# Find Your Perfect Park – An NLP-based Personalised Recommendation System

**LI Tianyi (A0262755X)  |  QIU Guanrong (A0198358J)  |  YANG Muzi (A0144173L)** [1]

## Abstract

This project paper presents a recommendation system for parks in Singapore that aims to provide personalized park recommendations based on user preferences. However, the effectiveness and accuracy of the system are limited by several factors, including the lack of user preference data and potential biases in the data. To address these limitations, future work should focus on collecting user preference data, reducing biases in the data, incorporating context into recommendations, evaluating the system using user feedback, and incorporating machine learning techniques. Despite the limitations, the recommendation system has the potential to be a valuable tool for individuals looking to explore parks in Singapore, and future improvements can make the system even more effective and personalized.

## 1. Introduction

Singapore is a densely populated city-state that faces the challenge of balancing urban development and environmental conservation. To address this challenge, the government has adopted a "**City in a Garden**" vision that aims to integrate greenery and nature into the urban landscape, enhancing the liveability and attractiveness of the city. One of the key components of this vision is the development and management of natural parks, which are public spaces that offer recreational, educational, and ecological benefits to residents and visitors. Since the 1960s, the government has been investing in creating and maintaining parks as part of its vision to improve the quality of life for its residents and visitors. Natural parks not only provide opportunities for people to enjoy nature and outdoor activities, but also contribute to the conservation of biodiversity and ecosystem services.

Despite the abundance and diversity of natural parks in Singapore, many people may not be aware of their existence or suitability for their preferences. This lack of knowledge is particularly relevant in the current post-COVID-19 pandemic era, where countries are lifting travel restrictions and promoting local tourism to support the recovery of the tourism industry. The Singapore government is no exception and has been actively encouraging locals to explore and appreciate the natural attractions within their own country. As a result, there is a need for innovative solutions that can help to promote domestic tourism and encourage locals to explore and appreciate the natural attractions within their own country. In this project, we aim to address this need by developing a **personalised recommendation system** for natural parks in Singapore that can help users to discover and visit the most suitable parks for them based on their personal preferences and needs.

Recommendation systems are valuable as they aid users in discovering products and services they may not have otherwise encountered independently. They are designed to analyse the interactions and behaviour of individuals and products, by gathering

---

data such as clicks, impressions, likes, and purchases. This information allows these systems to comprehend the preferences and characteristics of users. Recommendation systems are highly valued by content and product providers due to their ability to personalize recommendations and anticipate the interests and desires of each individual. These systems can guide users towards a wide range of products and services catering to their specific preferences.

This project holds immense significance, as it can create a win-win situation for both customers and natural parks. The personalised recommendation system will allow customers to find the most suitable natural parks based on their preferences, which can help improve customer satisfaction, engagement, and loyalty. It can also increase their awareness and appreciation of the natural beauty and attractions within Singapore. On the other hand, for natural parks, the system can help them to increase their visibility and reputation in the market by providing accurate and personalised recommendations to potential visitors. This can lead to increased visitation and revenue for the parks. It can also help them to improve their offerings and services by understanding and catering to the customers' needs.

Therefore, this project aims to develop an accurate and personalised recommendation system that can benefit both customers and natural parks. By providing a comprehensive understanding of customers' preferences, natural parks can improve their offerings and services, and customers can explore the beauty of Singapore's natural parks with ease. The success of this project has the potential to significantly impact the tourism industry and natural parks in Singapore.

**1.1 Objective**

The main aim of this project is to develop a personalised recommendation system for natural parks in Singapore that can assist customers in selecting the most appropriate park based on their preferences and interests. The system will utilise data from Google Map, to collect and analyse

reviews and ratings of parks. Natural Language Processing (NLP) will be used to extract park features from the reviews, which will be matched with customer input to provide tailored recommendations. The objective of the system is to enhance customer satisfaction, engagement, and loyalty while simultaneously increasing the revenue and visibility of the parks.

## 2. Methodology

The recommender function is a crucial element of recommendation systems as it uses information and user preferences to predict how a user might rate a park. This ability to predict user satisfaction or rating, even before receiving one, makes recommendation systems a powerful tool. To build a typical recommendation system, three types of data are required: user behaviour data (such as ratings, clicks, and browsing history), user demographic data (including age, education, income, and location), and product attribute data (such as park features). In our project, we can only gather data on natural parks using NLP techniques to extract information from reviews. While this is not a complete recommendation system, it can be integrated with other data sources to create a more comprehensive system or combined with other similar engines to create a sophisticated recommendation system. Besides, the features extracted from the NLP models can be add to the feature warehouse of larger recommendation system.

