

1. 非线性拟合实验报告

Author : Xinyu Zhao, Beihang University

Email: zxy2021@buaa.edu.cn

1. 非线性拟合实验报告

摘要

1.1 问题描述

1.2 数据描述

1.3 模型建立与拟合

1.3.1 多项式回归

1.3.2 神经网络拟合

1.4 结论

1.5 附录

摘要

本报告详细探讨了使用多项式回归、神经网络等多种方法进行非线性拟合的过程。我们的目标是找到一个能够最佳拟合给定训练集的模型，同时评估其在独立测试集上的性能。通过对不同非线性模型的比较、模型训练过程的优化、以及深入的结构分析，我们旨在提高模型的预测准确性和泛化能力。

1.1 问题描述

给定范围为0~10、均匀分布的100组二维数据点组成的训练集和测试集，我们采用**均方误差(MAE)**作为评价指标，旨在通过非线性拟合方法捕捉数据中的复杂关系。本实验首先尝试了多项式回归，并随后探索了使用神经网络进行拟合。

均方误差(MAE)是机器学习中较为常见的一种误差指标，对应于二次损失函数，使得计算和数学推导变得更加简单。在优化问题中，使用均方误差可以导出解析解，这是因为对于线性模型，二次损失函数的导数是线性的，这使得找到最小值变得相对容易。因此在本次实验中定义的评价指标为**均方误差**。

1.2 数据描述

数据集和测试集中分别包含范围为0~10且均匀分布的100组数据，作图如下，可以看出数据均为复杂的非线性分布。

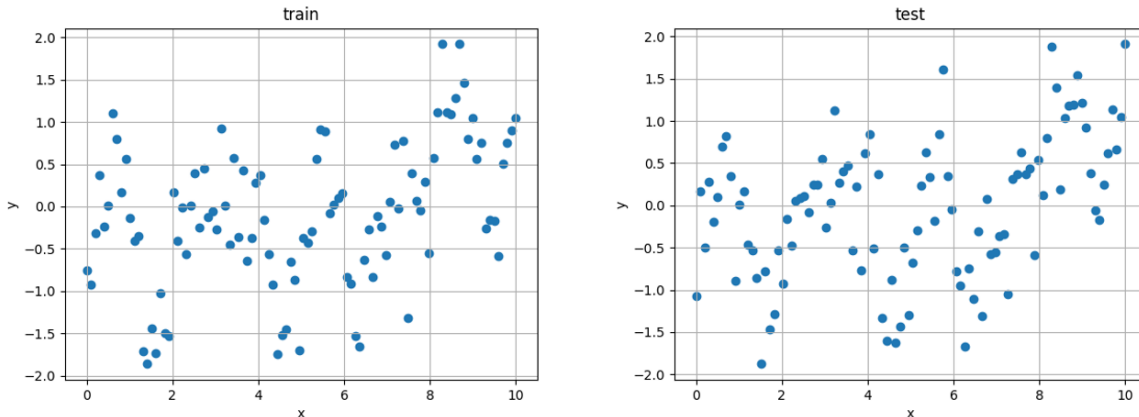


图1 原始训练/测试集的散点图


```

4     poly_model_best.fit(X_train, y_train)
5     y_pred_best = poly_model_best.predict(X_test)
6
7     # 计算 MAE
8     loss_best = mean_absolute_error(y_test, y_pred_best)
9     print(f'Best Model (Degree {degree_best}) MAE:', loss_best)
10    loss_list.append(loss_best)
11    # 创建用于绘图的预测范围
12    x_range = np.linspace(X_train['x'].min(), X_train['x'].max(), 500)
13    y_range_pred = poly_model_best.predict(x_range.reshape(-1, 1))
14
15    if degree_best in [5,8,10,15,20,35,50,65]:
16        # 训练集绘图
17        plt.figure(figsize=[10, 6])
18        plt.scatter(X_train['x'], y_train, color='blue', alpha=0.7,
19label='Training Data')
20        plt.plot(x_range, y_range_pred, color='green', label=f'Polynomial
Degree {degree_best} Fit')
21        plt.title('Polynomial Regression Fit')
22        plt.xlabel('x')
23        plt.ylabel('y')
24        plt.legend()
25        plt.grid(True)
26        plt.show()
27
28        # 测试集绘图
29        plt.figure(figsize=[10, 6])
30        plt.scatter(X_test['x'], y_test, color='red', alpha=0.5,
31label='Testing Data')
32        plt.plot(x_range, y_range_pred, color='green', label=f'Polynomial
Degree {degree_best} Fit')
33        plt.title('Polynomial Regression Fit')
34        plt.xlabel('x')
35        plt.ylabel('y')
36        plt.legend()
37        plt.grid(True)
38        plt.show()

