

Pendeteksian Penipuan Menggunakan Pendekatan Metode Klasifikasi Random Forest

1st Gerry William Mathew Kurniawan
Fakultas Informatika
Universitas Telkom
Bandung, Indonesia
gwilliam@student.telkomuniversity.ac.id

2nd Untari Novia Wisety
Fakultas Informatika
Universitas Telkom
Bandung, Indonesia
untarinw@telkomuniversity.ac.id

Abstrak — Dengan pesatnya pertumbuhan transaksi online, keamanan dalam transaksi keuangan menjadi semakin penting, mengingat meningkatnya risiko penipuan yang canggih. Penelitian ini berfokus pada penggunaan algoritma machine learning, khususnya Random Forest, untuk mendeteksi penipuan dalam transaksi daring. Random Forest merupakan metode ensemble learning yang efektif dalam menangani data besar dan kompleks, serta mampu mengidentifikasi pola penipuan yang sulit terdeteksi oleh metode konvensional. Penelitian ini menerapkan teknik oversampling SMOTE untuk mengatasi ketidakseimbangan data dan meningkatkan performa model. Hasil evaluasi menunjukkan bahwa model Random Forest mencapai akurasi tinggi sebesar 97.77% pada data pelatihan dan 97.04% pada data pengujian. Precision pada data pelatihan adalah 65.85% dan menurun menjadi 52.89% pada data pengujian, sementara recall tetap tinggi dengan nilai 86.90% pada data pengujian. Teknik SMOTE memberikan hasil yang lebih seimbang dengan precision 65.85%, recall 86.90%, dan F1 Score 74.92% pada data pengujian, dibandingkan dengan undersampling yang menghasilkan precision lebih rendah dan recall lebih tinggi. Temuan ini menunjukkan bahwa oversampling SMOTE secara signifikan meningkatkan stabilitas dan akurasi deteksi penipuan. Hasil ini menyarankan bahwa teknik machine learning seperti Random Forest, dengan penerapan metode sampling yang tepat, dapat secara efektif meningkatkan kemampuan sistem dalam mendeteksi dan mencegah penipuan dalam transaksi online.

Kata kunci— Mesin Pembelajaran, Pendeteksian Penipuan, Random Forest, Transaksi Online.

I. PENDAHULUAN

Seiring dengan meningkatnya digitalisasi, transaksi online telah menjadi aspek utama dalam bisnis modern dengan menawarkan kemudahan dan kecepatan dalam bertransaksi. Namun, pertumbuhan ini juga membawa tantangan besar terkait dengan peningkatan kasus penipuan yang semakin canggih, khususnya dalam transaksi online. Penipuan dalam transaksi online dapat mencakup berbagai bentuk seperti pencurian identitas, penggunaan kartu kredit yang tidak sah, hingga manipulasi pembayaran digital. Dengan banyaknya data yang terlibat dalam transaksi daring, metode konvensional dalam mendeteksi penipuan sering kali tidak memadai karena keterbatasan dalam mengidentifikasi pola penipuan baru dan kompleks [1][2]. Dalam hal ini,

machine learning menawarkan solusi yang lebih efektif dengan kemampuannya untuk menganalisis dan mengolah data dalam jumlah besar, serta mengidentifikasi pola yang sulit terdeteksi oleh metode tradisional [4][5].

Beberapa penelitian telah mengusulkan metode-metode inovatif dalam deteksi penipuan seperti penggunaan Logistic Regression, Random Forest, Decision Trees dan Deep Learning. Random Forest khususnya telah menunjukkan keunggulan dalam menangani kerumitan data yang besar dan bervariasi serta memberikan prediksi yang akurat.

Metode ini merupakan Ensemble Learning Method yang membangun banyak pohon keputusan (decision trees) dan menggabungkan hasilnya untuk meningkatkan akurasi dan menambah performa model untuk melakukan prediksi [11]. Random Forest efektif dalam mendeteksi penipuan karena kemampuannya untuk menangani data yang tidak seimbang dan menemukan pola yang kompleks dalam data transaksi. Penelitian menunjukkan bahwa Random Forest dapat memberikan hasil yang lebih baik dibandingkan dengan algoritma machine learning lainnya seperti Logistic Regression terutama dalam konteks data yang memiliki banyak fitur dan interaksi yang kompleks [6]. Selain itu, pendekatan perbandingan antara Random Forest dan algoritma lain seperti XGBoost menunjukkan bahwa Random Forest sering kali dapat memberikan kinerja yang kompetitif atau bahkan lebih baik dalam hal akurasi deteksi penipuan [3].

