

BAB II

LANDASAN TEORI

2.1 Penelitian Terdahulu

Dalam mengembangkan kerangka penelitian untuk analisis komparatif algoritma *Random Forest* dan *Gradient Boosting* pada prediksi lama menginap tamu hotel, diperlukan tinjauan mendalam terhadap penelitian-penelitian sebelumnya yang telah mengeksplorasi penerapan *Machine Learning* dalam industri perhotelan. Kajian literatur ini mencakup delapan penelitian kunci yang dipilih berdasarkan relevansi metodologi, konteks aplikasi, dan kualitas publikasi, meliputi:

Tabel 2. 1 Penelitian Terdahulu

Referensi ke-1	
Judul	<i>Hotel Guest Length of Stay Prediction Using Random Forest Regressor</i>
Nama Penulis	Yerik Afrianto Singgalen [26].
Sumber Jurnal	<i>Journal of Information Systems and Informatics</i> (SINTA 3)
Tahun	2024
Permasalahan	Memprediksi <i>length of stay</i> (LoS) hotel berdasarkan fitur <i>country</i> , <i>guest type</i> , <i>room type</i> , dan <i>rating</i> dengan akurasi tinggi.
Framework	<i>Random Forest Regression</i> , <i>Feature Engineering</i> , <i>Data Preprocessing</i> dengan <i>encoding categorical variables</i> .
Temuan	$R^2 = 0.85$, $MAE = 1.06$. <i>Feature importance</i> : <i>Country</i> (0.5), <i>Guest Type</i> (0.2), <i>Room Type</i> (0.15), <i>Rating</i> (0.15)
Pembahasan	<i>Random Forest</i> efektif menangani data <i>non-linear</i> dengan performa tinggi. <i>Country</i> menjadi prediktor terkuat dalam menentukan <i>length of stay</i> (LoS).
Relevansi	Memberikan <i>baseline</i> kuat untuk implementasi <i>Random Forest</i> dalam prediksi <i>length of stay</i> (LoS) hotel serta <i>benchmark</i> metrik evaluasi (R^2 , MAE) yang dapat dibandingkan dengan <i>Gradient Boosting</i> pada penelitian ini
Referensi ke-2	
Judul	<i>Big Data in Tourism and Hospitality Industry: Predictive Analytics of Hotel Room Trend</i>
Nama Penulis	Yerik Afrianto Singgalen [27].
Sumber Jurnal	<i>Indonesian Journal of Tourism and Leisure</i> (SINTA 3)
Tahun	2025

Permasalahan	Implementasi <i>XGBoost</i> untuk prediksi <i>room type</i> dan <i>trend okupansi</i> hotel berbasis <i>big data analytics</i> .
Framework	<i>XGBoost</i> , <i>Gradient Boosting</i> , <i>SMOTE</i> untuk <i>class imbalance</i> , <i>Grid Search hyperparameter optimization</i> .
Temuan	Akurasi prediksi 85% untuk <i>room selection patterns</i> . <i>Feature importance</i> : <i>Rating</i> (0.35), <i>Length of Stay</i> (0.28), <i>Guest Type</i> (0.22).
Pembahasan	<i>XGBoost</i> unggul dalam menangani dataset kompleks dan <i>imbalanced data</i> . <i>Gradient boosting</i> efektif untuk identifikasi pola tersembunyi.
Relevansi	Validasi keunggulan algoritma <i>Gradient Boosting (XGBoost)</i> pada dataset hotel dan memberikan <i>insight</i> tentang <i>feature engineering</i> serta <i>handling imbalanced data</i> yang relevan dan sejalan dengan penelitian ini.

Referensi ke-3

Judul	<i>Booker Prediction from Requests for Quotation via Machine Learning Techniques</i>
Nama Penulis	Samuel Runggaldier, Gabriele Sottocornola, Andrea Janes, Fabio Stella, Markus Zanker [28].
Sumber Jurnal	<i>Tourism and Hospitality Management</i> . (Scopus Q3)
Tahun	2023
Permasalahan	Memprediksi apakah permintaan penawaran (<i>request for quotation</i>) akan berubah menjadi pemesanan aktual menggunakan teknik <i>Machine Learning</i> untuk mengoptimalkan prioritas manajemen korespondensi hotel.
Framework	<i>Random Forest</i> , <i>Extra Tree</i> , <i>Gaussian Naive Bayes</i> , <i>Multi-layer Perceptron (MLP)</i> , <i>Support Vector Classifier (SVC)</i> , <i>10-fold Cross-Validation</i> , <i>Class Imbalance Subsampling</i> , <i>Feature Selection</i> dengan <i>Entropy-based Scoring</i> .
Temuan	<i>Random Forest</i> mencapai <i>F1-score</i> tertinggi (35%), <i>precision</i> 48%, <i>recall</i> 27%. <i>MLP</i> mencapai <i>precision</i> tertinggi (66%) dan <i>accuracy</i> 92%. <i>Naive Bayes</i> mencapai <i>recall</i> tertinggi (65%). <i>Feature</i> terpenting: <i>CR_RequestedDaysBeforeArrival</i> (0.31), <i>CG_CountryCode</i> (0.20), <i>CR_Duration</i> (0.12).
Pembahasan	Algoritma <i>Random Forest</i> dan <i>Extra Tree</i>) menunjukkan performa superior dalam prediksi <i>booking</i> hotel. <i>Class imbalance subsampling</i> dengan ratio 1:3 menghasilkan <i>trade-off</i> optimal <i>precision-recall</i> . <i>MLP</i> dan <i>Random Forest</i> paling menjanjikan dengan <i>F1-score</i> yaitu 42% setelah optimasi.
Relevansi	Memberikan validasi empiris tentang efektivitas <i>Random Forest</i> dalam prediksi hotel <i>booking</i> dan metodologi <i>handling class imbalance</i> . <i>Framework feature selection</i> dan <i>comparative analysis multiple algorithms</i> sangat relevan untuk penelitian komparasi model.

