

TUGAS BESAR II IF2123
ALJABAR LINEAR DAN GEOMETRI
2020/2021

Oleh

Muhammad Dzaki Razaan Faza 13519033

Joel Triwira 13519073

Rizky Anggita Syarbaini Siregar 13519132



SEKOLAH TEKNIK ELEKTRO DAN INFORMATIKA
INSTITUT TEKNOLOGI BANDUNG
2020

BAB I

Deskripsi Masalah

Hampir semua orang pernah menggunakan *search engine*, seperti *google*, *bing* dan *yahoo! search*. Setiap hari, bahkan untuk sesuatu yang sederhana orang-orang menggunakan mesin pencarian. Pada tugas besar Aljabar Linear dan Geometri kali ini penulis akan menjelaskan bagaimana cara *search engine* tersebut mendapatkan semua dokumen kita berdasarkan apa yang ingin pengguna cari.

Sebagaimana yang telah diajarkan di dalam kuliah pada materi vektor di ruang Euclidean, temu-balik informasi (*information retrieval*) merupakan proses menemukan kembali (*retrieval*) informasi yang relevan terhadap kebutuhan pengguna dari suatu kumpulan informasi secara otomatis. Biasanya, sistem temu balik informasi ini digunakan untuk mencari informasi pada informasi yang tidak terstruktur, seperti laman *web* atau dokumen.

Ide utama dari sistem temu balik informasi adalah mengubah *search query* menjadi ruang vektor. Setiap dokumen maupun *query* dinyatakan sebagai vektor $w = (w_1, w_2, \dots, w_n)$ di dalam R^n , dimana nilai w_i dapat menyatakan jumlah kemunculan kata tersebut dalam dokumen (*term frequency*). Penentuan dokumen mana yang relevan dengan *search query* dipandang sebagai pengukuran kesamaan (*similarity measure*) antara *query* dengan dokumen. Semakin sama suatu vektor dokumen dengan vektor *query*, semakin relevan dokumen tersebut dengan *query*. Kesamaan tersebut dapat diukur dengan *cosine similarity* dengan rumus:

$$\text{sim}(\mathbf{Q}, \mathbf{D}) = \cos \theta = \frac{\mathbf{Q} \cdot \mathbf{D}}{\|\mathbf{Q}\| \|\mathbf{D}\|}$$

Pada tugas besar kali ini dengan menggunakan model ruang vektor dan *cosine similarity* penulis membuat sebuah *search engine* sederhana.

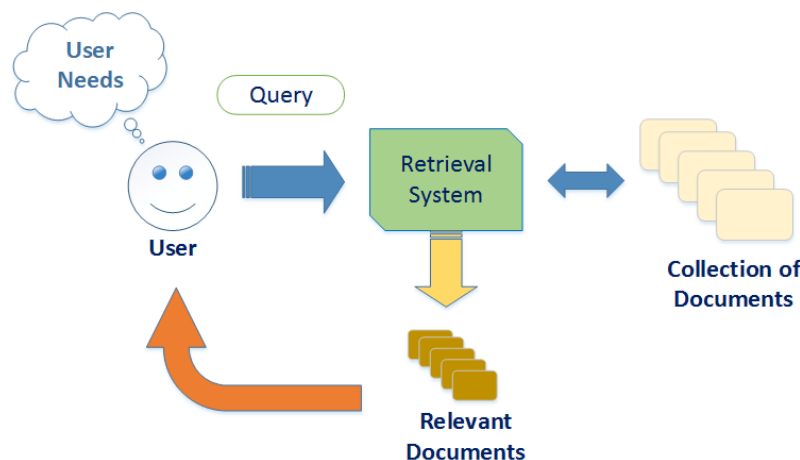
BAB II

Teori Singkat

2.1 Information Retrieval

Sistem *information retrieval* (IR) system adalah sistem yang digunakan untuk menemukan kembali (*retrieve*) informasi-informasi yang relevan terhadap kebutuhan pengguna dari suatu kumpulan informasi secara otomatis [Bunyamin, 2005].

