

TUGAS AKHIR

**OTOMATISASI PELABELAN KARYA TULIS ILMIAH
KOLEKSI PERPUSTAKAAN UNPAR**



Josep Cliff Jonathan

NPM: 2017730046

**PROGRAM STUDI TEKNIK INFORMATIKA
FAKULTAS TEKNOLOGI INFORMASI DAN SAINS
UNIVERSITAS KATOLIK PARAHYANGAN
2024**

FINAL PROJECT

THE AUTOMATION OF SCIENTIFIC PAPER LABELING FOR UNPAR LIBRARY COLLECTION



Josep Cliff Jonathan

NPM: 2017730046

**DEPARTMENT OF INFORMATICS
FACULTY OF INFORMATION TECHNOLOGY AND SCIENCES
PARAHYANGAN CATHOLIC UNIVERSITY
2024**

LEMBAR PENGESAHAN

OTOMATISASI PELABELAN KARYA TULIS ILMIAH KOLEKSI PERPUSTAKAAN UNPAR

Josep Cliff Jonathan

NPM: 2017730046

Bandung, 12 Januari 2024

Menyetujui,

Pembimbing
Digitally signed
by Veronica Sri
Moertini

Dr. Veronica Sri Moertini

Ketua Tim Penguji
Digitally signed
by Gede Karya

Gede Karya, M.T.

Anggota Tim Penguji
Digitally signed
by Cecilia Esti
Nugraheni

Dr.rer.nat. Cecilia Esti Nugraheni

Mengetahui,

Ketua Program Studi

Digitally signed
by Lionov

Lionov, Ph.D.

PERNYATAAN

Dengan ini saya yang bertandatangan di bawah ini menyatakan bahwa tugas akhir dengan judul:

OTOMATISASI PELABELAN KARYA TULIS ILMIAH KOLEKSI PERPUSTAKAAN UNPAR

adalah benar-benar karya saya sendiri, dan saya tidak melakukan penjiplakan atau pengutipan dengan cara-cara yang tidak sesuai dengan etika keilmuan yang berlaku dalam masyarakat keilmuan.

Atas pernyataan ini, saya siap menanggung segala risiko dan sanksi yang dijatuhan kepada saya, apabila di kemudian hari ditemukan adanya pelanggaran terhadap etika keilmuan dalam karya saya, atau jika ada tuntutan formal atau non-formal dari pihak lain berkaitan dengan keaslian karya saya ini.

Dinyatakan di Bandung,
Tanggal 12 Januari 2024



Josep Cliff Jonathan
NPM: 2017730046

ABSTRAK

Penelitian ini mengkaji tentang otomatisasi pelabelan karya ilmiah di perpustakaan Universitas Katolik Parahyangan (UNPAR) di Bandung, Indonesia. Latar belakang penelitian ini adalah tantangan yang dihadapi oleh perpustakaan dalam memberikan label atau tanda khusus pada setiap dokumen, terutama di era teknologi yang berkembang pesat. Tujuan penelitian ini adalah untuk mengotomatisasi pelabelan karya ilmiah berdasarkan jurusan dan topik yang terkait. Data yang digunakan dalam penelitian ini diambil dari situs Repozitori UNPAR dengan menggunakan teknik *web scraping*. Data kemudian diproses dengan berbagai teknik pra-pemrosesan, seperti pembersihan data, tokenisasi, penghapusan stop words, dan lemmatisasi. Untuk pelabelan jurusan, digunakan algoritma klasifikasi *Multinomial Naive Bayes*, sedangkan untuk analisis topik, digunakan algoritma pemodelan topik *Latent Dirichlet Allocation* (LDA). Hasil penelitian menunjukkan bahwa algoritma klasifikasi *Multinomial Naive Bayes* dapat menghasilkan label jurusan yang cukup baik, dan algoritma pemodelan topik LDA dapat mengidentifikasi topik yang berbeda dalam setiap jurusan. . Namun, hasil yang didapatkan tidak begitu spesifik dan rinci karena satu kata kunci dapat memiliki lebih dari satu makna untuk jurusan yang berbeda. Kesimpulan dari penelitian ini adalah bahwa otomatisasi pelabelan karya ilmiah dapat membantu perpustakaan dalam memudahkan perpustakaan dalam pelabelan dan pengguna dalam menemukan topik-topik yang sesuai dengan karya tulis yang dicari.

