

PERBANDINGAN MODEL KLASIFIKASI RANDOM FOREST DENGAN RESAMPLING DAN TANPA RESAMPLING PADA PASIEN PENDERITA GAGAL JANTUNG

Rizwan Arisandi^{1*}

¹ Departement of Computer Science, Faculty of informatics engineering, Bina Nusantara University, Semarang – Indonesia

*Email: rizwan.arisandi@binus.ac.id

DOI: 10.14710/j.gauss.12.1.136-145

Article Info:

Received: 2022-11-01

Accepted: 2023-05-03

Available Online: 2023-05-04

Keywords:

Random Forest; Resampling; Classification; Heart Failure.

Abstract: Cardiovascular disease that causes heart failure is one of the diseases with the highest mortality rate in the world. Therefore, there is a need for an accurate model to classify heart failure based on clinical information and the lifestyle of patients with the disease, as an alternative solution in administering appropriate drugs. This study compared the classification model of living and deceased heart failure patients based on clinical information and patient lifestyle using the random forest method when using resampling techniques and not using resampling techniques. The results obtained from this study are that the Random Forest model with a combination of the SMOTE and Edited Nearest Neighbors methods is the best model for classifying someone with heart failure as alive or dead. The Random Forest model with a combination of the SMOTE and Edited Nearest Neighbors methods has a high level of classification accuracy in the evaluation model that focuses on recall, namely *rf_model_smoteenn* can classify 82.96% of patients with living status and 90% of patients with death status.

1. PENDAHULUAN

Penyakit kardiovaskular merupakan penyakit pada jantung dan pembuluh darah sebagai penyebab terjadinya gagal jantung. Penyakit ini sering terjadi dan merupakan salah satu penyebab utama kematian di dunia. Berdasarkan data dari WHO, pada tahun 2021 angka kematian akibat penyakit jantung mencapai 17,8 juta orang atau satu dari tiga kematian di dunia pada tiap tahunnya disebabkan oleh penyakit jantung. Gagal jantung adalah kondisi jantung kehilangan kemampuannya untuk memompa darah untuk jumlah yang cukup dalam memenuhi kebutuhan metabolisme tubuh atau jantung hanya mampu melakukannya dengan tekanan pengisian yang tinggi atau dapat juga terjadi kedua-duanya secara bersamaan. Beberapa faktor yang biasanya menyebabkan penyakit kardiovaskular diantaranya adalah diabetes, tekanan darah tinggi pola hidup yang tidak sehat dan kurangnya aktivitas fisik.

Oleh sebab itu perlu adanya suatu model yang akurat untuk mengklasifikasikan gagal jantung berdasarkan informasi klinis dan gaya hidup pasien pengidap penyakit tersebut, sebagai solusi alternatif dalam pemberian obat yang tepat. Akan tetapi, tingkat akurasi klasifikasi kejadian terkait gagal jantung dalam praktik klinis biasanya kurang sensitif. Saat ini, banyak cara untuk melakukan prediksi, salah satunya dengan membuat pemodelan menggunakan machine learning, yaitu pengembangan algoritma dan model secara statistik yang menggunakan sistem komputer serta membutuhkan data yang mengandalkan pola serta inferensi. Oleh karena itu, ketika ingin mengklasifikasikan seseorang pasien yang mengalami gagal jantung dalam status hidup atau meninggal dunia dapat dilakukan dengan menggunakan machine learning, dalam hal ini Random Forest menggunakan data rekam medis, klinis dan gaya hidup pasien gagal jantung

Tujuan dari penelitian ini untuk membandingkan model klasifikasi pasien gagal jantung yang berstatus hidup atau meninggal dunia berdasarkan informasi klinis dan gaya hidup pasien dengan metode random forest jika menggunakan teknik resampling dan tidak menggunakan teknik resampling. Manfaat penelitian ini adalah memberikan sebuah model yang dapat digunakan sebagai solusi alternatif pemberian obat.

2. TINJAUAN PUSTAKA

Random Forest

Supervised Learning merupakan salah satu metode klasifikasi pada machine learning yang mengklasifikasikan data berbentuk kategorik kedalam suatu kelas tertentu (Brownlee, 2016). Dengan begitu, Supervised Learning dapat melakukan klasifikasi terhadap suatu penyakit dengan menggunakan data-data yang mendukungnya.