```

将每次的损失值存入列表loss_list中，绘制趋势图如下：

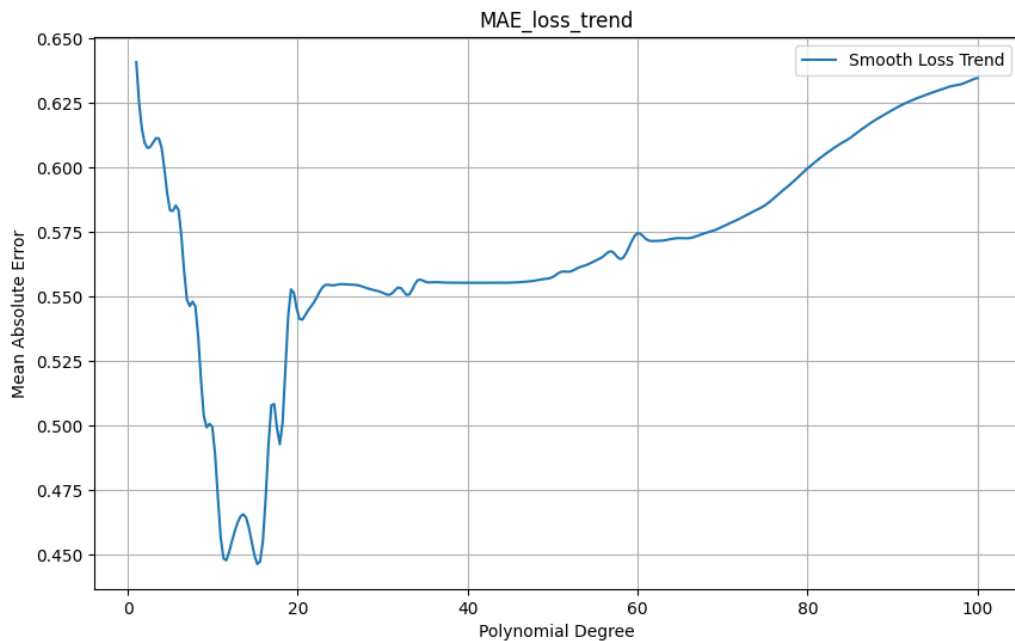


图4 损失值与多项式阶数之间的关系

代码如下：

```

1 # 创建一个索引列表，用作x轴数据
2 x_values = np.array(range(1, len(loss_list) + 1))
3 y_values = np.array(loss_list)
4
5 # 创建一个平滑的曲线
6 x_smooth = np.linspace(x_values.min(), x_values.max(), 300) # 生成足够多的x值
   以使曲线平滑
7 spl = make_interp_spline(x_values, y_values, k=3) # k是样条曲线的阶数，这里使用
   三次样条
8 y_smooth = spl(x_smooth)
9
10 # 使用matplotlib绘制平滑的曲线
11 plt.figure(figsize=[10, 6])
12 plt.plot(x_smooth, y_smooth, label='Smooth Loss Trend') # 平滑曲线
13 plt.title('MAE_loss_trend')
14 plt.xlabel('Polynomial Degree')
15 plt.ylabel('Mean Absolute Error')
16 plt.legend()
17 plt.grid(True)
18 plt.show()

```

1.3.2 神经网络拟合

为了进行拟合模型之间的性能比较，本文提出利用神经网络进行拟合，以提高模型的灵活性和准确性，关键代码如下：

```

1 # 数据读取
2 train_data = pd.read_excel('train.xlsx')
3 test_data = pd.read_excel('test.xlsx')
4 X_train = train_data['x'].values.reshape(-1, 1)
5 Y_train = train_data['y'].values.reshape(-1, 1)

