

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

Familiar
Novel

s - Maximum Logit Score (MLS)

