

Neural Contextual Bandits for Personalized Recommendation

Yikun Ban
University of Illinois at
Urbana-Champaign
yikunb2@illinois.edu

Yunzhe Qi
University of Illinois at
Urbana-Champaign
yunzheq2@illinois.edu

Jingrui He
University of Illinois at
Urbana-Champaign
jingrui@illinois.edu

ABSTRACT

In the dynamic landscape of online businesses, recommender systems are pivotal in enhancing user experiences. While traditional approaches have relied on static supervised learning, the quest for adaptive, user-centric recommendations has led to the emergence of the formulation of contextual bandits. This tutorial investigates the contextual bandits as a powerful framework for personalized recommendations. We delve into the challenges, advanced algorithms and theories, collaborative strategies, and open challenges and future prospects within this field. Different from existing related tutorials, (1) we focus on the exploration perspective of contextual bandits to alleviate the “Matthew Effect” in the recommender systems, i.e., the rich get richer and the poor get poorer, concerning the popularity of items; (2) in addition to the conventional linear contextual bandits, we will also dedicated to neural contextual bandits which have emerged as an important branch in recent years, to investigate how neural networks benefit contextual bandits for personalized recommendation both empirically and theoretically; (3) we will cover the latest topic, collaborative neural contextual bandits, to incorporate both user heterogeneity and user correlations customized for recommender system; (4) we will provide and discuss the new emerging challenges and open questions for neural contextual bandits with applications in the personalized recommendation, especially for large neural models.

Compared with other greedy personalized recommendation approaches, Contextual Bandits techniques provide distinct ways of modeling user preferences. We believe this tutorial can benefit researchers and practitioners by appreciating the power of exploration and the performance guarantee brought by neural contextual bandits, as well as rethinking the challenges caused by the increasing complexity of neural models and the magnitude of data.

1 INTRODUCTION

Recommender systems are indispensable in various online businesses, including e-commerce platforms and online streaming services. They leverage user correlations to assist the perception of user preferences, a field of study spanning several decades. In the past, considerable effort has been directed toward supervised-learning-based collaborative filtering methods within relatively static environments [20, 33]. However, the ideal recommender systems should adapt over time to consistently meet user interests. Consequently, it is natural to formulate the recommendation process as a sequential decision-making process. In this paradigm, the recommender

engages with users, observes their online feedback (i.e., rewards), and optimizes the user experience for long-term benefits, rather than fitting a model on the collected static data based on supervised learning [12, 17, 36]. Based on this idea, this tutorial focuses on the formulation of contextual bandits [1, 4, 6, 10, 18, 19, 25, 26, 28–30]. In particular, the emerging neural bandit techniques [5, 7, 9, 13, 21, 22, 31, 38, 39] have shown their superiority over linear or kernel [1, 3, 15, 27, 34] methods due to the representation power of neural networks. As a result, they have been considered the ideal solution to deal with the non-linear problem settings, in terms of the Contextual Bandits research and application.

To be specific, let us dive into a scenario involving a total of T rounds of recommendations. In the t -th round, the learner (model) is presented with K arms (items), each of which yields a reward based on an unknown function that captures the users’ preferences. The learner’s task is to select one arm in each round and recommend this arm to the serving user. Consequently, the learner observes the resulting rewards and updates its recommendation policy. The ultimate goal is to maximize the cumulative rewards across the T rounds, i.e., to minimize the cumulative regrets incurred over the course of these interactions. However, the fundamental challenge of balancing exploitation and exploration inherently arises in the context of sequential recommendation. In a round, the learner faces a dilemma. On one hand, it must ‘exploit’ the knowledge gleaned from the previous rounds to prioritize popular items. On the other hand, the model is also tasked with ‘exploring’ potential value from new or under-explored items in seek of long-term benefits. As a result, striking the right balance between exploitation and exploration is imperative for constructing a robust recommender system and contextual bandits provide powerful tools to tackle this dilemma.

