

```
In [10]: 1 import random
2 import math
3
4 import pandas as pd
5 pd.options.mode.chained_assignment = None # default='warn'
6
7 import numpy as np
8 import scipy as sp
9 import matplotlib.pyplot as plt
10 from matplotlib.lines import Line2D
11 from sklearn.preprocessing import KBinsDiscretizer
12 from sklearn.tree import DecisionTreeClassifier
13
14
15 import seaborn as sns
16 from scipy.stats import pearsonr
17 from itertools import cycle
```

Step 1:

Which dataset did you select?

We chose the Student Performance Data Set (Kaggle Description: This data approach student achievement in secondary education of two Portuguese schools) from

<https://www.kaggle.com/larsen0966/student-performance-data-set?select=student-por.csv>
(<https://www.kaggle.com/larsen0966/student-performance-data-set?select=student-por.csv>)

Which regulated domain does your dataset belong to?

The dataset belongs to the Education domain.

How many observations are in the dataset?

There are 649 observations.

How many variables in the dataset?

There are 33 variables in the dataset.

Which variables did you select as your dependent variables?

We selected the grades related to Portuguese: G1 - first-period grade (numeric: from 0 to 20) G2 - second-period grade (numeric: from 0 to 20) G3 - final grade (numeric: from 0 to 20, output target)

How many and which variables in the dataset are associated with a legally recognized protected class?

There are two variables associated with legally recognized protected classes: age and sex.

Which legal precedence/law (as discussed in the lectures) does each protected class fall under?

Age: Age Discrimination in Employment Act of 1967 Sex: Equal Pay Act of 1963; Civil Rights Act of 1964, 1991

In [4]:

```
1
2 def STEP1():
3     df = pd.read_csv('dataset/student-por.csv', delimiter=',')
4     print(df.head(5))
5     print("Number of records port: ", len(df))
6     r, c = df.shape
7     print("rows port: ", r)
8     print("columns port: ", c)
9     print()
10    for col in df.columns:
11        print(col)
12 STEP1()
```

	school	sex	age	address	famsize	Pstatus	Medu	Fedu	Mjob	Fjob
...	\									
0	GP	F	18	U	GT3	A	4	4	at_home	teacher
...										
1	GP	F	17	U	GT3	T	1	1	at_home	other
...										
2	GP	F	15	U	LE3	T	1	1	at_home	other
...										
3	GP	F	15	U	GT3	T	4	2	health	services
...										
4	GP	F	16	U	GT3	T	3	3	other	other
...										

	famrel	freetime	goout	Dalc	Walc	health	absences	G1	G2	G3
0	4	3	4	1	1	3	4	0	11	11
1	5	3	3	1	1	3	2	9	11	11
2	4	3	2	2	3	3	6	12	13	12
3	3	2	2	1	1	5	0	14	14	14
4	4	3	2	1	2	5	0	11	13	13

[5 rows x 33 columns]
Number of records port: 649
rows port: 649
columns port: 33

school
sex
age
address
famsize
Pstatus
Medu
Fedu
Mjob
Fjob
reason
guardian
traveltime
studytime
failures
schoolsup
famsup
paid
activities

nursery
 higher
 internet
 romantic
 famrel
 freetime
 goout
 Dalc
 Walc
 health
 absences
 G1
 G2
 G3

Step 2:

2.1:

Protected Class	Variable	Raw Values	Subsets
Gender	sex	F, M	Female, Male
Age	age	15, 16, 17, 18, 19, 20, 21, 22	[15-17] [18-22]

2.2:

Min Grade = 0

Avg Grade = 11

Max Grade = 19

bins = [0, 11, 19]

labels = [0, 1]

```
df['G1'] = pd.cut(df.G1, bins=bins, labels=labels, include_lowest=True)
```

```
df['G2'] = pd.cut(df.G2, bins=bins, labels=labels, include_lowest=True)
```

```
df['G3'] = pd.cut(df.G3, bins=bins, labels=labels, include_lowest=True)
```

Dependent Variables	Raw Categories	Discrete Categories
G1	0, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19	[0-11] represented as 0 [12-19] represented as 1
G2	0, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19	[0-11] represented as 0 [12-19] represented as 1
G3	0, 1, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19	[0-11] represented as 0 [12-19] represented as 1

Age	Bin	G1	G2	G3
15-17	0	227	225	205
15-17	1	241	243	263
18-22	0	116	106	96
18-22	1	65	75	85

Age	Bin	G1	G2	G3
15-17	0	227	225	205
15-17	1	241	243	263
18-22	0	116	106	96
18-22	1	65	75	85

2.3:

Age Group	Gender	Total
15-17	Female	275
15-17	Male	193
18-22	Female	108
18-22	Male	73

In [5]:

```
1
2 def STEP2():
3     df = pd.read_csv('dataset/student-por.csv', delimiter=',')
4
5     df.loc[df['sex'] == "M", ['sex']] = 'Male'
6     df.loc[df['sex'] == "F", ['sex']] = 'Female'
7
8     df.loc[df['age'].between(15, 17, inclusive=True), ['age_group']] =
9     df.loc[df['age'].between(18, 22, inclusive=True), ['age_group']] =
10
11     age_group_sex_freq = df.groupby(['age_group', 'sex']).size()
12     age_group_sex_freq.to_csv('out/age_group_sex_freq.csv')
13
14     age_unique = []
15     g1_unique = []
16     g2_unique = []
17     g3_unique = []
18
19     for x in df['age']:
20         if x not in age_unique:
21             age_unique.append(x)
22
23     for x in df['G1']:
24         if x not in g1_unique:
25             g1_unique.append(x)
26
27     for x in df['G2']:
28         if x not in g2_unique:
29             g2_unique.append(x)
30
31     for x in df['G3']:
32         if x not in g3_unique:
33             g3_unique.append(x)
34
35     age_unique.sort()
36     g1_unique.sort()
37     g2_unique.sort()
38     g3_unique.sort()
39
40     print("Age: ", age_unique)
41     print("G1: ", g1_unique)
42     print("G2: ", g2_unique)
43     print("G3: ", g3_unique)
44
45     print("G1 Max: ", df['G1'].max())
46     print("G1 Min: ", df['G1'].min())
47     print("G1 Avg: ", df['G1'].mean())
48
49     print("G2 Max: ", df['G2'].max())
50     print("G2 Min: ", df['G2'].min())
51     print("G2 Avg: ", df['G2'].mean())
52
53     print("G3 Max: ", df['G3'].max())
54     print("G3 Min: ", df['G3'].min())
55     print("G3 Avg: ", df['G3'].mean())
56
```

```

57 bins = [0, 11, 19]
58 labels = [0, 1]
59
60 df['G1'] = pd.cut(df.G1, bins=bins, labels=labels, include_lowest=True)
61 df['G2'] = pd.cut(df.G2, bins=bins, labels=labels, include_lowest=True)
62 df['G3'] = pd.cut(df.G3, bins=bins, labels=labels, include_lowest=True)
63
64 age_g1_freq = df.groupby(['age_group', 'G1']).size()
65 age_g1_freq.to_csv('out/age_g1_freq.csv')
66
67 age_g2_freq = df.groupby(['age_group', 'G2']).size()
68 age_g2_freq.to_csv('out/age_g2_freq.csv')
69
70 age_g3_freq = df.groupby(['age_group', 'G3']).size()
71 age_g3_freq.to_csv('out/age_g3_freq.csv')
72
73 sex_g1_freq = df.groupby(['sex', 'G1']).size()
74 sex_g1_freq.to_csv('out/sex_g1_freq.csv')
75
76 sex_g2_freq = df.groupby(['sex', 'G2']).size()
77 sex_g2_freq.to_csv('out/sex_g2_freq.csv')
78
79 sex_g3_freq = df.groupby(['sex', 'G3']).size()
80 sex_g3_freq.to_csv('out/sex_g3_freq.csv')
81
82 age_g1_freq.plot.bar(stacked=False, color="salmon")
83 plt.title('Age - First Period Grade')
84 plt.grid(True, axis='y', alpha=0.2, color='#999999')
85 plt.xlabel('Age - G1 Discretized')
86 plt.savefig('out/age_g1_freq.png', bbox_inches='tight')
87 plt.show()
88
89 age_g2_freq.plot.bar(stacked=False, color="salmon")
90 plt.title('Age - Second Period Grade')
91 plt.grid(True, axis='y', alpha=0.2, color='#999999')
92 plt.xlabel('Age - G2 Discretized')
93 plt.savefig('out/age_g2_freq.png', bbox_inches='tight')
94 plt.show()
95
96 age_g3_freq.plot.bar(stacked=False, color="salmon")
97 plt.title('Age - Final Grade')
98 plt.grid(True, axis='y', alpha=0.2, color='#999999')
99 plt.xlabel('Age - G3 Discretized')
100 plt.savefig('out/age_g3_freq.png', bbox_inches='tight')
101 plt.show()
102
103 sex_g1_freq.plot.bar(stacked=False, color="lightblue")
104 plt.title('Sex - First Period Grade')
105 plt.grid(True, axis='y', alpha=0.2, color='#999999')
106 plt.xlabel('Sex - G1 Discretized')
107 plt.savefig('out/sex_g1_freq.png', bbox_inches='tight')
108 plt.show()
109
110 sex_g2_freq.plot.bar(stacked=False, color="lightblue")
111 plt.title('Sex - Second Period Grade')
112 plt.grid(True, axis='y', alpha=0.2, color='#999999')
113 plt.xlabel('Sex - G2 Discretized')

