

The investigation related to the influence of dimension manipulation on regret performance based on the upper confidence bound algorithm

Yifei Mao

The Department of Computer Science, Purdue University, West Lafayette, Indiana
47906, United States of America

mao154@purdue.edu

Abstract. This paper aims to investigate the effects of dimension manipulation on the performance of a recommendation algorithm applied to a dataset of restaurant reviews. In this paper, the Zomato dataset which contains restaurant reviews and other relevant information in Bangalore City was used. This paper extracted ratings that each user gave for each restaurant from the list of user feedback for each restaurant. These different ratings were stored as a core feature that was used for the restaurant recommendation algorithm to estimate the true mean rating of each restaurant. Upper Confidence Bound bandit algorithm was used as the restaurant recommendation algorithm to find the restaurant with the highest average rating in the dataset. Dimension raising and dimension reduction were used as ways to manipulate dimension in this paper. Principal Component Analysis was used as the technique to reduce feature dimension in the dataset. It reduced features into two principal components and the first principal component was used in place of the original core feature. Dimension raising was implemented based on the original core feature and another feature that correlated with it the most. The product of these two features was used in place of the original core feature. The experimental results suggest that dimension manipulation leads to decreased regret performance when employing the Upper Confidence Bound algorithm for recommendation. Intriguingly, within the realm of dimension manipulation, dimension reduction exhibited a more adverse impact on regret performance compared to dimension raising.

Keywords: upper confidence bound, dimension manipulation, bandit problem.

1. Introduction

Picking an appropriate restaurant to dine in or order online delivery from can prove to be a challenging issue for some individuals since there are plenty of factors affecting one's decision, including the type of food provided, the environment of the restaurant, the tastefulness of the restaurant, etc. Acquiring pertinent information pertaining to these aspects is vital in facilitating informed choices. Online Opinion Platforms such as Zomato, TripAdvisor, and Yelp, provide customers an opportunity to express their opinions freely about the restaurants they have dined in [1]. These reviews could help potential customers to have better knowledge of prospective restaurants before dining in and therefore help them make a better decision. Moreover, after the outbreak of COVID-19, the interest and demand for food delivery services increased sharply [2]. With the increase in food delivery services and the advancement

in mobile technology, food delivery apps have become more popular. The recommendation systems that are widely used in such apps are essential as they are closely related to people's daily diets. The restaurants that appear on the customers' screen first should be tightly relevant to their potential decisions [3].

In general, the previous studies regarding restaurant recommendation mainly focus on proposing new algorithms or modifying existing recommendation algorithms in order to achieve a better-personalized recommendation result. For example, Gomathi et al. proposed a sentimental score measure Natural Language Processing (NLP) algorithm to examine user comments for hotel restaurants and compute the positive and negative percent of those comments and choose the restaurant with the highest ranking [4]. Furthermore, Chen et al. also proposed a context-aware recommendation system with different embedded feature selection methods to deal with the issue of data sparsity [5]. Other articles focus on combining additional background data with a recommendation system to produce a better user-based recommendation. In the study [6], Rajendran et al. combined browsing history and Wikipedia data with a recommendation system. They generated topics from Wikipedia by using the Latent Dirichlet Allocation model and analysed user preference from browsing history based on these topics. However, few previous studies have put the scope on the user review dataset themselves. Not only recommendation algorithms themselves are important, but datasets for user reviews and their related information might also play an important role in affecting recommendation decisions. If considering the ratings of a restaurant are enough to produce a satisfactory result, whether combining the ratings and the number of ratings or the number of dishes liked in a restaurant will bring an entirely different performance is worth studying. By taking a closer look at the dataset and by making use of the information that is given in the dataset already, for example by doing manipulation of features of the dataset, the present paper hopes to achieve a better result using a well-studied recommendation algorithm compares to using the same algorithm on the unmodified dataset.

