

Wide & Deep Book Recommender System

Author: Paul McSlarrow

Original paper: <https://arxiv.org/abs/1606.07792>.

This is not an implementation by the authors of the paper

Imports and Constants

```

1 import os
2 import random
3 import numpy as np
4 import pandas as pd
5 import torch
6 import torch.nn as nn
7 import torch.optim as optim
8 from torch.utils.data import Dataset, DataLoader, random_split
9 from torch.utils.tensorboard import SummaryWriter
10 from sklearn.model_selection import train_test_split
11 from sklearn.preprocessing import StandardScaler, OneHotEncoder
12 from collections import Counter
13 from google.colab import drive
14
15 # torch.manual_seed(1)

1 RATINGS_PATH = '/content/drive/MyDrive/ratings.csv'
2 BOOKS_PATH = '/content/drive/MyDrive/books.csv'
3 BOOK_TAGS_PATH = '/content/drive/MyDrive/book_tags.csv'
4 TAGS_PATH = '/content/drive/MyDrive/tags.csv'
5 TRAIN_SIZE = 0.8

```

Dataset creation

This section joins multiple tables to produce a unified dataset. Some tables contain genre information, while others provide the count of each genre. The join creates a single table that attaches the **three most popular genres** for each `goodreads_book_id`.

While I am also using a different dataset from the paper, to try my best to keep consistency between the two, I will turn this into a binary classification problem rather than a 5-class classification problem, and binarize the rating label to be 0 or 1. If the rating was greater than or equal to 4, I'll say that the user liked it, otherwise, the user didn't particularly love it and will be set to 0. This will help me utilize binary cross entropy loss as the paper did, which also allows me to use the Sigmoid activation function to enforce confidence in predictions rather than turning it into a sort of regression problem.

Information

- **Shape:** (5,976,479, 25)
- **Columns:**

```
['user_id', 'book_id', 'rating', 'goodreads_book_id', 'best_book_id', 'work_id', 'books_count', 'isbn', 'isbn13', 'authors', 'original_publication_year', 'original_title', 'title', 'language_code', 'average_rating', 'ratings_count', 'work_ratings_count', 'work_text_reviews_count', 'ratings_1', 'ratings_2', 'ratings_3', 'ratings_4', 'ratings_5', 'image_url', 'small_image_url']
```

Example Row (via `__getitem__`)

```
{
  'user_id': np.int64(34197),
  'book_id': np.int64(155),
  'rating': np.int64(5),
  'categorical_genre': 'fantasy',
  'genre_idxs': 11305,
  'categorical_book_title': 'The Two Towers (The Lord of the Rings, #2)',
  'title_idx': 8779
}
```

```

1 drive.mount('/content/drive/')
2
3 #
4 # PyTorch Dataset
5 #
6
7 class BookRecommenderDataset(Dataset):
8     def __init__(self, ratings_path, books_path, book_tags_path, tags_path):
9         #
10        # Raw datasets
11        #
12        self.reviews_df = pd.read_csv(ratings_path)
13        self.books_df = pd.read_csv(books_path)
14        self.book_tags_df = pd.read_csv(book_tags_path)
15        self.tags_df = pd.read_csv(tags_path)
16
17        #
18        # Joining
19        #
20        reviews_books_df = pd.merge(self.reviews_df, self.books_df, on="book_id").reset_index(drop=True)
21        tag_joined = pd.merge(self.book_tags_df, self.tags_df, on='tag_id')
22        top_tags = tag_joined.sort_values(['goodreads_book_id', 'count'], ascending=[True, False])
23        genre_per_goodreads_book_id = top_tags.groupby('goodreads_book_id')['tag_name'].apply(self.select_genres).reset_index()
24
25        #
26        # Master dataset
27        #
28        self.data = pd.merge(reviews_books_df, genre_per_goodreads_book_id, on='goodreads_book_id').iloc[:, [0, 1, 2, 25, 1
29
30        #
31        # Vocabularies
32        #
33        self.user_vocab = self.generate_vocab(self.data, 'user_id')
34        self.book_vocab = self.generate_vocab(self.books_df, 'book_id')
35        self.genre_vocab = self.generate_vocab(self.tags_df, 'tag_name')
36        self.cross_vocab = self.generate_cross_vocab()
37
38    def __len__(self):
39        return len(self.data)
40
41    def __getitem__(self, idx):
42        #
43        # IDs
44        #
45        user_id = self.data.iloc[idx]['user_id']
46        user_id_one_hot = self.one_hot(user_id, self.user_vocab)
47
48        book_id = self.data.iloc[idx]['book_id']
49        book_id_one_hot = self.one_hot(book_id, self.book_vocab)
50
51        #
52        # Genre
53        #
54        categorical_genre = self.data.iloc[idx]['tag_name']
55        genre_idx = self.genre_vocab.get(categorical_genre)
56
57        # Cross-features (user_id, book_id) combined
58        # cross_one_hot = torch.zeros(len(self.cross_vocab))
59        # cross_index = self.cross_vocab.get((user_id, book_id))
60        # if cross_index is not None:
61        #     cross_one_hot[cross_index] = 1.0
62
63        #
64        # Multi-hot vector for wide
65        #
66        wide_input = torch.cat([
67            user_id_one_hot,
68            book_id_one_hot,
69            # cross_one_hot
70        ], dim=-1)
71
72        #
73        # Label
74        #
75        rating = self.data.iloc[idx]['rating']
76        rating = torch.tensor(1.0 if rating >= 4 else 0.0)

