

PENGELOMPOKAN DOKUMEN BERBASIS ALGORITMA GENETIKA

CORNELIUS DAVID HERIANTO—2015730034

1 Data Skripsi

Pembimbing utama/tunggal: **Thomas Anung Basuki**

Pembimbing pendamping: -

Kode Topik : **TAB4502**

Topik ini sudah dikerjakan selama : **1 semester**

Pengambilan pertama kali topik ini pada : Semester **45 - Ganjil 18/19**

Pengambilan pertama kali topik ini di kuliah : **Skripsi 1**

Tipe Laporan : **C - Dokumen pendukung untuk pengambilan ke-2 di Skripsi 2**

2 Latar Belakang

Pengelompokan (*clustering*) merupakan teknik yang digunakan untuk membagi kumpulan data (*dataset*) ke dalam kelompok objek serupa. Setiap kelompok (*cluster*) terdiri dari objek-objek yang serupa satu sama lain berdasarkan suatu ukuran dan tidak serupa dengan objek-objek lain di luar kelompoknya.

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.

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.

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 berdasarkan prinsip *survival of the fittest*. Algoritma ini dinamakan demikian karena menggunakan konsep-konsep dalam genetika sebagai model pemecahan masalahnya.

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. Oleh karena itu, diperlukan fungsi *fitness* yang cocok sehingga GA dapat menghasilkan keluaran yang optimal. Pada masalah *clus-*

tering menggunakan GA, maka fungsi *fitness* yang digunakan harus bisa menggambarkan bahwa seluruh elemen sudah berada dalam *cluster* yang terbaik dan sudah sesuai.

3 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?

4 Tujuan

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.

5 Detail Perkembangan Pengerjaan Skripsi

Detail bagian pekerjaan skripsi sesuai dengan rencan kerja/laporan perkembangan terakhir :

1. **Melakukan studi literatur mengenai *Information Retrieval* (temu kembali informasi).**

Status : Baru ditambahkan semester ini

Hasil : Temu kembali informasi adalah bidang ilmu yang berurusan dengan representasi, penyimpanan, pengolahan, dan akses terhadap informasi. Pada penelitian ini, temu kembali informasi memiliki peran untuk mengubah dokumen yang semula berbentuk teks menjadi sesuatu yang memiliki nilai informasi agar nantinya bisa diproses lebih lanjut (dalam kasus ini dengan pengelompokan dokumen).

Dalam penelitian ini, temu kembali informasi digunakan untuk merepresentasikan dokumen ke dalam model ruang vektor agar bisa dilakukan pengelompokan. Tahap pertama yang dilakukan adalah membentuk kosa kata (*vocabulary*) yang berisi seluruh istilah berbeda yang ada di setiap dokumen yang akan diindeks. Kemudian, untuk setiap dokumen akan dibentuk suatu indeks yang terdiri dari pasangan antara istilah di dokumen tersebut dan jumlah kemunculannya (*term-document incidence matrix*). Tabel 1 merupakan contoh dari *term-document incidence matrix* dengan baris pada tabel menunjukkan istilah (*term*) dan kolom menunjukkan nama dokumen. Sel bernilai 1 apabila dokumen mengandung term tertentu dan 0 jika tidak.

Selanjutnya, setiap jumlah kemunculan suatu istilah dalam sebuah dokumen akan diubah menjadi suatu bobot tertentu berdasarkan teknik pemodelannya. Dalam penelitian ini, digunakan 2 teknik pembobotan yaitu bobot frekuensi dan bobot tf-idf.

Bobot frekuensi

Bobot frekuensi merupakan teknik pembobotan yang sangat sederhana karena bobotnya merupakan jumlah kemunculan istilah tersebut dalam dokumen. Bobot frekuensi dapat digambarkan dalam

rumus:

$$w_i = tf_i \quad (1)$$

dengan w_i merupakan bobot istilah ke- i dan tf_i merupakan frekuensi kemunculan istilah ke- i pada dokumen.

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth	...
Anthony	1	1	0	0	0	1	
Brutus	1	1	0	1	0	0	
Caesar	1	1	0	1	1	1	
Calpurnia	0	1	0	0	0	0	
Cleopatra	1	0	0	0	0	0	
mercy	1	0	1	1	1	1	
worser	1	0	1	1	1	0	
...							

Tabel 1: *Term-document incidence matrix*

Bobot TF-IDF

Pada bobot frekuensi, bobot hanya dihitung berdasarkan kemunculan istilah dalam dokumen itu sendiri. Namun dalam bobot TF-IDF (*term frequency-inverse document frequency*), bobot juga dihitung berdasarkan kemunculan istilah pada himpunan dokumen. Metode ini sangat populer digunakan oleh sistem rekomendasi berbasis teks. Rumus dari TF-IDF adalah sebagai berikut:

$$w_i = tf_i \times \log \frac{N}{N_i} \quad (2)$$

dengan w_i merupakan bobot istilah ke- i , tf_i merupakan frekuensi kemunculan istilah ke- i pada dokumen, N menyatakan banyaknya anggota himpunan dokumen, dan N_i menunjukkan frekuensi dokumen dari istilah ke- i (banyaknya dokumen pada himpunan dokumen yang memuat istilah ke- i).

Berdasarkan rumus tersebut, maka dapat ditarik dua kesimpulan yaitu:

- Semakin sering suatu istilah muncul di suatu dokumen, maka semakin representatif istilah tersebut terhadap isi dokumen.
- Semakin banyak dokumen yang memuat suatu istilah, maka nilai informasi istilah tersebut semakin kecil.

Metode penetapan bobot TF-IDF dianggap sebagai metode yang berkinerja baik karena mempertimbangkan frekuensi kemunculan istilah baik secara lokal (TF) maupun global (IDF).

Model ruang vektor

Model ruang vektor adalah representasi dari koleksi dokumen sebagai vektor dalam ruang vektor yang umum. Model ruang vektor ini biasanya digunakan dalam sejumlah operasi pencarian informasi mulai dari penilaian dokumen pada *query*, klasifikasi dokumen dan pengelompokan dokumen.

Dalam pengolahan model ruang vektor, dibutuhkan cara untuk menghitung kesamaan antara dua ruang vektor. Dalam pengelompokan dokumen, apabila kesamaan antara dua buah vektor semakin besar, maka peluang kedua vektor tersebut berada dalam sebuah kelompok yang sama akan semakin besar. Dalam penelitian ini, digunakan dua cara untuk menghitung kesamaan antara dua ruang

vektor yaitu menggunakan Jarak Euclidean (*Euclidean distance*) dan persamaan cosinus (*cosine similarity*).

Jarak Euclidean

Jarak Euclidean atau biasa disebut jarak garis lurus merupakan metode paling banyak digunakan untuk menghitung jarak antara dua buah vektor. Secara umum, rumus dari Jarak Euclidean adalah sebagai berikut:

$$d_{ij} = \sqrt{\sum_{v=1}^N (x_{vi} - x_{vj})^2} \quad (3)$$

dengan d_{ij} adalah jarak antara vektor ke- i dengan vektor ke- j , N adalah jumlah dimensi pada vektor, dan x_{vi} adalah nilai dimensi ke- v dari vektor ke- i .

Persamaan Cosinus

Persamaan cosinus merupakan normalisasi hasil kali titik dengan panjang masing-masing vektor. Persamaan cosinus memiliki persamaan sebagai berikut:

$$s_{ij} = \frac{i \cdot j}{\|i\| \times \|j\|} \quad (4)$$

dengan s_{ij} adalah kesamaan antara vektor ke- i dengan vektor ke- j , i adalah vektor ke- i , dan j adalah vektor ke- j . Persamaan ini menjelaskan bahwa semakin kecil sudut antara dua vektor, maka tingkat kemiripannya semakin besar.

2. Melakukan studi literatur mengenai *Document Clustering* (pengelompokan dokumen).

Status : Ada sejak rencana kerja skripsi

Hasil : Tugas pengelompokan (*clustering*) merupakan salah satu bagian dari pembelajaran mesin. Pengelompokan terdiri dari dua jenis berdasarkan metode pembelajarannya. *Classification* merupakan jenis pembelajaran terarah karena sudah diberikan label sejak awal, lalu dilakukan pengelompokan berdasarkan label untuk setiap kelompoknya. Sedangkan *clustering* merupakan jenis pembelajaran tak terarah karena tidak diberikan label sejak awal sehingga pengelompokan dilakukan berdasarkan suatu kesamaan tertentu.

