

Deep Learning-Based Personalized Recommendation Using Customer Behavior and Purchase History in E-Commerce

Da Young Hong[†] · Ga Yeong Kim^{**} · Hyon Hee Kim^{***}

ABSTRACT

In this paper, we present VAE-based recommendation using online behavior log and purchase history to overcome data sparsity and cold start. To generate a variable for customers' purchase history, embedding and dimensionality reduction are applied to the customers' purchase history. Also, Variational Autoencoders are applied to online behavior and purchase history. A total number of 12 variables are used, and nDCG is chosen for performance evaluation. Our experimental results showed that the proposed VAE-based recommendation outperforms SVD-based recommendation. Also, the generated purchase history variable improves the recommendation performance.

Keywords : Online Behavior Log, Purchase History, VAE-based Recommendation, Extracting Latent Space

전자상거래에서 고객 행동 정보와 구매 기록을 활용한 딥러닝 기반 개인화 추천 시스템

홍 다 영[†] · 김 가 영^{**} · 김 현 희^{***}

요 약

본 논문은 고객의 온라인 행동 정보와 구매 기록을 활용하여 기존의 추천 시스템이 갖는 데이터 희소성의 문제와 콜드 스타트 문제를 해결하고자 VAE 기반 추천 시스템을 제시하였다. 고객의 구매 기록을 임베딩하고 차원 축소하여 단일 변수로 생성하였으며, 온라인 행동 정보를 활용하여 학습을 통해 고객의 잠재 요인을 추출하는데 Variational Autoencoders를 적용하였다. VAE 기반 추천 시스템은 총 12개의 고객의 특성 변수에 VAE를 적용하여 저차원의 벡터를 생성한 뒤 이를 통해 유사 사용자를 찾는 다음, 유사 사용자들이 구매했던 상품들을 고객에게 추천한다. 이렇게 추천한 상품들의 점수를 매겨 nDCG로 성능을 평가하였다. 실험 결과 구매 기록 변수 그리고 온라인 행동 로그 변수를 활용한 VAE 기반의 추천시스템이 SVD 기반의 추천 시스템보다 성능이 좋다는 것을 알 수 있었다. 따라서 고객의 온라인 행동 로그 및 구매 기록을 사용하여 상품을 추천하면 정보 수집에 발생하는 비용과 시간을 줄일 수 있을 뿐만 아니라 기존 추천 시스템보다 더욱 효율적으로 상품을 추천할 수 있다는 것을 보여주었다.

키워드 : 온라인 행동 정보, 구매 기록 정보, VAE 기반 추천, 잠재 요인 추출

1. 서 론

최근 전자 상거래가 활발해짐에 따라 개인화 추천 시스템의 중요도가 부각되고 있다. 현대 사회에 수많은 상품과 서비스가 존재하기 때문에 각 고객에게 적절한 상품과 서비스를 효과적으로 제공하기 위해서는 추천 시스템이 필요하다. 통상적으로 많은 기업들에서는 평점을 기반으로 한 추천 시스템인 협업 필터링을 사용하고 있다. 하지만 현대 사회는 평점 데이터가 존재하지 않는 상품과 서비스도 존재하기 때문에 개인화 추

천시스템 개발에 많은 기업들이 관심을 갖고 있는 것이 사실이다. 또한 평점은 고객의 단순한 선호도를 표현한 정보이기 때문에 상품과 서비스를 적절하게 추천하기에는 한계가 있다. 예를 들어 음악 추천 시스템의 경우 단순히 고객이 음악에 대한 선호도를 나타낸 것으로 음악을 추천하기에는 한계가 있다. 음악의 장르, 가수, 음악을 듣는 시간대 등 평점으로는 나타내기 어려운 요소들을 활용하였을 때 더욱 적절한 음악을 고객에게 추천할 수 있을 것이다. 이렇듯 평점이라는 단편적인 선호도보다 다른 변수들을 활용하였을 때 더욱 효과적으로 추천 시스템을 활용할 수 있을 것으로 보인다.

앞서 언급한 바와 같이 기존의 추천시스템은 주로 고객의 구매 기록, 각 상품에 대해 매긴 평점 등 사용자가 상품에 대한 선호도를 직접 나타내는 정보를 사용하는 것이 일반적이다. 하지만 평점이나 좋아요 수와 같은 외부적 선호도 정보는

※ 이 논문은 2020년도 동덕여자대학교 학술연구비 지원에 의하여 수행된 것임.

† 준 회 원 : 고려대학교 의학통계학협동과정 석사과정

** 준 회 원 : 성균관대학교 인공지능학과 석사과정

*** 종신회원 : 동덕여자대학교 정보통계학과 부교수

Manuscript Received : October 15, 2021

Accepted : November 18, 2021

* Corresponding Author : Hyon Hee Kim(heekim@dongduk.ac.kr)

상품 전체에 비해 그 양이 매우 적은 데이터 희소성의 문제를 안고 있다. 또한 새로운 고객에게 추천을 할 때에도 해당 고객에 대한 구매 기록이 존재하지 않아 어떤 상품을 선호하는지 알 수 없기 때문에 상품 추천이 어렵고, 새롭게 등록된 상품 또한 고객의 외부적 선호도가 축적되기 전까지는 추천 상품 리스트에 진입하는 것이 어렵다.

