

Web 信息处理与应用实验 2 报告

翁屹禾 PB20000017

傅申 PB20000051

侯博文 PB20000054

目录

1. 阶段一: 图谱抽取	1
1.1. 获取豆瓣电影对应的实体	1
1.2. 抽取一跳子图	1
1.3. 抽取二跳子图	2
1.4. 运行结果	3
2. 阶段二: 图谱推荐	4
2.1. 映射知识图谱中的实体和关系	4
2.2. 基于图嵌入的模型	5
2.2.1. 数据加载部分	5
2.2.2. 模型搭建部分	5
2.3. 实验结果	7

1. 阶段一: 图谱抽取

该阶段的代码结构如下:

```
1 stage-1
2 |--- data
3 |   |--- douban2fb.txt
4 |   |--- freebase_douban.gz
5 |   |--- Movie_id.csv
6 |   |--- Movie_tag.csv
7 |   |--- kg_extracted.txt    # 提取的图谱
8 |--- kg_extraction.py        # 图谱抽取代码
```

1.1. 获取豆瓣电影对应的实体

在给出的链接信息文件 `douban2fb.txt` 中, 提供了豆瓣电影 ID 到图谱实体 ID 之间的映射关系. 因此, 可以通过该文件获取豆瓣电影对应的实体.

```
1 def load_movie_entities() -> set[str]:
2     # load movie entities
3     entities = set()
4     with open(LINK_INFO_PATH, "r") as f:
5         for line in f:
6             line = line.strip()
7             entities.add(f"{ENTITY_PREFIX}{line.split()[-1]}>")
8     return entities
```

1.2. 抽取一跳子图

首先, 获取所有以电影实体作为头实体的三元组, 作为一跳子图:

```

1 movie_entities = load_movie_entities()
2 first_hop = extract_subgraph(movie_entities)

```

其中 `extract_subgraph()` 函数的实现如下, 它约束了实体的前缀, 以保证抽取的知识图谱更加精简:

```

1 def extract_subgraph(entities: set[str]) -> kg_t:
2     # load knowledge graph with only provided entities
3     kg = set()
4     with gzip.open(KG_PATH, "rb") as f:
5         for line in process(f, KG_SIZE):
6             line = line.strip()
7             triplet = tuple(line.decode().split()[:3])
8             if triplet[0] in entities and triplet[2].startswith(ENTITY_PREFIX):
9                 kg.add(triplet)
10    return kg

```

然后过滤一跳子图, 只保留至少出现在 20 个三元组中的实体, 同时只保留出现超过 50 次的关系:

```

1 entity_count, relation_count = kg_info(first_hop)
2 first_hop = filter(
3     first_hop,
4     lambda triplet: entity_count[triplet[0]] > 20
5                     and entity_count[triplet[2]] > 20
6                     and relation_count[triplet[1]] > 50,
7 )

```

其中 `kg_info()` 函数用于统计实体和关系的出现次数, `filter()` 函数用于过滤子图:

```

1 def kg_info(kg: kg_t):
2     """get the entity and relation count of a knowledge graph"""
3     entity_count: dict[str, int] = {}
4     relation_count: dict[str, int] = {}
5     for start, relation, end in kg:
6         entity_count[start] = entity_count.get(start, 0) + 1
7         entity_count[end] = entity_count.get(end, 0) + 1
8         relation_count[relation] = relation_count.get(relation, 0) + 1
9     return entity_count, relation_count
10
11
12 def filter(kg: kg_t, filter_fn: Callable[[kg_entry_t], bool]):
13     filtered = set()
14     for triplet in process(kg, len(kg)):
15         if filter_fn(triplet):
16             filtered.add(triplet)
17     return filtered

```

在实际运行中, 得到的子图包含了 771 个实体和 31 个关系, 共 24851 个三元组。

1.3. 抽取二跳子图

首先, 使用 `second_hop = hop(first_hop)` 来获取二跳子图, 其中 `hop()` 通过将子图中所有出现过的实体作为头实体, 重新抽取子图:

```

1 def hop(kg: kg_t):

```

```
2 entities = {triplet[0] for triplet in kg}
3 entities.update({triplet[2] for triplet in kg})
4 return extract_subgraph(entities)
```

然后,对二跳子图进行处理,先过滤掉出现超过两万次的实体和出现少于 50 次的关系:

```
1 entity_count, relation_count = kg_info(second_hop)
2 second_hop = filter(
3     second_hop,
4     lambda triplet: entity_count[triplet[0]] < 20000
5     and entity_count[triplet[2]] < 20000
6     and relation_count[triplet[1]] > 50,
7 )
```