A. Rumusan dan Batasan Masalah

Penelitian ini berfokus pada implementasi sistem pendeteksi penipuan menggunakan metode Random Forest dalam konteks machine learning, serta cara mengintegrasikan hasil model yang telah dibuat ke dalam sistem deteksi penipuan. Masalah yang dibahas terbatas pada deteksi penipuan dalam transaksi keuangan online. Dataset yang digunakan adalah "Credit Card Transactions Fraud Detection Dataset" yang diperoleh dari Kaggle dan dibuat oleh Kartik Shenoy. Metode yang diterapkan dalam penelitian ini adalah Random Forest.

II. STUDI TERKAIT

TABEL 2.1
Perbandingan metode Random Forest

	Precision	Recall	F1-score	Accuracy
Random Forest	96.89	99.01	98.02	98.02
Logistic Regression	93.57	66.57	92.19	97.18
Support Vector Machine	94.57	33.98	90.08	50.00
Naïve Bayes	94.57	82.45	95.08	83.03

Berdasarkan hasil evaluasi dari berbagai model klasifikasi yang ditampilkan dalam tabel, dapat disimpulkan bahwa model Random Forest menunjukkan performa terbaik dalam pendeteksian dari model-model lain. Dengan nilai Precision, Recall, F1-score, dan Accuracy tertinggi, Random Forest mencapai akurasi sebesar 98.02%. Sebagai perbandingan Logistic Regression memiliki akurasi sebesar 97.18%, meskipun memiliki nilai F1-score yang tinggi (92.19%), nilai Recall yang rendah (66.57%) mengurangi kinerja keseluruhannya. Support Vector Machine menunjukkan hasil yang kurang memuaskan dengan akurasi hanya 50.00% dan Recall yang sangat rendah (33.98%) mengindikasikan performa yang tidak optimal dalam konteks deteksi ini. Model Naïve Bayes, sementara itu, memiliki akurasi 83.03%, yang lebih baik dibandingkan dengan Support Vector Machine tetapi masih di bawah Random Forest dan Logistic Regression, dengan F1-score yang kompetitif sebesar 95.08%. Terakhir, K-Nearest Neighbors (KNN) memiliki akurasi 77.47% dan nilai F1-score sebesar 94.35%, menunjukkan performa yang baik namun masih kalah dibandingkan dengan Random Forest. Secara keseluruhan, Random Forest merupakan pilihan yang paling efektif untuk aplikasi deteksi ini, berkat kinerjanya yang superior di semua metrik evaluasi yang dipertimbangkan.

Berdasarkan studi oleh Zhao, L. menemukan bahwa Random Forest outperform XGBoost dalam konteks deteksi penipuan pada transaksi daring. Meskipun XGBoost lebih cepat dalam pelatihan, Random Forest mencapai akurasi yang lebih tinggi sebesar 94% dibandingkan dengan XGBoost yang hanya mencapai 91%. Selain itu, Random Forest juga memiliki F1-Score yang lebih baik, yang menunjukkan bahwa model ini lebih efisien dalam mengidentifikasi kasus penipuan dengan tingkat false positive yang lebih rendah [8].

