ISSN: 2355-9365

Klasifikasi Kualitas Sungai Air Menggunakan Metode Pembelajaran Mesin k-Nearest Neighbour

1st Ahmad Fauzan Muhtar
Fakultas Teknik Elektro
Universitas Telkom
Bandung, Indonesia
afauzann@student.telkomuniversity.ac.

2nd IG. Prasetya Dwi Wibawa
Fakultas Teknik Elektro
Universitas Telkom
Bandung, Indonesia
prasdwibawa@telkomuniversity.ac. id

3rd Meta Kallista.

Fakultas Teknik Elektro

Universitas Telkom

Bandung, Indonesia

metakallista@telkomuniversity.ac.i d

id

Abstrak — Kualitas air Sungai Citarum di Jawa Barat masih belum memenuhi baku mutu air sepanjang tahun, terutama pada musim kemarau. Sungai ini masih tercemar berat akibat tingginya pencemar yang berasal dari berbagai aktivitas manusia seperti pertanian, peternakan, perikanan, industri, dan kegiatan domestik. Selain itu, sungai ini juga menjadi tempat pembuangan sampah rumah tangga, meningkatkan risiko masukan bahan pencemar yang mempengaruhi kualitas air. Kondisi tersebut diindikasikan oleh kekeruhan fisik sungai dan adanya pabrik limbah industri, pembuangan tinja perumahan, serta banyaknya sampah di aliran sungai. . Fokus utama dalam penelitian ini membahas tentang pemetaan kualitas air yang kualitas air di beberapa titik aliran sungai citarum menggunakan metode k-Nearest Neighbour (KNN). Parameter yang digunakan sebagai masukan dalam pembelajaran KNN seperti kekeruhan, pH, zat padat terlarut, dan suhu, diperoleh dengan instalasi sensor di beberapa titik di sepanjang aliran sungai citarum dan komunikasi jarak jauh menggunakan IoT. Dataset untuk pelatihan model pembelajaran KNN diperoleh dengan mengumpulkan beberapa sampel air dan telah melakukan pengukuran dengan alat ukur standar untuk selanjutnya dilakukan proses kalibrasi dengan sensor yang digunakan. Hasil dari model pembelajaran KNN menunjukkan akurasi sebesar 85%, dengan jumlah tetangga terdekat di k = 9 menggunakan 300 dataset. Hasil pengujian di beberapa titik sungai citarum dalam beberapa kali percobaan menunjukkan variasi indeks kualitas aliran air sungai Citarum.

Kata kunci— k-Nearest Neighbour, Dataset, Sensor, Akurasi, Air Sungai Citarum

I. PENDAHULUAN

Air merupakan kebutuhan yang paling penting bagi makhluk hidup khususnya manusia itu sendiri, salah satu sumber air tersebut adalah air sungai yang mengalir di sekitar kehidupan manusia. Sumber air tersebut tidak luput dari zat zat pengotor yang mencemari air tersebut [1]. kurangnya penelitian tentang distribusi rinci dari bahanbahan ini menyebabkan keterbatasan pengetahuan tentang dampak dan persebarannya secara detail. laporan yang diterbitkan oleh Organisasi Kesehatan Dunia (WHO), sekitar 1,8 miliar orang bergantung pada sumber air yang

terkontaminasi. Polusi air saat ini merupakan ancaman ekologis yang paling serius di dunia [2].

Setiap negara memiliki memiliki indeks kualitas air yang berbeda dengan parameter dan metode penilaian yang berbeda. Jepang memakain Indeks Kualitas Air (Water Quality Index, WQI), WQI adalah indeks yang digunakan untuk mengevaluasi kualitas air permukaan di berbagai lokasi. Indeks ini dikembangkan oleh Kementerian Lingkungan Hidup Jepang dan mempertimbangkan parameter fisik dan kimia seperti oksigen terlarut, kekeruhan, pH, nitrat, fosfat, bahan organik terlarut, logam berat, dan bahan organik total [3]. Pengukuran kualitas air di beberapa negara mengambil sampel air dikumpulkan dalam jumlah besar dari jarak minimal lima kaki di dalam botol plastik steril yang telah dicuci sebelumnya. Di lokasi pengambilan sampel, dilakukan analisis terhadap parameter pH, konduktivitas listrik (EC), total padatan terlarut (TDS), suhu, dan salinitas (dinyatakan sebagai NaCl) menggunakan perangkat HQ40-D Multi model meter buatan HACH (USA) (Hossain et al., 2010) [8].

