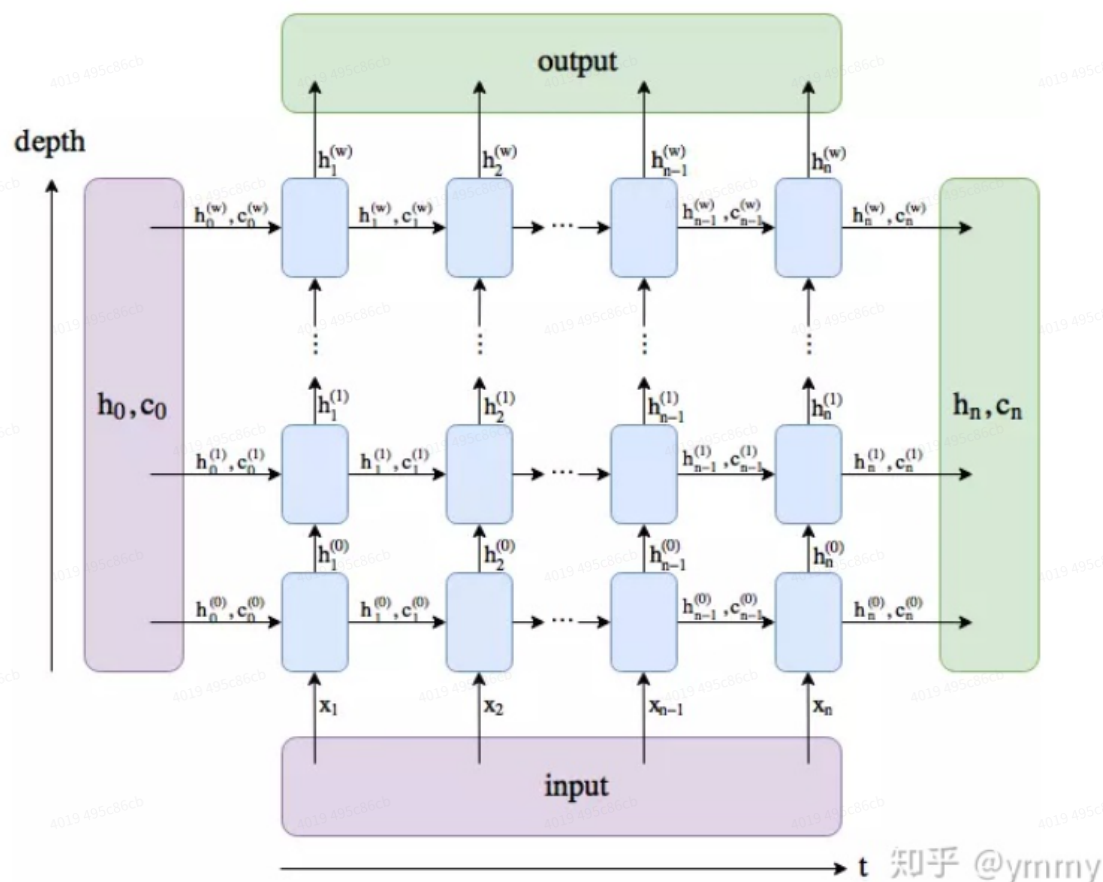


4.LSTM实操

1.LSTM网络结构与pytorch

LSTM中将整个网络看成这样：



蓝色的模块是指前面的单位LSTM结构。横向连接各个LSTM结构就是序列化。不同状态（时间）下的LSTM结构。纵向连接的是同一状态下LSTM结构。或者我们可以简单看作是多LSTM细胞结构串联，其个数为num_layers, 将其视为一个新的LSTM结构，并且按照时序连接起来。图中每个参数都是自带维度的。

Pytorch中的库：

nn.LSTM(input_size, hidden_size, num_layers)

input_size: 输入的X的维度。

hidden_size: 输出和输入的 h_i 的维度。

num_layers: depth, 即串联的LSTM个数。

2.时间序列实战

2.1. 数据

数据集，是从1949-1960,12年，每个月的乘客数量。即一共有144条数据，表示了144个月的乘客数量。

	year	month	passengers
0	1949	January	112
1	1949	February	118
2	1949	March	132
3	1949	April	129
4	1949	May	121
139	1960	August	606
140	1960	September	508
141	1960	October	461
142	1960	November	390
143	1960	December	432

2.2 思路

采用的肯定是利用前几期的数据来预测当前期的数据。(具体的方法是：利用前N（N= 3）期的数据为输入，当前期的数据为标签计算误差)。个人觉得可以N是一个超参，可以慢慢调。

由于时间序列数据每个时期的数据样本只有一个，那么X为(time_step, 1, input_size), Y为(time_step, 1, output_size)。