The methodology for this project involves a multi-step process, including data collection, exploratory data analysis, data pre-processing, park feature extraction, recommendation methods. Firstly, reviews and ratings of natural parks are crawled and collected from Google Map. The data will then be cleaned and prepared for further analysis. Park features extraction will be performed using NLP models to categorize the comments into different categories, and each park will be marked in different categories based on the sentiment of the corresponding comments. Customer preferences

will be analysed and summarized using NLP models, which will generate an overall score based on the matched features of parks and personalized requirements. The recommendation method will then suggest the park with the highest overall mark to the customer based on their requirements. The proposed methodology aims to provide a personalized recommendation system that captures user preferences and matches them with the features of natural parks in Singapore, ultimately boosting customer satisfaction and engagement with these parks.

## 3. Data Preparation

Collecting reviews and ratings data from Google Maps through API keys is a practical and efficient way to gather a large amount of data about parks in Singapore. This approach ensures that the recommendation system reflects the most realistic features of each park. Google Maps is a widely used platform where park-goers can provide feedback and share their experiences, making the data collected representative of the overall population of park-goers in Singapore. Using API keys to crawl data from reviews and ratings of parks allows the researchers to collect information about a large number of parks quickly and with minimal cost. Moreover, collecting data in a structured format through API keys enables researchers to process and analyse it more efficiently. This facilitates identifying trends and patterns in park preferences, which can be used to improve the accuracy of the recommendation system.

### 3.1 Data Collection

This dataset contains information about parks in Singapore that is necessary for building a recommendation system. The National Parks Board provided the full list of parks in Singapore, which served as the basis for the selection of the top 22 popular parks. Due to limited time and resources, only these parks were included in the dataset. Reviews and ratings data for the top 22 parks were collected using Google Maps API keys. Each park has 1,000 records, except for some parks that have

less than 1,000 records, in which case, the available records were taken. The dataset contains a total of 22,000 records, with each record consisting of the park name, overall rating, total number of reviews, individual review and rating, and the timestamp of the review. The dataset serves as the basis for building a recommendation system that can provide personalized park recommendations based on user preferences. The collected reviews and ratings are used to calculate the dimensions of the parks, including purpose, time, family children, and crowdedness. These dimensions can be used to generate a score for each park, which is used to match parks with user preferences. However, the dataset has missing data due to some parks having fewer than 1,000 records, which may limit the accuracy of the recommendation system.

### 3.2 Exploratory Data Analysis (EDA)

EDA is a crucial step in any data science project, as it provides insights into the underlying data and informs subsequent data pre-processing and modelling. In this project, we performed EDA on the park reviews using word cloud visualization to identify any patterns or trends in the data. Word cloud is a popular tool for visualizing textual data, which displays the most frequent words in the data in a visual and intuitive manner. By analysing the word cloud, we were able to gain a better understanding of the most used words and themes in the park reviews, which could inform our feature engineering and model selection. For instance, we find that the most frequent words are related to family or children, which would suggest that these attributes are important to visitors and should be included in our feature set. The insights gained from EDA are crucial for ensuring that subsequent modelling and analysis are based on a solid understanding of the data and its underlying patterns.

Word Cloud of Park Reviews



Word Cloud of Park Reviews: Remove Adjectives

The word cloud showed that the visitors were primarily concerned about the purpose of their visit, the time of their visit, whether the park was family- and children-friendly, and the level of crowding. These observations suggest that visitors place significant importance on these factors when reviewing a park. By identifying these important factors through EDA, we can focus on extracting and engineering features that are most relevant to visitors and improve the effectiveness of our recommendation system.

### 3.3 Feature Engineering

With the observations made in EDA, we have manually labelled 3 parks for the following features: purpose, time, family-friendliness, children-friendliness, and crowdedness based on the review text. The features are labelled as follows:

| Features | Tags |
|---|---|
| *purpose* | 1 – leisure<br>2 – social<br>3 – exercise<br>4 – photography<br>5 – other |
| *time* | 1 – morning<br>2 – afternoon<br>3 – evening |
| *family* | 0 – none<br>1 – yes<br>2 – no |
| *children* | 0 – none<br>1 – yes<br>2 – no |
| *crowded* | 0 – none<br>1 – yes<br>2 – no |

The purpose of manually labeling data for features such as purpose, time, family-friendliness, children-friendliness, and crowdedness is to use these labeled data to build and train a model that can predict these features for the remaining parks based on the reviews. This approach can help automate the process of extracting information from the reviews, potentially reducing the time and effort required for manual work. By using a trained model, we can quickly and accurately identify the features of parks based on the reviews, which can be beneficial for park managers, policymakers, and tourists. In the

4

following sections, we will explain in detail how this model is built and trained using the labeled data.

