

Content-Based Filtering untuk Sistem Rekomendasi Produk E-Commerce

Fiki Pratama^{1, a)}, Sabib Prastio^{2, b)}, Erni Seniwati^{3, c)}, Ninik Tri Hartanti^{4, d)}
Sudarmanto^{5, e)}

^{1,2,3,4)} Program Studi Sistem Informasi,

⁵⁾ Program Studi Rekayasa Perangkat Lunak Aplikasi
Fakultas Ilmu Komputer

Universitas Amikom Yogyakarta, Jl. Padjajaran, Yogyakarta 55283, Indonesia

Universitas Teknologi Digital Indonesia, Jl. Raya Janti Karang Jambe, Yogyakarta 55198, Indonesia

Author Emails

^{c)} Corresponding author: erni.s@amikom.ac.id

^{a)} fikipratama@students.amikom.ac.id

^{b)} sabibprastio@students.amikom.ac.id

^{d)} ninik.t@amikom.ac.id

^{e)} darmanto@utdi.ac.id

Abstract. The rapid growth of e-commerce platforms has increased the demand for recommendation systems that can effectively assist users in discovering relevant products. Although recommendation systems have been extensively studied, empirically validated implementations using Indonesian e-commerce product datasets remain limited. This study aims to develop a content-based filtering e-commerce product recommendation system by leveraging content similarity among products. The dataset used in this research is the Tokopedia Products 2025 dataset obtained from Kaggle, which contains product information such as product names, categories, and descriptions. Text feature representation is performed using the Term Frequency–Inverse Document Frequency (TF-IDF) method to extract significant characteristics from each product. Subsequently, product similarity is calculated using Cosine Similarity to generate the most relevant product recommendations based on user preferences. The system is implemented using the Python programming language with the support of the pandas, NumPy, and scikit-learn libraries. Evaluation results using various query sample sizes and K values show that, for a representative sample of 50 query products, the system achieves a Precision@10 of 88.40%, Recall@10 of 46.81%, and F1-Score@10 of 50.88%, with the optimal F1-Score of 66.79% achieved at K = 30. These results indicate that the combination of TF-IDF and Cosine Similarity is capable of producing relevant product recommendations with good accuracy, demonstrating its feasibility as a content-based recommendation solution for e-commerce platforms.

Keywords:

Confusion Matrix, Content-Based Filtering, Cosine Similarity, E-Commerce, Recommendation System, TF-IDF

Abstraksi. Perkembangan pesat platform e-commerce mendorong kebutuhan akan sistem rekomendasi yang mampu membantu pengguna menemukan produk yang relevan secara efektif. Meskipun sistem rekomendasi telah banyak dikembangkan, implementasi yang teruji secara empiris pada dataset produk e-commerce Indonesia masih terbatas. Penelitian ini bertujuan untuk mengembangkan sistem rekomendasi produk e-commerce berbasis *content-based filtering* dengan memanfaatkan kemiripan konten antar produk. Dataset yang digunakan bersumber dari Tokopedia Products 2025 yang diperoleh melalui Kaggle, berisi informasi produk seperti nama, kategori, dan deskripsi. Representasi fitur teks dilakukan menggunakan metode *Term Frequency–Inverse Document Frequency* (TF-IDF) untuk mengekstraksi karakteristik penting dari setiap produk. Selanjutnya, tingkat kemiripan antar produk dihitung menggunakan *Cosine Similarity* guna menghasilkan rekomendasi produk yang paling relevan berdasarkan preferensi pengguna. Implementasi

dilakukan menggunakan bahasa pemrograman Python dengan pustaka *pandas*, *NumPy*, dan *scikit-learn*. Hasil evaluasi menggunakan berbagai ukuran sampel query dan nilai K yang terdiri dari sample size 50 query yang paling representatif, sistem mencapai *Precision@10* sebesar 88.40%, *Recall@10* sebesar 46.81%, dan *F1-Score@10* sebesar 50.88%, dengan nilai *F1-Score* optimal dicapai pada $K=30$ sebesar 66.79%. Hasil penelitian menunjukkan bahwa kombinasi TF-IDF dan Cosine Similarity mampu memberikan rekomendasi produk yang relevan dengan tingkat akurasi yang baik, sehingga metode ini layak diterapkan sebagai solusi sistem rekomendasi pada platform e-commerce berbasis konten

Kata Kunci:

Confusion Matrix, Content-Based Filtering, Cosine Similarity, E-Commerce, Sistem Rekomendasi, TF-IDF

PENDAHULUAN

Perkembangan platform digital dan e-commerce telah mendorong peningkatan jumlah produk dan informasi yang tersedia bagi pengguna. Kondisi ini memberikan kemudahan dalam proses pencarian produk, namun sekaligus menimbulkan permasalahan *information overload*, yaitu kesulitan pengguna dalam menentukan pilihan produk yang sesuai dengan kebutuhan dan preferensinya. Untuk mengatasi permasalahan tersebut, sistem rekomendasi banyak dikembangkan sebagai solusi yang mampu membantu pengguna dalam proses pengambilan keputusan secara lebih efektif.

Berbagai penelitian sebelumnya menunjukkan bahwa sistem rekomendasi berbasis *content-based filtering* merupakan pendekatan yang efektif dalam memberikan rekomendasi berdasarkan karakteristik item. Pebrianti *et al.* [1] mengombinasikan *content-based filtering* dengan algoritma Apriori pada sistem belanja cerdas, sementara Sulami *et al.* [2] menerapkan *content-based filtering* untuk pemilihan produk skincare. Pendekatan serupa juga diterapkan pada berbagai domain lain, seperti rekomendasi resep makanan [3], berita berbasis RSS [4], pariwisata [5][6], film [7], dan restoran [8]. Hasil penelitian-penelitian tersebut menunjukkan bahwa *content-based filtering* mampu memberikan rekomendasi yang relevan tanpa bergantung pada data pengguna lain.

Dalam konteks representasi data berbasis teks, beberapa penelitian menerapkan metode *Term Frequency-Inverse Document Frequency* (TF-IDF) yang dikombinasikan dengan *Cosine Similarity* untuk mengukur tingkat kemiripan antar item. Pendekatan ini terbukti efektif pada sistem rekomendasi film berbasis sinopsis [7] serta berbagai sistem rekomendasi berbasis deskripsi teks lainnya. Dibandingkan dengan pendekatan *collaborative filtering* [9], *content-based filtering* memiliki keunggulan dalam mengatasi permasalahan *cold-start* serta lebih sesuai untuk dataset yang tidak memiliki riwayat interaksi pengguna yang lengkap. Beberapa penelitian terkini bahkan telah mulai mengeksplorasi pendekatan berbasis *deep learning* untuk meningkatkan performa sistem rekomendasi [10], namun pendekatan tersebut memiliki kompleksitas komputasi yang lebih tinggi dan membutuhkan sumber daya yang besar.

