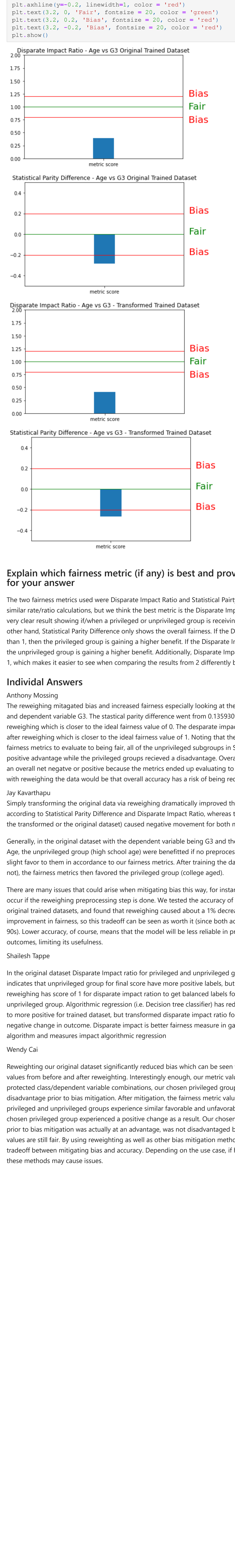


Step 3.4 Metrics

```
In [16]: # 3.4
fig, ax = plt.subplots(1, 1)
ax.bar('metric score', di_34, align='center')
ax.set_xlim(-3, 3)
ax.set_ylim(0, 2)
ax.set_title('Disparate Impact Ratio - Age/Sex vs G3')
plt.axhline(y=1, linewidth=1, color='green')
plt.axhline(y=0.8, linewidth=1, color='red')
plt.axhline(y=1.2, linewidth=1, color='red')
plt.text(3.2, 0.95, 'Fair', fontsize=20, color='green')
plt.text(3.2, 0.7, 'Bias', fontsize=20, color='red')
plt.text(3.2, 1.2, 'Bias', fontsize=20, color='red')
plt.show()

fig, ax = plt.subplots(1, 1)
ax.bar('metric score', spd_34, align='center')
ax.set_xlim(-3, 3)
ax.set_ylim(-0.5, 0.5)
ax.set_title('Statistical Parity Difference - Age/Sex vs G3')
plt.axhline(y=0, linewidth=1, color='green')
plt.axhline(y=0.2, linewidth=1, color='red')
plt.axhline(y=-0.2, linewidth=1, color='red')
plt.text(3.2, 0, 'Fair', fontsize=20, color='green')
plt.text(3.2, 0.2, 'Bias', fontsize=20, color='red')
plt.text(3.2, -0.2, 'Bias', fontsize=20, color='red')
plt.show()

print('Note: Graph looks empty because statistical parity difference is zero for Age/Sex vs G3')
```



Note: Graph looks empty because statistical parity difference is zero for Age/Sex vs G3

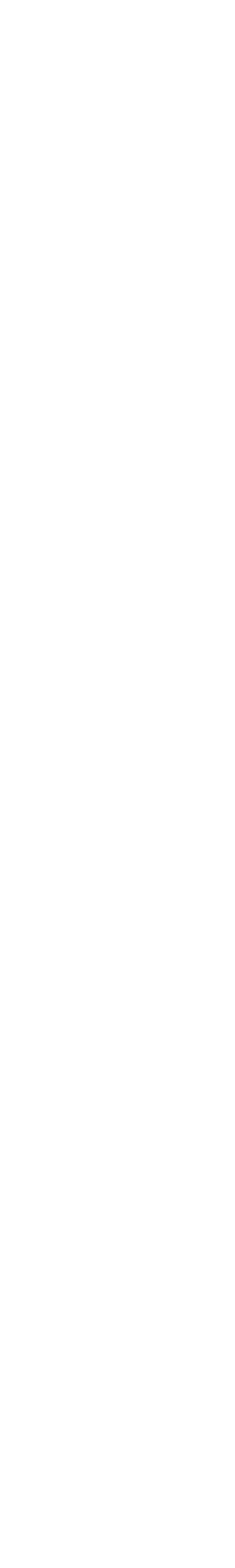
Step 4.5 Metrics

```
In [17]: # 4.5 - 1
fig, ax = plt.subplots(1, 1)
ax.bar('metric score', di_45_1, align='center')
ax.set_xlim(-3, 3)
ax.set_ylim(0, 2)
ax.set_title('Disparate Impact Ratio - Age vs G3 Original Trained Dataset')
plt.axhline(y=1, linewidth=1, color='green')
plt.axhline(y=0.8, linewidth=1, color='red')
plt.axhline(y=1.2, linewidth=1, color='red')
plt.text(3.2, 0.95, 'Fair', fontsize=20, color='green')
plt.text(3.2, 0.7, 'Bias', fontsize=20, color='red')
plt.text(3.2, 1.2, 'Bias', fontsize=20, color='red')
plt.show()

fig, ax = plt.subplots(1, 1)
ax.bar('metric score', spd_45_1, align='center')
ax.set_xlim(-3, 3)
ax.set_ylim(-0.5, 0.5)
ax.set_title('Statistical Parity Difference - Age vs G3 Original Trained Dataset')
plt.axhline(y=0, linewidth=1, color='green')
plt.axhline(y=0.2, linewidth=1, color='red')
plt.axhline(y=-0.2, linewidth=1, color='red')
plt.text(3.2, 0, 'Fair', fontsize=20, color='green')
plt.text(3.2, 0.2, 'Bias', fontsize=20, color='red')
plt.text(3.2, -0.2, 'Bias', fontsize=20, color='red')
plt.show()

# 4.5 - 2
fig, ax = plt.subplots(1, 1)
ax.bar('metric score', di_45_2, align='center')
ax.set_xlim(-3, 3)
ax.set_ylim(0, 2)
ax.set_title('Disparate Impact Ratio - Age vs G3 - Transformed Trained Dataset')
plt.axhline(y=1, linewidth=1, color='green')
plt.axhline(y=0.8, linewidth=1, color='red')
plt.axhline(y=1.2, linewidth=1, color='red')
plt.text(3.2, 0.95, 'Fair', fontsize=20, color='green')
plt.text(3.2, 0.7, 'Bias', fontsize=20, color='red')
plt.text(3.2, 1.2, 'Bias', fontsize=20, color='red')
plt.show()

fig, ax = plt.subplots(1, 1)
ax.bar('metric score', spd_45_2, align='center')
ax.set_xlim(-3, 3)
ax.set_ylim(-0.5, 0.5)
ax.set_title('Statistical Parity Difference - Age vs G3 - Transformed Trained Dataset')
plt.axhline(y=0, linewidth=1, color='green')
plt.axhline(y=0.2, linewidth=1, color='red')
plt.axhline(y=-0.2, linewidth=1, color='red')
plt.text(3.2, 0, 'Fair', fontsize=20, color='green')
plt.text(3.2, 0.2, 'Bias', fontsize=20, color='red')
plt.text(3.2, -0.2, 'Bias', fontsize=20, color='red')
plt.show()
```



