

Topic-Enhanced LightGCN (TE-LGCN): Addressing Data Sparsity in Movie Recommendation Systems

Lee Donghun
ldong0308@hanyang.ac.kr

Yun Juchan
juchan563@gmail.com

Jang Yeonguk
bike481227@gmail.com

November 24, 2025

Abstract

Recommender systems have become essential tools in the era of information overload, yet data sparsity and cold-start problems remain fundamental challenges. Graph neural network (GNN)-based collaborative filtering models, particularly LightGCN, demonstrate superior performance on dense datasets but struggle in sparse user-item interaction environments. The long-tail problem exacerbates these issues, with over 80% of users interacting with fewer than 10 items and more than 70% of items receiving insufficient exposure. To address these limitations, we propose a topic-enhanced heterogeneous graph approach that augments the standard user-item bipartite graph with topic nodes extracted through Latent Dirichlet Allocation (LDA) from movie overviews and descriptions. Our method incorporates a dual enhancement strategy: (1) Doc2Vec-based item node initialization that leverages textual content for better initial representations, and (2) explicit topic nodes that create additional pathways for information propagation through LDA topic modeling. Unlike existing methods that either abandon LightGCN’s lightweight architecture or rely solely on interaction data, our approach maintains computational efficiency while incorporating content-based signals through structured topic representations. This enables effective recommendations for light users with minimal interaction history and promotes niche items in the long tail. Experimental results in MovieLens datasets demonstrate that our method achieves a 26.4% improvement in Recall@10 compared to vanilla LightGCN. Our work bridges collaborative and content-based filtering in a unified graph framework, offering a practical solution for real-world sparse recommendation scenarios.

1 Introduction

Recommender systems are one of the core technologies in modern digital platforms, yet fundamental challenges that remain difficult to solve still exist. In particular, data sparsity poses a significant constraint on recommender system performance in real-world service environments.

Analysis of real-world recommender system data reveals a typical Pareto distribution. In the MovieLens dataset, over 80% of users interact with fewer than 10 movies, while more than 70% of movies receive insufficient user ratings, exhibiting a long-tail phenomenon. This sparsity creates two critical problems:

1. It becomes difficult to provide personalized recommendations to new users or cold-start users with limited interaction history.
2. Items with low popularity lose opportunities for discovery and become increasingly marginalized, leading to “popularity bias.”

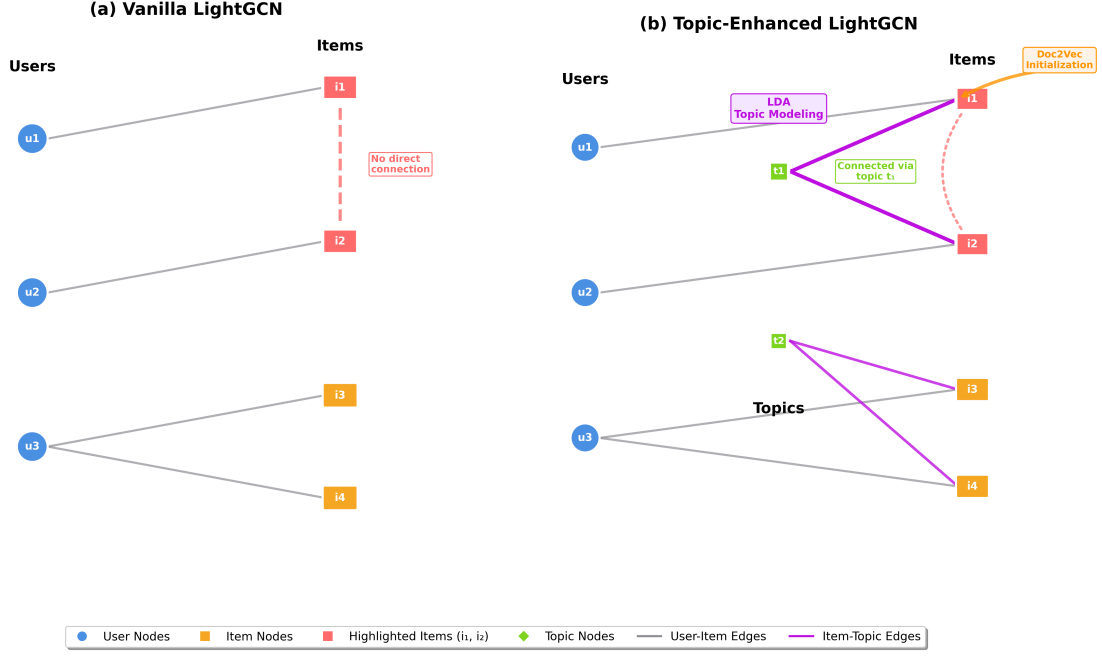


Figure 1: **Comparison of graph connectivity.** (a) In Vanilla LightGCN, sparse items like i_1 and i_2 lack direct connections and remain isolated. (b) Topic-Enhanced LightGCN introduces explicit Topic Nodes, creating semantic bridges between unconnected items and promoting broader discovery.

To address these challenges, Graph Neural Network (GNN)-based collaborative filtering methods have recently gained attention. LightGCN, in particular, proposed a lightweight architecture optimized for recommender systems by removing complex feature transformations and nonlinear activations from existing GCN-based methods. As shown in Figure 1(a), LightGCN demonstrated superior performance over existing complex methods using only standard user-item bipartite graphs.

However, LightGCN is fundamentally a collaborative filtering approach that relies solely on interaction data. Therefore, its performance degrades in sparse environments with insufficient interactions. When connections between users and items are scarce, information propagation through graphs alone cannot learn meaningful representations.

To overcome these limitations, approaches utilizing item content information have been explored. For movies, rich information such as plot summaries, genres, and cast details is available, providing important clues for understanding user preferences. For instance, if a user prefers science fiction movies, they are likely to be interested in other movies with similar SF themes. However, existing content-based approaches exhibit several limitations:

1. Most methods abandon LightGCN’s lightweight architecture and introduce complex structures.
2. Simply combining textual information as features fails to fully leverage the advantages of graph structures.
3. Systematic research on effective integration of content information with collaborative filtering signals is lacking.

