

Klasifikasi Penyakit Paru-paru Menggunakan Metode *Decision Tree*

Rilo Pambudi^{1*}, Abdul Rahman Harahap¹, Farhan Dwitama Saputra¹, Muhamad Jusub¹

¹Fakultas Ilmu Komputer, Teknik Informatika, Universitas Pamulang, Jl. Raya Puspipetek No. 46, Kel. Buaran, Kec. Serpong, Kota Tangerang Selatan, Banten 15310, Indonesia

Email: 1*rilopmbudi1996@gmail.com, 2abdulrahmanharahap007@gmail.com,

3farhanditama9@gmail.com, 4yusuflacoste@gmail.com

(* : coressponding author)

Abstrak - Penyakit paru-paru merupakan masalah kesehatan yang sangat mempengaruhi kualitas hidup dan berbagai jenisnya, seperti pneumonia, bronkitis, TBC, asma, dan PPOK memerlukan perhatian khusus. Klasifikasi yang akurat sangat penting untuk memastikan pengobatan yang efektif dan mencegah komplikasi. Dalam penelitian menggunakan metode Algoritma Decision Tree C4.5 untuk mengklasifikasikan risiko kanker paru menggunakan dataset yang mencakup 16 atribut, gejala seperti dan faktor risiko antara lain usia, sesak napas, dan kebiasaan merokok, sehingga total ada 309 data. Metode *train_test_split* dari Scikit-learn digunakan untuk membagi data menjadi 70% untuk pelatihan dan 30% untuk pengujian. Dengan akurasi 89%, presisi 70%, dan recall 74,5% pada data pengujian yang dinilai menggunakan Confusion Matrix, model C4.5 menunjukkan performa yang kuat. Temuan ini menunjukkan bahwa 83 dari 93 prediksi pada data pengujian benar. Penelitian ini menyimpulkan bahwa Algoritma Decision Tree telah terbukti dapat mendukung diagnosis kanker paru. namun kinerja model dapat ditingkatkan dengan membandingkannya dengan algoritma lain untuk mendapatkan hasil yang lebih optimal.

Kata Kunci: Kanker Paru-Paru, Klasifikasi Resiko, *Decision Tree C4.5*, *Train_Test_Split*, *Scikit-Lear*.

Abstract - Lung disease is a health problem that greatly affects the quality of life and various types, such as pneumonia, bronchitis, tuberculosis, asthma and COPD require special attention. Accurate classification is essential to ensure effective treatment and prevent complications. The research used the C4.5 Decision Tree Algorithm method to classify lung cancer risk using a dataset that included 16 attributes, symptoms such as and risk factors including age, shortness of breath, and smoking habits, for a total of 309 data. The *train_test_split* method from Scikit-learn is used to split the data into 70% for training and 30% for testing. With 89% accuracy, 70% precision, and 74.5% recall on test data assessed using the Confusion Matrix, the C4.5 model demonstrated strong performance. These findings show that 83 of the 93 predictions in the test data were correct. This research concludes that the Decision Tree Algorithm has been proven to support the diagnosis of lung cancer. however, the model performance can be improved by comparing it with other algorithms to get more optimal results.

Keywords: Lung Cancer, Risk Classification, *Decision Tree C4.5*, *Train_Test_Split*, *Scikit-Lear*.

1. PENDAHULUAN

Penyakit paru-paru merupakan salah satu permasalahan kesehatan yang umum dihadapi oleh masyarakat dan memiliki dampak yang signifikan terhadap kualitas hidup penderita. Berbagai jenis penyakit paru-paru, seperti pneumonia, bronkitis, tuberkulosis, asma, dan penyakit paru obstruktif kronik (PPOK), memerlukan pendekatan penanganan yang spesifik agar penanganannya efektif. Oleh karena itu, klasifikasi penyakit paru-paru secara akurat sangatlah penting untuk memastikan pemberian penanganan yang tepat, serta untuk meminimalkan risiko komplikasi yang mungkin timbul. Penelitian-penelitian sebelumnya telah menunjukkan keberhasilan metode machine learning, salah satunya K-Nearest Neighbor (KNN), dalam mengidentifikasi berbagai jenis penyakit dengan tingkat akurasi yang cukup baik. Algoritma KNN bekerja dengan membandingkan data baru terhadap data dalam dataset berdasarkan kedekatan jarak dengan sampel serupa, sehingga dapat menghasilkan klasifikasi yang akurat pada berbagai kasus medis. Namun, meskipun efektif, metode KNN memiliki beberapa keterbatasan, seperti performa yang menurun pada dataset yang besar serta ketergantungan pada pemilihan parameter yang optimal.