Referensi ke-4

Judul	<i>Development of Machine Learning Model to Predict Hotel Room Reservation Cancellations</i>
-------	--

Nama Penulis	Eka Rahmawati, Galih Setiawan Nurohim, Candra Agustina, Denny Irawan, Zainal Muttaqin [29].
Sumber Jurnal	Jurnal Teknologi Informasi Dan Terapan (J-TIT) (SINTA 2)
Tahun	2024
Permasalahan	Prediksi pembatalan reservasi hotel menggunakan <i>multiple Machine Learning algorithms</i> pada <i>dataset</i> Hotel Borobudur.
Framework	<i>Random Forest, Logistic Regression, Bayesian Networks, Cross-validation.</i>
Temuan	<i>Random Forest</i> : Akurasi 86.36%, Precision 88.06%, Recall 93.65%, <i>F1-score</i> 90.77%. <i>Random Forest outperforms Logistic Regression dan Bayesian Networks.</i>
Pembahasan	<i>Random Forest</i> sangat <i>robust</i> untuk <i>non-linear relationships</i> dan <i>imbalanced</i> hotel data. <i>Bayesian Networks underperform</i> pada dataset kompleks.
Relevansi	Konfirmasi superioritas <i>Random Forest</i> pada prediksi hotel-related tasks dalam konteks Indonesia, serta memberikan <i>benchmark</i> performa untuk evaluasi komparatif.
Reverensi ke-5	
Judul	<i>Predicting Hotel Booking Cancellations Using Machine Learning for Revenue Optimization</i>
Nama Penulis	Andy Hermawan, Aji Saputra, Nabila Lailinajma, Reska Julianti, Timothy Hartanto, Troy Kornelius Daniel [30].
Sumber Jurnal	<i>Router</i> : Jurnal Teknik Informatika dan Terapan (SINTA 3).
Tahun	2025
Permasalahan	Optimasi model <i>Machine Learning</i> untuk prediksi dan mitigasi pembatalan <i>booking</i> hotel dengan fokus <i>revenue optimization</i> .
Framework	<i>XGBoost, SMOTE resampling, Threshold tuning, SHAP interpretation.</i>
Temuan	<i>F2-score</i> 0.79 (<i>XGBoost</i> terbaik). Finansial: penghematan \$286,800. <i>Key features</i> : <i>deposit type, special requests</i> , dan <i>market segment</i> .
Pembahasan	<i>XGBoost</i> dengan <i>hyperparameter tuning</i> dan <i>threshold adjustment</i> memberikan ROI tertinggi dalam hotel <i>revenue management</i> .
Relevansi	Demonstrasi aplikasi praktis <i>Gradient Boosting</i> dalam konteks bisnis hotel Indonesia, serta metodologi <i>threshold optimization</i> yang dapat diadaptasi.
Referensi ke-6	
Judul	<i>Performance Comparison of Random Forest (RF) and Classification and Regression Trees (CART) for Hotel Star Rating Prediction</i>
Nama Penulis	Annisaa Utami, Dimas Fanny Hebrasianto Permadi, Yesy Diah Rosita, Jumanto Unjung [31].
Sumber Jurnal	<i>SJI: Scientific Journals of Informatics</i> (SINTA 3)
Tahun	2024

Permasalahan	Evaluasi komparatif <i>Random Forest</i> vs <i>CART</i> untuk klasifikasi rating hotel berbasis <i>customer reviews</i> .
Framework	<i>Random Forest</i> , <i>CART</i> , <i>Hyperparameter Tuning</i> , <i>Cross-validation</i> .
Temuan	<i>Random Forest</i> : Akurasi >99%, <i>superior robustness</i> . <i>CART</i> : <i>prone to overfitting</i> pada dataset besar.
Pembahasan	<i>Random Forest</i> konsisten unggul dibanding <i>CART</i> pada berbagai ukuran dataset hotel. <i>Ensemble method</i> lebih stabil.
Relevansi	Validasi keunggulan <i>Random Forest</i> dalam <i>domain</i> hotel <i>classification tasks</i> dan pentingnya analisis yang komparatif antara <i>ensemble</i> vs <i>single tree algorithms</i> .
Referensi ke-7	
Judul	<i>Personalized Restaurant Recommendations: A Hybrid Filtering Approach for Mobile Applications</i>
Nama Penulis	Christopher Matthew Marvelio, Alexander Waworuntu [32].
Sumber Jurnal	<i>International Journal of New Media Technology</i> UMN (IJMT UMN)
Tahun	2025
Permasalahan	Pengembangan sistem rekomendasi restaurant menggunakan <i>hybrid filtering</i> (<i>content-based</i> dan <i>collaborative filtering</i>).
Framework	<i>Hybrid Filtering</i> , <i>Content-Based Filtering</i> , <i>Collaborative Filtering</i> , <i>React Native</i> , <i>Flask</i> .
Temuan	<i>EUCS satisfaction</i> 93.9%. <i>The hybrid approach outperforms single filtering methods significantly</i> .
Pembahasan	Kombinasi <i>multiple filtering techniques</i> menghasilkan rekomendasi lebih akurat dan personal dibanding <i>single approach</i> .
Relevansi	Memberikan <i>insight</i> tentang <i>hybrid Machine Learning approaches</i> dan metodologi evaluasi <i>user satisfaction</i> yang dapat diadaptasi untuk penelitian prediksi hotel LoS.
Referensi ke-8	
Judul	<i>Adverse Media Classification: A New Era of Risk Management with XGBoost and Gradient Boosting Algorithms</i>
Nama Penulis	Reza Juliandri, Monika Evelin Johan, Jansen Wiratama, Samuel Ady Sanjaya [33].
Sumber Jurnal	2024 5th <i>International Conference on Big Data Analytics and Practices</i> (IBDAP) – IEEE Xplore
Tahun	2024
Permasalahan	Klasifikasi berita negatif (<i>adverse media</i>) untuk mendukung sistem <i>customer due diligence</i> dalam industri <i>fintech</i> , mengurangi risiko <i>fraud</i> dan <i>money laundering</i> .
Framework	<i>XGBoost</i> , <i>Gradient Boosting</i> , <i>Back Translation</i> , <i>Grid Search Hyperparameter Tuning</i> , <i>10-fold Cross-Validation</i> , <i>Class Imbalance Subsampling</i> .

