

Pengaruh *N-Gram* pada Klasifikasi Dokumen menggunakan  
*Algoritma Naïve Bayes Classifier*

*Diajukan Sebagai Syarat Untuk Menyelesaikan  
Pendidikan Program Strata-1 Pada  
Jurusan Teknik Informatika*



Oleh :

Fitria Khoirunnisa  
09021181520130

**JURUSAN TEKNIK INFORMATIKA  
FAKULTAS ILMU KOMPUTER UNIVERSITAS SRIWIJAYA  
2019**

## LEMBAR PENGESAHAN SKRIPSI

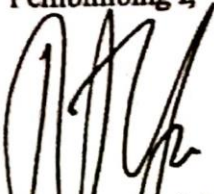
### PENGARUH *N-GRAM* PADA KLASIFIKASI DOKUMEN MENGUNAKAN ALGORITMA *NAÏVE BAYES CLASSIFIER*

Oleh :

**FITRIA KHOIRUNNISA**  
NIM : 09021181520130

Palembang, Desember 2019

Pembimbing I,



**Nowi Yuskani, M.T.**  
NIP. 198211082012122001

Pembimbing II,



**Desty Rodiah, M.T.**  
NIK. 1671016112890005

Mengetahui,  
Ketua Jurusan Teknik Informatika



**Rifkie Primartha, MT**  
NIP. 197706012009121004

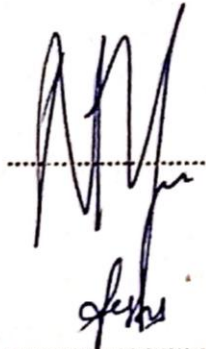
## TANDA LULUS UJIAN SIDANG SKRIPSI

Pada hari Kamis tanggal 26 Desember 2019 telah dilaksanakan ujian sidang skripsi oleh Jurusan Teknik Informatika Fakultas Ilmu Komputer Universitas Sriwijaya

Nama : Fitria Khoirunnisa  
NIM : 09021181520130  
Judul : Pengaruh *N-Gram* pada Klasifikasi Dokumen Menggunakan Algoritma *Naïve Bayes Classifier*

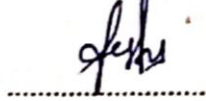
### 1. Pembimbing I

Novi Yusliani, M.T.  
NIP. 198211082012122001



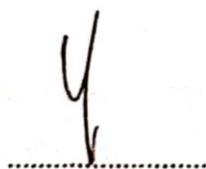
### 2. Pembimbing II

Desty Rodiah, M.T.  
NIK. 1671016112890005



### 3. Penguji I

Yunita, M.Cs  
NIP. 198306062015042002



### 4. Penguji II

Kanda Januar Miraswan, M.T.  
NIP. 199001092019031012



Mengetahui,  
Ketua Jurusan Teknik Informatika



Rifkie Primartha, M.T.  
NIP. 19770601200912004

## HALAMAN PERNYATAAN

Yang bertanda tangan dibawah ini :

Nama : Fitria Khoirunnisa  
NIM : 09021181520130  
Program Studi : Teknik Informatika  
Judul Skripsi : Pengaruh *N-Gram* Pada Klasifikasi Dokumen  
Menggunakan Algoritma *Naïve Bayes Classifier*  
Hasil Pengecekan Software *iThenticate/Turnitin* : 14%

Menyatakan bahwa Laporan Proyek saya merupakan hasil karya sendiri dan bukan hasil penjiplakan/plagiat. Apabila ditemukan unsur penjiplakan/plagiat dalam laporan proyek ini, maka saya bersedia menerima sanksi akademik dari Universitas Sriwijaya sesuai dengan ketentuan yang berlaku.

Demikian pernyataan ini saya buat dengan sebenarnya dan tidak ada paksaan oleh siapapun.



Palembang, Desember 2019



Fitria Khoirunnisa  
NIM. 09021181520130

## MOTTO DAN PERSEMBAHAN

Motto :

- “Hai orang-orang yang beriman, jadikanlah sabar dan shalat sebagai penolongmu, sesungguhnya Allah beserta orang-orang yang sabar”  
**(QS. Al-Baqarah : 153)**
  
- “Sesungguhnya bersama kesulitan ada kemudahan. Maka apabila engkau telah selesai (dari sesuatu urusan), tetaplah bekerja keras (untuk urusan yang lain). Dan hanya kepada Tuhanmulah engkau berharap”  
**(QS : Al-Insyirah : 6-8)**

**Kupersembahkan karya tulis ini kepada :**

- Allah Subhanahuwa ta'ala
- Kedua orang tua tercinta
- Ketiga adik tercinta
- Keluarga besarku
- Sahabat dan teman seperjuangan
- Seluruh dosen dan staff pengajar
- Almamater Universitas Sriwijaya

# **THE EFFECT OF N-GRAM ON DOCUMENT CLASSIFICATION USING NAÏVE BAYES CLASSIFIER ALGORITHM**

**By :**  
**Fitria Khoirunnisa**  
**09021181520130**

## **ABSTRACT**

News has become a major need for everyone, with news we can get the information needed. News can be distributed in the form of print mass media, electronic mass media and online media. The means of spreading the news now has grown very rapidly, making the amount of information being managed more and more and word management classified as not small. Therefore, we need a system in the classification of documents that are not structured. In this study, word processing in a document is done by N-Gram as a feature generation. The document classification process is carried out using the Naïve Bayes Classifier algorithm. This study examines the effect of N-Gram on document classification using the Naïve Bayes Classifier algorithm. The results of the classification accuracy of documents by applying N-Gram of 32.68% and without applying N-Gram of 84.97%. A decrease in the classification results occurs the number of features that result from solving N-Gram that is unique or dominant to another category. The accuracy of the results obtained shows that the application of N-Gram in the classification of documents using the Naïve Bayes Classifier algorithm gives a decreased effect on the performance of the classification.

