

40 lines - 37 Removals

```

1 import torch
2 from binaryLoader import BinaryDataset, pad_collate_func
3
4
5 def load_dataset(batch_size, max_seq_len=256, shuffle=True, num_workers=0):
6     train_set = BinaryDataset(good_dir='/data1/ember2018/train/benign/',
7                               bad_dir='/data1/ember2018/train/malware/',
8                               max_len=max_seq_len)
9
10    train_loader = torch.utils.data.DataLoader(train_set,
11                                                batch_size=batch_size,
12                                                shuffle=True,
13                                                num_workers=num_workers,
14                                                collate_fn=pad_collate_func,
15                                                drop_last=True)
16
17    test_set = BinaryDataset(good_dir='/data1/ember2018/test/benign/',
18                             bad_dir='/data1/ember2018/test/malware/',
19                             max_len=max_seq_len)
20
21    test_loader = torch.utils.data.DataLoader(test_set,
22                                              batch_size=batch_size,
23                                              shuffle=False,
24                                              num_workers=num_workers,
25                                              collate_fn=pad_collate_func,
26                                              drop_last=True)
27
28    return train_loader, test_loader
29

```

124 lines + 121 Additions

```

1 import datetime
2
3 import torch
4
5 from initializations import *
6 from binaryLoader import BinaryDataset, pad_collate_func
7
8
9 def printf(message):
10     timestamp = datetime.datetime.now()
11     timestamp_str = timestamp.strftime("%Y-%m-%d %H:%M:%S")
12     print(f"[{timestamp_str}] {message}")
13
14
15 def load_dataset(batch_size, max_seq_len=256, shuffle=True, num_workers=0):
16     if TRAIN:
17         train_set = BinaryDataset(
18             good_dir=TRAIN_DATA_PATH_GOOD,
19             bad_dir=TRAIN_DATA_PATH_BAD,
20             max_len=max_seq_len)
21
22     train_loader = torch.utils.data.DataLoader(
23         train_set,
24         batch_size=batch_size,
25         shuffle=True,
26         num_workers=num_workers,
27         collate_fn=pad_collate_func,
28         drop_last=True)
29
30     validation_set = BinaryDataset(
31         good_dir=VALIDATION_DATA_PATH_GOOD,
32         bad_dir=VALIDATION_DATA_PATH_BAD,
33         max_len=max_seq_len)
34
35     validation_loader = torch.utils.data.DataLoader(
36         validation_set,
37         batch_size=batch_size,
38         shuffle=True,
39         num_workers=num_workers,
40         collate_fn=pad_collate_func,
41         drop_last=True)
42
43     elif TEST:
44

```

```

30
31 if __name__ == '__main__':
32     train, test = load_dataset(batch_size=32, max_seq_len=20, shuffle=True, num_workers=0)
33
34     for x, y in train:
35         print(type(x))
36         print(x.shape)
37         print(y.shape)
38         print(x)
39         print(y)
40         break

```

```

44     # Testing: load test data set only
45     test_set = BinaryDataset(
46         good_dir=TEST_DATA_PATH_GOOD,
47         bad_dir=TEST_DATA_PATH_BAD,
48         max_len=max_seq_len)
49
50     test_loader = torch.utils.data.DataLoader(test_set,
51                                               batch_size=batch_size,
52                                               shuffle=shuffle,
53                                               num_workers=num_workers,
54                                               collate_fn=pad_collate_func,
55                                               drop_last=True)
56
57     return test_loader, len(test_set)
58
59 elif PREDICT:
60     test_set = BinaryDataset(
61         good_dir=TEST_DATA_PATH_GOOD + "temp",
62         bad_dir=TEST_DATA_PATH_BAD + "temp",
63         max_len=max_seq_len)
64
65     test_loader = torch.utils.data.DataLoader(test_set,
66                                               batch_size=batch_size,
67                                               shuffle=shuffle,
68                                               num_workers=num_workers,
69                                               collate_fn=pad_collate_func,
70                                               drop_last=True)
71
72     return test_loader, len(test_set)
73
74
75 try:
76     if TRAIN:
77         printf(f"Parameter selected: TRAIN")
78         train, validation, num_train_samples, num_validation_samples = load_dataset(batch_size=256, max_seq_len=256,
79
80         shuffle=True, num_workers=0)
81
82         printf(f"Number of training samples: {num_train_samples}")
83         printf(f"Number of validation samples: {num_validation_samples}\n")
84
85         printf("Printing training set data...")
86         for x, y in train:
87             printf(type(x))
88             printf(x.shape)
89             printf(y.shape)
90             printf(x)
91             printf(y)
92             break
93
94         print()

```

```

94         printf("Printing validation set data...")
95         for x, y in validation:
96             printf(type(x))
97             printf(x.shape)
98             printf(y.shape)
99             printf(x)
100            printf(y)
101            break
102
103     elif TEST:
104
105         printf(f"Parameter selected: TEST")
106         test, num_test_samples = load_dataset(batch
107         _size=256, max_seq_len=256, shuffle=True,
108         num_w
109         orkers=0)
110
111         printf("Printing test set data...")
112         for x, y in test:
113             printf(type(x))
114             printf(x.shape)
115             printf(y.shape)
116             printf(x)
117             printf(y)
118             break
119
120 except Exception as ex:
121     if ex.args[0] == 'num_samples should be a posit
122     ive integer value, but got num_samples=0':
123         printf("No file or folder found for one or
124         all of the specified dataset paths")
125         raise Exception(FileNotFoundError)
126     else:
127         raise ex

```