

# 네트워크 표현 학습에 기반한 다이버시티를 활용한 추천 시스템 설계

곽창수\*, 서창원\*, 신원용<sup>o</sup>

## Design of Recommender Systems Exploiting Diversity Based on Network Representation Learning

Changsoo Kwak\*, Changwon Seo\*,  
Won-Yong Shin<sup>o</sup>

### 요 약

최근 네트워크 임베딩 기술을 사용하여 추천 시스템을 개발하는 연구가 활발히 이루어지고 있다. 하지만 대부분의 연구는 추천 정확도를 개선하는 방향으로 수행되었으며, 정확도 이외에 사용자 만족도를 높이는 데 있어 중요한 척도로 고려될 수 있는 다이버시티를 개선하는 방향으로의 접근은 활발히 이루어지지 않은 편이다. 본 레터에서는 네트워크 임베딩에 기반을 둔 다이버시티를 개선한 협업 필터링 모델을 제안한다. 제안한 방법이 추천 정확도를 거의 떨어뜨리지 않으며 intra-list distance (ILD) 및 aggregate diversity 측면에서 우수성을 검증한다.

**Key Words** : Collaborative filtering, Diversity, Network embedding, Network representation learning, Recommender system

### ABSTRACT

Recently, it has been actively studied how to

develop recommender systems using network embedding. However, most studies were carried out in the sense of improving the recommendation accuracy, and thus approaches on exploiting diversity, thought of as another important measure in enhancing the quality of experiences, have been underexplored. In this letter, we propose a collaborative filtering model that improves the diversity based on network embedding. It is demonstrated that the proposed method is beneficial in terms of intra-list distance and aggregate diversity at the cost of slight accuracy reduction.

### 1. 서 론

추천 시스템 연구는 사용자로부터 누적된 평점, 행동 이력 등을 모아 새로운 아이템에 대한 선호도를 예측하는 협업 필터링을 중심으로 이루어지고 있다. 이러한 협업 필터링은 메모리 기반 협업 필터링과 모델 기반 협업 필터링 두 가지로 분류될 수 있는데, 메모리 기반 협업 필터링은 비슷한 성향을 가진 사용자나 유사한 아이템의 과거 선호도 (예: 영화 추천 시스템에서의 평점)를 이용하여 예측하는 반면, 모델 기반 협업 필터링은 과거 선호도의 패턴으로부터 사용자와 아이템의 latent representation을 추출하고, 이를 바탕으로 사용자-아이템 간 선호도를 예측하는 모델이다. 전자의 경우 간단히 구현할 수 있는 반면, 후자의 경우 행렬 분할, 차원 축소 기법 등을 이용하여 저장 공간을 절약하면서 메모리 기반 방식 대비 sparsity를 극복할 수 있고 scalable한 것으로 알려져 있다.

하지만 이러한 기존 협업 필터링 방식은 정확도에만 초점을 맞춰 사용자에게 편향된 정보만 제공하게 되는 필터 버블 현상이 발생할 수 있다. 이를 해결하기 위해 다이버시티 (diversity)가 가미된 추천 시스템에 대한 연구에 관심을 가지고 있다. 다양한 아이템을 추천함으로써 사용자는 다양한 정보를 선택할 수 있는 기회를 얻게 되고, 아이템 제공자는 사용자들의 만

\* 이 논문은 2017년도 정부(교육부)의 재원으로 한국연구재단의 지원을 받아 수행된 기초연구사업(2017R1D1A1A09000835)이고, 2019년도 연세대학교 연구비의 지원을 받아 수행된 것임(2019-22-0019).

• First Author : (ORCID:0000-0003-4566-8859)Department of Physics, Chung-Ang University & Department of Computational Science and Engineering, Yonsei University, cieske@hanmail.net, 학생(학사), 학생회원

o Corresponding Author : (ORCID:0000-0002-6533-3469)Department of Computational Science and Engineering, Yonsei University, wy.shin@yonsei.ac.kr, 부교수, 중신회원

\* (ORCID:0000-0003-3499-0578)Department of Computational Science and Engineering, Yonsei University, changwoni@yonsei.ac.kr, 학생(석사)

논문번호 : 202003-047-C-LU, Received March 2, 2020; Revised March 17, 2020; Accepted March 18, 2020

속도를 높이고 더 많은 사용자들이 그들의 서비스를 이용하도록 할 수 있다<sup>11</sup>.

