

Introduction & Motivation

- Automated ML/DL classification algorithms useful for Medicaid Eligibility Determination, but suffer from limitation and algorithmic bias due to a variety of factors (e.g. training data, algorithmic design)
- Fairgroup Construction to reduce unfairness in classification outcome. Fairness boosted through pre-processing the testing data before running actual classification model.
- Model agnostic to the specifics of classifier; can be generalized to other social decision problems such as Credit Card Approval and College Admission.

Definition of Fairness and preliminaries

Notion of Fairness: derived from legal doctrine of **Disparate Impact**, which calls for balanced representation of different classes. Here **balance** is simply the ratio of smaller class to larger class, and ranges from 0 to 1.

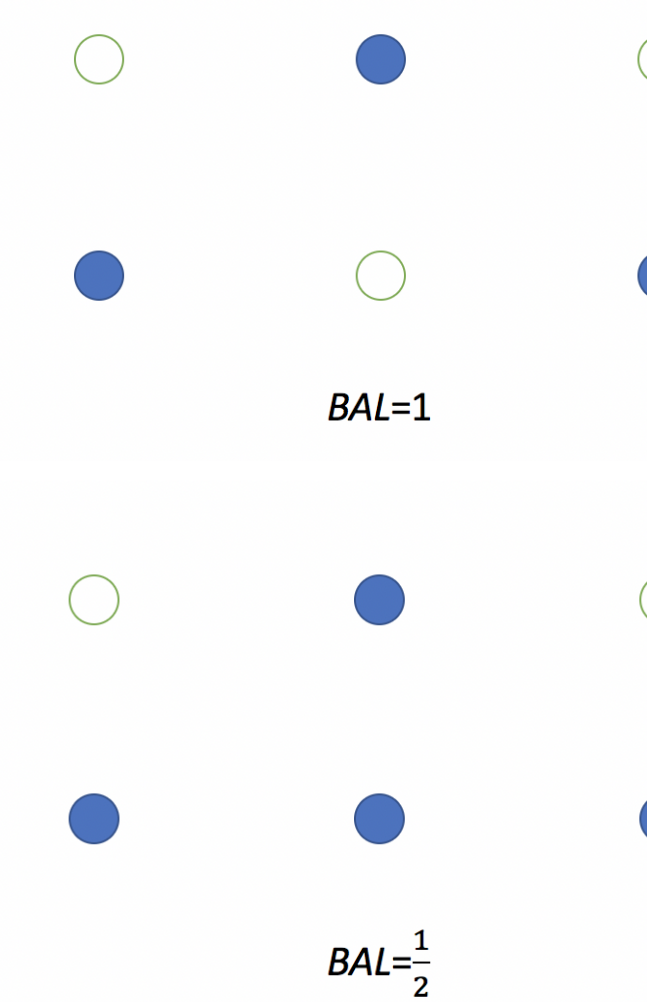


Fig. 1: Illustration of balance

We also notice that different features carry different levels of significance, and we can determine the importance of each feature in each data point from this observation. We construct the feature importance vector by computing the correlation between each numerical feature vector and the final decision vector. **The feature importance vector** encodes all such importance vectors, and will be used for subsequent models.