K-Means

K-means merupakan algoritma pengelompokan yang paling populer digunakan saat ini. algoritma ini membagi data ke dalam K cluster. Setiap cluster direpresentasikan dengan titik tengahnya (*centroid*). Setiap iterasi, titik tengah akan dihitung sebagai rata-rata dari semua titik data dari cluster tersebut. Rumus untuk menghitung *centroid* adalah sebagai berikut:

$$\mu_i = \frac{1}{N_i} \sum_{q=1}^{N_i} x_q \quad (5)$$

dengan μ_i merupakan *centroid* ke- i , N_i merupakan jumlah titik data pada cluster ke- i , dan x_q merupakan titik ke- q pada cluster ke- i .

Algoritma K-means diawali dengan penentuan *centroid* secara acak. Pada setiap iterasi, setiap data ditetapkan pada cluster yang memiliki *centroid* dengan jarak terdekat dari titik data tersebut, kemudian posisi *centroid* dari setiap cluster akan dihitung ulang dengan Persamaan 5. Iterasi akan terus diulang sampai posisi dari semua cluster tidak berubah.

Algorithm 1 K-Means**Input:** S (himpunan titik data), K (Jumlah *cluster*)**Output:** himpunan *cluster*

- 1: Pilih K titik data sebagai himpunan awal *centroid*.
- 2: **repeat**
- 3: Bentuk K *cluster* dengan menempatkan setiap titik data ke *cluster* dengan *centroid* terdekat.
- 4: Hitung ulang *centroid* untuk setiap *cluster*.
- 5: **until** *Centroid* tidak berubah.

3. Melakukan studi literatur mengenai Genetic Algorithm (algoritma genetik).**Status :** Ada sejak rencana kerja skripsi**Hasil :** Algoritma genetika (GA) adalah suatu algoritma pencarian yang terinspirasi dari proses seleksi alam yang terjadi secara alami dalam proses evolusi. Ada beberapa istilah yang digunakan dalam algoritma genetika di antaranya kromosom, seleksi, persilangan, mutasi, dan fungsi *fitness*.**Kromosom**

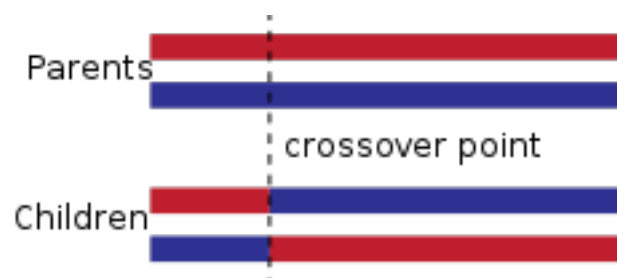
Dalam GA, kromosom adalah himpunan parameter yang mendefinisikan suatu solusi yang diusulkan. Kromosom biasanya direpresentasikan sebagai string yang berisi kumpulan nilai (gen), meskipun berbagai struktur data lainnya juga digunakan.

Seleksi

Seleksi dalam algoritma genetika bertugas untuk memilih kromosom dari populasi untuk proses persilangan. Salah satu teknik yang populer digunakan dalam seleksi adalah *roulette-wheel selection*. Teknik ini memungkinkan terjadinya pemilihan berdasarkan fungsi *fitness* sehingga kromosom yang memiliki fungsi *fitness* lebih besar memiliki peluang lebih besar untuk terpilih dan menjadi induk dari generasi selanjutnya.

Persilangan

Persilangan adalah operasi genetik yang digunakan untuk menggabungkan informasi genetik dari dua induk untuk menghasilkan keturunan baru. Teknik persilangan yang digunakan dalam penelitian ini adalah *Single-point crossover*. Dalam teknik ini, sebuah titik pada kedua induk dipilih untuk menjadi titik persilangan (*crossover point*). Bit yang berada di sebelah kanan titik persilangan bertukar antara kedua kromosom induk seperti yang ditunjukkan pada Gambar 1.

Gambar 1: *Single-point crossover*

Mutasi

Mutasi adalah suatu operator genetik yang digunakan untuk mempertahankan keragaman genetik dari satu generasi populasi dalam algoritma genetika. Mutasi mengubah satu atau beberapa nilai dalam gen. Mutasi terjadi berdasarkan probabilitas mutasi yang sudah ditentukan sebelumnya. Probabilitas ini seharusnya bernilai kecil, karena jika terlalu besar maka akan menjadi sama dengan algoritma pencarian acak primitif (*primitive random search*). Pada penelitian ini, kromosom memiliki gen yang bernilai non-biner. Oleh karena itu, mutasi dilakukan dengan memilih sebuah bilangan acak antara batas bawah dan batas atas nilai suatu gen.

Fungsi *Fitness*

Fungsi *fitness* adalah fungsi objektif yang digunakan untuk mengukur seberapa baik suatu kromosom sebagai calon solusi. Sebuah fungsi *fitness* harus bisa memperkirakan seberapa dekat sebuah calon solusi dengan solusi yang optimal. Dalam algoritma genetika, fungsi *fitness* juga digunakan dalam proses seleksi (*roulette-wheel selection*) untuk menentukan seberapa baik suatu kromosom untuk menjadi induk dari generasi berikutnya.

Algorithm 2 Algoritma Genetika

function Algoritma-Genetika(*populasi*, Fung-Fitness) **returns** solusi berupa individu

inputs: *populasi*, himpunan individu
Fung-Fitness, fungsi yang mengukur *fitness* dari tiap individu

```

1: repeat
2:   populasi_baru  $\leftarrow$  himpunan kosong
3:   for  $i=1$  to Size(populasi) do
4:      $x \leftarrow$  Seleksi-acak(populasi, Fung-Fitness)
5:      $y \leftarrow$  Seleksi-acak(populasi, Fung-Fitness)
6:     anak  $\leftarrow$  Persilangan( $x, y$ )
7:     if probabilitas mutasi then
8:       Mutasi(anak)
9:     end if
10:    tambahkan anak ke populasi_baru
11:  end for
12:  populasi  $\leftarrow$  populasi_baru
13: until kriteria berakhir dipenuhi atau sampai jumlah generasi tertentu
14: return individu terbaik dalam populasi, berdasarkan Fung-Fitness

```

function Persilangan(x, y) **returns** anak berupa individu

inputs: x, y , individu induk

```

1:  $n \leftarrow$  Length( $x$ )
2:  $c \leftarrow$  angka acak antara 1 sampai  $n$ 
3: return Append(Substring( $x, 1, c$ ), Substring( $y, c+1, n$ ))

```

Proses pencarian dalam algoritma genetika

Seperti yang disebutkan dalam Algoritma 2, algoritma genetika dimulai dengan menginisialisasi suatu populasi (biasanya dibangkitkan secara acak jika tidak diketahui heuristik masalahnya). Lalu, akan dibangkitkan populasi baru untuk generasi berikutnya dengan cara mengambil dua kromosom secara acak dengan teknik *roulette-wheel selection*. Setelah itu akan dilakukan persilangan dengan

teknik *single-point crossover*. Teknik ini dilakukan dengan cara menentukan sebuah titik potong c yang diambil secara acak antara angka 1 sampai panjang kromosom n . Keturunan dari kedua induk tersebut akan memiliki kromosom induk pertama dari gen ke-1 sampai gen ke- c dan dari induk kedua mulai dari gen ke- $(c+1)$ sampai gen ke- n . Setelah itu, untuk apabila terjadi mutasi, maka salah satu gen dari anak akan diubah nilainya. Iterasi ini akan dilakukan terus-menerus hingga kriteria berhenti tercapai atau sudah mencapai batas jumlah generasi tertentu.

