

# 趋势预测报告（LSTM模型）

## [参考代码链接](#)

该报告重点关注LSTM模型的参数调优问题，通过使用不同的参数组合来构建模型，比较不同模型下的预测结果，比较结果，得出最优参数。最后我们使用得到的最优参数组合来构建模型，并使用该模型进行趋势预测。

我们首先需要通过测试来得出，趋势预测效果最好的时间步个数。我们依次将时间步设置为1~30，分别用这些时间步个数构建模型，并进行训练和预测，通过比较预测的结果来得出最优的时间步个数，用于后续的趋势预测。

```
1  import os
2  import numpy as np
3  import pandas as pd
4  from datetime import datetime
5
6  from keras import Sequential
7  from keras.layers import Dropout, LSTM, Dense
8  from matplotlib import pyplot as plt
9  from sklearn.preprocessing import MinMaxScaler
10
11  os.environ['TF_CPP_MIN_LOG_LEVEL'] = '2'
12
13
14  def lstm_model_test(csv_file, steps, score1):
15      # Load the data
16      custom_date_parser = lambda x: datetime.strptime(x, "%Y-%m-%d %H:%M:%S")
17      df = pd.read_csv(csv_file, parse_dates=['trendTime'],
18      date_parser=custom_date_parser, index_col='trendTime')
19      # Rename the columns
20      last_data_time = df.index[-1]
21
22      start_time = pd.to_datetime(last_data_time) - pd.Timedelta(days=1) -
23      pd.Timedelta(days=7)
24      start_time = start_time.strftime('%Y-%m-%d %H:%M:%S')
25      print(start_time)
26
27      # last_data_time为%Y-%m-%d %H:%M:%S格式
28      mid_time = pd.to_datetime(last_data_time) - pd.Timedelta(days=1)
29      mid_time = mid_time.strftime('%Y-%m-%d %H:%M:%S')
30      print(mid_time)
31
32      end_time = pd.to_datetime(mid_time) + pd.Timedelta(days=1)
33      end_time = end_time.strftime('%Y-%m-%d %H:%M:%S')
34      print(end_time)
35
36      df = df.iloc[:, :300, :]
37
38      # 划分训练集和测试集
```

```

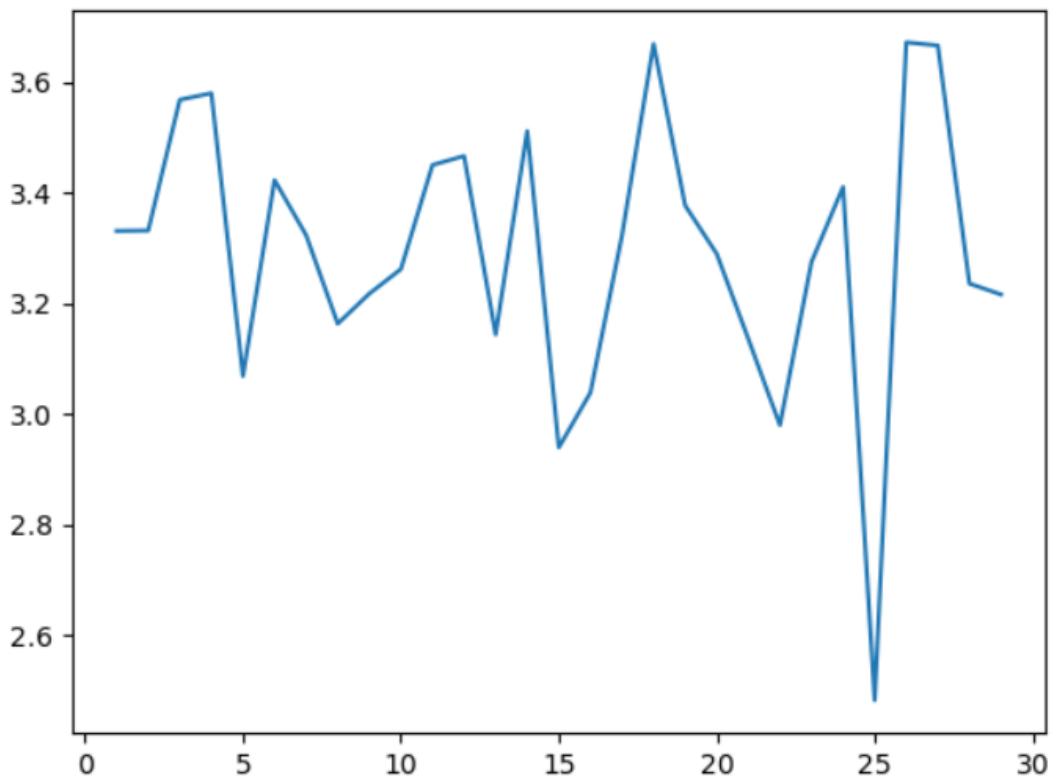
37     train_df = df[(df.index < mid_time) & (df.index > start_time)].values
