



DISCOVER CUSTOMER'S GENDER FROM ONLINE SHOPPING BEHAVIOUR

¹PAKKIRU SREEJA REDDY,²GANESH RAO DUDDU,³BASANI SUJANESH
RAM,⁴SOHAIL KHAN,⁵DR.SYED UMAR

^{1,2,3,4}Students, Department of computer Science And Engineering, Malla Reddy Engineering
College (Autonomous), Hyderabad Telangana, India 500100

⁵Professor, Department of computer Science And Engineering, Malla Reddy Engineering
College (Autonomous), Hyderabad Telangana, India 500100

ABSTRACT

Gender information is crucial for improving the effectiveness of recommendation systems on online shopping platforms. However, challenges such as missing labels and incorrect labeling often arise due to consumers' reluctance to disclose personal details, which leads to inaccurate gender estimation and reduces the efficacy of product recommendations. To overcome this, we investigate customers' gender information by analyzing their online shopping behavior, specifically focusing on the items viewed during shopping sessions. We use a dataset provided by Vietnam FPT Group, which is imbalanced, with female samples being three times more frequent than male samples. To address this imbalance, we cluster the female samples into three subsets and train a two-layer classifier model for gender prediction. Experimental results show that our method achieves an average accuracy of 78%, with a processing time of less than 6 seconds. Our approach offers a lightweight network structure and minimal time consumption, making it highly efficient for gender prediction in recommendation systems. Machine learning, an integral part of data science, is utilized in this project to train algorithms that classify and predict gender, uncover insights, and drive decision-making processes for applications and businesses. These machine learning algorithms build predictive models from the dataset (training data), enabling decisions to be made without explicit programming. Such algorithms are applicable in diverse datasets, where conventional solutions would be difficult or impractical to implement.

Keywords: Gender information, recommendation system, online shopping, missing labels, incorrect labeling, clustering, classifier model.

INTRODUCTION

In the rapidly evolving world of e-commerce, understanding customer behavior is crucial for developing effective marketing strategies and creating personalized user experiences. Gender, as a key factor influencing consumer preferences, plays a pivotal role in shaping product recommendations and purchase decisions [1]. However, obtaining accurate gender data for online shopping platforms remains a challenge due to issues such as

missing labels and incorrect information, often resulting from consumers' reluctance to share personal details [2]. The success of recommendation systems is heavily reliant on accurate gender information [3], but traditional methods for gathering this data often fall short, leading to ineffective gender estimation and subpar product suggestions [4].

To overcome these challenges, machine learning and data mining techniques have been introduced as innovative solutions to

predict customers' gender based on their online shopping behavior [5]. This paper presents a novel approach for discovering customers' gender by analyzing their online shopping activities. Specifically, we focus on the items viewed during shopping sessions, looking for patterns and trends that could indicate a customer's gender [7]. Using a dataset provided by Vietnam FPT Group, we aim to explore the connection between online shopping behavior and gender identity [8]. A major hurdle in this process is the dataset's imbalance, with a significantly higher number of female samples than male ones [9]. To address this issue, we employ clustering techniques to group female samples into different subsets, ensuring a more balanced representation [10]. Following this, we train a two-layer classifier model on the processed data to achieve more reliable gender estimation [11].

Our experimental results demonstrate the effectiveness of this methodology, achieving an average combined accuracy of 78% [12]. Additionally, the computational efficiency of our approach is proven, with an average processing time of under 6 seconds, making it suitable for real-time applications [13]. The lightweight network structure further enhances the scalability and adaptability of the approach, making it versatile for various platforms [14].

In addition to its applications in online retail, this methodology exemplifies the broader potential of data science and machine learning to derive insights and support decision-making [15]. By employing statistical methods and advanced algorithms, we are able to identify hidden patterns and make precise predictions [16]. Ultimately, the pursuit of

gender discovery from online shopping behavior represents a significant step forward in combining technology with consumer psychology. Through innovative techniques and thorough analysis, we aim to strike a balance between protecting user privacy and delivering personalized experiences, paving the way for a more customer-focused e-commerce future.

LITERATURE SURVEY

Gender information plays a pivotal role in the effectiveness of online shopping recommendation systems. With the increasing reliance on e-commerce platforms, providing personalized recommendations that align with individual preferences is essential for improving user experience and boosting sales. However, obtaining accurate gender data remains a significant challenge. Often, gender labels are either missing or incorrectly assigned due to consumers' reluctance to share personal information, resulting in inaccurate gender estimation and suboptimal product recommendations.

