

Algoritma Random Forest Untuk Prediksi Kelangsungan Hidup Pasien Gagal Jantung Menggunakan Seleksi Fitur Bestfirst

Yuri Yuliani^{1*}

¹Program Studi Sistem Informasi, Universitas Bina Sarana Informatika

*yuri.yyi@bsi.ac.id

Abstrak

Gagal jantung yang merupakan masalah kesehatan yang global yang tidak hanya menimbulkan masalah fisik, dampak lain seperti psikologis, social dan ekonomi, serta mengalami depresi, yang mempengaruhi dalam pengobatan, memperburuk status fungsional dan meningkatkan tingkat rawat inap hingga kematian. Menurut World Health Organization (WHO) hampir 17,5 juta orang meninggal yang diakibatkan oleh penyakit kardiovaskuler yang mewakili dari 31% kematian yang ada di dunia. Menggunakan machine learning untuk memprediksi kelangsungan hidup pasien penderita gagal jantung agar dapat melakukan pencegahan dari awal. Tahapan penelitian yang dilakukan meliputi tahapan pemahaman bisnis, tahapan pemahaman data, tahapan persiapan data, tahap pemodelan dan tahap evaluasi. Pada penelitian kali ini menggunakan seleksi fitur menggunakan bestfirst menghasilkan 4 fitur yang sangat berpengaruh yaitu age, enjection_fraction, serum_creatinene dan time, serta penanganan imbalance class menggunakan model class balancer. Algoritma random forest dengan metode percentage split 80% yang menghasilkan accuracy 91,45%, mean absolute error 0.1874, incorrectly classified instances 8.55%, precision 0.915, recall 0.914, AUC 0.953.

Kata kunci: Gagal Jantung, Random Forest, Bestfirst, Class Balancer

Abstract

Heart failure is a global health problem that not only causes physical problems, other impacts such as psychological, social, and economic, as well as depression, which affects treatment, worsens functional status, and increases hospitalization rates to death. According to the World Health Organization (WHO), nearly 17.5 million people die from cardiovascular disease, which represents 31% of deaths in the world. Using machine learning to predict the survival of patients with heart failure so that they can take precautions from the start. The stages of the research carried out include the business understanding stage, the data understanding stage, the data preparation stage, the modeling stage, and the evaluation stage. In this study, using feature selection using best-first resulted in 4 very influential features, namely age, injection_fraction, serum_creatinene and time, and handling imbalance class using the class balancer model. Random forest algorithm with 80% percentage split method which produces 91.45% accuracy, mean absolute error 0.1874, incorrectly classified instances 8.55%, precision 0.915, recall 0.914, AUC 0.953.

Keywords: Heart failure, Random Forest, Bestfirst, Class Balancer

1. Pendahuluan

Gagal jantung merupakan masalah global dalam kesehatan yang mempengaruhi jutaan orang [1].

Terjadi karena sindrom klinis kompleks yang diakibatkan oleh kerusakan structural atau fungsional akibat pengisian vantrikel atau

pemompa darah [2]. Serta gagal jantung itu sendiri yang dimana jantung gagal untuk memompa yang cukup dalam memenuhi kebutuhan tubuh [3]. Tidak hanya menimbulkan masalah fisik, dampak lain seperti psikologis, social dan ekonomi, serta mengalami depresi bagi pasien yang sedang menjalani rawat inap maupun rawat jalan, yang pastinya akan mempengaruhi dalam pengobatan, memperburuk status fungsional dan meningkatkan tingkat rawat inap hingga kematian [4]. Menurut World Health Organization (WHO) hampir 17,5 juta orang meninggal yang diakibatkan oleh penyakit kardiovaskuler yang mewakili dari 31% kematian yang ada di dunia. Di Amerika Serikat penyakit gagal jantung terjadi 550.000 kasus/tahun. Menurut Riset Kesehatan Dasar (2013)prevalensi penyakit gagal jantung di Indonesia tahun 2013 sebesar 0,13% atau diperkirakan sekitar 229.696 orang, sedangkan berdasarkan diagnosis dokter/ gejala sebesar 0,3% atau diperkirakan sekitar 530.068 orang. Penyakit gagal jantung Provinsi Jawa Tengah sebanyak 43.361 orang (0,18%) [5].