38     test_df = df[(df.index >= mid_time) & (df.index < end_time)].values
39
40     print(train_df.shape)
41     print(test_df.shape)
42
43     sc = MinMaxScaler(feature_range=(0, 1))
44     train_df_scaled = sc.fit_transform(train_df)
45
46     x_train = []
47     y_train = []
48     for i in range(steps, train_df_scaled.shape[0]):
49         x_train.append(train_df_scaled[i-steps:i, 0])
50         y_train.append(train_df_scaled[i, 0])
51
52     x_train, y_train = np.array(x_train), np.array(y_train)
53
54     x_train = np.reshape(x_train, (x_train.shape[0], x_train.shape[1], 1))
55
56     # 构建只有一层的LSTM模型
57     model = Sequential()
58     model.add(LSTM(units=128, return_sequences=True, input_shape=
(x_train.shape[1], 1)))
59     model.add(Dropout(0.2))
60
61     model.add(LSTM(units=128))
62     model.add(Dropout(0.2))
63
64     model.add(Dense(units=1))
65
66     model.compile(optimizer='rmsprop', loss='mse')
67
68     model.fit(x_train, y_train, epochs=20, batch_size=32)
69
70     # 获取train_df的最后steps个数据，用于预测
71     df1 = df.tail(steps)
72     test_df1 = df[(df.index >= mid_time) & (df.index < end_time)]
73     inputs = pd.concat([df1, test_df1], axis=0).values
74
75     inputs = inputs.reshape(-1, 1)
76     inputs = sc.transform(inputs)
77     x_test = []
78     for i in range(steps, inputs.shape[0]):
79         x_test.append(inputs[i - steps:i, 0])
80     x_test = np.array(x_test)
81
82     x_test = np.reshape(x_test, (x_test.shape[0], x_test.shape[1], 1))
83
84     predict_test = model.predict(x_test) # 预测
85     predict_test = sc.inverse_transform(predict_test) # 反归一化
86
87     score = model.evaluate(x_test, test_df1.values)

