

```

1 # =====
2 # SET-UP AND IMPORTS
3 # =====
4 import os
5 import torch
6 import torch.nn as nn
7 import numpy as np
8 import pandas as pd
9 import matplotlib.pyplot as plt
10 from torch.utils.data import DataLoader, TensorDataset, ConcatDataset
11 from sklearn.model_selection import KFold
12 from sklearn.preprocessing import StandardScaler
13
14 # Set directory and device setup
15 os.chdir('/home/dan/FUNDRAISING')
16 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
17 print(f"Using device: {device}")
18
19 # =====
20 # PREPARE DATA
21 # =====
22 df = pd.read_csv('./data/df.csv', index_col='Unnamed: 0')
23 df['zipconvert'] = df[['zipconvert2', 'zipconvert3', 'zipconvert4', 'zipconvert5']].apply(
24     lambda row: 'zc' + str(row.idxmax()[-1]) if pd.notna(row.idxmax()) else 'unknown', axis=1
25 )
26 test = df[df['type']=='test']
27
28 df = df[df['type']!='test']
29
30 df.drop(['zipconvert2', 'zipconvert3', 'zipconvert4', 'zipconvert5', 'type'], axis=1, inplace=True)
31 test.drop(['zipconvert2', 'zipconvert3', 'zipconvert4', 'zipconvert5', 'type'], axis=1, inplace=True)
32
33 cat_cols = ['homeowner', 'female', 'zipconvert', 'wealth', 'income', 'num_child']
34 cont_cols = [col for col in df.columns if col not in cat_cols + ['target']]
35 emb_sizes = [(df[col].astype('category').cat.codes.max() + 1, min(50, (df[col].nunique() + 1) // 2)) for col in cat_cols]
36
37 # Function to process datasets
38 def process_data(data):
39     cats = np.stack([data[col].astype('category').cat.codes.values for col in cat_cols], axis=1)
40     conts = np.stack([data[col].values for col in cont_cols], axis=1)
41     scaler = StandardScaler()
42
43     conts = scaler.fit_transform(conts)
44     y = data['target'].map({'Donor': 1, 'No Donor': 0}).values
45     return torch.tensor(cats, dtype=torch.int64), torch.tensor(conts, dtype=torch.float), torch.tensor(y, dtype=torch.long)
46
47 cats, conts, targets = process_data(df)
48 test_cats, test_conts, test_targets = process_data(test)
49
50 dataset = TensorDataset(cats, conts, targets)
51
52 test_dataset = TensorDataset(test_cats, test_conts, test_targets)
53 test_loader = DataLoader(test_dataset, batch_size=2048, shuffle=False)
54
55 # =====
56 # DEFINE A TABULAR MODEL
57 # =====
58
59 class TabularModel(nn.Module):
60     def __init__(self, emb_sizes, n_cont, out_sz, layers, p=0.5):
61         super().__init__()
62         self.embs = nn.ModuleList([nn.Embedding(ni, nf) for ni, nf in emb_sizes])
63         self.emb_drop = nn.Dropout(p)
64         self.bn_cont = nn.BatchNorm1d(n_cont)
65
66         # Calculate the total embedding output size
67         total_emb_size = sum(nf for _, nf in emb_sizes)
68         n_in = total_emb_size + n_cont # Total input size for the first linear layer
69
70         layerlist = nn.ModuleList()
71         for output_size in layers:
72             layerlist.append(nn.Linear(n_in, output_size))
73             layerlist.append(nn.ReLU())
74             layerlist.append(nn.BatchNorm1d(output_size))
75             layerlist.append(nn.Dropout(p))
76             n_in = output_size # Update n_in to the output size of the current layer
77
78         layerlist.append(nn.Linear(layers[-1], out_sz)) # Final output layer
79         self.layers = nn.Sequential(*layerlist)
80
81     def forward(self, x_cat, x_cont):
82         embeddings = [e(x_cat[:, i]) for i, e in enumerate(self.embs)]
83         x = torch.cat(embeddings, 1)
84         x = self.emb_drop(x)