Kata-kata kunci: Otomatisasi, Perpustakaan, Pemodelan topik, LDA, Naive Bayes Classifier, Klasifikasi

ABSTRACT

This research examines the automation of labeling scientific papers in the library of Parahyangan Catholic University (UNPAR) in Bandung, Indonesia. The background of this research is the challenge faced by the library in providing labels or special marks on each document, especially in the era of rapid technological development. The purpose of this research is to automate the labeling of scientific papers based on the department and related topics. The data used in this research were taken from the UNPAR Repository site using web scraping techniques. The data were then processed with various pre-processing techniques, such as data cleaning, tokenization, stop words removal, and lemmatization. For department labeling, the Multinomial Naive Bayes classification algorithm was used, while for topic analysis, the Latent Dirichlet Allocation (LDA) topic modeling algorithm was used. The results of the research show that the Multinomial Naive Bayes classification algorithm can produce good department labels, and the LDA topic modeling algorithm can identify different topics in each department. However, the results obtained are not very specific and detailed because one keyword can have more than one meaning for different departments. The conclusion of this research is that the automation of labeling scientific papers can help the library in facilitating the library in labeling and users in finding topics that match the papers they are looking for.

Keywords: Automation, Library, Topic modeling, LDA (Latent Dirichlet Allocation), Naive Bayes Classifier, Classification.

untuk saya

KATA PENGANTAR

Puji syukur kehadirat Tuhan Yang Maha Esa, atas berkat dan rahmat-Nya, saya dapat menyelesaikan tugas akhir dengan judul "OTOMATISASI PELABELAN KARYA TULIS ILMIAH KOLEKSI PERPUSTAKAAN UNPAR" sebagai salah satu syarat untuk memperoleh gelar Sarjana Komputer di Universitas Katolik Parahyangan.

Tugas akhir ini merupakan hasil dari pembelajaran dan pengalaman yang saya peroleh selama menempuh pendidikan di Jurusan Informatika Unpar angkatan 2017. Saya menyadari bahwa tugas akhir ini tidak akan terwujud tanpa bantuan dan dukungan dari berbagai pihak. Oleh karena itu, saya ingin menyampaikan ucapan terima kasih yang sebesar-besarnya kepada:

- Bapak, Ibu, dan Cici yang telah memberikan kasih sayang, doa, dan dukungan moral dan materi yang tidak pernah putus selama saya menempuh pendidikan.
- Ibu Maria Veronica, S.T., M.T. sebagai dosen pembimbing utama yang telah memberikan bimbingan, arahan, saran, dan kritik yang sangat berharga dalam penyusunan tugas akhir ini.
- Ibu Dr. Veronica Sri Moertini, Ir., M.T. sebagai dosen pembimbing kedua yang juga telah memberikan bimbingan, arahan, saran, dan kritik yang sangat berharga dalam penyusunan tugas akhir ini.
- Bapak Gede Karya, S.T., M.T. dan Ibu Dr. rer. nat Cecilia Esti Nugraheni, S.T., M.T. sebagai dosen penguji yang telah memberikan masukan dan evaluasi yang sangat berguna untuk penyempurnaan tugas akhir ini.
- Ka Putri, Ka Istoko, Mba Ratna, Katt, Rico, Monic, Sam, Andi, Alma, Jennison, Yani, Nadia, Yusuf, Jason, dan teman-teman lainnya yang telah memberikan semangat, motivasi, dan bantuan dalam berbagai hal selama saya menyelesaikan tugas akhir ini.
- Pihak-pihak lain yang tidak dapat saya sebutkan satu per satu, yang telah membantu saya dalam menyelesaikan tugas akhir ini.

Saya menyadari bahwa tugas akhir ini masih jauh dari sempurna dan masih banyak kekurangan. Oleh karena itu, saya sangat mengharapkan kritik dan saran yang membangun dari pembaca untuk perbaikan di masa depan. Semoga tugas akhir ini dapat bermanfaat bagi saya sendiri dan bagi orang lain yang berkepentingan.

