

Supplemental Material for Fair Feature Importance Scores for Interpreting Decision Trees

Reviewed on OpenReview: <https://openreview.net/forum?id=XXXX>

A Proof of Proposition 1

Proof. By the Total Law of Expectation, we have that

$$\begin{aligned} Bias^{DP}(lev(t)) = & \left| \left(E(\hat{y}_i | i \in lev_\ell(t), z_i = 1) P(i \in lev_\ell(t) | z_i = 1) \right. \right. \\ & \left. \left. + E(\hat{y}_i | i \in lev_r(t), z_i = 1) P(i \in lev_r(t) | z_i = 1) \right) \right. \\ & - \left(E(\hat{y}_i | i \in lev_\ell(t), z_i = 0) P(i \in lev_\ell(t) | z_i = 0) \right. \\ & \left. \left. + E(\hat{y}_i | i \in lev_r(t), z_i = 0) P(i \in lev_r(t) | z_i = 0) \right) \right|. \end{aligned}$$

Notice that the expectation of $\hat{y}_i \in lev_\ell(t) = \pi_{lev_\ell(t)}$. Replacing the expectation terms with $\pi_{lev_\ell(t)}$ and $\pi_{lev_r(t)}$, we can see that

$$\begin{aligned} Bias^{DP}(lev(t)) = & \left| \pi_{lev_\ell(t)} P(i \in lev_\ell(t) | z_i = 0) + \pi_{lev_r(t)} P(i \in lev_r(t) | z_i = 0) \right. \\ & \left. - \pi_{lev_\ell(t)} P(i \in lev_\ell(t) | z_i = 1) - \pi_{lev_r(t)} P(i \in lev_r(t) | z_i = 1) \right| \\ = & \left| \frac{\sum_i \mathbb{1}_{\{z_i=1, i \in lev_\ell(t)\}} * \pi_{lev_\ell(t)} + \sum_i \mathbb{1}_{\{z_i=1, i \in lev_r(t)\}} * \pi_{lev_r(t)}}{\sum_i \mathbb{1}_{\{z_i=1, i \in lev(t)\}}} \right. \\ & \left. - \frac{\sum_i \mathbb{1}_{\{z_i=0, i \in lev_\ell(t)\}} * \pi_{lev_\ell(t)} + \sum_i \mathbb{1}_{\{z_i=0, i \in lev_r(t)\}} * \pi_{lev_r(t)}}{\sum_i \mathbb{1}_{\{z_i=0, i \in lev(t)\}}} \right|. \end{aligned}$$

Combining similar terms and simplifying, we have that

$$\begin{aligned} Bias^{DP}(lev(t)) = & \left| \pi_{lev_\ell(t)} \left(\frac{\sum_i \mathbb{1}_{\{z_i=1, i \in lev_\ell(t)\}}}{\sum_i \mathbb{1}_{\{z_i=1, i \in lev(t)\}}} - \frac{\sum_i \mathbb{1}_{\{z_i=0, i \in lev_\ell(t)\}}}{\sum_i \mathbb{1}_{\{z_i=0, i \in lev(t)\}}} \right) \right. \\ & \left. + \pi_{lev_r(t)} \left(\frac{\sum_i \mathbb{1}_{\{z_i=1, i \in lev_r(t)\}}}{\sum_i \mathbb{1}_{\{z_i=1, i \in lev(t)\}}} - \frac{\sum_i \mathbb{1}_{\{z_i=0, i \in lev_r(t)\}}}{\sum_i \mathbb{1}_{\{z_i=0, i \in lev(t)\}}} \right) \right|. \end{aligned}$$

□

The proof for $E[Bias^{EQOP}(lev(t))]$ is analogous to the proof for demographic parity.

For the multiclass classification case, let $\pi_{lev_\ell(t)}^m$ and $\pi_{lev_r(t)}^m$ denote a vector of length K , where π_k denotes the proportion of that class in the node.

