

HW-1: Regression

姓名：罗威 学号：SA24218095

1 概述

本次作业旨在构建一个全连接神经网络模型，对波士顿房价进行预测。通过对波士顿房价数据集进行特征筛选、数据预处理，利用PyTorch搭建并训练神经网络模型，最终评估模型在测试集上的性能。

1.1 训练结果总结

- 首先，对13组数据进行了标签相关性计算，并在代码中对相关性较低的数据进行筛选，以此降低不相关数据产生的噪声影响。
- 对相关性较低的数据组合进行剔除过后，训练结果在测试集上的均方误差为11.9592。
- 对先前的代码加入早停策略，使模型在验证集上均方误差不再下降时保存模型并停止训练。通过早停策略训练得到的模型在测试集上的均方误差进一步降低到6.0042。

1.2 软硬件环境

(1) 系统环境：Ubuntu 22.04 LTS

(2) 语言及框架版本：

- Python: 3.11.11
- Pytorch: 2.5.0+cu124
- Pandas: 2.2.3
- Sklearn: 1.6.0

2 标签相关性分析

利用 pandas 库自带的 `corr()` 函数进行相关性分析。该函数可以计算出多种相关性系数（pearson、kendall、spearman），默认计算的是皮尔逊系数。

```
1 import pandas as pd
2
3 data = pd.read_excel('./BostonHousingData.xlsx', sheet_name='Sheet1')
4
5 correlation_matrix = data.corr()
6 correlation =
7 correlation_matrix['MEDV'].sort_values(ascending=False).drop('MEDV')
8
9 print("The correlation coefficients with MEDV:")
10 print(correlation)
```

```
1 The correlation coefficients with MEDV:
2 RM          0.695360
3 ZN          0.360445
4 B           0.333461
5 DIS         0.249929
6 CHAS        0.175260
7 AGE        -0.376955
```

```
8 RAD -0.381626
9 CRIM -0.388305
10 NOX -0.427321
11 TAX -0.468536
12 INDUS -0.483725
13 PTRATIO -0.507787
14 LSTAT -0.737663
15 Name: MEDV, dtype: float64
```

相关系数的取值范围为 $[-1, 1]$ ，其中系数越接近-1和1的标签，与房价的相关性越高。可以看到，计算出的相关系数中，DIS、CHAS与房价的相关性较低；同时，ZN、B、AGE、RAD、CRIM等的相关性相对也较低。因此，可以逐渐剔除这些相关性低下的标签，让模型获得更好的拟合效果。

3 全连接神经网络模型设计

构建了一个全连接神经网络模型 `Regression`，具体结构如下：

- (1) 包含5个有参层，分别为3个全连接层（`nn.Linear`）和2个批量归一化层（`nn.BatchNorm1d`）其中的2个隐藏层均使用 `nn.Linear` 进行线性变换，接着使用 `nn.BatchNorm1d` 进行批量归一化，`nn.ReLU` 作为激活函数，`nn.Dropout(0.2)` 进行正则化，防止过拟合；
 - (2) 使用常用的 `Adam` 优化器来进行参数更新，学习率设置为0.001，权重衰减设置为0.001，有助于模型在训练过程中更快地收敛并防止过拟合；
 - (3) 使用均方误差损失函数 `nn.MSELoss()`，用于衡量模型预测值与真实值之间的差异。
- 如下结果所示，最佳训练结果在测试集上的均方误差为11.9592。

```
1 import torch
2 import pandas as pd
3 import torch.nn as nn
4 import torch.optim as optim
5 from torch.utils.data import DataLoader, TensorDataset
6 from sklearn.preprocessing import StandardScaler
7
8 # 1.从.xlsx文件中读取数据集
9 data = pd.read_excel('BostonHousingData.xlsx', sheet_name='Sheet1')
10
11 # 2.计算相关系数并通过设置的阈值来筛选数据
12 threshold = 0.5
13 corre = data.corr()
14 correlation = corre['MEDV'].sort_values(ascending=False)
15 selected_features = correlation[abs(correlation) >=
16 threshold].index.tolist()
17
18 selected_data = data[selected_features]
19 x = selected_data.drop('MEDV', axis=1).values
20 y = selected_data['MEDV'].values.reshape(-1, 1)
21
22 # 3.划分训练集和测试集
23 x_train, x_test = x[:450], x[450:]
24 y_train, y_test = y[:450], y[450:]
25
26 # 4.数据标准化处理、转换格式为pytorch张量
27 scaler = StandardScaler()
28 x_train_scaled = scaler.fit_transform(x_train)
29 x_test_scaled = scaler.transform(x_test)
```

```

30 X_train_tensor = torch.tensor(X_train_scaled, dtype=torch.float32)
31 y_train_tensor = torch.tensor(y_train, dtype=torch.float32)
32 X_test_tensor = torch.tensor(X_test_scaled, dtype=torch.float32)
33 y_test_tensor = torch.tensor(y_test, dtype=torch.float32)
34
35 # 5.使用库函数加载数据集
36 train_dataset = TensorDataset(X_train_tensor, y_train_tensor)
37 train_loader = DataLoader(train_dataset, batch_size=32, shuffle=True)
38 test_dataset = TensorDataset(X_test_tensor, y_test_tensor)
39 test_loader = DataLoader(test_dataset, batch_size=32, shuffle=False)
40
41 # 定义输入数据维度，与阈值设置有关
42 input_size = selected_data.shape[1] - 1
43
44 # 6.全连接神经网络模型的定义
45 class Regression(nn.Module):
46     def __init__(self):
47         super().__init__()
48         self.fc = nn.Sequential(
49             nn.Linear(input_size, 128),
50             nn.BatchNorm1d(128),
51             nn.ReLU(),
52             nn.Dropout(0.2),
53             nn.Linear(128, 64),
54             nn.BatchNorm1d(64),
55             nn.ReLU(),
56             nn.Dropout(0.2),
57             nn.Linear(64, 1)
58         )
59
60     def forward(self, x):
61         return self.fc(x)
62
63 model = Regression()
64
65 # 7. 定义损失函数和优化器
66 criterion = torch.nn.MSELoss()
67 optimizer = optim.Adam(model.parameters(), lr=0.001, weight_decay=0.001)
68
69 # 8. 训练模型：自定义训练次数
70 for epoch in range(100):
71     model.train()
72     for inputs, labels in train_loader:
73         optimizer.zero_grad()
74         outputs = model(inputs)
75         loss = criterion(outputs, labels)
76         loss.backward()
77         optimizer.step()
78     if (epoch+1) % 10 == 0:
79         print(f'Epoch [{epoch+1}/{100}], Loss: {loss.item():.4f}')
80
81 # 9. 模型评估
82 model.eval()
83 total_loss = 0
84 with torch.no_grad():