It is well known that research works on “User Modeling and Recommendation” have always been an indispensable part of the Web Conference. Different from other greedy personalized recommendation approaches (e.g., the conventional collaborative filtering methods [20, 33]), Contextual Bandits techniques provide distinct ways of modeling user preferences [6, 18, 19, 23, 25, 26, 35, 37], by balancing exploiting the process knowledge and exploring the unrevealed benefits. Therefore, all of these aspects make this tutorial a good fit for the Web Conference and its community members.

In this lecture-style tutorial, we will provide a comprehensive review of contextual bandits algorithms under the personalized recommendation settings. We will begin by delving into the challenges of personalized recommendation and why contextual bandits emerge as a powerful tool. The subsequent section will provide an in-depth review of advanced algorithms and theories in

This paper is published under the Creative Commons Attribution 4.0 International (CC-BY 4.0) license. Authors reserve their rights to disseminate the work on their personal and corporate Web sites with the appropriate attribution.

WWW ’24, May 13–17, 2024, Singapore

© 2024 IW3C2 (International World Wide Web Conference Committee), published under Creative Commons CC-BY 4.0 License.

contextual bandits, ranging from linear to neural contextual bandits, emphasizing their role in enhancing exploration in personalized recommendation. The third part is the collaborative contextual bandits, extending beyond individual bandits to account for user correlations. We will explore two intriguing problem formulations: clustering of bandits, which is to formulate the scenarios in which we leverage the knowledge of a cluster of users to facilitate the decision-making of the serving user, and graph bandits learning, which is to formulate fine-grained user correlations represented by the strength of edges in the constructed graph. In conclusion, we will shed light on the open challenges and future directions in contextual bandits, particularly within personalized recommendation scenarios. Key areas of focus will include the scalability of contextual bandits in large recommender systems and ensuring the trustworthiness of personalized decision-making. This tutorial promises to be a captivating review and lecture on the evolving contextual bandits that holds immense potential for revolutionizing how we perform content recommendation to users.

2 TARGET AUDIENCE

The tutorial is developed for all the people who are interested in the related research areas, e.g., multi-armed bandits, reinforcement learning, information retrieval, data mining, recommender systems, etc. The audience is expected to have basic knowledge of machine learning and data mining. This tutorial balances the introductory and advanced material, i.e., 50 % for beginners and 50 % for intermediate and advanced researchers.

In addition, as this tutorial will be an organic combination of (1) the application of bandit techniques in recommendation, as well as (2) the theoretical insights of bandit algorithms, we believe that both the business practitioner as well as researchers focusing on algorithmic aspects will benefit from our contents.

3 SHORT BIO

In this section, we briefly introduce the three presenters of this tutorial.

Yikun Ban is a final-year Ph.D. student in the Department of Computer Science at the University of Illinois at Urbana-Champaign. He is a member of DAIS (Data and Information Systems) Research Lab. He received his M.CS. degree from Peking University in 2019 and B.Eng. degree from Wuhan University in 2016. His research interests lie in multi-armed bandits/Reinforcement Learning to design and develop principled exploration strategies in sequential decision-making. He has published more than 11 papers at top conferences in Machine Learning and Data Mining (e.g., WWW, KDD, NeurIPS, ICLR, AAAI) and has been a reviewer or program committee member of mainstream machine learning journals and conferences. He was an applied scientist intern at Amazon Web Service, and his research works have been powering primary applications in Amazon and Instacart.

Yunzhe Qi is a Ph.D. candidate in the School of Information Sciences at the University of Illinois at Urbana-Champaign (UIUC). He received his Master's degree from UIUC in 2021, and his Bachelor's degree in Beijing University of Posts and Telecommunications in 2019 respectively. His research interests mainly focus on utilizing Contextual Bandit methods to solve the exploitation-exploration

dilemma for machine learning tasks, such as online recommendation. He has published several papers at top machine learning / data mining conferences (e.g., KDD, NeurIPS), and has been serving as the reviewer as well as PC member for multiple machine learning / data mining conferences and journals. He was a Machine Learning Engineer Intern at Instacart, who designed and implemented Contextual Bandit frameworks for personalized recommendation that have been generating actual business growth.