4. Mempelajari penggunaan algoritma genetik dalam pengelompokan dokumen.

Status : Ada sejak rencana kerja skripsi

Hasil : Meskipun pada umumnya algoritma *K-means* digunakan dalam pengelompokan, tetapi ternyata algoritma *K-means* masih memiliki kekurangan yaitu masih dapat terjebak pada *local optimum*. Oleh karena itu, pada penelitian ini digunakan algoritma genetika sebagai solusi dari permasalahan *local optimum*. Algoritma genetika memiliki kemampuan untuk menghindari terjadinya *local optimum* karena algoritma genetika merupakan algoritma pencarian yang bersifat stokastik karena memperbolehkan terjadinya variasi acak dalam proses pencarian solusinya. Oleh karena itu, diharapkan pengelompokan berbasis GA dapat menghasilkan solusi yang lebih baik dibandingkan dengan pengelompokan pada umumnya yang menggunakan algoritma *K-means*.

5. Mencari dokumen yang akan dijadikan *training* dan *test datasets*. **Status :** Ada sejak rencana kerja skripsi

Hasil : Belum dikerjakan.

6. Melakukan analisis masalah dan mencari solusinya.

Status : Baru ditambahkan semester ini

Hasil : Untuk menyelesaikan masalah pengelompokan dokumen menggunakan algoritma genetika, maka perlu dibangun sebuah model yang dapat diterapkan ke dalam algoritma tersebut.

Representasi Kromosom

Setiap *string* kromosom merupakan deretan bilangan riil yang merepresentasikan K titik pusat *cluster* (*centroid*). Dalam ruang N dimensi, panjang dari kromosom akan menjadi $N \times K$ gen. N kata pertama merepresentasikan N dimensi dari *centroid* pertama, N kata selanjutnya merepresentasikan N dimensi dari *centroid* kedua, dan seterusnya.

Fungsi *Fitness*

Perhitungan *fitness* dalam penelitian ini terdiri dari dua tahap. Pada tahap pertama, terjadi pembentukan *cluster* berdasarkan titik pusat yang terkandung dalam kromosom. Hal ini dilakukan dengan menetapkan setiap titik $x_i, i = 1, 2, \dots, n$ ke dalam sebuah *cluster* C_j dengan *centroid* z_j sehingga

$$\|x_i - z_j\| < \|x_i - z_p\|, p = 1, 2, \dots, K, \text{ dan } p \neq j. \quad (6)$$

Setelah proses pengelompokan selesai, titik pusat yang terkandung dalam kromosom diganti dengan rata-rata titik dari tiap *cluster*. Dengan kata lain, untuk *cluster* C_i , *centroid* baru z_i^* dapat dihitung dengan rumus:

$$z_i^* = \frac{1}{n_i} \sum_{x_j \in C_i} x_j, i = 1, 2, \dots, K. \quad (7)$$

dengan z_i^* merupakan titik pusat *cluster* ke- i , n_i merupakan jumlah anggota *cluster* ke- i , dan x_j merupakan titik ke- j yang merupakan anggota dari *cluster* ke- i . z_i^* ini akan menggantikan z_i sebelumnya di kromosom. Ada dua metode perhitungan fungsi fitness yang diimplementasikan dalam penelitian ini yaitu *euclidean distance* dan *cosine similarity*.

Euclidean Distance

Pada perhitungan dengan *euclidean distance*, akan dihitung sebuah *clustering metric* M dengan rumus:

$$M = \sum_{i=1}^K M_i, \quad (8)$$

$$M_i = \sum_{x_j \in C_i} \|x_j - z_i\|$$

Lalu, fungsi *fitness* akan didefinisikan sebagai $f = 1/M$, sehingga maksimalisasi terhadap nilai f akan meminimalkan nilai M .

Cosine Similarity

Perhitungan *fitness* menggunakan *cosine similarity* dapat dilakukan dengan rumus:

$$f = \sum_{i=1}^K f_i, \quad (9)$$

$$f_i = \sum_{x_j \in C_i} \frac{x_j \cdot z_i}{\|x_j\| \times \|z_i\|}$$

semakin besar nilai dari fungsi *fitness* f , maka kromosom tersebut semakin mendekati solusi yang optimal.

Operasi Genetik

Ada beberapa operasi genetik yang akan dibahas, di antaranya: inisialisasi populasi, seleksi, persilangan, dan mutasi.

Inisialisasi Populasi

K *centroid* yang terkandung dalam kromosom pada mulanya dipilih secara acak sebanyak K titik dari keseluruhan himpunan data. Lalu proses ini diulang sebanyak P kali di mana P merupakan ukuran populasi yang diinginkan.

Seleksi

Proses seleksi ini terjadi berdasarkan konsep *survival of the fittest* yang diadaptasi dari sistem genetika alami. Dalam penelitian ini, calon induk dari generasi selanjutnya dipilih dengan menggunakan teknik *roulette-wheel selection* yang merupakan salah satu teknik yang digunakan karena mengimplementasikan *proportional selection strategy*.

Persilangan

Persilangan dalam penelitian ini terjadi terhadap dua induk dengan satu titik potong. Misalkan kromosom memiliki panjang l , sebuah angka acak akan diambil sebagai titik potong dalam batas $[1, l - 1]$. Bagian kromosom sebelah kanan titik potong akan ditukar antara kedua induk sehingga menghasilkan dua individu keturunan.

Mutasi

Setiap kromosom mengalami mutasi dengan probabilitas mutasi tetap μ_c . Setelah itu, ditentukan gen mana yang akan mengalami mutasi dengan mengambilnya secara acak. Perubahan nilai gen tersebut juga merupakan angka acak yang diambil mulai dari angka 0 sampai dengan jumlah kemunculan kata keseluruhan dalam *vocabulary*.

7. Membuat rancangan perangkat lunak yang menggunakan algoritma genetik sebagai algoritma pengelompokan dokumen.

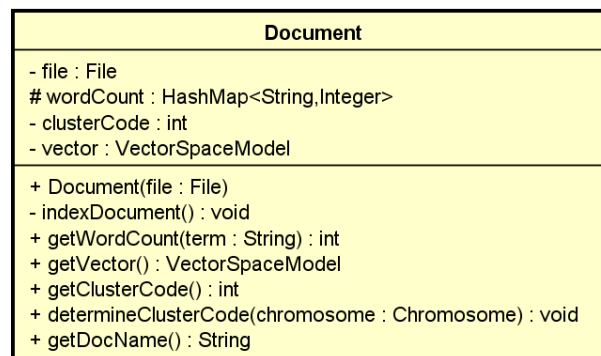
Status : Ada sejak rencana kerja skripsi

Hasil : Ada dua rancangan yang sudah dibuat yaitu rancangan kelas dan rancangan antarmuka pengguna.

Rancangan Kelas

Berdasarkan hasil analisis dari masalah yang dihadapi, dibentuklah diagram kelas pada gambar 12 sebagai gambaran dari perangkat lunak yang akan dibuat.

Document



Gambar 2: Kelas *Document*

Kelas ini merupakan representasi dari dokumen yang akan diproses dalam pengelompokan. Kelas ini berfungsi untuk menyimpan informasi yang dibutuhkan dari sebuah dokumen selama proses pengelompokan. Atribut yang dimiliki oleh kelas *Document* adalah:

- *wordCount*: bertipe *HashMap* dengan *key* bertipe *String* dan *value* bertipe *Integer*. Atribut ini menyimpan pasangan kata yang dimiliki oleh dokumen tersebut dan frekuensinya.
- *file*: atribut ini bertipe *File* milik *package java.io* yang berfungsi untuk merepresentasikan *file* dari dokumen yang akan diproses.
- *vector*: atribut bertipe *VectorSpaceModel* ini merepresentasikan model ruang vektor pada sebuah dokumen.

Method yang terdapat dalam kelas ini adalah:

- *Document*: merupakan *constructor* dengan sebuah parameter bertipe *File* yaitu *file* dari dokumen yang akan dikelompokkan.
- *indexDocument*: *method* tanpa kembalian (*void*) yang berfungsi untuk mengindeks dokumen untuk mengisi atribut *wordCount*.
- *getWordCount*: berfungsi untuk mengembalikan banyaknya istilah *term* muncul dalam dokumen.
- *getVector*: merupakan *getter* dari variabel *vector*.
- *determineClusterCode*: berfungsi untuk menentukan *cluster* dari dokumen.
- *getDocName*: mengembalikan nama *file* dari dokumen.

