







Hak Cipta Dilindungi Undang-Undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber:

- a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah.
- b. Pengutipan tidak merugikan kepentingan yang wajar UIN Suska Riau.

2. Dilarang mengemukakan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin UIN Suska Riau.

sistem dan temu kembali merupakan gabungan dari *user interface* dan *look-up-table*. Sistem temu kembali informasi didesain untuk menemukan dokumen atau informasi yang diperlukan oleh *user*.

Sistem Temu Kembali Informasi bertujuan untuk menjawab kebutuhan informasi *user* dengan sumber informasi yang tersedia dalam kondisi seperti sebagai berikut (Salton, 1989);

- a. Mempresentasikan sekumpulan ide dalam sebuah dokumen menggunakan sekumpulan konsep.
- b. Terdapat beberapa pengguna yang memerlukan ide, tapi tidak dapat mengidentifikasi dan menemukannya dengan baik.
- c. Sistem temu kembali informasi bertujuan untuk mempertemukan ide yang dikemukakan oleh penulis. Dalam dokumen dengan kebutuhan informasi pengguna yang dinyatakan dalam bentuk *keyword query*/istilah penelusuran.

Fungsi utama sistem temu kembali informasi (Salton, 1989):

- a. Mengidentifikasi sumber informasi yang relevan dengan minat masyarakat pengguna yang ditargetkan.
- b. Menganalisis isi sumber informasi (dokumen).
- c. Merepresentasikan isi sumber informasi dengan cara tertentu yang memungkinkan untuk dipertemukan dengan pertanyaan pengguna.
- d. Merepresentasikan pertanyaan (*query*) *user* dengan cara tertentu yang memungkinkan untuk dipertemukan sumber informasi yang terdapat dalam basis data.
- e. Mempertemukan pernyataan pencarian dengan data yang tersimpan dalam basis data.
- f. Menemu kembalikan informasi yang relevan.
- g. Menyempurnakan unjuk kerja sistem berdasarkan umpan balik yang diberikan oleh *user*.

Informasi atau data yang dicari dapat berupa berupa teks, *image*, *audio*, video dan lain-lain. Koleksi data teks yang dapat dijadikan sumber pencarian juga

**Hak Cipta Dilindungi Undang-Undang**

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber:

- a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah.
- b. Pengutipan tidak merugikan kepentingan yang wajar UIN Suska Riau.

2. Dilarang mengumumkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin UIN Suska Riau.

dapat berupa pesan teks, seperti *e-mail*, *fax*, dan dokumen berita, bahkan dokumen yang beredar di internet. Dengan jumlah dokumen koleksi yang besar sebagai sumber pencarian, maka dibutuhkan suatu sistem yang dapat membantu *user* menemukan dokumen yang relevan dalam waktu yang singkat dan tepat.

Dalam teknologi informasi terdapat istilah *data retrieval*, selain *information retrieval*. Dua hal ini sangatlah berbeda. *Data retrieval* secara umum menentukan dokumen yang tepat dari suatu koleksi data, yang isi dokumen tersebut mengandung *keyword* di dalam *query user*, tidak akan pernah cukup untuk memenuhi kebutuhan informasi *user*. Berbeda dengan *data retrieval*, *user* dari sistem *Information Retrieval* lebih memperhatikan dalam mendapatkan (*retrieve*) informasi melalui subyek, daripada *retrieve* data berdasarkan *query* yang diberikan, karena *user* tidak mau tahu bagaimana proses yang sedang berlangsung.

**Tabel 2.1 Perbedaan *Information Retrieval* dan *Data Retrieval***

<i>Information Retrieval</i>	<i>Data Retrieval</i>
Berhubungan dengan teks bahasa umum yang tidak selalu terstruktur dan ada kemungkinan memiliki kerancuan arti	Berhubungan dengan data, yang mana semantik strukturnya sudah terdefiniskan
Informasi yang diambil mengenai subjek atau topik	Isi dokumen/data mengandung bagian dari <i>keyword</i>
Semantik sering kali hilang	Semantik terdefinisi dengan baik
Kesalahan kecil masih bisa ditorensi	Kesalahan kecil/tunggal dari suatu obyek menunjukkan kegagalan



Model yang terdapat dalam *Information Retrieval* terbagi dalam 3 model besar, yaitu:

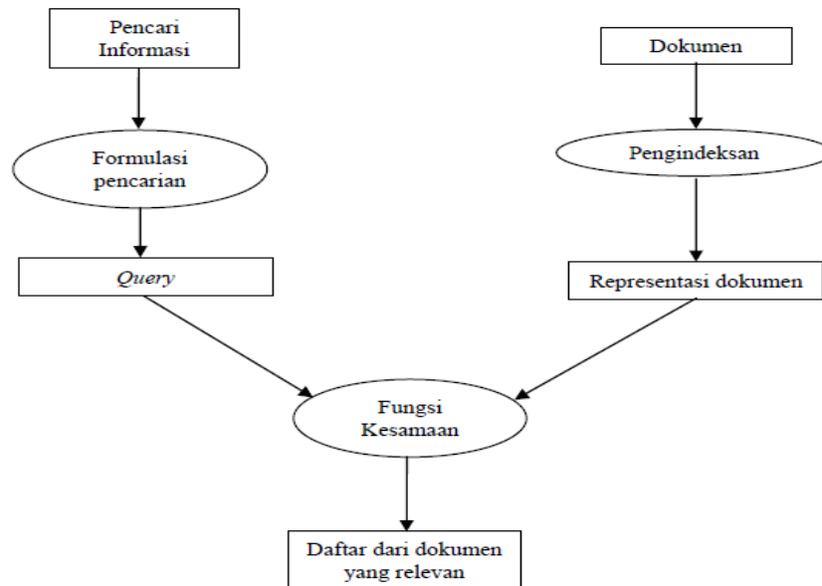
1. *Set-theoretic models*, model merepresentasikan dokumen sebagai himpunan kata atau frase. Contoh model ini ialah *standard Boolean model* dan *extended Boolean model*.
2. *Algebraic model*, model merepresentasikan dokumen dan *query* sebagai vektor atau matriks *similarity* antara vektor dokumen dan vektor *query* yang direpresentasikan sebagai sebuah nilai skalar. Contoh model ini ialah *vector space model* dan *latent semantic indexing (LSI)*.
3. *Probabilistic model*, model memperlakukan proses pengembalian dokumen sebagai sebuah *probabilistic inference*. Contoh model ini ialah penerapan teorema bayes dalam model probabilistik.

Proses dalam *Information Retrieval* dapat digambarkan sebagai sebuah proses untuk mendapatkan *relevant documents* dari *collection documents* yang ada melalui pencarian *query* yang diinputkan *user*.

### 2.2.1 Arsitektur Sistem Temu Kembali Informasi

Sistem Temu Kembali Informasi terdapat proses yang berjalan yang terbagi dalam 2 proses. Proses pertama adalah pencarian informasi oleh *user* dengan memasukkan suatu kata kunci yang nantinya akan diterjemahkan oleh sistem dalam bentuk *query*. Proses kedua adalah proses pengindeksan dokumen di dalam *database*. Pada bagian ini dokumen akan direpresentasikan ke dalam bentuk indeks yang nantinya akan dibandingkan dengan *query* dari pengguna menggunakan fungsi kesamaan untuk mendapatkan dokumen yang relevan.

Berikut adalah arsitektur sederhana dari sistem temu kembali informasi dapat dilihat pada gambar 2.1:



**Gambar 2.1** Arsitektur Sederhana *Information Retrieval* (Ingwersen, 1992)

Perlu diingat bahwa pencarian sebuah informasi di dalam sistem temukembali informasi belum tentu mengembalikan seluruh dokumen yang relevan. Bisa hanya sebagian atau tidak sama sekali. Sistem temu kembali informasi mungkin tidak memberikan hasil apapun jika memang tidak ditemukan dokumen yang relevan.

### 2.2.2 Tahapan Proses dalam Sistem Temu Kembali Informasi

Adapun beberapa tahapan proses dari Sistem temu kembali informasi dalam menemukan dokumen yang relevan terhadap *query* pengguna adalah sebagai berikut:

#### 1. Pengindeksan

Indeks adalah bahasa yang digunakan di dalam sebuah buku konvensional untuk mencari informasi berdasarkan kata atau istilah yang mengacu ke dalam suatu halaman. Dengan menggunakan indeks si pencari informasi dapat dengan mudah menemukan informasi yang diinginkannya. Pada sistem temu kembali

Hak Cipta Diindungi Undang-Undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber:

a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah.

b. Pengutipan tidak merugikan kepentingan yang wajar UIN Suska Riau.

2. Dilarang mengummumkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin UIN Suska Riau.

© Hak cipta milik UIN Suska Riau

Stie Islamic University of Sultan Syarif Kasim Riau



**Hak Cipta Dilindungi Undang-Undang**

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber:

- a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah.
- b. Pengutipan tidak merugikan kepentingan yang wajar UIN Suska Riau.

2. Dilarang mengumumkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin UIN Suska Riau.

informasi, indeks ini nantinya yang digunakan untuk merepresentasikan informasi di dalam sebuah dokumen.

Elemen dari indeks adalah istilah indeks (*index term*) yang didapatkan dari teks yang dipecah di dalam sebuah dokumen. Elemen lainnya adalah bobot istilah (*term weighting*) sebagai penentuan ranking dari kriteria relevan sebuah dokumen yang memiliki istilah yang sama.