Berdasarkan celah penelitian tersebut, permasalahan yang dikaji dalam penelitian ini adalah bagaimana menerapkan metode *content-based filtering* menggunakan TF-IDF dan *Cosine Similarity* pada dataset produk e-commerce lokal berskala besar serta bagaimana kinerja sistem rekomendasi yang dihasilkan apabila dievaluasi menggunakan pendekatan *confusion matrix*. Penelitian ini bertujuan untuk membangun sistem rekomendasi produk e-commerce berbasis konten dengan memanfaatkan TF-IDF sebagai metode representasi fitur dan *Cosine Similarity* sebagai metode pengukuran kemiripan antar produk, serta mengevaluasi kinerjanya menggunakan metrik akurasi, presisi, *recall*, dan F1-score.

Kontribusi (novelty) dari penelitian ini terletak pada tiga aspek utama. Pertama, penelitian ini menggunakan dataset Tokopedia Products 2025 yang merepresentasikan data e-commerce Indonesia yang aktual dan berskala besar, sehingga memperluas kajian *content-based filtering* pada konteks lokal. Kedua, penelitian ini menerapkan skema evaluasi berbasis *confusion matrix* pada sistem rekomendasi berbasis konten, yang masih relatif jarang digunakan dalam penelitian-penelitian sebelumnya. Ketiga, penelitian ini menyajikan analisis empiris kinerja sistem rekomendasi berbasis TF-IDF dan *Cosine Similarity* pada data e-commerce nyata, sehingga dapat menjadi referensi metodologis bagi penelitian dan pengembangan sistem rekomendasi di bidang ilmu komputer dan sistem informasi.

TINJAUAN PUSTAKA

Sistem rekomendasi merupakan salah satu bidang penelitian yang berkembang pesat dalam ilmu komputer, khususnya dalam konteks pengolahan data dan sistem pendukung keputusan. Sistem ini bertujuan untuk memberikan rekomendasi item yang relevan kepada pengguna berdasarkan data dan karakteristik tertentu. Secara umum, pendekatan sistem rekomendasi dapat diklasifikasikan menjadi *content-based filtering*, *collaborative filtering*, dan pendekatan hibrida yang mengombinasikan keduanya.

Pendekatan *content-based filtering* merekomendasikan item berdasarkan kesamaan karakteristik atau konten antar item. Beberapa penelitian menunjukkan bahwa pendekatan ini efektif dalam berbagai domain. Pebrianti *et al.* [1] mengombinasikan *content-based filtering* dengan algoritma Apriori dalam sistem belanja cerdas, sedangkan Sulami *et al.* [2] menerapkan metode ini untuk rekomendasi produk skincare. Pendekatan serupa juga digunakan dalam sistem rekomendasi resep makanan [3], berita berbasis RSS [4], pariwisata [5][6], film [7], dan restoran [8]. Hasil penelitian-penelitian tersebut menunjukkan bahwa *content-based filtering* mampu memberikan rekomendasi yang relevan tanpa bergantung pada preferensi pengguna lain.

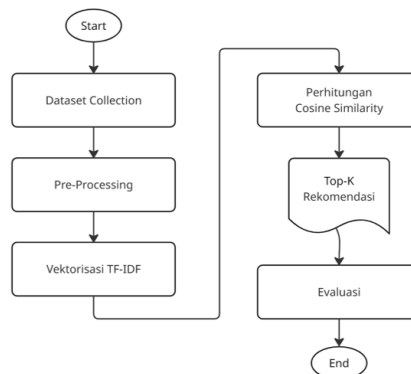
Dalam sistem rekomendasi berbasis teks, proses kombinasi TF-IDF dengan *Cosine Similarity* telah terbukti efektif dalam mengukur tingkat kemiripan antar item, sebagaimana diterapkan pada sistem rekomendasi film berbasis sinopsis [7] serta berbagai sistem rekomendasi berbasis deskripsi teks lainnya.

Selain *content-based filtering*, pendekatan *collaborative filtering* juga banyak digunakan dalam sistem rekomendasi. Pendekatan ini memanfaatkan pola perilaku dan preferensi pengguna lain untuk menghasilkan rekomendasi, sebagaimana diterapkan pada sistem rekomendasi tempat wisata [9]. Perkembangan terkini dalam sistem rekomendasi juga mengarah pada pemanfaatan metode pembelajaran mesin dan *deep learning* untuk meningkatkan kinerja rekomendasi. Maulana dan Setiawan [10] menerapkan algoritma *deep feed-forward neural network* berbasis *content-based filtering* dan menunjukkan peningkatan performa dibandingkan metode konvensional.

Berdasarkan kajian pustaka tersebut, dapat disimpulkan bahwa *content-based filtering* dengan TF-IDF dan *Cosine Similarity* merupakan pendekatan yang telah banyak digunakan dan terbukti efektif dalam berbagai domain. Pada penelitian ini difokuskan untuk mengisi celah tersebut dengan menyajikan analisis empiris kinerja sistem rekomendasi produk e-commerce berbasis konten.

METODE PENELITIAN

Pada penelitian ini memiliki alur yang dapat dilihat pada Gambar 1.



GAMBAR 1. Alur Penelitian

Dataset

Metode penelitian diawali dengan pencarian dataset yang relevan dari dataset publik di platform *Kaggle*. Tahap dataset collection ini bertujuan untuk mengumpulkan data produk e-commerce yang akan digunakan sebagai objek penelitian. Data yang dikumpulkan berisi informasi produk yang merepresentasikan karakteristik konten setiap item dan menjadi dasar utama dalam proses analisis serta pengembangan sistem rekomendasi.

Pre-processing

Tahap *Pre-Processing* dilakukan untuk membersihkan dan mempersiapkan data sebelum diproses lebih lanjut. Proses ini mencakup seleksi atribut yang relevan, penghapusan data yang tidak diperlukan, serta normalisasi data teks agar memiliki format yang konsisten dan siap untuk tahap vektorisasi.

Vektorisasi TF-IDF

Tahap Vektorisasi *TF-IDF* bertujuan untuk mengubah data teks menjadi bentuk numerik menggunakan metode *Term Frequency-Inverse Document Frequency (TF-IDF)*. Proses ini menghasilkan representasi vektor yang menggambarkan tingkat kepentingan setiap kata dalam suatu produk terhadap keseluruhan dataset. Pada tahap ini akan menggunakan 3 persamaan yaitu *Term-Frequency (TF)*, *Inverse Document Frequency (IDF)*, dan *TF-IDF*.

$$TF(t, d) = n_{t,d} \quad 1)$$

$$IDF(t) = \log_{10} \left(\frac{N}{df_t} \right) + 1 \quad 2)$$

$$TF-IDF(t, d) = TF(t, d) \times IDF(t) \quad 3)$$

Perhitungan Cosine Similarity

Tahap *Perhitungan Cosine Similarity* dilakukan untuk mengukur tingkat kemiripan antar produk berdasarkan vektor *TF-IDF* yang telah dihasilkan. Nilai kemiripan ini digunakan untuk menentukan seberapa dekat karakteristik konten antar produk dalam ruang vektor. Pada tahapan ini akan menggunakan persamaan 4, yaitu:

$$||D1|| = \sqrt{D1_1^2 + D1_2^2 + \dots + D1_n^2} \quad 4)$$

Rekomendasi Top-K Produk

Tahap *Top-K Rekomendasi* merupakan proses pemilihan sejumlah *Top-K* item produk dengan nilai kemiripan tertinggi terhadap produk acuan, Tahap ini merupakan sebuah hasil implementasi sistem rekomendasi produk e-commerce dalam bentuk prototype website. Produk-produk tersebut kemudian ditetapkan/ditampilkan dalam interface ui website sebagai hasil rekomendasi yang paling relevan bagi pengguna berdasarkan perhitungan nilai kemiripan dari metode cosine similarity.