Our Solution: Topic-Enhanced LightGCN

This study proposes Topic-Enhanced LightGCN to address these challenges. As illustrated in Figure 1(b), our core idea is to extend the traditional user-item bipartite graph with topic nodes to construct a

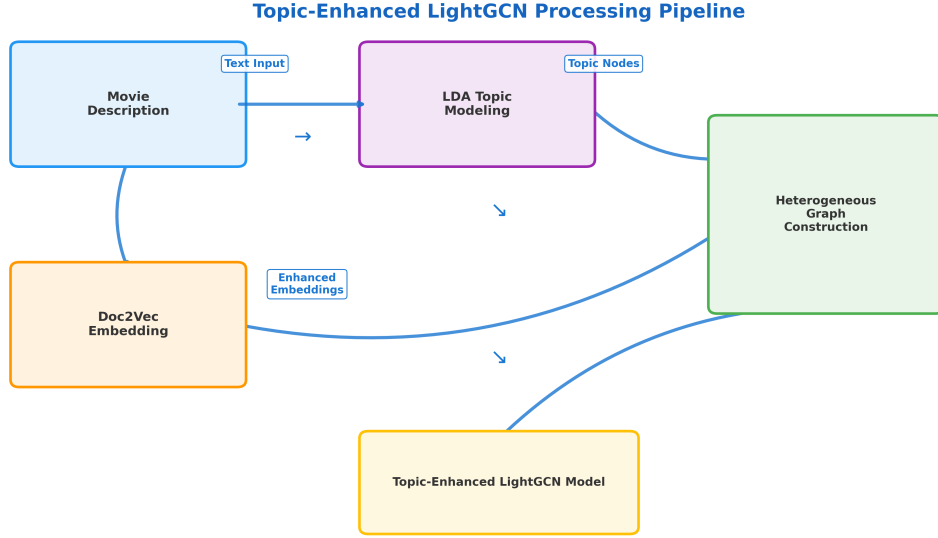


Figure 2: **The overall processing pipeline of Topic-Enhanced LightGCN.** Movie descriptions undergo two parallel processes: LDA Topic Modeling constructs explicit topic nodes for graph extension, while Doc2Vec generates initial semantic embeddings to mitigate cold-start issues.

heterogeneous graph. This enables natural integration of textual content information while maintaining LightGCN’s simplicity and efficiency.

Figure 2 shows our complete processing pipeline. Specifically, we present a dual enhancement strategy. First, we initialize item nodes with semantic vectors extracted from movie plot summaries using Doc2Vec. This naturally injects semantic information from textual content into the graph learning process. Second, we add topics extracted through Latent Dirichlet Allocation (LDA) as explicit nodes to the graph. This structurally models semantic connections between items, enabling effective information propagation even in sparse interaction environments.

The introduction of topic nodes provides indirect information propagation pathways mediated by topics, in addition to the existing direct user-item connections, as shown in Figure 1(b). For example, item i_4 , which user u_2 has not directly interacted with, can be connected through common topic t_2 , enabling recommendations based on semantic similarity

The main contributions of this research are as follows:

1. **Heterogeneous Graph Structure** We extend LightGCN with topic nodes, going beyond simple user-item graphs to structurally model semantic connections.
2. **Dual Utilization of Text** We effectively utilize textual information through two complementary strategies: Doc2Vec-based initialization and LDA topic nodes, enriching the graph structure itself.
3. **Practical Approach** We present a method that addresses data sparsity without complex architectural changes, enhancing applicability in real service environments.
4. **Performance on Sparsity** We demonstrate effectiveness in sparse data environments through systematic performance analysis on long-tail items.

2 Related Work

2.1 Graph Neural Networks for Recommendation

Graph-based collaborative filtering has shown remarkable success in capturing user-item interactions through message passing mechanisms. He et al. proposed LightGCN [1], which simplifies GCN architectures by removing feature transformation and nonlinear activation, focusing solely on neighborhood aggregation. LightGCN demonstrated an average performance improvement of over 16% compared to NGCF on benchmark datasets such as Gowalla, Yelp2018, and Amazon-Book, but suffers from performance degradation in extremely sparse environments where user-item interactions are severely limited. Recent studies have attempted to enhance LightGCN through contrastive learning [2] and embedding norm scaling [3], but these approaches do not fundamentally address the cold-start problem caused by insufficient user-item interactions.

2.2 Addressing Cold-Start and Data Sparsity

Cold-start recommendation remains a fundamental challenge in recommender systems. Meta-learning approaches [4] address this by learning to generalize from similar users or items, while recent TMAG [5] proposed graph augmentation techniques through task-aligned meta-learning. Content-based methods leverage item metadata to infer embeddings for new items. Qian et al. proposed AGNN [6], which constructs attribute graphs to handle cold-start scenarios, demonstrating the effectiveness of incorporating side information such as item attributes. However, these methods have limited integration of collaborative signals with attribute information or require complex meta-learning frameworks. Our work bridges this gap by augmenting the collaborative graph structure with content-derived topic nodes while maintaining the simplicity of LightGCN.

2.3 Topic Modeling and Content Integration

Traditional topic modeling approaches have been integrated with collaborative filtering to address sparsity issues. Wang and Blei proposed Collaborative Topic Regression (CTR) [7], which combines probabilistic matrix factorization with LDA-based topic modeling to leverage both user-item interactions and item textual content. CTR demonstrated that topic distributions extracted from item descriptions can effectively bridge the gap between collaborative and content-based signals, particularly proving effective in alleviating cold-start problems by utilizing content information. However, CTR is based on a probabilistic framework and structurally differs from recent graph neural network-based approaches. Our research, similar to CTR, utilizes topic modeling but is differentiated by integrating topics as explicit nodes within a graph neural network framework instead of probabilistic matrix factorization.

2.4 Heterogeneous Graph and Knowledge-Enhanced Approaches

Heterogeneous graph neural networks have emerged as a powerful paradigm for integrating multi-modal information in recommendation systems. Zhang et al. introduced Light Heterogeneous Graph CF [8], which incorporates textual descriptions as nodes and uses SBERT embeddings to initialize text representations. Similarly, [knowledge graph-based approaches](<https://arxiv.org/pdf/2403.18667>) leverage semantic relationships between items to enhance recommendation diversity and address cold-start users.

Wang et al. proposed KGAT (Knowledge Graph Attention Network) [9], which constructs collaborative knowledge graphs (CKG) by linking items with their attributes such as genres, directors, and actors,

and employs graph attention networks to explicitly model high-order connectivity. KGAT overcomes the manual path selection problems of existing path-based methods and the implicit relation modeling limitations of regularization-based methods, demonstrating significant performance improvements. However, KGAT relies on pre-defined structured knowledge graphs, which have practical limitations requiring manual curation by domain experts or external knowledge bases. The increased model complexity and computational costs associated with knowledge graph construction and maintenance are also major challenges. Moreover, KGAT does not include mechanisms to automatically discover and utilize latent topics or semantic nuances inherent in unstructured text such as movie overviews or plot descriptions.

Recently, Lee et al. proposed KGMC (Keyword-enhanced Graph Matrix Completion) [10], which extracts keywords from user reviews using TF-IDF or BERT to construct additional edges in the interaction graph. KGMC alleviates data sparsity problems by strengthening user-user, item-item, and user-item connections through keyword co-occurrence patterns. While this shows a similar approach of explicitly integrating keyword nodes into graph structures, its dependence on user reviews limits applicability in situations where review data is insufficient or unavailable.

