

**ANALISIS DATA ULASAN PENGUNJUNG MENGGUNAKAN  
*LEXICON BASED, SUPPORT VECTOR MACHINE, RANDOM  
FOREST* DALAM MENENTUKAN SKALA PRIORITAS  
PEMBANGUNAN OBJEK WISATA LABUAN BAJO**

**Tesis**

**untuk memenuhi sebagian persyaratan  
mencapai derajat Sarjana S-2 Program Studi  
Magister Sistem Informasi**



**SEKOLAH PASCASARJANA**

**Arnoldus Janssen Dahur**

**30000322410011**

**SEKOLAH PASCASARJANA  
UNIVERSITAS DIPONEGORO**

HALAMAN PENGESAHAN

SEMARANG

ANALISIS DATA ULASAN PENGGUNJUNG MENGGUNAKAN  
LEXICON BASED, SUPPORT VECTOR MACHINE, RANDOM  
FOREST DALAM MENENTUKAN SKALA PRIORITAS  
PEMBANGUNAN OBJEK SATELABUAN BAJO

Oleh:  
Arnoldus Janssen Dahur  
30000322410011

Telah diujikan dan dinyatakan lulus ujian tesis pada tanggal 23 Januari 2024 oleh tim penguji Program Studi Magister Sistem Informasi Sekolah Pascasarjana Universitas Diponegoro.

Semarang, 23 Januari 2024  
Mengetahui,

Penguji I



Dr. Budi Warsito, S.Si., M.Si.  
NIP 197508241999031003

Penguji II



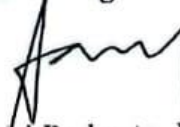
Dr. Drs. Catur Edi Widodo, M.T.  
NIP 196405181992031002

Pembimbing I



Dr. Eng. Wahyu Amien Syafei, S.T., M.T.  
NIP 197112181995121001


Pembimbing II



Ir. Toni Prahasto, MAsc., Ph.D.  
NIP 196208091988031001

Mengetahui :

Dekan Sekolah Pascasarjana  
Universitas Diponegoro



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

Ketua Program Studi  
Magister Sistem Informasi



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 : Arnoldus Janssen Dahur  
NIM : 30000322410011  
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 DATA ULASAN PENGUNJUNG MENGGUNAKAN *LEXICON BASED, SUPPORT VECTOR MACHINE, RANDOM FOREST* DALAM MENENTUKAN SKALA PRIORITAS PEMBANGUNAN OBJEK WISATA LABUAN BAJO**

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 : 23 Januari 2024  
Yang menyatakan



Arnoldus Janssen Dahur  
NIM 30000322410011

## 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, 23 Januari 2023

  
Arnoldus Janssen Dahur

SEKOLAH PASCASARJANA

## KATA PENGANTAR

Segala Puji dan syukur penulis panjatkan kehadirat Tuhan Yang Maha Esa yang telah melimpahkan segala rahmatNya sehingga penulis dapat menyelesaikan tesis dengan judul “Analisis Data Ulasan Pengunjung Menggunakan *Lexicon Based, Support Vector Machine, Random Forest* dalam Menentukan Skala Prioritas Pembangunan Objek Wisata Labuan Bajo” guna memenuhi sebagian persyaratan untuk memperoleh gelar magister pada program studi Magister Sistem Informasi, Sekolah Pascasarjana Universitas Diponegoro.