```

```

6 X_test = test_data['x'].values.reshape(-1, 1)
7 Y_test = test_data['y'].values.reshape(-1, 1)
8 # 定义神经网络结构
9 class SimpleNN(nn.Module):
10     def __init__(self):
11         super(SimpleNN, self).__init__()
12         self.fc1 = nn.Linear(1, 64) # 输入层到第一个隐藏层
13         self.fc2 = nn.Linear(64, 64) # 第二个隐藏层
14         self.fc3 = nn.Linear(64, 1) # 输出层
15
16     def forward(self, x):
17         x = torch.relu(self.fc1(x))
18         x = torch.relu(self.fc2(x))
19         return self.fc3(x)
20
21 # 初始化模型、损失函数和优化器
22 model = SimpleNN()
23 criterion = nn.MSELoss()
24 optimizer = optim.Adam(model.parameters(), lr=0.001)
25
26 # 转换数据为torch tensors
27 X_train_tensor = torch.tensor(X_train, dtype=torch.float32)
28 Y_train_tensor = torch.tensor(Y_train, dtype=torch.float32)
29 X_test_tensor = torch.tensor(X_test, dtype=torch.float32)
30
31 # 训练模型
32 epochs = 1000
33 for epoch in range(epochs):
34     optimizer.zero_grad() # 清除之前的梯度
35     output = model(X_train_tensor) # 前向传播
36     loss = criterion(output, Y_train_tensor) # 计算损失
37     loss.backward() # 反向传播, 计算梯度
38     optimizer.step() # 更新权重
39
40     if epoch % 100 == 0:
41         print(f'Epoch [{epoch+1}/{epochs}], Loss: {loss.item():.4f}')
42
43 # 使用训练好的模型进行预测
44 model.eval() # 设置模型为评估模式
45 with torch.no_grad():
46     Y_pred_nn = model(X_test_tensor).view(-1).numpy()
47
48 # 计算平均绝对误差
49 mae_nn = mean_absolute_error(Y_test, Y_pred_nn)
50 print(f'Test MAE of the neural network: {mae_nn}')
51
52 '''
53 Epoch [1/1000], Loss: 0.7944
54 Epoch [101/1000], Loss: 0.5523
55 Epoch [201/1000], Loss: 0.5039
56 Epoch [301/1000], Loss: 0.5000
57 Epoch [401/1000], Loss: 0.4995
58 Epoch [501/1000], Loss: 0.4992
59 Epoch [601/1000], Loss: 0.4988
60 Epoch [701/1000], Loss: 0.4983
61 Epoch [801/1000], Loss: 0.4978

```

```
62 Epoch [901/1000], Loss: 0.4973
63 Test MAE of the neural network: 0.5818524205386759
64 ...
```

发现运行结果得到的误差为0.5818524205386759，不如15阶多项式模型的拟合结果，推测原因如下：

- **神经网络结构不当**：如果神经网络的规模（层数或每层的神经元数量）太小，它可能没有足够的能力（表达能力）来捕捉数据中的复杂关系；15阶多项式由于其高度的非线性和复杂性，能够很好地拟合或过拟合某些类型的数据。
- **迭代次数不足**：如果训练神经网络的迭代次数（epoch数量）不够，网络可能还没有足够的时间来学习数据中的模式。
- **数据量**：相比于多项式拟合，神经网络通常需要更多的数据来有效学习。
- **优化器和学习率**：选择不恰当的优化算法或设置不合适的学习率导致神经网络训练不充分。
- **过拟合**：神经网络可能过度学习训练数据中的噪声，而不是数据的真实模式，导致在测试数据上的表现不佳。
- **欠拟合**：相反，如果神经网络太简单，不能捕捉数据中的所有相关模式，也会导致拟合效果不佳。