Adapun tahap-tahap dari proses pengindeksan dokumen adalah sebagai berikut:

a. *Tokenizing*

Pada tahap ini isi dokumen akan dipecah menjadi unit-unit yang lebih kecil berupa kata, frasa atau kalimat dengan cara menghilangkan seluruh tanda baca yang ada pada dokumen sehingga menghasilkan kata yang berdiri sendiri. Unit tersebut biasanya disebut sebagai token. Sedangkan algoritma untuk memecahkan kumpulan kalimat atau frasa menjadi token disebut *tokenizer*. Token seringkali disebut sebagai istilah (*term*) atau kata, sebagai contoh sebuah token merupakan suatu urutan karakter dari dokumen tertentu yang dikelompokkan sebagai unit semantik yang berguna untuk diproses (Salton, 1989).

b. *Filtering*

Tahap *filtering* adalah tahap pengambilan kata-kata yang penting dari hasil *tokenizing*. Pada tahap ini terdapat dua proses yaitu eliminasi *stopwords* dan pengambilan *wordlist*. Eliminasi *stopwords* yaitu penyaringan (*filtering*) terhadap kata-kata yang tidak memiliki arti atau tidak layak untuk dijadikan sebagai kata kunci dalam pencarian dokumen sehingga kata-kata tersebut dapat dihilangkan dari dokumen. Sedangkan *wordlist* adalah daftar kata yang mungkin digunakan sebagai kata kunci dalam pencarian dokumen. Eliminasi *stopwords* memiliki banyak keuntungan, yaitu akan mengurangi *space* pada tabel *term index* hingga 40% atau lebih (Baeza, 1999).



#### Hak Cipta Dilindungi Undang-Undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber:
  - a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah.
  - b. Pengutipan tidak merugikan kepentingan yang wajar UIN Suska Riau.
2. Dilarang mengemukakan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin UIN Suska Riau.

**Tabel 2.4 Aturan untuk *First Order Derivational Prefix***

Akhiran	<i>Replacement</i>	<i>Additional Condition</i>
-meng	<i>Null</i>	<i>Null</i>
-meny	S	V...*
-men	<i>Null</i>	<i>Null</i>
-mem	P	V...
-mem	<i>Null</i>	<i>Null</i>
-me	<i>Null</i>	<i>Null</i>
-peng	<i>Null</i>	<i>Null</i>
-peny	S	V...
-pen	<i>Null</i>	<i>Null</i>
-pem	P	V...
-pem	<i>Null</i>	<i>Null</i>
-di	<i>Null</i>	<i>Null</i>
-ter	<i>Null</i>	<i>Null</i>
-ke	<i>Null</i>	<i>Null</i>

**Tabel 2.5 Aturan untuk *Second Order Derivational Prefix***

Akhiran	<i>Replacement</i>	<i>Additional Condition</i>
-ber	<i>Null</i>	<i>Null</i>
-bel	<i>Null</i>	<i>Ajar</i>
-be	<i>Null</i>	k*er
-per	<i>Null</i>	<i>Null</i>
-pel	<i>Null</i>	<i>Ajar</i>
-pe	<i>Null</i>	<i>Null</i>

**Tabel 2.6 Aturan untuk *Derivational Suffix***

Akhiran	<i>Replacement</i>	<i>Additional Condition</i>
---------	--------------------	-----------------------------



#### Hak Cipta Dilindungi Undang-Undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber:

- a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah.
- b. Pengutipan tidak merugikan kepentingan yang wajar UIN Suska Riau.

2. Dilarang mengumumkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin UIN Suska Riau.

-kan	<i>Null</i>	<i>Prefix</i> bukan anggota {ke, peng}
-an	<i>Null</i>	<i>Prefix</i> bukan anggota {di, meng, ter}
-i	<i>Null</i>	<i>Prefix</i> bukan anggota {ber, ke, peng}

#### d. *Inverted Index*

*Inverted index* adalah salah satu mekanisme untuk pengindeksan sebuah koleksi teks yang digunakan untuk mempercepat proses pencarian. Struktur dari *inverted index* terdiri dari dua elemen yaitu kosakata dan posisinya di dalam sebuah dokumen (Baeza-Yates dan Ribeiro-Neto, 1999). Posisi dari sebuah istilah di dalam indeks pada sebuah buku, diterjemahkan dalam bentuk nomor halaman.

Pada *inverted index*, setiap istilah di masukan ke dalam *inverted list* yang menyimpan daftar dari istilah yang menunjuk ke sejumlah dokumen yang memiliki istilah tersebut. *Inverted list* juga kadang-kadang disebut *posting list* (Witten et al, 1999).

Penggunaan *inverted index* di dalam sistem temu-kembali informasi memiliki kelemahan yaitu lambat di dalam pengindeksan, tetapi cepat di dalam proses pencarian informasi.

