

基于RNN的行为评分卡模型实战

原创 梅子行 大数据风控与机器学习 2019-11-07



背景：因风控场景下有大量的具有时间先后顺序的数据。近期使用的较为广泛的一种建模方式，就是先使用一个针对时序数据进行建模的模型，输出时序预测分数，然后将分数应用于评分卡中作为一个特征。

应用场景：主要应用于时序数据较多的 B 卡中。

数据：复贷超过3次的B卡客户，以月份为切片构造基础特征向量，历史数据超过2年。

数据扩充：由于负样本不足5万。为得到更好的模型效果，使用pd5代替M1+。将观察期介于2~3年的客户首尾相连，进行随机截断，以获取更多的24维向量客户。再将观察期介于1~2年之间的客群首尾相连，通过prophet算法前向延伸，作为辅助训练样本。

选用模型：2层LSTM模型（示例，目前最深不超过6层）

深度学习框架：pytorch

使用方法：特征带入 LSTM 中输出 score 或隐层参数，分别作为单维度特征，一同带入逻辑回归中进行建模。



环境加载。

```
1 import torch
2 import torch.nn as nn
3 import random
4 from sklearn.model_selection import train_test_split
5 import torchvision.transforms as transforms
6 import torchvision.datasets as datasets
7 from torch.autograd import Variable
```

数据加载。

```
1 random_st = random.choice(range(10000))
2 trainB, testB = train_test_split(trainB, test_size=0.25,
3                                   random_state=random_st)
4
5 train_data = MyDataset(trainB)
6 test_data = MyDataset(testB)
7
8 train_loader = torch.utils.data.DataLoader(train_data, batch_size=50,
9                                              shuffle=True, num_workers=0)
10 test_loader = torch.utils.data.DataLoader(test_data, batch_size=25,
11                                            shuffle=False, num_workers=0)
```

定义网络。

```
1 #搭建LSTM网络
2 class Rnn(nn.Module):
3     def __init__(self, in_dim, hidden_dim, n_layer, n_class):
4         super(Rnn, self).__init__()
5         self.n_layer = n_layer
6         self.hidden_dim = hidden_dim
7         self.LSTM = nn.LSTM(in_dim, hidden_dim,
8                               n_layer, batch_first=True)
9         self.linear = nn.Linear(hidden_dim, n_class)
10        self.sigmoid = nn.Sigmoid()
11
12
13    def forward(self, x):
```

```
14         x = x.sum(dim = 1)
15         out, _ = self.LSTM(x)
16         out = out[:, -1, :]
17         out = self.linear(out)
18         out = self.sigmoid(out)
19         return out
```

指定参数。

```
1  #指定网络参数。
2
3  #28个特征, 42个月切片, 2个隐层, 2分类
4  model = Rnn(28,42,2,2)
5  device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
6  model = model.to(device)
7
8  #使用二分类对数损失函数
9
10 criterion = nn.SoftMarginLoss(reduction='mean')
11 opt = torch.optim.Adam(model.parameters())
12 total_step = len(train_loader)
13 total_step_test = len(test_loader)
14
15 num_epochs = 50
```

