

Klastering Penyakit Diabetes dengan Metode K-Means

Muhammad Khalidin Basyir^{1*}

¹Fakultas Ilmu Komputer, Universitas Islam Negeri Sumatera Utara, Indonesia

Email: muhammadkhalidinbasyir@gmail.com

(* : coresponding author)

Abstrak—Penelitian ini bertujuan untuk mengklasifikasikan penyakit diabetes menggunakan metode *K-Means* berdasarkan karakteristik klinis dan demografis pasien. Data yang digunakan diambil dari platform *Kaggle*, terdiri dari 768 entri dengan delapan fitur numerik yang relevan. Proses analisis melibatkan normalisasi data dengan *StandardScaler* dan penentuan jumlah kluster optimal menggunakan metode *Elbow*. Hasil penelitian menunjukkan bahwa tiga kluster optimal dapat diidentifikasi, masing-masing mencerminkan karakteristik unik yang berkaitan dengan profil risiko diabetes. Kluster pertama menunjukkan pasien dengan kadar glukosa tinggi dan BMI tinggi, yang berkaitan dengan diabetes tipe 2. Kluster kedua mencerminkan pasien dengan BMI rendah dan kadar glukosa normal, sedangkan kluster ketiga mengelompokkan pasien dengan kadar insulin rendah dan usia muda. Temuan ini menggarisbawahi potensi algoritma *K-Means* untuk mendukung pengambilan keputusan medis yang lebih personal dan terarah.

Kata Kunci: Diabetes, Klastering, *K-Means*, Analisis Data, Visualisasi

Abstract—This study aims to classify diabetes using the *K-Means Method* based on the clinical and demographic characteristics of patients. The data used was taken from the *Kaggle platform*, consisting of 768 entries with eight relevant numerical features. The analysis process involved normalizing the data with *StandardScaler* and determining the optimal number of clusters using the *Elbow Method*. The results showed that three optimal clusters could be identified, each reflecting unique characteristics relating to diabetes risk profiles. The first cluster represents patients with high Glucose levels and high BMI, which is associated with type 2 diabetes. The second cluster reflects patients with low BMI and normal Glucose levels, while the third cluster groups patients with low insulin levels and young Age. These findings underscore the potential of the *K-Means algorithm* to support more personalized and targeted medical decision-making.

Keywords: Diabetes, Clustering, *K-Means*, Data Analysis, Visualization

1. PENDAHULUAN

Diabetes adalah salah satu penyakit kronis yang semakin meningkat prevalensinya di seluruh dunia, dengan angka yang mencapai lebih dari 463 juta orang pada tahun 2019 dan diperkirakan akan terus meningkat (Yan et al., 2018). Penyakit ini tidak hanya berdampak pada kesehatan individu, tetapi juga memberikan beban yang signifikan pada sistem kesehatan global. Oleh karena itu, pendekatan yang lebih efektif dalam pengelolaan dan pencegahan diabetes sangat diperlukan. Salah satu metode yang menjanjikan dalam analisis data kesehatan adalah klastering, klastering merupakan metode untuk mengelompokkan titik data menjadi dua kelompok atau lebih, sehingga titik-titik dalam kelompok yang sama memiliki kesamaan yang lebih tinggi dibandingkan dengan titik di kelompok lainnya, dan semua ini dilakukan berdasarkan informasi yang ada pada setiap titik data (Herlinda & Darwis, 2021). Tujuan dari klastering adalah mengelompokkan karakteristik yang serupa ke dalam satu area, sementara data dengan karakteristik yang berbeda akan tergabung dalam kelompok objek yang memiliki kesamaan (Dewi et al., 2022). Klastering memungkinkan pengelompokan pasien berdasarkan karakteristik klinis dan demografis mereka, sehingga memfasilitasi intervensi yang lebih terarah dan personalisasi dalam perawatan (Yang et al., 2024).