Perbandingan diagnosa yang menderita gagal jantung, terjadi dari satu hingga dua dari setiap 100 orang dewasa dalam populasi umum dan lebih dari satu dari 10 orang berusia diatas 70 tahun [6]. Kemudian 6 – 10% individu diatas 60 tahun memiliki prevalensi penyakit gagal jantung yang menunjukkan pola meningkat seiring

bertambahnya usia [7]. Gagal jantung bagian dari proses patologis rumit yang bermula dari cedera jantung yang sering muncul dari penyakit jantung iskemik yang berkaitan dengan penyakit pembuluh darah koroner. Serta gagal jantung dapat diakibatkan dari kemoterapi ex posure dan bentuk kardiomiopati lainnya, seperti kardiomiopati bawaan dan virus [3].

Dari penjelasan terkait permasalahan tersebut, terdapat beberapa penelitian yang memprediksi terkait dari kelangsungan hidup, salah satunya yaitu penelitian yang dilakukan Minh Tuan Le, Minh Thanh Vo, Nhat Tan Pham, Son V.T Dao berjudul *Predicting Heart Failure Using A Wrapper-Based Feature Selection* pada tahun 2021 dengan hasil akurasi 83%. Dilihat dari hasil akurasi dapat dikembangkan lagi untuk mendapatkan akurasi terbaik, agar dapat membantu pencegahan dalam kelangsungan hidup bagi pasien penderita gagal jantung.

Penelitian yang akan dilakukan dengan melakukan analisis terhadap algoritma dan fitur yang dapat memberikan hasil yang optimal. Klasifikasi algoritma *random forest*, *random subspace* dan *logitboost* dengan metode bestfirst untuk seleksi fitur untuk menentukan prediksi kelangsungan hidup pada pasien penderita gagal jantung.

2. Tinjauan Pustaka

Prediksi kelangsungan hidup bagi pasien penderita gagal jantung, pastinya sudah ada para peneliti yang melakukan penelitian tersebut, berikut penelitian terkait yang relevan dengan penelitian kali ini. Terdapat 5 penelitian terkait diantaranya:

1. Penelitian yang dilakukan oleh Tariq Ahmad, Lars H. Lund, Pooja Rao, Rohit Ghosh, Prashant Warier, Benjamin Vaccaro, Ulf Dahlström, Christopher M. O'Connor, G. Michael Felker, dan Nihar R. Desai dengan judul "Machine Learning Methods Improve Prognostication, Identify Clinically Distinct Phenotypes, and Detect Heterogeneity in Response to Therapy in a Large Cohort of Heart Failure Patients" yang terbit pada tahun 2018. Penelitian tersebut menggunakan algoritma random forest yang menghasilkan AUC 0,83 [8].
2. Penelitian selanjutnya yaitu dari Faisal Maqbool Zahid, Shakeela Ramzan, Shahla Faisal dan Ijaz Hussain dengan judul "Gender based survival prediction models for heart failure patients: A case study in Pakistan" yang terbit pada tahun 2019. Penelitian tersebut menggunakan approach dengan hasil fitur untuk laki-laki, merokok, diabetes, dan anemia, sedangkan untuk perempuan, fraksi ejeksi, natrium, dan trombosit [9].
3. Penelitian selanjutnya yaitu dari Z. Kucukakcali, I. Balikci Cicek, E. Guldogan, and C. Colak dengan judul "Assessment Of Associative Classification Approach For Predicting Mortality By Heart Failure" yang terbit pada tahun 2020. Penelitian tersebut menggunakan Associative Classification Model dengan hasil akurasi Akurasi 0,866 [10].
4. Penelitian selanjutnya yaitu dari Davide Chicco and Giuseppe Jurman dengan judul "Machine Learning Can Predict Survival Of Patients With Heart Failure From Serum Creatinine And Ejection Fraction Alone" yang terbit pada 2020. Penelitian tersebut menggunakan algoritma random forest yang menghasilkan akurasi 83.8% [11], dan
5. Penelitian selanjutnya yaitu dari Minh Tuan Le, Minh Thanh Vo, Nhat Tan Pham, Son V.T Dao dengan judul "Predicting Heart Failure Using A Wrapper-Based Feature Selection" yang terbit pada tahun 2021. Penelitian tersebut menggunakan algoritma random forest dengan hasil akurasi 85% [12].