본 논문에서는 이와 같은 추천 시스템의 문제점을 극복하기 위해서 전자 상거래 고객들의 온라인 행동 로그를 이용한 새로운 추천시스템을 제안하고자 한다. 온라인 행동 로그는 전자상거래에서 발생하는 고객의 행동 정보를 기록한 것으로 고객이 상품을 조회한 시각이나 고객이 어떤 채널로 유입했는지, 상품을 사기까지 웹 사이트에 머문 시간 등 외부적으로 정보를 제공하지 않아도 암시적인 고객의 온라인 행동을 수집할 수 있다. 이러한 고객의 전자상거래 상에서 일어난 행동 정보는 온라인 행동 로그의 형태로 자동으로 수집되기 때문에 데이터 수집의 시간과 비용을 단축할 수 있으며, 상품 추천에 효과적인 변수로 활용하여 추천시스템의 성능을 높일 수 있다.

최근 딥러닝 기반의 추천 시스템이 많은 관심을 받고 있다 [1-3]. 딥러닝 기반의 추천 시스템은 기존의 행렬 분해 (Matrix Factorization)[4] 기반의 추천 시스템에서 잠재 공간을 찾아내는 과정에 딥러닝을 적용하면 추천 성능을 향상시킬 수 있다는 결과를 보여주었다[5,6]. 특히 딥러닝 기법 중에서도 행렬 분해 과정에 Variational Autoencoders(VAE)[7]를 적용한 연구[8,9]가 좋은 추천 성능을 보여주었으며, 본 연구에서 활용한 바와 같이 고객의 암시적 정보를 활용하고 딥러닝으로 학습하면 성능이 향상된다는 것을 보여주었다[10].

본 논문에서는 7개의 범주형 변수와 6개의 연속형 변수로 구성된 고객 행동 정보를 사용하여 사용자 특성의 잠재 공간을 추출하기 위해 VAE모형을 사용하였다. 고객의 구매 내역 정보를 변수로 활용하기 위하여 상품 아이템을 임베딩하여 차원 축소를 실시한 뒤 단일 변수로 사용하였다. 이렇게 추출한 잠재 공간을 통해 고객 별 유사 사용자를 찾아내어 이들이 구매한 상품 목록을 통해 해당 고객에게 상품을 추천하여 성능 평가를 진행했다. 성능 평가 과정에서 모델, 변수의 변화, 그리고 고객 별 유사 사용자 수에 따라 추천 시스템의 성능을 평가하고자 하였다.

제안하는 VAE 기반 추천 시스템의 성능을 평가하고자 온라인 쇼핑몰 고객 정보인 L.Point 데이터를 사용하였으며, 특이값 분해(Singular Value Decomposition, SVD)를 행렬 분해에 적용한 SVD 모델과 VAE를 행렬 분해에 적용한 VAE 모델을 비교하였다. 또한 추천 시스템에 활용하는 변수의 중요도를 파악하기 위해서 구매 내역 변수의 유무에 따른 성능의 향상 정도를 알아보려고 하였다. 마지막으로 고객 별 유사 사용자의 수에 따른 성능의 변화도 알아보기 위해 유사 사용자를 20명, 30명, 40명으로 설정하여 상품을 추천하여 실험하였다. 성능 평가 방식으로는 추천 순서에 따른 추천 정

확도를 평가할 수 있는 nDCG(normalized Discounted Cumulative Gain)[11]를 사용하였다.

실험 결과, 모든 상황에서 VAE 모델을 사용하여 상품을 추천했을 때의 성능이 SVD 모델을 활용했을 때보다 높았다. 또한 변수의 변화 측면에서는 기본변수와 온라인 행동 로그 변수만을 사용했을 때보다 기본변수와 온라인 행동 로그 변수와 구매 내역 변수 모두를 활용하여 상품을 추천하는 경우에 높은 성능을 보였다. 유사 사용자의 수에 따른 변화로 보면 고객 당 20명의 유사 사용자들의 상품을 추천 했을 때의 성능이 가장 높았다. 이를 종합하였을 때 VAE 모델을 기반으로 하여 기본변수, 온라인 행동 로그 변수, 구매 내역 변수를 활용하여 고객 별 20명의 유사 사용자를 찾아 상품을 추천했을 때 가장 높은 성능을 보인다.

본 연구의 공헌은 현재 통상적으로 사용되고 있는 추천 시스템이 가지고 있는 데이터 희소성 문제나 콜드 스타트 문제 등의 문제를 해결할 수 있다는 데에 있다. 또한 평점뿐만이 아닌 고객의 암시적 정보를 잠재 공간을 활용하기 때문에 상품을 추천하는 데 눈으로 확인하기 어렵지만 유의미한 변수들의 중요성을 알 수 있다. 이는 상품 추천 서비스뿐만이 아닌 유사 사용자를 찾아 활용하는 다른 서비스의 추천 시스템에도 활용이 가능하다는 것을 의미한다.

본 연구의 구성은 다음과 같다. 2장에서는 온라인 행동 로그를 사용한 추천 시스템이나 VAE를 사용한 추천 시스템에 관한 연구들에 대해서 다룬다. 3장에서는 데이터의 전처리 과정과 latent space를 통해 유사 사용자를 찾아 상품을 추천하며 이에 대한 성능 평가를 진행한다. 4장에서는 실험 결과를 분석하고 5장에서는 결론 및 향후 연구 방향을 제시하고자 한다.