Kajian literatur menunjukkan bahwa berbagai teknik klastering, termasuk *K-Means*, telah diterapkan dalam konteks kesehatan untuk mengidentifikasi pola dalam data pasien dan mendukung pengambilan keputusan klinis (Chen et al., 2017). Sebagai contoh, penelitian oleh Hutchins menunjukkan bahwa teknik klastering dapat digunakan untuk mengidentifikasi tren yang tidak terlihat dalam populasi pasien, yang pada gilirannya dapat membantu dalam merancang intervensi yang lebih sesuai dengan kebutuhan kelompok tertentu. Selain itu, penelitian oleh Yadav menekankan pentingnya teknik klastering dalam pengembangan model prediksi yang dapat meningkatkan layanan kesehatan (Yadav et al., 2023). Namun, meskipun ada kemajuan yang signifikan dalam penerapan teknik ini, masih terdapat kekurangan dalam pemahaman tentang

bagaimana klustering dapat dioptimalkan untuk spesifik penyakit diabetes, yang menjadi dasar pernyataan kebaruan ilmiah dari artikel ini.

Pernyataan kebaruan ilmiah dari penelitian ini adalah penerapan metode *K-Means* dalam klustering pasien diabetes dengan mempertimbangkan variabel klinis dan non-klinis yang lebih luas, serta pengembangan model yang dapat memberikan wawasan lebih dalam tentang karakteristik kelompok pasien yang berbeda. Penelitian ini bertujuan untuk mengisi celah dalam literatur yang ada dengan memberikan analisis yang lebih mendalam tentang bagaimana klustering dapat digunakan untuk meningkatkan pemahaman tentang diabetes dan mendukung pengembangan intervensi yang lebih efektif.

Permasalahan penelitian yang diangkat dalam artikel ini adalah bagaimana metode *K-Means* dapat dioptimalkan untuk mengidentifikasi kluster pasien diabetes yang berbeda berdasarkan karakteristik klinis dan demografis mereka. Hipotesis yang diajukan adalah bahwa penerapan metode *K-Means* yang disesuaikan dengan variabel spesifik diabetes akan menghasilkan kluster yang lebih relevan dan informatif, yang pada gilirannya dapat meningkatkan efektivitas intervensi kesehatan yang dirancang untuk kelompok pasien tersebut.

