similarity* kurang dari 40% [8].

Mochammad Iran Dary [9] melakukan percobaan untuk mengetahui keterkaitan antar ayat didalam al-Quran menggunakan metode *latent semantic analysis*. Menggunakan algoritma *Cosine Similarity* untuk membandingkan jarak dari 2 buah *vector*. Hasil dari penelitian ini adalah metode *latent semantic analysis* cukup efektif untuk mencari kesamaan antar ayat pada alquran yang berupa *short text*. Terlihat dari nilai *F-Measure* yang paling tinggi adalah 14% [9].

Radiant Victor, dkk [14] melakukan mengembangkan sebuah aplikasi yang mengimplementasikan *Cosine Similarity* yang berguna untuk mengukur kesamaan teks berdasarkan kemunculan kata-kata dalam teks tersebut dan algoritma *Smith-Waterman* yang berfungsi untuk menghitung kemiripan teks berdasarkan urutan kata [14]. Dalam proses mengubah kata-kata berbahasa indonesi menjadi kata-kata dasar (*stemming*) menggunakan algoritma Nazief-Adriani. Hasil dari penelitian ini dapat mendeteksi tingkat kemiripan teks dari sangat mirip hingga sangat tidak mirip

berdasarkan kemunculan kata di dalamnya dengan menggunakan *Cosine Similarity* dan algoritma *Smith-Waterman*.

Sedangkan Kadek Bayu, dkk [10] Membangun sebuah sistem yang dapat menentukan prioritas opini dari suatu fitur produk berdasarkan nilai kemiripan atau *similarity* suatu kata. sistem yang akan dibangun menggunakan pendekatan *Word2vec* dan pendekatan *WordNet* untuk proses klasifikasi opini dimana *Word2vec* merupakan model representasi kata dalam bentuk *vector*, dan pada penelitian ini *Word2vec* digunakan untuk melakukan klasifikasi opini pada fitur suatu produk [10]. Hasil dari penelitian ini yaitu Penerapan pendekatan *Word2vec* pada proses klasifikasi memiliki nilai akurasi yang lebih tinggi dibandingkan penerapan pendekatan *WordNet*, yaitu dengan rata-rata persentase akurasi pada *Word2vec* sebesar 43,55% sedangkan rata-rata persentase akurasi pada *WordNet* sebesar 41,48%, dengan selisih 2,07% [10].

Tabel 2. 1 Penelitian Terdahulu

No	Penulis	Penulis	Langkah Perancangan Sistem	Objek Penelitian	Perbedaan dengan penelitian yang akan dilakukan
1	A. W. Z. Nasution dkk (2017)	Analisis dan Implementasi Perhitungan <i>Semantics Similarity</i> Pada Ayat AL-Quran Dengan Pendekatan <i>Word Alignment</i> Berdasarkan <i>Support Vector Regression</i>	<i>Preprocessing alignment, contextual Similarity</i> , metode ekstraksi TF-IDF, <i>preprocessing Word2vec</i> , perhitungan <i>semantic Word2vec</i> , SVR	Terjemahan ayat suci Al-Quran berbahasa Inggris	penelitian ini data yang digunakan berbahasa Inggris dan korpus yang dibuat juga Bahasa Inggris. Sedangkan penulis menggunakan Bahasa Indonesia
2	Sekarwati dkk (2013)	Pengukuran Kemiripan Dokumen Dengan <i>Tools Gensim</i>	<i>Preprocessing</i> , metode ekstraksi TF-IDF, <i>Latent Semantic Indexing</i> , perhitungan kemiripan <i>Cosine Similarity, Dice's Similarity</i>	Dokumen pribadi	Penelitian ini dilakukan menggunakan model <i>latent semantic analysis</i> untuk mendapatkan nilai kemiripan. Sedangkan penelitian yang dilakukan penulis menggunakan <i>Word2vec</i> dan <i>cosine similarity</i>
3	Sugiamta (2015)	Sistem Deteksi Kemiripan Dokumen Dengan Algoritma <i>Cosine Similarity</i> Dan <i>Single Pass Clustering</i>	<i>Preprocessing, Term indexing, term similarity</i> metode <i>cosine, Clustering</i>	550 Abstrak Skripsi	Penelitian ini melakukan pengelompokan nilai similaritas paling tinggi menggunakan <i>single pass clustering</i> . Sedangkan penelitian yang dilakukan melakukan pengelompokan dari nilai tertinggi ke terendah secara langsung