#### 2. *Term Weighting*

Pembobotan istilah (*term weighting*) adalah proses pembobotan pada istilah. Istilah di dalam suatu indeks harus bisa membedakan kepentingan dari sebuah dokumen pada sebuah informasi. Caranya yaitu dengan pemberian bobot kepada sebuah istilah terhadap suatu dokumen. Semakin tinggi bobot dari sebuah istilah maka semakin penting istilah tersebut dibandingkan dengan istilah lainnya di dalam sebuah dokumen. Bobot dari istilah ini dicantumkan pada *inverted index* untuk digunakan dalam proses penemu-kembali dokumen.



#### Hak Cipta Dilindungi Undang-Undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber:

- a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah.
- b. Pengutipan tidak merugikan kepentingan yang wajar UIN Suska Riau.

2. Dilarang mengumumkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin UIN Suska Riau.

Pembobotan dasar dilakukan dengan menghitung frekuensi kemunculan *term* dalam dokumen. Frekuensi kemunculan (*term frequency*) merupakan petunjuk sejauh mana *term* tersebut mewakili isi dokumen. Semakin besar kemunculan suatu *term* dalam dokumen akan memberikan nilai kesesuaian yang semakin besar.

Metode TF-IDF (*Term Frequency Inverse Document Frequency*) merupakan suatu cara untuk memberikan bobot hubungan suatu kata (*term*) terhadap dokumen (Robertson, 2005). Metode ini menggabungkan dua konsep untuk perhitungan bobot yaitu, frekuensi kemunculan sebuah kata didalam sebuah dokumen tertentu dan *inverse* frekuensi dokumen yang mengandung kata tersebut.

Terdapat beberapa cara atau metode dalam melakukan pembobotan istilah pada metode TF-IDF, yaitu melalui skema pembobotan *query* dan dokumen. Pada teknik pembobotan ini, bobot istilah telah dinormalisasi. Dalam menentukan bobot suatu istilah tidak hanya berdasarkan frekuensi kemunculan istilah di satu dokumen, tetapi juga memperhatikan frekuensi terbesar pada suatu istilah yang dimiliki oleh dokumen bersangkutan. Hal ini untuk menentukan posisi relatif bobot dari istilah dibanding dengan istilah-istilah lain didokumen yang sama. Selain itu teknik ini juga memperhitungkan jumlah dokumen yang mengandung istilah yang bersangkutan dan jumlah keseluruhan dokumen.

Hal ini berguna untuk mengetahui posisi relatif bobot istilah bersangkutan pada suatu dokumen dibandingkan dengan dokumen-dokumen lain yang memiliki istilah yang sama. Sehingga jika sebuah istilah mempunyai frekuensi kemunculan yang sama pada dua dokumen belum tentu mempunyai bobot yang sama.

### 2.3 Vektor Space Model

*Vector Space Model* (VSM) adalah metode untuk melihat tingkat kedekatan atau kesamaan (*similarity*) *term* dengan cara pembobotan *term*. Dokumen dipandang sebagai sebuah vektor yang memiliki *magnitude* (jarak) dan *direction* (arah). Pada *Vector Space Model*, sebuah istilah direpresentasikan



**Hak Cipta Dilindungi Undang-Undang**

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber:

- a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah.
- b. Pengutipan tidak merugikan kepentingan yang wajar UIN Suska Riau.

2. Dilarang mengumumkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin UIN Suska Riau.

dengan sebuah dimensi dari ruang vektor. Relevansi sebuah dokumen ke sebuah *query* didasarkan pada similaritas diantara vektor dokumen dan vektor *query* (Baeza, 1999).

Pada *Information Retrieval System* terdapat beberapa metode yang digunakan dalam *Searching* salah satunya adalah dengan merepresentasikan proses *Searching* menggunakan Model Ruang Vektor. Model Ruang Vektor dibuat berdasarkan pemikiran bahwa isi dari dokumen ditentukan oleh kata-kata yang digunakan dalam dokumen tersebut. Model ini menentukan kemiripan (*similarity*) antara dokumen dengan *query* dengan cara merepresentasikan dokumen dan *query* masing-masing ke dalam bentuk vektor. Tiap kata yang ditemukan pada dokumen dan *query* diberi bobot dan disimpan sebagai salah satu elemen vektor.

Bobot istilah yang akhirnya digunakan untuk menghitung tingkat kesamaan antara setiap dokumen yang tersimpan dalam sistem dan permintaan *user*. Dokumen yang terambil disortir dalam urutan yang memiliki kemiripan, model vektor dan memperhitungkan pertimbangan dokumen yang relevan dengan permintaan *user*. Hasilnya adalah himpunan dokumen yang terambil jauh lebih akurat (dalam arti sesuai dengan informasi yang dibutuhkan oleh *user*).

