

PERBANDINGAN ALGORITMA KLASIFIKASI UNTUK PREDIKSI PENYAKIT LIVER

Muhammad RizkiFahdia¹

Program Studi Sistem Informasi STMIK Nusa Mandiri
www.nusamandiri.ac.id
rizki.muz@nusamandiri.ac.id

Abstract—The liver is the largest organ in the human body, weighing around 3 lbs. The liver is located on the right side of the abdominal cavity below the diaphragm. Where the liver itself acts as a blood filter in the body. Almost all of the blood passes through the liver. liver damage causes complications for other diseases. Several studies were conducted to diagnose patients affected by liver disease. From the process above, it is possible to do research by finding correlations between attributes that influence liver disease. Data mining processing is one of the solutions in research to find linkages between data attributes to gain knowledge and patterns. In addition, it must be explained how far the level of cursation obtained from the model. By comparing the five classification algorithms specified, including Decision Tree (C4.5), Naive Bayes (NB), K-Nearest Neighbor (kNN), Log-R, Deep learning. one of the best algorithms is used for decision making in determining which patients are affected by liver disease. In this study we used a dataset obtained from the UCI repository consisting of 583 patients with 11 attributes. comparisons obtained from the five algorithms using algorithms with Rapidminer software to find out which is the most accurate in predicting liver disease. To determine the accuracy indicator, use the Area Under Curve (AUC) method and a different test (t-Test). The comparison of the five Decision Tree (C4.5) algorithms is the algorithm with the best results, with accuracy (72.56%) and AUC (0.594), to get better performance with Feature Extraction and Feature Selection, there is an increase to 73.24% (accuracy) and 0.602 (AUC).

Keywords - Liver, AlgoritmaKlasifikasi

Abstrak - Hati merupakan organ terbesar dalam tubuh manusia, beratnya sekitar 3 lbs. Hati letaknya berada di sisi kanan rongga perut di bawah diafragma. Dimana hati itu sendiri bertindak sebagai filter darah dalam tubuh. Hampir semua darah melewati liver. kerusakan liver menimbulkan komplikasi terhadap penyakit yang lain. Beberapa penelitian dilakukan untuk mendiagnosis pasien yang terkena penyakit liver. Dari permasalahan diatas hal yang memungkinkan adalah melakukan penelitian dengan cara menemukan korelasi antar atribut yang berpengaruh terhadap timbulnya penyakit liver. Pengolahan data mining menjadi salah satu solusi dalam penelitian untuk menemukan keterkaitan antar atribut data untuk mendapatkan pengetahuan dan pola. Selain itu harus diketahui sejauh mana tingkat kaurasi yang didapat dari model. Dengan cara membandingkan lima algoritma klasifikasi yang ditentukan, diantaranya Decision Tree (C4.5), Naive Bayes(NB), K-Nearest Neighbor(kNN), Log-R, Deep learning. salah satu algoritma terbaik dijadikan untuk pengambilan keputusan dalam menentukan pasien yang terkena atau tidaknya terhadap penyakit liver. Dalam penelitian ini kami menggunakan dataset yang didapat dari repositori UCI yang terdiri dari 583 pasien catatan dengan 11 atribut. perbandingan yang didapat dari ke lima algoritma menggunakan algoritma dengan perangkat lunak Rapidminer untuk mengetahui mana yang paling akurat dalam memprediksi penyakit liver. Untuk menentukan indikator akurasi menggunakan metode Area Under Curve(AUC) dan uji beda (t-Test). Dari perbandingan kelima algoritma Decision Tree (C4.5) merupakan algoritma dengan hasil paling baik, dengan tingkat accuracy (72.56%) dan AUC (0.594), untuk mendapatkan performa yang lebih baik dilakukan Feature Extraction dan Feature Selection, ada kenaikan menjadi 73.24% (accuracy) dan 0.602(AUC).