2. METODE PENELITIAN

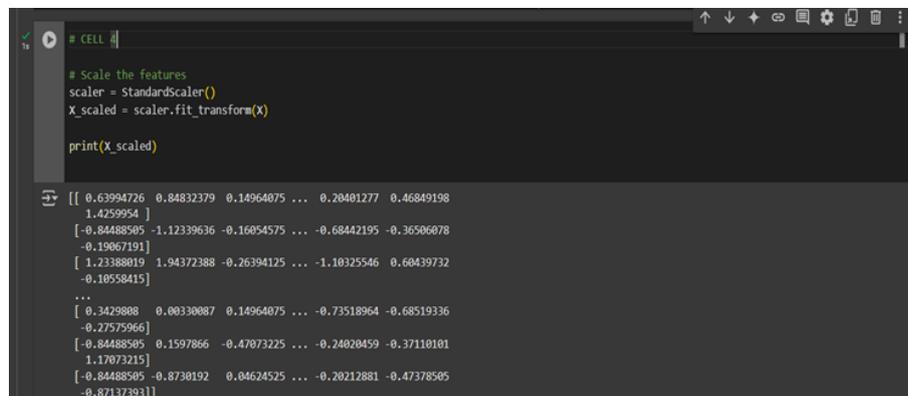
2.1 Persiapan Data

Dataset ini diperoleh dari *Kaggle*, *Kaggle* merupakan sebuah platform yang digunakan oleh berbagai perusahaan untuk menyajikan kompetisi di bidang *data science* (Angkasa & Pangaribuan, 2022). Dataset terdiri dari 768 entri serta 8 fitur numerik, yang digunakan untuk menganalisis penyakit diabetes. Fitur-fitur tersebut mencakup jumlah kehamilan (*Pregnancies*), kadar glukosa dalam darah (*Glucose*), tekanan darah diastolik (*BloodPressure*), ketebalan kulit pada triceps (*SkinThickness*), kadar insulin serum (*Insulin*), indeks massa tubuh (BMI), fungsi silsilah diabetes (*DiabetesPedigreeFunction*), dan usia pasien (*Age*). Dataset ini tidak mengandung nilai yang hilang, sehingga sangat cocok untuk penelitian pembelajaran mesin, seperti klustering atau prediksi risiko diabetes.

2.2 Preprocessing

Tahap *preprocessing* melibatkan normalisasi data menggunakan metode *StandardScaler* dari pustaka *Scikit-learn*. *Standard Scaler* adalah teknik *preprocessing* yang berfungsi untuk menstandarisasi fitur dengan cara menghilangkan rata-rata dan menskalakan berdasarkan varians (Prasetyo et al., 2022). Normalisasi diperlukan karena fitur dalam dataset memiliki rentang nilai yang berbeda-beda, seperti kadar glukosa yang memiliki nilai jauh lebih besar dibandingkan indeks massa tubuh (BMI). Tanpa normalisasi, fitur dengan skala besar dapat mendominasi hasil klustering.

Setelah normalisasi, setiap fitur memiliki rata-rata 0 dan deviasi standar 1, sehingga memastikan setiap variabel memiliki kontribusi yang seimbang dalam algoritma *K-Means*. Hasil dari proses ini adalah dataset yang siap untuk dianalisis dengan metode klustering.



```
# Scale the features
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

print(X_scaled)
```

```
[[ 0.63994726  0.84832379  0.14964075 ...  0.20401277  0.46849198
  1.4259954 ]
 [-0.84488505 -1.12339636 -0.16054575 ... -0.68442195 -0.36506078
 -0.19067191]
 [ 1.23388019  1.94372388 -0.26394125 ... -1.10325546  0.60439732
 -0.10558413]
 ...
 [ 0.3429888  0.00330087  0.14964075 ... -0.73518964 -0.68519336
 -0.27575966]
 [-0.84488505  0.1597866 -0.47073225 ... -0.24020459 -0.37110101
  1.17073213]
 [-0.84488505 -0.8730192  0.04624525 ... -0.20212881 -0.47378505
 -0.87137393]]
```

Gambar 1. *Preprocessing* Menggunakan Metode *Standard Scaler*

2.3 Penentuan Jumlah Cluster (*Elbow Method*)

Metode *Elbow* merupakan salah satu pendekatan yang berfungsi sebagai evaluasi dan digunakan untuk menentukan nilai k optimal dari serangkaian percobaan untuk menguji nilai k (Fuadah et al., 2021). Metode *Elbow* digunakan untuk menentukan jumlah *cluster* (k) yang optimal. Dalam penelitian ini, nilai k diuji dari 1 hingga 10, dan nilai *inertia* dihitung untuk setiap k . *Inertia* adalah total jarak kuadrat antara data dalam *cluster* dan *centroid*-nya.

Grafik *Elbow Method* dianalisis untuk menemukan "titik siku" yang menunjukkan penurunan signifikan pada nilai *inertia*. Titik ini mencerminkan jumlah *cluster* optimal, yang merupakan keseimbangan antara variabilitas data yang dijelaskan dan kompleksitas model.