Penelitian terkait tersebut memiliki dataset yang sama dengan pengolahan data pada penelitian kali ini. Kemudian terdapat beberapa landasan teori yang relevan untuk penelitian kali ini diantaranya:

1. Gagal Jantung

Gagal jantung adalah suatu kondisi di mana jantung menjadi lemah dan tidak mampu memompa cukup darah ke seluruh tubuh. Kondisi ini dapat menyerang siapa saja, tetapi lebih sering terjadi pada orang yang berusia di atas 65 tahun [13]. Penyebab gagal jantung dikarenakan muncul karena katup jantung yang rusak, akibat melemahnya ruang jantung atau ventrikel kiri yang bertugas memompa darah ke seluruh tubuh dan akibat kakunya ventrikel sebelah kiri, sehingga jantung sulit terisi darah [14].

2. Data Mining

Data mining digunakan untuk menentukan informasi ataupun pengetahuan didalam sebuah database yang memiliki tugas seperti deskripsi, estimasi, prediksi, klasifikasi, clustering dan asosiasi [15]. Data mining menjadi salah satu solusi pembelajaran untuk menjelaskan proses penambangan informasi didalam basis data yang bersekala besar [16].

3. Random Forest

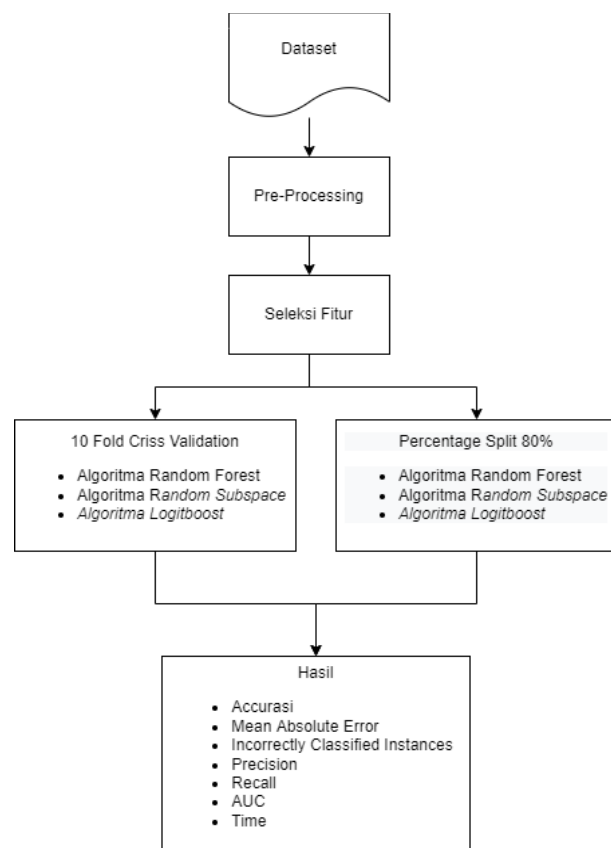
Random forest merupakan algoritma esamble learning yang menggunakan serta membangun struktur tree secara bertahap. Dalam penggunaan pohon keputusan dibangun berdasarkan memilih ataupun mengambil data secara acak. Penentuan kelas dalam suatu data menggunakan

sistem voting berdasarkan hasil dari decision tree [17].

4. Algoritma Best First Search

Algoritma Best First Search adalah metode yang meningkatkan simpul berikutnya dari suatu simpul terbaik yang diantara semia leaf nodes yang sudah pernah ditingkatkan [18]. Depth-first search dan breadth-first search menjadi kelebihan dari algoritma pencarian best first [19].

Selanjutnya terdapat tahapan-tahapan pada penelitian yaitu sebagai berikut:



Gambar 1. Tahapan Penelitian

a. Dataset

Dataset merupakan dataset publik yang didapat melalui website kaggle dengan nama dataset *heart failure*

b. *Preprocessing*

Preprocessing dilakukan untuk mengetahui data duplikat, *null* dan visualisasi data

c. Seleksi Fitur

Seleksi Fitur dilakukan untuk mengetahui fitur mana saja yang sangat berpengaruh, agar dapat menghasilkan prediksi yang maksimal.

d. Penerapan Algoritma

Algoritma yang digunakan yaitu perbandingan dari algoritma *random forest*, *random subspace* dan *logitboost*, untuk melihat performa mana yang paling baik.

e. Hasil Evaluasi

Hasil dari evaluasi pengolahan algoritma didapatkan *accurasi*, *mean absolute error*, *incorrectly classified instances*, *precision*, *recall*, *AUC*, dan *time*.