Evaluasi

Tahap *Evaluasi* bertujuan untuk menilai kinerja sistem rekomendasi yang telah dibangun. Proses ini dilakukan dengan menganalisis kualitas dan relevansi rekomendasi yang dihasilkan berdasarkan pada metode evaluasi *confusion matrix*.

HASIL DAN PEMBAHASAN

Dataset

Dataset yang digunakan dalam penelitian ini merupakan dataset publik yang bersumber dari platform Kaggle dengan judul "*Tokopedia Products 2025*" yang dapat diakses melalui url <https://www.kaggle.com/datasets/nsmlehq/tokopedia-products-2025/>. Dataset ini berisi informasi produk e-commerce yang diambil dari platform Tokopedia, salah satu marketplace terbesar di Indonesia. Total dataset terdiri dari 1.200 record produk yang tersimpan dalam format csv.

Dataset memiliki struktur yang terdiri dari 36 atribut yang merepresentasikan berbagai karakteristik produk. Atribut-atribut tersebut meliputi informasi dasar produk seperti *id (ID)*, *name (Nama Produk)*, *description (Deskripsi)*, *price (Harga)*, *condition (Kondisi)*, dan *status*. Informasi kategori produk mencakup *category_id (ID Kategori)*, *category_name (Kategori)*, dan *category_url (URL Kategori)*. Informasi toko penjual terdiri dari *shop_id (ID Toko)*, *shop_name (Nama Toko)*, *shop_city (Kota Toko)*, *shop_url (URL Shop)*, dan *shop_reputation (Reputasi Toko)*. Atribut visual produk mencakup *thumbnail* dan *images* yang berisi url gambar produk. Selain itu, dataset juga menyediakan informasi statistik performa produk seperti *sold count*, *view count*, *review count*, *talk count*, dan *rating* yang merepresentasikan popularitas dan tingkat kepercayaan konsumen terhadap produk. Informasi diskon tersedia melalui atribut *discount_percent*, *price_text*, dan *price_slash_text*. Atribut tambahan seperti *min_order*, *max_order*, *weight*,

dan *weight_unit* memberikan informasi teknis terkait pemesanan dan pengiriman produk. Berikut pada Gambar 2 adalah contoh beberapa data yang ada dalam dataset.

TABEL 1. Struktur Dataset

ID	Nama Produk	Deskripsi	Harga	Berat	Kondisi	Status	Kategori	URL Kategori	Sold Count	View Count	Rating
110296338456	ONLINE EXCLUSIVE VIVO Y100	Vivo Y100 Smartphone	3.077.000	0.6 kg	Baru	Aktif	Handphone	tokopedia.com	17rb+	988	4.9
110429538485	Paket Bundle Your Coffee	Kopi & Tumbler	117.000	0.5 kg	Baru	Aktif	Kopi	tokopedia.com	1rb+	398	4.9
110264351	Sweaty Bronze Dry T-Shirt	Kaos Olahraga	110.000	0.2 kg	Baru	Aktif	Pakaian	tokopedia.com	3rb+	684	4.9

Pada penelitian ini, atribut yang digunakan sebagai fitur utama untuk sistem rekomendasi adalah *name*, *description*, *category_name*, dan *shop_name*. Keempat atribut ini dipilih karena mengandung informasi tekstual yang relevan untuk menganalisis karakteristik konten produk dan mengidentifikasi kesamaan antar produk menggunakan metode *Content-Based Filtering*.

1. Pre-processing

Pada penelitian ini, proses pre-processing dilakukan melalui 5 tahapan berurutan yang diimplementasikan menggunakan library *Natural.js* dan *Stopword*. Proses tersebut meliputi:

A. Feature Selection

Tahap pertama adalah pemilihan fitur yang relevan dari dataset. Dari total 36 atribut yang tersedia dalam dataset, dipilih 4 atribut utama yang mengandung informasi tekstual paling relevan untuk analisis konten produk. Atribut-atribut tersebut adalah *name* yang merepresentasikan nama produk, *description* yang berisi deskripsi detail produk, *category_name* yang menunjukkan kategori produk, dan *shop_name* yang mengidentifikasi nama toko penjual. Keempat atribut ini digabungkan menjadi satu teks yang merepresentasikan karakteristik konten produk.

B. Cleaning & Normalization

Tahap cleaning dan normalisasi dilakukan untuk membersihkan data dari elemen-elemen yang tidak diperlukan dalam analisis teks. Proses ini dilakukan secara berurutan melalui beberapa tahapan. Pertama, dilakukan penghapusan *html tags* untuk menghilangkan elemen markup html. Kedua, penghapusan emoji dilakukan dengan mendeteksi dan menghilangkan karakter emoji unicode yang tidak relevan untuk analisis konten tekstual. Ketiga, penghapusan simbol dan karakter seperti tanda baca dan simbol, hanya menyisakan huruf, angka, dan spasi. Keempat, normalisasi spasi dilakukan dengan mengganti multiple spaces menjadi single space untuk konsistensi format. Kelima, case folding diterapkan dengan mengkonversi seluruh teks menjadi huruf kecil (*lowercase*) untuk memastikan konsistensi dalam pencocokan term dan menghindari duplikasi term yang sama dengan kapitalisasi berbeda.

C. Tokenization

Pada penelitian ini, tokenisasi dilakukan untuk memecah teks menjadi unit-unit kata individual. Tokenisasi ini memisahkan teks berdasarkan delimiter seperti spasi dan tanda baca dengan pemisahan kata. Hasil dari proses tokenisasi adalah sebuah *array of strings* yang berisi kata-kata individual yang telah terpisah yang siap untuk tahap preprocessing selanjutnya.

D. Stopword Removal

Stopword removal adalah proses penghapusan kata-kata umum seperti “yang”, “dan”, “di”, “the”, “a”, “is” yang tidak memberikan nilai diskriminatif. Pada penelitian ini, penghapusan stopwords dilakukan menggunakan kombinasi daftar stopwords bahasa Indonesia dan bahasa Inggris mengingat dataset produk *e-commerce* mengandung istilah dari kedua bahasa. Setiap token dibandingkan dengan kumpulan stopwords yang telah didefinisikan dan dihilangkan jika termasuk di dalamnya. Hasil akhir adalah array token bermakna dengan nilai informasi tinggi untuk analisis konten dan mengurangi dimensi vektor.