Penulis menyadari kelemahan serta keterbatasan yang ada sehingga dalam menyelesaikan tesis ini memperoleh bantuan dari berbagai pihak, dalam kesempatan ini penulis menyampaikan ucapan terimakasih kepada :

1. Bapak Dr. R.B. Sularto, S.H., M.Hum. selaku dekan Sekolah Pascasarjana Universitas Diponegoro.
2. Bapak Drs. Bayu Surarso, M.Sc., Ph.D. selaku Ketua Program Studi Magister Sistem Informasi Universitas Diponegoro.
3. Bapak Dr. Eng.Wahyul Amien Syafei, S.T.,MT. Selaku Dosen Pembimbing I, yang telah memberikan bimbingan, arahan, serta kritik kepada penulis hingga bisa selesai dalam penyusunan tesis.
4. Bapak Ir. Toni Prahasto, MAsc., Ph.D. selaku dosen Pembimbing II, yang sudah membimbing, memberikan arahan, dan masukan yang baik kepada penulis sehingga dapat penulis menyelesaikan laporan tesis ini.
5. Bapak Dr. Budi Warsito, S.Si., M.Si. selaku Dosen Ketua Penguji yang telah memberikan banyak saran dan masukan yang baik kepada penulis.
6. Bapak Dr. Drs. Catur Edi Widodo, M.T. sebagai Dosen Penguji yang telah memberikan banyak saran dan masukan dalam penyusunan tesis.
7. Semua pihak yang tidak dapat disebutkan satu persatu, yang telah membantu sampai dengan terselesaikannya tesis ini.

Penulis menyadari bahwa tesis ini masih banyak kekurangan baik isi maupun susunannya. Semoga tesis ini dapat bermanfaat tidak hanya bagi penulis juga bagi para pembaca.

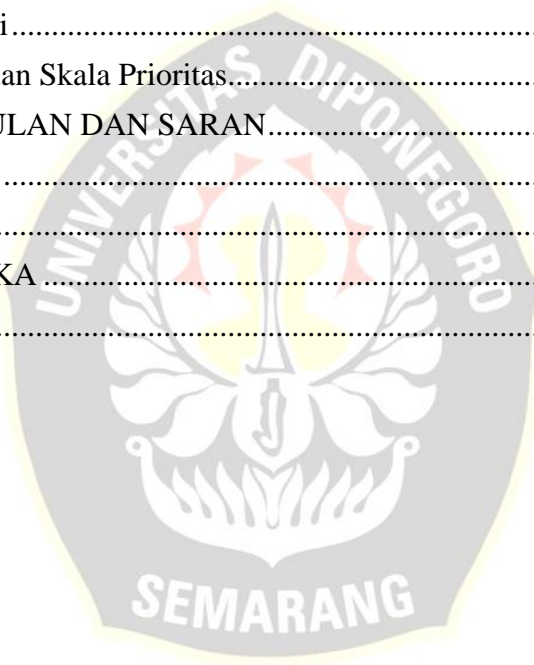
Semarang, 30 Januari 2024

Arnoldus Janssen Dahur

## DAFTAR ISI

Halaman Judul.....	i
Halaman Pengesahan .....	ii
PERNYATAAN PERSETUJUAN .....	iii
PERNYATAAN.....	iv
KATA PENGANTAR .....	v
DAFTAR ISI.....	vi
DAFTAR GAMBAR .....	viii
DAFTAR TABEL.....	ix
Daftar Arti Lambang dan Singkatan.....	ix
ABSTRAK .....	xi
ABSTRACT.....	xii
BAB I PENDAHULUAN.....	1
1.1. Latar Belakang.....	1
1.2. Tujuan Penelitian.....	5
1.3. Manfaat Penelitian.....	5
BAB II KERANGKA TEORI.....	6
2.1 Tinjauan Pustaka.....	6
2.2 Dasar Teori .....	8
2.2.1 <i>Web Scraping</i> .....	8
2.2.2 <i>Text Mining</i> .....	9
2.2.3 Analisis Sentimen .....	11
2.2.4 Pembobotan Kata.....	12
2.2.5 Lexicon Based.....	14
2.2.6 Support Vector Machine .....	16
2.2.7 Random Forest.....	19
2.2.8 Klasifikasi .....	20
BAB III METODE PENELITIAN.....	23
3.1. Bahan dan Alat Penelitian.....	23
3.1.1. Bahan Penelitian .....	23
3.1.2. Alat Penelitian.....	24