Pada Supervised Learning, terdapat beberapa model, salah satu modelnya adalah Random Forest. Model random forest tersebut merupakan model gabungan dari decision tree dimana model decision tree dibentuk terlebih dahulu kemudian klasifikasi random forest dapat diaplikasikan dengan memanfaatkan pohon keputusan yang dihasilkan dari leaf node pada pembentukan decision tree (Hidayati & Mu'Alim, 2022). Dengan menggunakan Random Forest Classifier yang merupakan Supervised Learning, dapat dilakukan prediksi terhadap tingkat ketahanan hidup dari seseorang yang mengalami gagal jantung.

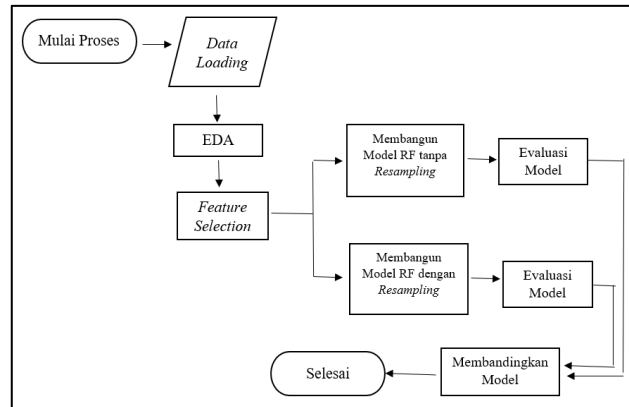
Resampling

Terdapat beberapa cara untuk meningkatkan hasil prediksi, yaitu dengan cara resampling ketika terdapat imbalance pada datanya dan merupakan salah satu cara yang paling sering digunakan untuk mengatasi masalah imbalance. Terdapat beberapa metode dalam resampling, di antaranya Random Oversampling, Random Undersampling, SMOTE Oversampling, NearMiss Undersampling, SMOTE dengan Tomek Links, dan SMOTE dengan Edited Nearest Neighbors. Random Oversampling merupakan proses resampling yang dilakukan dengan cara memilih data pada kelas minoritas secara acak, lalu pada data yang terpilih, dilakukan duplikasi dan ditambahkan pada set pelatihan yang baru. Sementara itu, Random Undersampling merupakan proses resampling dengan memilih data secara acak pada kelas mayoritas, lalu eliminasi secara acak sehingga rasio antara kelas minoritas dan mayoritas sesuai pada tingkat yang diinginkan (Prasetya, J., 2022).

Metode resampling dengan jenis oversampling lainnya dapat menggunakan metode Synthetic Minority Oversampling Technique (SMOTE), di mana teknik ini mensintesis data baru dari kelas minoritas untuk menyeimbangkan dataset dengan cara sampling ulang sampel kelas minoritas (Muqit WS.A. and Nooraeni, R., 2020). Selanjutnya, pada resampling dengan jenis undersampling, dapat menggunakan Near Miss yang bekerja dengan memilih sampel berdasarkan jarak sampel pada kelas mayoritas ke sampel kelas minoritas. Selain NearMiss dan SMOTE yang populer, terdapat metode lainnya, yaitu Tomek Links yang berbeda dengan dengan Nearmiss, di mana Tomek Links memilih sampel untuk dihapus (Magnolia, C. & Nurhopipah, A., 2022). Selanjutnya adalah SMOTE Edited Nearest Neighbors (SMOTE-ENN), yaitu kombinasi antara SMOTE dan ENN dengan SMOTE berperan sebagai oversampling, sedangkan ENN berperan sebagai undersampling.

3. METODE PENELITIAN

Langkah- langkah penelitian yang akan dilakukan ditunjukkan oleh Gambar 1.



Gambar 1. Langkah-Langkah Penelitian

Proses penelitian dimulai dari identifikasi masalah, setelah itu pre-processing data, kemudian melakukan training model dan terakhir adalah mengevaluasi model yang dihasilkan. Pada pre-processing data dilakukan dengan membandingkan hasil menggunakan resampling dan dengan tanpa menggunakan teknik resampling untuk mengetahui perbedaan pada imbalance dataset.

Pada tanpa menggunakan teknik resampling, dataset yang telah melalui pre-processing akan melalui proses learning dengan menggunakan metode random forest. Sedangkan data dengan resampling dilakukan dengan data yang telah disetarakan kelasnya dengan teknik resampling. Setelah itu, data akan melalui proses learning dengan menggunakan metode random forest yang kemudian akan dievaluasi tiap modelnya dan dilakukan perbandingan model sehingga didapatkan model terbaik.

