

**ANALISIS ULASAN HOTEL DI SITUS TRIPADVISOR  
MENGUNAKAN METODE *TERM FREQUENCY-INVERSE*  
*DOCUMENT FREQUENCY* DAN *K-NEAREST NEIGHBOR***

**Tesis untuk Tesis S-2  
Program Studi Magister Sistem Informasi**



**Khairul Huda  
30000318410014**

**PROGRAM PASCASARJANA  
UNIVERSITAS DIPONEGORO  
SEMARANG  
2022**

**HALAMAN PENGESAHAN**

**TESIS**

**ANALISIS ULASAN HOTEL DI SITUS TRIPADVISOR MENGGUNAKAN METODE  
*TERM FREQUENCY-INVERSE DOCUMENT FREQUENCY*  
(TF-IDF) DAN *K-NEAREST NEIGHBOR* (KNN)**

**Oleh:**

**Khairul Huda  
30000318410014**

Telah diujikan dan dinyatakan lulus ujian tesis pada tanggal 15 Agustus 2022 oleh tim penguji Program Studi Magister Sistem Informasi Sekolah Pascasarjana Universitas Diponegoro.

Semarang, 15 Agustus 2022  
Mengetahui,

**Penguji I**

**Penguji II**

Prof. Dr. Kusworo Adi, S.Si., M.T.  
NIP: 197203171998021001

Drs. Bayu Surarso, M.Sc., Ph.D  
NIP. 196311051988031001

**Pembimbing I**

**Pembimbing II**

Dr. Catur Edi Widodo, MT.  
NIP: 196405181992031002

Vincencius Gunawan S.K., M.Si., Ph.D  
NIP: 197105221997021001

**Mengetahui :**  
**Dekan Sekolah Pascasarjana  
Universitas Diponegoro**

**Ketua Program Studi  
Magister Sistem Informasi**

Dr. R.B. Sularto, S.H., M.Hum  
NIP. 196701011991031005

Drs. Bayu Surarso, M.Sc., Ph.D  
NIP. 196311051988031001

**PERNYATAAN PERSETUJUAN  
PUBLIKASI TESIS UNTUK KEPENTINGAN AKADEMIS**

Sebagai sivitas akademik Universitas Diponegoro, saya yang bertanda tangan di bawah ini:

Nama : Khairul Huda  
NIM : 30000318410014  
Program Studi : Magister Sistem Informasi  
Program : Sekolah Pascasarjana  
Jenis Karya : Tesis

Demi pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada Universitas Diponegoro Hak Bebas Royalti Noneksklusif atas karya ilmiah saya yang berjudul:

**ANALISIS ULASAN HOTEL DI SITUS TRIPADVISOR  
MENGUNAKAN METODE *TERM FREQUENCY-INVERSE DOCUMENT  
FREQUENCY* DAN *K-NEAREST NEIGHBOR***

beserta perangkat yang ada. Dengan Hak bebas Royalti Noneksklusif ini Program Studi Magister Sistem Informasi Sekolah Pascasarjana Universitas Diponegoro berhak menyimpan, mengalihmedia/formatkan, mengelola dalam bentuk pangkalan data (*database*) merawat, dan mempublikasikan tesis saya selama tetap mencantumkan nama saya sebagai penulis/pencipta dan sebagai pemilik hak cipta.

Dibuat di: Semarang  
Pada tanggal: 12 Juli 2022  
Yang menyatakan

Khairul Huda  
NIM. 30000318410014

## PERNYATAAN

Dengan ini saya menyatakan bahwa dalam tesis ini tidak terdapat karya yang pernah diajukan untuk memperoleh gelar akademik di suatu perguruan tinggi, dan sepanjang pengetahuan saya juga tidak terdapat karya atau pendapat yang pernah ditulis atau diterbitkan oleh orang lain, kecuali secara tertulis diacu dalam naskah ini dan disebutkan dalam daftar pustaka.