3.2.    Prosedur Penelitian .....	25
<b>BAB IV HASIL PENELITIAN DAN PEMBAHASAN .....</b>	<b>31</b>
4.1 Hasil Penelitian .....	31
4.1.1 Pengumpulan data .....	31
4.1.2 Pelabelan Data.....	32
4.1.3 Preprocessing Data.....	35
4.1.4 Klasifikasi .....	37
4.1.5 Sistem Informasi .....	39
4.2 Pembahasan.....	46
4.2.1 Evaluasi.....	46
4.2.2 Penentuan Skala Prioritas.....	50
<b>BAB V KESIMPULAN DAN SARAN.....</b>	<b>54</b>
5.1 Kesimpulan .....	54
5.2 Saran.....	55
<b>DAFTAR PUSTAKA .....</b>	<b>56</b>
<b>LAMPIRAN.....</b>	<b>59</b>



**SEKOLAH PASCASARJANA**

## DAFTAR GAMBAR

Gambar 2. 1 Proses SVM dalam Menemukan Hyperline .....	17
Gambar 2. 2 Klasifikasi ruang data.....	17
Gambar 2. 3 Jarak D dari titik X.....	19
Gambar 2. 4 Representasi Tahapan Random Forest. ....	20
Gambar 2. 5 Proses Klasifikasi. ....	21
Gambar 3. 1 Flowchart Prosedur Penelitian .....	26
Gambar 3. 2 Flowchart Prosedur Penelitian .....	27
Gambar 4. 1 Proses Scrapping data menggunakan Web Scrapper .....	32
Gambar 4.2 Hasil Pelabelan menggunakan Lexicon Based.....	34
Gambar 4.3 Tampilan Hasil Preprocessing Data .....	37
Gambar 4.4 Hasil klasifikasi tanpa menggunakan teknik undersampling .....	38
Gambar 4.5 Hasil klasifikasi tanpa menggunakan teknik undersampling .....	39
Gambar 4.6 Tampilan menu sistem informasi .....	40
Gambar 4.7 Tampilan Halaman Dashboard.....	41
Gambar 4.8 Tampilan dataset dan tampilan translate data .....	41
Gambar 4.9 Tampilan Hasil Pelabelan.....	42
Gambar 4.10 Tampilan Preprocessing .....	42
Gambar 4.11 Tampilan Klasifikasi .....	43
Gambar 4.12 Prediksi kalimat positif dan kalimat negatif.....	44
Gambar 4.13 Tampilan Skala Prioritas .....	45
Gambar 4. 14 Hasil ROC AUC tanpa menggunakan teknik Undersampling .....	48
Gambar 4.15 Hasil ROC AUC menggunakan teknik Undersampling.....	49

SEKOLAH PASCASARJANA



## DAFTAR TABEL

Tabel 2. 1 Confusion Matrix. ....	21
Tabel 3. 1 Jumlah data ulasan pada website TripAdvisor.....	24
Tabel 3. 2 Perangkat keras dan perangkat lunak yang digunakan .....	24
Tabel 4.1 Contoh Poscore dan NegScore masing-masing kata pada kamus Sentiwordnet .....	33
Tabel 4.2 Contoh kata-kata dalam kamus stopwords .....	35
Tabel 4.3 Contoh kata-kata sebelum normalizing dan setelah normalizing .....	36
Tabel 4.4 Contoh kata-kata sebelum proses stemming dan sesudah stemming.....	36
Tabel 4. 5 Confusion Matrix tanpa menggunakan Teknik Undersampling .....	46
Tabel 4. 6 Confusion Matrix dengan menggunakan Teknik Undersampling .....	47
Tabel 4.7 Hasil confusion matrix teknik undersampling dan tanpa undersampling .....	47
Tabel 4.8 Nilai rentang AUC dan tingkat klasifikasi menurut Gorunescu (2011).49	
Tabel 4.9 10 Kata Positif teratas dan jumlah sentiment setiap kata.....	50
Tabel 4. 10 10 Kata negatif teratas dan jumlah sentiment setiap kata .....	51
Tabel 4.11 Wordcloud sentiment positif setiap destinasi .....	52
Tabel 4.12 Wordcloud sentiment negatif setiap destinasi.....	52

