

# 超大规模数据集类的创建

在前面的学习中我们只接触了数据可全部储存于内存的数据集，这些数据集对应的数据集类在创建对象时就将所有数据都加载到内存。然而在一些应用场景中，**数据集规模超级大，我们很难有足够大的内存完全存下所有数据**。因此需要一个**按需加载样本到内存的数据集类**。在此上半节内容中，我们将学习为一个包含上千万个图样本的数据集构建一个数据集类。

## Dataset基类简介

在PyG中，我们通过继承`torch_geometric.data.Dataset`基类来自定义一个按需加载样本到内存的数据集类。此基类与Torchvision的Dataset类的概念密切相关，这与第6节中介绍的`torch_geometric.data.InMemoryDataset`基类是一样的。

**继承`torch_geometric.data.InMemoryDataset`基类要实现的方法，继承此基类同样要实现，此外还需要实现以下方法：**

- `len()`：返回数据集中的样本的数量。
- `get()`：实现加载单个图的操作。注意：在内部，`__getitem__()`返回通过调用`get()`来获取Data对象，并根据`transform`参数对它们进行选择转换。

下面让我们通过一个简化的例子看**继承`torch_geometric.data.Dataset`基类的规范**：

```
1 import os.path as osp
2
3 import torch
4 from torch_geometric.data import Dataset, download_url
5
6 class MyOwnDataset(Dataset):
7     def __init__(self, root, transform=None, pre_transform=None):
```

```

8         super(MyOwnDataset, self).__init__(root, transform,
pre_transform)
9
10     @property
11     def raw_file_names(self):
12         return ['some_file_1', 'some_file_2', ...]
13
14     @property
15     def processed_file_names(self):
16         return ['data_1.pt', 'data_2.pt', ...]
17
18     def download(self):
19         # Download to `self.raw_dir`.
20         path = download_url(url, self.raw_dir)
21         ...
22
23     def process(self):
24         i = 0
25         for raw_path in self.raw_paths:
26             # Read data from `raw_path`.
27             data = Data(...)
28
29             if self.pre_filter is not None and not
self.pre_filter(data):
30                 continue
31
32             if self.pre_transform is not None:
33                 data = self.pre_transform(data)
34
35             torch.save(data, osp.join(self.processed_dir,
'data_{}.pt'.format(i)))
36             i += 1
37
38     def len(self):
39         return len(self.processed_file_names)
40
41     def get(self, idx):
42         data = torch.load(osp.join(self.processed_dir,
'data_{}.pt'.format(idx)))
43         return data
44

```

其中，每个Data对象在process()方法中单独被保存，并在get()中通过指定索引进行加载。

## *跳过download/process*

对于无需下载数据集原文件的情况，我们不重写（override）download方法即可跳过下载。对于无需对数据集做预处理的情况，我们不重写process方法即可跳过预处理。

## *无需定义Dataset类*

通过下面的方式，我们可以不用定义一个Dataset类，而直接生成一个DataLoader对象，直接用于训练：

```
1 from torch_geometric.data import Data, DataLoader
2
3 data_list = [Data(...), ..., Data(...)]
4 loader = DataLoader(data_list, batch_size=32)
```

我们也可以通过下面的方式将一个列表的Data对象组成一个batch：

```
1 from torch_geometric.data import Data, Batch
2
3 data_list = [Data(...), ..., Data(...)]
4 loader = Batch.from_data_list(data_list, batch_size=32)
```

# 图样本封装成批（BATCHING）与DataLoader类

---

内容来源：[ADVANCED MINI-BATCHING](#)

## 合并小图组成大图

图可以有任意数量的节点和边，它不是规整的数据结构，因此对图数据封装成批的操作与对图像与序列等数据封装成批的操作不同。PyTorch Geometric中采用的将多个图封装成批的方式是，将小图作为连通组件（connected component）的形式合并，构建一个大图。于是小图的邻接矩阵存储在大图邻接矩阵的对角线上。大图的邻接矩阵、属性矩阵、预测目标矩阵分别为：

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_1 & & \\ & \ddots & \\ & & \mathbf{A}_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \vdots \\ \mathbf{X}_n \end{bmatrix}, \quad \mathbf{Y} = \begin{bmatrix} \mathbf{Y}_1 \\ \vdots \\ \mathbf{Y}_n \end{bmatrix}.$$

**此方法有以下关键的优势：**

- 依靠消息传递方案的GNN运算不需要被修改，因为消息仍然不能在属于不同图的两个节点之间交换。
- 没有额外的计算或内存的开销。例如，这个批处理程序的工作完全不需要对节点或边缘特征进行任何填充。请注意，邻接矩阵没有额外的内存开销，因为它们是以稀疏的方式保存的，只保留非零项，即边。

