

```
1 # imports
2 import pandas as pd
3 import numpy as np
4 np.random.seed(7)
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 sns.set(style="white", palette="muted", color_codes=True, context="talk")
8 from IPython import display
9 %matplotlib inline
10
11 import sklearn as sk
12 from sklearn.model_selection import train_test_split
13 from sklearn.preprocessing import StandardScaler
14 from sklearn.metrics import accuracy_score, roc_auc_score
15 from sklearn.utils.class_weight import compute_class_weight
16
17
18 import torch
19 from torch.utils.data import TensorDataset
20 from torch.utils.data import DataLoader
21 import torch.optim as optim
22 import torch.nn as nn
23 import torch.nn.functional as F
24
25
26 print("pythoch version: {}".format(torch.__version__))
27
28 device = "cpu"
```

```
⇒ pythoch version: 1.2.0
```

New Section

FAIRER MACHINE LEARNING: IS IT POSSIBLE?

ML FOR DECISION MAKING

From credit ratings to online dating, machine learning models are increasingly used to automate 'everyday' decision making processes.

GROWING IMPACT ON SOCIETY

MACHINE LEARNING & DISCRIMINATION

We need to step-up and look for ways to mitigate emergent discrimination in our models. Make sure that our predictions do not hurt people unfairly with certain sensitive characteristics (e.g., gender, ethnicity, etc.).

▼ Making income predictions

Let's start by training a basic classifier that can predict whether or not a person's income is larger than 50K dollar a year. We'll be working with some California Census Data, we'll be trying to use various features of an individual to predict what class of income they belong in ($>50k$ or $\leq 50k$). It is not hard to imagine that financial institutions train models on similar data sets and use them to decide whether or not someone is eligible for a loan, or to set the height of an insurance premium.

Before training a model, we first parse the data into three datasets: features, targets and sensitive attributes. The set of features X contains the input attributes that the model uses for making the predictions, with attributes like age, education level and occupation.

The targets y contain the binary class labels that the model needs to predict. These labels are $y \in \{income > 50K, income \leq 50K\}$.

Finally, the set of sensitive attributes Z contains the attributes for which we want the prediction to fair. These are $z_{race} \in \{\text{black}, \text{white}\}$ and $z_{gender} \in \{\text{male}, \text{female}\}$. It is important to note that datasets are non-overlapping, so the sensitive attributes race and sex are not part of the features used for training the model.

Here is some information about the data:

| Column Name | Type | Description |
|----------------|-------------|---|
| age | Continuous | The age of the individual |
| workclass | Categorical | The type of employer the individual has (government, military, private, etc.). |
| fnlwgt | Continuous | The number of people the census takers believe that observation represents (sample weight). This variable will not be used. |
| education | Categorical | The highest level of education achieved for that individual. |
| education_num | Continuous | The highest level of education in numerical form. |
| marital_status | Categorical | Marital status of the individual. |
| occupation | Categorical | The occupation of the individual. |
| relationship | Categorical | Wife, Own-child, Husband, Not-in-family, Other-relative, Unmarried. |
| race | Categorical | White, Asian-Pac-Islander, Amer-Indian-Eskimo, Other, Black. |
| gender | Categorical | Female, Male. |
| capital_gain | Continuous | Capital gains recorded. |
| capital_loss | Continuous | Capital Losses recorded. |
| hours_per_week | Continuous | Hours worked per week. |
| native_country | Categorical | Country of origin of the individual. |
| income | Categorical | " $>50K$ " or " $\leq 50K$ ", meaning whether the person makes more than \$50,000 annually. |

THE DATA

Downloading the data

```
1 !wget https://archive.ics.uci.edu/ml/machine-learning-databases/adult/adult.data
→ --2019-10-16 11:49:01-- https://archive.ics.uci.edu/ml/machine-learning-databases/adult/adult.data
Resolving archive.ics.uci.edu (archive.ics.uci.edu)... 128.195.10.252
Connecting to archive.ics.uci.edu (archive.ics.uci.edu)|128.195.10.252|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 3974305 (3.8M) [application/x-httpd-php]
Saving to: 'adult.data.1'

adult.data.1      100%[=====] 3.79M 16.0MB/s   in 0.2s

2019-10-16 11:49:01 (16.0 MB/s) - 'adult.data.1' saved [3974305/3974305]
```