Corollary 1. *Consider multiclass classification with probabilistic trees:*

$$\begin{aligned} Bias^{DP}(lev(t)) = & \left| \pi_{lev_\ell(t)}^m \left(\frac{\sum_i \mathbb{1}_{\{z_i=1, i \in lev_\ell(t)\}}}{\sum_i \mathbb{1}_{\{z_i=1, i \in lev(t)\}}} - \frac{\sum_i \mathbb{1}_{\{z_i=0, i \in lev_\ell(t)\}}}{\sum_i \mathbb{1}_{\{z_i=0, i \in lev(t)\}}} \right) \right. \\ & \left. + \pi_{lev_r(t)}^m \left(\frac{\sum_i \mathbb{1}_{\{z_i=1, i \in lev_r(t)\}}}{\sum_i \mathbb{1}_{\{z_i=1, i \in lev(t)\}}} - \frac{\sum_i \mathbb{1}_{\{z_i=0, i \in lev_r(t)\}}}{\sum_i \mathbb{1}_{\{z_i=0, i \in lev(t)\}}} \right) \right| \end{aligned}$$

B Additional Simulation Results

B.1 Classification Results

We evaluate our method on the same experiments as Figure 2 in the main paper but in Figure A1, $N = 500$ and in Figures A2 and A3, we use Equality of Opportunity as the fairness metric. In Figure A4, we consider a simulation with a large number of features with $p = 250$. In the large p simulation, there are 5 features in each group and we otherwise follow the same setting as our other classification simulations. In the correlated simulations, shown in Figure A5, we use an autoregressive design with $\Sigma_{j,j+1}^{-1} = 0.5$ instead of $\Sigma = \mathbf{I}$ in the uncorrelated p simulations. Similar to the results in the main paper, we see the correct magnitude and direction of the scores in all of the simulation scenarios.

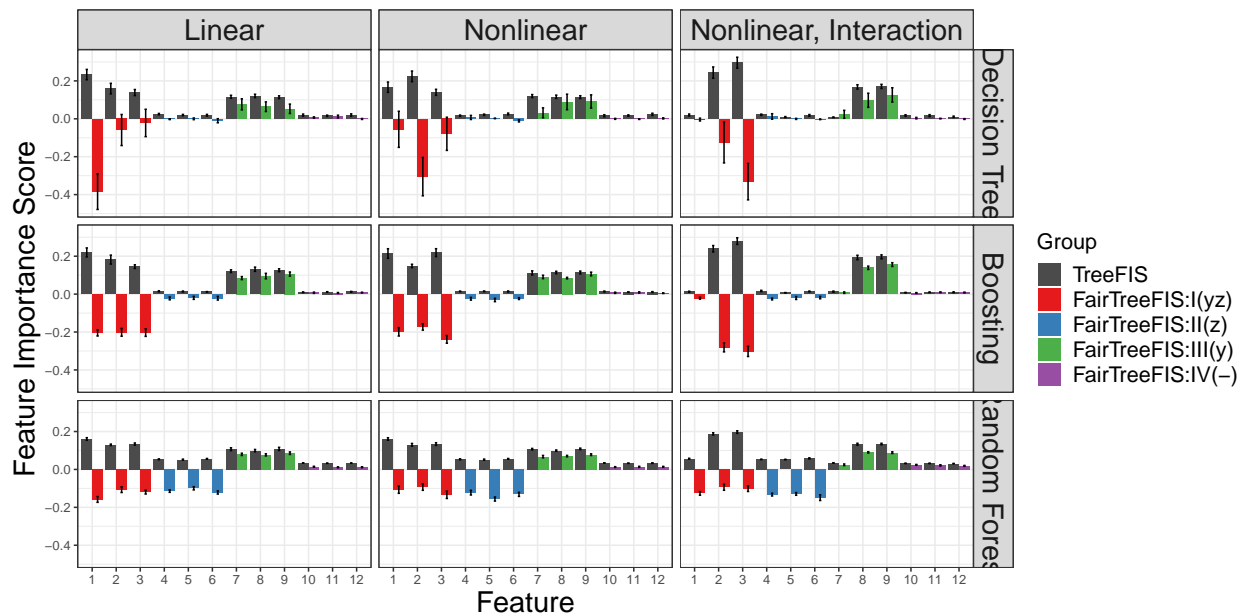


Figure A1: Classification *TreeFIS* and *FairTreeFIS* results for accuracy and Demographic Parity on three major simulation types that include a linear model (left), a non-linear additive model (middle), and a non-linear additive model with pairwise interactions (right), with $N = 500$ and $p = 12$. We examine a decision tree classifier, a boosting classifier, and a random forest classifier.