```

```

77     return {
78         'user_id': user_id,
79         'book_id': book_id,
80         'genre_idx': genre_idx,
81         'rating' : rating,
82         'wide_input': wide_input,
83     }
85
86     def __str__(self):
87         print("Retrieving a single random instance.", end='\n\n')
88         return str(self.__getitem__(735)) # favorite book <3
89
90     def one_hot(self, id, vocab):
91         a = torch.zeros(len(vocab))
92         a[vocab[id]] = 1.0
93         return a
94
95     def head(self, n=5):
96         return self.data.head(n)
97
98     def generate_vocab(self, df, column_name):
99         unique = set(df[column_name])
100        vocab = {key: idx for idx, key in enumerate(sorted(unique))}
101        vocab["<UNK>"] = len(vocab)
102        return vocab
103
104    def generate_cross_vocab(self):
105        pairs = list(zip(self.data['user_id'], self.data['book_id']))
106
107        unique_pairs = sorted(set(pairs))
108        cross_vocab = {pair: idx for idx, pair in enumerate(unique_pairs)}
109        return cross_vocab
110
111    def select_genres(self, tags, n=1):
112        l = tags[:5].tolist()
113
114        for tag in tags:
115            if tag not in {'to-read', 'favorites', 'currently-reading', 'books-i-own', 'book-club'}:
116                return tag
117
118 #
119 # Creating the dataset and train/test split
120 #
121 dataset = BookRecommenderDataset(RATINGS_PATH, BOOKS_PATH, BOOK_TAGS_PATH, TAGS_PATH)
122 train_size = int(TRAIN_SIZE * len(dataset))
123 test_size = len(dataset) - train_size
124
125 #
126 # Custom collate fn to design batches
127 #
128 def custom_collate_fn(batch):
129     batch_dict = {}
130     batch_dict['user_id'] = torch.tensor([item['user_id'] for item in batch], dtype=torch.int64)
131     batch_dict['book_id'] = torch.tensor([item['book_id'] for item in batch], dtype=torch.int64)
132     batch_dict['genre_idx'] = torch.tensor([item['genre_idx'] for item in batch], dtype=torch.int64)
133     batch_dict['rating'] = torch.tensor([item['rating'] for item in batch], dtype=torch.float64)
134     wide_inputs = []
135     for item in batch:
136         wide_input = item['wide_input']
137         if isinstance(wide_input, torch.Tensor):
138             wide_inputs.append(wide_input)
139         else:
140             wide_inputs.append(torch.tensor(wide_input, dtype=torch.float32))
141     batch_dict['wide_input'] = torch.stack(wide_inputs)
142     return batch_dict
143
144 train_data, test_data = random_split(dataset, [train_size, test_size])
145 train_dataloader = DataLoader(train_data, batch_size=64, shuffle=True, collate_fn=custom_collate_fn)
146 test_dataloader = DataLoader(test_data, batch_size=64, shuffle=True, collate_fn=custom_collate_fn)