Penelitian yang diambil dalam makalah ini adalah memantau dan melakukan klasifikasi air sungai Citarum. Menurut Kementerian Lingkungan Hidup dan Kehutanan (KLHK) di Indonesia, salah satu masalah utama pencemaran sungai terjadi di DAS Citarum, Cisadane, dan Ciujung. Sungai Citarum, yang mengalir dari Kabupaten Bandung ke Kabupaten Bekasi, telah menarik perhatian internasional karena kondisinya yang sangat tercemar. Data menunjukkan bahwa 54% air sungai Citarum tercemar berat, 23% tercemar sedang, 20% tercemar ringan, dan hanya 3% yang memenuhi standar kualitas air. Meskipun masih termasuk dalam daftar lima sungai paling tercemar di dunia, diperlukan inovasi dalam pengembangan alat pengukuran kualitas air di sepanjang daerah aliran sungai ini.[4]. Metode machine learning k-NN digunakan untuk membantu dalam penelitian kualitas air sungai di Citarum dengan mengambil sampel air dan melakukan klasifikasi. K-Nearest Neighbor merupakan salah satu algoritma yang paling sederhana yang digunakan dalam machine learning untuk regresi dan klasifikasi [5].

Penggunaan kNN yang mudah seperti alokasi memori yang tidak besar untuk menyimpan model pembelajaran,

waktu pembelajaran dan pengklasifikasian yang efisien, implementasi system yang sangat mudah di ubah ke dalam bahasa C dan terakhir kecepatan komputasi relative cepat dalam melakukan prediksi atau klasifikasi menjadikan penggunaan metode kNN cocok dalam pembelajaran mesin di system pada makalah ini [6]. Hal ini memungkinkan data yang dihasilkan dapat diperoleh secara real-time dan langsung dikirim ke platform Internet of Things (IoT) serta aplikasi seluler yang memudahkan akses pengguna. Penelitian tentang Sungai Citarum masih melibatkan pengambilan sampel air secara manual di lapangan, diikuti sampel di laboratorium dengan pengujian mendapatkan hasil yang dianalisis oleh pakar yang berkompeten [7].

Dalam makalah ini proses pemantauan dan klasifikasi kualitas air Sungai Citarum. Menggunakan Metode yang berbeda dengan penelitian sebelumnya karena melibatkan pengambilan sampel langsung di lapangan menggunakan perangkat keras yang dirancang khusus untuk efisiensi penggunaan. Setelah itu, klasifikasi kualitas air dilakukan menggunakan metode K-Nearest Neighbors (KNN) dalam pembelajaran mesin menggunakan perangkat lunak.

II. KAJIAN TEORI

K-Nearest Neighbors (KNN) adalah algoritma pembelajaran mesin yang serbaguna yang digunakan untuk tugas klasifikasi dan regresi. Konsep dasar KNN didasarkan pada asumsi bahwa data dengan atribut yang mirip cenderung memiliki label atau nilai target yang mirip juga. Dalam pendekatan KNN, objek yang akan diklasifikasikan diatributkan ke kelas mayoritas dari K tetangga terdekatnya dalam ruang fitur.Proses KNN melibatkan tiga langkah utama. Pertama, algoritma menghitung jarak antara objek yang akan diklasifikasikan dengan semua objek dalam dataset pelatihan menggunakan metrik jarak seperti jarak Euclidean. Selanjutnya, algoritma memilih K tetangga terdekat dengan objek berdasarkan jarak terkecil yang dihitung sebelumnya. Terakhir, klasifikasi ditentukan melalui proses pemungutan suara di antara K tetangga, dengan kelas mayoritas menjadi kelas yang diprediksi untuk objek tersebut. Misalnya, jika sebagian besar tetangga termasuk dalam kelas A, maka objek akan diklasifikasikan sebagai milik kelas A. KNN memiliki beberapa keunggulan yang membuatnya menjadi pilihan populer dalam berbagai aplikasi.