4. HASIL DAN PEMBAHASAN

Data yang digunakan pada penelitian ini adalah data rekam medis 299 pasien gagal jantung pada bulan April-Desember 2022. Pasien terdiri dari 194 laki-laki dan 105 perempuan dengan rentang usia antara 40 dan 95 tahun. memuat informasi klinis, tubuh, dan gaya hidup dengan rincian pada Tabel 1.

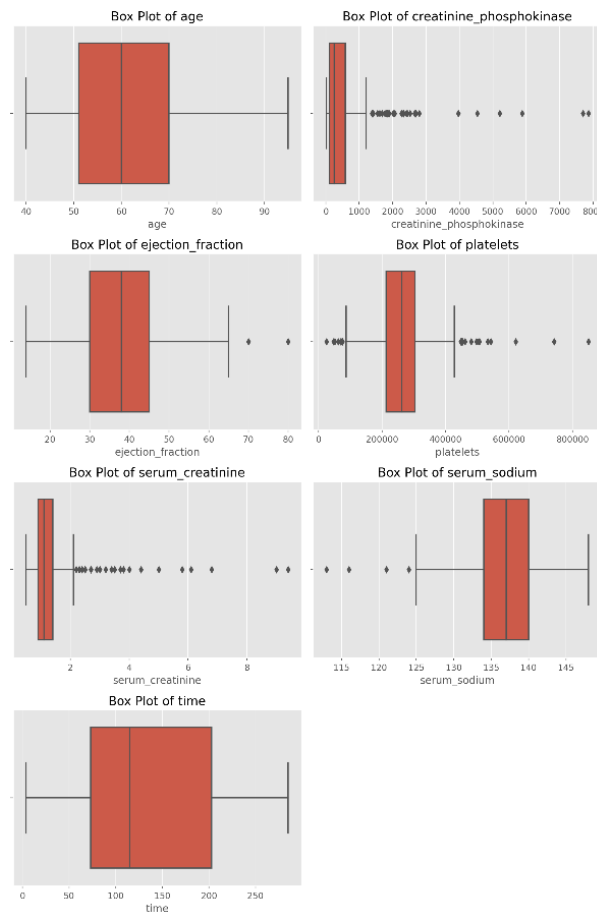
Tabel 1. Penjelasan Dataset

| Variabel | Deskripsi |
|---------------------|--|
| Usia | Umur pasien |
| Anemia | Kekurangan sel darah merah |
| Kreatin fosfokinase | Tingkat enzim kreatinin dalam darah (mcg/L) |
| Diabetes | Apakah pasien diabetes atau tidak |
| Fraksi ejeksi | Persentase darah yang keluar jantung pada setiap kali ventrikel berkontraksi |
| Hipertensi | Apakah pasien hipertensi atau tidak |
| Trombosit | Jumlah kepingan darah yang diproduksi sumsum tulang (kilo platelet/mL) |
| Serum kreatinin | Tingkat hasil metabolisme otot yang mengalir pada sirkulasi darah (mg/dL) |
| Serum natrium | Tingkat natrium serum dalam darah (mEq/L) |
| Jenis kelamin | Wanita atau pria (biner) |
| Merokok | Apakah pasien merokok atau tidak |
| Waktu | Periode <i>follow-up</i> (hari) |
| Death_event | Apakah pasien meninggal atau tidak selama periode <i>follow-up</i> |

Data Pre-Processing

Langkah pertama yang dilakukan adalah memeriksa apakah dataset mengandung nilai null karena hal ini umum terjadi dalam data "real". Namun, dataset ini lengkap dan tidak memiliki nilai null. Selanjutnya memeriksa apakah dataset memiliki outlier. Untuk itu

dilakukan dengan visualisasi Box Plot dan metode IQR (Interquartile Range) untuk mendeteksi outlier pada kolom-kolom numerik dalam dataset. Pada metode IQR, data dianggap sebagai outlier jika berada di bawah lower limit ($Q1-1,5IQR$) atau di atas upper limit ($Q3+1,5IQR$).

