

Pengamatan Tren Ulasan Hotel Menggunakan Pemodelan Topik Berbasis *Latent Dirichlet Allocation*

Suparyati ^{1*}, Emma Utami ^{2**} Agus Fathurahman ^{3***}

* Magister Teknik Informatika, Universitas Amikom Yogyakarta

** Magister Teknik Informatika, Universitas Amikom Yogyakarta

*** Magister Teknik Informatika, Universitas Amikom Yogyakarta

suparyati@students.amikom.ac.id¹, ema.u@amikom.ac.id², agusagusfatkhurohman@amikom.ac.id³

Article Info

Article history:

Received 2022-01-01

Revised 2022-05-13

Accepted 2022-05-18

Keyword:

*Latent Dirichlet Allocation,
Topic Modeling,
Machine Learning.*

ABSTRACT

The accuracy in extracting and summarizing thousands of reviews into several topics is the key in the implementation of data processing and further information. The hotel industry is no exception, where a review is an asset which, when processed, can produce information that will later be used for business expansion and business continuity. This hotel review topic modeling research uses Latent Dirichlet Allocation as a means to summarize the document. Latent Dirichlet Allocation is proven to be effective in the processing of summarizing words and many studies have used this method. The purpose of this research is to get a summary of words that make up a topic that represents the whole review which can produce data for hotel management to maintain their existence in the business and expand by considering the results of modeling the topic. The results showed that the words location, service, hotel, breakfast, resort and beach were the terms that most often appeared among the dominant topics.



This is an open access article under the [CC-BY-SA](https://creativecommons.org/licenses/by-sa/4.0/) license.

I. PENDAHULUAN

Kesulitan dalam pengubahan data tidak terstruktur menjadi data terstruktur merupakan tantangan tersendiri dewasa ini. Informasi dalam bentuk teks yang ada saat ini umumnya dalam bentuk teks yang tidak berlabel yang tidak dapat kita atribusikan secara ketat ke domain tematik berjumlah sangat banyak. Manajemen kategorisasi manual terhadap dokumen teks memerlukan sumber daya keuangan dan sumber daya manusia yang besar sehingga diperlukan permodelan topik untuk pengklasifikasian dokumen teks tersebut.

Penerapan *Latent Dirichlet Allocation* (LDA) dan *Latent Semantic Index* (LSA), yang mampu menangani data diskrit. Selain itu, dilakukan perbandingan untuk menguji divergensi, throughput, kualitas, dan waktu respons, karena keduanya dapat mengklasifikasikan data berdasarkan konten dan dengan memberi label pada masing-masing kategori. Algoritma dengan divergensi yang lebih baik diimplementasikan yang dapat menangani persyaratan organisasi dengan menghadirkan area teratas yang perlu ditingkatkan/konsentrasi tergantung pada analitik yang dibuat

oleh algoritma pada data diskrit yang tersedia, dan dengan menerapkan teknik visualisasi, hasilnya akan ditampilkan secara merata dalam format grafis [1]. Studi perbandingan pada klasifikasi dokumen teks ilmiah tidak terstruktur berdasarkan teks lengkap yang paling banyak diterapkan pendekatan pemodelan topik populer (LDA, LSA) untuk mengelompokkan kata-kata menjadi satu set topik sebagai kata kunci penting untuk klasifikasi [2].

Metode yang diusulkan membantu dalam mengonversi data yang tidak dapat dibaca ke format terstruktur yang dapat dibaca dengan bantuan Pembelajaran mesin adalah klasifikasi, dan pengelompokan memainkan peran penting dalam mengubah data operasional menjadi model data dan memvisualisasikan informasi yang diproses ke pengguna akhir. Penelitian ini bertujuan mendapatkan gambaran dan ringkasan yang berupa kata-kata yang mengarah pada suatu topik yang dapat membantu memberikan pertimbangan-pertimbangan terhadap pihak manajemen hotel untuk memberikan suatu kebijakan dan keputusan dalam hal memajukan bisnis dan ekspansi usahanya.