VectorSpaceModel

<i>VectorSpaceModel</i>
termsWeight : HashMap<String,Double> - distanceCalculator : DistanceCalculator
VectorSpaceModel(wordCount : HashMap<String,Integer>, distanceMode : String) # generateVector(wordCount : HashMap<String,Integer>) : void + getWeight(term : String) : double + calculateDistance(otherVSM : VectorSpaceModel) : double + mutate() : void + setTermsWeight(termsWeight : HashMap<String,Double>) : void

Gambar 3: Kelas *VectorSpaceModel*

Kelas ini merupakan kelas abstrak yang merepresentasikan model ruang vektor. kelas ini memiliki fungsi-fungsi yang umum dimiliki oleh sebuah model ruang vektor. Atribut yang dimiliki kelas ini adalah:

- *termsWeight*: bertipe *Hashmap* dengan *key* berupa *String* dan *value* berupa *Double*. atribut ini menyimpan pasangan istilah dan bobotnya sesuai dengan metode pembobotan.
- *distanceCalculator*: bertipe *DistanceCalculator* dan merupakan objek yang akan digunakan untuk menghitung jarak antar vektor.

Method yang terdapat dalam kelas ini adalah:

- *VectorSpaceModel*: merupakan *constructor* dengan dua parameter yaitu *wordCount* bertipe *HashMap<String,Integer>* yang merupakan pasangan kata dan banyak kemunculannya dalam dokumen serta *distanceMode* bertipe *String* yang akan menentukan tipe perhitungan jarak antar vektor.
- *generateVector*: merupakan *method* tanpa parameter untuk mengubah banyak kemunculan kata menjadi berat.
- *getWeight*: berfungsi untuk mengembalikan berat dari istilah *term*.
- *calculateDistance*: berfungsi untuk menghitung jarak antara vektor ini dengan *otherVSM* menggunakan metode yang dipilih pada parameter di *constructor*.
- *mutate*: merupakan *method* untuk melakukan mutasi pada sebuah dimensi dalam vektor.
- *setTermsWeight*: merupakan *setter* dari atribut *termsWeight*.

BagOfWordVSM
+ BagOfWordVSM(wordCount : HashMap<String,Integer>, distanceMode : String) + getWeight(term : String) : double # generateVector(wordCount : HashMap<String,Integer>) : void

Gambar 4: Kelas *BagOfWordVSM****BagOfWordVSM***

Kelas ini merupakan kelas yang mengimplementasikan kelas *VectorSpaceModel*. Kelas ini menggunakan metode pembobotan *BagOfWord* di mana bobot tiap istilah merupakan banyaknya kemunculan istilah itu sendiri. Atribut yang dimiliki kelas ini seluruhnya merupakan atribut yang diturunkan dari kelas *VectorSpaceModel* dan hanya melakukan *override* pada *method* yang masih bersifat abstrak.

DistanceCalculator

DistanceCalculator
+ calculateDistance(vsm1 : VectorSpaceModel, vsm2 : VectorSpaceModel) : double

Gambar 5: Kelas *DistanceCalculator*

Kelas ini merupakan kelas abstrak yang berfungsi untuk menghitung jarak antara dua buah objek bertipe *VectorSpaceModel*. Kelas ini tidak memiliki atribut dan hanya memiliki sebuah *method* abstrak yaitu *calculateDistance* yang memiliki dua buah parameter *vsm1* dan *vsm2*. Hasil yang dikembalikan oleh *method* ini adalah jarak dari kedua vektor tersebut sesuai dengan metode perhitungan jaraknya.

EuclideanDistanceCalculator

EuclidianDistanceCalculator
+ calculateDistance(vsm1 : VectorSpaceModel, vsm2 : VectorSpaceModel) : double

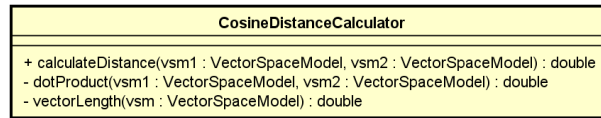
Gambar 6: Kelas *EuclideanDistanceCalculator*

Kelas ini mengimplementasikan kelas abstrak *DistanceCalculator*. Kelas ini hanya melakukan *override* pada *method* *calculateDistance* dengan melakukan perhitungan menggunakan jarak euclidean untuk menghitung jarak antara dua buah vektor (Subbab 1).

CosineDistanceCalculator

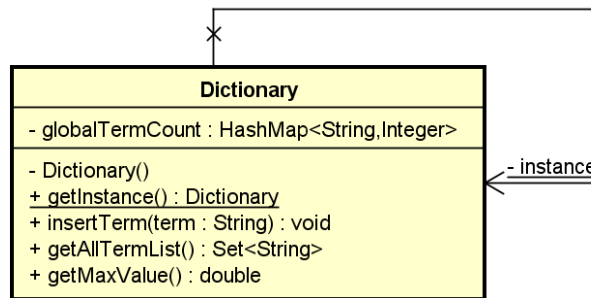
Kelas ini juga mengimplementasikan kelas abstrak *DistanceCalculator*. Kelas ini memiliki dua *method* tambahan selain melakukan *override* pada *method* *calculateDistance*. *Method* yang ada pada kelas ini adalah:

- *calculateDistance*: merupakan *method* yang diturunkan dari kelas *DistanceCalculator*. *Method* ini mengembalikan jarak dari *vsm1* dan *vsm2* yang dihitung menggunakan persamaan cosinus (Subbab 1).

Gambar 7: Kelas *CosineDistanceCalculator*

- *dotProduct*: berfungsi untuk menghitung hasil perkalian titik (*dot product*) antara *vsm1* dan *vsm 2*.
- *vectorLength*: berfungsi untuk menghitung panjang dari vektor *vsm*.

Dictionary

Gambar 8: Kelas *Dictionary*

Kelas ini merepresentasikan sebuah kamus yang menangani seluruh kebutuhan dalam proses pengelompokan yang membutuhkan akses global untuk keseluruhan koleksi dokumen. Atribut yang ada dalam kelas ini adalah:

- *globalTermCount*: bertipe *HashMap* dengan *key* bertipe *String* dan *value* bertipe *Integer*. Atribut ini berfungsi untuk menyimpan seluruh istilah yang muncul dan banyak kemunculannya dalam keseluruhan koleksi dokumen.
- *instance*: merupakan objek bertipe *Dictionary* sebagai instansiasi satu-satunya dari kelas *Dictionary* karena kelas ini bersifat *singleton*.

Method yang ada pada kelas ini adalah:

- *Dictionary*: merupakan *constructor private* untuk menjamin tidak akan ada lebih dari satu *instance* selama perangkat lunak berjalan.
- *getInstance*: merupakan *method static* yang berfungsi sebagai *getter* dari atribut *instance*.
- *insertTerm*: berfungsi untuk memasukkan istilah *term* ke dalam variabel *globalTermCount*.
- *getAllTermList*: bertugas mengembalikan daftar seluruh istilah yang pernah muncul di seluruh koleksi dokumen.
- *getValue*: bertugas mengembalikan banyaknya kata *term* muncul dalam seluruh koleksi dokumen.

Gene

Kelas ini merepresentasikan gen dalam algoritma genetika. Kelas ini hanya memiliki sebuah atribut *value* bertipe *VectorSpaceModel*. Atribut ini menyimpan model ruang vektor yang menjadi titik pusat *cluster* (*centroid*). *Method* yang ada pada kelas ini adalah:

Gene
- value : VectorSpaceModel
+ Gene(value : VectorSpaceModel) + getValue() : VectorSpaceModel + mutate() : void + setTermsWeight(termsWeight : HashMap<String,Double>) : void

Gambar 9: Kelas *Gene*

- *Gene*: merupakan *constructor* dari kelas *Gene* yang membutuhkan sebuah parameter bertipe *VectorSpaceModel* untuk mengisi variabel *value*.
- *getValue*: merupakan *getter* dari atribut *value*.
- *mutate*: berfungsi untuk melakukan mutasi pada gen. *Method* ini sebenarnya hanya bertugas memanggil fungsi *mutate()* dari atribut *value*.
- *setTermsWeight*: berfungsi untuk mengubah nilai atribut *termsWeight* milik atribut *value*.

