

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

Familiar
Novel

0.2 0.4 0.6 0.8 1.0 1.2 1.4 1.6 1.8

s - Maximum Logit Score (MLS)