## 4. NLP Text Labelling

Auto-labelling with NLP models is an efficient and effective technique that can save time and resources while still providing accurate labels for text data. In our project, it helps us quickly label a large number of reviews to train our model for predicting the features of parks. In our project, we used a pre-trained BERT model to automatically label the reviews of the parks based on the same features that we manually labelled, such as *purpose*, *time*, *family*, *children*, and *crowdedness*. By using a pre-trained NLP model, we can leverage the power of machine learning to accurately predict the features of the parks based on the text in the reviews. This saved us a significant amount of time and resources compared to manually labelling each review.

This process is also known as transfer learning. We leverage pretrained language model, BERT, to train a new model to predict the features of parks based on reviews. The pre-trained model has already been trained on a large amount of text data, which allows it to learn general language representations. By using transfer learning, we can benefit from the knowledge learned by the pre-trained models, which can reduce the amount of data needed for training and improve the performance of our model. Additionally, transfer learning can save us a lot of time and resources compared to training a model from scratch, as training large language models on large amounts of text data can be computationally intensive and time-consuming. Furthermore, we can leverage the knowledge from the pre-trained model to improve the accuracy and efficiency of our model in predicting the features of parks based on the reviews.

### 4.1 BERT-based NLP Labeling Model

We combined park review data from multiple csv files into a comprehensive Data Frame. This dataset was then split into training and testing sets to facilitate model evaluation. To process the raw review data and make it compatible with the pre-trained BERT model, we implemented a custom dataset class called Review_Dataset, which is inherited from PyTorch's Dataset class. We employed Auto Tokenizer for text processing and generated an encoded dictionary for each review, containing input features and target labels.

After the pre-processing pipeline, we developed a BERT-based multi-task classification model to predict five attributes of park reviews: *purpose*, *time*, *family*, *children*, and *crowded*. During the training phase, multiple epochs were performed. In each batch, data was sent to the GPU, and predictions were made using the model. Losses for each attribute were calculated and summed up to obtain the total loss. This total loss was used for gradient backpropagation and optimizer updates.

### 4.2 Performance Evaluation

In the evaluation phase, predictions were made on the test set. The predicted results were compared with the actual labels to calculate the accuracy of each attribute as follows:

**Accuracy = [0.615 0.99 0.96 0.945 0.97]**

Our model demonstrated excellent performance in predicting the attributes of *time*, *family*, *children*, and *crowded*, with accuracies above 94%. However, the accuracy for the *purpose* attribute was comparatively lower at 61.5%. The possible reason could be that the purpose of the park visits could be more subjective (there could be some reviews which contains irrelevant contents) and context-dependent (the aspect of the review may not be single, but we only choose one of them to label) than the other attributes, leading to a higher degree of complexity and ambiguity in the data, making it difficult for the model to capture the patterns effectively.

Despite a few minor issues, we remain confident in our model's overall predictive capabilities. Consequently, we decided to employ this model to automatically label the remaining *10,000+* reviews. After labeling, we inspected some of the generated labels, which showed that the accuracy was indeed quite high.

### 4.3 Auto Labelling Pipeline

We want to automate the process by building an Auto Labelling Pipeline using NLP transfer learning. For our project, we have chosen *__bert_base_uncased__*, which has been pretrained on a large corpus of text data, to serve as the foundation of our pipeline. It will allow us to extract information of the required features from raw text data. Next, we have fine-tuned our pre-trained model using our manually labelled dataset of parks, so that it can learn the patterns and features of parks that we are interested in. After fine-tuning, we can use the model to auto label any new text data related to parks. This means that we can feed raw review text into the pipeline, and the model will automatically extract and categorize information related to those five features. More data from any new parks can be processed efficiently and effectively, without the need for manual labelling. We can easily expand the size of training data for recommendation system to improve the accuracy of recommendations with this pipeline in the future.

## 5. Recommendation Systems

### 5.1 Logistic Regression Based System

The first recommendation system we applied is based on the logistic regression model. We managed to use the five attributes as input features and predict the park associated with each review as output. The data, consisting of labelled and predicted instances, is consolidated into a single dataframe, which is then divided into training and test sets. The training set accounts for 80% of the data, while the test set comprises the remaining 20%.