Chromosome

Chromosome
+ <u>MUTATION_PROBABILITY : double = .20</u> - rand : Random - genes : List<Gene>
+ Chromosome() + addGene(gene : Gene) : void + mutate() : void + crossover(otherChromosome : Chromosome) : Chromosome + computeFitness() : double + getAllGenes() : List<Gene>

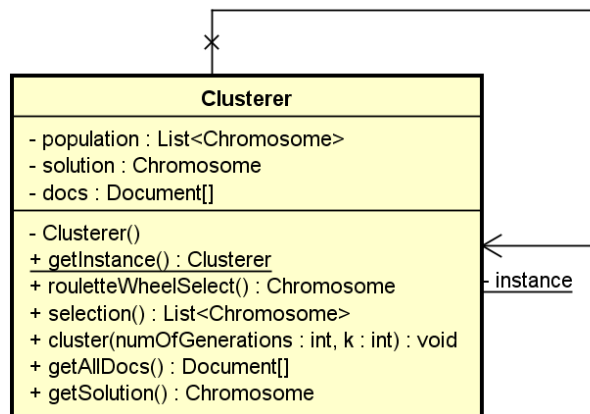
Gambar 10: Kelas *Chromosome*

Kelas ini merepresentasikan kromosom dalam algoritma genetika (Subbab 3). Atribut yang terdapat dalam kelas ini adalah:

- *genes*: bertipe *List of Gene* dan merupakan kumpulan gen yang terdapat dalam kromosom.
- *MUTATION_PROBABILITY*: merupakan atribut yang bersifat *static* dan *final* yang berisi probabilitas terjadinya mutasi dalam proses pembangkitan keturunan.

Method yang terdapat dalam kelas ini adalah:

- *Chromosome*: merupakan *constructor* tanpa parameter untuk membentuk objek dari kelas *Chromosome*.
- *addGene*: bertugas untuk menambahkan satu gen ke dalam kromosom (ke dalam atribut *genes*).
- *mutate*: berfungsi untuk melakukan mutasi pada kromosom dengan cara melakukan mutasi pada sebuah gen secara acak (Subbab 3).
- *crossover*: bertugas untuk melakukan persilangan dengan kromosom lain untuk menghasilkan keturunan (Subbab 3).
- *computeFitness*: mengembalikan nilai *fitness* dari kromosom (Subbab 3).
- *getAllGenes*: merupakan *getter* dari atribut *genes*.

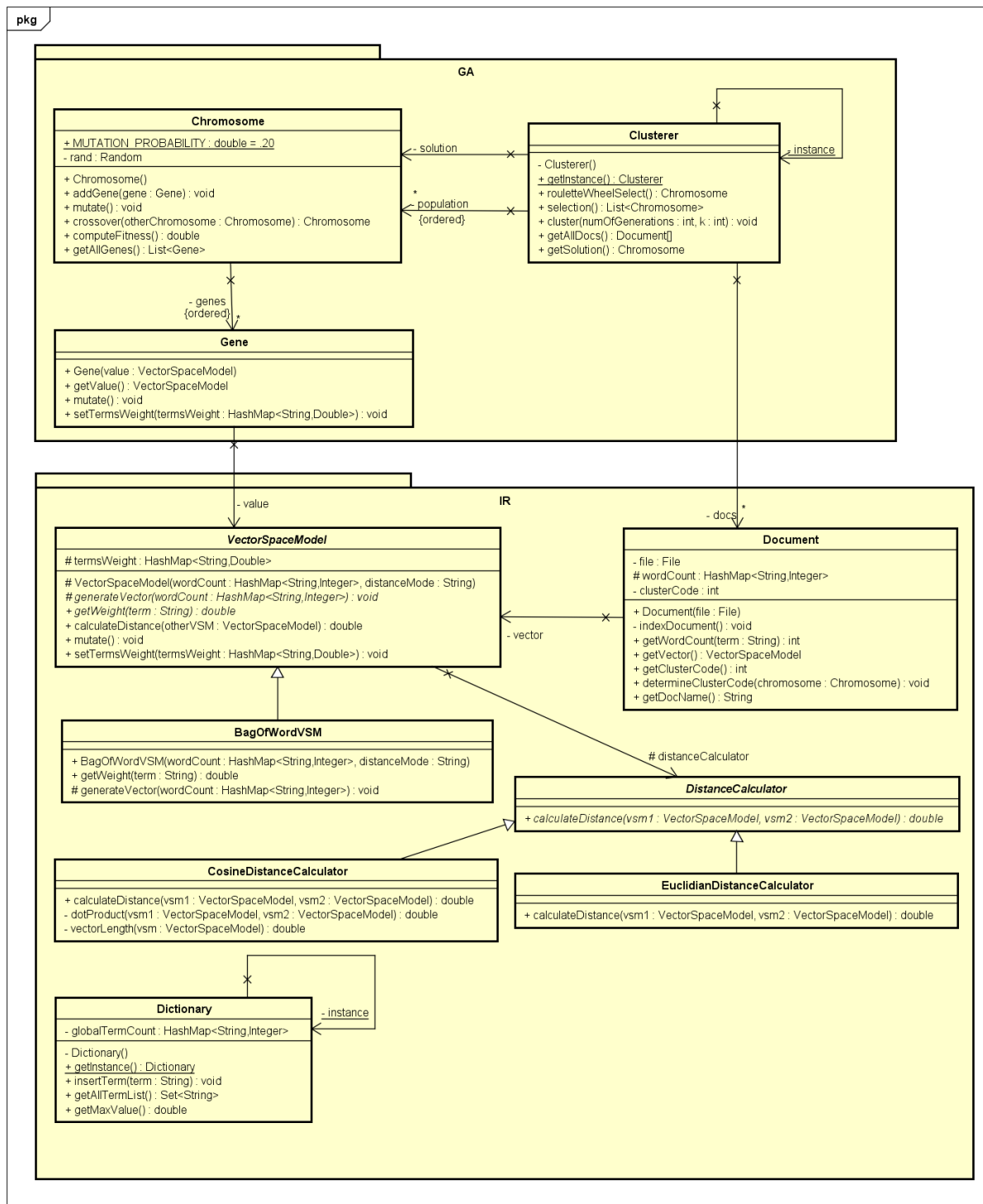
ClustererGambar 11: Kelas *Clusterer*

Kelas ini merupakan kelas utama yang akan mengatur jalannya proses pengelompokan. Kelas ini merupakan kelas singleton. Atribut yang terdapat dalam kelas ini adalah:

- *docs*: bertipe *List of Document* yang berfungsi untuk menyimpan seluruh koleksi dokumen.
- *instance*: variabel *static* ini berfungsi untuk menyimpan *instance* dari kelas *Clusterer*.
- *population*: bertipe *List of Chromosome* yang merepresentasikan populasi pada generasi saat ini.
- *solution*: bertipe *Chromosome* yang mencatat kromosom dengan nilai *fitness* terbaik.

Method yang terdapat dalam kelas ini adalah:

- *Clusterer*: merupakan *constructor private* yang berfungsi untuk menjamin tidak akan ada *instance* dibuat diluar dari kelas ini.
- *getInstance*: merupakan *getter* dari variabel *instance*.
- *rouletteWheelSelect*: bertugas untuk memilih dua kromosom dan melakukan persilangan untuk menghasilkan sebuah keturunan selanjutnya.
- *selection*: bertugas untuk melakukan *roulette wheel selection* sebanyak populasi untuk menghasilkan populasi dari generasi selanjutnya.
- *cluster*: merupakan *method* utama yang bertugas melakukan pengelompokan dokumen dengan dua parameter yaitu jumlah generasi dan nilai *k*.
- *getAllDocs*: merupakan *getter* dari atribut *docs*.
- *getSolution*: merupakan *getter* dari atribut *solution*.



Gambar 12: Diagram kelas pada tahap perancangan

Perancangan Antarmuka Pengguna

Antarmuka yang dirancang untuk perangkat lunak ini hanya terdiri dari satu jendela utama dan dua jendela *pop-up*. Pada penelitian ini, perancangan antarmuka dibuat menggunakan perangkat lunak *balsamiq*¹. Setiap objek dan *field* akan diberi label unik agar dapat disesuaikan dengan tabel keterangan. Berikut akan dibahas rancangan antarmuka pengguna dari perangkat ini.