While these works demonstrate the value of heterogeneous structures and knowledge integration, they have constraints such as using complex modified RGCN architectures, knowledge graph construction costs, or dependence on user review data. In contrast, our approach maintains LightGCN’s lightweight design while incorporating topic nodes extracted directly from movie overviews and descriptions through LDA-based topic modeling. This enables automatic discovery and utilization of latent semantic structures from unstructured textual content of items without requiring external knowledge bases or user reviews.

3 Methodology

3.1 Problem Formulation

3.1.1 Sparsity Challenge in Collaborative Filtering

We consider a standard recommendation setting with users $\mathcal{U} = \{u_1, \dots, u_N\}$, items $\mathcal{I} = \{i_1, \dots, i_M\}$, and implicit feedback matrix $\mathbf{R} \in \{0, 1\}^{N \times M}$. Real-world datasets exhibit severe sparsity, creating fundamental challenges for collaborative filtering approaches.

Traditional methods optimize the BPR objective:

$$\mathcal{L}_{BPR} = - \sum_{(u,i,j) \in \mathcal{D}_s} \ln \sigma(\hat{y}_{ui} - \hat{y}_{uj}) \quad (1)$$

where $\hat{y}_{ui} = \mathbf{e}_u^T \mathbf{e}_i$ and $\mathcal{D}_s = \{(u, i, j) : R_{ui} = 1, R_{uj} = 0\}$.

However, in sparse environments, the limited interactions make it difficult to learn meaningful embeddings \mathbf{e}_u and \mathbf{e}_i , leading to poor prediction quality, particularly for cold-start users and long-tail items.

3.1.2 Content-Enhanced Problem Formulation

To address this limitation, we reformulate the problem to incorporate textual content $\mathcal{D} = \{d_1, \dots, d_M\}$ (item descriptions) into the recommendation process. Our goal is to enhance the prediction function \hat{y}_{ui} such that it captures both collaborative signals and semantic relationships derived from content.

We propose to:

1. Extract semantic structure from item descriptions using topic modeling
2. Initialize embeddings with content-aware representations
3. Extend graph connectivity through topic-mediated paths

This enables BPR optimization with richer embeddings that remain effective even when direct user-item interactions are sparse, while preserving the proven ranking objective framework.

3.2 Methodology

3.2.1 Topic-Enhanced Graph Construction

Our approach transforms the traditional user-item bipartite graph into a heterogeneous graph that incorporates semantic relationships derived from item textual content. This transformation involves three key steps: topic extraction, content-aware initialization, and graph extension.

Topic Extraction via LDA We apply Latent Dirichlet Allocation (LDA) to extract K latent topics from the collection of item descriptions $\mathcal{D} = \{d_1, \dots, d_M\}$. LDA produces topic distributions $\theta_i \in R^K$ for each item i , where θ_{ik} represents the probability that item i belongs to topic k . We establish item-topic connections when $\theta_{ik} > \alpha$, creating edges $\mathcal{E}_{it} = \{(i, t_k) : \theta_{ik} > \alpha\}$ that link items to their dominant topics.

Content-Aware Initialization To inject semantic information from the textual content, we utilize Doc2Vec to generate dense vector representations $\mathbf{v}_i = \text{Doc2Vec}(d_i)$ for each item description. These vectors serve as initial embeddings for item nodes: $\mathbf{e}_i^{(0)} = \mathbf{v}_i$, ensuring that the model starts with meaningful content-based representations rather than random initialization.

Heterogeneous Graph Construction The enhanced graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ extends the standard user-item structure with topic nodes, where $\mathcal{V} = \mathcal{U} \cup \mathcal{I} \cup \mathcal{T}$ and $\mathcal{E} = \mathcal{E}_{ui} \cup \mathcal{E}_{it}$. This creates additional pathways for information propagation: users can now discover items through shared topic preferences, enabling better recommendations for sparse users and promoting long-tail item visibility through semantic similarity rather than popularity alone.

3.2.2 Enhanced Message Passing

Our model extends LightGCN’s message passing mechanism to operate on the heterogeneous graph with three node types. The enhanced propagation enables information flow between users, items, and topics, creating richer embeddings that capture both collaborative and semantic signals.

Multi-Type Node Updates The layer-wise embedding updates are formulated as:

$$\mathbf{e}_u^{(l+1)} = \sum_{i \in \mathcal{N}(u)} \frac{1}{\sqrt{|\mathcal{N}(u)||\mathcal{N}(i)|}} \mathbf{e}_i^{(l)} \quad (2)$$

$$\begin{aligned} \mathbf{e}_i^{(l+1)} &= \sum_{u \in \mathcal{N}(i)} \frac{1}{\sqrt{|\mathcal{N}(i)||\mathcal{N}(u)|}} \mathbf{e}_u^{(l)} \\ &\quad + \sum_{t \in \mathcal{N}(i)} \frac{1}{\sqrt{|\mathcal{N}(i)||\mathcal{N}(t)|}} \mathbf{e}_t^{(l)} \end{aligned} \quad (3)$$

$$\mathbf{e}_t^{(l+1)} = \sum_{i \in \mathcal{N}(t)} \frac{1}{\sqrt{|\mathcal{N}(t)||\mathcal{N}(i)|}} \mathbf{e}_i^{(l)} \quad (4)$$

The key innovation lies in the item update equation, where items now aggregate information from both connected users (collaborative signal) and associated topics (semantic signal).

Topic-Mediated Information Flow This architecture enables multi-hop information propagation through topic nodes. For instance, a user who watched “Interstellar” can receive recommendations for “Blade Runner” through their shared “Science Fiction” topic, even without direct collaborative signals between these items. This semantic pathway is particularly valuable for cold-start scenarios and long-tail item discovery.

Final Embedding Computation Following LightGCN’s design, we obtain final embeddings by aggregating across all layers:

$$\begin{aligned} \mathbf{e}_u^{(\text{final})} &= \frac{1}{L+1} \sum_{l=0}^L \mathbf{e}_u^{(l)}, & \mathbf{e}_i^{(\text{final})} &= \frac{1}{L+1} \sum_{l=0}^L \mathbf{e}_i^{(l)}, \\ \mathbf{e}_t^{(\text{final})} &= \frac{1}{L+1} \sum_{l=0}^L \mathbf{e}_t^{(l)} \end{aligned} \quad (5)$$

The prediction score remains $\hat{y}_{ui} = \mathbf{e}_u^{(\text{final})T} \mathbf{e}_i^{(\text{final})}$, preserving computational simplicity while benefiting from semantically enriched embeddings.

3.2.3 Learning Objective

Our learning framework maintains the BPR ranking objective while incorporating regularization terms to handle the extended graph structure and content information.