No	Penulis	Penulis	Langkah Perancangan Sistem	Objek Penelitian	Perbedaan dengan penelitian yang akan dilakukan
4	N. Ariyani dkk (2016)	Aplikasi Pendeteksi Kemiripan Isi Teks Dokumen Menggunakan Metode <i>Levenshtein Distance</i>	<i>Preprocessing, Metode Levenshtein Distance, term similarity</i>	Dokumen Pribadi	Penelitian ini menggunakan <i>Levenshtein Distance</i> untuk mendapatkan nilai <i>similarity</i> dari dokumen. Data yang digunakan dokumen berplagiat yang diambil dari artikel/berita. Sedangkan penelitian yang akan dilakukan menggunakan <i>cosine similarity</i> untuk mendapatkan nilai kemiripan
5	M. Dary (2017)	Analisis dan Implementasi <i>Short Text Similarity</i> dengan Metode <i>Latent Semantic Analysis</i> Untuk Mengetahui Kesamaan Ayat al-Quran	<i>Preprocessing, Text Preprocessing, Term Document, Implementasi Latent Semantic Analysis, Reconstructed matrix, evaluasi</i>	Terjemahan ayat suci Al-Quran berbahasa Indonesia	Penelitian ini menerapkan <i>latent semantic analysis</i> untuk mencari kesamaan antar ayat pada al-Quran yang berupa <i>short text</i> . Sedangkan penelitian yang akan dilakukan menggunakan dokumen abstrak.
6	R. Imbar dkk (2014)	Implementasi <i>Cosine Similarity</i> dan Algoritma <i>Smith-Waterman</i> untuk Mendeteksi Kemiripan Teks	<i>Preprocessing, Term Similarity Cosine Similarity, Term Similarity algoritma Smith Waterman</i>	Dokumen Pribadi	penelitian ini dapat mendeteksi tingkat kemiripan teks dari sangat mirip hingga sangat tidak mirip berdasarkan urutan kata pembentuknya dengan menggunakan algoritma <i>Smith-Waterman</i> . Sedangkan penelitian yang akan dilakukan tidak berdasarkan urutan kata pembentuknya.

No	Penulis	Penulis	Langkah Perancangan Sistem	Objek Penelitian	Perbedaan dengan penelitian yang akan dilakukan
7	B Surarso dkk (2016)	Analisa Performa Metode <i>Cosine</i> dan <i>Jaccard</i> Pada Pengujian Dokumen	<i>Preprocessing, Indexing, Hitung Similaritas metode Jaccard dan cosine, klustering</i>	Abstrak Skripsi	penelitian ini membandingkan bandingkan performa dari metode <i>Cosine</i> dan <i>Jaccard</i> untuk menguji tingkat kemiripan dokumen dalam bentuk abstrak
8	I.Bayu dkk (2017)	Klasifikasi Opini Pada Fitur Produk Berbasis <i>Graph</i>	<i>Dataset, Ekstraksi, Klasifikasi Word2vec dan WordNet</i>	Data review produk dari <i>paper</i> “ <i>Mining and Summarizing Customers Reviews</i>	Penelitian ini menganalisis opini yang berupa opini positif atau negatif secara otomatis. Sedangkan penelitian yang akan dilakukan mendeteksi kemiripan teks dalam dokumen.

## 2.2 Landasan Teori

### 2.2.1 Plagiarisme

Plagiarisme merupakan tindakan kriminal yang sering terjadi dalam dunia akademis. Plagiarisme itu sendiri berasal dari kata latin “Plagiarus” yang berarti penculik dan “Plagiare” yang berarti mencuri. Jadi, secara sederhana plagiat berarti mengambil ide, kata-kata, dan kalimat seseorang dan memposisikannya sebagai hasil karyanya sendiri atau menggunakan ide, kata-kata, dan kalimat tanpa mencantumkan sumber dimana seorang penulis mengutipnya [1].

Menurut Peraturan Menteri Pendidikan Nasional Republik Indonesia Nomor 17 Tahun 2010 mengatakan bahwa plagiat adalah perbuatan sengaja atau tidak sengaja dalam memperoleh atau mencoba memperoleh kredit atau nilai untuk suatu karya ilmiah, dengan mengutip sebagian atau seluruh karya dan atau karya ilmiah pihak lain yang diakui sebagai karya ilmiahnya, tanpa menyatakan sumber secara tepat dan memadai [15]. Dalam Kamus Besar Bahasa Indonesia (2008) Plagiat adalah

pengambilan karangan (pendapat dan sebagainya) orang lain dan menjadikannya seolah-olah karangan (pendapat) sendiri”[16].

