

**SKRIPSI**

**PENGELOMPOKAN DOKUMEN BERBASIS ALGORITMA  
GENETIKA**



**Cornelius David Herianto**

**NPM: 2015730034**

**PROGRAM STUDI TEKNIK INFORMATIKA  
FAKULTAS TEKNOLOGI INFORMASI DAN SAINS  
UNIVERSITAS KATOLIK PARAHYANGAN  
2019**

**UNDERGRADUATE THESIS**

**GA-BASED DOCUMENT CLUSTERING**



**Cornelius David Herianto**

**NPM: 2015730034**

**DEPARTMENT OF INFORMATICS  
FACULTY OF INFORMATION TECHNOLOGY AND SCIENCES  
PARAHYANGAN CATHOLIC UNIVERSITY  
2019**

# LEMBAR PENGESAHAN

## PENGELOMPOKAN DOKUMEN BERBASIS ALGORITMA GENETIKA

Cornelius David Herianto

NPM: 2015730034

Bandung, 22 Mei 2019

Menyetujui,

Pembimbing

Kristopher David Harjono, M.T.

Ketua Tim Penguji

Anggota Tim Penguji

Husnul Hakim, M.T.

Rosa De Lima, M.Kom.

Mengetahui,

Ketua Program Studi

Mariskha Tri Adithia, P.D.Eng

## **PERNYATAAN**

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

### **PENGELOMPOKAN DOKUMEN BERBASIS ALGORITMA GENETIKA**

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 dijatuhkan 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 22 Mei 2019

Meterai Rp. 6000
---------------------

Cornelius David Herianto  
NPM: 2015730034

## ABSTRAK

Pengelompokan (*clustering*) merupakan sebuah metode untuk menggabungkan himpunan objek ke dalam kelompok-kelompok sedemikian rupa sehingga objek dalam kelompok (*cluster*) lebih mirip (karena suatu hal) satu sama lain daripada objek di kelompok lain [1]. *Document clustering* (pengelompokan dokumen) merupakan proses pengelompokan yang dilakukan terhadap suatu koleksi dokumen. Pengelompokan dokumen diterapkan dalam beberapa bidang seperti penambangan web, mesin pencari (*search engine*), dan temu kembali informasi (*information retrieval*) [3]. Hal yang dilakukan dalam pengelompokan dokumen adalah mengukur kemiripan (*similarity*) antar dokumen dan mengelompokkan dokumen yang serupa. Salah satu algoritma pengelompokan yang paling sering digunakan adalah *K-means*. Namun, algoritma *K-means* memiliki kekurangan yaitu dapat terjebak dalam *local optimum*. *Local optimum* adalah suatu solusi yang optimal (baik maksimal maupun minimal) diantara kandidat solusi yang berdekatan dalam masalah optimasi. Dikatakan lokal karena solusi ini hanya optimal apabila dibandingkan dengan kandidat solusi yang berdekatan, tidak optimal secara keseluruhan (*global optimum*).

Algoritma genetika atau biasa disebut *Genetic Algorithm* (GA) adalah suatu algoritma pencarian yang terinspirasi dari proses seleksi alam yang terjadi secara alami dalam proses evolusi. GA merupakan metode penyelesaian masalah yang menggunakan genetika sebagai pemodelannya. Dalam penelitian ini, GA akan digunakan sebagai solusi dari masalah *local optimum*. *Local optimum* dapat diatasi oleh GA yang sudah terbukti efektif dalam masalah pencarian dan optimasi. GA dapat digunakan untuk mengelompokkan dokumen dengan beberapa adaptasi terhadap representasi kromosom, fungsi *fitness*, seleksi, persilangan, dan mutasi.

Algoritma genetika dan algoritma *K-means* diuji menggunakan suatu *dataset* berlabel untuk membandingkan waktu dan hasil pengelompokan dari kedua algoritma tersebut. Berdasarkan hasil eksperimen menggunakan *dataset* dalam penelitian ini, rata-rata nilai *purity* dari hasil pengelompokan menggunakan algoritma genetika adalah sebesar 0.799, lebih baik 56% dibandingkan dengan menggunakan algoritma *K-means*. Hal ini membuktikan bahwa algoritma genetika sudah dapat mengelompokkan dokumen dengan hasil yang memuaskan. Namun dari segi waktu, algoritma genetika membutuhkan waktu 4365% lebih lama dibandingkan dengan algoritma *K-means*. Hal ini disebabkan oleh proses komputasi yang dilakukan pada algoritma genetika jauh lebih banyak dan kompleks dibandingkan dengan algoritma *K-means*.

