

IMPLEMENTASI *FOCUSED CRAWLING* DENGAN
TOPIC-SPECIFIC, MULTINOMIAL NAIVE
BAYES, DAN *BREADTH FIRST SEARCH* GUNA
PENGUMPULAN DATA *MEDIA MONITORING*
GEOPARK CILETUH

TUGAS AKHIR



CLARA VELITA PRANOLO

1162001016

PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNIK DAN ILMU KOMPUTER
UNIVERSITAS BAKRIE
JAKARTA
2021

**IMPLEMENTASI *FOCUSED CRAWLING* DENGAN
*TOPIC-SPECIFIC, MULTINOMIAL NAIVE
BAYES*, DAN *BREADTH FIRST SEARCH* GUNA
PENGUMPULAN DATA *MEDIA MONITORING*
GEOPARK CILETUH**

TUGAS AKHIR

Diajukan sebagai salah satu syarat untuk memperoleh gelar
Sarjana Komputer



CLARA VELITA PRANOLO

1162001016

**PROGRAM STUDI INFORMATIKA
FAKULTAS TEKNIK DAN ILMU KOMPUTER
UNIVERSITAS BAKRIE
JAKARTA**

2021

Halaman Pernyataan Orisinalitas

Tugas Akhir ini adalah hasil karya saya sendiri, dan semua sumber baik yang dikutip maupun dirujuk telah saya nyatakan dengan benar.

Nama : Clara Velita Pranolo

NIM : 1162001016

Tanda Tangan :

A handwritten signature in black ink, appearing to read 'Clara Velita', written over two horizontal lines. The signature is stylized and cursive.

Tanggal : 15 April 2021

Halaman Pengesahan

Tugas akhir ini diajukan oleh :

Nama : Clara Velita Pranolo

NIM : 1162001016

Program Studi : Teknik Informatika

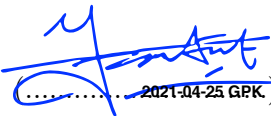
Fakultas : Teknik dan Ilmu Komputer

Judul Tugas Akhir : Implementasi Focused Crawling dengan Topic-Specific, Multinomial Naive Bayes, dan Breadth First Search guna Pengumpulan Data Media Monitoring Geopark Ciletuh

Telah berhasil dipertahankan dihadapan Dewan Penguji sebagai persyaratan yang diperlukan untuk memperoleh gelar Sarjana Komputer pada Program Studi Informatika Fakultas Teknik dan Ilmu Komputer, Universitas Bakrie.


DEWAN PENGUJI

Pembimbing : Guson P. Kuntarto, S.T., M.Sc.



(..... 2021-04-25 GPK)

Penguji 1 : Yusuf Lestanto, S.T., M.Sc.



(..... 15.04.2021)

Penguji 2 : Irwan Prasetya Gunawan, ST, MEng, PhD.



(..... 202104252501161910376)

Ditetapkan : Jakarta

Tanggal : 15 April 2021

Ungkapan Terima Kasih

Puji Syukur dipanjatkan kepada Allah SWT atas petunjuk, rahmat dan hidayah-Nya, sehingga terselesaikannya Tugas Akhir ini yang berjudul "Implementasi *Focused Crawling* dengan *Topic-Specific, Multinomial Naïve Bayes*, dan *Breadth First Search* guna Pengumpulan Data *Media Monitoring* Geopark Ciletuh" sebagai salah satu syarat dalam mencapai gelar Sarjana Komputer di Program Studi Informatika, Fakultas Teknik dan Ilmu Komputer, Universitas Bakrie.