为此，接下来将代码进行了改进，优化神经网络结构，调整训练次数和学习率。

```
1  # 定义改进的神经网络结构
2  class ImprovedNN(nn.Module):
3      def __init__(self):
4          super(ImprovedNN, self).__init__()
5          self.fc1 = nn.Linear(1, 128) # 输入层到第一个隐藏层，增加神经元数量
6          self.dropout1 = nn.Dropout(0.25) # 添加Dropout层以减少过拟合
7          self.fc2 = nn.Linear(128, 128) # 第二个隐藏层
8          self.dropout2 = nn.Dropout(0.25)
9          self.fc3 = nn.Linear(128, 1) # 输出层
10
11      def forward(self, x):
12          x = torch.relu(self.fc1(x))
13          x = self.dropout1(x)
14          x = torch.relu(self.fc2(x))
15          x = self.dropout2(x)
16          return self.fc3(x)
17
18  # 初始化模型、损失函数和优化器
19  model = ImprovedNN()
20  criterion = nn.MSELoss()
21  optimizer = optim.Adam(model.parameters(), lr=0.01)
22  scheduler = torch.optim.lr_scheduler.StepLR(optimizer, step_size=1000,
23  gamma=0.9) # 学习率调度器
24
25  # 转换数据为torch tensors
26  X_train_tensor = torch.tensor(X_train, dtype=torch.float32)
27  Y_train_tensor = torch.tensor(Y_train, dtype=torch.float32)
28  X_test_tensor = torch.tensor(X_test, dtype=torch.float32)
29
30  # 训练模型
31  epochs = 10000
```

```

31 for epoch in range(epochs):
32     optimizer.zero_grad() # 清除之前的梯度
33     output = model(X_train_tensor) # 前向传播
34     loss = criterion(output, Y_train_tensor) # 计算损失
35     loss.backward() # 反向传播, 计算梯度
36     optimizer.step() # 更新权重
37     scheduler.step() # 更新学习率
38
39     if epoch % 1000 == 0:
40         print(f'Epoch [{epoch+1}/{epochs}], Loss: {loss.item():.4f}, LR:
{scheduler.get_last_lr()[0]}')
41
42 # 使用训练好的模型进行预测
43 model.eval() # 设置模型为评估模式
44 with torch.no_grad():
45     Y_pred_nn = model(X_test_tensor).view(-1).numpy()
46
47 # 计算平均绝对误差
48 mae_nn = mean_absolute_error(Y_test, Y_pred_nn)
49 print(f'Test MAE of the neural network: {mae_nn}')
50 '''
51 Epoch [1/10000], Loss: 2.4662, LR: 0.01
52 Epoch [1001/10000], Loss: 0.3971, LR: 0.009000000000000001
53 Epoch [2001/10000], Loss: 0.4269, LR: 0.008100000000000001
54 Epoch [3001/10000], Loss: 0.3931, LR: 0.007290000000000001
55 Epoch [4001/10000], Loss: 0.4117, LR: 0.006561000000000002
56 Epoch [5001/10000], Loss: 0.3081, LR: 0.005904900000000002
57 Epoch [6001/10000], Loss: 0.2981, LR: 0.005314410000000002
58 Epoch [7001/10000], Loss: 0.3212, LR: 0.004782969000000002
59 Epoch [8001/10000], Loss: 0.3076, LR: 0.004304672100000002
60 Epoch [9001/10000], Loss: 0.3468, LR: 0.003874204890000002
61 Test MAE of the neural network: 0.4521120688577277
62 '''

```

我们可以明显地发现，拟合效果要好于优化前的神经网络，误差值仅为0.45。

通过深入的模型比较和结构分析，我们得到了关于非线性拟合最有效方法的见解。我们发现，尽管多项式回归简单且易于实现，但在处理复杂数据时，神经网络提供了更高的灵活性和准确性。此外，模型优化和正则化技术的应用对于防止过拟合至关重要。

1.4 结论

本报告展示了非线性拟合在复杂数据集上的应用，并比较了多种方法的性能，结果强调了选择合适模型、优化训练过程和深入理解模型结构的重要性。以后将探索更多的非线性拟合技术，并在更广泛的数据集上验证这些方法的有效性。

1.5 附录

实验所用代码、图片以及数据集后续会上传到github上。