```

```

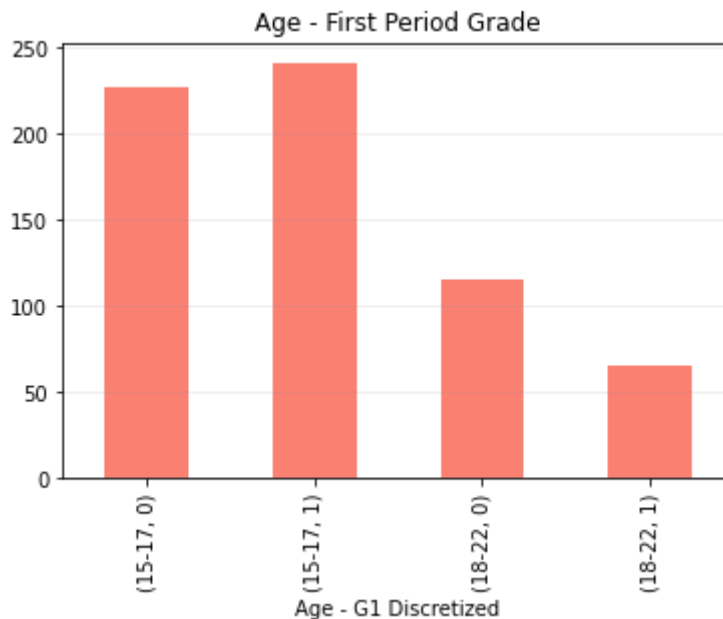
114 plt.savefig('out/sex_g2_freq.png', bbox_inches='tight')
115 plt.show()
116
117 sex_g3_freq.plot.bar(stacked=False, color="lightblue")
118 plt.title('Sex - Final Grade')
119 plt.grid(True, axis='y', alpha=0.2, color='#999999')
120 plt.xlabel('Sex - G3 Discretized')
121 plt.savefig('out/sex_g3_freq.png', bbox_inches='tight')
122 plt.show()
123
124 STEP2()

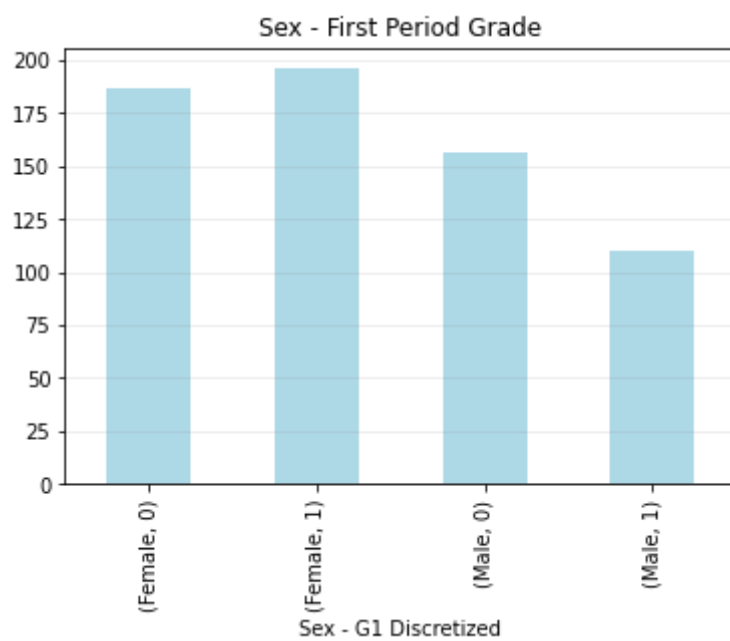
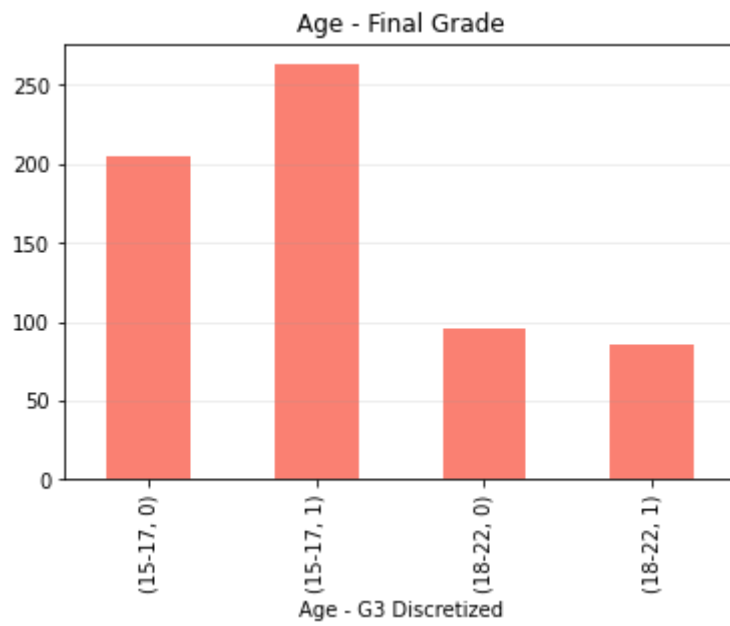
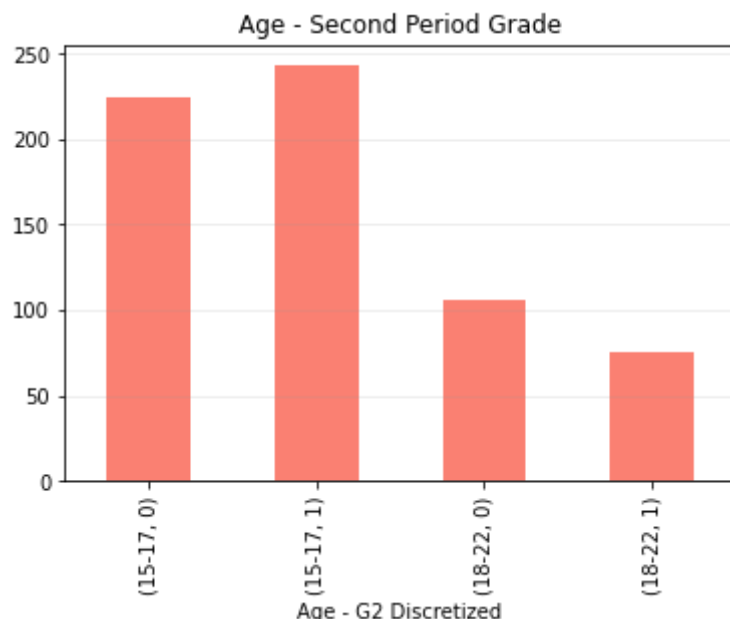
```

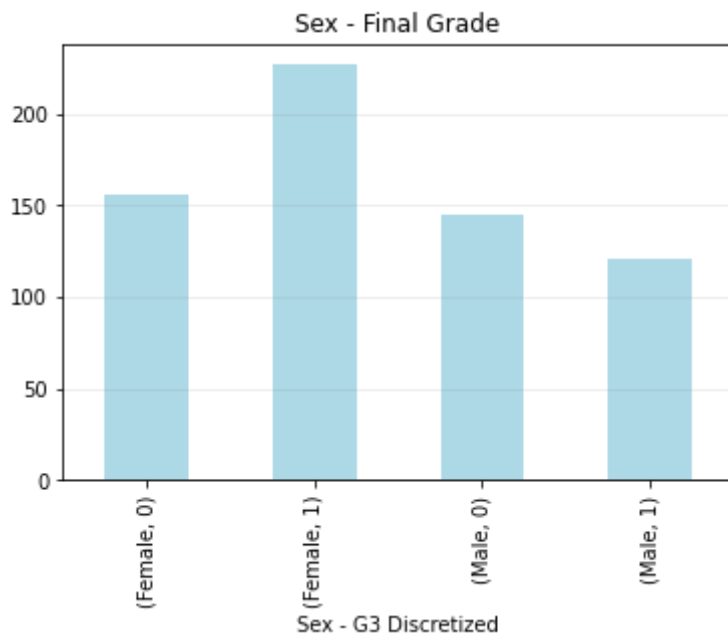
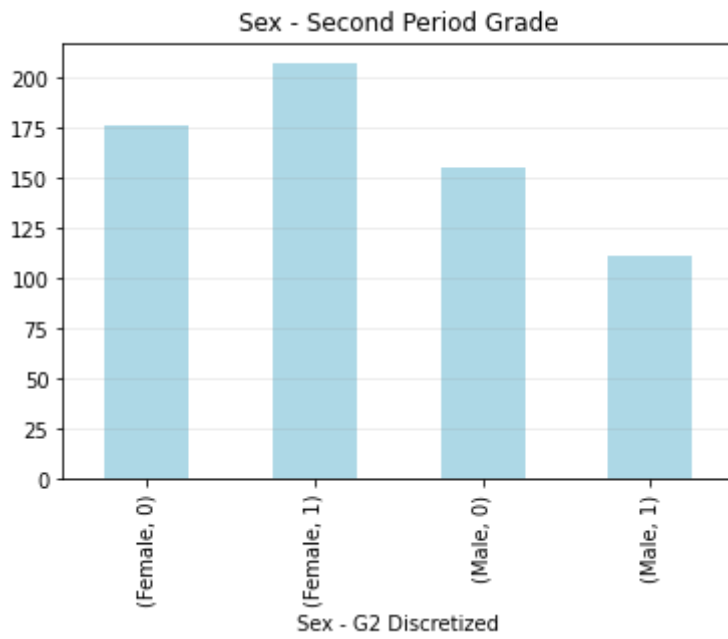
```

Age:  [15, 16, 17, 18, 19, 20, 21, 22]
G1:  [0, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19]
G2:  [0, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19]
G3:  [0, 1, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19]
G1 Max:  19
G1 Min:  0
G1 Avg:  11.399075500770415
G2 Max:  19
G2 Min:  0
G2 Avg:  11.570107858243452
G3 Max:  19
G3 Min:  0
G3 Avg:  11.906009244992296

```







Step 3:

Protected Class	Variable	Privileged Group	Unprivileged Group
Gender	sex	Female	Male

Protected Class	Variable	Privileged Group	Unprivileged Group
Age	age	[15-17]	[18-22]

The fairness metrics selected are: 1. Statistical Parity Difference 2. Disparate Impact The threshold chosen to calculate the fairness metric is the grade value of 12 (maximum grade value is 19).

Fairness Metrics Calculated:

Dependent Variable	Protected Class Variable	Statistical Parity Difference	Disparate Impact
G1	sex	-0.0982155	0.8080788
G1	age	-0.1558412	0.6973705
G2	sex	-0.1231767	0.7720932
G2	age	-0.1048661	0.7980356
G3	sex	-0.1378020	0.7674969
G3	age	-0.0923525	0.8356616

Bias Mitigation Strategy selected: **Reweighting**.

The weights were calculated for each Dependent Variable and Protected Class Variable Combination, and used these weights to calculate the new fairness metrics. (Different weights were used for each row, actual weight values can be found in the out folder after running the application).

Fairness Metrics Calculated after Reweighting:

Dependent Variable	Protected Class Variable	Statistical Parity Difference	Disparate Impact
G1	sex	0	1.0000000000000002
G1	age	0	1.0000000000000002
G2	sex	0	1
G2	age	0	1
G3	sex	0	1
G3	age	0	1.0000000000000002