a. Jenis plagiarisme berdasarkan aspek yang dicuri

Dalam jenis plagiarisme ini terdiri dari 4 kategori yaitu plagiarisme ide, isi, kata, kalimat, paragraph, dan plagiarisme total. Plagiarisme ide sering dihubungkan dengan laporan hasil penelitian replikatif. Penelitian replikatif adalah penelitian yang secara garis besar mengulang penelitian orang lain, dengan maksud untuk menambah data, menguji hipotesis apakah hasil yang sudah ditemukan dalam suatu populasi berlaku pula untuk populasi yang lain [1]. Plagiarisme isi fabrikasi dan falsifikasi. Fabrikasi adalah tindakan membuat data yang tidak ada menjadi ada. Sedangkan falsifikasi adalah mengubah data dengan maksud agar sesuai dengan yang dikehendaki oleh peneliti tersebut [17]. Plagiarisme kata, kalimat, paragraph merupakan jenis plagiarisme yang mengubah sebagian kecil dalam sebuah karya tulis, seperti mengganti sebagian kalimat atau paragraf atau meniru seluruh isinya dengan bahasa yang berbeda [1]. Plagiarisme total merupakan plagiarisme yang dibuat tanpa merubah total isi dari penelitian tersebut dan sama persis dengan karya orang lain. plagiarisme ini merupakan plagiarisme yang paling berat [1].

b. Klafikasi berdasarkan proporsi atau persentasi kata, kalimat, paragraf yang dibajak

Plagiarisme ringan: <30%

Plagiarisme sedang: 30-70%

Plagiarisme berat atau total : >70%

(angka-angka tersebut tentu dibuat secara arbitrer berdasarkan “kepantasan” tanpa dasar kuantitatif yang definitif) [1].

### **2.2.2 Dokumen**

Menurut Kamus Besar Bahasa Indonesia (KBBI) menyebutkan dokumen adalah sesuatu yang tertulis atau tercetak yang dapat dipergunakan sebagai bukti atau keterangan. Dokumen merupakan salah satu hal yang sangat penting karena merupakan sumber informasi yang diperlukan oleh suatu instansi, organisasi, atau Negara [16]. Tanpa dokumen kita akan kehilangan data-data yang diperlukan untuk

kegiatan kantor/organisasi masa yang akan datang. Jenis-jenis format dokumen adalah sebagai berikut:

a. DOC (.doc)

*File* Doc merupakan dokumen pengolah kata yang diciptakan oleh program Microsoft Word. Jenis *file* ini bisa berisi susunan teks, image, tabel, grafik, *chart*, format halaman, dan pengaturan cetak. Setiap *file* DOC selalu diikuti dengan ekstensi.doc [18]. Jenis *file* DOC secara otomatis tercipta ketika membuat dokumen pada program MS Word dan menyimpannya. Ekstensi .doc memang tidak selalu muncul di akhir nama *file*. Namun, untuk mengetahui identitas atau jenis *file* tersebut, anda bisa melihat kolom *Type* dimana pada kolom tersebut diinformasikan jenis *file* adalah Microsoft Word Document [18]

b. DOCX (.docx)

*File* DOCX dirancang untuk membuat dokumen yang kontennya dapat diakses dengan mudah. Hal ini dimaksudkan bahwa sebuah dokumen dengan format DOCX bisa memiliki format *file* yang beragam. Sebagai contoh, misalnya dokumen teks disimpan menggunakan *file plain text* dan dokumen gambar disimpan sebagai *file* gambar individual dalam *file* DOCX. Jenis *file* ini memiliki ekstensi .docx [18].

Berbeda dengan *file* DOC yang menyimpan dokumen data dalam *file* binary tunggal, *file* DOCX diciptakan menggunakan format Open XML yang menyimpan dokumen sebagai sebuah kumpulan *file* dan folder yang terpisah dalam paket kompresi zip. *File* DOCX berisi *file* XML dan tiga folder, yakni *docProps*, *Word*, dan *\_rels* yang menjaga dokumen properties, konten, dan hubungan di antara *file*. *File* DOCX sebenarnya merupakan kumpulan *file* XML yang dikompres. *File* XML dapat ditampilkan secara individual dengan mengubah ulang nama *file* .docx ke dalam .zip dan memperlakukan *file* tersebut sebagai *file* .zip [18].

c. TXT (.txt)

*File* TXT merupakan standar dokumen teks yang berisi rangkaian teks yang tidak terformat. Jenis *file* ini diakui oleh sembarang program pengolah kata. *File* jenis ini umumnya dihasilkan oleh sebuah *software* atau program yang



bernama *Notepad* [18]. Jenis *file* TXT disebut juga sebagai *file plain text*. Berbeda dengan *file* RTF, pada *file* TXT tidak didukung oleh format teks seperti pengaturan *style* teks, baik tebal, miring, ataupun garis bawah dan pengaturan *style font*, baik jenis maupun ukuran *font* untuk teks tertentu [18].

