

PERBANDINGAN ALGORITMA MACHINE LEARNING UNTUK ANALISIS SENTIMEN PADA ULASAN HOTEL

Viktor Handrianus Pranatawijaya ^{a,1,*}, Efrans Christian ^{b,2}

^{ab} Universitas Palangka Raya, Kampus Tunjung Nyaho Jalan Yos Sudarso, Palangka Raya, Kalimantan Tengah, Indonesia

¹ viktorhp@it.upr.ac.id*; ² efrans@it.upr.ac.id; ³ Email penulis ketiga (9pt)

* corresponding author

ARTICLE INFO

ABSTRACT

Keywords

Machine Learning, Gradient Boosting, Support Vector Machine (SVM), K-Nearest Neighbor (KNN)

The paper extensively explores machine learning algorithms for evaluating sentiments in hotel reviews, particularly within the tourism and hospitality industry. It underscores the importance of precise reviews in utilizing artificial intelligence for improved operational efficiency, revenue optimization, and heightened customer satisfaction. Notably, supervised machine learning algorithms like Gradient Boosting, Support Vector Machine, and K-Nearest Neighbor are highlighted for offering recommendations based on reviews to predict user preferences. The research methodology involves data scraping, cleaning, preprocessing, and labeling, followed by training and testing the chosen machine learning algorithms. Results indicate that the Support Vector Machine algorithm demonstrated superior performance with accuracy 0.8553, precision 0.8433, recall 0.8553, dan F1-score 0.8424, suggesting its appropriateness for sentiment analysis in hotel reviews. The paper concludes by recommending the implementation of the Support Vector Machine model for sentiment analysis in hotel reviews in Palangka Raya, Indonesia, and proposes avenues for further industry development and enhancement.

1. Pendahuluan

Industri pariwisata dan perhotelan, yang menarik jutaan orang di seluruh dunia, telah berkembang menjadi salah satu subsektor yang sangat penting dalam industri jasa. Menyusun dan merencanakan perjalanan wisata melibatkan banyak pertimbangan, terutama dalam memilih destinasi yang menarik dan memilih akomodasi yang tepat. Komentar wisatawan juga sangat penting untuk memilih akomodasi terbaik [1].

Keakuratan penilaian yang dibuat melalui ulasan ini sangat penting karena berdampak pada keputusan yang dibuat oleh wisatawan. Industri pariwisata dan perhotelan menggunakan kecerdasan buatan untuk mengatasi kompleksitas ini dengan tujuan meningkatkan efisiensi operasional, mengoptimalkan pendapatan, dan memberikan pengalaman wisatawan yang lebih memuaskan [2]. Namun, masalah dengan menilai kecerdasan buatan dengan benar masih menjadi masalah besar [3].

Salah satu yang dapat digunakan dengan *supervised Machine Learning* (SML) adalah untuk memberikan rekomendasi kepada pengguna berdasarkan ulasan. Hal ini mencakup pengembangan ide untuk memilih algoritma terbaik setelah pelatihan dan pengujian. Dalam hal ini, algoritma seperti *Gradient Boosting* (GB), *K-Nearest Neighbor* (KNN), dan *Support Vector Machine* (SVM) sangat membantu mengantisipasi preferensi pengguna terhadap ulasan.

KNN, yang sering digunakan dalam masalah klasifikasi dan regresi, adalah salah satu algoritma yang sederhana namun berhasil [4]. Sangat populer untuk mendukung keputusan berbasis opini karena kemampuan untuk memprediksi berdasarkan data terdekat. Selain itu, terbukti bahwa SVM [5]. membantu pengguna melakukan transaksi informasi, terutama dalam hal rekomendasi berbasis ulasan.

Dalam teori pembelajaran, algoritma GM menghasilkan pengklasifikasi kombinasi yang kuat dengan menggabungkan sekelompok pengklasifikasi yang kurang akurat dengan pohon keputusan [6].

GM adalah pilihan yang bagus untuk meningkatkan akurasi dalam hal rekomendasi berbasis ulasan karena keunggulan ini.