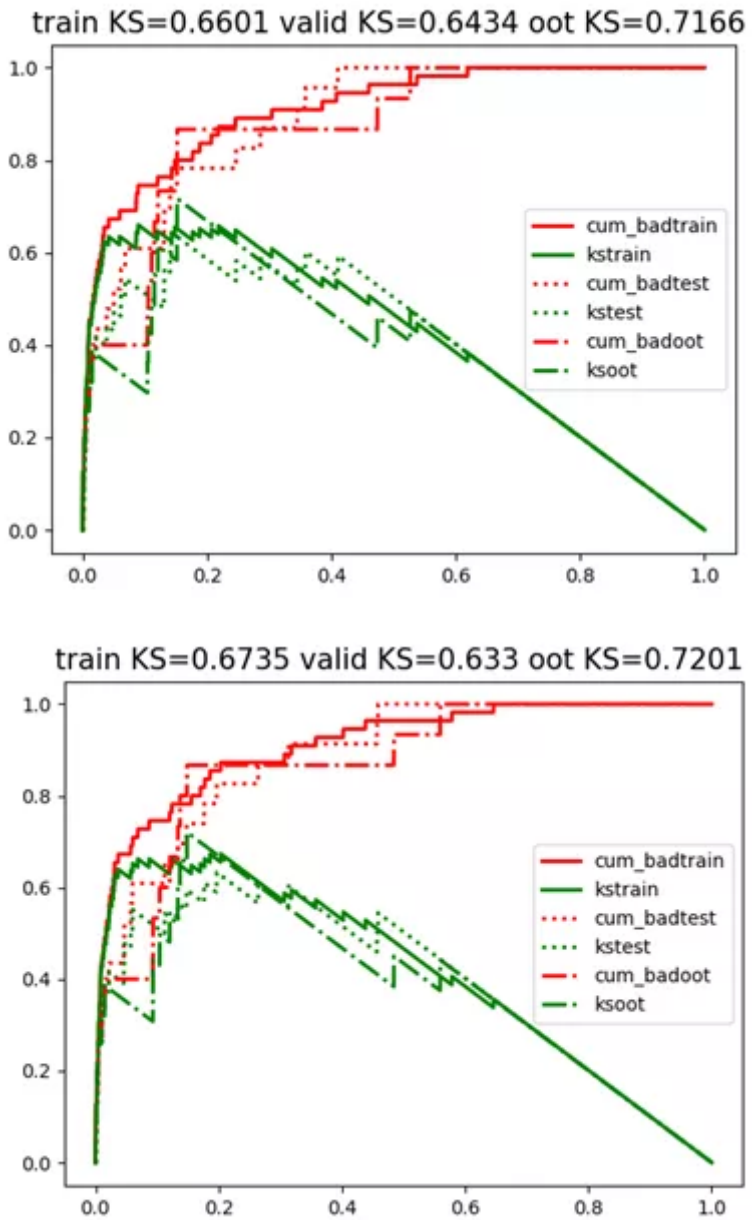
训练网络。

```
1  #指定网络参数。
2
3  #28个特征, 42个月切片, 2个隐层, 2分类
4  model = Rnn(28,42,2,2)
5  device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
6  model = model.to(device)
7
8  #使用二分类对数损失函数
9
10 criterion = nn.SoftMarginLoss(reduction='mean')
11 opt = torch.optim.Adam(model.parameters())
12 total_step = len(train_loader)
13 total_step_test = len(test_loader)
14
15 num_epochs = 50
```

```
14 训练网络。
15
16 #训练得到LSTM模型并计算单模型的KS值和AUC值。
17 for epoch in range(num_epochs):
18     train_label = []
19     train_pred = []
20     model.train()
21     for i, (B_features, labels) in enumerate(train_loader):
22         B_features = B_features.to(device)
23         labels = labels.to(device)
24         #网络训练
25         out = model(B_features)
26         loss = criterion(out, labels)
27         opt.zero_grad()
28         loss.backward()
29         opt.step()
30         #每一百轮打印一次
31         if i%100 == 0:
32             print('train epoch: {}/{}'.format(epoch + 1, num_epochs),
33                   loss: {}'.format(loss),
34                   i + 1, total_step, loss))
35         #真实标记和预测值
36         train_label.extend(labels.cpu().numpy().flatten().tolist())
37         train_pred.extend(out.detach().cpu().numpy().flatten().tolist())
38     #计算真正率和假正率
39     fpr_lm_train, tpr_lm_train, _ = roc_curve(np.array(train_label),
40                                               np.array(train_pred))
41     #计算KS和AUC
42     print('train epoch: {}/{}'.format(epoch + 1, num_epochs),
43           KS: {}, ROC: {}'.format(
44             epoch + 1, num_epochs, abs(fpr_lm_train - tpr_lm_train).max(),
45             metrics.auc(fpr_lm_train, tpr_lm_train)))
46
47     test_label = []
48     test_pred = []
49
50     model.eval()
51     #计算测试集上的KS值和AUC值
52
53     for i, (B_features, labels) in enumerate(test_loader):
54
55         B_features = B_features.to(device)
```

```
54     labels = labels.to(device)
55     out = model(B_features)
56     loss = criterion(out, labels)
57
58     # 计算KS和AUC
59     if i%100 == 0:
60         print('test epoch: {}/{}, round: {}/{},
61               loss: {}'.format(epoch + 1, num_epochs,
62                                i + 1, total_step_test, loss))
63         test_label.extend(labels.cpu().numpy().flatten().tolist())
64         test_pred.extend(out.detach().cpu().numpy().flatten().tolist())
65
66         fpr_lm_test, tpr_lm_test, _ = roc_curve(np.array(test_label),
67                                                  np.array(test_pred))
68
69         print('test epoch: {}/{}, KS: {}, ROC: {}'.format(
70               epoch + 1, num_epochs,
71               abs(fpr_lm_test -
```

最终模型结果。



[阅读原文](#)

喜欢此内容的人还喜欢

教育部答复“防止男性青少年女性化”提案：加强体育教师配备

体育老师

世袭罔替“作二代”？中国作协副主席之女，以屎尿作诗誉满文坛！

红色文化网