


<div><div>Exploratory Analysis and Sentiment Mining of Destination Reviews in Sri Lanka Hannah Cinderella L Kishan V Dr. Pattabiraman V School of Computer Science Engineering</div></div>																																											
INTRODUCTION	RESULTS																																										
<p>Tourism in Sri Lanka is a thriving industry fueled by the country’s scenic beauty and cultural heritage. With the increasing influence of digital platforms, tourists leave reviews that contain rich, unstructured data.</p> <p>Project Objectives:</p> <ul style="list-style-type: none">Understand travel preferences via review analysis.Group destinations by sentiment using clustering.Recommend travel paths based on co-visitation patterns.Build predictive models to classify review sentiment.	<p>Top 20 Destinations Identified: Horton Plains, Sigiriya, Ella, Nuwara Eliya, Mirissa, Galle Fort, Yala National Park, Anuradhapura, Polonnaruwa, and others.</p> <p>Word Cloud showed frequent mentions of “beautiful”, “nature”, “peaceful”, “waterfall”, “temple”, “historic”.</p> <p>Association Rules: "Horton Plains" often co-mentioned with "Bakers Falls", and "Nuwara Eliya" – indicating popular travel circuits.</p> <p>“Galle” with “Mirissa” and “Unawatuna” for beach tourism</p> <p>DBSCAN Clusters: 3 major clusters indicating positive, neutral, and negative sentiment zones.</p> <p>Predictive Modeling:</p> <p>Class Imbalance Handling:</p> <ul style="list-style-type: none">Applied upsampling to balance sentiment classes. <p>Model Performance</p> <table><tr><th>Model</th><th>Accuracy</th><th>F1 Score</th><th>Precision</th><th>Best performer in</th></tr><tr><td>Logistic Regression</td><td>91.4%</td><td>0.91</td><td>0.92</td><td>Precision <input checked="" type="checkbox"/></td></tr><tr><td>Naive Bayes</td><td>86.3%</td><td>0.86</td><td>0.86</td><td>---</td></tr><tr><td>Random Forest</td><td>95.5%</td><td>0.95</td><td>0.96</td><td>All metrics <input checked="" type="checkbox"/></td></tr></table>	Model	Accuracy	F1 Score	Precision	Best performer in	Logistic Regression	91.4%	0.91	0.92	Precision <input checked="" type="checkbox"/>	Naive Bayes	86.3%	0.86	0.86	---	Random Forest	95.5%	0.95	0.96	All metrics <input checked="" type="checkbox"/>																						
Model	Accuracy	F1 Score	Precision	Best performer in																																							
Logistic Regression	91.4%	0.91	0.92	Precision <input checked="" type="checkbox"/>																																							
Naive Bayes	86.3%	0.86	0.86	---																																							
Random Forest	95.5%	0.95	0.96	All metrics <input checked="" type="checkbox"/>																																							
SCOPE OF THE PROJECT																																											
<p>This project focuses on extracting meaningful insights from user - generated reviews of Sri Lankan tourist destinations. By leveraging Exploratory Data Analysis, sentiment mining, clustering, and predictive modelling, the project aims to:</p> <ul style="list-style-type: none">Identify popular tourist attractions based on review volume and sentiment.Classify tourist experiences into positive and negative groups using machine learning.Discover frequently co-visited destination pairs to enhance travel planning.Provide data-backed recommendations for tourism boards and travel service providers.																																											
METHODOLOGY																																											
<p>Step 1: Data Cleaning & Sentiment Extraction (TextBlob)</p> <p>Step 2: TF-IDF Vectorization (converts text to numerical features) & Variance Threshold Feature Selection (removes low-variance features)</p> <p>Step 3: EDA — Visualize Top 20 Reviewed Destinations, Word Frequency (Word Cloud), and print frequent destination pairs using association-style logic</p> <p>Step 4: DBSCAN Clustering on TF-IDF + SVD-Reduced Features</p> <p>Step 5: Association Rule Mining — Use co-occurrence of destinations per time span to extract frequent destination pairs</p> <p>Step 6: Predictive Modeling — Train and compare: Logistic Regression, Naive Bayes, Random Forest</p> <p>Step 7: Evaluation using Accuracy & F1-score</p> <p>Step 8: Resampling Techniques — Apply Random Oversampling to balance class distribution</p>	CONCLUSION																																										
	<p>Sentiment analysis reveals - 70% of reviews are positive.</p> <p>Top attractions are mainly nature parks, waterfalls, beaches, and historical sites.</p> <p>Random Forest achieved best sentiment prediction (95.5%)</p> <p>Co-visit patterns can help plan better tour packages.</p> <p>DBSCAN effectively grouped sentiment-based clusters.</p> <p>Feature extraction significantly improved predictive performance.</p>																																										