Penerapan algoritma ini sangat terasa, terutama dalam meningkatkan kemampuan pengklasifikasi yang tidak begitu tepat, seperti pohon keputusan. Dengan demikian, kecerdasan buatan dalam SML membantu meningkatkan kualitas rekomendasi yang diberikan kepada pengguna berdasarkan ulasan, sehingga pengalaman pengguna menjadi lebih personal dan memuaskan. Ada kemungkinan bahwa implementasi yang cermat dari algoritma-algoritma ini akan berdampak positif pada industri, tergantung pada pendapat dan ulasan pengguna.

Dalam meningkatkan pembobotan dari sisi pengguna pada analisis ulasan hotel, analisis sentimen menjadi elemen kunci untuk memahami nuansa dan pendapat yang terkandung dalam ulasan tersebut. Beberapa penelitian terkait, seperti yang dijelaskan dalam literatur [7–10], memberikan kontribusi berharga terkait metode analisis sentimen yang dapat diadopsi.

Penelitian [7] mengenai analisis sentimen pada ulasan pengguna hotel menyajikan pendekatan yang dapat diintegrasikan untuk memahami apakah pendapat yang disampaikan cenderung positif, negatif, atau netral. Selain itu, penelitian [8] menyoroti teknik-teknik pemrosesan bahasa alami yang efektif dalam mengekstrak sentimen dari teks ulasan, memungkinkan identifikasi aspek-aspek yang paling memengaruhi persepsi pengguna terhadap sebuah hotel. Implementasi metode tersebut dijadikan dasar untuk meningkatkan pembobotan dari sudut pandang pengguna.

Selanjutnya, penelitian [9] menitikberatkan pada penggunaan teknologi deep learning dalam analisis sentimen. Pendekatan ini membuka peluang untuk lebih mendalam dan akurat dalam menangkap kompleksitas sentimen yang mungkin tersembunyi dalam ulasan pengguna. Begitu pula, penelitian [10] mengeksplorasi penggunaan metode-metode yang inovatif dalam analisis sentimen, memberikan perspektif tambahan dalam meningkatkan ketepatan pembobotan.

Untuk menyelesaikan masalah ini, kita perlu mengambil pendekatan yang cermat dan solusi kreatif untuk menerapkan kecerdasan buatan di industri pariwisata dan perhotelan. Dengan demikian, untuk memastikan kemajuan industri dan memberikan pengalaman wisata yang lebih memuaskan bagi seluruh pengunjung, upaya untuk meningkatkan keakuratan penilaian ulasan sangat penting.

2. Metodologi Penelitian

Langkah awal dalam penelitian ini adalah melakukan pengumpulan data dari situs web www.agoda.com yang berfokus pada informasi hotel di kota Palangka Raya, termasuk ulasan pengguna. Proses scraping data dilakukan untuk mendapatkan dataset yang komprehensif, mencakup informasi penting seperti atribut hotel dan ulasan yang diberikan oleh pengguna. Langkah ini menjadi krusial karena dataset yang baik merupakan fondasi utama dalam melaksanakan pelatihan dan pengujian terhadap algoritma pada SML.

Setelah pengumpulan data, langkah berikutnya dalam penelitian ini adalah melakukan serangkaian proses untuk memastikan bahwa dataset yang diperoleh bersih, terstruktur, dan siap untuk diolah menggunakan algoritma supervised machine learning. Langkah-langkah ini mencakup *data cleaning*, *preprocessing*, dan *labelling* pada ulasan.

Pertama, *data cleaning* melibatkan identifikasi dan penanganan data yang tidak lengkap, duplikat, atau tidak konsisten. Ini mencakup penghapusan entri data yang tidak relevan atau tidak lengkap, penanganan nilai yang hilang, dan pengidentifikasian serta penanganan duplikasi data. Proses ini penting untuk memastikan integritas dan kualitas dataset sebelum dilibatkan dalam analisis lebih lanjut.