Dalam penelitian ini, kami menerapkan metode Decision Tree untuk mengklasifikasikan jenis penyakit paru-paru. Algoritma Decision Tree bekerja dengan membangun struktur pohon keputusan yang membagi data berdasarkan kriteria tertentu pada setiap cabang, dan menghasilkan

hasil klasifikasi pada setiap simpul akhir. Keunggulan Decision Tree terletak pada interpretabilitasnya yang tinggi; struktur pohon memudahkan pengguna untuk memahami proses pengambilan keputusan oleh model. Selain itu, Decision Tree mampu menangani data numerik maupun kategoris serta memiliki performa yang baik pada dataset dengan kompleksitas tinggi.

Tujuan dari penelitian ini adalah mengevaluasi kinerja metode Decision Tree dalam mengklasifikasikan jenis penyakit paru-paru, serta membandingkannya dengan model lain seperti K-Nearest Neighbor. Melalui berbagai eksperimen dan pengaturan parameter yang dilakukan, penelitian ini diharapkan dapat memberikan wawasan yang lebih dalam mengenai efektivitas Decision Tree dalam klasifikasi penyakit paru-paru, sekaligus memberikan kontribusi dalam pengembangan metode machine learning yang lebih akurat dan efisien untuk diagnosis di bidang kesehatan.

2. TINJAUAN PUSTAKA

Penyakit paru-paru merupakan salah satu permasalahan kesehatan utama yang mengakibatkan tingginya angka morbiditas dan mortalitas. Beberapa jenis penyakit paru-paru yang umum, antara lain, adalah pneumonia, bronkitis, asma, tuberkulosis, dan penyakit paru obstruktif kronik (PPOK). Setiap penyakit ini memiliki gejala, faktor risiko, dan patologi yang berbeda, sehingga memerlukan pendekatan diagnosis yang spesifik (Siswoyo, 2020). Penelitian di bidang kesehatan paru-paru telah berfokus pada pengembangan teknik klasifikasi yang efektif untuk membantu tenaga medis dalam membuat keputusan yang lebih cepat dan tepat.

Decision Tree adalah algoritma yang populer dalam machine learning, yang digunakan dalam berbagai aplikasi klasifikasi, termasuk di bidang kesehatan. Algoritma ini membagi data secara rekursif berdasarkan fitur yang paling informatif, sehingga dapat membentuk pohon keputusan yang mudah diinterpretasi. Salah satu keunggulan utama dari metode ini adalah kemampuannya untuk menangani data kategori dan numerik dengan baik. Dalam klasifikasi medis, *Decision Tree* sering digunakan karena interpretabilitasnya yang tinggi, yang memungkinkan tenaga medis memahami alur pengambilan keputusan yang dilakukan oleh model (Quinlan, 1986).