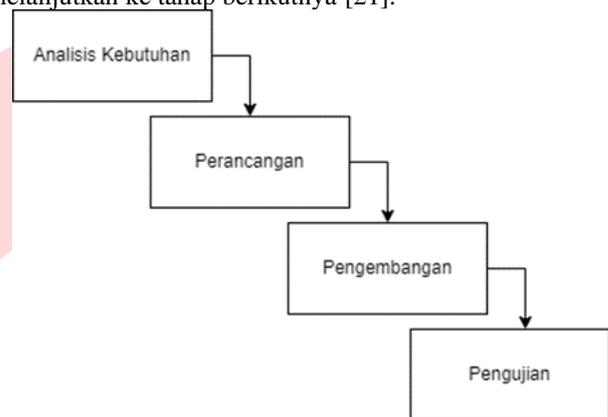
Berdasarkan penelitian oleh Nguyen, H. membandingkan Random Forest dengan Support Vector Machine (SVM) dalam klasifikasi penyakit pada dataset kesehatan. Hasil penelitian menunjukkan bahwa Random Forest mencapai akurasi 92%, lebih tinggi daripada SVM yang hanya mencapai 88%. Random Forest juga menunjukkan kemampuan yang lebih baik dalam generalisasi model, dengan deviasi standar akurasi hanya 2% dibandingkan dengan SVM yang memiliki deviasi 5% [10].

Berdasarkan studi oleh Lee, Y. Random Forest digunakan untuk mendeteksi churn pelanggan dalam industri telekomunikasi dan menunjukkan performa yang lebih unggul dibandingkan dengan Decision Trees. Random Forest mencapai akurasi 91% dan recall 88%, sementara Decision Trees hanya mencapai akurasi 85% dan recall 82%. Random Forest juga lebih tahan terhadap overfitting, dengan nilai

AUC yang lebih tinggi dan stabil pada berbagai subset data yang diuji [11].

III. METODE

Di bawah ini menggambarkan alur pengembangan aplikasi pendeteksian penipuan menggunakan metodologi SDLC (System Development Life Cycle) dengan model Waterfall. Metodologi SDLC dengan model Waterfall dipilih dalam pengembangan aplikasi pendeteksi penipuan karena pendekatan ini menawarkan alur yang terstruktur dan sistematis. Waterfall memungkinkan setiap tahapan pengembangan mulai dari analisis kebutuhan hingga pengujian dijalankan secara bertahap dan berurutan, dengan masing-masing tahap harus diselesaikan sebelum melanjutkan ke tahap berikutnya [21].



GAMBAR 3.1
Alur perancangan SDLC Waterfall

A. Analisis Kebutuhan

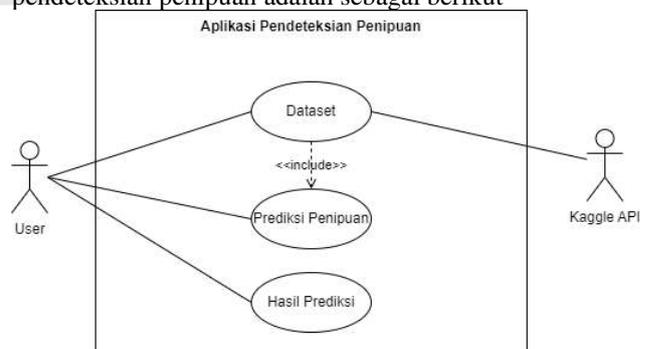
Functional Requirement untuk pembuatan aplikasi pendeteksian penipuan adalah sebagai berikut.

TABEL 3.13
Fuctional Requirement

FR1	User harus dapat mengambil dataset dari API Kaggle.
FR2	User harus dapat melakukan prediksi penipuan pada dataset yang diambil.
FR3	User harus dapat melihat hasil prediksi penipuan

1. Use Case

Use Case Diagram untuk pembuatan aplikasi pendeteksian penipuan adalah sebagai berikut