Selanjutnya, *preprocessing* melibatkan serangkaian langkah untuk menyiapkan data agar sesuai dengan kebutuhan algoritma machine learning. Ini termasuk normalisasi data, konversi format, dan pemilihan atribut yang relevan. Proses ini bertujuan untuk meningkatkan efisiensi dan performa algoritma, serta memastikan konsistensi dan keakuratan analisis.

Setelah itu, langkah *labelling* melibatkan pengkategorian sentimen pada ulasan pengguna sebagai variabel target. Proses ini melibatkan penentuan label atau kategori yang sesuai dengan setiap ulasan, seperti positif, negatif, atau netral. Labeling ini menjadi dasar bagi algoritma supervised machine learning dalam memahami dan memprediksi sentimen pada ulasan.

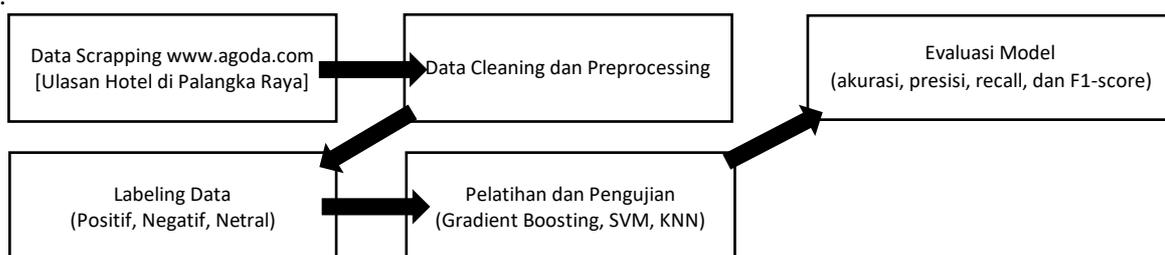
Dengan menyelesaikan langkah-langkah ini, dataset yang sudah dibersihkan, diproses, dan diberi label dapat digunakan untuk melanjutkan ke tahap pelatihan dan pengujian menggunakan algoritma machine learning yang telah dipilih. Proses ini memastikan bahwa analisis sentimen pada ulasan hotel dapat dilakukan dengan akurat dan relevan, dengan hasil yang dapat digunakan untuk meningkatkan pemahaman tentang preferensi pengguna dan kualitas pelayanan hotel di kota Palangka Raya.

Setelah mendapatkan dataset yang cukup representatif, penelitian dilanjutkan dengan pelatihan dan pengujian menggunakan tiga algoritma supervised ML utama, yaitu GB, KNN, dan SVM. Pada tahap ini, penelitian merujuk pada metodologi yang telah diterapkan dalam penelitian sebelumnya [11] [12], dengan fokus pada aspek klasifikasi dan regresi. Pelatihan dan pengujian tersebut melibatkan penerapan algoritma-algoritma tersebut pada dataset hotel dan ulasan yang telah diperoleh, mengamati kemampuan mereka dalam mengklasifikasikan dan meramalkan berbagai variabel yang relevan.

Hasil pelatihan dan pengujian menjadi penentu algoritma terbaik untuk diterapkan pada dataset khusus ini. Assessment melibatkan penilaian akurasi, presisi, recall, dan metrik evaluasi lainnya yang sesuai dengan konteks penelitian. Algoritma yang memberikan performa tertinggi dari segi klasifikasi dan regresi pada dataset hotel dan ulasan kota Palangka Raya dipilih sebagai yang terbaik. Kesimpulan dari penelitian ini memberikan arah yang jelas untuk implementasi algoritma supervised machine learning yang paling efektif dalam mengoptimalkan analisis ulasan pengguna terkait hotel di kota tersebut.

Dengan merujuk pada kontribusi-kontribusi ini, penelitian dapat memperkaya analisis sentimen pada ulasan pengguna hotel di kota Palangka Raya. Dengan mengadopsi dan menggabungkan elemen-elemen yang efektif dari penelitian terdahulu, diharapkan dapat tercipta sistem pembobotan yang lebih canggih dan sensitif terhadap nuansa sentimen dalam ulasan. Pendekatan ini diharapkan mampu memberikan pemahaman yang lebih mendalam terhadap preferensi pengguna, dengan hasil pembobotan yang lebih akurat dan relevan dalam meningkatkan kualitas analisis ulasan hotel.

