

EXTRA 3 — post_mlp: $\text{logit}(\text{' Smith'}) - \text{logit}(\text{' Jones'})$