```

```

85         x_cont = self.bn_cont(x_cont)
86         x = torch.cat([x, x_cont], 1)
87         return self.layers(x)
88
89 # Initialize model
90 model = TabularModel(emb_sizes, len(cont_cols), 2, [200, 100], p=0.4)
91 model = model.to(device)
92
93
94 #endregion
95 #region # EARLY STOPPING
96 # =====
97 # EARLY STOPPING
98 # =====
99 class EarlyStopping:
100     def __init__(self, patience=5, verbose=False, delta=0):
101         self.patience = patience
102         self.verbose = verbose
103         self.delta = delta
104         self.best_score = None
105         self.early_stop = False
106         self.counter = 0
107
108     def __call__(self, val_loss, model):
109         score = -val_loss
110
111         if self.best_score is None:
112             self.best_score = score
113         elif score < self.best_score + self.delta:
114             self.counter += 1
115             if self.verbose:
116                 print(f'EarlyStopping counter: {self.counter} out of {self.patience}')
117             if self.counter >= self.patience:
118                 self.early_stop = True
119         else:
120             self.best_score = score
121             self.counter = 0
122
123
124 # =====
125 # CROSS-VALIDATION SETUP
126 # =====
127 def calculate_accuracy(y_pred, y_true):
128     y_pred_classes = torch.argmax(y_pred, dim=1)
129     correct = (y_pred_classes == y_true).float() # convert into float for division
130     acc = correct.sum() / len(correct)
131     return acc
132
133 kf = KFold(n_splits=20, shuffle=True, random_state=42)
134 early_stopping = EarlyStopping(patience=5, verbose=True)
135 results = []
136
137 for fold, (train_idx, val_idx) in enumerate(kf.split(dataset)):
138     train_subsampler = torch.utils.data.SubsetRandomSampler(train_idx)
139     val_subsampler = torch.utils.data.SubsetRandomSampler(val_idx)
140     train_loader = DataLoader(dataset, batch_size=32, sampler=train_subsampler)
141     val_loader = DataLoader(dataset, batch_size=32, sampler=val_subsampler)
142
143     model = TabularModel(emb_sizes, len(cont_cols), 2, [100, 200], p=0.4).to(device)
144     criterion = nn.CrossEntropyLoss()
145     optimizer = torch.optim.Adam(model.parameters(), lr=1e-5)
146
147     for epoch in range(50): # Adjust as needed
148         model.train()
149         total_loss, total_acc = 0, 0
150         for cats, conts, y in train_loader:
151             cats, conts, y = cats.to(device), conts.to(device), y.to(device)
152             optimizer.zero_grad()
153             outputs = model(cats, conts)
154             loss = criterion(outputs, y)
155             acc = calculate_accuracy(outputs, y)
156             total_loss += loss.item()
157             total_acc += acc.item()
158             loss.backward()
159             optimizer.step()
160         avg_train_loss = total_loss / len(train_loader)
161         avg_train_acc = total_acc / len(train_loader)
162
163         model.eval()
164         val_loss, val_acc = 0, 0
165         with torch.no_grad():
166             for cats, conts, y in val_loader:
167                 cats, conts, y = cats.to(device), conts.to(device), y.to(device)
168                 outputs = model(cats, conts)
169                 loss = criterion(outputs, y)
170                 acc = calculate_accuracy(outputs, y)
171                 val_loss += loss.item()

```

```

172         val_acc += acc.item()
173     avg_val_loss = val_loss / len(val_loader)
174     avg_val_acc = val_acc / len(val_loader)
175
176     print(f'Fold {fold+1}, Epoch {epoch+1}, Train Loss: {avg_train_loss:.4f}, Train Acc: {avg_train_acc:.4f}, Val Loss: {avg_val_loss:.4f}, Val
Acc: {avg_val_acc:.4f}')
177
178     # Call early stopping
179     early_stopping(avg_val_loss, model)
180     if early_stopping.early_stop:
181         print("Early stopping")
182         break
183
184     results.append((avg_val_loss, avg_val_acc))
185
186 average_val_loss = sum(x[0] for x in results) / len(results)
187 average_val_acc = sum(x[1] for x in results) / len(results)
188 print(f'Average Validation Loss: {average_val_loss:.4f}, Average Validation Accuracy: {average_val_acc:.4f}')
189
190 preds = []
191 model.eval()
192
193 with torch.no_grad():
194     for cats, conts, y in test_loader:
195         cats, conts, y = cats.to(device), conts.to(device), y.to(device)
196         output = model(cats, conts)
197         predicted = output.argmax(dim=1) # Ensure you use dim=1
198         # print(predicted)
199         preds.append(predicted.cpu()) # Move predictions to CPU
200 # print(type(preds[0]))
201 # Concatenate all batch predictions into a single tensor
202 # preds = torch.cat(preds)
203 preds = preds[0].tolist()
204 # print(preds)
205 final_preds = []
206 donor=1
207 no_donor=1
208 for i in preds:
209     if i == 1:
210         final_preds.append('Donor')
211         donor +=1
212     else:
213         final_preds.append('No Donor')
214         no_donor+=1
215 print(f'% Donor = {donor / (donor+no_donor)}')
216
217
218 save_df = pd.DataFrame(final_preds, columns=['values'])
219 save_df.to_csv('./preds/preds.csv', index=False)

```