

Komparasi Algoritma K-Means dan *Hierarchical Clustering* Untuk Mengetahui Data Customer Dalam Layanan Internet

Afri Yudha¹, Shofi Khaerunnisa², Timor Setyaningsih³

^{1,2,3}Program Studi Teknologi Informasi Universitas Darma Persada

ibnugazali@gmail.com

Abstract

This study aims to compare two popular clustering algorithms, namely K-Means and Hierarchical Clustering, in grouping customer data of internet services. The data used includes information such as date, purchase order number, type of internet, product description, quantity in Mbps, region, city, and rate in IDR. Clustering is a technique in data mining used to group data based on similar characteristics, helping companies understand usage patterns, market segmentation, and improve services to customers. In this study, both algorithms were tested using the Sum of Squared Errors (SSE) metric to evaluate clustering performance. The results showed that the K-Means algorithm outperformed Hierarchical Clustering in terms of efficiency and accuracy when evaluating SSE values. K-Means was able to cluster data more quickly and produce clearer clusters, while Hierarchical Clustering required longer computation time and yielded less optimal clustering results in some cases. In conclusion, the K-Means algorithm is more suitable for analyzing customer data of internet services in this context due to its advantages in speed and accuracy of clustering based on SSE results. This study provides guidance for companies in choosing the appropriate clustering algorithm to group customer data to enhance business strategies and service quality.

Keywords: K-Means, Hierarchical Clustering, clustering, customer data, internet services, data mining, market segmentation, Sum of Squared Errors (SSE).

Abstrak

Penelitian ini bertujuan untuk membandingkan dua algoritma clustering yang populer, yaitu K-Means dan Hierarchical Clustering, dalam mengelompokkan data pelanggan layanan internet. Data yang digunakan mencakup informasi seperti tanggal, nomor pesanan, jenis internet, deskripsi produk, kuantitas Mbps, wilayah, kota, dan tarif IDR. Clustering adalah teknik dalam data mining yang digunakan untuk mengelompokkan data berdasarkan karakteristik yang serupa, membantu perusahaan dalam memahami pola penggunaan, segmentasi pasar, dan meningkatkan layanan kepada pelanggan. Pada penelitian ini, kedua algoritma diuji menggunakan metrik *Sum of Squared Errors (SSE)* untuk mengevaluasi kinerja clustering. Hasil penelitian menunjukkan bahwa algoritma K-Means lebih unggul dibandingkan dengan Hierarchical Clustering dalam hal efisiensi dan akurasi ketika mengevaluasi nilai SSE. K-Means mampu mengelompokkan data dengan lebih cepat dan menghasilkan cluster yang lebih jelas, sementara Hierarchical Clustering membutuhkan waktu komputasi yang lebih lama dan hasil clustering yang kurang optimal dalam beberapa kasus. Kesimpulannya, algoritma K-Means lebih cocok digunakan untuk analisis data pelanggan layanan internet dalam konteks ini karena keunggulannya dalam kecepatan dan akurasi clustering berdasarkan hasil pengecekan SSE. Penelitian ini memberikan panduan bagi perusahaan dalam memilih algoritma clustering yang tepat untuk mengelompokkan data pelanggan guna meningkatkan strategi bisnis dan kualitas layanan.

Kata Kunci: K-Means, Hierarchical Clustering, clustering, data pelanggan, layanan internet, data mining, segmentasi pasar, *Sum of Squared Errors (SSE)*.

