

Melon Playlist Continuation

CUAI 4기 메로나팀

권송아(소프트웨어학), 김소은(응용통계학), 김태운(소프트웨어학), 김효민(소프트웨어학), 최은서(소프트웨어학), 홍지호(기계공학)

2. 본 론

[요약]

이 프로젝트는 카카오 아레나의 Melon Playlist Continuation 참가를 가정하고 진행하였으며, 프로젝트의 목표는 음악 플레이리스트에 수록된 곡과 태그의 절반 또는 전부가 숨겨져 있을 때, 주어지지 않은 곡들과 태그를 예측하는 것이다. 숨겨진 곡과 태그를 예측할 수 있는 모델을 만든다면, 이를 주어진 플레이리스트에 대해 그 플레이리스트와 어울리는 곡을 추천해주는 용도로 사용할 수 있다. 멜론 플레이리스트 데이터 EDA를 통해 장르, 곡, 태그에 대한 데이터의 특징을 파악하였다. 이후 오토인코더를 사용한 협업 필터링을 기반으로 모델을 설계하여 평가 지표로서 nDCG와 Recall을 사용하였다. 또한 사전 훈련된 모델을 통해 코사인 유사도를 기반으로 플레이리스트 제목과 유사한 태그를 추천하는 시스템을 설계하였다.

1. 서 론

국내 음원 스트리밍 시장에서 독보적인 점유율을 유지하던 ‘멜론’은 ‘지니’, ‘플로’ 등 국내 음악 서비스는 물론이고 ‘애플 뮤직’, ‘스포티파이’, 유튜브에서 음악 관련 기능만 집중한 ‘유튜브 뮤직’이 주요 경쟁자로 부상함에 따라 그 지위가 흔들리고 있다.

이처럼 음원 시장에서의 각 플랫폼의 생존 경쟁도 치열하다. 사용자의 취향이나 분위기에 맞는 곡을 제안하는 추천시스템이 중요하게 부각되고 있다. 음원 스트리밍 플랫폼에는 수천만 개가 넘는 곡을 서비스 하는데, 이 수많은 곡들 중 자신의 취향에 맞는 음악을 사용자가 일일이 찾는 것은 시간이 많이 걸리는 작업이다. 따라서 사용자의 취향에 맞는 곡을 최대한 많이 효율적으로 탐색하는 것을 도와주는 시스템을 구축해야 한다.

이에 본 참가팀은 협업 필터링(Collaborative filtering) 기법을 적용하여 플레이리스트의 태그와 곡을 추천하는 모델을 설계하고자 한다. 예측 곡과 예측 태그를 nDCG의 가중평균값으로 점수를 매겨 모델의 성능을 점검하고자 한다. 감추어진 곡과 태그를 정확히 예측하는 모델을 완성하게 된다면, 음원 스트리밍 플랫폼이 사용자의 플레이리스트에 어울리는 곡을 효율적으로 추천하는 것에 도움을 줄 수 있을 것으로 예상된다.

1) EDA

장르 데이터인 genre_gn_all.json을 통해 총 254개의 장르코드가 존재함을 확인하였다. 곡 별 메타 데이터인 song_meta.json을 통해 대부분의 곡들은 한 개의 대분류 장르와 매핑되어 있고 2014~2019년도에 발매된 곡의 비중이 높았다. 플레이리스트의 태그는 1개부터 11개까지 존재하고, 노래는 1개부터 200개까지 존재한다.