Complete Objective Function The total loss combines three components:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{BPR}} + \lambda_1 \mathcal{L}_{\text{reg}} + \lambda_2 \mathcal{L}_{\text{content}} \quad (6)$$

BPR Loss We preserve the standard ranking loss to maintain compatibility with implicit feedback:

$$\mathcal{L}_{\text{BPR}} = - \sum_{(u,i,j) \in \mathcal{D}_s} \ln \sigma \left(\mathbf{e}_u^{(\text{final})T} \mathbf{e}_i^{(\text{final})} - \mathbf{e}_u^{(\text{final})T} \mathbf{e}_j^{(\text{final})} \right) \quad (7)$$

where $\mathcal{D}_s = \{(u, i, j) : R_{ui} = 1, R_{uj} = 0\}$ represents preference triplets. The key advantage is that while the loss formulation remains unchanged, the embeddings now encode richer collaborative and semantic information.

Extended Regularization To prevent overfitting in the expanded parameter space:

$$\mathcal{L}_{\text{reg}} = \|\mathbf{E}_{\mathcal{U}}\|_F^2 + \|\mathbf{E}_{\mathcal{I}}\|_F^2 + \|\mathbf{E}_{\mathcal{T}}\|_F^2 \quad (8)$$

The inclusion of topic embedding regularization $\|\mathbf{E}_{\mathcal{T}}\|_F^2$ is essential as topic nodes have sparser connections and are more prone to overfitting.

Content Consistency An optional term to maintain semantic coherence:

$$\mathcal{L}_{\text{content}} = \sum_{i=1}^M \|\mathbf{e}_i^{(\text{final})} - \mathbf{v}_i\|_2^2 \quad (9)$$

This encourages learned item embeddings to preserve the semantic structure captured by Doc2Vec initialization, though setting $\lambda_2 = 0$ allows pure collaborative adaptation from content-aware starting points.

Training Overview The model is optimized using standard SGD with negative sampling, where

the enhanced graph structure enables more effective gradient propagation to sparse users and items through topic-mediated paths, improving convergence in sparse scenarios.

4 Experimental Setup

4.1 Dataset and Setup

We conduct experiments using a preprocessed MovieLens dataset adapted for our evaluation. The preprocessing pipeline involves converting explicit ratings to binary interactions by treating ratings ≥ 4 as positive feedback, followed by k -core filtering with $k = 10$ to ensure minimum connectivity for both users and items. The final dataset consists of 605 users interacting with 1,234 movies, resulting in 38,405 total interactions.

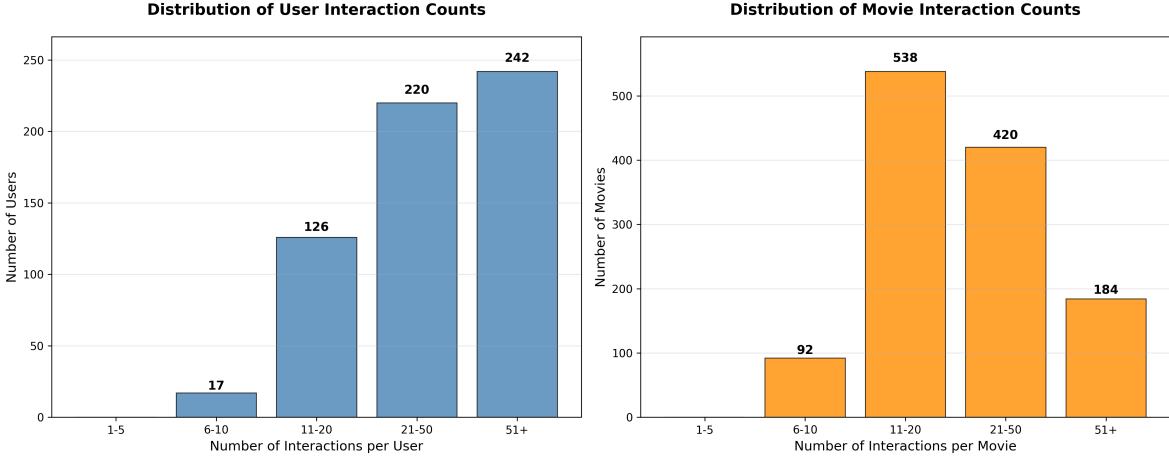
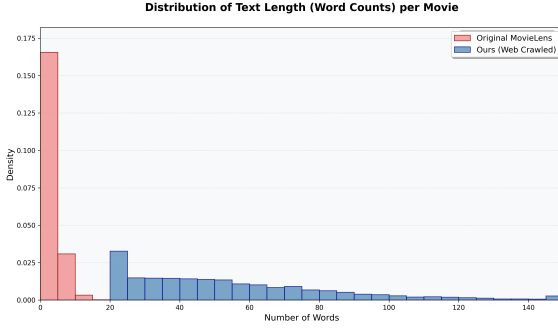


Figure 3: **Distribution of interactions.** Even with $k = 10$ core filtering, the dataset exhibits a typical long-tail distribution where a significant number of movies lack interactions.

Data Sparsity Analysis The processed dataset exhibits sparsity of 94.86%, where users interact with an average of 63.5 movies out of 1,234 available items. As illustrated in Figure 3, the k -core filtering ensures that all entities maintain minimum connectivity (≥ 10 interactions), yet the dataset still demonstrates clear long-tail characteristics. User interactions are distributed with 40.0% having 51+ interactions and 20.8% in the 11-20 range, while movies show a more balanced distribution with 43.6% receiving 11-20 interactions and 34.0% receiving 21-50 interactions. Despite the filtering, 7.5% of movies still receive only 10 or fewer interactions, representing the long-tail segment that benefits from our content-enhanced approach.

Evaluation Setup We adopt a temporal split strategy where interactions are sorted by timestamp, with 70% for training, 15% for validation, and 15% for testing. This setup ensures that the model is evaluated on future interactions, simulating real-world deployment scenarios. For comprehensive analysis, we evaluate performance across different user activity levels and focus particularly on long-tail movie recommendations.

Textual Content The original dataset contained only keywords, lacking rich textual data for content analysis. To address this limitation, we collected plot summary information for each movie through web crawling to construct an additional column. The collected textual data underwent preprocessing steps including tokenization, stopword removal, and lemmatization. After preprocessing, the average plot summary contains 53.95 words. As shown in Figure 4, our dataset demonstrates a substantial



(a) Expansion of vocabulary range



(b) Qualitative shift in vocabulary

Figure 4: **Comparison of Textual Content.** (a) The enhanced dataset demonstrates a substantial increase in vocabulary range. (b) While the original keywords focused on simple genres (e.g., Comedy, Sci-Fi), the crawled plot summaries capture deeper semantic meanings (e.g., life, story).

increase in both text length (Figure 4a) and vocabulary diversity compared to the original MovieLens data, providing rich semantic content for our topic modeling and embedding initialization procedures.