I. Pendahuluan

PT XYZ adalah perusahaan telekomunikasi terkemuka di Indonesia, yang bergerak dalam sektor telekomunikasi. PT XYZ memiliki jaringan yang luas dan canggih di seluruh Indonesia, termasuk infrastruktur serat optik, jaringan nirkabel, satelit, dan kabel bawah laut. Sebagai BUMN, PT XYZ berkomitmen untuk mendukung pembangunan masyarakat dan perekonomian Indonesia. Mengingat tingginya permintaan pelanggan, PT XYZ perlu mengolah data untuk menyusun strategi bisnis yang efektif guna mempertahankan hubungan baik dengan pelanggan, terutama pelanggan setia. Salah satu cara untuk mencapai hal ini adalah dengan melakukan segmentasi pelanggan, yaitu mengelompokkan pelanggan dengan karakteristik serupa ke dalam satu kelompok. Segmentasi ini dapat dilakukan menggunakan metode clustering, yang merupakan salah satu teknik data mining untuk mengelompokkan data menjadi kelompok-kelompok tertentu. Metode clustering memiliki berbagai algoritma, di antaranya algoritma K-Means dan Hierarchical. Algoritma K-Means adalah teknik partisi yang membagi atau memisahkan objek ke dalam wilayah yang terpisah. Pada algoritma ini, setiap objek harus dimasukkan ke dalam kelompok atau cluster tertentu. Sebaliknya, Hierarchical Clustering mengelompokkan data berdasarkan hierarki atau tingkatan tertentu, memungkinkan analisis yang lebih mendalam terhadap struktur data.

Penelitian ini memilih algoritma K-Means dan Hierarchical Clustering dengan beberapa pertimbangan penting. Pertama, algoritma K-Means dikenal efisien dalam menangani dataset besar dan mampu menghasilkan cluster yang jelas dan terpisah. K-Means juga relatif mudah diimplementasikan dan memberikan hasil clustering yang mudah diinterpretasikan. Di sisi lain, Hierarchical Clustering dipilih karena kemampuannya membangun hierarki atau dendrogram yang memberikan wawasan mendalam tentang struktur data dan hubungan antar cluster. Kombinasi kedua algoritma ini memungkinkan peneliti mendapatkan gambaran yang lebih komprehensif tentang data pelanggan. Dalam penelitian ini, peneliti melakukan studi kasus pada perusahaan distributor yang bergerak di bidang layanan, dengan target pemasaran yang dapat mencakup berbagai segmen pelanggan. Penelitian ini tidak hanya fokus pada pemilihan algoritma clustering yang tepat, tetapi juga bertujuan memberikan rekomendasi strategi bisnis yang lebih efektif berdasarkan hasil clustering tersebut. Dengan demikian, PT XYZ dapat lebih memahami

kebutuhan pelanggan, meningkatkan kualitas layanan, dan mempertahankan loyalitas pelanggan yang tinggi.

Selain itu, penelitian ini juga bertujuan untuk menjawab tantangan yang dihadapi oleh banyak perusahaan di era digital saat ini, yaitu bagaimana mengelola dan memanfaatkan data pelanggan yang besar dan kompleks untuk mendukung keputusan bisnis strategis. Dengan menerapkan metode clustering, perusahaan dapat mengidentifikasi pola-pola tersembunyi dalam data pelanggan, mengoptimalkan kampanye pemasaran, dan meningkatkan efisiensi operasional. Penelitian ini diharapkan dapat memberikan kontribusi signifikan bagi PT XYZ dalam upaya mereka tetap menjadi pemimpin di industri telekomunikasi nasional, serta memberikan wawasan berharga bagi perusahaan lain yang menghadapi tantangan serupa dalam pengelolaan data pelanggan.

II. Metodologi penelitian

1) Data Mining

Data mining merupakan “proses pengolahan data untuk menemukan pola-pola baru dan bermanfaat dari kumpulan data besar” (Efori Buulolo, 2020). Menurut Han dan Kamber data mining adalah “bagian dari proses Knowledge Discovery in Databases (KDD) yang terdiri dari beberapa tahap, mulai dari pembersihan data hingga interpretasi hasil” (Han et al., 2012). Teknik-teknik dalam data mining digunakan untuk berbagai tujuan seperti klasifikasi, regresi, clustering, dan asosiasi. *Output* dalam *data mining* dapat dipergunakan sebagai alternatif dalam pengambilan keputusan atau untuk memperbaiki keputusan di masa yang akan datang (Efori Buulolo, 2020).