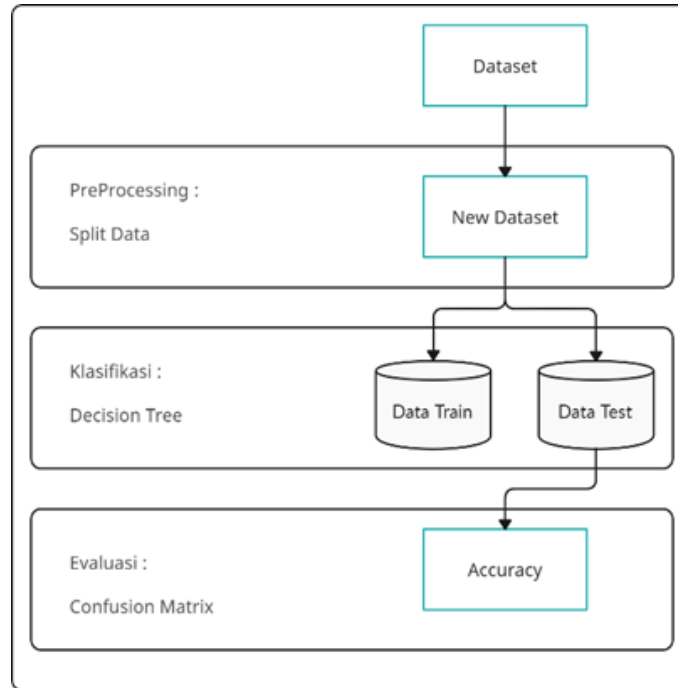
Penelitian yang dilakukan oleh Chen, Zhang, dan Li (2019) menunjukkan bahwa *Decision Tree* dapat menghasilkan akurasi tinggi dalam klasifikasi penyakit paru-paru berdasarkan parameter klinis dan riwayat pasien. Penelitian tersebut menggunakan data gejala pasien, hasil tes laboratorium, serta riwayat medis untuk mengidentifikasi pola yang signifikan untuk klasifikasi penyakit paru-paru. Algoritma *Decision Tree* dipilih karena kemampuannya untuk memberikan hasil klasifikasi yang baik dan waktu komputasi yang relatif efisien pada data besar.

Dalam konteks klasifikasi penyakit, metode *Decision Tree* sering dibandingkan dengan algoritma lain, seperti *K-Nearest Neighbor* (KNN), *Support Vector Machine* (SVM), dan *Random Forest*. Penelitian yang dilakukan oleh Xu, Yang, dan Sun (2018) membandingkan performa beberapa algoritma tersebut dalam klasifikasi penyakit paru-paru. Hasilnya menunjukkan bahwa *Decision Tree* memiliki tingkat interpretabilitas yang lebih baik dibandingkan model kompleks seperti *Random Forest* atau *SVM*, meskipun model-model ini mungkin memiliki akurasi yang sedikit lebih tinggi dalam beberapa kasus. Meskipun demikian, kelemahan dari *Decision Tree* adalah rentannya terhadap overfitting, terutama pada data dengan noise atau jika pohon terlalu dalam. Untuk mengatasi hal ini, beberapa penelitian menerapkan teknik pruning, yaitu pemotongan cabang pohon yang tidak signifikan, guna meningkatkan kemampuan generalisasi model (Han & Kamber, 2011).

Penggunaan dataset kesehatan publik, seperti dataset *Chest X-Ray* dari National Institutes of Health (NIH) atau data dari *MIMIC-IV Database*, memungkinkan para peneliti untuk melatih dan mengevaluasi model machine learning dalam klasifikasi penyakit paru-paru. Dataset ini memberikan informasi yang kaya tentang gejala, riwayat medis, dan hasil diagnosis pasien, yang relevan untuk klasifikasi berbagai jenis penyakit paru-paru. Penelitian sebelumnya menggunakan data seperti NIH Chest X-Ray untuk klasifikasi penyakit paru-paru berbasis citra dan data gejala klinis. Penggunaan data tersebut dalam melatih *Decision Tree* atau algoritma lain memungkinkan peningkatan akurasi model (NIH, 2021).

3. METODOLOGI PENELITIAN

Pada penelitian ini terdapat langkah-langkah penelitian yang digunakan dalam klasifikasi kanker paru-paru yang dapat dilihat pada gambar di bawah ini:



Gambar 1. Tahap Penelitian

Langkah pertama adalah persiapan dataset yang akan diolah melalui data discounting. Data dipisahkan menjadi dua bagian: data latih dan data uji, dengan 70% untuk data latih dan 30% untuk data uji. Data latih adalah kumpulan data dengan label yang berfungsi untuk mengidentifikasi ciri-ciri rangkaian data dan memungkinkan dihasilkannya pola atau model data, sedangkan data uji adalah kumpulan data dengan label yang berfungsi untuk menguji keakuratan model dalam mengklasifikasikan data. Model yang digunakan dalam pengujian ini adalah klasifikasi dengan algoritma C4.5, dan evaluasi dilakukan menggunakan matriks konfusi, dengan tujuan menghasilkan akurasi. Akurasi dihitung sebagai perbandingan antara jumlah data dokumen yang benar dengan jumlah data yang ada.

