

Improving LightGCN based recommendation system with NLP metadata embeddings

Pratik Acharya
achary17@purdue.edu
collaboration with Vidhi Jain
Thanks to Prof. Bruno Ribeiro

Abstract

1 In this study we explored various methodologies relating to recommendation
2 systems using Deep Learning algorithms. LightGCN is one of the state of the art
3 matching systems for recommendation based on collaborative filtering. We apply
4 this to a book recommendation problem which has been growing in popularity
5 recently. LightGCN only uses the user book interactions and not any user or book
6 metadata. Along the way we devise a new technique to initialize embeddings of
7 LightGCN nodes using sentence embeddings from an NLP transformer model
8 called MiniLM. This also means we have a content based recommendation along
9 with collaborative filtering. This improves the performance significantly.

10 1 Motivation

11 Shopping, social media, streaming websites and apps actively invest in recommendation systems to
12 provide curated content for each user. This state-of-the-art methodology is used less in the publishing
13 category because reading as a habit had been declining over the past few decades. Over the past 2
14 years (the pandemic years), we have seen a sudden rise in the number of books people read. We offer
15 2 pieces of evidence to support our claim: a. The annual revenue of the eBook industry has risen to
16 more than a billion dollars in the US alone. b. The subreddit r/suggestmeabook, where users provide
17 book recommendations to each other, has gained 1 million followers since 2020.

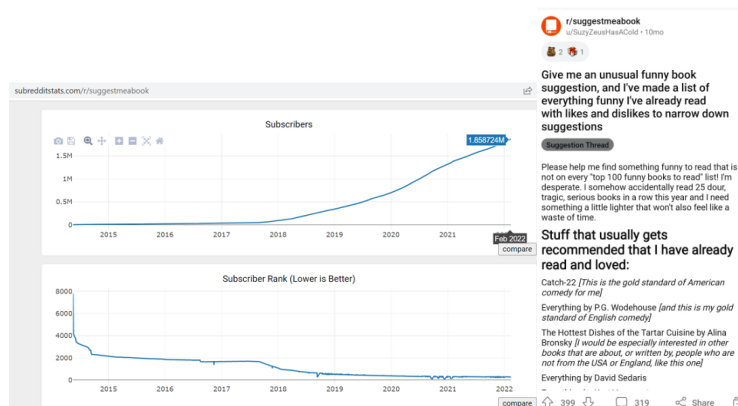


Figure 1: r/suggestmeabook

2 Dataset

We chose the Goodreads books dataset provided by University of California, San Diego. The dataset can be found here. This dataset was collected in the year 2017 and consists of 3 files – 1. Books: Contains information on each book in JSON format. 2. Shelves: Captures the user-book interactions in a CSV format. 3. Reviews: Captures any text-based reviews shared by users. Some specific statistics of the dataset: This dataset contains data pertaining to 2.3M books including 1.5M works, 400k series and 830k authors. It also has dataset pertaining to 870k users and 230M user-book interactions including 112M reading statuses and 104M ratings. An example entry from this dataset is shown in the figure below.