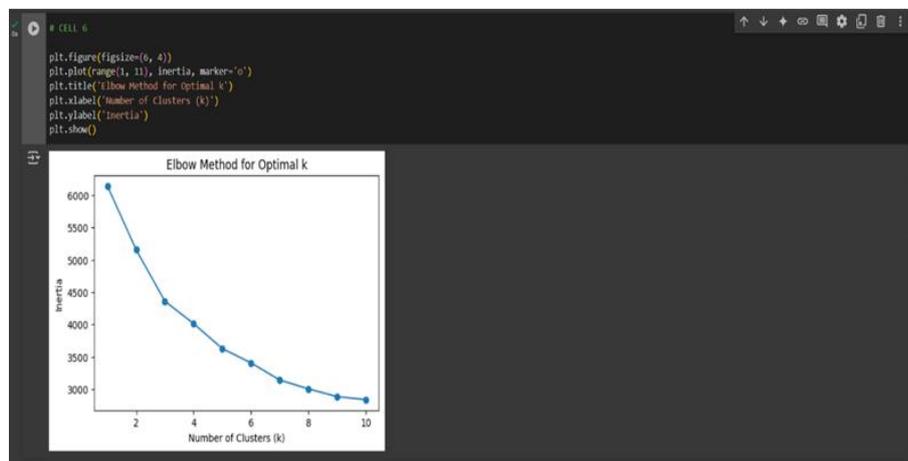
```
# CELL 5
# Elbow method to find optimal k
inertia = []
for k in range(1, 11): # Test k from 1 to 10
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(X_scaled)
    inertia.append(kmeans.inertia_)

inertia
[6143.9999999999994,
 5156.2502953401345,
 4360.470411418354,
 4013.437467578209,
 3625.7434665428546,
 3405.1387568451079,
 3139.810055595146,
 3002.471914445965,
 2882.34033583747,
 2836.913676999315]
```

Gambar 2. Penentuan Jumlah Cluster (*Elbow Method*)

2.3 Penerapan *K-Means* dan Visualisasi

Setelah jumlah *cluster* (k) ditentukan, algoritma *K-Means* diterapkan pada dataset yang telah dinormalisasi. Setiap data kemudian dikelompokkan ke dalam *cluster* berdasarkan jaraknya ke *centroid* terdekat, yang diperbarui secara iteratif hingga konvergen. Proses ini menghasilkan pembagian data menjadi *cluster-cluster* dengan karakteristik yang berbeda.



Gambar 3. Visualisasi *Elbow Method*

3. ANALISA DAN PEMBAHASAN

3.1 Penentuan Optimal K

Berdasarkan hasil analisis *Elbow Method*, jumlah *cluster* optimal yang ditemukan adalah $k=3$. Pada grafik, terlihat penurunan signifikan pada nilai *inertia* antara $k=1$ hingga $k=3$, sementara setelah $k=3$ penurunannya cenderung melandai. Hal ini menunjukkan bahwa $k=3$ merupakan jumlah

cluster yang memberikan keseimbangan terbaik antara variabilitas data yang menjelaskan dan kesederhanaan model.

Pemilihan $k=3$ juga didukung oleh karakteristik dataset, di mana pembagian ke dalam tiga *cluster* dianggap relevan untuk menggambarkan variasi pasien dengan penyakit diabetes berdasarkan faktor risiko utama seperti kadar glukosa, indeks massa tubuh, dan usia.

```
# CELL 7
# Apply KMeans clustering with the chosen k (e.g., k=3 based on the elbow method)
optimal_k = 3 # Replace with your optimal k value from the elbow method
kmeans = KMeans(n_clusters=optimal_k, random_state=42)
data['cluster'] = kmeans.fit_predict(X_scaled)
print(data['cluster'])
```

```
0    1
1    2
2    1
3    2
4    0
..
763  1
764  2
765  2
766  1
767  2
Name: cluster, Length: 768, dtype: int32
```

Gambar 4. Penentuan Optimal K

3.2 Klastering Penyakit Diabetes

Hasil klastering menunjukkan bahwa setiap *cluster* memiliki karakteristik yang unik. *Cluster* pertama mencakup pasien dengan kadar glukosa tinggi, BMI tinggi, dan usia di atas 40 tahun, yang sesuai dengan profil pasien diabetes tipe 2. *Cluster* kedua didominasi oleh pasien dengan BMI rendah, kadar glukosa normal, dan usia muda, yang mencerminkan kondisi pasien tanpa risiko diabetes yang signifikan. *Cluster* ketiga mencakup pasien dengan kadar insulin rendah dan usia di bawah 30 tahun, yang lebih cenderung menggambarkan pasien diabetes tipe 1 atau diabetes gestasional.