4.2 Baseline and Evaluation

Baseline Method We compare our approach against vanilla LightGCN as the primary baseline. This choice is strategically motivated as LightGCN represents the lightweight graph-based collaborative filtering, which directly aligns with our goal of maintaining computational efficiency while incorporating content information. LightGCN’s simplicity and proven effectiveness on sparse datasets make it an ideal foundation for demonstrating the incremental benefits of our topic-enhanced extensions.

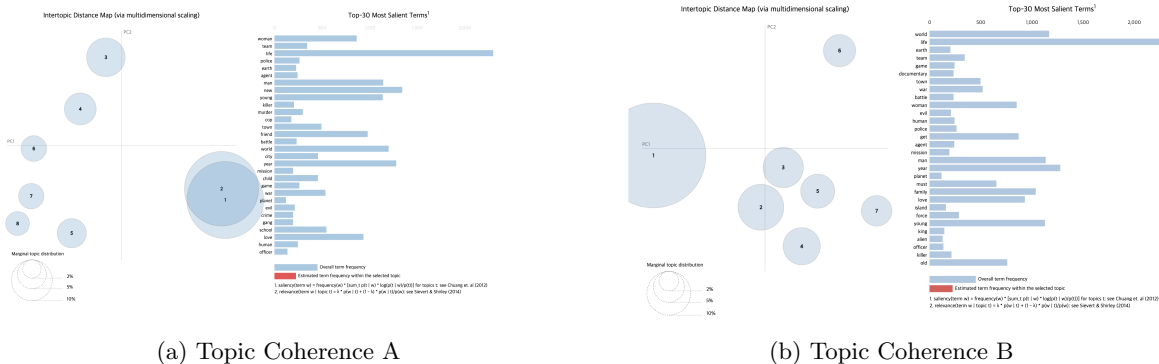
Evaluation Metrics We employ standard ranking-based metrics commonly used in implicit feedback recommendation systems:

1. **Precision@10** Measures the proportion of recommended items that are actually relevant, providing insight into recommendation accuracy.
2. **Recall@10** Measures the proportion of relevant items successfully retrieved in the top-10 recommendations.
3. **NDCG@10** Normalized Discounted Cumulative Gain that accounts for the ranking position of relevant items, providing a more nuanced evaluation of recommendation quality.

Evaluation Protocol For each user in the test set, we generate recommendations by ranking all unobserved items and computing metrics based on held-out positive interactions.

4.3 Implementation Details

Hyperparameter Settings Our model uses an embedding dimension of $d = 64$ for all node types (users, items, and topics). We employ $L = 3$ propagation layers in the graph neural network, which provides sufficient depth for multi-hop information propagation while maintaining computational efficiency. For topic modeling, we set the number of LDA topics to $K = 7$, selected based on topic coherence analysis during preprocessing. As shown in Figure 3, setting the number of topics to 7 achieves an optimal distribution where all objects maintain non-overlapping characteristics. The topic connection threshold is set to $\alpha = 0.1$, meaning items are connected to topics only when their topic probability exceeds this threshold.



Loss Function Parameters The regularization weight $\lambda_1 = 1e - 5$ is set to control the balance between model expressiveness and overfitting prevention. This value was determined by independently applying values ranging from $1e - 2$ to $1e - 6$ to vanilla LightGCN and selecting the parameter that achieved the highest recall value, with these results verifiable through Table 1. For the content consistency term, we experiment with $\lambda_2 \in \{0, 1e - 4, 1e - 3, 1e - 2, 1e - 1\}$ to evaluate the impact of preserving Doc2Vec initialization versus allowing pure collaborative adaptation. The results of this experiment are presented in Appendix B.

Table 1: **Performance comparison of Vanilla LightGCN** with different regularization weights (λ_1). The optimal value ($1e-5$) is highlighted.

Parameter	Evaluation Metrics			
	HitRate	Prec@10	Recall@10	NDCG@10
1e-6	0.1547	0.0155	0.1547	0.0818
1e-5	0.1622	0.0162	0.1622	0.0871
1e-4	0.1592	0.0159	0.1592	0.0850
1e-3	0.1532	0.0153	0.1532	0.0819
1e-2	0.1532	0.0153	0.1532	0.0792

Training Configuration We use the Adam optimizer with a learning rate of $1e-3$ and batch size of 1024. Using these parameters, we perform BPR loss computation for each positive interaction. Training is conducted for a maximum of 50 epochs with early stopping based on validation performance. All experiments are implemented in PyTorch and performed on Google Colab T4 GPU.

Reproducibility For consistent results, we fix random seeds across all experiments. Detailed hyperparameter configurations and additional implementation specifics are provided in Appendix A to ensure full reproducibility of our results.

5 Results and Analysis

5.1 Overall Performance

Table 2 presents a comprehensive performance comparison between our Topic-Enhanced LightGCN (TE-LGCN) and the vanilla LightGCN baseline across all evaluation metrics. Our method demonstrates consistent and substantial improvements across all metrics, with particularly notable enhancements in recall and ranking quality.

Table 2: **Overall Performance Comparison.** Our method achieves over 25% improvement across all metrics compared to LightGCN.

Model	Precision@10	Recall@10	NDCG@10
LightGCN	0.0162	0.1622	0.0871
TE-LGCN (Ours)	0.0205	0.2050	0.1091
<i>Improvement</i>	+26.5%	+26.4%	+25.3%

The experimental results show that our Topic-Enhanced LightGCN achieves consistent improvements across all evaluation metrics. Notably, our method demonstrates substantial performance gains exceeding 25% across all metrics. This indicates that our approach significantly enhances recommendation accuracy while successfully identifying a greater number of relevant items for users.

Key Observations

1. **Enhanced Recall Performance** The remarkable improvements in recall metrics empirically demonstrate that our topic-enhanced approach enables the model to capture a broader spectrum of user preferences. This is particularly significant in sparse data environments where traditional collaborative filtering methods exhibit limitations in identifying relevant niche items.
2. **Improved Ranking Quality** The sustained improvements in NDCG scores indicate that our method possesses the capability to not merely recommend more relevant items, but also to rank them more effectively. This suggests that the topic-based information propagation mechanism helps the model achieve a more accurate understanding of item relevance and the intensity of user preferences.

The overall performance results firmly establish that the integration of topic information through our heterogeneous graph construction and enhanced message passing mechanism provides substantial and meaningful improvements over the strong LightGCN baseline. These performance enhancements are achieved while fully preserving the computational efficiency of the original LightGCN architecture, making our approach highly practical for real-world industrial applications.

5.2 Ablation Study

To systematically analyze the contribution of each component, we conduct a comprehensive ablation study. Table 3 demonstrates the performance changes as components are added step-by-step, providing a quantitative evaluation of how each element of our methodology impacts overall performance.