Pertama, kesederhanaan dan intuitifnya membuatnya mudah diimplementasikan bahkan bagi pemula dalam pembelajaran mesin. Selain itu, KNN adalah algoritma non-parametrik, artinya tidak berasumsi tentang distribusi data tertentu, sehingga cocok untuk menangani data yang kompleks dan tidak terstruktur. Selain itu, KNN beradaptasi dengan baik pada dataset yang dinamis karena tidak melibatkan proses pembelajaran, sehingga memungkinkan pembaruan data langsung tanpa harus mengulangi tahap pelatihan. Namun, perlu diakui bahwa KNN juga memiliki beberapa kelemahan. Algoritma ini sensitif terhadap nilai K yang dipilih, dan nilai yang tidak tepat dapat menyebabkan overfitting atau underfitting. Selain itu, data outlier dalam dataset dapat signifikan mempengaruhi kinerja KNN dan menyebabkan kesalahan klasifikasi. Selain itu, KNN

memerlukan penyimpanan data pelatihan yang besar karena harus mengakses seluruh dataset saat tahap prediksi. Untuk memanfaatkan kelebihan KNN sambil mengatasi kelemahannya, pemilihan nilai K yang tepat dan langkahlangkah pra-pemrosesan data menjadi penting. Secara keseluruhan, KNN adalah algoritma yang berharga dan banyak digunakan dalam domain pembelajaran mesin, menyediakan pendekatan yang sederhana namun efektif dalam tugas klasifikasi dan prediksi..

A. Klasifikasi kualitas air menggunakan k-NN

Klasifikasi menggunakan K-Nearest Neighbors (KNN) merupakan salah satu metode populer dalam machine learning yang sering digunakan untuk mengklasifikasikan objek berdasarkan atribut-atributnya dan tingkat kemiripan dengan objek-objek yang sudah ada dalam dataset. KNN adalah algoritma klasifikasi non-parametrik yang memiliki keunggulan tidak memerlukan asumsi tertentu tentang distribusi data, sehingga dapat digunakan dalam berbagai <mark>kasus dengan fleksibilitas y</mark>ang tinggi. Idea dasar dari algoritma KNN adalah mencari K titik terdekat (neighbors) dari objek yang akan diklasifikasikan berdasarkan jarak Euclidean atau metrik lainnya. Dalam proses ini, KNN melakukan perhitungan jarak antara data yang akan diklasifikasikan dengan data lain dalam dataset. Data yang memiliki jarak terdekat dengan data yang diklasifikasikan akan dianggap sebagai tetangga terdekat. Setelah mendapatkan K tetangga terdekat, proses klasifikasi dilakukan dengan memeriksa mayoritas kelas dari tetanggatetangga tersebut. Jika mayoritas tetangga termasuk dalam suatu kelas, maka objek yang akan diklasifikasikan juga akan diklasifikasikan ke dalam kelas tersebut. Pendekatan ini sangat sederhana dan intuitif, sehingga KNN sering digunakan sebagai langkah pertama dalam analisis data dan pemahaman masalah klasifikasi.

Meskipun sederhana, KNN dapat memberikan hasil yang baik dalam kasus-kasus di mana data memiliki pola yang mudah dipisahkan atau ketika jumlah data pelatihan yang terbatas. Namun, perlu dicatat bahwa pemilihan nilai K yang tepat sangat penting dalam mengoptimalkan kinerja algoritma KNN, serta perlu diwaspadai terhadap masalah overfitting atau underfitting yang mungkin terjadi. Dengan pemilihan nilai K dan pemrosesan data yang tepat, KNN dapat menjadi alat yang efektif untuk melakukan klasifikasi dan memahami karakteristik data dengan lebih baik..

III. METODE

Penggunaan metode yang tepat dalam uji kualitas air sungai Citarum sangat tergantung pada parameter yang ingin diuji dan tujuan pengujian. Untuk mendapatkan gambaran yang lebih lengkap tentang kualitas air, seringkali digunakan kombinasi metode fisik, kimia, dan biologi. Oleh karena itu, beberapa instansi kampus telah melakukan penelitian dengan metode dan pendekatan yang berbeda-beda. Dalam studi pendahuluan yang dilakukan oleh beberapa penelitian, langkah awal adalah mengambil sampel air secara langsung di aliran Sungai Citarum. Kemudian, sampel-sampel tersebut diuji di laboratorium kesehatan daerah (Labkesda) untuk menentukan kualitas air dan mencari tahu pencemar dominan yang mungkin ada. Namun, penelitian yang diangkat dalam