Sistem IR terutama berhubungan dengan pencarian informasi yang isinya tidak memiliki struktur. Demikian pula ekspresi kebutuhan pengguna yang disebut *query*, juga tidak memiliki struktur. Hal ini yang membedakan sistem IR dengan sistem basis data. Dokumen adalah contoh informasi yang tidak terstruktur. Isi dari suatu dokumen sangat tergantung pada pembuat dokumen tersebut.



Gambar 2.1 Gambaran Information Retrieval

Sumber: <https://ir.cs.ui.ac.id/new/images/ir.png>

Sebagai suatu sistem, sistem IR memiliki beberapa bagian yang membangun sistem secara keseluruhan. Berikut adalah bagian-bagian dari sistem IR:

1. *Text Operations* (operasi terhadap teks) yang meliputi pemilihan kata-kata dalam *query* maupun dokumen (*term selection*) dalam pentransformasian dokumen atau *query* menjadi term *index* (indeks dari kata-kata).

2. *Query formulation* (formulasi terhadap *query*) yaitu memberi bobot pada indeks kata-kata *query*.
3. *Ranking* (perangkingan), mencari dokumen-dokumen yang relevan terhadap *query* dan mengurutkan dokumen tersebut berdasarkan kesesuaiannya dengan *query*.

Sistem IR menerima *query* dari pengguna, kemudian melakukan perangkingan terhadap dokumen pada koleksi berdasarkan kesesuaiannya dengan *query*. Hasil perangkingan yang diberikan kepada pengguna merupakan dokumen yang menurut sistem relevan dengan *query*.

2.2 Vektor

Vektor adalah kuantitas fisik yang memiliki besar dan arah. Salah satu model IR adalah model menggunakan ruang vektor. Model ini menggunakan teori di dalam aljabar vektor. Misalkan terdapat n kata berbeda sebagai kamus kata (*vocabulary*) atau indeks kata (*term index*). Kata-kata tersebut membentuk ruang vektor berdimensi n . Setiap dokumen maupun *query* dinyatakan sebagai vektor $w = (w_1, w_2, \dots, w_n)$ di dalam R^n . w_i = bobot setiap kata i di dalam *query* atau dokumen. Nilai w_i dapat menyatakan jumlah kemunculan kata tersebut dalam dokumen (*term frequency*).

Ada beberapa langkah atau proses untuk mendapatkan hasil dari *query* yang dimasukkan dengan menggunakan model ruang vektor:

1. Membuang kata depan dan kata penghubung.
2. Menggunakan *stemmer* pada kumpulan dokumen dan *query*, yaitu aplikasi yang digunakan untuk menghilangkan imbuhan (awalan, akhiran).
3. Menghitung banyaknya frekuensi atau kemunculan kata dalam kumpulan dokumen yang sesuai dengan *query*.
4. Mengubah dokumen dan *query* menjadi vektor
5. Mengurutkan dokumen berdasarkan similaritas, dengan menghitung perkalian titik vektor (*cosine similarity*).

2.3 Cosine similarity

Penentuan dokumen mana yang relevan dengan *query* dipandang sebagai pengukuran kesamaan (*similarity measure*) antara *query* dengan dokumen. Semakin sama suatu vektor dokumen dengan vektor *query*, semakin relevan dokumen tersebut dengan

query. Kesamaan (*sim*) antara dua vektor $Q = (q_1, q_2, \dots, q_n)$ dan $D = (d_1, d_2, \dots, d_n)$ diukur dengan rumus *cosine similarity* yang merupakan bagian dari rumus perkalian titik (*dot product*) dua buah vektor:

$$\text{sim}(\mathbf{Q}, \mathbf{D}) = \cos \theta = \frac{\mathbf{Q} \cdot \mathbf{D}}{\|\mathbf{Q}\| \|\mathbf{D}\|}$$

dengan $\mathbf{Q} \cdot \mathbf{D}$ adalah perkalian titik yang didefinisikan sebagai:

$$\mathbf{Q} \cdot \mathbf{D} = q_1 d_1 + q_2 d_2 + \dots + q_n d_n$$

Jika $\cos \theta = 1$, berarti $\theta = 0$, vektor Q dan D berimpit, yang berarti dokumen D sesuai dengan *query* Q . Jadi, nilai cosinus yang besar (mendekati 1) mengindikasikan bahwa dokumen cenderung sesuai dengan *query*. Setiap dokumen di dalam koleksi dokumen dihitung kesamaannya dengan *query* dengan rumus cosinus di atas.

