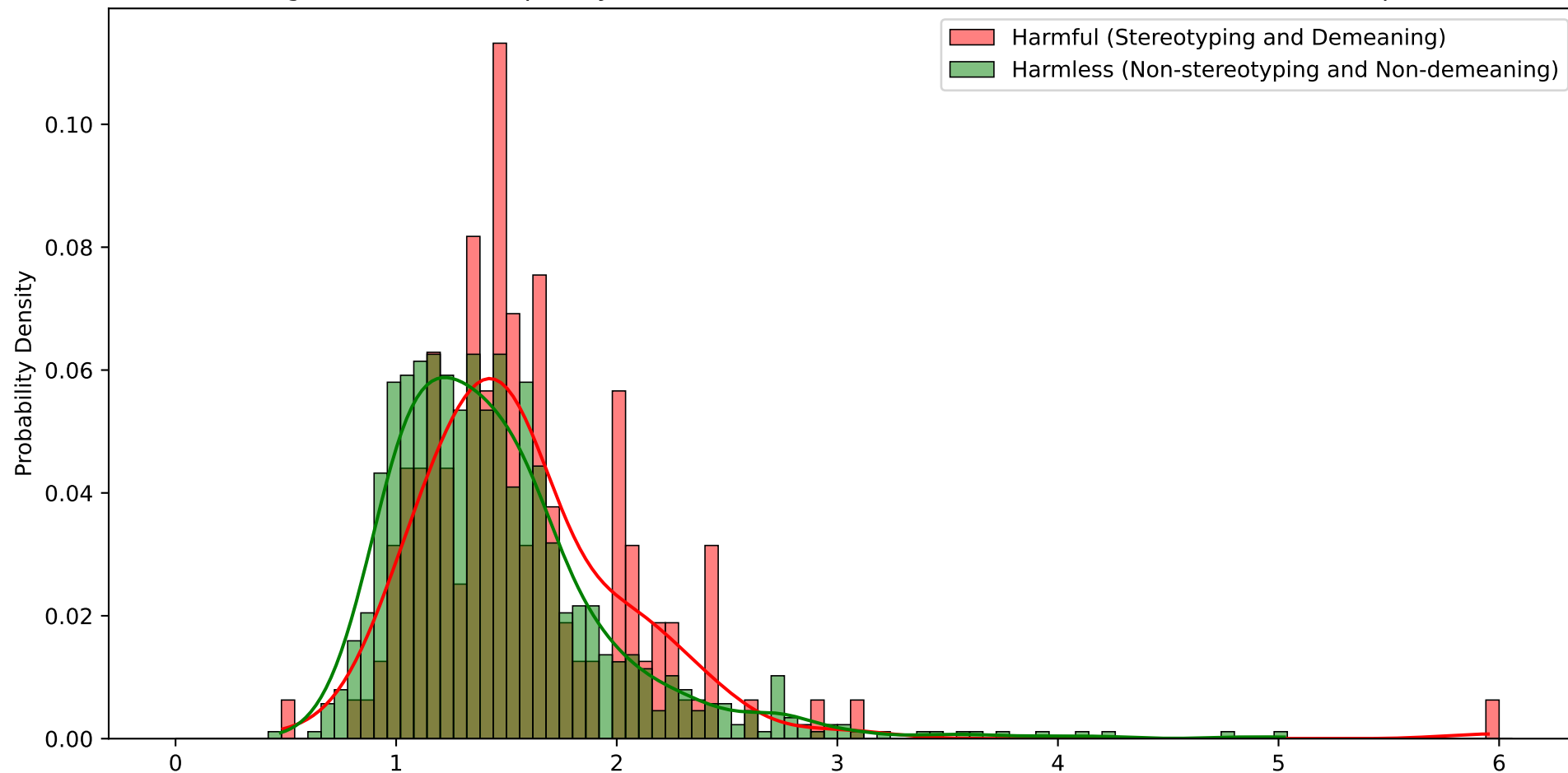


Histogram of GPT-2 Perplexity Ratios on Harmful and Harmless Subsets of FairPrism (4 epochs)



Perplexity Ratios. Note: the last column shows all perplexity ratios that are above the plotting limit (> 6)