Semarang, 12 Juli 2022

Khairul Huda



## KATA PENGANTAR

Puji syukur ke hadirat Allah SWT, atas limpahan rahmat dan karunia-Nya sehingga penulis dapat menyelesaikan penelitian tesis dengan judul “Analisis Ulasan Hotel Di Situs TripAdvisor Menggunakan Metode *Term Frequency-Inverse Document Frequency* (TF-IDF) Dan *K-Nearest Neighbor* (KNN)”. Penulis menyadari tidak dapat menyelesaikan Tesis ini tanpa bantuan dan dorongan dari berbagai pihak. Untuk itu pada kesempatan ini penulis ingin mengucapkan terima kasih kepada:

1. Bapak Dr. Catur Edi Widodo, MT selaku pembimbing pertama penulis yang telah memberikan wawasan, masukan, arahan dan nasehat dalam penulisan tesis.
2. **Bapak Vincencius Gunawan S.K., M.Si., Ph.D** selaku pembimbing kedua yang juga telah memberikan bimbingan, koreksi dan dukungan dalam penulisan tesis.
3. Bapak Drs. Bayu Surarso, M.Sc., Ph.D **Selaku Ketua Program Studi, Magister Sistem Informasi, Sekolah Pascasarjana, Universitas Diponegoro, Semarang.**
4. Bapak **Dr. R.B. Sularto, S.H., M.Hum, Selaku Dekan Sekolah Pascasarjana Universitas Diponegoro Semarang.**
5. Keluargaku tercinta, Bapak Drs. Pardi dan Ibu Sri Wahyuti, serta istri dan kakak-kakak saya.
6. Seluruh Dosen Magister Sistem Informasi, Universitas Diponegoro yang tidak dapat penulis sebutkan satu persatu, terimakasih atas ilmu yang diberikan selama perkuliahan.
7. Rekan-rekan Magister Sistem Informasi 2018 dan semua staf karyawan Pascasarjana, Universitas Diponegoro, terimakasih atas kerja sama yang telah terjalin selama ini, sehingga penulis menyelesaikan administrasi hingga menyelesaikan tesis ini.

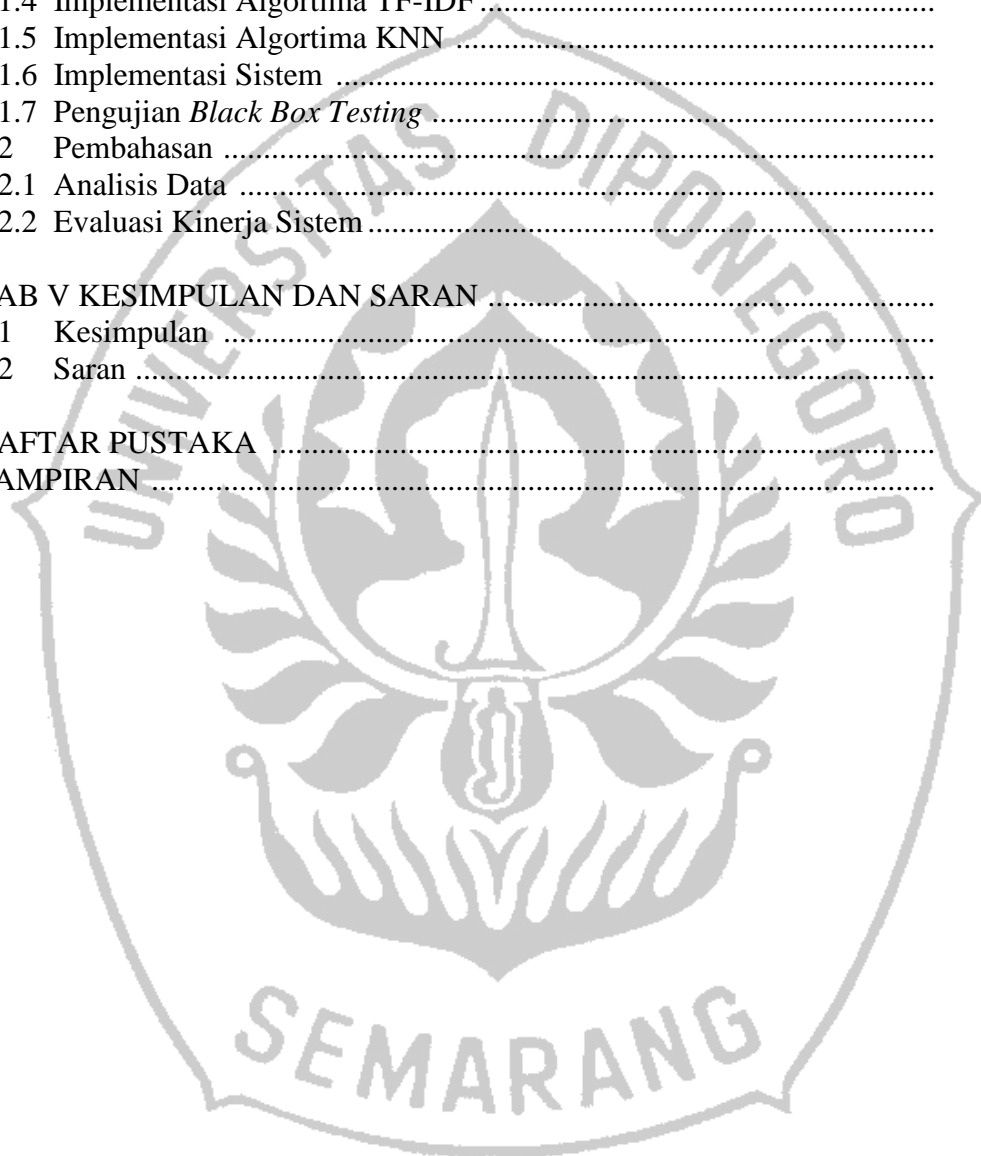
Penulis menyadari bahwa tesis ini masih jauh dari sempurna, untuk itu saran dan kritik yang membangun sangat penulis harapkan. Semoga tesis ini dapat bermanfaat bagi pihak-pihak yang membutuhkan.