Temuan	<i>Gradient Boosting</i> (753 records) mencapai akurasi 82.31% pada test data dan 84.93% pada <i>validation data</i> , <i>XGBoost</i> (1,281 records) mencapai akurasi 78.54% pada test data dan 82.52% pada <i>validation data</i> .
Pembahasan	<i>Back translation</i> efektif menangani <i>class imbalance</i> , <i>Grid Search</i> yang membantu optimasi pada <i>parameter</i> , dan <i>ensemble methods</i> unggul dalam klasifikasi <i>news headline</i> berbasis teks.
Relevansi	Metodologi <i>comparative analysis</i> dan teknik <i>preprocessing</i> teks (<i>back translation</i> , <i>TF-IDF/CountVectorizer</i>) dapat diadaptasi untuk prediksi lama menginap tamu hotel menggunakan <i>Random Forest</i> dan <i>Gradient Boosting</i> .

Seluruh penelitian terdahulu yang dirangkum dalam Tabel 2.1 memiliki area kajian yang sejenis, yaitu penerapan algoritma *Machine Learning* pada industri perhotelan dan sektor terkait, dengan fokus utama pada prediksi perilaku pelanggan, okupansi, hingga pembatalan reservasi. Berdasarkan hasil dari delapan jurnal tersebut, penelitian pertama dan kedua [26] [27], menjadi acuan utama dalam penelitian ini karena keduanya membahas penerapan *Random Forest* dan *XGBoost* dalam konteks prediksi hotel. Pada penelitian pertama, algoritma *Random Forest* berhasil mencapai nilai R^2 sebesar 0,85 dan MAE 1,06, menunjukkan performa tinggi dalam memprediksi *length of stay* (LoS) berdasarkan fitur seperti *country* dan *guest type* [26]. Sementara pada penelitian keduanya, model *XGBoost* memperoleh akurasi prediksi sebesar 85%, membuktikan kemampuannya dalam menangani data hotel berskala besar.

Penelitian lain turut memberikan kontribusi penting terhadap pengembangan model prediktif. Penelitian yang membandingkan berbagai algoritma seperti *Random Forest*, *MLP*, dan *Naïve Bayes*, dengan hasil *F1-score* 42% setelah optimasi, menunjukkan kinerja lebih baik dalam prediksi *hotel booking conversion* [28]. Selanjutnya penelitian lain juga memperkuat temuan tersebut dengan menunjukkan bahwa *Random Forest* mencapai akurasi 86,36%, *precision* 88,06%, dan *recall* 93,65%, mengungguli *Logistic Regression* dan *Bayesian Networks* [29]. Penelitian lain juga membuktikan superioritas *XGBoost* dengan hasil *F2-score* 0,79 serta dampak finansial berupa penghematan sebesar \$286.800, menegaskan potensi algoritma ini dalam optimasi pendapatan hotel melalui prediksi pembatalan.

Selain itu, penelitian yang menunjukkan bahwa *Random Forest* memiliki akurasi di atas 99% pada klasifikasi *rating* hotel, mengungguli metode *CART* yang rentan terhadap *overfitting* [31]. Pendekatan serupa juga tampak pada penelitian terdahulu yang ke tujuh yang menggunakan metode *hybrid filtering* dan menghasilkan tingkat kepuasan pengguna (EUCS) sebesar 93,9%, membuktikan efektivitas kombinasi pendekatan *Machine Learning* dalam sistem rekomendasi restoran [32]. Adapun penelitian terakhir menjelaskan terkait performa *Gradient Boosting* dan *XGBoost* dalam klasifikasi teks, dengan akurasi masing-masing 84,93% dan 82,52%, serta menekankan pentingnya teknik *feature engineering* seperti *back translation* untuk meningkatkan kualitas model.

Berdasarkan hasil dari seluruh penelitian terdahulu, dapat disimpulkan bahwa algoritma *Random Forest* dan *Gradient Boosting* memiliki keunggulan signifikan dalam hal akurasi, stabilitas, dan kemampuan menangani data kompleks maupun tidak seimbang. Oleh karena itu, dalam penelitian ini, kedua algoritma tersebut diadopsi untuk membangun model prediksi *length of stay* (LoS) tamu hotel. Analisis dilakukan untuk membandingkan performa keduanya dengan menggunakan metrik regresi standar yaitu RMSE, MAE, dan R^2 , tanpa menerapkan metrik klasifikasi.

