Jurnal Ilmiah Dan Karya Mahasiswa Vol. 1 No. 6 Desember 2023





e-ISSN :2985-7732, p-ISSN :2985-6329, Hal 382-396 DOI: https://doi.org/10.54066/jikma.v1i6.1162

Implementasi Model Machine Learning dalam Mengklasifikasi Kualitas Air

Stacyana Jesika

Universitas Negeri Medan Email: stacyanajs@mhs.unimed.ac.id

Suci Ramadhani

Universitas Negeri Medan Email: suciramadhani122002@gmail.com

Yohanna Permata Putri

Universitas Negeri Medan Email: <u>yohannapasaribu6@gmail.com</u>

Jalan Willem Iskandar, Pasar V Medan Estate, Percut Sei Tuan, Deli Serdang

Abstract. Water quality is an important factor in maintaining human health and environmental sustainability. Water pollution is a major problem in Indonesia, so it is important to monitor and classify water quality effectively. Implementation of machine learning models in classifying water quality can provide important benefits in the environmental and health fields. This research uses two machine learning algorithms, namely KNN and SVM, to classify water quality. The water quality data used comes from the website www.kaggle.com, which was uploaded by MsSmartyPants in 2021 with the title "Water Quality (Dataset for water quality classification)". Implementation of this machine learning model involves the steps of data collection, data pre-processing, selection of relevant attributes, algorithm selection, model training, evaluation, and model implementation for real-time water quality classification.

Keywords: Classification, K-NN, SVM, and Water Quality.

Abstrak. Kualitas air merupakan faktor penting dalam menjaga kesehatan manusia dan keberlanjutan lingkungan. Pencemaran air menjadi masalah utama di Indonesia, sehingga penting untuk memantau dan mengklasifikasikan kualitas air dengan efektif. Implementasi model machine learning dalam mengklasifikasikan kualitas air dapat memberikan manfaat penting dalam bidang lingkungan dan kesehatan. Penelitian ini menggunakan dua algoritma machine learning, yaitu KNN dan SVM, untuk mengklasifikasikan kualitas air. Data kualitas air yang digunakan bersumber dari website www.kaggle.com, yang telah diunggah oleh MsSmartyPants pada tahun 2021 dengan judul "Water Quality (Dataset for water quality classification)". Implementasi model machine learning ini melibatkan langkah-langkah pengumpulan data, pra-pemrosesan data, pemilihan atribut yang relevan, pemilihan algoritma, pelatihan model, evaluasi, dan implementasi model untuk klasifikasi kualitas air secara real-time.

Kata kunci: Klasifikasi, K-NN, SVM, dan Kualitas Air.

LATAR BELAKANG

Kualitas air merupakan faktor penting dalam menjaga kesehatan manusia dan keberlanjutan lingkungan. Air memiliki standarisasi tersendiri dalam layak atau tidaknya bagi tubuh manusia. Pencemaran air merupakan isu utama di Indonesia, dimana sumber utamanya berasal dari limbah domestik dan rumah tangga. Dampaknya adalah menurunnya kualitas dan ketersediaan air bersih di berbagai wilayah di negara ini. Oleh karena itu, penting untuk memantau dan mengklasifikasikan kualitas air dengan efektif (Said et al., 2022).

Implementasi model machine learning dalam mengklasifikasikan kualitas air adalah penerapan algoritma dan teknik machine learning untuk menganalisis data terkait kualitas air dan menghasilkan model yang dapat memprediksi kualitas air berdasarkan atribut-atribut yang relevan (Savitri & Nursalim, 2023). Algoritma machine learning yang digunakan pada penelitian ini yaitu KNN dan SVM untuk mengklasifikasi serta memprediksi tingkat akurasi pada klasifikasi kualitas air.

Penelitian ini menggunakan algoritma KNN dan SVM karena beberapa alasan yaitu diantaranya adalah, SVM merupakan algoritma pembelajaran yang digunakan untuk masalah klasifikasi dan regresi. Prinsip dasar SVM adalah mencari hyperplane (bidang pemisah) yang optimal untuk memisahkan dua kelas data dengan jarak maksimum di antara kelas-kelas tersebut. SVM memiliki keunggulan dalam kemampuan generalisasi yang tinggi, kemampuan menghasilkan model klasifikasi yang baik dengan data pelatihan terbatas, dan kemudahan implementasi melalui formulasi support vector dalam masalah Quadratic Programming (QP) (Suyanto, 2017). Sementara itu untuk metode KNN merupakan algoritma pembelajaran yang berdasarkan prinsip "tetangga terdekat". KNN melakukan klasifikasi dengan cara mencari k-tetangga terdekat dari data baru yang akan diprediksi, kemudian menggunakan mayoritas kelas dari tetangga-tetangga tersebut untuk memprediksi kelas data baru. KNN memiliki beberapa keunggulan, seperti pelatihan yang cepat, sederhana dan mudah dipahami, ketahanan terhadap data pelatihan yang bising, serta efektivitasnya dalam menangani data pelatihan yang besar (Mutrofin et al., 2014).