Penelitian dan penyusunan Tugas Akhir ini tidak akan terwujud tanpa dukungan serta bantuan dari berbagai pihak. Oleh sebab itu, penulis ingin menyampaikan terima kasih yang sebesar-besarnya kepada :

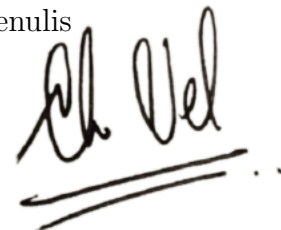
1. Kedua Orang Tua, Ibu Amalia Virnawati dan Bapak Tato Pranolo terima kasih untuk kasih sayang, kerja keras, pengorbanan, pengertian, semangat, serta dukungan yang tidak pernah ada hentinya diberikan untuk penulis dalam menggapai cita-cita sehingga terselesaikannya Tugas Akhir ini.
2. Bapak Guson P. Kuntarto, S.T., M.Sc. selaku dosen Pembimbing Tugas Akhir, yang selalu dengan sabar membimbing saya, selalu dapat meluangkan waktu, memberikan nasihat serta dukungan selama proses penelitian dan penyusunan Tugas Akhir ini.
3. Bapak Yusuf Lestanto, S.T., M.Sc. selaku dosen pembahas Tugas Akhir yang telah memberikan masukan dan perbaikan dalam penyusunan Tugas Akhir ini
4. Bapak Irwan Prasetya Gunawan, ST, MEng, PhD. selaku dosen pembahas Tugas Akhir serta dosen Pembimbing Akademik yang telah memberikan masukan, nasihat dan membimbing penulis dalam melangkah selama di perkuliahan sampai dengan di Tugas Akhir ini.

5. Bapak Berkah I. Santoso, S.T., M.T.I. yang pernah menjadi dosen Pembimbing Akademik yang telah memberikan masukan dan membimbing di semester awal perkuliahan.
6. Seluruh Bapak dan Ibu dosen Program Studi Informatika yang telah mendidik dan memberikan ilmu sehingga bertambahnya pengetahuan penulis yang tentunya sangat membantu dalam proses Tugas Akhir ini.
7. Mutiara Julia Ifra, Pilipus Delevia Vegas, Dezan Andhika, Hafiz Kurnia Aji, dan Fitrah Cahya yang selalu memberikan dukungan, pengertian dan baik selama diperkuliahan. Khususnya Mutiara Julia Ifra yang selalu sabar mendengarkan keluh kesah, memberikan dukungan supaya penulis dapat percaya diri menyelesaikan Tugas Akhir ini.
8. Teman-teman Informatika 2016, yang telah bersama selama 4 tahun.
9. Seluruh pihak Program Studi Informatika Universitas Bakrie yang telah memberikan ilmu dan pembelajaran serta pengalaman yang sangat bermanfaat bagi peneliti selama di perkuliahan.

Meskipun menyadari bahwa penelitian ini masih jauh dari kata sempurna, peneliti berharap bahwa penelitian dan Tugas Akhir ini dapat memberikan manfaat dan berguna bagi kalangan pendidikan, khususnya Informatika.

Jakarta, 15 April 2021

Penulis

A handwritten signature in black ink, appearing to read 'Clara Velita Pranolo', written over a double horizontal line.

Clara Velita Pranolo

Halaman Pernyataan Persetujuan Publikasi

Sebagai sivitas akademik Universitas Bakrie, saya yang bertanda tangan dibawah ini:

Nama : Clara Velita Pranolo
NIM : 1162001016
Program Studi : Teknik Informatika
Fakultas : Teknik dan Ilmu Komputer

demi pengembangan ilmu pengetahuan, menyetujui untuk memberikan kepada Universitas Bakrie **Hak Bebas Royalti Noneksklusif (*Non-exclusive Royalty-Free Right*)** atas karya ilmiah saya yang berjudul:

Implementasi Focused Crawling dengan Topic-Specific, Multinomial Naive Bayes, dan Breadth First Search guna Pengumpulan Data Media Monitoring Geopark Ciletuh

beserta perangkat yang ada (jika diperlukan). Dengan Hak Bebas Royalti Noneksklusif ini Universitas Bakrie berhak menyimpan, mengalihmedia/formatkan, mengelola dalam bentuk pangkalan data (database), merawat, dan mempublikasikan tugas akhir saya selama tetap mencantumkan nama saya sebagai penulis/pencipta dan sebagai pemilik Hak Cipta untuk kepentingan akademis.