E. Stemming

Stemming merupakan tahap akhir dari pre-processing dalam penelitian ini yang mentransformasi kata menjadi bentuk dasar dengan menghilangkan imbuhan bahasa Indonesia seperti “me-”, “ber-”, “di-”, “-kan”, “-an”, dan kombinasinya. Proses ini memastikan variasi bentuk kata yang secara semantik sama dipetakan ke satu term konsisten. Sebagai contoh, “membeli”, “dibeli”, “pembelian”, dan “belian” akan dikonversi menjadi “beli”. Hasil akhir dari seluruh tahapan pre-processing adalah array of string tokens yang telah dibersihkan, dinormalisasi, dan siap untuk vektorisasi *TF-IDF*.

Vektorisasi TF-IDF

Pada bagian ini akan diuraikan contoh perhitungan dari algoritma *TF-IDF* dan kemudian pada tahapan berikutnya proses *Cosine Similarity* melalui perhitungan manual menggunakan data sampel produk *e-commerce*. Tahapan cara kerja model dan algoritma ini dilakukan secara sistematis mulai dari tahap pemilihan data sampel, tahap preprocessing data, vektorisasi menggunakan *TF-IDF*, hingga pengukuran kemiripan menggunakan *Cosine Similarity* untuk menghasilkan rekomendasi produk yang relevan.

Data sampel yang digunakan dalam perhitungan ini terdiri dari 5 item produk *e-commerce* yang dipilih dari dataset *Tokopedia Products 2025*. Informasi mengenai atribut yang digunakan dalam data sampel dapat dilihat pada TABEL 2.

TABEL 2. Data Sampel

	<i>Nama</i>	<i>Deskripsi</i>	<i>Kategori</i>	<i>Toko</i>
D1	Teko Listrik Stainless 2L Pemanas Air Kettle Listrik Electric Kettle	Electric Kettle - Teko Listrik Model : Electric Kettle / Teko Listrik Kapasitas : 1,8L Daya : 500 Watt Voltage : 220-240V Kelengkapan : 1 Teko x 1 Induksi x 1 Manual x 1 Electric Kettle - Teko Listrik yang sangat berguna untuk Rumah Tangga Anda memudahkan dan dengan waktu yang cepat untuk memasak air Dengan Tenaga Listrik sangat memudahkan Anda untuk memperoleh air panas saat dibutuhkan tanpa kompor Dilengkapi Overheat protector sangat aman digunakan	Teko Listrik	AchillesCollection
D2	Tolak Angin Madu Sido Muncul Box isi 12 sachet Tolak angin	Tolak Angin Madu Sido Muncul Box isi 12 Sachet adalah produk herbal yang mengandung bahan-bahan alami seperti ekstrak jahe, daun mint, dan madu, yang dirancang untuk membantu menjaga kesehatan tubuh dan meningkatkan daya tahan tubuh. Dengan tambahan madu, produk ini memiliki rasa yang lebih manis dan menyegarkan, membuatnya lebih mudah diminum. Setiap box berisi 12 sachet yang praktis, sehingga memudahkan Anda untuk mengonsumsinya kapan saja, baik di rumah, di kantor, atau saat bepergian.	Multivitamin	Tobaking Murah Online
D3	Kipas Angin Sekai HDO 615-S	Type HDO 615-S - Kipas angin high velocity 15 cm (6") - Baling baling dari Plastik - 2 pengatur kecepatan - Hembusan angin kencang - Hemat listrik - Tidak berisik - Dilengkapi sekering pengaman - Garansi 3 tahun - Daya 20 W, 220 VAC, 50 Hz	Kipas Angin Listrik	amac electronic store
D4	Kopi LIONG BULAN Plus Gula 28 GR (Isi 20 Sachet)	HALAL - Isi : 20 Sachet (20 x 28 Gr) Kopi Liong Bulan Plus Gula - Terbuat dari 100% Biji kopi murni dan gula tanpa campuran dan pengawet - THE LOCAL LEGEND - Kopi LIONG BULAN merupakan merek kopi asal Bogor yang sudah hadir sejak tahun 1945 - Cita RASA yang UNIK serta AROMA nya yang KHAS membuat Kopi LIONG BULAN menjadi cerita dari mulut ke mulut secara turun temurun - Barang selalu fresh, produksi terbaru dari	Kopi Kemasan	Kopi Liong Bulan

		pabrik.\nExp Date / Kadaluarsa : 3 bulan sejak tanggal pemesanan (Tercantum di kode produksi di belakang)\n		
D5	saringan kopi teh portable V60 / coffee filter stainless double mesh	saringan kopi portable double mesh bahan terbuat dari stainless 304 dengan pinggiran ABS\ntidak perlu menggunakan kertas lagi\nBrand: vesper\ndimensi\ntinggi 7,5cm\ndiameter atas 6,5cm\ndiameter atas dengan kuping 12cm\npremium coffee filter\n	Saringan Masak	Vesperindo _NEW

Dari data sampel yang terdapat pada TABEL 2 dilakukan tahap pre-processing yaitu tahap feature selection, cleaning & normalization, tokenization, stopwords removal, dan stemming yang kemudian menghasilkan kumpulan token yang bersih, konsisten, dan relevan untuk analisis konten produk. Hasil pre-processing dapat dilihat pada TABEL 3.

TABEL 3. Token Data Sampel

	Tokens
D1	[listrik, teko, kettle, 1, electric, air, x, panas, lengkap, mudah, 220, 500, stainless, 2l, model, kapasitas, 8l, daya, watt, voltase, 240v, induksi, manual, guna, rumah, tangga, cepat, masak, tenaga, oleh, butuh, kompor, overheat, protector, aman, achillescollection]
D2	[madu, 12, box, isi, sachet, tolak, angin, sido, muncul, produk, bahan, tubuh, mudah, tolak angin, herbal, kandung, alami, ekstrak, jahe, daun, mint, rancang, bantu, jaga, sehat, tingkat, tahan, tambah, milik, manis, segar, buat, minum, praktis, konsumsi, baik, kantor, pergi, multivitamin, tobaking, murah, online, daya, rumah]
D3	[kipas, angin, 615, hdo, s, baling, listrik, 2, 3, 6, 15, 20, 50, seka, type, high, velocity, cm, plastik, atur, hembus, kencang, hemat, berisik, sekring, kam, garansi, w, vac, hz, amac, electronic, store, 220, cepat, lengkap, daya]
D4	[kopi, liong, gula, 20, 28, plus, gr, mulut, turun, produksi, isi, sachet, 100, 1945, halal, biji, murni, campur, awet, local, legend, merek, bogor, hadir, cita, unik, aroma, nya, khas, cerita, barang, fresh, baru, pabrik, exp, date, kadaluarsa, tanggal, mesan, cantum, kode, kemas, 3, x, buat]
D5	[saring, portable, coffee, filter, double, mesh, 5cm, kopi, stainless, 7, 304, teh, v60, pinggir, abs, kertas, brand, vesper, dimensi, diameter, diamater, kuping, 12cm, premium, vesperindo, new, 6, bahan, masak, buat]

Setelah melewati tahapan preprocessing, selanjutnya dilakukan konversi/ekstraksi fitur menggunakan algoritma *TF-IDF* untuk mengkonversi token-token teks menjadi representasi vektor numerik. Proses ini dimulai dengan menghitung frekuensi kemunculan setiap term dalam dokumen (*Term Frequency/TF*), kemudian mengalikannya dengan nilai kebalikan frekuensi dokumen (*Inverse Document Frequency/IDF*) untuk menghasilkan bobot *TF-IDF*. Bobot *TF-IDF* ini mencerminkan tingkat kepentingan setiap term dalam merepresentasikan karakteristik unik suatu produk terhadap keseluruhan dataset.

