

实验四：Pytorch 实现玉米基因表达量预测模型

Hu Lab

2025 年 4 月 10 日

1. 实验背景及数据

1.1 实验背景

基因表达是生物学中的核心问题之一。基因在表达过程中，先转录成 mRNA，再翻译成蛋白质，这里基因转录出 mRNA 的量称为基因表达量。基因的转录受多种因素的影响，称为转录调控，包括顺式调控，反式调控及环境的影响等等。反式调控是指细胞中对基因表达有影响的蛋白质对基因的调控，顺式调控是指与基因在同一条 DNA 链上的 DNA 序列对基因的调控。顺式调控是影响基因表达的主要因素，也是我们的研究对象。这些调控序列 (称为调控元件) 根据其位置及功能的不同，可以分为启动子，增强子，沉默子等。启动子是转录起始位点 (TSS) 上游约 1kb 的 DNA 序列，其具有起始转录的能力，增强子可能位于基因上游远端，基因下游或者内含子区，对基因的表达起增强作用，沉默子对基因表达起抑制作用。另外，不同的调控元件有其特有的序列特征，这是 DNA 序列可以用来深度学习建模的基础。

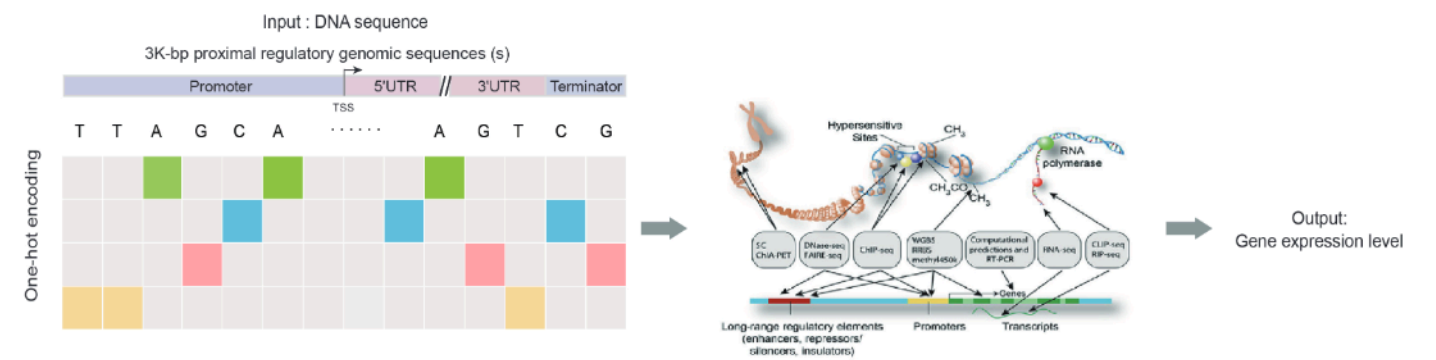


图 1: 实验背景

1.2 实验数据

本实验以玉米基因 TSS 附近上下游的四个调控区共 3kb 的 DNA 序列的 DNA 序列作经过 one-hot 编码为输入，one-hot 编码是将每种碱基 (有 A、T、C 和 G 四种) 映射成一个四维向量，其中一个维度是 1，剩余三个维度是 0。本实验中使用的标签基因表达量，是通过 RNA-seq 实验数据处理得到。output(标签) 为每个基因通过实验测得的表达量，单位为 TPM，代表着基因转录出 mRNA 的量。我们使用的基因数为 37979 个，本实验将数据划分为训练集和测试集，分别占 80% 和 20%，每一个样本对应一个基因。本实验拟通过搭建深度残差卷积网络，来实现对玉米基因表达量的预测。

2. 模型介绍

用 python 的 Pytorch 模块实现空洞残差卷积网络。网络结构包含 Convolution block、Convolution tower 和 Dilation residual block 及一个全连接层。此模型结构 (Basenji) 来源于文献：

<https://genome.cshlp.org/content/28/5/739.full?sid=3ecf5a07-edd6-4fe8-b183-10b7d1e269f4>