Semarang, Juni 2022  
Penulis

## DAFTAR ISI

	Halaman
Halaman Judul .....	i
Halaman Persetujuan .....	ii
Halaman Pernyataan Persetujuan Publikasi .....	iii
Halaman Pernyataan .....	iv
Kata Pengantar .....	v
Daftar Isi .....	vi
Daftar Gambar .....	viii
Daftar Tabel .....	ix
Daftar Lampiran .....	x
Daftar Arti Lambang dan Singkatan .....	xi
Abstrak .....	xii
<i>Abstract</i> .....	xiii
<b>BAB I PENDAHULUAN</b> .....	<b>1</b>
1.1 Latar Belakang .....	1
1.2 Tujuan Penelitian .....	3
1.3 Manfaat Penelitian .....	3
<b>BAB II TINJAUAN PUSTAKA DAN DASAR TEORI</b> .....	<b>4</b>
2.1 Tinjauan Pustaka .....	4
2.2 Dasar Teori .....	5
2.2.1 <i>Data Mining</i> .....	5
2.2.2 Fungsi Penerapan <i>Data Mining</i> .....	6
2.2.3 Metode Pengembangan <i>Data Mining</i> .....	7
2.2.4 <i>Text Mining</i> .....	8
2.2.5 Analisis Sentimen .....	8
2.2.6 Prapengolahan Teks .....	9
2.2.7 Metode <i>Resampling</i> .....	10
2.2.8 Algoritma TF-IDF .....	11
2.2.9 Algoritma KNN ( <i>K-Nearest Neighbor</i> ) .....	12
2.2.10 Pengukuran Kinerja Klasifikasi .....	14
<b>BAB III METODE PENELITIAN</b> .....	<b>15</b>
3.1 Bahan dan Alat Penelitian .....	15
3.2 Prosedur Penelitian .....	15
3.3 Perancangan Sistem .....	19
3.3.1 Kerangka Sistem .....	19
3.3.2 Diagram Alir Sistem .....	21
3.4 Desain Sistem .....	22
3.4.1 Perancangan Diagram Alir Data .....	22
3.4.2 Pemodelan Sistem .....	24