To create the recommendation model, we employed a linear SVM classifier using the extracted features (purpose, time, family, children, and crowded). After training the model on the training set, we evaluated its performance on the test set. Unfortunately, the recommendation system yielded a low accuracy.

There could be several reasons for the low accuracy. One possible explanation is that the features may not sufficiently capture the nuances needed for an effective recommendation system, as the parks share similar features. For instance, it may be challenging for even humans to correctly recognize over 90% of the reviews concerning West Coast Park and East Coast Park based solely on the content of the reviews. Additionally, the linear SVM classifier might not be the optimal choice for this problem, as the data could have a more complex, non-linear relationship.

Despite the low accuracy, we decided to proceed with the current approach for a few reasons. Firstly, this is an initial attempt, and further improvements can be made in the feature selection and model choice. Secondly, the model might still provide some insights into park recommendations, even if the accuracy is not optimal. Finally, as more data becomes available, the model's performance could improve over time, and the low accuracy could serve as a valuable learning experience for future iterations.

### 5.2 Similarity Matching Based System

We implemented a similarity matching-based recommendation system to better consider the specific preferences of park visitors. This approach uses an algorithm that calculates the similarity between parks based on their average ratings for different attributes, in our case *purpose*, *time*, *family*, *children*, and *crowdedness*, and the users' preference. The main idea behind this approach is to recommend parks that are similar to the visitor's own five aspects' preferences which they could choose in our system. The following is a detailed explanation of the algorithm and its implementation:

Load labelled and predicted park review data, and merge them into a single dataframe. Retain only the relevant columns, including purpose, time, family, children, crowded, park name, and rating. Map the numerical values in the dataframe to their corresponding category names (e.g., purpose map, time map, etc.). Use one-hot encoding to transform the categorical variables into binary features, creating an expanded dataframe. Calculate the

average ratings for each park across all binary features by grouping the expanded dataframe by park name. Standardize the park ratings using the MinMaxScaler from the Sklearn library. Define a personal preference weight vector, which represents the importance of each attribute for the visitor (e.g., leisure, afternoon, family, children, no crowded). Calculate a recommendation score for each park by multiplying the standardized park ratings with the personal preference weight vector. Sort the parks by their recommendation scores in descending order.

By using this similarity matching-based recommendation system, we can provide personalized park recommendations to visitors based on their preferences. This approach differs from the previous recommendation system, as it directly utilizes the average ratings of park attributes rather than trying to predict the park based on the review attributes.

## 6. Results & Analysis

To illustrate the effectiveness of the Similarity Matching Based System, we have performed multiple experiments using different orders of preferences. Specifically, we have tested the system using different weight vectors to represent visitor preferences, in order to examine how the system responds to changes in the relative importance of different attributes (such as leisure, time of day, and family-friendliness). By varying the weight vector, we are able to demonstrate how the system can provide recommendations that are tailored to the individual needs and interests of each visitor, while also ensuring that the recommendations are consistent with the visitor's overall preferences. The results of these experiments show that the Similarity Matching Based System is highly effective at identifying parks that are likely to be of interest to visitors and can provide a more satisfying and enjoyable experience for park-goers.

```
Assume the preference is: leisure > afternoon > family > children > no crowded
Recommendation Result:
```

|    | park_name | score |
|----|-----------|-------|
| 21 | west_coast_park(1) | 0.623198 |
| 2 | bukit_timah_nature_reserve | 0.507307 |
| 0 | admiralty_park | 0.432289 |
| 15 | sembawang_park | 0.277135 |
| 3 | changi_beach_park | 0.254151 |
| 6 | east_coast_park | 0.254030 |
| 13 | pasir_ris_park | 0.232805 |
| 11 | lower_seletar_reservoir_park | 0.202875 |
| 17 | singapore_botanic_gardens | 0.200063 |
| 20 | upper_peirce_reservoir_park | 0.169539 |
| 4 | chinese_garden | 0.111582 |
| 16 | sengkang_riverside_park | 0.105405 |
| 9 | kallang_riverside_park | 0.097977 |
| 14 | punggol_waterway_park | 0.073340 |
| 1 | bukit_batok_nature_park | 0.066103 |
| 7 | hindhede_nature_park | 0.064702 |
| 8 | japanese_garden | 0.062808 |
| 10 | kranji_marshes | 0.050954 |
| 18 | sungei_buloh_wetland_reserve | 0.048558 |
| 19 | telok_blangah_hill_park | 0.044968 |