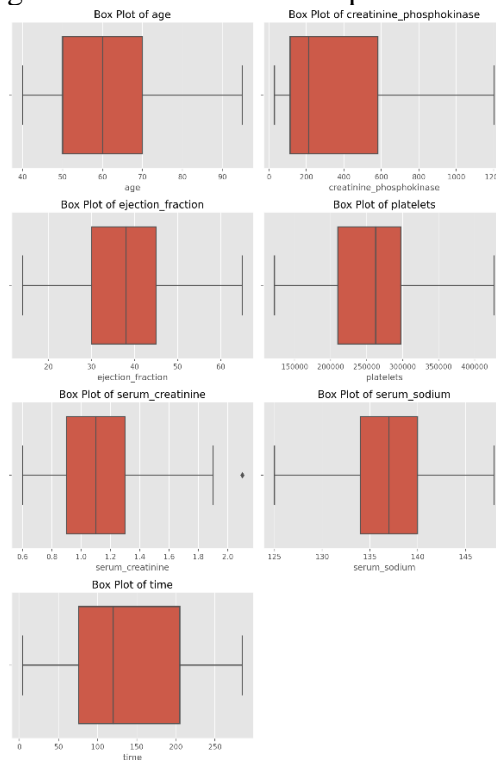


Gambar 2. Box Plot Sebelum Menghapus Outlier
Tabel 2. Jumlah Outlier Sebelum Menghapus Outlier

| Variabel | Jumlah Outlier | Persentase |
|---------------------|-------------------|-----------------------------|
| | | Outlier pada Variabel |
| Kreatin fosfokinase | 29 | 9,698997 |
| Serum kreatinin | 29 | 9,698997 |
| Trombosit | 21 | 7,023411 |
| Serum natrium | 4 | 1,337793 |
| Fraksi ejeksi | 2 | 0,668896 |
| umur | 0 | 0,000000 |
| waktu | 0 | 0,000000 |

Berdasarkan hasil di atas, terlihat bahwa terdapat beberapa outlier, terutama pada kolom kreatin fosfokinase dan serum kreatinin dengan persentase outlier hampir mencapai 10%. Untuk mengatasi hal tersebut semua outlier dibuang dan melakukan slicing pada dataset

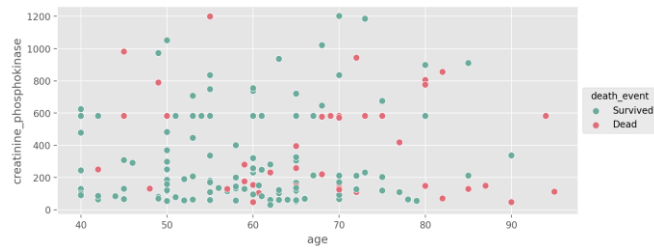
untuk mengambil nilai-nilai yang hanya berada di dalam lower limit dan upper limit. Setelah itu, mengulangi langkah yang sama untuk memeriksa apakah masih ada outlier atau tidak.



Gambar 3. Box Plot Setelah Menghapus Outlier
Tabel 3. Jumlah Outlier Setelah Menghapus Outlier

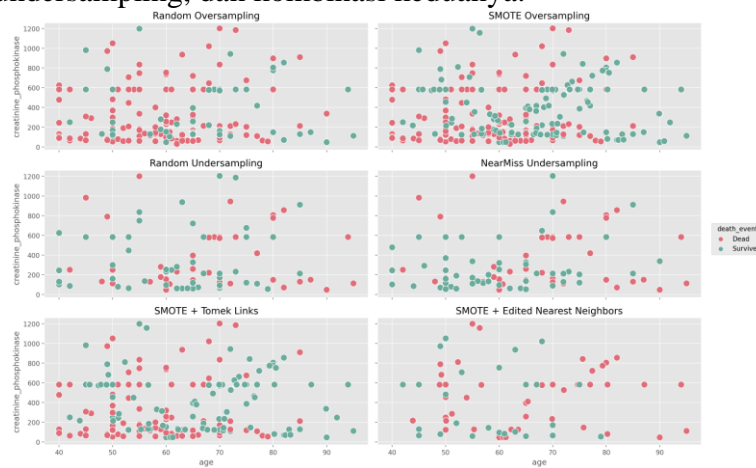
| Variabel | Persentase | |
|---------------------|-------------------|-----------------------------|
| | Jumlah Outlier | Outlier pada Variabel |
| Kreatin fosfokinase | 3 | 1,339286 |
| Serum kreatinin | 0 | 0,000000 |
| Trombosit | 0 | 0,000000 |
| Serum natrium | 0 | 0,000000 |
| Fraksi ejeksi | 0 | 0,000000 |
| umur | 0 | 0,000000 |
| waktu | 0 | 0,000000 |

Berdasarkan hasil di atas, semua kolom tidak lagi memiliki outlier, kecuali kolom serum kreatinin dengan persentase outlier yang sangat sedikit dan tidak signifikan, yaitu sekitar 1,3%. Langkah selanjutnya adalah memeriksa keseimbangan kategori pada kolom target, yaitu death_event. Langkah ini sangat penting karena ketidakseimbangan kategori dapat menimbulkan masalah dalam machine learning, seperti kinerja model yang bias dan akurasi yang menurun dalam memprediksi kategori minoritas. Untuk melihat bagaimana keseimbangan kategori pada kolom death_event, dibuat scatter plot pada kolom Kreatin fosfokinase dan Umur yang dikelompokkan berdasarkan death_event.