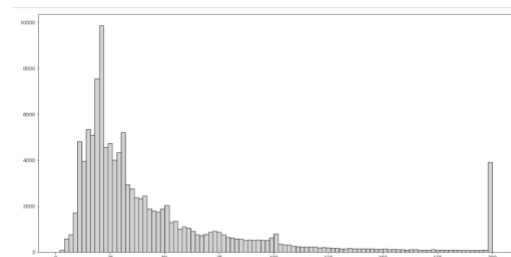


그림 1. 플레이리스트에 수록된 곡

	min	25%	50%	mean	75%	max	std
songs	1.0	19.0	30.0	45.94	54.0	200.0	43.95

표 1. 플레이리스트에 수록된 곡 통계

train.json을 통해 플레이리스트에는 평균 약 46개의 곡이 수록되어 있고 약 4.1개의 태그가 수록되어 있다. 수록곡의 약 51%, 태그의 약 40.2%는 두 개 이상의 플레이리스트에 중복 수록되어 있다. 플레이리스트 아티스트는 다르지만 같은 이름을 갖는 플레이리스트도 약 2% 존재했고 대부분 곡 수와 태그 수도 비슷한 편이었다. 기분전환, 계절, 드라이브 등의 태그가 매핑 기준 상위권에 있었으며 다른 태그와 조합이 없는 장르 관련 태그가 상위권을 차지하고 있었다.

2) 노래 추천

2-1) 데이터 세트

학습에 이용할 수 있는 전체 데이터 세트에는 총 115,071개의 플레이리스트와 707,989개의 노래가 존재했다. 이 중, 10개 미만의 노래가 수록된 플레이리스트는 학습에서 제외시켰다. 그 이후 검증 용도로 30,000개, 테스트 용도로 10,000개의 플레이리스트를 분리시켰다.

2-2) MultiDAE를 이용한 협업필터링 추천 시스템

노래 파트는 협업 필터링을 이용해 플레이리스트에

어울리는 새로운 노래를 예측하였다. 플레이리스트에 수록된 노래들이 비슷하게 수록된 다른 플레이리스트들을 찾아 새로운 노래를 추천하는 방식이다.

데이터의 차원의 수가 크기 때문에 데이터를 잠재공간으로 압축해 학습을 진행하는 오토인코더(AutoEncoder)를 사용해 협업 필터링을 구현했다.

오토인코더 모델에는 DAE(Denoising AutoEncoder)와 VAE(Variational AutoEncoder)가 있다. DAE는 기본 오토인코더 모델보다 Robust한 결과를 내기 위해 입력 데이터에 노이즈를 추가해 학습을 진행한다. 그리고 VAE는 은닉층이 입력 데이터와 최대한 같은 분포를 가지도록 학습한다. 베이지안 추론(Bayesian Inference)을 이용해 입력 데이터와 가까운 분포를 가지는 은닉층에서 원본 데이터와 가까운 값을 출력하는 것이 VAE의 목표이다.

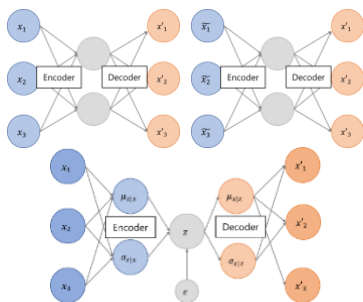


그림 3. AE, DAE, VAE

후보 모델로 MultiDAE(Multinomial DAE)와 MultiVAE(Multinomial VAE) [1]를 고려했다. 추천시스템에서의 ranking loss는 직접적으로 최적화되기 어렵고 보통 근사값을 이용한다. 참고한 논문에서 다항가능도 함수(Multinomial Likelihood)가 ranking loss를 감소시키는데 적합하다는 것과 implicit한 feedback에 적합하다는 것을 알 수 있었다.

두 모델을 MovieLens-20M 데이터 세트를 이용한 baseline 코드를 돌려보았을 때, MultiVAE가 근소하게 더 높은 성능을 보였지만 연산량이 더 많았다. 우리가 사용한 데이터 세트는 이보다 더 큰 데이터 세트였기 때문에 성능을 조금 포기하는 대신 MultiDAE를 택했다.

입력값으로는 플레이리스트에 들어있는 노래를 1로, 들어있지 않은 노래를 0으로 표현한 행렬을 사용한다. MultiDAE의 layer들의 차원은 [노래 수, 200, 노래 수]이다. 참고한 논문에서 가장 높은 성능을 보였던 구조를 택했다.

학습 과정에서는 인코더에 해당하는 가중치와 디코더에 해당하는 가중치를 학습한다. 각 layer를 통과할 때마다 입력 값과 가중치를 곱한 값에 bias가 더해진다. 활성화 함수로는 tanh를 사용했고 12 정규화를 적용했다. 또, 가중치는 xavier 초기화를 사용했고, bias는 잘린 정규분포에서 샘플링을 해 초기화를 했다.