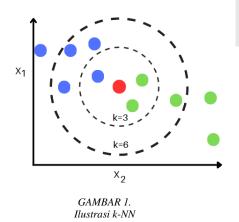
penelitian ini berbeda. Dalam konteks pengukuran kualitas air, metode machine learning menggunakan algoritma k-Nearest Neighbors (k-NN) diadopsi sebagai pengklasifikasi untuk menentukan kualitas air berdasarkan data yang telah diambil sebelumnya. Proses pembelajaran k-NN dilakukan dengan menggunakan data pelatihan yang dikumpulkan dari pengukuran air sebelumnya. Setelah proses pembelajaran selesai. model k-NN dapat secara otomatis mengklasifikasikan kualitas air berdasarkan data masukan yang baru. Dengan memanfaatkan teknik machine learning, metode k-NN memberikan pendekatan yang inovatif dalam menganalisis kualitas air sungai Citarum. menggunakan data historis dan pengukuran sebelumnya, model k-NN dapat memberikan hasil klasifikasi secara realtime dan membantu dalam pemantauan dan pemahaman yang lebih efektif tentang kualitas air sungai tersebut.. Secara umum untuk mengetahui rumus formula Euclidean distance pada 1 – dimensional space seperti berikut.

$$dis(x_1, x_2) = \sqrt{\sum_{i=0}^{n} (x_{1i} - x_{2i})^2}$$

Tetapi jika variabelnya tidak independent dan jumlah varibel lebih dari satu maka rumus formula yang digunakan sebagai berikut.

$$dis = \sqrt{\sum_{i=0}^{n} (x_{1i} - x_{2i})^2 + (y_{1i} - y_{2i})^2 + \cdots}$$

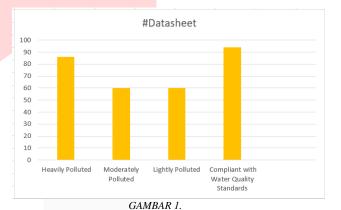
Kelebihan dari KNN adalah kemampuannya dalam menangani data non-linear dan kemampuan adaptasinya terhadap perubahan pada dataset. Namun, KNN juga memiliki beberapa kelemahan, seperti ketergantungan terhadap skala atribut, sensitivitas terhadap data yang tidak relevan, dan biaya komputasi yang tinggi pada dataset besar.Dalam prakteknya, KNN digunakan dalam berbagai aplikasi seperti klasifikasi data, pengenalan pola, rekomendasi, dan sistem pengenalan suara. Penting untuk melakukan pemilihan nilai K yang optimal, normalisasi data, dan validasi model untuk mendapatkan hasil yang akurat dan baik dalam klasifikasi menggunakan KNN.



IV. HASIL DAN PEMBAHASAN

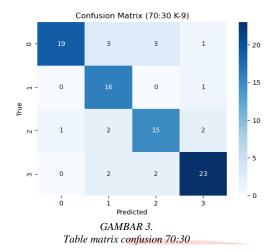
Pengklasifikasian menggunakan algoritma KNN sangat penting dalam memproses dan mengklasifikasikan data dengan akurat. Persiapan data yang tepat, normalisasi data untuk menghilangkan bias, pembagian data menjadi data pelatihan dan pengujian, serta pemilihan parameter K yang optimal merupakan langkah awal yang krusial dalam proses ini. Selanjutnya, melatih model KNN dengan menggunakan data pelatihan dan mengklasifikasikan data pengujian berdasarkan mayoritas label dari K tetangga terdekat memainkan peran penting dalam menentukan kelas klasifikasi. Dengan langkah-langkah yang tepat, model KNN dapat digunakan untuk mengklasifikasikan data baru tanpa label dengan hasil yang dapat diandalkan

K yang terbaik karena nilai akurasinya tinggi dan ratarata matriks evaluasi tinggi diantara K yang lain, dan berwarna merah juga adalah K yang ganjil tertinggi nilainya.



GAMBAR 2. Table confusion matrix 80:20 dan 75:25

Dengan data di atas, kita dapat membentuk confusion matrix untuk mengevaluasi performa model KNN dengan perbandingan 80:20 k-9. Dengan menggunakan nilai TP, TN, FP, dan FN yang diperoleh, kita dapat menghitung berbagai metrik evaluasi seperti akurasi, presisi, recall, dan F1-score untuk memahami sejauh mana model tersebut berhasil dalam melakukan klasifikasi.