```

```

88     score1.append(score)
89     print(score)
90
91 if __name__ == "__main__":
92     score1 = []
93     for steps in range(1, 31):
94         lstm_model_test("cs4.csv", steps, score1)
95         # 输出最大值对应的索引
96         print(score1.index(max(score1)) + 1)
97     plt.plot(range(1, 30), score1)
98     plt.show()

```



结果显示当时时间步个数为26时，预测效果最佳，所以后续的预测过程，我们将把26作为时间步个数用于构建模型。

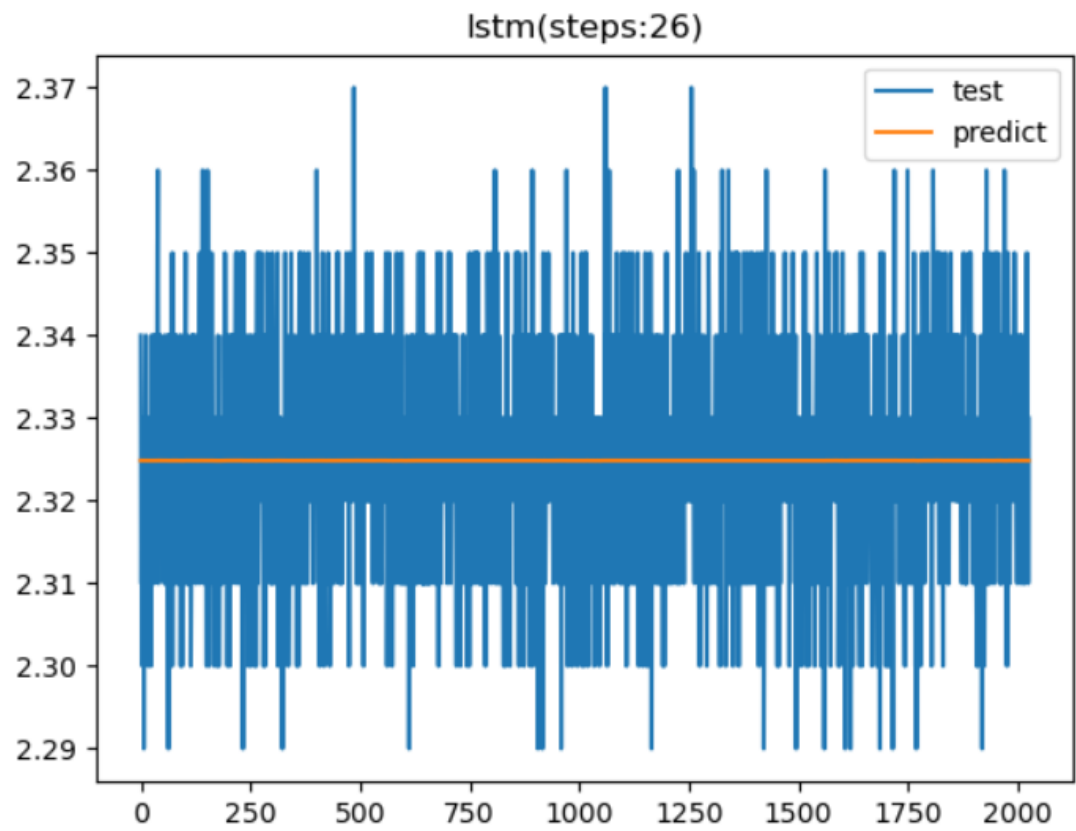
### 数据参数

参数	值
开始时间	2023-04-26 16: 04: 14
分隔时间	2023-05-03 16: 04: 14
结束时间	2023-05-04 16: 04: 14
训练数据集数据量	133436