3.4.3 Perancangan Basis Data .....	26
3.4.4 Perancangan Antar Muka .....	27
<b>BAB IV HASIL PENELITIAN DAN PEMBAHASAN .....</b>	<b>29</b>
4.1 Hasil Penelitian .....	29
4.1.1 Pemilihan Data .....	29
4.1.2 Prapengolahan Teks .....	30
4.1.3 Pelabelan Kelas Sentimen .....	30
4.1.4 Implementasi Algoritma TF-IDF .....	31
4.1.5 Implementasi Algoritma KNN .....	33
4.1.6 Implementasi Sistem .....	34
4.1.7 Pengujian <i>Black Box Testing</i> .....	37
4.2 Pembahasan .....	38
4.2.1 Analisis Data .....	38
4.2.2 Evaluasi Kinerja Sistem .....	39
<b>BAB V KESIMPULAN DAN SARAN .....</b>	<b>47</b>
5.1 Kesimpulan .....	47
5.2 Saran .....	47
<b>DAFTAR PUSTAKA .....</b>	<b>48</b>
<b>LAMPIRAN .....</b>	<b>51</b>



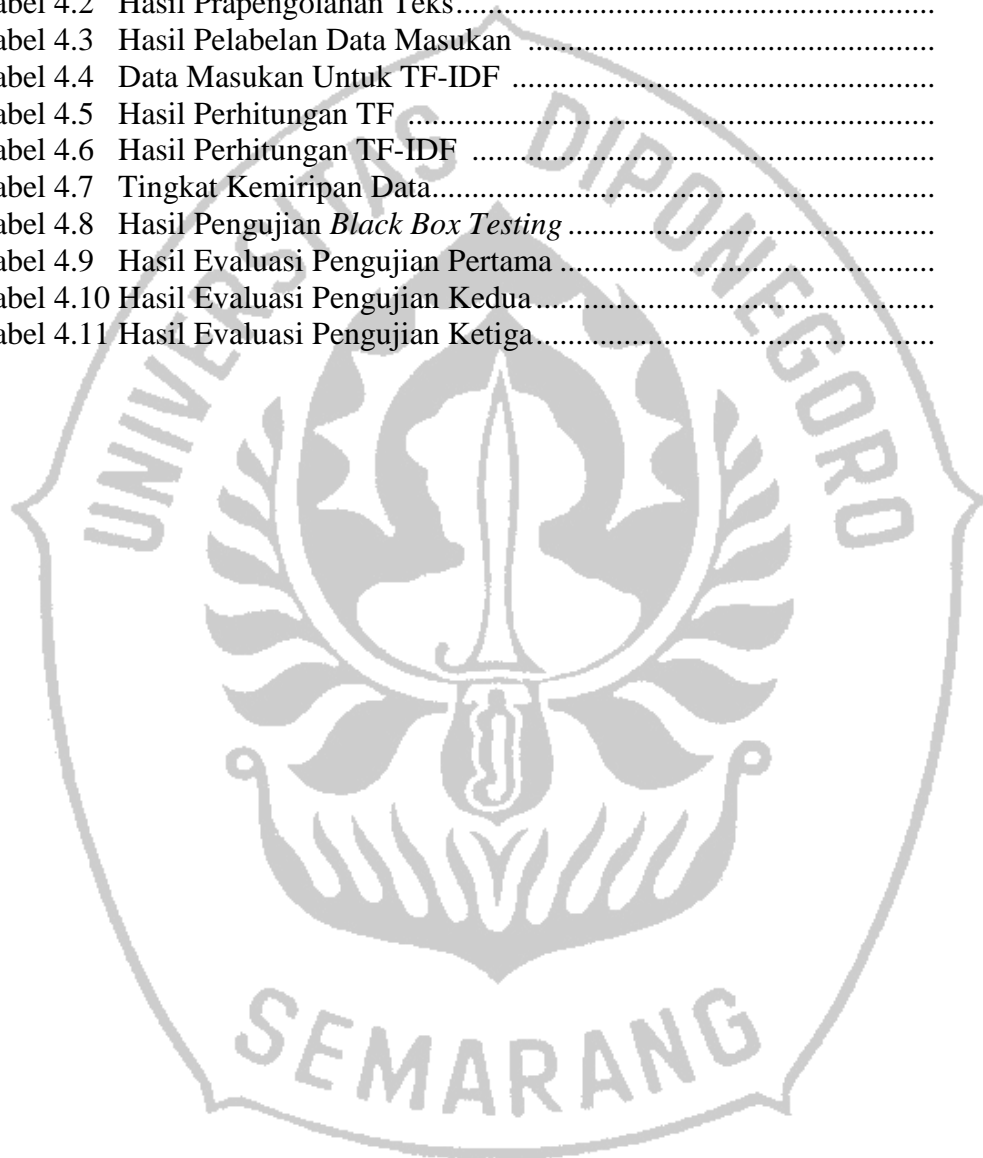
## DAFTAR GAMBAR

	Halaman
Gambar 3.1 SDLC dengan model <i>waterfall</i> .....	16
Gambar 3.2 Prosedur Penelitian.....	18
Gambar 3.3 Kerangka Sistem yang akan dibangun .....	20
Gambar 3.4 Diagram Alir Sistem.....	21
Gambar 3.5 DFD Level 0 .....	22
Gambar 3.6 DFD Level 1.....	23
Gambar 3.7 <i>Use Case</i> Diagram.....	24
Gambar 3.8 <i>Sequence</i> Diagram .....	25
Gambar 3.9 <i>Activity</i> Diagram .....	26
Gambar 3.10 Desain <i>Mockup</i> Sistem.....	27
Gambar 3.11 Desain Tampilan Utama Sistem.....	28
Gambar 4.1 Tampilan Modul Masukan Data <i>Training</i> .....	35
Gambar 4.2 Tampilan Modul Masukan Data uji .....	35
Gambar 4.3 Tampilan Modul Proses .....	36
Gambar 4.4 Tampilan Keluaran Sistem .....	36
Gambar 4.5 Statistik Data masukan .....	38
Gambar 4.6 Statistik Hasil <i>Rebalance</i> Data.....	39