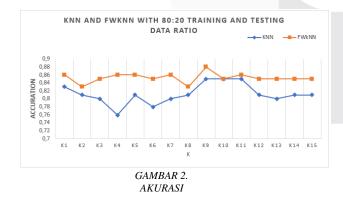


Dengan data di atas, kita dapat membentuk confusion matrix untuk mengevaluasi performa model KNN dengan data dan uji 70:30 k-9. Berdasarkan nilai TP, TN, FP, dan FN yang diperoleh, kita dapat menghitung berbagai metrik evaluasi seperti akurasi, presisi, recall, dan F1-score untuk memahami sejauh mana model tersebut berhasil dalam melakukan klasifikasi pada data pengujian.

TABEL 1. kNN vs FWkNN rasio latih uji 80:20

К	ACURATION		PRESISSION		RECALL		F1-SCORE	
	KNN	EXYLNN	KNN	FWkNN	KNN	FWkNN	KNN	EWKNN
K1	0,83	0,86	0,84	0,89	0,84	0,87	0,83	0,87
K2	0,81	0,83	0,82	0,84	0,83	0,85	0,81	0,83
K3	0,8	0,85	0,8	0,86	0,79	0,84	0,79	0,84
K4	0,76	0,86	0,77	0,87	0,76	0,87	0,75	0,86
K5	0,81	0,86	0,81	0,87	0,81	0,85	0,81	0,85
K6	0,78	0,85	0,77	0,85	0,78	0,85	0,77	0,84
K7	0,8	0,86	0,8	0,87	0,79	0,85	0,79	0,85
K8	0,81	0,83	0,82	0,84	0,82	0,82	0,81	0,82
K9	0,85	0,88	0,86	0,89	0,86	0,87	0,85	0,87
K10	0,85	0,85	0,86	0,86	0,86	0,84	0,85	0,84
K11	0,85	0,86	0,85	0,87	0,85	0,84	0,84	0,86
K12	0,81	0,85	0,82	0,86	0,82	0,84	0,81	0,84
K13	0,8	0,85	0,80	0,85	0,8	0,84	0,79	0,83
K14	0,81	0,85	0,82	0,85	0,81	0,84	0,80	0,83
K15	0,81	0,85	0,82	0,86	0,81	0,83	0,79	0,83

DISTRIBUSI JUMLAH DATA BERDASARKAN UJI TEST PERBANDINGAN

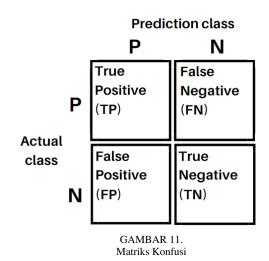






GAMBAR 4.

Dalam menentukan k terbaik dari grafik diatas pembelajaran mesin menggunakan metode k-Nearest Neighbors (KNN), kita dapat menggunakan teknik validasi silang (cross-validation). Dalam konteks ini, kita akan membagi dataset menjadi subset pelatihan dan pengujian, dan mengulangi proses ini beberapa kali dengan menggunakan nilai k yang berbeda untuk setiap literasi.



Rumus untuk menghitung metrik evaluasi yang umum digunakan dalam KNN adalah sebagai berikut:

- Akurasi (Accuracy): Akurasi = (TP + TN) / (TP + TN + FP + FN)
- Presisi (Precision): Presisi = TP / (TP + FP)
- Sensitivitas atau Recall: Sensitivitas = TP / (TP + FN)
- F1-score:

F1-score = 2 * (Presisi * Sensitivitas) / (Presisi + Sensitivitas)

Dalam validasi silang, kita menghitung metrik evaluasi untuk setiap nilai k pada setiap iterasi, dan akhirnya memilih nilai k yang memberikan performa terbaik dalam hal metrik evaluasi yang diinginkan. Hal ini membantu kita menentukan k terbaik untuk model KNN dalam pembelajaran mesin. Confusion matrix adalah tabel yang digunakan untuk membandingkan hasil prediksi model dengan label yang sebenarnya pada data pengujian dalam evaluasi kinerja model klasifikasi. Dalam confusion matrix, terdapat empat istilah yang masing-masing memiliki arti dan peranan penting. True Positives (TP) mengacu pada jumlah data yang benar-benar diklasifikasikan dengan benar oleh model sebagai positif (kelas yang diinginkan). True Negatives (TN) adalah jumlah data yang benar-benar diklasifikasikan dengan benar oleh model sebagai negatif (kelas yang bukan diinginkan). Di sisi lain, False Positives (FP) merujuk pada data yang sebenarnya adalah negatif tetapi salah diklasifikasikan sebagai positif oleh model, sedangkan False Negatives (FN) menggambarkan jumlah data yang sebenarnya adalah positif tetapi diklasifikasikan sebagai negatif oleh model. Confusion matrix memberikan informasi yang berguna dalam menilai performa model secara lebih rinci, dan dengan nilai TP, TN, FP, dan FN, kita dapat menghitung berbagai metrik evaluasi seperti akurasi, presisi, sensitivitas (recall), dan F1-score untuk mendapatkan pemahaman yang lebih komprehensif tentang kinerja model klasifikasi.