2. 관련 연구

최근 고객의 온라인 행동 정보를 활용하여 추천 성능을 향상시키고자 하는 연구가 관심을 받고 있다[12]. [13]에서는 고객의 온라인 구매 행동을 고려한 토픽 모델링을 통해 도서를 추천하였다. 토픽 모델링은 구조화된 텍스트 데이터에서 각 문서의 주제를 탐색하는 방법으로 최근 텍스트 마이닝에서 대표적으로 사용되는 기법이다. 책 내용과 목차에서 주제를 추출하여 아이템의 특성을 파악하였으며, 온라인 사용자의 구매 행태를 활용하여 도서를 추천하였다. 이 때 온라인 사용자의 구매 행태는 사용자의 인구 통계학적 특성과 도서의 가격이나 구매 계절, 그리고 사용자의 검색 기록 등을 활용하여 최신성과 가격을 추천시스템에 반영하고자 하였다. 그 결과 기존의 협업 필터링보다 높은 성과를 보이며 기존의 협업 필터링이 가지는 한계를 일부 해결하였다.

음악 추천 연구[14]에서는 본 연구에서 구매 내역을 임베딩하여 사용한 것과 유사한 방식으로 음악 플레이리스트 데이터를 임베딩하여 새로운 음악을 추천하는 연구를 진행했

다. 음악은 다른 콘텐츠와는 다르게 상품에 대한 특정 정보 이외에도 플레이리스트라는 개념이 존재하기 때문에 이를 통해 음악 간 관계를 분석하여 유사도를 측정해 새로운 음악을 추천하였다. 플레이리스트에 있는 곡들을 곡의 태그, 장르, 세부 장르 등을 Word2Vec 알고리즘으로 알려진 SGNS에 적용하여 다차원 벡터 공간에 임베딩하였다. 임베딩 후에는 각 곡의 특징 벡터를 생성하였다. 이렇게 생성된 임베딩된 곡들을 코사인 유사도를 기반으로 새로운 곡을 추천하는 방식으로 연구를 진행하였다. 이 때 비교 대상으로 Item2Vec 방식으로 선정하여 nDCG를 통해 성능을 평가했다. Item2Vec 방식의 nDCG는 0.1850, SGNS 방식의 nDCG는 0.2996으로 SGNS 방식이 더욱 높은 성능을 보였다.

최근에 딥러닝 기법들이 다양해지면서 이를 활용한 다양한 추천 시스템에 관한 연구들이 이루어지고 있다. 고객의 구매 히스토리, 이미지 및 텍스트 정보를 활용하여 개인화 추천에 적용한 연구[15]가 진행되었으며, 그 중에도 VAE를 활용한 협업필터링이 많은 관심을 받고 있다. [16]의 연구는 VAE 모델은 사용자의 선호도에 영향을 미치는 요인을 더 정확하게 집어낼 수 있었으며, 잠재인자의 개수에 따라 VAE의 결과가 달라짐을 보여주었다. 또한 협업 VAE 추천시스템이 잠재공간에서 확률론적 분포를 추론하기 때문에 추천을 위해 다양한 멀티미디어 양식을 통합할 수 있는 모델을 제시하였다.

본 연구에서는 기초적인 고객 정보, 온라인 행동 정보, 그리고 구매 내역 정보를 통합한 VAE 기반의 추천 시스템을 제시하였다. 온라인 행동 정보를 활용하기 위해서 VAE를 사용하여 잠재 요인을 찾아내고, 구매 내역을 활용하기 위해서 아이템을 임베딩하고 차원 축소를 실시하였다. 추천 순위를 고려한 nDCG 성능 평가 결과 제안하는 추천 시스템이 추천 성능을 향상시킬 수 있다는 것을 보여주었다.

3. 실험 설계

3.1 데이터 소개 및 전처리

본 논문에서는 롯데멤버스, 제6회 L.POINT Big Data Competition에서 제공하는 온라인 행동 정보, 거래 정보, 고객 정보, 그리고 구매 기록 데이터를 사용했다. 해당 데이터는 세션별로 되어 있는 3,196,362개의 데이터로, 이 중 구매 확정을 한 22,239개의 고객의 세션 데이터로 추출한 후 이를 8,851명의 고객 별 데이터로 변환하여 연구를 진행하였다.

Table 1은 본 연구에서 사용된 변수들을 보여주고 있다. 기본 변수, 온라인 행동 로그 변수, 그리고 구매 내역 변수의 세 범주로 구분하였으며, 기본 변수에는 고객의 성별, 나이, 총 구매 금액 변수를 포함하고 있다. 고객의 성별과 나이는 범주형 변수로 전처리를 하였고 총 구매 금액 변수는 구매 금액과 구매 수량이라는 변수를 통해 생성하였다. 온라인 행동 로그 변수는 세션 경과 시간, 총 페이지 조회 건수, 총 세션

Table. 1. Categories of Variables

Category	Variable name	Variable description
Base variables	clnt_gender	client gender
	clnt_age	client age
	tot_buy_am	total purchase amount
Online behavior log variables	trfc_src	access channel
	dvc_ctg_nm	device type
	b_unit	business unit
	action_type	action type
	de_dt	demand date
	de_tm	demand time
	tot_pag_view_ct	total page view case
Purchase variable	tot_sess_hr_v	total session hour variable
	hit_pss_tm	hit pass time
	em	purchase history

시간 값, 유입 채널, 기기 유형, 업종 유형, 구매 월 및 구매 시각을 사용하였다.