**Kata-kata kunci:** Algoritma genetika, Pengelompokan dokumen, Algoritma *K-means*, TF-IDF, *Local optimum*

## ABSTRACT

Clustering is a method of creating groups of objects, or clusters, in such a way that objects in one cluster are very similar and objects in different clusters are quite distinct [1]. Document clustering is an organization of documents into clusters. It has been studied intensively because of its wide applicability in various areas such as web mining, search engines, and information retrieval [3]. It is measuring similarity between documents and grouping similar documents together. One of the most frequently used algorithm in clustering is K-means. However, K-means can easily stuck in local optimum. Local optimum of an optimization problem is a solution that is optimal (either maximal or minimal) within a neighboring set of candidate solutions. This is in contrast to a global optimum, which is the optimal solution among all possible solutions, not just those in a particular neighborhood of values.

Genetic Algorithm (GA) is a search algorithm inspired by the natural selection process that occurs naturally in the evolutionary process. GA is a problem solving method that used genetics as its model. A solution candidate is modeled as an individual in GA. Set of these individuals are called population. Every individual in a population is represented by a chromosome. Chromosome is a collection of parameters which formed a solution. That parameters is called gene. Every individual in the population is assigned, by means of a fitness function, a measure of its goodness. In this study, GA will be used as a solution to local optimum, which have proven to be effective in search and optimization problems. GA can be used to cluster documents by adapting chromosome representation, fitness function, selection, crossover, and mutation.

Genetic algorithm and K-means algorithm will be tested using a labeled datasets to compare the running time and clustering result. Based on the experimental results using the datasets in this study, the average purity value of clustering results using genetic algorithms is 0.799, 56% greater than using the K-means algorithm. This proves that the genetic algorithm is able to cluster documents with satisfactory results. But in terms of running time, genetic algorithms take 4365% more time than the K-means algorithm. This is caused by the computational process carried out on the genetic algorithm is far more complex than the K-means algorithm.

**Keywords:** Genetic algorithm, Document clustering, K-means algorithm, TF-IDF, Local optimum

*Dipersembahkan kepada Tuhan YME, keluarga tercinta,  
dan diri sendiri*

## KATA PENGANTAR

Puji syukur kepada Tuhan Yang Maha Esa atas berkat yang diberikan kepada penulis sehingga dapat menyelesaikan skripsi dengan judul **Pengelompokan Dokumen Berbasis Algoritma Genetika** dengan baik dan tepat waktu. Selama menjalani proses perkuliahan dan penyusunan skripsi, penulis telah mendapat banyak bantuan dan dukungan dalam menghadapi hambatan yang ada. Oleh karena itu, penulis ingin menyampaikan rasa terima kasih kepada:

1. Keluarga penulis yaitu Papa dan Mama yang selalu mendukung penulis secara moral dan materiil sehingga penulis dapat menyelesaikan proses perkuliahan dan skripsi ini dengan baik. Serta kepada Michelle dan Vincent sebagai adik penulis yang selalu mendukung penulis dalam menyelesaikan skripsi ini.
2. Bapak Kristopher David Harjono, M.T. sebagai dosen pembimbing yang telah membimbing penulis hingga dapat menyelesaikan skripsi ini.
3. Bapak Husnul Hakim, M.T. dan Ibu Rosa De Lima, M.Kom. sebagai dosen penguji yang telah membantu dalam menguji dan memperbaiki skripsi ini.
4. Ibu Mariskha Tri Adithia, P.D.Eng selaku Ketua Program Studi Teknik Informatika Fakultas Teknologi Informasi dan Sains Universitas Katolik Parahyangan.
5. Arlin Sasqia Puspa Shiffa sebagai sahabat yang selalu mendampingi penulis dalam penyusunan skripsi dari awal hingga akhir bahkan membantu penulis mempersiapkan presentasi sidang skripsi.
6. Vania Stephanie dan Khezia Josephine sebagai sahabat penulis yang senantiasa memberikan dukungan, semangat, masukan yang berguna untuk penulis, serta menghibur penulis selama proses penyusunan skripsi ini terutama saat sedang mengalami kesulitan.
7. Teman-teman dari Grup MANAYGBILANGWGAGUNA yaitu AK, Otung, Devie, Gilbert, Hoshea, Jeane, Khen, Mark, Matthew, Oyeng, Bebe, Rifo, Rizky, Sean, Vinny, dan WM yang telah menghibur dan mendukung penulis dalam penyusunan skripsi ini.
8. Teman-teman dari Grup Korea yaitu Ario, Dandy, Fakhry, Felis, Hima, Hizkia, Irvan, Acong, Joshua, Kezia, Edrick, Ocín, Momon, Thoby, Victor, Yona, Yudhis, dan Matthew Ariel sebagai teman seperjuangan di Teknik Informatika UNPAR angkatan 2015 yang telah memberi semangat dan dukungan kepada penulis.
9. Teman-teman penulis lain yang tidak dapat disebutkan satu persatu. Terima kasih untuk segala dukungannya sehingga skripsi ini dapat diselesaikan dengan baik.