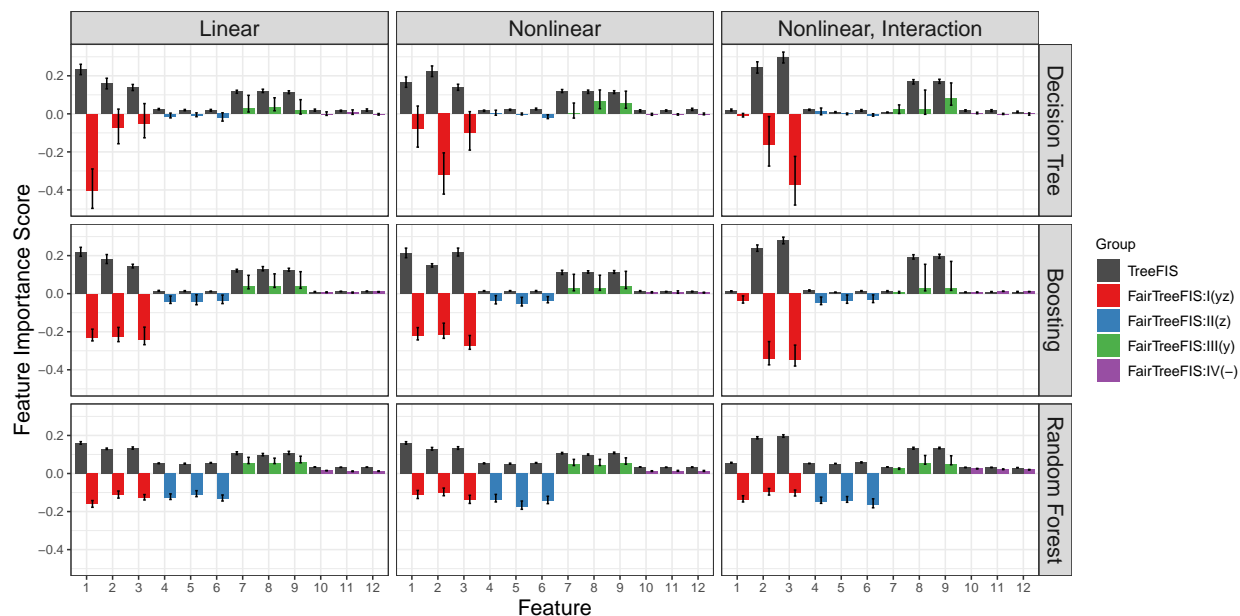


Figure A2: Classification *TreeFIS* and *FairTreeFIS* results for accuracy and Equality of Opportunity on three major simulation types that include a linear model (left), a non-linear additive model (middle), and a non-linear additive model with pairwise interactions (right), with $N = 500$ and $p = 12$. We examine a decision tree classifier, a boosting classifier, and a random forest classifier.

