

Correlations between bias features and preference labels