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

The two fairness metrics used were Disparate Impact Ratio and Statistical Parity Difference. They have very similar rate/ratio calculations, but we think the best metric is the Disparate Impact. Disparate Impact gives a very clear result showing if when a privileged or unprivileged group is receiving a higher benefit. On the other hand, Statistical Parity Difference only shows the overall fairness. If the Disparate Impact is higher than 1, then the privileged group is gaining a higher benefit. If the Disparate Impact is lower than 1, then the unprivileged group is gaining a higher benefit. Additionally, Disparate Impact is not bounded by -1 and 1, which makes it easier to see when comparing the results from 2 differently biased groups.

Individual Answers

Anthony Mossing

The reweighting mitigated bias and increased fairness especially looking at the protected class variable Sex and dependent variable G3. The statistical parity difference went from 0.13593091 to -5.55e-17 after reweighting which is closer to the ideal fairness value of 1. Noting that the reweighting resulted in our fairness metrics to evaluate to being fair, all of the unprivileged subgroups in Sex and G3 received a positive advantage while the privileged groups received a disadvantage. Overall, no group came out with an overall net negative or positive because the metrics ended up evaluating to fair. An issue that may arise with reweighting the data would be that overall accuracy has a risk of being reduced.

Jay Kavarthapu

Simply transforming the original data via reweighting dramatically improved the fairness of the data according to Statistical Parity Difference and Disparate Impact Ratio, whereas training the data (on either the transformed or the original dataset) caused negative movement for both metrics.

Generally, in the original dataset with the dependent variable being G3 and the protected variable being Age, the unprivileged group (high school age) were benefited if no preprocessing step was done, with slight favor to them in accordance to our fairness metrics. After training the data (either transformed or not), the fairness metrics then favored the privileged group (college aged).

There are many issues that could arise when mitigating bias this way, for instance a loss in accuracy could occur if the reweighting preprocessing step is done. We tested the accuracy of both the transformed and original trained datasets, and found that reweighting caused about a 1% decrease in score, with a marginal improvement in fairness, so this tradeoff can be seen as worth it (since both accuracy scores were in the 90s). Lower accuracy, of course, means that the model will be less reliable in predicting real world outcomes, limiting its usefulness.

Shalesh Tappe

In the original dataset Disparate Impact ratio for privileged and unprivileged group is over 4, which indicates that unprivileged group for final score have more positive labels, but transferred dataset after reweighting has score of 1 for disparate impact ratio to get balanced labels for both privileged and unprivileged group. Algorithmic regression (i.e. Decision tree classifier) has reduced disparate impact ratio to more positive for trained dataset, but transformed disparate impact ratio for trained dataset has negative change in outcome. Disparate impact is better fairness measure in gauging fairness with different algorithm and measures impact algorithmic regression

Wendy Cai

Reweighting our original dataset significantly reduced bias which can be seen from the fairness metric values from before and after reweighting. Interestingly enough, our metric values show that for most of our protected class/dependent variable combinations, our chosen privileged group are actually at a disadvantage prior to bias mitigation. After mitigation, the fairness metric values suggest that both the privileged and unprivileged groups experience similar favorable and unfavorable outcomes, therefore our chosen privileged group experienced a positive change as a result. Our chosen unprivileged group, which prior to bias mitigation was actually at an advantage, was not disadvantaged because the resulting metric values are still fair. By using reweighting as well as other bias mitigation methods, there is sometimes a tradeoff between mitigating bias and accuracy. Depending on the use case, if high accuracy is critical, then these methods may cause issues.