

**Tosan Iqbal Kurniawan**

**A11.2023.15409**

**UTS STKI**

## **Pendahuluan**

Proyek ini bertujuan membangun sebuah mini project sistem temu kembali informasi (STKI). Saya memfokuskan untuk mencari dan mengklasifikasi berita dari dataset berita CNN. STKI sendiri merupakan sistem yang dirancang khusus untuk menangani data tidak terstruktur seperti teks berita. Dengan memproses koleksi dokumen dan mengembalikannya berdasarkan query pengguna. Elemen kunci dalam sistem ini adalah penggunaan indeks dan mekanisme ranking untuk mengurutkan dokumen yang paling relevan.

Untuk mengimplementasikannya, proyek ini menggunakan dua pendekatan utama: Boolean Retrieval Model dan Vector Space Model (VSM). Model Boolean bekerja dengan mencocokkan kata kunci menggunakan operator logika (AND, OR, NOT), sedangkan VSM menggunakan representasi vektor untuk menghitung kesamaan (similarity) antara query dan dokumen. Seluruh alur kerja sistem, mulai dari preprocessing data hingga tahap ranking dan presentasi hasil, akan divisualisasikan menggunakan diagram arsitektur search engine klasik.

## **Data & Preprocessing**

Untuk tahap preprocessing, saya menggunakan data berupa kumpulan berita yang diambil dari CNN. Sebanyak 9 dokumen dipilih untuk diproses dalam proyek ini. Setiap dokumen berita berisi teks yang cukup panjang dengan berbagai informasi terkait topik terkini. Dokumen-dokumen ini kemudian diubah ke dalam format .txt untuk mempermudah pemrosesan lebih lanjut. Pada langkah pertama dalam preprocessing, dilakukan tokenisasi untuk memecah setiap dokumen menjadi token-token yang lebih kecil, seperti kata atau frasa. Langkah selanjutnya adalah case-folding, yang dilakukan untuk mengubah seluruh teks menjadi huruf kecil. Hal ini bertujuan untuk menghindari perbedaan antara kata-kata yang memiliki kapitalisasi yang berbeda. Setelah case-folding, langkah berikutnya adalah stopword removal, di mana kata-kata yang tidak membawa makna signifikan dalam pencarian, seperti "dan", "di", "yang", dihapus. Setelah menghapus stopwords, dilakukan stemming untuk mengubah kata-kata turunan menjadi bentuk dasar. Langkah terakhir dalam preprocessing adalah normalisasi angka dan tanda baca. Semua angka dan tanda baca yang tidak relevan dihapus atau diganti untuk memudahkan analisis lebih lanjut. Setelah kesembilan dokumen tersebut melalui lima tahapan preprocessing, data teks mentah telah berhasil ditransformasi menjadi kumpulan token yang bersih, seragam, dan siap untuk dianalisis lebih lanjut. Untuk menunjukkan secara konkret efektivitas dari proses ini, berikut disajikan perbandingan antara kondisi dokumen asli (before preprocessing) dengan hasil akhir dokumen (after preprocessing):

The screenshot shows a terminal window with four sections of text, each representing a document. The first section is 'Dokumen 2', the second is 'Dokumen 3', and the third is 'Dokumen 4'. Each section has an 'Original:' part at the top and a 'Processed:' part below it. The 'Original:' text is the raw news clipping, while the 'Processed:' text is the same text after being run through the preprocessing pipeline, showing how it has been tokenized, converted to lowercase, and stop words removed.

```
=====
--- Dokumen 2 ---
Original:
Mike Tyson menyebut legenda tinju dunia,Muhammad Ali, sebagai petinju yang bisa mengalahkannya saat berada di puncak kariernya. Tyson merupakan salah satu petinju paling ganas di atas ring. Dia juga dikenal karena tekniknya yang agresif dan gagah. Tyson pernah mengalahkan petinju legendaris lainnya seperti Muhammad Ali dan Joe Frazier. Meskipun karier profesionalnya berakhir pada tahun 1988, Tyson tetap menjadi ikon dalam dunia tinju.

Processed:
mike tyson menyebut legenda tinju duniamuhammad ali petinju mengalahkannya puncak kariernya tyson salah petinju gana ring merebut gelar juara dunia mengalahkan trevor berwick tyson

=====

--- Dokumen 3 ---
Original:
Anggota Komisi IX DPR RI Ashabul Kahfi mengatakan perusahaan di Surabaya yang diduga memotong gaji karyawannya yang salat Jumat dan menahan ijazah karyawannya, telah melanggar hukum dan tidak adil.

Processed:
anggota komisi ix dpr ri ashabul kahfi perusahaan surabaya diduga memotong gaji karyawannya yang salat jumat dan menahan ijazah karyawannya melanggar hukum ditoleransi ashabul menerangkan

=====

--- Dokumen 4 ---
Original:
Pesta diskon bertajuk Transmart Full Day Sale kembali lagi, Minggu (20/4). Kulkasside by side(SBS) banting harga gila-gilaan, dari harga Rp9 juta jadi Rp6 jutaan saja! Khusus untuk pembelian pulau jawa bali lampung kulka sb 1 trans

Processed:
pesta diskon bertajuk transmart full day salekembali minggu kulkassid by sidesb banting harga gilagilaan harga rp 9 juta rp 6 jutaan aja khusus pembelian pulau jawa bali lampung kulka sb 1 trans
```