Bandung, Januari 2024

Penulis

DAFTAR ISI

KATA PENGANTAR	xv
DAFTAR ISI	xvii
DAFTAR GAMBAR	xxi
DAFTAR TABEL	xxiii
DAFTAR KODE PROGRAM	xxvii
1 PENDAHULUAN	1
1.1 Latar Belakang	1
1.2 Rumusan Masalah	2
1.3 Tujuan	2
1.4 Batasan Masalah	2
1.5 Metodologi	2
1.6 Sistematika Pembahasan	3
2 LANDASAN TEORI	5
2.1 NLP (<i>Natural Language Processing</i>)	5
2.2 Klasifikasi	10
2.2.1 Algoritma klasifikasi	10
2.3 Tahap Evaluasi Model	14
2.4 Latent Dirichlet Allocation (LDA)	16
2.4.1 Istilah dalam LDA	16
2.4.2 Proses LDA	16
2.4.3 Contoh Kasus LDA	17
3 EKSPLORASI TEKNOLOGI	19
3.1 Eksplorasi NLTK	19
3.2 Eksplorasi <i>CountVectorizer</i>	23
3.3 Eksplorasi <i>TFIDFVectorizer</i>	24
3.4 Eksplorasi <i>Max_features CountVectorizer</i> dan <i>TFIDFVectorizer</i>	25
3.5 Eksplorasi fitur pembeda karya tulis ilmiah	26
3.6 Eksplorasi model-model klasifikasi untuk penebakan jurusan	27
3.7 Eksplorasi LDA <i>Gensim</i>	28
3.8 Eksplorasi LDA Sklearn	31
3.9 Eksplorasi perbandingan LDA <i>Gensim</i> dan <i>Sklearn</i>	33
3.10 Eksplorasi Model Terbaik	34
3.11 Eksplorasi Uji Anova	35
3.12 Eksplorasi Pemilihan fitur dengan KBBI	35
3.13 Eksplorasi Tkinter Designer	37
3.14 Eksplorasi PyQt5	37

4 PENAMBANGAN DATA	39
4.1 Perpustakaan UNPAR	39
4.2 Analisis Masalah	42
4.3 Pengumpulan Data	42
4.4 Deskripsi Data	44
4.5 Fakta Data	51
4.6 Pra-pemrosesan Data	53
4.6.1 Membersihkan data	53
4.6.2 Tokenisasi menggunakan <i>library</i> NLTK	54
4.6.3 Penghapusan kata umum menggunakan <i>library</i> NLTK	54
4.6.4 <i>Lemmatization</i> menggunakan <i>library</i> Sastrawi	55
4.6.5 Penyamarataan nilai <i>class2</i> atau sub kategori jurusan	55
4.6.6 Pembuangan description	55
4.6.7 Penghapusan bahasa yang dideteksi selain dari bahasa indonesia	56
4.6.8 Melakukan pembersihan ulang untuk setiap jurusan	56
4.7 Analisis Data Secara Manual	57
4.8 Pemilihan Fitur	59
4.8.1 Pemilihan fitur dengan menghapus 100 kata umum dengan IDF	59
4.8.2 Pemilihan fitur dengan IDF dan information gain berdasarkan jurusan	60
4.8.3 Pemilihan fitur dengan IDF untuk tiap Jurusan Tahap 1	61
4.8.4 Pemilihan fitur dengan IDF untuk tiap Jurusan Tahap 2	71
4.9 Analisis Data untuk Pemodelan Topik	73
4.9.1 Pemodelan Topik Tahap Satu	73
4.9.2 Pemodelan Topik Tahap Dua	74
4.9.3 Pemodelan Topik Tahap Tiga	75
4.9.4 Pemodelan Topik Tahap Empat	76
4.9.5 Pemodelan Topik Tahap Lima	77
5 PERANGKAT LUNAK UNTUK MELUNCURKAN MODEL KLASIFIKASI	79
5.1 Implementasi Model Topik	79
5.1.1 Hasil Pemodelan Topik untuk Tiap Jurusan	79
5.2 Implementasi Model Klasifikasi untuk Jurusan	118
5.2.1 Data yang digunakan	118
5.2.2 Model Klasifikasi yang digunakan	118
5.2.3 Hasil Metriks Model	119
5.2.4 Model yang digunakan	119
5.3 Implementasi Perangkat Lunak	120
5.3.1 Use Case Scenario	120
5.3.2 Lingkungan Perangkat Lunak	122
5.3.3 Diagram Kelas	122
5.3.4 Antar Muka Pengguna	124
5.4 Pengujian Perangkat Lunak	128
5.4.1 Pengujian Fungsional	128
5.4.2 Pengujian Eksperimental Pertama	129
5.4.3 Pengujian Eksperimental Kedua	132
5.4.4 Pengujian Eksperimental Ketiga	133
5.4.5 Pengujian Eksperimental Keempat	134
6 KESIMPULAN DAN SARAN	135
6.1 Kesimpulan	135
6.2 Saran	135