Kata kunci - Liver, AlgoritmaKlasifikasi

PENDAHULUAN

Hati merupakan organ terbesar dalam tubuh manusia, beratnya sekitar 3 lbs. Hati letaknya berada di sisi kanan rongga perut di bawah diafragma (Kumar &

Katyal, 2018). Setiap gangguan hati dapat menyebabkan peradangan akut atau peradangan kronis, kelemahan hati, dan bahkan dapat merusak organ lain di dalam tubuh. Pada



berbagai gangguan hati perlu membutuhkan perawatan khusus oleh praktisi medis atau profesional dalam perawatan kesehatan (Kumar & Katyal, 2018)

Peneliti sebelumnya sudah banyak yang melakukan penelitian tentang penyakit liver dengan menggunakan berbagai algoritma, diantaranya adalah penelitian tentang memprediksi kematanjangkapendek yang disebabkan salah satunya oleh penakit liver pada pasien yang kerumahsakit untuk dirawat dengan sirosis menggunakan algoritma regresi logistik dan jaringan saraf memori jangka pendek (Harrison et al., 2018).

Peneliti yang memprediksi virus hepatitis C terkait penyakit liver dengan menggunakan jaringan saraf polynomial (Lara et al., 2015). Prediksikomorbidity terkait penyakit liver menggunakan Jaringan Bayesian yang Resampling dan Dinamis dengan Variabel Laten (Yousefi et al., 2017). Peneliti lain melakukan penelitian tentang diagnosis non-invasif penyakit hati berlemak non-alkohol (NAFLD) menggunakan echogenisitas gambar ultrasound (Benjamin et al., 2017). Peneliti berikutnya melakukan penelitian tentang optimisasi aturan klasifikasi C5.0 Boosted Menggunakan Algoritma Genetika untuk Prediksi penyakit hati (Hassoon et al., 2017).

Penelitian tentang perbandingan pendekatan pembelajaran mesin untuk prediksi fibrosis hati lanjut pada Pasien Hepatitis C Kronik (Hashem et al., 2018). Dari keseluruhan penelitian tersebut semuanya menggunakan dataset yang beragam dan tentu saja hasil dari penelitian menyelesaikan permasalahan dari dataset masing-masing.

Pada penelitian ini menggunakan dataset liver yang didapat dari repositori UCI (Universal Child Immunization) yang terdiri dari 583 pasien catat dengan 11 atribut.

Target penelitian ini juga mengacu pada atribut yang digunakan untuk menentukan apakah pasien terkenapa penyakit liver atau tidak. Ketika dokter memutuskan untuk memeriksa pasien maka hal yang pertama dilakukan adalah mengajukan tes bilirubin bersam dengan lainnya (*alkaline phosphatase, aspartate aminotransferase, alanine aminotransferase*). Biasanya liver di cek dengan melakukan beberapa tes, yang meliputi *tes bilirubin, kemudian tes alanine transaminase (ALT), tes aspartate transaminase (AST), dan juga tes alkaline phosphatase (ALP), albumin total protein dan lain-lain* (Lika Aprilia Samadi, 2018). Informasi yang didapat dari berbagai literatur menunjukkan adanya hubungan yang terkait dengan sejumlah atribut dataset yang digunakan untuk memprediksi adanya penyakit liver pada pasien.

Salah satu cara untuk menyelesaikan masalah di atas adalah dengan mengolah data sehingga informasi yang terkandung di dalamnya dapat diambil untuk mendapatkan pola/pengetahuan yang diperlukan.

Dalam hal ini peran Data Mining digunakan untuk mengatasi masalah tentang penyakit liver pada pasien, analisa yang akurat akan memberikan solusi terbaik bagi peneliti dalam mengambil keputusan, inisiatif yang baik untuk memilih suatu metode yang dapat digunakan untuk membantu dalam memprediksi penyakit liver dengan menggunakan teknik CRISP-DM. Ada beberapa hal yang dilakukan dalam penelitian ini yaitu melakukan analisis faktor guna mendapatkan korelasi antar atribut hal ini di lakukan sebagai acuan dasar dalam menemukan pola yang tepat, Melakukan perbandingan lima algoritma klasifikasi untuk menghasilkan algoritma terbaik sebagai tindak lanjut penelitian yang lebih relevan, juga melakukan perbandingan dengan *feature extraction* dan *feature selection*. Data yang digunakan untuk penelitian ini bersumber dari alamat web: <http://archive.ics.uci.edu/ml/>. Data tersebut merupakan hasil pemeriksaan dari 583 orang yang tersebar di wilayah Andhra Pradesh, India jumlah atribut sebanyak 11, Model dibuat dengan menggunakan Rapidminer 8.