This paper mainly focuses on manipulating the dimensionality of features to observe how the recommendation result is affected by the dimension of the feature space. The paper uses a dataset of Zomato that collected data in Bangalore city. Upper Confidence Bound Bandit Algorithm is used as the recommendation algorithm for picking the best restaurant. The paper modifies the dimensionality of features of the dataset by using Principal Component Analysis and other techniques. The result of the original dataset is compared with the results from the modified datasets to see if a significantly different regret is achieved. By performing dimension reduction, the total of 17 features of the dataset is first to be reduced to only two principal "features". By performing dimension raising, two selected features of the dataset (overall rating and different user rating) are multiplied together to produce a new feature. By comparing the result of applying the Upper Confidence Bound bandit algorithm (UCB) on the original dataset and the results on the modified datasets after dimension reduction and dimension raising, significant differences are shown in the performance between these three different scenarios. Both dimension reduction and dimension raising result in worse performance and dimension reduction results in an even worse performance compared to dimension raising.

2. Method

2.1 Data preprocessing

This paper employs a dataset from Zomato called "Zomato Bangalore Dataset" which was posted on Kaggle. This dataset records Zomato data for analyzing restaurant conditions in Bangalore City, consisting of a dimension of 51, 717*17, which means it contains 17 features and 51, 717 records of data in total.

In the current data preprocessing stage, the initial thing entails feature selection: wherein the identification and subsequent removal of features presumed inconsequential to model performance transpire. A feature called "reviews_list" in the dataset includes a list of reviews for each restaurant. Every review includes a rating and written feedback that each user gives for the restaurant. In this stage, the user ratings for each restaurant are extracted from the list of reviews and are stored as a new feature

“reward”. This feature “reward” is the core feature that is used to estimate the true average rating for each restaurant in the recommendation algorithm.

In addition, during this stage, the duplicate data and any row that has missing value are removed from the dataset. There are features in the dataset that are recorded in texts instead of integer values. The “Yes” and “No” in the features “online_order” and “book_table” are replaced by 1 and 0 correspondingly. Each restaurant has also been assigned a unique ID. Since some restaurants have the same name, the unique ID for each restaurant is assigned based on both its name and its location. Each restaurant also has one or more modes that the restaurant is running under (Café, Pub, Casual Dining, etc.). Since this feature is also recorded in text form, each different mode is created as a new column in the dataset. If a restaurant belongs to one or more specific modes, the corresponding cell will be assigned 1 otherwise 0. These are all preparation for dimension manipulation that will be done on the dataset later.

In order to reduce the number of restaurants that are compared to reach a better demonstration performance in the paper, only restaurants with more than 3000 reviews are kept in the dataset. In other words, 50 unique restaurants are left in the dataset.

2.2 Upper confidence bound algorithm

This paper mainly uses the basic UCB shown in Figure 1 as the model to choose the best restaurant [7].

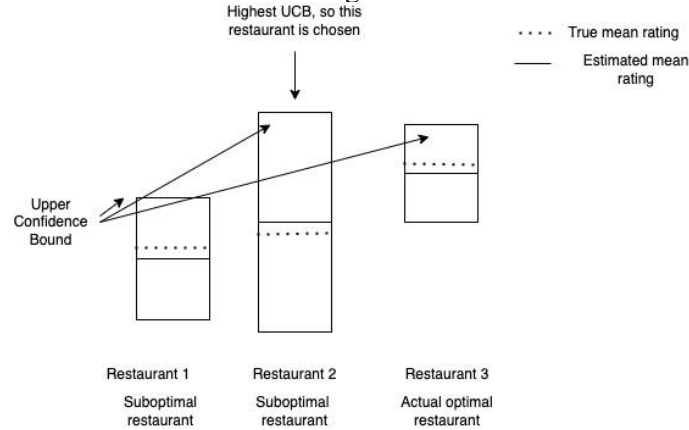


Figure 1. A demonstration of how UCB is built for different arms (restaurants) (Photo/Picture credit: Original).

The main logic of the model is to build an upper confidence bound based on the estimated average rating. At each round t , the restaurant with the highest UCB value will be chosen even though its estimated average rating or true mean rating might not be the highest. The UCB value for each arm i at each round t is calculated using equation (1):

$$UCB_{i,t} = \hat{\mu}_{i,t} + \left(\frac{\ln t}{n_{i,t}} \right)^{\frac{1}{2}} \quad (1)$$

where n is the number of times the arm i has been chosen until round t and $\hat{\mu}_i$ is the estimated mean reward for arm i . The estimated average rating is calculated by averaging n input user rating for the restaurant i .