Topic Modeling merupakan tugas *natural language processing* yang mengekstrak topik yang relevan dari

dokumen tekstual. *Topic modelling* yang juga teks *mining* yang menggunakan teknik *supervised* dan *unsupervised machine learning* yang digunakan dalam mengidentifikasi pola dalam korpus yang mengandung sejumlah besar teks yang struktur dan tidak terstruktur dengan cara mengelompokkan istilah korpus ke dalam kelompok istilah sehingga menghasilkan topik menggunakan pemrosesan kesamaan. Topik modeling sendiri memiliki masalah karena dapat menghasilkan topik yang tidak dapat diandalkan yang dapat menyesatkan konsumen berita, tetapi topic modeling tetap dapat menjadi metode yang ampuh untuk menganalisis dan mengelompokkan artikel secara efisien dalam skala besar oleh organisasi secara internal [3].

Permodelan topik dapat dilakukan dengan beberapa teknik seperti keluhan pelanggan berbasis teks online yang dianalisis dengan *Latent Dirichlet Allocation* (LDA), GenSim LDA, Mallet LDA dan Gibbs *Sampling* untuk model Dirichlet Multinomial Mixture (GSDMM)[4]. Algoritma *Latent Dirichlet Allocation* (LDA) diterapkan untuk mengekstraksi topik-topik penting seperti penjualan dan konsumsi alkohol, tinggal di rumah, penelusuran statistik harian, kebrutalan polisi, teori konspirasi 5G dan vaksin menjadi beberapa topik yang dibahas dan seputar sikap dan persepsi yang dibentuk oleh warga [5]. Dalam analisis sentimen dan *topic modeling* pada data Twitter yang berisi “#IndiaFightsCorona” untuk menganalisis opini publik, LDA berkinerja lebih baik untuk sentimen positif, dan untuk sentimen negatif, LSA berkinerja lebih baik. Eksperimen dilakukan pada kumpulan data berbasis Twitter yang dibuat menggunakan twit dengan kata kunci sungai Cauvery, tagihan Lokpal dan Rahul Gandhi menggunakan metode permodelan topik LDA, LSI dan NMF yang mengevaluasi keakuratan topik yang dibentuk dengan menggunakan langkah-langkah *confusion*, kemungkinan log, dan koherensi topik dimana topik terbaik yang terbentuk kemudian diumpankan ke model regresi Logistik dan diketahui model yang dibuat menunjukkan akurasi yang lebih baik dengan LDA [6]. Model LDA menghasilkan ringkasan yang singkat, jelas, dan koheren [7]. LDA ditemukan sebagai metode komputasi yang efisien dan dapat ditafsirkan dalam mengadopsi bahasa Inggris Kumpulan data Alkitab Versi

Internasional Baru yang belum dibuat daripada metode LSA [8]. Dalam sentiment analisis mengenai data twitter tentang kendaraan listrik menunjukkan bahwa LDA memberikan yang lebih baik wawasan tentang topik, serta akurasi yang lebih baik daripada LSA [9]. Penggunaan algoritma *Latent Semantic Index* (LSA) dalam peringkasan dokumen mampu mencapai skor lebih besar dari peringkasan teks menggunakan pemodelan topik *Latent Dirichlet Allocation* [10]. Sedangkan dalam pencarian kembali informasi terutama untuk Pengambilan Informasi Ad-hoc seperti mengklasifikasikan dokumen dan memodelkan hubungannya antara berbagai topik menggunakan LDA. Studi perbandingan empiris antara dua pendekatan pemodelan topik penting menggunakan LSA dan LDA pada penggunaan korpus publikasi ilmiah telah dilakukan untuk mengetahui bahwa dokumen ilmiah menggunakan kosakata yang sangat khusus [11]. Penggunaan secara bersamaan teknik *Latent Dirichlet Allocation* (LDA) dan *Latent Semantic Index* (LSA) saling melengkapi dibandingkan dengan penggunaan metode tunggal [12]. Dari penelitian-penelitian yang telah ada tersebut diatas maka dalam penelitian ini akan menggunakan metode *Latent Dirichlet Allocation* untuk memperkuat atas hasil penelitian-penelitian sebelumnya mengenai penggunaan teknik *Latent Dirichlet Allocation* dalam permodelan topik.