Jendela Utama

Gambar 13 merupakan jendela yang pertama kali ditampilkan saat perangkat lunak dijalankan. Jendela tersebut berisi berbagai macam hal yang dibutuhkan pengguna dalam melakukan pengelompokan dokumen.

Gambar 13: Rancangan antarmuka jendela utama

Penjelasan setiap *field* dalam jendela utama adalah sebagai berikut:

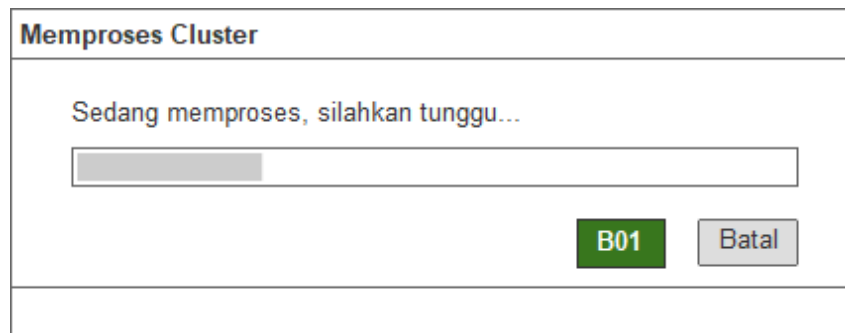
Kode	Nama	Jenis	Default value	Wajib	Aturan validasi
A01	Direktori sumber	<i>file chooser</i>	-	ya	Harus berupa direktori yang berisi dokumen (tidak boleh kosong)
A02	Direktori tujuan	<i>file chooser</i>	-	ya	Harus berupa direktori yang kosong
A03	Jumlah cluster	<i>spinner</i>	2	ya	Nilai minimum 2
A04	Jumlah generasi	<i>spinner</i>	1	ya	Nilai minimum 1
A05	Metode perhitungan jarak	<i>dropdown</i>	<i>Cosine Similarity</i>	ya	-

Tabel 2: Rincian *field* pada jendela utama

Jendela ini hanya memiliki sebuah tombol dengan kode **A06** yang berfungsi untuk memulai proses pengelompokan dan membuka jendela *loading*. Selain tombol, pada jendela ini juga terdapat lima *field* seperti yang tertera pada Tabel 7.

¹<https://balsamiq.com/>

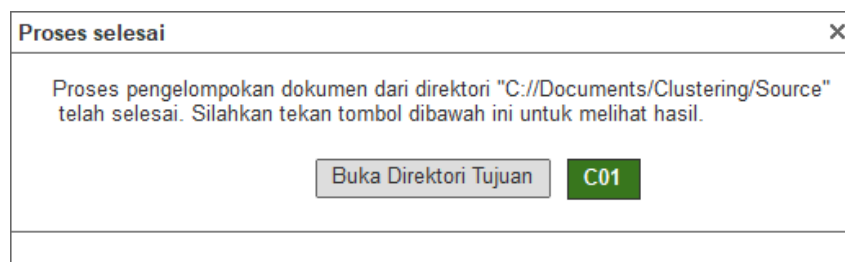
Jendela *Loading*



Gambar 14: Rancangan antarmuka jendela *loading*

Gambar 14 merupakan jendela yang akan muncul setelah pengguna menekan tombol "Mulai Pengelompokan". Jendela ini berisi informasi perkembangan proses pengelompokan. Informasi ini disajikan dalam bentuk *progress bar*. Hanya ada sebuah tombol pada jendela ini yaitu tombol dengan kode **B01**. Sesuai dengan labelnya, tombol ini berfungsi untuk membatalkan proses pengelompokan. Apabila tombol batal ditekan, maka pengguna akan dikembalikan ke jendela utama (Gambar 13).

Jendela Proses Berhasil



Gambar 15: Rancangan antarmuka jendela proses berhasil

Gambar 15 akan ditampilkan setelah proses pengelompokan berhasil, yaitu saat *progress bar* di jendela *loading* sudah terisi penuh. Jendela ini hanya memiliki satu buah tombol yaitu tombol dengan kode **C01** yang berfungsi untuk membuka *Windows Explorer* pada direktori hasil yang sudah dipilih pengguna pada jendela utama untuk menampilkan hasil dari proses pengelompokan.

8. Mengimplementasikan hasil rancangan menjadi perangkat lunak.

Status : Ada sejak rencana kerja skripsi

Hasil : Hasil implementasi dilampirkan.

6 Pencapaian Rencana Kerja

Langkah-langkah kerja yang berhasil diselesaikan dalam Skripsi 1 ini adalah sebagai berikut:

1. Melakukan studi literatur mengenai *Information Retrieval* (temu kembali informasi).
2. Melakukan studi literatur mengenai *Document Clustering* (pengelompokan dokumen).

3. Melakukan studi literatur mengenai *Genetic Algorithm* (algoritma genetik).
4. Mempelajari penggunaan algoritma genetik dalam pengelompokan dokumen.
5. Melakukan analisis masalah dan mencari solusinya.
6. Membuat rancangan perangkat lunak yang menggunakan algoritma genetik sebagai algoritma pengelompokan dokumen.
7. Mengimplementasikan hasil rancangan menjadi perangkat lunak.

7 Kendala yang Dihadapi

Kendala - kendala yang dihadapi selama mengerjakan skripsi :

- Skripsi 1 diambil bersamaan dengan Skripsi 2 sehingga waktu yang tersedia sangat terbatas sedangkan banyak hal yang perlu dikerjakan.
- Ada beberapa tugas besar dari mata kuliah lain yang menyita cukup banyak waktu sehingga mengurangi waktu untuk mengerjakan skripsi.

Bandung, 15/11/2018

Cornelius David Herianto

Menyetujui,

Nama: Thomas Anung Basuki
Pembimbing Tunggal

Lampiran Kode Program

Listing 1: Chromosome.java

```

1 package GA;
2
3
4 import IR.Dictionary;
5 import IR.Document;
6 import java.util.Arrays;
7 import java.util.HashMap;
8 import java.util.LinkedList;
9 import java.util.List;
10 import java.util.Random;
11
12 /*
13  * To change this license header, choose License Headers in Project Properties.
14  * To change this template file, choose Tools | Templates
15  * and open the template in the editor.
16  */
17
18 /**
19  *
20  * @author CorneliusDavid
21  */
22 public class Chromosome {
23     public static final double MUTATION_PROBABILITY=.20;
24     private List<Gene> genes;
25     private Random rand;
26
27     public Chromosome() {
28         this.genes = new LinkedList<>();
29         this.rand = new Random();
30         // this.fitnessValue=-1;
31     }
32
33     public void addGene(Gene gene){
34         this.genes.add(gene);
35     }
36
37     public void mutate(){
38         if(rand.nextDouble()<MUTATION_PROBABILITY){
39             int mutatedGeneIdx=rand.nextInt(genes.size());
40             genes.get(mutatedGeneIdx).mutate();
41         }
42     }
43
44     public Chromosome crossover(Chromosome otherChromosome){
45         int breakPoint=rand.nextInt(genes.size());
46         Chromosome result=new Chromosome();
47         for (int i = 0; i < breakPoint; i++) {
48             result.addGene(this.genes.get(i));
49         }
50
51         for (int i = breakPoint; i < this.genes.size(); i++) {
52             result.addGene(otherChromosome.genes.get(i));
53         }
54
55         return result;
56     }
57
58     public double computeFitness(){
59         // if(fitnessValue!=-1)return fitnessValue;
60         int fitnessValue=0;
61         Document docs[]=Clusterer.getInstance().getAllDocs();
62         HashMap<String, Double> sumOfCentroid[]=new HashMap[genes.size()];
63         for (int i = 0; i < genes.size(); i++) {
64             sumOfCentroid[i]=new HashMap<>();
65         }
66         int pointCount[]=new int[genes.size()];
67         for (int i = 0; i < docs.length; i++) {
68             docs[i].determineClusterCode(this);
69             int clusterCode=docs[i].getClusterCode();
70             double tmp=docs[i].getVector().calculateDistance(this.genes.get(clusterCode).getValue());
71             fitnessValue+=tmp;
72             System.out.println(tmp);
73             for (String term:Dictionary.getInstance().getAllTermList()){
74                 double nextValue=sumOfCentroid[clusterCode].containsKey(term)?sumOfCentroid[clusterCode].get(term):0;
75                 nextValue+=tmp;
76                 sumOfCentroid[clusterCode].put(term, nextValue);
77             }
78             pointCount[clusterCode]++;
79         }
80
81         for (int i = 0; i < genes.size(); i++) {
82             for (String term:Dictionary.getInstance().getAllTermList()){
83                 double nextValue=sumOfCentroid[i].containsKey(term)?sumOfCentroid[i].get(term):0;
84                 nextValue=pointCount[i]==0.0?0.0:nextValue/pointCount[i];
85                 sumOfCentroid[i].put(term, nextValue);
86             }
87             System.out.println(sumOfCentroid[i]);
88             genes.get(i).setTermsWeight(sumOfCentroid[i]);
89         }