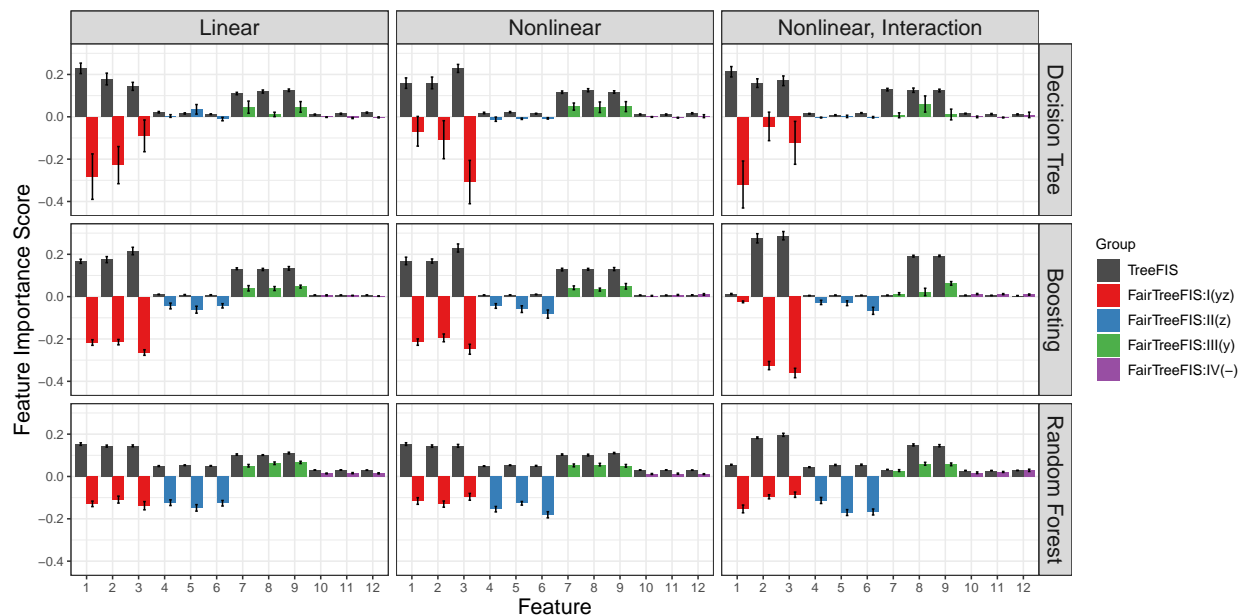


Figure A3: Classification *TreeFIS* and *FairTreeFIS* results for accuracy and Equality of Opportunity on three major simulation types that include a linear model (left), a non-linear additive model (middle), and a non-linear additive model with pairwise interactions (right), with $n = 1000$ and $p = 12$. We examine a decision tree classifier, a boosting classifier, and a random forest classifier.

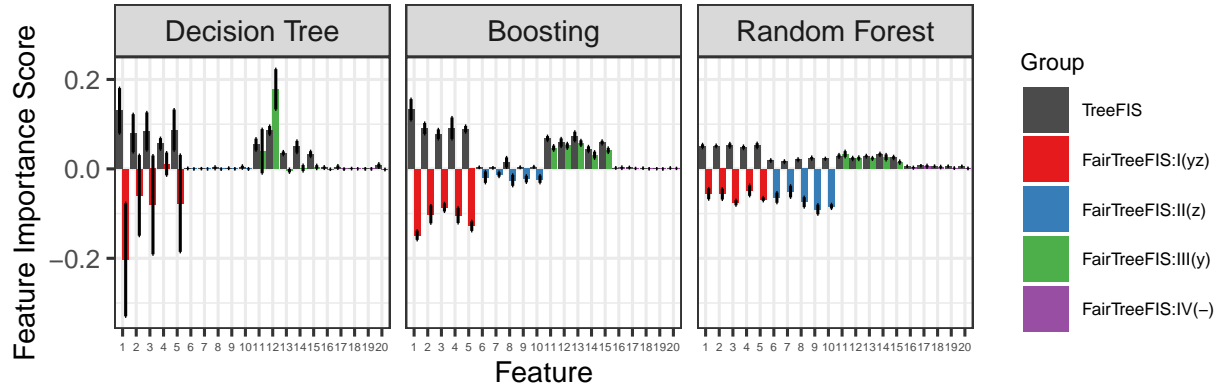


Figure A4: Large p classification *TreeFIS* and *FairTreeFIS* results for accuracy and Demographic Parity for a decision tree classifier, a boosting classifier, and a random forest classifier, with $N = 1000$ and $p = 250$. We show the *TreeFIS* and *FairTreeFIS* scores for the first 20 features.