세션은 고객이 일정 기간 내에 사이트를 방문하여 활동하는 영역으로 사용된 세션 경과 시간(hit_pss_tm)은 세션이 시작된 이후 해당 조회까지 경과한 시간으로 분단위로 변환하여 사용하였다. 총 세션 시간 값(tot_sess_hr_v)은 세션 머문 총 시간을 의미한다. 범주형 변수는 유입 채널(trfs_src)과 기기 유형(dvc_ctg_nm)의 두 가지로, 유입 채널의 경우 DIRECT/ PUSH/ WEBSITE/ PORTAL/ unknown으로 구성된 범주형 변수이며, 유입 채널은 모바일 웹/모바일 앱/PC의 세 가지 범주로 구성되었다. 총 페이지 조회 건수(tot_pag_view_ct) 및 조회 경과 시간(hit_pss_tm)은 세션 별 정보이므로 이를 고객 당 세션의 평균으로 대체하여 사용하였다.

3.2 변수 생성 및 변수 선정

전자상거래에서 고객이 상품을 구매한 구매 내역은 물품 추천에 있어 유용한 정보를 제공할 수 있다. 그러나 한 명의 고객이 장기간에 걸쳐 구매한 물품은 매우 방대할 뿐 아니라, 고객마다 품목이 다르기 때문에 한 개의 변수로 활용하는 것이 어렵다. 본 연구에서는 임베딩과 차원축소를 적용하여 고객의 구매 품목 리스트를 변수 em 값으로 생성하여 사용하였다.

Fig. 1은 고객의 구매 물품 리스트를 한 개의 변수로 생성하는 과정을 나타낸다. 먼저, 날짜별로 분산되어 있는 고객의 구매 물품 리스트를 Fig. 1의 왼쪽 위의 테이블에서와 같이 고객별로 구매 품목 리스트를 생성한다. 구매 품목 리스트에서 사용된 숫자는 물품번호이다. 즉, B 고객은 186번, 151번, 351번, 그리고 189번 물품을 구매했음을 나타낸다. 이렇

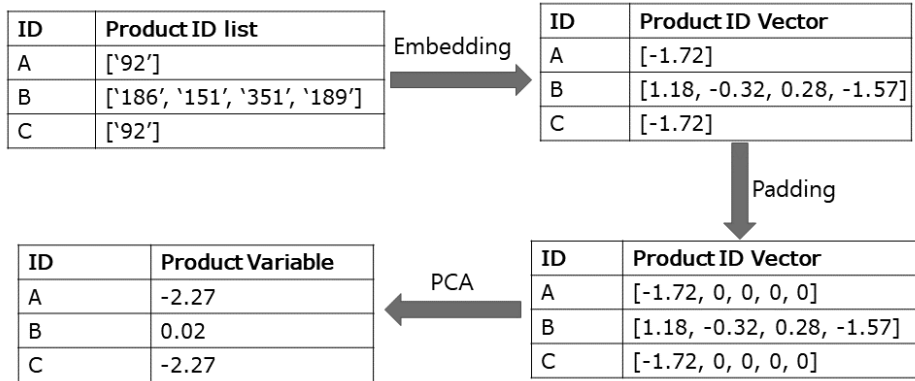


Fig. 1. Process of Generating a Variable of Purchase Lists

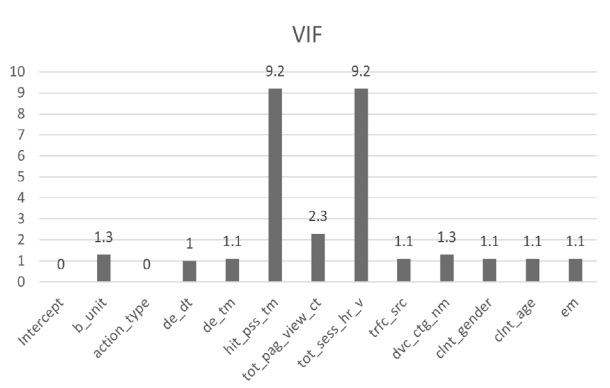


Fig. 2. VIF Value of Each Variable

계 고객마다 리스트의 길이 및 품목 번호가 다르기 때문에 그림의 오른쪽 위에서 보는 바와 같이 임베딩을 통하여 각 물품의 벡터값을 생성한다. 변수로 사용하려면 한 개의 값으로 차원을 축소해야 하고 이를 위해서 주성분분석 (Principle Component Analysis, PCA)을 적용하였다. 주성분분석을 적용하기 위해 모든 리스트의 길이가 동일하여야 하므로 '0'값으로 패딩하였으며, 주성분분석을 통해 최종적으로 Fig. 1의 왼쪽 아래에서 보는 바와 같이 단일값을 생성하여 사용하였다.