Berbeda dengan penelitian terdahulu yang banyak berfokus pada prediksi pembatalan reservasi, konversi pemesanan, atau rekomendasi layanan, penelitian ini memiliki ruang kajian yang lebih spesifik. Penelitian ini berfokus pada pembangunan model prediksi lama menginap tamu (*length of stay*) dengan memanfaatkan data transaksi nyata dari operasional hotel.

2.2 Tinjauan Teori

2.2.1 Machine Learning

Machine Learning adalah bidang ilmu komputer yang berfokus pada pengembangan algoritma dan teknik yang memungkinkan sistem untuk belajar dari data dan meningkatkan kinerjanya seiring bertambahnya pengalaman tanpa pemrograman eksplisit [34]. Konsep ini dipopulerkan oleh Tom M. Mitchell dengan menegaskan bahwa pembelajaran berorientasi pada peningkatan

performa model yang diukur dengan metrik tertentu ketika diberikan lebih banyak data [35]. Algoritma pembelajaran mesin pada dasarnya dapat diklasifikasikan ke dalam tiga kategori utama, yakni *supervised learning*, *unsupervised learning*, dan *reinforcement learning*. Pada *supervised learning*, proses pelatihan dilakukan dengan memanfaatkan data yang telah dilabeli untuk membimbing model dalam melakukan prediksi. Sebaliknya, *unsupervised learning* berfokus pada penggalian pola atau struktur yang tersembunyi tanpa adanya bantuan label pada data. Adapun *reinforcement learning* menekankan pada interaksi agen dengan lingkungannya, di mana agen tersebut memperoleh pengalaman melalui umpan balik berupa penghargaan (*reward*) atau hukuman (*punishment*) guna mencapai tujuan tertentu secara optimal.

Paradigma *Machine Learning* terbagi menjadi *supervised learning* untuk tugas klasifikasi dan regresi dengan data berlabel, *unsupervised learning* untuk eksplorasi struktur data tanpa label, serta *reinforcement learning* untuk pengambilan keputusan berbasis *reward* [36]. Dalam penelitian ini, pendekatan *supervised learning* dipilih karena tersedia data historis lengkap beserta target *variable* (lama menginap). Data awal divalidasi, kemudian dilakukan *preprocessing* seperti penanganan *missing values*, *encoding categorical variables*, dan *feature scaling* sebelum dimasukkan ke dalam model regresi.

2.2.2 Feature Engineering

Feature engineering adalah proses transformasi data mentah menjadi fitur yang lebih informatif dan representatif untuk meningkatkan performa model *Machine Learning* [37]. Proses ini meliputi seleksi fitur, ekstraksi fitur, transformasi, dan pembuatan fitur baru yang dapat menangkap pola penting dalam data. *Feature engineering* sangat krusial karena kualitas fitur secara langsung mempengaruhi kemampuan model dalam belajar dan melakukan prediksi.

Penelitian terkini menekankan pentingnya automasi *feature engineering* dengan metode seperti *feature selection* berbasis statistik dan algoritma pembelajaran untuk mengurangi beban manual dan meningkatkan

efisiensi pengembangan model [38]. Dalam praktiknya, *feature engineering* dapat melibatkan teknik seperti *encoding* variabel kategorikal, normalisasi, pembuatan interaksi antar fitur, serta penggunaan *domain knowledge* untuk menambah fitur yang relevan.

2.2.3 Evaluasi Regresi

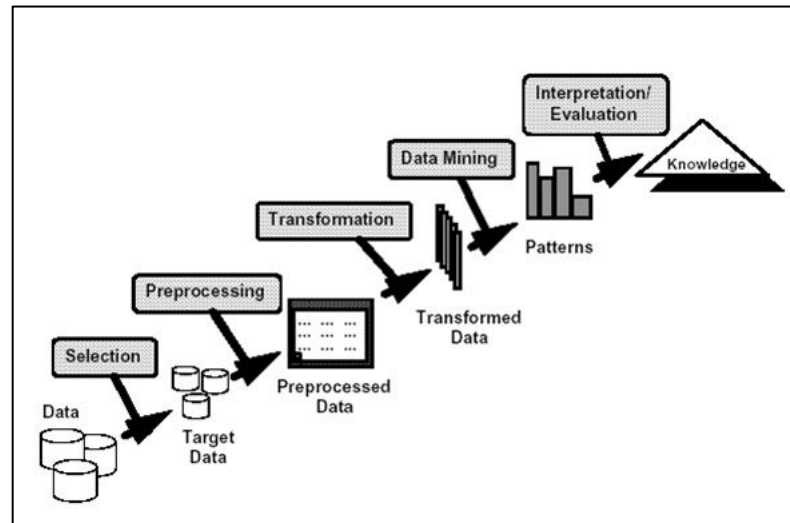
Evaluasi regresi adalah proses pengukuran kinerja model regresi dalam memprediksi nilai kontinu. Metrik evaluasi yang umum digunakan meliputi *Root Mean Squared Error* (RMSE), *Mean Absolute Error* (MAE), *Mean Absolute Percent Error* (MAPE), dan *Coefficient of Determination* (R^2). RMSE mengukur akar kuadrat rata-rata dari selisih kuadrat antara nilai prediksi dan aktual, memberikan penalti lebih besar pada *error* besar. MAE menghitung rata-rata nilai mutlak selisih, memberikan gambaran *error* rata-rata yang lebih stabil terhadap *outlier*. R^2 mengukur proporsi variansi data yang dapat dijelaskan oleh model, memberikan ukuran kecocokan model secara keseluruhan [39]. Pemilihan metrik evaluasi harus disesuaikan dengan tujuan dan karakteristik data. Kombinasi ketiga metrik ini memberikan gambaran komprehensif tentang akurasi, *robustness*, dan kemampuan generalisasi model. Selain itu, teknik validasi silang (*cross-validation*) dan uji statistik seperti *paired t-test* digunakan untuk memastikan keandalan dan signifikansi perbedaan performa antar model.

