Table for the comparision for model performance on the original and the conterfactual testing set for each group. "Or_" represents original testing set results, "CF_" represents counterfactual testing set results, "F" for female, "M" for male, "Diff" for the difference between two groups.

| | | | SP | TPR | FPR | Acc | F1 | AUC |
|---|---|---|---|---|---|---|---|---|
| CF | Prompt1 | Or_F | 0.6606 | 0.7451 | 0.5215 | 0.6443 | 0.7227 | 0.6118 |
| | | CF_F | 0.6992 | 0.7974 | 0.5376 | 0.6707 | 0.7508 | 0.6299 |
| | | Diff | 0.0386 | 0.0523 | 0.0161 | 0.0264 | 0.0281 | 0.0181 |
| | | Or_M | 0.6446 | 0.7978 | 0.4762 | 0.6673 | 0.7152 | 0.6608 |
| | | CF_M | 0.6144 | 0.7726 | 0.4405 | 0.6711 | 0.7110 | 0.6660 |
| | | Diff | 0.0302 | 0.0252 | 0.0357 | 0.0038 | 0.0042 | 0.0052 |
| | Prompt2 | Or_F | 0.6504 | 0.7418 | 0.5000 | 0.6504 | 0.7252 | 0.6209 |
| | | CF_F | 0.6911 | 0.7908 | 0.5269 | 0.6707 | 0.7492 | 0.6320 |
| | | Diff | 0.0407 | 0.0490 | 0.0269 | 0.0203 | 0.0240 | 0.0111 |
| | | Or_M | 0.6163 | 0.7690 | 0.4484 | 0.6654 | 0.7065 | 0.6603 |
| | | CF_M | 0.6011 | 0.7726 | 0.4127 | 0.6843 | 0.7193 | 0.6799 |
| | | Diff | 0.0152 | 0.0036 | 0.0357 | 0.0189 | 0.0128 | 0.0196 |
| | Prompt3 | Or_F | 0.7988 | 0.8627 | 0.6935 | 0.6524 | 0.7553 | 0.5846 |
| | | CF_F | 0.8455 | 0.9085 | 0.7419 | 0.6626 | 0.7700 | 0.5832 |
| | | Diff | 0.0467 | 0.0458 | 0.0484 | 0.0102 | 0.0147 | 0.0014 |
| | | Or_M | 0.7996 | 0.8953 | 0.6944 | 0.6144 | 0.7086 | 0.6004 |
| | | CF_M | 0.7618 | 0.8664 | 0.6468 | 0.6219 | 0.7059 | 0.6098 |
| | | Diff | 0.0378 | 0.0289 | 0.0476 | 0.0075 | 0.0027 | 0.0094 |
| | Prompt4 | Or_F | 0.6728 | 0.7680 | 0.5161 | 0.6606 | 0.7378 | 0.6259 |
| | | CF_F | 0.7297 | 0.8300 | 0.5645 | 0.6809 | 0.7639 | 0.6328 |
| | | Diff | 0.0569 | 0.0620 | 0.0484 | 0.0203 | 0.0261 | 0.0069 |
| | | Or_M | 0.6541 | 0.8087 | 0.4841 | 0.6692 | 0.7191 | 0.6623 |
| | | CF_M | 0.6125 | 0.7906 | 0.4167 | 0.6919 | 0.7288 | 0.6870 |
| | | Diff | 0.0416 | 0.0181 | 0.0674 | 0.0227 | 0.0097 | 0.0247 |
| | Prompt5 | Or_F | 0.5752 | 0.6471 | 0.4570 | 0.6077 | 0.6723 | 0.5950 |
| | | CF_F | 0.6179 | 0.7026 | 0.4785 | 0.6341 | 0.7049 | 0.6121 |
| | | Diff | 0.0427 | 0.0555 | 0.0215 | 0.0264 | 0.0326 | 0.0171 |