SEKOLAH PASCASARJANA

## Daftar Arti Lambang dan Singkatan

### Daftar Arti Lambang

Lambang	Arti Lambang
$W^T$	Parameter Bobot
$x$	Vector Input
$b$	bias
$\ w\ $	Norma Euclidean dari vektor bobot $w$
$D$	Jarak support vector dengan hyperplane
$x_1, x_2$ ... $x_n$	Data miror acak
$tf_{t,d}$	Jumlah yang muncul di <i>term</i> $t$ pada dokumen $d$
$N$	Jumlah yang terdapat di dokumen teks
$df_t$	Jumlah dokumen yang terdapat term $t$
$W_{tf_{t,d}}$	Nilai term frequency
$df_t$	Jumlah dokumen yang terdapat term
<i>Recall</i>	Rasio prediksi benar positif dibandingkan dengan keseluruhan data yang benar positif
<i>Precision</i>	Rasio prediksi benar positif dibandingkan dengan keseluruhan hasil yang diprediksi positif
$\log_{10} \frac{N}{df_t}$	Log dari jumlah keseluruhan dokumen dibagi dengan jumlah kata yang sama pada semua dokumen

### Daftar Singkatan

Singkatan	Kepanjangan Singkatan
TP	<i>True Positive</i>
TN	<i>True Negative</i>
FP	<i>False Positive</i>
FN	<i>False Negative</i>
HTML	<i>HyperText Markup Language</i>
HTML	<i>Extensible HyperText Markup Language</i>
CSV	<i>Comma-separated- values</i>
SVM	<i>Support Vector Machine</i>
RF	<i>Random Forest</i>
ROC	<i>Receiver Operating Characteristic</i>
AUC	<i>Area Under the Curve</i>
CRISP-DM	<i>Cross-Industry Standard Process for Data Mining</i>
TF-IDF	<i>Term Frequency-Inverse Document Frequency</i>

# **ANALISIS DATA ULASAN PENGUNJUNG MENGGUNAKAN LEXICON BASED, SUPPORT VECTOR MACHINE, RANDOM FOREST DALAM MENENTUKAN SKALA PRIORITAS PEMBANGUNAN OBJEK WISATA LABUAN BAJO**