DAFTAR REFERENSI	137
A KODE PROGRAM	139

DAFTAR GAMBAR

2.1 Grafik harmonic average	15
3.1 Abstrak Mentah	19
3.2 Hasil Keluaran dari word_tokenize()	20
3.3 Hasil Keluaran dari WhitespaceTokenizer().tokenize(s)	20
3.4 Hasil Keluaran dari penghapusan stopword	21
3.5 Hasil Keluaran dari stemming	22
3.6 Hasil Keluaran dari lemmatization	22
3.7 Sebelum dan sesudah pemilihan fitur dengan kbbi	36
4.1 merupakan gambar dari perpustakaan UNPAR di ciumbuleuit berserta kegiatan yang dilakukan	39
4.2 UNPAR-IR atau Repository UNPAR	41
4.3 Metadata tiap dokumen pada repositori UNPAR	43
4.4 Dataset yang digunakan	44
4.5 Book	45
4.6 Book Chapter	45
4.7 Journal Article	45
4.8 Non-Journal Article	46
4.9 Diploma	46
4.10 Undergraduate Theses	46
4.11 Master Theses	47
4.12 Dissertation	47
4.13 Conference Paper	48
4.14 Research Report	48
4.15 Community Service Report	49
4.16 Scientific Oration	49
4.17 Unpublished Lecturer Paper	50
4.18 Unpublished Student Paper	50
4.19 Poster	51
4.20 Fakta pertama	51
4.21 Bahasa yang terdeteksi	52
4.22 Contoh Description	52
4.23 Pembersihan data	53
4.24 Tokenisasi	54
4.25 Stop Word Removal	54
4.26 Stemming	55
4.27 Penyamarataan nilai class 2	55
4.28 Pemisahan abstract dan description	55
4.29 Hasil pembuangan description	56
4.30 Hasil data	56
4.31 Pembersihan ulang untuk setiap jurusan	56
4.32 Grafik abstrak rata-rata jumlah kata dalam tiap jenis dokumen	58

4.33 Histogram yang diberi tanda silang dibuang	62
4.34 Perbandingan distribusi IDF Economics sebelum dan setelah menghapus pencilan	62
4.35 Perbandingan distribusi IDF Theology sebelum dan setelah menghapus pencilan .	63
4.36 Perbandingan distribusi IDF Accounting sebelum dan setelah menghapus pencilan	63
4.37 Perbandingan distribusi IDF Architecture sebelum dan setelah menghapus pencilan	63
4.38 Perbandingan distribusi IDF Business Administration sebelum dan setelah menghapus pencilan	64
4.39 Perbandingan distribusi IDF Chemical Engineering sebelum dan setelah menghapus pencilan	64
4.40 Perbandingan distribusi IDF Civil Engineering sebelum dan setelah menghapus pencilan	65
4.41 Perbandingan distribusi IDF Developmental Economy sebelum dan setelah menghapus pencilan	65
4.42 Perbandingan distribusi IDF Electronics Engineering sebelum dan setelah menghapus pencilan	66
4.43 Perbandingan distribusi IDF Industrial Engineering sebelum dan setelah menghapus pencilan	66
4.44 Perbandingan distribusi IDF Informatics Engineering sebelum dan setelah menghapus pencilan	67
4.45 Perbandingan distribusi IDF International Relations sebelum dan setelah menghapus pencilan	67
4.46 Perbandingan distribusi IDF Law sebelum dan setelah menghapus pencilan	68
4.47 Perbandingan distribusi IDF Management sebelum dan setelah menghapus pencilan	68
4.48 Perbandingan distribusi IDF Mathematics sebelum dan setelah menghapus pencilan	69
4.49 Perbandingan distribusi IDF Philosophy sebelum dan setelah menghapus pencilan	69
4.50 Perbandingan distribusi IDF Physics sebelum dan setelah menghapus pencilan	70
4.51 Perbandingan distribusi IDF Public Administration sebelum dan setelah menghapus pencilan	70
4.52 Perbandingan distribusi IDF Social Science sebelum dan setelah menghapus pencilan	71
4.53 Perbandingan pemilihan fitur sebelumnya dan sekarang (arsir merah berarti dihapus)	71
4.54 Perbandingan distribusi IDF <i>Social Science</i> sebelum dan setelah menghapus pencilan	72
 5.1 Diagram Usecase	120
5.2 Diagram Kelas Perangkat Lunak	122
5.3 Tampilan halaman awal	124
5.4 Tampilan halaman opsi prediksi satu dokumen	124
5.5 Hasil keluaran dari prediksi menggunakan satu dokumen	125
5.6 Tampilan halaman Prediksi Topik dengan PDF	125
5.7 Tampilan halaman Prediksi Topik dan Jurusan dengan PDF	126
5.8 Tampilan halaman Daftar File PDF	126
5.9 Tampilan halaman Proses Prediksi	127
5.10 Tampilan halaman Sudah Selesai	127
5.11 Abstrak dari karya tulis milik 2017730056	131