**Key Word :** N-Gram, Naïve Bayes Classifier, Text Mining

# **PENGARUH *N-GRAM* PADA KLASIFIKASI DOKUMEN MENGUNAKAN ALGORITMA *NAÏVE BAYES CLASSIFIER***

**Oleh :  
Fitria Khoirunnisa  
09021181520130**

## **ABSTRAK**

Berita sudah menjadi kebutuhan utama bagi setiap orang, dengan berita kita dapat memperoleh informasi yang dibutuhkan. Berita dapat disebar dalam bentuk media massa cetak, media massa elektronik dan media online. Sarana penyebaran berita sekarang ini telah berkembang sangat pesat, membuat jumlah informasi yang dikelola semakin banyak dan pengolahan kata tergolong tidak sedikit. Oleh karena itu, dibutuhkan suatu sistem dalam klasifikasi dokumen yang tidak terstruktur. Pada penelitian ini, pemrosesan kata dalam suatu dokumen dilakukan dengan *N-Gram* sebagai pembangkitan fitur. Proses klasifikasi dokumen dilakukan dengan menggunakan algoritma *Naïve Bayes Classifier*. Penelitian ini menguji pengaruh *N-Gram* pada klasifikasi dokumen menggunakan algoritma *Naïve Bayes Classifier*. Hasil akurasi klasifikasi dokumen dengan menerapkan *N-Gram* sebesar 32,68% dan tanpa menerapkan *N-Gram* sebesar 84,97%. Penurunan hasil klasifikasi terjadi banyaknya fitur yang dihasilkan dari pemecahan *N-Gram* yang unik atau dominan ke kategori lain. Dari hasil akurasi yang didapatkan menunjukkan bahwa penerapan *N-Gram* pada klasifikasi dokumen menggunakan algoritma *Naïve Bayes Classifier* memberikan pengaruh yang menurun terhadap kinerja klasifikasi.

**Kata Kunci :** *N-Gram, Naïve Bayes Classifier, Text Mining*

## KATA PENGANTAR

Puji dan syukur atas kehadiran Allah ta'ala atas segala nikmat, rahmat, dan karunia-Nya sehingga penulis dapat menyelesaikan Tugas Akhir ini. Tugas Akhir yang berjudul **“Pengaruh *N-Gram* pada Klasifikasi Dokumen menggunakan Algoritma *Naïve Bayes Classifier*”** disusun untuk memenuhi salah satu persyaratan kelulusan tingkat sarjana (S1) pada Jurusan Teknik Informatika Fakultas Ilmu Komputer Universitas Sriwijaya.

Pada kesempatan ini, penulis mengucapkan banyak terima kasih kepada pihak-pihak yang telah memberikan bantuan dan dukungan baik secara langsung maupun secara tidak langsung dalam penyelesaian Tugas Akhir ini, yaitu :

1. Bapak Joharudin dan ibu Sobatini selaku orang tua penulis, Andre Cevtiansyah, Iffat Vajriansyah, dan Siti Ai'syah selaku saudara penulis serta keluarga besar yang telah memberikan dukungan baik moril maupun materil kepada penulis.
2. Bapak Jaidan Jauhari, M.T. selaku Dekan Fakultas Ilmu Komputer Universitas Sriwijaya, Bapak Rifkie Primartha, M.T. selaku Ketua Jurusan Teknik Informatika Universitas Sriwijaya dan Ibu Hardini Novianti, M.T. selaku Sekretaris Jurusan Teknik Informatika Universitas Sriwijaya.
3. Bapak Hadipurnawan Satria, Ph.D dan Ibu Rifka Widyastuti S.Kom., M.Ti, M.Im., selaku pembimbing akademik, yang telah membimbing, mengarahkan dan memberikan motivasi penulis dalam proses perkuliahan dan pengerjaan Tugas Akhir.
4. Ibu Novi Yusliani, M.T. selaku pembimbing I Tugas Akhir dan ibu Desty Rodiah, M.T. selaku pembimbing II Tugas Akhir yang telah memberikan pengarahan, bimbingan, bantuan, serta masukan kepada penulis sehingga Tugas Akhir ini dapat diselesaikan dengan baik.
5. Ibu Yunita, M.Cs. selaku penguji I Tugas Akhir dan Bapak Kanda Januar Miraswan, M.T. selaku penguji II Tugas Akhir yang telah memberikan saran dan kritik serta arahan selama penulisan Tugas Akhir.



6. Seluruh bapak dan ibu Dosen Jurusan Teknik Informatika Fakultas Ilmu Komputer Universitas Sriwijaya yang telah memberikan ilmu pengetahuan kepada penulis selama perkuliahan.
7. Seluruh karyawan dan karyawanati Fakultas Ilmu Komputer Universitas Sriwijaya yang telah membantu penulis dalam urusan administrasi selama kegiatan akademik.
8. Mbak Mida dan Lupina sahabat penulis dari awal perkuliahan yang selalu ada sebagai tempat berbagi cerita, curhat, memberi semangat, do'a dan dukungan kepada penulis.
9. Keluarga Besar Asrama Palembang dan Rumah Kosan43 Bapak Feri, Ibu Linda, Septi, Sandira, Alma, Okta, Yulik, Teteh, Amel, Indah, Mak Tika, Marina, dan Dolor-Dolor lainnya yang selalu memberikan asupan semangat dan doa kepada penulis.
10. Teman seperjuangan penulis seluruh mahasiswa/mahasiswa BEM, DPM, FASCO, group bimbingan TA dan Teknik Informatika Reguler/Bilingual terutama Reguler Angkatan 2015.
11. Seluruh pihak yang telah membantu dalam penyusunan dan penyempurnaan tugas akhir ini yang tidak dapat disebutkan satu persatu.