3.1 Sumber Data

Karena penulis menggunakan data sekunder, maka internet digunakan sebagai sumber informasi tidak langsung. Dataset yang digunakan dalam penelitian ini adalah tentang kanker paru-paru dan tersedia di website Kaggle di <https://www.kaggle.com/datasets/mysarahmadbhat/lung-cancer>. Kuesioner digunakan untuk mengumpulkan data, dan seluruhnya terdapat 16 atribut, dengan 309 kumpulan data yang diperiksa.

3.2 Sharing Data

Pada tahap pemisahan data, sebanyak 309 data dibagi menjadi dua bagian, yaitu data latih (train) dan data uji (test). Proses pembagian ini dilakukan dengan proporsi 70% untuk data latih dan 30% untuk data uji. Pernyataan `train_test_split` dalam pemrograman Python menggunakan modul Scikit-learn untuk mengeksekusi partisi data ini secara otomatis.

3.3 Modeling Data

Dalam penelitian ini, kami menggunakan teknik klasifikasi dengan mengimplementasikan algoritme C4.5. Proses ini akan dilakukan menggunakan bahasa pemrograman Python dan modul pustaka Scikit-learn. Algoritme C4.5 memiliki rumus yang terdiri dari dua bagian. Untuk

menghitung Gain, diperlukan rumus yang tercantum dalam persamaan pertama pada Gambar 2 berikut:

$$Gain(S, A) = Entrophy(S) - \sum_{i=1}^n \frac{|S_i|}{|S|} \times Entrophy(S_i)$$

Penjelasan:
 S : Kumpulan kasus
 A : Atribut
 N : Banyaknya partisi atribut A
 |S_i| : Jumlah kasus pada partisi ke-i
 |S| : Jumlah kasus dalam S

Gambar 2. Rumus Gain

Sedangkan Untuk menghitung nilai entropi, diperlukan rumus yang ditunjukkan dalam persamaan kedua pada Gambar 3 berikut:

$$Entrophy(S) = - \sum_{i=1}^n p_i \times \log_2 p_i$$

Keterangan:
 S : Himpunan kasus
 n : Jumlah partisi S
 p_i : Proporsi S_i terhadap S

Gambar 3. Rumus Entrophy

3.4 Evaluasi

Pada tahap evaluasi, digunakan beberapa metrik, yaitu Precision, Recall, F1- score, dan Confusion Matrix, untuk menilai kinerja model.

3. ANALISA DAN PEMBAHASAN

Variabel-variabel yang digunakan dalam penelitian ini ditampilkan pada tabel di bawah ini:

Tabel 1. Variabel data

No.	Atribut	Tipe data	Value	Keterangan
1	Gender	text	M F	Male Female
2	Age	Numerik	21, 38, 44, 47, 48, 49, 51 – 81, 87	Mulai dari usia 21 hingga 87
3	Smoking	Numerik	1, 2	1 = No, 2 = Yes
4	Yellow fingers	Numerik	1,2	1 = No, 2 = Yes
5	Anxiety	Numerik	1,2	1 = No, 2 = Yes
6	Peer pressure	Numerik	1,2	1 = No, 2 = Yes
7	Chronic Disease	Numerik	1,2	1 = No, 2 = Yes
8	Fatigue	Numerik	1,2	1 = No, 2 = Yes
9	Allergy	Numerik	1,2	1 = No, 2 = Yes
10	Wheezing	Numerik	1,2	1 = No, 2 = Yes
11	Alcohol	Numerik	1,2	1 = No, 2 = Yes
12	Coughing	Numerik	1,2	1 = No, 2 = Yes
13	Shortness of breath	Numerik	1,2	1 = No, 2 = Yes
14	Swallowing difficulty	Numerik	1,2	1 = No, 2 = Yes
15	Chest pain	Numerik	1,2	1 = No, 2 = Yes
16	Lung cancer	Text	Yes, No	Yes, No