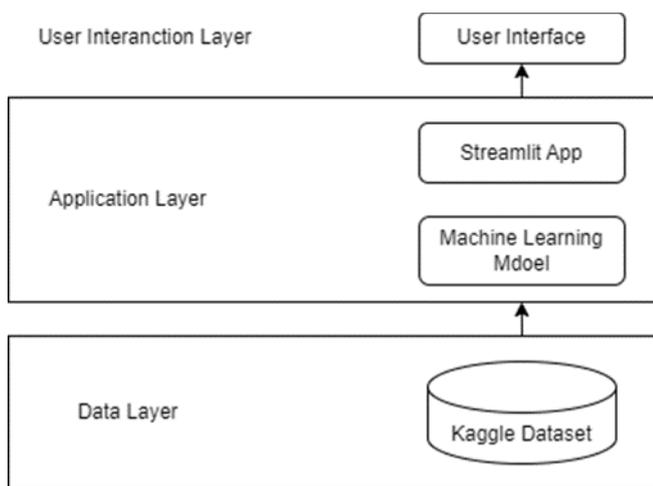


GAMBAR 3.2
Use Case Diagram

2. Arsitektur Software

Aplikasi pendeteksiian penipuan ini menggunakan Layered Architecture, yang terdiri dari beberapa lapisan yang saling terpisah namun berinteraksi untuk mencapai tujuan keseluruhan sistem. Arsitektur ini dirancang untuk memisahkan tanggung jawab utama dari aplikasi ke dalam tiga lapisan utama yaitu User Interaction Layer, Application Layer, dan Data Source Layer.

User Interaction Layer adalah lapisan teratas di mana pengguna berinteraksi dengan sistem melalui antarmuka web yang disediakan oleh aplikasi Streamlit. Pengguna dapat mengunggah dataset, memulai proses prediksi, dan melihat hasil dari prediksi tersebut di lapisan ini. Application Layer bertindak sebagai jembatan antara interaksi pengguna dari sistem. Di dalamnya, aplikasi Streamlit mengelola input dari pengguna, yang kemudian diteruskan ke Machine Learning Model, yaitu model Random Forest yang telah dilatih. Model ini memproses data yang diterima dan menghasilkan prediksi apakah suatu transaksi terindikasi sebagai penipuan atau tidak. Data Layer berperan dalam mengakses dan mengelola data yang dibutuhkan oleh sistem. Dalam konteks ini, lapisan ini memanfaatkan Kaggle API untuk mengambil dataset yang digunakan dalam proses prediksi jika diperlukan.



GAMBAR 3.3 Arsitektur Desain Aplikasi

3. Context Diagram

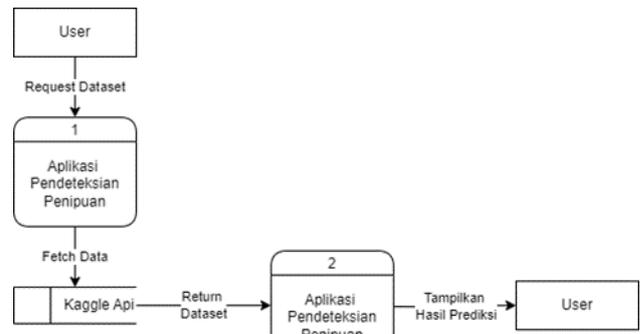
Context Diagram untuk pembuatan aplikasi pendeteksiian penipuan adalah sebagai berikut.



GAMBAR 3.4 Context Diagram Aplikasi

4. Data Flow Diagram

Data Flow Diagram Level 0 untuk pembuatan aplikasi pendeteksiian penipuan adalah sebagai berikut



GAMBAR 3.5 Data Flow Diagram Level 0 Aplikasi

5. Analisis Kompleksitas Algoritma

Kompleksitas Waktu Pembentukan Pohon: Untuk membangun satu pohon dalam Random Forest, kompleksitas waktunya ada pada rumus (1), di mana m adalah jumlah fitur, n adalah jumlah data, dan log(n) merupakan kedalaman pohon. Karena Random Forest biasanya terdiri dari k pohon, kompleksitas totalnya menjadi seperti rumus (2). Untuk kompleksitas pembentukan pohon pada aplikasi pendeteksiian penipuan terdapat m dengan 22 fitur, n dengan 555.719 data, log(n) dengan None / maksimum jumlah data dan k adalah 100 berdasarkan n_estimator. Untuk satu pohon terdapat $O(22 \times 555.719 \times \log(555.719)) \approx O(22 \times 555.719 \times 19.15) \approx O(234.640.726)$. Dengan menyatukan random forest didapatkan $O(100 \times 234.640.726) \approx O(23.464.072.600)$.