Selanjutnya hasil perhitungan diurutkan berdasarkan nilai cosinus dari besar ke kecil sebagai proses pemilihan dokumen yang “dekat” dengan *query*. Pengurutan tersebut menyatakan dokumen yang paling relevan hingga yang kurang relevan dengan *query*. Nilai cosinus yang besar menyatakan dokumen yang relevan, nilai cosinus yang kecil menyatakan dokumen yang kurang relevan dengan *query*.

BAB III

Implementasi Program

Program yang dibuat menggunakan bahasa *Python* dan menggunakan framework *Flask* untuk *backend* nya. Pada bagian *frontend*, kami menggunakan *HTML*, *CSS*, dan *Bootstrap 4*.

Terdapat 4 buah program modular (*backend*), 3 buah file *HTML*, dan 1 buah file *CSS* yang kami buat pada Tugas Besar II ini.

1. *App.py*

Kami menggunakan *Flask* sebagai *framework* untuk membuat *web* aplikasi ini. *Library* yang digunakan pada *app.py* adalah *flask*, *werkzeug*, *os*, dan *shutil*. Berikut adalah penjelasan mengenai fungsi yang terdapat pada *app.py*.

Tabel 3.1 *App.py*

Fungsi/Prosedur	Keterangan
<i>home</i>	Menampilkan <i>home.html</i> .
<i>show</i>	Menampilkan <i>file txt</i> yang di-input.
<i>search</i>	Menerima <i>file .txt</i> dan <i>query</i> yang ingin dicari, kemudian menampilkan hasil pencarian. <i>Method = GET, POST</i>
<i>allowed_files</i>	Mengecek apakah <i>file</i> yang diinput memiliki format yang benar (<i>.txt</i>).
<i>request_txt</i>	Menerima <i>multiple files input</i> dari <i>user</i> .

2. *Text.py*

Kumpulan fungsi yang melakukan *text-processing* pada dokumen yang di-input oleh pengguna. Dokumen yang diinput adalah berbahasa Indonesia. Kami menggunakan *library Sastrawi* untuk melakukan beberapa proses yaitu *filtering* *stopword* dan *stemming*. Berikut beberapa fungsi yang terdapat di *text.py*.

Tabel 3.2 *Text.py*

Fungsi/Prosedur	Keterangan
<i>stemming_doc</i>	Menyiapkan objek <i>Stemmer</i> dan melakukan <i>stemming</i> pada dokumen/ <i>query</i> .
<i>filtering_stopword</i>	Menyiapkan objek <i>StopWordRemover</i> dan melakukan <i>filtering stopwords</i> pada dokumen/ <i>query</i> .
<i>stemming_filtering_doc</i>	Melakukan <i>filtering stopwords</i> dan <i>stemming</i> pada kumpulan dokumen dengan menggunakan fungsi <i>stemming_doc</i> dan <i>filtering_stopword</i> di atas secara sekaligus.
<i>stemming_filtering_query</i>	Melakukan <i>filtering stopwords</i> dan <i>stemming</i> pada <i>query</i> dengan menggunakan fungsi <i>stemming_doc</i> dan <i>filtering_stopword</i> di atas secara sekaligus.
<i>first_sentence</i>	Melakukan pengambilan kalimat pertama dari setiap dokumen yang di- <i>input user</i> .

3. *Docs.py*

Kumpulan fungsi yang melakukan proses lebih detail terhadap dokumen yang di-*input* oleh pengguna. Berikut fungsi-fungsi yang terdapat pada *docs.py*.