Penulis menyadari bahwa Tugas Akhir ini jauh dari kata sempurna. Untuk itu, penulis mengharapkan kritik dan sarab yang membangun untun kesempurnaan Tugas Akhir ini dimasa mendatang, harapannya semoga Tugas Akhir ini dapat bermanfaat untuk semua.

Palembang, Desember 2019

Penyusun,

## DAFTAR ISI

	Halaman
HALAMAN JUDUL.....	i
HALAMAN PENGESAHAN SKRIPSI.....	ii
HALAMAN TANDA LULUS UJIAN SIDANG SKRIPSI .....	iii
HALAMAN PERNYATAAN .....	iv
HALAMAN MOTTO DAN PERSEMBAHAN.....	v
ABSTRACT.....	vi
ABSTRAK .....	vii
KATA PENGANTAR .....	viii
DAFTAR ISI.....	x
DAFTAR TABEL.....	xiii
DAFTAR GAMBAR .....	xv
DAFTAR LAMPIRAN.....	xvii
<b>BAB I PENDAHULUAN</b>	
1.1 Pendahuluan .....	I-1
1.2 Latar Belakang .....	I-1
1.3 Rumusan Masalah .....	I-4
1.4 Tujuan Penelitian.....	I-4
1.5 Manfaat Penelitian.....	I-5
1.6 Batasan Masalah.....	I-5
1.7 Sistematika Penulisan.....	I-6
1.8 Kesimpulan.....	I-7
<b>BAB II KAJIAN LITERATUR</b>	
2.1 Pendahuluan .....	II-1
2.2 <i>Text Mining</i> .....	II-1
2.3 <i>Text Preprocessing</i> .....	II-2
2.4 <i>Naïve Bayes Classifier</i> .....	II-5
2.5 <i>N-Gram</i> .....	II-6
2.6 <i>Term Frequency</i> .....	II-9
2.7 Berita .....	II-9
2.8 <i>K-Fold Cross Validation</i> .....	II-10
2.9 <i>Rational Unified Process</i> .....	II-12
2.10 Penelitian Lain Yang Relevan.....	II-13
2.11 Kesimpulan.....	II-16
<b>BAB III METODOLOGI PENELITIAN</b>	
3.1 Pendahuluan .....	III-1
3.2 Pengumpulan Data .....	III-1
3.2.1 Jenis dan Sumber Data.....	III-1

3.2.2 Metode Pengumpulan Data.....	III-1
3.3 Tahap Penelitian .....	III-2
3.3.1 Menetapkan Kerangka Kerja .....	III-2
3.3.2 Menetapkan Kriteria Pengujian .....	III-5
3.3.3 Menetapkan Format Data Pengujian .....	III-6
3.3.4 Menentukan Alat yang Digunakan dalam Pelaksanaan Penelitian .....	III-7
3.3.5 Melakukan Pengujian Penelitian .....	III-7
3.3.6 Menentukan Analisis Hasil Pengujian dan Membuat Kesimpulan Penelitian .....	III-8
3.4 Metode Pengembangan Perangkat Lunak.....	III-10
3.4.1 Fase Insepsi.....	III-10
3.4.2 Fase Elaborasi .....	III-11
3.4.3 Fase Konstruksi .....	III-11
3.4.4 Fase Transisi .....	III-12
3.5 Manajemen Proyek Penelitian .....	III-12

#### BAB IV PENGEMBANGAN PERANGKAT LUNAK

4.1 Pendahuluan .....	IV-1
4.2 Fase Insepsi .....	IV-1
4.2.1 Pemodelan Bisnis.....	IV-1
4.2.2 Kebutuhan Sistem .....	IV-2
4.2.3 Analisis dan Desain .....	IV-3
4.2.3.1 Analisis Kebutuhan Perangkat Lunak.....	IV-3
4.2.3.2 Analisis Data .....	IV-4
4.2.3.3 Analisis <i>Text Preprocessing</i> .....	IV-5
4.2.3.4 Analisis <i>N-Gram</i> .....	IV-8
4.2.3.5 Analisis <i>Term Frequency (TF)</i> .....	IV-10
4.2.3.6 Analisis <i>Naïve Bayes Classifier</i> .....	IV-15
4.2.3.7 Analisis <i>K-Fold Cross Validation</i> .....	IV-24
4.2.3.8 Analisis Persentase Akurasi.....	IV-25
4.2.3.9 Desain Perangkat Lunak .....	IV-27
4.3 Fase Elaborasi .....	IV-37
4.3.1 Pemodelan Bisnis.....	IV-37
4.3.2 Perancangan Data .....	IV-38
4.3.3 Perancangan Antarmuka.....	IV-38
4.3.4 Kebutuhan Sistem.....	IV-39
4.3.5 Diagram Aktivitas.....	IV-40
4.3.6 Diagram <i>Sequence</i> .....	IV-45
4.4 Fase Konstruksi .....	IV-47
4.4.1 Kebutuhan Sistem.....	IV-47
4.4.2 Diagram Kelas .....	IV-47
4.4.3 Implementasi.....	IV-49
4.4.3.1 Implementasi Kelas.....	IV-49
4.4.3.2 Implementasi Antarmuka.....	IV-51
4.5 Fase Transisi.....	IV-52