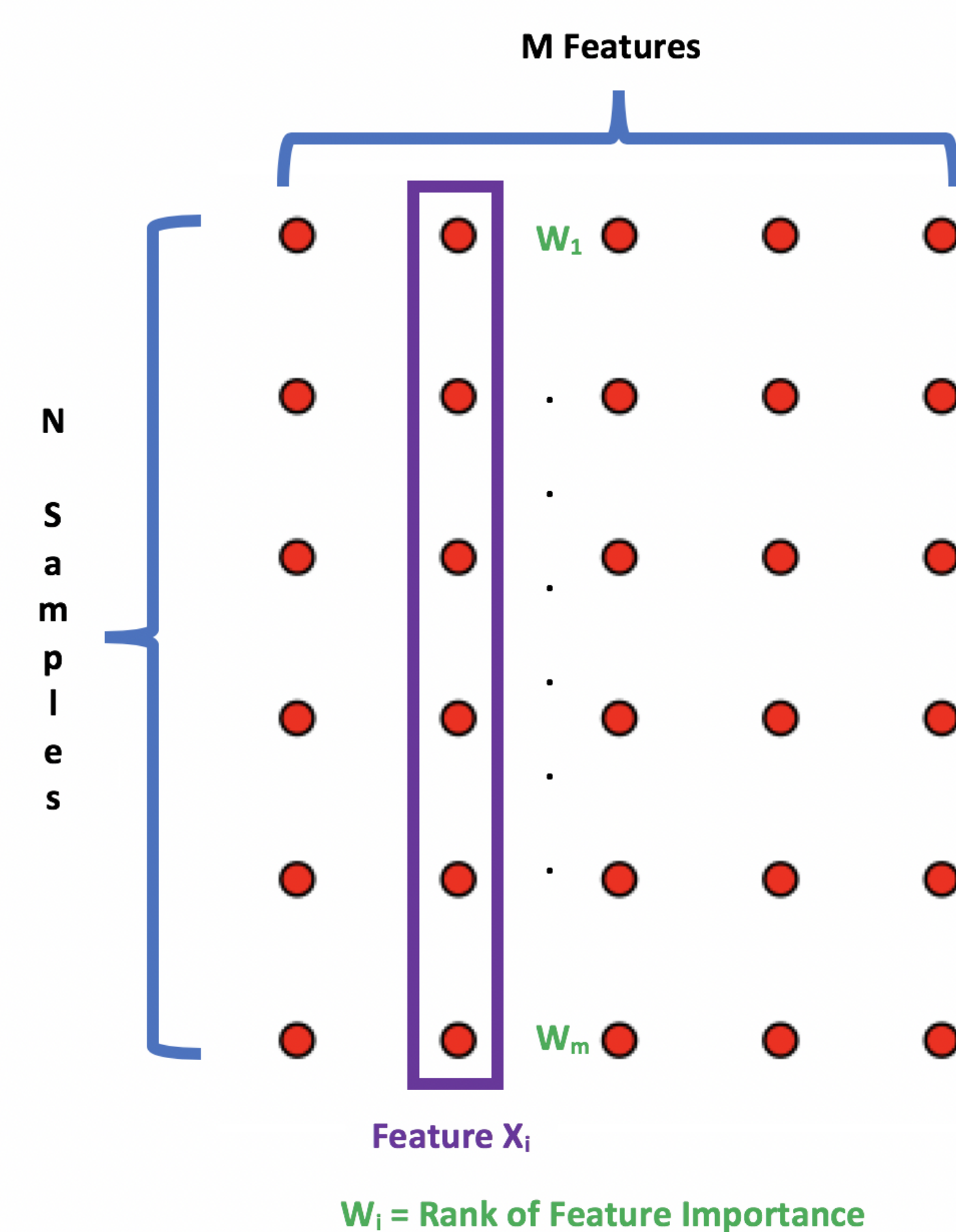


Fig. 2: Feature importance vectors

Fairness Model Demonstration

Our algorithm consists of three steps:

- K-clustering to ensure similarity;
- Intra-cluster fairgroup construction to ensure fairness;
- Actual classification to note the properties of original classifier.

