

热点挖掘：LSTM模型

数据集

2022年5月16日到9月12日全国新增确诊人数与新增无症状感染者数量

思路

初步想法是直接将前60%的数据制作成训练集，后40%用于预测，偏离预测值较大的即为当天出现热点信息。但由于数据集较小，且通过抽取前几日确诊人数特征预测当日确诊人数，缺乏一定的逻辑合理性，因此，改换思路，将新增无症状感染者数据作为训练集，将新增确诊人数作为测试集，然后再将偏离值较大的选出，作为本次热点分析结果，这是一个朴素的想法，无症状感染者与确诊者有较强的关联性，通过挖掘无症状感染者的数量变化规律可以较好的预测出确诊人数，如若不然，就说明当天有其他因素影响了确诊人数的变化，需要作为热点重点分析。

模型

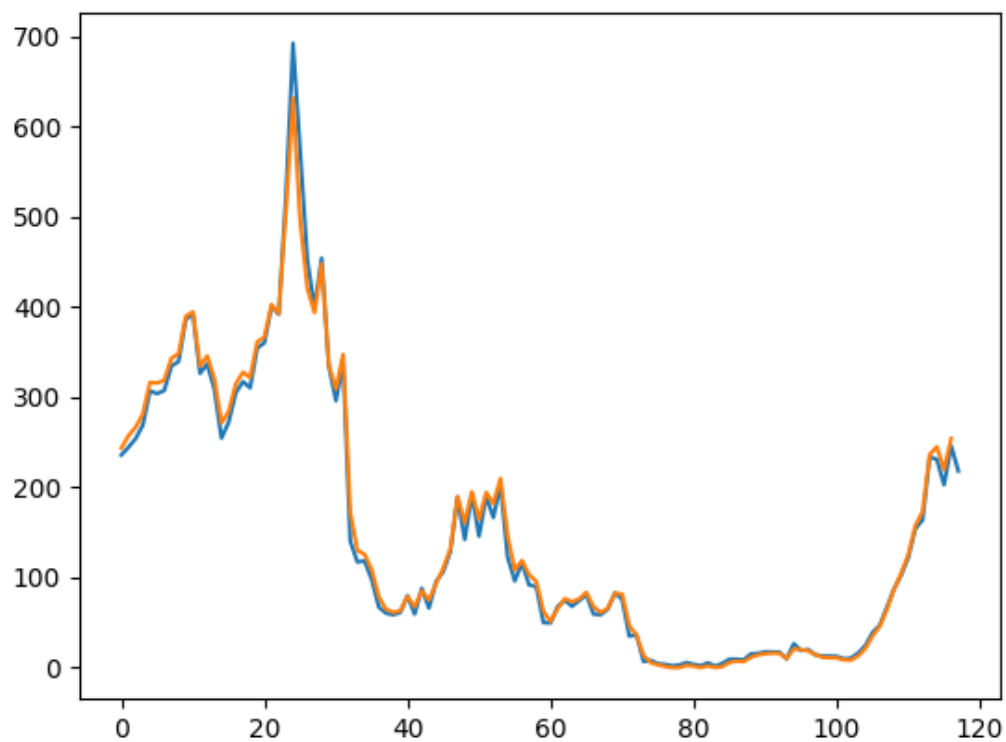
```
1 model = Sequential()
2 model.add(LSTM(4, input_shape=(None,1)))
3 model.add(Dense(1))
4 adam = optimizers.Adam(learning_rate=0.01, beta_1=0.9, beta_2=0.999,
5   amsgrad=False)
6 model.compile(loss='mean_squared_error', optimizer=adam)
```