d. Pdf (.pdf)

*File* Pdf atau *Portable Document Format* merupakan standar untuk distribusi elektronik yang dibuat oleh Adobe. *File* ini menyediakan *layout*, tipografi, gambar *bitmap*, transparansi, dan gambar *vector* dengan hasil yang sangat baik. PDF juga mempunyai kualitas dan presisi *layout* untuk proses yang baik dan sangat ideal untuk distribusi secara *online*. *File* PDF merupakan format yang sangat umum digunakan pada *World Wide Web* [18].

### 2.2.3 Text Mining

*Text mining* adalah proses penemuan akan informasi atau *trend* baru yang sebelumnya tidak terungkap dengan memproses dan menganalisis data dalam jumlah besar. Dalam menganalisa sebagian atau keseluruhan *unstructured text*, *text mining* mencoba untuk mengasosiasikan satu bagian *text* dengan yang lainnya berdasarkan aturan tertentu. Hasil yang di harapkan adalah informasi baru atau “*insight*” yang tidak terungkap jelas sebelumnya [19].

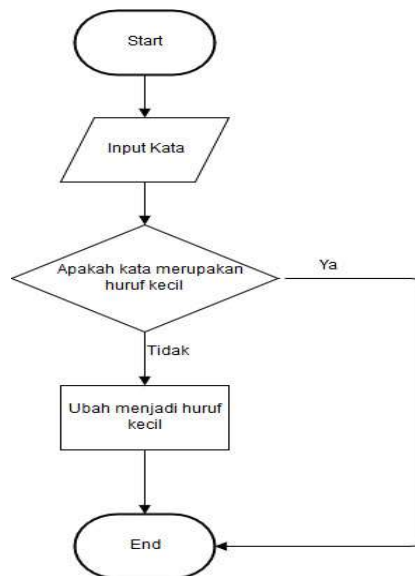
Seperti halnya data mining, *text mining* juga menghadapi masalah yang sama, termasuk jumlah data yang besar, dimensi yang tinggi, data dan struktur yang terus berubah, dan data “*noise*.” Berbeda dengan data *mining* yang utamanya memproses *structured* data, data yang digunakan *text mining* pada umumnya dalam bentuk *unstructured*, atau minimal semi-*structured*, *text*. Akibatnya, *text mining* mempunyai tantangan tambahan yang tidak ditemukan dalam data *mining*, seperti struktur *text* yang complex dan tidak lengkap, arti yang tidak jelas dan tidak standard, dan bahasa yang berbeda ditambah translasi yang tidak akurat. Dikarenakan *structured* data ditujukan agar mudah di proses komputer secara *automatic*, *pre-process* data di data *mining* jauh lebih mudah dilakukan dari pada *unstructured text* [19]. *Text* di ciptakan bukan untuk di gunakan oleh mesin, tapi untuk dikonsumsi manusia langsung. Karena itu, pada umumnya “*Natural Language Processing*” digunakan untuk memproses *unstructured text*.

#### 2.2.4 Text Processing

Struktur data yang baik dapat memudahkan proses komputerisasi secara otomatis. Pada *text mining*, informasi yang akan digali berisi informasi-informasi yang strukturnya sembarang [20]. Oleh karena itu, diperlukan proses perubahan bentuk menjadi data yang terstruktur sesuai kebutuhannya untuk proses dalam data *mining*, yang biasanya akan menjadi nilai-nilai numerik. Proses ini sering disebut *text preprocessing* [21]. Setelah data menjadi data terstruktur dan berupa nilai numerik maka data dapat dijadikan sebagai sumber data yang dapat diolah lebih lanjut. Berberapa proses yang dilakukan adalah sebagai berikut:

a. *Case folding*

*Case folding* merupakan proses pertama dari rangkaian *preprocessing* dokumen. Dalam proses ini mengubah semua huruf dalam dokumen menjadi huruf kecil. Hanya huruf a sampai dengan z yang diterima [21]. Gambar 2.1 menunjukkan *flowchart* dari proses *case folding*.



Gambar 2. 1 Tahapan proses *case folding*.

b. *Tokenizing*

Tahap *Tokenizing* adalah tahap pemotongan string input berdasarkan tiap kata yang menyusunnya [20]. Karakter selain huruf akan dianggap *delimiter* dan akan dihilangkan atau dihapus untuk proses mendapat kata-kata penyusun teks. Gambar 2.2 menunjukkan *flowchart* tahapan *tokenizing*.