In [13]:

```
1
2 def STEP3():
3     df = pd.read_csv('dataset/student-por.csv', delimiter=',')
4
5     df.loc[df['sex'] == "M", ['sex']] = 'Male'
6     df.loc[df['sex'] == "F", ['sex']] = 'Female'
7
8     df.loc[df['age'].between(15, 17, inclusive=True), ['age_group']] =
9     df.loc[df['age'].between(18, 22, inclusive=True), ['age_group']] =
10
11     bins = [0, 11, 19]
12     labels = [0, 1]
13
14     df['G1'] = pd.cut(df.G1, bins=bins, labels=labels, include_lowest=True)
15     df['G2'] = pd.cut(df.G2, bins=bins, labels=labels, include_lowest=True)
16     df['G3'] = pd.cut(df.G3, bins=bins, labels=labels, include_lowest=True)
17
18     age_g1_freq = df.groupby(['age_group', 'G1']).size()
19     age_g2_freq = df.groupby(['age_group', 'G2']).size()
20     age_g3_freq = df.groupby(['age_group', 'G3']).size()
21     sex_g1_freq = df.groupby(['sex', 'G1']).size()
22     sex_g2_freq = df.groupby(['sex', 'G2']).size()
23     sex_g3_freq = df.groupby(['sex', 'G3']).size()
24
25     up_sex = 'Male'
26     up_age = '18-22'
27     p_sex = 'Female'
28     p_age = '15-17'
29
30     # statistical parity difference G1 and age
31     spd_g1_age = (age_g1_freq[up_age][1] / (age_g1_freq[up_age][1] + age_g1_freq[up_age][0]) -
32                  age_g1_freq[p_age][1] / (age_g1_freq[p_age][1] + age_g1_freq[p_age][0]))
33
34     # statistical parity difference G2 and age
35     spd_g2_age = (age_g2_freq[up_age][1] / (age_g2_freq[up_age][1] + age_g2_freq[up_age][0]) -
36                  age_g2_freq[p_age][1] / (age_g2_freq[p_age][1] + age_g2_freq[p_age][0]))
37
38     # statistical parity difference G3 and age
39     spd_g3_age = (age_g3_freq[up_age][1] / (age_g3_freq[up_age][1] + age_g3_freq[up_age][0]) -
40                  age_g3_freq[p_age][1] / (age_g3_freq[p_age][1] + age_g3_freq[p_age][0]))
41
42     # statistical parity difference G1 and sex
43     spd_g1_sex = (sex_g1_freq[up_sex][1] / (sex_g1_freq[up_sex][1] + sex_g1_freq[up_sex][0]) -
44                  sex_g1_freq[p_sex][1] / (sex_g1_freq[p_sex][1] + sex_g1_freq[p_sex][0]))
45
46     # statistical parity difference G1 and sex
47     spd_g2_sex = (sex_g2_freq[up_sex][1] / (sex_g2_freq[up_sex][1] + sex_g2_freq[up_sex][0]) -
48                  sex_g2_freq[p_sex][1] / (sex_g2_freq[p_sex][1] + sex_g2_freq[p_sex][0]))
49
50     # statistical parity difference G1 and sex
51     spd_g3_sex = (sex_g3_freq[up_sex][1] / (sex_g3_freq[up_sex][1] + sex_g3_freq[up_sex][0]) -
52                  sex_g3_freq[p_sex][1] / (sex_g3_freq[p_sex][1] + sex_g3_freq[p_sex][0]))
53
54     spd_data = [['G1', 'Age', spd_g1_age], ['G2', 'Age', spd_g2_age],
55                ['G1', 'Sex', spd_g1_sex], ['G2', 'Sex', spd_g2_sex],
```

```

57 pd.DataFrame(spd_data,
58               columns=['Dependent Variable', 'Protected Class Variable',
59                       'out/spd.csv', index=False)
60
61 # disparate impact G1 and age
62 di_g1_age = (age_g1_freq[up_age][1] / (age_g1_freq[up_age][1] + age_g1_freq[p_age][1]) -
63             age_g1_freq[p_age][1] / (age_g1_freq[p_age][1] + age_g1_freq[up_age][1]))
64
65 # disparate impact G2 and age
66 di_g2_age = (age_g2_freq[up_age][1] / (age_g2_freq[up_age][1] + age_g2_freq[p_age][1]) -
67             age_g2_freq[p_age][1] / (age_g2_freq[p_age][1] + age_g2_freq[up_age][1]))
68
69 # disparate impact G3 and age
70 di_g3_age = (age_g3_freq[up_age][1] / (age_g3_freq[up_age][1] + age_g3_freq[p_age][1]) -
71             age_g3_freq[p_age][1] / (age_g3_freq[p_age][1] + age_g3_freq[up_age][1]))
72
73 # disparate impact G1 and sex
74 di_g1_sex = (sex_g1_freq[up_sex][1] / (sex_g1_freq[up_sex][1] + sex_g1_freq[p_sex][1]) -
75             sex_g1_freq[p_sex][1] / (sex_g1_freq[p_sex][1] + sex_g1_freq[up_sex][1]))
76
77 # disparate impact G2 and sex
78 di_g2_sex = (sex_g2_freq[up_sex][1] / (sex_g2_freq[up_sex][1] + sex_g2_freq[p_sex][1]) -
79             sex_g2_freq[p_sex][1] / (sex_g2_freq[p_sex][1] + sex_g2_freq[up_sex][1]))
80
81 # disparate impact G3 and sex
82 di_g3_sex = (sex_g3_freq[up_sex][1] / (sex_g3_freq[up_sex][1] + sex_g3_freq[p_sex][1]) -
83             sex_g3_freq[p_sex][1] / (sex_g3_freq[p_sex][1] + sex_g3_freq[up_sex][1]))
84
85 di_data = [['G1', 'Age', di_g1_age], ['G2', 'Age', di_g2_age], ['G3', 'Age', di_g3_age],
86            ['G1', 'Sex', di_g1_sex], ['G2', 'Sex', di_g2_sex], ['G3', 'Sex', di_g3_sex]]
87
88 pd.DataFrame(di_data,
89               columns=['Dependent Variable', 'Protected Class Variable', 'Outcome',
90                       'out/di.csv', index=False)
91
92 # re-weighting
93
94 # weights for G1 and age
95 w_pp_age_g1 = (age_g1_freq[p_age].sum() * (age_g1_freq[p_age][1] + age_g1_freq[up_age][1]) -
96               age_g1_freq.values.sum() * age_g1_freq[p_age][1]) /
97               (age_g1_freq[p_age].sum() * (age_g1_freq[p_age][1] + age_g1_freq[up_age][1]) -
98               age_g1_freq.values.sum() * age_g1_freq[p_age][1])
99 w_pu_age_g1 = (age_g1_freq[up_age].sum() * (age_g1_freq[p_age][1] + age_g1_freq[up_age][1]) -
100               age_g1_freq.values.sum() * age_g1_freq[up_age][1]) /
101               (age_g1_freq[up_age].sum() * (age_g1_freq[p_age][1] + age_g1_freq[up_age][1]) -
102               age_g1_freq.values.sum() * age_g1_freq[up_age][1])
103
104 weights = [[w_pp_age_g1, w_pu_age_g1, w_np_age_g1, w_nu_age_g1]]
105 pd.DataFrame(weights,
106               columns=['Positive outcome - Privileged Group', 'Positive outcome - Protected Group',
107                       'Negative outcome - Privileged Group', 'Negative outcome - Protected Group',
108                       'out/weights_g1_age.csv', index=False)
109
110 # applying the weights to calculate spd and di
111 w_spd_g1_age = (w_pu_age_g1 * age_g1_freq[up_age][1] / (
112               w_pu_age_g1 * age_g1_freq[up_age][1] + w_nu_age_g1 * age_g1_freq[up_age][0] +
113               w_pp_age_g1 * age_g1_freq[p_age][1] / (

```

```

114         w_pp_age_g1 * age_g1_freq[p_age][1] + w_np_
115 w_di_g1_age = (w_pu_age_g1 * age_g1_freq[up_age][1] / (
116         w_pu_age_g1 * age_g1_freq[up_age][1] + w_nu_age_g1 * age_g
117         w_pp_age_g1 * age_g1_freq[p_age][1] / (
118         w_pp_age_g1 * age_g1_freq[p_age][1] + w_np_a
119
120 # weights for G2 and age
121 w_pp_age_g2 = (age_g2_freq[p_age].sum() * (age_g2_freq[p_age][1] +
122         age_g2_freq.values.sum() * age_g2_freq[p_age][1])
123 w_pu_age_g2 = (age_g2_freq[up_age].sum() * (age_g2_freq[p_age][1] +
124         age_g2_freq.values.sum() * age_g2_freq[up_age][1])
125 w_np_age_g2 = (age_g2_freq[p_age].sum() * (age_g2_freq[p_age][0] +
126         age_g2_freq.values.sum() * age_g2_freq[p_age][0])
127 w_nu_age_g2 = (age_g2_freq[up_age].sum() * (age_g2_freq[p_age][0] +
128         age_g2_freq.values.sum() * age_g2_freq[up_age][0])
129
130 weights = [[w_pp_age_g2, w_pu_age_g2, w_np_age_g2, w_nu_age_g2]]
131 pd.DataFrame(weights,
132         columns=['Positive outcome - Privileged Group', 'Posi
133         'Negative outcome - Privileged Group', 'Nega
134         'out/weights_g2_age.csv', index=False)
135
136 # applying the weights to calculate spd and di
137 w_spd_g2_age = (w_pu_age_g2 * age_g2_freq[up_age][1] / (
138         w_pu_age_g2 * age_g2_freq[up_age][1] + w_nu_age_g2 * age_g
139         w_pp_age_g2 * age_g2_freq[p_age][1] / (
140         w_pp_age_g2 * age_g2_freq[p_age][1] + w_np_
141 w_di_g2_age = (w_pu_age_g2 * age_g2_freq[up_age][1] / (
142         w_pu_age_g2 * age_g2_freq[up_age][1] + w_nu_age_g2 * age_g
143         w_pp_age_g2 * age_g2_freq[p_age][1] / (
144         w_pp_age_g2 * age_g2_freq[p_age][1] + w_np_a
145
146 # weights for G3 # and age
147 w_pp_age_g3 = (age_g3_freq[p_age].sum() * (age_g3_freq[p_age][1] +
148         age_g3_freq.values.sum() * age_g3_freq[p_age][1])
149 w_pu_age_g3 = (age_g3_freq[up_age].sum() * (age_g3_freq[p_age][1] +
150         age_g3_freq.values.sum() * age_g3_freq[up_age][1])
151 w_np_age_g3 = (age_g3_freq[p_age].sum() * (age_g3_freq[p_age][0] +
152         age_g3_freq.values.sum() * age_g3_freq[p_age][0])
153 w_nu_age_g3 = (age_g3_freq[up_age].sum() * (age_g3_freq[p_age][0] +
154         age_g3_freq.values.sum() * age_g3_freq[up_age][0])
155
156 weights = [[w_pp_age_g3, w_pu_age_g3, w_np_age_g3, w_nu_age_g3]]
157 pd.DataFrame(weights,
158         columns=['Positive outcome - Privileged Group', 'Posi
159         'Negative outcome - Privileged Group', 'Nega
160         'out/weights_g3_age.csv', index=False)
161
162 # applying the weights to calculate spd and di
163 w_spd_g3_age = (w_pu_age_g3 * age_g3_freq[up_age][1] / (
164         w_pu_age_g3 * age_g3_freq[up_age][1] + w_nu_age_g3 * age_g
165         w_pp_age_g3 * age_g3_freq[p_age][1] / (
166         w_pp_age_g3 * age_g3_freq[p_age][1] + w_np_
167 w_di_g3_age = (w_pu_age_g3 * age_g3_freq[up_age][1] / (
168         w_pu_age_g3 * age_g3_freq[up_age][1] + w_nu_age_g3 * age_g
169         w_pp_age_g3 * age_g3_freq[p_age][1] / (
170         w_pp_age_g3 * age_g3_freq[p_age][1] + w_np_a