超参数	值
time_step(look_back)	2
epoch	30
batch_size	3
learning_rate	0.01

结果

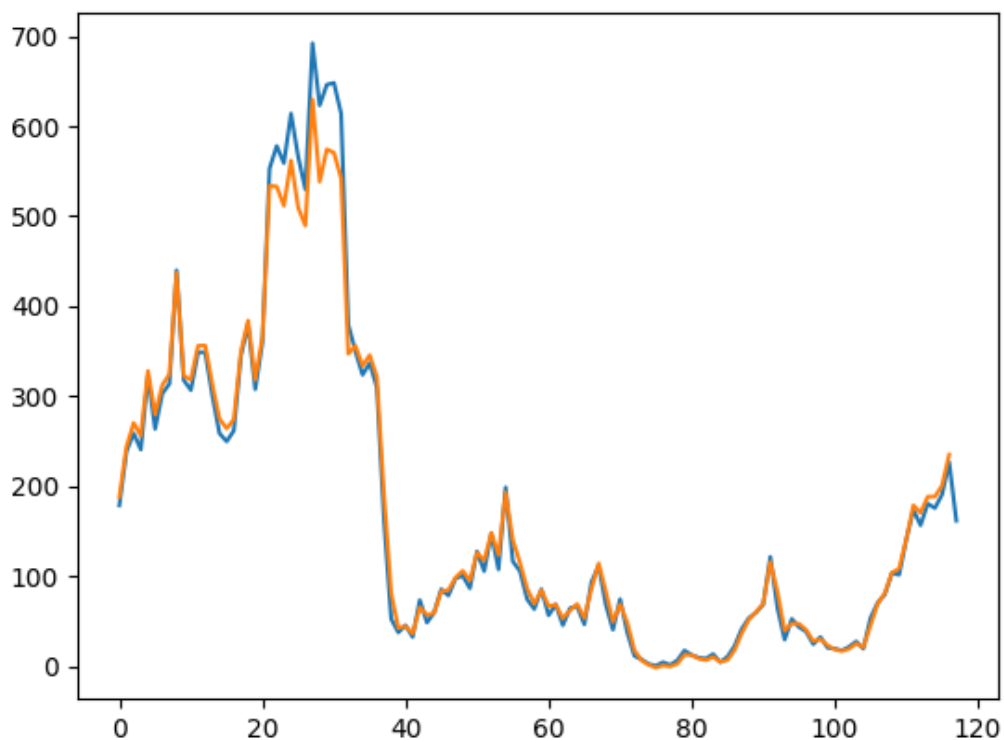
- 蓝线为标签
- 橙线为预测

训练结果



```
1 //loss 收敛
2 40/40 - 0s - loss: 0.0036 - 343ms/epoch - 9ms/step
3 Epoch 26/30
4 40/40 - 0s - loss: 0.0039 - 379ms/epoch - 9ms/step
5 Epoch 27/30
6 40/40 - 0s - loss: 0.0035 - 362ms/epoch - 9ms/step
7 Epoch 28/30
8 40/40 - 0s - loss: 0.0037 - 335ms/epoch - 8ms/step
9 Epoch 29/30
10 40/40 - 0s - loss: 0.0036 - 346ms/epoch - 9ms/step
11 Epoch 30/30
12 40/40 - 0s - loss: 0.0039 - 361ms/epoch - 9ms/step
```

预测结果



热点挖掘

结果

1	日期下标	真实值	预测值	偏差百分比
2	25	566	[509]	[10.07067138]
3	28	623	[538]	[13.64365971]
4	29	646	[574]	[11.14551084]
5	30	648	[570]	[12.03703704]
6	31	614	[541]	[11.88925081]
7	37	162	[189]	[16.66666667]
8	38	53	[81]	[52.83018868]
9	39	38	[42]	[10.52631579]
10	42	74	[65]	[12.16216216]
11	43	49	[57]	[16.32653061]
12	51	106	[117]	[10.37735849]
13	53	108	[123]	[13.88888889]
14	55	117	[140]	[19.65811966]
15	57	75	[86]	[14.66666667]
16	60	56	[67]	[19.64285714]
17	62	46	[53]	[15.2173913]
18	65	47	[53]	[12.76595745]
19	68	69	[83]	[20.28985507]
20	69	41	[49]	[19.51219512]
21	71	38	[48]	[26.31578947]
22	72	12	[18]	[50.]
23	73	8	[7]	[12.5]
24	74	3	[2]	[33.33333333]
25	75	1	[-1]	[200.]
26	76	5	[1]	[80.]
27	77	2	[0]	[100.]

```
28      78 6      [3] [50.]
29      79 18     [13] [27.77777778]
30      81 10     [8] [20.]
31      82 9      [7] [22.22222222]
32      83 13     [10] [23.07692308]
33      85 11     [7] [36.36363636]
34      86 23     [18] [21.73913043]
35      87 42     [37] [11.9047619]
36      92 65     [82] [26.15384615]
37      93 30     [39] [30.]
38      97 25     [28] [12.]
39     102 22     [19] [13.63636364]
40     105 54     [45] [16.66666667]
```