Tabel 3.3 *Docs.py*

Fungsi	Keterangan
<i>open_doc</i>	Melakukan <i>opening file</i> dan mengembalikan <i>txt</i> dari <i>file</i> tersebut.
<i>count_kata_doc</i>	Menghitung jumlah kata pada kumpulan dokumen sebelum dilakukan <i>text-processing</i> apapun pada kumpulan dokumen tersebut.
<i>count_word</i>	Menghitung jumlah kata per dokumen

<i>remove_duplicate</i>	Menghilangkan duplikasi pada kumpulan kata (<i>term</i>) dari setiap dokumen dan <i>query</i> .
<i>makeTerm</i>	Membentuk <i>term</i> dari kumpulan kata yang berasal dari dokumen dan <i>query</i> .

4. Similarity.py

Kumpulan fungsi yang melakukan proses terhadap perhitungan kemiripan (similarity) antara query dan dokumen. Terdapat pustaka(library) *math* yang digunakan untuk menghitung akar pangkat 2 (square root). Berikut fungsi-fungsi yang terdapat pada similarity.py

Fungsi	Keterangan
<i>dot_product</i>	Menghitung hasil kali titik (dot product) dari dua buah vektor.
<i>besar</i>	Menghitung norm(besar/panjang) dari sebuah vektor.
<i>sim</i>	Menghitung cosine similarity pada query dan dokumen dengan menggunakan fungsi <i>dot_product</i> dan <i>besar</i> .
<i>count_found_term</i>	Menghitung kemunculan term pada pada dokumen dan query dan mengembalikan sebuah array yaitu <i>arrayHasil</i> .
<i>sortHasil</i>	Melakukan pengurutan(sorting) pada beberapa array. Pengurutan berdasarkan <i>arrayHasil</i> yang merupakan hasil perhitungan dari <i>count_found_term</i> .

5. Layout.html

Merupakan *template* yang dibuat untuk di-load pada file *html* lainnya.

6. *Index.html*

Merupakan *file html* yang digunakan pada saat pengguna mengklik *Search*.

7. *Search.html*

Merupakan *file html* yang digunakan pada saat hasil pencarian ditampilkan.

8. *View.html*

Merupakan *file html* yang digunakan pada saat menampilkan isi dari *file txt* yang di-*upload* oleh pengguna.

9. *Home.html*

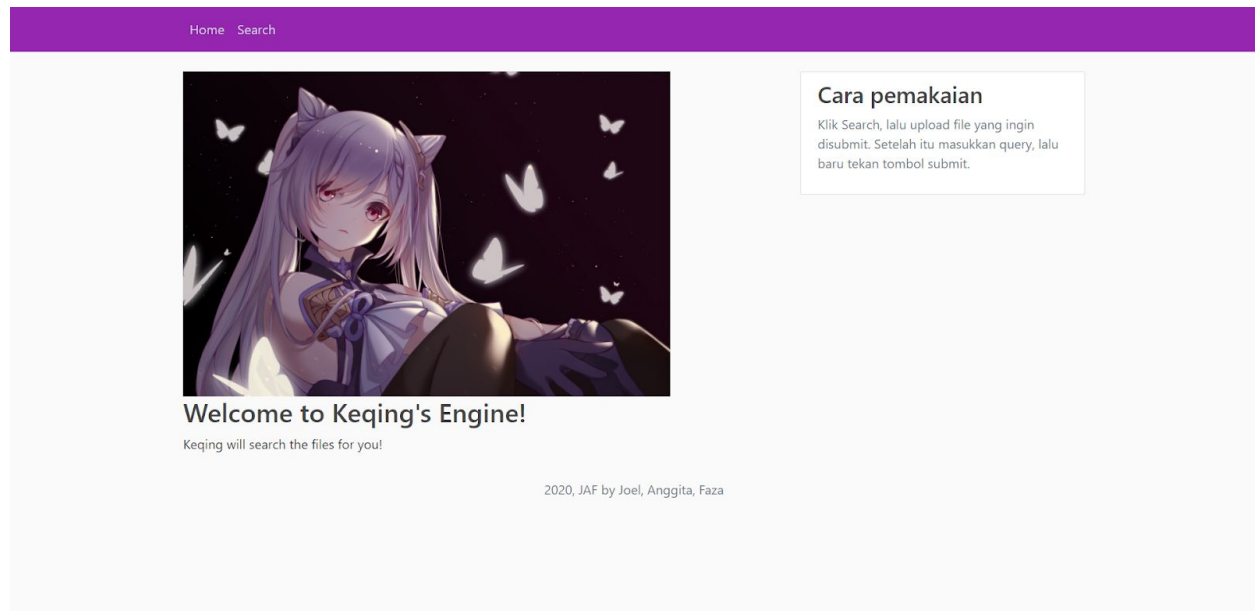
Merupakan *file html* yang digunakan untuk *landing page*.