In Figure 1, Restaurant 2 will be chosen even though the optimal restaurant is Restaurant 3. As one restaurant is chosen more times, the upper confidence bound of that restaurant’s estimated rating will become smaller and its estimated average rating will become closer to its actual average rating. Eventually, the UCB value for Restaurant 2 will be smaller than it is shown in Figure 1 and closer to its actual average rating. By that time, Restaurant 3 will be chosen repeatedly, and the model has found the actual optimal restaurant.

2.3. Dimension reduction using principal component analysis

This paper mainly uses Principal Component Analysis (PCA) to achieve dimension reduction. Principal Component Analysis is a dimension reduction technique proposed by Pearson and developed by Hotelling to deal with the curse of dimension [8]. The principal components in PCA refer to the linear combination of the original features [9]. The number of principal components can be set manually to different values. This paper prefers to explain variance for more than 90 percent. In order to achieve that, the number needs to be set to 2. In addition, fewer principal components mean that each principal component will preserve more information about the original dataset. Only one feature can be treated as the core feature in this work. Therefore, the case where each principal component preserves more information about the original dataset is favored. By considering these reasons, the number of principal components is set to 2. Since the principal components are ordered in descending order of importance, the paper chooses to use the first principal component in place of the original core feature as the measurement to calculate the estimated average reward in the modified dataset.

The present paper first brings all features to the same scale. In PCA, eigenvalues can be calculated using equation (2):

$$|A - \lambda I| = 0 \quad (2)$$

where A is an $n \times n$ matrix. After finding eigenvalues, the corresponding eigenvectors of A could be found using equation (3):

$$AX = \lambda X \quad (3)$$

The principal components in PCA are also projections of data instances onto the eigenvectors of the covariance matrix of data. The projection $\text{Proj}_{p_i}(\vec{\mu})$ can be done using equation (4):

$$\text{Proj}_{p_i}(\vec{\mu}) = \frac{p_i \cdot \vec{\mu}}{|\mu|} \quad (4)$$

Dimension reduction is achieved by only keeping these principal components that explain most of the variance in the data.

2.4. Dimension raising algorithm

Dimension raising is to increase the number of features (dimensions) in a dataset while obtaining more relevant information. In this paper, “reward” is considered to be the core feature that is used to estimate the true mean rating for each restaurant in the original dataset. According to Figure 2, “rate” is the feature that correlates with “reward” the most. The correlation of the other features to “reward” is not comparable with the correlation of “rate” to reward. Therefore, this work chooses to apply dimension raising on these two features.

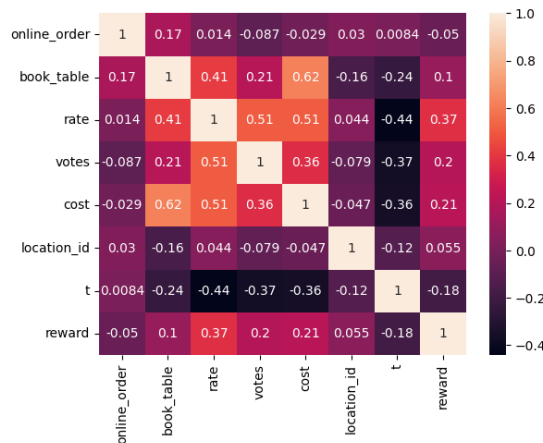


Figure 2. Heat map of correlation between each variable (Photo/Picture credit: Original).

After dimension raising, the two features x_1 “rate” and x_2 “reward” are raised to 5 features: $x_1, x_2, x_1^2, x_1x_2, x_2^2$ etc.). The current paper then chooses x_1x_2 as the new core feature used to calculate the estimated mean rating.