METODE PENELITIAN

Metode Correlation Matrix digunakan untuk mengetahui hubungan antar faktor dalam mendeteksi pasien yang terkenapa penyakit liver. Kemudian penelitian dibandingkan dengan 5 algoritma umum, yaitu Decision Tree (C4.5), Naive Bayes (NB), K-Nearest Neighbor (kNN), Logistic Regression (LogR) dan Deep Learning, yang kemudian akan didapatkan algoritma terbaik untuk menentukan penyakit liver. Selanjutnya digunakan perbandingan 4 metode Feature Extraction (PCA) Feature Selection (FS) antara lain Forward Selection, Information Gain, dan Backward Elimination yang kemudian diperoleh metode Feature Selection yang paling baik dalam meningkatkan performa algoritma yang akan digunakan tersebut. Penelitian dilakukan dengan metode CRISP-DM yang terdiri dari 6 fase.

Eksperimen dilakukan pada laptop berbasis Core i3 2.53 GHz CPU, 8 GB RAM dan sistem operasi Windows 10 Enterprise 32-bit. Aplikasi yang digunakan untuk penelitian yaitu Rapid Miner 8.0. Metode CRISP-DM meliputi

- a. Business Understanding
- b. Data Understanding
- c. Data Preparation
- d. Modeling
- e. Evaluasi

HASIL DAN PEMBAHASAN

A. Business Understanding

- Motivasi:
Industri kesehatan memerlukan analisa yang akurat dan efektif untuk menghasilkan prediksi penyakit pada pasien, agar masalah dan kebutuhan pasien teratasi dengan baik. Ingin mengetahui bagaimana data yang ada diolah dapat menghasilkan pengetahuan atau model yang akurat.
- Objektif:
dalam hal ini akan dilakukan pencarian relasi antara faktor yang saling berpengaruh

B. Data Understanding

Dataset yang digunakan dalam penelitian ini diambidar UCI Machine Learning Repository dengan alamat: [https://archive.ics.uci.edu/ml/datasets/ILPD+\(Indian+Liver+Patient+Dataset\)](https://archive.ics.uci.edu/ml/datasets/ILPD+(Indian+Liver+Patient+Dataset)).

Kumpulan data ini berisi 416 catatan pasien liver dan 167 catatan pasien non liver. Kumpulan data diambil dari utara timur Andhra Pradesh India. selector adalah label kelas yang digunakan untuk membagi ke dalam kelompok pasien liver atau tidak. Kumpulan data ini berisi 441 catatan pasien laki-laki dan 142 catatan pasien perempuan. Setiap pasien yang umurnya 89 tahun tercatat berusia 90 tahun. Dataset dengan atribut sebagai berikut:

- Usia (umur penderita)
- Gender dari pasien
- TB (Total Bilirubin)
- DB (Bilirubin Langsung)
- Alkphos (alkaline phosphotase)
- SGPT alamine aminotransferase
- Sulfataspertat aminotransferase
- TP (total Protein)
- ALB (albumin)
- Rasio Albumin dan Globulin Rasio A/G
- Pasien terkena liver atau tidak (1:ya, 2:no)

C. Data Preparation

Atribut akan disesuaikan dengan korelasi atributnya, bisa berupa numerik maupun nominal dan menggunakan label. dalam hal ini keterangannya adalah 1=(terkena liver) atau 2=(tidak terkena).

Kemudian setelah dievaluasi kualitas data ternyata datasetnya missing (kosong) pada salah satu atribut yaitu Rasio Albumin dan Globulin Rasio A/G sudah konsisten sehingga dapat dimodelkan. Selain itu juga pada beberapa atribut menunjukkan nilai yang abnormal (tidak wajar)

Name	Type	Missing	Statistics
direct_bilirubin	Real	0	Min: 0.100, Max: 19.700, Average: 1.486
alkaline_phosphatase	Integer	0	Min: 63, Max: 2110, Average: 290.576
sgpt_ait	Integer	0	Min: 10, Max: 2000, Average: 80.714
sgot_ait	Integer	0	Min: 10, Max: 4929, Average: 109.911
protein_total	Real	0	Min: 2.700, Max: 9.600, Average: 6.483
albumin	Real	0	Min: 0.900, Max: 5.500, Average: 3.142
albumin_globulin_ratio	Real	4	Min: 0.300, Max: 2.800, Average: 0.947

Gambar 1. Statistik Dataset Pasien Liver masiherdapat missing dan noisy

Kemudian melakukan Preparation data dengan operator *Replace Missing value* (mengganti nilai yang hilang dengan nilai rata-rata, dan Operator *Filter Example* (menghapus nilai yang tidak wajar).