To ensure experimental fairness, we configure the hyperparameters of the objective function as follows. The regularization weight λ_1 is fixed at $1e-5$, which achieved the best results in the baseline, to observe pure performance improvements from each variant. In contrast, the content consistency weight λ_2 is individually tuned to achieve optimal Recall@10 values for each model variant. Detailed

specifications of these hyperparameter settings and sensitivity analysis results are comprehensively covered in Appendix B.

Table 3: **Ablation Study Results.** Step-by-step performance gains by adding components.

Model Variant	Prec@10	Recall@10	NDCG@10	Improv.
LightGCN (baseline)	0.0162	0.1622	0.0871	-
+ Doc2Vec init	0.0193	0.1934	0.1086	+19.2%
+ Topic nodes	0.0198	0.1983	0.1125	+22.3%
+ Complete (TE-LGCN)	0.0205	0.2050	0.1091	+26.4%

Component-wise Contribution Analysis

1. **Strong Effect of Doc2Vec Initialization** Applying only Doc2Vec-based item embedding initialization yielded a substantial 19.2% improvement in Recall@10. This indicates that semantic representations extracted from textual content effectively capture intrinsic relationships between items that are difficult to learn from sparse interaction data alone. Particularly in environments with extreme sparsity like our dataset (94.86%), this empirically demonstrates that content-based prior knowledge provides a robust foundation for collaborative filtering learning.
2. **Additional Enhancement Through Topic Nodes** Adding topic nodes to the heterogeneous graph increased the overall improvement to 22.3%, achieving an additional 3.1 percentage point gain. This shows that topics extracted through LDA can model deeper levels of semantic associations beyond simple textual similarity. Topic-mediated message passing provides information propagation pathways through thematic similarity even when direct user-item connections are absent, significantly enhancing the discoverability of long-tail items.
3. **Synergistic Effects of Complete System** The complete method with all components combined ultimately achieved a 26.4% improvement, showing an additional 4.1 percentage point enhancement. This gradual yet consistent improvement indicates that each component not only contributes independently but also creates complementary synergistic effects. The combination of Doc2Vec’s semantic initialization and topic-based structural enhancement enables the learning of richer and more expressive embedding spaces.
4. **Performance Stability and Consistency** Notably, consistent improvement patterns are observed across all metrics (Precision@10, Recall@10, NDCG@10). This indicates that our method does not bias toward specific metrics but enhances overall recommendation quality in a balanced manner. Particularly, the improvement in NDCG@10 from 0.0871 to 0.1091 demonstrates substantial enhancement in ranking quality.

5.3 Qualitative Recommendation Analysis: Real Case Study

To demonstrate the practical effectiveness of our method, we conduct qualitative analysis of recommendation results for a specific user. The analyzed subject is an active user (User ID 564) who has interacted with 487 movies and shows preferences for diverse genres (Adventure, Drama, Thriller, Comedy, etc.).

5.3.1 Model-wise Recommendation Characteristics Analysis

Limitations of Basic Collaborative Filtering Base LightGCN shows a popularity-based safe recommendation pattern. While the score range is broadly distributed (5.58 \sim 7.09), most recommendations consist of well-known classic films, limiting opportunities for personalized discovery.

Table 4: **Detailed Comparison of Top-6 Recommended Movies by Model.** Movie titles are in bold, followed by score and genres.

Rank	Base LightGCN	Doc2Vec LightGCN	LDA LightGCN	TE-LGCN (Complete)
1st	The African Queen (7.09, Adv/War/Rom)	Shakespeare in Love (7.94, Rom/Hist)	The African Queen (1.54, Adv/War/Rom)	Amadeus (1.56, Dra/Hist/Mus)
2nd	Shakespeare in Love (6.47, Rom/Hist)	Amadeus (7.76, Dra/Hist/Mus)	Amadeus (1.50, Dra/Hist/Mus)	Shakespeare in Love (1.55, Rom/Hist)
3rd	Amadeus (6.36, Dra/Hist/Mus)	The African Queen (7.56, Adv/War/Rom)	Shakespeare in Love (1.46, Rom/Hist)	The Postman (1.55, Com/Dra/Rom)
4th	The Maltese Falcon (6.00, Mys/Cri/Thr)	The Maltese Falcon (7.52, Mys/Cri/Thr)	Hoop Dreams (1.46, Documentary)	Four Weddings and a Funeral (1.55, Com/Dra/Rom)
5th	The Bridge on the River Kwai (5.74, Dra/Hist/War)	The Bridge on the River Kwai (7.26, Dra/Hist/War)	Babe (1.44, Fan/Dra/Com)	The African Queen (1.52, Adv/War/Rom)
6th	Bonnie and Clyde (5.58, Cri/Dra)	Singin' in the Rain (7.10, Com/Mus/Rom)	The Maltese Falcon (1.43, Mys/Cri/Thr)	The English Patient (1.50, Dra/Rom/War)

Semantic Enhancement Effects of Doc2Vec Initialization Interesting changes are observed in the Doc2Vec-enhanced model. Overall score increases (7.10 ~ 7.94) demonstrate confidence improvement effects. “Shakespeare in Love” rises to first place, and the emergence of “Singin’ in the Rain” (6th place) indicates more sophisticated capture of users’ emotional/artistic preferences through textual content.

Enhanced Discoverability Through Topic-based Approach The most notable result in the LDA-only model is the 4th place entry of “Hoop Dreams” (Documentary). Despite being a genre difficult to discover from existing interaction patterns, this connection through the latent topic of “fact-based narrative” that users prefer demonstrates that topic-mediated recommendation effectively captures semantic associations across genre boundaries.

5.3.2 Synergistic Effects of the Complete Model

Complex Preference Modeling In the TE-LGCN complete model, “Amadeus” takes first place (1.56 points), which is interpreted as the best capture of users’ latent preferences for music/art films. Particularly, the 6th place entry of “The English Patient” (1.50 points) is evaluated as a recommendation that sophisticatedly combines users’ complex preferences (narrative depth of war films + emotional elements of romance films).

Balance of Genre Diversity The positioning of romantic comedies like “The Postman” and “Four Weddings and a Funeral” side by side in 3rd-4th places in the complete model’s recommendation list indicates personalized recommendations that reflect users’ multifaceted tastes beyond simple popularity bias.

Enhanced Musical Elements The appearance of “Singin’ in the Rain” in the Doc2Vec model and “Amadeus” taking first place in the complete model both demonstrate effective capture of users’ preferences for musical/artistic elements that were not explicitly expressed in interaction history.

5.3.3 Implications of Recommendation Score Distribution

Confidence Improvement The overall score increase in the Doc2Vec model (Base: 5.74 ~ 7.09 → Doc2Vec: 7.26 ~ 7.94) indicates that content-based signals enhanced the confidence of collaborative

filtering. This signifies that textual similarity enables evidence-based recommendations beyond simple interaction patterns.

Normalized Scoring System The compression of scores in LDA-based models to the $1.44 \sim 1.56$ range is a natural phenomenon due to the characteristics of probabilistic topic distributions. What matters is not the absolute values of scores but the improvement in relative ranking and recommendation diversity.