Geneid	sequence	TPM	dataset
Zm00001d027230	TGGAGGAAGCGAGCACATGTCCGGCCGG	34.621006	train
Zm00001d027235	TTTGCTCCTGCCATTCCCCAAAGCTTA	1.673876	valid
Zm00001d027250	TAATGTGTCCAAAAAACGGTTCCTCTA	4.509122	train
Zm00001d027253	AAATAGACACAAAAAGCAAGTCTACTTG	5.616972	train
Zm00001d027254	TCAGCTTCAGTCTAGTCTCTTCTTGTG	3.272217	valid
Zm00001d027256	TAGGAAATTAACCCACCGTCCCTGGTT	12.515516	valid
Zm00001d027258	CTTTAAAAAGTGGGGAATTTTATGCTA	291.908051	train
Zm00001d027265	GTTCTCCTTTGAGGGCGGAGGTAAGT	27.537571	test
Zm00001d027266	GCGGCTTCGAGCACTGCACCGGCAAGC	152.018448	valid
Zm00001d027267	TCATCATCTGCTGTCAAAATTTGGTGG	140.49733	train
Zm00001d027269	GAACCGGTATGGAATCTTAGCCGTTGG	0.882457	train
Zm00001d027273	TCTTACCGAATGGGTGGGACTCGTCTC	5.569419	train
Zm00001d027275	TCAGGTCCGTCATCACCTAAACCACTT	6.147862	train
Zm00001d027278	TTGGGTAAACATTTATATAAGTTGTAT	93.774559	train
Zm00001d027279	GAAATGTGAAGTTGTGACTAAATCATC	2.54835	train
Zm00001d027283	GATGATCCAGACTTACCTTTATGAGTG	196.878952	train
Zm00001d027285	TTATCTAAAAAGGTTAAAAAGTCTTAT	6.53042	test
Zm00001d027286	CGGTCGCGCCCTACTGGGCCTAGGAGC	10.948885	valid
Zm00001d027289	TTGAATGCCAGGTCATCAAAAAAAGT	2.951355	test
Zm00001d027290	GCAAGTGTACATTTCGTCAATTTTAA	3768.132568	train
Zm00001d027292	CTTCACATCATCAACTTCATAATAATT	21.867928	test
Zm00001d027295	TTCCAGAAGACATAAAGTCTGTGTTAA	2193.95166	train
Zm00001d027304	CATCGACGACGTGGGAGGAGGGTGGG	100.480644	test
Zm00001d027307	GCACTGGACAGTGTCCGGTGGGCCAAC	73.999771	train
Zm00001d027311	ACTTGCAATAAGCATGTCTTTTAAATT	9.518116	train
Zm00001d027312	CTGATATACGAACCCCGTGTTCAGCA	20.452984	train
Zm00001d027314	TACGACCCAAACACCGCACAAAGCTCG	187.788651	train
Zm00001d027317	TCATACACGCATGACACACATGTGCTT	191.618698	test
Zm00001d027319	TGTCCACACATCAAATAATCTGGCAGA	52.433243	train
Zm00001d027321	TAGCCTCCAATCATCCCTCAACTCTAT	163.306503	train
Zm00001d027322	TGTTTGGAGCTGATGGGAATCGTCGTA	96.274788	train
Zm00001d027323	GTGAGACCTCTCTTGGCCATTCTTGGC	382.299866	train
Zm00001d027324	AGTGGGGCGCGAAAAATTCATTTGATGA	89.222183	train
Zm00001d027325	GAGCTGCCGCTCTCCACATGCCCGCC	37.097458	train
Zm00001d027329	AGCGAAAGAGGGCAAGCAAGCAGAGAA	91.757378	train
Zm00001d027330	TCATCTTTTCTTAGAAGCGCTATTCA	260.086792	train

图 2: 数据集示意图

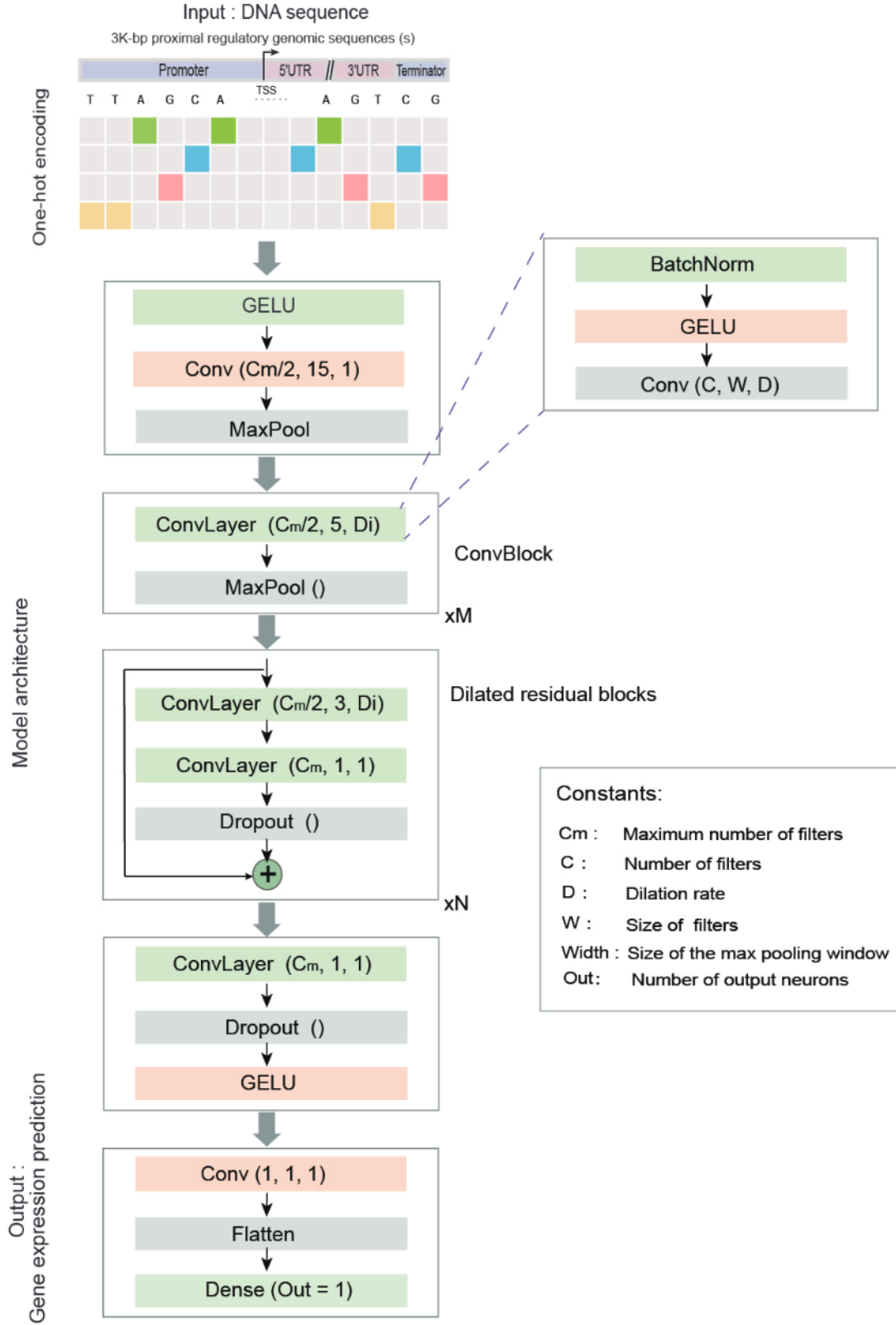


图 3: 网络结构示意图

3. 实验代码

3.1 数据预处理代码

```
1 # %%
2 import os
3 import pandas as pd
4 import numpy as np
5 from sklearn.model_selection import train_test_split
6 # %%
7 def one_hot_encode_along_channel_axis(sequence):
8     to_return = np.zeros((4, len(sequence)), dtype=np.int8)
9     seq_to_one_hot_fill_in_array(zeros_array=to_return, sequence=sequence,
10                                   one_hot_axis=0)
11
12     return to_return
13
14 def seq_to_one_hot_fill_in_array(zeros_array, sequence, one_hot_axis):
15     assert one_hot_axis==0 or one_hot_axis==1
16     if (one_hot_axis==0):
17         assert zeros_array.shape[1] == len(sequence)
18     elif (one_hot_axis==1):
19         assert zeros_array.shape[0] == len(sequence)
20     #will mutate zeros_array
21     for (i, char) in enumerate(sequence):
22         if (char=="A" or char=="a"):
23             char_idx = 0
24         elif (char=="C" or char=="c"):
25             char_idx = 1
26         elif (char=="G" or char=="g"):
27             char_idx = 2
28         elif (char=="T" or char=="t"):
29             char_idx = 3
30         elif (char=="N" or char=="n"):
31             continue #leave that pos as all 0's
32         else:
33             raise RuntimeError("Unsupported character: "+str(char))
34     if (one_hot_axis==0):
35         zeros_array[char_idx, i] = 1
36     elif (one_hot_axis==1):
37         zeros_array[i, char_idx] = 1
38 # %%
39 # X = range(len(y)) # 37979
40 # x_train, x_test, y_train, y_test = train_test_split(X,
```

```

40 # y,
41 # train_size=0.8,
42 # random_state=2023)
43 # x_train, x_valid, y_train, y_valid = train_test_split(x_train,
44 # y_train,
45 # train_size=0.9,
46 # random_state=2023)
47 # x_train:27344, x_valid:3039, x_test:7596
48 ###
49 #为了统一，这里提前为大家划分好了训练集，验证集和测试集，见df['dataset']列
50 #导入划分好的数据集，并对DNA序列one-hot编码处理
51 df = pd.read_excel('/mnt/cgshare/dataset.xlsx')
52 print('genes number:', df.shape[0])
53 df_train = df[df['dataset']=='train']
54 y_train = np.log2(df_train['TPM'].values + 1)
55 train_data = np.array([one_hot_encode_along_channel_axis(i.strip()) for i in df_train
    ['sequence'].values])
56 df_valid = df[df['dataset']=='valid']
57 y_valid = np.log2(df_valid['TPM'].values + 1)