```

```

90         return fitnessValue;
91     }
92
93     public List<Gene> getAllGenes() {
94         return genes;
95     }
96
97
98 }

```

Listing 2: Clusterer.java

```

1 package GA;
2
3 import IR.Document;
4 import java.io.File;
5 import java.util.Arrays;
6 import java.util.Collections;
7 import java.util.LinkedList;
8 import java.util.List;
9 import java.util.Random;
10
11 /*
12  * To change this license header, choose License Headers in Project Properties.
13  * To change this template file, choose Tools | Templates
14  * and open the template in the editor.
15  */
16 /**
17  *
18  * @author CorneliusDavid
19  */
20 public class Clusterer {
21
22     private List<Chromosome> population;
23     private Chromosome solution;
24     private Document docs[];
25     private static Clusterer instance;
26
27     private Clusterer() {
28         population = new LinkedList<>();
29     }
30
31     public static Clusterer getInstance() {
32         if (instance == null) {
33             instance = new Clusterer();
34         }
35         return instance;
36     }
37
38     public Chromosome rouletteWheelSelect() {
39         double totalWeight = 0;
40         for (int i = 0; i < population.size(); i++) {
41             totalWeight += population.get(i).computeFitness();
42         }
43
44         double selectedValue = (new Random()).nextDouble() * totalWeight;
45
46         for (int i = 0; i < population.size(); i++) {
47             selectedValue -= population.get(i).computeFitness();
48             if (selectedValue <= 0) {
49                 return population.get(i);
50             }
51         }
52         return population.get(population.size() - 1);
53     }
54
55     public List<Chromosome> selection() {
56         List<Chromosome> nextGen = new LinkedList<>();
57         for (int i = 0; i < population.size(); i++) {
58             Chromosome temp = this.rouletteWheelSelect();
59             nextGen.add(temp);
60         }
61
62         return nextGen;
63     }
64
65     public void cluster(int numOfGenerations, int k) {
66         int popSize = 10;
67
68         //document indexing
69         File folder = new File("testdoc");
70         File[] files = folder.listFiles();
71         docs = new Document[files.length];
72         for (int i = 0; i < docs.length; i++) {
73             docs[i] = new Document(files[i]);
74         }
75
76         List<Document> list = Arrays.asList(docs);
77
78         //generate initial population
79         for (int i = 0; i < popSize; i++) {
80             Chromosome temp = new Chromosome();
81             Collections.shuffle(list);
82             for (int j = 0; j < k; j++) {

```

```

83         temp.addGene(new Gene(list.get(j).getVector()));
84     }
85     population.add(temp);
86 }
87
88
89     for (int i = 0; i < numOfGenerations; i++) {
90         //selection
91         population = selection();
92         for (int j = 0; j < popSize; j++) {
93             Chromosome temp=population.get(j);
94             if (solution == null || solution.computeFitness()< temp.computeFitness()) {
95                 solution = temp;
96             }
97         }
98         // System.out.println(solution.computeFitness());
99         //mutation
100         for (Chromosome c : population) {
101             c.mutate();
102         }
103     }
104 }
105
106 public Document[] getAllDocs() {
107     return docs;
108 }
109
110 public Chromosome getSolution() {
111     return solution;
112 }
113 }

```

Listing 3: Gene.java

```

1 package GA;
2
3 import IR.VectorSpaceModel;
4 import java.util.HashMap;
5
6 /*
7  * To change this license header, choose License Headers in Project Properties.
8  * To change this template file, choose Tools | Templates
9  * and open the template in the editor.
10 */
11
12 /**
13  *
14  * @author CorneliusDavid
15  */
16 public class Gene {
17     private VectorSpaceModel value;
18
19     public Gene(VectorSpaceModel value) {
20         this.value = value;
21     }
22
23     public VectorSpaceModel getValue() {
24         return value;
25     }
26
27     public void mutate(){
28         value.mutate();
29     }
30
31     @Override
32     public String toString() {
33         return value.toString();
34     }
35
36     public void setTermsWeight(HashMap<String, Double> termsWeight) {
37         this.value.setTermsWeight(termsWeight);
38     }
39 }

```

Listing 4: BagOfWordVSM.java

```

1 /*
2  * To change this license header, choose License Headers in Project Properties.
3  * To change this template file, choose Tools | Templates
4  * and open the template in the editor.
5  */
6 package IR;
7
8 import java.util.HashMap;
9 import java.util.Map.Entry;
10 import java.util.Set;
11 /**
12  *
13  * @author CorneliusDavid
14  */
15 public class BagOfWordVSM extends VectorSpaceModel {
16

```

```

17 public BagOfWordVSM(HashMap<String, Integer> wordCount, String distanceMode) {
18     super(wordCount, distanceMode);
19 }
20
21 @Override
22 public double getWeight(String term) {
23     if(!termsWeight.containsKey(term)) return 0.0;
24     return termsWeight.get(term);
25 }
26
27
28 @Override
29 protected void generateVector(HashMap<String, Integer> wordCount) {
30     Set<Entry<String,Integer>> entrySet=wordCount.entrySet();
31     for(Entry<String,Integer> entry:entrySet){
32         this.termsWeight.put(entry.getKey(), entry.getValue().doubleValue());
33     }
34 }
35
36 }

```

Listing 5: CosineDistanceCalculator.java

```

1  /*
2  * To change this license header, choose License Headers in Project Properties.
3  * To change this template file, choose Tools | Templates
4  * and open the template in the editor.
5  */
6  package IR;
7
8  import java.util.Iterator;
9  import java.util.Set;
10
11 /**
12  *
13  * @author CorneliusDavid
14  */
15 public class CosineDistanceCalculator extends DistanceCalculator{
16
17     @Override
18     public double calculateDistance(VectorSpaceModel vsm1, VectorSpaceModel vsm2) {
19         // System.out.println("l1 "+vectorLength(vsm1));
20         // System.out.println("l2 "+vectorLength(vsm2));
21         // System.out.println("dot "+dotProduct(vsm1, vsm2));
22         if(vectorLength(vsm1)==0 || vectorLength(vsm2)==0){
23             System.out.println("asdasad");
24             return 0;
25         }
26         return dotProduct(vsm1, vsm2)/(vectorLength(vsm1)*vectorLength(vsm2));
27     }
28
29     private double dotProduct(VectorSpaceModel vsm1, VectorSpaceModel vsm2){
30         Set<String> terms=Dictionary.getInstance().getAllTermList();
31         Iterator<String> it=terms.iterator();
32         double result=0;
33         while(it.hasNext()){
34             String term=it.next();
35             double weight1=vsm1.getWeight(term);
36             double weight2=vsm2.getWeight(term);
37             result+=(weight1*weight2);
38         }
39         return result;
40     }
41
42     private double vectorLength(VectorSpaceModel vsm){
43         Set<String> terms=Dictionary.getInstance().getAllTermList();
44         Iterator<String> it=terms.iterator();
45         double result=0;
46         while(it.hasNext()){
47             String term=it.next();
48             double weight=vsm.getWeight(term);
49             result+=(weight*weight);
50         }
51         return Math.sqrt(result);
52     }
53 }