5.3.4 Practical Implications

Advancement in Personalization Level This case study demonstrates that our method achieves substantial recommendation quality improvement beyond simple performance metric enhancement. Particularly, the ability to meaningfully recommend genres that users have not explicitly interacted with (documentary) or films with complex characteristics will directly translate to improved user satisfaction in real services.

Balance Between Discovery and Safety The complete model achieves an appropriate balance between safe recommendations (4 out of top 5 are common across all models) and new discoveries (“The Postman”, “The English Patient”). This is evaluated as effectively resolving the trade-off between accuracy and novelty required in commercial recommendation systems.

Cross-genre Connectivity The recommendation of “Hoop Dreams” through topic modeling demonstrates our method’s ability to connect different genres through semantic similarity, providing users with serendipitous discoveries while maintaining relevance to their underlying preferences.

These qualitative analysis results empirically demonstrate that our Topic-Enhanced LightGCN provides meaningful value not only in quantitative performance improvement but also in actual user experience aspects.

6 Conclusion and Future Work

In this study, we propose Topic-Enhanced LightGCN (TE-LGCN) to address the data sparsity problem in recommender systems. Our method extends the standard user-item bipartite graph with topic nodes extracted through LDA to construct a heterogeneous graph, significantly improving recommendation performance in sparse environments through Doc2Vec-based initialization and enhanced message passing.

Key Contributions

1. **Heterogeneous Graph Architecture** We propose a novel architecture that effectively integrates content-based signals while maintaining LightGCN’s lightweight structure. Additional information propagation pathways through topic nodes enable connections between users and items without direct interactions, fundamentally mitigating sparsity issues.
2. **Empirical Performance Improvements** Through extensive experiments on the MovieLens dataset, we achieve consistent and substantial performance improvements across all evaluation metrics. Particularly, the 26.4% improvement in Recall@10 demonstrates that our method significantly expands users’ preference coverage.
3. **Collaborative-Content Fusion Paradigm** We present a new paradigm that integrates the advantages of collaborative filtering and content-based filtering. This provides a practical solution that mutually compensates for the limitations of each approach while maintaining computational efficiency.

4. **Systematic Analysis Methodology** Through comprehensive ablation studies and component-wise contribution analysis, we clearly identify the impact of each element of our methodology on performance. This provides important guidelines for future research directions for improving each component.

Future Research Directions Future research requires scalability verification on larger datasets and exploration of generalization possibilities across various domains. Extensions using multi-modal content such as images and videos, and the application of more sophisticated topic modeling techniques would also be interesting research topics. Developing dynamic topic update mechanisms in real-time streaming environments would be a research direction with high practical value.

Our research provides a practical and effective solution for improving recommender system performance in sparse data environments, and demonstrates the potential for new research paradigms through the fusion of collaborative filtering and content-based methods.

References

- [1] X. He, K. Deng, X. Wang, Y. Li, Y. Zhang, and M. Wang, “LightGCN: Simplifying and Powering Graph Convolution Network for Recommendation,” in *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020, pp. 639–648.
- [2] X. Li et al., “Graph-based Collaborative Filtering,” *arXiv preprint arXiv:2302.08191*, 2023.
- [3] J. Xu, J. Zhang, Q. Liu, Y. Wang, X. Li, and Q. Yao, “Hypergraph-based Dynamic Recommender System,” in *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR ’24)*, 2024.
- [4] A. Vaswani et al., “Attention Is All You Need,” in *Advances in Neural Information Processing Systems 30 (NeurIPS 2017)*, 2017.
- [5] J. Yang, H. Du, Z. Wang, C. Xu, and X. Li, “Recommendation as Language Modeling (RecLM),” *arXiv preprint arXiv:2208.05716*, 2022.
- [6] H. Tang, X. Xu, H. Chen, R. Zhang, and X. Lin, “Neural Collaborative Filtering with Contextualized Word Representations,” *arXiv preprint arXiv:1912.12398*, 2019.
- [7] C. Wang and D. M. Blei, “Collaborative Topic Modeling for Recommending Scientific Articles,” in *Proceedings of the 17th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD ’11)*, 2011.
- [8] M. Soltanizadeh, M. Zabihi, M. Analoui, and M. Yaghmaeizadeh, “A novel approach for improving personalized recommendation with knowledge-based features and collaborative filtering,” *Expert Systems with Applications*, vol. 189, Art. no. 116175, Feb. 2022.
- [9] H. Wang et al., “Multimodal Retrieval-Augmented Generation for Recommendation,” *arXiv preprint arXiv:2403.18667*, 2024.
- [10] J. Shin, “Keyword-enhanced recommender system based on inductive graph matrix completion,” Master’s thesis, Hansung Univ., Seoul, South Korea, 2024.

Appendix

A Dataset Construction and Preprocessing Methodology

A.1 Data Filtering and Graph Construction

To ensure meaningful graph connectivity while maintaining computational efficiency, we apply systematic data filtering procedures to the original MovieLens dataset. Our filtering strategy consists of two primary steps designed to balance data quality with sparsity constraints.

K-Core Filtering We implement k -core filtering with $k = 10$ to guarantee minimum connectivity in the user-item interaction graph. This filtering ensures that each user has interacted with at least 10 items and each item has received interactions from at least 10 users, thereby eliminating extremely sparse nodes that would contribute minimal information to the graph neural network learning process.

Positive Feedback Selection To construct binary implicit feedback suitable for collaborative filtering, we retain only positive interactions, defined as ratings ≥ 4 on the original 5-point scale. This threshold selection is based on the assumption that ratings of 4 and 5 represent genuine user preference and positive engagement, while lower ratings may indicate neutral or negative sentiment that could introduce noise in the recommendation learning process.

Final Dataset Statistics After applying these filtering procedures, our final dataset comprises 605 users, 1,234 movies, and 38,405 total interactions, resulting in a sparsity level of 94.86%. This configuration provides a realistic sparse recommendation scenario while maintaining sufficient connectivity for meaningful graph-based learning.

A.2 Textual Content Collection and Preprocessing

Content Acquisition For the 1,234 movies in our filtered dataset, we systematically collect movie overviews from IMDb using the official MovieLens-IMDb linking provided in the MovieLens dataset. This approach ensures data authenticity and consistency with the original dataset curation standards.