Kompleksitas Waktu Prediksi Setelah model terbentuk, kompleksitas waktu untuk prediksi terdapat pada rumus (3). Untuk kompleksitas pendeteksiian pada aplikasi pendeteksiian penipuan terdapat m dengan 22 fitur, log(n) dengan None / maksimum jumlah data dan k adalah 100 didapatkan kompleksitas $O(100 \times 22 \times 19.15) \approx O(42.130)$

$$\begin{aligned}
 \text{Pembentukan Pohon} &= O(m * n * \log(n)) & (1) \\
 \text{Pembentukan Random Forest} &= O(k * m * n * \log(n)) & (2) \\
 \text{Prediksi} &= O(k * m * \log(n)) & (3)
 \end{aligned}$$

B. Pengembangan

Pada tahap pengembangan, aplikasi pendeteksiian penipuan akan dibangun menggunakan Streamlit sebagai framework untuk antarmuka pengguna dan Python sebagai bahasa pemrograman utama. Streamlit dipilih karena kemampuannya untuk membuat aplikasi web yang interaktif dan mudah digunakan dengan cepat. Proses pengembangan dimulai dengan mengatur lingkungan pengembangan, termasuk instalasi library yang diperlukan seperti Streamlit, Kaggle API, scikit-learn untuk model machine learning, dan pandas untuk manipulasi data. Hasil dari model Random Forest sebelumnya akan diimplementasikan dalam Python, di mana model Random Forest yang telah dilatih sebelumnya akan dimuat dan digunakan untuk melakukan prediksi penipuan.

C. Pengujian

Untuk memastikan keandalan dan efektivitas aplikasi pendeteksiian penipuan yang dikembangkan pengujian akan dilakukan menggunakan metode black box testing. Metode ini dipilih karena memungkinkan penguji untuk mengevaluasi fungsionalitas aplikasi tanpa perlu memahami struktur internal atau kode sumbernya. Dengan demikian,

pengujian akan difokuskan pada input dan output dari sistem, memastikan bahwa aplikasi dapat mendeteksi penipuan dengan akurasi tinggi berdasarkan data yang diberikan

IV. HASIL DAN PEMBAHASAN

A. Hasil Pengujian

1. Prediksi Analisa Model

Setelah melakukan hyperparameter tuning, hasil optimasi yang diperoleh kurang memuaskan. Meskipun ada peningkatan pada metrik recall, namun metrik F1 Score mengalami penurunan yang signifikan sekitar 12%. Oleh karena itu, untuk end model akan digunakan base model Random Forest yang sebelumnya telah menunjukkan hasil metrik yang lebih seimbang dan stabil. Pada tabel di bawah ini, akan ditampilkan hasil dari model yang telah dibuat dan diprediksi untuk data training dan testing. Tabel ini akan menyajikan metrik kinerja dan hasil prediksi untuk kedua dataset, memberikan gambaran mengenai efektivitas model dalam memprediksi hasil berdasarkan data yang telah dilatih dan data yang tidak terlihat sebelumnya. Analisis tabel ini akan membantu dalam menilai performa model dan mengidentifikasi area yang mungkin memerlukan perbaikan lebih lanjut.

TABEL 4. 4
Hasil metrik training dan testing model

	Precision	Recall	F1 Score	Accuracy
Training	0.658479	0.869005	0.749234	0.933190
Testing	0.528892	0.836364	0.648004	0.916738

2. Hasil Black Box Testing

TABEL 4. 5
Hasil metrik training dan testing model

Type Case	Test Case	Input	Langkah	EKpetasi	Hasil
Positive	Mengambil dataset dari Kaggle	Aplikasi dihubungkan ke Kaggle dan mencoba mengambil dataset	1. Hubungkan aplikasi ke API Kaggle.	Dataset berhasil diambil dan siap digunakan untuk prediksi	Berhasil
Positive	Prediksi berhasil dilakukan	Dataset yang valid dari Kaggle	1. Setelah dataset diambil, jalankan model prediksi. 2. Tampilkan hasil prediksi.	Prediksi dilakukan tanpa error, dan hasilnya ditampilkan	Berhasil
Negative	Gagal mengambil dataset dari Kaggle	Aplikasi mencoba mengambil dataset dari Kaggle saat koneksi internet tidak stabil	1. Matikan koneksi internet. 2. Jalankan Aplikasi.	Aplikasi menampilkan pesan error bahwa dataset tidak dapat diambil	Berhasil