D. Modeling

Merupakan proses memilih teknik data mining dengan menentukan algoritma yang akan digunakan. Tool yang dipakai yaitu Rapid Miner versi 8. Untuk mendapatkan hubungan antar atribut digunakan Korelasi Matrix yang mampu mendeskripsikan bentuk dan kekuatan hubungan antar atribut tersebut.

Untuk metode Klasifikasi, dengan membandingkan 5 Algoritma yaitu: Decision Tree, Naive Bayes, K-NN, Logistic Regression dan deep learning. untuk mengetahui nilai accuracy dan AUC.

Digunakan Uji beda (T-Test) Untuk dilakukan perbandingan kinerja (performa) dari kelima algoritma tersebut. Dengan uji beda (T-test) dapat diketahui Akurasi untuk Klasifikasi dan perbedaan signifikan dari kelima algoritma dari masing-masing metode Klasifikasi tersebut untuk dianalisis lebih lanjut.

Untuk memperbaiki kinerja (performa) dari masing-masing metode dapat digunakan metode Feature Extraction, Feature Selection yang terdiri dari beberapa jenis, sedangkan yang digunakan dalam penelitian ini adalah :

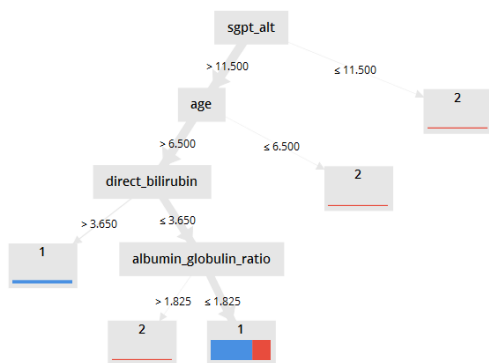
- Feature Exctration menggunakan Principal Component analysis
- Feature Selection

Filter (information gain), Wrapper (Forward Selection dan Backward Elimination).
Feature Selection dapat memperbaiki performanya dan dapat meningkatkan akurasi. Untuk membandingkan metode Feature Selection tersebut juga digunakan uji beda (t-Test). Hasil dari uji beda tersebut digunakan untuk mengetahui metode feature Selection yang terbaik.

E. Evaluasi

Dari hasil pemodelan tersebut dapat diketahui bahwa yang paling berpengaruh dalam menentukan Pasien yang terkena penyakit liver adalah sgpt_alt (tes darah yang dilakukan dokter terhadap alanine aminotransferase untuk memeriksa kerusakan liver)

Jika sgpt_ald dengan kadar kurang dari 11.500 pasien tidak terkena penyakit hati, dan jika kadar sgpt_ald lebih besar dari 11.500 maka selanjutnya memeriksa umur pasien, direct bilirubin dan albumin_globulin_ratio. Pohon keputusan (decision tree) yang digunakan untuk mengetahui hubungan antar atribut :



Gambar 2. pola klasifikasi pasien terkena liver

Pola/pengetahuan yang didapatkan adalah sebagai berikut

1. sgpt_alt <= 11.500 tidak terkena penyakit liver
2. sgpt_ald >= 11.5, umur <= 6.5 tahun tidak terkena liver
3. sgpt_ald >= 11.5, umur >= 6.5 tahun dan direct bilirubin > 3.65 terkena liver
4. sgpt_ald >= 11.5, umur >= 6.5 tahun dan direct bilirubin <= 3.65 dan ratio albumin globulin > 1.825 tidak terkena Liver
5. sgpt_ald >= 11.5, umur >= 6.5 tahun dan direct bilirubin <= 3.65 dan ratio albumin globulin <= 1.825 sebagian besar terkena liver

Grafik dari Hasil Correlation Matrix dapat diketahui pada masing-masing atribut yang memiliki korelasi terhadap pasien yang terkena penyakit liver atau pun tidak.

Tabel yang didapat dari hasil pemodelan dapat juga diketahui hubungan antar faktor di antaranya adalah hubungan positif (berbanding lurus) seperti hubungan antar liver_patient dan albumin. Jika nilai correlation negatif maka hubungannya negatif (berbanding terbalik)

Lalu untuk hubungan negatif (berbanding terbalik) yaitu antara Liver_patient dan direct-bilirubin dimana semakin rendah tingkat direct_bilirubin maka kemungkinan tidak terkena penyakit liver.