II. METODOLOGI PENELITIAN

A. Pra-pemrosesan Data

Dataset yang digunakan adalah data ulasan hotel pada Tripadvisor data sebanyak 20.491 baris yang terdiri dari dua kolom yaitu Review dan Rating. Preprocessing data sangat penting dan memengaruhi secara substansial hasil percobaan[11]. *Preprocessing* data dilakukan dengan menggunakan *lemmatizer* dan *stopwords* dengan menghapus semua kata bahasa Inggris atau kata-kata paling umum dalam bahasa Inggris yang tidak menambahkan banyak arti pada suatu kalimat. Selain itu dilakukan juga penghapusan kata-kata yang memiliki panjang kurang dari 3 karakter yang tidak memiliki arti penting dalam sebuah kalimat seperti hm,at,ab,cc,er,ww,zc,nm,dll.

	Review	Rating	Review_cleaned_text
0	nice hotel expensive parking got good deal sta...	4	nice hotel expensive parking good deal stay ho...
1	ok nothing special charge diamond member hilito...	2	nothing special charge diamond member hilton d...
2	nice rooms not 4* experience hotel monaco seat...	3	nice room experience hotel monaco seattle good...
3	unique, great stay, wonderful time hotel monac...	5	unique great stay wonderful time hotel monaco ...
4	great stay great stay, went seahawk game aweso...	5	great stay great stay went seahawk game awesom...

Gambar 1. Hasil pra-pemrosesan data

Gambar 1 menjelaskan adanya perbedaan setelah dilakukan penghapusan *stopwords*, tanda baca serta lemmatisasi dimana dapat dilihat contohnya pada baris ke tiga

```
['nice hotel expensive parking good deal stay hotel anniversary arrived late evening took advice previous review valet
'nothing special charge diamond member hilton decided chain shot 20th anniversary seattle start booked suite paid extr
'nice room experience hotel monaco seattle good hotel level.positives large bathroom mediterranean suite comfortable p
'unique great stay wonderful time hotel monaco location excellent short stroll main downtown shopping area friendly rc
'great stay great stay went seahawk game awesome downfall view building complain room huge staff helpful booked hotel
'love monaco staff husband stayed hotel crazy weekend attending memorial service best friend husband celebrating 12th
'cozy stay rainy city husband spent night monaco early january 2008. business trip chance come ride.we booked monte ca
'excellent staff housekeeping quality hotel chocked staff make feel home experienced exceptional service desk staff cc
'hotel stayed hotel monaco cruise room generous decorated uniquely hotel remodeled pacific bell building charm sturdin
'excellent stayed hotel monaco past delight reception staff friendly professional room smart comfortable particularly
```

Gambar 2. Korpus dari ulasan hotel Tripadvisor

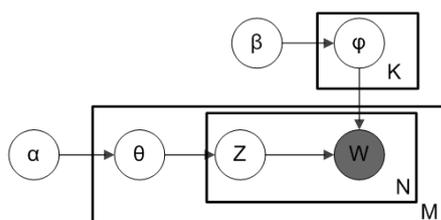
B. Permodelan Topik

Teknik *information retrieval* yang diimplementasikan pada algoritma *Latent Dirichlet Allocation* (LDA) untuk memperoleh informasi dari ulasan hotel pada Tripadvisor. Penggunaan algoritma LDA dapat menunjukkan bagaimana algoritma tersebut menghasilkan topik atau pengelompokkan kata pada korpus teks secara otomatis.