Berikut ini merupakan gambar yang menjelaskan langkah-langkah yang dilakukan dalam penelitian ini.



Gambar 1. Alur Metodologi Penelitian

3. Hasil dan Pembahasan

Bagian ini akan membahas hasil dan pembahasan penelitian berdasarkan tahapan metodologi penelitian yang telah dilakukan, mulai dari pengumpulan data melalui scrapping ulasan hotel di Palangka Raya, proses data cleaning dan preprocessing, hingga tahapan labelling, pelatihan, dan pengujian algoritma Gradient Boosting (GB), Support Vector Machine (SVM), dan K-Nearest Neighbor (KNN). Penjelasan berikut memberikan gambaran menyeluruh tentang temuan dan kesimpulan yang dapat diambil dari setiap langkah metodologi tersebut.

A. 3.1. Data

Scraping dari www.agoda.com untuk ulasan hotel di Palangka Raya

Proses pertama dalam penelitian ini melibatkan pengumpulan data dari situs web www.agoda.com yang difokuskan pada ulasan hotel di kota Palangka Raya. Untuk mencapai ini, alat yang digunakan adalah Web Scraper dari www.webscraper.io/. Melalui alat ini, data yang terkumpul mencakup informasi detail mengenai berbagai hotel dan ulasan pengguna yang terkait. Keseluruhan, setelah proses *scrapping* selesai, berhasil diperoleh sebanyak 376 data yang akan digunakan untuk analisis lebih lanjut.

Gambar 2 memberikan gambaran contoh dari dataset yang dikumpulkan melalui proses *scrapping* ini. Data yang terdapat dalam tabel mencakup beragam informasi, termasuk atribut-atribut hotel seperti nama hotel, rating, identitas konsumen, dan ulasan, bersama dengan ulasan-ulasan pengguna yang mencakup sentimen terhadap layanan dan pengalaman mereka di hotel tersebut. Data yang kaya dan beragam ini menjadi dasar untuk penelitian analisis sentimen selanjutnya.

nama_hotel	nilai	nilaikata2	nama	negara	identitas	kanar	kesan	ulasan	tanggal	review_english
0	Rungan Sari	8.4	Excellent	Tisha	Indonesia	Couple	Standard Double or Twin	Escaping the cityâ€¦ If you look for a staycation close to forest, ...	Reviewed December 22, 2021	If you are looking for a staycation close to l...
1	Rungan Sari	7.6	Very good	Tisha	Indonesia	Couple	Standard Double or Twin	A calming placeâ€¦ If you look for a short getaway from city hust...	Reviewed July 01, 2019	If you are looking for a short getaway from th...
2	Rungan Sari	9.2	Exceptional	Nigel	United Kingdom	Couple	Standard Double or Twin	Great unpretentious and great value hotelâ€¦ We have stayed here now on two occasions. A lo...	Reviewed November 06, 2018	We have stayed here now on two occasions. A lo...
3	Rungan Sari	9.6	Exceptional	Merv	New Zealand	Couple	Standard Double or Twin	Great valueâ€¦ Resort was in great location, fantastic staff...	Reviewed July 23, 2016	Resort was in great location, fantastic staff...
4	Rungan Sari	9.7	Exceptional	Erwin	Indonesia	Business traveler	Standard Double or Twin	good for vacation tripâ€¦ I love very much	Reviewed September 11, 2015	I love very much
...
371	Best Western	10.0	Exceptional	Nurbaya	Indonesia	Group	1 Double Bed, Smoking, Deluxe	Okâ€¦ Ok	Reviewed January 08, 2023	OK
372	Best Western	7.2	Very good	Emi	Indonesia	Business traveler	1 Double Bed Non-Smoking, Deluxe	Very goodâ€¦ Tempatnya cukup bersih dan sarapannya bagus say...	Reviewed November 21, 2022	The place is quite clean and the breakfast is ...
373	Best Western	10.0	Exceptional	Vhelinda	Indonesia	Couple	1 Double Bed Non-Smoking, Deluxe	Exceptionalâ€¦ Tempat menginap yang sangat rekomen dan wajib ...	Reviewed October 26, 2022	A highly recommended place to stay and a must ...
374	Best Western	9.2	Exceptional	Yovanda	Indonesia	Business traveler	1 Double Bed Non-Smoking, Deluxe	Hotel okee sihâ€¦ Tempatnya bersih dan lokasi strategis	Reviewed October 19, 2022	The place is clean and strategic location
375	Best Western	10.0	Exceptional	Difa	Indonesia	Couple	1 Double Bed, Smoking, Deluxe	Enak banget â€¦ Masakannya enak, bersih dan nyaman banget	Reviewed July 21, 2022	The food is delicious, clean and very comfortable