3. Results and discussion

To compare how UCB performs in three different cases (on the original dataset, the dataset after dimension reduction, and the dataset after dimension raising), cumulative regret is used to evaluate performance using equation (5):

$$R = E\left[\sum_{t=1}^T (r_{a^*} - r_{a,t})\right] \quad (5)$$

where r_{a^*} is the reward of the optimal action and $r_{a,t}$ is the reward of the action chosen at round t .

3.1. The performance of the models

Figure 3 demonstrated that the regret of the UCB on the unmodified dataset successfully reached logarithmic regret, which matches the expectation.

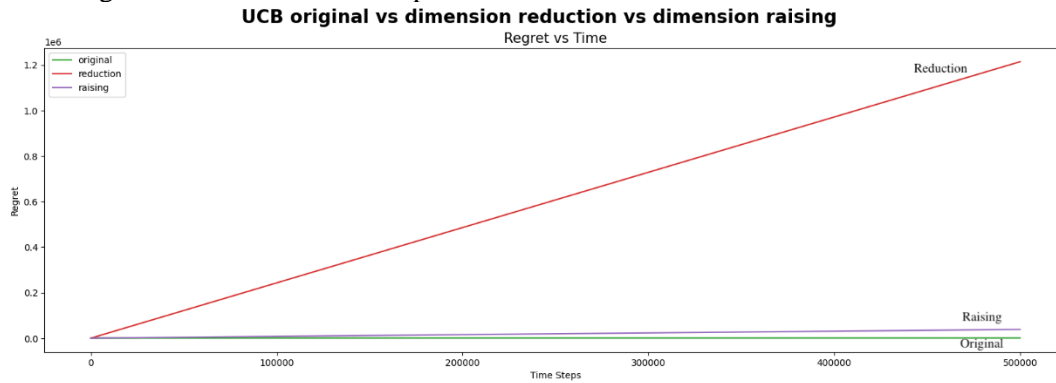


Figure 3. Cumulative regret of UCB in three different cases (Photo/Picture credit: Original).

However, when applying UCB on the dataset in which was modified by applying dimension reduction, and the other one was modified by applying dimension raising, it can be seen in Figure 3 that the regret is linear instead of logarithmic. The regret of dimension reduction is too different from the others, which makes it seems like the regret of the unmodified dataset and dimension raising are two close lines. This distinction is accentuated in the graphical representation presented in Figure 4, enabling a clear discrimination of performance nuances between the unmodified dataset and the dimension-raised counterpart. Furthermore, a linear growth trajectory characterizes the regret associated with UCB application on the dimension-raised dataset.

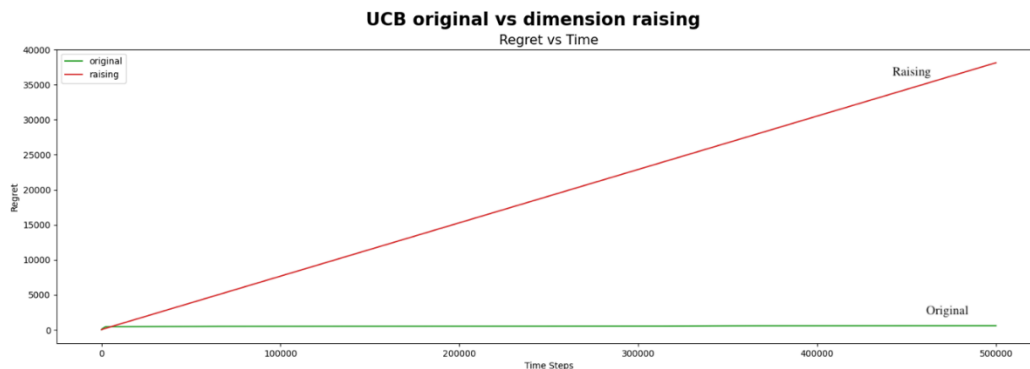


Figure 4. Cumulative regret of UCB on original dataset and dataset after dimension raising (Photo/Picture credit: Original).

From Figure 3 and Figure 4 it can be concluded that UCB performs worse on the modified dataset than the original dataset. It can also be concluded that the performance of UCB on the dataset whose feature space is reduced is worse than the performance of UCB on the dataset whose feature space is raised.