Demikian pernyataan ini saya buat dengan sebenarnya.

Dibuat di : Jakarta

Pada tanggal : 15 April 2021

Yang Menyatakan,



Clara Velita Pranolo

IMPLEMENTASI *FOCUSED CRAWLING* DENGAN *TOPIC-SPECIFIC, MULTINOMIAL NAIVE* *BAYES*, DAN *BREADTH FIRST SEARCH* GUNA PENGUMPULAN DATA *MEDIA MONITORING* GEOPARK CILETUH

Clara Velita Pranolo

ABSTRAK

Sebagai bagian dari *UNESCO Global Geopark* (UGG), wisata Geopark Ciletuh memiliki praktisi *Public Relations* (PR) yang mempunyai tujuan dalam hal membangun, mengembangkan, dan mempertahankan reputasi serta citra dari kawasan wisata Geopark Ciletuh. Untuk mencapai tujuan tersebut, praktisi PR melakukan *media monitoring* dengan mencari dan memilih berita dari berbagai sumber *media online*, lalu mengumpulkannya kedalam klipng berita untuk diidentifikasi dan dianalisis. Dalam *media monitoring*, proses ini disebut *data back-end*. Penelitian ini mengusulkan *focused crawling* untuk diimplementasikan pada *data back-end media monitoring* Geopark Ciletuh supaya proses pengumpulan data lebih cepat. *Focused crawling* diimplementasikan dengan menggunakan tiga metode yaitu metode *crawler* dengan Algoritma *Breadth First Search* (BFS) untuk mendapatkan URL berita yang lebih banyak, metode *distiller* dengan *Topic-Specific Weight Table* dan *Page Relevance* untuk fitur parameter dataset, serta metode klasifikasi dengan *Multinomial Naïve Bayes* untuk menentukan berita yang relevan. Hasil penelitian dengan algoritma BFS dapat melakukan *crawling* sebanyak 470 URL untuk Detik dan 290 URL untuk Kompas. Sedangkan dalam menentukan berita yang relevan akurasi yang didapatkan model *Multinomial Naïve Bayes* dengan *Page Relevance* yaitu 83.46% untuk *dataset* Detik, 89% untuk *dataset* Kompas dan diatas 88.16% untuk kedua gabungan *dataset* Detik dan Kompas.

Kata Kunci : *Web Crawling, Focused Crawling, Breadth First Search, Topic-Specific Weight Table, Page Relevance, Multinomial Naïve Bayes*

IMPLEMENTATION OF FOCUSED CRAWLING USING TOPIC-SPECIFIC, MULTINOMIAL NAIVE BAYES, AND BREADTH FIRST SEARCH FOR MEDIA MONITORING DATA COLLECTION GEOPARK CILETUH

Clara Velita Pranolo

ABSTRACT

As part of the UNESCO Global Geopark (UGG), Ciletuh Geopark has Public Relations (PR) practitioners who has the purpose to building, developing, and maintaining the reputation and image of the Ciletuh Geopark tourist area. To achieve this purpose, PR practitioners conduct media monitoring by searching and selecting news from various online media sources, then collecting them into news clippings to be identified and analyzed. In media monitoring, this process is called backend data. This study proposes focused crawling to be implemented in the back-end data of the Ciletuh Geopark media monitoring so the data collection process more effectively. Focused crawling is implemented using three methods, the crawler method with the Breadth First Search (BFS) algorithm to get more news URLs, the distiller method with the Topic-Specific Weight Table and Page Relevance for the feature parameter dataset, and the classification method with Multinomial Naïve Bayes for determine relevant news. The results of this study that BFS Algorithm crawl 470 URLs on Detik dataset and 290 URLs on Kompas dataset. The result also showed that Multinomial Naïve Bayes using the Page Relevance feature. The result also showed that Multinomial Naïve Bayes classify URLs based on Page Relevance with accuracy 83.46% on Detik dataset, 89% on Kompas Dataset and above 88.16% on combination of Detik and Kompas dataset.