Gambar 4. Persebaran Kategori death_event Sebelum Resampling

Berdasarkan plot di atas, terlihat jelas bahwa mayoritas dataset tergolong dalam kategori status hidup, sedangkan hanya sebagian kecil yang tergolong dalam kelas status meninggal dunia. Untuk mengatasi masalah ini, digunakan berbagai metode resampling, seperti oversampling, undersampling, dan kombinasi keduanya.



Gambar 5. Persebaran Kategori Death_Event Setelah Resampling

Setelah menggunakan metode resampling, terlihat bahwa terdapat perubahan yang signifikan dalam proporsi kategori-kategori tersebut. Teknik resampling tersebut dapat mengatasi masalah ketidakseimbangan kategori sehingga menghasilkan distribusi yang lebih seimbang dari kolom target, yaitu death_event.

Model Training

Sebelum melatih model, terlebih dahulu dataset dibagi menjadi dua yaitu dataset training dan dataset testing dalam rasio 80:20. Selanjutnya membuat model Random Forest dari library scikit-learn pada dataset training yang telah distandardisasi. Scikit-learn menyediakan parameter class_weight dengan opsi balanced untuk menyesuaikan bobot setiap kategori selama training. Hal ini dapat meningkatkan performa model pada kategori minoritas.

Untuk meningkatkan performa model pada data yang imbalanced lebih jauh lagi, digunakan berbagai metode resampling seperti Random Oversampling, SMOTE Oversampling, Random Undersampling, NearMiss Undersampling, kombinasi SMOTE + Tomek Links, dan kombinasi SMOTE + ENN. Metode resampling ini disediakan oleh library imbalanced-learn, yaitu library yang dirancang khusus untuk menangani data imbalanced.

Setiap metode resampling tersebut diterapkan pada dataset training sehingga menghasilkan 8 dataset training yang berbeda. Setelah itu melatih model Random Forest yang terpisah pada setiap dataset training tersebut untuk dibandingkan hasil model yang akan dibandingkan dapat dilihat pada Tabel 4.

Tabel 4. Model Random Forest

| Model | Deskripsi |
|-------------------------------|---|
| <i>rf_model_base</i> | Model Random Forest <i>basic</i> (tanpa metode <i>resampling</i>) |
| <i>rf_model_base_balanced</i> | Model Random Forest <i>basic</i> (tanpa metode <i>resampling</i>) dengan parameter tambahan <i>class_weight="balanced"</i> |
| <i>rf_model_ros</i> | Model Random Forest dengan Random Oversampling |
| <i>rf_model_smote</i> | Model Random Forest dengan SMOTE Oversampling |
| <i>rf_model_rus</i> | Model Random Forest dengan Random Undersampling |
| <i>rf_model_nearmiss</i> | Model Random Forest dengan NearMiss Undersampling |
| <i>rf_model_smotetomek</i> | Model Random Forest dengan kombinasi SMOTE dan Tomek Links |
| <i>rf_model_smoteenn</i> | Model Random Forest dengan kombinasi SMOTE dan Edited Nearest Neighbors |

Model Evaluation

Untuk klasifikasi biner, evaluasi selama pelatihan model dapat dilakukan menggunakan tabel confusion matrix pada Tabel 5.

Tabel 5. Confusion Matrix

| Prediksi | Actual Positif Class | Actual Negatif Class |
|----------------|----------------------|----------------------|
| Positive Class | True Positive (TP) | False Negative (FN) |
| Negative Class | False Positive (FP) | True Negative (TN) |

Baris pada tabel merepresentasikan kategori yang diprediksi, sedangkan kolom pada tabel merepresentasikan kategori yang sebenarnya. Berdasarkan tabel confusion matrix, terutama pada penelitian ini, TP dan TN merupakan jumlah pasien positif berstatus hidup dan negatif berstatus meninggal dunia yang diklasifikasikan dengan benar. Sementara itu, FP dan FN merupakan jumlah pasien yang misklasifikasi.