参数	值
测试数据集数据量	20254

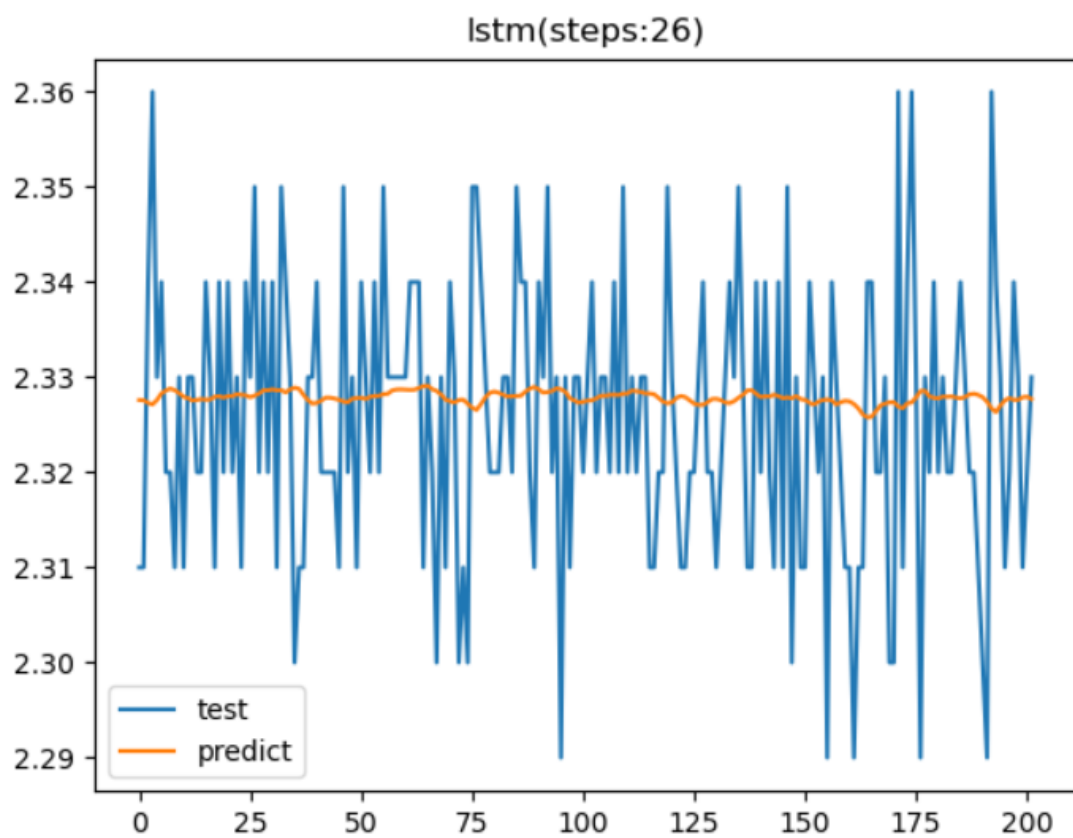
我们尝试每10个数据取一个进行训练和测试

此时用于训练和测试的数据量分别为13343和2026



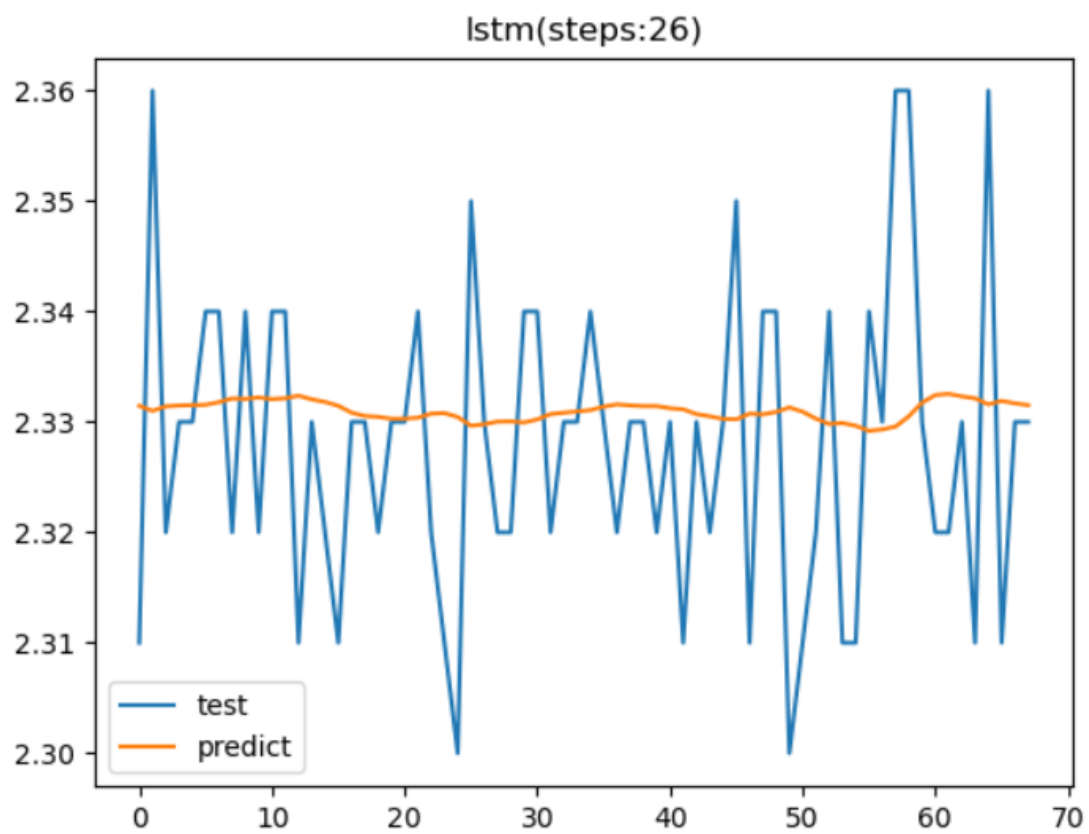
我们尝试每100个数据取一个进行训练和测试

此时用于训练和测试的数据量分别为1335和202



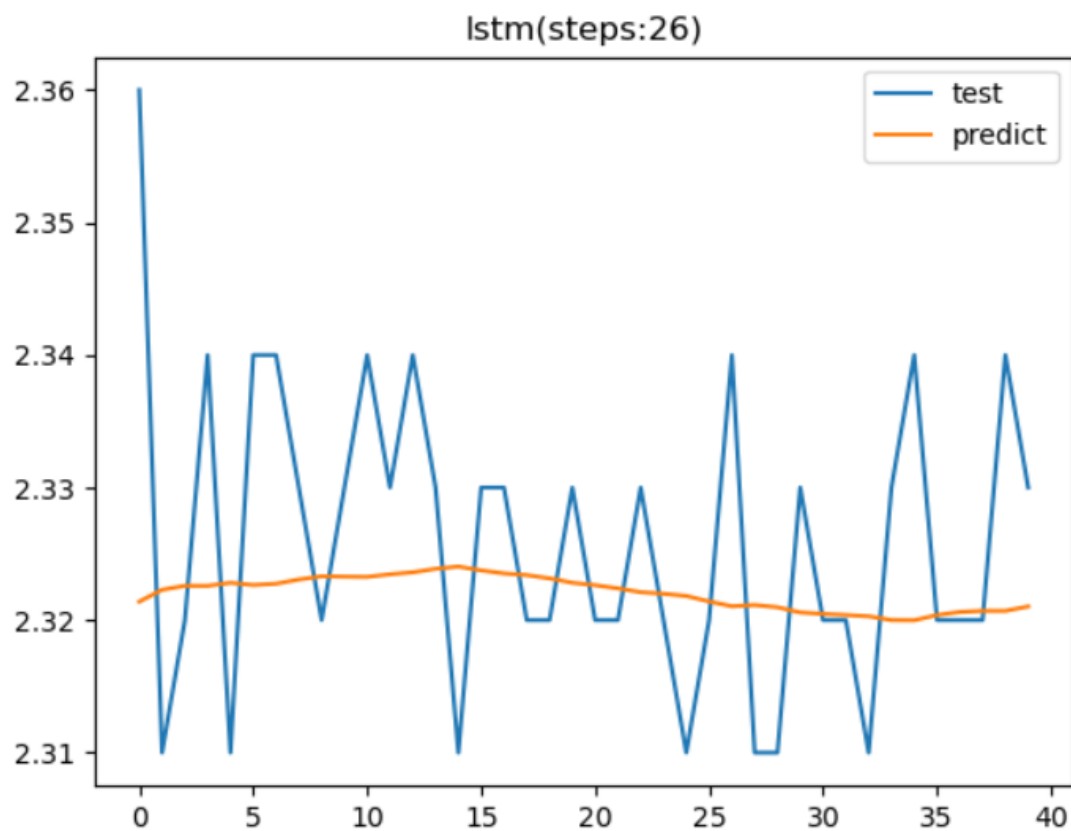
我们尝试每300个数据取一个进行训练和测试

此时用于训练和测试的数据量分别为445和68



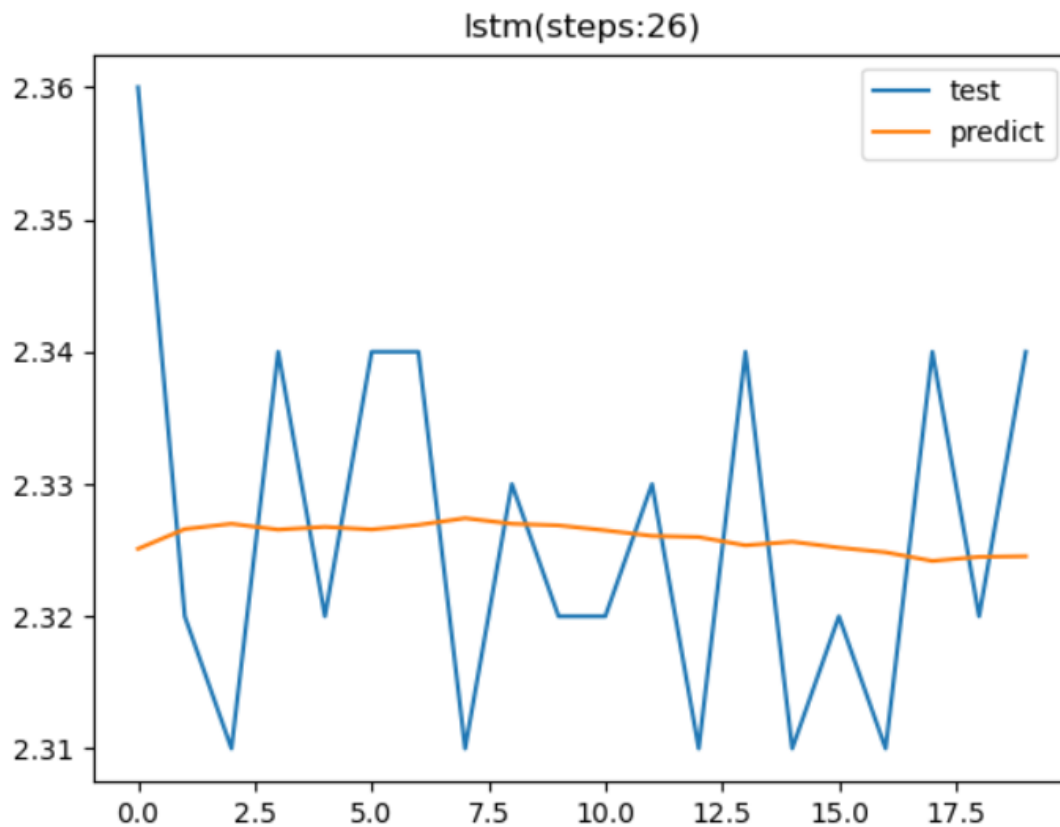
我们尝试每500个数据取一个进行训练和测试

此时用于训练和测试的数据量分别为267和40



我们尝试每1000个数据取一个进行训练和测试

此时用于训练和测试的数据量分别为133和20



每100条数据取一个数据时，模型预测得到的趋势曲线的走势能够更好地呈现出来，所以在后续的测试过程中我们将继续使用每100条数据取一个数据的形式来，得到更好的预测结果。

```

1 param_grid = {
2     'batch_size': [32, 64, 128], # batch_size是指进行梯度下降时每个batch包含的样
    本数
3     'epochs': [20, 50, 100], # epochs是指训练过程中数据将被“轮”多少次
4     'optimizer': ['rmsprop', 'adam', 'sgd'] # optimizer是指优化器
5 }
6
7 all_params = [dict(zip(param_grid.keys(), v)) for v in
8               itertools.product(*param_grid.values())]
9 score1 = [] # 用来记录不同参数组合对应的rmse
10 for params in all_params:
11     m = Sequential()
12     m.add(LSTM(units=128, return_sequences=True, input_shape=
13           (x_train.shape[1], 1)))
14     m.add(Dropout(0.2))
15
16     m.add(LSTM(units=128))
17     m.add(Dropout(0.2))
18
19     m.add(Dense(units=1))
20
21     m.compile(optimizer=params['optimizer'], loss='mse')