2) Clustering

Clustering adalah “salah satu teknik dalam data mining yang digunakan untuk mengelompokkan data menjadi beberapa kelompok atau cluster berdasarkan kemiripan karakteristik”. Menurut Jain clustering adalah “proses pengelompokan objek-objek ke dalam kelompok yang anggotanya memiliki kesamaan tertentu”. Clustering membantu dalam memahami struktur dan hubungan dalam data serta mengidentifikasi pola-pola yang tersembunyi” (Jain et al., 1999).

3) K-Means

Metode K-Means adalah “salah satu metode populer dalam analisis data dan merupakan teknik yang umum digunakan untuk melakukan klausterisasi (clustering) data”. Algoritma ini mengelompokkan data

berdasarkan jarak antara setiap data dengan pusat kelompok yang mewakili kelompok tersebut. Rumus penggunaan K- Means melibatkan dua komponen utama, yaitu jarak Euclidean dan pembaruan pusat kelompok (Kaur et al., 2020).

4) Hierarchical Clustering

Hierarchical Clustering adalah "Metode klausterisasi yang berusaha untuk membentuk hierarki kluster berdasarkan jarak antara data point". Hierarchical Clustering adalah "Metode clustering yang membangun hierarki atau dendrogram dari data. Metode ini bisa bersifat agglomerative (bottom-up) atau divisive (top-down)". Menurut Johnson *Agglomerative Clustering* dimulai dengan setiap data point sebagai cluster terpisah dan menggabungkannya secara bertahap, sedangkan *divisive clustering* dimulai dengan satu cluster dan memecahnya secara bertahap (Johnson, 1967).

5) Sum Of Squared Error (SSE)

Sum of Squared Errors (SSE) adalah "metrik yang digunakan untuk mengevaluasi kinerja clustering dengan mengukur jumlah kesalahan kuadrat antara setiap data point dan centroid terdekatnya". SSE yang lebih rendah menunjukkan clustering yang lebih baik. Menurut Hartigan SSE digunakan untuk menentukan kualitas clustering dan membantu dalam memilih jumlah cluster yang optimal (Jollyta et al., 2019).

6) Segmentasi Pelanggan

Segmentasi pelanggan adalah "proses mengelompokkan pelanggan ke dalam kelompok-kelompok berdasarkan karakteristik atau perilaku mereka". Segmentasi ini membantu perusahaan dalam memahami kebutuhan pelanggan, mengidentifikasi peluang pasar, dan merancang strategi pemasaran yang lebih efektif. Menurut Kotler dan Keller, segmentasi pelanggan adalah "langkah kunci dalam manajemen hubungan pelanggan (CRM) yang berfokus pada meningkatkan kepuasan dan loyalitas pelanggan" (Kotler & Keller, 2016)

7) UML

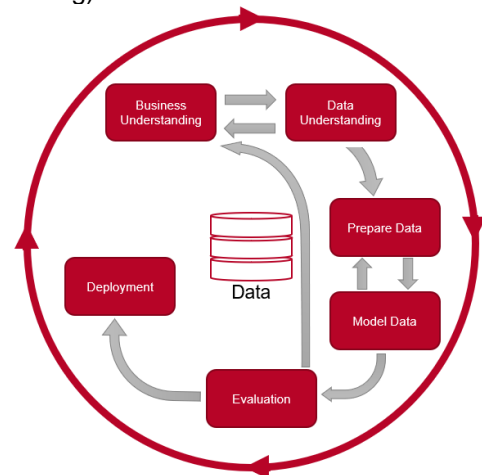
Unified Modeling Language (UML) adalah "bahasa standar yang digunakan untuk spesifikasi, visualisasi, pengembangan, dan pendokumentasian artefak dari sistem perangkat lunak". UML menyediakan serangkaian notasi untuk menggambarkan berbagai aspek dari sistem, termasuk struktur, perilaku, dan interaksi antar komponen (Sudarman, 2019).