3. Metode Penelitian

Data yang digunakan merupakan data publik yang di peroleh dari website kaggle. Yang di upload oleh Larex yang merupakan ilmuwan data senior di rumah sakit Israelta Albert Einstein dengan nama dataset "Heart Failure" yang didapat dari aisalabad Institute of Cardiology and at the Allied Hospital in Faisalabad (Punjab,

Pakistan) dengan total data 299 penderita gagal jantung [20]. Metode penelitian yang digunakan pada penelitian kali ini terdiri dari tahapan pemahaman bisnis, pemahaman data, tahap persiapan data, pemodelan dan evaluasi. Adapun penjelasannya sebagai berikut:

1. Tahap Pemahaman Bisnis

Berdasarkan hasil pemeriksaan pasien penyakit gagal jantung di aisalabad Institute of Cardiology and at the Allied Hospital in Faisalabad (Punjab, Pakistan) yang diperoleh dari Kaggle dengan jumlah 299 orang , untuk mengetahui kelangsungan hidup pasien agar dapat dapat diketahui lebih awal dan dapat mendapatkan penanganan yang lebih tepat.

2. Tahap Pemahaman Data

Data tersebut terdiri dari 12 atribut predictor dan 1 atribut hasil yang dapat diketahui status pasien yang bertahan hidup (1) dan tidak bertahan hidup (0).

3. Tahap Persiapan Data

Persiapan data dengan melakukan pengecekan duplikat, null, visualisasi data kemudian penanganan seleksi fitur menghasilkan 4 fitur yang berpengaruh menggunakan metode bestfirst yaitu age, enjection_fraction, serum_creatinene dan time, serta penanganan imbalance class menggunakan model *class balancer*.

4. Tahapan Pemodelan

Tahapan pemodelan menggunakan 3 algoritma yaitu algoritma *random forest*, *random subspace* dan *logitboost* dengan metode cross validation dengan 10 fold dan percentage split 80%.

5. Tahap Evaluasi

Tahapan dari evaluasi ini akan melihat hasil dari *accurasi*, *mean absolute error*, *incorrectly classified instances*, *precision*, *recall*, *AUC*, *time* dari perbandingan algoritma *random forest*, *random subspace* dan *logitboost*.

4. Hasil dan Pembahasan

Perbandingan algoritma *random forest*, *random subspace* dan *logitboost* dengan menggunakan 10 fold cross validation dan percentage split 80%, maka dihasilkan akurasi yang disajikan dalam tabel 2 sebagai berikut:

Tabel 1. Perbandingan Accuracy

Algoritma	Mode Tes	Accuracy
Random Forest	Cross Validation	82.47%
	Percentage Split	91.45%
Random Subspace	Cross Validation	82.69%
	Percentage Split	83.94%
Logitboost	Cross Validation	82.44%
	Percentage Split	83.94%

Accuracy tertinggi dihasilkan oleh algoritma *random forest* dengan metode percentage split dengan hasil accuracy 91,45%, berikutnya hasil

dari mean absolute error yang disajikan dalam tabel 3 sebagai berikut:

Tabel 2. Perbandingan Mean Absolute Error

Algoritma	Mode Tes	Mean Absolute Error
Random Forest	Cross Validation	0.2309
	Percentage Split	0.1874
Random Subspace	Cross Validation	0.3314
	Percentage Split	0.332
Logitboost	Cross Validation	0.2372
	Percentage Split	0.2222

Mean absolute error yang paling rendah dihasilkan oleh algoritma *random forest* dengan metode percentage split dengan hasil 0.1874, berikutnya hasil dari *incorrectly classified instances* yang disajikan dalam tabel 4 sebagai berikut:

Tabel 3. Perbandingan Incorrectly Classified Instances

Algoritma	Mode Tes	Incorrectly Classified Instances
Random Forest	Cross Validation	17.53%
	Percentage Split	8.55%
Random Subspace	Cross Validation	17.31%
	Percentage Split	16.06%
Logitboost	Cross Validation	17.56%
	Percentage Split	16.06%

Incorrectly classified instances yang merupakan kesalahan dalam pemrosesan dengan persentase terendah dihasilkan oleh algoritma

random forest dengan metode percentage split dengan hasil 8.55%, kemudian berikut untuk hasil dari precision, recall dan AUC yang disajikan dalam tabel 5 sebagai berikut:

Tabel 4. Perbandingan Precision, Recall dan AUC

Algoritma	Mode Tes	Precision	Recall	AUC
Random Forest	Cross Validation	0.825	0.825	0.904
	Percentage Split	0.915	0.914	0.953
Random Subspace	Cross Validation	0.828	0.827	0.906
	Percentage Split	0.853	0.839	0.950
Logitboost	Cross Validation	0.826	0.824	0.900
	Percentage Split	0.853	0.839	0.906

Precision, recall dan AUC yang paling tinggi dihasilkan oleh algoritma random forest dengan metode percentage split dengan hasil 0.915, 0.914 dan 0.953, dan untuk waktu berdasarkan detik dengan hasil yang disajikan dalam tabel 6 sebagai berikut:

Tabel 5. Perbandingan Time

Algoritma	Mode Tes	Waktu
Random Forest	Cross Validation	0.23
	Percentage Split	0.01
Random Subspace	Cross Validation	0.13
	Percentage Split	0.01
Logitboost	Cross Validation	0.11
	Percentage Split	0

Waktu yang dihasilkan paling cepat dalam pemrosesan data yaitu saat memproses

algoritma logitboost dengan metode percentage split menghasilkan waktu 0 detik.

5. Kesimpulan

Berdasarkan hasil penelitian yang menggunakan aplikasi weka dengan melakukan seleksi fitur dengan metode bestfirst serta metode class balancer untuk menangani class yang tidak balance dan perbandingan terhadap 3 algoritma yang menunjukkan performa terbaik yaitu pada algoritma random forest dengan metode percentage split 80% yang menghasilkan *accuracy* 91,45%, mean absolute error 0.1874, *incorrectly classified instances* 8.55%, precision 0.915, recall 0.914, AUC 0.953. Sedangkan untuk waktu tercepat dihasilkan algoritma logitboost dengan metode percentage split dengan waktu 0 detik.

Performa algoritma random forest dengan metode percentage split 80% menghasilkan terbaik dengan dibandingkan dari penelitian sebelumnya yang memiliki *accurasi* 83%. Jika dibandingkan dari penelitian sebelumnya mendapatkan kenaikan akurasi sebesar 8,45%. Dengan kata lain penelitian yang dilakukan kali ini mengalami peningkatan akurasi sehingga permasalahan untuk mencari performa algoritma yang lebih baik dapat terselesaikan.