然后再次过滤,只保留出现大于 15 次的实体和出现大于 50 次的关系:

```
1 entity_count, relation_count = kg_info(second_hop)
2 second_hop = filter(
3     second_hop,
4     lambda triplet: entity_count[triplet[0]] > 15
5                     and entity_count[triplet[2]] > 15
6                     and relation_count[triplet[1]] > 50,
7 )
```

最后将得到的二跳子图写入文件:

```
1 with open(OUTPUT_KG_PATH, "w") as f:
2     for triplet in process(second_hop):
3         f.write(" ".join(triplet) + "\n")
```

1.4. 运行结果

程序的运行过程如下图所示,最终得到的子图包含 1938 个实体,56 个关系,共 43711 个三元组:

```
Loading movie entities ...  
100%|██████████████████████████████████████████████████████████████████████████████| 578/578 [00:00<00:00, 1245150.34it/s]  
Done, size: 578  
  
First hop ...  
100%|██████████████████████████████████████████████████████████████████████████████| 395577070/395577070 [13:06<00:00, 503176.11it/s]  
Done, size: 125372  
Filtering first hop ...  
100%|██████████████████████████████████████████████████████████████████████████████| 125372/125372 [00:00<00:00, 1246916.27it/s]  
Done, 771 entities, 31 relations, size: 24851  
  
Second hop ...  
100%|██████████████████████████████████████████████████████████████████████████████| 395577070/395577070 [14:52<00:00, 443100.84it/s]  
Done, size: 104742330  
Filtering second hop ...  
100%|██████████████████████████████████████████████████████████████████████████████| 104742330/104742330 [01:16<00:00, 1361481.01it/s]  
100%|██████████████████████████████████████████████████████████████████████████████| 762594/762594 [00:01<00:00, 671881.79it/s]  
Done, 1938 entities, 56 relations, size: 43711  
  
Saving extracted knowledge graph ...  
43711it [00:00. 637306.76it/s]
```

图 1 阶段一运行结果

2. 阶段二: 图谱推荐

该阶段的代码结构如下:

```

1 stage-2
2 |--- data
3 |   |--- Douban
4 |   |   |--- entity_map.txt    # 实体映射
5 |   |   |--- kg_final.txt      # 映射后的知识图谱
6 |   |   |--- relation_map.txt  # 关系映射
7 |   |   |--- test.txt
8 |   |   |--- train.txt
9 |   |--- douban2fb.txt
10 |   |--- kg_extracted.txt      # 阶段一中提取的知识图谱
11 |   |--- movie_id_map.txt
12 |   |--- user_id_map.txt
13 |--- data_loader/*
14 |--- model/*
15 |--- parser/*
16 |--- trained_model/*
17 |--- utils/*
18 |--- main_Embedding_based.py
19 |--- main_GNN_based.py
20 |--- main_KG_free.py
21 |--- mapping.py                # 实体和关系映射代码

```

2.1. 映射知识图谱中的实体和关系

首先, 加载知识图谱 `kg_extracted.txt`、链接信息 `douban2fb.txt` 和电影 ID 到索引值的映射关系 `movie_id_map.txt`.

```

1 raw_kg = load_raw_kg()
2 id_to_entity = load_link_info()
3 id_to_index = load_id_map()

```

然后, 先将电影 ID 对应的实体映射到相应的索引值:

```

1 entity_to_index = {}
2 relation_to_index = {}
3 for id, entity in id_to_entity.items():
4     entity_to_index[entity] = id_to_index[id]

```

再将剩余的实体和关系映射到相应的索引值:

```

1 entity_index = max(id_to_index.values()) + 1
2 relation_index = 0
3 for start, relation, end in raw_kg:
4     if start not in entity_to_index:
5         entity_to_index[start] = entity_index
6         entity_index += 1
7     if relation not in relation_to_index:
8         relation_to_index[relation] = relation_index
9         relation_index += 1
10    if end not in entity_to_index:
11        entity_to_index[end] = entity_index
12        entity_index += 1

```

最后, 将实体映射、关系映射和映射后的知识图谱写入文件:

```
1 with open(ENTITY_MAP_PATH, "w") as f:
2     for entity, index in entity_to_index.items():
3         f.write(entity + " " + str(index) + "\n")
4 with open(RELATION_MAP_PATH, "w") as f:
5     for relation, index in relation_to_index.items():
6         f.write(relation + " " + str(index) + "\n")
7
8 with open(FINAL_KG_PATH, "w") as f:
9     for triplet in map_kg(raw_kg, entity_to_index, relation_to_index):
10        f.write(" ".join(triplet) + "\n")
```

2.2. 基于图嵌入的模型

2.2.1. 数据加载部分

数据加载部分位于 stage-2/dataloader/loader_Embedding_based.py 文件内, 需要补全 `construct_data()` 函数.