Ide awal *Latent Dirichlet Allocation* (LDA) adalah asumsi bahwa dokumen dianggap sebagai kombinasi dari beberapa topik, dimana karakteristik suatu topik ditentukan oleh distribusi kata[1]. Dalam LDA, kata tersebut disebut sebagai istilah. Kumpulan kata disebut dokumen, dan kumpulan dokumen disebut corpus. Kumpulan semua istilah dalam korpus disebut sebagai kosakata. Generatif dasar LDA mekanismenya sangat mirip dengan pLSI. Asumsi yang digunakan dalam LDA adalah dokumen dianggap sebagai kantong kata-kata. Beberapa langkah yang diperlukan dalam LDA untuk menentukan topik adalah sebagai berikut.

- a) Menentukan jumlah topik
- b) Memberikan inisialisasi topik acak dalam kata-kata yang merupakan proses generatif untuk setiap dokumen pada korpus
- c) Menghitung nilai probabilitas topik pada dokumen dan probabilitas kata pada topik untuk melihat prevalensi topik pada dokumen dan kemungkinan kata-kata pada topik.
- d) Memperbarui topik pada setiap kata berdasarkan nilai probabilitas tertinggi.

Model graf probabilistik dari LDA pada langkah kedua dapat dilihat pada gambar 1 sebagai *Directed Acyclic Graph* (DAG).



Gambar 3. Model grafik probabilistik dari LDA

kata “rooms” berubah menjadi “room”. Setelah itu akan disiapkan korpus untuk kebutuhan analisis serta memeriksa 10 entri pertama dari review seperti terlihat di bawah ini.

Berdasarkan gambar 1, α dan β merupakan parameter pada tingkat corpus. θ adalah variabel di tingkat dokumen. Sedangkan z dan w adalah variabel dalam tingkat kata (istilah). Variabel adalah simbol objek yang dapat diisi dengan konten bergantian tetapi hanya jenis yang memiliki distribusi. Dalam hal ini, isi z dan w dapat diubah. Karena perubahan nilai yang mengikuti distribusi, variabel memiliki parameter tertentu.

Setelah didapatkan korpus dari hasil *pre-processing* data, kemudian dilakukan pembuatan folder *temporary* untuk menampung kamus data sementara dan korpus. Kemudian akan dilakukan penghapusan kata-kata yang umum yang ada dalam korpus serta melakukan tokenisasi. Langkah selanjutnya mengekstrak fitur dan menciptakan *Document-Term-Matrix* (DTM) yang mana didalam *Document-Term-Matrix* nilainya adalah nilai dari Term Frequency-Inverse Document Frequency (TF-IDF) yang dimanfaatkan untuk mengetahui frekuensi suatu kata berapa kali muncul dalam suatu dokumen.

(20491, 1000)	
(0, 314)	0.09762007796425641
(0, 613)	0.10218520564050582
(0, 771)	0.11399237618097087
(0, 261)	0.11429820501621801
(0, 957)	0.10054130749159866
(0, 383)	0.05282654657047758
(0, 496)	0.06181750173261944
(0, 499)	0.16031068782605226
(0, 838)	0.10647601847069599
(0, 22)	0.1683203249198376
(0, 867)	0.14254066929574236
(0, 895)	0.13523975875057437
(0, 75)	0.1277610609193184
(0, 583)	0.1347704070277188
(0, 536)	0.13946410109734583
(0, 396)	0.16402785793218935

Gambar 4. Hasil *Term Frequency-Inverse Document Frequency* (TF-IDF)

Pada gambar 4 diatas dapat dilihat bahwa kata-kata yang paling sering dan yang jarang muncul di ulasan hotel tripadvisor berdasarkan skor *Inverse Document Frequency* dimana semakin kecil nilainya maka kata-kata tersebut adalah yang paling banyak muncul dalam ulasan hotel tripadvisor dalam semua dokumen. Demikian juga sebaliknya, nilai *skor Inverse Document Frequency* yang semakin besar nilainya

menunjukkan bahwasanya kata-kata tersebut adalah kata-kata yang jarang muncul dalam semua dokumen ulasan hotel tripadvisor.

```
room iberostar
1.1691546692024446
5.702043171345892
```