构建一个LSTM网络，输入的数据就是X，hidden_size可任意取，num_layers也可视情况取，即用多少层LSTM串联(同一时期内)。

最后用一层线性层Linear(hidden_size, output_size)进行输出。比较输出与Y的误差。不断迭代对参数进行优化。

2.3代码

Python

```
1 #step 1. 加载飞行数据
2 flight_data = pd.read_csv('flights.csv')
3 # 数据归一化
4 maxPassenger = flight_data['passengers'].max()
5 minPassenger = flight_data['passengers'].min()
6 flight_data['passengers'] = (flight_data['passengers'] - minPassenger) \
7     / (maxPassenger - minPassenger)
8
9 dataset = flight_data['passengers'].values.tolist()
10
11 #step 2. 划分数据集
12 # 数据集目标函数值赋值, 其中dataset为数据, look_back为以几行数据为特征数目
13 # look_back表示3期回头, 即使用前三期的数据预测下一期
14 # 用前3期数据预测下1期
15 def createDataset(dataset, look_back):
16     dataX = []
17     dataY = []
18     for i in range(len(dataset)-look_back):
19         dataX.append(dataset[i:i+look_back])
20         dataY.append(dataset[i+look_back])
21
22     dataX = torch.tensor(dataX)
23     dataX = dataX.reshape(-1, 1, look_back)
24     dataY = torch.tensor(dataY)
25     dataY = dataY.reshape(-1, 1, 1)
26     return dataX, dataY
27
28 data = createDataset(dataset=dataset, look_back=3) # 划分数据集, 3个月为一组
29
30 # step3. 划分训练集和测试集
31 # 由于是时间序列数据, 不适合这样随机打乱
32 def splitData(data, rate=0.7): #默认是0.7的训练集, 0.2的测试集
33     # 默认训练集比例为0.7
34     dataX, dataY = data
35     nSample = dataX.shape[0]
36     nTrain = int(nSample*rate)
37     trainData = (dataX[:nTrain], dataY[:nTrain])
38     testData = (dataX[nTrain:], dataY[nTrain:])
39     return trainData, testData
40
41
42 # 获取训练集和测试集, 用80%的数据来训练拟合, 20%的数据来预测
43 rate = 0.8
44 trainData, testData = splitData(data, rate=rate)
```

```

1  # step4: 定义模型
2  class LstmModel(nn.Module):
3      def __init__(self, inputSize=5, hiddenSize=6):
4          super().__init__()
5          # LSTM层-> 两个LSTM单元叠加
6          self.lstm = nn.LSTM(input_size = inputSize,
7                               hidden_size = hiddenSize,
8                               num_layers = 2)
9          self.output = nn.Linear(6,1) # 线性输出
10
11
12  def forward(self,x):
13      # x: input->(time_step, batch, input_size)
14      # x的维度是【数量量: 整批数量量: 输入特征维度】
15      # x是【112 ,1, 3】
16      # lstm两层, 目标是从 3->6
17
18      x1, (h ,c)= self.lstm(x)
19      # x1: output->(time_step, batch, output_size)
20
21      a, b, c = x1.shape
22      out = self.output(x1.view(-1,c)) # 只有三维数据转化为二维才能作为输入
23      # 重新将结果转化为三维
24      out = out.view(a,b,-1)
25      return out
26
27  # 定义模型
28  lstm = LstmModel(inputSize=3) # inputSize与look_back保持一致
29
30  # step5.模型训练
31  def training_loop(nEpochs, model, optimizer, lossFn, trainData,
32                    testData=None):
33      trainX, trainY = trainData
34      if testData is not None:
35          testX, testY = testData
36      for epoch in range(1, nEpochs+1):
37          optimizer.zero_grad() # 梯度清0
38          trainP = model(trainX)
39          loss = lossFn(trainP, trainY)
40          loss.backward() # 反向传播
41          optimizer.step()
42          if epoch % 100 == 0:
43              print(f"Epoch: {epoch}, Loss: {loss.item()}")
44      return model

```

```

45
46 # 使用优化器Adam比SGD更好
47 optimizer = optim.Adam(lstm.parameters(), lr=0.1)
48 loss_func = nn.MSELoss()
49
50 # 训练模型
51 lstm = training_loop(nEpochs=1000, model= lstm,
52                     optimizer=optimizer, lossFn=loss_func,
53                     trainData=trainData)

```

Dockerfile

```

1 #Step6: 可视化
2 dataX, dataY = data # 原始数据 -> (time_step, batch, input_size)
3 dataY = dataY.view(-1).data.numpy() # 展开为1维
4 dataY = dataY * (maxPassenger - minPassenger) + minPassenger
5 dataP = lstm(dataX) # 进行拟合
6 dataP = dataP.view(-1).data.numpy() # 展开为1维
7 dataP = dataP * (maxPassenger - minPassenger) + minPassenger
8
9 nTrain = int(dataX.shape[0] * rate) # 拟合的数量
10 nData = dataX.shape[0] # 预测的数量
11
12 # 绘制对比图
13 plt.rcParams['font.sans-serif'] = 'KaiTi' # 正常显示中文
14 fig = plt.figure(dpi=400)
15 ax = fig.add_subplot(111)
16 ax.plot(dataY, color='blue', label="实际值")
17 ax.plot(np.arange(nTrain), dataP[:nTrain], color='green', \
18         linestyle='--', label = '拟合值')
19 ax.plot(np.arange(nTrain, nData), dataP[nTrain:], \
20         linestyle='--', color = 'red', label='预测值')
21 ax.legend()
22 fig.savefig('test.png', dpi=400)
23

```