Keywords : Web Crawling, Focused Crawling, Breadth First Search, Topic-Specific Weight Table, Page Relevance, Multinomial Naïve Bayes

Daftar Isi

| | |
|--|----------|
| Halaman Pernyataan Orisinalitas | i |
| Halaman Pengesahan | ii |
| Ungkapan Terima Kasih | iii |
| Halaman Pernyataan Persetujuan Publikasi | v |
| Abstrak | vi |
| Abstract | vii |
| Daftar Isi | vii |
| Daftar Tabel | xi |
| Daftar Gambar | xiii |
| Daftar Skrip | xv |
| Daftar Singkatan | xvii |
| I Pendahuluan | 1 |
| 1.1 Latar Belakang | 1 |
| 1.2 Rumusan Masalah | 5 |
| 1.3 Tujuan Penelitian | 6 |
| 1.4 Manfaat Penelitian | 6 |
| 1.5 Ruang Lingkup Penelitian | 6 |

| | | |
|------------|--|-----------|
| 1.6 | Kontribusi Penelitian | 7 |
| 1.7 | Sistematika Penulisan | 8 |
| II | Tinjauan Pustaka | 10 |
| 2.1 | Penelitian Terkait | 10 |
| 2.2 | <i>Media Monitoring</i> dan Pengumpulan Informasi | 15 |
| 2.3 | <i>Web Crawling</i> | 16 |
| 2.3.1 | <i>Focused Crawling</i> | 17 |
| 2.4 | <i>Web Scraping</i> | 18 |
| 2.5 | <i>Topic-Specific Weight Table (Topic-Specific / TSWT)</i> | 19 |
| 2.6 | <i>Page Relevance</i> | 20 |
| 2.7 | <i>Breadth First Search (BFS)</i> | 20 |
| 2.8 | <i>Multinomial Naïve Bayes</i> | 21 |
| 2.9 | <i>Pengukuran Performa Keakuratan</i> | 22 |
| III | Metodologi Penelitian | 24 |
| 3.1 | Tahapan Penelitian | 25 |
| 3.1.1 | Studi Literatur | 26 |
| 3.1.2 | Analisis Masalah dan Tujuan | 26 |
| 3.1.3 | Pengumpulan Data | 27 |
| 3.1.4 | Pelaksanaan Penelitian | 27 |
| 3.1.5 | Hasil Implementasi dan Pembahasan | 27 |
| 3.1.6 | Penyusunan Hasil Laporan Penelitian | 28 |
| 3.2 | Kerangka Kerja | 30 |
| 3.2.1 | <i>Dataset</i> | 30 |
| 3.2.2 | <i>Training Data</i> | 31 |
| 3.2.3 | <i>Testing Data</i> | 31 |
| 3.2.4 | Implementasi <i>Breadth First Search (BFS)</i> | 32 |
| 3.2.5 | Implementasi <i>Topic-Specific Weight Table</i> | 33 |
| 3.2.6 | Implementasi <i>Page Relevance</i> | 34 |

| | | |
|-----------|--|-----------|
| 3.2.7 | Implementasi <i>Multinomial Naïve Bayes</i> | 35 |
| 3.2.8 | <i>Scraping</i> data dari URL | 35 |
| 3.2.9 | Uji Skenario dan Hasil Pengukuran Performa | 40 |
| 3.3 | Instrumen atau Alat Penelitian | 40 |
| IV | Implementasi dan Hasil Penelitian | 41 |
| 4.1 | Pengumpulan Data | 41 |
| 4.1.1 | <i>Breadth First Search</i> (BFS) | 43 |
| 4.2 | <i>Topic-Specific Weight Table</i> (<i>Topic-Specific</i> / TSWT) | 46 |
| 4.3 | <i>Page Relevance</i> | 48 |
| 4.3.1 | <i>URL Word</i> | 48 |
| 4.3.2 | <i>Parent Page</i> | 50 |
| 4.3.3 | <i>Anchor Text</i> | 51 |
| 4.3.4 | <i>Surrounding Text</i> | 52 |
| 4.4 | <i>Labeling Dataset</i> | 52 |
| 4.5 | Pembentukan <i>Multinomial Naïve Bayes</i> | 53 |
| 4.6 | <i>Scraping</i> data dari URL | 57 |
| 4.7 | Eksperimen | 57 |
| 4.7.1 | Hasil Uji Coba <i>Multinomial Naïve Bayes</i> | 57 |
| 4.7.2 | Hasil Uji Coba <i>Breadth First Search</i> (BFS) | 61 |
| 4.7.3 | Hasil Uji Coba <i>Multithread Scraping</i> | 63 |
| 4.8 | Pembahasan | 64 |
| V | Simpulan dan Saran | 66 |
| 5.1 | Simpulan | 66 |
| 5.2 | Saran | 67 |
| A | Dataset <i>Topic-Specific Weight Table</i> | 72 |
| B | Data Training | 73 |
| C | Data Testing Detik | 96 |