Atribut...	liver_pa...	age	gender	tot_bilir...	direct_...	alkaline...	sgpt_ald	spot_ald	protein...	albumin	albumin...
liver_gsl...	1	-0.136	-0.089	-0.231	-0.240	-0.212	-0.202	-0.219	0.031	0.158	0.160
age	-0.136	1	0.058	0.014	0.018	0.031	-0.076	-0.029	-0.189	-0.263	-0.209
gender	-0.089	0.058	1	0.116	0.126	0.045	0.100	0.115	-0.083	-0.098	-0.014
tot_bilru...	-0.231	0.014	0.116	1	0.979	0.227	0.200	0.308	0.004	-0.230	-0.200
direct_bil...	-0.240	0.018	0.126	0.979	1	0.235	0.200	0.294	0.006	-0.224	-0.189
alkaline_...	-0.212	0.031	0.045	0.227	0.235	1	0.160	0.128	0.005	-0.141	-0.224
sgpt_ald	-0.202	-0.076	0.100	0.200	0.200	0.160	1	0.802	-0.008	-0.036	-0.064
spot_ald	-0.219	-0.029	0.115	0.308	0.294	0.128	0.802	1	-0.025	-0.128	-0.130
protein_...	0.031	-0.189	-0.083	0.004	0.006	0.005	-0.008	-0.025	1	0.785	0.233
albumin	0.158	-0.263	-0.098	-0.230	-0.224	-0.141	-0.036	-0.128	0.785	1	0.682
albumin...	0.160	-0.209	-0.014	-0.200	-0.189	-0.224	-0.064	-0.130	0.233	0.682	1

Gambar 3. Tabel hasil Correlation matrix hubungan antar faktor atribut

Selain itu dapat diketahui juga kekuatan hubungannya semakin besar nilai correlationnya maka semakin kuat/banyak hubungannya begitu juga

sebaliknya semakin kecil nilai correlationnya maka semakin lemah/sedikit hubungannya contohnya adalah hubungan antara sgpt_ald dengan sgpt_ald kedua serum ini memiliki korelasi yang kuat sehingga nilai correlationnya besar. Pada

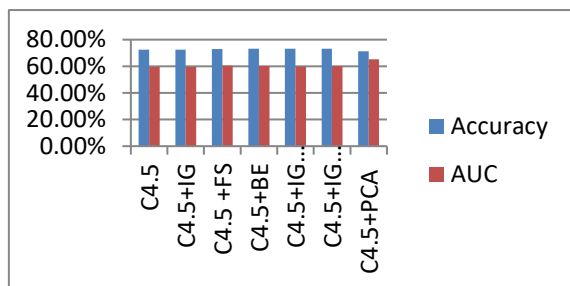
kenyataannya jika dokter memeriksa kadar sgottinggima ka untuk menentukan penyakit liver maka sgpt juga harus diketahui kadarnya. Sedangkan untuk yang hubungannya lemah/sedikit contohnya adalah hubungan antara Pasien yang terkena liver dengan tidak selaludisebabkan oleh salah satu faktor penyebabnya.

dengan tekanan udara karena memiliki nilai correlation yang kecil. Sedangkan yang nilai correlationnya sangat kecil atau bisadikatakan tidak berhubungan karena nilai correlationnya kurang dari 0.4 seperti yang ditunjukkan pada gambar hasil correlation matrix berikut ini :

-1 ↔ -0.8	-0.8 ↔ -0.6	-0.6 ↔ -0.4	-0.4 ↔ -0	0 ↔ 0.4	0.4 ↔ 0.6	0.6 ↔ 0.8	0.8 ↔ 1.0
Very Strong	Strong	Some	No	No	Some	Strong	Very strong
Correlation	Correlation	Correlation	correlation	correlation	correlation	correlation	correlation

Gambar 4. Kekuatan hubungan antar faktor

Dari kelima algoritma yang dipilih menunjukkan hasil dari uji beda (t-Test) diketahui bahwa pada Klasifikasi model yang terbaik adalah decision Tree (C4.5) karena memiliki tingkat Akurasi yang lebih tinggi dan memiliki perbandingan yang signifikan terhadap algoritma yang lain.



