

Signed Maximum Logit Change; $S_{adv} - S$

2.0
1.5
1.0
0.5
0.0

Familiar

Novel