Pada sistem IR, kemiripan antar dokumen didefinisikan berdasarkan representasi *bag of words* dan dikonversi ke suatu model ruang vektor (*vector space model*, VSM). Model ini diperkenalkan oleh (Salton, 1983) dan telah digunakan secara luas. Dimensi sesuai dengan jumlah *term* dalam dokumen yang terlibat. Pada model ini:

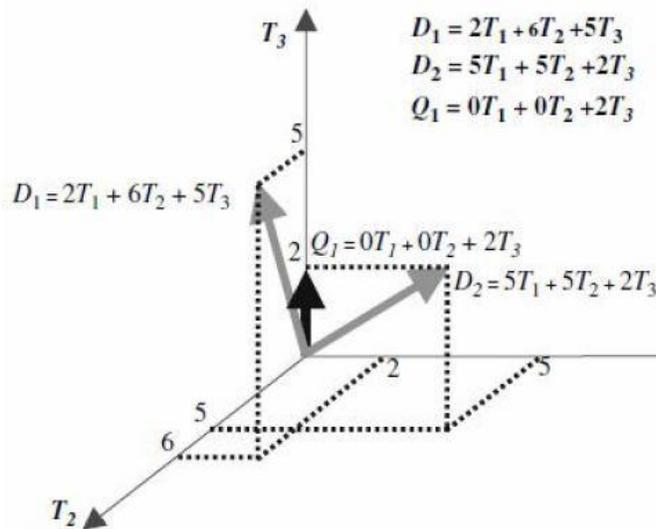
- a. *Vocabulary* merupakan kumpulan semua term berbeda yang tersisa dari dokumen setelah *preprocessing* dan mengandung *t term index*. Term-term ini membentuk suatu ruang vektor.
- b. Setiap *term*  $i$  di dalam dokumen atau *query*  $j$ , diberikan suatu bobot (*weight*) bernilai *real*  $w_{ij}$ .

## Hak Cipta Dilindungi Undang-Undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber:
  - a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah.
  - b. Pengutipan tidak merugikan kepentingan yang wajar UIN Suska Riau.
2. Dilarang mengemukakan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin UIN Suska Riau.

c. Dokumen dan *query* diekspresikan sebagai vektor  $t$  dimensi  $d_j = (w_1, w_2, \dots, w_t)$  dan terdapat  $n$  dokumen di dalam koleksi, yaitu  $j = 1, 2, \dots, n$ .

Contoh dari model ruang vektor tiga dimensi untuk dua dokumen  $D_1$  dan  $D_2$ , satu *query* pengguna  $Q_1$ , dan tiga term  $T_1$ ,  $T_2$  dan  $T_3$  diperlihatkan pada Gambar 2.2 berikut ini :



**Gambar 2.2 Representasi Dokumen dan Query pada Model Ruang Vektor**

Dalam model ruang vektor, koleksi dokumen direpresentasikan oleh matriks term document (atau matriks *term-frequency*). Setiap sel dalam matriks bersesuaian dengan bobot yang diberikan dari suatu *term* dalam dokumen yang ditentukan. Nilai nol berarti bahwa *term* tersebut tidak hadir di dalam dokumen.

Keberhasilan dari model VSM ini ditentukan oleh skema pembobotan terhadap suatu term baik untuk cakupan local maupun global, dan faktor normalisasi. Pembobotan lokal hanya berpedoman pada frekuensi munculnya *term* dalam suatu dokumen dan tidak melihat frekuensi kemunculan *term* tersebut di dalam dokumen lainnya. Pendekatan dalam pembobotan lokal yang paling banyak diterapkan adalah *term frequency* (tf) meskipun terdapat skema lain



Hak Cipta Dilindungi Undang-Undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber:

- a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah.
- b. Pengutipan tidak merugikan kepentingan yang wajar UIN Suska Riau.

2. Dilarang mengemukakan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin UIN Suska Riau.

seperti pembobotan biner, *augmented normalized tf*, logaritmik *tf* dan logaritmik alternatif.

Pembobotan global digunakan untuk memberikan tekanan terhadap *term* yang mengakibatkan perbedaan dan berdasarkan pada penyebaran dari *term* tertentu di seluruh dokumen. Banyak skema didasarkan pada pertimbangan bahwa semakin jarang suatu *term* muncul di dalam total koleksi maka *term* tersebut menjadi semakin berbeda. Pemanfaatan pembobotan ini dapat menghilangkan kebutuhan *stopword removal* karena *stopword* mempunyai bobot global yang sangat kecil.

Namun pada prakteknya lebih baik menghilangkan *stopword* di dalam fase *pre processing* sehingga semakin sedikit *term* yang harus ditangani. Pendekatan terhadap pembobotan global mencakup *inverse document frequency (idf)*, *squared idf*, *probabilistic idf*, TF-IDF, *entropy*. Pendekatan IDF merupakan pembobotan yang paling banyak digunakan saat ini.

Faktor normalisasi digunakan untuk menormalkan vektor dokumen sehingga proses *retrieval* tidak terpengaruh oleh panjang dari dokumen. Normalisasi ini diperlukan karena dokumen panjang biasanya mengandung perulangan *term* yang sama sehingga menaikkan frekuensi *term* (*tf*). Dokumen panjang juga mengandung banyak *term* yang berbeda sehingga menaikkan ukuran kemiripan antara *query* dengan dokumen tersebut, meningkatkan peluang di-retrievenya dokumen yang lebih panjang. Beberapa pendekatan normalisasi adalah normalisasi cosinus, penjumlahan bobot, normalisasi ke-4, normalisasi bobot maksimal dan normalisasi *pivoted unique*.