一、添加逆向三元组: 使用 `rename()` 函数将深拷贝的三元组中的头实体和尾实体交换, 再给它们的关系加上 `n_relations`, 最后将源三元组和逆向三元组连接起来.

```
1 inverted_kg_data = copy.deepcopy(kg_data)
2 inverted_kg_data = inverted_kg_data.rename({"h": "t", "t": "h"}, axis=1)
3 inverted_kg_data["r"] += max(kg_data["r"]) + 1
4 self.kg_data = pd.concat(
5     [kg_data, inverted_kg_data],
6     axis=0,
7     ignore_index=True
8 )
```

二、计算关系数, 实体数和三元组数量.

```
1 self.n_relations = self.kg_data["r"].max() + 1
2 self.n_entities = self.kg_data["h"].max() + 1
3 self.n_kg_data = self.kg_data.shape[0]
```

三、构建字典 `kg_dict` 和 `relation_dict`:

```
1 self.kg_dict = collections.defaultdict(list)
2 self.relation_dict = collections.defaultdict(list)
3 for _, (h, r, t) in self.kg_data.iterrows():
4     self.kg_dict[h].append((t, r))
5     self.relation_dict[r].append((h, t))
```

2.2.2. 模型搭建部分

模型搭建部分位于 stage-2/model/Embedding_based.py 文件内, 需要补全相关代码, 实现 TransE 和 TransR 算法.

`calc_kg_loss_TransE()` 函数和 `calc_kg_loss_TransR()` 函数分别根据 TransE 算法和 TransR 算法计算嵌入的损失函数, 需要进行补全.

TransE 算法

TransE 算法的基本思想是: 将实体和关系分别映射到对应空间中, 使得头实体和关系的和接近尾实体, 即 $\mathbf{h} + \mathbf{r} \approx \mathbf{t}$. 具体而言, 需要最小化下面的损失函数:

$$\mathcal{L} = \sum_{t_r \in T_r} \sum_{t_r' \in T_r'} \max(0, \gamma + f(t_r) - f(t_r'))$$

其中, 得分 $f(\cdot)$ 计算了 $\mathbf{h} + \mathbf{r}$ 和 \mathbf{t} 的距离. 换言之, 我们需要尽可能地使得负样例的得分大于正样例的得分. 如果采用 BPR 损失函数, 则有:

$$\mathcal{L}_{\text{BPR}} = - \sum_{t_r \in T_r} \sum_{t_r' \in T_r'} \log(\sigma(f(t_r') - f(t_r)))$$

基于上述思想, 可以完成 `calc_kg_loss_TransE()` 函数.

一、计算头实体、关系和正负样例尾实体在空间中的嵌入:

```
1 r_embed = self.relation_embed(r) # (kg_batch_size, relation_dim)
2
3 h_embed = self.entity_embed(h) # (kg_batch_size, embed_dim)
4 pos_t_embed = self.entity_embed(pos_t) # (kg_batch_size, embed_dim)
5 neg_t_embed = self.entity_embed(neg_t) # (kg_batch_size, embed_dim)
```

二、对这些嵌入进行 L2 归一化:

```
1 r_embed = r_embed / torch.norm(r_embed, dim=1, keepdim=True)
2 h_embed = h_embed / torch.norm(h_embed, dim=1, keepdim=True)
3 pos_t_embed = pos_t_embed / torch.norm(pos_t_embed, dim=1, keepdim=True)
4 neg_t_embed = neg_t_embed / torch.norm(neg_t_embed, dim=1, keepdim=True)
```

三、计算正负样例的得分 $f(\cdot)$, 采用 L2 距离:

```
1 pos_diff = h_embed + r_embed - pos_t_embed # (kg_batch_size, embed_dim)
2 neg_diff = h_embed + r_embed - neg_t_embed # (kg_batch_size, embed_dim)
3 pos_score = torch.norm(pos_diff, dim=1) ** 2 # (kg_batch_size)
4 neg_score = torch.norm(neg_diff, dim=1) ** 2 # (kg_batch_size)
```

四、计算 BPR 损失函数:

```
1 kg_loss = -F.logsigmoid(neg_score - pos_score).mean()
```

TransR 算法

TransR 算法的基本思想与 TransE 算法类似. TransR 算法在实体空间和多个关系空间中建模实体和关系. 对于每个三元组 (h, r, t) , 将实体空间中的实体 \mathbf{h} 和 \mathbf{r} 通过矩阵 \mathbf{W}_r 投影到 r 关系空间中, 即

$$\mathbf{h}_r = \mathbf{h}\mathbf{W}_r, \mathbf{t}_r = \mathbf{t}\mathbf{W}_r$$

然后, 在 r 关系空间中, 尽可能使得 $\mathbf{h}_r + \mathbf{r} \approx \mathbf{t}_r$. 即尽可能使负样例得分大于正样例的得分.