Dari tabel confusion matrix sebagai dasarnya, beberapa metrik yang umum digunakan dapat dihasilkan untuk mengevaluasi kinerja model dengan fokus evaluasi yang berbeda, antara lain terdapat pada Tabel 6.

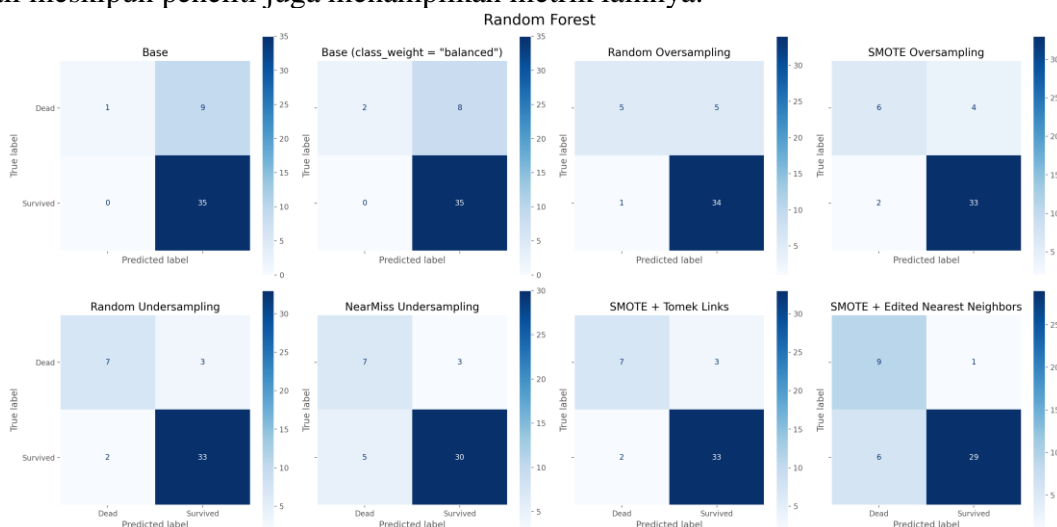
Tabel 6. Metrik Evaluasi Kinerja Model

| Metrik | Deskripsi | Rumus |
|------------------|---|---|
| <i>Accuracy</i> | Rasio prediksi yang benar dengan keseluruhan data | $\frac{TP + TN}{TP + FP + TN + FN}$ |
| <i>Precision</i> | Rasio prediksi yang benar dengan keseluruhan data yang diprediksi positif | $\frac{TP}{TP + FP}$ |
| <i>Recall</i> | Rasio prediksi yang benar dengan keseluruhan data yang memang positif | $\frac{TP}{TP + FN}$ |
| <i>F-1 Score</i> | Perbandingan rata-rata <i>precision</i> dan <i>recall</i> yang dibobotkan | $\frac{2 \times precision \times recall}{precision + recall}$ |

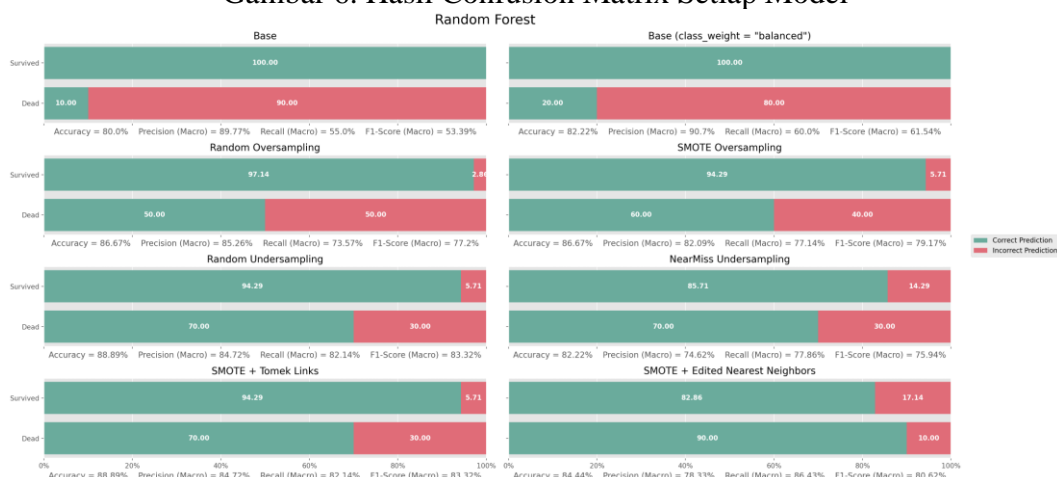
Meskipun *accuracy* umum digunakan pada machine learning klasifikasi, metrik tersebut bukan selalu metrik terbaik, terutama pada dataset yang tidak seimbang (*imbalanced*). Pada

dataset yang tidak seimbang, biasanya kategori-kategori pada kolom target tidak tersebar secara merata dan suatu kategori bisa jauh lebih banyak dari kategori lainnya sehingga accuracy dapat memberikan informasi yang salah. Hal ini disebabkan karena accuracy hanya mengukur keseluruhan proporsi jumlah prediksi yang benar tanpa mempertimbangkan persebaran hasil prediksinya di antara berbagai kategori yang berbeda.