## Boolean Retrieval Model

Saya mengimplementasikan Boolean Retrieval Model dengan memproses query menggunakan operator logika AND, OR, dan NOT. Langkah fundamental untuk model ini adalah membangun Vocabulary, yaitu kamus yang berisi semua kata unik dari 9 dokumen berita yang telah diproses. Dari *vocabulary* tersebut, saya membangun dua struktur data inti:

1. Incidence Matrix: Sebuah matriks biner (sparse) yang menunjukkan ada (1) atau tidaknya (0) sebuah kata pada setiap dokumen.
  2. Inverted Index: Struktur data yang lebih efisien yang memetakan setiap kata (*term*) ke daftar dokumen yang mengandung kata tersebut.

Berikut adalah cuplikan dari *Vocabulary* (10 kata pertama), *Incidence Matrix* (untuk beberapa dokumen dan kata pertama), dan *Inverted Index* (5 term pertama) yang berhasil dibangun dari korpus berita CNN:

Selain itu, saya juga menguji dan mendemonstrasikan model pencari Boolean serta melakukan uji wajib menghitung precision/recall pada Boolean result set dengan gold set dengan hasil berikut:

```
--- Query 3: dpr ri gaji karyawan ---  
Query diproses: ['dpr', 'ri', 'gaji', 'karyawan']  
Ditemukan di dokumen ID: [3]  
Cuplikan Dokumen 3: anggota komisi ix dpr ri ashabul kahfi perusahaan surabaya diduga memotonggajikaryawan yangsalat jumatdan menahan ijazah karyawannya melalui  
  
--- Query 4: transmart diskon sepeda ---  
Query diproses: ['transmart', 'diskon', 'sepeda']  
Ditemukan di dokumen ID: [6]  
Cuplikan Dokumen 6: transmart full day salebalik kasih diskon gede pelanggan berbelanja gelaran bawa pulang sepeda diobral rp jutaan aneka sepeda gunung  
  
--- Query 5: liverpool gelar liga ---  
Query diproses: ['liverpool', 'gelar', 'liga']  
Ditemukan di dokumen ID: [5, 9]  
Cuplikan Dokumen 5: moham salahberpeluang menciptakan rekor unik saatliverpoolmelawan leicest citi laga krusial perebutan gelar liga inggri stadion king power  
Cuplikan Dokumen 9: arsen peluang merebut gelar juaraliga championsmusim berat mengejar liverpool liga inggri rotasi pemain arsen arseen tertigg angka liverpool  
  
--- Evaluating Boolean Retrieval Model Performance ---  
  
--- Scenario 1: Query = 'Transmart AND diskon' ---  
Gold Relevant Docs (IDs): [3, 6]  
Retrieved Docs (IDs): [4, 6]  
Precision: 0.5000  
Recall: 0.5000  
  
--- Scenario 2: Query = 'Liverpool OR tinju' ---  
Gold Relevant Docs (IDs): [2, 5, 7]  
Retrieved Docs (IDs): [2, 5, 9]  
Precision: 0.6667  
Recall: 0.6667  
  
--- Scenario 3: Query = 'korupsi NOT Surabaya' ---  
Gold Relevant Docs (IDs): [1, 4]  
Retrieved Docs (IDs): [1, 7]  
Precision: 0.5000  
Recall: 0.5000
```

## Vector Space Model (VSM)

Pada bagian ini, kami mengimplementasikan Vector Space Model (VSM) untuk perankingan dokumen berdasarkan kesamaan antara query dan dokumen. VSM merepresentasikan dokumen dan query sebagai vektor, di mana setiap elemen vektor mewakili suatu kata dalam koleksi dokumen. Untuk menghitung Term Frequency (TF), kami pertama-tama menggunakan CountVectorizer untuk menghitung frekuensi kemunculan setiap kata dalam dokumen.

