dmlc / **dgl**   Public

| `<>` **Code** | ⊙ Issues   337 | ⇄ Pull requests   110 | ▷ Actions | ⊞ Projects   1 | ⊘ Security | ⌇ In |

**dgl** / examples / pytorch / rgcn / **link.py** ⧉        •••

🐻 **frozenbugs** and **Ubuntu** examples (#5323)   •••   ✕      10 months ago   •••   ↺

352 lines (305 loc) · 11.1 KB

| Code | Blame |         Raw   ⧉   ⤓   `<>` |

```python
1   import dgl
2   import numpy as np
3   import torch
4   import torch.nn as nn
5   import torch.nn.functional as F
6   import tqdm
7   from dgl.data.knowledge_graph import FB15k237Dataset
8   from dgl.dataloading import GraphDataLoader
9   from dgl.nn.pytorch import RelGraphConv
10
11
12  # for building training/testing graphs
13  def get_subset_g(g, mask, num_rels, bidirected=False):
14      src, dst = g.edges()
15      sub_src = src[mask]
16      sub_dst = dst[mask]
17      sub_rel = g.edata["etype"][mask]
18
19      if bidirected:
20          sub_src, sub_dst = torch.cat([sub_src, sub_dst]), torch.cat(
21              [sub_dst, sub_src]
22          )
23          sub_rel = torch.cat([sub_rel, sub_rel + num_rels])
24
25      sub_g = dgl.graph((sub_src, sub_dst), num_nodes=g.num_nodes())
26      sub_g.edata[dgl.ETYPE] = sub_rel
27      return sub_g
28
29
30  class GlobalUniform:
31      def __init__(self, g, sample_size):
32          self.sample_size = sample_size
33          self.eids = np.arange(g.num_edges())
34
35      def sample(self):
36          return torch.from_numpy(np.random.choice(self.eids, self.sample_size))
37
38
39  class NegativeSampler:
```

```python
40          def __init__(self, k=10):  # negative sampling rate = 10
41              self.k = k
42
43   v      def sample(self, pos_samples, num_nodes):
44              batch_size = len(pos_samples)
45              neg_batch_size = batch_size * self.k
46              neg_samples = np.tile(pos_samples, (self.k, 1))
47
48              values = np.random.randint(num_nodes, size=neg_batch_size)
49              choices = np.random.uniform(size=neg_batch_size)
50              subj = choices > 0.5
51              obj = choices <= 0.5
52              neg_samples[subj, 0] = values[subj]
53              neg_samples[obj, 2] = values[obj]
54              samples = np.concatenate((pos_samples, neg_samples))
55
56              # binary labels indicating positive and negative samples
57              labels = np.zeros(batch_size * (self.k + 1), dtype=np.float32)
58              labels[:batch_size] = 1
59
60              return torch.from_numpy(samples), torch.from_numpy(labels)
61
62
63   v  class SubgraphIterator:
64   v      def __init__(self, g, num_rels, sample_size=30000, num_epochs=6000):
65              self.g = g
66              self.num_rels = num_rels
67              self.sample_size = sample_size
68              self.num_epochs = num_epochs
69              self.pos_sampler = GlobalUniform(g, sample_size)
70              self.neg_sampler = NegativeSampler()
71
72          def __len__(self):
73              return self.num_epochs
74
75   v      def __getitem__(self, i):
76              eids = self.pos_sampler.sample()
77              src, dst = self.g.find_edges(eids)
78              src, dst = src.numpy(), dst.numpy()
79              rel = self.g.edata[dgl.ETYPE][eids].numpy()
80
81              # relabel nodes to have consecutive node IDs
82              uniq_v, edges = np.unique((src, dst), return_inverse=True)
83              num_nodes = len(uniq_v)
84              # edges is the concatenation of src, dst with relabeled ID
85              src, dst = np.reshape(edges, (2, -1))
86              relabeled_data = np.stack((src, rel, dst)).transpose()
87
88              samples, labels = self.neg_sampler.sample(relabeled_data, num_nodes)
89
90              # use only half of the positive edges
91              chosen_ids = np.random.choice(
92                  np.arange(self.sample_size),
93                  size=int(self.sample_size / 2),
94                  replace=False,
```

```python
 95            )
 96            src = src[chosen_ids]
 97            dst = dst[chosen_ids]
 98            rel = rel[chosen_ids]
 99            src, dst = np.concatenate((src, dst)), np.concatenate((dst, src))
100            rel = np.concatenate((rel, rel + self.num_rels))
101            sub_g = dgl.graph((src, dst), num_nodes=num_nodes)
102            sub_g.edata[dgl.ETYPE] = torch.from_numpy(rel)
103            sub_g.edata["norm"] = dgl.norm_by_dst(sub_g).unsqueeze(-1)
104            uniq_v = torch.from_numpy(uniq_v).view(-1).long()
105
106            return sub_g, uniq_v, samples, labels
107
108
109    class RGCN(nn.Module):
110        def __init__(self, num_nodes, h_dim, num_rels):
111            super().__init__()
112            # two-layer RGCN
113            self.emb = nn.Embedding(num_nodes, h_dim)