10. *Main.css*

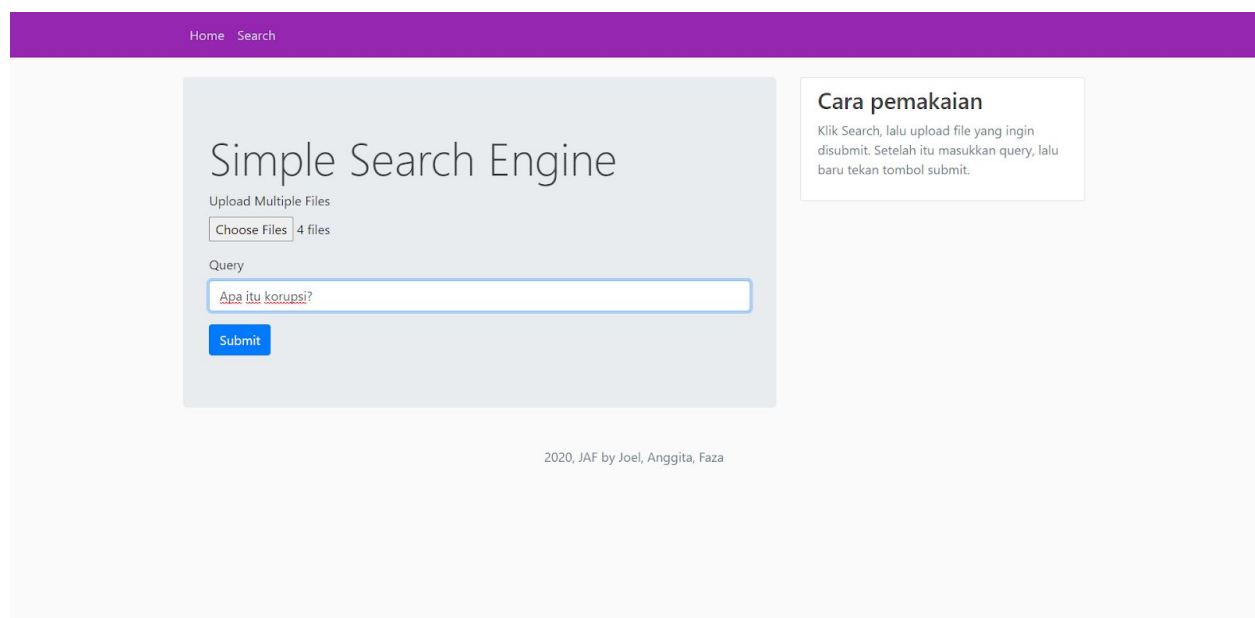
Merupakan *file CSS* yang digunakan untuk *styling* pada tampilan *web* ini.

BAB IV

Eksperimen



Gambar 4.1 Halaman Depan *Website*



Gambar 4.2 Tampilan Input Dokumen dan Query dari Pengguna

[Home](#)
[Search](#)

Upload Multiple Files
 No file chosen

Query

Cara pemakaian
 Klik Search, lalu upload file yang ingin di-submit. Setelah itu masukkan query, lalu baru tekan tombol submit.

Query yang dicari
 Apa itu korupsi?

Hasil Pencarian

1. Dokumen Apa Itu Korupsi.txt

Jumlah kata : 57
 Tingkat kemiripan: 19.069%
 Kalimat pertama: Korupsi atau rasuah berasal dari bahasa Latin: corruptio dari kata kerja corrumpere yang bermakna busuk, rusak, menggoyahkan, memutarbalik, menyogok.