Distribusi data dalam *cluster* menunjukkan bahwa metode *K-Means* mampu memisahkan kelompok pasien dengan pola karakteristik yang jelas. Temuan ini memberikan wawasan tentang bagaimana faktor risiko tertentu dapat membantu mengidentifikasi tipe diabetes yang berbeda dalam populasi pasien.

```
# CELL 8
# Assign diabetes labels
data['Diabetes_Label'] = data.apply(assign_diabetes_label, axis=1)
data
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Cluster	Diabetes_Label
0	6	148	72	35	0	33.6	0.627	50	1	Diabetes Tipe 2
1	1	85	66	29	0	26.6	0.351	31	2	Diabetes Gestasional
2	8	183	64	0	0	23.3	0.672	32	1	Diabetes Gestasional
3	1	89	66	23	94	28.1	0.167	21	2	Diabetes Gestasional
4	0	137	40	35	168	43.1	2.288	33	0	Diabetes Tipe 2
..
763	10	101	76	48	180	32.9	0.171	63	1	Diabetes Tipe 2
764	2	122	70	27	0	36.8	0.340	27	2	Diabetes Tipe 1
765	5	121	72	23	112	26.2	0.245	30	2	Diabetes Gestasional
766	1	128	60	0	0	30.1	0.349	47	1	Diabetes Tipe 2
767	1	93	70	31	0	30.4	0.315	23	2	Diabetes Tipe 1

768 rows x 10 columns

Gambar 5. Klastering Penyakit Diabetes

3.3 Analisis Centorid dan Hubungan dengan Tipe Diabetes

Centorid dari setiap *cluster* dianalisis untuk memahami karakteristik utama dari masing-masing kelompok. Sebagai contoh, *centorid cluster* pertama menunjukkan nilai rata-rata glukosa dan BMI yang tinggi, mendukung hipotesis bahwa *cluster* ini didominasi oleh pasien diabetes tipe 2. Sebaliknya, *centorid cluster* kedua memiliki nilai rata-rata glukosa yang rendah, menunjukkan kondisi metabolisme yang lebih sehat.

Hasil ini menunjukkan bahwa klustering dapat membantu dalam proses pengelompokan pasien berdasarkan risiko dan karakteristik penyakit diabetes. Jika digabungkan dengan label tipe diabetes, analisis ini dapat memberikan informasi tambahan untuk pengambilan keputusan medis.

```
centroids = pd.DataFrame(scaler.inverse_transform(means.cluster_centers_), columns=features)
print(centroids)
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin		BMI	DiabetesPedigreeFunction	Age
0	2.953852	141.544681	74.192488	34.779343	197.417840		36.972778	0.688216	31.769953
1	7.276818	128.411765	76.859729	12.882353	27.678733		32.188090	0.437557	45.787338
2	2.143713	102.751497	68.739539	16.517964	39.278443		28.692515	0.407638	25.877246