Proses perhitungan VSM melalui tahapan pembobotan kata *perhitungan term frequency (tf)* menggunakan persamaan (1):

$$tf_{ij} = \frac{f_{ij}}{\max_i(f_{ij})} \quad (2.1)$$



Dimana  $f_{ij}$  adalah jumlah berapa kali *term*  $i$  muncul di dalam dokumen  $j$ . Frekuensi tersebut dinormalisasi dengan frekuensi dari *most common term* di dalam dokumen tersebut. Bobot global dari suatu *term*  $i$  pada pendekatan *inverse document frequency* ( $idf_i$ ) dapat didefinisikan sebagai,

$$idf_i = \log \frac{N}{df_i} \quad (2.2)$$

Dengan  $idf_i$  adalah *inverse document frequency*,  $N$  adalah jumlah dokumen yang diambil oleh sistem, dan  $df_i$  adalah banyaknya dokumen dalam koleksi dimana *term*  $t_i$  muncul di dalamnya, maka Perhitungan  $idf_i$  digunakan untuk mengetahui banyaknya *term* yang dicari ( $df_i$ ) yang muncul dalam dokumen lain yang ada pada database (korpus).

Perhitungan *Term Frequency Inverse Document Frequency* (TFIDF), menggunakan persamaan (3):

$$W_{ij} = tf_{ij} \times idf_i \quad (2.3)$$

Dengan  $W_{ij}$  adalah bobot dokumen,  $N$  adalah Jumlah dokumen yang diambil oleh sistem,  $tf_{ij}$  adalah banyaknya kemunculan *term*  $t_i$  pada dokumen  $d_j$ , dan  $df_i$  adalah banyaknya dokumen dalam koleksi dimana *term*  $t_i$  muncul di dalamnya. Bobot dokumen ( $W_{ij}$ ) dihitung untuk didapatkannya suatu bobot hasil perkalian atau kombinasi antara *term frequency* ( $tf_i, j$ ) dan *Inverse Document Frequency* ( $idf_i$ ).

Pengukuran *Cosine Similarity* (menghitung nilai kosinus sudut antara dua vector) menggunakan persamaan (7):

$$sim(q, dj) = \frac{\sum_{i=1}^t (W_{iq} \cdot W_{ij})}{\sqrt{\sum_{j=1}^t (W_{iq})^2} \cdot \sqrt{\sum_{i=1}^t (W_{ij})^2}} \quad (2.4)$$





## Hak Cipta Dilindungi Undang-Undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber:

- a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah.
- b. Pengutipan tidak merugikan kepentingan yang wajar UIN Suska Riau.

2. Dilarang mengumumkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin UIN Suska Riau.

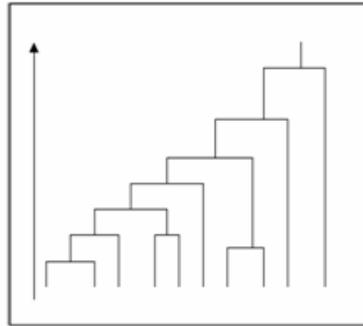
*Hierarchical Agglomerative Clustering* adalah salah satu algoritma klustering yang dapat digunakan untuk mengkluster dokumen (*document clustering*). Dari teknik *Hierarchical Agglomerative Clustering*, dapat dihasilkan suatu kumpulan partisi yang berurutan, dimana dalam kumpulan tersebut terdapat:

1. Klaster-klaster yang mempunyai poin-poin individu. Klaster-klaster ini berada di level yang paling bawah.
2. Sebuah klaster yang didalamnya terdapat poin-poin yang dipunyai semua klaster didalamnya. Single klaster ini berada di level yang paling atas.

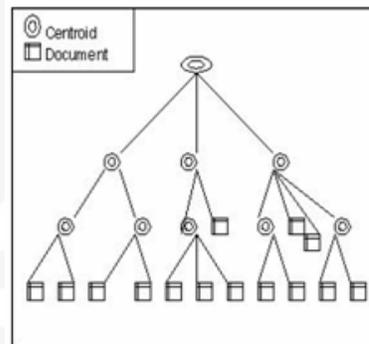
Pembentukan klaster dokumen dalam sistem temu kembali informasi dengan metode hirarkhis adalah sebagai berikut:

1. Mengidentifikasi dua dokumen yang paling mirip dan menggabungkannya menjadi sebuah klaster.
2. Mengidentifikasi dan menggabungkan dua dokumen yang paling mirip berikutnya menjadi sebuah klaster sampai semua dokumen tergabung dalam klaster-klaster yang terbentuk.
3. Proses penelusuran dokumen dilakukan dengan cara mencocokkan *query* dengan *centroid*. *Centroid* merupakan dokumen *parent* pada masing-masing klaster dokumen. Berikutnya dokumen yang berada dalam satu klaster dengan *centroid* akan ditampilkan sebagai hasil *query*.

