

KLASIFIKASI KANKER PARU-PARU MENGGUNAKAN METODE *NAIVE BAYES*

¹Eva Wulandari ²Andreas Perdana

¹Teknik Informatika STMIK Dharma Wacana Metro,

²STMIK Dharma Wacana Metro

¹wulandarieva22@gmail.com , ² andreas.perdana@dharmawacana.ac.id

ABSTRAK

Kanker paru-paru merupakan penyebab utama kematian akibat kanker dibandingkan dengan jenis kanker lainnya. Penyakit ini disebabkan oleh perubahan sel yang terus tumbuh di luar kendali. Melalui klasifikasi, pola kanker paru-paru ditemukan. Metode klasifikasi yang umum digunakan adalah metode Naive Bayes. Naive Bayes adalah algoritma klasifikasi sederhana dengan kinerja dan akurasi tinggi. Proses klasifikasi naif Bayesian terdiri dari menghitung probabilitas sebelumnya dengan menjumlahkan kasus secara rinci dengan atribut data. Kemudian menghitung probabilitas posterior untuk menentukan kelas kasus baru saat memproses data. Selanjutnya hitung peluang keseluruhan untuk setiap kelas. Dan akurasi dapat diperoleh dengan perhitungan menggunakan confusion matrix. Confusion matrik yang akan dihitung adalah accuracy, precision dan recall.

Keyword: Algoritma, Kematian,, Penyakit

1.PENDAHULUAN

Kanker paru-paru merupakan penyebab utama kematian akibat kanker dibandingkan dengan jenis kanker lainnya. Data yang diperoleh menunjukkan bahwa 32% kematian akibat kanker pada pria dan 25% pada wanita disebabkan oleh kanker paru-paru. Sebagian besar kasus kanker paru-paru menyerang orang berusia antara 35 dan 75 tahun, dengan insiden tertinggi antara 55 dan 65 tahun (Aliya dkk., 2016). Dalam sebuah penelitian, Tommy dkk. (2022) menjelaskan bahwa kanker adalah penyakit yang disebabkan oleh perubahan sel yang menyebabkan pertumbuhan dan pembelahan sel tidak terkendali. Berdasarkan permasalahan tersebut, penulis tertarik untuk menerapkan metode klasifikasi Naive Bayes pada data mining dalam penelitian ini. Memuat latar belakang atau konteks penelitian, landasan teori jika diperlukan) atau hasil kajian pustaka yang menunjukkan adanya kesenjangan temuan penelitian, wawasan rencana pemecahan masalah dan potensi kontribusinya bagi bidang ilmu dan juga memuat rumusan tujuan penelitian.

Berdasarkan latar belakang tersebut, maka rumusan masalah dalam penelitian ini adalah penerapan metode Naive Bayes untuk klasifikasi kanker paru-paru. Penerapan metode Naive Bayes bertujuan untuk mengklasifikasikan data secara akurat dan dengan akurasi yang tinggi.

Data mining adalah teknologi yang memungkinkan pengguna untuk menggunakan hubungan data sebagai dasar pengambilan keputusan dengan menemukan hubungan dari data yang tidak diketahui pengguna dan menyajikannya dalam format yang mudah dipahami. Data mining dapat dibagi menjadi beberapa kelompok berdasarkan tugas yang dapat mereka lakukan. mendeskripsikan, memperkirakan, memprediksi, mengklasifikasikan, mengelompokkan, dan mengasosiasikan (Larose, 2005). Berisi tentang teori utama yang menjadi fokus utamapenelitian anda.

2.2 Naive Bayes

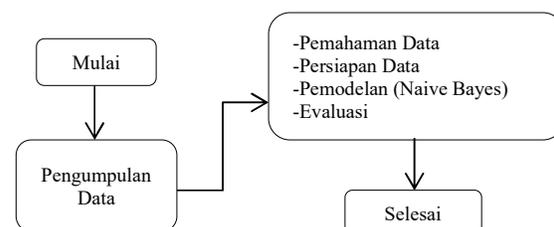
Metode Klasifikasi adalah proses menemukan model atau fungsi yang menggambarkan atau membedakan kelas data, dengan membentuk model yang dapat membagi data ke dalam kelas yang berbeda berdasarkan aturan atau fungsi tertentu, sehingga memungkinkan objek yang tidak ditentukan. (Fadlan dkk., 2018). Dalam penelitian ini akan mengklasifikasikan kanker paru-paru berdasarkan data yang diperoleh dan mencari keakuratan data. Metode yang digunakan dalam pengklasifikasian pada penelitian ini adalah Naive Bayes. Metode Naive Bayes dipilih karena menghasilkan akurasi tertinggi dengan data latih yang lebih sedikit (Handayani, 2015).