8) PHP-MYSQL-Bootstrap

PHP adalah "bahasa skrip server-side yang dirancang khusus untuk pengembangan web. PHP memungkinkan pembuatan halaman web yang dinamis dengan kemampuan untuk berinteraksi dengan basis data dan memproses input pengguna". Bootstrap adalah "framework front-end yang digunakan untuk pengembangan web responsif dan mobile-first". Bootstrap menyediakan komponen HTML, CSS, dan JavaScript yang sudah siap pakai, yang memudahkan pengembang dalam membuat desain web yang konsisten dan modern. MySQL adalah "sistem manajemen basis data relasional (RDBMS) yang menggunakan SQL (*Structured Query Language*) sebagai bahasa pengelolaan data". (Wibowo P, 2018).

III. HASIL PEMBAHASAN

Dalam penelitian ini, langkah-langkah penelitian yang penulis lakukan menggunakan CRISP-DM (Cross-Industry Standard Process for Data Mining).



Gambar 1. CRISP-DM

Sumber: Michael Fuchs Python (2024)

1. *Business Understanding* (Pemahaman Bisnis)
Pada fase ini penulis ingin memahami tujuan bisnis PT XYZ terkait dengan segmentasi pelanggan layanan internet dan Mengidentifikasi pola penggunaan layanan internet oleh pelanggan.
2. *Data Understanding* (Pemahaman Data)
Pada fase ini penulis Mengumpulkan data pelanggan yang relevan, menganalisis data untuk memahami distribusi, korelasi, dan anomali. Mengidentifikasi data yang hilang atau tidak konsisten dan memeriksa kelengkapan data dan mengatasi data yang hilang atau anomali
3. *Data Preparation* (Persiapan Data)

Gambar 5. Hasil Algoritma K-Means dan SSE

Penulis menampilkan hasil algoritma *Hierarchical Clustering* setelah melakukan proses perhitungan dari data yang telah di masukkan.

Kesimpulan Akhir:

Klaster terdekat adalah Cluster 4 dan Cluster 6 dengan jarak

Jarak antara Cluster 1 dan Cluster 3: 60
Jarak antara Cluster 1 dan Cluster 4: 10
Jarak antara Cluster 1 dan Cluster 5: 110
Jarak antara Cluster 1 dan Cluster 6: 10
Jarak antara Cluster 1 dan Cluster 7: 130
Jarak antara Cluster 1 dan Cluster 8: 10
Jarak antara Cluster 2 dan Cluster 3: 80
Jarak antara Cluster 2 dan Cluster 4: 30
Jarak antara Cluster 2 dan Cluster 5: 130
Jarak antara Cluster 2 dan Cluster 6: 30
Jarak antara Cluster 2 dan Cluster 7: 130
Jarak antara Cluster 2 dan Cluster 8: 10
Jarak antara Cluster 3 dan Cluster 4: 50
Jarak antara Cluster 3 dan Cluster 5: 50
Jarak antara Cluster 3 dan Cluster 6: 50
Jarak antara Cluster 3 dan Cluster 7: 50
Jarak antara Cluster 3 dan Cluster 8: 70
Jarak antara Cluster 4 dan Cluster 5: 100
Jarak antara Cluster 4 dan Cluster 6: 0
Jarak antara Cluster 4 dan Cluster 7: 100
Jarak antara Cluster 4 dan Cluster 8: 20
Jarak antara Cluster 5 dan Cluster 6: 100
Jarak antara Cluster 5 dan Cluster 7: 0
Jarak antara Cluster 5 dan Cluster 8: 120
Jarak antara Cluster 6 dan Cluster 7: 100
Jarak antara Cluster 6 dan Cluster 8: 20
Jarak antara Cluster 7 dan Cluster 8: 120

Sum of Squared Errors (SSE):

SSE: 155100

SSE (Sum of Squared Errors) adalah metrik yang mengukur j:

Gambar 6. Hasil Algoritma *Hierarchical Clustering* dan SSE

Penulis melakukan komparasi algoritma k-means dan *Hierarchical Clustering* dengan menggunakan SSE sebagai standart yang dapat dilihat di bawah ini.