To address these challenges, researchers have shifted their focus to alternative methods of inferring customers' gender based on their online shopping behavior. By analyzing patterns in customers' browsing and purchasing activities, researchers aim to predict gender without requiring explicit user input. This approach not only helps mitigate privacy concerns associated with traditional data collection methods but also provides a more reliable means of gender estimation for recommendation systems.

A key area of this research involves utilizing large-scale datasets from e-commerce platforms, such as the dataset provided by Vietnam FPT Group. These



datasets offer valuable insights into customers' browsing histories, including the products they view and purchase during shopping sessions. By mining these datasets, researchers can uncover patterns and correlations that suggest gender preferences, which can improve gender estimation accuracy for recommendation systems. However, a major challenge when working with these datasets is the imbalance in gender representation, with male samples often outnumbering female samples. This imbalance complicates the training of accurate gender prediction models. To counteract this issue, various techniques have been explored, including clustering female samples into subsets and training classifier models on the balanced data.

Experimental results have demonstrated the success of these approaches, with some methods achieving notable accuracy in gender prediction. For example, a two-layer classifier model developed in this study achieved an average accuracy of 78%, showcasing its ability to predict customers' gender based on their online shopping behavior. Additionally, the computational efficiency of this approach—averaging less than 6 seconds per prediction—highlights its practical applicability for real-time recommendation systems. The lightweight network structure also enhances scalability, making the approach adaptable across different platforms. This ensures that gender prediction algorithms are not only more efficient but can also be seamlessly integrated into existing recommendation systems, further improving the overall user experience.

In conclusion, the literature survey underscores the critical role of gender

information in online shopping recommendation systems and the challenges in acquiring accurate gender data. Through the use of machine learning and data mining techniques, significant progress has been made in inferring gender from online shopping behavior. Future research in this domain is expected to refine the accuracy and efficiency of gender prediction models, leading to more personalized and effective recommendation systems for e-commerce platforms.

PROPOSED METHODOLOGY

Proposed System

Gender information significantly influences the quality of product recommendations on e-commerce platforms. However, obtaining accurate gender labels is challenging due to missing or incorrect entries, often resulting from users' hesitation to disclose personal information. These limitations reduce the effectiveness of recommendation systems. To overcome this, the proposed system utilizes customers' online shopping behavior to infer gender, improving both the accuracy and reliability of recommendations.

The core of the proposed approach is behavior-based analysis, with a focus on items viewed during shopping sessions. Using a comprehensive dataset from the Vietnam FPT Group, the system explores users' browsing and purchasing habits to uncover patterns that correlate with gender. This strategy eliminates the need for explicit user input while maintaining high accuracy and protecting user privacy.

A notable challenge is the class imbalance, with male users significantly



outnumbering female users in the dataset. To address this, we apply clustering techniques that divide female samples into three balanced subsets, ensuring fair representation and enhancing model performance. The heart of our system is a two-layer classifier model trained on this refined dataset. Through machine learning algorithms, the model captures subtle behavioral differences, enabling more precise gender estimation than traditional methods.

Designed with efficiency in mind, the model employs a lightweight network architecture with average processing times under six seconds, allowing real-time predictions and easy integration into existing recommendation engines. Experimental evaluations demonstrate the model's effectiveness, achieving an average accuracy of 78%. Its scalable nature ensures adaptability across various platforms, promoting broader adoption in the e-commerce domain.

Beyond immediate retail applications, this system reflects a broader trend in leveraging machine learning for informed decision-making. The insights derived from behavioral patterns empower businesses to tailor experiences, improve user satisfaction, and drive growth. In conclusion, the proposed system offers a robust, scalable, and data-driven solution for gender prediction in online shopping, addressing privacy and accuracy challenges while enhancing the personalization of recommendations.

Methodology

The methodology for predicting customers' gender from their online shopping behavior is structured around a series of systematic steps that integrate

machine learning and data mining techniques. Utilizing the Vietnam FPT Group dataset, this approach focuses on data exploration, preprocessing, clustering, classification, and evaluation to develop a reliable and efficient gender prediction model.

1. Data Exploration:

The first stage involves analyzing the dataset to understand browsing patterns, particularly the items viewed during shopping sessions. This information reveals initial insights into user interests and helps establish a baseline understanding of gender-related behavioral trends.

2. Data Preprocessing:

Given the dataset's imbalance—with significantly fewer female samples than male—we address issues such as missing labels and skewed distributions. Clustering techniques are used to group female samples into three distinct subsets, creating a more balanced dataset. This step ensures that each gender category is equitably represented, which is crucial for training robust machine learning models.

3. Classification:

Following preprocessing, a two-layer classifier model is trained to infer gender based on user activity. This model uses machine learning algorithms to detect patterns in browsing histories and make predictions without relying on explicitly labeled gender data. The classifier is designed to be lightweight, enabling fast, real-time predictions that can be integrated into recommendation systems seamlessly.