```

```

171
172 # weights for G1 and sex
173 w_pp_sex_g1 = (sex_g1_freq[p_sex].sum() * (sex_g1_freq[p_sex][1] +
174 sex_g1_freq.values.sum() * sex_g1_freq[p_sex][1])
175 w_pu_sex_g1 = (sex_g1_freq[up_sex].sum() * (sex_g1_freq[p_sex][1] +
176 sex_g1_freq.values.sum() * sex_g1_freq[up_sex][1])
177 w_np_sex_g1 = (sex_g1_freq[p_sex].sum() * (sex_g1_freq[p_sex][0] +
178 sex_g1_freq.values.sum() * sex_g1_freq[p_sex][0])
179 w_nu_sex_g1 = (sex_g1_freq[up_sex].sum() * (sex_g1_freq[p_sex][0] +
180 sex_g1_freq.values.sum() * sex_g1_freq[up_sex][0])
181
182 weights = [[w_pp_sex_g1, w_pu_sex_g1, w_np_sex_g1, w_nu_sex_g1]]
183 pd.DataFrame(weights,
184 columns=['Positive outcome - Privileged Group', 'Positive outcome - Disadvantaged Group',
185 'Negative outcome - Privileged Group', 'Negative outcome - Disadvantaged Group'],
186 'out/weights_g1_sex.csv', index=False)
187
188 # applying the weights to calculate spd and di
189 w_spd_g1_sex = (w_pu_sex_g1 * sex_g1_freq[up_sex][1] / (
190 w_pu_sex_g1 * sex_g1_freq[up_sex][1] + w_nu_sex_g1 * sex_g1_freq[up_sex][0]) +
191 (w_pp_sex_g1 * sex_g1_freq[p_sex][1] / (
192 w_pp_sex_g1 * sex_g1_freq[p_sex][1] + w_np_sex_g1 * sex_g1_freq[p_sex][0])
193 w_di_g1_sex = (w_pu_sex_g1 * sex_g1_freq[up_sex][1] / (
194 w_pu_sex_g1 * sex_g1_freq[up_sex][1] + w_nu_sex_g1 * sex_g1_freq[up_sex][0]) -
195 (w_pp_sex_g1 * sex_g1_freq[p_sex][1] / (
196 w_pp_sex_g1 * sex_g1_freq[p_sex][1] + w_np_sex_g1 * sex_g1_freq[p_sex][0])
197
198 # weights for G2 and sex
199 w_pp_sex_g2 = (sex_g2_freq[p_sex].sum() * (sex_g2_freq[p_sex][1] +
200 sex_g2_freq.values.sum() * sex_g2_freq[p_sex][1])
201 w_pu_sex_g2 = (sex_g2_freq[up_sex].sum() * (sex_g2_freq[p_sex][1] +
202 sex_g2_freq.values.sum() * sex_g2_freq[up_sex][1])
203 w_np_sex_g2 = (sex_g2_freq[p_sex].sum() * (sex_g2_freq[p_sex][0] +
204 sex_g2_freq.values.sum() * sex_g2_freq[p_sex][0])
205 w_nu_sex_g2 = (sex_g2_freq[up_sex].sum() * (sex_g2_freq[p_sex][0] +
206 sex_g2_freq.values.sum() * sex_g2_freq[up_sex][0])
207
208 weights = [[w_pp_sex_g2, w_pu_sex_g2, w_np_sex_g2, w_nu_sex_g2]]
209 pd.DataFrame(weights,
210 columns=['Positive outcome - Privileged Group', 'Positive outcome - Disadvantaged Group',
211 'Negative outcome - Privileged Group', 'Negative outcome - Disadvantaged Group'],
212 'out/weights_g2_sex.csv', index=False)
213
214 # applying the weights to calculate spd and di
215 w_spd_g2_sex = (w_pu_sex_g2 * sex_g2_freq[up_sex][1] / (
216 w_pu_sex_g2 * sex_g2_freq[up_sex][1] + w_nu_sex_g2 * sex_g2_freq[up_sex][0]) +
217 (w_pp_sex_g2 * sex_g2_freq[p_sex][1] / (
218 w_pp_sex_g2 * sex_g2_freq[p_sex][1] + w_np_sex_g2 * sex_g2_freq[p_sex][0])
219 w_di_g2_sex = (w_pu_sex_g2 * sex_g2_freq[up_sex][1] / (
220 w_pu_sex_g2 * sex_g2_freq[up_sex][1] + w_nu_sex_g2 * sex_g2_freq[up_sex][0]) -
221 (w_pp_sex_g2 * sex_g2_freq[p_sex][1] / (
222 w_pp_sex_g2 * sex_g2_freq[p_sex][1] + w_np_sex_g2 * sex_g2_freq[p_sex][0])
223
224 # weights for G3 and sex
225 w_pp_sex_g3 = (sex_g3_freq[p_sex].sum() * (sex_g3_freq[p_sex][1] +
226 sex_g3_freq.values.sum() * sex_g3_freq[p_sex][1])
227 w_pu_sex_g3 = (sex_g3_freq[up_sex].sum() * (sex_g3_freq[p_sex][1] +

```

```

228     sex_g3_freq.values.sum() * sex_g3_freq[up_sex][1])
229 w_np_sex_g3 = (sex_g3_freq[p_sex].sum() * (sex_g3_freq[p_sex][0] +
230     sex_g3_freq.values.sum() * sex_g3_freq[p_sex][0])
231 w_nu_sex_g3 = (sex_g3_freq[up_sex].sum() * (sex_g3_freq[p_sex][0] +
232     sex_g3_freq.values.sum() * sex_g3_freq[up_sex][0])
233
234 weights = [[w_pp_sex_g3, w_pu_sex_g3, w_np_sex_g3, w_nu_sex_g3]]
235 pd.DataFrame(weights,
236     columns=['Positive outcome - Privileged Group', 'Posi
237     'Negative outcome - Privileged Group', 'Nega
238     'out/weights_g3_sex.csv', index=False)
239
240 # applying the weights to calculate spd and di
241 w_spd_g3_sex = (w_pu_sex_g3 * sex_g3_freq[up_sex][1] / (
242     w_pu_sex_g3 * sex_g3_freq[up_sex][1] + w_nu_sex_g3 * sex_g
243     w_pp_sex_g3 * sex_g3_freq[p_sex][1] / (
244     w_pp_sex_g3 * sex_g3_freq[p_sex][1] + w_np_s
245 w_di_g3_sex = (w_pu_sex_g3 * sex_g3_freq[up_sex][1] / (
246     w_pu_sex_g3 * sex_g3_freq[up_sex][1] + w_nu_sex_g3 * sex_g
247     w_pp_sex_g3 * sex_g3_freq[p_sex][1] / (
248     w_pp_sex_g3 * sex_g3_freq[p_sex][1] + w_np_s
249
250 w_spd_data = [['G1', 'Age', w_spd_g1_age], ['G2', 'Age', w_spd_g2_
251     ['G1', 'Sex', w_spd_g1_sex], ['G2', 'Sex', w_spd_g2_
252
253 pd.DataFrame(w_spd_data,
254     columns=['Dependent Variable', 'Protected Class Variab
255     'out/w_spd.csv', index=False)
256
257 w_di_data = [['G1', 'Age', w_di_g1_age], ['G2', 'Age', w_di_g2_age
258     ['G1', 'Sex', w_di_g1_sex], ['G2', 'Sex', w_di_g2_sex
259
260 pd.DataFrame(w_di_data,
261     columns=['Dependent Variable', 'Protected Class Variab
262     'out/w_di.csv', index=False)
263
264 # apply weights to original data set TODO clean this up
265 weighted_df = pd.read_csv('dataset/student-por.csv', delimiter=',')
266 weighted_df['G1_Weighted'] = weighted_df['G1'] + 0.0
267 weighted_df.loc[weighted_df['sex'] == "M", ['sex']] = 'Male'
268 weighted_df.loc[weighted_df['sex'] == "F", ['sex']] = 'Female'
269 weighted_df.loc[weighted_df['age'].between(15, 17, inclusive=True)
270 weighted_df.loc[weighted_df['age'].between(18, 22, inclusive=True)
271 weighted_df['G1_PassFail'] = pd.cut(weighted_df.G1, bins=bins, lab
272 weighted_df['G2_PassFail'] = pd.cut(weighted_df.G2, bins=bins, lab
273 # weighted_df['G3_PassFail'] = pd.cut(weighted_df.G3, bins=bins, l
274
275 # add weight columns as floats
276 weighted_df['G1_Weighted'] = weighted_df['G1'] + 0.0
277 weighted_df['G2_Weighted'] = weighted_df['G2'] + 0.0
278 # weighted_df['G3_Weighted'] = weighted_df['G3'] + 0.0
279
280 # outcome: // PO PG == female 15-17 e.g. row 2
281 p_pos = (weighted_df['sex'] == p_sex) & (weighted_df['age_group']
282 weighted_df.loc[p_pos, 'G1_Weighted'] = weighted_df.loc[p_pos, 'G1
283 p_pos = (weighted_df['sex'] == p_sex) & (weighted_df['age_group']
284 weighted_df.loc[p_pos, 'G2_Weighted'] = weighted_df.loc[p_pos, 'G2