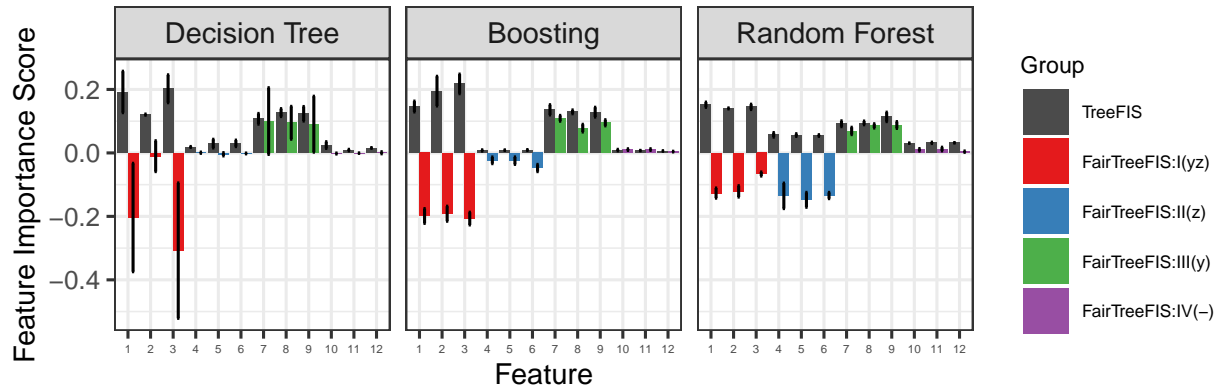


Figure A5: Correlated feature classification *TreeFIS* and *FairTreeFIS* results for accuracy and Demographic Parity for a decision tree classifier, a boosting classifier, and a random forest classifier, with $N = 1000$ and $p = 12$.

B.2 Regression Results

In Figure A6, we evaluate *FairTreeFIS* results for Demographic Parity in the regression setting. Here, $\beta_j = 3$ for $j \in G_1$ or G_3 and $\beta_j = 0$ for $j \in G_2$ or G_4 and $\alpha_j = 0.4$ for $j \in G_1$ or G_2 and $\alpha_j = 0$ for $j \in G_3$ or G_4 . All other aspects of the base simulation as described in the main paper remain the same. Similar to the results in the main paper and the additional classification results, the magnitudes and directions of the scores are as expected from the simulation design.