## **ABSTRAK**

Objek wisata Labuan Bajo merupakan salah satu destinasi wisata super prioritas di Indonesia. Pentingnya mendapatkan dan menganalisis ulasan pengunjung wisata untuk mengetahui preferensi berupa pandangan pengunjung terhadap fasilitas dan pelayanan yang ada saat ini. Oleh karena itu penelitian ini dilakukan untuk mendapatkan dan menganalisis data ulasan pengunjung yang didapat dari website TripAdvisor dan Google Maps. Adapun metode yang digunakan dalam analisis ulasan pengunjung ini yaitu Lexicon Based untuk melakukan pelabelan, metode Support Vector Machine (SVM) dan Random Forest untuk klasifikasi. Hasil pelabelan menggunakan metode Lexicon Based didapat sentiment positif sejumlah 4187 ulasan, sentiment negatif sejumlah 1796 ulasan dan sentiment netral sejumlah 1774 ulasan. Klasifikasi dilakukan dengan menggunakan teknik *undersampling* dan tanpa menggunakan teknik *undersampling* karena ketidak seimbangan data. Hasil menggunakan teknik *undersampling* dengan SVM yaitu *accuracy* 0.89 *precisi* 0.95 *recall* 0.85 dan *f1-measure* 0.90 serta nilai ROC AUC menggunakan teknik *undersampling* yaitu 0.94 dan tanpa menggunakan teknik *undersampling* *accuracy* 0.79 *presisi* 0.80 *recall* 0.94 dan *f1-measure* 0.86 serta nilai ROC AUC yaitu 0.83. Hasil menggunakan teknik *undersampling* dengan Random Forest yaitu *accuracy* 0.87 *precisi* 0.91 *recall* 0.86 dan *f1-measure* 0.88 serta nilai ROC AUC menggunakan teknik *undersampling* yaitu 0.93 dan tanpa menggunakan SMOTE *accuracy* 0.77 *presisi* 0.78 *recall* 0.94 dan *f1-measure* 0.85 serta nilai ROC AUC yaitu 0.81. Penentuan skala prioritas dilakukan dengan mendapatkan 10 kata teratas dan jumlah sentiment dari masing-masing kata yang berkaitan dengan pembangunan didapat kata-kata sentimen positif yang sering muncul yaitu 'indah', 'alami', 'eksotik', 'pandang', 'bersih', 'purba', 'takjub', 'sejarah'. Pelestarian aset alami dan aset sejarah tentunya harus dijaga dan terus dipertahankan. Sebaliknya kata-kata negatif yang sering muncul yaitu 'mahal', 'biaya', 'pandu', 'jalan', 'sampah', 'panas'. Berdasarkan kata tersebut pembangunan transportasi dan infrastruktur tentunya sangat diperlukan dalam peningkatan daya tarik wisata Labuan Bajo.

Kata kunci: Analisis Ulasan, Labuan Bajo, Lexicon Based, Support Vector Machine, Random Forest

**ANALYSIS OF VISITOR REVIEW DATA USING *LEXICON BASED*,  
*SUPPORT VECTOR MACHINE*, *RANDOM FOREST* IN DETERMINING  
THE PRIORITY SCALE OF BUILDING LABUAN BAJO TOURISM  
OBJECTS**

**ABSTRACT**

Labuan Bajo tourist destination is one of the super priority tourist destinations in Indonesia. The importance of obtaining and analyzing tourists' reviews is to understand their preferences and views on the existing facilities and services. Therefore, this research is conducted to obtain and analyze visitor review data obtained from TripAdvisor and Google Maps. The methods used in analyzing these visitor reviews are Lexicon-Based for labeling, Support Vector Machine (SVM), and Random Forest for classification. The labeling results using the Lexicon-Based method showed 4187 positive reviews, 1796 negative reviews, and 1774 neutral reviews. The classification was performed using undersampling technique and without using undersampling technique due to data imbalance. Results using undersampling technique with SVM showed an accuracy of 0.89, precision of 0.95, recall of 0.85, and f1-measure of 0.90, with an ROC AUC value of 0.94. Without using undersampling technique, the accuracy was 0.79, precision was 0.80, recall was 0.94, and f1-measure was 0.86, with an ROC AUC value of 0.83. Results using undersampling technique with Random Forest showed an accuracy of 0.87, precision of 0.91, recall of 0.86, and f1-measure of 0.88, with an ROC AUC value of 0.93. Without using undersampling technique, the accuracy was 0.77, precision was 0.78, recall was 0.94, and f1-measure was 0.85, with an ROC AUC value of 0.81. The determination of priority scale was done by obtaining the top 10 words and the number of sentiments related to development. The frequently occurring positive sentiment words were 'beautiful,' 'natural,' 'exotic,' 'scenic,' 'clean,' 'ancient,' 'amazed,' and 'historical.' The preservation of natural and historical assets must be maintained and continuously preserved. On the other hand, the frequently occurring negative words were 'expensive,' 'cost,' 'guide,' 'road,' 'garbage,' and 'hot.' Based on these words, the development of transportation and infrastructure is undoubtedly needed to enhance the attractiveness of Labuan Bajo as a tourist destination.

*Keywords:* Review Analysis, Labuan Bajo, Lexicon Based, Support Vector Machine, Random Forest.