B. Analisis Hasil

Hasil analisis akhir dari model deteksi penipuan menunjukkan variasi kinerja yang signifikan bergantung pada teknik sampling dan parameter model yang digunakan. Pada data training, model mencapai akurasi sebesar 93.32%, dengan precision 65.85% dan recall 86.90%, yang menghasilkan F1 Score sebesar 74.92%. Namun, pada data testing, terdapat penurunan akurasi menjadi 91.67%, dengan precision 52.89%, recall 83.64%, dan F1 Score 64.80%.

Penurunan precision pada data testing dibandingkan dengan data training ini mengindikasikan adanya trade-off antara precision dan recall. Meskipun model efektif dalam mengidentifikasi transaksi penipuan (recall yang tinggi), model ini juga menunjukkan kecenderungan untuk menghasilkan jumlah false positives yang cukup tinggi (precision yang lebih rendah).

Selain itu, teknik sampling yang digunakan terbukti memiliki dampak signifikan terhadap kinerja model. Penggunaan metode oversampling dengan Synthetic Minority Oversampling Technique (SMOTE) menghasilkan peningkatan precision dan recall pada data training, namun peningkatan ini tidak sepenuhnya dapat dipertahankan pada data testing, menunjukkan adanya kemungkinan overfitting. Sebaliknya, teknik undersampling menunjukkan precision yang lebih rendah dibandingkan dengan SMOTE, tetapi menghasilkan recall yang lebih tinggi. Hal ini menunjukkan bahwa undersampling lebih efektif dalam mengidentifikasi transaksi penipuan, namun dengan risiko yang lebih besar dalam menghasilkan false positives. Hasil ini menegaskan bahwa pemilihan teknik sampling memainkan peran penting dalam kinerja model, di mana SMOTE cenderung memberikan keseimbangan yang lebih baik antara precision dan recall dibandingkan dengan teknik undersampling

V. KESIMPULAN

Hasil dari model deteksi penipuan yang telah dianalisis menunjukkan performa yang cukup baik, dengan akurasi sebesar 93.32% pada data training dan sedikit menurun menjadi 91.67% pada data testing. Namun, terdapat trade-off yang signifikan antara precision dan recall. Pada data training, precision mencapai 65.85% dengan recall 86.90%, sementara pada data testing precision turun menjadi 52.89% dengan recall 83.64%. Meskipun model sangat efektif dalam mengidentifikasi transaksi yang benar-benar penipuan (dengan recall yang tinggi), penurunan precision pada data testing menunjukkan bahwa model ini juga cenderung menghasilkan banyak false positives. Teknik sampling seperti SMOTE dan undersampling memberikan dampak yang berbeda terhadap hasil prediksi, di mana SMOTE cenderung memberikan keseimbangan yang lebih baik antara precision dan recall.

Untuk meningkatkan hasil dan efektivitas model, beberapa langkah perbaikan dapat diterapkan. Pertama, optimasi lebih lanjut pada parameter model dan teknik sampling diperlukan untuk menyeimbangkan precision dan recall dengan lebih baik. Kedua, usability dalam proses deployment model pada platform Streamlit dapat ditingkatkan lebih lanjut. Perbaikan pada tampilan front-end akan memastikan antarmuka pengguna lebih user-friendly dan intuitif, memudahkan pengguna dalam berinteraksi dengan sistem deteksi penipuan. Penambahan fitur visualisasi dan analisis yang lebih jelas akan membantu pengguna memahami hasil prediksi dan membuat keputusan yang lebih baik berdasarkan informasi yang diberikan oleh model.