```



```

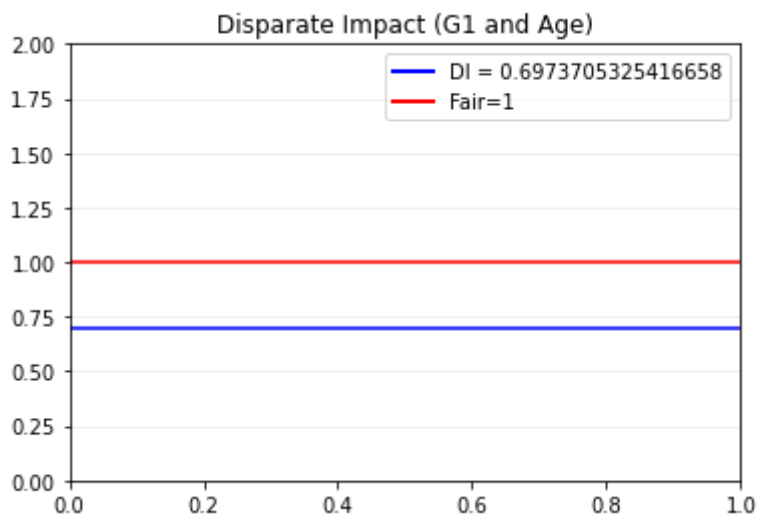
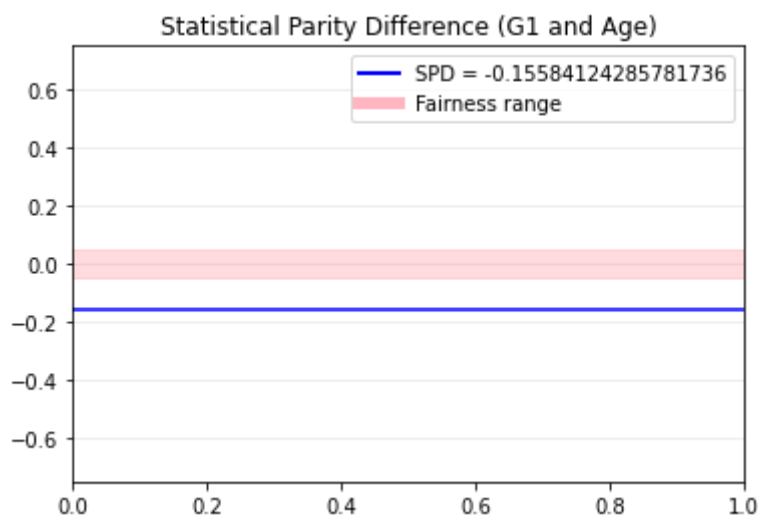
285 # p_pos = (weighted_df['sex'] == p_sex) & (weighted_df['age_group']
286 # weighted_df.loc[p_pos, 'G3_Weighted'] = weighted_df.loc[p_pos, '
287
288 # outcome: // PO UG == male 18-22 e.g. row 229
289 up_pos = (weighted_df['sex'] == up_sex) & (weighted_df['age_group']
290 weighted_df.loc[up_pos, 'G1_Weighted'] = weighted_df.loc[up_pos, '
291 up_pos = (weighted_df['sex'] == up_sex) & (weighted_df['age_group']
292 weighted_df.loc[up_pos, 'G2_Weighted'] = weighted_df.loc[up_pos, '
293 # up_pos = (weighted_df['sex'] == up_sex) & (weighted_df['age_group']
294 # weighted_df.loc[up_pos, 'G3_Weighted'] = weighted_df.loc[up_pos,
295
296 # outcome: // NO PG == female 15-17 e.g. row
297 p_neg = (weighted_df['sex'] == p_sex) & (weighted_df['age_group']
298 weighted_df.loc[p_neg, 'G1_Weighted'] = weighted_df.loc[p_neg, 'G1_
299 p_neg = (weighted_df['sex'] == p_sex) & (weighted_df['age_group']
300 weighted_df.loc[p_neg, 'G2_Weighted'] = weighted_df.loc[p_neg, 'G2_
301 # p_neg = (weighted_df['sex'] == p_sex) & (weighted_df['age_group']
302 # weighted_df.loc[p_neg, 'G3_Weighted'] = weighted_df.loc[p_neg, '
303
304 # outcome: // NO UG == male 18-22 e.g. row 165
305 up_neg = (weighted_df['sex'] == up_sex) & (weighted_df['age_group']
306 weighted_df.loc[up_neg, 'G1_Weighted'] = weighted_df.loc[up_neg, '
307 up_neg = (weighted_df['sex'] == up_sex) & (weighted_df['age_group']
308 weighted_df.loc[up_neg, 'G2_Weighted'] = weighted_df.loc[up_neg, '
309 # up_neg = (weighted_df['sex'] == up_sex) & (weighted_df['age_group']
310 # weighted_df.loc[up_neg, 'G3_Weighted'] = weighted_df.loc[up_neg,
311
312 weighted_df.to_csv('out/weighted.csv', index=False)
313
314 for i in range(0, 6):
315     spd_i = spd_data[i]
316     plt.axhline(y=spd_i[2], color='blue')
317     plt.axhspan(-0.05, 0.05, alpha=0.5, color='#ffb6c1')
318     axes = plt.gca()
319     axes.set_ylim([-0.75, 0.75])
320     plt.grid(True, axis='y', alpha=0.2, color='#999999')
321     plt.title('Statistical Parity Difference (' + spd_i[0] + ' and
322     plt.legend([Line2D([0], [0], color='blue', lw=2), Line2D([0],
323         ['SPD = ' + str(spd_i[2]), 'Fairness range'])
324     plt.savefig('out/spd_' + spd_i[0] + '_' + spd_i[1] + '.png')
325     plt.show()
326
327     di_i = di_data[i]
328     plt.axhline(y=di_i[2], color='blue')
329     plt.axhline(y=1, color='red')
330     axes = plt.gca()
331     axes.set_ylim([0, 2])
332     plt.grid(True, axis='y', alpha=0.2, color='#999999')
333     plt.title('Disparate Impact (' + di_i[0] + ' and ' + di_i[1] +
334     plt.legend([Line2D([0], [0], color='blue', lw=2), Line2D([0],
335         ['DI = ' + str(di_i[2]), 'Fair=1'])
336     plt.savefig('out/di_' + di_i[0] + '_' + di_i[1] + '.png')
337     plt.show()
338
339     w_spd_i = w_spd_data[i]
340     plt.axhline(y=w_spd_i[2], color='blue')
341     plt.axhspan(-0.05, 0.05, alpha=0.5, color='#ffb6c1')

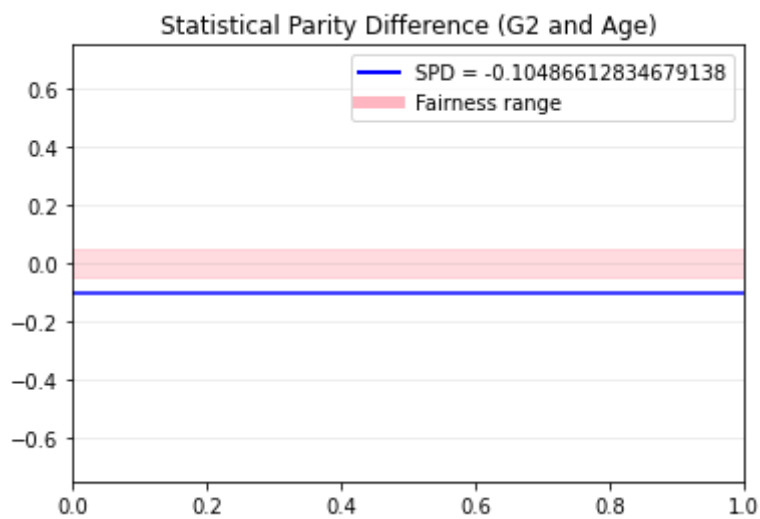
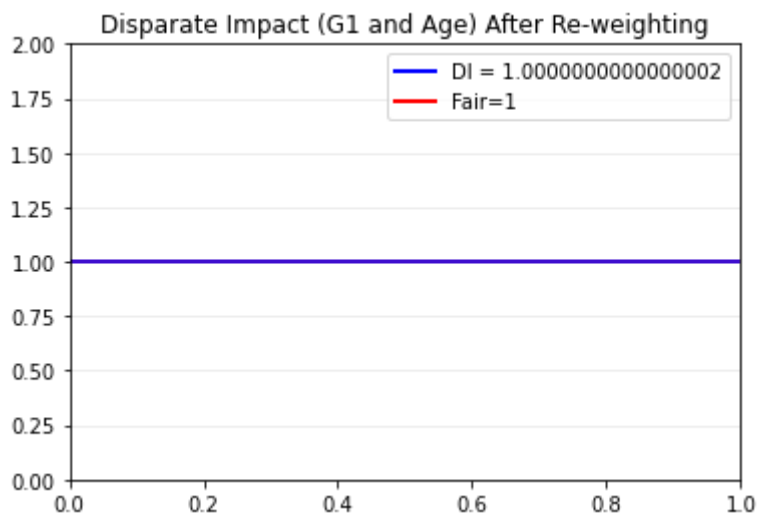
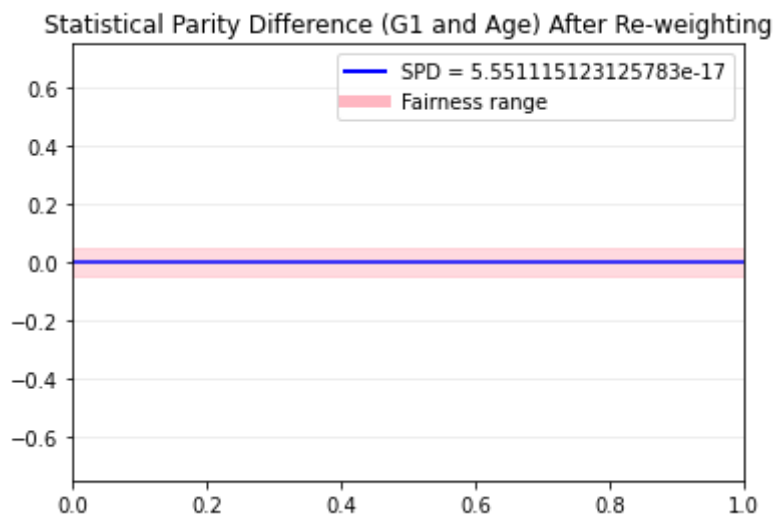
```

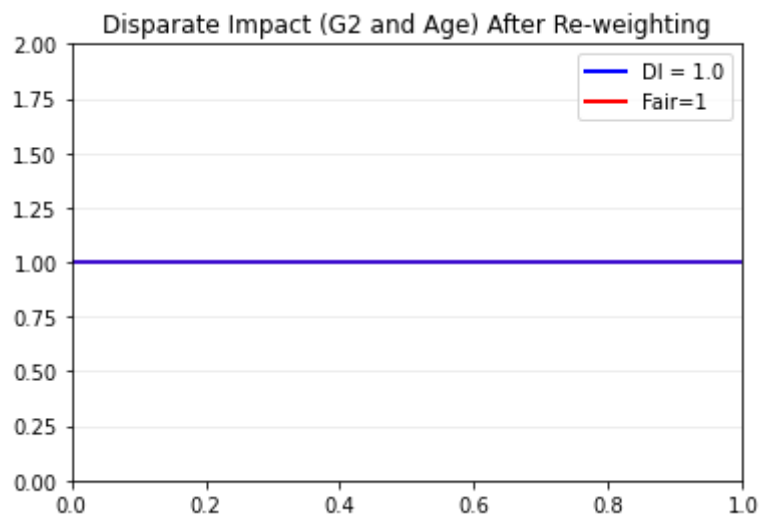
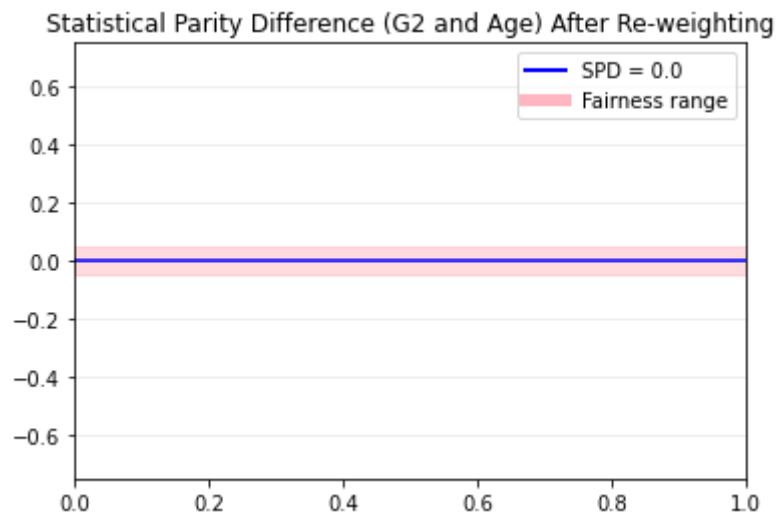
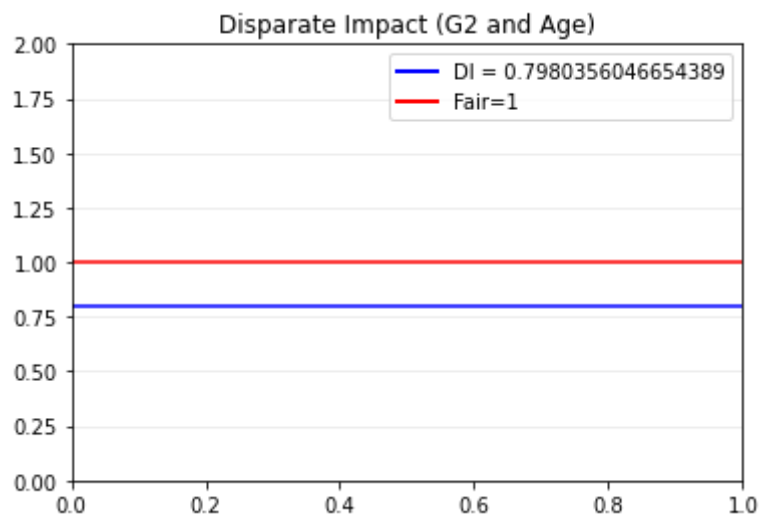
```

342     axes = plt.gca()
343     axes.set_ylim([-0.75, 0.75])
344     plt.grid(True, axis='y', alpha=0.2, color='#999999')
345     plt.title('Statistical Parity Difference (' + w_spd_i[0] + ' and ' + w_spd_i[1] + ')')
346     plt.legend([Line2D([0], [0], color='blue', lw=2), Line2D([0], [0],
347         ['SPD = ' + str(w_spd_i[2]), 'Fairness range'])
348     plt.savefig('out/w_spd_' + w_spd_i[0] + '_' + w_spd_i[1] + '.png')
349     plt.show()
350
351     w_di_i = w_di_data[i]
352     plt.axhline(y=1, color='red')
353     plt.axhline(y=w_di_i[2], color='blue')
354     axes = plt.gca()
355     axes.set_ylim([0, 2])
356     plt.grid(True, axis='y', alpha=0.2, color='#999999')
357     plt.legend([Line2D([0], [0], color='blue', lw=2), Line2D([0], [0],
358         ['DI = ' + str(w_di_i[2]), 'Fair=1'])
359     plt.title('Disparate Impact (' + w_di_i[0] + ' and ' + w_di_i[1] + ')')
360     plt.savefig('out/w_di_' + w_di_i[0] + '_' + w_di_i[1] + '.png')
361     plt.show()
362
363 STEP3()

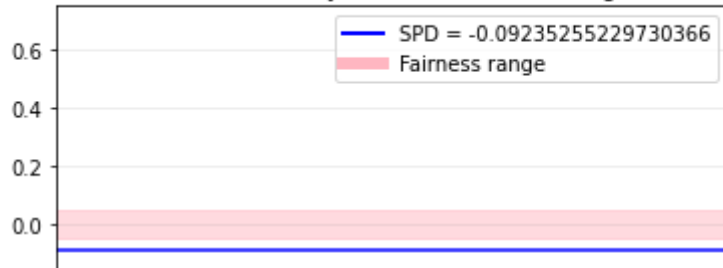
```