2.3 Framework dan Algoritma yang digunakan

2.3.1 Knowledge Discovery in Database (KDD)

Knowledge Discovery in Database (KDD) adalah proses yang digunakan untuk menemukan pola, informasi, dan pengetahuan yang bermanfaat dari kumpulan data besar. Tujuan utama KDD adalah mentransformasikan data mentah menjadi pengetahuan yang dapat digunakan untuk pengambilan Keputusan [40]. Proses ini terdiri dari beberapa tahapan yang sistematis, dimulai dari pemilihan data, *preprocessing*, transformasi, penambangan data (*data mining*), hingga interpretasi dan evaluasi hasil. KDD sangat penting dalam konteks analisis data besar dan *Machine Learning* karena

menyediakan kerangka kerja yang terstruktur guna mengolah dan mengekstrak informasi penting dari data yang tersebar dan kompleks.



Gambar 2. 1 Tahapan Proses Knowledge Discovery in Database (KDD)

Sumber: [41]

Gambar 2.1 menunjukkan tahapan-tahapan dalam metodologi *Knowledge Discovery in Database* (KDD) yang menggambarkan alur transformasi data dari bentuk awal hingga menghasilkan pengetahuan. Kerangka kerja KDD banyak digunakan dalam penelitian berbasis data *mining* dan *Machine Learning* karena kemampuannya dalam menyediakan struktur analisis yang jelas, sistematis, dan berorientasi pada eksplorasi pola dari data berskala besar. Dalam penelitian ini, *framework* KDD dipilih karena sesuai dengan karakteristik data transaksi hotel yang bersifat kompleks dan memiliki banyak variabel operasional. Penelitian diawali dengan tahap *business understanding* untuk memahami permasalahan operasional Hotel XYZ, proses pengolahan data selanjutnya mengikuti tahapan inti KDD sebagai kerangka teknis analisis data. Pada tahap akhir, hasil evaluasi difokuskan pada *business recommendation*, yaitu penerjemahan hasil analisis dan prediksi menjadi rekomendasi keputusan operasional yang relevan bagi manajemen hotel. Penerapan KDD dalam penelitian ini juga bertujuan untuk menghasilkan pengetahuan yang dapat mendukung pengambilan keputusan bisnis secara

terukur dan berbasis data, tanpa memasuki tahap implementasi sistem secara langsung. Berikut ini merupakan tahapan-tahapan dalam KDD:

1) *Selection*

Tahap *selection* merupakan proses awal dalam KDD yang berfokus pada pemilihan *subset* data yang relevan dari keseluruhan *dataset* yang tersedia. Pada tahap ini, data yang dianggap penting dan sesuai dengan tujuan analisis dipilih untuk diproses lebih lanjut. Pemilihan data ini sangat krusial karena akan menentukan kualitas dan fokus dari proses penemuan pengetahuan selanjutnya. Data yang tidak relevan atau berlebihan dapat diabaikan agar proses analisis menjadi lebih efisien dan hasilnya lebih akurat.

2) *Pre-processing*

Selanjutnya terdapat tahapan *preprocessing*, yang dimana tahap ini adalah tahap dimana data dilakukan persiapan dengan eksplorasi data, pengecekan tipe data dan pembersihan data yang bertujuan menghilangkan berbagai masalah seperti *noise*, nilai yang hilang, duplikasi, dan inkonsistensi dalam data. Proses ini memastikan bahwa data yang digunakan dalam analisis sudah dalam kondisi yang bersih dan konsisten sehingga tidak mengganggu hasil pemodelan. Selain itu, *preprocessing* juga dapat meliputi pengisian nilai yang hilang dan penghapusan data yang tidak valid, sehingga meningkatkan kualitas data secara keseluruhan.

3) *Transformation*

Pada tahap ini, penelitian ini melakukan data dinormalisasi untuk menyamakan skala variabel numerik, dilakukan rekayasa/pemilihan fitur, dan integrasi data sehingga dataset siap untuk pemodelan. Transformasi ini bertujuan mencegah dominasi fitur berkisar besar dan memperkenalkan variabel turunan yang menangkap dinamika hotel lebih baik.

4) *Data Mining*

Tahap ini dilakukannya pemodelan dengan membagi data menjadi set latih dan uji lalu menerapkan dua algoritma: *Random Forest* dan *Gradient Boosting*, dengan konfigurasi *hyperparameter* terkait yang

dikontrol di *Python*. Tujuan tahap ini adalah membangun model regresi untuk memprediksi LoS dan membandingkan performa kedua pendekatan model *Machine Learning*.

5) *Evaluation*

Pada tahap akhir, model dievaluasi menggunakan metrik regresi RMSE, MAE, MAPE, dan R^2 untuk menilai akurasi dan kemampuan penjelasan varians, selain metrik numerik, dilakukan analisis *feature importance* untuk mengidentifikasi variabel paling berpengaruh terhadap lama menginap tamu, lalu hasil teknis diinterpretasikan menjadi rekomendasi bisnis.