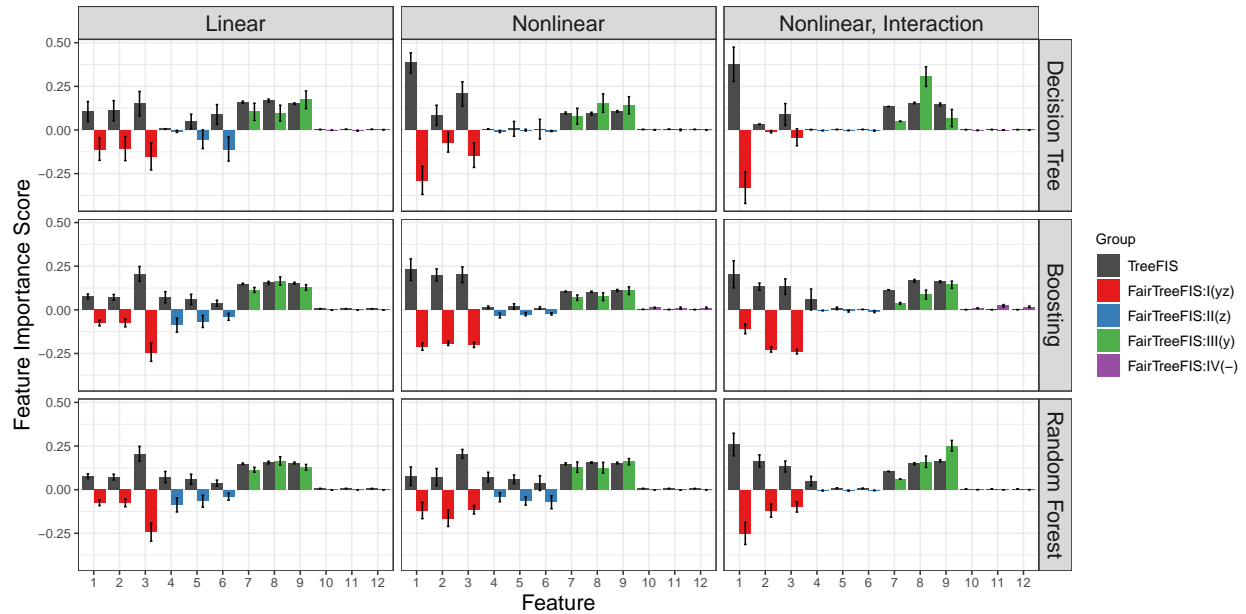


Figure A6: Regression *TreeFIS* and *FairTreeFIS* results for accuracy and Demographic Parity on three major simulation types that include a linear model (left), a non-linear additive model (middle), and a non-linear additive model with pairwise interactions (right), with $N = 500$ and $p = 12$. We examine a decision tree regressor, a boosting regressor, and a random forest regressor.

C Additional Results on Benchmark Datasets

We include the same experiment as Figure 6 from the main paper for the C & C dataset with Race as the protected attribute and the German dataset with Gender as the protected attribute in order to validate the use of global surrogates. We see that the magnitudes and the directions between the scores of the boosting classifier and the tree-based surrogate of the boosting classifier are similar.