```

```

85     for inputs, labels in test_loader:
86         outputs = model(inputs)
87         total_loss += criterion(outputs, labels).item() * inputs.size(0)
88
89 mse = total_loss / len(test_dataset)
90
91 print(f'\nmSE on test dataset: {mse:.4f}')

```

```

1  Epoch [10/100], Loss: 95.0033
2  Epoch [20/100], Loss: 25.6044
3  Epoch [30/100], Loss: 28.5615
4  Epoch [40/100], Loss: 13.7481
5  Epoch [50/100], Loss: 125.9341
6  Epoch [60/100], Loss: 72.2138
7  Epoch [70/100], Loss: 29.9736
8  Epoch [80/100], Loss: 95.5638
9  Epoch [90/100], Loss: 20.1719
10 Epoch [100/100], Loss: 176.2024
11
12 MSE on test dataset: 11.9592

```

4 使用早停与模型保存的训练方法

对于模型训练，固定的 `epoch` 使得最终训练得到的模型可能并不是最佳模型。`epoch` 过多可能导致模型记忆噪声，造成过拟合。`epoch` 不足时模型未充分学习，造成欠拟合。因此，采用较大的 `epoch` 结合早停策略，可以在模型性能表现最佳时将模型保存下来。

早停策略的参数设置如下：

`best_val_loss`：初始化为正无穷大，用于记录验证集上的最小损失。

`patience`：设定为20，表示当验证集损失在连续20个 `epoch` 中没有下降时，触发早停机制。

`counter`：用于记录验证集损失没有下降的连续 `epoch` 数量，初始化为0。

如下结果所示，最佳训练结果在测试集上的均方误差进一步降低到6.0042。

```

1  import torch
2  import pandas as pd
3  import torch.nn as nn
4  import torch.optim as optim
5  from torch.utils.data import DataLoader, TensorDataset
6  from sklearn.preprocessing import StandardScaler
7  from sklearn.model_selection import train_test_split
8
9  # 1.从.xlsx文件中读取数据集
10 data = pd.read_excel('BostonHousingData.xlsx', sheet_name='Sheet1')
11
12 # 2.计算相关系数并通过设置的阈值来筛选数据
13 threshold = 0
14 corre = data.corr()
15 correlation = corre['MEDV'].sort_values(ascending=False)
16 selected_features = correlation[abs(correlation) >=
17     threshold].index.tolist()
18
19 selected_data = data[selected_features]
20 x = selected_data.drop('MEDV', axis=1).values
21 y = selected_data['MEDV'].values.reshape(-1, 1)