Preprocessing Pipeline The raw textual content requires extensive preprocessing to ensure suitability for natural language processing tasks. We implement a comprehensive six-stage preprocessing pipeline:

1. **Case Normalization** All text is converted to lowercase to ensure consistent token matching and reduce vocabulary size while preserving semantic meaning.
2. **Markup and Symbol Removal** HTML tags, special characters, and non-alphabetic symbols are systematically removed using regular expression patterns to eliminate formatting artifacts and focus on semantic content.
3. **Alphabetic Character Retention** Only alphabetic characters are preserved, with punctuation and numerical content removed. For example, “he’s 25-years-old!!!” is transformed to “hes years old”, maintaining semantic structure while eliminating noise.
4. **Tokenization and Stopword Elimination** Text is tokenized into individual words, followed by removal of common stopwords (the, a, an, of, in, on, with, etc.) that occur frequently but carry minimal semantic significance for topic modeling and content analysis.
5. **Lemmatization** Words are reduced to their canonical forms using lemmatization techniques. This process normalizes different inflected forms to their base representation (e.g., “movies,” “movie,”

“movied” \rightarrow “movie”; “running,” “ran” \rightarrow “run”), reducing vocabulary redundancy while preserving semantic meaning.

6. **Short Token Filtering** Tokens with length ≤ 2 characters are eliminated as they typically carry limited semantic information. Only tokens with length ≥ 3 characters are retained for subsequent analysis.

Content Statistics Following preprocessing, we obtain an average of 53.95 tokens per movie description, providing substantial textual content for semantic analysis while maintaining computational tractability. This preprocessing approach ensures high-quality textual features suitable for Doc2Vec embedding generation and LDA topic modeling.

A.3 Data Quality Validation

1. **Coverage Analysis** We verify that 100% of movies in our filtered dataset have corresponding textual content, ensuring no missing data issues in our content-enhanced approach.
2. **Text Length Distribution** The processed text lengths follow a reasonable distribution with sufficient variance to capture diverse movie descriptions while avoiding extremely sparse or overly verbose content that could bias the topic modeling process.
3. **Vocabulary Statistics** The final vocabulary consists of meaningful semantic terms relevant to movie content, with effective removal of noise and redundant tokens through our systematic preprocessing pipeline.

This comprehensive dataset construction and preprocessing methodology ensures high-quality data suitable for evaluating our Topic-Enhanced LightGCN approach while maintaining realistic sparse recommendation scenarios representative of real-world applications.

B Hyperparameter Sensitivity Analysis for Content Consistency Weight λ_2

B.1 Experimental Setup

To ensure fair comparison across different model variants in our ablation study, we conduct systematic hyperparameter tuning for the content consistency weight λ_2 while keeping the regularization weight λ_1 fixed at $1e-5$ (the optimal value determined for the baseline LightGCN). This approach allows us to isolate the effect of content integration strategies without confounding factors from different regularization strengths.

For each model variant, we evaluate λ_2 values from the set $\{0, 1e-4, 1e-3, 1e-2, 1e-1\}$ and select the configuration that achieves the highest Recall@10 performance, as this metric best reflects our objective of improving recommendation coverage in sparse data environments.

B.2 Doc2Vec-Only Model Results

Analysis For the Doc2Vec-only model, optimal performance is achieved at $\lambda_2 = 1e-1$. This relatively high value suggests that enforcing consistency between learned embeddings and Doc2Vec initializations is crucial when topic-based structural enhancement is absent. The content consistency term compensates for the lack of additional graph structure by maintaining semantic relationships encoded in the pre-trained Doc2Vec representations.

Table 5: **Performance of Doc2Vec Initialization** with Different λ_2 Values. Optimal performance ($\lambda_2 = 1e - 1$) is highlighted.

λ_2	HitRate@10	Prec@10	Recall@10	NDCG@10
0	0.1868	0.0187	0.1868	0.1058
1e-4	0.1868	0.0187	0.1868	0.1059
1e-3	0.1901	0.0190	0.1901	0.1084
1e-2	0.1884	0.0188	0.1884	0.1081
1e-1	0.1934	0.0193	0.1934	0.1086

Table 6: **Performance of LDA Topic Nodes** ($\lambda_2 = 0$ by design).

Model Component	HitRate@10	Prec@10	Recall@10	NDCG@10
LDA Topic Nodes	0.1983	0.0198	0.1983	0.1125

B.3 LDA Topic Nodes Model Results

Analysis The LDA-only model does not utilize content consistency loss ($\lambda_2 = 0$ by design) as it relies purely on structural enhancements through topic nodes without Doc2Vec initialization. The strong performance (Recall@10 = 0.1983) demonstrates the effectiveness of topic-mediated information propagation as a standalone enhancement.

B.4 Complete Model (LDA + Doc2Vec) Results

Table 7: **Performance of Complete Model** with Different λ_2 Values. Optimal performance ($\lambda_2 = 1e - 2$) is highlighted.

λ_2	HitRate@10	Prec@10	Recall@10	NDCG@10
0	0.2000	0.0200	0.2000	0.1098
1e-4	0.1901	0.0190	0.1901	0.1013
1e-3	0.1983	0.0198	0.1983	0.1077
1e-2	0.2050	0.0205	0.2050	0.1091
1e-1	0.1934	0.0193	0.1934	0.1065

Analysis For the complete model combining both LDA topic nodes and Doc2Vec initialization, optimal performance is achieved at $\lambda_2 = 1e - 2$. This moderate value indicates that when structural topic-based enhancements are present, a balanced approach to content consistency is most effective. Too high values ($\lambda_2 = 1e - 1$) overly constrain the learning process, while too low values ($\lambda_2 = 1e - 4$) underutilize the semantic initialization.

B.5 Cross-Model Comparison and Insights

Key Observations

1. **Complementary Enhancement Pattern** The optimal λ_2 decreases as more structural enhancements (topic nodes) are added, suggesting that content consistency and structural enhancement serve complementary roles.

Table 8: **Optimal Configuration Summary.**

Model Variant	Optimal λ_2	Best Recall@10	Gain*	* Relative to baseline
Doc2Vec Only	1e-1	0.1934	+19.2%	
LDA Only	N/A (0)	0.1983	+22.3%	
Complete Model	1e-2	0.2050	+26.4%	
LightGCN (Recall@10 = 0.1622)				

2. **Diminishing Returns Effect** When topic nodes provide structural semantic connections, the need for explicit content consistency constraints diminishes, as evidenced by the lower optimal λ_2 in the complete model compared to the Doc2Vec-only variant.
3. **Robustness Analysis** The complete model shows more robust performance across different λ_2 values, with less dramatic performance drops compared to the Doc2Vec-only model, indicating that structural enhancements provide stability.
4. **Synergistic Effects** The complete model’s best performance (0.2050) exceeds the sum of individual improvements, confirming synergistic effects between Doc2Vec initialization and topic-based structural enhancement.

B.6 Practical Implications

These results provide important guidance for practical deployment:

1. When computational resources limit the use of topic nodes, Doc2Vec initialization with $\lambda_2 = 1e - 1$ provides substantial improvement.
2. For maximum performance, the complete model with $\lambda_2 = 1e - 2$ is recommended.
3. The robustness of the complete model makes it suitable for production environments where hyper-parameter sensitivity is a concern.