6. Daftar Pustaka

- [1] L. Ali *et al.*, "A Feature-Driven Decision Support System for Heart Failure Prediction Based on χ^2 Statistical Model and Gaussian Naive Bayes," *Comput. Math. Methods Med.*, vol. 2019, 2019, doi: 10.1155/2019/6314328.
- [2] Nursyamsiah and R. Hasan, "High-sensitivity c-reactive protein (hs-CRP) value with 90 days mortality in patients with heart failure," *IOP Conf. Ser. Earth Environ. Sci.*, vol. 125, no. 1, 2018, doi: 10.1088/1755-1315/125/1/012124.
- [3] T. R. Heallen, Z. A. Kadow, J. H. Kim, J. Wang, and J. F. Martin, "Stimulating Cardiogenesis as a Treatment for Heart Failure," *Circ. Res.*, vol. 124, no. 11, pp. 1647–1657, 2019, doi: 10.1161/CIRCRESAHA.118.313573.
- [4] S. E. Awan, M. Bennamoun, F. Sohel, F. M. Sanfilippo, and G. Dwivedi, "Machine learning-based prediction of heart failure readmission or death: implications of choosing the right model and the right metrics," *ESC Hear. Fail.*, vol. 6, no. 2, pp. 428–435, 2019, doi: 10.1002/ehf2.12419.
- [5] A. P. Utami, Fitri, "Gambaran Karakteristik Personal pada Pasien Gagal Jantung: A Narrative Review Article," *J. Ilm. Keperawatan Indones.*, vol. 5, no. 1, pp. 45–57, 2022, [Online]. Available: <https://www.google.com/search?q=jurnal+ilmiah+keperawatan&ie=utf-8&oe=utf-8>
- [6] N. R. Jones, A. K. Roalfe, I. Adoki, F. D. R. Hobbs, and C. J. Taylor, "Survival of patients with chronic heart failure in the community: a systematic review and meta-analysis," *Eur. J. Heart Fail.*, vol. 21, no. 11, pp. 1306–1325, 2019, doi: 10.1002/ehf.1594.
- [7] M.- Kamal, "Potential effectiveness of sleep hygiene and relaxation Benson in improving the quality of sleep in patients with heart failure: Literature review," *Int. J. Nurs. Heal. Serv.*, vol. 2, no. 1, pp. 101–107, 2019, doi: 10.35654/ijnhs.v2i1.69.
- [8] T. Ahmad *et al.*, "Machine learning methods improve prognostication, identify clinically distinct phenotypes, and detect heterogeneity in response to therapy in a large cohort of heart failure patients," *J. Am. Heart Assoc.*, vol. 7, no. 8, pp. 1–14, 2018, doi: 10.1161/JAHA.117.008081.
- [9] F. M. Zahid, S. Ramzan, S. Faisal, and I. Hussain, "Gender based survival prediction models for heart failure patients: A case study in Pakistan," *PLoS One*, vol. 14, no. 2, pp. 1–10, 2019, doi: 10.1371/journal.pone.0210602.
- [10] Z. Kucukakcali, I. B. Cicek, E. Guldogan, and C. Colak, "ASSESSMENT OF ASSOCIATIVE CLASSIFICATION APPROACH FOR PREDICTING MORTALITY BY HEART FAILURE," *J. Cogn. Syst.*, vol. 5, no. 2, pp. 41–45, 2020.
- [11] D. Chicco and G. Jurman, "Machine learning can predict survival of patients with heart failure from serum creatinine and ejection fraction alone," *BMC Med. Inform. Decis. Mak.*, vol. 20, no. 1, pp. 1–16, 2020, doi: 10.1186/s12911-020-1023-5.
- [12] M. T. Le, M. T. Vo, N. T. Pham, and S. V. T. Dao, "Predicting heart failure using a wrapper-based feature selection," *Indones. J. Electr. Eng. Comput. Sci.*, vol. 21, no. 3, pp. 1530–1539, 2021, doi: 10.11591/ijeecs.v21.i3.pp1530-1539.
- [13] R. Fadli, "Gagal Jantung."
- [14] P. Pittara, "Gagal Jantung," *Alodokter*, 2022. <https://www.alodokter.com/gagal-jantung>
- [15] B. A. Candra Permana and I. K. Dewi Patwari, "Komparasi Metode Klasifikasi Data Mining Decision Tree dan Naive Bayes Untuk Prediksi Penyakit Diabetes," *Infotek J. Inform. dan Teknol.*, vol. 4, no. 1, pp. 63–69, 2021, doi: 10.29408/jit.v4i1.2994.
- [16] S. Suhartini and R. Yuliani, "Penerapan Data Mining untuk Mengcluster Data Penduduk Miskin Menggunakan Algoritma

- K-Means di Dusun Bagik Endep Sukamulia Timur,” *Infotek J. Inform. dan Teknol.*, vol. 4, no. 1, pp. 39–50, 2021, doi: 10.29408/jit.v4i1.2986.
- [17] A. Wiraguna, S. Al Faraby, and Adiwijaya, “Klasifikasi Topik Multi Label pada Hadis Bukhari dalam Terjemahan Bahasa Indonesia Menggunakan Random Forest,” *e-Proceeding Eng.*, vol. 6, no. 1, pp. 2144–2153, 2019.
- [18] I. Maulana, M. Irawan Padli Nasution, and A. Ikhwan, “Aplikasi Pendaftaran Siswa Baru Menggunakan Algoritma Best First Search pada SMP Negeri 1 Medab,” *J. Chem. Inf. Model.*, vol. 53, no. 9, pp. 1689–1699, 2020.
- [19] L. I. Liana and S. R. Nudin, “Implementasi Algoritma Best-First Search untuk Aplikasi Mesin Pencari Handphone pada E-commerce (Apenphone),” *J. Informatics Comput. Sci.*, vol. 2, no. 01, pp. 67–73, 2020, doi: 10.26740/jinacs.v2n01.p67-73.
- [20] Larxel, “Heart Failure Prediction,” *Kaggle*, 2020.
<https://www.kaggle.com/datasets/andrewmvd/heart-failure-clinical-data>