<pre>{ "isbn": "", "text_reviews_count": "7", "series": ["189911"], "country_code": "US", "language_code": "eng", "# top user-generated shelves for a book used to define genres by goodreads": { "popular_shelves": [{"count": "58", "name": "to-read"}, {"count": "15", "name": "fantasy"}, {"count": "6", "name": "fiction"}, {"count": "5", "name": "owned"}, ...] }, "asin": "B00071IKUY", "is_ebook": "false", "average_rating": "4.03", "kindle_asin": "", "# a list of books that users who like the current book also like": { "similar_books": [{"id": "19997", "id2": "828466", "id3": "1569323", "id4": "425389", "id5": "1176674", "id6": "262748", "id7": "3743837", "id8": "888461", "id9": "2292726", "id10": "1883810", "id11": "1888197", "id12": "625158", "id13": "1988846", "id14": "399178", "id15": "2620131", "id16": "383106", "id17": "1597281"}, {"id": "7327624", "id2": "8948723", "id3": "148", "id4": "148", "id5": "148", "id6": "148", "id7": "148", "id8": "148", "id9": "148", "id10": "148", "id11": "148", "id12": "148", "id13": "148", "id14": "148", "id15": "148", "id16": "148", "id17": "148"}], "description": "Omnibus book club edition containing the Ladies of Madrigyn and the Witches of Wenshar.", "format": "Hardcover", "link": "https://www.goodreads.com/book/show/7327624-the-unschooled-wizard", "authors": [{"author_id": "10333", "role": "I"}], "publisher": "Nelson Doubleday, Inc.", "num_pages": "680", "publication_day": "", "isbn13": "", "publication_month": "", "edition_information": "Book Club Edition", "publication_year": "1987", "url": "https://www.goodreads.com/book/show/7327624-the-unschooled-wizard", "image_url": "https://images.gr-assets.com/books/1304100136m/7327624.jpg", "book_id": "7327624", "ratings_count": "148", "work_id": "8948723", "title": "The Unschooled Wizard (Sun Wolf and Starhawk, #1-2)", "title_without_series": "The Unschooled Wizard (Sun Wolf and Starhawk, #1-2)" }</pre>	<table><tr><th>user_id</th><th>book_id</th><th>is_read</th><th>rating</th><th>is_reviewed</th></tr><tr><td>678442</td><td>7381</td><td>1</td><td>5</td><td>0</td></tr><tr><td>135202</td><td>6166</td><td>1</td><td>3</td><td>0</td></tr><tr><td>126691</td><td>262380</td><td>1</td><td>5</td><td>0</td></tr><tr><td>3224</td><td>237810</td><td>1</td><td>4</td><td>0</td></tr><tr><td>274050</td><td>6955</td><td>1</td><td>5</td><td>0</td></tr></table>	user_id	book_id	is_read	rating	is_reviewed	678442	7381	1	5	0	135202	6166	1	3	0	126691	262380	1	5	0	3224	237810	1	4	0	274050	6955	1	5	0
user_id	book_id	is_read	rating	is_reviewed																											
678442	7381	1	5	0																											
135202	6166	1	3	0																											
126691	262380	1	5	0																											
3224	237810	1	4	0																											
274050	6955	1	5	0																											

Figure 2: Goodreads book graph by UCSD

3 Implementation

3.1 Training data generation

To get genuine items, we remove NaNs, winsorize top and bottom 5% of items (books and users) with unusual values, mean, std, count of different attributes, and analysis is done only on English books. Further, remove books with <10 reviews, description length <50 characters, <100 ratings. After this we get around 10 million user-book interactions left, then we draw a random sample to get 2000 unique users, 4500 unique books to add to our graph data structure with 500,000 edges. To make the graph we used a library called snap in python <https://snap.stanford.edu/snappy/doc/reference/graphs.html>. This was useful to return a k-connected-graph for pytorch geometric. To do the train test split, we use a function called RandomLinkSplit in the pytorch geometric library, which does an inductive split as shown in the figure, the nodes are same in train, val and test splits. Only the edges are distributed.

- Suppose we have a dataset of 3 graphs. Each inductive split will contain an independent graph
- In **train** or **val** or **test** set, each graph will have 2 types of edges: message edges + supervision edges
 - Supervision edges are not the input to GNN

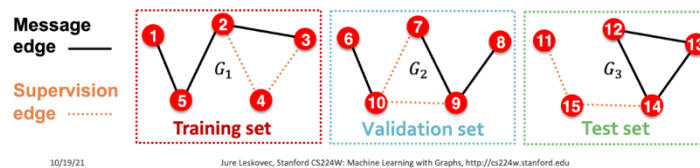


Figure 3: inductive train split

3.2 LightGCN

LightGCN is a type of GCN which is not as computationally expensive. A comparison of the propagation rules is given in the figure below. In LightGCN, we remove activations and other weight matrices except layer 0 embeddings which are learnt, and simplify the formula for aggregation. This means the model is much faster than usual GNN and will suit our large samples from the dataset. For each node, we take the average of embeddings of all the layers. Number of layers here represent maximum degrees of separation between users/books that are averaged to generate the final embeddings. We didn't see a significant change in the performance on adding more layers so we chose to keep 3 layers for our model which aids computation. It is a very simple model where we just make lightgcn layers and take mean after propagation.

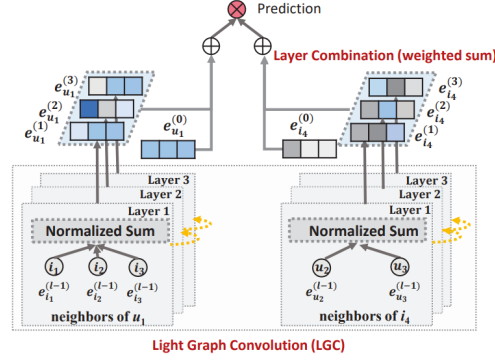


Figure 4: LightGCN Model Architecture from the paper

$$e_u^{(k+1)} = \sigma \left(W_1 e_u^{(k)} + \sum_{i \in N_u} \frac{1}{\sqrt{|N_u|} \sqrt{|N_i|}} (W_1 e_i^{(k)} + W_2 (e_i^{(k)} \odot e_u^{(k)})) \right)$$

$$e_i^{(k+1)} = \sigma \left(W_1 e_i^{(k)} + \sum_{u \in N_i} \frac{1}{\sqrt{|N_u|} \sqrt{|N_i|}} (W_1 e_u^{(k)} + W_2 (e_u^{(k)} \odot e_i^{(k)})) \right)$$