Hasil keseluruhan dari algoritma *Hierarchical Clustering* secara grafik dapat digambarkan sebagai *tree*, yang disebut dengan dendogram. *Tree* ini secara grafik menggambarkan proses penggabungan dari klaster-klaster yang ada, sehingga menghasilkan klaster dengan level yang lebih tinggi. Cabang-cabang dalam pohon menyajikan *cluster*. Kemudian cabang-cabang bergabung pada *node* yang posisinya sepanjang sumbu jarak (similaritas) menyatakan tingkat di mana penggabungan terjadi. Gambar 2.3 dan 2.4 memperlihatkan struktur dendogram dan diagram pohon untuk klustering hirarkhis.



**Gambar 2.3 Struktur Dendrogram dan Diagram Pohon Untuk Klastering Hirarkhis**



**Gambar 2.4 Dendrogram dan Struktur Pohon dari *Hierarchical Clustering* (Salton, 1989)**

Kemiripan antar dokumen ditentukan dengan mengukur jarak antar dokumen. Dua dokumen yang mempunyai jarak paling kecil dikatakan mempunyai kemiripan paling tinggi dan dikelompokkan ke dalam satu klaster yang sama. Sebaliknya dua dokumen yang mempunyai jarak paling besar dikatakan mempunyai kemiripan paling rendah, dan dimasukkan ke dalam klaster yang berbeda.

#### 2.4.1 Metode *Hierarchical Agglomerative Clustering*

Metode *Hierarchical Agglomerative Clustering* adalah metode yang menggunakan strategi desain *Bottom-Up* yang dimulai dengan meletakkan setiap obyek sebagai sebuah klaster tersendiri (atomic klaster) dan selanjutnya



menggabungkan *atomic* kluster-*atomic* kluster tersebut menjadi kluster yang lebih besar dan lebih besar lagi sampai akhirnya semua obyek menyatu dalam sebuah kluster atau proses dapat pula berhenti jika telah mencapai batasan kondisi tertentu. Langkah-langkah dalam algoritma *Hierarchical Agglomerative Clustering*:

1. Mulai dengan N kluster, setiap kluster mengandung entiti tunggal dan sebuah matriks simetrik dari jarak (similarities)  $D = \{d_{ik}\}$  dengan tipe matrik adalah  $N \times N$ .
2. Cari matriks jarak untuk pasangan kluster yang terdekat (paling mirip), yaitu dengan mencari similaritas terbesar. Misalkan jarak antara kluster U dan V yang paling mirip adalah  $d_{uv}$ .
3. Gabungkan kluster U dan V. Label kluster yang baru dibentuk dengan (UV). Update entries pada matrik jarak dengan cara :
  - a. Hapus baris dan kolom yang bersesuaian dengan kluster U dan V
  - b. Tambahkan baris dan kolom yang memberikan jarak-jarak antara kluster (UV) dan kluster-kluster yang tersisa

#### 2.4.2 Metode *Single Linkage*

Input untuk algoritma single linkage bisa berwujud jarak atau similarities antara pasangan-pasangan dari objek-objek. Kelompok-kelompok dibentuk dari entities individu dengan menggabungkan jarak paling pendek atau similarities (kemiripan) yang paling besar. Pada awalnya, kita harus menemukan jarak terpendek dalam  $D = \{d_{uv}\}$  dan menggabungkan objek-objek yang bersesuaian misalnya U dan V, untuk mendapatkan (UV). Untuk langkah (3) dari algoritma di atas jarak-jarak antara (UV) dan *cluster* W yang lain dihitung dengan cara,

$$d_{(UV)W} = \min\{d_{UW}, d_{VW}\} \quad (2.5)$$

Di sini besaran-besaran  $d_{UW}$  dan  $d_{VW}$  berturut-turut adalah jarak terpendek antara *cluster-cluster* U dan W dan juga *cluster-cluster* V dan W.



Hak Cipta Dilindungi Undang-Undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber:

- a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah.
- b. Pengutipan tidak merugikan kepentingan yang wajar UIN Suska Riau.

2. Dilarang mengemukakan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin UIN Suska Riau.

## 2.5 Precision & Recall

Evaluasi dari sistem temu-kembali informasi dipengaruhi oleh dua parameter utama yaitu *recall* dan *precision*. *Recall* adalah rasio antara dokumen relevan yang berhasil ditemu kembalikan dari seluruh dokumen relevan yang ada di dalam sistem, sedangkan *precision* adalah rasio dokumen relevan yang berhasil ditemu kembalikan dari seluruh dokumen yang berhasil ditemu kembalikan (Grossman, 2002).