Jingrui He (Corresponding Tutor) is a professor in the School of Information Sciences at the University of Illinois Urbana-Champaign. She received her Ph.D. from Carnegie Mellon University in 2010. Her research focuses on heterogeneous machine learning, rare category analysis, active learning and semi-supervised learning, with applications in security, social network analysis, healthcare, and manufacturing processes. Dr. He is the recipient of the 2016 NSF CAREER Award, the 2020 OAT Award, three-time recipient of the IBM Faculty Award in 2018, 2015, and 2014, and was selected as IJCAI 2017 Early Career Spotlight. Dr. He has more than 100 publications at major conferences (e.g., WWW, IJCAI, AAAI, KDD, ICML, NeurIPS) and journals (e.g., TKDE, TKDD, DMKD), and is the author of two books. Her papers have received the Distinguished Paper Award at FAccT 2022, as well as Bests of the Conference at ICDM 2016, ICDM 2010, and SDM 2010. She has several years of course teaching experience as an instructor and has offered several tutorials at major conferences, e.g., KDD, AAAI, IJCAI, SDM, IEEE BigData, etc. in the past few years. For more information, please refer to her homepage at <https://www.hejingrui.org/>.

4 OUTLINE

This will be a **3-hour lecture-style** tutorial to cover the state-of-the-art research for neural contextual bandit algorithms and theories with applications in personalized recommendation.

• Introduction

Duration: 15 minutes, *Presenter:* Yikun Ban

We start by motivating the formulation of contextual bandits for personalized recommendation and its existing challenges.

- Motivation
- Formulation of Sequential Decision-Making [32]
- Challenges

• Part I: Linear Contextual Bandits

Duration: 30 minutes, *Presenter:* Yikun Ban

In this part, we introduce the problem definition of linear contextual bandits and the existing linear exploration strategy with applications in personalized recommendation.

- Linear Upper Confidence Bound (UCB) [1]
- Linear Thompson Sampling (TS) [2]
- Linear Personalized Recommendation [24]

• Part II: Neural Contextual Bandits

Duration: 40 minutes, *Presenter:* Yunzhe Qi and Yikun Ban

This is the core part. We introduce the problem definition, latest algorithms, and theoretical guarantee for neural contextual bandits, with elaboration on how the neural networks benefit contextual bandits for personalized recommendation.

- Neural Upper Confidence Bound (UCB) [39]
- Neural Thompson Sampling (TS) [38]

- Exploitation-Exploration Networks [9, 11]
- Neural Inverse Gap Weighting [14, 16]
- Neural Personalized Recommendation

• Part III Collaborative Contextual Bandits

Duration: 60 minutes, *Presenter:* Yunzhe Qi and Yikun Ban
This is the core part. We introduce the motivation of collaborative Contextual Bandits and focus on two problems: clustering of bandits and graph bandits learning.

- Motivation and Challenges
- Clustering of Bandits
 - * Linear Clustering of Bandits [6, 19]
 - * Neural Clustering of Bandits [8]
- Graph Bandits Learning [29, 30]
- Applications in Recommender Systems [7]

• Part IV: Open Questions and Future Trends

Duration: 35 minutes, *Presenter:* Yikun Ban and Jingrui He
In this part, we will discuss the emerging challenges and open questions for neural contextual bandits in applications of recommender systems.

- Large Search Space: Arm and User Space
- Transparency: Rationales and Explanation for Exploration
- Fairness: Exploit-Explore Fairly
- Privacy: Privacy-preserving Decision-making

5 RELATED TUTORIALS OR TALKS

• *Practical Bandits: An Industry Perspective*

The WebConf 2023 Tutorial

Event date: May 3, 2023. **Location:** Austin, Texas, USA

Differences: In this WebConf tutorial, they present bandits algorithms from a practical perspective to facilitate practitioners to determine non-contextual or contextual approaches, on- or off-policy setups, delayed or immediate feedback, etc. Instead, our tutorial focuses on the neural contextual bandits that have attained increasing attention recently year from both practical and theoretical aspects, and our backbone is for SOTA exploration strategies with applications in the recommender system.

• *Bridging Learning and Decision Making*

The ICML 2022 Tutorial

Event Date: July 18, 2022. **Location:** Baltimore, Maryland, USA

Difference: This ICML tutorial gives an overview of the theoretical foundations of contextual bandits and reinforcement learning. In particular, they focus on the unified bandit framework to build the connection between supervised estimation and sequential decision-making. However, they lack the coverage of neural contextual bandits and collaborative contextual bandits tailored for the recommender system, which is the main gap we try to fill in this tutorial.

6 PREVIOUS EDITIONS

This will be the first edition of this tutorial, but the presenters have experience in teaching material covering similar topics in the past. We anticipate to present (an extended version of) this tutorial at similar conferences in the future.

7 TUTORIAL MATERIALS

We will set up a website to release all the related materials, including presentation slides, references, and open-source data & code, in order to get the audience familiarized with the tutorial content before the tutorial begins. Besides, all the researchers and practitioners are encouraged to have Q&A interaction during the tutorial or further discussions offline.

8 LINK TO VIDEO TEASER

The video teaser file can be found at the following Dropbox link:

<https://www.dropbox.com/scl/fo/32idrzywjz6gc9qp3iixi/h?rlkey=0xagsm4m4hbf>

REFERENCES

- [1] Y. Abbasi-Yadkori, D. Pál, and C. Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320, 2011.
- [2] S. Agrawal and N. Goyal. Thompson sampling for contextual bandits with linear payoffs. In *International Conference on Machine Learning*, pages 127–135. PMLR, 2013.
- [3] P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- [4] Y. Ban and J. He. Generic outlier detection in multi-armed bandit. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 913–923, 2020.
- [5] Y. Ban and J. He. Convolutional neural bandit: Provable algorithm for visual-aware advertising. *arXiv preprint arXiv:2107.07438*, 2021.
- [6] Y. Ban and J. He. Local clustering in contextual multi-armed bandits. In *Proceedings of the Web Conference 2021*, pages 2335–2346, 2021.
- [7] Y. Ban, J. He, and C. B. Cook. Multi-facet contextual bandits: A neural network perspective. In *The 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, Singapore, August 14-18, 2021*, pages 35–45, 2021.
- [8] Y. Ban, Y. Qi, T. Wei, and J. He. Neural collaborative filtering bandits via meta learning. *arXiv preprint arXiv:2201.13395*, 2022.
- [9] Y. Ban, Y. Yan, A. Banerjee, and J. He. EE-net: Exploitation-exploration neural networks in contextual bandits. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=X_ch3VrNSRg.
- [10] Y. Ban, Y. Zhang, H. Tong, A. Banerjee, and J. He. Improved algorithms for neural active learning. *Advances in Neural Information Processing Systems*, 35: 27497–27509, 2022.
- [11] Y. Ban, Y. Yan, A. Banerjee, and J. He. Neural exploitation and exploration of contextual bandits. *arXiv preprint arXiv:2305.03784*, 2023.
- [12] M. Chen, C. Xu, V. Gatto, D. Jain, A. Kumar, and E. Chi. Off-policy actor-critic for recommender systems. In *Proceedings of the 16th ACM Conference on Recommender Systems*, pages 338–349, 2022.
- [13] Z. Dai, Y. Shu, A. Verma, F. X. Fan, B. K. H. Low, and P. Jaillet. Federated neural bandit. *arXiv preprint arXiv:2205.14309*, 2022.
- [14] R. Deb, Y. Ban, S. Zuo, J. He, and A. Banerjee. Contextual bandits with online neural regression. *arXiv preprint arXiv:2312.07145*, 2023.
- [15] A. A. Deshmukh, U. Dogan, and C. Scott. Multi-task learning for contextual bandits. In *Advances in neural information processing systems*, pages 4848–4856, 2017.
- [16] D. Foster and A. Rakhlin. Beyond ucb: Optimal and efficient contextual bandits with regression oracles. In *International Conference on Machine Learning*, pages 3199–3210. PMLR, 2020.
- [17] C. Gao, K. Huang, J. Chen, Y. Zhang, B. Li, P. Jiang, S. Wang, Z. Zhang, and X. He. Alleviating matthew effect of offline reinforcement learning in interactive recommendation. *arXiv preprint arXiv:2307.04571*, 2023.
- [18] C. Gentile, S. Li, and G. Zappella. Online clustering of bandits. In *International Conference on Machine Learning*, pages 757–765, 2014.
- [19] C. Gentile, S. Li, P. Kar, A. Karatzoglou, G. Zappella, and E. Etrud. On context-dependent clustering of bandits. In *Proceedings of the 34th International Conference on Machine Learning—Volume 70*, pages 1253–1262. JMLR. org, 2017.
- [20] X. He, L. Liao, H. Zhang, L. Nie, X. Hu, and T.-S. Chua. Neural collaborative filtering. In *Proceedings of the 26th international conference on world wide web*, pages 173–182, 2017.
- [21] Y. Jia, W. Zhang, D. Zhou, Q. Gu, and H. Wang. Learning neural contextual bandits through perturbed rewards. In *International Conference on Learning Representations*, 2022.
- [22] P. Kassarai and A. Krause. Neural contextual bandits without regret. In *International Conference on Artificial Intelligence and Statistics*, pages 240–278. PMLR, 2022.

- [23] N. Korda, B. Szörényi, and L. Shuai. Distributed clustering of linear bandits in peer to peer networks. In *Journal of machine learning research workshop and conference proceedings*, volume 48, pages 1301–1309. International Machine Learning Society, 2016.
- [24] L. Li, W. Chu, J. Langford, and R. E. Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670, 2010.
- [25] S. Li, A. Karatzoglou, and C. Gentile. Collaborative filtering bandits. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, pages 539–548, 2016.
- [26] S. Li, W. Chen, S. Li, and K.-S. Leung. Improved algorithm on online clustering of bandits. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pages 2923–2929. AAAI Press, 2019.
- [27] W. Li, X. Wang, R. Zhang, Y. Cui, J. Mao, and R. Jin. Exploitation and exploration in a performance based contextual advertising system. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 27–36, 2010.
- [28] T. M. McDonald, L. Maystre, M. Lalmas, D. Russo, and K. Ciosek. Impatient bandits: Optimizing recommendations for the long-term without delay. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1687–1697, 2023.
- [29] Y. Qi, Y. Ban, and J. He. Neural bandit with arm group graph. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1379–1389, 2022.
- [30] Y. Qi, Y. Ban, and J. He. Graph neural bandits. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1920–1931, 2023.
- [31] Y. Qi, Y. Ban, T. Wei, J. Zou, H. Yao, and J. He. Meta-learning with neural bandit scheduler. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [32] A. Slivkins et al. Introduction to multi-armed bandits. *Foundations and Trends® in Machine Learning*, 12(1-2):1–286, 2019.
- [33] X. Su and T. M. Khoshgoftaar. A survey of collaborative filtering techniques. *Advances in artificial intelligence*, 2009, 2009.
- [34] M. Valko, N. Korda, R. Munos, I. Flaounas, and N. Cristianini. Finite-time analysis of kernelised contextual bandits. *arXiv preprint arXiv:1309.6869*, 2013.
- [35] J. Wu, C. Zhao, T. Yu, J. Li, and S. Li. Clustering of conversational bandits for user preference learning and elicitation. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pages 2129–2139, 2021.
- [36] W. Xue, Q. Cai, R. Zhan, D. Zheng, P. Jiang, and B. An. Resact: Reinforcing long-term engagement in sequential recommendation with residual actor. *arXiv preprint arXiv:2206.02620*, 2022.
- [37] L. Yang, B. Liu, L. Lin, F. Xia, K. Chen, and Q. Yang. Exploring clustering of bandits for online recommendation system. In *Fourteenth ACM Conference on Recommender Systems*, pages 120–129, 2020.
- [38] W. Zhang, D. Zhou, L. Li, and Q. Gu. Neural thompson sampling. In *International Conference on Learning Representations*, 2021.
- [39] D. Zhou, L. Li, and Q. Gu. Neural contextual bandits with ucb-based exploration. In *International Conference on Machine Learning*, pages 11492–11502. PMLR, 2020.