N_u : the set of all neighbors of user u (items liked by u)
 N_i : the set of all neighbors of item i (users who liked i)
 $e_u^{(k)}$: k -th layer user embedding
 $e_i^{(k)}$: k -th layer item embedding
 W_1, W_2 : trainable weight matrices

Figure 5: GCN propagation rule

$$e_u^{(k+1)} = \sum_{i \in N_u} \frac{1}{\sqrt{|N_u|} \sqrt{|N_i|}} e_i^{(k)} \quad e_i^{(k+1)} = \sum_{u \in N_i} \frac{1}{\sqrt{|N_u|} \sqrt{|N_i|}} e_u^{(k)}$$

Figure 6: LightGCN propagation rule

The loss function used here is Bayesian Personalized Loss. It is commonly used for recommendation problems. Note: here we are also sampling negative edges randomly, since in our training edges exist only between users and items that are rated positively. Negative edges are sampled such that they equal positive edges. BPR loss in PyG is shown in figure below, we use similar score but without regularization.

$$L_{BPR} = - \sum_{u=1}^M \sum_{i \in N_u} \sum_{j \in N_u} \ln \sigma(\hat{y}_{ui} - \hat{y}_{uj}) + \lambda \|E^{(0)}\|^2$$

\hat{y}_{ui} : predicted score of a positive sample
 \hat{y}_{uj} : predicted score of a negative sample
 λ : hyperparameter which controls the L2 regularization strength

Figure 7: BPR Loss

```
reg_loss = lambda_val * (users_emb_0_norm(2).pow(2) +
                        pos_items_emb_0_norm(2).pow(2) +
                        neg_items_emb_0_norm(2).pow(2)) * L2_loss

pos_scores = torch.nn.functional.matmul(users_emb_final, pos_items_emb_final)
neg_scores = torch.nn.functional.matmul(users_emb_final, neg_items_emb_final)
neg_scores = torch.nn.functional.matmul(users_emb_final, neg_items_emb_final)

loss = -torch.mean(torch.nn.functional.softplus(pos_scores - neg_scores)) + reg_loss
```

Figure 8: BPR Loss in PyG

3.3 MiniLMv2

MiniLM is trained to mimic the self attention module of a transformer model. We have chosen this model specifically over other transformers because our dataset has 2 million books and to generate embeddings for all descriptions, we need a fast model. As we can see in the benchmark comparison, MiniLM is near the top and is one of the best in terms of accuracy.

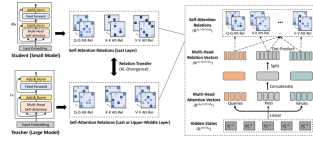


Figure 9: MiniLM Model Architecture from the paper

Model Name	Performance Sentence Embeddings (14 Datasets)	Performance Semantic Search (14 Datasets)	Avg. Performance	Speed	Model Size
paraphrase-MiniLM-L3-v2	62.25	39.19	50.72	19000	61 MB
all-MiniLM-L6-v2	68.06	49.54	58.80	14200	80 MB
multi-qa-MiniLM-L6-cos-v1	64.31	51.83	58.08	14200	80 MB
paraphrase-multilingual-MiniLM-L12-v2	64.25	39.19	51.72	7500	420 MB
all-MiniLM-L12-v2	68.70	50.82	59.76	7500	120 MB
paraphrase-distil-multilingual-v1	64.46	40.04	52.25	5000	43 MB
distil-base-multilingual-cased-v1	61.30	29.87	45.59	4000	480 MB
distil-base-multilingual-cased-v2	60.18	27.35	43.77	4000	480 MB
all-distilbert-v1	68.73	50.94	59.84	4000	200 MB

Figure 10: MiniLM benchmark with other models via huggingface

59 To check whether the embeddings we get are genuine, we cluseter this using kmeans and DBSCAN.
 60 Each of the clusters gave a detailed genre. Examples are shown in figure below. This verifies that the
 61 embeddings are making sense.

