

1 lines 0 Removals

124 lines + 123 Additions

1

```
1 import os
2 import os.path as osp
3 from typing import Callable, List, Optional
4
5 import torch
6
7 from torch_geometric.data import (
8     Data,
9     InMemoryDataset,
10     download_url,
11     extract_tar,
12     extract_zip,
13 )
14
15
16 class MalNetTiny(InMemoryDataset):
17     r"""The MalNet Tiny dataset from the
18     `A Large-Scale Database for Graph Representation Learning
19     <https://openreview.net/pdf?id=1xDTDk3XPW>`_ paper.
20     :class:`MalNetTiny` contains 5,000 malicious and benign
21     software function call graphs across 5 different types. Each graph
22     contains at most 5k nodes.
23
24     Args:
25         root (str): Root directory where the dataset should be saved.
26         split (str, optional): If :obj:`"train"`, loads the training dataset.
27             If :obj:`"val"`, loads the validation dataset.
28             If :obj:`"trainval"`, loads the training and validation dataset.
29             If :obj:`"test"`, loads the test dataset.
30             If :obj:`None`, loads the entire dataset.
31             (default: :obj:`None`)
32         transform (callable, optional): A function/transform that takes in an
33             :obj:`torch_geometric.data.Data` object and returns a transformed
34             version. The data object will be transformed before every access.
35             (default: :obj:`None`)
36         pre_transform (callable, optional): A function/transform that takes in
37             an :obj:`torch_geometric.data.Data` object and returns a
38             transformed version. The data object will be transformed before
39             being saved to disk. (default: :obj:`None`)
40         pre_filter (callable, optional): A function that takes in an
41             :obj:`torch_geometric.data.Data` object and returns a boolean
42             value, indicating whether the data object should be included in the
43             final dataset. (default: :obj:`None`)
44
45     """
46     data_url = ('http://malnet.cc.gatech.edu/'
47                 'graph-data/malnet-graphs-tiny.tar.gz')
```

```

47     split_url = 'http://malnet.cc.gatech.edu/split-info/s
plit_info_tiny.zip'
48     splits = ['train', 'val', 'test']
49
50     def __init__(
51         self,
52         root: str,
53         split: Optional[str] = None,
54         transform: Optional[Callable] = None,
55         pre_transform: Optional[Callable] = None,
56         pre_filter: Optional[Callable] = None,
57     ):
58         if split not in {'train', 'val', 'trainval', 'tes
t', None}:
59             raise ValueError(f'Split "{split}" found, but
expected either '
60                             f'"train", "val", "trainva
l", "test" or None')
61         super().__init__(root, transform, pre_transform,
pre_filter)
62         self.data, self.slices = torch.load(self.processe
d_paths[0])
63
64         if split is not None:
65             split_slices = torch.load(self.processed_path
s[1])
66             if split == 'train':
67                 self._indices = range(split_slices[0], sp
lit_slices[1])
68             elif split == 'val':
69                 self._indices = range(split_slices[1], sp
lit_slices[2])
70             elif split == 'trainval':
71                 self._indices = range(split_slices[0], sp
lit_slices[2])
72             elif split == 'test':
73                 self._indices = range(split_slices[2], sp
lit_slices[3])
74
75         @property
76         def raw_file_names(self) -> List[str]:
77             return ['malnet-graphs-tiny', osp.join('split_inf
o_tiny', 'type')]
78
79         @property
80         def processed_file_names(self) -> List[str]:
81             return ['data.pt', 'split_slices.pt']
82
83         def download(self):
84             path = download_url(self.data_url, self.raw_dir)
85             extract_tar(path, self.raw_dir)
86             os.unlink(path)
87
88             path = download_url(self.split_url, self.raw_dir)
89             extract_zip(path, self.raw_dir)
90             os.unlink(path)
91
92         def process(self):
93             y_map = {}
94             data_list = []
95             split_slices = [0]
96
97             for split in ['train', 'val', 'test']:
98                 with open(osp.join(self.raw_paths[1], f'{spli
t}.txt'), 'r') as f:
99                     filenames = f.read().split('\n')[:-1]

```

```

100         split_slices.append(split_slices[-1] + len(
101         n(filenamees))
102
103         for filename in filenamees:
104             path = osp.join(self.raw_paths[0], f'{filename}.edgelist')
105             malware_type = filename.split('/')[0]
106             y = y_map.setdefault(malware_type, len(y_map))
107
108             with open(path, 'r') as f:
109                 edges = f.read().split('\n')[5:-1]
110
111                 edge_index = [[int(s) for s in edge.split()] for edge in edges]
112                 edge_index = torch.tensor(edge_index).t().contiguous()
113                 num_nodes = int(edge_index.max()) + 1
114
115                 data = Data(edge_index=edge_index, y=y, num_nodes=num_nodes)
116                 data_list.append(data)
117
118                 if self.pre_filter is not None:
119                     data_list = [data for data in data_list if self.pre_filter(data)]
120
121                 if self.pre_transform is not None:
122                     data_list = [self.pre_transform(data) for data in data_list]
123
124                 torch.save(self.collate(data_list), self.processed_paths[0])
125
126                 torch.save(split_slices, self.processed_paths[1])

```