Langkah pertama dalam *TF-IDF* adalah dengan mengukur *Term Frequency (TF)*, yakni mengukur seberapa sering suatu *term* muncul dalam dokumen tertentu menggunakan rumus pada persamaan 1. Dari rumus tersebut, berikut adalah contoh perhitungan *TF* untuk term “daya”.

Term *daya* muncul 1 kali dalam *D1*:

$$TF(daya, D1) = 1$$

Term *daya* muncul 1 kali dalam *D2*:

$$TF(daya, D2) = 1$$

Lakukan proses perhitungan yang sama untuk memperoleh $TF(daya, D3)$, $TF(daya, D4)$, dan $TF(daya, D5)$. Selanjutnya adalah perhitungan *Inverse Document Frequency (IDF)* yang menggunakan persamaan 2.

$$IDF(daya) = \log_{10} \left(\frac{5}{3} \right) + 1$$

$$IDF(daya) = \log_{10}(1.667) + 1$$

$$IDF(daya) = 0.2218 + 1 = 1.2218$$

Kemudian menghitung *TF-IDF* dengan menggunakan persamaan 3 maka perhitungan *TF-IDF* term *daya* pada dokumen *D1*, *D2*, *D3*, *D4*, dan *D5* adalah sebagai berikut:

$$TF-IDF(daya, D1) = 1 \times 1.2218 = 1.2218$$

Lakukan proses perhitungan yang sama untuk memperoleh $TF-IDF(daya, D2)$ sampai $TF-IDF(daya, D5)$.

Dari hasil perhitungan *TF-IDF* semua term pada dokumen *D1*, *D2*, *D3*, *D4*, dan *D5*. Hasil perhitungan *TF-IDF* untuk seluruh term yang muncul dalam kelima dokumen sampel disajikan pada TABEL 4.

TABEL 4. Vector Hasil Perhitungan TF-IDF

[illegible]

Cosine Similarity

Setelah menghitung *TF-IDF*, selanjutnya ialah menghitung similaritas antar sudut kosinus vector *TF-IDF* menggunakan algoritma *Cosine Similarity*. Namun sebelum itu, disini pada penelitian ini kami hanya menggunakan beberapa data vector hasil perhitungan *TF-IDF* sebagai sampel perhitungan agar mempermudah ilustrasi perhitungan *Cosine Similarity*. Data vector hasil perhitungan *TF-IDF* yang kami gunakan dapat dilihat pada TABEL 5.

TABEL 5. Matrix Bobot Sampel TF-IDF

<i>Term</i>	<i>D1</i>	<i>D2</i>	<i>D3</i>	<i>D4</i>	<i>D5</i>
<i>listrik</i>	10.5758	0	3.0217	0	0
<i>buat</i>	0	1.2231	0	1.2231	1.2231
<i>panas</i>	3.8326	0	0	0	0
<i>lengkap</i>	3.0217	0	1.5108	0	0
<i>mudah</i>	3.0217	3.0217	0	0	0
<i>sachet</i>	0	4.5325	0	3.0217	0
<i>aman</i>	1.9163	0	0	0	0
<i>stainless</i>	1.5108	0	0	0	3.0217
<i>model</i>	1.9163	0	0	0	0
<i>kapasitas</i>	1.9163	0	0	0	0
<i>daya</i>	1.2231	1.2231	1.2231	0	0

Pada contoh perhitungan berikut ini akan menghitung similaritas antar kedua dokumen menggunakan *Cosine Similarity* untuk dokumen $D1$ dengan dokumen $D3$. Langkah awal ialah menghitung *Dot Product* $D1 \cdot D3$, Dot product dihitung dengan mengalikan nilai *TF-IDF* setiap term yang bersinggungan pada kedua dokumen, kemudian menjumlahkan seluruh hasil perkalian.

$$\begin{aligned}
D1 \cdot D3 &= (D1_1 \times D3_1) + (D1_2 \times D3_2) + \dots + (D1_n \times D3_n) \\
D1 \cdot D3 &= (10.5758 \times 3.0217) + (0 \times 0) + (3.8326 \times 0) + (3.0217 \times 1.5108) + \\
&\quad (3.0217 \times 0) + (0 \times 0) + (1.9163 \times 0) + (1.5108 \times 0) + \\
&\quad (1.9163 \times 0) + (1.9163 \times 0) + (1.2231 \times 1.2231) \\
D1 \cdot D3 &= 31.9620 + 0 + 0 + 4.5644 + 0 + 0 + 0 + 0 + 0 + 0 + 1.4960 \\
D1 \cdot D3 &= \mathbf{38.0224}
\end{aligned}$$

Kemudian dilanjutkan dengan menghitung magnitude vektor $|D1|$ yang dimana dihitung dengan akar kuadrat dari jumlah kuadrat setiap elemen dengan menggunakan persamaan 4.

$$\begin{aligned} ||D1|| &= \sqrt{10.5758^2 + 0^2 + 3.8326^2 + 3.0217^2 + 3.0217^2 + 0^2 + 1.9163^2 + 1.5108^2 + 1.9163^2 + 1.2231^2} \\ ||D1|| &= \sqrt{111.8475 + 0 + 14.6888 + 9.1307 + 9.1307 + 0 + 3.6722 + 2.2825 + 3.6722 + 1.4960} \\ ||D1|| &= \sqrt{159.5928} = \mathbf{12.6332} \end{aligned}$$

Lalu dilanjutkan dengan menghitung magnitude vektor $|D3|$ yang dimana dihitung dengan cara yang sama seperti menghitung magnitude vektor $|D1|$:

$$\begin{aligned} ||D3|| &= \sqrt{3.0217^2 + 0^2 + 0^2 + 1.5108^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 1.2231^2} \\ ||D3|| &= \sqrt{9.1307 + 0 + 0 + 2.2825 + 0 + 0 + 0 + 0 + 0 + 1.4960} = \sqrt{12.9092} = \mathbf{3.5930} \end{aligned}$$

Setelah mendapatkan nilai dot product dan magnitude kedua vektor, dapat dilanjutkan dengan menghitung *Cosine Similarity*:

$$\begin{aligned} \text{Cos}(D1, D3) &= \frac{D1 \cdot D3}{||D1|| \times ||D3||} \\ \text{Cos}(D1, D3) &= \frac{38.0224}{12.6332 \times 3.5930} \\ \text{Cos}(D1, D3) &= \mathbf{0.8377} \end{aligned}$$

Dari hasil perhitungan nilai *Cosine Similarity* diatas antara produk $D1$ dan produk $D3$ dapat disimpulkan bahwa kemiripannya adalah sebesar **0.8377** atau **83.77%** similar/mirip. Nilai tersebut mengindikasikan bahwa kedua produk memiliki tingkat kemiripan yang sangat tinggi $>80\%$. Hal ini dapat dijelaskan karena keduanya merupakan produk elektronik yang menggunakan listrik, sehingga berbagi term-term yang sama seperti “*listrik*”, “*lengkap*”, dan “*daya*” yang berkontribusi signifikan terhadap nilai similarity.