Gambar 5. Contoh hasil *Term Frequency-Inverse Document Frequency* (TF-IDF)

Dari gambar 5 bahwa berdasarkan nilai *Inverse Document Frequency*, 'room' adalah kata yang paling sering muncul sedangkan 'iberostar' adalah kata yang paling kurang kemunculannya dalam ulasan hotel tripadvisor. Secara keseluruhan *dataframe* dari *Latent Dirichlet Allocation* yang terlihat pada gambar 6 berikut ini.

	room	hotel	night	told	time	stay	desk	service	like	said	staff
0	0.026180	0.014856	0.007182	0.006673	0.006413	0.005777	0.00538	0.005136	0.005121	0.005076	0.004892
1	0.008526	0.000000	0.000000	0.000000	0.010848	0.000000	0.000000	0.005222	0.006274	0.000000	0.005291
2	0.035447	0.037187	0.006905	0.000000	0.007386	0.012659	0.000000	0.010281	0.000000	0.000000	0.010740
3	0.030224	0.055009	0.008919	0.000000	0.000000	0.011157	0.000000	0.000000	0.000000	0.000000	0.013551

Gambar 6. *Dataframe Latent Dirichlet Allocation*

Pada gambar 6 menunjukkan hasil permodelan topik yang menampilkan kata-kata pada setiap topik dimana setiap nilai

menunjukkan prosentase kontribusi topik yang sesuai dalam dokumen.

```
[(0,
'0.026*room" + 0.015*hotel" + 0.007*night" + 0.007*told" + 0.006*time"'),
(1,
'0.018*beach" + 0.016*resort" + 0.013*pool" + 0.013*food" + 0.011*time"'),
(2,
'0.037*hotel" + 0.035*room" + 0.020*great" + 0.013*stay" + 0.011*nice"'),
(3,
'0.055*hotel" + 0.030*room" + 0.014*staff" + 0.013*location" + 0.011*stay"'),
(4,
'0.022*hotel" + 0.020*room" + 0.015*good" + 0.012*great" + 0.011*nice"')]
```

Gambar 7. Nilai koherensi kata pada setiap topik

Pelabelan pada topik merujuk pada penyebaran kata yang diperoleh dari permodelan topik dan dapat dilihat bahwa

setiap topik dapat membentuk suatu makna sehingga bisa disimpulkan persebaran katanya baik. Dapat terlihat dari hasil pelabelan pada topik di bawah ini.

```
Topic 0:
resort beach room pool food time hotel service great people

Topic 1:
hotel room great location staff stay good breakfast excellent stayed

Topic 2:
amsterdam madrid smoking smoke non europe tram princess puerto rico

Topic 3:
hilton london sydney harbour westin pleasantly executive marriott surprised traveller

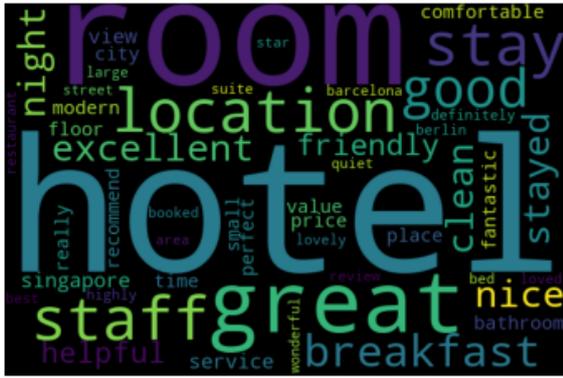
Topic 4:
hotel room great good location stay staff nice breakfast night
```