DAFTAR TABEL

2.1	Tabel contoh <i>data train</i> dan <i>data test</i>	11
2.2	Tabel contoh untuk hitung jumlah kemunculan kata Industri di dalam kelas FTI	12
2.3	Tabel contoh untuk hitung jumlah kemunculan seluruh kata di dalam kelas FTI	12
2.4	Tabel contoh untuk hitung jumlah kata lain di dalam seluruh kelas selain kata Industri	12
2.5	Kemiripan kata yang dihasilkan dengan kelas FTI dan FTIS	13
2.6	Visualisasi dari kondisi model evaluasi pada matriks dua dimensi	14
3.1	Tabel perbandingan model klasifikasi dengan max features = 10	27
3.2	Tabel Perbandingan model klasifikasi dengan max features 50-100	27
4.1	Tabel jumlah <i>record</i> dari tiap jenis dokumen	57
4.2	Tabel abstrak rata-rata jumlah kata dalam tiap jenis dokumen <i>class2</i>	58
4.3	Tabel abstrak rata-rata jumlah kata dalam tiap jenis dokumen <i>class2</i>	59
5.1	Contoh menentukan deskripsi topik	79
5.2	Dataset Topik 0 - Accounting	81
5.3	Dataset Topik 1 - Accounting	81
5.4	Dataset Topik 2 - Accounting	82
5.5	Dataset Topik 3 - Accounting	82
5.6	Dataset Topik 4 - Accounting	83
5.7	Dataset Topik 0 - Architechture	83
5.8	Dataset Topik 1 - Architechture	84
5.9	Dataset Topik 2 - Architechture	84
5.10	Dataset Topik 3 - Architechture	85
5.11	Dataset Topik 4 - Architechture	85
5.12	Dataset Topik 0 - Business Administration	86
5.13	Dataset Topik 1 - Business Administration	86
5.14	Dataset Topik 2 - Business Administration	86
5.15	Dataset Topik 3 - Business Administration	87
5.16	Dataset Topik 4 - Business Administration	87
5.17	Dataset Topik 0 - Chemical Engineering	87
5.18	Dataset Topik 1 - Chemical Engineering	88
5.19	Dataset Topik 2 - Chemical Engineering	88
5.20	Dataset Topik 3 - Chemical Engineering	88
5.21	Dataset Topik 4 - Chemical Engineering	89
5.22	Dataset Topik 0 - Civil Engineering	89
5.23	Dataset Topik 1 - Civil Engineering	89
5.24	Dataset Topik 2 - Civil Engineering	90
5.25	Dataset Topik 3 - Civil Engineering	90
5.26	Dataset Topik 4 - Civil Engineering	90
5.27	Dataset Topik 0 - Developmental Economy	91
5.28	Dataset Topik 1 - Developmental Economy	91
5.29	Dataset Topik 2 - Developmental Economy	91