총 13개의 변수들이 추천 시스템에 유용한 변수인지 확인하기 위하여 회귀 분석을 통해 다중 공선성을 확인하였다. 총 구매금액(tot_buy_am)이 실제 고객의 선호도를 직접적으로 반영한 변수이므로 총 구매금액을 종속변수로 두고 다중 선형 회귀 분석을 실시한 다음 Variance Inflation Factors (VIF) 값을 계산하였다. Fig. 2는 각 변수의 VIF값을 나타낸다. 조회 경과 시간(hit_pss_tm)과 총 세션 시간(tot_sess_hr_v) 사이의 다중 공선성이 매우 높게 나타났으므로 조회 경과 시간(hit_pss_tm) 변수를 제거하였으며, 조회 경과 시간(hit_pss_tm) 변수를 제거하고 나면 총 세션 시간(tot_sess_hr_v)의 VIF 값은 5 이하로 감소하므로 추가로 다른 변수들을 제거하지 않았다. 따라서 조회 경과 시간 (hit_pss_tm) 변수를 제외한 총 12개의 변수를 실험에 사용하였다.

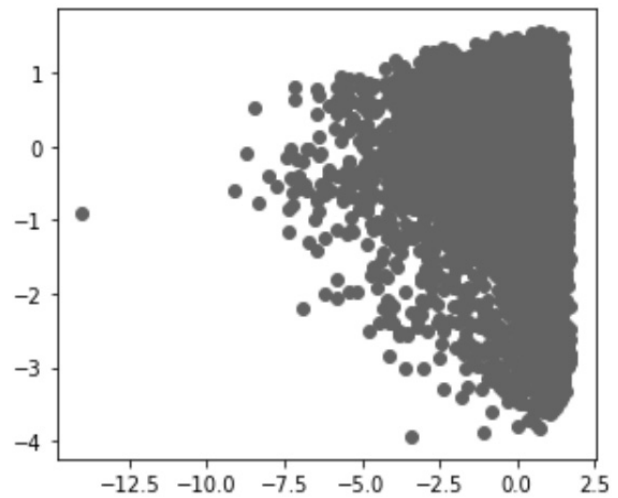


Fig. 3. Visualization of Latent Space

3.3 VAE 기반 추천 알고리즘

선정된 변수들로부터 잠재 요인을 추출하고자 행렬 분해를 하는 과정에서 VAE[17]를 적용하였다. VAE는 기존의 Autoencoder 모델에 평균과 분산을 구할 수 있는 층을 추가하여 보다 확률적인 모델을 구현할 수 있으며, 학습이 진행되면 잠재 공간은 고객의 주요한 특성을 2차원 형태로 함축적으로 표현하게 된다. Fig. 3은 사용자의 특징을 갖는 z벡터들을 시각화하여 표현한 것이다. Fig. 3에서 각각의 점들은 고객을 나타내고, 점사이의 거리가 가까울수록 유사한 고객을 나타낸다.

유사한 고객을 찾기 위해서 코사인 유사도를 사용하였다. 코사인 유사도 결과로 유사 고객을 찾은 후, 유사 고객이 구매했던 상품들의 추천 횟수를 계산하고 추천수가 높은 순서대로 상품을 정렬한 뒤 n개의 상품을 추천하였다. Algorithm 1은 VAE 기반 추천 시스템의 유사 코드이다.

본 논문에서 사용된 데이터 및 코드는 <https://github.com/YOUNG0308/VAE-recommendation> 에서 다운로드하여 확인할 수 있다.

Algorithm 1. Pseudo-Code for VAE-based recommendation

Input : X : clients' online behavior data
 Y : clients' information data
 Z : transaction data

Output : A : preprocessed online behavior data of clients' purchase history
 B : data of selected variables by VIF
 C : latent space of predicted X by VAE
 D : list of K similar users
 E : top-N recommended products' score list per user

STEP 1. Preprocessing
Read X, Y, Z
Merge X, Y, Z with clnt_id and trans_id
Preprocess merged data
 A ← preprocessed data
Return A

STEP 2. Selecting variables by VIF
Read A
Apply VIF algorithm into variables of A
If variable's VIF of A > 5 **then** drop the variable
 B ← data of selected variables
Return B

STEP 3. Find latent space by VAE
Read B
Perform VAE
 C ← latent space of predicted B
Return C

STEP 4. K-NearestNeighbors
Read C
Decide K, which is the optimal number of similar users
Apply K-NearestNeighbors algorithm into A using the K
 D ← list of K similar users per user
Return D

STEP 5. Score of recommended products
Read D
Create list of all products which is similar users purchased
Count how many recommendations for each product per person
Divide the number of times by average recommendation of product
 E ← Score of recommended products
Return E

4. 실험 결과

제안하는 VAE 기반 추천 시스템의 성능을 평가하기 위해서, 동일한 변수를 사용하는 추천 시스템에 잠재 공간을 찾기 위해 SVD를 적용한 추천 시스템 [18]과 비교 평가를 수행하였다. 본 연구를 위해 생성하여 사용한 구매 내역 변수가 추천 시스템의 효율성을 향상시켰는지 확인하기 위하여 모든 실험은 기본변수와 행동로그 변수로 구성된 추천 시스템과 기본변수, 행동 로그 변수, 그리고 생성된 구매 변수로 구성된 추천 시스템에 대해 성능 평가를 실시하였다. 마지막으로 추천 상품의 개수를 10, 20, 30, 40, 그리고 50개로 증가시키면서 성능 평가를 실시하였으며, 이때 유사 사용자의 수는 20, 30, 그리고 40명으로 달리하여 실험을 실시하였다. 실험은 Intel Core i5-7200U CPU, RAM 8GB, windows 10 환경에서 진행되었다.