DAE의 objective 함수는 아래와 같다. 인코더 네트워크의 결과 값이 주어졌을 때, 입력 값이 나올 확률에 log를 취한 것이다.

$$\mathcal{L}_u(\theta, \phi) = \log p_{\theta}(\mathbf{x}_u | g_{\phi}(\mathbf{z}_u))$$

학습을 최적화하기 위해 adam optimizer를 사용했고, Learning rate는 $1e-3$, regularization rate는 $1e-2$, epoch 6회로 학습을 진행했다.

2-3) 결과

X_t = np.transpose(np.nonzero(X)) song_meta.loc[X_t[:,1]]['song_name']		[126] song_meta.loc[ids_topk_part]['song_name']	
17889	사랑, 참...	0	Feelings
25280	그때로	471988	두근두근
46162	복국성	471989	춤추
51612	오늘은 가자	471990	Last Dance
88503	윤금성	471991	Memories
97412	상록수		
109500	오빠는		
117184	그리움	259897	떠나갈래
118150	사랑할 것 같아서	259898	Gang Bang
121846	올지마	259899	On My Way
138319	옛날여자	259901	춤추는데
142777	보날	707988	Queen 맘 알로디
171689	가난한 사랑		

그림 4. 노래 학습 데이터와 예측값

이별을 담은 플레이리스트를 모델에 넣었을 때, 이별에 관련한 '떠나갈래', '춤추는데'라는 노래가 담긴 노래를 예측한 것을 볼 수 있었다.

```
Test NDCG@3000=0.00521 (0.00016)
Test Recall@3000=0.01782 (0.00038)

Test NDCG@1000=0.00365 (0.00015)
Test Recall@1000=0.00999 (0.00027)

Test NDCG@500=0.00285 (0.00014)
Test Recall@500=0.00636 (0.00022)

Test NDCG@100=0.00192 (0.00014)
Test Recall@100=0.00275 (0.00016)
```

그림 5. 평가지표 nDCG, Recall

평가지표로는 관련성이 높은 아이템을 순서대로 예측했는지 평가하는 nDCG와, 실제 추천되어야 하는 것이 얼마나 포함되어 있는지 나타내는 지표인 Recall을 사용했다. @뒤의 값은 평가의 대상이 상위 top-N개를 뜻한다.

이 모델은 보통 추천시스템에서 수많은 데이터 중 후보군을 추려내는데 사용하는 모델이다. 따라서 절대적인 수치가 작아 보이지만 의미가 없는 값이 아니다. Recall 3000 값을 보면 1000개 중 17개가 바르게 추천된 것을 볼 수 있는데, 3000개 중에서는 51개가 바르게 추천이 되었음을 의미한다. 큰 데이터 세트에서 후보군을 만든 다음 랭킹 모델에 넣어서 바르게 추천되어야 하는 51개의 순서가 앞으로 오도록 학습을 진행하면 더 나은 성능의 추천시스템을 만들 수 있을 것이다.

3) 태그 추천

3-1) 데이터 세트

데이터 전처리는 비슷한 분야의 태그를 합치는 방향으로 진행되었다. 먼저 태그 데이터에서 한글이 아닌 영어, 숫자로 이루어진 데이터를 삭제하였고 이를 바탕으로 태그간 유사도 분석을 진행하였다. 유사도 분석은 FastText[2]에서 사전 훈련된 모델을 통해 이루어졌다. 비슷한 태그끼리 매칭이 이루어진 후 가장 많은 태그와 매칭이 된 태그를 상위 태그로 선정하였고 상위 태그와 비슷한 태그를 같은 라벨로 라벨링하였다. 그 이후 결측치를 제거 하였고 태그, 플레이리스트, 노래로 이루어진 새로운 데이터프레임으로 바꾸는 작업을 거쳤다.

데이터 전처리를 거쳐 만들어진 111,133개의 플레이리스트와 16,451개의 태그 중 20,000개의 플레이리스트는 학습에 10,000가 테스트에 이용되었다.

3-2) 플레이리스트 제목을 이용한 태그 추천시스템

태그의 경우 텍스트 기반제목 벡터를 이용한 추천