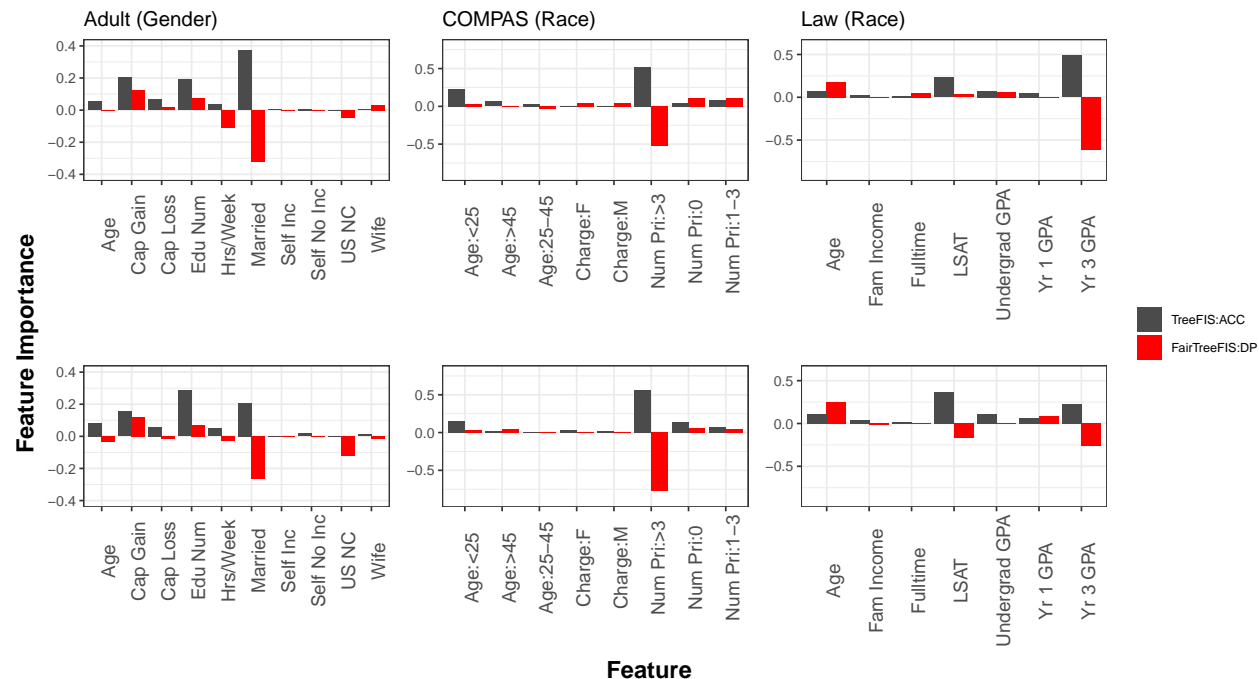


Figure A7: Global surrogate validation. The top row shows *TreeFIS* and *FairTreeFIS* results on a boosting classifier for the C & C dataset with Race as the protected attribute and the German dataset with Gender as the protected attribute. The bottom row shows *TreeFIS* and *FairTreeFIS* results for a tree-based surrogate of a boosting classifier. The scores between the top and bottom rows are similar in magnitude and direction, indicating that our scores are effective when used to interpret tree-based global surrogates.

In Figure A8, we explore the quality of *FairTreeFIS* interpretations of tree-based surrogates of a deep learning model (multi-layer perceptron with two hidden layers each with p units and ReLU activation) on the German dataset with Gender as the protected attribute and the Law School dataset with Race as the protected attribute. As shown in the main paper when discussing Figure 4, the *FairTreeFIS* results provide reasonable feature interpretations in terms of fairness.

In order to validate using trees for interpretation of deep learning models versus model-specific interpretation, we compare *TreeFIS* scores of a tree-based surrogate of a deep learning model (multi-layer perceptron with two hidden layers each with p units and ReLU activation) to scores from Layerwise Relevance Propagation (LRP) Montavon et al. (2019) of the same deep learning model for the Adult dataset with Gender as the protected attribute, the Law School dataset with Race as the protected attribute, the COMPAS dataset with Race as the protected attribute, and the German dataset with Gender as the protected attribute as shown in Figure A9. We implement LRP using the DeepExplain package with “elrp” set as the method name. We set the first layer of the MLP as the input layer and the last layer as the output. For all the datasets, we see that in general the magnitude of the importance scores for the tree surrogate and LRP surrogate are comparable. Specifically, both methods identify the same features as highly predictive, as reflected in the magnitude of the scores. These results validate that we can reasonably use trees for interpretation versus model-specific validation Lundberg & Lee (2017).

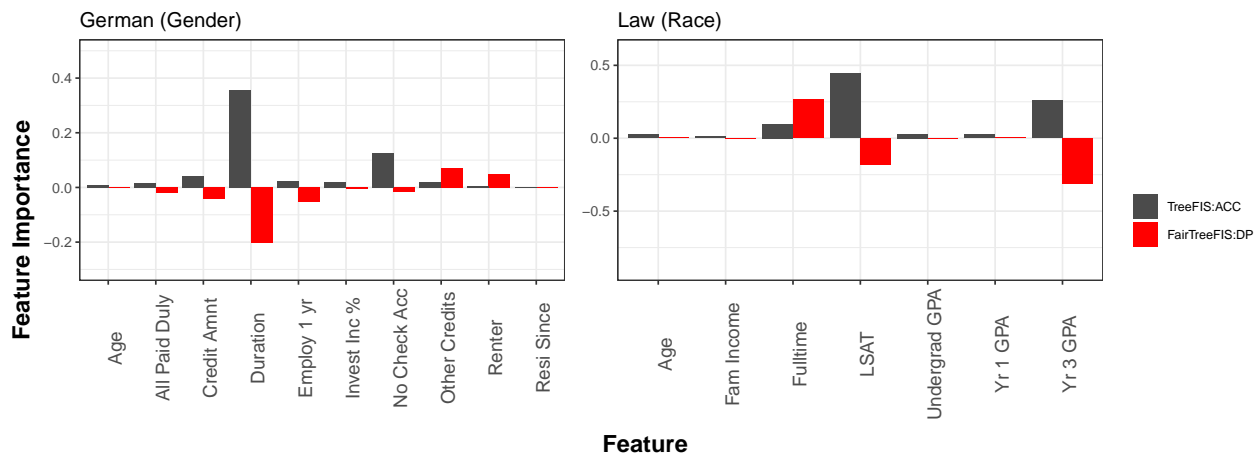


Figure A8: Importance scores for a tree-based surrogate of a deep learning model for the German dataset with Gender as the protected attribute (left) and Law School dataset with Race as the protected attribute (right).