通过 `torch_geometric.data.DataLoader` 类，多个小图被封装成一个大图。`torch_geometric.data.DataLoader` 是 PyTorch 的 `DataLoader` 的子类，它覆盖了 `collate()` 函数，改函数定义了一列表的样本是如何封装成批的。因此，所有可以传递给 PyTorch `DataLoader` 的参数也可以传递给 PyTorch Geometric 的 `DataLoader`，例如，`num_workers`。

## 小图的属性增值与拼接

将小图存储到大图中时需要对小图的属性做一些修改，一个最显著的例子就是要对节点序号增值。在最一般的形式中，PyTorch Geometric 的 `DataLoader` 类会自动对 `edge_index` 张量增值，增加的值是当前被处理图的前面的图的累积节点数量。比方说，现在对第  $k$  个图的 `edge_index` 张量做增值，前面  $k - 1$  个图的累积节点数量为  $n$ ，那么对第  $k$  个图的 `edge_index` 张量的增值  $n$ 。增值后，对所有图的 `edge_index` 张量（其形状为  $[2, \text{num\_edges}]$ ）在第二维中连接起来。

然而，有一些特殊的场景中（如下所述），基于需求我们希望能修改这一行为。PyTorch Geometric 允许我们通过覆盖 `torch_geometric.data.__inc__()` 和 `torch_geometric.data.__cat_dim__()` 函数来实现我们希望的行为。在未做修改的情况下，它们在 `Data` 类中的定义如下。

```
1 def __inc__(self, key, value):
2     if 'index' in key or 'face' in key:
3         return self.num_nodes
4     else:
5         return 0
6
7 def __cat_dim__(self, key, value):
8     if 'index' in key or 'face' in key:
9         return 1
10    else:
11        return 0
```

我们可以看到，`__inc__()` 定义了两个连续的图的属性之间的增量大小，而 `__cat_dim__()` 定义了同一属性的图形张量应该在哪个维度上被连接起来。PyTorch Geometric 为存储在 `Data` 类中的每个属性调用此二函数，并以它们各自的 key 和值 item 作为参数。

在下面的内容中，我们将学习一些对 `__inc__()` 和 `__cat_dim__()` 的修改可能是绝对必要的案例。

## 图的匹配 (Pairs of Graphs)

如果你想在 `Data` 对象中存储多个图，例如用于图匹配等应用，我们需要确保所有这些图的正确封装成批行为。例如，考虑将两个图，一个源图  $G_s$  和一个目标图  $G_t$ ，存储在一个 `Data` 类中，即

```
1 class PairData(Data):
2     def __init__(self, edge_index_s, x_s, edge_index_t, x_t):
3         super(PairData, self).__init__()
4         self.edge_index_s = edge_index_s
5         self.x_s = x_s
6         self.edge_index_t = edge_index_t
7         self.x_t = x_t
8
```

在这种情况下，`edge_index_s`应该根据源图 $G_s$ 的节点数做增值，即`x_s.size(0)`，而`edge_index_t`应该根据目标图 $G_t$ 的节点数做增值，即`x_t.size(0)`。

```
1 class PairData(Data):
2     def __init__(self, edge_index_s, x_s, edge_index_t, x_t):
3         super(PairData, self).__init__()
4         self.edge_index_s = edge_index_s
5         self.x_s = x_s
6         self.edge_index_t = edge_index_t
7         self.x_t = x_t
8
9     def __inc__(self, key, value):
10        if key == 'edge_index_s':
11            return self.x_s.size(0)
12        if key == 'edge_index_t':
13            return self.x_t.size(0)
14        else:
15            return super().__inc__(key, value)
```

我们可以通过设置一个简单的测试脚本来测试我们的PairData批处理行为。

```
1 edge_index_s = torch.tensor([
2     [0, 0, 0, 0],
3     [1, 2, 3, 4],
4 ])
5 x_s = torch.randn(5, 16) # 5 nodes.
6 edge_index_t = torch.tensor([
7     [0, 0, 0],
8     [1, 2, 3],
9 ])
10 x_t = torch.randn(4, 16) # 4 nodes.
11
12 data = PairData(edge_index_s, x_s, edge_index_t, x_t)
13 data_list = [data, data]
14 loader = DataLoader(data_list, batch_size=2)
15 batch = next(iter(loader))
16
17 print(batch)
18 # Batch(edge_index_s=[2, 8], x_s=[10, 16], edge_index_t=[2, 6],
19 #        x_t=[8, 16])
20 print(batch.edge_index_s)
```

```

21 # tensor([[0, 0, 0, 0, 5, 5, 5, 5], [1, 2, 3, 4, 6, 7, 8, 9]])
22
23 print(batch.edge_index_t)
24 # tensor([[0, 0, 0, 4, 4, 4], [1, 2, 3, 5, 6, 7]])