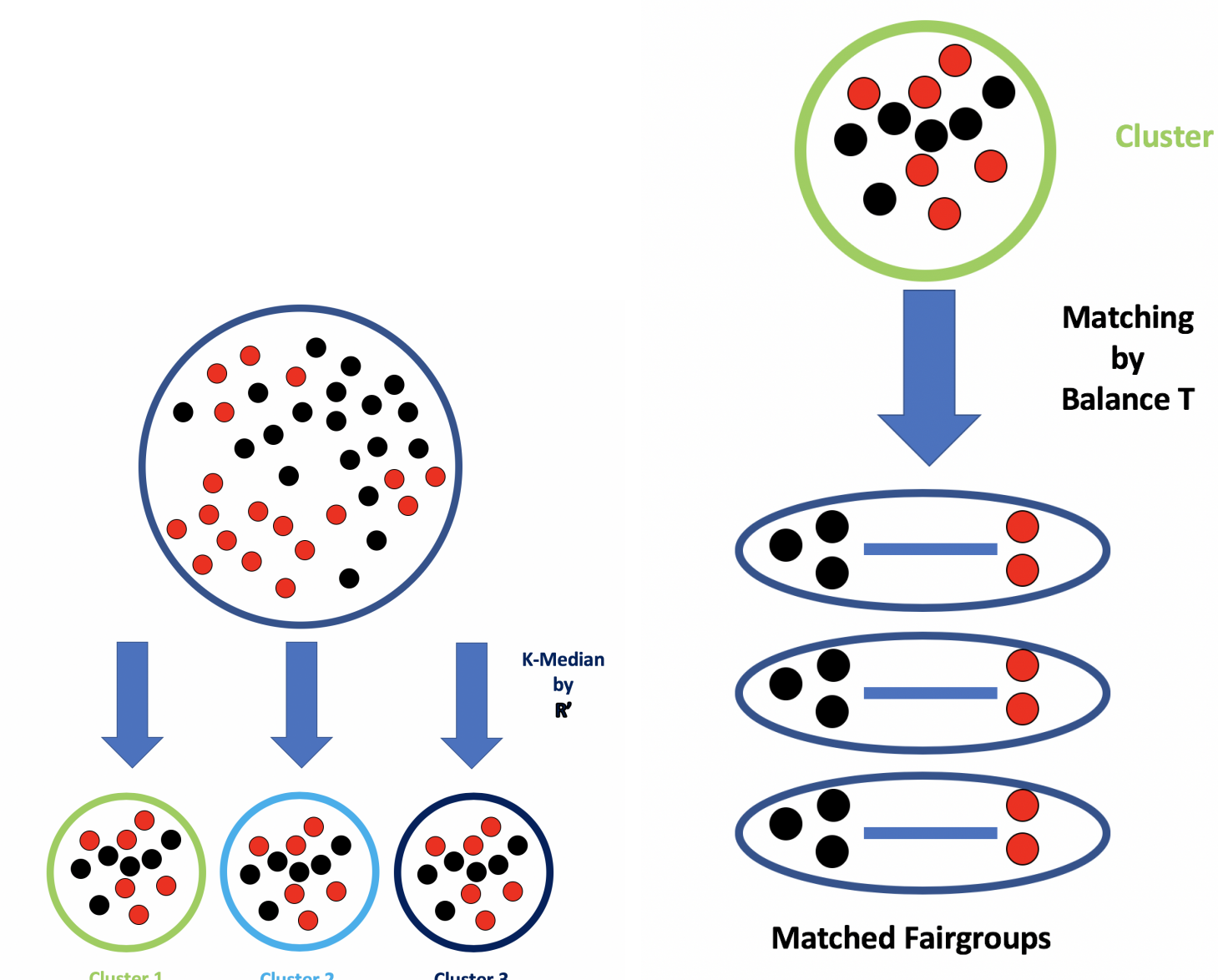


Fig. 3: K-Clustering and Fairgroup Construction

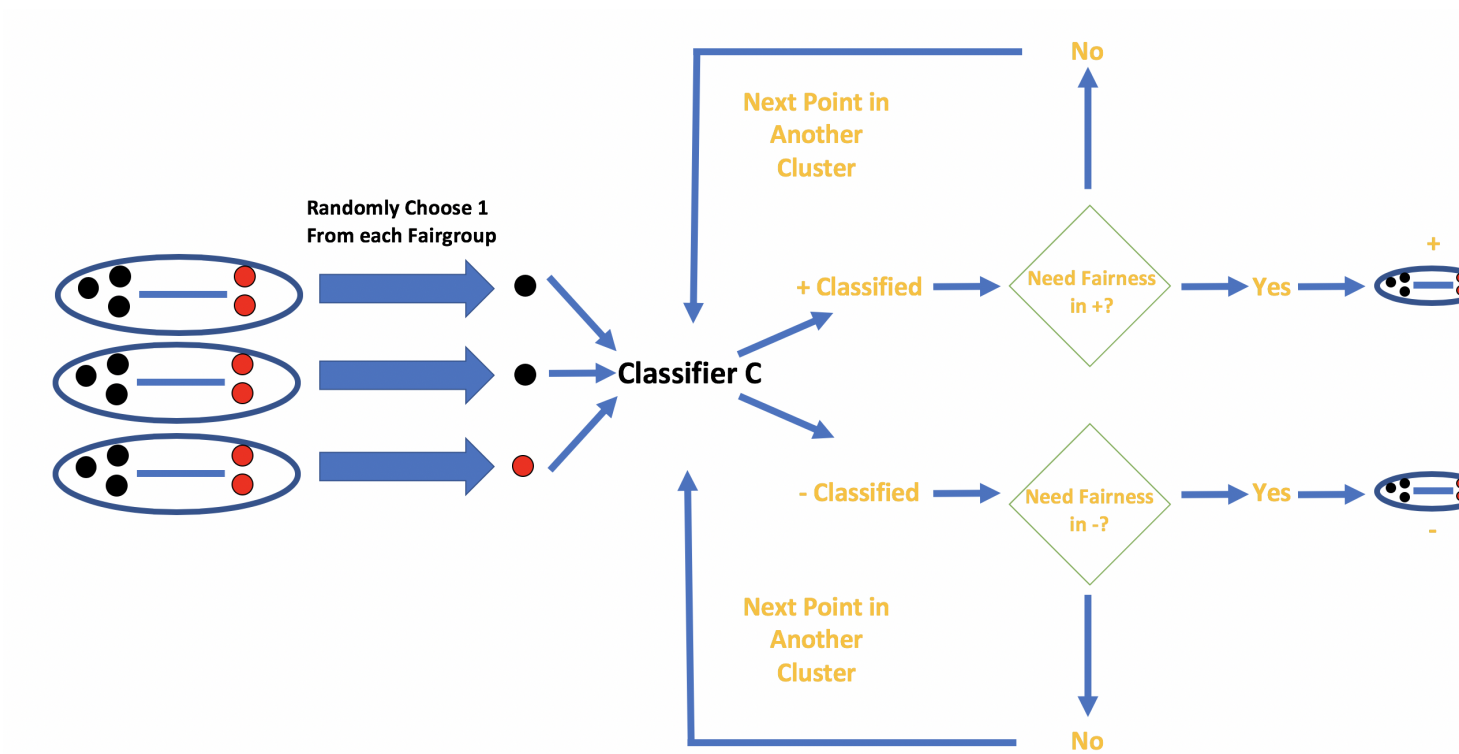


Fig. 4: Actual Classification

Experimental Results

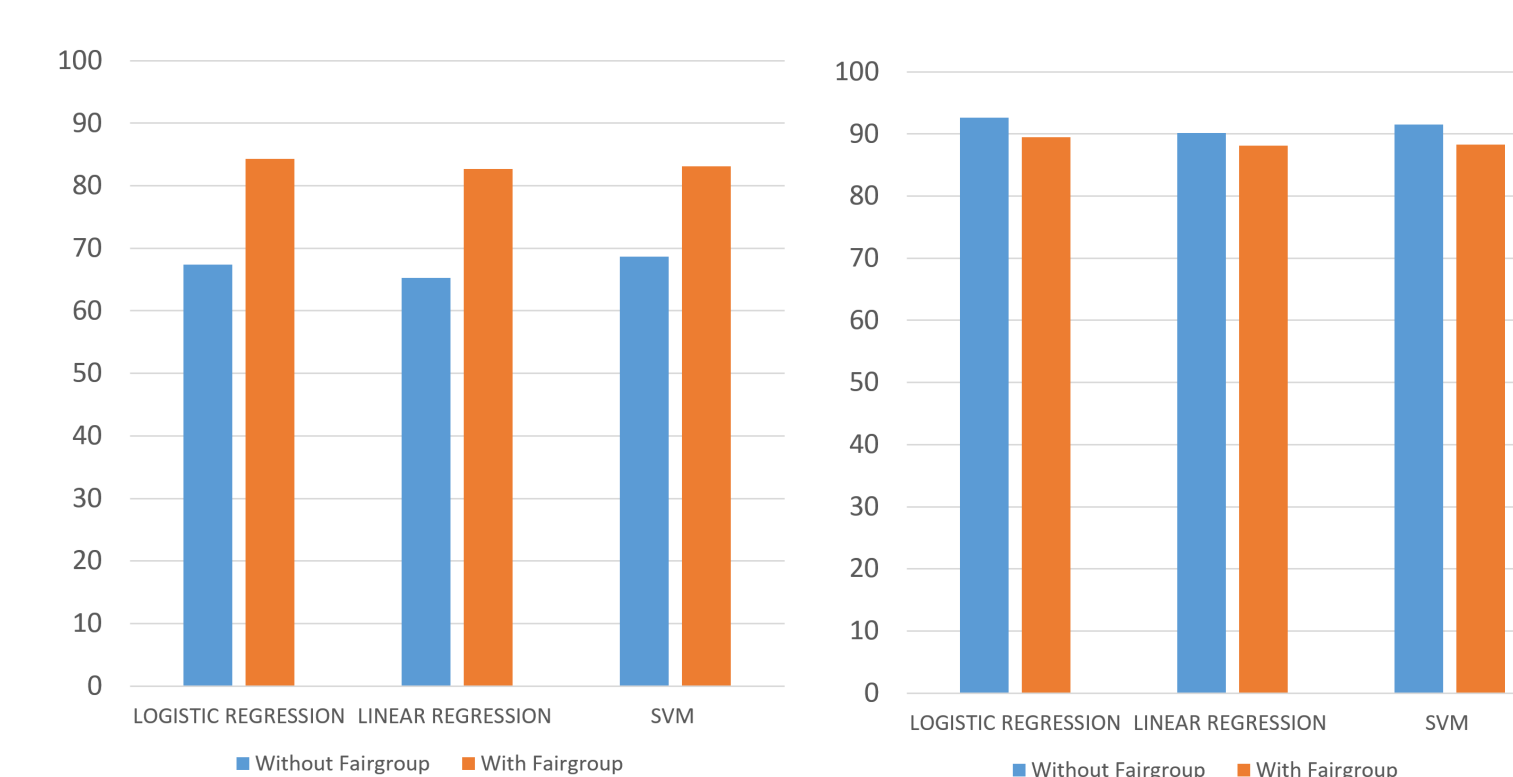


Fig. 5: Fairness and accuracy comparison - Poverty

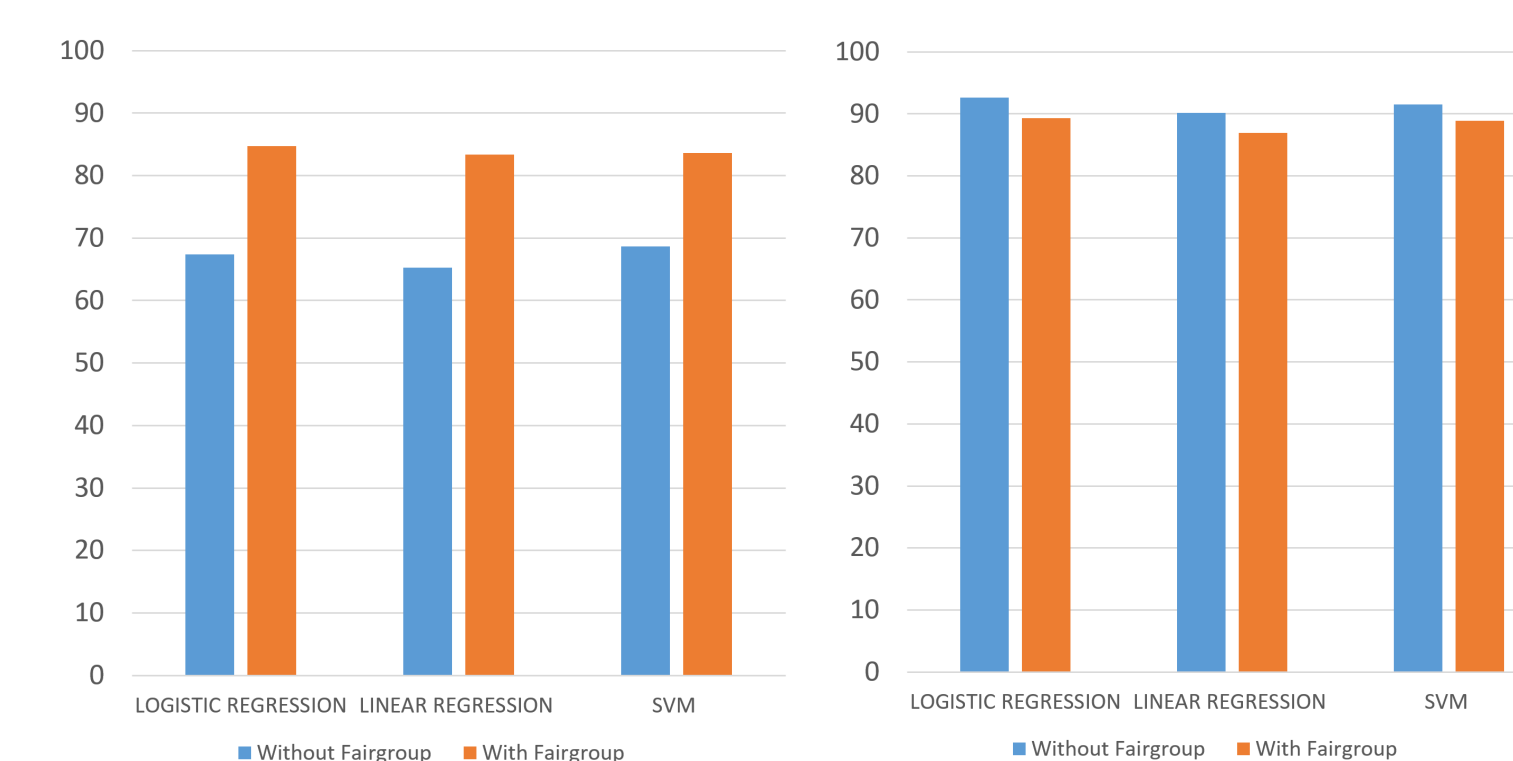


Fig. 6: Fairness and accuracy comparison - Income

Algorithm

Algorithm 1: Fairness machine learning algorithm

Result: Predicted decisions of data points
 Construct the feature importance vector r_i for each data point;