Gambar 4.7 Hasil Evaluasi Tanpa <i>Rebalance</i> Data .....	39
Gambar 4.8 Hasil Evaluasi Setelah <i>Rebalance</i> Data .....	41
Gambar 4.9 Distribusi Kelas Pengujian Pertama.....	42
Gambar 4.10 Distribusi Kelas Pengujian Kedua .....	43
Gambar 4.11 Distribusi Kelas Pengujian Ketiga .....	45



## DAFTAR TABEL

	Halaman
Tabel 1.1 Penerapan Prapengolahan Teks .....	10
Tabel 4.1 Sampel data yang digunakan.....	29
Tabel 4.2 Hasil Prapengolahan Teks.....	30
Tabel 4.3 Hasil Pelabelan Data Masukan .....	31
Tabel 4.4 Data Masukan Untuk TF-IDF .....	32
Tabel 4.5 Hasil Perhitungan TF .....	32
Tabel 4.6 Hasil Perhitungan TF-IDF .....	33
Tabel 4.7 Tingkat Kemiripan Data.....	34
Tabel 4.8 Hasil Pengujian <i>Black Box Testing</i> .....	37
Tabel 4.9 Hasil Evaluasi Pengujian Pertama .....	42
Tabel 4.10 Hasil Evaluasi Pengujian Kedua.....	44
Tabel 4.11 Hasil Evaluasi Pengujian Ketiga.....	45



## DAFTAR LAMPIRAN

	Halaman
Lampiran 1. Tabel perhitungan TF .....	51
Lampiran 2. Tabel perhitungan TF-IDF.....	52



## DAFTAR ARTI LAMBANG DAN SINGKATAN

### DAFTAR ARTI LAMBANG

Lambang	Arti Lambang
$D_{xy}$	Jarak antara objek x dan y
$Cos_{ij}$	Nilai kemiripan data uji dengan data latih
$d_{ik}$	Panjang vektor data uji
$d_{jk}$	Panjang vektor data latih
$x$	Data uji yang akan diklasifikasi
$y$	Data latih
$n$	Jumlah data latih
$i$	Data ke- $i$
$Tf$	Banyaknya kata dalam dokumen
$Tf_{t,d}$	Banyaknya kata $t$ pada dokumen $d$
$d$	Jumlah dokumen
$df$	Jumlah dokumen yang mengandung suatu <i>term</i>
$W_{t,d}$	Nilai bobot akhir kata
$tf_{t,d}$	Banyaknya kata dalam dokumen
$idf_t$	Pembobotan keseluruhan