Metode Naive Bayes telah banyak digunakan oleh para peneliti sebelumnya dengan berbagai dataset. Dalam sebuah studi oleh Sapto dkk. (2019) menjelaskan algoritma pengujian untuk memprediksi waktu belajar siswa. Pengujian dalam penelitian ini dilakukan di AMIKOM Universitas Yogyakarta menggunakan metode algoritma Naive Bayes pada dataset 300 data dengan akurasi 68%, precision 61,3%, dan recall 65,3%. Sebuah studi oleh Annur (2018) menjelaskan klasifikasi penduduk miskin dengan metode Naive Bayes. Data kecamatan Tibawa menggunakan atribut umur, pendidikan, pekerjaan, pendapatan, kerabat dan status. Hasil pengujian mencapai akurasi 73%, precision 92%, dan recall 86%.

2.METODOLOGI

3.1 Tahapan Penelitian

Tahapan penelitian untuk klasifikasi kanker paru-paru dapat dilihat pada gambar berikut :



Gambar.1 Tahapan Penelitian

Tahapan penelitian dapat diuraikan sebagai berikut:

1. Pengumpulan Data

Dalam penelitian ini penulis menggunakan dataset lung cancer

yang didapat dari web kaggle, yang mana pengumpulan datanya dilakukan dengan menggunakan kuesioner. Data berupa file .csv sebanyak 306 data dan terdapat 16 atribut yaitu: *gender, anxiety, age, smoking, fatigue, yellow fingers, peer pressure, chronic disease, wheezing, alcohol consuming, coughing, shortness of breath, allergy, swallowing difficulty, chest paint, dan lung cancer.*

2. Pemahaman Data

Penelitian ini menggunakan data kaggle dengan 16 atribut. Data kaggle berupa data nominal dan data numerik. Data nominal adalah data yang mempunyai nilai berupa lambang atau nama dan tidak mempunyai urutan nilai. Dan data numerik adalah data kuantitatif yang dapat dihitung dan direpresentasikan sebagai bilangan bulat atau bilangan real. Yang termasuk data nominal yaitu *gender, smoking, anxiety, yellow fingers, peer pressure, chronic disease, wheezing, alcohol consuming, fatigue, coughing, shortness of breath, allergy, swallowing difficulty*, dan *chest paint*. Sedangkan atribut *Age* adalah numerik dan *lung cancer* atribut *class*. Penyelesaian menggunakan metode Naive Bayes dan diuji menggunakan Aplikasi RapidMiner.

3. Persiapan Data

Data yang digunakan dalam format .csv. Data yang akan diproses dimasukkan pada tabel data awal dan kemudian akan dilakukan perubahan pada beberapa jenis data.

4. Pemodelan

Langkah pertama dalam perhitungan Naive Bayes yaitu menghitung probabilitas prior, kemudian menghitung hipotesis class dengan melakukan perhitungan terhadap probabilitas pada kondisi tertentu atau probabilitas prior dengan menjumlahkan kasus secara terperinci tiap atribut data. Selanjutnya yaitu menghitung probabilitas posterior yang berguna dalam penentuan kelas terhadap kasus baru dalam mengelola data. Setelah diketahui nilai tiap atribut probabilitasnya maka dapat dirumuskan dengan $P(C_i|X)$. Langkah selanjutnya yaitu dilakukan perhitungan terhadap jumlah keseluruhan probabilitas tiap kelas.

Teorema bayes dapat ditulis menggunakan persamaan sebagai berikut :

$$P(C_i|X) = \frac{P(X|C_i) \cdot P(C_i)}{P(X)} \dots \dots \dots \text{Persamaan (1)}$$

Keterangan :

$P(C_i|X)$: Probabilitas posterior dari C_i pada kondisi X (posterior probability)

$P(X|C_i)$: Probabilitas posterior dari X pada kondisi C_i

$P(C_i)$: Probabilitas prior C_i (prior probability)

$P(X)$: Probabilitas prior X

X : Data dengan class yang belum diketahui

C_i : Hipotesis data X merupakan class spesifik

5. Evaluasi

Evaluasi performa menggunakan *confusion matrix* untuk menghitung persentase *accuracy, precision* dan *recall*. *Confusion matrix* yaitu tabel yang menampilkan hasil klasifikasi jumlah data uji yang terprediksi benar dan jumlah data uji yang terprediksi salah dapat ditulisdengan rumus sebagai berikut:

$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FN+FP} \dots \dots \dots \text{Persamaan (2)}$$

$$\text{Recall} = \frac{TP}{TP+FN} \dots \dots \dots \text{Persamaan (3)}$$

$$\text{Presisi} = \frac{TP}{TP+FP} \dots \dots \dots \text{Persamaan (4)}$$

Keterangan :

TP = True Positive

TN = True Negative

FN = False Negative

FP = False Positif

3.HASIL DAN PEMBAHASAN

Dari 309 data kanker paru-paru akan dibagi menjadi dua data, yaitu: persentase 70% data train dan persentase 30% data test, dapat dilihat pada tabel 1 dan tabel 2 dibawah ini:

Tabel.1 Data Train

	A 1	A 2	A 3	A 4	A 5	A 6	A 7	A 8	A 9	A 10	A 11	A 12	A 13	A 14	A 15	Pr edic. A
M 1	6	1	2	1	2	2	2	2	2	2	1	2	2	2	2	YES
M 2	7	4	2	1	1	1	2	1	1	1	2	2	2	2	2	YES
F 1	5	9	1	1	2	2	2	1	2	1	1	2	2	1	2	NO
M 1	6	3	2	2	1	1	1	2	1	2	1	2	1	2	1	NO
F 1	6	3	1	2	1	2	1	1	2	1	1	1	1	2	1	NO
F 2	7	5	1	2	1	2	2	1	2	1	2	1	2	1	2	YES
M 1	5	2	2	1	1	2	2	1	2	2	1	2	2	1	2	YES
F 1	5	1	2	2	2	1	2	2	1	1	2	1	2	2	2	YES
F 1	6	8	2	1	1	1	2	2	1	1	1	1	1	1	1	NO
M 2	5	3	2	2	2	1	1	2	1	2	2	2	2	1	2	YES
F 2	6	1	2	2	2	2	2	2	2	1	1	1	1	2	2	YES
M 2	7	2	1	1	1	2	2	1	2	2	2	2	2	2	1	YES
F 1	6	0	2	1	1	1	2	1	1	1	1	1	1	2	1	NO
M 1	5	8	2	1	1	2	2	1	2	2	2	2	2	2	1	YES
M 1	6	9	2	1	1	2	1	1	2	2	2	2	2	1	1	NO
F 2	4	8	1	2	2	2	2	2	2	1	2	1	2	2	2	YES
M 2	7	5	2	1	1	2	1	1	2	2	2	2	2	2	1	YES
M 2	5	7	2	2	2	1	1	2	1	2	1	2	1	2	2	YES
F 2	6	8	2	2	2	2	2	2	1	1	1	1	1	2	1	YES
F 2	6	1	1	1	1	1	2	1	1	1	1	1	1	2	1	NO

Keterangan:

- A1 = Gender
- A2 = Chronic Disease
- A3 = Age
- A4 = Smoking
- A5 = Yellow Fingers
- A6 = Peer Pressure
- A7 = Coughing
- A8 = Fatigue
- A9 = Anxiety
- A10 = Wheezing
- A11 = Alcohol Consuming
- A12 = Allergy
- A13 = Chest Pain
- A14 = Shortness Of Breath
- A15 = Swallowing Difficulty

Predic. A = Lung Cancer

Data train berguna untuk mengklasifikasikan kanker paru-paru pada kelas sentimennya (probabilitas class) dan jumlah datanya lebih banyak. Pada penelitian ini menggunakan data train sebanyak 216 data (70%) dari total data keseluruhan.