Mengklasifikasikan kualitas air dengan menggunakan model machine learning dapat memberikan manfaat penting dalam bidang lingkungan dan kesehatan. Dengan memanfaatkan teknologi machine learning, kita dapat mengidentifikasi pola-pola kompleks dalam data air, mengklasifikasikan air berdasarkan standar kualitas yang

ditetapkan, dan mengambil tindakan yang tepat untuk menjaga dan meningkatkan kualitas air.

Data kualitas air yang didapatkan sendiri bersumber dari website www.kaggle.com yang diupload oleh MsSmartyPants pada tahun 2021 dengan judul Water Quality (Dataset for water quality classification). Implementasi model machine learning melibatkan langkah-langkah seperti pengumpulan data kualitas air, prapemrosesan data, pemilihan atribut yang relevan, pemilihan algoritma machine learning yang sesuai, pelatihan model menggunakan data latih, evaluasi dan pengujian model, serta implementasi dan penggunaan model untuk mengklasifikasikan kualitas air secara real-time. Oleh karena itu penelitian ini dibuat untuk membangun model machine learning yang mampu mengklasifikasikan kualitas air berdasarkan atribut-atribut yang tersedia, memprediksi tingkat akurasi dalam klasifikasi kualitas air dengan menggunakan algoritma K-NN dan SVM, serta meningkatkan pemahaman tentang kualitas air dan faktor-faktor yang mempengaruhinya melalui analisis data menggunakan model machine learning.

KAJIAN TEORITIS

1. Klasifikasi

Klasifikasi merupakan tugas krusial dalam machine learning dan data mining yang bertujuan untuk mengategorikan data ke dalam kelas atau kategori yang telah ditentukan berdasarkan atribut-atribut atau fitur-fitur tertentu. Dalam konteks klasifikasi, kita memiliki kumpulan data yang sudah diberi label sebelumnya, yang berarti setiap data di dalamnya sudah memiliki label yang sesuai dengan kelas atau kategori tertentu. Tujuan utama dari klasifikasi adalah untuk mengembangkan model yang mampu memahami pola-pola atau relasi-relasi antara atribut-atribut tersebut, lalu menggunakannya untuk memprediksi kelas atau kategori yang tepat untuk data-data baru yang belum pernah terlihat sebelumnya (Han et al., 2012).

2. Super Vector Mechine (SVM)

a. Pengertian SVM

SVM, atau Mesin Vektor Dukungan, adalah metode pembelajaran mesin yang terawasi yang digunakan untuk tugas seperti klasifikasi dan regresi. SVM pada dasarnya adalah sebuah pendekatan klasifikasi biner yang membagi titik-

titik data ke dalam dua kelas. Saat melatih SVM, tujuannya adalah menemukan hiperplane margin maksimum yang dapat memisahkan kedua kelas dengan jarak terbesar dari titik data pelatihan terdekat (McCartney, 2015).

SVM, atau Mesin Vektor Dukungan, adalah metode pembelajaran mesin yang terawasi yang digunakan untuk tugas seperti klasifikasi dan regresi. SVM pada dasarnya adalah sebuah pendekatan klasifikasi biner yang membagi titiktitik data ke dalam dua kelas. Saat melatih SVM, tujuannya adalah menemukan hiperplane margin maksimum yang dapat memisahkan kedua kelas dengan jarak terbesar dari titik data pelatihan terdekat (Han et al., 2012).

3. K-Nearest Neighbor (KNN)

K-Nearest Neighbor atau KNN, merupakan algoritma yang digunakan untuk mengklasifikasikan data dengan memanfaatkan data latih yang berasal dari K tetangga terdekat dalam sekitarnya (Maulida, 2020). k-NN dilakukan dengan mengidentifikasi kelompok k objek pada data pelatihan yang memiliki karakteristik yang serupa dengan objek pada data baru atau data pemeriksaan (Bachtiar et al., 2019). Salah satu teknik lalai belajar, K-Nearest Neighbor (K-NN) mencari sekelompok k objek dalam data pelatihan yang memiliki kemiripan paling besar dengan objek dalam data baru atau data pengujian (Maulida, 2020).