完整代码

```
1  from pickletools import optimize
2  import numpy
3  import matplotlib.pyplot as plt
4  from keras.models import Sequential
5  from keras.layers import Dense
6  from keras.layers import LSTM
7  import pandas as pd
8  import os
9  from keras.models import Sequential, load_model
10 from sklearn.preprocessing import MinMaxScaler
11 from keras import optimizers
12
13 def create_dataset(dataset, look_back):
14     #这里的look_back与timestep相同
15     dataX, dataY = [], []
16     for i in range(len(dataset)-look_back):
17         a = dataset[i:(i+look_back)]
18         dataX.append(a)
19         dataY.append(dataset[i + look_back])
20     return numpy.array(dataX), numpy.array(dataY)
21
22 dataframe1 = pd.read_excel('data.xlsx', usecols=[2])
23 dataframe2 = pd.read_excel('data.xlsx', usecols=[1])
24
25
26 dataset1 = dataframe1.values
27 dataset2 = dataframe2.values
28
29 # 将整型变为float
30 dataset1 = dataset1.astype('float32')
31 dataset2 = dataset2.astype('float32')
32 #归一化
33 scaler = MinMaxScaler(feature_range=(0, 1))
34 dataset1 = scaler.fit_transform(dataset1)
35 dataset2 = scaler.fit_transform(dataset2)
36
37 # train_size = int(len(dataset) * 0.65)
38 trainlist = dataset1
39 testlist = dataset2
40
41
```

```

42
43 #训练数据太少 look_back并不能过大
44 look_back = 2
45 trainX,trainY = create_dataset(trainlist,look_back)
46
47 testX,testY = create_dataset(testlist,look_back)
48
49 trainX = numpy.reshape(trainX, (trainX.shape[0], trainX.shape[1], 1))
50 testX = numpy.reshape(testX, (testX.shape[0], testX.shape[1], 1))
51
52 # # create and fit the LSTM network
53 # model = Sequential()
54 # model.add(LSTM(4, input_shape=(None,1)))
55 # model.add(Dense(1))
56 # adam = optimizers.Adam(learning_rate=0.01, beta_1=0.9, beta_2=0.999,
57 # amsgrad=False)
58 # model.compile(loss='mean_squared_error', optimizer=adam)
59
60 # model.fit(trainX, trainY, epochs=30, batch_size=3, verbose=2)
61 # model.save(os.path.join("DATA","Test" + ".h5"))
62
63 # make predictions
64 model = load_model(os.path.join("DATA","Test" + ".h5"))
65 trainPredict = model.predict(trainX)
66 testPredict = model.predict(testX)
67
68 #反归一化
69 trainPredict = scaler.inverse_transform(trainPredict)
70 trainY = scaler.inverse_transform(trainY)
71 testPredict = scaler.inverse_transform(testPredict)
72 testY = scaler.inverse_transform(testY)
73
74
75 testY = testY.astype(int)
76 testPredict = testPredict.astype(int)
77 testY = testY.reshape(-1)
78
79 print("选出偏差百分比大于20%的数据")
80 print("日期下标\t真实值\t预测值\t偏差百分比")
81
82 df1 = pd.DataFrame({'predict': testY})
83 df1.to_excel('predict.xlsx', sheet_name='Sheet1', index=False) # index false
84 # 为不写入索引
85 for i in range(len(testY)-1):
86     delta = (testY[i] - testPredict[i+1])/testY[i] * 100
87     delta = abs(delta)
88     if delta > 10:
89         print(i,testY[i],testPredict[i+1],delta)
90
91 # plt.plot(trainY)
92 # plt.plot(trainPredict[1:])
93 # plt.show()
94 # plt.plot(testY)
95 # plt.plot(testPredict[1:])
96 # plt.show()

```