```

→ Drive already mounted at /content/drive/; to attempt to forcibly remount, call drive.mount("/content/drive/", force_remount=

```

1 print(dataset)
2 item = dataset.__getitem__(735)
3

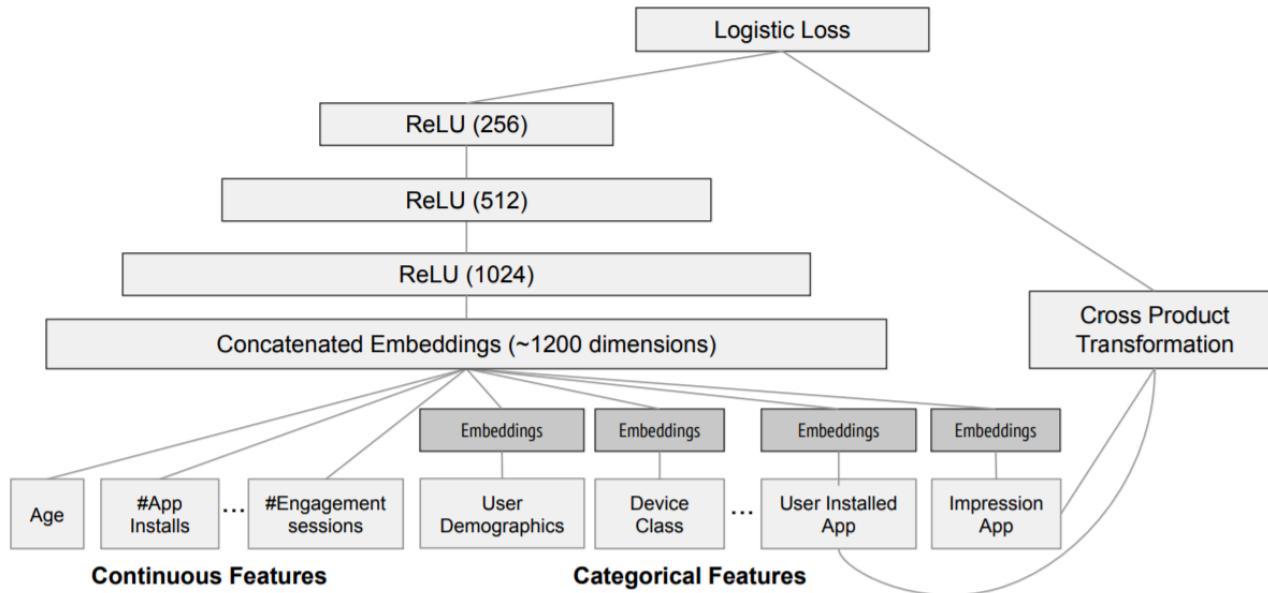
→ Retrieving a single random instance.

{'user_id': np.int64(65), 'book_id': np.int64(155), 'genre_idx': 11305, 'rating': tensor(1.), 'wide_input': tensor([0., 0.,