Sum of Squared Errors (SSE) Dari Algoritma Hierarchy:

SSE: 155100

Sum of Squared Errors (SSE) Dari Algoritma Kmeans:

SSE: 98935

Untuk mengetahui Algoritma clusterisasi terbaik dari K-means dan Hierarchy maka menggunakan sum of Squared Errors (SSE)

Jumlah SSE yang lebih rendah menunjukkan bahwa data lebih terkelompok dan klaster lebih terdefinisi dengan baik

Kesimpulan: Algoritma Kmeans Lebih baik dari Algoritma Hierarchy dikarenakan SSE algoritma Kmeans 98935 lebih kecil dari algoritma Hierarchy

Gambar 7. Hasil Komparasi Algoritma K-Means dan *Hierarchical Clustering*

Hasil penerapan Algoritma K-Means dan *Hierarchical Clustering* sudah dilakukan dari tahap pemodelan hingga tahap perhitungan. Metode K-Means dan Hierarchical Clustering sesuai dengan yang diharapkan yaitu mengcluster data customer dalam layanan internet PT. XYZ.

5. Kesimpulan dan Saran

1. Kesimpulan

- Pengembangan sistem informasi menggunakan analisis clustering berhasil dilakukan dengan menerapkan metode K-Means dan Hierarchical Clustering.
- Hasil analisis menunjukkan bahwa model K-Means memiliki Sum of Squared Errors (SSE) yang lebih rendah (SSE 98935) dibandingkan dengan model Hierarchical (SSE 155100), menunjukkan prediksi yang lebih baik dalam mengclusterisasi layanan internet.
- Kedua metode clustering ini efektif dalam mengorganisir data pelanggan layanan internet berdasarkan pola penggunaan yang berbeda-beda.

2. Saran

- Gunakan dataset yang lebih luas, seperti data dari beberapa tahun terakhir, untuk meningkatkan akurasi dan keakuratan clustering.
- Uji performa model dalam situasi yang lebih kompleks untuk memvalidasi keandalan dan keefektifan prediksi pasar.
- Eksplorasi metode clustering lain atau variasi dari K-Means dan Hierarchical Clustering untuk memperluas pemahaman terhadap data pelanggan dalam layanan internet.
- Fokuskan penelitian pada pengembangan strategi penjualan yang lebih efektif berdasarkan hasil clustering untuk meningkatkan kepuasan pelanggan dan pertumbuhan bisnis.

Daftar Pustaka

- Efori Buulolo. (2020). *Data Mining Untuk Perguruan Tinggi - Efori Buulolo*. Google Book.
- Han, J., Kamber, M., & Pei, J. (2012). *Data Mining: Concepts and Techniques*: Concepts and Techniques (3rd Edition). In *Data Mining*. Morgan Kaufmann. <http://linkinghub.elsevier.com/retrieve/pii/B9780123814791000010>
- Jain, A. K., Murty, M. N., & Flynn, P. J. (1999). *Data Clustering ACM Computing Surveys*. *Intelligent Multidimensional Data Clustering and Analysis*, 31(3), 265–323.
- Johnson, S. C. (1967). *Hierarchical Clustering Schemes*. *Psychometrika*.
- Jollyta, D., Efendi, S., Zarlis, M., & Mawengkang, H. (2019). *Optimasi Cluster Pada Data Stunting: Teknik Evaluasi Cluster Sum of Square Error dan Davies Bouldin Index*.

- Prosiding Seminar Nasional Riset Information Science (SENARIS)*, 1(September), 918.
<https://doi.org/10.30645/senaris.v1i0.100>
- Kaur, A., Pal, S. K., & Singh, A. P. (2020). Quantum-inspired ant lion optimized hybrid K-Means for cluster analysis and intrusion detection. *Knowledge-Based Syst*, 97. <https://doi.org/10.1016/j.asoc.2019.105523>
- Kotler, P., & Keller, K. L. (2016). *Manajemen Pemasaran* (G. A. Pratama (ed.); 13th ed.). Gelora Aksara Pratama.
- Sudarman. (2019). *Pemodelan Sistem Informasi Menggunakan Unified Modeling Language (UML)*. Elex Media Komputindo.
- Wibowo P. (2018). *Pemrograman Web Dinamis Menggunakan PHP dan MySQL*. Andi.