Akhir kata, semoga skripsi ini dapat bermanfaat bagi pembaca dan dapat menjadi dasar untuk penelitian yang terkait dengan skripsi ini.

Bandung, Mei 2019

Penulis

# DAFTAR ISI

<b>KATA PENGANTAR</b>	<b>xv</b>
<b>DAFTAR ISI</b>	<b>xvii</b>
<b>DAFTAR GAMBAR</b>	<b>xix</b>
<b>DAFTAR TABEL</b>	<b>xxi</b>
<b>1 PENDAHULUAN</b>	<b>1</b>
1.1 Latar Belakang . . . . .	1
1.2 Rumusan Masalah . . . . .	2
1.3 Tujuan Penelitian . . . . .	2
1.4 Batasan Masalah . . . . .	2
1.5 Metodologi . . . . .	2
1.6 Sistematika Pembahasan . . . . .	3
<b>2 LANDASAN TEORI</b>	<b>5</b>
2.1 Pengelompokan . . . . .	5
2.1.1 Definisi Pengelompokan . . . . .	5
2.1.2 Aplikasi Pengelompokan . . . . .	6
2.1.3 <i>Local Optimum</i> . . . . .	7
2.2 <i>K-Means</i> . . . . .	7
2.3 Algoritma Genetika . . . . .	8
2.3.1 <i>Fitness</i> . . . . .	9
2.3.2 Seleksi . . . . .	10
2.3.3 Persilangan . . . . .	11
2.3.4 Mutasi . . . . .	11
2.3.5 Proses Pencarian Dalam Algoritma Genetika . . . . .	11
2.3.6 GA dalam Pengelompokan . . . . .	14
2.4 Model Ruang Vektor . . . . .	14
2.5 Pembobotan <i>Term</i> ( <i>Term Weighting</i> ) . . . . .	14
2.5.1 Bobot frekuensi . . . . .	15
2.5.2 Bobot TF-IDF . . . . .	15
2.6 Metrik <i>Intrachuster</i> untuk Mengukur Kinerja Metode <i>Clustering</i> . . . . .	16
2.7 <i>Clustering Purity</i> . . . . .	17
<b>3 ANALISIS</b>	<b>19</b>
3.1 Analisis <i>Dataset</i> . . . . .	19
3.2 Representasi Dokumen . . . . .	19
3.3 Model Ruang Vektor . . . . .	20
3.3.1 Bobot Frekuensi . . . . .	20
3.3.2 Bobot TF-IDF . . . . .	20
3.4 Representasi Kromosom . . . . .	22