Algoritma ini beroperasi dengan menggunakan perhitungan jarak terpendek dari sampel uji ke sampel latihan untuk menentukan K-Nearest Neighbors (KNN). Setelah mengidentifikasi KNN, langkah berikutnya adalah mengambil mayoritas dari KNN tersebut sebagai prediksi untuk sampel uji. Biasanya, tingkat kedekatan atau jarak diukur menggunakan metrik jarak Euclidean (Maulida, 2020).

METODE PENELITIAN

Penelitian ini dimulai dengan mengumpulkan dan memahami data yang diperoleh. Selanjutnya, dilakukan persiapan data melalui proses preprocessing, yang meliputi langkah-langkah seperti membersihkan data, memilih data yang relevan, dan menyeimbangkan data. Setelah itu, data yang telah diolah akan dibagi menjadi bagianbagian yang sesuai untuk memulai proses klasifikasi menggunakan algoritma Machine Learning. Akhirnya, dilakukan evaluasi model dan perbandingan tingkat akurasi dari setiap model yang diuji untuk memprediksi kualitas air. Berikut adalah tahapan-tahapan

yang dilakukan pada penelitian ini yaitu:

1. Pengumpulan dan pemahaman data

Pengumpulan dan pemahaman data merupakan tahap awal sangat penting dalam mengelolah data. Data mining ialah metode untuk mengekstrak atau mengumpulkan informasi bernilai dari kumpulan data besar. Prosesnya sering melibatkan pemanfaatan teknik statistik dan matematika dengan menggunakan teknologi kecerdasan buatan (Said et al., 2022). Melalui analisis data, kita dapat memperoleh informasi yang beragam, seperti jumlah rekaman data, jenis data yang ada, dan jumlah kolom data yang tersedia. Data yang digunakan pada penelitian ini besifat universal (didapatkan dari berbagai tempat) yang bersumber dari website www.kaggle.com diunggah oleh MsSmartyPants pada tahun 2021 dengan judul Water Quality (Dataset for water quality classification), adapun rincian data yang digunakan sebagai berikut:

no	variabel	rentang nilai
1	Aluminium	0-5
2	Ammonia	-0.08 - 29.84
3	Arsenic	0 - 1.05
4	Barium	0-4.49
5	Cadmium	0-0.13
6	Chloramine	0 - 8.68
7	Chromium	0 - 0.90
8	Copper	0-2
9	Flouride	0 - 1.50
10	Lead	0 - 0.20
11	Nitrat	0-19.83
12	Nitrit	0-2.93
13	Mercury	0 - 0.01
14	Perchlorate	0-60.01
15	Radium	0 - 7.99
16	Selenium	0 - 0.100
17	Silver	0 - 0.5
18	Uranium	0 - 0.09
19	Bacteria	0-1
20	Virus	0-1
21	is_safe	0-1

Gambar 1: Tabel atribut dataset

2. Preprocessing Data

Data preprocessing merupakan teknik awal dalam proses penambangan data yang bertujuan untuk mengubah data mentah menjadi format yang lebih efisien dan bermanfaat. Format data mentah yang berasal dari berbagai sumber seringkali mengandung kesalahan, nilai yang hilang, dan inkonsistensi. Oleh karena itu, perlu dilakukan pembenahan format data agar hasil dari proses penambangan data menjadi

tepat dan akurat (Said et al., 2022). Dengan melakukan preprocessing data yang baik, dapat mempengaruhi hasil dari penambangan data dengan meningkatkan akurasi dan menghasilkan informasi yang lebih jelas. Tahapan dalam data preprocessing terdiri dari data cleansing, data selection, dan data balancing.