58 valid_data = np.array([one_hot_encode_along_channel_axis(i.strip()) for i in df_valid
    ['sequence'].values])
59 df_test = df[df['dataset']=='test']
60 y_test = np.log2(df_test['TPM'].values + 1)
61 test_data = np.array([one_hot_encode_along_channel_axis(i.strip()) for i in df_test[
    'sequence'].values])

```

3.2 搭建模型及训练、测试代码

```

1 # %%
2 import numpy as np
3 import os
4 import torch
5 import torch.nn as nn
6 from torch import Tensor
7 import torch.optim as optim
8 from torch.utils.data import Dataset, DataLoader, ConcatDataset
9 import torch.nn.functional as F
10 #from torch.cuda.amp.grad_scaler import GradScaler
11 #from torch.cuda.amp import autocast
12 #from torchsummary import summary
13 from tqdm import tqdm
14 from scipy.stats import pearsonr
15 import matplotlib.pyplot as plt
16 from einops.layers.torch import Rearrange

```

```

17 from torchsummary import summary
18
19 # define blocks
20 class upd_GELU(nn.Module):
21     def __init__(self):
22         super(upd_GELU, self).__init__()
23         self.constant_param = nn.Parameter(torch.Tensor([1.702]))
24         self.sig = nn.Sigmoid()
25
26     def forward(self, input: Tensor) -> Tensor:
27         outval = torch.mul(self.sig(torch.mul(self.constant_param, input)), input)
28         return outval
29
30 # 由于torch的padding操作没有padding='same',因此重新定义了maxpooling类
31 class KerasMaxPool1d(nn.Module):
32     def __init__(
33         self,
34         pool_size=2,
35         padding="valid",
36         dilation=1,
37         return_indices=False,
38         ceil_mode=False,
39     ):
40         """Hard coded to only be compatible with kernel size and stride of 2."""
41         super().__init__()
42         self.padding = padding
43         _padding = 0
44         if pool_size != 2:
45             raise NotImplementedError("MaxPool1D with kernel size other than 2.")
46         self.pool = nn.MaxPool1d(
47             kernel_size=pool_size,
48             padding=_padding,
49             dilation=dilation,
50             return_indices=return_indices,
51             ceil_mode=ceil_mode,
52         )
53
54     def forward(self, x):
55         # (b c h)
56         if self.padding == "same" and x.shape[-1] % 2 == 1:
57             x = F.pad(x, (0, 1), value=float("inf"))
58         return self.pool(x)
59