376 rows x 11 columns

Gambar 2. Dataset Ulasan Hotel

3.2. Proses Data Cleaning dan Preprocessing

Setelah pengumpulan data, dilakukan proses data cleaning untuk menangani data yang tidak lengkap, duplikat, atau tidak konsisten. Langkah ini memberikan kebersihan pada dataset, memastikan integritas data. Selanjutnya, proses preprocessing dilakukan untuk menormalkan data, mengubah format, dan memilih atribut yang relevan. Hal ini bertujuan untuk meningkatkan kualitas data sebelum dilibatkan dalam pembelajaran mesin. Gambar 3 menunjukkan hasil dari proses ini.

sentimen dari ulasan hotel, memungkinkan evaluasi yang lebih terarah. Proses *labelling* menggunakan alence Aware Dictionary and sEntiment Reasoner (VADER). Untuk lebih jelasnya dapat dilihat pada gambar 4. Distribusi yang seimbang pun dilakukan berdasarkan kolom *sentiment_label* dengan *oversampling*.

stemmed_tokens	sentiment_vader	sentiment_positive	sentiment_negative	sentiment_neutral	sentiment_compound	sentiment_label	sentiment_score
[look, staycat, close, forest, right, place]	{'neg': 0.0, 'neu': 1.0, 'pos': 0.0, 'compound': ...}	0.000	0.000	1.000	0.0000	neutral	2
[look, short, getaway, citi, hustle, bustle, wou...]	{'neg': 0.092, 'neu': 0.787, 'pos': 0.121, 'co...}	0.121	0.092	0.787	0.2263	positive	1
[stay, two, occas, long, way, palangkaraya, pe...]	{'neg': 0.0, 'neu': 0.656, 'pos': 0.344, 'comp...}	0.344	0.000	0.656	0.9403	positive	1
[resort, great, local, fantast, staff, good, p...]	{'neg': 0.0, 'neu': 0.481, 'pos': 0.519, 'comp...}	0.519	0.000	0.481	0.9349	positive	1
[love, much]	{'neg': 0.0, 'neu': 0.192, 'pos': 0.808, 'comp...}	0.808	0.000	0.192	0.6369	positive	1