3.2. Discussion

Since the regret is growing linearly in both cases where the dataset is modified, it can be concluded that in neither modified dataset case the algorithm is able to find the best restaurant. The reason that the performance of UCB on the dataset after dimension reduction is performing worse might be that the features have lost their physical meaning when being reduced into only two principal components [10]. The reason for performance becoming worse on the dataset after dimension raising might follow the same logic. In addressing the comparatively favorable performance outcome of dimension raising in contrast to dimension reduction, a pertinent aspect to contemplate is the feature selection mechanism. Specifically, the dimension-raising methodology pursued in this study entails the incorporation of the core feature in conjunction with a secondary feature exhibiting the highest degree of correlation with said core feature. It is plausible that this established correlation nexus assumes a consequential role in influencing performance dynamics, potentially facilitating the retention of discernible portions of the original core feature's upward or downward trend, thereby contributing to the observed performance distinctions. In the future, it could also be considered to employ neural networks for extracting features after dimension reduction and investigate their impact on the performance of multi-armed bandit algorithms, given their remarkable performance in other tasks [11, 12].

4. Conclusion

This study deliberated upon the impact induced by dimension manipulation on the efficacy of the recommendation algorithm when applied to a dataset comprising restaurant reviews. Upper Confidence Bound was used as the recommendation algorithm and the dataset was modified in two different ways: dimension reduction and dimension raising. The performance for both ways of dimension manipulation was worse compared to the performance on the unmodified dataset. Dimension manipulation might overly cause the features to lose their physical meaning which hurt the regret performance of UCB. In addition, the performance of the dataset whose feature space was reduced was worse than the dataset whose feature space was raised. In the future, the performance of different degrees of dimension raising or dimension reduction can be evaluated to see whether these different ways will bring significant changes.

References

- [1] Meek S Wilk V and Lambert C 2021 A big data exploration of the informational and normative influences on the helpfulness of online restaurant reviews *Journal of Business Research* vol 125 (Elsevier) pp 354-367.
- [2] Sánchez C N Domínguez-Soberanes J Arreola A and Graff M 2023 Recommendation System for a Delivery Food Application Based on Number of Orders *Applied Sciences* vol 13 (MDPI) p 2299.
- [3] Jung H-H Yoon H-H and Song M-K 2021 A Study on Dining-Out Trends Using Big Data: Focusing on Changes since COVID-19 *Sustainability* vol 13 (MDPI) p 11489.
- [4] Gomathi R M Ajitha P Krishna G H S and Pranay I H 2019 Restaurant Recommendation System for User Preference and Services Based on Rating and Amenities *International Conference on Computational Intelligence in Data Science (IEEE)* pp 1-6.
- [5] Chen L and Xia M 2021 A context-aware recommendation approach based on feature selection *Applied Intelligence* vol 51 (Springer Nature) pp 865-875.
- [6] Rajendran D P D and Sundarraj R P 2021 Using topic models with browsing history in hybrid collaborative filtering recommender system: Experiments with user ratings *International Journal of Information Management Data Insights* vol 1 (Elsevier).

- [7] Auer P 2002 Using Confidence Bounds for Exploitation-Exploration Trade-offs *Journal of Machine Learning Research* vol 3, ed P M Long pp 397-422.
- [8] Jolliffe I T 2002 *Principal Component Analysis* (New York: Springer-Verlag).
- [9] Jolliffe I T and Cadima J 2016 Principal component analysis: a review and recent developments *Philosophical Transactions of the Royal Society A* (Royal Society Publishing).
- [10] Jia W Sun M Lian J and Hou S 2022 Feature dimensionality reduction: a review *Complex and Intelligent Systems* vol 8 (Springer Nature) pp 2663-2693.
- [11] Yu Q Wang J Jin Z et al 2022 Pose-guided matching based on deep learning for assessing quality of action on rehabilitation training *Biomedical Signal Processing and Control* 72: 103323.
- [12] Malhotra P Gupta S Koundal D et al 2022 Deep neural networks for medical image segmentation *Journal of Healthcare Engineering* 2022.