```

到目前为止，一切看起来都很好! `edge_index_s`和`edge_index_t`被正确地封装成批了，即使在为 $G_s$ 和 $G_t$ 含有不同数量的节点时也是如此。然而，由于PyTorch Geometric无法识别PairData对象中实际的图，所以batch属性（将大图每个节点映射到其各自对应的小图）没有正确工作。此时就需要DataLoader的follow\_batch参数发挥作用。在这里，我们可以指定我们要为哪些属性维护批信息。

```

1 loader = DataLoader(data_list, batch_size=2, follow_batch=['x_s',
  'x_t'])
2 batch = next(iter(loader))
3
4 print(batch)
5 # Batch(edge_index_s=[2, 8], x_s=[10, 16], x_s_batch=[10],
6         edge_index_t=[2, 6], x_t=[8, 16], x_t_batch=[8])
7 print(batch.x_s_batch)
8 # tensor([0, 0, 0, 0, 0, 1, 1, 1, 1, 1])
9
10 print(batch.x_t_batch)
11 # tensor([0, 0, 0, 0, 1, 1, 1, 1])

```

可以看到，`follow_batch=['x_s', 'x_t']`现在成功地为节点特征`x_s`和`x_t`分别创建了名为`x_s_batch`和`x_t_batch`的赋值向量。这些信息现在可以用来在一个单一的'Batch'对象中对多个图进行聚合操作，例如，全局池化。

## 二部图 (Bipartite Graphs)

二部图的邻接矩阵定义两种类型的节点之间的连接关系。一般来说，不同类型的节点数量不需要一致，于是二部图的邻接矩阵 $A \in \{0, 1\}^{N \times M}$ 可能为平方矩阵，即可能有 $N \neq M$ 。对二部图的封装成批过程中，`edge_index`中边的源节点与目标节点做的增值操作应是不同的。我们将二部图中两类节点的特征张量分别存储为`x_s`和`x_t`。

```

1 class BipartiteData(Data):
2     def __init__(self, edge_index, x_s, x_t):
3         super(BipartiteData, self).__init__()
4         self.edge_index = edge_index
5         self.x_s = x_s
6         self.x_t = x_t

```

为了对二部图实现正确的封装成批，我们需要告诉PyTorch Geometric，它应该在edge\_index中独立地为边的源节点和目标节点做增值操作。

```

1 def __inc__(self, key, value):
2     if key == 'edge_index':
3         return torch.tensor([[self.x_s.size(0)],
4                               [self.x_t.size(0)]])
5     else:
6         return super().__inc__(key, value)

```

其中，edge\_index[0]（边的源节点）根据x\_s.size(0)做增值运算，而edge\_index[1]（边的目标节点）根据x\_t.size(0)做增值运算。我们可以再次通过运行一个简单的测试脚本来测试我们的实现。

```

1 edge_index = torch.tensor([
2     [0, 0, 1, 1],
3     [0, 1, 1, 2],
4 ])
5 x_s = torch.randn(2, 16) # 2 nodes.
6 x_t = torch.randn(3, 16) # 3 nodes.
7
8 data = BipartiteData(edge_index, x_s, x_t)
9 data_list = [data, data]
10 loader = DataLoader(data_list, batch_size=2)
11 batch = next(iter(loader))
12
13 print(batch)
14 # Batch(edge_index=[2, 8], x_s=[4, 16], x_t=[6, 16])
15
16 print(batch.edge_index)
17 # tensor([[0, 0, 1, 1, 2, 2, 3, 3], [0, 1, 1, 2, 3, 4, 4, 5]])

```

可以看到我们得到我们期望的结果。



## 在新的维度上做拼接

有时，Data对象的属性需要在一个新的维度上做拼接（如经典的封装成批），例如，图级别属性或预测目标。具体来说，形状为[num\_features]的属性列表应该被返回为[num\_examples, num\_features]，而不是[num\_examples \* num\_features]。PyTorch Geometric通过在\_\_cat\_dim\_\_()中返回一个None的连接维度来实现这一点。

```
1 class MyData(Data):
2     def __cat_dim__(self, key, item):
3         if key == 'foo':
4             return None
5         else:
6             return super().__cat_dim__(key, item)
7
8 edge_index = torch.tensor([
9     [0, 1, 1, 2],
10    [1, 0, 2, 1],
11 ])
12 foo = torch.randn(16)
13
14 data = MyData(edge_index=edge_index, foo=foo)
15 data_list = [data, data]
16 loader = DataLoader(data_list, batch_size=2)
17 batch = next(iter(loader))
18
19 print(batch)
20 # Batch(edge_index=[2, 8], foo=[2, 16])
```