3.5	Fungsi <i>Fitness</i> . . . . .	23
3.6	Operasi Genetik Dalam Pengelompokan Dokumen . . . . .	24
3.6.1	Inisialisasi Populasi . . . . .	24
3.6.2	Seleksi . . . . .	24
3.6.3	Persilangan . . . . .	24
3.6.4	Mutasi . . . . .	25
3.7	Evaluasi Hasil Pengelompokan Menggunakan <i>Purity</i> . . . . .	26
<b>4</b>	<b>PERANCANGAN</b> . . . . .	<b>27</b>
4.1	Kebutuhan Masukan dan Keluaran . . . . .	27
4.2	Rancangan Kelas . . . . .	29
4.2.1	<i>Document</i> . . . . .	30
4.2.2	<i>Vector</i> . . . . .	31
4.2.3	<i>SimilarityCalculator</i> . . . . .	32
4.2.4	<i>CosineSimilarityCalculator</i> . . . . .	32
4.2.5	<i>TermWeighting</i> . . . . .	33
4.2.6	<i>FrequencyWeighting</i> . . . . .	33
4.2.7	<i>TFIDFWeighting</i> . . . . .	33
4.2.8	<i>Lexicon</i> . . . . .	34
4.2.9	<i>Gene</i> . . . . .	35
4.2.10	<i>Chromosome</i> . . . . .	35
4.2.11	<i>GAClusterer</i> . . . . .	37
4.2.12	<i>Params</i> . . . . .	39
4.2.13	<i>KMeans</i> . . . . .	41
4.2.14	<i>FXMLLDocumentController</i> . . . . .	43
4.3	Perancangan Antarmuka Pengguna . . . . .	45
4.3.1	Halaman Algoritma Genetika . . . . .	46
4.3.2	Halaman <i>K-Means</i> . . . . .	48
<b>5</b>	<b>PENGUJIAN DAN EKSPERIMEN</b> . . . . .	<b>51</b>
5.1	Skenario Pengujian Eksperimental . . . . .	51
5.2	Eksperimen Algoritma Genetika . . . . .	53
5.3	Eksperimen <i>K-Means</i> . . . . .	62
5.4	Analisis Hasil Eksperimen . . . . .	64
<b>6</b>	<b>KESIMPULAN DAN SARAN</b> . . . . .	<b>67</b>
6.1	Kesimpulan . . . . .	67
6.2	Saran . . . . .	68
	<b>DAFTAR REFERENSI</b> . . . . .	<b>69</b>
	<b>A KODE PROGRAM</b> . . . . .	<b>71</b>
	<b>B HASIL EKSPERIMEN</b> . . . . .	<b>91</b>
	<b>C CONTOH DATASET</b> . . . . .	<b>97</b>
	C.1 Topik <i>Business</i> . . . . .	97

## DAFTAR GAMBAR

2.1	Contoh <i>cluster</i> hasil pengelompokan . . . . .	5
2.2	Hasil pencarian ( <i>clustered search result</i> ) di Yippy . . . . .	6
2.3	<i>Local Optimum</i> dan <i>Global Optimum</i> . . . . .	7
2.4	Alur algoritma genetika dasar . . . . .	9
2.5	Ilustrasi <i>roulette-wheel</i> . . . . .	10
2.6	<i>Single-point crossover</i> . . . . .	11
2.7	<i>Ilustrasi untuk jarak intracuster</i> . . . . .	17
3.1	Formula representasi kromosom . . . . .	22
3.2	Contoh representasi <i>centroid</i> ke dalam kromosom . . . . .	23
3.3	Ilustrasi kromosom untuk persilangan . . . . .	24
3.4	Ilustrasi persilangan . . . . .	25
3.5	Ilustrasi mutasi kromosom . . . . .	25
4.1	Diagram kelas . . . . .	29
4.2	Kelas <i>Document</i> . . . . .	30
4.3	Kelas <i>Vector</i> . . . . .	31
4.4	Kelas <i>SimilarityCalculator</i> . . . . .	32
4.5	Kelas <i>CosineSimilarityCalculator</i> . . . . .	32
4.6	Kelas <i>TermWeighting</i> . . . . .	33
4.7	Kelas <i>FrequencyWeighting</i> . . . . .	33
4.8	Kelas <i>TFIDFWeighting</i> . . . . .	33
4.9	Kelas <i>Lexicon</i> . . . . .	34
4.10	Kelas <i>Gene</i> . . . . .	35
4.11	Kelas <i>Chromosome</i> . . . . .	35
4.12	Kelas <i>GAClusterer</i> . . . . .	37
4.13	Kelas <i>Params</i> . . . . .	39
4.14	Kelas <i>KMeans</i> . . . . .	41
4.15	Kelas <i>FXMLDocumentController</i> . . . . .	43
4.16	Rancangan antarmuka halaman algoritma genetika . . . . .	46
4.17	Rancangan antarmuka halaman algoritma genetika . . . . .	48
5.1	Grafik hubungan banyaknya populasi dengan waktu pengelompokan . . . . .	53
5.2	Diagram hubungan banyaknya populasi dengan <i>intracuster similarity</i> . . . . .	54
5.3	Diagram hubungan banyaknya populasi dengan banyaknya iterasi . . . . .	54
5.4	Diagram hubungan banyaknya populasi dengan nilai <i>purity</i> . . . . .	55
5.5	Diagram hubungan metode pembobotan dengan waktu pengelompokan . . . . .	55
5.6	Diagram hubungan metode pembobotan dengan <i>intracuster similarity</i> . . . . .	56
5.7	Diagram hubungan metode pembobotan dengan banyaknya iterasi . . . . .	56
5.8	Diagram hubungan metode pembobotan dengan nilai <i>purity</i> . . . . .	57
5.9	Diagram hubungan probabilitas mutasi dengan waktu pengelompokan . . . . .	57
5.10	Diagram hubungan probabilitas mutasi dengan <i>intracuster similarity</i> . . . . .	58
5.11	Diagram hubungan probabilitas mutasi dengan banyaknya iterasi . . . . .	58