```

```

60 class Residual(nn.Module):
61     def __init__(self, module):
62         super().__init__()
63         self.module = module
64
65     def forward(self, x):
66         return self.module(x) + x
67
68 class ConvBlock(nn.Module):
69     def __init__(self, in_channels, filters, kernel_size=1, padding='same',
70                 stride=1, dilation_rate=1, pool_size=1, dropout=0, bn_momentum
71                 =0.1):
72         super().__init__()
73         block = nn.ModuleList()
74         block.append(upd_GELU())
75         block.append(nn.Conv1d(in_channels=in_channels,
76                               out_channels=filters,
77                               kernel_size=kernel_size,
78                               stride=stride,
79                               padding=padding,
80                               dilation=int(round(dilation_rate)),
81                               bias=False))
82         block.append(nn.BatchNorm1d(filters, momentum=bn_momentum, affine=True))
83         if dropout > 0:
84             block.append(nn.Dropout(p=dropout))
85         if pool_size > 1:
86             block.append(KerasMaxPool1d(pool_size=pool_size, padding=padding))
87         self.block = nn.Sequential(*block)
88         self.out_channels = filters
89
90     def forward(self, x):
91         return self.block(x)
92
93 class ConvTower(nn.Module):
94     def __init__(
95         self,
96         in_channels,
97         filters_init,
98         filters_end=None,
99         filters_mult=None,
100         divisible_by=1,
101         repeat=2,
102         **kwargs,

```



```

102 ):
103     super().__init__()
104
105     def _round(x):
106         return int(np.round(x / divisible_by) * divisible_by)
107
108     # determine multiplier
109     if filters_mult is None:
110         assert filters_end is not None
111         filters_mult = np.exp(np.log(filters_end / filters_init) / (repeat - 1))
112
113     rep_filters = filters_init
114     in_channels = in_channels
115     tower = nn.ModuleList()
116     for _ in range(repeat):
117         tower.append(
118             ConvBlock(
119                 in_channels=in_channels, filters=_round(rep_filters), **kwargs
120             )
121         )
122         in_channels = _round(rep_filters)
123         rep_filters *= filters_mult
124
125     self.tower = nn.Sequential(*tower)
126     self.out_channels = in_channels
127
128     def forward(self, x):
129         return self.tower(x)
130
131 class DilatedResidual(nn.Module):
132     def __init__(
133         self,
134         in_channels,
135         filters,
136         kernel_size=3,
137         rate_mult=2,
138         dropout=0,
139         repeat=1,
140         **kwargs,
141     ):
142         super().__init__()
143         dilation_rate = 1 # 初始化为1, 后面累乘
144         in_channels = in_channels

```

```

145     block = nn.ModuleList()
146     for _ in range(repeat):
147         inner_block = nn.ModuleList()
148
149         inner_block.append(
150             ConvBlock(
151                 in_channels=in_channels,
152                 filters=filters,
153                 kernel_size=kernel_size,
154                 dilation_rate=int(np.round(dilation_rate)),
155                 **kwargs,
156             )
157         )
158
159         inner_block.append(
160             ConvBlock(
161                 in_channels=filters,
162                 filters=in_channels,
163                 dropout=dropout,
164                 **kwargs,
165             )
166         )
167
168         block.append(Residual(nn.Sequential(*inner_block)))
169
170         dilation_rate *= rate_mult
171         dilation_rate = np.round(dilation_rate)
172     self.block = nn.Sequential(*block)
173     self.out_channels = in_channels
174
175     def forward(self, x):
176         return self.block(x)
177
178 class BasenjiFinal(nn.Module):
179     def __init__(
180         self, in_features, units=1, activation='linear', **kwargs):
181         super().__init__()
182         block = nn.ModuleList()
183         #if flatten:
184         block.append(Rearrange('b ... -> b (...)'))
185
186         # Rearrange('b ... -> b (...)') 代表将batch size维度保留，剩下的维度变成一维
187         # 这里相当于flatten