3. Split Data

Pembagian data merupakan salah satu metode yang digunakan untuk mengevaluasi kinerja model, dengan menyesuaikan persentase data dalam set pelatihan dan pengujian. Hal ini bertujuan untuk memperoleh performa maksimal pada model machine learning yang akan digunakan. Split data ini melibatkan pemisahan dataset menjadi dua bagian, yang akan digunakan sebagai data pelatihan dan data pengujian dengan proporsi tertentu. Dalam penelitian ini, variasi persentase dan jumlah pembagian data akan dilakukan dalam beberapa percobaan untuk mengamati hasil akurasi tertinggi.(Said et al., 2022)

4. Algortima Machine Learning

Adapun beberapa algortima Machine leraning yang digunakan dalam memprediksi klasifikasi kualitas air diantaranya:

a. K-Nearest Neighbors (KNN)

Adalah suatu metode machine learning yang memanfaatkan k-nearest neighbors (k tetangga terdekat) untuk melakukan klasifikasi. Dalam konteks prediksi kualitas air, KNN dapat digunakan untuk mengklasifikasikan air berdasarkan atribut-atribut yang relevan. Cara kerja algoritma ini adalah dengan menghitung jarak antara sampel data yang baru dengan tetangga terdekatnya. Kualitas air baru kemudian diklasifikasikan berdasarkan mayoritas label dari tetangga terdekatnya. Dengan demikian, KNN memanfaatkan informasi dari tetangga terdekat untuk membuat prediksi klasifikasi.

b. SVM

SVM adalah algoritma yang beroperasi dengan cara berikut: ia mengaplikasikan pemetaan non-linear untuk mengganti data train asal ke dalam dimensi yang lebih tinggi. Di dalam dimensi baru tersebut, SVM mencari hiperplane pemisah linear yang optimal (dikenal sebagai "batas keputusan") yang memisahkan tuple dari satu kelas dengan yang lain. Melalui pemetaan non-linear yang sesuai ke dimensi yang cukup tinggi, data dari kedua kelas selalu bisa

dibedakan oleh hiperplane. SVM mengidentifikasi hiperplane ini dengan memanfaatkan vektor dukungan (yang merupakan tuple pelatihan yang sangat penting) dan margin (yang ditentukan oleh vektor dukungan).(Han et al., 2012)

5. Evaluasi Model

Tahapan ini dilakukan dalam menilai Kemampuan atau hasil kerja dari suatu algoritma dalam tingkatan akurasi pada model, diperlukan penghitungan matriks kebingungan (confusion matrix). Confusion matrix adalah teknik yang dipakai dalam menilai hasil kerja dari model yang diuji terutama dalam kasus klasifikasi pada machine learning (Savitri & Nursalim, 2023). Confusion matrix memiliki empat kombinasi yang berbeda berdasarkan nilai prediksi dan nilai aktual. Tabel evaluasi model dapat ditemukan pada tabel berikut:

Confusion	Matrix	Nilai Aktual	
Conjusion	11141114	Positif	Negatif
Nilai Prediksi	Positif	True Positives	False Positives
	Negatif	False Negatives	True Negatives

Gambar 2: Tabel confusion matrix

Confusion matrix bertujuan untuk memvisualisasikan prediksi dan keadaan aktual dari data yang dihasilkan oleh algoritma machine learning. Hal ini dilakukan dengan menghitung akurasi, presisi, recall, dan F1-score. Metrik-metrik ini sangat penting dalam mengevaluasi kinerja klasifikasi atau algoritma machine learning yang digunakan untuk melakukan prediksi. Rumus untuk keempat metrik ini ditunjukkan dalam tabel berikut:

No.	Pengukuran	Rumus
1.	Akurasi	$Accuracy = \frac{TP + TN}{TP + FP + FN + TN}$
2.	Presisi	$Precision = \frac{TP}{TP + FP}$
3.	Recall	$Recall = \frac{TP}{TP + FN}$
4.	F1-Score	$F1 - Score = \frac{2 \times Recall \times Precision}{Recall + Precision}$

Gambar 3: Tabel rumus confusion matrix

HASIL DAN PEMBAHASAN

Berikut ini merupakan tahapan-tahapan dalam mengklasifikasikan dataset mengenai kualitas air pada penelitian ini dengan menggunakan serta membandingkan antara 2 algoritma model mechine learning yaitu algoritma KNN dan algoritma SVM:

1. Exploratory Data Analysis (EDA)

Pada tahapan ini bertujuan untuk melakukan analisis pada data yang akan diuji, seperti memahami struktur, karakteristik dan informasi yang ada pada dataset yang dianalisis. Selain itu tahapan ini juga digunakan dalam melakukan identifikasi pola dan korelasi antar variable sebelum melakukan analisis statistik yang lebih mendalam lagi pada Pembangunan model yang akan diuji.