2.3.2 Algoritma Random Forest

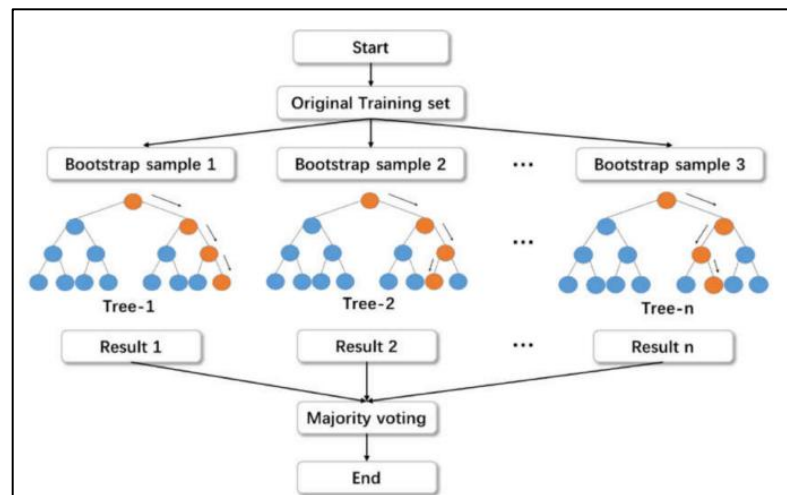
Random Forest adalah salah satu metode yang terdiri dari kumpulan *decision trees* yang dibangun secara acak (*bootstrap sampling*) dan independen, kemudian hasilnya digabungkan (*majority voting*) untuk klasifikasi atau *averaging* untuk regresi) untuk menghasilkan prediksi akhir yang lebih akurat dan stabil. Algoritma ini dikenal karena kemampuannya dalam menangani data dengan dimensi tinggi dan kompleksitas *non-linear*, serta ketahanannya terhadap *overfitting* [42]. Secara umum, *Random Forest* terdiri dari dua elemen utama: pembentukan pohon keputusan secara acak dan penggabungan hasilnya menjadi *output* akhir. Formulasi matematis dari *Random Forest* melibatkan proses *bootstrap sampling* yang menghasilkan *dataset subset* acak D_1, D_2, \dots, D_n , masing-masing digunakan untuk membangun pohon keputusan independen. Prediksi akhir γ pada model regresi dihitung sebagai rata-rata prediksi dari n pohon keputusan:

$$\gamma = \frac{1}{n} \sum_{i=1}^n f_i(x) \quad (2.1)$$

Berdasarkan rumus *Random Forest* di (2.1) menjelaskan bahwa hasil prediksi akhir diperoleh dari rata-rata semua prediksi yang dihasilkan oleh setiap pohon keputusan. Setiap pohon dibangun dari *subset* data acak yang berbeda melalui proses *bootstrap*. Setelah semua pohon menghasilkan

prediksinya masing-masing, seluruh nilai prediksi dijumlahkan lalu dibagi dengan jumlah pohon untuk menghasilkan nilai akhir yang lebih stabil dan akurat dibandingkan prediksi dari satu pohon saja.

Berikut ini merupakan struktur pelatihan algoritma *Random Forest*:



Gambar 2. 2 Struktur Algoritma *Random Forest*

Sumber: [43]

Gambar 2.2 di atas memperlihatkan proses pelatihan *Random Forest* yang dimulai dengan *bootstrap sampling* data, membangun pohon keputusan pada setiap sampel, dan menggabungkan hasil prediksi dari seluruh pohon dengan *majority voting*. Struktur ini memungkinkan *Random Forest* mengurangi varians model sekaligus mempertahankan bias rendah karena kombinasi pohon yang bervariasi. Dalam aplikasi seperti prediksi lama menginap tamu hotel, *Random Forest* dapat menangkap pola yang kompleks dan beragam dengan stabilitas yang baik meski data berisik atau tidak seimbang.

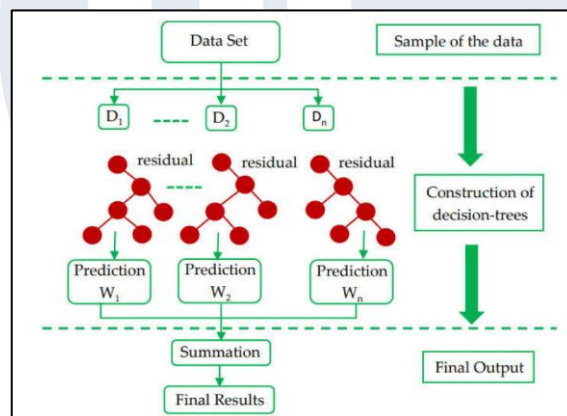
2.3.2 Algoritma Gradient Boosting

Gradient Boosting adalah teknik membangun model secara berurutan. Pada setiap iterasi, model baru ditujukan untuk memperbaiki kesalahan (*residual*) dari model sebelumnya dengan menggunakan pendekatan penurunan gradien pada fungsi *loss*. Metode ini menghasilkan model kombinasi yang kuat

dari sejumlah *weak learner* seperti pohon keputusan yang sederhana [44]. Rumus umum *Gradient Boosting* adalah:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x) \quad (2.2)$$