4.5.1	Pemodelan Bisnis.....	IV-52
4.5.2	Kebutuhan Sistem.....	IV-52
4.5.3	Rencana Pengujian.....	IV-53
4.5.3.1	Rencana Pengujian <i>Use Case</i> Muat Data .....	IV-53
4.5.3.2	Rencana Pengujian <i>Use Case</i> Mempraproses Data .....	IV-54
4.5.3.3	Rencana Pengujian <i>Use Case</i> Melatih Data.....	IV-54
4.5.3.4	Rencana Pengujian <i>Use Case</i> Menguji Data .....	IV-55
4.5.3.5	Rencana Pengujian <i>Use Case</i> Melatih Data <i>N-Gram</i> .....	IV-55
4.5.3.6	Rencana Pengujian <i>Use Case</i> Menguji Data <i>N-Gram</i> .....	IV-56
4.5.3.7	Rencana Pengujian <i>Use Case</i> Menghitung Akurasi .....	IV-56
4.5.3.8	Rencana Pengujian <i>Use Case</i> Menghitung Akurasi <i>N-Gram</i> .....	IV-56
4.5.4	Implementasi.....	IV-57
4.5.4.1	Pengujian <i>Use Case</i> Muat Data .....	IV-58
4.5.4.2	Pengujian <i>Use Case</i> Mempraproses Data .....	IV-59
4.5.4.3	Pengujian <i>Use Case</i> Melatih Data .....	IV-60
4.5.4.4	Pengujian <i>Use Case</i> Menguji Data .....	IV-61
4.5.4.5	Pengujian <i>Use Case</i> Melatih Data <i>N-Gram</i> .....	IV-62
4.5.4.6	Pengujian <i>Use Case</i> Menguji Data <i>N-Gram</i> .....	IV-64
4.5.4.7	Pengujian <i>Use Case</i> Menghitung Akurasi .....	IV-64
4.5.4.8	Pengujian <i>Use Case</i> Menghitung Akurasi <i>N-Gram</i> .....	IV-64
4.6	Kesimpulan.....	IV-66

65

## BAB V HASIL DAN ANALISIS PENELITIAN

5.1	Pendahuluan .....	V-1
5.2	Data Hasil Penelitian .....	V-1
5.2.1	Konfigurasi Percobaan.....	V-1
5.2.2	Data Hasil Konfigurasi I.....	V-2
5.2.3	Data Hasil Konfigurasi II.....	V-6
5.3	Analisis Hasil Penelitian .....	V-6
5.4	Kesimpulan .....	V-8

## BAB VI KESIMPULAN DAN SARAN

6.1	Pendahuluan .....	VI-1
6.2	Kesimpulan.....	VI-1
6.3	Saran.....	VI-2

## DAFTAR PUSTAKA

## DAFTAR TABEL

	Halaman
II-1. Tabel Pembentukan <i>Bi-Gram</i> (Chandra et al., 2016) .....	II-8
II-2. Tabel Skema 10 <i>Fold Cross Validation</i> .....	II-11
III-1. Tabel Rancangan Evaluasi Sistem Hasil Klasifikasi Dokumen .....	III-7
III-2. Tabel Rancangan Hasil Pengujian pada Percobaan K-Fold.....	III-8
III-3. Tabel Rancangan Hasil Analisis Klasifikasi Dokumen .....	III-9
III-4. Tabel Rancangan Perbandingan Hasil Klasifikasi Dokumen Terhadap Penelitian Sebelumnya .....	III-10
III-5. Tabel Manajemen Proyek Penelitian dalam Bentuk <i>Work Breakdown Structure</i> (WBS).....	III-13
IV-1. Tabel Kebutuhan Fungsional .....	IV-3
IV-2. Tabel Kebutuhan Non-Fungsional .....	IV-3
IV-3. Tabel Contoh Data Berita Online .....	IV-5
IV-4. Tabel Hasil Tahapan Case Folding .....	IV-6
IV-5. Tabel Hasil Tahapan Tokenizing .....	IV-7
IV-6. Tabel Hasil Tahapan Filtering .....	IV-7
IV-7. Tabel Hasil Tahapan Stemming .....	IV-8
IV-8. Tabel Pemodelan N-Gram.....	IV-9
IV-9. Tabel Hasil Perhitungan <i>Term Frekuensi</i> (TF) dari Data Berita .....	IV-11
IV-10. Tabel Hasil Perhitungan <i>Term Frekuensi</i> (TF) dari Pemodelan N-Gram .....	IV-12
IV-11. Tabel Perhitungan Probabilitas <i>Likelihood</i> dari Data Berita .....	IV-16
IV-12. Tabel Perhitungan Probabilitas <i>Likelihood</i> dari Pemodelan <i>N-Gram</i> .....	IV-18
IV-13. Tabel Proses <i>K-Fold Cross Validation</i> .....	IV-25
IV-14. Tabel Contoh Hasil Klasifikasi .....	IV-26
IV-15. Tabel Definisi Aktor .....	IV-28
IV-16. Tabel Definisi Diagram <i>Use Case</i> .....	IV-28
IV-17. Tabel Definisi Muat Data .....	IV-29
IV-18. Tabel Definisi Memproses Data .....	IV-30
IV-19. Tabel Definisi Proses Melatih Data .....	IV-31
IV-20. Tabel Definisi Proses Menguji Data .....	IV-32
IV-21. Tabel Definisi Melatih Data <i>N-Gram</i> .....	IV-33
IV-22. Tabel Definisi Menguji Data <i>N-Gram</i> .....	IV-34
IV-23. Tabel Definisi Menghitung Akurasi.....	IV-35
IV-24. Tabel Definisi Menghitung Akurasi <i>N-Gram</i> .....	IV-36
IV-25. Tabel Implementasi Kelas .....	IV-49
IV-26. Tabel Rencana Pengujian <i>Use Case</i> Muat Data .....	IV-53
IV-27. Tabel Rencana Pengujian <i>Use Case</i> Memproses Data .....	IV-54
IV-28. Tabel Rencana Pengujian <i>Use Case</i> Melatih Data.....	IV-54
IV-29. Tabel Rencana Pengujian <i>Use Case</i> Menguji Data.....	IV-55
IV-30. Tabel Rencana Pengujian <i>Use Case</i> Melatih Data <i>N-Gram</i> .....	IV-55