Gambar 8. Hasil pemodelan topik

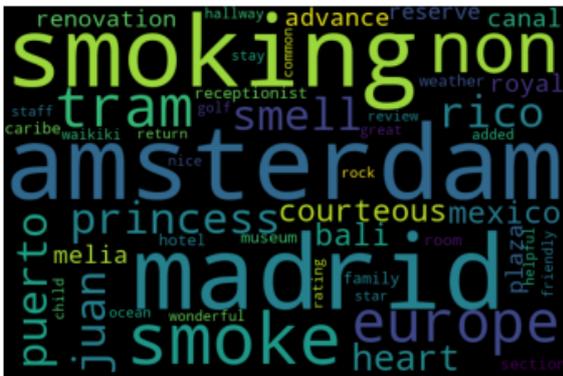
Pemrosesan dalam LDA dilakukan dengan mengkonversi dokumen ke dalam kamus data serta mengkonversinya ke dalam corpus (matriks dokumen) kemudian akan algoritma LDA digunakan untuk membentuk permodelan topik. Parameter input yang digunakan dalam pemrosesan LDA berupa jumlah topik sebanyak lima buah.

III. HASIL DAN PEMBAHASAN

Gambaran mengenai kata-kata penting yang ada di korpus sebanyak tiga puluh buah. Visualisasi tersebut dapat terlihat pada gambar 9 tentang visualisasi LDA.



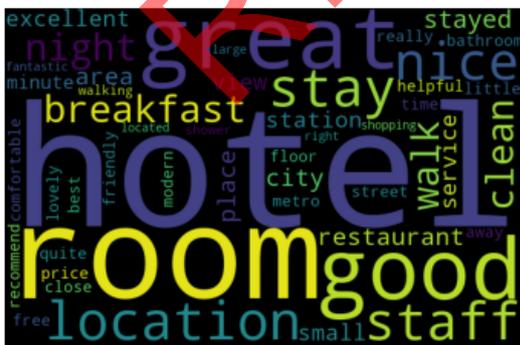
Gambar 11. Wordcloud topik kedua



Gambar 12. Wordcloud topik ketiga



Gambar 13. Wordcloud topik keempat



Gambar 14. Wordcloud topik kelima

IV. KESIMPULAN

Dari hasil pemodelan topik *Latent Dirichlet Allocation* yang telah dilakukan terhadap dataset review tripadvisor dapat disimpulkan bahwa tren ulasan lebih banyak membahas mengenai *location*, *service*, *hotel*, *breakfast*, *resort* dan *beach*. Tren ulasan yang didapat tersebut bisa dipergunakan bagi pemangku kepentingan dalam hal ini manajemen hotel untuk lebih memaksimalkan pelayanan serta pertimbangan lokasi pembangunan hotel yang baru untuk dapat meningkatkan performa dan keuntungan finansial yang memadai. Pemodelan topik yang dilakukan adalah *unsupervised learning* yang memerlukan ketepatan Analisa dalam mengartikan topik yang dihasilkan oleh model yang mana terlihat dari beberapa sebaran kata yang kurang sesuai dengan bahasan topik. Kedepannya untuk bisa memberikan nilai lebih metode ini dapat dikembangkan ataupun di gabungkan menggunakan pemodelan topik yang lain.