'The Christmas wedding', 'Merry Ex-Hus', 'A Christmas Present',
 'Silver Bells', 'Snow Globe',
 'Sleigh Ride (Homesick Christmas, #2)',
 'Gus Was A Friendly Ghost', 'The Fifth Day of Christmas',
 'Big Girls Do It on Christmas (Big Girls Do It, #6)',
 '12 Terrors of Christmas', 'The Not Very Merry Hoot-Noot Fish',
 'A Scandal in Battersea (Elemental Pastors, #12)',
 'Speters Closes', 'The Holiday Bride (Holiday Brides, #2)',
 'Silent Night 3 (Fear Street Supercollider, #13)',
 'Christmas with an Angel (Honey Honey, #1.5)',
 'Christmas in Harmony (Harmony, #2.5)',
 'Choosing Christmas (Piper Anderson, #2.5)',
 'Christmas in Paradise (T Jensen Paradise Lake Mysteries #4)',
 'The Christmas Bouquet / Bayside Retreat (Chesapeake Promise, #11 & 12.5)',
 'Twelve Days of Christmas',
 'Christmas at Leo's (Memoirs of a Houseboy, #5)',
 'Merry's Christmas: a love story (Redding Romance Series)',
 'A Regency Christmas: Scarlet Ribbons / Christmas Promise / A Little Christmas',
 'Silent Belle (Celebration Bay, #2)',
 'Love Finds You in Holly Beach, South Carolina',
 'Sprinkles on Top (Sugar Springs, #1)',
 'Ruth's first Christmas Tree (Ruth Galloway, #4.5)',
 'I Heart Christmas', 'Waiting in Vain', dtype='<U486>'

Figure 11: Books with christmas theme

[['Honey So Sweet, Vol. 2', 'White Crane (Samurai Kids, #1)',
 'Tsubasa: REServoir CHRONICLE, Vol. 13', 'Samsir',
 'Maid-samurai Vol. 01 (Maid-samurai, #1)', 'La Corda d'Oro, Volume 1',
 'Naruto, Vol. 66: The New Three (Naruto, #66)',
 'Strobe Edge, Vol. 5 (Strobe Edge, #5)', 'Skip Beat!, Vol. 12',
 'Close The Last Door, Volume 01',
 'Twittering Birds Never Fly, Vol.1',
 'The Wallflower, Vol. 3 (The Wallflower, #3)',
 'キズナとちどち (Kisuu Yorio Hayaku #1) (Faster than a Kiss #4)',
 'Your Lie in April, Vol. 4', 'Noragami: Stray God, Vol. 12',
 'The Ring of Sky (Young Samurai, #3)',
 'The Fox Woman (Love/Har/Death, #1)',
 'Punch!, Vol. 3 (Punch!, #2)', 'Gochikid, Volume 04',
 'Bleach, Volume 87', 'The Decay of the Angel',
 'Dragon Ball Z, Vol. 21: Tournament of the Heavens (Dragon Ball Z, #21)',
 'Get Well, Kikakuhi (Polkman Reader #9)',
 'What Did You Eat Yesterday, Volume 4',
 'Your Lie in April, Vol. 4',
 'A Silent Voice, Vol. 7 (A Silent Voice, #7)',
 'Night of Cake & Puppets (Daughter of Smoke & Bone, #2.5)',
 'Sand Chronicles, Vol. 9', 'A Pale View of Hills',
 'Dorazaru! Saika Arc, Vol. 1', dtype='<U486>']

Figure 12: Japanese Fiction

62 Using the book descriptions, we get sentence embeddings of size 384 for each book. Now we use
 63 dimensionality reduction to get 10 dimensional embeddings for books. To make embeddings for
 64 users, we take a rating weighted average of book embeddings, of books read by that user. Now we
 65 standardize these embeddings and this gives initial user and book embeddings for all nodes. The
 66 process is shown below.

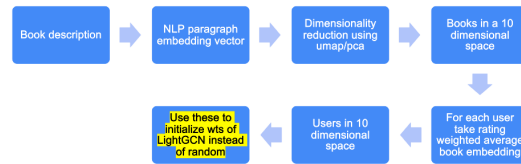


Figure 13: Process of making the embeddings

67 4 Results

68 We find that the model trains quickly 30 minutes for 30 epochs considering that the dataset is large.
 69 The validation recall quickly reaches the asymptotic value. With increase in embedding size of each
 70 node, the validation recall increases slightly. So we fix a size of 10 here in the 3 layers deep model.