TABEL 6 merupakan hasil lengkap dari nilai *Cosine Similarity* antar semua dokumen dalam sampel yang digunakan.

TABEL 6. Matrix Similarity

	<i>D1</i>	<i>D2</i>	<i>D3</i>	<i>D4</i>	<i>D5</i>
<i>D1</i>	1.0000	0.1472	0.8377	0.0000	0.1201
<i>D2</i>	0.1472	1.0000	0.1030	0.8154	0.0000
<i>D3</i>	0.8377	0.1030	1.0000	0.0000	0.0000
<i>D4</i>	0.0000	0.8154	0.0000	1.0000	0.0000
<i>D5</i>	0.1201	0.0000	0.0000	0.0000	1.0000

Rekomendasi Top-K Product

Tahapan ini dilakukan untuk memvalidasi efektivitas algoritma *Content-Based Filtering* yang di implementasikan dengan variasi parameter *Top-K* dengan menggunakan dataset beberapa item produk. Evaluasi menggunakan 50 query produk yang disampling dari berbagai kategori dengan similarity threshold 0.85 untuk menentukan relevansi. Gambar 2 dibawah ini menunjukkan hasil evaluasi komprehensif sistem rekomendasi pada semua variasi sample size dan nilai K.

SUPPORT EVALUATE SIMPLER SAMPLE SIZE							
	SampleSize	K	Precision	Recall	FL-Score	MAP	2003
0	10	10	98.00	26.30	51.44	38.38	2.033
1	10	10	97.33	51.37	63.48	51.37	2.033
2	10	10	98.00	61.62	78.05	61.62	1.921
3	10	10	98.67	82.27	87.88	82.27	1.878
4	10	10	97.68	100.00	88.73	100.00	1.219
5	10	10	98.00	100.00	90.00	100.00	0.728
6	10	10	98.00	58.65	41.84	58.65	1.289
7	10	10	97.33	33.33	42.63	33.33	0.688
8	10	10	97.33	67.97	58.85	67.97	1.087
9	10	10	98.00	100.00	62.63	100.00	0.853
10	10	10	98.00	98.98	50.27	98.98	1.285
11	10	10	97.33	99.35	54.42	99.35	0.866
12	10	10	97.33	99.35	54.42	99.35	0.866
13	10	10	97.33	57.85	64.29	57.85	2.356
14	10	10	97.33	68.68	67.29	68.68	1.933
15	10	10	97.33	79.00	81.33	81.33	1.763
16	10	10	97.33	97.33	81.33	97.33	1.000
17	10	10	98.00	97.86	47.72	97.86	0.938
18	10	10	98.00	98.98	65.44	98.98	1.285
19	10	10	98.00	98.98	65.44	98.98	1.285
20	10	10	98.00	72.13	67.47	72.13	2.281
21	10	10	99.11	69.11	68.24	69.11	2.080
22	10	10	94.87	91.80	58.87	91.80	1.044
23	10	10	97.87	97.87	58.87	97.87	1.044
24	10	10	92.75	91.80	50.27	91.80	2.338
25	10	10	92.75	91.80	50.27	91.80	2.338
26	10	10	92.75	91.80	50.27	91.80	2.338
27	10	10	92.75	91.80	50.27	91.80	2.338
28	10	10	92.75	91.80	50.27	91.80	2.338
29	10	10	92.75	91.80	50.27	91.80	2.338
30	10	10	92.75	91.80	50.27	91.80	2.338
31	10	10	92.75	91.80	50.27	91.80	2.338
32	10	10	92.75	91.80	50.27	91.80	2.338
33	10	10	92.75	91.80	50.27	91.80	2.338
34	10	10	92.75	91.80	50.27	91.80	2.338
35	10	10	92.75	91.80	50.27	91.80	2.338
36	10	10	92.75	91.80	50.27	91.80	2.338
37	10	10	92.75	91.80	50.27	91.80	2.338
38	10	10	92.75	91.80	50.27	91.80	2.338
39	10	10	92.75	91.80	50.27	91.80	2.338
40	10	10	92.75	91.80	50.27	91.80	2.338
41	10	10	92.75	91.80	50.27	91.80	2.338
42	10	10	92.75	91.80	50.27	91.80	2.338
43	10	10	92.75	91.80	50.27	91.80	2.338
44	10	10	92.75	91.80	50.27	91.80	2.338
45	10	10	92.75	91.80	50.27	91.80	2.338
46	10	10	92.75	91.80	50.27	91.80	2.338
47	10	10	92.75	91.80	50.27	91.80	2.338
48	10	10	92.75	91.80	50.27	91.80	2.338
49	10	10	92.75	91.80	50.27	91.80	2.338
50	10	10	92.75	91.80	50.27	91.80	2.338
51	10	10	92.75	91.80	50.27	91.80	2.338
52	10	10	92.75	91.80	50.27	91.80	2.338
53	10	10	92.75	91.80	50.27	91.80	2.338
54	10	10	92.75	91.80	50.27	91.80	2.338
55	10	10	92.75	91.80	50.27	91.80	2.338
56	10	10	92.75	91.80	50.27	91.80	2.338
57	10	10	92.75	91.80	50.27	91.80	2.338
58	10	10	92.75	91.80	50.27	91.80	2.338
59	10	10	92.75	91.80	50.27	91.80	2.338
60	10	10	92.75	91.80	50.27	91.80	2.338
61	10	10	92.75	91.80	50.27	91.80	2.338
62	10	10	92.75	91.80	50.27	91.80	2.338
63	10	10	92.75	91.80	50.27	91.80	2.338
64	10	10	92.75	91.80	50.27	91.80	2.338
65	10	10	92.75	91.80	50.27	91.80	2.338
66	10	10	92.75	91.80	50.27	91.80	2.338
67	10	10	92.75	91.80	50.27	91.80	2.338
68	10	10	92.75	91.80	50.27	91.80	2.338
69	10	10	92.75	91.80	50.27	91.80	2.338
70	10	10	92.75	91.80	50.27	91.80	2.338
71	10	10	92.75	91.80	50.27	91.80	2.338
72	10	10	92.75	91.80	50.27	91.80	2.338
73	10	10	92.75	91.80	50.27	91.80	2.338
74	10	10	92.75	91.80	50.27	91.80	2.338
75	10	10	92.75	91.80	50.27	91.80	2.338
76	10	10	92.75	91.80	50.27	91.80	2.338
77	10	10	92.75	91.80	50.27	91.80	2.338
78	10	10	92.75	91.80	50.27	91.80	2.338
79	10	10	92.75	91.80	50.27	91.80	2.338
80	10	10	92.75	91.80	50.27	91.80	2.338
81	10	10	92.75	91.80	50.27	91.80	2.338
82	10	10	92.75	91.80	50.27	91.80	2.338
83	10	10	92.75	91.80	50.27	91.80	2.338
84	10	10	92.75	91.80	50.27	91.80	2.338
85	10	10	92.75	91.80	50.27	91.80	2.338
86	10	10	92.75	91.80	50.27	91.80	2.338
87	10	10	92.75	91.80	50.27	91.80	2.338
88	10	10	92.75	91.80	50.27	91.80	2.338
89	10	10	92.75	91.80	50.27	91.80	2.338
90	10	10	92.75	91.80	50.27	91.80	2.338
91	10	10	92.75	91.80	50.27	91.80	2.338
92	10	10	92.75	91.80	50.27	91.80	2.338
93	10	10	92.75	91.80	50.27	91.80	2.338
94	10	10	92.75	91.80	50.27	91.80	2.338
95	10	10	92.75	91.80	50.27	91.80	2.338
96	10	10	92.75	91.80	50.27	91.80	2.338
97	10	10	92.75	91.80	50.27	91.80	2.338
98	10	10	92.75	91.80	50.27	91.80	2.338
99	10	10	92.75	91.80	50.27	91.80	2.338