2. Dokumen Analisis Kompetitor.txt

Jumlah kata : 201
 Tingkat kemiripan: 0.0%
 Kalimat pertama: Kompetitor yang dimiliki Bepi terdiri dari dua yaitu kompetitor langsung dan tidak langsung.

3. Dokumen Analisis Masalah.txt

Jumlah kata : 238
 Tingkat kemiripan: 0.0%
 Kalimat pertama: Menurut WHO (World Health Organization) corona merupakan infeksi penyakit yang disebabkan oleh virus baru yang dinamakan COVID-19.

4. Dokumen arduino.txt

Jumlah kata : 84
 Tingkat kemiripan: 0.0%
 Kalimat pertama: Arduino adalah pengendali mikro single-board yang bersat sumber terbuka, diturunkan dari Wiring platform, dirancang untuk memudahkan penggunaan elektronik dalam berbagai bidang.

No.	Term	Q	D1	D2	D3	D4
1	apa	1	0	0	0	0
2	korupsi	1	2	0	0	0

2020, JAF by Joel, Anggita, Faza

Gambar 4.3 Hasil Pencarian *Terms Query* pada Dokumen

BAB V

Simpulan

5.1 Simpuln

Search engine sederhana dapat dibuat dengan menerapkan *cosine similarity* yang memanfaatkan model ruang vektor. *Query* dan dokumen yang didapat dari pengguna dapat diolah menggunakan *stemming* dan penghapusan *stopword* lalu dimasukkan ke dalam sebuah array yang sudah menampung *term* yang bersih. Selanjutnya, kemunculan *term* di *query* dalam dokumen dapat dihitung sehingga didapat vektor yang dapat dimanfaatkan untuk mencari *cosine similarity* dari *query* dan dokumen. Lalu, dokumen diurutkan berdasarkan nilai kemiripannya dengan *query*. Dengan bantuan *framework flask* dan sedikit *html* serta *css*, *search engine* dapat diaplikasikan dalam bentuk *website* lokal sederhana.

5.2 Saran

Penulis menyadari bahwa dalam pengerjaan program masih banyak yang dapat dikembangkan, sebagai berikut:

1. Memperindah desain *website* supaya lebih menarik bagi pengguna. Beberapa diantaranya adalah desain tombol, *typography*, *color pallete*, dan membuat tampilan yang responsif.
2. Membuat kode lebih rapi dan memberikan komentar yang lebih detail serta mudah dipahami.
3. Memberikan konvensi dalam menentukan standar nama fungsi atau variabel.
4. Menerapkan *web scraping* agar bisa menerima input *html*.

5.3 Refleksi

Melalui tugas ini, penulis belajar untuk bekerja sama dalam kelompok, tidak menunda - nunda pekerjaan, berbagi tugas, dan juga membuat sebuah *website* baik dari segi *frontend* maupun *backend* menggunakan *flask*.

Daftar Pustaka

- Bunjamin, H., & Negara, C. P. (2008). Aplikasi Information Retrieval (IR) CATA Dengan Metode Generalized Vector Space Model. *Jurnal Informatika*, 4(1), 29-38.
<https://media.neliti.com/media/publications/219301-none.pdf>
- Flask Documentation. (n.d.). *Flask Table*. <https://flask-table.readthedocs.io/en/stable/>
- Munir, R. (2020). Aplikasi Dot Product pada Sistem Temu-balik Informasi (Information Retrieval System). In *Bahan kuliah IF2123 Aljabar Linier dan Geometri* (pp. 8-14). Program Studi Teknik Informatika, STEI-ITB.
- Schafer, C. (2018, May Friday). *Python Flask Tutorial: Full-Featured Web App*. Retrieved November, 2020, from
https://www.youtube.com/watch?v=QnDWIZuWYW0&list=PL-osiE80TeTs4UjLw5MM6OjgkjFeUxCYH&index=2&ab_channel=CoreySchafer
- Visual Studio Code. (n.d.). *Flask Tutorial in Visual Studio Code*.
<https://code.visualstudio.com/docs/python/tutorial-flask>