성능 평가 기준으로는 nDCG(normalized Discounted Cumulative Gain)[19]를 사용하였다. nDCG는 DCG 평가 지표를 정규화한 것으로, DCG[20]는 수식 (1)에서 보는 바와 같이 추천 결과들을 동일한 비중으로 계산한 CG에서 랭킹 순서대로 비중을 줄여 관련도를 계산한 것이다. 하지만 nDCG는 정규화 과정을 거쳤기 때문에 사용자마다 추천된 상품의 개수가 달라도 평가를 할 수 있다. nDCG는 수식 (2)에서 볼 수 있는 것처럼 순위가 있는 추천 시스템에 사용되는 평가지표로 관련성이 높은 결과를 상위권에 노출시켰는지를 기반으로 한다. 이는 하위권의 결과 예측보다 상위권의 결과 예측이 중요한 추천 시스템에 사용될 때 유용하다.

$$DCG_p = \sum_{i=1}^p \frac{rel_i}{\log_2(i+1)} \quad (1)$$

$$nDCG_p = \frac{DCG_p}{IDCG_p} \quad (2)$$

본 논문에서는 고객 별 유사 사용자들이 많이 구매한 상품을 기준으로 순위를 매겨 상품을 추천하였기 때문에 랭킹이 있는 추천 시스템에 적용하기 용이한 nDCG를 활용하였다. 다음 세 가지 그래프를 통해 각 상품 추천 모델의 추천 상품의 개수에 따른 nDCG의 변화를 볼 수 있다.

Fig. 4는 상위 20명의 사용자에게 대해 추천 상품 갯수를 10에서 50까지 증가시키면서 추천한 결과이다. VAE 기반의 추천은 파란색과 주황색으로 나타나고, SVD 기반 추천이 초록색과 붉은 색으로 나타난다. VAE 기반 추천이 SVD 기반 추천 시스템보다 좋은 성능을 보이는 것을 알 수 있으며, 두 가지 경우 모두 구매 내역 변수를 추가로 사용하였을 때 성능이 좋아지는 것을 볼 수 있다. 또한 추천하는 상품의 개수가 늘

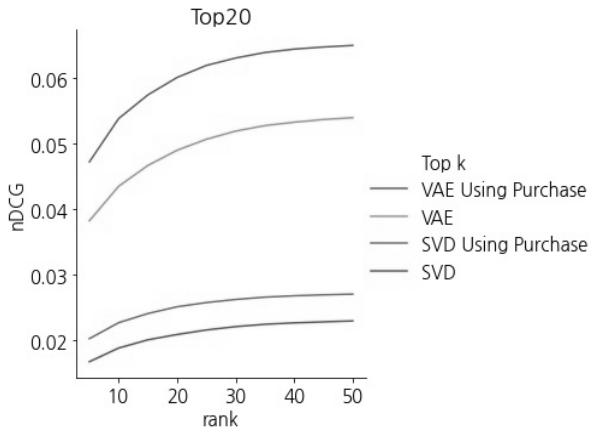


Fig. 4. Performance Evaluation using Top 20 Similar Customers

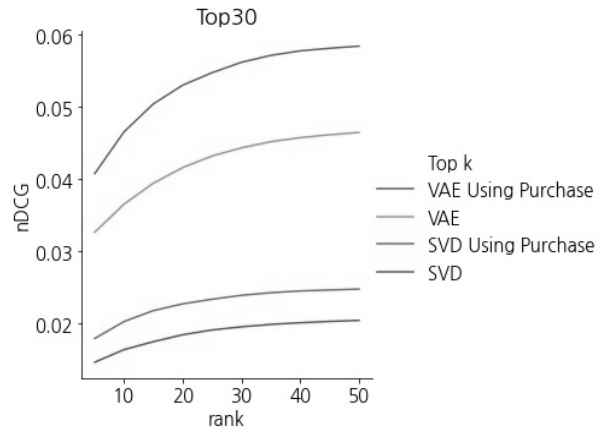


Fig. 5. Performance Evaluation using Top 30 Similar Customers

어남에 따라 성능이 향상되는 모습을 볼 수 있다.

Fig. 5에서 나타난 바와 같이 유사 사용자를 30명으로 설정한 경우도 구매 내역 변수를 포함한 경우가 기본 변수와 온라인 행동 로그 변수만을 이용했을 때보다 높았으며, 역시 추천 상품의 개수가 늘어남에 따라 추천 성능이 향상하는 모습을 보인다. 또한 Fig. 4와 Fig. 5 모두 같은 변수들을 사용하더라도 SVD 기반 추천을 했을 때보다 VAE 기반 추천을 이용하여 상품을 추천한 경우의 성능이 월등히 높다는 것을 알 수 있다.

마지막으로 Fig. 6은 40명의 유사 사용자를 찾아 상품을 추천한 다음 성능 평가를 실시한 경우이다. 이 경우에도 12가지의 모든 변수를 사용한 VAE 기반 추천이 모든 경우에서 가장 성능이 높은 것을 알 수 있다. 이 경우에도 역시 추천하는 상품의 개수가 많아질수록 성능이 향상한다.