### DAFTAR SINGKATAN

Singkatan	Kepanjangan Singkatan
<i>SVM</i>	<i>Support Vector Machine</i>
<i>ANN</i>	<i>Artificial Neural Network</i>
<i>KNN</i>	<i>K-Nearest Neighbor</i>
<i>CNN</i>	<i>Convolutional Neural Network</i>
<i>SQL</i>	<i>Structured Query Language</i>
<i>TF-IDF</i>	<i>Term Frequency-Inverse Document Frequency</i>
<i>KDD</i>	<i>Knowledge Discovery in Database</i>
<i>NLTK</i>	<i>Natural Language Toolkit</i>
<i>API</i>	<i>Application Programming Interface</i>
<i>NLP</i>	<i>Natural Language Processing</i>
<i>UML</i>	<i>Unified Modelling Language</i>
<i>SDLC</i>	<i>Software Development Language Cycle</i>

**ANALISIS ULASAN HOTEL DI SITUS TRIPADVISOR  
MENGUNAKAN METODE *TERM FREQUENCY-INVERSE DOCUMENT  
FREQUENCY* DAN *K-NEAREST NEIGHBOR***

**ABSTRAK**

Penelitian ini dilatarbelakangi oleh evaluasi produk dan jasa menggunakan metode konvensional melalui wawancara, survei dan kuisioner yang berakibat pada hasil analisis menjadi tidak akurat dan tidak konsisten. Penelitian ini bertujuan untuk menerapkan algoritma *Term Frequency-Inverse Document Frequency* (TF-IDF) dan *K-Nearest Neighbor* serta mengevaluasi hasil dari sistem yang dibangun dengan tingkat akurasi yang paling optimal. Salah satu upaya untuk menanggulangi permasalahan tersebut yaitu membangun sistem untuk analisis ulasan pelanggan hotel di situs TripAdvisor yang bernilai positif, negatif dan netral menggunakan teknik *text mining*. Algoritma klasifikasi yang digunakan dalam penelitian adalah *K-Nearest Neighbor* karena memiliki kelebihan dalam hal komputasi berkinerja tinggi, tahan terhadap berbagai karakteristik data yang besar, dan memiliki kompleksitas algoritma yang relatif kecil. Hasil penelitian menunjukkan bahwa sistem dapat melakukan klasifikasi terhadap ulasan hotel di situs TripAdvisor yang bernilai positif, negatif dan netral dengan nilai performa paling baik pada  $K=31$  dan memiliki tingkat akurasi mencapai 76% untuk data *training*, dan menghasilkan peningkatan akurasi mencapai 84% dengan menerapkan metode *random over-sampling* untuk *rebalance* data.

Kata kunci : *text mining, term frequency-inverse document frequency, k-nearest neighbor, random over-sampling*

# **ANALYSIS OF HOTEL REVIEWS ON THE TRIPADVISOR SITE USING THE TERM FREQUENCY-INVERSE DOCUMENT FREQUENCY AND K-NEAREST NEIGHBOR METHODS**

## **ABSTRACT**

This research is based on analysis using conventional methods through interviews, surveys and questionnaires which resulted in the analysis being inaccurate and inconsistent. This study aims to apply the Term Frequency-Inverse Document Frequency (TF-IDF) and K-Nearest Neighbor algorithms and evaluate the results of the system built with the most optimal level of accuracy. To solve these problems is to build a system for analyzing hotel customer reviews on the TripAdvisor site which have positive, negative and neutral values using text mining techniques. The classification algorithm used in the research is K-Nearest Neighbor because it has advantages in terms of high-performance computing, is resistant to various characteristics of large data, and has a relatively small algorithm complexity. The results show that the system can classify hotel reviews on the TripAdvisor site which are positive, negative and neutral with the best performance value at  $K = 31$  and has an accuracy of 76% for training data and increase in accuracy of up to 84% by applying the random over-sampling method for data rebalance.

*Keywords: text mining, term frequency-inverse document frequency, k-nearest neighbor, random over-sampling*