| | | Or_M | 0.6011 | 0.7437 | 0.4444 | 0.6541 | 0.6924 | 0.6496 |
| | | CF_M | 0.5595 | 0.7040 | 0.4008 | 0.6541 | 0.6806 | 0.6516 |
| | | Diff | 0.0416 | 0.0397 | 0.0436 | 0.0000 | 0.0118 | 0.0020 |
| | Prompt6 | Or_F | 0.5874 | 0.6765 | 0.4409 | 0.6321 | 0.6958 | 0.6178 |
| | | CF_F | 0.6809 | 0.7712 | 0.5323 | 0.6565 | 0.7363 | 0.6195 |
| | | Diff | 0.0935 | 0.0947 | 0.0914 | 0.0244 | 0.0405 | 0.0017 |
| | | Or_M | 0.6673 | 0.8051 | 0.5159 | 0.6522 | 0.7079 | 0.6446 |
| | | CF_M | 0.5917 | 0.7292 | 0.4405 | 0.6484 | 0.6847 | 0.6444 |
| | | Diff | 0.0756 | 0.0759 | 0.0754 | 0.0038 | 0.0232 | 0.0002 |
| | Prompt7 | Or_F | 0.5427 | 0.6307 | 0.3978 | 0.6199 | 0.6736 | 0.6164 |
| | | CF_F | 0.6199 | 0.6961 | 0.4946 | 0.6240 | 0.6972 | 0.6007 |
| | | Diff | 0.0772 | 0.0654 | 0.0968 | 0.0041 | 0.0236 | 0.0157 |
| | | Or_M | 0.6163 | 0.7509 | 0.4683 | 0.6465 | 0.6899 | 0.6413 |
| | | CF_M | 0.5350 | 0.6679 | 0.3889 | 0.6408 | 0.6607 | 0.6395 |
| | | Diff | 0.0813 | 0.0830 | 0.0794 | 0.0057 | 0.0292 | 0.0018 |
| | Prompt8 | Or_F | 0.5102 | 0.5980 | 0.3656 | 0.6118 | 0.6571 | 0.6162 |
| | | CF_F | 0.5000 | 0.5980 | 0.3387 | 0.6220 | 0.6630 | 0.6297 |
| | | Diff | 0.0102 | 0.0000 | 0.0269 | 0.0102 | 0.0059 | 0.0135 |
| | | Or_M | 0.4650 | 0.6317 | 0.2817 | 0.6730 | 0.6692 | 0.6750 |
| | | CF_M | 0.4934 | 0.6426 | 0.3294 | 0.6560 | 0.6617 | 0.6566 |
| | | Diff | 0.0284 | 0.0109 | 0.0477 | 0.0170 | 0.0075 | 0.0184 |
| | LR | Or_F | 0.8211 | 0.8889 | 0.7097 | 0.6626 | 0.7662 | 0.5896 |
| | | CF_F | 0.6159 | 0.7353 | 0.4194 | 0.6768 | 0.7389 | 0.6580 |
| | | Diff | 0.2052 | 0.1536 | 0.2903 | 0.0142 | 0.0273 | 0.0684 |
| | | Or_M | 0.4858 | 0.6426 | 0.3135 | 0.6635 | 0.6667 | 0.6646 |
| | | CF_M | 0.6994 | 0.8303 | 0.5556 | 0.6465 | 0.7110 | 0.6374 |
| | | Diff | 0.2136 | 0.1877 | 0.2421 | 0.0170 | 0.0443 | 0.0272 |
| | MLP | Or_F | 0.6199 | 0.7092 | 0.4631 | 0.6402 | 0.7103 | 0.6180 |
| | | CF_F | 0.5508 | 0.6634 | 0.3656 | 0.6524 | 0.7036 | 0.6489 |
| | | Diff | 0.0691 | 0.0458 | 0.0975 | 0.0122 | 0.0067 | 0.0309 |
| | | Or_M | 0.4405 | 0.5451 | 0.3254 | 0.6068 | 0.5922 | 0.6099 |
| | | CF_M | 0.5539 | 0.6715 | 0.4246 | 0.6257 | 0.6526 | 0.6234 |
| | | Diff | 0.1134 | 0.1264 | 0.0992 | 0.0189 | 0.0604 | 0.0135 |