Assuming the preference of the visitor is leisure > afternoon > family > children > no crowd, we have generated a list of recommended parks based on this preference, as shown in the screenshot above. Each park is given a recommendation score, which is calculated by multiplying the standardized park rating with the corresponding weight of each attribute in the preference vector. In this case, the system has recommended "West Coast Park" as the top choice, with a recommendation score of 0.623. This suggests that "West Coast Park" is an excellent choice for visitors who prioritize leisure activities, prefer to visit parks in the afternoon, are looking for family-friendly options, and want to avoid crowded spaces. The next recommended parks in the list are "Bukit Timah Nature Reserve" and "Admiralty Park" with recommendation scores of 0.507 and 0.432, respectively. These parks are also well-suited for visitors who value leisure activities and family-friendly options, but they may be slightly less attractive to those who prioritize visiting parks in the afternoon or avoiding crowded spaces.

```
Assume the preference is: Children > family > leisure > afternoon > no
Recommendation Result:
```

| | park_name | score |
|---|---|---|
| 0 | admiralty_park | 0.768358 |
| 21 | west_coast_park(1) | 0.621372 |
| 15 | sembawang_park | 0.316450 |
| 6 | east_coast_park | 0.313081 |
| 13 | pasir_ris_park | 0.288881 |
| 2 | bukit_timah_nature_reserve | 0.242694 |
| 3 | changi_beach_park | 0.237964 |
| 17 | singapore_botanic_gardens | 0.221217 |
| 11 | lower_seletar_reservoir_park | 0.175388 |
| 18 | sungei_buloh_wetland_reserve | 0.107771 |
| 14 | punggol_waterway_park | 0.105916 |
| 16 | sengkang_riverside_park | 0.102682 |
| 4 | chinese_garden | 0.102007 |
| 20 | upper_peirce_reservoir_park | 0.089137 |
| 1 | bukit_batok_nature_park | 0.063827 |
| 7 | hindhede_nature_park | 0.060667 |
| 10 | kranji_marshes | 0.052287 |
| 9 | kallang_riverside_park | 0.040567 |
| 19 | telok_blangah_hill_park | 0.029176 |

When we performed the experiment with the preference order of children > family > leisure > afternoon > no crowd, the recommendation system generated a different ranking of parks. In this case, the system prioritized the presence of activities for children, followed by a family-friendly atmosphere, leisure options, afternoon availability, and the absence of crowds. As a result, the top recommended parks included those that had high ratings for children's activities and family atmosphere, even if they had slightly lower ratings for leisure or were more crowded. This experiment demonstrates how changing the order of preference can significantly affect the outcome of the recommendation system, highlighting the importance of understanding and defining personal preferences before using such a system.

The Similarity Matching Based System we developed proved to be an effective method for recommending parks to visitors based on their preferences. Through our experiments with different orders of preferences, we showed that the system can generate customized recommendations by taking into account the relative importance of various attributes. The use of one-hot encoding and MinMaxScaler allowed us to standardize the park ratings and weight them according to the visitor's

preference, respectively. The results demonstrated that changing the preference order can significantly affect the ranking of recommended parks, underscoring the importance of having well-defined and accurate preferences. Overall, our recommendation system can be used as a valuable tool for visitors to make informed decisions about which parks to visit based on their personal preferences. Future research could explore how to incorporate additional attributes into the recommendation system and optimize the weighting of preferences to enhance its performance.

## 7. Limitations & Future Research

Our recommendation system may provide a limited degree of personalisation. Without data on user preferences, the recommendation system can only provide generic recommendations based on the overall characteristics of each park. This can limit the system's ability to provide personalized recommendations that match the specific interests and needs of each user. For example, if a user is looking for a park with a playground for their children, the system may recommend a park that has a playground, but it may not be the best match for the user's other preferences, such as the location or the level of crowding. To improve personalisation, it is essential to collect data on user preferences. This can be done through user surveys, online reviews, or by tracking user behaviour on the website or app. By analysing this data, the recommendation system can better understand each user's individual interests and needs and provide more personalized recommendations.