Dengan menggunakan nilai dari parameter *recall* dan *precision* akan dicari nilai dari *average precision* untuk menghitung keefektifan dan keakuratan dari algoritma sistem temu-kembali informasi. *Average precision* adalah suatu ukuran evaluasi sistem temu-kembali informasi yang diperoleh dengan cara menghitung rata-rata *precision* pada seluruh tingkat *recall* (Grossman, 2002).

Sedangkan untuk menentukan nilai dari *recall* dan *precision* harus didapatkan jumlah dokumen yang relevan terhadap suatu topik informasi. Satu-satunya cara untuk mendapatkannya yaitu dengan membaca dokumen itu satu persatu.

Menurut Rijsbergen (1979) relevansi merupakan sesuatu yang sifatnya subyektif. Setiap orang mempunyai perbedaan untuk mengartikan sesuatu dokumen tersebut relevan terhadap sebuah topik informasi.

Menurut Mizzaro (1998), evaluasi pada sebuah sistem temu-kembali informasi dengan menggunakan *recall* dan *precision* sudah cukup baik untuk menjadi ukuran dari sistem tersebut.

Adapun rumus dari *Precision* dan *Recall* yang akan digunakan akan dijelaskan pada tabel 2.7 berikut:

**Tabel 2.7 Precision dan Recall (Manning dkk, 2009)**

	<i>Relevant</i>	<i>Not Relevant</i>
<i>Retrieved</i>	<i>true positif (tp)</i>	<i>false positif (fp)</i>
<i>Not Retrieved</i>	<i>false negative (fn)</i>	<i>true negative (tn)</i>



$$\text{Maka: Precision } (P) = tp/(tp+fp) \quad (2.6)$$

$$\text{Recall } (R) = tp/(tp+fn) \quad (2.7)$$

Berdasarkan pada Tabel 2.7, ketika kondisi dokumen ditemukembali (retrieved), maka terdapat dua jenis kesesuaian yang muncul yaitu apakah dokumen tersebut relevan atau tidak relevan dengan *query* pengguna. Adapun keterangan penjelasan pada tabel di atas yaitu:

*true positif (tp)*: jumlah dokumen yang ditemu kembalikan sesuai dengan *query* pengguna

*false negative (fn)*: jumlah dokumen yang relevan dengan *query* namun tidak ditemukembali oleh sistem

*false positif (fp)*: jumlah dokumen yang ditemu kembalikan tidak sesuai dengan *query*

*true negative (tn)*: jumlah dokumen tidak relevan dengan *query* yang tidak ditemu kembalikan oleh sistem

## 2.6 Penelitian Terkait

Berdasarkan penelitian sebelumnya tentang perbandingan metode *single linkage clustering* dan *k-means clustering* (Rendy et all, 2014) dengan menggunakan parameter perhitungan *Silhouette Coefficient* dan *Purity* sebagai pembanding performasi kedua metode tersebut, didapatkan bahwa metode perhitungan *Single Linkage* memiliki performansi temu kembali informasi yang lebih baik dibandingkan dengan algoritma *k-means*. Hal ini dikarenakan proses saat pembentukan *cluster* awal pada algoritma *k-means* dilakukan secara acak. Pada penelitian tersebut data yang diujicobakan adalah dokumen berita *online* yang disimpan di dalam *notepad*. Proses pengujiannya adalah menguji performansi *clustering* yang dibangun dengan *dataset* 4 kategori (25, 50, 75 dan 100 *dataset*) dengan jumlah *cluster* adalah 3, 4, 5, 6, 7, 8, 9, dan 10 yang selanjutnya akan dicari masing-masing nilai *Silhouette Coefficient* dan *Purity* dari 4 bagian *dataset* yang diujicobakan. Hasil penelitian yang didapat dalam penelitian tersebut bahwa



melakukan pengujian, setiap *path* yang ada dalam tiap proses dapat dijalankan dengan baik.

Pengujian terhadap proses integrasi dilakukan dengan membandingkan proses integrasi dengan perhitungan manual dan dengan sistem. Dokumen yang dipergunakan adalah dokumen dengan judul Data Mining 1 dan Data Mining 2. Pada akhir pengujian didapatkan hasil bahwa hasil integrasi secara manual menghasilkan dokumen yang sama dengan dokumen yang diintegrasikan oleh sistem. *Agglomerative*

Hak Cipta Dilindungi Undang-Undang

1. Dilarang mengutip sebagian atau seluruh karya tulis ini tanpa mencantumkan dan menyebutkan sumber:
  - a. Pengutipan hanya untuk kepentingan pendidikan, penelitian, penulisan karya ilmiah, penyusunan laporan, penulisan kritik atau tinjauan suatu masalah.
  - b. Pengutipan tidak merugikan kepentingan yang wajar UIN Suska Riau.
2. Dilarang mengumumkan dan memperbanyak sebagian atau seluruh karya tulis ini dalam bentuk apapun tanpa izin UIN Suska Riau.