```

Model Building

As a summary for those who have not read the paper, the **Wide** component uses a linear transformation using cross-products for better **memorization** of co-occurrence between feature items and target variables. Whereas the **Deep** component uses a deep neural network to learn deep connections between features to hopefully become more **generalizable**, that is, it makes better recommendations for unseen examples. The deep component also utilizes embedding layers for categorical features, to learn dense representations of the otherwise sparse inputs.



```

1 USER_VOCAB_SIZE = len(dataset.user_vocab)
2 BOOK_VOCAB_SIZE = len(dataset.book_vocab)
3 GENRE_VOCAB_SIZE = len(dataset.genre_vocab)
4 CROSS_VOCAB_SIZE = len(dataset.cross_vocab)
5
6 USER_EMBEDDING_DIM = 32
7 BOOK_EMBEDDING_DIM = 32
8 GENRE_EMBEDDING_DIM = 20
9
10 class WideAndDeepModel(nn.Module):
11     def __init__(self):
12         super().__init__()
13
14         #
15         # Embedding layers
16         #
17         self.user_embedding = nn.Embedding(USER_VOCAB_SIZE, USER_EMBEDDING_DIM)
18         self.book_embedding = nn.Embedding(BOOK_VOCAB_SIZE, BOOK_EMBEDDING_DIM)
19         self.genre_embedding = nn.Embedding(GENRE_VOCAB_SIZE, GENRE_EMBEDDING_DIM)
20
21         #
22         # Learned parameters
23         #
24         self.deep_weight = nn.Parameter(torch.tensor(1.0))
25         self.wide_weight = nn.Parameter(torch.tensor(1.0))
26
27         #
28         # Deep feed-forward layers
29         #
30         self.ffnn = nn.Sequential(

```

```

1      nn.Linear(USER_EMBEDDING_DIM + BOOK_EMBEDDING_DIM + GENRE_EMBEDDING_DIM, 1024),
2      nn.ReLU(),
3      nn.Linear(1024, 512),
4      nn.ReLU(),
5      nn.Linear(512, 256),
6      nn.ReLU(),
7      nn.Linear(256, 1)
8  )
9
10 #
11 # Wide Layers
12 #
13 self.wide = nn.Linear(USER_VOCAB_SIZE + BOOK_VOCAB_SIZE, 1)
14
15 def forward(self, data):
16     #
17     # DEEP
18     #
19     embedded_user_ids = self.user_embedding(data['user_id']) # [B, E]
20     embedded_book_ids = self.book_embedding(data['book_id']) # [B, E]
21     embedded_genre = self.genre_embedding(data['genre_idx']) # [B, E]
22     deep_input = torch.cat([embedded_user_ids, embedded_book_ids, embedded_genre], dim=1)
23     deep_output = self.ffnn(deep_input) # [64 predictions, 1]
24
25     #
26     # WIDE
27     #
28     wide_output = self.wide(data['wide_input']) # [64 predictions, 1]
29
30     #
31     # Combine outputs (w/ learnable parameters)
32     #
33     logits = self.deep_weight * deep_output + self.wide_weight * wide_output
34
35     return logits.view(-1) # [64,]
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64

```

▼ Training

▼ Training Parameters & GPU

```

1 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
2
3 label_counts = Counter(dataset.data['rating'].apply(lambda x: 1.0 if x >= 4 else 0))
4 pos_weight = torch.tensor([label_counts[0] / label_counts[1]], device=device)
5
6 model = WideAndDeepModel()
7 loss_fn = nn.BCEWithLogitsLoss(pos_weight=pos_weight)
8 optimizer = torch.optim.Adam(model.parameters(), lr=0.001)
9
10 writer = SummaryWriter('./logs/')
11 best_test_loss = float('inf')
12
13 SAVE_PATH = "/content/drive/MyDrive/book_recommender_best_model.pt"
14 EPOCHS = 5
15 GLOBAL_STEP = 0

```

▼ Loop

```

1 #
2 # Training loop
3 #
4 for epoch in range(EPOCHS):
5     running_loss = 0.0
6
7     for batch_idx, data in enumerate(train_dataloader):
8         model.train()
9         optimizer.zero_grad()
10
11     preds = model(data) # [B]
12     true_labels = data['rating']
13
14     loss = loss_fn(preds, true_labels)