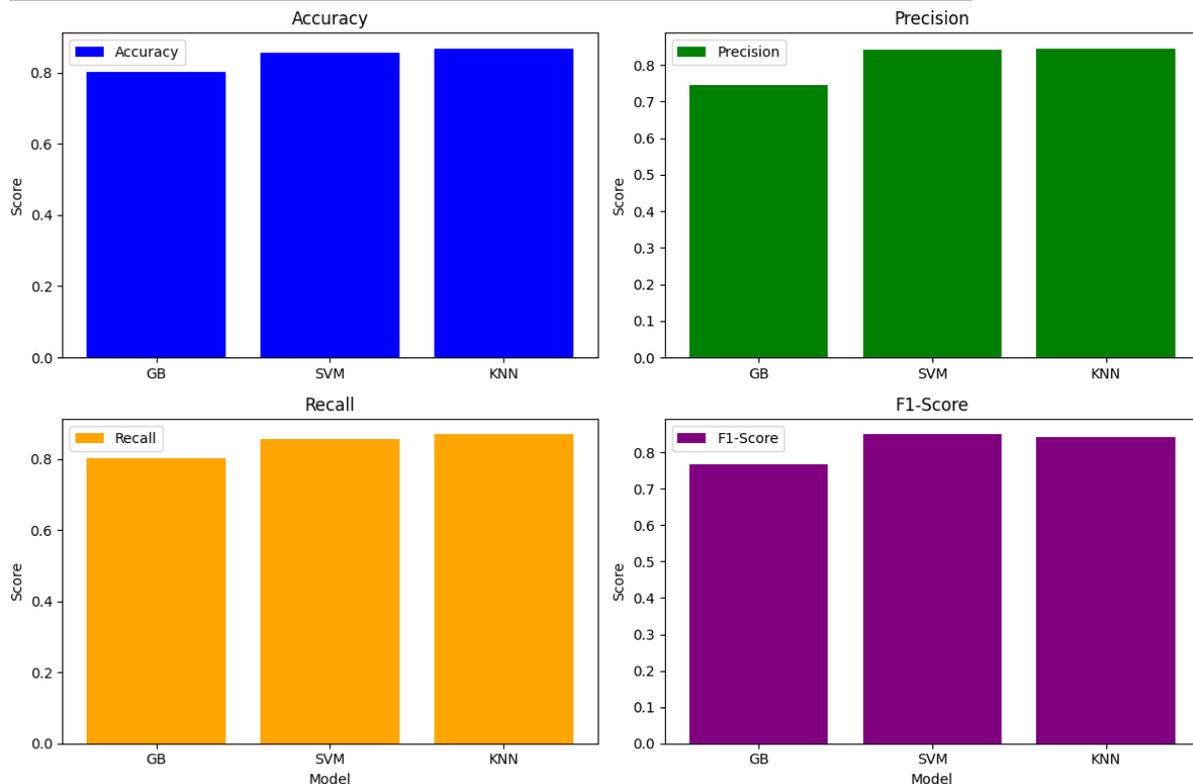
Gambar 5. Data Labelling

3.5 Pelatihan dan Pengujian Algoritma GB, SVM, dan KNN:

Tahap selanjutnya adalah pelatihan dan pengujian menggunakan tiga algoritma utama, yaitu GB, SVM, dan KNN. Melalui dataset yang sudah dipersiapkan, model-model ini dilatih untuk mengenali pola dan hubungan dalam data serta diuji pada dataset yang belum pernah dilihat sebelumnya. Langkah awalnya dengan membagi dataset menjadi data pelatihan dan pengujian. Selain itu melakukan *feature engineering* dengan menggunakan Term Frequency-Inverse Document Frequency (TF-IDF). Setelah itu baru dilatih dan diuji modelnya.

3.6 Evaluasi Model

Hasil pengujian diukur menggunakan metrik evaluasi seperti akurasi, presisi, recall, dan F1-score. Evaluasi ini memberikan gambaran sejauh mana setiap algoritma dapat mengklasifikasikan sentimen ulasan hotel dengan akurat. Hasil evaluasi menjadi dasar untuk mengidentifikasi algoritma yang paling efektif dalam konteks analisis sentimen ulasan hotel di Palangka Raya. Berdasarkan hasil uji pada model GB, Accuracy: 0.8026, Precision: 0.7445, Recall: 0.8026, dan F1-score: 0.7673. Pada SVM hasil ujinya adalah Accuracy: 0.8553, Precision: 0.8433, Recall: 0.8553, dan F1-score: 0.8492. Pada KNN pengujiannya menghasilkan Accuracy: 0.8684, Precision: 0.8459, Recall: 0.8684, dan F1-score: 0.8424. Disini dapat dilihat bahwa model SVM yang memiliki akurasi tertinggi. Untuk lebih jelasnya dapat dilihat pada gambar 6 yaitu Grafik Evaluasi dari ketiga model yang dilatih dan diuji.