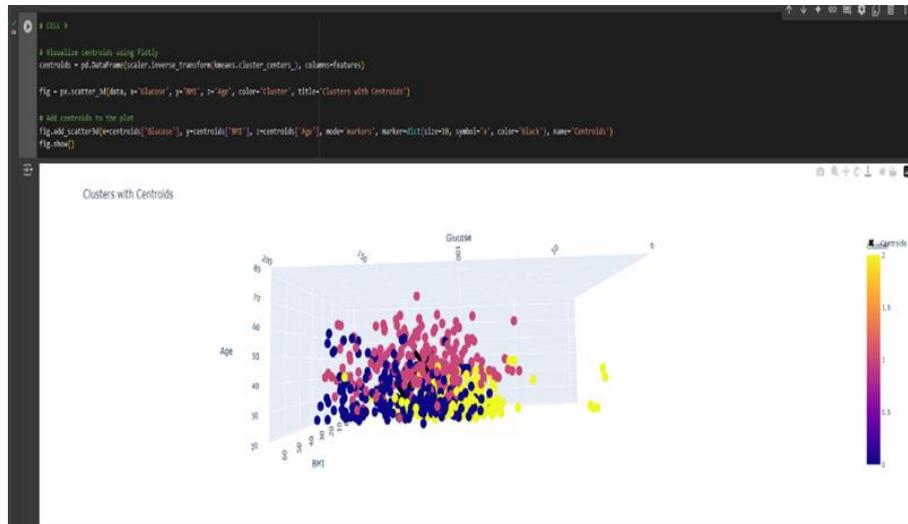
Gambar 6. Analisis Centroid

3.4 Visualisasi Cluster

Hasil klustering divisualisasikan dalam ruang 3D menggunakan Plotly. Plotly merupakan salah satu pustaka Python yang digunakan untuk menghasilkan grafik dan diagram dalam format tiga dimensi (Rakhmad, 2021). Visualisasi ini memanfaatkan fitur seperti *Glucose*, *BMI*, dan *Age* untuk memudahkan interpretasi hasil. *Centorid* dari setiap *cluster* juga divisualisasikan untuk memberikan gambaran tentang karakteristik pusat masing-masing *cluster*.

Centorid ditampilkan sebagai tanda silang hitam untuk menandai pusat dari masing-masing *cluster*. Distribusi ini menunjukkan bagaimana algoritma *K-Means* berhasil mengelompokkan data berdasarkan kesamaan fitur.

Melalui visualisasi, dapat diamati bahwa *cluster* dengan nilai glukosa dan BMI yang tinggi cenderung terpisah jauh dari *cluster* dengan nilai glukosa rendah. Hal ini mencerminkan pengaruh signifikan dari faktor-faktor tersebut terhadap pengelompokan pasien diabetes.



Gambar 7. Visualisasi

4. KESIMPULAN

Penelitian ini berhasil menerapkan metode *K-Means* untuk mengelompokkan pasien diabetes berdasarkan karakteristik klinis dan demografis. Berdasarkan analisis *Elbow Method*, jumlah kluster optimal yang ditemukan adalah tiga, dengan setiap kluster menunjukkan pola karakteristik pasien yang berbeda. Kluster pertama mencakup pasien dengan kadar glukosa tinggi, BMI tinggi, dan usia di atas 40 tahun, yang sesuai dengan profil diabetes tipe 2. Kluster kedua mendominasi pasien dengan BMI rendah, kadar glukosa normal, dan usia muda, sedangkan kluster ketiga mencerminkan pasien dengan kadar insulin rendah dan usia di bawah 30 tahun.

Hasil penelitian ini menunjukkan bahwa metode *K-Means* dapat digunakan secara efektif untuk mengidentifikasi pola dalam data pasien diabetes. Dengan pemahaman yang lebih baik terhadap kelompok risiko, metode ini memiliki potensi besar dalam mendukung pengembangan intervensi kesehatan yang lebih personal dan tepat sasaran. Visualisasi kluster juga memberikan wawasan tambahan yang memperkuat relevansi teknik ini dalam konteks analisis data kesehatan. Penelitian lebih lanjut diperlukan untuk mengintegrasikan klustering dengan label tipe diabetes guna meningkatkan akurasi dan manfaat aplikasinya dalam pengambilan keputusan medis.