DAFTAR PUSTAKA

- [1] Y. Kalepalli, S. Tasneem, P. D. P. Teja, and S. Manne, "Effective Comparison of LDA with LSA for Topic Modelling," Proc. Int. Conf. Intell. Comput. Control Syst. ICICCS 2020, no. Iciccs, pp. 1245–1250, 2020, doi: 10.1109/ICICCS48265.2020.9120888.
- [2] S. H. Mohammed and S. Al-Augby, "LSA & LDA topic modeling classification: Comparison study on E-books," Indones. J. Electr. Eng. Comput. Sci., vol. 19, no. 1, pp. 353–362, 2020, doi: 10.11591/ijeecs.v19.i1.pp353-362.
- [3] J. Blad, K. Svensson, J. Blad, and K. Svensson, "Exploring NMF and LDA Topic Models of Swedish News Articles News Articles," no. December, 2020.
- [4] S. İLHAN OMURCA, E. EKİNCİ, E. YAKUPOĞLU, E. ARSLAN, and B. ÇAPAR, "Automatic Detection of the Topics in Customer Complaints with Artificial Intelligence," Balk. J. Electr. Comput. Eng., vol. 9, no. 3, pp. 268–277, 2021, doi: 10.17694/bajece.832274.
- [5] M. B. Mutanga and A. Abayomi, "Tweeting on COVID-19 pandemic in South Africa: LDA-based topic modelling approach," African J. Sci. Technol. Innov. Dev., vol. 0, no. 0, pp. 1–10, 2020, doi: 10.1080/20421338.2020.1817262.
- [6] P. Tijare and P. J. Rani, "Exploring popular topic models," J. Phys. Conf. Ser., vol. 1706, no. 1, 2020, doi: 10.1088/1742-6596/1706/1/012171.
- [7] R. Rani and D. K. Lobiyal, "An extractive text summarization approach using tagged-LDA based topic modeling," Multimed. Tools Appl., vol. 80, no. 3, pp. 3275–3305, 2021, doi: 10.1007/s11042-020-09549-3.
- [8] V. K. Garbhapu, "A comparative analysis of Latent Semantic analysis and Latent Dirichlet allocation topic modeling methods using Bible data," Indian J. Sci. Technol., vol. 13, no. 44, pp. 4474–4482, 2020, doi: 10.17485/ijst/v13i44.1479.
- [9] H. P. Suresha and K. Kumar Tiwari, "Topic Modeling and Sentiment Analysis of Electric Vehicles of Twitter Data," Asian J. Res. Comput. Sci., no. October, pp. 13–29, 2021, doi: 10.9734/ajrcos/2021/v12i230278.
- [10] H. Gupta and M. Patel, "Method of Text Summarization Using Lsa and Sentence Based Topic Modelling with Bert," Proc. - Int. Conf. Artif. Intell. Smart Syst. ICAIS 2021, pp. 511–517, 2021, doi: 10.1109/ICAIS50930.2021.9395976.
- [11] S. Bellaouar, M. M. Bellaouar, and I. E. Ghada, "Topic modeling: Comparison of LSA and LDA on scientific publications," ACM Int. Conf. Proceeding Ser., pp. 59–64, 2021, doi: 10.1145/3456146.3456156.
- [12] T. Williams and J. Betak, "A Comparison of LSA and LDA for the Analysis of Railroad Accident Text," J. Ubiquitous Syst.

- Pervasive Networks, vol. 11, no. 1, pp. 11–15, 2019, doi: 10.5383/juspn.11.01.002.
- [13] J. C. Campbell, A. Hindle, and E. Stroulia, “Latent Dirichlet Allocation: Extracting Topics from Software Engineering Data,” *Art Sci. Anal. Softw. Data*, vol. 3, pp. 139–159, 2015, doi: 10.1016/B978-0-12-411519-4.00006-9.
- [14] S. Nikou, H. Bin Selemat, R. C. M. Yussoff, and M. M. Khiabani, “Identifying the impact of hotel image on customer loyalty: a case study from four star hotels in Kuala Lumpur, Malasia,” *Int. J. Soc. Sci. Econ. Res.*, vol. 2, no. 3, pp. 2786–2812, 2017.
- [15] G. Tovmasyan, “Evaluating the quality of hotel services based on tourists’ perceptions and expectations: The case study of Armenia,” *J. Int. Stud.*, vol. 13, no. 1, pp. 93–107, 2020, doi: 10.14254/2071-8330.2020/13-1/6.
- [16] A. Mandić and L. Petrić, *The impacts of location and attributes of protected natural areas on hotel prices : implications for sustainable tourism development*, no. 0123456789. Springer Netherlands, 2020.
- [17] A. Gelbman and A. Gelbman, “Seaside hotel location and environmental impact : land use dilemmas dilemmas,” *J. Tour. Cult. Chang.*, vol. 0, no. 0, pp. 1–21, 2021, doi: 10.1080/14766825.2021.1961797.

RETRACTED