Selanjutnya, untuk menghitung TF-IDF (Term Frequency-Inverse Document Frequency), kami menggunakan TfidfVectorizer, yang memperhitungkan bobot kata berdasarkan frekuensinya dalam dokumen dan seberapa jarang kata tersebut muncul di seluruh koleksi dokumen. Proses ini memungkinkan model untuk menilai kata-kata yang lebih relevan dalam konteks pencarian, dengan memberi bobot lebih pada kata-kata yang sering muncul dalam dokumen tetapi jarang muncul di dokumen lainnya. Dengan menggunakan TF-IDF, model dapat memberikan ranking yang lebih akurat berdasarkan kesamaan antara query dan dokumen.

Sebagai langkah awal, CountVectorizer menghasilkan matriks Term Frequency (TF) yang berisi frekuensi mentah (raw count) dari setiap kata dalam setiap dokumen. Berikut adalah cuplikan dari matriks TF untuk 5 dokumen pertama dan 10 term pertama:

--- Term Frequency (TF) Matrix (Raw Counts - Cuplikan) ---										
	abdul	advertis	agam	agama	agamanya	aguero	agung	aja	akar	akibat
Doc 1	1	1	1	0	0	0	2	0	0	0
Doc 2	0	1	0	0	0	0	0	0	0	0
Doc 3	0	1	0	1	2	0	0	0	0	0
Doc 4	0	1	0	0	0	0	0	1	0	0
Doc 5	0	1	0	0	0	1	0	0	0	0

Penjelasan: Setiap angka dalam tabel menunjukkan berapa kali suatu 'term' muncul di 'dokumen' tersebut.

Setelah mendapatkan frekuensi mentah (TF), langkah selanjutnya adalah menghitung Inverse Document Frequency (IDF). IDF memberikan bobot lebih tinggi pada kata-kata yang unik (jarang muncul di seluruh dokumen) dan bobot rendah pada kata-kata yang umum. IDF dihitung dari Document Frequency (DF), yaitu jumlah dokumen di mana sebuah term muncul.

Untuk menganalisis korpus saya, saya mengidentifikasi term yang paling umum (DF tinggi) dan yang paling unik (IDF tinggi). Seperti yang terlihat di bawah, term seperti 'advertis' dan 'scroll' sangat umum (muncul di 8 dari 9 dokumen), sementara term seperti 'zona' dan 'abdul' sangat unik (IDF tinggi).

... --- Top 10 Terms by IDF (paling unik/jarang) ---	
Term	IDF
zona	2.6894
abdul	2.6894
zinchenko	2.6894
agam	2.6894
virgil	2.6894
villa	2.6894
van	2.6894
utama	2.6894
uskup	2.6894
usaha	2.6894

  

... --- Top 10 Terms by DF (paling sering muncul di dokumen berbeda) ---	
Term	DF
scroll	8
advertis	8
content	8
continu	8
minggu	5
orang	5
memiliki	4
rp	4
pidana	3
khusu	3

Tabel berikut memberikan ringkasan lebih rinci dari nilai DF dan IDF untuk beberapa term yang dihitung dari korpus saya:

... --- Ringkasan DF dan IDF untuk Beberapa Term (Tabel) ---		
Term	DF	IDF
advertis	8	1.1054
content	8	1.1054
continu	8	1.1054
scroll	8	1.1054
minggu	5	1.5108
orang	5	1.5108
memiliki	4	1.6931
rp	4	1.6931
berat	3	1.9163
dikutip	3	1.9163
gelar	3	1.9163
juara	3	1.9163
juta	3	1.9163
kalah	3	1.9163
khusu	3	1.9163

Terakhir, kami menggunakan TfidfVectorizer untuk menggabungkan nilai TF dan IDF menjadi satu matriks pembobotan TF-IDF. Matriks inilah yang menjadi representasi VSM untuk korpus 9 berita CNN kami. Ukuran (shape) dari matriks *sparse* yang dihasilkan menunjukkan bahwa korpus kami terdiri dari 9 dokumen dan 1015 term unik (kata):

... TF-IDF Sparse Matrix Shape: (9, 1015)

Dengan matriks TF-IDF ini, sistem siap menerima *query*, mengubah *query* tersebut menjadi vektor TF-IDF, dan menghitung Cosine Similarity untuk menemukan dokumen yang paling relevan.

Setelah matriks TF-IDF siap, kami menghitung Cosine Similarity antara vektor *query* dan vektor dokumen untuk mendapatkan *ranking* relevansi. Berikut adalah hasil retrieve top k\_documents (dengan k=3) untuk dua query pengujian.

```
... Functions 'calculate_cosine_similarity' and 'retrieve_top_k_documents' have been defined.

Top 3 documents for query: 'diskusi kejaksaan agung dan uang'
Rank 1: Document ID 1 (Similarity: 0.3275)
  Snippet: kejaksaan agung kejagung menyebut uang dititip tersangkahakim djuyamto kepada satpam pengadilan negeri jakarta selatan mencapai rp juta kepala pusat
Rank 2: Document ID 7 (Similarity: 0.0227)
  Snippet: uskup agung jakarta kardin mgr ignatius suharyo mengatakan paska momen keduluan upaya membantu lemah dilemahkan indik bakti tuhan bangsa damai se
Rank 3: Document ID 9 (Similarity: 0.0000)
  Snippet: arsen peluang merebut gelar juaraliga championsmusim berat mengejar liverpool liga inggris rotasi pemain arsen arseen tertinggi angka liverpool siswa

Top 3 documents for query: 'sepak bola liga inggris liverpool'
Rank 1: Document ID 9 (Similarity: 0.3380)
  Snippet: arsen peluang merebut gelar juaraliga championsmusim berat mengejar liverpool liga inggris rotasi pemain arsen arseen tertinggi angka liverpool siswa
Rank 2: Document ID 5 (Similarity: 0.2445)
  Snippet: moham salahberpeluang menciptakan rekor unik saat liverpool melawan leicester city laga krusial perebutan gelar liga inggris stadion king power minggu
Rank 3: Document ID 8 (Similarity: 0.0000)
  Snippet: serangan udara as kami melumpuhkan pelabuhan bahan bakar hodeidah yaman serangan menyebabkan orang tewa orang terluka staf pelabuhan mengutuk serang
```

Hasilnya menunjukkan:

- Untuk *query* 'diskusi kejaksaan agung dan uang', Dokumen 1 adalah yang paling relevan (Skor: 0.3275).
- Untuk *query* 'sepak bola liga inggris liverpool', Dokumen 9 adalah yang paling relevan (Skor: 0.3380).
- Skor 0.0000 berarti tidak ada *term* relevan yang ditemukan antara *query* dan dokumen tersebut.

Saya juga membandingkan VSM dengan Gold Set yang saya buat dengan hasil seperti berikut:

```
... Evaluating TF-IDF Retrieval Performance ...
...
--- Scenario 1: Query = 'Transmart AND diskon' ---
Gold Relevant Docs (IDs): [3, 6]

Top 9 documents for query: 'Transmart AND diskon'
Rank 1: Document ID 6 (Similarity: 0.4138)
  Snippet: transmart full day salebalik kasih diskon gede pelanggan berbelanja gelaran bawa pulang sepeda diobral rp jutaan aneka sepeda gunung sepeda
Rank 2: Document ID 4 (Similarity: 0.3142)
  Snippet: pesta diskon bertajuk transmart full day salekembali minggu kulkassid by sidesb bant harga gilagila harga rp juta rp jutaan aja khusus pemb
Rank 3: Document ID 9 (Similarity: 0.0000)
  Snippet: arsen peluang merebut gelar juaraliga championsmusim berat mengejar liverpool liga inggris rotasi pemain arsen arseen tertinggi angka liverpool siswa
Rank 4: Document ID 7 (Similarity: 0.0000)
  Snippet: uskup agung jakarta kardin mgr ignatius suharyo mengatakan paska momen keduluan upaya membantu lemah dilemahkan indik bakti tuhan bangsa damai se
Rank 5: Document ID 8 (Similarity: 0.0000)
  Snippet: serangan udara as kami melumpuhkan pelabuhan bahan bakar hodeidah yaman serangan menyebabkan orang tewa orang terluka staf pelabuhan mengutuk serang
Rank 6: Document ID 5 (Similarity: 0.0000)
  Snippet: moham salahberpeluang menciptakan rekor unik saat liverpool melawan leicester city laga krusial perebutan gelar liga inggris stadion king power minggu
Rank 7: Document ID 3 (Similarity: 0.0000)
  Snippet: anggota komisi ix dpr ri ashabul kahfi perusahaan surabaya diduga memotong gaji karyawannya yang salat jumat dan menahan ijazah karyawannya melanggar hu
Rank 8: Document ID 2 (Similarity: 0.0000)
  Snippet: mike tyson menyebut legenda tinju duniamuhammad ali petinju mengalahkannya puncak kariernya tyson salah petinju gana ring merebut gelar juara dunia
Rank 9: Document ID 1 (Similarity: 0.0000)
  Snippet: kejaksaan agung kejagung menyebut uang dititip tersangkahakim djuyamto kepada satpam pengadilan negeri jakarta selatan mencapai rp juta kepala pusat
Retrieved Docs (IDs for P@3): [np.int64(6), np.int64(4), np.int64(9)]
Precision@3: 0.3333
Average Precision (AP): 0.6429
```