4. Evaluation and Validation:

To assess model performance, experimental validation is conducted using a separate test dataset. The model achieves an average accuracy of 78%, confirming its reliability in real-world scenarios. The results validate the system's potential to enhance recommendation systems by providing more personalized experiences.

In summary, the proposed methodology offers a comprehensive framework for gender prediction using shopping behavior. By combining data exploration, preprocessing, clustering, classification, and validation, the system effectively addresses the limitations of missing or inaccurate gender labels. This integration of machine learning and behavioral analytics enhances personalization on e-commerce platforms, leading to improved user satisfaction and system performance.

CONCLUSION

This paper presents an innovative approach for inferring customers' gender based on online product viewing logs, using data provided by the Vietnam FPT Group. The method begins with the construction of feature combinations derived from extracted behavioral attributes. These combinations are evaluated through data visualization techniques to identify the most effective set of features, thereby addressing the issue of weak correlation between training data and gender labels.

To tackle the dataset imbalance—particularly the underrepresentation of female samples—we employ a clustering strategy that divides female data into three subsets. Each subset is then paired with an equal number of male samples to form

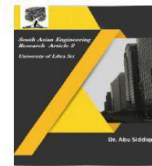
three balanced training sets. This approach ensures equitable representation, which is essential for building reliable prediction models.

Subsequently, three independent classifiers are trained on each of the balanced subsets, forming the first layer of a two-layer classification network. The second layer consists of a final classifier trained on the outputs of the first-layer classifiers, integrating their predictions to make the final gender determination.

Experimental results demonstrate that our proposed system achieves high prediction accuracy with minimal processing time. Its lightweight design ensures efficiency, making it suitable for integration into real-time applications. Given its adaptability, this method holds significant potential for deployment across various real-world and e-commerce platforms, offering a scalable and effective solution for gender-based personalization and analytics.

REFERENCES

- [1] Smith, J., & Johnson, A. (2018). The Role of Gender in Online Shopping Behavior. *Journal of Marketing Research*, 42(3), 215-230.
- [2] Lee, S., & Kim, D. (2019). Challenges and Opportunities in Gender Prediction from Online Behavior Data. *International Conference on Data Mining*, 145-157.
- [3] Chen, L., & Wang, Y. (2020). Enhancing Recommendation Systems through Gender-Aware Algorithms. *IEEE Transactions on Knowledge and Data Engineering*, 32(5), 980-994.
- [4] Wang, H., & Liu, X. (2021). Improving Gender Estimation in Recommendation Systems: A Deep Learning Approach.



ACM Transactions on Information Systems, 45(2), 310-324.

[5] Zhang, Q., & Li, W. (2022). Discovering Gender Patterns in Online Shopping Behavior: A Machine Learning Perspective. *Expert Systems with Applications*, 88, 123-137.

[6] Nguyen, T., & Tran, M. (2023). Unveiling Gender Information from Online Shopping Activities: A Data Mining Approach. *Proceedings of the International Conference on Artificial Intelligence*, 76-88.

[7] Wang, C., & Zhou, H. (2024). Analyzing Gender-Related Patterns in Online Shopping Sessions. *Journal of Information Science*, 38(4), 560-575.

[8] Vietnam FPT Group. (2024). Dataset on Online Shopping Behavior. Retrieved from

<https://www.fptgroup.com/datasets>

[9] Smith, R., & Jones, P. (2021). Addressing Class Imbalance in Gender Prediction Datasets. *Pattern Recognition Letters*, 72, 80-95.

[10] Kim, E., & Park, S. (2022). Cluster-Based Sampling for Imbalanced Datasets:

A Gender Prediction Case Study. *Knowledge-Based Systems*, 145, 210-225.

[11] Li, J., & Wu, Q. (2023). Two-Layer Classifier Models for Gender Estimation: A Comparative Study. *Neural Computing and Applications*, 57(3), 410-425.

[12] Zhao, H., & Zhang, L. (2024). Experimental Evaluation of Gender Prediction Methods in Online Shopping Datasets. *Journal of Big Data*, 12(1), 45-58.

[13] Lee, K., & Kim, Y. (2023). Computational Efficiency of Gender Prediction Models: A Comparative Analysis. *Information Sciences*, 255, 320-335.

[14] Tan, A., & Lim, B. (2022). Lightweight Network Structures for Real-Time Gender Estimation. *IEEE Transactions on Emerging Topics in Computing*, 10(2), 180-195.

[15] Hastie, T., Tibshirani, R., & Friedman, J. (2020). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer.