



Unofficial template Master SSE Thesis

by

Studentname Middlename Lastname

GRADUATION REPORT

Submitted to

Hanze University of Applied Sciences

in partial fulfillment of the requirements
for the degree of

Master Smart Systems Engineering

Confidential document

Groningen

2024

ABSTRACT

Insert your abstract here. Original Hanze format states one can use a maximum of 200 words. Though, it is advised is to discuss with your supervisor(s) if you need more.

Keywords: Keyword, Keyword, Keyword.

Contents

1	Introduction	7
2	Literature Review	8
2.1	Explainable AI In Knowledge Graph Recommender System	8
2.1.1	Path-based Modeling	8
2.1.2	Integration of User Profiles and Behavior	9
2.1.3	Explainability through Symbolic and Logical Reasoning	9
2.2	Counterfactual Methods in Different Types of Recommender Systems . . .	10
2.3	Application of Counterfactuals for Fairness and Bias Mitigation	12
2.3.1	Mitigating Bias Across Different AI Frameworks	12
2.3.2	Enriching Data and Ensuring Equitable Outcomes	13
2.3.3	Leveraging Knowledge Graphs for Fair Recommendations	13
3	Research Design and Methodology	14
3.1	Recommender System Details	15
3.1.1	Data and Implementation	15
3.1.2	Knowledge Graph Composition	15
3.1.3	Path-Based Recommendation Mechanics	15
3.2	Counterfactual Analysis	16
3.2.1	Attribute Selection	16
3.2.2	Performing Counterfactual Analysis	17
3.3	Case Study	17
4	Results	18
5	Conclusion	19
	References	20

Chapter 1

Introduction

Chapter 2

Literature Review

2.1 Explainable AI In Knowledge Graph Recommender System

Knowledge Graphs (KGs) are pivotal in enhancing the explainability and accuracy of recommender systems. These structured, relational frameworks capture complex interactions among users, items, and their attributes, allowing for more nuanced recommendations coupled with clear, logical explanations. This literature review synthesizes recent advancements in explainable artificial intelligence (XAI) that utilize KGs to illustrate how these technologies not only refine recommendation quality but also enhance user trust and understanding through transparency.

2.1.1 Path-based Modeling

Path-based modeling has emerged as a fundamental innovation in the utilization of KGs for recommender systems. Techniques such as the Knowledge-aware Path Recurrent Network (KPRN) and Path Language Modeling Recommendation (PLM-Rec) illustrate this trend's dynamic nature. KPRN leverages LSTM networks to interpret paths of entities and relationships, emphasizing those connections that are most influential in understanding user preferences. This method enriches the recommendation process by providing a temporal and semantic depth that traditional models lack, allowing for a better prediction of user behavior based on past interactions (Wang et al., 2019). On the other hand, PLM-Rec employs a novel approach by integrating natural language processing techniques to extend KG paths. This model treats paths as sentences, using

a language model to dynamically predict and extend these paths within the KG. Such extensions help the system explore new, potentially uncharted areas of the KG, thereby enhancing the system’s ability to recommend items that were previously unreachable. This approach addresses the inherent limitations of static KG structures and improves the system’s recall capabilities, making it particularly valuable for discovering long-tail items (Geng et al., 2022). Together, these path-based methods signify a shift towards more dynamic and exploratory use of KGs, expanding both the depth and breadth of what recommender systems can achieve.

2.1.2 Integration of User Profiles and Behavior

The integration of user behavioral data into KGs has significantly refined the personalization capabilities of recommender systems. The "Cafe" model by Xian et al. (2020) represents a sophisticated application of this concept, employing a coarse-to-fine strategy where initially broad user profiles help to narrow down and guide the path-finding algorithms in KGs. These profiles are crafted from historical data and are instrumental in focusing the recommendation process on paths most relevant to individual users, thus enhancing both the relevance and personalization of the recommendations. This method mirrors strategies used in other models that combine knowledge-base embeddings (KBE) with collaborative filtering. By embedding user behaviors and item characteristics into a unified representation, these models achieve a granular understanding of user-item relationships. This integration allows for a tailored recommendation experience, where the system’s outputs are closely aligned with individual preferences and behaviors, as demonstrated in the work by Ai et al. (2018).

2.1.3 Explainability through Symbolic and Logical Reasoning

The demand for explainability in AI has driven the adoption of models that incorporate transparent, logical reasoning processes. Monotonic GNNs (MGNNs), introduced by Tena et al. (2022), exemplify this trend by ensuring that every transformation within the network adheres to a set of logical rules, akin to traditional rule-based systems. This adherence guarantees that the network’s operations are interpretable and justifiable, enhancing user trust by providing comprehensible explanations for the recommendations made. Similarly, the Policy-Guided Path Reasoning (PGPR) model uses reinforcement learning to navigate through the KG, selecting paths that not only lead to relevant recommendations but are also interpretable. This model provides explicit paths that detail

the reasoning behind each recommendation, fulfilling the dual requirements of accuracy and transparency in the recommendation process (Xian et al., 2019).

The convergence of these methodologies highlights a crucial trend towards enhancing both the predictive accuracy and the interpretability of KG-based systems. Through the integration of dynamic path exploration, personalized user profile analysis, and logical reasoning, these approaches offer a more profound understanding of the intricacies involved in making recommendations. They collectively emphasize a shift towards recommender systems that are not only effective in their predictions but also provide transparent and understandable explanations, aligning with the growing user demand for transparency and accountability in AI systems.

2.2 Counterfactual Methods in Different Types of Recommender Systems

Counterfactual reasoning in recommender systems has emerged as a pivotal technique within the domain of explainable artificial intelligence (XAI), enhancing both the transparency and fairness of recommendations. By modeling alternative scenarios where specific variables are modified, this approach provides insights into the potential impacts of different data configurations, helping to elucidate the inner workings and dependencies within these systems.

The introduction of the KGCR model (Y. Wei et al., 2023) marks a significant advancement in embedding causal inference within graph-based recommender systems. Utilizing Graph Convolutional Networks, this model enriches user, item, and attribute embeddings, which allow for a more nuanced understanding of user preferences. By constructing a causal graph and applying do-calculus interventions, the KGCR model effectively mitigates biases introduced by previous user interactions, offering a refined approach to understanding how bias influences recommendation outcomes.

In a similar vein, Tran et al. (2021) developed the ACCENT framework, which facilitates the generation of actionable counterfactual explanations in neural recommender systems. This framework leverages extended influence functions to explore how changes in user-item interactions could affect recommendation outputs, significantly enhancing computational efficiency through Fast Influence Analysis. This methodology underscores the minimal adjustments in user behavior that could lead to different recommendations, thereby aiding in the creation of more transparent recommendation mechanisms.

Addressing selection bias, (Liu et al., 2022) implemented counterfactual policy learning to recalibrate recommendation fairness and effectiveness. Their approach utilizes Inverse Propensity Scoring to weigh observed interactions, allowing the system to simulate outcomes under different recommendation policies. By integrating these counterfactual outcomes into the learning process, the model achieves an improved balance, enhancing both the performance and equity of recommendations across various user groups and item categories.

The Prince method (Ghazimatin et al., 2020), emphasizes the importance of trust and understanding in recommendation systems through counterfactual reasoning within heterogeneous information networks. By identifying key user actions and employing Personalized PageRank, Prince efficiently predicts the impact of these actions on recommendation outcomes. This approach not only avoids exhaustive computations but also outperforms traditional heuristic methods in providing understandable and trust-enhancing explanations.

Yang et al. (2021) utilize causal inference through Structural Equation Models (SEMs) to address data sparsity in recommender systems. By generating counterfactual training samples, they enrich the dataset with diverse user responses that are otherwise not observed but plausible. This approach not only enhances the performance of the recommender systems but also strengthens their capacity to handle scenarios marked by data imbalance.

Finally, the Counterfactual Explainable Recommendation (CountER) model (Tan et al., 2021) focuses on identifying minimal attribute changes that could reverse a recommendation decision. Through a structured optimization process, CountER iteratively adjusts item attributes to discover the least extensive yet impactful changes required for altering outcomes. This model utilizes novel metrics to evaluate the necessity and sufficiency of these changes, demonstrating enhanced precision in providing actionable insights into recommendation decisions.

In conclusion, counterfactual reasoning offers a robust framework for enhancing the explainability and fairness of recommender systems by providing a deeper understanding of the implications of various data interactions and policies. These innovative approaches not only clarify the decision-making processes but also foster more equitable and user-centric recommendation practices.

2.3 Application of Counterfactuals for Fairness and Bias Mitigation

Counterfactual reasoning plays a pivotal role in the domain of explainable artificial intelligence (XAI), especially for mitigating biases in automated decision-making systems. This method involves hypothesizing alternative scenarios where key variables are altered, allowing for the exploration of how such changes impact outcomes. This not only uncovers hidden biases but also ensures fairness in AI operations. Broadly applied in various AI frameworks, from graph neural networks to recommender systems, counterfactual reasoning enhances transparency and equity in AI outcomes, establishing it as an essential tool for ethical AI development.

2.3.1 Mitigating Bias Across Different AI Frameworks

The use of counterfactual reasoning in graph-based models like those studied Guo et al. (2023) demonstrates a rigorous approach to maintaining consistency in model predictions across varying sensitive attributes. By implementing Graph Variational Autoencoders (GraphVAE), they not only perturb attributes but also train the network to minimize discrepancies in outputs between the original and counterfactual nodes. This methodology effectively addresses biases at a fundamental level, ensuring the fairness of the model's outcomes. Medda et al. (2024) extend this approach within graph neural network-based recommender systems. Their innovative use of counterfactual reasoning to adjust user-item interactions on a bipartite graph includes strategically adding or removing connections, which serves to simulate various scenarios where demographic disparities can be analyzed and mitigated, ensuring a more equitable distribution of utility among users. The field of recommender systems frequently grapples with biases such as popularity and exposure, which can distort user preferences. T. Wei et al. (2021) dissect these issues through the Model-Agnostic Counterfactual Reasoning (MACR) framework, which explicitly separates the influence of item popularity from actual user preferences. By adjusting input data to simulate a scenario where item popularity is neutralized, MACR provides a recalibrated basis for recommendation, aligning more closely with unbiased user preferences. Meanwhile, Xu et al. (2020) focus on exposure bias by employing a counterfactual approach that involves a minimax adversarial model. This model simulates worst-case scenarios to test the resilience of the recommendation system, ensuring that it can withstand and adapt to a range of user exposure conditions, thus promoting

a more fair and balanced recommendation landscape.

2.3.2 Enriching Data and Ensuring Equitable Outcomes

Addressing data sparsity and imbalance, Yang et al. (2021) utilize causal inference via Structural Equation Models (SEMs) to generate counterfactual scenarios that enrich training datasets. This not only addresses the immediate issue of insufficient data but also simulates a broader spectrum of user interactions, which helps in developing a more robust and responsive recommender system. On a more focused level, Chiappa (2019) pioneers the use of Path-Specific Counterfactual Fairness (PSCF) within decision-making processes. This approach manipulates causal pathways, particularly those that might be influenced by sensitive attributes such as race or gender, to ensure that resulting decisions are free from the undue influence of these attributes, thus promoting fairness in critical decision-making contexts.

2.3.3 Leveraging Knowledge Graphs for Fair Recommendations

Expanding the utility of knowledge graphs, Balloccu et al. (2022) integrate counterfactual reasoning within the Policy-Guided Path Reasoning (PGPR) model to optimize recommendation systems. By re-ranking items and explanations based on various fairness-oriented criteria, such as recency, popularity, and diversity, PGPR enhances the quality and equity of recommendations. This approach not only improves the relevance of the recommendations but also significantly increases user trust and satisfaction by ensuring that recommendations cater equitably to diverse user groups.

Chapter 3

Research Design and Methodology

This chapter delineates the methodology employed to achieve the primary objective of this study: extending the explainability of path-based knowledge graph recommender systems to explore "what if" scenarios. The methodology is structured as follows:

1. **Initial Recommendation:** The process commences with the system generating a product recommendation for a user. In a path-based recommendation system, the inherent explanation typically comprises a sequence of entities and relationships that link the user to the recommended product.
2. **Counterfactual Analysis:**
 - **Extraction of Relevant Information:** Initially, the analysis extracts pertinent information related to the entities along the recommendation path. This includes the specific attributes that influenced the selection of the final product, as well as other potentially relevant attributes associated with the recommended product.
 - **Scenario Construction:** Utilizing the extracted information, a collection of hypothetical scenarios is constructed. These scenarios are crafted to test various alterations in attributes and their impact on the recommendation outcome.
3. **Recommendation System Utilization:** For the recommendation engine, we employ CAFE (Coarse-to-Fine Neural Symbolic Reasoning for Explainable Recommendation). This system is particularly suitable for our purposes due to its path-based nature and its capability to evaluate the plausibility of different paths by assigning probability scores to the steps that connect users to products.

This methodology both facilitates a deeper understanding of the decision-making processes inherent in the recommender system and also allows us to simulate and evaluate

how changes in product attributes or user-product relationships might alter the system’s recommendations.

3.1 Recommender System Details

CAFE (Coarse-to-Fine Neural Symbolic Reasoning for Explainable Recommendation) is used as the foundational framework for our recommender system. This section provides an overview of its implementation and core functionalities.

3.1.1 Data and Implementation

The CAFE model is implemented using the Amazon review dataset, the beauty category, which includes comprehensive user and product interactions. It leverages predefined embeddings train in the model developed by Ai et al. (2018) described in the previous section, as input to their symbolic network.

3.1.2 Knowledge Graph Composition

The knowledge graph at the heart of this recommendation system is intricately structured, comprising several types of entities and relationships:

- **Users** are linked to the words they have used and the products they have purchased.
- **Products** are associated with descriptive words, their brand, category, and other related products. Relationships with related products include those that have been 'bought together', 'also viewed', and 'also bought'.
- **Brands and Categories** form additional nodes, creating multiple pathways that connect different aspects of the data.
- **Related categories** mentioned above.

3.1.3 Path-Based Recommendation Mechanics

The system operates on predefined metapaths that represent meaningful relationships leading a user to a product. These metapaths are substantial for understanding the logic behind the recommendations. Todo: That paths . The recommender system assigns a probability score to each step along the path, determining the strength of the connection between the user and the potential product recommendations. It then selects the top 10 paths with the highest cumulative scores, and the products associated with these

paths are recommended to the user. This scoring and selection process ensures that the recommendations are both relevant and tailored to the user’s preferences and behavior patterns. This structured approach allows the CAFE system to recommend products effectively but also provide insights into the reasons behind each recommendation, providing explainability to the system.

3.2 Counterfactual Analysis

subsectionCommunity Detection and Graph Analysis Node FilteringThe counterfactual analysis begins with the detection of communities within the knowledge graph of interactions. Communities are identified using Louvain Method. This helps to cluster entities that share significant similarities and interactions. The Louvain Method is an efficient algorithm designed for detecting communities in large-scale networks by optimizing modularity, a measure that quantifies the density of links within communities relative to those between them introduced in (Blondel et al., 2008). The algorithm operates in two iterative phases: initially, it optimizes modularity locally, evaluating potential gains by moving individual nodes into different communities. Nodes are shifted to the community that maximizes this gain, and the process is repeated until no further improvement is possible, achieving a local maximum of modularity. In the second phase, the method aggregates these identified communities into new nodes of a reduced network, and the process is reapplied. This hierarchical approach allows the algorithm to uncover community structures at multiple levels effectively. Notable for its speed, the Louvain Method can handle networks with up to millions of nodes efficiently, making it well-suited for modern datasets of substantial size. To refine the analysis further, we calculate the degree centrality for each type of pair within the graph. Nodes that do not provide significant insight are filtered out based on their z-score; specifically, nodes with a z-score exceeding [specific threshold] are removed.

3.2.1 Attribute Selection

Following the predictions provided by the recommender system, a threshold is set for the minimum score required for a product path to be recommended to a user, which is the path score of the last product recommended in the top 10 recommended products. For the analysis of a recommended path, we retrieve the first-level attributes and their

related products. This forms the initial layer of attribute selection, predicated on the hypothesis that first-level connected items possess more relevant attributes. These selected attributes are then evaluated to determine whether they fall within the community of the recommended product. If they do, and their z-score is within an acceptable range, they are considered for further counterfactual analysis.

3.2.2 Performing Counterfactual Analysis

For each attribute, an appropriate metapath is selected based on the type of the attribute. Using the recommender engine, we calculate the score for a user-product combination, which incorporates all the products previously purchased by the user and the counterfactual attribute in question. This approach is grounded in the assumption that the system discerns the user's preferences through their purchase history. If the recalculated score for a product, when considering a counterfactual attribute, exceeds the set minimum score, the attribute is considered a positive influence for a product similar to the recommended one. This isolated attribute analysis not only aids in understanding the influence of specific attributes on product recommendations but also provides marketing insights. Such insights can further be used to enhance the diversity and precision of the recommender system.

3.3 Case Study

Chapter 4

Results

Chapter 5

Conclusion

References

- Ai, Q., Azizi, V., Chen, X., & Zhang, Y. (2018). Learning heterogeneous knowledge base embeddings for explainable recommendation. *Algorithms*, 11(9), 137. <https://doi.org/10.3390/a11090137>
- Balloccu, G., Boratto, L., Fenu, G., & Marras, M. (2022). Post processing recommender systems with knowledge graphs for recency, popularity, and diversity of explanations. *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 646–656. <https://doi.org/10.1145/3477495.3532041>
- Blondel, V. D., Guillaume, J.-L., Lambiotte, R., & Lefebvre, E. (2008). Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2008(10), P10008. <https://doi.org/10.1088/1742-5468/2008/10/P10008>
- Chiappa, S. (2019). Path-specific counterfactual fairness. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(1), 7801–7808. <https://doi.org/10.1609/aaai.v33i01.33017801>
- Geng, S., Fu, Z., Tan, J., Ge, Y., De Melo, G., & Zhang, Y. (2022). Path language modeling over knowledge graphs for explainable recommendation. *Proceedings of the ACM Web Conference 2022*, 946–955. <https://doi.org/10.1145/3485447.3511937>
- Ghazimatin, A., Balalau, O., Saha Roy, R., & Weikum, G. (2020). PRINCE: Provider-side interpretability with counterfactual explanations in recommender systems. *Proceedings of the 13th International Conference on Web Search and Data Mining*, 196–204. <https://doi.org/10.1145/3336191.3371824>
- Guo, Z., Li, J., Xiao, T., Ma, Y., & Wang, S. (2023). Towards fair graph neural networks via graph counterfactual. *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, 669–678. <https://doi.org/10.1145/3583780.3615092>

- Liu, Y., Yen, J.-N., Yuan, B., Shi, R., Yan, P., & Lin, C.-J. (2022). Practical counterfactual policy learning for top-k recommendations. *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 1141–1151. <https://doi.org/10.1145/3534678.3539295>
- Medda, G., Fabbri, F., Marras, M., Boratto, L., & Fenu, G. (2024). GNNUERS: Fairness explanation in GNNs for recommendation via counterfactual reasoning. *ACM Transactions on Intelligent Systems and Technology*, 3655631. <https://doi.org/10.1145/3655631>
- Tan, J., Xu, S., Ge, Y., Li, Y., Chen, X., & Zhang, Y. (2021). Counterfactual explainable recommendation. *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 1784–1793. <https://doi.org/10.1145/3459637.3482420>
- Tran, K. H., Ghazimatin, A., & Saha Roy, R. (2021). Counterfactual explanations for neural recommenders. *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 1627–1631. <https://doi.org/10.1145/3404835.3463005>
- Wang, X., Wang, D., Xu, C., He, X., Cao, Y., & Chua, T.-S. (2019). Explainable reasoning over knowledge graphs for recommendation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(1), 5329–5336. <https://doi.org/10.1609/aaai.v33i01.33015329>
- Wei, T., Feng, F., Chen, J., Wu, Z., Yi, J., & He, X. (2021). Model-agnostic counterfactual reasoning for eliminating popularity bias in recommender system. *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 1791–1800. <https://doi.org/10.1145/3447548.3467289>
- Wei, Y., Wang, X., Nie, L., Li, S., Wang, D., & Chua, T.-S. (2023). Causal inference for knowledge graph based recommendation. *IEEE Transactions on Knowledge and Data Engineering*, 35(11), 11153–11164. <https://doi.org/10.1109/TKDE.2022.3231352>
- Xian, Y., Fu, Z., Muthukrishnan, S., De Melo, G., & Zhang, Y. (2019). Reinforcement knowledge graph reasoning for explainable recommendation. *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 285–294. <https://doi.org/10.1145/3331184.3331203>
- Xian, Y., Fu, Z., Zhao, H., Ge, Y., Chen, X., Huang, Q., Geng, S., Qin, Z., De Melo, G., Muthukrishnan, S., & Zhang, Y. (2020). CAFE: Coarse-to-fine neural symbolic

- reasoning for explainable recommendation. *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, 1645–1654. <https://doi.org/10.1145/3340531.3412038>
- Xu, D., Ruan, C., Korpeoglu, E., Kumar, S., & Achan, K. (2020). Adversarial counterfactual learning and evaluation for recommender system [Publisher: arXiv Version Number: 1]. <https://doi.org/10.48550/ARXIV.2012.02295>
- Yang, M., Dai, Q., Dong, Z., Chen, X., He, X., & Wang, J. (2021). Top-n recommendation with counterfactual user preference simulation [Publisher: [object Object] Version Number: 2]. <https://doi.org/10.48550/ARXIV.2109.02444>

Appendix

You are encouraged to put in appendices in your final report. In an appendix you can include things such as large tables or background information. Anything that is useful to know for the reader, but prevents the reader to read your main text in a fluent manner. Each appendix should have a number and a self-explanatory title.