IV-31.	Tabel Rencana Pengujian <i>Use Case</i> Menguji Data <i>N-Gram</i> .....	IV-56
IV-32.	Tabel Rencana Pengujian <i>Use Case</i> Menghitung Akurasi .....	IV-56
IV-33.	Tabel Rencana Pengujian <i>Use Case</i> Menghitung Akurasi <i>N-Gram</i> ...	IV-56
IV-34.	Tabel Pengujian <i>Use Case Muat Data</i> .....	IV-58
IV-35.	Tabel Pengujian <i>Use Case</i> Mempraproses Data .....	IV-59
IV-36.	Tabel Pengujian <i>Use Case</i> Melatih Data .....	IV-60
IV-37.	Tabel Pengujian <i>Use Case</i> Menguji Data .....	IV-61
IV-38.	Tabel Pengujian <i>Use Case</i> Melatih Data <i>N-Gram</i> .....	IV-62
IV-39.	Tabel Pengujian <i>Use Case</i> Menguji Data <i>N-Gram</i> .....	IV-64
IV-40.	Tabel Pengujian <i>Use Case</i> Menghitung Akurasi .....	IV-64
IV-41.	Tabel Pengujian <i>Use Case</i> Menghitung Akurasi <i>N-Gram</i> .....	IV-64
V-1.	Tabel Hasil Pengujian Sistem Klasifikasi Dokumen .....	V-2
V-2.	Tabel Evaluasi pada Percobaan K-Fold Cross Validation .....	V-4
V-3.	Tabel Perbandingan Hasil Pengujian Klasifikasi Dokumen .....	V-6

## DAFTAR GAMBAR

	Halaman
II-1. Gambar Contoh tahap <i>case folding</i> (Mooney,2006) .....	II-3
II-2. Gambar Contoh tahap <i>tokenizing</i> (Mooney,2006) .....	II-3
II-3. Gambar Contoh tahap <i>filtering</i> (Mooney,2006) .....	II-4
II-4. Gambar Contoh tahap <i>stemming</i> (Mooney,2006) .....	II-4
II-5. Gambar Struktur berita (Musthafa, 2009) .....	II-10
II-6. Gambar Arsitektur <i>Rational Unified Process</i> (Kruchten, 2003) .....	II-12
III-1. Gambar Diagram Tahapan Proses Perangkat Lunak dengan <i>N-Gram</i>	III-3
III-2. Gambar Diagram Tahapan Proses Perangkat Lunak tanpa <i>N-Gram</i> ..	III-3
III-3. Gambar Penjadwalan untuk Tahap Menentukan Ruang Lingkup dan Unit Penelitian .....	III-19
III-4. Gambar Penjadwalan untuk Tahap Menentukan Dasar Teori yang Berkaitan dengan Penelitian .....	III-20
III-5. Gambar Penjadwalan untuk Tahap Menentukan Kriteria Pengujian .	III-20
III-6. Gambar Penjadwalan untuk Tahap Menentukan Alat yang Digunakan untuk Pelaksanaan Penelitian Fase Insepsi .....	III-21
III-7. Gambar Penjadwalan untuk Tahap Menentukan Alat yang Digunakan untuk Pelaksanaan Penelitian Fase Elaborasi .....	III-21
III-8. Gambar Penjadwalan untuk Tahap Menentukan Alat yang Digunakan untuk Pelaksanaan Penelitian Fase Konstruksi .....	III-22
III-9. Gambar Penjadwalan untuk Tahap Menentukan Alat yang Digunakan untuk Pelaksanaan Penelitian Fase Transisi .....	III-22
III-10. Gambar Penjadwalan untuk Tahap Melakukan Pengujian Penelitian	III-23
III-11. Gambar Penjadwalan untuk Tahap Analisa Hasil Pengujian Penelitian dan Membuat Kesimpulan .....	III-23
IV-1. Gambar Perhitungan Probabilitas <i>Priori</i> .....	IV-16
IV-2. Gambar Perhitungan Probabilitas <i>Posterior</i> dari Data Berita .....	IV-23
IV-3. Gambar Perhitungan Probabilitas <i>Posterior</i> dari Pemodelan <i>N-Gram</i>	IV-23
IV-4. Gambar Klasifikasi Data <i>Testing</i> dari Data Berita .....	IV-23
IV-5. Gambar Klasifikasi Data <i>Testing</i> dari Pemodelan <i>N-Gram</i> .....	IV-24
IV-6. Gambar Contoh Perhitungan Persentase Akurasi .....	IV-26
IV-7. Gambar Diagram <i>Use Case</i> .....	IV-27
IV-8. Gambar Rancangan Antarmuka Perangkat Lunak .....	IV-39
IV-9. Gambar Diagram Aktivitas Mempraproses Data .....	IV-41
IV-10. Gambar Diagram Aktivitas Melatih Data .....	IV-42
IV-11. Gambar Diagram Aktivitas Menguji Data .....	IV-42
IV-12. Gambar Diagram Aktivitas Mempraproses Data <i>N-Gram</i> .....	IV-43
IV-13. Gambar Diagram Aktivitas Melatih Data <i>N-Gram</i> .....	IV-44
IV-14. Gambar Diagram Aktivitas Menguji Data <i>N-Gram</i> .....	IV-44
IV-15. Gambar Diagram Aktivitas Menghitung Akurasi .....	IV-45
IV-16. Gambar Diagram <i>Sequence</i> Perangkat Lunak .....	IV-46

IV-17.	Gambar Diagram Kelas Perangkat Lunak.....	IV-48
IV-18.	Gambar Antarmuka Perangkat Lunak.....	IV-52
V-1.	Gambar Perhitungan Persentase Akurasi Hasil Klasifikasi .....	V-4
V-2.	Gambar Perhitungan Persentase Rata-Rata Akurasi Hasil Klasifikasi .....	V-4
V-3.	Gambar Grafik Hasil Pengujian 10 <i>Fold</i> pada Klasifikasi Dokumen dengan Algoritma <i>Naïve Bayes Classifier</i> tanpa Menggunakan Pemodelan <i>N-Gram</i> .....	V-5
V-4.	Gambar Grafik Hasil Pengujian 10 <i>Fold</i> pada Klasifikasi Dokumen dengan Algoritma <i>Naïve Bayes Classifier</i> dengan Menggunakan Pemodelan <i>N-Gram</i> .....	V-5
V-5.	Gambar Perbandingan Hasil Pengujian Klasifikasi Dokumen.....	V-7



## DAFTAR LAMPIRAN

Halaman

1. Hasil Perhitungan *Term Frequency N-Gram*
2. Hasil Perhitungan *Posterior*
3. *Source Code* Program

# **BAB I**

## **PENDAHULUAN**

### **1.1 Pendahuluan**

Bab ini membahas latar belakang masalah, rumusan masalah, tujuan penelitian, manfaat penelitian dan batasan masalah yang menjadi dasar dalam penelitian ini serta memberikan penjelasan umum mengenai keseluruhan penelitian. Pendahuluan dimulai dengan penjelasan singkat tentang dokumen yang diteliti. Selanjutnya, bab ini akan membahas kebutuhan dan tantangan dalam melakukan klasifikasi dokumen berdasarkan kategori yang menjadi latar belakang masalah penelitian ini, serta pembangkit fitur sebagai solusinya.

### **1.2 Latar Belakang Masalah**

Berita sudah menjadi kebutuhan utama bagi setiap orang. Hal ini dikarenakan berita dapat menghasilkan informasi yang dibutuhkan setiap orang. Dengan adanya berita, suatu informasi tentang fakta atau sedang terjadi dapat diketahui. Berita adalah bentuk laporan tentang suatu kejadian yang sedang terjadi baru-baru ini atau keterangan terbaru dari suatu peristiwa (Juditha, 2013). Berita dapat disebar dalam bentuk media massa cetak seperti surat kabar, *tabloid*, *newsletter*, majalah, *bulletin*, dan buku. Media massa elektronik seperti radio, televisi, dan film. Maupun media online yakni website internet yang berisikan informasi aktual layaknya media cetak.

Sarana penyebaran berita sekarang ini telah berkembang sangat pesat, membuat jumlah informasi yang disebar semakin meningkat. Maka jumlah informasi yang dikelola semakin banyak dan akan membuat pengolahan kata tergolong tidak sedikit. Sehingga, pengolahan kata yang dibutuhkan dalam mengategorikan berita perlu mengetahui inti dari isi dokumen. Proses yang dipakai dalam mengategorikan dokumen yang tidak terstruktur adalah proses klasifikasi dokumen dengan algoritma *Naïve Bayes Classifier* (Indriani, 2014).

Algoritma *Naïve Bayes Classifier* adalah algoritma yang digunakan dalam penelitian *text mining*. Algoritma *Naïve Bayes Classifier* merupakan metode klasifikasi berdasarkan probabilitas dan *Teorema Bayesian*. Tahapan klasifikasi ditentukan berdasarkan nilai kategori dari suatu dokumen yakni *term* yang muncul dalam dokumen yang diklasifikasi (Hamzah, 2012).

Salah satu kelebihan algoritma *Naïve Bayes Classifier* dibandingkan algoritma lain menurut Setiawan dan Nursantika (2017) adalah untuk mengetahui kategori pada sebuah dokumen dengan pengukuran kemiripan dokumen terkait yakni melalui proses pengenalan teks dan dokumen. Algoritma *Naïve Bayes Classifier* merupakan klasifikasi statistik yang bisa memprediksi probabilitas sebuah kelas dan tingkat akurasi yang tinggi. Menurut Narayanan, Arora, dan Bhatia (2013) mengatakan algoritma *Naïve Bayes Classifier* dapat meningkatkan hasil akurasi klasifikasi dan menghilangkan *noise* dengan cara memilih jenis fitur yang tepat dan sesuai. Salah satu model yang tepat dalam pengolahan data berupa teks di *data mining* dan pemrosesan kata (Kumar, Sc, Phil, & Ph, 2018) untuk membangkitkan fitur adalah *N-Gram* (Sugianto, Liliana, dan Rostianingsih, 2013).

*N-Gram* adalah model probabilistik yang dapat digunakan untuk membangkitkan karakter dan kata serta memprediksi kata berikutnya dalam urutan kata tertentu. Dalam pembangkitan karakter, *N-Gram* terdiri dari *substring* sepanjang  $n$  karakter dari sebuah string. *N-Gram* digunakan untuk mengambil potongan-potongan karakter huruf sejumlah  $n$  dari sebuah kata yang secara kontinuitas dibaca dari teks sumber hingga akhir dari dokumen (Sugianto et al., 2013). Semakin besar nilai  $n$  dari sebuah kata maka berbanding terbalik dengan jumlah frekuensi keluar yang didapat, yaitu semakin kecil atau lebih jarang keluar. Penggunaan model Bi-Gram dan Tri-Gram untuk model bahasa masih memungkinkan, karena hasil dari jumlah frekuensi keluar pada suku *N-Gram*-nya masih cukup besar dan datanya masih valid apabila diproses lebih lanjut (Rostianingsih, Sugianto, dan Pustaka, 2012).

Penelitian yang menerapkan proses *N-Gram* adalah Pang, Lee, Rd, dan Jose (2002). *N-Gram* yang digunakan yakni Uni-Gram, yang dipakai dalam melakukan klasifikasi sentimen terhadap *review* film dengan menggunakan berbagai teknik pembelajaran mesin, dimana hasil penelitian ini menunjukkan bahwa metode *N-Gram* sebagai pendekatan untuk melakukan ekstraksi fitur. Dimana hasil klasifikasi dengan Uni-Gram dalam penelitian ini menghasilkan akurasi sebanyak 82,9%. Penelitian selanjutnya menerapkan prediksi kata untuk membantu mempercepat proses pengetikan dengan menerapkan *N-Gram* adalah Sugianto et al. (2013). Hasil penelitian ini menunjukkan bahwa *N-Gram* sebagai metode dasar dalam proses prediksi sangatlah membantu pemilahan kata. Sehingga proses prediksi menjadi

lebih efektif, mampu menghasilkan prediksi efektif di atas 20% dari total prediksi yang terjadi.

Berdasarkan dari hasil uraian tersebut, maka penelitian ini akan menguji pengaruh *N-Gram* pada hasil akurasi klasifikasi dokumen menggunakan algoritma *Naïve Bayes Classifier*.

### **1.3 Rumusan Masalah**

Rumusan masalah dalam penelitian ini adalah bagaimana pengaruh *N-Gram* pada klasifikasi dokumen menggunakan algoritma *Naïve Bayes Classifier*. Untuk menyelesaikan permasalahan tersebut, maka penelitian ini dibagi dalam beberapa *research question* antara lain:

1. Bagaimana mekanisme *N-Gram* dalam membangkitkan fitur?
2. Bagaimana mekanisme algoritma *Naïve Bayes Classifier* dalam mengklasifikasikan dokumen?
3. Bagaimana hasil akurasi klasifikasi dokumen menggunakan algoritma *Naïve Bayes Classifier*?
4. Bagaimana pengaruh *N-Gram* pada klasifikasi dokumen menggunakan algoritma *Naïve Bayes Classifier*?

### **1.4 Tujuan Penelitian**

Tujuan dari penelitian ini adalah:

1. Mengetahui mekanisme *N-Gram* dalam membangkitkan fitur.

2. Mengetahui mekanisme algoritma *Naïve Bayes Classifier* dalam mengklasifikasikan dokumen.
3. Mengetahui hasil akurasi klasifikasi dokumen menggunakan algoritma *Naïve Bayes Classifier*.
4. Mengetahui pengaruh *N-Gram* pada klasifikasi dokumen menggunakan algoritma *Naïve Bayes Classifier*.

### **1.5 Manfaat Penelitian**

Manfaat penelitian ini adalah:

1. Memahami mekanisme pembangkitan fitur pada *N-Gram* dan algoritma *Naïve Bayes Classifier* dalam proses klasifikasi dokumen.
2. Hasil penelitian dapat digunakan sebagai rujukan penelitian dalam pengaruh *N-Gram* pada klasifikasi dokumen menggunakan algoritma *Naïve Bayes Classifier*.

### **1.6 Batasan Masalah**

Adapun batasan masalah penelitian ini adalah sebagai berikut:

1. Metode *N-Gram* yang digunakan adalah *Bi-Gram*.
2. Dokumen yang digunakan merupakan dokumen berita online.
3. Berita online yang digunakan dalam klasifikasi yakni berita bahasa Indonesia.
4. Sumber data yang akan digunakan berasal dari situs <http://news.kompas.com> dan <http://www.republika.co.id>.

5. Sistem menerima masukkan berupa file dokumen bertipe teks dengan format teks (\*.txt).
6. Banyak kata dalam setiap file dokumen yakni lebih kurang 300 buah kata.
7. Kategori terdiri dari 5 kategori yakni berita olahraga, berita ekonomi, berita kesehatan, berita teknologi, dan berita politik.

### **1.7 Sistematika Penulisan**

Sistematika penulisan proposal skripsi ini adalah sebagai berikut:

#### **BAB I. PENDAHULUAN**

Bab I menguraikan latar belakang masalah, rumusan masalah, tujuan dan manfaat penelitian, batasan masalah, dan sistematika penulisan penelitian ini.