위의 실험을 통해 얻은 결과를 요약하면 다음과 같다. 먼저, VAE 기반 추천 시스템이 SVD 기반 추천 시스템보다 성능이 높다. 또한, VAE 기반 추천 시스템과 SVD 기반 추천 시스템 모두 기본변수, 온라인 행동 로그 변수 그리고 구매 내역 변수까지 모든 변수를 활용했을 때의 성능이 가장 높다는 것을 알 수 있었다. 세 가지 실험을 종합적으로 보면 고객 별 유사 사용자의 수가 20명일때 추천 시스템의 성능이 가장 좋았으며, 유사 사용자의 수를 증가시킬수록 추천 시스템의 성능이 저하됨을 알 수 있었다. 이를 통해 유사 사용자의 수가 많을수록 추천 시스템의 성능이 향상되는 것은 아니며 적절한 수의 유사 사용자를 통해 상품을 추천하는 것이 추천 시스템의 성능 향상에 도움이 될 것이라는 결론을 내릴 수 있다.

5. 결론 및 향후 연구

현재 많은 곳에서 사용하고 있는 추천 시스템은 평점 데이터를 사용하여 고객에게 상품을 추천한다. 하지만 이러한 기존의 추천 시스템은 평점 데이터를 기반으로 하기 때문에 콜

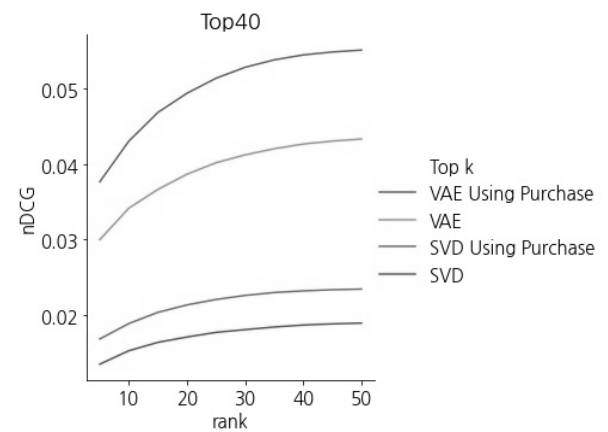


Fig. 6. Performance Evaluation using Top 40 Similar Customers

드스타트 문제와 데이터 희소성 문제를 가지고 있다. 이러한 문제가 발생할 경우 고객에게 적절한 상품을 추천하기가 어려운 상황을 극복하기 위하여 본 논문에서는 고객의 온라인 행동 로그 데이터와 구매 기록 데이터를 활용한 추천 시스템을 제시하였다. 고객 특성의 잠재 요인을 유사한 고객을 찾는 데 활용하기 위해 VAE를 적용하였으며, 구매 기록 데이터를 변수로 생성하기 위해서 임베딩과 차원 축소를 진행하여 구매기록 변수를 생성하였다.

제안하는 추천 시스템의 성능을 평가하기 위해 기본 변수, 온라인 행동 로그 변수, 구매 내역 변수로 변수들을 나누었고 기본 변수를 기반으로 온라인 행동 로그 변수와 구매 내역 변수를 활용했을 때의 추천 시스템 성능의 변화를 보고자 실험을 진행하였다. 실험 결과, 고객의 기본 변수와 온라인 행동 로그 변수, 구매 내역 변수를 활용하여 상품을 추천하는 경우가 기본 변수와 온라인 행동 변수만을 활용했을 때보다 높은 성능을 보였다. 이를 통해 구매 내역 변수가 상품을 추천하는 데에 상당한 영향을 미친다는 것을 알 수 있다. 즉, 고객이 구매했던 상품들이 또 다른 상품을 구매하는 데에 중요한 요소임을 알 수 있다.

또한 같은 변수들을 사용하더라도 SVD 기반 추천이 아닌 VAE 기반 추천을 사용하여 상품을 추천했을 때의 성능이 높다는 것을 알 수 있었으며 이를 통해 제안하는 모델의 성능이 우수함을 증명하였다. 또한 고객 별 유사 사용자의 수에 따라서도 추천 시스템의 성능이 달라지는 것을 알 수 있다. 이는 적절한 고객 별 유사 사용자의 수를 활용하여 상품을 추천한다면 더욱 효율적인 추천 시스템을 개발할 수 있다는 것을 의미한다.

위 실험 결과를 바탕으로 고객의 온라인 행동 로그와 VAE 모델을 통해 상품 추천을 효율적으로 할 수 있다는 것을 알 수 있다. 온라인 행동 로그는 정보 수집에 용이하기 때문에 데이터 수집에 소요되는 비용과 시간을 줄일 수 있으며 더욱 효율적으로 상품을 추천할 수 있다. 또한 본 논문에서 활용한 변수들 이외에도 다양한 온라인 행동 로그를 수집하여 고객의 구매에 영향을 미치는 변수들을 사용한다면 이는 더욱 높은 성능의 상품 추천이 가능할 것으로 기대된다. 향후에는 고객의 검색 키워드 등 더 많은 변수를 활용하여 더욱 개인화된 추천 시스템을 개발시킬 연구를 진행해보고자 한다.