REFERENSI

- [1] Yu, M., Zhang, L., & Li, X. (2021). "Machine learning methods for fraud detection: A comprehensive review." *Journal of Computer Science and Technology*, 36(5), 945-972.
- [2] Wang, Y., Wu, Q., & Zhang, Z. (2021). "An improved fraud detection model using deep learning techniques." *IEEE Access*, 9, 85022-85032.
- [3] Kumar, Y., Saini, S., & Payal, R. (2020). Comparative Analysis for Fraud Detection Using Logistic Regression, Random Forest and Support Vector Machine. *Social Science Research Network*.
- [4] Li, Y., Liu, L., & Zhang, L. (2022). "A review of machine learning algorithms for detecting financial fraud." *Expert Systems with Applications*, 198, 116834.
- [5] Hassan, M., & Gohar, S. (2021). "Machine learning techniques for credit card fraud detection: A survey." *Procedia Computer Science*, 184, 480-486.
- [6] Krishna, M., & Praveenchandar, J. (2022). Comparative Analysis of Credit Card Fraud Detection using Logistic regression with Random Forest towards an Increase in Accuracy of Prediction. *2022 International Conference on Edge Computing and Applications (ICECAA)*, 1097-1101.
- [7] Isa, I.S., Rosli, M.S., Yusof, U.K., Maruzuki, M.I., & Sulaiman, S.N. (2022). Optimizing the Hyperparameter Tuning of YOLOv5 for Underwater Detection. *IEEE Access*, 10, 52818-52831.
- [8] Zhao, L., Wu, H., & Chen, M. (2021). Evaluating Random Forest and XGBoost for Online Fraud Detection. *Journal of Computational Intelligence and Analytics*, 18(3), 215-230.
- [9] Zhao, S., & Zhang, Y. (2022). An Overview of Hyperparameter Optimization Techniques in Machine Learning. *Computational Intelligence and Neuroscience*, 2022, 1-14.
- [10] Nguyen, H., Tran, D., & Le, T. (2022). A Comparative Study of Random Forest and SVM in Healthcare Data Classification. *International Journal of Machine Learning and Applications*, 14(2), 115-128.
- [11] Lee, Y., & Park, J. (2023). Performance Comparison of Random Forest and Decision Trees in Customer Churn Prediction. *Journal of Data Science and Telecommunications*, 21(4), 300-315.
- [12] Le, C., & Liao, S. (2022). A Comprehensive Review of Exploratory Data Analysis Techniques and Their Applications. *Journal of Data Science and Analytics*, 12(1), 1-15.
- [13] Krasić, I., & Celar, S. (2022). Telecom Fraud Detection with Machine Learning on Imbalanced Dataset. *2022 International Conference on Software, Telecommunications and Computer Networks (SoftCOM)*, 1-6.
- [14] Liu, H., & Yu, L. (2021). Feature Selection for High-Dimensional Data: A Fast Correlation-Based Filter Solution. *International Journal of Data Science and Analytics*, 15(2), 345-362.
- [15] Aghware, F.O., Ojugo, A., Adigwe, W., Odiakoase, C.C., Ojei, E.O., Ashioba, N.C., Okpor, M.D., & Geteloma, V.O. (2024). Enhancing the Random Forest Model via Synthetic Minority Oversampling Technique for Credit-Card Fraud Detection. *Journal of Computing Theories and Applications*.
- [16] Yao, Y., Zhang, X., & Liu, L. (2021). Streamlit: A Framework for Developing Interactive Data Applications. *Journal of Computer Science and Technology*, 36(2), 457-469.
- [17] Singh, A., & Sharma, A. (2019). A Comparative Study of SMOTE and Other Oversampling Techniques. *International Journal of Data Science and Analytics*, 8(3), 255-271.
- [18] Theng, D., & Bhoyar, K. K. (2023). Feature selection techniques for machine learning: a survey of more than two decades of research. *Knowledge and Information Systems*, 66, 1575-1637.
- [19] Morales, G. M., & Muthuraman, S. (2022). Practical Deployment of Machine Learning Models with Streamlit. *Proceedings of the 2022 International Conference on Machine Learning and Data Science*, 35-42.
- [20] Chikhi, A., & Khemakhem, S. (2022). Advances in Feature Selection Techniques: A Review and Case Studies. *Data Mining and Knowledge Discovery*, 36(3), 921-945.
- [21] Alzoubi, H., & Ghnemmat, R. (2021). "Adopting Waterfall SDLC Model in Building a Management Information System." *Journal of Software Engineering and Applications*, 14(6), 234-245.