Read in the census_data.csv data with pandas

```
1 def load_ICU_data(path):
2     column_names = ['age', 'workclass', 'fnlwgt', 'education', 'education_num',
3                      'marital_status', 'occupation', 'relationship', 'race', 'sex',
4                      'capital_gain', 'capital_loss', 'hours_per_week', 'country', 'target']
5     input_data = (pd.read_csv(path, names=column_names,
6                               na_values=?, sep=r'\s*,\s*', engine='python')
7                   .loc[lambda df: df['race'].isin(['White', 'Black'])])
8
9     # sensitive attributes; we identify 'race' and 'sex' as sensitive attributes
10    sensitive_attribs = ['race', 'sex']
11    Z = (input_data.loc[:, sensitive_attribs]
12         .assign(race=lambda df: (df['race'] == 'White').astype(int),
13                 sex=lambda df: (df['sex'] == 'Male').astype(int)))
14
15    # targets; 1 when someone makes over 50k , otherwise 0
16    y = (input_data['target'] == '>50K').astype(int)
17
18    # features; note that the 'target' and sentive attribute columns are dropped
19    X = (input_data
20          .drop(columns=['target', 'race', 'sex'])
21          .fillna('Unknown')
22          .pipe(pd.get_dummies, drop_first=True))
23
24    print("features X: {} samples, {} attributes".format(X.shape[0], X.shape[1]))
25    print("targets y: {} samples".format(y.shape[0]))
26    print("sensitive Z: {} samples, {} attributes".format(Z.shape[0], Z.shape[1]))
```

```
26 print('sensitives Z: {} samples, {} attributes'.format(Z.shape[0], Z.shape[1]))  
27 return X, y, Z
```

```
1 # load ICU data set  
2 X, y, Z = load_ICU_data('adult.data')  
3  
4 n_features = X.shape[1]  
5 n_sensitive = Z.shape[1]
```

```
→ features X: 30940 samples, 94 attributes  
targets y: 30940 samples  
sensitives Z: 30940 samples, 2 attributes
```

Our dataset contains the information of just over 30K people. Next, we split the data into train and test sets, where the split is 50/50, and scale the features X using standard scaling.

```
1 # split into train/test set  
2 X_train, X_test, y_train, y_test, Z_train, Z_test = train_test_split(X, y, Z, test_size=0.5,  
3                                         stratify=y,  
4                                         random_state=7)  
5  
6 # standardize the data  
7 scaler = StandardScaler().fit(X_train)  
8 scale_df = lambda df, scaler: pd.DataFrame(scaler.transform(df), columns=df.columns, index=df.index)  
9 X_train = X_train.pipe(scale_df, scaler)  
10 X_test = X_test.pipe(scale_df, scaler)
```

```
1 class PandasDataSet(TensorDataset):  
2  
3     def __init__(self, *dataframes):  
4         tensors = (self._df_to_tensor(df) for df in dataframes)  
5         super(PandasDataSet, self).__init__(*tensors)  
6  
7     def _df_to_tensor(self, df):  
8         if isinstance(df, pd.Series):  
9             df = df.to_frame()  
10        return torch.from_numpy(df.values).float()  
11  
12  
13 train_data = PandasDataSet(X_train, y_train, Z_train)  
14 test_data = PandasDataSet(X_test, y_test, Z_test)
```

```
1 train_loader = DataLoader(train_data, batch_size=32, shuffle=True, drop_last=True)  
2 print('# features:', n_features)  
3 print('# training samples:', len(train_data))
```

```
4 print('# batches:', len(train_loader))
```

```
⇒ # features: 94  
# training samples: 15470  
# batches: 483
```