#### **BAB II. KAJIAN LITERATUR**

Bab II berisi landasan teori yang digunakan pada penelitian ini, antara lain definisi-definisi *text mining*, *text preprocessing*, *Naïve Bayes Classifier*, *N-Gram*, *term frequency*, berita dan *k-fold cross validation*. Selain itu bab II juga membahas penelitian-penelitian lain yang relevan dengan penelitian ini.

#### **BAB III. METODOLOGI PENELITIAN**

Bab III berisi pembahasan mengenai tahapan yang akan dilaksanakan pada penelitian ini. Rencana tahapan penelitian akan dideskripsikan dengan rinci dengan mengacu pada suatu kerangka kerja.

Di akhir bab III akan dijabarkan perancangan manajemen proyek untuk pelaksanaan penelitian ini.

## **1.8 Kesimpulan**

Bab ini telah dibahas latar belakang masalah penelitian dalam melakukan klasifikasi dokumen. Karena itu, penelitian ini akan menguji pengaruh *N-Gram* pada klasifikasi dokumen menggunakan algoritma *Naïve Bayes Classifier*. Dan hasil klasifikasinya akan dibandingkan, yakni hasil pengujian pada algoritma *Naïve Bayes Classifier* dengan *N-Gram* dan algoritma *Naïve Bayes Classifier* tanpa *N-Gram* dalam klasifikasi dokumen.



## DAFTAR PUSTAKA

- Fathan Hidayatullah, A., & Rifqi Ma, M. (2016). Penerapan Text Mining dalam Klasifikasi Judul Skripsi. *Seminar Nasional Aplikasi Teknologi Informasi (SNATI)*, 1907–5022.
- Hamzah, A. (2012). Klasifikasi Teks dengan Naïve Bayes Classifier (NBC) untuk Pengelompokan Teks Berita dan Abstract Akademis Amir. *Prosiding Seminar Nasional Aplikasi Sains & Teknologi (SNAST) Periode III*, 3(2011), 269–277. <https://doi.org/1979-911X>
- Indriani, A. (2014). Klasifikasi Data Forum dengan menggunakan Metode Naive Bayes Classifier. *Seminar Nasional Aplikasi Teknologi Informasi (SNATI)*, 28–33. <https://doi.org/10.1155/2011/172853>
- Juditha, C. (2013). News Accuracy in Online Journalism ( News of Alleged Corruption The Constitutional Court in Detiknews ). *Jurnal Pekomnas*. 16(3): 145–154.
- Kumar, S. S., & Rajini, A. (2018). An Efficient Sentimental Analysis for Twitter Using Neural Network based on Rmsprop. *IOSR Journal of Engineering (IOSRJEN)*, (Iccids), 17-25.
- Ginting, S., L., Br. & Trinanda, R. P. (2016). Teknik Data Mining Menggunakan Metode Bayes Classifier pada Aplikasi Perpustakaan. *Jurnal Teknologi dan Informasi (JATI)*. 3(2).
- Maghfira, T. N., Cholissodin, I., & Widodo, A. W. (2017). Deteksi Kesalahan Ejaan dan Penentuan Rekomendasi Koreksi Kata yang Tepat Pada Dokumen Jurnal JTIK Menggunakan Dictionary Lookup dan Damerau-Levenshtein Distance. *Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer*. 1(6): 498–506.
- Mooney, J. 2006. CS 391L : Machine Learning Text Categorization. *Austin : University of Texas*.
- Musthafa, A.. 2009. Klasifikasi Otomatis Dokumen Berita Kejadian Berbahasa Indonesia. Skripsi. Fakultas Sains dan Teknologi universitas Islam Negeri (UIN) Maulana Malik Ibrahim Malang
- Narayanan, V., Arora, I., & Bhatia, A. (2013b). Fast and Accurate Sentiment Classification Using an Enhanced Naive Bayes Model. H. *Yin et al. (Eds.): IDEAL 2013*, LNCS 8206, pp. 194–201.
- Sammut, C., & Webb, G. I. (2016). *Encyclopedia of Machine Learning and Data Mining*. Springer Publishing Company, Incorporated. <https://doi.org/10.1007/978-1-4899-7502-7>
- Setiawan, I., & Nursantika, D. (2017). Klasifikasi Artikel Berita Menggunakan

Metode Text Mining Dan Naive Bayes Classifier. *Prosiding SENIATI*. 1–6.

Sugianto, S. A., Liliana, L., & Rostianingsih, S. (2013). Pembuatan Aplikasi Predictive Text Menggunakan Metode N-gram-based. *Jurnal Infra*. 1(2). Retrieved from <https://www.neliti.com/id/publications/105718/pembuatan-aplikasi-predictive-text-menggunakan-metode-n-gram-based>

Wibowo, A. T., & Septiana, G. (2015). Pembobotan Fitur Ekstraksi Pada Peringkasan Teks Bahasa Indonesia Menggunakan Algoritma Genetika. *E-Proceeding of Engineering*. 2(2):6481–6489.

Wijaya, A. P. (2012). Klasifikasi Dokumen dengan Naive Bayes Classification ( NBC ) Untuk Mengatahui Konten. *Jurnal Datamining Indonesia*. 222D, 1–6.

Zaman, B., Hariyanti, E., Purwanti, E., & Bahasa, A. D. (2015). Sistem Deteksi Bahasa pada Dokumen menggunakan N-Gram. *Jurnal Multinetics*. 1(2): 21–26.