| | |
|---|------------|
| D Data Testing Kompas | 119 |
| E Hasil Uji Coba Multinomial Naïve Bayes | 134 |
| F Source Code | 139 |

Daftar Tabel

| | | |
|-----|---|-----|
| 2.1 | Rangkuman Penelitian Terkait | 14 |
| 4.1 | <i>Topic-Specific Weight Table</i> yang digunakan pada penelitian | 48 |
| 4.2 | Tahap <i>Text Preprocessing</i> pada <i>URL Word</i> | 49 |
| 4.3 | Potongan Hasil <i>Page Relevance Training Data</i> | 52 |
| 4.4 | Potongan Hasil <i>Page Relevance Data Detik</i> | 52 |
| 4.5 | Potongan Hasil <i>Page Relevance Data Kompas</i> | 52 |
| 4.6 | Pembagian Dataset untuk Pengujian Gambar 4.5 | 59 |
| 4.7 | Hasil Uji Coba Breadth First Search (BFS) pada <i>Crawling URL</i> | 61 |
| D.1 | Testing Data Kompas | 119 |
| E.1 | Skenario ke-1 MNB Dataset Gabungan - Training Seluruh Tahun, Testing Seluruh Tahun | 134 |
| E.2 | Skenario ke-1 Improved MNB Dataset Gabungan - Training Seluruh Tahun, Testing Seluruh Tahun | 134 |
| E.3 | Skenario ke-1 MNB Dataset Detik dan Kompas - Training Seluruh Tahun, Testing Seluruh Tahun | 135 |
| E.4 | Skenario ke-1 Improved MNB Dataset Detik dan Kompas - Training Seluruh Tahun, Testing Seluruh Tahun | 135 |
| E.5 | Skenario ke-2 MNB - Training Pertahun, Testing Pertahun | 136 |
| E.6 | Skenario ke-2 Improved MNB - Training Pertahun, Testing Pertahun | 136 |
| E.7 | Skenario ke-3 MNB - Training Seluruh Tahun, Testing Pertahun | 137 |
| E.8 | Skenario ke-3 Improved MNB - Training Seluruh Tahun, Testing Pertahun | 137 |
| E.9 | Skenario ke-4 MNB - Training Pertahun, Testing Seluruh Tahun | 138 |

| | |
|---|-----|
| E.10 Skenario ke-4 Improved MNB - Training Pertama, Testing Seluruh Tahun | 138 |
|---|-----|