GAMBAR 2. Hasil Top-K

Berdasarkan hasil evaluasi pada berbagai sample size dan nilai K yang terlihat seperti pada Gambar 2, sistem rekomendasi yang diimplementasikan menunjukkan performa yang konsisten dan robust untuk produk e-commerce. Untuk implementasi praktis, nilai $K=20$ hingga $K=30$ direkomendasikan sebagai ukuran optimal yang memberikan keseimbangan antara kualitas rekomendasi (precision 70.67%-87.50%) dan coverage produk relevan (recall 71.12%-83.72%) pada sample size representatif. Precision tinggi pada $K=10$ (rentang 88.40%-98.00% untuk sample size 20-50) mengindikasikan bahwa pengguna akan mendapatkan rekomendasi yang sangat relevan pada posisi teratas, meningkatkan user experience dan potensi konversi pembelian. Recall yang mencapai 90.65%-92.58% pada $K=50$ menunjukkan bahwa sistem mampu mengidentifikasi hampir seluruh produk relevan, memberikan variasi pilihan yang cukup bagi pengguna tanpa mengorbankan relevansi secara signifikan. Konsistensi performa pada berbagai sample size memvalidasi robustness algoritma *Content-Based Filtering* dengan *TF-IDF* dan *Cosine Similarity* untuk kategori produk item e-commerce pada penelitian ini.

Evaluasi

Evaluasi sistem rekomendasi dilakukan menggunakan *metric calculation/confussion matrix* standar yang umum digunakan dalam penelitian sistem rekomendasi yaitu *Precision@K*, *Recall@K*, dan *F1-Score@K*. Proses evaluasi diimplementasikan dengan menggunakan bahasa pemrograman Python yang instruksi evaluasi tertampil pada Gambar 3. Pada Gambar 3 terlihat pengaturan nilai K yang digunakan dalam evaluasi adalah 10, 15, 20, 30, 50, dan 100 dengan menggunakan threshold sebesar 0,05.

```

2 // @ts-ignore
3 import { getProducts, type ProductWithTranslations } from "src/lib/data";
4 import { recommender } from "../lib/recommender";
5
6 // --- Parameter defaults ---
7 const K_SIZE = 10; // 10, 20, 30, 50, 100;
8 const THRESHOLD = 0.5;
9 const SAMPLE_SIZE = [5, 10, 20, 30, 50, 60, 70, 80, 90, 100];
10 const VERSION = "v1.0.0";
11
12 // --- Utilities ---
13
14 function precision(recommended: number[], relevant: number[]): number {
15     const hit = recommended.filter((id) => id === relevant).length;
16     return hit / recommended.length > 0 ? hit / recommended.length : 0;
17 }
18
19 function recall(recommended: number[], relevant: number[]): number {
20     const hit = recommended.filter((id) => id === relevant).length;
21     return hit / relevant.length > 0 ? hit / relevant.length : 0;
22 }
23
24 function f1p(number: number, precision: number, recall: number): number {
25     return 2 * precision * recall / (precision + recall);
26 }
27
28 function mean(arr: number[]): number {
29     return arr.length > 0 ? arr.reduce((a, b) => a + b) / arr.length : 0;
30 }
31
32 function meanAveragePrecision(recommended: number[], relevant: number[]): number {
33     let hit = 0;
34     sum = 0;
35     recommended.forEach((id, idx) => {
36         if (relevant.includes(id)) {
37             hit++;
38             sum = hit / (idx + 1);
39         }
40     });
41     return relevant.length > 0 ? sum / relevant.length : 0;
42 }

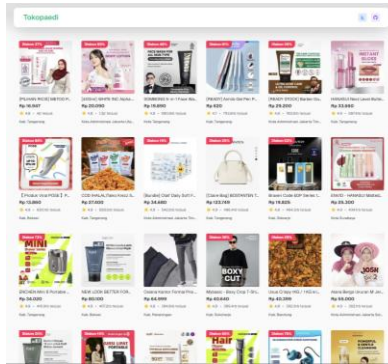
```

GAMBAR 3. Evaluation File

Proses evaluasi dilakukan dengan menjalankan script evaluasi yang mengiterasi seluruh query produk dalam ground truth, menghitung rekomendasi untuk setiap query, dan mengakumulasi metrik evaluasi. Evaluasi dilakukan untuk berbagai nilai K (5, 10, 15, 20, 30, ..., 100) untuk menganalisis performa sistem pada berbagai ukuran rekomendasi.

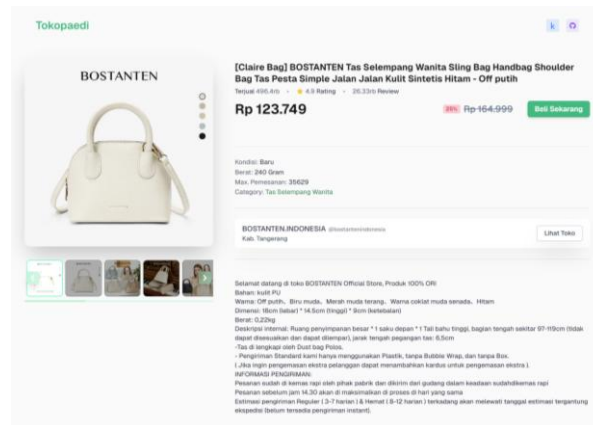
Hasil Implementasi

Hasil implementasi dari penelitian ini ialah berupa halaman website yang terdiri dari homepage, toko, product detail serta rekomendasinya. Gambar 4 menampilkan halaman homepage website yang berisi daftar produk atau berbagai macam produk. Produk-produk yang tampil pada Gambar 4 dapat dipilih oleh pengguna dengan bebas.