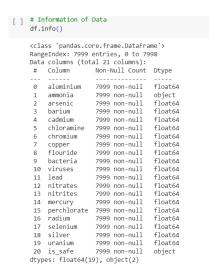
A. Jumlah kolom dan baris

```
[ ] # Dataset Dimension
df.shape
(7999, 21)
```

Gambar 4: Jumlah baris dan kolom

Dataset pada penelitian ini terdiri dari 7999 baris dan 21 kolom

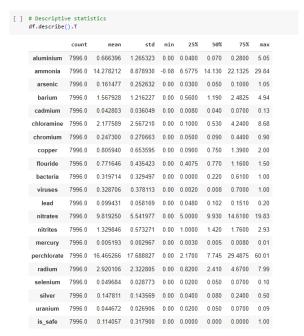
B. Mencari tipe data



Gambar 5: Tipe data

Dataset pada penelitian ini memiliki 2 tipe data yaitu 19 atribut bertipe data float dan atribut lainya bertipe data object.

C. Ringkasan statistic



Gambar 6: Ringkasan statistik

Ringkasan statistik ini berguna dalam membantu analisis statistik secara keseluruhan.

Keterangan:

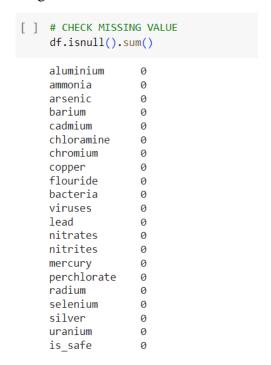
- ➤ Count = jumlah data
- ➤ Mean = rata-rata
- ➤ Std = standar deviasi
- ➤ Min = nilai minimum
- > 25 % = kuartil pertama
- > 50 % = kuartil kedua atau median
- > 75 % = kuartil ketiga
- ➤ Max = nilai maksimum

2. Preprocessing Data

Sebelum melakukan pengujian model maka sangatlah diperlukan tahapan pra pemrosesan data yaitu dengan melakukan serangkaian langkah untuk merapikan data sebelum memulai analisis. Sebelum menuju ke tahap pemprosesan. Data mentah akan diolah terlebih dahulu dengan tujuan untuk memastikan data dalam keadaan siap untuk diuji.

A. Data cleaning (Check Missing Value)

Tahap ini sangatlah penting dilakukan dalam melakukan prepoccessing data, karena memastikan bahwa data tidak memiliki missing value atau tidak berdampak pada kualitas data itu sendiri. Data dengan kualitas yang baik adalah data yang terbebas dari adanya missing value. Missing value ini juga memiliki dampak terhadap pengujian model karena missing value dapat berdampak kepada nilainilai yang diperlukan seperti mean, modus, media dan standar deviasi yang berdampak kepada tingkat keakuratan model tersebut.



Gambar 7: Check missing value

Dari gambar diatas dapat disimpulkan bahwa data yang diperoleh bebas dari adanya missing value.

B. Feature scaling

Gambar 8: Feature scaling

Tahap ini bertujuan untuk memastikan bahwasanya skala pada data ada di tingkatan yang sama, tahapan ini juga penting dilakukan pada preprocessing data sebelum melakukan pengujian model

C. Split data

Proses spilt data dilakukan sebelum memasuki pengujian model yaitu dengan cara membagi data menjadi 2 yaitu data testing dan data training. Split data ini sangat berguna dalam melakukan pengujian model, seperti membantu mengevaluasi kinerja model, mencegah overfitting dan lain-lain.

```
[ ] #SPLIT DATA

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.2, stratify = Y, random_state = 42)
```

Gambar 9: Split data

Berdasarkan gambar diatas data data dibagi menjadi 2 yaitu 20 % data testing dan 80 % data training.

D. Uji KNN

Gambar 10: Library K-NN

Berdasarkan gambar diatas dilakukan pemanggilan library KNN terlebih dahulu sebelum melakukan ngeujian model yang ingin diuji, dan kemudian pada pengujian model KNN ini parameter k yang mewakili jumlah tetangga terdekat yang akan dipertimbangkan adalah 4 tertangga terdekat dengan menggunakan data train (data latih)

Gambar 11: Uji K-NN

Kode pada gambar diatas digunakan untuk menampilkan hasil prediksi dengan model KNN, dengan menggunakan model K-NN yang telah dilatih sebelumnya untuk memprediksi label data pengujian ('X_test'). Hasil prediksi disimpan dalam DataFrame 'Pred', yang mencakup kolom 'Prediction' (hasil prediksi numerik).