▼ Making income predictor

Let's train our basic income level predictor. So it predicts "the probability that this person's income is larger than 50K".

```
1 class Classifier(nn.Module):  
2  
3     def __init__(self, n_features, n_hidden=32, p_dropout=0.2):  
4         super(Classifier, self).__init__()  
5         self.network = nn.Sequential(  
6             nn.Linear(n_features, n_hidden),  
7             nn.ReLU(),  
8             nn.Dropout(p_dropout),  
9             nn.Linear(n_hidden, n_hidden),  
10            nn.ReLU(),  
11            nn.Dropout(p_dropout),  
12            nn.Linear(n_hidden, n_hidden),  
13            nn.ReLU(),  
14            nn.Dropout(p_dropout),  
15            nn.Linear(n_hidden, 1),  
16        )  
17  
18    def forward(self, x):  
19        return torch.sigmoid(self.network(x))
```

```
1 clf = Classifier(n_features=n_features)  
2 clf_criterion = nn.BCELoss()  
3 clf_optimizer = optim.Adam(clf.parameters())
```

```
1 N_CLF_EPOCHS = 2  
2  
3 for epoch in range(N_CLF_EPOCHS):  
4     for x, y, _ in train_loader:  
5         clf.zero_grad()  
6         p_y = clf(x)  
7         loss = clf_criterion(p_y, y)  
8         loss.backward()  
9         clf_optimizer.step()
```

```
1 with torch.no_grad():
2     pre_clf_test = clf(test_data.tensors[0])
3
4 y_pre_clf = pd.Series(pre_clf_test.data.numpy().ravel(), index=y_test.index)

1 iter = 1
2 val_metrics = pd.DataFrame(columns=['ROC AUC', 'Accuracy'])
3 val_metrics.loc[iter,'ROC AUC'] = roc_auc_score(y_test, y_pre_clf)
4 val_metrics.loc[iter,'Accuracy'] = 100*accuracy_score(y_test, y_pre_clf > 0.5)
5 print("ROC AUC: {:.2f}".format(val_metrics.loc[iter,'ROC AUC']))
6 print("Accuracy: {:.1f}%".format(val_metrics.loc[iter, 'Accuracy']))
```

ROC AUC: 0.90
Accuracy: 84.8%

With a ROC AUC larger than 0.90 and a prediction accuracy of 85% we can say that our basic classifier performs pretty well! However, if it is also fair in its predictions, that remains to be seen.

▼ Qualitative model fairness

We start the investigation into the fairness of our classifier by analysing the predictions it made on the test set. The plots in the figure below show the distributions of the predicted $P(\text{income} > 50K)$ given the sensitive attributes.

```

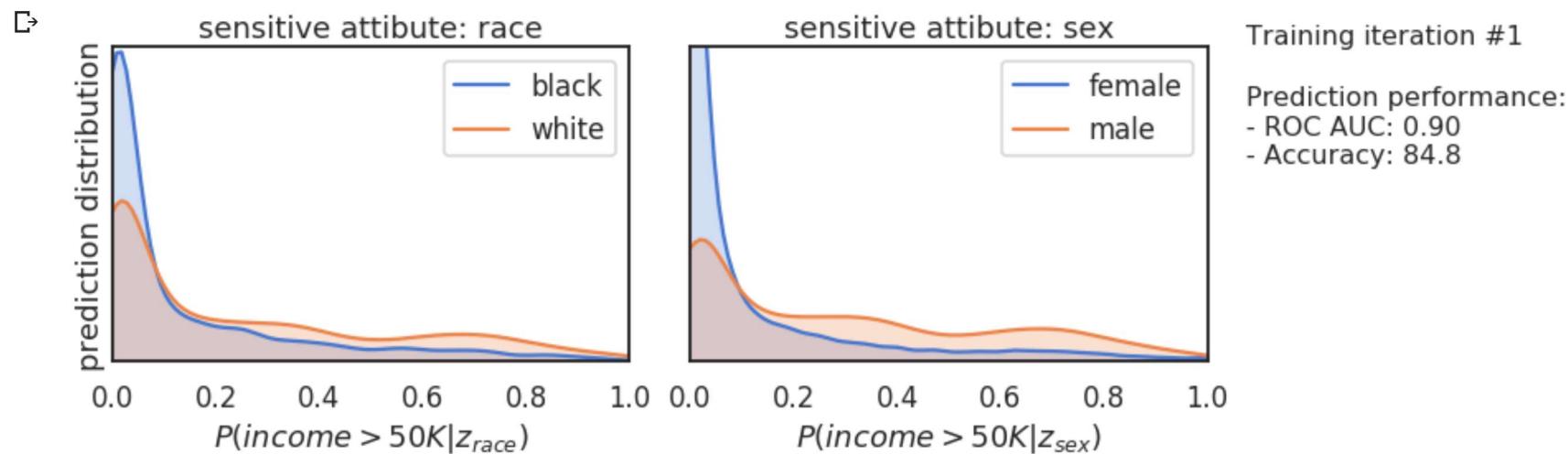
1 def plot_distributions(y, Z, iteration=None, val_metrics=None, p_rules=None, fname=None)
2     fig, axes = plt.subplots(1, 2, figsize=(10, 4), sharey=True)
3     legend={'race': ['black','white'],
4             'sex': ['female','male']}
5     for idx, attr in enumerate(Z.columns):
6         for attr_val in [0, 1]:
7             ax = sns.distplot(y[Z[attr] == attr_val], hist=False,
8                                kde_kws={'shade': True,},
9                                label='{}{}'.format(legend[attr][attr_val]),
10                               ax=axes[idx])
11             ax.set_xlim(0,1)
12             ax.set_ylim(0,7)
13             ax.set_yticks([])
14             ax.set_title("sensitive attribute: {}".format(attr))
15             if idx == 0:
16                 ax.set_ylabel('prediction distribution')
17                 ax.set_xlabel(r'$P(\{\text{income}>50K\} | z_{\{\cdot\}})$'.format(attr))
18             if iteration:
19                 fig.text(1.0, 0.9, "Training iteration #{}".format(iteration), fontsize='16')
20             if val_metrics is not None:
21                 fig.text(1.0, 0.65, '\n'.join(["Prediction performance:",
```

```

22             "- ROC AUC: {:.2f}".format(val_metrics.loc[iteration, 'ROC AUC']),
23             "- Accuracy: {:.1f}".format(val_metrics.loc[iteration, 'Accuracy'])),
24         fontsize='16')
25     if p_rules is not None:
26         fig.text(1.0, 0.4, '\n'.join(["Satisfied p%-rules:"]
27                                     [- '{0}: {1:.0f}% rule".format(attr, p_rules[attr])
28                                      for attr in p_rules.keys()]),
29         fontsize='16')
30 fig.tight_layout()
31 if fname is not None:
32     plt.savefig(fname, bbox_inches='tight')
33 return fig

```

```
1 fig = plot_distributions(y_pre_clf, z_test, iteration= iter, val_metrics= val_metrics, fname='biased_training.png')
```



▼ Quantitative model fairness

p-rules

The rule states that the ratio between the probability of a positive outcome given the sensitive attribute being true and the same probability given the sensitive attribute being false is no less than p/100

$$\min\left(\frac{P(\hat{y} = 1|z = 1)}{P(\hat{y} = 1|z = 0)}, \frac{P(\hat{y} = 1|z = 0)}{P(\hat{y} = 1|z = 1)}\right) \geq \frac{p}{100}$$

```

1 def p_rule(y_pred, z_values, threshold=0.5):
2     y_z_1 = y_pred[z_values == 1] > threshold if threshold else y_pred[z_values == 1]
3     y_z_0 = y_pred[z_values == 0] > threshold if threshold else y_pred[z_values == 0]
4     odds = y_z_1.mean() / y_z_0.mean()

```

```

5     return np.min([odds, 1/odds]) * 100

1 print("The classifier satisfies the following %p-rules:")
2 print("\tgiven attribute race; {:.0f}%‐rule".format(p_rule(y_pre_clf, z_test['race'], threshold= 0.5)))
3 print("\tgiven attribute sex;  {:.0f}%‐rule".format(p_rule(y_pre_clf, z_test['sex'], threshold= 0.5)))