```

```

188     # block.append(nn.Conv1d(in_features, units, kernel_size=1, stride=1, padding
189                               =0))
190
191     block.append(nn.Linear(in_features=in_features, out_features=units))
192     self.block = nn.Sequential(*block)
193
194     def forward(self, x):
195         return self.block(x)
196
197     class BasenjiModel(nn.Module):
198         def __init__(self, conv1_filters=8, conv1_ks=15,
199             # 第一个 conv block 参数
200             conv1_pad=7, conv1_pool=2, conv1_pdrop=0.4, conv1_bn_momentum=0.1,
201             # conv tower 参数
202             convt_filters_init=16, filters_end=32, convt_repeat=2, convt_ks=5,
203             convt_pool=2,
204             # dil res block 参数
205             dil_in_channels=32, dil_filters=16, dil_ks=3, rate_mult=2, dil_pdrop=0.3,
206             dil_repeat=2,
207             conv2_in_channels=32, conv2_filters=32, # 第二个 conv block 参数
208             conv3_in_channels=32, conv3_filters=1, # 第三个 conv block (1*1 conv) 参数
209             final_in_features=int(3000/(2**3)) # final block 参数, 2**3表示经过3次
210             pooling
211             ):
212         super().__init__()
213         block = nn.ModuleList()
214         block.append(ConvBlock(in_channels=4,
215             filters=conv1_filters,
216             kernel_size=conv1_ks,
217             padding=conv1_pad,
218             pool_size=conv1_pool,
219             dropout=conv1_pdrop,
220             bn_momentum=conv1_bn_momentum))
221
222         block.append(ConvTower(in_channels=conv1_filters,
223             filters_init=convt_filters_init,
224             # filters_mult=convt_filters_mult,
225             filters_end=filters_end,
226             divisible_by=1,
227             repeat=convt_repeat,
228             kernel_size=convt_ks,
229             pool_size=convt_pool,
230             ))

```

```

227
228     block.append(DilatedResidual(
229         in_channels=dil_in_channels,
230         filters=dil_filters,
231         kernel_size=dil_ks,
232         rate_mult=rate_mult,
233         dropout=dil_pdrop,
234         repeat=dil_repeat,))
235
236     block.append(ConvBlock(in_channels=conv2_in_channels,
237         filters=conv2_filters,
238         kernel_size=1))
239     block.append(ConvBlock(in_channels=conv3_in_channels,
240         filters=conv3_filters,
241         kernel_size=1))
242     block.append(BasnjiFinal(final_in_features))
243     self.block = nn.Sequential(*block)
244
245     def forward(self, x):
246         return self.block(x)
247
248     # %%
249     # define dataset class
250     class MyDataset(Dataset):
251         def __init__(self, input, label):
252             inputs = torch.tensor(input, dtype=torch.int8)
253             labels = torch.tensor(label, dtype=torch.float32)
254             self.inputs = inputs
255             self.labels = labels
256
257         def __getitem__(self, index):
258             return self.inputs[index], self.labels[index]
259
260         def __len__(self):
261             return len(self.labels)
262
263     # %%
264     BATCH_SIZE = 32
265     EPOCHS = 10
266     lr=1e-3
267
268     # valid data可以用于早停，这里没有加入早停机制，你们可以自己尝试加入
269     train_dataset = MyDataset(train_data, y_train)