Gambar 5. Accuracy dan AUC pada 5 Algoritma

Dari hasil tersebut dapat diketahui bahwa urutan model Algoritma terbaik : 1. C4.5 2. Log R 3. kNN 4. DP 5. NB.

Pohon Keputusan (*Decision Tree*) merupakan metode klasifikasi dan prediksi yang sangat kuat dan terkenal. Metode pohon keputusan mengubah fakta yang sangat besar menjadi pohon keputusan yang merepresentasikan aturan.

Aturan dapat dengan mudah dipahami dengan bahasa alami. Aturan ini juga dapat diekspresikan dalam bentuk bahasa basis data seperti SQL untuk mencari *record* pada kategori tertentu. Pohon keputusan juga berguna untuk mengeksplorasi data, menemukan hubungan tersembunyi antara sejumlah variabel input dengan sebuah variabel target. Karena pohon keputusan memadukan antara eksplorasi data dan pemodelan, pohon keputusan ini sangat bagus sebagai langkah awal dalam proses pemodelan bahkan ketika dijadikan sebagai model akhir dari beberapa teknik lain (J R Quinlan, 1993) Theorema untuk menghitung nilai pada Decision Tree adalah sebagai berikut :

Menghitung nilai Entropy :

$$Entropy(S) = \sum_{j=1}^k -p_j \log_2 p_j$$

S : himpunan kasus
 k : jumlah partisi S
 p_j : probabilitas yang didapat dari jumlah (ya/tidak) dibagi

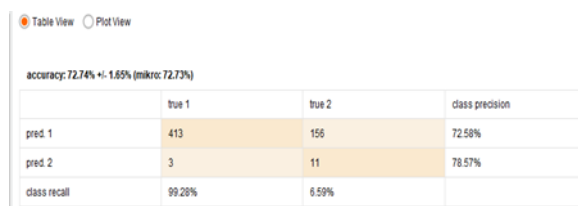
total kasus.

Untuk menghitung gain digunakan rumus berikut :

$$Gain(S, A) = Entropy(S) - \sum_{i=1}^n \left| \frac{s_i}{S} \right| * entropy(s_i)$$

S : himpunan kasus
 A : atribut
 n : jumlah partisi atribut A
 |s_i| : jumlah kasus pada partisi ke i
 |S| : jumlah kasus dalam S

Untuk mengetahui akurasi dari Decision Tree ini dapat dilihat dengan Confusion Matrix berikut ini :



Gambar 6. Confusion Matrix Decision tree



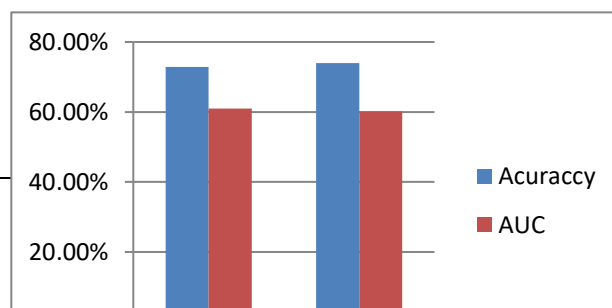
Gambar 7. Kurva ROC-AUC untuk Decision tree

Kategori Klasifikasi AUC:

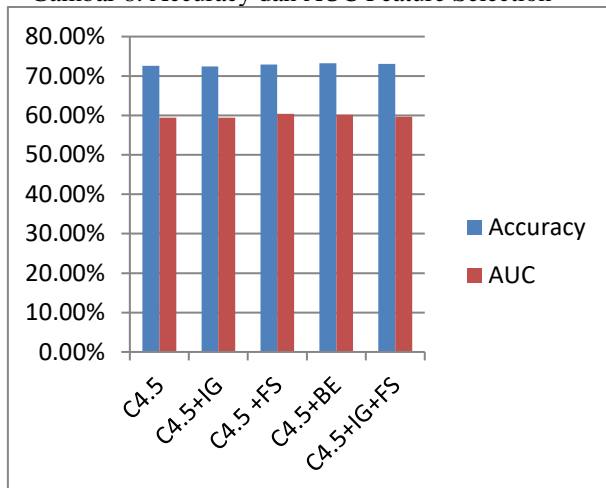
1. 0.90 - 1.00 = excellent classification
2. 0.80 - 0.90 = good classification
3. 0.70 - 0.80 = fair classification
4. 0.60 - 0.70 = poor classification
5. 0.50 - 0.60 = failure

Dari kurva ROC-AUC model Decision tree memiliki AUC sebesar 0.596 ini berarti termasuk dalam kategori klasifikasi rendah.