반면, 데이터 마이닝 문제를 해결하기 위해 네트워크 임베딩 (또는 네트워크 표현 학습) 방법이 연구되고 있다. 기존 네트워크 연구에서는 입력으로 인접 행렬을 이용했다면, 네트워크 임베딩 방식은 네트워크를 임베딩 차원으로 투영시키는 벡터화 방법을 이용한다. 자연어 처리에서의 워드 임베딩 아이디어를 적용한 랜덤 워크 (random walk) 기반 네트워크 임베딩 방식인 DeepWalk<sup>12</sup>를 시작으로 네트워크 내 노드 간 1차 및 2차 proximity 정보 학습을 통한 LINE<sup>13</sup>, 그래프 신경망을 이용한 딥 러닝<sup>14</sup>, 기댓값 최대화 알고리즘 기반 네트워크 임베딩 방식<sup>15</sup> 등 다양한 네트워크 임베딩 모델들이 연구되고 있다.

네트워크 임베딩은 각 노드의 latent representation을 사용하기 때문에 기존 방식 보다 메모리 공간 측면에서 이점을 가지고 있다. 또한 평점 이력을 가지고 이분 (bipartite) 그래프화하여 네트워크 임베딩 기법을 사용하면 추천 시스템 문제도 유사하게 해결할 수 있다. 네트워크 임베딩을 이용한 추천 시스템 모델 연구<sup>16</sup>도 진행되고 있지만 다이버시티를 개선시키려는 시도는 거의 진행되지 않았다.

본 레터에서는 네트워크 임베딩에 기반을 둔 다이버시티를 개선한 협업 필터링 모델을 제안한다. 일반적으로 랜덤 워크는 연관성이 높은 노드를 샘플링하기 위해 사용되지만<sup>2,5</sup>, 본 연구에서는 다양한 아이템을 샘플링 하기 위해 랜덤 워크에 다이버시티를 높이는 방향으로 텔레포트 (teleport)를 추가하는 방법을 제시한다. 실험 결과로서, 다이버시티 성능 척도로 intra-list distance (ILD) 및 aggregate diversity (AggDiv)를 사용할 때, 텔레포트 확률이 증가함에 따라 추천 정확도를 거의 떨어뜨리지 않으며 다이버시티가 개선됨을 확인한다.

## II. 방법론

이분 그래프  $G = (U, V, E)$ 에서  $U$ 를 사용자의 집합,  $V$ 를 아이템의 집합이라 하고,  $u_i, v_j$ 는 두 집합의 각각  $i$ 번째,  $j$ 번째 노드로 정의한다.  $E$ 는 이분 그래프 내 사용자와 아이템 간 링크이고, 가중치  $w_{u_i, v_j}$ 는 사용자  $u_i$ 가 아이템  $v_j$ 에 남긴 평점을 의미한다.

랜덤 워크를 수행하기 위해서는 먼저 이분 그래프에서 특정 unipartite 그래프 (예: 아이템-아이템 그래프) 내의 노드들의 숨겨진 관계를 끌어내야 한다. 즉,

랜덤 워크를 통해 아이템-아이템 그래프에서 샘플링을 수행하기 위해서는 아이템 간 관계를 필요로 한다. 이를 위해 본 연구에서는 2차 proximity  $w_{v_i, v_j} = \sum_{k \in U} w_{u_k, v_i} w_{u_k, v_j}$ 를 사용한다<sup>6</sup>. 이 때, 동시에 아이템  $v_i, v_j$ 에 높은 평점을 남기는 사용자가 많을수록 두 아이템 간 링크의 가중치가 커지게 된다. 따라서 가중치가 높은 두 아이템은 연관성이 높다고 할 수 있다.