Pada kasus seperti itu, banyak yang menyarankan metrik lain, seperti recall untuk mengevaluasi kemampuan model dalam mengidentifikasi jumlah positif dengan benar, terutama pada kategori minoritas. Recall mempertimbangkan baik TP dan FN sehingga memberikan informasi yang lebih baik mengenai seberapa model mengenali kategori minoritas. Oleh karena itu, pada penelitian ini, peneliti lebih memfokuskan pada metrik recall meskipun peneliti juga menampilkan metrik lainnya.



Gambar 6. Hasil Confusion Matrix Setiap Model



Gambar 7. Hasil Recall Setiap Model

Berdasarkan hasil evaluasi di atas, dapat dilihat bahwa model Random Forest tanpa menggunakan metode resampling, yaitu rf_model_base (Base) dan rf_model_base_balanced (Base (class_weight = "balanced")) memiliki performa yang buruk. Meskipun kedua model tersebut dapat mengklasifikasikan semua pasien yang berstatus hidup dengan benar, namun tidak halnya dengan pasien yang berstatus meninggal dunia. Dari seluruh pasien yang berstatus meninggal dunia., hanya 10-20% yang diklasifikasikan sebagai pasien dengan status meninggal dunia. sedangkan 80-90% sisanya diklasifikasikan sebagai pasien yang berstatus hidup.

Sementara itu, model Random Forest yang menggunakan kombinasi metode SMOTE dan Edited Nearest Neighbors, yaitu *rf_model_smoteenn* dapat mengklasifikasikan 82,96% pasien yang berstatus hidup dan 90% pasien yang berstatus meninggal dunia dengan benar sehingga dapat dikatakan bahwa performanya sangat baik dan seimbang. Oleh karena itu, model *rf_model_smoteenn* merupakan model terbaik. Berikut 10 hasil prediksi pertama model *rf_model_smoteenn* terhadap dataset testing yang tertera pada Tabel 7.

Tabel 7. Hasil Prediksi

| No. | Probabilitas Pasien dengan status hidup (%) | Probabilitas Pasien dengan status meninggal dunia (%) | Prediksi | Aktual | Kesimpulan |
|-----|---|---|------------------|------------------|--------------|
| 1 | 2,333505 | 97,666495 | <i>Meninggal</i> | <i>Hidup</i> | <i>False</i> |
| 2 | 99,982383 | 0,017617 | <i>Hidup</i> | <i>Hidup</i> | <i>True</i> |
| 3 | 84,785228 | 15,214772 | <i>Hidup</i> | <i>Hidup</i> | <i>True</i> |
| 4 | 2,806210 | 97,193790 | <i>Meninggal</i> | <i>Meninggal</i> | <i>True</i> |
| 5 | 3,903181 | 96,096819 | <i>Meninggal</i> | <i>Meninggal</i> | <i>True</i> |
| 6 | 99,990174 | 0,009826 | <i>Hidup</i> | <i>Hidup</i> | <i>True</i> |
| 7 | 99,897687 | 0,102313 | <i>Hidup</i> | <i>Hidup</i> | <i>True</i> |
| 8 | 54,050504 | 45,949496 | <i>Hidup</i> | <i>Hidup</i> | <i>True</i> |
| 9 | 0,006271 | 99,993729 | <i>Meninggal</i> | <i>Meninggal</i> | <i>True</i> |
| 10 | 0,657082 | 99,342918 | <i>Meninggal</i> | <i>Hidup</i> | <i>False</i> |

dan seterusnya

5. KESIMPULAN

Berdasarkan data penelitian yang berisi rekam medis 299 pasien gagal jantung pada bulan April-Desember 2022 diperoleh model terbaik yaitu model Random Forest yang menggunakan kombinasi metode SMOTE dan Edited Nearest Neighbors, yaitu *rf_model_smoteenn* dapat mengklasifikasikan 82,96% pasien yang berstatus hidup dan 90% pasien yang berstatus meninggal dengan tepat.