5.12	Diagram hubungan probabilitas mutasi dengan nilai <i>purity</i> . . . . .	59
5.13	Diagram hubungan individu elitisme dengan waktu pengelompokan . . . . .	60
5.14	Diagram hubungan individu elitisme dengan <i>intracluster similarity</i> . . . . .	60
5.15	Diagram hubungan individu elitisme dengan banyaknya iterasi . . . . .	61
5.16	Diagram hubungan individu elitisme dengan nilai <i>purity</i> . . . . .	61
5.17	Diagram hubungan algoritma dengan waktu tempuh . . . . .	62
5.18	Diagram hubungan algoritma dengan <i>intracluster similarity</i> . . . . .	63
5.19	Diagram hubungan algoritma dengan banyaknya iterasi . . . . .	63
5.20	Diagram hubungan algoritma dengan nilai <i>purity</i> . . . . .	64

## DAFTAR TABEL

2.1	<i>Term-document incidence matrix</i> . . . . .	15
2.2	Tabel hasil pengelompokan . . . . .	18
3.1	Hasil perhitungan bobot frekuensi . . . . .	20
3.2	Hasil perhitungan TF . . . . .	21
3.3	Hasil perhitungan IDF . . . . .	21
3.4	Hasil perhitungan bobot TF-IDF . . . . .	22
4.1	Contoh keluaran dalam bentuk <i>file CSV</i> . . . . .	28
4.2	Rincian <i>field</i> pada halaman algoritma genetika . . . . .	47
4.3	Rincian <i>field</i> pada halaman algoritma <i>K-means</i> . . . . .	49
5.1	Variasi nilai variabel bebas . . . . .	52
5.2	Rata-rata hasil pengelompokan dengan variasi variabel banyaknya populasi . . . . .	53
5.3	Rata-rata hasil pengelompokan dengan variasi variabel metode pembobotan . . . . .	55
5.4	Rata-rata hasil pengelompokan dengan variasi variabel probabilitas mutasi . . . . .	57
5.5	Rata-rata hasil pengelompokan dengan variasi variabel individu elitisme . . . . .	60
5.6	Rata-rata hasil pengelompokan dengan menggunakan algoritma <i>K-means</i> . . . . .	62
5.7	Hasil perhitungan statistika terhadap hasil eksperimen algoritma genetika dan <i>K-means</i> . . . . .	64
B.1	Hasil eksperimen kasus uji 1 (parameter ideal) . . . . .	91
B.2	Hasil eksperimen kasus uji 2 (Populasi=50) . . . . .	92
B.3	Hasil eksperimen kasus uji 3 (Populasi=150) . . . . .	92
B.4	Hasil eksperimen kasus uji 4 (Bobot frekuensi) . . . . .	93
B.5	Hasil eksperimen kasus uji 5 (Probabilitas mutasi=0) . . . . .	93
B.6	Hasil eksperimen kasus uji 6 (Probabilitas mutasi=0.25) . . . . .	94
B.7	Hasil eksperimen kasus uji 7 (Individu elitisme=0) . . . . .	94
B.8	Hasil eksperimen kasus uji 8 (Individu elitisme=5) . . . . .	95
B.9	Hasil eksperimen algoritma <i>K-means</i> . . . . .	95

# BAB 1

## PENDAHULUAN

### 1.1 Latar Belakang

Pengelompokan (*clustering*) merupakan prosedur untuk mencari struktur alami dari suatu kumpulan data. Proses ini melibatkan pemilihan data atau objek ke dalam kelompok (*cluster*) sehingga objek-objek dalam cluster yang sama akan lebih mirip satu sama lain dibandingkan dengan objek yang berada di *cluster* lain. *Clustering* berguna untuk mereduksi data (mereduksi data dengan volume besar ke dalam kelompok-kelompok dengan karakteristik tertentu), mengembangkan skema klasifikasi (juga dikenal sebagai taksonomi), dan memberikan masukan atau dukungan terhadap hipotesis mengenai struktur suatu data.