Tabel.2 Data Test

A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	Pr edi c.A
1	2	3	4	5	6	7	8	9	0	1	1	1	1	1	1	YE S
F	2	6	7	2	2	2	1	2	2	1	1	1	1	2	2	YE S
F	2	5	6	2	2	2	1	1	2	2	2	2	2	1	2	YE S
F	2	7	0	1	1	2	2	1	1	1	2	2	1	2	1	YE S
M	2	7	0	1	1	1	2	2	1	1	2	2	2	2	1	YE S
F	2	5	7	1	1	2	1	2	2	2	2	2	2	2	2	YE S
M	2	6	1	1	1	2	1	1	2	1	2	1	1	1	1	N O
F	2	7	7	1	1	2	1	2	1	2	2	2	2	1	1	YE S
M	2	6	3	2	2	1	1	2	2	2	1	2	1	2	1	YE S
F	1	6	2	2	1	2	2	2	1	2	2	2	2	1	2	YE S
M	2	5	9	2	1	2	2	2	1	2	2	2	2	2	2	YE S
F	2	7	0	1	2	2	2	2	1	1	2	2	1	1	1	YE S
M	2	7	1	1	2	1	2	1	2	2	2	2	1	1	2	YE S
F	1	5	6	1	2	2	2	2	1	1	1	1	2	2	1	YE S
M	2	5	7	1	1	1	2	1	1	1	2	2	2	2	2	YE S
M	2	7	8	1	2	2	1	2	1	2	1	1	2	2	2	YE S
M	1	6	4	2	2	2	2	2	2	1	2	2	2	2	1	YE S
M	2	6	2	1	2	2	2	2	1	2	2	2	2	1	2	YE S
F	1	4	9	1	2	1	2	1	2	1	1	1	1	1	1	YE S
M	1	7	7	1	2	2	2	2	1	2	2	2	1	1	1	YE S
F	2	6	4	1	1	1	1	1	2	2	2	2	1	1	2	YE S

Tabel.3 Probabilitas

KATEGORI/TRIBUT	SUBSET	YES	NO
GENDER	M	0,540	0,407
	F	0,460	0,444
SMOKING	1	0,418	0,519
	2	0,582	0,481
YELLOW_FINGERS	1	0,402	0,593
	2	0,598	0,407
ANXIETY	1	0,476	0,704
	2	0,524	0,296
PEER_PRESSURE	1	0,476	0,704
	2	0,524	0,296
CHRONIC DISEASE	1	0,476	0,630
	2	0,524	0,370
FATIGUE	1	0,312	0,593
	2	0,688	0,407
ALLERGY	1	0,402	0,926
	2	0,598	0,074
WHEEZING	1	0,418	0,778
	2	0,582	0,222
ALCOHOL CONSUMING	1	0,386	0,815
	2	0,614	0,185
COUGHING	1	0,386	0,704
	2	0,614	0,296
SHORTNESS OF BREATH	1	0,365	0,444
	2	0,635	0,556
SWALLOWING DIFFICULTY	1	0,481	0,852
	2	0,519	0,148
CHEST PAINT	1	0,392	0,704
	2	0,608	0,296

Data test atau data uji berguna untuk mengetahui seberapa baik classifier melakukan pengklasifikasian dengan benar. Jumlah data yang digunakan dalam penelitian ini sebanyak 93 data (30%) dari data keseluruhan.

- I. (2021). Penerapan Algoritma Naive Bayes Dalam Memprediksi Penyakit Lambung. *Journal of Information System, Informatics and Computing*, 5(2), 524–531. <https://doi.org/10.52362/JISICOM.V5I2.659>
- 6) Liantoni, F., & Nugroho, H. (2015). KLASIFIKASI DAUN HERBAL MENGGUNAKAN METODE NAÏVE BAYES CLASSIFIER DAN KNEAREST NEIGHBOR. *Jurnal Simantec*, 5(1). <https://doi.org/10.21107/SIMANTEC.V5I1.1009>
- 7) Novianti, D., Nusa, S., Jakarta, M., & Sitasi, C. (2019). Implementasi Algoritma Naive Bayes Pada Data Set Hepatitis Menggunakan Rapid Miner. 21(1). <https://doi.org/10.31294/p.v20i2>
- 8) Purnama, I., Saputra, R., & Wibowo, A. (n.d.). IMPLEMENTASI DATA MINING MENGGUNAKAN CRISP-DM PADA SISTEM INFORMASI EKSEKUTIF DINAS KELAUTAN DAN PERIKANAN PROVINSI JAWA TENGAH.
- 9) Sari, R. W., Wanto, A., & Windarto, A. P. (2018). IMPLEMENTASI RAPIDMINER DENGAN METODE K-MEANS (STUDY KASUS: IMUNISASI CAMPAK PADA BALITA BERDASARKAN PROVINSI). *KOMIK (Konferensi Nasional Teknologi Informasi dan Komputer)*, 2(1). <https://doi.org/10.30865/KOMIK.V2I1.930>
- 10) Utomo, D. P., & Purba, B. (2019). Penerapan Datamining pada Data Gempa Bumi Terhadap Potensi Tsunami di Indonesia. *Prosiding Seminar Nasional Riset Information Science (SENARIS)*, 1(0), 846–853. <https://doi.org/10.30645/SENARIS.V1I0.91>