```

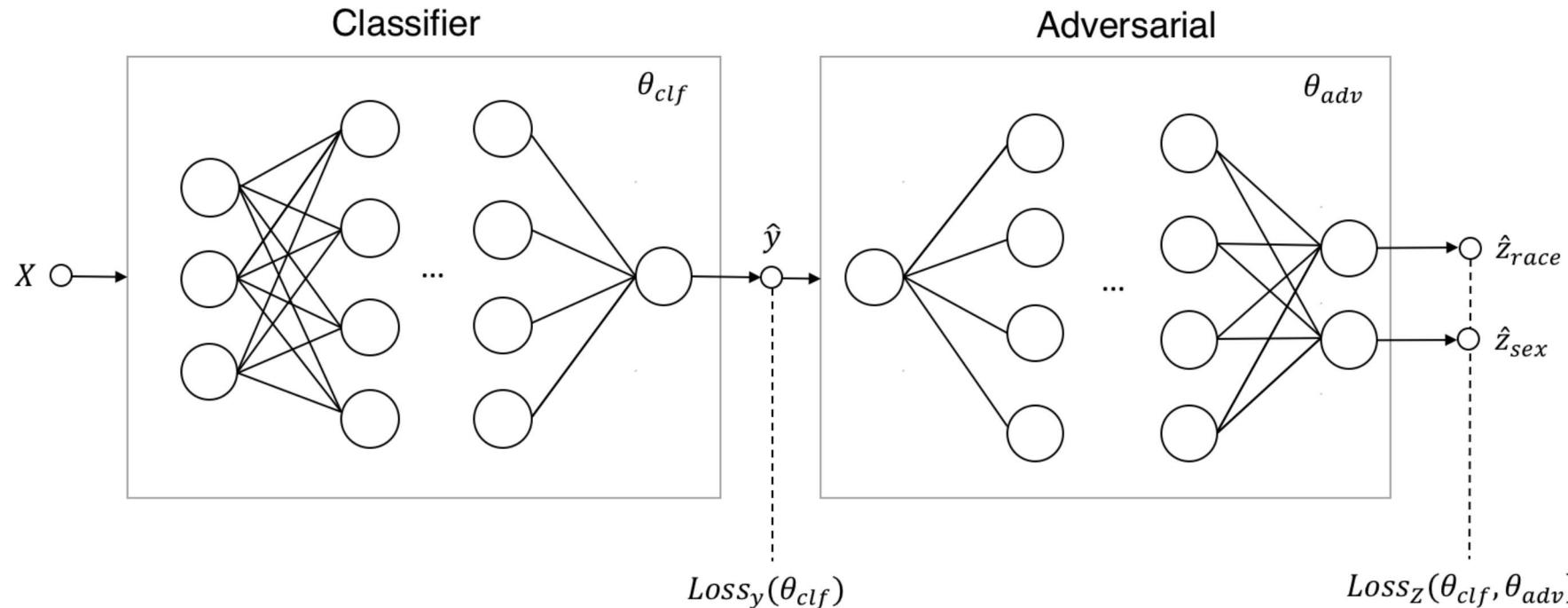
⇒ The classifier satisfies the following %p-rules:
 given attribute race; 42%-rule
 given attribute sex; 33%-rule

Both sensitive attributes have p%-rule significantly lower than 80%: The trained classifier is **unfair** in making its predictions

▼ Fair income predictions

ADVERSARIAL TRAINING PROCEDURE

- First generative model generates actual predictions \hat{y} (the probability a person's income is larger than 50K) based on the input X .
- Second adversarial classifier predicts the sensitive attribute values ($\text{race} \in \{\text{black, white}\}$ and $\text{sex} \in \{\text{male, female}\}$) from the predicted \hat{y} of the classifier.
- Finally, the objectives that both nets try to optimize are based on the prediction losses of the target and sensitive attributes



▼ FAIRCLASSIFIER

For the classifier the objective of twofold: make the best possible income level predictions at the same time as ensuring that sensitive attributes cannot be derived from them. This is captured by the following objective function:

$$\min_{\theta_{clf}} [Loss_y(\theta_{clf}) - \lambda Loss_Z(\theta_{clf}, \theta_{adv})].$$