```

```

269 trainDataLoader = torch.utils.data.DataLoader(dataset = train_dataset , batch_size =
    BATCH_SIZE, shuffle = True)
270 valid_dataset = MyDataset(valid_data , y_valid)
271 validDataLoader = torch.utils.data.DataLoader(dataset = valid_dataset , batch_size =
    BATCH_SIZE)
272 test_dataset = MyDataset(test_data , y_test)
273 testDataLoader = torch.utils.data.DataLoader(dataset = test_dataset , batch_size =
    BATCH_SIZE)
274 # %%
275 device = "cuda:0" if torch.cuda.is_available() else "cpu"
276
277 net = BasenjiModel()
278 summary(net, input_size=[(4, 3000)], batch_size=BATCH_SIZE, device="cpu")
279 optimizer = torch.optim.Adam(params= net.parameters() , lr=lr)
280 lossF = nn.MSELoss()
281 print(net.to(device))
282 #存储训练过程
283 history = { 'Valid Loss':[] , 'Valid pcc':[] }
284 for epoch in range(1,EPOCHS + 1):
285     progressBar = tqdm(trainDataLoader , unit = 'step')
286     net.train(True)
287     epoch_loss_all = 0
288     for step, (inputs, labels) in enumerate(progressBar):
289         #将序列和标签传输进 device 中
290         inputs = inputs.to(device)
291         labels = labels.to(device)
292         #清空模型的梯度
293         net.zero_grad()
294         #对模型进行前向推理
295         outputs = net(inputs)
296         #计算本轮推理的 Loss 值
297         loss = lossF(outputs.reshape(labels.shape) , labels)
298         epoch_loss_all += loss*len(labels)
299         #进行反向传播求出模型参数的梯度
300         loss.backward()
301         #使用迭代器更新模型权重
302         optimizer.step()
303         #将本 step 结果进行可视化处理
304         progressBar.set_description( "[%d/%d] Loss: %.4f" %
305                                     (epoch, EPOCHS, loss.item()) )
306     if step == len(progressBar)-1:
307         valid_totalLoss = 0
308         y_valid_pred_all = []

```

```

309     y_valid_true_all = []
310     train_loss = epoch_loss_all / len(train_data)
311     #关闭模型的训练状态
312     net.train(False)
313     with torch.no_grad():
314         #对验证集的 DataLoader 进行迭代
315         for x_valid, y_valid in validDataLoader:
316             x_valid = x_valid.to(device)
317             y_valid = y_valid.to(device)
318             y_valid_pred = net(x_valid)
319             y_valid_pred_all.extend(y_valid_pred.flatten().tolist())
320             y_valid_true_all.extend(y_valid.tolist())
321             valid_batch_loss = lossF(y_valid_pred.reshape(y_valid.shape),
322                                     y_valid)
323             valid_totalLoss += valid_batch_loss * len(y_valid)
324         validLoss = valid_totalLoss/len(valid_data)
325         pcc_valid, _ = pearsonr(y_valid_pred_all, y_valid_true_all)
326         history['Valid Loss'].append(validLoss.item())
327         history['Valid pcc'].append(pcc_valid.item())
328         progressBar.set_description("[%d/%d] Train Loss: %.4f, Valid Loss: %.4f, Valid pcc: %.4f" %
329                                     (epoch, EPOCHS, train_loss.item(), validLoss.item(), pcc_valid))
330
331     progressBar.close()
332
333 # %%
334 # calculate test pcc
335 # calculate test pcc and R2
336 net.train(False)
337 y_test_pred_all = []
338 y_test_true_all = []
339 with torch.no_grad():
340     #对测试集的 DataLoader 进行迭代
341     for x_test, y_test in testDataLoader:
342         x_test = x_test.to(device)
343         y_test = y_test.to(device)
344         y_test_pred = net(x_test)
345         y_test_pred_all.extend(y_test_pred.flatten().tolist())
346         y_test_true_all.extend(y_test.tolist())
347 # 计算测试集上的 pcc 和 R2，并画测试集上真实值和预测值之间的散点图
348 pcc, _ = pearsonr(y_test_pred_all, y_test_true_all)
349 print('pcc: ', pcc)

```

```
349 plt.xlabel('y test predict ')
350 plt.ylabel('y test true ')
351 plt.scatter(y_test_pred_all, y_test_true_all, s=0.1)
```