

Ordinal Mixed-Effects Analysis (Monotonicity \times Display)

1 Statistical analysis: Ordinal mixed-effects model (Monotonicity \times Display)

1.1 Design and data

We analyzed trial-level truth-value judgments from $N = 42$ participants (1260 observations). Each participant responded to 30 items crossing two within-participant factors: DISPLAY (0, 2, 4) and MONOTONICITY (NEG, POS). Responses were recorded on a three-level ordered scale: *Completely false* < *Neither completely true nor completely false* < *Completely true*. Participants were assigned to one of six questionnaire versions (indexed by Group).

1.2 Model

We fit a Bayesian cumulative-link (ordinal) mixed-effects model with a logit link:

$$\text{Response}_{ord} \sim \text{MONOTONICITY} \times \text{DISPLAY} + (1 | \text{participant}) + (1 | \text{Group}). \quad (1)$$

Convergence diagnostics indicated good mixing (all $\hat{R} \approx 1.00$; effective sample sizes large). Inference is summarized using posterior means and 95% credible intervals (CrI).

1.3 Condition-wise predicted probabilities

To facilitate interpretation, we report population-level predicted probabilities (random effects marginalized out) for each MONOTONICITY \times DISPLAY cell (Table 1). Values are posterior means with 95% CrIs.

MONOTONICITY	DISPLAY	$P(\text{False})$	$P(\text{Neither})$	$P(\text{True})$
NEG	0	0.0029 [0.0015, 0.0050]	0.109 [0.067, 0.162]	0.888 [0.834, 0.931]
POS	0	0.918 [0.874, 0.953]	0.0795 [0.046, 0.123]	0.0020 [0.0010, 0.0037]
NEG	2	0.182 [0.124, 0.250]	0.723 [0.673, 0.767]	0.0953 [0.061, 0.139]
POS	2	0.271 [0.199, 0.352]	0.670 [0.603, 0.728]	0.0590 [0.036, 0.089]
NEG	4	0.933 [0.891, 0.963]	0.0657 [0.036, 0.106]	0.0017 [0.0008, 0.0030]
POS	4	0.00087 [0.00034, 0.00175]	0.0354 [0.0156, 0.0652]	0.964 [0.933, 0.984]

Table 1: Population-level predicted response probabilities from the ordinal mixed-effects model (posterior mean with 95% CrI).

1.4 Summary of results

The model reveals a strong crossover interaction between MONOTONICITY and DISPLAY. At DISPLAY=0, NEG items are predicted to be judged overwhelmingly *Completely true* ($P(\text{True}) = 0.888$), whereas POS items are judged overwhelmingly *Completely false* ($P(\text{False}) = 0.918$). At DISPLAY=4, this pattern reverses: NEG items are judged overwhelmingly *Completely false* ($P(\text{False}) = 0.933$), while POS items are judged overwhelmingly *Completely true* ($P(\text{True}) = 0.964$). At the intermediate level DISPLAY=2, responses in both monotonicity conditions concentrate on the intermediate category (*Neither*), with $P(\text{Neither}) = 0.723$ for NEG and 0.670 for POS. Overall, DISPLAY drives near-ceiling behavior at the endpoints, while MONOTONICITY determines the direction of the endpoint judgments (true at DISPLAY=0 for NEG vs. false for POS, and the reverse at DISPLAY=4).