Hasil metode Feature Extraction, Feature selection pada klasifikasi yang terbaik adalah Backward Elimination (BE) sehingga model terbaik yang digunakan adalah C4.5+BE. Dengan Feature Selection ini performanya lebih baik daripada sebelumnya karena akurasi meningkat dari 72.74% menjadi 73.24% dan AUC meningkat dari 0.596 menjadi 0.602.



Gambar 8. Accuracy dan AUC Feature Selection



Gambar 9. Grafik Accuracy dan AUC pada C4.5 sebelum dan sesudah menggunakan BE

A	B	C	D
	0.720 +/- 0.026	0.718 +/- 0.030	0.724 +/- 0.020
0.720 +/- 0.026		0.689	0.743
0.718 +/- 0.030			0.646
0.724 +/- 0.020			

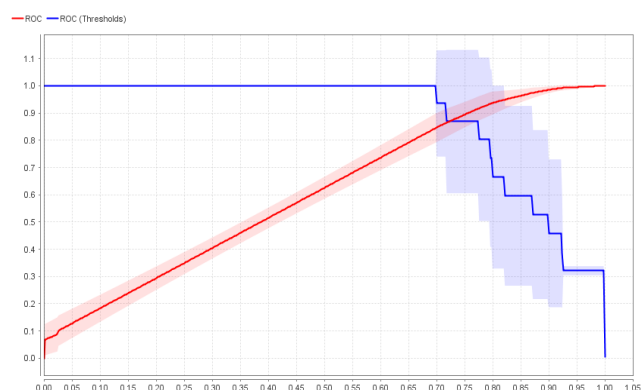
Gambar 10. Hasil Uji beda (T-Test) Feature Selection untuk C4.5

Dari hasil tersebut dapat diketahui bahwa urutan model terbaiknya adalah: 1. C4.5+BE 2. C4.5+FS 3. C4.5+FILTER

accuracy: 72.37% +/- 1.98% (mikro: 72.37%)

	true 1	true 2	class precision
pred 1	393	153	71.98%
pred 2	2	13	86.67%
class recall	99.49%	7.83%	

Gambar 11. Confusion Matrix C4.5+BE



Gambar 12. Kurva ROC-AUC untuk C4.5+BE

Dari kurva ROC-AUC model NB+BE memiliki AUC sebesar 0.606 ini berarti termasuk dalam kategori Poor Classification.

F. Deployment

Pola sebagai pengetahuan yang dihasilkan untuk dijadikan informasi barudalam proses data mining. Pola pengetahuan tersebut didapat dari metode Korelasi, Klasifikasi, dan Feature Selection untuk menentukan Penyakit Liver berdasarkan dataset ILPD (India Liver Patient Dataset). Atribut yang ada pada dataset tersebut korelasinya sangat kecil. Dalam menentukan klasifikasi penyakit liver yang paling menentukan adalah sgt_alt (serum yang terkandung dalam darah). Keakuratan Klasifikasi dapat ditingkatkan dengan menggunakan Backward Elimination sehingga dapat menghasilkan keputusan klasifikasi yang lebih akurat. Pengetahuan yang diperoleh dapat digunakan sebagai dasar dalam mengambil keputusan untuk menentukan pasien yang terkena penyakit liver.

KESIMPULAN

Penelitian dapat menggunakan beberapa peran data mining yaitu : Korelasi, Klasifikasi, dan Feature Extraction, Feature Selection (Filter dan Wrapper)

1. Hubungan Antar Faktor Pasien yang terkena Penyakit liver. Faktor yang paling mempengaruhi pasien yang terkena penyakit liver adalah $sgt_alt \geq 11.500$. Jika $sgt_alt \leq 11.500$ tidak terkena penyakit liver. Apabila $sgt_alt \geq 11.500$, $umur \leq 6.5$ tahun tidak terkena liver. $sgt_alt \geq 11.5$, $umur \geq 6.5$ tahun dan kadar direct bilirubin > 3.65 terkena liver. $sgt_alt \geq 11.5$, $umur \geq 6.5$ tahun dan direct bilirubin ≤ 3.65 dan ratio albumin globulin > 1.825 tidak terkena liver. $sgt_alt \geq 11.5$, $umur \geq 6.5$ tahun dan direct bilirubin ≤ 3.65 dan ratio albumin globulin ≤ 1.825 sebagian besar terkena liver
2. Perbandingan 5 Algoritma Klasifikasi Pasien yang terkena penyakit liver. Algoritma terbaiknya adalah C4.5 dengan Akurasi 72.56 %, AUC 0.594, dan tidak memiliki perbedaan signifikan terhadap algoritma yang lain. Sehingga algoritma C4.5 dapat digunakan untuk klasifikasi dalam menentukan pasien yang terkena penyakit liver.
3. Peningkatan performadengan Feature Selection Feature Selection dapat mengurangi faktor/atribut