5.30 Dataset Topik 3 - Developmental Economy	92
5.31 Dataset Topik 4 - Developmental Economy	92
5.32 Dataset Topik 2 - Economics	93
5.33 Dataset Topik 3 - Economics	93
5.34 Dataset Topik 4 - Economics	93
5.35 Dataset Topik 0 - Electronics Engineering	94
5.36 Dataset Topik 1 - Electronics Engineering	94
5.37 Dataset Topik 2 - Electronics Engineering	94
5.38 Dataset Topik 3 - Electronics Engineering	95
5.39 Dataset Topik 4 - Electronics Engineering	95
5.40 Dataset Topik 0 - Industrial Engineering	96
5.41 Dataset Topik 1 - Industrial Engineering	96
5.42 Dataset Topik 2 - Industrial Engineering	96
5.43 Dataset Topik 3 - Industrial Engineering	97
5.44 Dataset Topik 4 - Industrial Engineering	97
5.45 Dataset Topik 0 - International Relations	98
5.46 Dataset Topik 1 - International Relations	98
5.47 Dataset Topik 3 - International Relations	99
5.48 Dataset Topik 4 - International Relations	99
5.49 Dataset Topik 0 - Law	100
5.50 Dataset Topik 1 - Law	100
5.51 Dataset Topik 2 - Law	101
5.52 Dataset Topik 3 - Law	101
5.53 Dataset Topik 4 - Law	102
5.54 Dataset Topik 0 - Management	102
5.55 Dataset Topik 1 - Management	103
5.56 Dataset Topik 2 - Management	103
5.57 Dataset Topik 3 - Management	104
5.58 Dataset Topik 4 - Management	104
5.59 Dataset Topik 0 - Mathematics	105
5.60 Dataset Topik 1 - Mathematics	105
5.61 Dataset Topik 2 - Mathematics	106
5.62 Dataset Topik 3 - Mathematics	106
5.63 Dataset Topik 4 - Mathematics	107
5.64 Dataset Topik 0 - Philosophy	108
5.65 Dataset Topik 1 - Philosophy	108
5.66 Dataset Topik 2 - Philosophy	108
5.67 Dataset Topik 3 - Philosophy	109
5.68 Dataset Topik 4 - Philosophy	109
5.69 Dataset Topik 0 - Physics	110
5.70 Dataset Topik 1 - Physics	110
5.71 Dataset Topik 2 - Physics	111
5.72 Dataset Topik 3 - Physics	111
5.73 Dataset Topik 4 - Physics	111
5.74 Dataset Topik 0 - Public Administration	112
5.75 Dataset Topik 1 - Public Administration	112
5.76 Dataset Topik 2 - Public Administration	113
5.77 Dataset Topik 3 - Public Administration	113
5.78 Dataset Topik 4 - Public Administration	114
5.79 Dataset Topik 0 - Social Science	114
5.80 Dataset Topik 1 - Social Science	114

5.81	Dataset Topik 2 - Social Science	115
5.82	Dataset Topik 3 - Social Science	115
5.83	Dataset Topik 4 - Social Science	115
5.84	Dataset Topik 0 - Theology	116
5.85	Dataset Topik 1 - Theology	116
5.86	Dataset Topik 2 - Theology	117
5.87	Dataset Topik 3 - Theology	117
5.88	Dataset Topik 4 - Theology	118
5.89	Hasil <i>testing</i> dengan prediksi untuk 20 baris teratas	119
5.90	Use Case 1: Prediksi Topik dengan PDF	121
5.91	Use Case 2: Prediksi Topik dan Jurusan dengan PDF	121
5.92	Hasil Pengujian Fungsional	128
5.93	Hasil Output Prediksi Jurusan dan Topik	129
5.94	<i>Confusion Matrix</i> prediksi jurusan dan topik kasus Informatika	129
5.95	Pemeriksaan apakah topik yang sudah diprediksi jurusannya benar dan dihasilkan sesuai dengan input	130
5.96	Hasil Output Prediksi Topik	130
5.97	Hasil Output Prediksi Jurusan dan Topik dengan abstrak Teknik Kimia	132
5.98	<i>Confusion Matrix</i> prediksi jurusan dan topik kasus 10 dokumen Cover - Bab 1 Teknik Kimia	132
5.99	Hasil Output Prediksi Jurusan dan Topik dengan abstrak Teknik Informatika	133
5.100	<i>Confusion Matrix</i> prediksi jurusan dan topik kasus 10 dokumen abstrak Informatika	133
5.101	Hasil Output Prediksi Jurusan dan Topik dengan jurusan yang berbeda-beda	134
5.102	<i>Confusion matrix</i> prediksi jurusan dan topik dengan jurusan yang berbeda-beda	134

DAFTAR KODE PROGRAM

3.1	Contoh kode stop word removal	21
3.2	Contoh kode stemming	21
3.3	Contoh kode lemmatization	22
3.4	Contoh kode CountVectorizer	23
3.5	Keluaran dari CountVectorizer	23
3.6	Contoh kode TFIDFVectorizer	24
3.7	LDA Gensim	28
3.8	Output dari bow_corpus	29
3.9	Output dari model LDA	30
3.10	Contoh kode import sklearn	31
3.11	Contoh kode hitung sparsitas	31
3.12	Contoh kode membuat model LDA scikit-learn	32
3.13	Contoh kode menampilkan output LDA	32
3.14	Contoh kode mengukur kinerja model	32
3.15	Contoh kode membuat gridsearch	33
3.16	Contoh kode KBBI	35
4.1	Contoh kode membuat model LDA scikit-learn	59
A.1	GUI.py	139
A.2	TebakJurusandanTopik.py	151

