

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.preprocessing import LabelEncoder
4 import torch
5 from torch.utils.data import Dataset, DataLoader
6 import torch
7 import torch.nn as nn
8 import torch.nn.functional as F

1 df = pd.read_csv('data_cleaned.csv')
2 df['platform_order_time_date']=pd.to_datetime(df['platform_order_time_date'])
3 df=df.sort_values(by='platform_order_time_date')

1 time_interval = '5T' # 'T' 代表分钟df['platform_order_time_date']
2 df['order_time_interval']=df['platform_order_time_date'].dt.floor(time_interval)

1 all_intervals = pd.date_range(start=df['order_time_interval'].min().floor('5T'),
2                               end=df['order_time_interval'].max().floor('5T'),
3                               freq='5T')
4 all_h3 = df['H3_Index'].unique()
5
6 # 创建所有可能的组合
7 full_index = pd.MultiIndex.from_product(
8     [all_intervals, all_h3],
9     names=['order_time_interval', 'H3_Index']
10 )
11 print(len(all_intervals))
12 print(len(all_h3))
13 print(len(full_index ))

↔ 2309
   1560
   3602040

1 grouped=df.groupby(['order_time_interval','H3_Index']).size().agg('count')
2 grouped = grouped.reindex(full_index , fill_value=0)
3 df = grouped.reset_index(name='count')

1 le = LabelEncoder()
2 df['H3_Index_encoded'] = le.fit_transform(df['H3_Index'])

1 class H3TimeSeriesDataset(Dataset):
2     def __init__(self, sequences, labels):
3         """
4         sequences: list of tuples (seq_counts, seq_h3)
5         labels: list of tuples (target_count, target_h3)
6         """
7         self.sequences = sequences
8         self.labels = labels
9
10    def __len__(self):
11        return len(self.sequences)
12
13    def __getitem__(self, idx):
14        seq_counts, seq_h3 = self.sequences[idx]
15        target_count= self.labels[idx]
16        # 转换为张量
17        seq_counts = torch.tensor(seq_counts, dtype=torch.float32).unsqueeze(-1) # (T, 1)
18        seq_h3 = torch.tensor(seq_h3, dtype=torch.long).unsqueeze(-1) # (T, 1)
19        # 合并特征, 形状 (T, 2)
20        input_seq = torch.cat((seq_counts, seq_h3), dim=1)
21        # 目标
22        target_count = torch.tensor(target_count, dtype=torch.float32) # 泊松目标为float
23        #target_h3 = torch.tensor(target_h3, dtype=torch.float32) # 二分类标签
24        return input_seq, target_count

1 def create_sequences(df, T):
2     sequences = []
3     labels = []
4     # 按 H3_Index_encoded 分组
5     for h3, group in df.groupby('H3_Index_encoded'):
6         group = group.sort_values('order_time_interval')

```

```

7         counts = group['count'].values
8         h3_encoded = group['H3_Index_encoded'].values
9         # 创建长度为 T+1 的序列
10        for i in range(len(counts) - T-1):
11            seq_counts = counts[i:i+T]
12            seq_h3 = h3_encoded[i:i+T]
13            target_counts = counts[i+1:i+T+1]
14            #target_h3 = h3_encoded[i+1:i+T+1]
15            sequences.append((seq_counts, seq_h3))
16            labels.append(target_counts)
17        return sequences, labels

1 T = 10 # 序列长度
2 sequences, labels = create_sequences(df, T)
3 dataset = H3TimeSeriesDataset(sequences, labels)
4 batch_size = 512
5 dataloader = DataLoader(dataset, batch_size=batch_size, shuffle=True)

1
2 class H3LSTMModel(nn.Module):
3     def __init__(self, h3_vocab_size, h3_embedding_dim, hidden_dim, num_layers, max_count):
4         super(H3LSTMModel, self).__init__()
5         self.h3_embedding = nn.Embedding(num_embeddings=h3_vocab_size, embedding_dim=h3_embedding_dim)
6         self.lstm = nn.LSTM(
7             input_size=1 + h3_embedding_dim,
8             hidden_size=hidden_dim,
9             num_layers=num_layers,
10            batch_first=True
11        )
12        self.fc_binary = nn.Linear(hidden_dim, 1) # 输出二分类概率
13        self.fc_count = nn.Linear(hidden_dim, max_count + 1) # 输出Softmax概率分布
14
15    def forward(self, x_counts, x_h3):
16        """
17        x_counts: (batch_size, T, 1)
18        x_h3: (batch_size, T, 1)
19        """
20        # 嵌入 H3_Index
21        x_h3 = x_h3.squeeze(-1) # (batch_size, T)
22        embedded_h3 = self.h3_embedding(x_h3) # (batch_size, T, h3_embedding_dim)
23
24        # 合并 count 和嵌入后的 H3_Index
25        x = torch.cat((x_counts, embedded_h3), dim=2) # (batch_size, T, 1 + h3_embedding_dim)
26
27        # LSTM
28        lstm_out, _ = self.lstm(x) # lstm_out: (batch_size, T, hidden_dim)
29
30        # 计算二分类概率
31        binary_logits = self.fc_binary(lstm_out).squeeze(-1) # (batch_size, T)
32        binary_probs = torch.sigmoid(binary_logits) # (batch_size, T)
33
34        # 计算Softmax概率分布
35        count_logits = self.fc_count(lstm_out) # (batch_size, T, max_count +1)
36        count_probs = F.softmax(count_logits, dim=-1) # (batch_size, T, max_count +1)
37
38        return binary_probs, count_probs
39