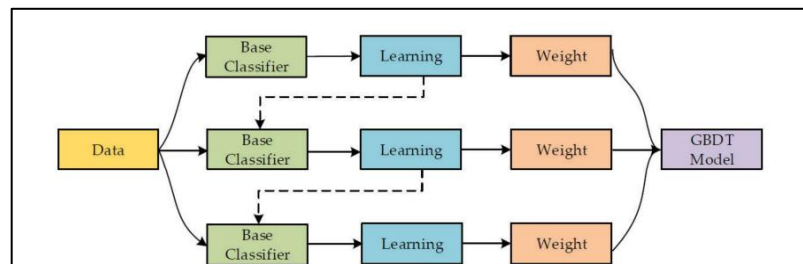
Rumus *Gradient Boosting* menunjukkan bahwa model pada iterasi ke- m merupakan hasil penambahan antara model sebelumnya dengan kontribusi model baru yang dikalikan oleh nilai *learning rate*. Model sebelumnya berfungsi sebagai dasar prediksi, sementara model baru (*weak learner*) berperan untuk memperbaiki kesalahan prediksi yang masih ada. *Learning rate* mengatur seberapa besar pengaruh model baru terhadap model akhir agar proses pembelajaran tetap stabil.



Gambar 2. 3 Struktur Algoritma *Gradient Boosting*

Sumber: [45]

Berdasarkan gambar 2.3 diatas, menggambarkan proses *algoritma Gradient Boosting* yang pembangunan model dilakukan secara berurutan berdasarkan *residual error*. Dataset awal dibagi menjadi beberapa *subset data* (D_1, D_2, \dots, D_n), di mana setiap *subset* digunakan untuk membangun *decision tree* yang fokus pada memprediksi residual dari model sebelumnya. Setiap pohon menghasilkan prediksi (W_1, W_2, \dots, W_n) yang kemudian digabungkan melalui proses *summation* untuk menghasilkan prediksi akhir (*final results*) [45]. Proses ini memungkinkan algoritma untuk secara progresif memperbaiki kesalahan prediksi dari model sebelumnya, sehingga menghasilkan model yang lebih akurat dan *robust* dalam menangkap pola kompleks dalam data.



Gambar 2. 4 Proses Training *Gradient Boosting*

Sumber: [46]

Selanjutnya, berdasarkan gambar 2.4 diatas mengilustrasikan terkait alur proses *training* pada algoritma *Gradient Boosting*, yang menunjukkan bagaimana data *input* diproses melalui serangkaian *base classifier* secara berurutan. Setiap *base classifier* menjalani tahap *learning* dan menghasilkan *weight* (bobot) yang menentukan kontribusinya terhadap model gabungan akhir. Garis putus-putus menunjukkan adanya *feedback mechanism* di mana *output* dari setiap *classifier* mempengaruhi *training classifier* berikutnya, memungkinkan setiap model baru untuk fokus pada kesalahan yang dibuat oleh model sebelumnya. Proses ini berlanjut secara iteratif hingga mencapai konvergensi atau kriteria penghentian yang ditetapkan, menghasilkan GBDT (*Gradient Boosting Decision Tree*) Model yang merupakan agregasi dari seluruh *base classifier* dengan pembobotan yang optimal untuk meminimalkan fungsi *loss* keseluruhan.

2.3.3 Evaluasi Model Prediksi

Evaluasi model prediksi dalam *Machine Learning* adalah langkah penting untuk mengetahui seberapa baik sebuah model dapat memperkirakan hasil pada data yang belum pernah ditemui sebelumnya. Pada prinsipnya, proses evaluasi ini menggunakan berbagai metrik, seperti *Root Mean Squared Error* (RMSE), *Mean Absolute Error* (MAE), *Mean Absolute Percent Error* (MAPE) dan *Coefficient of Determination* (R^2), yang semuanya memberikan sudut pandang berbeda mengenai kinerja model [47]. RMSE akan sangat menonjolkan *error* yang nilainya besar, MAE lebih menggambarkan rata-rata deviasi model terhadap data asli secara umum, dan R^2 memberikan gambaran umum tentang seberapa besar variasi data yang berhasil dijelaskan oleh model

[48]. Pada penelitian ini, proses evaluasi model prediksi akan menggunakan ketiga metrik tersebut agar hasil yang didapat benar-benar objektif, valid, dan mampu memberi gambaran menyeluruh tentang keandalan algoritma *Random Forest* maupun *Gradient Boosting*. Untuk memastikan evaluasi dilakukan secara objektif, sejumlah rumus statistik dipakai agar hasil pengujian model dapat dipertanggungjawabkan. Berikut ini merupakan ketiga rumus berdasarkan RMSE, MAE, MAPE, dan R²:

1) Rumus RMSE (*Root Mean Squared Error*):

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2} \quad (2.3)$$

Rumus ini menghitung akar dari rata-rata kuadrat *error* antara nilai aktual (y_i) dan prediksi model (\hat{y}_i), dengan n sebagai jumlah data, sehingga sangat peka terhadap *error* besar. Dengan tujuan rumus ini untuk menghitung akar dari rata-rata kuadrat *error* antara nilai aktual dan prediksi model, sehingga sangat peka terhadap *error* besar.

2) Rumus MAE (*Mean Absolute Error*):

$$MAE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i) \quad (2.4)$$

MAE menghitung nilai rata-rata dari selisih *absolut* antara nilai aktual (y_i) dan nilai prediksi (\hat{y}_i), serta n sebagai jumlah seluruh data, sehingga memberikan pengukuran *error* rata-rata yang stabil dan mudah diinterpretasikan. Dengan tujuan rumus MAE untuk menghitung rata-rata dari selisih *absolut* antara nilai aktual dan hasil prediksi.