Daftar Gambar

| | | |
|------|---|----|
| 1.1 | Arsitektur Media Monitoring | 2 |
| 1.2 | <i>Data Back-end</i> sebelumnya vs Solusi Penelitian dengan <i>Focused Crawling</i> | 3 |
| 2.1 | Arsitektur Media Monitoring | 15 |
| 2.2 | Arsitektur <i>Web Crawling</i> [18] | 17 |
| 2.3 | Arsitektur <i>Focused Crawling</i> [18] | 18 |
| 2.4 | Perbedaan <i>Web Crawling</i> dan <i>Web Scraping</i> [16] | 18 |
| 2.5 | Confusion Matrix | 23 |
| 3.1 | Tahap Penelitian | 26 |
| 3.2 | Kerangka Kerja Penelitian | 29 |
| 3.3 | Ilustrasi parameter <i>dataset</i> | 30 |
| 3.4 | Proses <i>crawling</i> URL | 32 |
| 3.5 | Proses Pembuatan Topic-Specific Weight Table | 33 |
| 3.6 | Hasil Topic-Specific "Geopark Ciletuh" | 34 |
| 3.7 | Proses Perhitungan <i>Page Relevance</i> untuk menghasilkan skor parameter berdasarkan <i>keyword</i> Geopark Ciletuh | 35 |
| 3.8 | Proses <i>Scraping</i> data dari URL detik dan Kompas | 35 |
| 3.9 | DOM <i>Structure</i> Detik kategori Finance | 36 |
| 3.10 | DOM <i>Structure</i> Detik kategori News | 37 |
| 3.11 | DOM <i>Structure</i> Detik kategori Oto | 37 |
| 3.12 | DOM <i>Structure</i> Detik kategori Travel | 38 |
| 3.13 | DOM <i>Structure</i> Kompas Subdomain Regional | 38 |
| 3.14 | DOM <i>Structure</i> Kompas Subdomain travel | 39 |

| | | |
|------|--|----|
| 3.15 | DOM <i>Structure</i> Kompas Subdomain Pesona Indonesia | 39 |
| 4.1 | Potongan Hasil <i>Scraping Training Data</i> | 42 |
| 4.2 | <i>Topic-Specific</i> Geopark Ciletuh | 47 |
| 4.3 | Contoh <i>Topic-Specific Weight Table</i> dan <i>URL Word Weight Table</i> | 50 |
| 4.4 | Hasil uji coba Model Skenario ke-1 : <i>Dataset</i> Gabungan | 58 |
| 4.5 | Hasil uji coba Model Skenario ke-1 : <i>Dataset</i> Detik dan Kompas | 59 |
| 4.6 | Hasil uji coba Model Skenario ke-2 | 59 |
| 4.7 | Hasil uji coba Model Skenario ke-3 | 60 |
| 4.8 | Hasil uji coba Model Skenario ke-4 | 60 |
| 4.9 | Hasil Uji Coba Breadth First Search (BFS) | 61 |
| 4.10 | Hasil uji coba scraping berita | 62 |
| 4.11 | Hasil uji coba scraping berita tanpa model klasifikasi | 63 |
| 4.12 | Hasil uji coba scraping dengan atau tanpa Multithread | 63 |

Daftar Skrip

| | | |
|-----|--|----|
| 4.1 | BFS pada Detik 1 : Pencarian <i>Keyword</i> (<i>Node</i> Awal) | 43 |
| 4.2 | BFS pada Kompas 1 : Pencarian <i>Keyword</i> (<i>Node</i> Awal) | 44 |
| 4.3 | BFS pada Detik 2 : <i>Crawling</i> URL Layer 1 | 44 |
| 4.4 | BFS pada Kompas 2 : <i>Crawling</i> URL Layer 1 | 45 |
| 4.5 | <i>labeling</i> dataset | 53 |
| 4.6 | Multinomial Naive Bayes : Perhitungan Prior Class | 54 |
| 4.7 | Multinomial Naive Bayes : Perhitungan <i>Conditional Probabilities</i> | 54 |
| 4.8 | Multinomial Naive Bayes : Menentukan <i>Class</i> pada <i>Testing data</i> | 56 |

Daftar Singkatan

| Singkatan | Penjelasan |
|------------------|---|
| PR | <i>Public Relations</i> |
| URL | <i>Uniform Resource Locator</i> |
| UGG | <i>UNESCO Global Geopark</i> |
| UNESCO | <i>United Nations of Educational, Scientific, and Cultural Organization</i> |
| TSWT | <i>Topic-Specific Weight Table / Topic-Specific</i> |
| BFS | <i>Breadth First Search</i> |
| WNRR | <i>Word that Not Represent Receipt</i> |