正如我们期望的，batch.foo现在由两个维度来表示，一个批维度，一个特征维度。

## 创建超大规模数据集类实践

**PCQM4M-LSC**是一个分子图的量子特性回归数据集，它包含了3,803,453个图。

注意以下代码依赖于ogb包，通过pip install ogb命令可安装此包。ogb文档可见于[Get Started | Open Graph Benchmark \(stanford.edu\)](#)。

我们定义的数据集类如下：

```
1 import os
2 import os.path as osp
3
4 import pandas as pd
5 import torch
6 from ogb.utils import smiles2graph
7 from ogb.utils.torch_util import replace_numpy_with_torchtensor
8 from ogb.utils.url import download_url, extract_zip
9 from rdkit import RDLogger
10 from torch_geometric.data import Data, Dataset
11 import shutil
12
13 RDLogger.DisableLog('rdApp.*')
14
15 class MyPCQM4MDataset(Dataset):
16
17     def __init__(self, root):
18         self.url = 'https://dgl-data.s3-
19 accelerate.amazonaws.com/dataset/OGB-LSC/pcqm4m_kddcup2021.zip'
20         super(MyPCQM4MDataset, self).__init__(root)
21
22         filepath = osp.join(root, 'raw/data.csv.gz')
23         data_df = pd.read_csv(filepath)
24         self.smiles_list = data_df['smiles']
25         self.homolumogap_list = data_df['homolumogap']
26
27     @property
28     def raw_file_names(self):
29         return 'data.csv.gz'
30
31     def download(self):
32         path = download_url(self.url, self.root)
33         extract_zip(path, self.root)
34         os.unlink(path)
35         shutil.move(osp.join(self.root,
36 'pcqm4m_kddcup2021/raw/data.csv.gz'), osp.join(self.root,
37 'raw/data.csv.gz'))
38
39     def len(self):
40         return len(self.smiles_list)
41
42     def get(self, idx):
```

```

40         smiles, homolumogap = self.smiles_list[idx],
self.homolumogap_list[idx]
41         graph = smiles2graph(smiles)
42         assert(len(graph['edge_feat']) ==
graph['edge_index'].shape[1])
43         assert(len(graph['node_feat']) == graph['num_nodes'])
44
45         x = torch.from_numpy(graph['node_feat']).to(torch.int64)
46         edge_index =
torch.from_numpy(graph['edge_index']).to(torch.int64)
47         edge_attr =
torch.from_numpy(graph['edge_feat']).to(torch.int64)
48         y = torch.Tensor([homolumogap])
49         num_nodes = int(graph['num_nodes'])
50         data = Data(x, edge_index, edge_attr, y,
num_nodes=num_nodes)
51         return data
52
53     # 获取数据集划分
54     def get_idx_split(self):
55         split_dict =
replace_numpy_with_torchtensor(torch.load(osp.join(self.root,
'pcqm4m_kddcup2021/split_dict.pt')))
56         return split_dict
57
58 if __name__ == "__main__":
59     dataset = MyPCQM4MDataset('dataset2')
60     from torch_geometric.data import DataLoader
61     from tqdm import tqdm
62     dataloader = DataLoader(dataset, batch_size=256, shuffle=True,
num_workers=4)
63     for batch in tqdm(dataloader):
64         pass
65

```

在生成一个该数据集类的对象时，程序首先会检查指定的文件夹下是否存在 data.csv.gz 文件，如果不在，则会执行 download 方法，这一过程是在运行 super 类的 \_\_init\_\_ 方法中发生的。然后程序继续执行 \_\_init\_\_ 方法的剩余部分，读取 data.csv.gz 文件，获取存储图信息的 smiles 格式的字符串，以及回归预测的目标 homolumogap。我们将由 smiles 格式的字符串转成图的过程在 get() 方法中实现，这样我们在生成一个 DataLoader 变量时，通过指定 num\_workers 可以实现并行执行生成多个图。

## 参考资料

---

- Dataset类官方文档: [torch\\_geometric.data.Dataset](#)
- 将图样本封装成批 (BATCHING) : [ADVANCED MINI-BATCHING](#)
- 分子图的量子特性回归数据集: [PCQM4M-LSC](#)
- [Get Started](#) | [Open Graph Benchmark \(stanford.edu\)](#)