114            self.conv1 = RelGraphConv(
115                h_dim,
116                h_dim,
117                num_rels,
118                regularizer="bdd",
119                num_bases=100,
120                self_loop=True,
121            )
122            self.conv2 = RelGraphConv(
123                h_dim,
124                h_dim,
125                num_rels,
126                regularizer="bdd",
127                num_bases=100,
128                self_loop=True,
129            )
130            self.dropout = nn.Dropout(0.2)
131
132        def forward(self, g, nids):
133            x = self.emb(nids)
134            h = F.relu(self.conv1(g, x, g.edata[dgl.ETYPE], g.edata["norm"]))
135            h = self.dropout(h)
136            h = self.conv2(g, h, g.edata[dgl.ETYPE], g.edata["norm"])
137            return self.dropout(h)
138
139
140    class LinkPredict(nn.Module):
141        def __init__(self, num_nodes, num_rels, h_dim=500, reg_param=0.01):
142            super().__init__()
143            self.rgcn = RGCN(num_nodes, h_dim, num_rels * 2)
144            self.reg_param = reg_param
145            self.w_relation = nn.Parameter(torch.Tensor(num_rels, h_dim))
146            nn.init.xavier_uniform_(
147                self.w_relation, gain=nn.init.calculate_gain("relu")
148            )
149
150        def calc score(self   embedding   triplets):
```

```python
150          def calc_score(self, embedding, triplets):
151              s = embedding[triplets[:, 0]]
152              r = self.w_relation[triplets[:, 1]]
153              o = embedding[triplets[:, 2]]
154              score = torch.sum(s * r * o, dim=1)
155              return score
156
157          def forward(self, g, nids):
158              return self.rgcn(g, nids)
159
160          def regularization_loss(self, embedding):
161              return torch.mean(embedding.pow(2)) + torch.mean(self.w_relation.pow(2))
162
163          def get_loss(self, embed, triplets, labels):
164              # each row in the triplets is a 3-tuple of (source, relation, destination)
165              score = self.calc_score(embed, triplets)
166              predict_loss = F.binary_cross_entropy_with_logits(score, labels)
167              reg_loss = self.regularization_loss(embed)
168              return predict_loss + self.reg_param * reg_loss
169
170
171      def filter(
172          triplets_to_filter, target_s, target_r, target_o, num_nodes, filter_o=True
173      ):
174          """Get candidate heads or tails to score"""
175          target_s, target_r, target_o = int(target_s), int(target_r), int(target_o)
176          # Add the ground truth node first
177          if filter_o:
178              candidate_nodes = [target_o]
179          else:
180              candidate_nodes = [target_s]
181          for e in range(num_nodes):
182              triplet = (
183                  (target_s, target_r, e) if filter_o else (e, target_r, target_o)
184              )
185              # Do not consider a node if it leads to a real triplet
186              if triplet not in triplets_to_filter:
187                  candidate_nodes.append(e)
188          return torch.LongTensor(candidate_nodes)
189
190
191      def perturb_and_get_filtered_rank(
192          emb, w, s, r, o, test_size, triplets_to_filter, filter_o=True
193      ):
194          """Perturb subject or object in the triplets"""
195          num_nodes = emb.shape[0]
196          ranks = []
197          for idx in tqdm.tqdm(range(test_size), desc="Evaluate"):
198              target_s = s[idx]
199              target_r = r[idx]
200              target_o = o[idx]
201              candidate_nodes = filter(
202                  triplets_to_filter,
203                  target_s,
204                  target_r,
205                  target_o,
```

```python
206                  num_nodes,
207                  filter_o=filter_o,
208              )
209              if filter_o:
210                  emb_s = emb[target_s]
211                  emb_o = emb[candidate_nodes]
212              else:
213                  emb_s = emb[candidate_nodes]
214                  emb_o = emb[target_o]
215              target_idx = 0
216              emb_r = w[target_r]
217              emb_triplet = emb_s * emb_r * emb_o
218              scores = torch.sigmoid(torch.sum(emb_triplet, dim=1))
219
220              _, indices = torch.sort(scores, descending=True)
221              rank = int((indices == target_idx).nonzero())
222              ranks.append(rank)
223      return torch.LongTensor(ranks)
224
225
226  def calc_mrr(
227      emb, w, test_mask, triplets_to_filter, batch_size=100, filter=True
228  ):
229      with torch.no_grad():
230          test_triplets = triplets_to_filter[test_mask]
231          s, r, o = test_triplets[:, 0], test_triplets[:, 1], test_triplets[:, 2]