또한, 사용자와 연관성이 있는 아이템들을 샘플링하기 위해 랜덤 워크 방식을 사용한다. 사용자가 평점을 매긴 아이템들을 시작 노드로 하는 랜덤 워크 수열을 2차 proximity를 이용하여 만든 아이템-아이템 그래프 상에서 생성한다. 이 때 사용자의 평점에 따른 선호도를 반영하기 위해 평점에 비례해서 랜덤 워크 횟수를 조절한다. 아이템-아이템 그래프에 존재하는 링크의 가중치는 아이템 간 연관성을 의미하기 때문에 해당 가중치에 비례하는 확률로 랜덤 워크를 수행하게 되면 연관성이 높은 아이템은 샘플링 되지만 다양한 종류의 아이템이 샘플링 되지 않을 수 있다. 따라서 다이버시티를 높이기 위해 랜덤 워크에 텔레포트 개념을 추가하여 수정된 샘플링을 수행한다. 구체적으로,  $0 \leq \lambda < 1$ 를 랜덤 워크에서의 매 샘플링마다 텔레포트를 수행할 확률이라 할 때,  $1 - \lambda$ 의 확률로 랜덤 워크의  $i$ 번째 노드 (현재 노드)의 주위 노드로 가중치를 준 랜덤 워크를 수행하여  $i+1$ 번째 노드 (다음 노드)를 방문하거나  $\lambda$ 의 확률로 이전까지 해당 워크 내 아이템들이 가지고 있는 모든 장르를 취합하고 한 번도 랜덤 워크 내에 등장하지 않은 장르를 가진 아이템들을 뽑아 그 중 하나로 텔레포트를 수행한다. 텔레포트 이후에는 랜덤 워크 시작 노드와의 연관성을 위해 시작 노드 주위의 노드 중 하나로 돌아온다.

다음으로, 매 랜덤 워크 수열이 완성될 때마다 사용자와 수열 내 아이템들의 임베딩 벡터를 업데이트한다. 수열 내에는 사용자가 평점을 남긴 아이템과 연관성이 높은 아이템, 새로운 장르를 가진 아이템이 포함되어 있는데 모두를 사용자와 연관성이 있는 아이템이라고 간주하고 1차 proximity에 기반한 두 노드 간 확률  $p(u_i, v_j) = \frac{1}{1 + \exp(x_i^T y_j)}$ 를 사용하여 임베딩 벡터를 학습한다. 여기에서  $x_i$ 와  $y_j$ 는 각각  $u_i$ 와  $v_j$ 의 임베딩 벡터이다. 사용자  $u_i$ 가 평점을 매긴 특정 아이템으로부터 시작하여 만들어진 랜덤 워크 수열을

$L_{u_i}$ 로 정의하자. 이 때, 손실 함수  $O = - \sum_{v_j \in L_{u_i}} \log p(u_i, v_j)$ 를 통해 임베딩 벡터의 학습을 반복하여 사용자 벡터 주위에 해당 사용자 입장에서 선호도가 높은 아이템 벡터들과 다양한 장르를 가진 아이템 벡터들이 위치하게 된다.

### III. 실험

데이터셋으로는 943명의 사용자와 1,682개의 아이템, 10만개의 평점 이력으로 구성된 MovieLens-100k를 사용하였다. 10만개의 평점 이력을 training set과 test set의 비율이 8:2가 되도록 나누고 5-fold validation을 수행한다. 먼저 전체 데이터셋을 이용해 이분 그래프와 아이템-아이템 그래프를 만들고 training set을 이용하여 임베딩을 진행한다. 변수로는 임베딩 벡터의 차원  $d = 128$ , 랜덤 워크의 길이  $l = 10$ 을 사용하였다.

성능 측정을 위해 정확도로는  $F_1$  score를, 다이버시티로는 ILD 및 AggDiv를 사용한다. 사용자  $u$ 에게 추천된 추천 리스트를  $R_u$ , ground-truth를  $G_u$ 라 할 때,  $F_1$  score는 precision과 recall의 조화평균으로 주어지고  $precision = \frac{|G_u \cap R_u|}{|R_u|}$  및  $recall = \frac{|G_u \cap R_u|}{|G_u|}$  이