 Form K Clusters for r_i 's with K-Median algorithm;
while \exists points unmatched **do**
 Make match for the groups by balance t ;
 if \exists no more unmatched points **then**
 break;
 else
 continue matching;
 end
end
for \forall fair group **do**
 Randomly pick a point;
 result=classification(random point);
 Plus-fair = False;
 Minus-fair = False;
 if Fairness Required for '+' **then**
 Plus-fair = True;
 else
 Fairness Required for '-' Minus-fair = True;
 if Plus-fair **then**
 if result = positive **then**
 for each point of the group **do**
 prediction result=positive;
 end
 else
 for each point of the group **do**
 prediction result=classification(point);
 end
 end
 else
 Minus-fair **if** result = negative **then**
 for each point of the group **do**
 prediction result=negative;
 end
 else
 for each point of the group **do**
 prediction result=classification(point);
 end
 end
end
return Decisions of each point
end

References(selected)

- [1] US Census Bureau. 2017. American Community Survey 2017 5-year Estimate. (2017). <https://www.census.gov/programs-surveys/acs/>
- [2] Flavio Chierichetti, Ravi Kumar, Silvio Lattanzi, and Sergei Vassilvitskii. 2017. Fair clustering through fairlets. In Advances in Neural Information Processing Systems. 5029–5037
- [3] Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. 2015. Certifying and removing disparate impact. In Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM, 259–268.
- [4] Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. 2016. Inherent trade-offs in the fair determination of risk scores. arXiv preprint arXiv:1609.05807