BAB 1

PENDAHULUAN

1.1 Latar Belakang

Pada era teknologi yang berkembang pesat, data yang dihasilkan juga semakin pesat. Kebanyakan dari data yang dihasilkan adalah teks yang tidak berstruktur dalam bentuk yang beraneka ragam, contohnya yaitu buku, halaman web, artikel, dan dokumen lainnya. Perpustakaan adalah tempat penyimpanan kumpulan buku, majalah, jurnal, dan bahan pustaka lainnya disimpan, diorganisir, dan dipinjamkan kepada pengguna. Perpustakaan biasanya berada di dalam suatu lembaga pendidikan seperti sekolah, universitas, kampus, dan lainnya. Perpustakaan di dalam universitas memiliki koleksi buku dan sumber daya informasi yang luas dan mendalam, dan terdiri dari berbagai macam ilmu yang terdapat di dalam universitas tersebut. Salah satu tantangan besar bagi organisasi seperti perpustakaan adalah untuk memberikan tanda khusus atau label pada setiap dokumen agar mempermudah pelanggan dalam mencari informasi yang sesuai.

Sebagai solusi dari tantangan tersebut penelitian ini bertujuan untuk mengotomatisasi pelabelan koleksi Perpustakaan Universitas Katolik Parahyangan (UNPAR) di Kota Bandung, Jawa Barat, Indonesia. Ratusan ribu koleksi Perpustakaan UNPAR terdiri dari buku, karya mahasiswa, karya dosen, penelitian, dan lainnya. Dalam studi kasus ini dikumpulkan sampel dari beberapa koleksi perpustakaan UNPAR dari berbagai jenis dan topik, dengan label jenis dan topik.

Koleksi pada Perpustakaan UNPAR biasanya memiliki label berdasarkan jenis dan topiknya atau yang biasa disebut katalog. Sebagai contoh, katalog dari buku "Filosofi Teras: Filsafat Yunani-Romawi Kuno untuk Mental Tangguh Masa Kini". Katalog untuk buku tersebut adalah "171.2 MAN f" yang menunjukkan bahwa buku tersebut termasuk dalam kategori filsafat umum dan diterbitkan oleh penerbit Mantra.

Meringkas teks dan mengkategorikan dokumen yang ada di perpustakaan dikenal sebagai katalog. Membuat ringkasan singkat tentang isi dokumen, seperti abstrak atau deskripsi, disebut meringkas teks. Proses ini menghasilkan abstrak yang dibuat secara manual. Pada penelitian ini, dokumen dikategorikan ke dalam jurusan tertentu dan penentuan topik atau *related terms* dari tiap dokumen

Ilmu Informasi merupakan ilmu yang mempelajari pengkodean informasi dengan menggunakan pendekatan matematika. Salah satu aplikasi dari ilmu ini adalah *citation index* yang berfungsi untuk mempelajari dan mengukur teks dan informasi menggunakan sebuah bilangan. *Natural Language Processing* (NLP) merupakan ilmu gabungan dari *Linguistic Science* dan Ilmu Informasi yang bertujuan untuk memodelkan bahasa alami manusia di komputer.

Dataset ini diperoleh dari situs web Repository UNPAR, juga dikenal sebagai UNPAR-IR. Tujuan dari penelitian ini adalah untuk melabeli karya tulis ilmiah dengan jurusan dan kata kunci, sehingga sivitas akademika dapat lebih mudah menemukan informasi yang mereka butuhkan. Hasil akhir dari penelitian ini adalah program yang dapat melakukan prediksi topik dan jurusan dengan input pdf. Program ini dievaluasi dengan pengujian fungsional dan eksperimental.

1.2 Rumusan Masalah

Rumusan masalah yang muncul berdasarkan deskripsi dan latar belakang yang sudah dibahas adalah sebagai berikut:

1. Bagaimana mengumpulkan *dataset* dari Perpustakaan UNPAR?
2. Bagaimana menyiapkan *dataset* yang diperoleh dari Perpustakaan UNPAR?
3. Bagaimana melakukan otomatisasi pelabelan jurusan dan topik untuk karya tulis ilmiah di Perpustakaan UNPAR?
4. Bagaimana cara membangun model klasifikasi jurusan dan model topik menggunakan LDA untuk pelabelan karya tulis ilmiah?
5. Bagaimana mengevaluasi keakuratan dari model klasifikasi dan model LDA yang digunakan untuk melabeli dokumen?
6. Bagaimana cara membangun perangkat lunak untuk meluncurkan model klasifikasi?