```



```

20     m.fit(x_train, y_train, epochs=params['epochs'],
batch_size=params['batch_size'])
21     predict_test = m.predict(x_test) # 预测
22     predict_test = sc.inverse_transform(predict_test) # 反归一化
23
24     score = m.evaluate(x_test, test_df1.values)
25     score1.append(score)
26
27     # 输出最佳参数组合
28     print('best params: ', all_params[np.argmin(score1)])

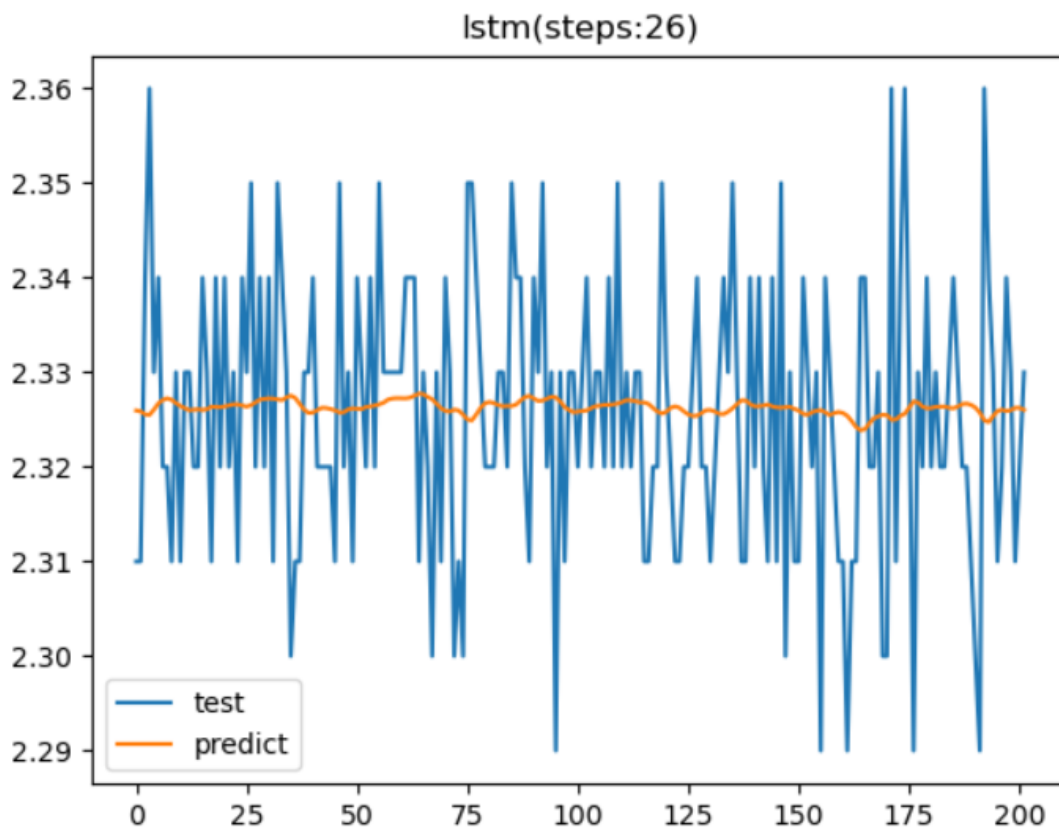
```

此时的最优参数为;

```

1 | best params: {'batch_size': 128, 'epochs': 50, 'optimizer': 'rmsprop'}

```



从模型的预测结果可以看出，此时预测得到的趋势变化曲线，已经能够较为清晰地展现出数据的变化趋势，当真实的数据点总体呈上升趋势时，趋势预测曲线也呈上升趋势，当真实数据点呈下降趋势时，趋势预测曲线也呈下降趋势，符合实验预期结果。但是，我们可以发现，该模型得到的趋势预测曲线只能在一定程度上反应数据的变化趋势，在数据的数值预测上的准确度还有待提升，这部分提升可以通过对不同的参数组合进行测试，并比较预测结果，不断缩小最优参数组合的范围，最终确定最优参数，用于模型构建。