*Clustering* merupakan salah satu teknik pembelajaran tak terarah (*unsupervised learning*). Pembagian kelompok dalam *clustering* tidak berdasarkan sesuatu yang telah diketahui sebelumnya, melainkan berdasarkan kesamaan tertentu menurut suatu ukuran tertentu [2].

*Document clustering* (pengelompokan dokumen) merupakan proses pengelompokan yang dilakukan terhadap suatu koleksi dokumen. Pengelompokan dokumen diterapkan dalam beberapa bidang seperti penambangan web, mesin pencari (*search engine*), dan temu kembali informasi (*information retrieval*) [3]. Hal yang dilakukan dalam pengelompokan dokumen adalah mengukur kemiripan (*similarity*) antar dokumen dan mengelompokkan dokumen yang serupa. Suatu dokumen dapat terdiri dari beberapa jenis informasi seperti teks, jenis tulisan, ukuran tulisan, warna tulisan, dan gambar.

Salah satu algoritma pengelompokan yang paling sering digunakan adalah *K-means* yang dilakukan dengan cara membagi data ke dalam  $K$  kelompok. Kelompok tersebut dibentuk dengan cara meminimalkan jarak antara titik pusat *cluster* (*centroid*) dengan setiap anggota *cluster* tersebut. Titik pusat *cluster* dicari dengan menggunakan rata-rata (*mean*) dari nilai setiap anggota *cluster*. Dalam hal ini, setiap anggota *cluster* dimodelkan sebagai vektor dalam  $n$  dimensi ( $n$  merupakan banyaknya atribut). *K-means* sudah terbukti efektif dalam melakukan pengelompokan dalam situasi apapun. Namun, cara tersebut tetap saja memiliki kekurangan yaitu dapat terjebak dalam *local optima* tergantung dengan pemilihan *centroid* awal [4].

Masalah *local optima* dapat ditangani menggunakan *Genetic Algorithm* (GA) yang telah terbukti efektif dalam menyelesaikan masalah pencarian dan optimasi. GA merupakan teknik pencarian heuristik tingkat tinggi yang menirukan proses evolusi yang secara alami terjadi [5] berdasarkan prinsip *survival of the fittest*. Algoritma ini dinamakan demikian karena menggunakan konsep-konsep dalam genetika sebagai model pemecahan masalahnya [6].

Dalam GA, parameter dari *search space* dikodekan dalam bentuk deretan objek yang disebut kromosom. Kumpulan kromosom tersebut lalu dikenal sebagai populasi. Pada awalnya, populasi dibangkitkan secara acak. Kemudian, akan dipilih beberapa kromosom menggunakan teknik *roulette wheel selection* berdasarkan fungsi *fitness*. Operasi dasar yang terinspirasi dari Ilmu Biologi seperti persilangan (*crossover*) dan mutasi (*mutation*) digunakan untuk membangkitkan generasi berikutnya. Proses seleksi, persilangan, dan mutasi ini berlangsung dalam jumlah generasi tertentu atau sampai kondisi akhir tercapai.

Fungsi *fitness* tidak hanya berfungsi untuk menentukan seberapa baik solusi yang dihasilkan

namun juga menentukan seberapa dekat solusi tersebut dengan hasil yang optimal [6]. Oleh karena itu, diperlukan fungsi *fitness* yang cocok sehingga GA dapat menghasilkan keluaran yang optimal. Pada masalah *clustering* menggunakan GA, maka fungsi *fitness* yang digunakan harus bisa menggambarkan bahwa seluruh elemen sudah berada dalam *cluster* yang terbaik dan sudah sesuai.

## 1.2 Rumusan Masalah

Berdasarkan latar belakang yang telah dipaparkan, rumusan masalah dari penelitian ini adalah sebagai berikut:

1. Bagaimana algoritma genetik dapat digunakan untuk mengelompokkan dokumen?
2. Bagaimana membangun perangkat lunak yang menggunakan algoritma genetik untuk dapat mengelompokkan dokumen?