UCAPAN TERIMA KASIH

Penulis mengucapkan terima kasih yang sebesar-besarnya kepada pihak-pihak yang telah berkontribusi dalam penelitian ini. Terima kasih kepada Universitas Islam Negeri Sumatera Utara, khususnya Program Studi Ilmu Komputer, yang telah memberikan dukungan selama proses penelitian ini berlangsung. Penulis juga mengucapkan terima kasih kepada platform Kaggle yang menyediakan dataset yang digunakan dalam penelitian ini. Selain itu, penghargaan juga disampaikan kepada rekan-rekan sejawat dan pembimbing yang telah memberikan masukan berharga selama proses penyusunan artikel ini. Semoga hasil penelitian ini dapat memberikan manfaat bagi pengembangan analisis data kesehatan, khususnya dalam pengelolaan penyakit diabetes.

REFERENCES

- Angkasa, V., & Pangaribuan, J. J. (2022). Information System Development Komparasi Tingkat Akurasi Random Forest Dan KNN Untuk Mendiagnosis Penyakit Kanker Payudara. *Journal of Information System Development*, 7(1), 37–38. Retrieved from <http://dx.doi.org/10.19166/xxxx>
- Chen, W., Chen, S., Zhang, H., & Wu, T. (2017). A hybrid prediction model for type 2 diabetes using K-Means and decision tree. In *2017 8th IEEE International Conference on Software Engineering and Service Science (ICSESS)* (pp. 386–390). <https://doi.org/10.1109/ICSESS.2017.8342938>
- Dewi, N. L. P. P., Purnama, I. N., & Utami, N. W. (2022). Penerapan Data Mining Untuk Clustering Penilaian Kinerja Dosen Menggunakan Algoritma K-Means (Studi Kasus: STMIK Primakara). *Jurnal Ilmiah Teknologi Informasi Asia*, 16(2), 105. <https://doi.org/10.32815/jitika.v16i2.761>
- Fuadah, A. W., Arifin, F. N., & Juwita, O. (2021). Optimasi K-Klasterisasi Ketahanan Pangan Kabupaten Jember Menggunakan Metode Elbow. *INFORMAL Informatics Journal*, 6(3), 136. <https://doi.org/10.19184/isj.v6i3.28363>
- Herlinda, V., & Darwis, D. (2021). Analisis Clustering Untuk Recredesialing Fasilitas Kesehatan Menggunakan Metode Fuzzy C-Means. *Darwis, Dartono*, 2(2), 94–99. Retrieved from <http://jim.teknokrat.ac.id/index.php/JTSI>
- Prasetyo, V. R., Mercifia, M., Averina, A., Sunyoto, L., & Budiarjo, B. (2022). Prediksi Rating Film Pada Website IMDb Menggunakan Metode Neural Network. *Network Engineering Research Operation*, 7(1), 1. <https://doi.org/10.21107/nero.v7i1.268>
- Rahmad, K. (2021). Rekonstruksi Tiga Dimensi Area Pegunungan Berdasarkan Citra Satelit Dua Dimensi Secara Otomatis. Tugas Akhir - Universitas Islam Indonesia. Retrieved from <https://dspace.uui.ac.id/handle/123456789/34042%25>
- Yadav, S., Singh, M. K., & Kumar, P. (2023). Data Mining Applications for Enhancing Healthcare Services: A Comprehensive Review. *International Journal of Engineering Technology and Management Sciences*, 7(5), 325–333. <https://doi.org/10.46647/ijetms.2023.v07i05.038>
- Yan, S., Kwan, Y. H., Tan, C. S., Thumboo, J., & Low, L. L. (2018). A systematic review of the clinical application of data-driven population segmentation analysis. *BMC Medical Research Methodology*, 18(1), 1–12. <https://doi.org/10.1186/s12874-018-0584-9>
- Yang, W. C., Lai, J. P., Liu, Y. H., Lin, Y. L., Hou, H. P., & Pai, P. F. (2024). Using Medical Data and Clustering Techniques for a Smart Healthcare System. *Electronics*, 13(1), 1–20. <https://doi.org/10.3390/electronics13010140>