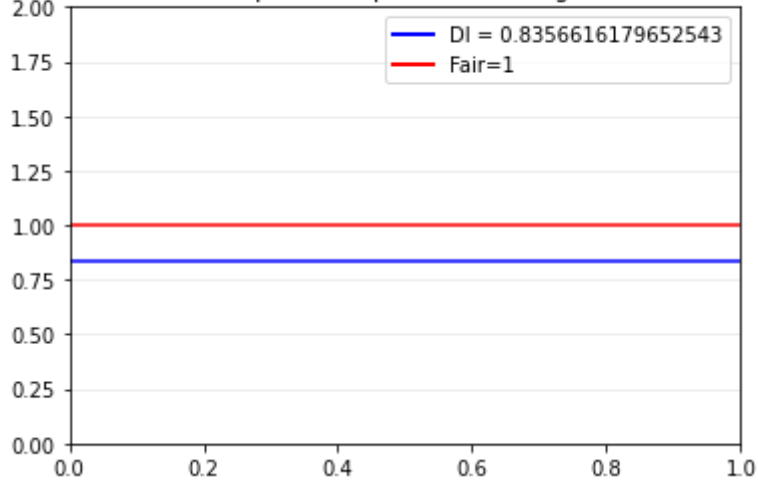




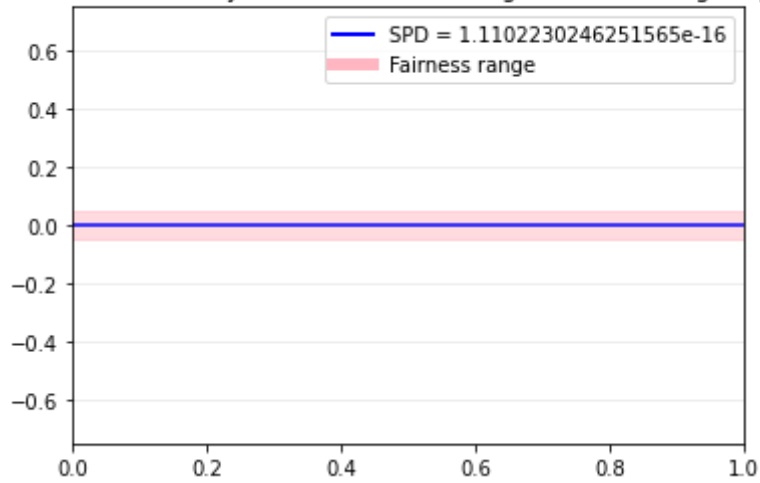
Statistical Parity Difference (G3 and Age)

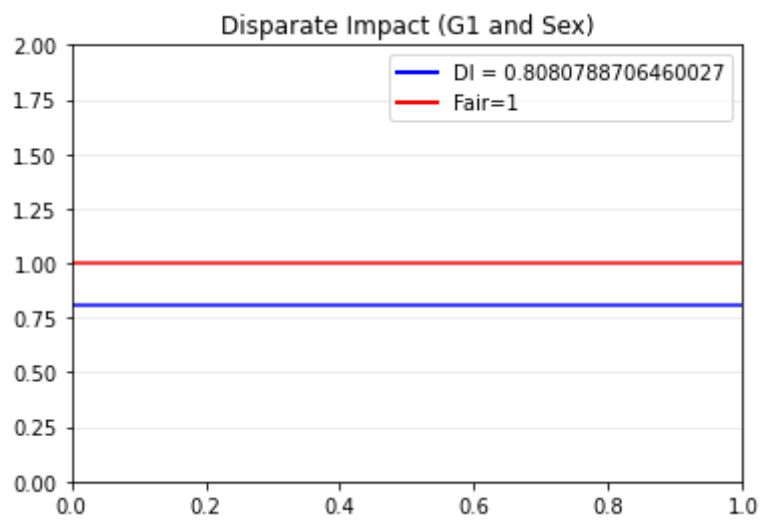
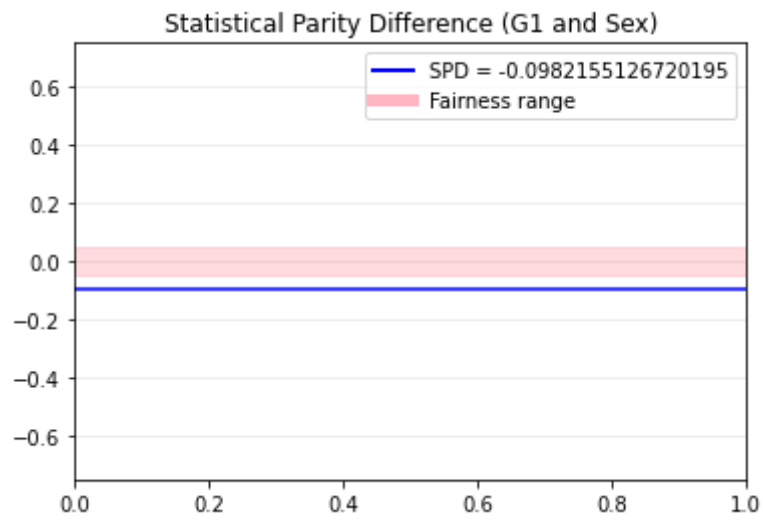
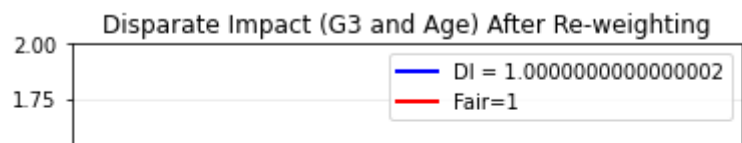


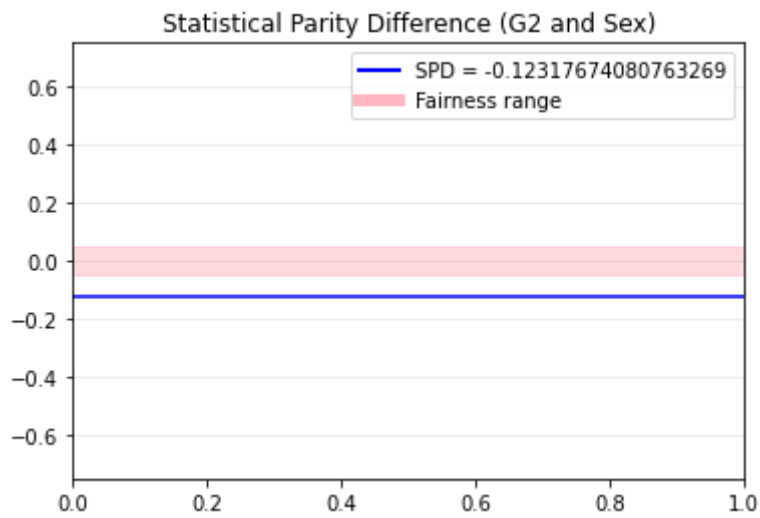
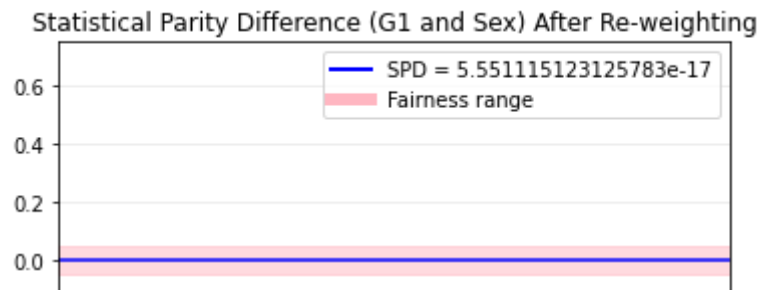
Disparate Impact (G3 and Age)

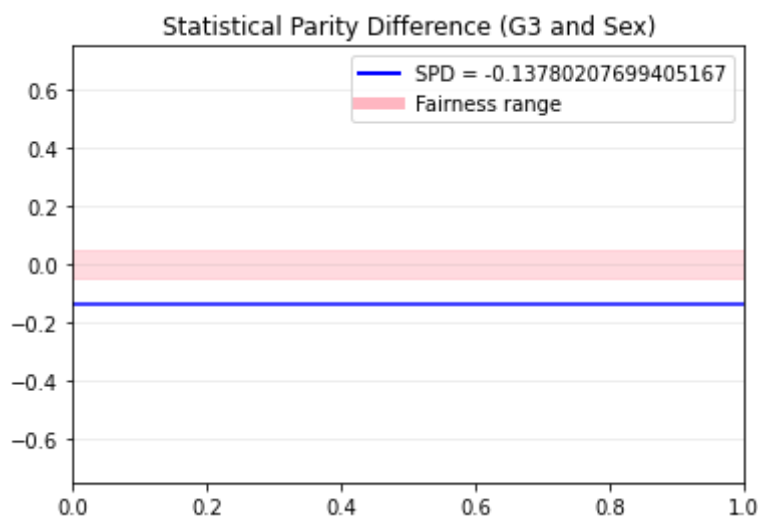
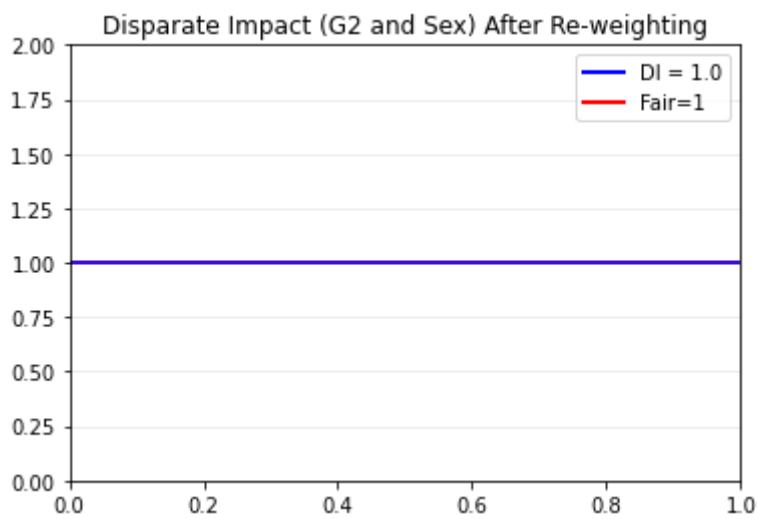
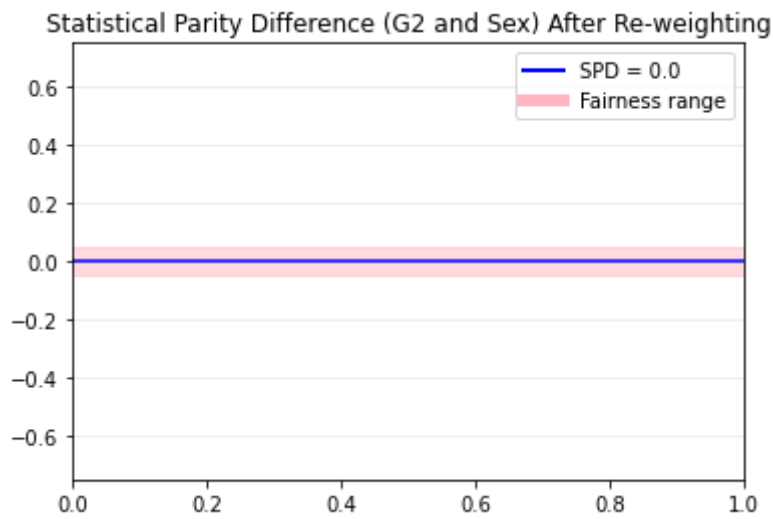
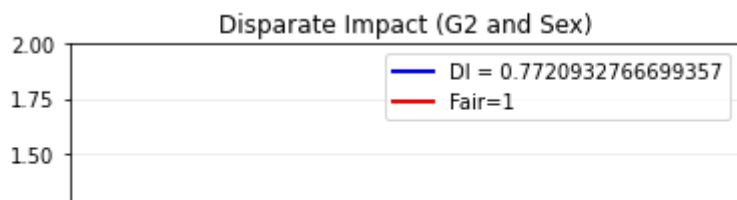