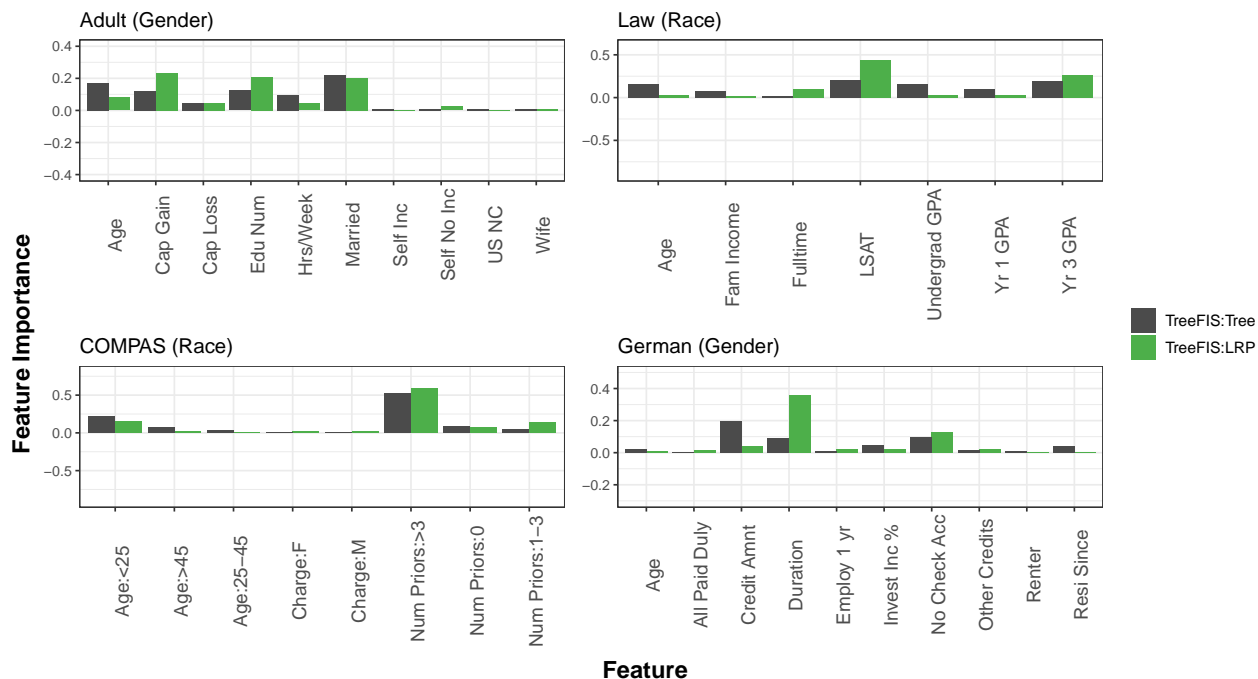


Figure A9: Validation for using trees as surrogates. For the Adult dataset with Gender as the protected attribute, the Law dataset with Race as the protected attribute, the COMPAS dataset with Race as the protected attribute, and the German dataset with Gender as the protected attribute, we show *TreeFIS* scores for a tree-based surrogate of an MLP and an LRP surrogate. The magnitudes between the two methods are similar, validating we can use trees for interpreting deep learning models.

References

- Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/file/8a20a8621978632d76c43dfd28b67767-Paper.pdf>.
- Grégoire Montavon, Alexander Binder, Sebastian Lapuschkin, Wojciech Samek, and Klaus-Robert Müller. Layer-wise relevance propagation: an overview. *Explainable AI: interpreting, explaining and visualizing deep learning*, pp. 193–209, 2019.