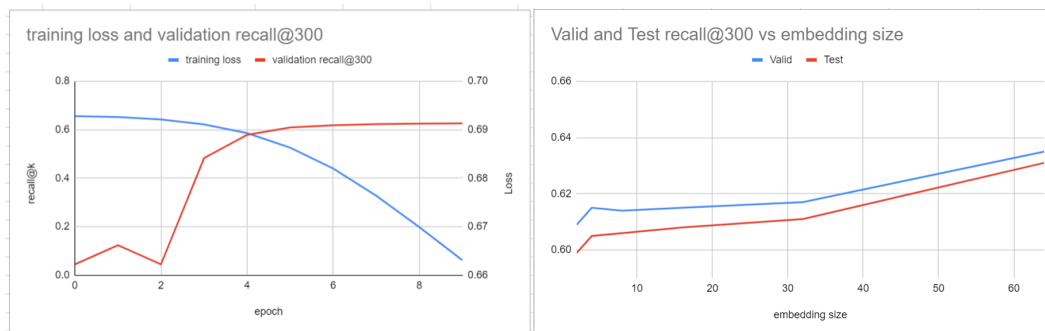


Figure 14: Training LightGCN for goodreads data

Now after initializing node embeddings with NLP embeddings from MiniLM(dimension reduced to 10), we can see 2 advantages immediately. The convergence is very quick, almost in 1 epoch vs 10-30 for normal sampled embeddings.

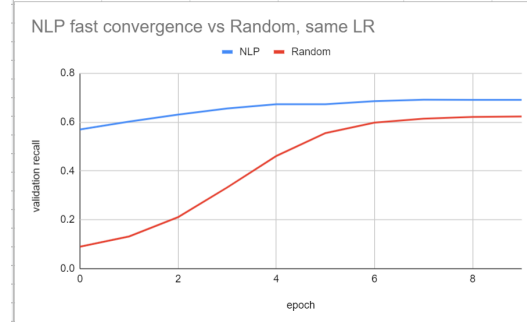


Figure 15: Fast convergence for NLP embeddings

Another very big advantage is the performance improvement in the recall score after convergence. We can see the results of 8 samples in the table below and the corresponding paired t tests. With a very low p value, the recall@50 has increased as much as 15%, which is huge considering that LightGCN is already state of the art. The performance has increased significantly for all recall scores.



Figure 16: Comparison of both embeddings performance

5 Conclusion and Discussion

Using LightGCN+NLP embeddings initialization we were able to recommend books with a recall@300 score of 66%, which means that 66% of the relevant books in test data are recommended in top 7% recommendations. This is a very good score. Using these embeddings as initial values in LightGCN model, we were able to improve the recall@50 performance by as much as 15%(statistically significant for 8 samples), which is a lot considering that this model is sota for recommendation tasks.

6 Comments from the presentation

Yes this is a very simple permutation invariant model. I haven't made the LightGCN architecture or trained the MiniLM model. It is pretrained, downloaded from huggingface. My idea is used to improve the performance of sota LightGCN for book recommendation purpose. For this i try a new idea to generate embeddings from the book description which are later used to initialize embeddings for the user and book nodes of LightGCN.

92 7 Future ideas

- 93 • Test this technique on other graph based recommendation networks
- 94 • Does the performance of the transformer change recall score
- 95 • Try different dimension reduction techniques
- 96 • Add metadata other than book description by adding sentences like this book is long/short
- 97 etc. based on number of pages, also add written by this author
- 98 • Include authors in the graph as nodes

99 References

- 100 [1] Wang, W., Bao, H., Huang, S., Dong, L., Wei, F. (n.d.). MINILMv2: Multi-Head Self-Attention
101 Relation Distillation for Compressing Pretrained Transformers ArXiv:2012.15828v2 [cs.CL] 27 Jun 2021
102 <https://arxiv.org/pdf/2012.15828.pdf>
103
- 104 [2] He, X., Deng, K., Wang, X., Li, Y., Zhang, Y., Wang, M. (2020, July 7). LIGHTGCN: Simplifying and
105 powering graph convolution network for recommendation. <https://arxiv.org/abs/2002.02126>
106
- 107 [3] Wang, X., He, X., Wang, M., Fung, F., Chua, T.-S. (n.d.). Neural graph collaborative filtering -
108 <https://arxiv.org/pdf/1905.08108.pdf>
109
- 110 [4] Gao, C., Zheng, Y., Li, N., Li, Y., Qin, Y., Piao, J., Quan, Y., Chang, J., Jin, D., He, X., Li, Y. (n.d.). Chen
111 Gao*, Yinfeng Li, Yingrong Qin ... - Graph Neural Networks for Recommender Systems: Challenges, Methods,
112 and Directions <https://arxiv.org/pdf/2109.12843.pdf>