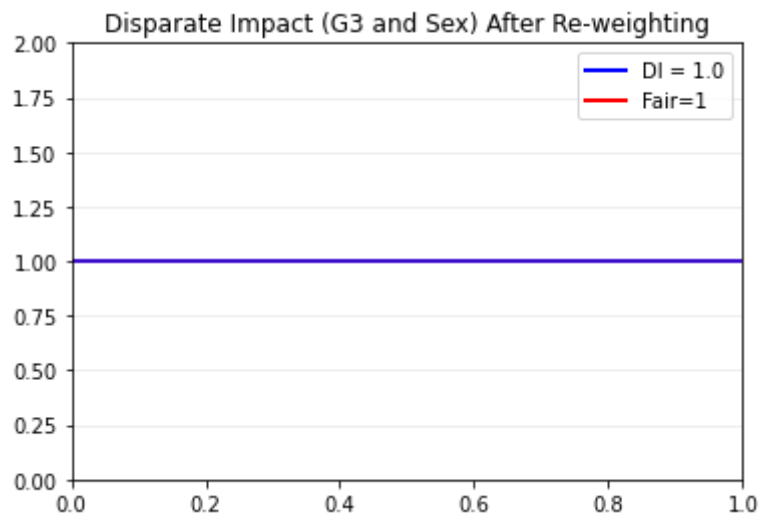
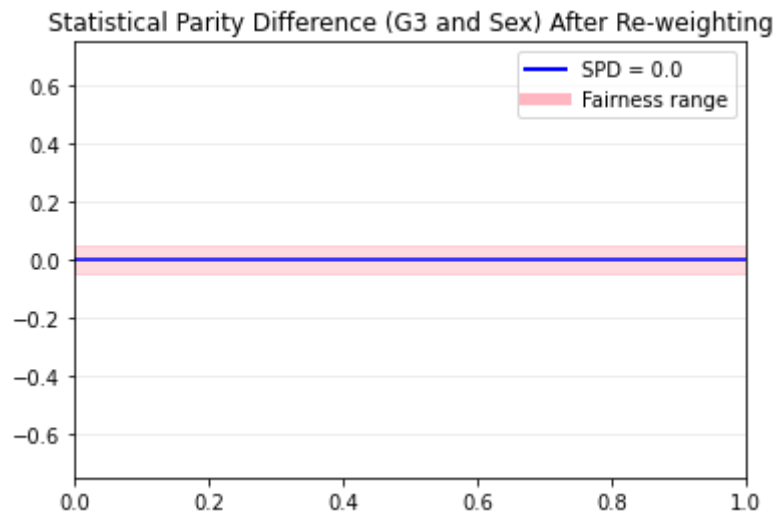
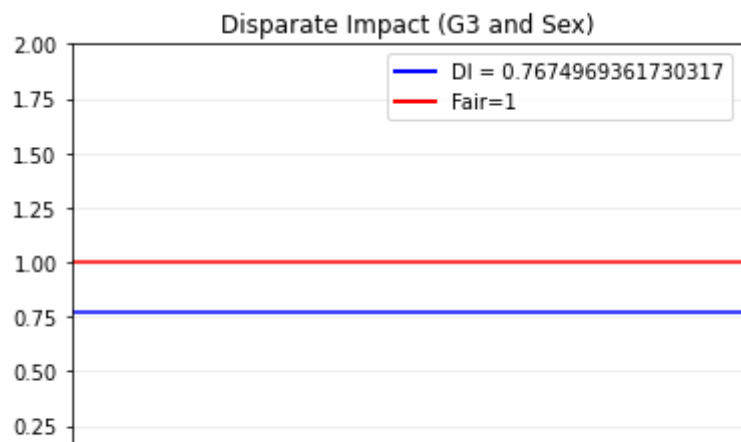
Statistical Parity Difference (G3 and Age) After Re-weighting











Step 4:

Option A: Independent Variable: 0.024299344 Privileged Group: Female Unprivileged Group: Male

Independent		Metric	Original	Transformed	Difference
Sex	Statistical Parity Difference		-0.119759532	-0.151282051	0.03152252
Sex	Disparate Impact		0.774299344	0.75	0.024299344

Was there a positive change, negative change, or no change on that fairness metric after transforming the dataset (from Step 3.4)?

After transforming the dataset, there was a slight negative change for both Statistical Parity Difference and Disparate Impact.

Was there a positive change, negative change, or no change on that fairness metric after training the classifier - with respect to the original testing dataset and the transformed testing dataset?

In our opinion, there was no change on the fairness metric after training the classifier. If trained on the transformed data, the testing results would result in the same or similar metrics as the training transformed data itself.

In [11]:

```
1
2 def STEP4():
3     raw_df = pd.read_csv(r"dataset/student-por.csv", delimiter=',')
4     df = raw_df.replace({'sex': {'M': 0, 'F': 1}})
5
6     df_num_rows = int(len(df))
7
8     # split data into train and test sets
9     shuffled = df.sample(frac=1)
10    df_train = shuffled.iloc[:int(df_num_rows / 2)]
11    df_test = shuffled.iloc[int(df_num_rows / 2):]
12
13    # classify data Y in this case is G3. Features are G1 and G2,
14
15    features = ['sex', 'G1', 'G2']
16    y_feature = 'G3'
17
18    # create training data
19    df_train_x = df_train[features]
20    df_train_y = df_train[y_feature]
21
22    # create and train classifier
23    orig_classifier = DecisionTreeClassifier()
24    orig_classifier.fit(df_train_x, df_train_y)
25
26    # run classifier against the test set
27    df_test_x = df_test[features]
28    predicted = orig_classifier.predict(df_test_x)
29
30    df_test_x['G3_Predicted'] = predicted
31    df_test_x['G3_Actual'] = df_test[y_feature]
32
33    # create a classifier for the data from Step 3.3
34    weighted_raw_df = pd.read_csv(r"out/weighted.csv", delimiter=',')
35    weighted_df = weighted_raw_df.replace({'sex': {'Male': 0, 'Female': 1}})
36
37    weighted_df_num_rows = int(len(weighted_df))
38
39    # split data into train and test sets
40    weighted_shuffled = weighted_df.sample(frac=1)
41    weighted_df_train = weighted_shuffled.iloc[:int(weighted_df_num_rows / 2)]
42    weighted_df_test = weighted_shuffled.iloc[int(weighted_df_num_rows / 2):]
43
44    # create training data
45    features = ['sex', 'G1_Weighted', 'G2_Weighted']
46    y_feature = 'G3'
47    weighted_df_train_x = weighted_df_train[features]
48    weighted_df_train_y = weighted_df_train[y_feature]
49
50    # create and train classifier
51    weighted_classifier = DecisionTreeClassifier()
52    weighted_classifier.fit(weighted_df_train_x, weighted_df_train_y)
53
54    # run classifier against the test set
55    weighted_df_test_x = weighted_df_test[features]
56    weighted_predicted = weighted_classifier.predict(weighted_df_test_x)
```

```

57 weighted_df_test_x['G3_Predicted'] = weighted_predicted
58 weighted_df_test_x['G3_Actual'] = weighted_df_test[y_feature]
59
60
61 df_test_x.to_csv('out/predicted_g3-unweighted.csv', index=False)
62 weighted_df_test_x.to_csv('out/predicted_g3-weighted.csv', index=False)
63
64 # calculate fairness based on sex
65 up_sex = 0 # male
66 p_sex = 1 # female
67 bins = [0, 11, 19]
68 labels = [0, 1]
69
70 # calculate fairness for original dataset
71 df_test_x['G3_Predicted_Discrete'] = pd.cut(df_test_x['G3_Predicted'],
72                                             include_lowest=True)
73 sex_g3_freq_og = df_test_x.groupby(['sex', 'G3_Predicted_Discrete'])
74
75 # statistical parity difference G3 and sex
76 spd_g3_sex_og = (sex_g3_freq_og[up_sex][1] / (sex_g3_freq_og[up_sex][1] +
77                                             sex_g3_freq_og[p_sex][1]) -
78                  sex_g3_freq_og[p_sex][1] / (sex_g3_freq_og[p_sex][1] +
79                                             sex_g3_freq_og[up_sex][1]))
79
80 # disparate impact G3 and sex
81 di_g3_sex_og = (sex_g3_freq_og[up_sex][1] / (sex_g3_freq_og[up_sex][1] +
82                                             sex_g3_freq_og[p_sex][1]) -
83                  sex_g3_freq_og[p_sex][1] / (sex_g3_freq_og[p_sex][1] +
84                                             sex_g3_freq_og[up_sex][1]))
85
86 # calculate fairness for transformed dataset
87 weighted_df_test_x['G3_Predicted_Discrete'] = pd.cut(weighted_df_test_x['G3_Predicted'],
88                                                         include_lowest=True)
89 sex_g3_freq_transformed = weighted_df_test_x.groupby(['sex', 'G3_Predicted_Discrete'])
90
91 # spd of weighted for weighted dataset
92 spd_g3_sex_transformed = (sex_g3_freq_transformed[up_sex][1] / (
93     sex_g3_freq_transformed[up_sex][1] + sex_g3_freq_transformed[p_sex][1]) -
94     sex_g3_freq_transformed[p_sex][1] / (sex_g3_freq_transformed[p_sex][1] +
95     sex_g3_freq_transformed[up_sex][1]))
96
97 # disparate impact G3 and sex for weighted dataset
98 di_g3_sex_transformed = (sex_g3_freq_transformed[up_sex][1] / (
99     sex_g3_freq_transformed[up_sex][1] + sex_g3_freq_transformed[p_sex][1]) -
100    sex_g3_freq_transformed[p_sex][1] / (sex_g3_freq_transformed[p_sex][1] +
101    sex_g3_freq_transformed[up_sex][1]))
102
103 print("SPD original: {}".format(spd_g3_sex_og))
104 print("DI original: {}".format(di_g3_sex_og))
105
106 print("SPD weighted: {}".format(spd_g3_sex_transformed))
107 print("DI weighted: {}".format(di_g3_sex_transformed))
108
109 outcomes = pd.DataFrame(columns=['Independent Variable', 'Metric',
110                                 'Value'])
111 outcomes.loc[len(outcomes)] = ['Sex', 'Statistical Parity Difference',
112                                spd_g3_sex_og - spd_g3_sex_transformed]
113 outcomes.loc[len(outcomes)] = ['Sex', 'Disparate Impact',
114                                di_g3_sex_og - di_g3_sex_transformed]
115
116 outcomes.to_csv('out/fairness_metrics_classified.csv', index=False)
117
118 # Plots for step 5 #TODO Make these pretty

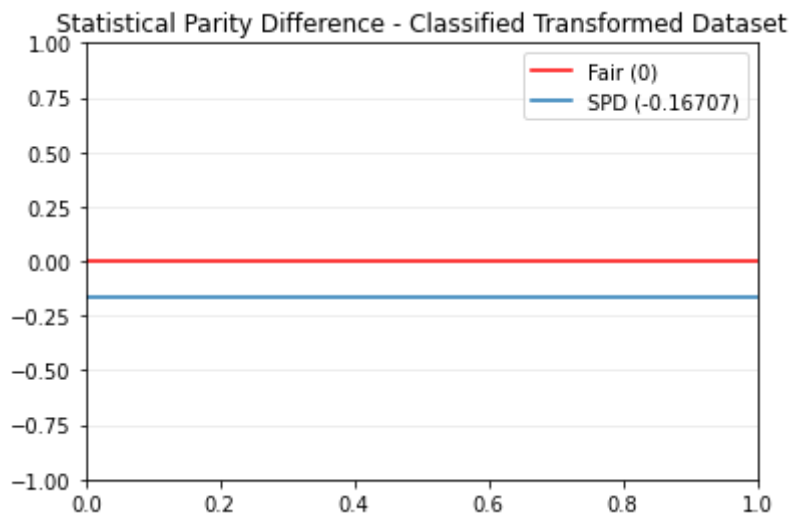
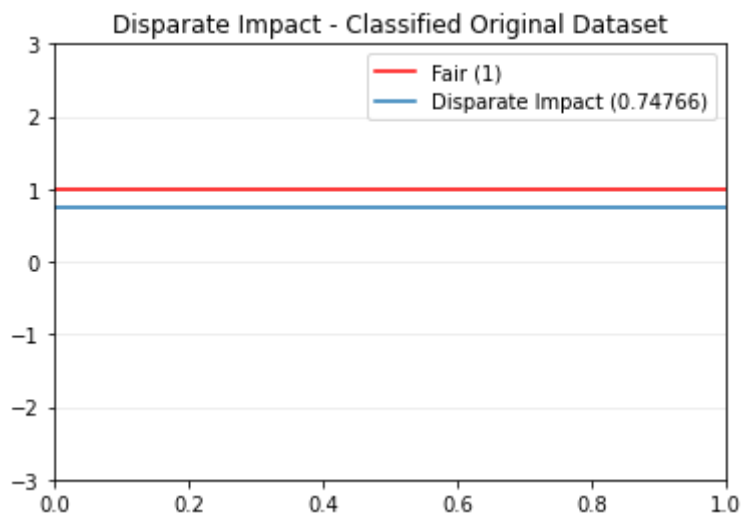
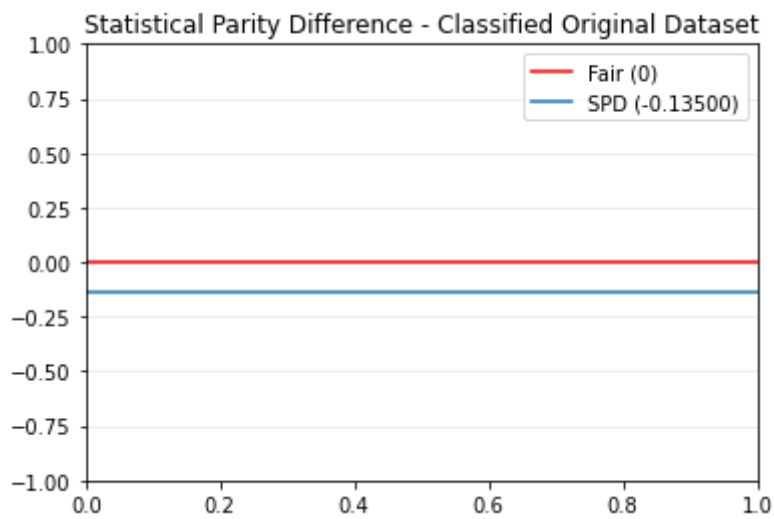
```