References

- [1] S. Zhang, L. Yao, A. Sun, and Y. Tay, "Deep learning based recommender system: A survey and new perspectives," *ACM Computing Surveys*, Vol.1, No.1, pp.1-35, 2018.
- [2] R. Mu, X. Zeng, and L. Han, "A survey of recommender systems based on deep learning," *IEEE Access*, Vol.6, pp. 69009-69022, 2018.
- [3] A. Da'u and N. Salim, "Recommendation system based on deep learning methods: A systematic review and new directions," *Artificial Intelligence Review*, Vol.53, No.4, pp.2709-2847, 2020.
- [4] Y. Koren, R. Bell, and C. Volinsky, "Matrix Factorization Techniques for Recommender Systems," *IEEE Computer*, Vol. 42, No.8, pp.30-37, 2009.
- [5] S. Sedhain, A. K. Menon, S. Scanner, and L. Xie, "AutoRec: Autoencoders meet collaborative filtering," In *Proceedings of the 24th International Conference on World Wide Web*, May 2015.
- [6] H. J. Xue, X. Y. Dai, J. Zhang, S. Huang, and J. Chen, "Deep matrix factorization models for recommender systems," In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pp.3203-3209, 2017.
- [7] D. P. Kingma, "Auto-encoding variational bayes," In *Proceedings of the 2nd International Conference on Learning Representations (ICLR)*, Apr. 2014.
- [8] D. Liang, R. G. Krishnan, M. D. Hoffman, and T. Jebara, "Variational autoencoders for collaborative filtering," In *Proceedings of the 24th International Conference on World Wide Web*, Apr. 2018.
- [9] I. Shenbin, A. Alekseev, E. Tutubalina, V. Malykh, and S. I. Nikolenko, "RecVAE: A new variational autoencoder for Top-N recommendations with implicit feedback," In *Proceedings of the Web Search and Data Mining*, Feb. 2020.
- [10] X. He, L. Liao, H. Zhang, L. Nie, X. Hu, and T-S Chua, "Neural collaborativ filtering," In *Proceedings of the International Conference on World Wide Web*, Apr. 2017.
- [11] Y. Wang, L. Wang, Y. Li, D. He, T. Y. Liu, and W. Chen, "A Theoretical Analysis of NDCG Ranking Measures," In *Proceedings of 26th Annual Conference on Learning Theory*, 2013.
- [12] G. Kim, J. Kwak, D. Hong, and H. H. Kim, "A variational autoencoders based recommendation using the online behavior log," In *Proceedings of the Korean Institute of Information Scientists and Engineers*, pp.1147-1149, 2020.
- [13] Y. Jung and Y. Cho, "Topic modeling-based book recommendations considering online purchase behavior," *Knowledge Management Research*, Vol.18, No.4, pp.97-118, 2017.
- [14] H. Lee, S. Hong, J. Bang, and H. Kim, "A study on the music recommendation based on user playlist using data embedding," *The Journal of Korean Institute of Information Technology*, Vol.18, No.9, pp.27-34, 2020.
- [15] Y. Guan, Q. Wei, and G. Chen, "Deep learning based personalized recommendation with multi-view information integration," *Decision Support Systems*, Vol.118, pp.58-69, 2019.
- [16] X. Li and J. She, "Collaborative variational autoencoder for recommender systems," In *Proceedings of ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, Aug. 2017.
- [17] D. P. Kingma and M. Welling, "An introduction to variational autoencoders," *Foundations and Trends in Machine Learning*, Vol.12, No.4, pp.15-32, 2019.
- [18] M. G. Kim and K. Kim, "Recommender systems using SVD with social network information," *Journal of Intelligent Information Systems*, Vol.22, No.4, pp.1-18, 2016.
- [19] Y. Wang, L. Wang, Y. Li, D. He, W. Chen, and T. Y. Liu, "A theoretical analysis of normalized discounted cumulative gain (NDCG) ranking measures," In *Proceedings of the 26th Annual Conference on Learning Theory (COLT)*, 2013.
- [20] K. Jaavelin and J. Kekaelaenen, "Cumulated gain-based evaluation of IR techniques," *ACM Transactions on Information Systems*, Vol.20, No.4, pp.422-446, 2022.



홍 다 영

<https://orcid.org/0000-0001-5157-2488>

e-mail : dayoung0308@daum.net

2022년 동덕여자대학교 정보통계학과
(학사)

2022년~현 재 고려대학교

의학통계학협동과정 석사과정

관심분야 : Deep Learning, Biostatistics, Data Analysis



김 가 영

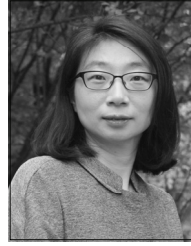
<https://orcid.org/0000-0001-8668-7963>

e-mail : gajil103@naver.com

2021년 동덕여자대학교 정보통계학과
(학사)

2021년~현 재 성균관대학교
인공지능학과 석사과정

관심분야 : Recommendation System, Machine Learning



김 현 희

<https://orcid.org/0000-0002-7507-8342>

e-mail : heekim@dongduk.ac.kr

1996년 이화여자대학교 컴퓨터학과(학사)

1998년 이화여자대학교 컴퓨터학과(석사)

2005년 이화여자대학교 컴퓨터공학과(박사)

2005년~2006년 LG전자디지털미디어
연구소 선임연구원

2006년~현 재 동덕여자대학교 정보통계학과 부교수

관심분야 : Machine Learning, Deep Learning, Big Data
Analysis