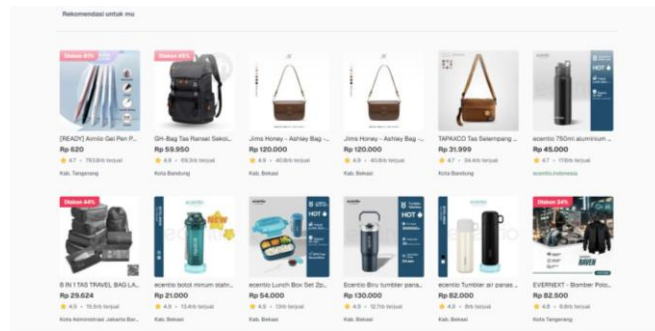
GAMBAR 4. Halaman Homepage

Ketika seorang pengguna memilih salah satu produk yang ditampilkan pada halaman homepage (Gambar 4), maka akan tampil halaman detail produk yang dipilih. Misal pengguna memilih produk Claire Bag dengan merek BOSTANTEN maka detail produk akan tampil seperti terlihat pada Gambar 5.



GAMBAR 5. Detail Produk

Pada Gambar 5, jika pengguna scrolling kebawah maka akan menampilkan hasil rekomendasi produk berdasarkan tingkat kemiripan produk terhadap produk referensinya (dalam case/kasus/contoh ini adalah jenis Claire Bag), mulai dari tingkat kemiripan yang tertinggi ke yang terendah.



GAMBAR 6. Halaman Detail Produk Section Rekomendasi Produk

Gambar 6 merupakan hasil rekomendasi dari hasil aksi yang dilakukan pada Gambar 5. Pada Gambar 5 dipilih produk berjenis Claire Bag dengan merek BOSTANTEN kemudian hasil produk lain yang direkomendasikan kepada user dapat terlihat pada Gambar 6 yang urutan produk rekomendasi nya telah diurutkan dari yang memiliki nilai kemiripan tertinggi sampai ke rendah.

KESIMPULAN

Penelitian ini berhasil mengimplementasikan sistem rekomendasi produk e-commerce menggunakan metode *Content-Based Filtering* dengan pendekatan *TF-IDF* dan *Cosine Similarity* pada dataset produk *Tokopedia* yang terdiri dari 1.199 record. Pada hasil evaluasi menggunakan berbagai ukuran sampel query dan nilai K yang terdiri dari sample size 50 query yang paling representatif, sistem mencapai *Precision@10* sebesar 88.40%, *Recall@10* sebesar 46.81%, dan *F1-Score@10* sebesar 50.88%, dengan nilai *F1-Score* optimal dicapai pada $K=30$ sebesar 66.79%. Analisis variasi sample size memvalidasi stabilitas algoritma, dimana sample size yang lebih besar (20-50 query) menghasilkan metrik evaluasi yang lebih representatif dengan variance rendah, dengan precision tetap berada di atas 70% pada $K=30$ untuk semua sample size representatif dan recall mencapai 90.65%-95.96% pada $K=50-100$, membuktikan bahwa algoritma *Cosine Similarity* berhasil mengurutkan produk berdasarkan tingkat kemiripan dengan akurat. Berdasarkan hasil evaluasi, nilai $K=20$ hingga $K=30$ direkomendasikan sebagai ukuran optimal untuk implementasi praktis pada aplikasi e-commerce, memberikan keseimbangan antara kualitas rekomendasi (precision 70.67%-80.70%) dan coverage produk relevan (recall 71.12%-83.72%).

TINJAUAN PUSTAKA

1. D. Pebrianti, D. Ahmad, L. Bayuaji, L. Wijayanti, and M. Mulyadi, "Using Content-Based Filtering and Apriori for Recommendation Systems in a Smart Shopping System," *Indonesian Journal of Computing, Engineering, and Design (IJoCED)*, vol. 6, no. 1, pp. 58–70, Apr. 2024, doi: 10.35806/ijoced.v6i1.393.
2. A. Sulami, V. Atina, and N. Nurmawati, "Penerapan Metode Content Based Filtering dalam Sistem Rekomendasi Pemilihan Produk Skincare," *Satuan Tulisan Riset dan Inovasi Teknologi (STRING)*, vol. 9, no. 2, p. 172, Dec. 2024, doi: 10.30998/string.v9i2.24066.
3. N. P. T. Widiyanti, I. M. A. D. Suarjaya, and N. K. D. Rusjajanti, "Food Recipe Recommendation System with Content-Based Filtering and Collaborative Filtering Methods," *Jurnal Dan Penelitian Teknik Informatika (Sinkron)*, vol. 9, no. 3, pp. 1167–1776, Aug. 2025, doi: 10.33395/sinkron.v9i3.14778.
4. I. A. W. Nandita, I. M. A. Dwi Suarjaya, and I. P. A. Bayupati, "News Recommendation System Using Content-Based Filtering through RSS Customization Service," *Journal of Applied Informatics and Computing (JAIC)*, vol. 9, no. 4, pp. 1858–1867, Aug. 2025, doi: 10.30871/jaic.v9i4.9807.
5. C. Salsabilla and D. W. Utomo, "Pekalongan Regency Tourism Recommendation System with Content based Filtering," *Jurnal Sistem Informasi (SISTEMASI)*, vol. 14, no. 1, p. 262, Jan. 2025, doi: 10.32520/stmsi.v14i1.4839.

6. I. Yulfihani and M. Zakariyah, “Optimization of Tourism Destination Recommendations in Batang Regency Using Content-Based Filtering,” *Journal of Applied Informatics and Computing (JAIC)*, vol. 8, no. 2, pp. 499–508, Nov. 2024, doi: 10.30871/jaic.v8i2.8618.
7. A. H. J. P. Juni Permana and Agung Toto Wibowo, “Movie Recommendation System Based on Synopsis Using Content-Based Filtering with TF-IDF and Cosine Similarity,” *International Journal on Information and Communication Technology (IJoICT)*, vol. 9, no. 2, pp. 1–14, Dec. 2023, doi: 10.21108/ijoi.v9i2.747.
8. I. Ilhamsyah, S. Rahmayuda, D. M. Midyanti, and S. Martha, “Pemodelan Sistem Rekomendasi Restoran berdasarkan Preferensi Pengguna dengan Pendekatan Content-Based Filtering,” *Jurnal Edukasi dan Penelitian Informatika (JEPIN)*, vol. 10, no. 1, p. 154, Apr. 2024, doi: 10.26418/jp.v10i1.74008.
9. S. R. Cholil, N. A. Rizki, and T. F. Hanifah, “Sistem Rekomendasi Tempat Wisata Di Kota Semarang Menggunakan Metode Collaborative Filtering,” *Jurnal Informatika dan Komputer (JIKO)*, vol. 7, no. 1, p. 118, Feb. 2023, doi: 10.26798/jiko.v7i1.727.
10. F. Maulana and E. B. Setiawan, “Performance of Deep Feed-Forward Neural Network Algorithm Based on Content-Based Filtering Approach,” *Jurnal Ilmiah Penelitian dan Penerapan Teknologi Sistem Informasi (INTENSIF)*, vol. 8, no. 2, pp. 278–294, Aug. 2024, doi: 10.29407/intensif.v8i2.22904.