```

Listing 6: Dictionary.java

```

1  /*
2  * To change this license header, choose License Headers in Project Properties.
3  * To change this template file, choose Tools | Templates
4  * and open the template in the editor.
5  */
6  package IR;
7
8  import java.io.BufferedReader;
9  import java.io.FileReader;
10 import java.util.HashMap;
11 import java.util.HashSet;
12 import java.util.Map.Entry;
13 import java.util.Set;

```

```

14
15 /**
16  *
17  * @author CorneliusDavid
18  */
19 public class Dictionary {
20
21     private HashMap<String, Integer> globalTermCount;
22     private HashSet<String> stopwords;
23     private static Dictionary instance;
24
25     private Dictionary() {
26         this.globalTermCount = new HashMap<>();
27         this.stopwords=new HashSet<>();
28         this.loadStopwords();
29     }
30
31     private void loadStopwords(){
32         try {
33             BufferedReader br=new BufferedReader(new FileReader("stopword-list.txt"));
34             String tmp=br.readLine();
35             while(tmp!=null && tmp.length()>0){
36                 this.stopwords.add(tmp);
37                 tmp=br.readLine();
38             }
39         } catch (Exception ex) {}
40     }
41
42     public static Dictionary getInstance() {
43         if (instance == null) {
44             instance = new Dictionary();
45         }
46         return instance;
47     }
48
49     /**
50     * memasukkan term ke dalam
51     *
52     * @param term String yang akan dimasukkan
53     */
54     public void insertTerm(String term) {
55         if (!globalTermCount.containsKey(term)) {
56             globalTermCount.put(term, 0);
57         }
58         globalTermCount.put(term, globalTermCount.get(term) + 1);
59     }
60
61     public Set<String> getAllTermList(){
62         return this.globalTermCount.keySet();
63     }
64
65     public double getMaxValue(){
66         double max=0;
67         String str="";
68         for(Entry<String,Integer> x:globalTermCount.entrySet()){
69             if(max<x.getValue()){
70                 max=x.getValue();
71                 str=x.getKey();
72             }
73         }
74         return max;
75     }
76
77     public boolean isStopWord(String term){
78         return stopwords.contains(term);
79     }
80 }

```

Listing 7: DistanceCalculator.java

```

1 /**
2  * To change this license header, choose License Headers in Project Properties.
3  * To change this template file, choose Tools | Templates
4  * and open the template in the editor.
5  */
6 package IR;
7
8 /**
9  *
10  * @author CorneliusDavid
11  */
12 public abstract class DistanceCalculator {
13     public abstract double calculateDistance(VectorSpaceModel vsm1,VectorSpaceModel vsm2);
14 }

```

Listing 8: Document.java

```

1 /**
2  * To change this license header, choose License Headers in Project Properties.
3  * To change this template file, choose Tools | Templates
4  * and open the template in the editor.
5  */

```

```

6 package IR;
7
8 import GA.Chromosome;
9 import java.io.File;
10 import java.io.FileInputStream;
11 import java.io.FileNotFoundException;
12 import java.io.InputStreamReader;
13 import java.util.HashMap;
14 import java.util.Scanner;
15
16 /**
17  *
18  * @author CorneliusDavid
19  */
20 public class Document {
21     private File file;
22     protected HashMap<String,Integer> wordCount;
23     private VectorSpaceModel vector;
24     private int clusterCode;
25
26     public Document(File file){
27         this.file=file;
28         wordCount=new HashMap<>();
29
30         try {
31             this.indexDocument();
32         } catch (FileNotFoundException ex) {}
33
34
35         this.vector=new BagOfWordVSM(wordCount,"COS");//TODO: ganti param
36     }
37
38     private void indexDocument() throws FileNotFoundException {
39
40
41         Scanner sc=new Scanner(new InputStreamReader(new FileInputStream(file)));
42         // StemmerIndo stemmer=new StemmerIndo();
43         while(sc.hasNext())
44         {
45             String temp=sc.next();
46             temp=temp.replaceAll("[^A-Za-z0-9]", "").toLowerCase();
47             if(temp.length()==0)continue;
48             // if(ClusteringConfig.getInstance().USE_STEMMER){
49             //     temp=stemmer.getRootWord(temp);
50             // }
51             if(ClusteringConfig.getInstance().STOPWORD_REMOVAL){
52                 if(Dictionary.getInstance().isStopWord(temp)){
53                     continue;
54                 }
55             }
56             Dictionary.getInstance().insertTerm(temp);
57             if (!wordCount.containsKey(temp)) {
58                 wordCount.put(temp, 0);
59             }
60             wordCount.put(temp, wordCount.get(temp) + 1);
61         }
62     }
63
64     public int getWordCount(String term) {
65         if(!wordCount.containsKey(term)){
66             return 0;
67         }
68         return wordCount.get(term).intValue();
69     }
70
71     @Override
72     public String toString() {
73         return vector.toString();
74     }
75
76     public VectorSpaceModel getVector() {
77         return vector;
78     }
79
80     public int getClusterCode() {
81         return clusterCode;
82     }
83
84     public void determineClusterCode(Chromosome chromosome) {
85         this.clusterCode=0;
86         double distance=Double.MAX_VALUE;
87         for (int i = 0; i < chromosome.getAllGenes().size(); i++) {
88             double tempDist=chromosome.getAllGenes().get(i).getValue().calculateDistance(vector);
89             if(tempDist<distance){
90                 distance=tempDist;
91                 clusterCode=i;
92             }
93         }
94     }
95
96     public String getDocName(){
97         return file.getName();
98     }
99 }

```


Listing 9: EuclidianDistanceCalculator.java

```

1  /*
2  * To change this license header, choose License Headers in Project Properties.
3  * To change this template file, choose Tools | Templates
4  * and open the template in the editor.
5  */
6  package IR;
7
8  import java.util.Iterator;
9  import java.util.Set;
10
11 /**
12  *
13  * @author CorneliusDavid
14  */
15 public class EuclidianDistanceCalculator extends DistanceCalculator{
16
17     @Override
18     public double calculateDistance(VectorSpaceModel vsm1, VectorSpaceModel vsm2) {
19         Set<String> terms=Dictionary.getInstance().getAllTermList();
20         Iterator<String> it=terms.iterator();
21         double result=0;
22         while(it.hasNext()){
23             String term=it.next();
24             double weight1=vsm1.getWeight(term);
25             double weight2=vsm2.getWeight(term);
26             result+=((weight1-weight2)*(weight1-weight2));
27         }
28         return Math.sqrt(result);
29     }
30 }
31 }

```

Listing 10: VectorSpaceModel.java

```

1  /*
2  * To change this license header, choose License Headers in Project Properties.
3  * To change this template file, choose Tools | Templates
4  * and open the template in the editor.
5  */
6  package IR;
7
8  import java.util.HashMap;
9  import java.util.Random;
10
11 /**
12  *
13  * @author CorneliusDavid
14  */
15 public abstract class VectorSpaceModel {
16
17     protected HashMap<String, Double> termsWeight;
18     protected DistanceCalculator distanceCalculator;
19
20     protected VectorSpaceModel(HashMap<String, Integer> wordCount, String distanceMode) {
21         termsWeight = new HashMap<>();
22
23         switch (distanceMode) {
24             case "EUC":
25                 distanceCalculator = new EuclidianDistanceCalculator();
26                 break;
27             case "COS":
28                 distanceCalculator = new CosineDistanceCalculator();
29                 break;
30         }
31
32         generateVector(wordCount);
33     }
34
35     protected abstract void generateVector(HashMap<String, Integer> wordCount);
36
37     public abstract double getWeight(String term);
38
39     public double calculateDistance(VectorSpaceModel otherVSM) {
40         return distanceCalculator.calculateDistance(this, otherVSM);
41     }
42
43     @Override
44     public String toString() {
45         return termsWeight.toString();
46     }
47
48     public void mutate(){
49         Random rand=new Random();
50         String key=(String) termsWeight.keySet().toArray()[rand.nextInt(termsWeight.size())];
51         termsWeight.put(key, rand.nextDouble()*termsWeight.get(key)*2); //mutasi dari 0 - nilai sendiri dikali 2
52     }
53
54     public void setTermsWeight(HashMap<String, Double> termsWeight) {
55         this.termsWeight = termsWeight;
56     }
57
58 }
59 }

```