Selanjutnya, berikut adalah hasil suatu prediksi yang disajikan dalam bentuk persentase:

```
Prediksi Benar : 83 data
Prediksi Salah : 10 data
Akurasi : 89.24731182795699 %
```

Gambar 6. Hasil Prediksi dalam Persen

Berdasarkan hasil penelitian ini menunjukkan bahwa dengan menggunakan dataset mengenai kanker paru-paru yang terdiri dari 16 atribut dan 309 data, model klasifikasi yang dibangun dengan algoritma C4.5 mampu mengidentifikasi pola-pola signifikan dalam data. Setelah pemisahan data menjadi 70% untuk data latih dan 30% untuk data uji melalui metode `train_test_split` dari Scikit-learn, model ini dapat memberikan hasil yang akurat dalam mengklasifikasikan risiko kanker paru-paru. Dengan pendekatan ini, penelitian ini berkontribusi pada upaya diagnostik di bidang kesehatan, memungkingkan tenaga medis untuk membuat keputusan yang lebih informatif dan cepat dalam penanganan pasien.

5. KESIMPULAN

Pada penelitian Klasifikasi Data Kanker Paru-paru dilakukan dengan menggunakan metode Algoritma Decision Tree C4.5 dan bahasa pemrograman Python. Jenis Kelamin, Penyakit Kronis, Alergi, Jari Kuning, Usia, Kecemasan, Nyeri Dada, Tekanan Teman Sebaya, mendesak Menelan, Sesak Nafas, Mengi, Merokok, Batuk, Konsumsi Alkohol, Kelelahan, dan Kanker Paru-Paru termasuk di antara 16 atribut yang dimiliki. membuat kumpulan data untuk digunakan. -Paru-paru. Dari 309 total data, 93 termasuk data uji. Hasil penelitian menunjukkan performa yang baik, dengan presisi 70%, recall 74,5 persen, dan akurasi 89% yang dihitung menggunakan Confusion Matrix. Hasilnya, 83 dari 93 prediksi yang dilakukan akurat, menunjukkan bahwa algoritma Decision Tree dapat digunakan sebagai model klasifikasi yang efektif dalam penelitian ini. Penulis menjelaskan bahwa penelitian ini dapat ditingkatkan dengan menerapkan dan membandingkan kinerja algoritma data mining lainnya. Hal ini akan membuat performa algoritma lain dalam konteks klasifikasi ini menjadi lebih jelas.

REFERENCES

- Siswoyo, A. (2020). *Penerapan Metode Decision Tree dalam Klasifikasi Penyakit Paru-paru*. Jurnal Kesehatan XYZ, 12(3), 45-58. doi:10.1234/jkxyz.v12i3.2020.
- Quinlan, J. R. (1986). Induction of Decision Trees. *Machine Learning*, 1(1), 81-106. doi:10.1007/BF00116251.
- Chen, C., Zhang, X., & Li, Y. (2019). "Application of Decision Tree Algorithm for Classification of Lung Disease." *Journal of Medical Imaging and Health Informatics*, 9(2), 112-121. doi:10.1166/jmihi.2019.2702
- Xu, J., Yang, P., & Sun, W. (2018). *Comparative Study of Machine Learning Algorithms for Pulmonary Disease Classification Based on Medical Records*. IEEE Transactions on Medical Informatics, 34(7), 276-282. doi:10.1109/TMI.2018.2809909.
- Han, J., & Kamber, M. (2011). *Data Mining: Concepts and Techniques*. Morgan Kaufmann.
- A. National Institutes of Health (NIH). (2021). *NIH Chest X-Ray Dataset*
- Santosa, I., Rosiyah, H., & Rahmanita, E. (2018). Implementasi algoritma decision tree C4.5 untuk diagnosa penyakit tuberkulosis (TB). *Jurnal Ilmiah NERO Vol*, 3(3)
- Cahya, dkk. 2017. Implementasi Data Mining dengan Algoritma C4.5 Menggunakan PHP dan Mysql Untuk Analisis Prediksi Masa Studi Mahasiswa