Gambar 12: Evaluasi model K-NN

Pada gambar diatas dapat dilihat hasil dari pengujian knn didaptkan nilai akurasinya adalah 89%. Yang artinya tingkatan akurasi pada model KNN ini sebesar 89%.

E. Uji SVM

```
[ ] # Calling the Support Vector Machine library
    from sklearn.svm import SVC
    model = SVC(kernel='rbf')
    model.fit(X_train, Y_train)
```

Gambar 13: Library SVM

Gambar diatas yaitu pemanggilan library SVM dalam melatih model SVM dengan kernel RBF.

```
[ ] # Report of Model Evaluation
    print(f"Report Model Evaluation: \n {classification report(Y test, y pred SVM1)}")
    print(f"Confusion Matrix: \n {confusion_matrix(Y_test, y_pred_SVM1)}")
    Report Model Evaluation:
                               recall f1-score
                  precision
                                                  support
             0.0
                       0.91
                               0.99
                                          0.95
                                                    1418
             1.0
                       0.75
                                0.28
                                          0.41
                                                    182
        accuracy
                                          0.91
                                                    1600
       macro avg
                      0.83
                                0.63
                                          0.68
                                                    1600
    weighted avg
                      0.90
                                0.91
                                          0.89
                                                    1600
    Confusion Matrix:
     [[1401 17]
     [ 131 51]]
```

Gambar 14: Evaluasi model SVM

Pada gambar diatas dapat dilihat hasil dari pengujian SVM didaptkan nilai akurasinya adalah 91%. Yang artinya tingkatan akurasi pada model KNN ini sebesar 91%.

KESIMPULAN DAN SARAN

Berdasarkan hasil dan pembahasan diatas dalam perbandingan antara model K-Nearest Neighbors (KNN) dan model Support Vector Machine (SVM), hasil menunjukkan bahwa model SVM memiliki tingkat akurasi yang sedikit lebih tinggi, yakni mencapai 91%, dibandingkan dengan model KNN yang mencapai 89%. Hal ini berarti bahwa model SVM lebih unggul sedikit dalam memodelkan hubungan dan pola dalam data yang digunakan dalam pengujian. Tingkat akurasi yang tinggi pada kedua model yaitu melebihi 85%, menunjukkan bahwa keduanya memiliki kemampuan yang baik dalam mengklasifikasikan data dengan benar.

DAFTAR REFERENSI

- Bachtiar, F., Syahputra, I., & Wicaksono, S. (2019). PERBANDINGAN ALGORITME MACHINE LEARNING UNTUK MEMPREDIKSI PENGAMBIL MATA KULIAH. *Jurnal Teknologi Informasi dan Ilmu Komputer (JTIIK)*, 6(5), 543–548. https://doi.org/10.25126/jtiik.2019611755
- Han, J., Kamber, M., & Pei, J. (2012). Techniques to Improve Classification Accuracy. In *Data Mining, Concepts and Techniques*.
- Maulida, A. (2020). Penerapan Metode Klasifikasi K-Nearest Neigbor pada Dataset Penderita Penyakit Diabetes. 1(2), 29–33.
- McCartney, P. R. (2015). Big Data Science. In *MCN The American Journal of Maternal/Child Nursing* (Vol. 40, Nomor 2). https://doi.org/10.1097/NMC.000000000000118
- Mutrofin, S., Izzah, A., Kurniawardhani, A., & Masrur, M. (2014). OPTIMASI TEKNIK KLASIFIKASI MODIFIED K NEAREST NEIGHBOR MENGGUNAKAN ALGORITMA GENETIKA. *GAMMA*, *10*, *Nomor*(September), 130–134. https://doi.org/10.1017/9781316534946.021
- Said, H., Matondang, N. H., & Irmanda, H. N. (2022). Penerapan Algoritma K-Nearest Neighbor Untuk Memprediksi Kualitas Air Yang Dapat Dikonsumsi. *Techno.Com*, 21(2), 256–267. https://doi.org/10.33633/tc.v21i2.5901
- Savitri, L., & Nursalim, R. (2023). Klasifikasi Kualitas Air Minum menggunakan Penerapan Algoritma Machine Learning dengan Pendekatan Supervised Learning. *Diophantine Journal of Mathematics and Its Applications*, 2(01), 30–36. https://doi.org/10.33369/diophantine.v2i01.28260
- Suyanto, D. (2017). *Data Mining Untuk Klasifikasi dan Klastering Data*. Informatika Bandung.