	FUTURE SCOPE																																										
	<p>Add confusion matrix & sentiment-wise feature importance.</p>																																										
	CONTACT DETAILS																																										
	<p>Hannah Cinderella L - hannahcinderella.l2023@vitstudent.ac.in</p> <p>Kishan V - kishan.v2023@vitstudent.ac.in</p> <p>GitHub Repository: https://github.com/Hannah-Cinderella/EDA_DestinationReviews_Srilanka</p>																																										
	REFERENCES																																										
	<p>Dataset: (https://www.kaggle.com/datasets/nethumdperera/travel-destinations-reviews-in-sir-lanka)</p> <p>Libraries Used: pandas, scikit-learn, TextBlob, matplotlib, seaborn, imbalanced-learn</p> <p>Tools: Jupyter Notebook, Python 3.9, WordCloud, DBSCAN, Apriori (mlxtend).</p> <p>Research article: https://www.mdpi.com/2071-1050/14/15/9572</p>																																										
<div><div>Cluster Distribution:</div><div><table><tr><th>Cluster</th><th>Count</th></tr><tr><td>0</td><td>1652</td></tr><tr><td>-1</td><td>252</td></tr><tr><td>5</td><td>29</td></tr><tr><td>2</td><td>20</td></tr><tr><td>1</td><td>17</td></tr><tr><td>3</td><td>9</td></tr><tr><td>4</td><td>7</td></tr><tr><td>6</td><td>7</td></tr><tr><td>7</td><td>7</td></tr></table></div><div>Name: count, dtype: int64</div></div>	Cluster	Count	0	1652	-1	252	5	29	2	20	1	17	3	9	4	7	6	7	7	7	<div><div>Top 10 Most Frequently Visited Destination Pairs:</div><div><table><tr><th>Destination Pair</th><th>Count</th></tr><tr><td>(Horton plains national park, Moon plains)</td><td>102013</td></tr><tr><td>(Bambarakiri ella, Riverston)</td><td>94811</td></tr><tr><td>(Riverston, Sembuwatta lake)</td><td>83207</td></tr><tr><td>(Pitawala pathana , riverston, Riverston)</td><td>76771</td></tr><tr><td>(Bambarakiri ella, Sembuwatta lake)</td><td>70201</td></tr><tr><td>(Riverston, Sera ella water falls)</td><td>65563</td></tr><tr><td>(Bambarakiri ella, Pitawala pathana , riverston)</td><td>64324</td></tr><tr><td>(Riverston, Sri muththumari amman kovil)</td><td>64144</td></tr><tr><td>(Kurunagela clock tower, Yapahuwa rock fortress)</td><td>63291</td></tr><tr><td>(Horton plains national park, Horton plains national park)</td><td>63257</td></tr></table></div></div>	Destination Pair	Count	(Horton plains national park, Moon plains)	102013	(Bambarakiri ella, Riverston)	94811	(Riverston, Sembuwatta lake)	83207	(Pitawala pathana , riverston, Riverston)	76771	(Bambarakiri ella, Sembuwatta lake)	70201	(Riverston, Sera ella water falls)	65563	(Bambarakiri ella, Pitawala pathana , riverston)	64324	(Riverston, Sri muththumari amman kovil)	64144	(Kurunagela clock tower, Yapahuwa rock fortress)	63291	(Horton plains national park, Horton plains national park)	63257
Cluster	Count																																										
0	1652																																										
-1	252																																										
5	29																																										
2	20																																										
1	17																																										
3	9																																										
4	7																																										
6	7																																										
7	7																																										
Destination Pair	Count																																										
(Horton plains national park, Moon plains)	102013																																										
(Bambarakiri ella, Riverston)	94811																																										
(Riverston, Sembuwatta lake)	83207																																										
(Pitawala pathana , riverston, Riverston)	76771																																										
(Bambarakiri ella, Sembuwatta lake)	70201																																										
(Riverston, Sera ella water falls)	65563																																										
(Bambarakiri ella, Pitawala pathana , riverston)	64324																																										
(Riverston, Sri muththumari amman kovil)	64144																																										
(Kurunagela clock tower, Yapahuwa rock fortress)	63291																																										
(Horton plains national park, Horton plains national park)	63257																																										
Figure 1: DBSCAN Clustering	Figure 2: Association Rule Mining																																										