KESIMPULAN V.

Penggunaan metode pembelajaran mesin menggunakan k-Nearest Neighbours (KNN) dalam melakukan klasifikasi kualitas air sungai memiliki beberapa keunggulan yang signifikan. Metode KNN telah terbukti efektif dalam mengklasifikasikan kualitas air sungai berdasarkan data pengukuran seperti kekeruhan, pH, zat padat terlarut, dan suhu. Dengan menghitung jarak antara objek yang akan diklasifikasikan dengan data pelatihan sebelumnya, KNN

dapat mengidentifikasi kelas mayoritas dari K tetangga terdekatnya dalam ruang fitur, sehingga memungkinkan penentuan kualitas air dengan cukup akurat. Selain itu, penggunaan sensor di beberapa titik aliran Sungai Citarum dan penerapan teknologi Internet of Things (IoT) dalam metode ini memungkinkan pengumpulan data pengukuran secara real-time dan pengiriman data ke platform IoT untuk analisis lebih lanjut. Dengan adanya implementasi sensor dan IoT, pemantauan kualitas air sungai dapat dilakukan secara efisien dan mendukung upaya pengelolaan lingkungan yang lebih baik.

Penelitian ini memberikan kontribusi signifikan dalam memahami tingkat pencemaran air Sungai Citarum dan memiliki implikasi positif untuk pengelolaan dan pemulihan kualitas air. Dengan metode KNN dan teknologi sensor-IoT yang digunakan, pihak berwenang dapat dengan cepat mengidentifikasi daerah yang tercemar dan mengambil tindakan korektif yang diperlukan. Selain itu, penggunaan KNN dan teknologi sensor-IoT juga memungkinkan pemantauan hasil upaya pemulihan secara lebih efisien. Dengan demikian, penggunaan metode pembelajaran mesin menggunakan KNN dalam melakukan klasifikasi kualitas air sungai memberikan potensi besar dalam meningkatkan pemantauan dan pengelolaan lingkungan sungai, serta berkontribusi dalam menjaga kualitas air yang lebih baik bagi masyarakat dan ekosistem yang bergantung pada Sungai Citarum.

REFERENSI

"Impact of Industrial Wastewater on Water Quality: A Case Study of XYZ River" - Jurnal: Water Research. Richard Helmer dan Ivanildo Hespanhol "Water Pollution Control: A Guide to the Use of Water Quality Management Principles" (2008)

Management Principles" (2008).

T. Widodo, M. T. S. Budiastuti, dan K. Komariah, "Water Quality and Pollution Index in Grenjeng River, Boyolali Regency, Indonesia," Caraka Tani J. Sustain. Agric., vol. 34, no. 2, hal. 150, 2019, doi: 10.20961/carakatani.v34i2.29186

"Development of a water quality index model for river systems: A case study in Japan" oleh Y. Li et al.

Flem, B. et al. Inorganic chemical quality of European tap-water: Geographical distribution. Appl.

tap-water: 2. Geographical distribution. Appl. Geochem. **59**, 211–224 (2015). Hossain F, Chang NB, Wanielista M, Xuan Z, Daranpob A (2010) Nitrification and denitrification in a passive on-site wastewater treatment system with a recirculation filtration tank. Water Qual Expo Health 2:31–46.

D. H. Jayani, "Indonesia pada tahun 2019 tercatat 314 anak balita meninggal dunia karena penyakit diare le," databoks, https://databoks.katadata.co.id/datapublish/2021/04/26/ diare-penyebab-utama-kematian-anak-di-indonesiapada-2019#:~:text=Kementerian Kesehatan (Kemenkes) mencatat penyebab,%2C dan malaria (22). (diakses 13 Oktober 2022).

T. Cover.,K-Nearest Neighbors - 25 years later".

H. Utomo et al, "Evaluation of Water Quality and

Potential Human Health Risks in the Citarum River, Indonesia" (2018).