Gambar 6. Grafik Evaluasi Model

4. Kesimpulan

Proses penelitian ini berhasil mengumpulkan data ulasan hotel di Palangka Raya dari situs web www.agoda.com menggunakan Web Scraper dari www.webscraper.io/. Melalui langkah-langkah data *cleaning*, *preprocessing*, dan *labelling* sentimen dengan VADER, dataset berhasil disiapkan untuk dilibatkan dalam pembelajaran mesin. Dari hasil evaluasi model, algoritma Support Vector Machine (SVM) menunjukkan kinerja terbaik dalam mengklasifikasikan sentimen ulasan hotel, dengan akurasi 0.8553, presisi 0.8433, recall 0.8553, dan F1-score 0.8424 yang lebih unggul dibandingkan dengan GB dan K-Nearest Neighbor KNN.

Untuk pengembangan lebih lanjut, disarankan untuk mengimplementasikan model SVM dalam pengelolaan dan analisis sentimen ulasan hotel di Palangka Raya. Dalam upaya mempertahankan ketepatan model, perlu dilakukan pembaruan dataset secara berkala dan evaluasi performa model. Selain itu, eksplorasi terhadap aspek pada ulasan yang mempengaruhi sentimen pengguna dapat memberikan wawasan tambahan. Meningkatkan interpretabilitas model dan mempertimbangkan integrasi feedback pelanggan juga dapat meningkatkan relevansi dan keefektifan analisis sentimen di industri perhotelan.

Daftar Pustaka

- [1] Lim KH, Chan J, Karunasekera S, Leckie C. Tour recommendation and trip planning using location-based social media: a survey. *Knowl Inf Syst.* 2019 Sep 1;60(3):1247–75.
- [2] Krishnan K. Travel and tourism industry applications and usage. In: *Building Big Data Applications.* Elsevier; 2020. p. 145–55.
- [3] Baizal ZKA, Tarwidi D, Adiwijaya, Wijaya B. Tourism Destination Recommendation Using Ontology-based Conversational Recommender System. *International Journal of Computing and Digital Systems.* 2021;10(1):829–38.

- [4] Al-Ghobari M, Muneer A, Fati SM. Location-aware personalized traveler recommender system (lapta) using collaborative filtering knn. *Computers, Materials and Continua*. 2021;69(2):1553–70.
- [5] Chidambarathanu K, Shunmuganathan KL. Predicting user preferences on changing trends and innovations using SVM based sentiment analysis. *Cluster Comput*. 2019 Sep 1;22:11877–81.
- [6] Lin X, Zhang X, Xu X. Efficient Classification of Hot Spots and Hub Protein Interfaces by Recursive Feature Elimination and Gradient Boosting. *IEEE/ACM Trans Comput Biol Bioinform*. 2020 Sep 1;17(5):1525–34.
- [7] Tey FJ, Wu TY, Lin CL, Chen JL. Accuracy improvements for cold-start recommendation problem using indirect relations in social networks. *J Big Data*. 2021 Dec 1;8(1).
- [8] Shokeen J, Rana C. A study on features of social recommender systems. *Artif Intell Rev*. 2020 Feb 1;53(2):965–88.
- [9] Osman NA, Noah SAM, Darwich M, Mohd M. Integrating contextual sentiment analysis in collaborative recommender systems. *PLoS One*. 2021 Mar 1;16(3 March).
- [10] Ziani A, Azizi N, Schwab D, Aldwairi M, Chekkai N, Zenakhra D, et al. Recommender System Through Sentiment Analysis [Internet]. 2017. Available from: <https://hal.science/hal-01683511>
- [11] Kuanr M, Mohapatra P. Assessment Methods for Evaluation of Recommender Systems: A Survey. *Foundations of Computing and Decision Sciences*. 2021 Dec 1;46(4):393–421.
- [12] Zheng X, Luo Y, Sun L, Zhang J, Chen F. A tourism destination recommender system using users' sentiment and temporal dynamics. *J Intell Inf Syst*. 2018 Dec 1;51(3):557–78.