

1. Distance correlation is superior to Maximum mean discrepancy, Mutual information and adversarial de-biasing

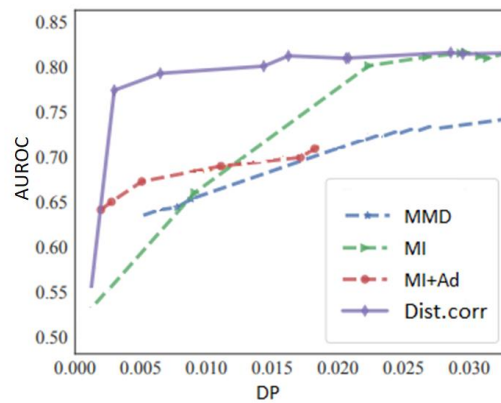


Figure 1 : Comparison of AUROC and DP for various measures

2. Experiment on CelebA data (for grey-scale CelebA with $n_channels=1$ and the results are shown from 600th epoch due to time constraints)



Case 1 – Fairness penalty computed for the whole 'gender' attribute.



Case 2 – Fairness penalty computed for only 'male' sensitive sub-groups



Case 3 – Fairness penalty computed for only 'female' sensitive sub-groups

Figure 2 : Fairness experiment on CelebA dataset

Remarks – It is interesting note that when we enforced fairness for female sub-groups (that is fairness penalty added for male, case 2 in above figure), all the male images are transformed to females and vice versa (case 3).