3) Rumus MAPE (*Mean Absolute Percentage Error*):

$$MAPE = \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{\hat{y}_i} \right| \times 100\% \quad (2.5)$$

MAPE menghitung rata-rata persentase kesalahan *absolut* antara nilai aktual (y_i) dan nilai prediksi (\hat{y}_i), dengan n sebagai jumlah data. Rumus ini menampilkan seberapa besar kesalahan model dalam bentuk persentase terhadap nilai aktual. Tujuan dari rumus MAPE adalah untuk mengukur model dalam satuan persen.

4) Rumus R^2 (Koefisien Determinasi):

$$R = \frac{\text{Covar}(x, y)}{\sqrt{\text{Var}(x)\text{Var}(y)}} \quad (2.6)$$

Rumus R^2 (Koefisien Determinasi) menunjukkan hubungan antara nilai prediksi (x) dan nilai aktual (y). Nilai kovarians antara keduanya, yaitu $\text{covar}(x, y)$, dibandingkan dengan hasil akar dari perkalian varians masing-masing variabel, $\sqrt{\text{Var}(x)\text{Var}(y)}$, untuk menghasilkan koefisien korelasi *Pearson* (R). Nilai R kemudian dikuadratkan menjadi R^2 untuk menunjukkan seberapa besar variasi nilai. Semakin tinggi nilai R^2 , semakin baik model dalam pada hubungan antara variabel.

Berikut merupakan tabel kriteria evaluasi metrik regresi yang akan digunakan untuk menilai hasil model.

Tabel 2. 2 Kriteria Evaluasi Metrik Regresi

Metrik Evaluasi	Rentang Nilai	Kategori
RMSE	Semakin kecil semakin baik	Akurasi tinggi jika nilai mendekati 0
MAE	Semakin kecil semakin baik	<i>Error</i> rata-rata kecil menunjukkan prediksi stabil
MAPE	< 10%	Akurasi sangat baik
	10% – 20%	Akurasi baik
	20% – 50%	Akurasi moderat
	> 50%	Akurasi rendah / tidak dapat diandalkan
R^2	0.90 – 1.00	<i>Excellent fit</i>
	0.80 – 0.90	<i>Good fit</i>
	0.70 – 0.80	<i>Fair fit</i>
	0.50 – 0.70	<i>Poor fit</i>
	< 0.50	<i>Failure fit</i>

Tabel 2.2 menjelaskan kriteria penilaian untuk metrik regresi yang digunakan dalam evaluasi model. RMSE dan MAE menunjukkan tingkat kesalahan *absolut* antara nilai aktual dan prediksi, di mana nilai yang semakin kecil menandakan model bekerja dengan baik. MAPE memberikan gambaran persentase kesalahan, sehingga memudahkan interpretasi prediksi dalam konteks operasional; nilai MAPE di bawah 10% umumnya dianggap sangat baik. Sementara itu, R^2 digunakan untuk menilai seberapa besar variansi data yang dapat dijelaskan oleh model, dengan nilai mendekati 1 menunjukkan kecocokan model yang sangat baik.

2.4 Tools yang digunakan

2.4.1 Python

Python merupakan bahasa pemrograman tingkat tinggi yang berorientasi objek dan bersifat sumber terbuka, dikenal karena kemudahannya untuk dipahami serta fleksibilitasnya dalam berbagai bidang pengembangan [49]. Dengan manajemen memori otomatis dan sintaks yang sederhana, *python* memungkinkan pengembang untuk membangun berbagai jenis proyek mulai dari pengembangan situs, pengelolaan data, hingga pembuatan permainan. Selain itu, *python* juga memiliki ekosistem *library* yang kuat untuk berbagai bidang komputasi modern seperti *Big Data*, *Deep Learning*, *Data Mining*, *Machine Learning*, *Deep Visualization*, dan *Data Science* [50]. Keunggulan inilah yang menjadikan *python* sebagai salah satu bahasa pemrograman utama dalam pengembangan kecerdasan buatan dan aplikasi berbasis data.

2.4.1 Jupyter Notebook

Jupyter Notebook merupakan sebuah aplikasi berbasis *web* bersifat *open source* yang dirancang untuk memfasilitasi pembuatan serta penyebaran dokumen interaktif yang mengintegrasikan kode pemrograman, visualisasi data, perhitungan matematis, dan penjelasan naratif dalam satu kesatuan [51]. Aplikasi ini banyak dimanfaatkan dalam berbagai kegiatan analisis data seperti pembersihan dan transformasi data, simulasi numerik, pemodelan statistik, serta pengembangan model pembelajaran mesin [52]. Dan *Jupyter Notebook*

mendukung penggunaan berbagai pustaka dasar *python* seperti *numpy*, *pandas*, *matplotlib*, dan *seaborn* yang berfungsi untuk pengolahan data, analisis statistik, serta visualisasi informasi secara efisien.

2.4.1 MariaDB

MariaDB merupakan sistem manajemen basis data relasional bersifat *open-source* yang dikembangkan sebagai turunan (*fork*) dari *MySQL* dengan tujuan mempertahankan kompatibilitas sekaligus meningkatkan kinerja dan keandalannya [53]. Sistem ini mendukung berbagai *storage engine* seperti *Aria*, *InnoDB*, dan *ColumnStore* yang memungkinkan pengelolaan data berskala besar dengan efisien. *MariaDB* juga dilengkapi fitur *replication*, *clustering*, serta optimasi *query* yang mendukung kebutuhan analisis dan pengolahan data secara cepat dan stabil [54]. Keunggulan-keunggulan ini menjadikan *MariaDB* sebagai salah satu pilihan utama dalam pengembangan aplikasi riset dan sistem informasi modern.