## 1.3 Tujuan Penelitian

Berdasarkan rumusan masalah yang telah disebutkan, tujuan dari penelitian ini adalah sebagai berikut:

1. Mempelajari algoritma genetik dan hubungannya dengan pengelompokan dokumen.
2. Membangun perangkat lunak yang mengimplementasikan algoritma genetik untuk dapat mengelompokkan dokumen.

## 1.4 Batasan Masalah

Rumusan masalah yang telah disebutkan memiliki ruang lingkup yang cukup luas. Dengan menyadari terbatasnya waktu serta kemampuan, penelitian ini difokuskan dengan memperlihatkan batasan masalah sebagai berikut:

1. Jenis dokumen yang dapat diproses dengan perangkat lunak yang dibuat hanyalah *Text Document* dengan ekstensi *TXT*.
2. Informasi dari dokumen yang diproses dalam pengelompokan hanya berasal dari teks yang menjadi isi dari dokumen tersebut. Gambar dan *metadata* (pemilik, tanggal modifikasi) tidak diperhitungkan.

## 1.5 Metodologi

Langkah-langkah yang dilakukan dalam penelitian ini adalah:

1. Melakukan studi literatur mengenai model ruang vektor, *Document Clustering* (pengelompokan dokumen), *Genetic Algorithm* (algoritma genetik), dan penggunaan algoritma genetik dalam pengelompokan dokumen.
2. Mencari dokumen yang dijadikan *datasets*.
3. Membuat rancangan perangkat lunak yang menggunakan algoritma genetik sebagai algoritma pengelompokan dokumen.
4. Mengimplementasikan hasil rancangan menjadi perangkat lunak dalam bahasa pemrograman Java.

5. Melatih dan menguji perangkat lunak dengan dokumen yang telah tersedia.
6. Mengevaluasi hasil pengujian lalu melakukan implementasi dan pengujian kembali sampai didapatkan hasil yang sudah sesuai dengan harapan.

## 1.6 Sistematika Pembahasan

Dokumentasi dari penelitian ini disajikan dalam enam bab dengan sistematika pembahasan sebagai berikut:

1. Bab 1 Pendahuluan  
Bab 1 berisi latar belakang pemilihan "Pengelompokan Dokumen berbasis Algoritma Genetika" sebagai judul dari penelitian ini. Selain itu, dibahas juga rumusan masalah, tujuan penelitian, batasan masalah, serta metodologi penelitian yang menjadi acuan dari penelitian ini.
2. Bab 2 Landasan Teori  
Bab 2 memuat landasan teori yang digunakan dalam penelitian ini. Konsep-konsep yang dibahas yaitu pengelompokan, *local optimum*, K-means, algoritma genetika beserta seluruh operasinya, model ruang vektor, pembobotan term yang terdiri dari bobot frekuensi dan bobot TF-IDF, metrik *Intracuster* untuk mengukur kinerja metode *clustering*, dan evaluasi hasil pengelompokan menggunakan *purity*.
3. Bab 3 Analisis  
Bab 3 memuat hasil analisis berdasarkan landasan teori. Hasil analisis yang ditulis pada bab 3 antara lain analisis *dataset*, representasi dokumen, modifikasi terhadap model ruang vektor beserta metode pembobotannya, representasi kromosom, fungsi *fitness*, dan operasi genetika lainnya.
4. Bab 4 Perancangan  
Bab 4 memuat hasil perancangan berdasarkan hasil analisis pada bab 3. Terdapat tiga bagian dalam bab perancangan yaitu Kebutuhan masukan dan keluaran, perancangan kelas, dan perancangan antarmuka pengguna.
5. Bab 5 Pengujian dan Eksperimen  
Bab 5 memuat hasil pengujian dan eksperimen yang telah dilakukan. Pada bab ini dibahas mengenai skenario pengujian, eksperimen pada algoritma genetika, dan eksperimen pada algoritma *K-means*.
6. Bab 6 Kesimpulan dan Saran  
Bab 6 memuat kesimpulan dari penulis berdasarkan hasil penelitian yang telah dilakukan dan saran untuk peneliti berikutnya agar dapat mengembangkan penelitian ini menjadi lebih baik lagi.