

Ariel ML Data Challenge

Fu-yun Wang¹ 191300051¹

¹ Nanjing University

摘要

简单介绍在 Ariel ML Data Challenge 的排名情况与代码设计。

1. 设备与计算资源

轻薄本 cpu。

2048 个训练样本。

1024 个测试样本。

训练时长 8 小时。

语言:python

深度学习框架: PyTorch.

2. 得分情况

如图1中所示, 我的个人队伍 ifoyooo 得分 9810, 排名 11 位。个人代码地址:

<https://github.com/ifoyooo/ml-network>

3. 代码设计

本次的代码设计主要参考了如下两个项目:

<https://github.com/ucl-exoplanets/ML-challenge-baseline>

<https://github.com/ifoyooo/pytorch-CycleGAN-and-pix2pix>

项目文件及各自的功能包括:

1. base_dataset.py 数据的抽象类。
2. challenge_dataset.py 本例中的数据类别具体实现, 主要实现了包括: `__getitem__` 函数, 用

Current Leader Board

Rank	Name	Score
1	Valis	9931
2	TheReturnOfBasel321	9916
3	DeepBlueAI	9909
4	major_tom	9899
5	Cantina_Band	9892
6	RozanskiT	9878
7	Christophe_pere	9872
8	Wabinab	9864
9	voidzero	9829
10	qcer	9822
11	ifoyooo	9810

图 1. 排名情况

于构建 pytorch 必备的 DataLoader; 数据与处理函数:simple_transform.

3. modules.py 神经网络的各种基本结构
4. networks.py 神经网络的具体实现
5. predict.py 利用训练好的模型预测测试集的结果。
6. score.py 返回模型在测试集上的得分。
7. test.py 测试训练好的模型在整个训练集上的 loss 与得分情况。
8. train.py 训练模型。

3.1. 数据处理与网络结构

受限于计算设备的限制, 仅仅在 2048 个样本上训练就用掉了 8 个小时, 所以实现的数据预处理与网络结构都非常简单, 并没有进行太多的尝试。在这两个方面做更多的尝试应该能较容易的有所提升。

经过对整个训练集读取分析可以得到该数据各个属性的分布的均值都越为 1, 标准差约为 0.04, 所以数据预处理设置为简单的标准化。如下所示

Listing 1. Transform

```
1 def simple_transform(x):
2     out = x.clone()
3     # out[:, :30]=1
4     # centering
5     out -= 1.
6     # rough rescaling
7     out /= 0.05
8     return out
```

神经网络结构主要为简单的卷积神经网络与全连接的组合。

Listing 2. NetWork

```
1 class SimpleConv(nn.Module):
2     def __init__(self):
3         super(SimpleConv, self).__init__()
4         self.conv1 = nn.Conv2d(1, 4, [16, 2],
5                                 [1, 1], padding=[4, 0])
6         self.pool1 = nn.MaxPool2d([2, 1], [2, 1])
7         self.conv21 = nn.Conv2d(4, 1, [16, 1],
8                                 [1, 1], padding=[4, 0])
9         self.pool2 = nn.MaxPool2d([2, 1], [1, 1])
10        self.mll = nn.Sequential(nn.Linear(4784, 1024),
11                                nn.ReLU(),
12                                nn.Linear(1024, 256),
13                                nn.ReLU(),
14                                nn.Linear(256, 55),
15                                )
16    def forward(self, input):
17        x = input.unsqueeze(1)
18        x = self.conv1(x)
```

```
x = self.pool1(x)
x = F.relu(x)
x = self.conv21(x)
x = self.pool2(x)
x = F.relu(x)
x = torch.flatten(x, start_dim=1)
x = self.mll(x)
return x
```

3.2. 训练与测试

为了便于调整训练参数, 代码中使用了 parser 库, 能够通过命令行直接控制模型的参数, 例如当需要使用已经保存过的参数继续训练时只需要执行

```
1 python train.py --continue_train=True
```

的命令就可以继续训练了。

Listing 3. Train

```
1 parser = argparse.ArgumentParser()
2 parser.add_argument("--lc_train_path",
3                     help="train path",
4                     type=str, default=pathlib.Path(__file__).
5                     parent.absolute()/
6                     "data/noisy_train/home/ucapats/Scratch
7                     /ml_data_challenge/training_set/noisy_train")
8 parser.add_argument("--lc_val_path",
9                     help="val path", type=str, default=
10                    pathlib.Path(__file__).parent.absolute()
11                    /"data/noisy_train/home/ucapats/
12                    Scratch/ml_data_challenge/
13                    training_set/noisy_train")
14 parser.add_argument("--params_train_path",
15                     help="params_train_path",
16                     type=str, default=pathlib.Path(__file__).
17                     parent.absolute()/ "data/params_train
18                     /home/ucapats/Scratch/ml_data_challenge
19                     /training_set/params_train")
20 parser.add_argument("--params_val_path",
21                     help="params_val_path",
```

```

22 type=str,default=pathlib.Path(__file__).
23 parent.absolute()/"data/params_train
24 /home/ucapats/Scratch/ml_data_challenge
25 /training_set/params_train")
26 parser.add_argument("--train_size",
27 type=int,default=2048)
28 parser.add_argument("--val_size",
29 type=int,default=1024)
30 parser.add_argument("--epochs",
31 type=int,default=10)
32 parser.add_argument("--save_from",
33 type=int,default=3)
34 parser.add_argument("--device",
35 type=str,default="cuda"
36 if torch.cuda.is_available() else "cpu")
37 # parser.add_argument("--device",
38 type=str,default="cpu")
39 parser.add_argument("--batch_size",
40 type=int,default=128)
41 parser.add_argument("--seed",
42 type=int,default=2048)
43 parser.add_argument("--input_dim",
44 type=int,default=55*300)
45 parser.add_argument("--output_dim",
46 type=int,default=55)
47 parser.add_argument("--model",
48 type=str,default="Conv2d")
49 parser.add_argument("--MLlinearlist",
50 type=list,default=[55*300,1024,256])
51 parser.add_argument("--continue_train",
52 type=bool,default=False)

```

此外，保存过的模型会保存在 outputs 文件夹下，根据模型的选择的不同，会保存在不同模型的子文件夹下，并保存训练过程的 loss 变化到 txt 文件并绘图。

最后当使用训练好的模型预测时，只需要执行

```
1 python predict.py
```

即可。由于一次性读取所有数据会导致内存溢出，所以我仍然使用 Dataloader 作为模型的输入，每次仅预测很少一部分的结果，然后追加保存在文件当中。

4. 如何测试我的代码

请将训练数据与测试数据放于 code 文件夹中，以训练集为例，其中的所有训练数据的地址应该为

```

1 code\data\noisy_train\home\ucapats\Scratch
2 \ml_data_challenge\training_set\noisy_train
3 \文件名

```

然后执行 main.py

致谢. 感谢老师的阅读。

参考文献

无