1 import torch
2 from torch.distributions import Poisson
3
4 def poisson_log_prob( lambda_, k):
5     """
6     计算泊松分布的对数概率。
7
8     Args:
9         k (torch.Tensor): 事件发生次数，形状可以是任意的非负整数张量。
10        lambda_ (torch.Tensor): 平均发生次数，形状与 k 相同，且所有元素均为正数。
11
12    Returns:
13        torch.Tensor: 对数概率，形状与 k 相同。
14    """
15    # 创建泊松分布对象
16    poisson_dist = Poisson(rate=lambda_)
17
18    # 计算对数概率

```

```

19         log_prob = poisson_dist.log_prob(k)
20
21     return log_prob

1 device='cpu'
2 # model
3 model=H3LSTMModel(h3_vocab_size=len(all_h3), h3_embedding_dim=3, hidden_dim=32, num_layers=1, max_count=df['count'].max()).to(device)
4 # 二分类损失
5 criterion_binary = nn.BCELoss()
6
7 # 分类交叉熵损失
8 criterion_ce = nn.CrossEntropyLoss()
9
10 #mse loss
11 criterion_mse = nn.MSELoss()
12 # 优化器
13 optimizer = torch.optim.Adam(model.parameters(), lr=0.001)
14

```

```

1 import torch
2 print(torch.cuda.is_available())

```

False

```

1 def softargmax(count_probs, max_count):
2     """
3     计算 Softargmax, 返回期望值作为连续预测。
4
5     Args:
6         count_probs (torch.Tensor): Softmax 输出概率, 形状为 (batch_size, T, max_count + 1)。
7         max_count (int): 最大计数 k 值。
8
9     Returns:
10         torch.Tensor: 预测的连续计数, 形状为 (batch_size, T)。
11     """
12     # 创建一个包含所有计数类别的张量 k = [0, 1, 2, ..., max_count]
13     k = torch.arange(0, max_count + 1, dtype=count_probs.dtype, device=count_probs.device).view(1, 1, -1)
14
15     # 计算 Softargmax: sum_k P(k) * k
16     soft_argmax = (count_probs * k).sum(dim=-1)
17
18     return soft_argmax

1 # 训练模型
2 num_epochs = 100
3 max_count = df['count'].max() # 从0到8
4
5 for epoch in range(num_epochs):
6     model.train()
7     epoch_loss = 0
8     i=0
9     for batch_inputs, batch_targets in dataloader:
10         # 分离输入特征
11         batch_counts = batch_inputs[:, :, 0].unsqueeze(-1).to(device) # (batch_size, T, 1)
12         batch_h3 = batch_inputs[:, :, 1].long().unsqueeze(-1).to(device) # (batch_size, T, 1)
13
14         # 目标
15         y = batch_targets # y_count: (batch_size,), y_binary: (batch_size,)
16         y_count = y.to(device) # (batch_size, T)
17         y_binary = ((y==1)*1.0).to(device) # (batch_size, T)
18         #print( y_binary.shape)
19         # 前向传播
20         optimizer.zero_grad()
21         binary_probs, count_probs = model(batch_counts, batch_h3) # binary_probs: (batch_size, T), count_probs: (batch_size, T, max_count+1)
22
23         # 计算二分类损失
24         loss_binary = criterion_binary(binary_probs, y_binary)
25
26         soft_argmax=softargmax(count_probs, max_count)
27
28
29         ##计算mse loss
30         loss_mse=criterion_mse(soft_argmax,y_count)
31         #print(soft_argmax.shape)
32         # 计算分类交叉熵损失

```

```
33     # CrossEntropyLoss expects input of shape (N, C) and target of shape (N)
34     # 这里将每个时间步视为一个独立的样本
35     #count_probs_reshaped = count_probs.view(-1, max_count +1)    # (batch_size * T, max_count +1)
36     #y_count_reshaped = y_count.view(-1)    # (batch_size * T,)
37     #loss_ce = criterion_ce(count_probs_reshaped, y_count_reshaped)
38
39     # 计算KL散度损失
40     log_possion=-poisson_log_prob(soft_argmax,y_count).mean()
41     #kl_div = kl_divergence_softmax_poisson(count_probs, y_count, max_count)    # (batch_size, T)
42     #loss_kl = kl_div.mean()
43
44     # 总损失
45     loss = loss_binary + loss_mse + log_possion
46     i+=1
47     # 反向传播和优化
48
49     loss.backward()
50     optimizer.step()
51
52     epoch_loss += loss.item()
53     if i>=1000:
54         break
55     # 每10个epoch打印一次
56     #if (epoch + 1) % 10 == 0:
57     avg_loss = epoch_loss / len(dataloader)
58     print(f"Epoch {epoch}, Loss: {avg_loss:.4f}")
59
```

↻ Epoch [0], Loss: 3.9421
Epoch [1], Loss: 0.1952
Epoch [2], Loss: 0.1600

1 len(dataloader)

↻ 14004