다. 또한, ILD는 추천된 리스트 내 아이тем들의 특성이 얼마나 비슷한지 pairwise하게 계산하여 역으로 추천 리스트의 다이버시티를 평가하는 방법이며, AggDiv는 모든 사용자에게 추천된 중복되지 않은 아이тем의 개수이다.  $s(v_i, v_j)$ 를 아이тем의 장르 기반 Jaccard index라 할 때,  $ILD = 1 - \frac{2}{|R_u|(|R_u| - 1)} \sum_{v_i \in R_u} \sum_{v_j \in R_u, i \neq j} s(v_i, v_j)$ 가 되고<sup>[7]</sup>,

$AggDiv = \left| \bigcup_{u \in U} R_u \right|$ 로 주어진다<sup>[1]</sup>.

표 1은 top-5 recommendation 수행 시 텔레포트 확률  $\lambda$  증가에 따른 정확도 및 다이버시티를 보여준다.

표 1. 텔레포트 확률에 따른 정확도 및 다이버시티  
Table. 1. Accuracy and diversity according to the teleport probability

	$\lambda = 0.0$	$\lambda = 0.2$	$\lambda = 0.4$	$\lambda = 0.6$
$F_1$ score	0.45601	0.45418	0.45148	0.44286
ILD	0.61395	0.62692	0.66392	0.70036
AggDiv	418	425	444	476

다.  $\lambda = 0$ 인 경우는 다이버시티 활용을 하지 않는 기존 방법에 대응한다. 그 결과 텔레포트 확률이  $\lambda = 0.6$ 까지 증가함에 따라  $\lambda = 0$ 인 경우 대비  $F_1$  score가 2.88% 감소하는 대신 다이버시티 metric인 ILD와 AggDiv이 각각 14.07%, 13.87% 증가함을 확인할 수 있다. 즉, 실험 결과로부터 텔레포트 확률이 증가할수록 정확도는 조금 감소하지만 다이버시티가 크게 증가하며 모든 사용자들에게 다양한 아이тем이 추천되면서 각 사용자에게 제공되는 추천 리스트마다 서로 다른 장르를 가진 아이тем들이 많이 들어있다고 판단할 수 있다.

### IV. 결론

본 레터에서는 텔레포트 개념을 추가한 랜덤 워크를 통해 샘플링을 하여 새로운 네트워크 표현 학습을 수행하고, 이를 통해 추천 시스템에서의 다이버시티를 크게 증가시키는 방법을 제안하였다. 실험 결과 제안한 방법은 기존 방법 대비 정확도를 거의 감소시키지 않으면서 두 가지 다이버시티 metric인 ILD 및 AggDiv를 증가시킬 수 있음을 검증하였다.

### References

- [1] G. Adomavicius and Y. Kwon, "Improving aggregate recommendation diversity using ranking-based techniques," *IEEE Trans. Knowledge and Data Eng.*, vol. 24, no. 5, pp. 896-911, May 2012.
- [2] B. Perrozi, R. Al-Rfou, and S. Skiena, "DeepWalk: Online learning of social representations," in *Proc. 20th ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, New York, USA, Aug. 2014.
- [3] J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, and Q. Mei, "LINE: Large-scale information network embedding," in *Proc. 24th Int. Conf. World Wide Web*, Florence, Italy, May 2015.
- [4] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional Networks," in *Proc. Int. Conf. Int. Conf. Learning Representations*, Toulon, France, Apr. 2017.
- [5] G.-T. Park, C. Tran, and W.-Y. Shin,

- “Network-embedding-based link prediction using the expectation maximization algorithm,” *J. KICS*, vol. 44, no. 11, pp. 2123-2126, Nov. 2019.
- [6] M. Gao, X. He, L. Chen, and A. Zhou, “BiNE: Bipartite network embedding,” in *Proc. 41<sup>st</sup> ACM SIGIR Conf. Research & Development in Information Retrieval*, pp. 715-724, Ann Arbor, MI, USA, Jul. 2018.
- [7] Y. Liu, Y. Zhang, Q. Wu, C. Miao, L. Cui, B. Zhao, Y. Zhao, and L. Guan, “Diversity-promoting deep reinforcement learning for interactive recommendation,” *arXiv preprint arXiv:1903.07826*, 2019..