1.3 Tujuan

Tujuan dari penelitian ini adalah sebagai berikut:

1. Mengumpulkan data-data yang tersedia di koleksi Perpustakaan UNPAR.
2. Menyiapkan data berupa karya tulis ilmiah yang tersedia di laman Repositori UNPAR.
3. Membuat otomatisasi pelabelan jurusan dan topik untuk karya tulis ilmiah di koleksi Perpustakaan UNPAR.
4. Membangun model klasifikasi jurusan dan model topik menggunakan LDA untuk karya tulis ilmiah yang tersedia di Perpustakaan UNPAR.
5. Melakukan evaluasi keakuratan dari model klasifikasi dan model LDA yang digunakan untuk melabeli dokumen.
6. Membuat perangkat lunak berbasis *desktop* untuk meluncurkan model klasifikasi dan model otomatisasi label.

1.4 Batasan Masalah

Batasan masalah pada penelitian ini adalah:

1. Data yang diambil hanya mencakup repositori UNPAR.
2. Data yang dikumpulkan hanya mencakup tahun publikasi 2017 sampai tahun 2022.
3. Data yang telah terambil memiliki beberapa bahasa dan tipe. Pada penelitian ini, data yang dipakai hanya abstrak berbahasa indonesia.

1.5 Metodologi

Metodologi yang digunakan dalam penyusunan tugas akhir ini terdiri dari langkah-langkah sebagai berikut:

1. Studi literatur terkait pengolahan bahasa alami, klasifikasi teks, dan LDA.
2. Studi literatur mengenai penggunaan NLTK Python.
3. Melakukan pra-pemrosesan text sebagai tahap awal pengolahan *dataset* karya tulis ilmiah yang tersedia di UNPAR-IR.
4. Analisis fitur-fitur yang dapat membedakan jenis-jenis karya tulis ilmiah.
5. Eksperimen pembuatan model klasifikasi untuk pelabelan jenis karya tulis dengan fitur-fitur sesuai hasil analisis.
6. Mengevaluasi hasil pelabelan topik dan sub topik di topik tertentu.
7. Memperbaiki pelabelan dokumen.
8. Membuat perangkat lunak model sub topik di topik tertentu.
9. Menampilkan hasil pelabelan topik.

10. Menyelesaikan perangkat lunak dan melakukan peluncuran.
11. Melakukan analisis dari hasil pengujian dan eksperimen yang telah dilakukan.
12. Menyusun dokumen skripsi.

1.6 Sistematika Pembahasan

1. Bab 1 Pendahuluan

Bagian pertama dari skripsi ini membahas mengenai latar belakang permasalahan, rumusan masalah, tujuan penelitian, batasan masalah, serta metodologi penelitian.

2. Bab 2 Landasan Teori

Bab 2 membahas beberapa konsep dan teknik yang digunakan dalam pengolahan teks. Konsep-konsep tersebut meliputi teknik *Natural Language Processing* (NLP) untuk mengolah dan memproses teks, seleksi fitur untuk mengekstrak fitur dari data teks, evaluasi model teks, penambangan data, klasifikasi, evaluasi klasifikasi, serta penggunaan *Latent Dirichlet Allocation* (LDA) untuk melakukan klasifikasi topik pada teks.

3. Bab 3 Eksplorasi Teknologi

Bab ini mencakup studi kasus kecil, hasil komputasi atau perhitungan manual yang dilakukan untuk mempelajari algoritma yang digunakan. Hasil eksplorasi penggunaan teknologi juga dicantumkan di sini, termasuk hasil pengujian menggunakan *library* yang diperlukan di lingkungan pemrograman Python.

4. Bab 4 Penambangan Data

Bab ini mencakup deskripsi masalah yang ingin diselesaikan, cara pengumpulan data dan langkah pra-pemrosesan data untuk tujuan klasifikasi, serta hasil yang diperoleh dari proses penambangan data.

5. Bab 5 Perangkat Lunak untuk Meluncurkan Model Klasifikasi

Bab ini berisi implementasi model, pengujian perangkat lunak, dan eksperimen berdasarkan rancangan yang telah dibuat.

6. Bab 6 Kesimpulan dan Saran

Bab ini berisi kesimpulan dan saran dari seluruh penelitian yang telah dilakukan.