```

```
15     loss.backward()
16     optimizer.step()
17
18     running_loss += loss.item()
19     GLOBAL_STEP += 1
20
21     if batch_idx % 100 == 99:
22         avg_train_loss = running_loss / 100
23
24         print(f"Epoch {epoch+1}, Batch {batch_idx+1}, Train Loss: {avg_train_loss:.4f}")
25
26         writer.add_scalar("Loss/train", avg_train_loss, GLOBAL_STEP)
27
28     #
29     # Test set
30     #
31
32     model.eval()
33     test_loss = 0.0
34     total_test_batches = 0
35
36     with torch.no_grad():
37         test_data = next(iter(test_dataloader))
38
39         test_preds = model(test_data)
40         test_labels = test_data['rating']
41
42         loss_val = loss_fn(test_preds, test_labels)
43         test_loss += loss_val.item()
44         total_test_batches += 1
45
46         avg_test_loss = test_loss / total_test_batches
47
48         writer.add_scalar("Loss/test", avg_test_loss, GLOBAL_STEP)
49
50     #
51     # Saving the best model based on test-loss (just a rough approx.)
52     #
53     if avg_test_loss < best_test_loss:
54         best_test_loss = avg_test_loss
55         torch.save(model.state_dict(), SAVE_PATH)
56
57     running_loss = 0.0
58
59 writer.close()
60
```

```
1 %load_ext tensorboard
2 %tensorboard --logdir ./logs/
```



TensorBoard

TIME SERIES SCALARS

INACTIVE

Filter tags (regex)

All Scalars Image Histogram

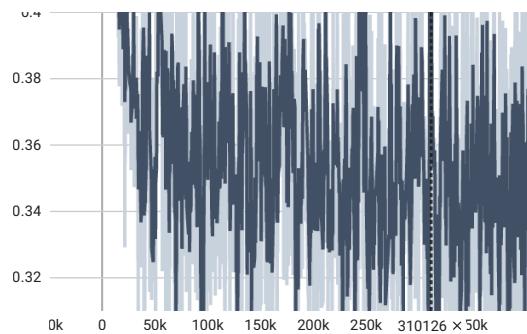
Settings

Pinned

Pin cards for a quick view and comparison

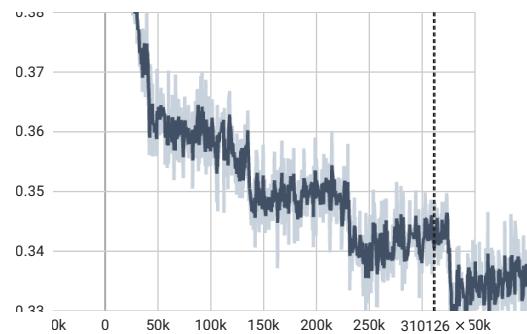
Loss 2 cards

Loss/test



Run ↑	Smoothed	Value	Step	Relative
.	0.3356	0.3296	310,126	9.49 hr

Loss/train



Run ↑	Smoothed	Value	Step	Relative
.	0.3352	0.3346	310,126	9.49 hr

Calculating AUC

With only three features (user_id, book_id, genre), the model was able to achieve 0.779 AUC score. While it could theoretically be stronger than this, I am pleased to see that with such simple features, the wide and deep network was able to learn a relationship between these features and the rating one would give.

```

1 from sklearn.metrics import roc_auc_score
2
3 model.eval()
4 all_labels = []
5 all_preds = []
6
7 with torch.no_grad():
8     for idx, batch in enumerate(test_dataloader):
9         preds = model(batch)
10        labels = batch['rating']
11
12        probs = torch.sigmoid(preds).squeeze().cpu().numpy()
13
14        all_preds.extend(probs)
15        all_labels.extend(labels.cpu().numpy())
16
17 auc = roc_auc_score(all_labels, all_preds)
18 positive_rate = sum(all_labels) / len(all_labels)
19
20 print(f"AUC: {auc:.4f}")
21 print(f"Positive rate: {positive_rate:.4f}")
22

```

AUC: 0.7789
Positive rate: 0.6898

```
1 # import json
2 # with open("book_recommender_vocabs.json", "w") as f:
3 #     json.dump({"user_vocab": dataset.user_vocab,
4 #                 "book_vocab": dataset.book_vocab,
5 #                 "genre_vocab": dataset.genre_vocab}, f)
```

▼ Conclusion