一、计算头实体、关系和正负样例尾实体在空间中的嵌入, 获取 r 关系空间的投影矩阵 \mathbf{W}_r :

```
1 r_embed = self.relation_embed(r) # (kg_batch_size, relation_dim)
2 W_r = self.trans_M[r] # (embed_dim, relation_dim)
3
4 h_embed = self.entity_embed(h) # (kg_batch_size, embed_dim)
```

```

5 pos_t_embed = self.entity_embed(pos_t) # (kg_batch_size, embed_dim)
6 neg_t_embed = self.entity_embed(neg_t) # (kg_batch_size, embed_dim)

```

二、将头实体和尾实体投影到 r 关系空间中:

```

1 def trans_r_mul(W_r, entity_embed):
2     """
3     计算 TransR 中的投影嵌入
4     """
5     return torch.matmul(entity_embed.unsqueeze(1), W_r).squeeze(1)
6
7 r_mul_h = trans_r_mul(W_r, h_embed) # (kg_batch_size, relation_dim)
8 r_mul_pos_t = trans_r_mul(W_r, pos_t_embed) # (kg_batch_size, relation_dim)
9 r_mul_neg_t = trans_r_mul(W_r, neg_t_embed) # (kg_batch_size, relation_dim)

```

三、对这些嵌入进行 L2 归一化:

```

1 r_embed = r_embed / torch.norm(r_embed, dim=1, keepdim=True)
2 r_mul_h = r_mul_h / torch.norm(r_mul_h, dim=1, keepdim=True)
3 r_mul_pos_t = r_mul_pos_t / torch.norm(r_mul_pos_t, dim=1, keepdim=True)
4 r_mul_neg_t = r_mul_neg_t / torch.norm(r_mul_neg_t, dim=1, keepdim=True)

```

四、计算正负样例的得分 $f(\cdot)$, 采用 L2 距离:

```

1 pos_diff = r_mul_h + r_embed - r_mul_pos_t # (kg_batch_size, relation_dim)
2 neg_diff = r_mul_h + r_embed - r_mul_neg_t # (kg_batch_size, relation_dim)
3 pos_score = torch.norm(pos_diff, dim=1) ** 2 # (kg_batch_size)
4 neg_score = torch.norm(neg_diff, dim=1) ** 2 # (kg_batch_size)

```

五、计算 BPR 损失函数:

```

1 kg_loss = -F.logsigmoid(neg_score - pos_score).mean()

```

在后面的计算中, 还需要为物品嵌入注入实体嵌入的语义信息, 可以采用相加、逐元素相乘和拼接等方法:

```

2 def inject(item_embed, item_kg_embed):
3     """
4     为 物品嵌入 注入 实体嵌入 的语义信息, 可以采用相加/逐元素相乘/拼接等方式
5     """
6     return item_embed + item_kg_embed
7     # return item_embed * item_kg_embed
8     # return torch.cat([item_embed, item_kg_embed], dim=1)

```

如果使用拼接的方式, 则还需要对 `calc_cf_loss()` 和 `calc_score()` 中的 `user_embed` 扩充维度, 以保证维度一致

```

9 user_embed = torch.cat([user_embed, user_embed], dim=1)

```

2.3. 实验结果

我们运行了图谱嵌入模型的训练代码, 对比了不同的算法和注入语义信息方式对模型性能的影响, 同时与 Baseline (MF) 进行比较, 结果如下表所示 (表中为训练过程中最好的结果):

算法 + 语义信息注入方式		Recall@5	Recall@10	NDCG@5	NDCG@10
图谱嵌入 w/ TransE	相加	0.0676	0.1156	0.3104	0.2904
	逐元素相乘	0.0608	0.0999	0.2824	0.2557
	拼接	0.0676	0.1159	0.3102	0.2905
图谱嵌入 w/ TransR	相加	0.0658	0.1126	0.3118	0.2849
	逐元素相乘	0.0627	0.1087	0.2916	0.2690
	拼接	0.0657	0.1127	0.3117	0.2853
Baseline(MF)		0.0660	0.1094	0.3110	0.2829

表 1 实验结果对比

根据结果可以看出, TransE 和 TransR 算法在不同的指标上各有优劣, 而不同的语义信息注入方式对结果的影响比较一致, 即“相加” \approx “拼接” $>$ “逐元素相乘”. 在与 Baseline 的比较中, TransE 和 TransR 在采用相加/拼接的语义信息注入方式时, 可以取得稍微好一点的结果.