232          test_size = len(s)
233          triplets_to_filter = {
234              tuple(triplet) for triplet in triplets_to_filter.tolist()
235          }
236          ranks_s = perturb_and_get_filtered_rank(
237              emb, w, s, r, o, test_size, triplets_to_filter, filter_o=False
238          )
239          ranks_o = perturb_and_get_filtered_rank(
240              emb, w, s, r, o, test_size, triplets_to_filter
241          )
242          ranks = torch.cat([ranks_s, ranks_o])
243          ranks += 1  # change to 1-indexed
244          mrr = torch.mean(1.0 / ranks.float()).item()
245      return mrr
246
247
248  def train(
249      dataloader,
250      test_g,
251      test_nids,
252      test_mask,
253      triplets,
254      device,
255      model_state_file,
256      model,
257  ):
258      optimizer = torch.optim.Adam(model.parameters(), lr=1e-2)
259      best_mrr = 0
260      for epoch, batch_data in enumerate(dataloader):  # single graph batch
```

```python
261                model.train()
262                g, train_nids, edges, labels = batch_data
263                g = g.to(device)
264                train_nids = train_nids.to(device)
265                edges = edges.to(device)
266                labels = labels.to(device)
267
268                embed = model(g, train_nids)
269                loss = model.get_loss(embed, edges, labels)
270                optimizer.zero_grad()
271                loss.backward()
272                nn.utils.clip_grad_norm_(
273                    model.parameters(), max_norm=1.0
274                )  # clip gradients
275                optimizer.step()
276                print(
277                    "Epoch {:04d} | Loss {:.4f} | Best MRR {:.4f}".format(
278                        epoch, loss.item(), best_mrr
279                    )
280                )
281            if (epoch + 1) % 500 == 0:
282                # perform validation on CPU because full graph is too large
283                model = model.cpu()
284                model.eval()
285                embed = model(test_g, test_nids)
286                mrr = calc_mrr(
287                    embed, model.w_relation, test_mask, triplets, batch_size=500
288                )
289                # save best model
290                if best_mrr < mrr:
291                    best_mrr = mrr
292                    torch.save(
293                        {"state_dict": model.state_dict(), "epoch": epoch},
294                        model_state_file,
295                    )
296                model = model.to(device)
297
298
299    if __name__ == "__main__":
300        device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
301        print(f"Training with DGL built-in RGCN module")
302
303        # load and preprocess dataset
304        data = FB15k237Dataset(reverse=False)
305        g = data[0]
306        num_nodes = g.num_nodes()
307        num_rels = data.num_rels
308        train_g = get_subset_g(g, g.edata["train_mask"], num_rels)
309        test_g = get_subset_g(g, g.edata["train_mask"], num_rels, bidirected=True)
310        test_g.edata["norm"] = dgl.norm_by_dst(test_g).unsqueeze(-1)
311        test_nids = torch.arange(0, num_nodes)
312        test_mask = g.edata["test_mask"]
313        subg_iter = SubgraphIterator(train_g, num_rels)  # uniform edge sampling
314        dataloader = GraphDataLoader(
315            subg_iter, batch_size=1, collate_fn=lambda x: x[0]
```

```python
316              )
317
318              # Prepare data for metric computation
319              src, dst = g.edges()
320              triplets = torch.stack([src, g.edata["etype"], dst], dim=1)
321
322              # create RGCN model
323              model = LinkPredict(num_nodes, num_rels).to(device)
324
325              # train
326              model_state_file = "model_state.pth"
327              train(
328                  dataloader,
329                  test_g,
330                  test_nids,
331                  test_mask,
332                  triplets,
333                  device,
334                  model_state_file,
335                  model,
336              )
337
338              # testing
339              print("Testing...")
340              checkpoint = torch.load(model_state_file)
341              model = model.cpu()  # test on CPU
342              model.eval()
343              model.load_state_dict(checkpoint["state_dict"])
344              embed = model(test_g, test_nids)
345              best_mrr = calc_mrr(
346                  embed, model.w_relation, test_mask, triplets, batch_size=500
347              )
348              print(
349                  "Best MRR {:.4f} achieved using the epoch {:04d}".format(
350                      best_mrr, checkpoint["epoch"]
351                  )
352              )
```