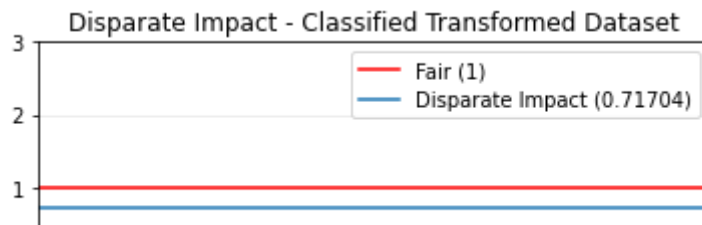
```

114 plt.clf()
115 plt.axhline(y=0, color='r', label='Fair (0)')
116 plt.axhline(spd_g3_sex_og, label="SPD ({:.5f})".format(spd_g3_sex_og))
117 plt.grid(True, axis='y', alpha=0.2, color='#999999')
118 plt.title("Statistical Parity Difference - Classified Original Dataset")
119 top = math.ceil(abs(spd_g3_sex_og) + .3)
120 plt.ylim(-top, top)
121 plt.savefig('out/spd_g3_sex_og.png', bbox_inches='tight')
122 plt.legend()
123 plt.show()
124
125 plt.clf()
126 plt.axhline(y=1, color='r', label='Fair (1)')
127 plt.axhline(di_g3_sex_og, label="Disparate Impact ({:.5f})".format(di_g3_sex_og))
128 plt.grid(True, axis='y', alpha=0.2, color='#999999')
129 plt.title("Disparate Impact - Classified Original Dataset")
130 top = math.ceil(abs(di_g3_sex_og) + .3 + 1)
131 plt.ylim(-top, top)
132 plt.legend()
133 plt.savefig('out/di_g3_sex_og.png', bbox_inches='tight')
134 plt.show()
135
136 plt.clf()
137 plt.axhline(y=0, color='r', label='Fair (0)')
138 plt.axhline(spd_g3_sex_transformed, label="SPD ({:.5f})".format(spd_g3_sex_transformed))
139 plt.grid(True, axis='y', alpha=0.2, color='#999999')
140 plt.title("Statistical Parity Difference - Classified Transformed Dataset")
141 top = math.ceil(abs(spd_g3_sex_transformed) + .3)
142 plt.ylim(-top, top)
143 plt.legend()
144 plt.savefig('out/spd_g3_sex_transformed.png', bbox_inches='tight')
145 plt.show()
146
147 plt.clf()
148 plt.axhline(y=1, color='r', label='Fair (1)')
149 plt.axhline(di_g3_sex_transformed, label="Disparate Impact ({:.5f})".format(di_g3_sex_transformed))
150 plt.grid(True, axis='y', alpha=0.2, color='#999999')
151 plt.title("Disparate Impact - Classified Transformed Dataset")
152 top = math.ceil(abs(di_g3_sex_transformed) + .3 + 1)
153 plt.ylim(-top, top)
154 plt.legend()
155 plt.savefig('out/di_g3_sex_transformed.png', bbox_inches='tight')
156 plt.show()
157
158 STEP4()

```

SPD original: -0.135
 DI original: 0.7476635514018691
 SPD weighted: -0.16706786768131704
 DI weighted: 0.7170382060893009





Step 5:

Team Members:

- Alex Carmona
- Sarah Hernandez
- Neesha Sinha
- Yaima Valdivia

Questions, Answered:

1) Explain which fairness metric (if any) is best and provide a justification for your answer:

Both the Statistical Parity Difference and the Disparate Impact metrics, used here, successfully mitigated bias in our dataset, and both are used in industry to mitigate bias on a grander scale. The difference between these two measures is subtle. The Statistical Parity Difference measures, the statistical difference in outcomes between privileged and underprivileged groups, while the disparate impact measures the ratio of outcomes. While both serve useful functions, ultimately I would pick the disparate impact measurement because it is 1) easier to explain, and 2) easier for the less statistical-minded to visualize. Because communication is key in any industry, I would pick the more communicable metric over the lesser one.

Individual Responses:

Alex Carmona:

Did any of these approaches seem to work to mitigate bias (or increase fairness)? Explain your reasoning.

Yes, after applying the reweighting strategy, the bias seemed to be mitigated. Both SPD and DI were moved closer towards being fair.

Did any group receive a positive advantage?

Yes. The weights provided an advantage towards underprivileged groups.

Was any group disadvantaged by these approaches?

No. These improvements helped provide a positive outcome for underprivileged groups without making an impact on privileged groups.

What issues would arise if you used these methods to mitigate bias?

While using the weighted values to mitigate bias provides a “prettier” picture, it does not necessarily reflect reality. Additionally, the algorithm we used to create the weights is very rudimentary and would need to be revisited if we were to expand the data set.

Sarah Hernandez:

Did any of these approaches seem to work to mitigate bias (or increase fairness)? Explain your reasoning.

Yes! Our reweighting strategy proved to be most effective. For both our metrics, Statistical Parity Difference and Disparate Impact, reweighting brought their values to their statistical ideal (0 and 1, respectively!).

Did any group receive a positive advantage?

Yes! Reweighting resulted in a more fair distribution of outcomes, resulting in a positive advantage for our underprivileged groups.

Was any group disadvantaged by these approaches?

For this particular reweighting, no group was disadvantaged by the result.

What issues would arise if you used these methods to mitigate bias?

For one, applying our reweighting algorithm to another set would not prove useful, as it is specifically designed for this dataset alone. Additionally, while our reweighting method only serves to increase fairness across multiple metrics, the broader public may not see it as such, and could result in public outcry.

Neesha Sinha:

Did any of these approaches seem to work to mitigate bias (or increase fairness)? Explain your reasoning.

Reweighting strategies did work to mitigate bias. After applying different weights to each population, the fairness metrics we considered (SPD and DI) both changed to match the fair threshold (SPD = 0 and DI = 1).

Did any group receive a positive advantage?

Since the reweighting technique used different weights for each dependent variable and protected class variable combination, each of the six groups would be deemed “fair” once the appropriate weights were applied.

Was any group disadvantaged by these approaches?

Similar to the above answer, since we used different weights for each combination, each group was fair, and no one was disadvantaged since each combination was weighed separately.

What issues would arise if you used these methods to mitigate bias?

If we use reweighting to mitigate bias, it results in an inaccurate representation of the actual data. Since different populations would be weighed differently, we’re essentially changing the data set to portray the populations on an equal footing.

Yaima Valdivia:

Did any of these approaches seem to work to mitigate bias (or increase fairness)? Explain your reasoning.

Reweighting did work to mitigate bias. The fairness metrics used achieved the fair threshold (Statistical Parity Difference = 0, Disparate Impact = 1) after the approach.

Did any group receive a positive advantage?

Both groups received a positive benefit after applying the custom reweighting techniques.

Was any group disadvantaged by these approaches?

There was no introduced disadvantage after reweighting the sets.

What issues would arise if you used these methods to mitigate bias?

The reweighting technique is ideal when we cannot change the values. Because the approach was customized to match the existing dataset, the same method could introduce biases in anything different.

In []: 1