```

```

21
22 # 3.划分训练集和验证集、测试集
23 X_train_val, X_test, y_train_val, y_test = train_test_split(X, y,
24     test_size=0.1, random_state=42)
25
26 # 4.数据标准化处理、转换格式为pytorch张量
27 scaler = StandardScaler()
28 X_train_scaled = scaler.fit_transform(X_train)
29 X_val_scaled = scaler.transform(X_val)
30 X_test_scaled = scaler.transform(X_test)
31
32 X_train_tensor = torch.tensor(X_train_scaled, dtype=torch.float32)
33 y_train_tensor = torch.tensor(y_train, dtype=torch.float32)
34 X_val_tensor = torch.tensor(X_val_scaled, dtype=torch.float32)
35 y_val_tensor = torch.tensor(y_val, dtype=torch.float32)
36 X_test_tensor = torch.tensor(X_test_scaled, dtype=torch.float32)
37 y_test_tensor = torch.tensor(y_test, dtype=torch.float32)
38
39 # 5.使用库函数加载数据集
40 train_dataset = TensorDataset(X_train_tensor, y_train_tensor)
41 train_loader = DataLoader(train_dataset, batch_size=64, shuffle=True)
42 val_dataset = TensorDataset(X_val_tensor, y_val_tensor)
43 val_loader = DataLoader(val_dataset, batch_size=64, shuffle=False)
44 test_dataset = TensorDataset(X_test_tensor, y_test_tensor)
45 test_loader = DataLoader(test_dataset, batch_size=64, shuffle=False)
46
47 # 定义输入数据维度，与阈值设置有关
48 input_size = selected_data.shape[1] - 1
49
50 # 6.全连接神经网络模型的定义
51 class Regression(nn.Module):
52     def __init__(self):
53         super().__init__()
54         self.fc = nn.Sequential(
55             nn.Linear(input_size, 128),
56             nn.BatchNorm1d(128),
57             nn.ReLU(),
58             nn.Dropout(0.2),
59             nn.Linear(128, 64),
60             nn.BatchNorm1d(64),
61             nn.ReLU(),
62             nn.Dropout(0.2),
63             nn.Linear(64, 1)
64         )
65
66     def forward(self, x):
67         return self.fc(x)
68
69 model = Regression()
70
71 # 7. 定义损失函数和优化器
72 criterion = torch.nn.MSELoss()
73 optimizer = optim.Adam(model.parameters(), lr=0.0005, weight_decay=0.0001)

```

```

74
75 # 8. 早停策略参数
76 best_val_loss = float('inf')
77 patience = 20
78 counter = 0
79
80 # 9. 训练模型：自定义训练次数
81 num_epochs = 500
82 for epoch in range(num_epochs):
83     model.train()
84     for inputs, labels in train_loader:
85         optimizer.zero_grad()
86         outputs = model(inputs)
87         loss = criterion(outputs, labels)
88         loss.backward()
89         optimizer.step()
90
91     # 验证集评估
92     model.eval()
93     val_loss = 0
94     with torch.no_grad():
95         for inputs, labels in val_loader:
96             outputs = model(inputs)
97             val_loss += criterion(outputs, labels).item() * inputs.size(0)
98     val_loss /= len(val_dataset)
99
100     if val_loss < best_val_loss:
101         best_val_loss = val_loss
102         torch.save(model.state_dict(), 'model.pth')
103         counter = 0
104     else:
105         counter += 1
106         if counter >= patience:
107             print(f'Early stopping at epoch {epoch+1}')
108             break
109
110     if (epoch+1) % 10 == 0:
111         print(f'Epoch [{epoch+1}/{num_epochs}], Train Loss:
112 {loss.item():.4f}, Val Loss: {val_loss:.4f}')
113
114 # 10. 加载最佳模型
115 model.load_state_dict(torch.load('model.pth', weights_only = True))
116
117 # 11. 模型评估
118 model.eval()
119 total_loss = 0
120 with torch.no_grad():
121     for inputs, labels in test_loader:
122         outputs = model(inputs)
123         total_loss += criterion(outputs, labels).item() * inputs.size(0)
124
125 mse = total_loss / len(test_dataset)
126
127 print(f'\nMSE on test dataset: {mse:.4f}')

```

```
1 Epoch [10/500], Train Loss: 417.8493, Val Loss: 545.7247
2 Epoch [20/500], Train Loss: 348.2777, Val Loss: 461.9928
3 Epoch [30/500], Train Loss: 365.9896, Val Loss: 386.1222
4 Epoch [40/500], Train Loss: 420.8316, Val Loss: 315.0311
5 Epoch [50/500], Train Loss: 207.3795, Val Loss: 231.4717
6 Epoch [60/500], Train Loss: 244.2944, Val Loss: 179.0083
7 Epoch [70/500], Train Loss: 193.6992, Val Loss: 133.3739
8 Epoch [80/500], Train Loss: 136.5826, Val Loss: 88.3605
9 Epoch [90/500], Train Loss: 106.8433, Val Loss: 60.4388
10 Epoch [100/500], Train Loss: 29.5513, Val Loss: 36.1021
11 Epoch [110/500], Train Loss: 33.9615, Val Loss: 20.6541
12 Epoch [120/500], Train Loss: 21.0119, Val Loss: 16.6938
13 Epoch [130/500], Train Loss: 15.2447, Val Loss: 10.4564
14 Epoch [140/500], Train Loss: 23.1641, Val Loss: 9.5385
15 Epoch [150/500], Train Loss: 61.7226, Val Loss: 8.8601
16 Epoch [160/500], Train Loss: 14.9008, Val Loss: 8.6332
17 Epoch [170/500], Train Loss: 15.9442, Val Loss: 8.9074
18 Epoch [180/500], Train Loss: 34.8720, Val Loss: 7.9747
19 Early stopping at epoch 186
20
21 MSE on test dataset: 6.0042
```