Analisis Singkat:

- Gold Set untuk query ini adalah dokumen [3, 6].
- Sistem saya berhasil menempatkan Dokumen 6 di Peringkat 1.
- Karena hanya 1 dari 3 dokumen teratas yang dikembalikan ([6, 4, 9]) yang ada di Gold Set, Precision@3 adalah 0.3333.
- Average Precision (AP) untuk query ini adalah 0.6429.

## DISKUSI

Disini saya ingin membahas refleksi atas keseluruhan proyek, mengidentifikasi kelebihan dan keterbatasan dari implementasi yang telah dilakukan, serta memberikan saran untuk pengembangan di masa depan.

Kelebihan :

- Implementasi Model Inti Berhasil: Proyek ini berhasil mengimplementasikan dua model STKI fundamental dari awal: Boolean Retrieval (dengan inverted index) dan Vector Space Model (dengan TF-IDF dan cosine similarity).
- Preprocessing Efektif: Pipeline preprocessing (tokenisasi, case folding, stopword removal, stemming) berhasil diterapkan untuk membersihkan 9 dokumen berita CNN, yang terbukti dari vocabulary dan matriks TF-IDF yang dihasilkan.
- Kemampuan Ranking: VSM yang dibangun terbukti mampu memberikan ranking relevansi (bukan hanya kecocokan Biner), yang divalidasi melalui skor cosine similarity dan evaluasi MAP.
- Evaluasi Kuantitatif: Proyek ini berhasil dievaluasi secara kuantitatif menggunakan metrik standar (P@3, AP, MAP), memberikan baseline kinerja yang jelas.

#### Keterbatasan:

1. Skala Korpus Sangat Kecil: Keterbatasan terbesar adalah Saya hanya mengambil sedikit korpus yang terdiri dari 9 dokumen. Ukuran sekecil ini membuat statistik TF-IDF (terutama IDF) kurang andal dan membuat evaluasi menjadi sangat sensitif.
2. Subjektivitas Gold Set: Evaluasi sangat bergantung pada Gold Set yang dibuat secara manual. Pada korpus yang kecil, penentuan relevansi bisa menjadi sangat subjektif.
3. Query Parser Sederhana: Implementasi parser untuk query Boolean tidak ditunjukkan, dan query VSM (seperti 'Transmart AND diskon') masih diproses sebagai bag-of-words sederhana oleh TfidfVectorizer tanpa logika Boolean yang eksplisit.

#### Saran Pengembangan:

Berdasarkan keterbatasan yang telah diidentifikasi, berikut adalah beberapa saran pengembangan untuk meningkatkan sistem ini di masa depan:

1. Peningkatan Skala Korpus: Langkah paling krusial adalah memperbesar dataset. Menggunakan ratusan atau ribuan dokumen berita akan menghasilkan bobot TF-IDF yang jauh lebih bermakna dan hasil evaluasi yang lebih stabil.
2. Integrasi Modul (Orkestrasi): Dalam pembuatan search engine saya mengalami kesulitan sehingga tidak menambahkan kedalam mini project UTS.
3. Implementasi Skema Pembobotan Lain: Membandingkan kinerja TF-IDF dengan skema pembobotan yang lebih canggih seperti BM25 (Okapi), yang seringkali memberikan hasil lebih baik untuk koleksi dokumen standar.