

Python使用Faiss库实现向量近邻搜索

原创 蚂蚁 蚂蚁学Python 8月2日

本文是Python应用于推荐系统领域的技术文章。

Embedding的近邻搜索是当前图推荐系统非常重要的一种召回方式，通过item2vec、矩阵分解、双塔DNN等方式都能够产出训练好的user embedding、item embedding，对于embedding的使用非常的灵活：

- 输入user embedding，近邻搜索item embedding，可以给user推荐感兴趣的items
- 输入user embedding，近邻搜索user embedding，可以给user推荐感兴趣的user
- 输入item embedding，近邻搜索item embedding，可以给item推荐相关的items

然而有一个工程问题，一旦user embedding、item embedding数据量达到一定的程度，对他们的近邻搜索将会变得非常慢，如果离线阶段提前搜索好在高速缓存比如redis存储好结果当然没问题，但是这种方式很不实时，如果能在线阶段上线几十MS的搜索当然效果最好。

Faiss是Facebook AI团队开源的针对聚类 and 相似性搜索库，为稠密向量提供高效相似度搜索和聚类，支持十亿级别向量的搜索，是目前最为成熟的近似近邻搜索库。

接下来通过jupyter notebook的代码，给大家演示下使用faiss的简单流程，内容包括：

1. 读取训练好的Embedding数据
2. 构建faiss索引，将待搜索的Embedding添加进去
3. 取得目标Embedding，实现搜索得到ID列表
4. 根据ID获取电影标题，返回结果

1、读取预训练好的Embedding数据

```
1 import pandas as pd
2 import numpy as np
3
```

```
4 df = pd.read_csv("./datas/movielens_sparkals_item_embedding.csv")
5
6 # 提取要使用的电影ID列表，注意要转换成int64
7 ids = df["id"].values.astype(np.int64)
8
9 # 记录ID列表的大小
10 ids_size = ids.shape[0]
11
12 # 读取内容embedding数据，转换成二维array
13 import json
14 import numpy as np
15 datas = []
16 for x in df["features"]:
17     datas.append(json.loads(x))
18 # 变成二维array
19 datas = np.array(datas).astype(np.float32)
20
21 # 记录数据维度
22 dimension = datas.shape[1]
```

2、使用faiss建立索引

```
1 import faiss
2
3 index = faiss.IndexFlatL2(dimension)
4 index2 = faiss.IndexIDMap(index)
5
6 index2.add_with_ids(datas, ids)
```

3、实现近邻搜索

```
1 # 读取user embedding数据
2 df_user = pd.read_csv("./datas/user_embedding.csv")
3 df_user.head()
4
5 # 挑选一条user的embedding，转换成1行N列的二维array
6 user_embedding = np.array(json.loads(df_user[df_user["id"] == 10]["features"]))
7 user_embedding = np.expand_dims(user_embedding, axis=0).astype(np.float32)
```

```
8 user_embedding
9
10 # 实现搜索, 这里的I就是近邻ID列表
11 tok = 30
12 D, I = index.search(user_embedding, topk)    # actual search
```

4、拼接得到内容列表

```
1 # 把搜索出来的ID变成Series
2 target_ids = pd.Series(I[0], name="MovieID")
3
4 # 读取电影数据
5 df_movie = pd.read_csv("./datas/ml-1m/movies.dat",
6                         sep="::", header=None, engine="python",
7                         names = "MovieID::Title::Genres".split("::"))
8 # 实现内容JOIN
9 df_result = pd.merge(target_ids, df_movie)
```

df_result就是最终可以返回给前端的推荐的电影ID列表了。

如果本文对你有帮助的话, 帮忙点个赞吧^_^

代码: <https://github.com/peiss/ant-learn-recsys>