Note that increasing the size of λ steers the classifier towards fairer predictions while sacrificing prediction accuracy.

```
1 class FairClassifier(nn.Module):
2     def __init__(self, n_features, n_hidden= 32):
3         super(FairClassifier, self).__init__()
4         self.nn= nn.Sequential(
5             nn.Linear(in_features= n_features, out_features= n_hidden),
6             nn.ReLU(),
7             nn.Dropout(p= 0.2),
8             nn.Linear(in_features= n_hidden, out_features= n_hidden),
9             nn.ReLU(),
10            nn.Dropout(p= 0.2),
11            nn.Linear(in_features= n_hidden, out_features= n_hidden),
12            nn.ReLU(),
13            nn.Dropout(p= 0.2),
14            nn.Linear(in_features= n_hidden, out_features= 1)
15        )
16        self.criterion = nn.BCELoss()
17        self.optimizer = optim.Adam(self.nn.parameters())
18
19    def forward(self, x):
20        return torch.sigmoid(self.nn(x))
21
22    def calculate_loss(self, adv_nn, p_y, y, p_z, z, lambdas):
23        return self.criterion(p_y, y) - (adv_nn.criterion(p_z, z) * lambdas).mean()
24
25    def pretrain(self, data_loader):
26        for x, y, _ in data_loader:
27            self.zero_grad()
28            p_y = self(x)
29            loss = self.criterion(p_y, y)
30            loss.backward()
31            self.optimizer.step()
32
33    def train_single_batch(self, data_loader, adv_nn, lambdas):
34        for x, y, z in data_loader:
35            pass
36            self.zero_grad()
37            p_y = self(x)
38            p_z = adv_nn(p_y).detach()
39            loss_adv = adv_nn.calculate_loss(p_z, z, lambdas)
```

```

39     loss_adv = adv_nn.calculate_loss(p_z, z, lambdas)
40     loss= self.calculate_loss(adv_nn, p_y, y, p_z, z, lambdas)
41     loss.backward()
42     self.optimizer.step()
43
44 def train(self, data_loader, adv_nn, lambdas):
45     for x, y, z in data_loader:
46         self.zero_grad()
47         p_y = self(x)
48         p_z = adv_nn(p_y)#.detach()
49         loss_adv = adv_nn.calculate_loss(p_z, z, lambdas)
50         loss= self.calculate_loss(adv_nn, p_y, y, p_z, z, lambdas)
51         loss.backward()
52         self.optimizer.step()
53

```

▼ FAIR ADVERSARIAL

For the adversarial: predict race and sex based on the income level predictions of the classifier. This is captured in the following objective function:

$$\min_{\theta_{adv}} [Loss_Z(\theta_{clf}, \theta_{adv})].$$

The adversarial does not care about the prediction accuracy of the classifier. It is only concerned with minimizing its own prediction losses of sensitive attributes

```

1 class FairAdversarial(nn.Module):
2     def __init__(self, n_sensitive, n_hidden= 32):
3         super(FairAdversarial, self).__init__()
4         self.nn= nn.Sequential(
5             nn.Linear(in_features= 1, out_features= n_hidden),
6             nn.ReLU(),
7             nn.Linear(in_features= n_hidden, out_features= n_hidden),
8             nn.ReLU(),
9             nn.Linear(in_features= n_hidden, out_features= n_hidden),
10            nn.ReLU(),
11            nn.Linear(in_features= n_hidden, out_features= n_sensitive)
12        )
13        self.criterion = nn.BCELoss(reduction= 'elementwise_mean')
14        self.optimizer = optim.Adam(self.nn.parameters())
15
16    def forward(self, x):
17        return torch.sigmoid(self.nn(x))
18
19    def calculate_loss(self, p_z, z, lambdas):
20        return (self.criterion(p_z, z) * lambdas).mean()
21
22    def pretrain(self, data_loader, clf_nn, lambdas):
23        for x, y, z in data_loader:

```

```
23     for x, _, z in data_loader:
24         self.zero_grad()
25         p_y = clf_nn(x)#.detach()
26         p_z = self(p_y)
27         loss = (self.criterion(p_z, z) * lambdas).mean()
28         loss.backward()
29         self.optimizer.step()
30
31 def train_single_batch(self, data_loader, clf_nn, lambdas):
32     for x, y, z in data_loader:
33         pass
34     self.zero_grad()
35     p_y = clf_nn(x)#.detach()
36     p_z = self(p_y)
37     loss = self.calculate_loss(p_z, z, lambdas)
38     loss.backward()
39     self.optimizer.step()
40
41 def train(self, data_loader, clf_nn, lambdas):
42     for x, y, z in data_loader:
43         self.zero_grad()
44         p_y = clf_nn(x)#.detach()
45         p_z = self(p_y)
46         loss = self.calculate_loss(p_z, z, lambdas)
47         loss.backward()
48         self.optimizer.step()
```

```
1 lambdas = torch.Tensor([30, 130])
2 fair_clf = FairClassifier(n_features, 32)
3 adv = FairAdversarial(n_sensitive, 32)
```

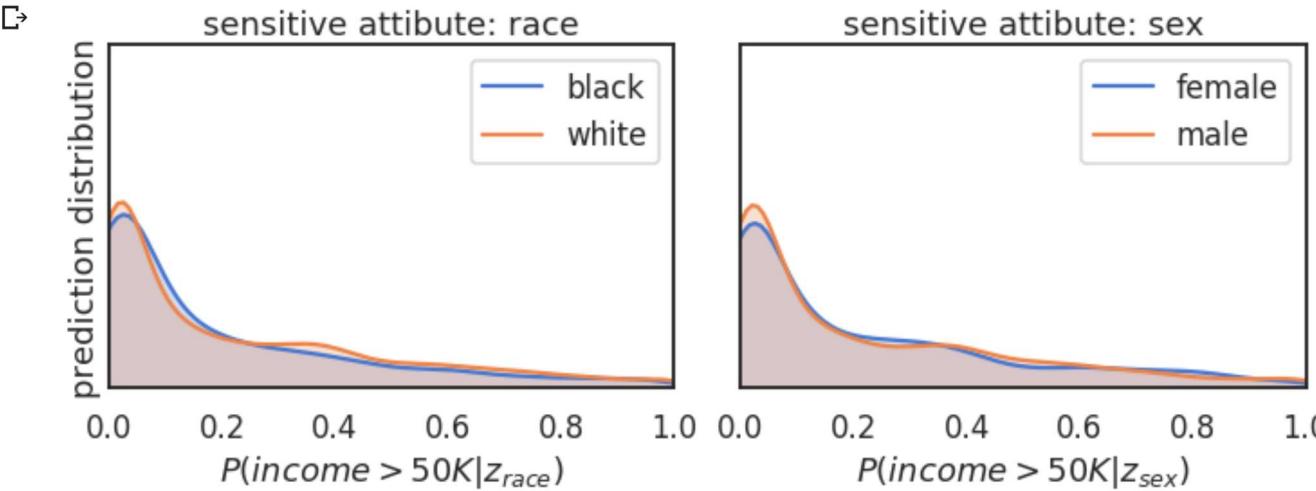
```
1 def train_GAN(clf, adv, data_loader, lambdas):
2     adv.train(data_loader, fair_clf, lambdas)
3     clf.train_single_batch(data_loader, adv, lambdas)
4     return clf, adv
```

```
1 N_CLF_EPOCHS = 2
2 for epoch in range(N_CLF_EPOCHS):
3     fair_clf.pretrain(train_loader)
4
5 N_ADV_EPOCHS = 5
6 for i in range(N_ADV_EPOCHS):
7     adv.pretrain(train_loader, clf_nn= fair_clf, lambdas= lambdas)
8
9 N_EPOCH_COMBINED = 80
```

```

10
11 performance_metrics = pd.DataFrame(columns=['ROC AUC', 'Accuracy'])
12 fairness_metrics = {}
13
14 for epoch in range(1, N_EPOCH_COMBINED):
15     fair_clf, adv = train_GAN(fair_clf, adv, train_loader, lambdas)
16     with torch.no_grad():
17         clf_pred = fair_clf(test_data.tensors[0])
18         adv_pred = adv(clf_pred)
19
20     y_post_clf = pd.Series(clf_pred.data.numpy().ravel(), index=y_test.index)
21     Z_post_adv = pd.DataFrame(adv_pred.data.numpy(), columns=Z_test.columns, index=Z_test.index)
22
23     performance_metrics.loc[epoch,'ROC AUC'] = roc_auc_score(y_test, y_post_clf)
24     performance_metrics.loc[epoch,'Accuracy'] = 100*accuracy_score(y_test, y_post_clf > 0.5 )
25     for sensitive_attr in Z_test.columns:
26         fairness_metrics[sensitive_attr] = p_rule(y_post_clf, Z_test[sensitive_attr], threshold= 0.5)
27
28     fig = plot_distributions(y_post_clf, Z_test, iteration= epoch,
29                             val_metrics= performance_metrics, p_rules= fairness_metrics)
30     display.clear_output(wait=True)
31     plt.savefig('torch_{:08d}.png'.format(epoch + 1), bbox_inches='tight')
32     plt.show(plt.gcf())

```



Training iteration #79

Prediction performance:
- ROC AUC: 0.82
- Accuracy: 81.0

Satisfied p%-rules:
- race: 78%-rule
- sex: 99%-rule