yang tidak teraluberpengaruh sehingga dapat meningkatkan Akurasi dan AUC. Sedangkan metode Feature Selection yang terbaik untuk C4.5 pada penelitian ini adalah Backward Elimination (BE) modelnya menjadi C4.5+BE. Metode dengan model C4.5+BE memiliki tingkat akurasi yang baik sehingga hasilnya cukup akurat dengan demikian metode ini dapat digunakan sebagai rekomendasi dalam membantu mengambil keputusan yang tepat untuk menentukan pasien yang terkena liver. Untuk penelitian mendatang dapat menggunakan algoritma yang lain agar akurasinya lebih baik lagi atau dataset yang ada atributnya bisa ditambahkan agar kinerja dari algoritma yang ada menunjukkan performa yang lebih baik lagi. Selanjutnya untuk lebih meningkatkan performansi dari algoritma tersebut dapat dilakukan dengan menggunakan metode Feature Selection atau Feature Extraction lain yang dapat meningkatkan performa algoritma menjadi lebih cepat dan lebih akurat.

REFERENSI

- Benjamin, A., Zubajlo, R., Thomenius, K., Dhyani, M., Kaliannan, K., Samir, A. E., & Anthony, B. W. (2017). Non-invasive diagnosis of non-alcoholic fatty liver disease (NAFLD) using ultrasound image echogenicity. *Proceedings of the Annual International Conference of the IEEE Engineering in Medicine and Biology Society, EMBS*, 2920–2923. <https://doi.org/10.1109/EMBC.2017.8037468>
- Harrison, E., Chang, M., Hao, Y., & Flower, A. (2018). Using machine learning to predict near-term mortality in cirrhosis patients hospitalized at the University of Virginia health system. *2018 Systems and Information Engineering Design Symposium, SIEDS 2018*, 112–117. <https://doi.org/10.1109/SIEDS.2018.8374719>
- Hashem, S., Esmat, G., Elakel, W., Habashy, S., Raouf, S. A., ElHefnawi, M., Eladawy, M., & ElHefnawi, M. (2018). Comparison of Machine Learning Approaches for Prediction of Advanced Liver Fibrosis in Chronic Hepatitis C Patients. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 15(3), 861–868. <https://doi.org/10.1109/TCBB.2017.2690848>
- Hassoon, M., Kouhi, M. S., Zomorodi-Moghadam, M., & Abdar, M. (2017). Rule Optimization of Boosted C5.0 Classification Using Genetic Algorithm for Liver disease Prediction. *2017 International Conference on Computer and Applications, ICCA 2017*, 299–305. <https://doi.org/10.1109/COMAPP.2017.8079783>
- Kumar, S., & Katyal, S. (2018). Effective Analysis and Diagnosis of Liver Disorder by Data Mining. *Proceedings of the International Conference on Inventive Research in Computing Applications, ICIRCA 2018*, 1047–1051. <https://doi.org/10.1109/ICIRCA.2018.8596817>
- Lara, J., Khudyakov, Y., Rossi, L., & Vaughan, G. (2015). Highlights: Predicting the cross-immunoreactivity of hepatitis C virus hyper-variable region 1 peptides using polynomial neural networks. *IEEE*, 1–1. <https://doi.org/10.1109/iccabs.2015.7344732>
- Yousefi, L., Saachi, L., Bellazzi, R., Chiovato, L., & Tucker, A. (2017). Predicting Comorbidities Using Resampling and Dynamic Bayesian Networks with Latent Variables. *Proceedings - IEEE Symposium on Computer-Based Medical Systems, 2017-June*, 205–206. <https://doi.org/10.1109/CBMS.2017.32>