Furthermore, without data on user preferences, it can be difficult to evaluate the accuracy and effectiveness of the recommendation system. It may be unclear whether the recommendations provided by the system are relevant or useful to users. Evaluating the system may also be challenging because it can be difficult to determine whether the recommendations are based on objective criteria or on the biases and limitations of the data. To counter this limitation, we can collect user feedback on the

recommendations provided. This can be done through user surveys, online reviews, or by tracking user behaviour on the website or app. By analysing this feedback, the recommendation system can identify areas for improvement and make changes to the recommendation algorithm accordingly.

Another potential limitation of our project is the lack of context. Without data on user preferences, it can be difficult to understand the context in which the recommendation system is being used. For instance, users may be looking for parks for different purposes, such as exercising, picnicking, or walking their dogs, and the recommendation system may not be able to differentiate between these different contexts. The system may also not be able to account for other factors that may affect the user's experience, such as the weather, time of day, or the user's mood. We can analyse the user's location, time of day, weather, and other relevant factors to improve the relevance of recommendations. By taking these factors into account, the recommendation system can provide more contextually relevant recommendations that reflect the user's immediate needs and preferences.

## 8. Conclusion

The recommendation system for parks in Singapore has the potential to provide a useful tool for individuals looking to explore parks in the area. By analysing data such as park ratings, reviews, and manually labelled dimensions, the system can generate a score for each dimension of the parks, such as purpose, time, family children, and crowdedness, to provide tailored recommendations for users. We have also solved the overall process of addressing all similar questions, from data acquisition to the construction of recommendation systems. Any organization or individual can easily build their own recommendation system in any field through our pipeline. They can also use the pipeline to extract relevant features and add them to the existing recommendation system.

In conclusion, the recommendation system for parks in Singapore has the potential to be a valuable tool

for individuals looking to explore parks in the area or any potential domains. However, to ensure its effectiveness and accuracy, it is important to address the limitations of the study, including the lack of user preference data and potential biases in the data. Future work should focus on collecting user preference data, reducing biases in the data, incorporating context into recommendations, evaluating the system using user feedback, and incorporating machine learning techniques. By addressing these areas, the recommendation system can provide accurate and personalized recommendations that reflect the individual interests and needs of each user.

## References

CSC. (2008). *A city in a garden*. Ethos, 4, 29-34. https://www.csc.gov.sg/articles/a-city-in-a-garden

Park, S., Lee, J., & Kim, J. (2007, November). *Location based recommendation system using bayesian user's preference model in mobile devices.* In International Conference on Ubiquitous Intelligence and Computing (pp. 1130-1139). Springer, Berlin, Heidelberg.

Urban Redevelopment Authority. (n.d.). *Environmental impact assessment*. https://www.ura.gov.sg/Corporate/Planning/Our-Planning-Process/Bringing-plans-to-Reality/Environmental-Impact-Assessment

Ang, S. L. (2021). *Nature conservation in Singapore.* BiblioAsia, 17(1), 4-11. https://biblioasia.nlb.gov.sg/vol-17/issue-1/apr-jun-2021/nature

Zheng, Y., Zhang, L., Xie, X., & Ma, W. Y. (2009). *Mining interesting locations and travel sequences from GPS trajectories.* In Proceedings of the 18th international conference on World wide web (pp. 791-800).

Ye, M., Yin, P., Lee, W. C., & Lee, D. L. (2010). *Exploiting geographical influence for collaborative point-of-interest recommendation.* In Proceedings of the 34th international ACM

SIGIR conference on Research and development in Information Retrieval (pp. 325-334).

Covington, P., Adams, J., & Sargin, E. (2016). *Deep neural networks for YouTube recommendations*. In Proceedings of the 10th ACM conference on recommender systems (pp. 191-198). https://dl.acm.org/doi/10.1145/2959100.2959190

Vatsal, V. (2021, July 12). *Recommendation systems explained*. Towards Data Science. https://towardsdatascience.com/recommendation-systems-explained-a42fc60591ed

Wikipedia contributors. (2021, December 13). *Recommender syste*m. In Wikipedia, The Free Encyclopedia. https://en.wikipedia.org/wiki/Recommender_system

Nan, X., Kanato, K., & Wang, X. (2022). *Design and implementation of a personalized tourism recommendation system based on the data mining and collaborative filtering algorithm*. Computational Intelligence and Neuroscience, 2022, Article ID 1424097. https://doi.org/10.1155/2022/1424097

Zhou, G., Liu, C., Liu, Y., Liu, Z., & Gao, J. (2018). *Personalized recommendation systems: Five hot research topics you must know*. Microsoft Research Blog. https://www.microsoft.com/en-us/research/lab/microsoft-research-asia/articles/personalized-recommendation-systems/