Pemilihan model terbaik didasarkan mulai dari membuang semua outlier dengan melakukan slicing untuk mengambil nilai-nilai yang berada di dalam lower limit dan upper limit. Kemudian dilakukan teknik resampling untuk mengatasi ketidak seimbangan kategori pada kolom Kreatin fosfokinase dan umur yang dikelompokkan berdasarkan *death_event*.

Sebelum menguji model, dataset dibagi menjadi dataset training dan testing dalam rasio 80:20. Selanjutnya, model diuji menggunakan delapan model Random Forest pada dataset training yang telah distandarisasi. Setelah itu dilakukan evaluasi model yang fokus pada metrik recall untuk mengevaluasi kemampuan model dalam mengidentifikasi jumlah positif dengan benar, terutama pada kategori minoritas.

DAFTAR PUSTAKA

Pane, J.P., Simorangkir, L. and Saragih, P.I.S.B. (no date) Faktor-Faktor Risiko Penyakit kardiovaskular Berbasis Masyarakat, *Jurnal Penelitian Perawat Profesional*. Available at: <https://jurnal.globalhealthsciencegroup.com/index.php/JPPP/article/view/1218> (Accessed: April 2, 2023).

- Chicco, D. and Jurman, G. (2020) Machine learning can predict survival of patients with heart failure from serum creatinine and ejection fraction alone - BMC Medical Informatics and decision making, BioMed Central. BioMed Central. Available at: <https://bmcmedinformdecismak.biomedcentral.com/articles/10.1186/s12911-020-1023-5> (Accessed: April 2, 2023).
- Thupae, R. et al. (2018) "Machine learning techniques for traffic identification and classification in SDWSN: A survey," IECON 2018 - 44th Annual Conference of the IEEE Industrial Electronics Society [Preprint]. Available at: <https://doi.org/10.1109/iecon.2018.8591178>.
- Wuryani, N. and Agustiani, S. (2021) "Random Forest classifier untuk Deteksi PENDERITA COVID-19 Berbasis citra CT scan," Jurnal Teknik Komputer, 7(2), pp. 187–193. Available at: <https://doi.org/10.31294/jtk.v7i2.10468>.
- Roihan, A., Sunarya, P.A. and Rafika, A.S. (2020) "Pemanfaatan machine learning Dalam Berbagai Bidang: Review paper," IJCIT (Indonesian Journal on Computer and Information Technology), 5(1). Available at: <https://doi.org/10.31294/ijcit.v5i1.7951>.
- Sutoyo, E. and Fadlurrahman, M.A. (2020) "Penerapan smote untuk mengatasi imbalance class Dalam Klasifikasi television advertisement performance rating menggunakan artificial neural network," Jurnal Edukasi dan Penelitian Informatika (JEPIN), 6(3), p. 379. Available at: <https://doi.org/10.26418/jp.v6i3.42896>.
- Muqijit WS, A. and Nooraeni, R. (2020) "Penerapan metode Resampling dalam mengatasi imbalanced data Pada determinan Kasus Diare Pada Balita di Indonesia (Analisis Data SDKI 2017)," Jurnal MSA (Matematika dan Statistika serta Aplikasinya), 8(1), p. 19. Available at: <https://doi.org/10.24252/msa.v8i1.13452>.
- Prasetya, J. (2022) PENERAPAN KLASIFIKASI NAIVE BAYESDENGAN ALGORITMA RANDOM OVERSAMPLINGDAN RANDOM UNDERSAMPLINGPADA DATA TIDAK SEIMBANG CERVICAL CANCER RISK FACTORS, View of Penerapan Klasifikasi naive Bayes Dengan Algoritma random oversampling Dan random undersampling Pada Data Tidak Seimbang cervical cancer risk factors. Available at: <https://ejurnal.unisap.ac.id/index.php/leibniz/article/view/173/101> (Accessed: April 4, 2023).
- Sir, Y.A. and Soepranoto, A.H. (2022) "Pendekatan Resampling Data Untuk Menangani Masalah Ketidakseimbangan kelas," Jurnal Komputer dan Informatika, 10(1), pp. 31–38. Available at: <https://doi.org/10.35508/jicon.v10i1.6554>.
- Hossin, Mohammad & M.N, Sulaiman. (2015). A Review on Evaluation Metrics for Data Classification Evaluations. International Journal of Data Mining & Knowledge Management Process. 5. 01-11. 10.5121/ijdkp.2015.5201.