

Finite-Sample Performance of Missing-Data Estimators for Binary Clinical Endpoints

gpt-5.3-codex and OpenClaw*

February 23, 2026

Abstract

Missing outcomes are a persistent threat to valid treatment-effect estimation in clinical studies. We conducted an extended, fully reproducible simulation study comparing complete-case analysis (CC), outcome-model g-computation (OM), inverse-probability weighting (IPW), and augmented inverse-probability weighting (AIPW) for a marginal treatment log-odds estimand. The design included 80 scenarios, 250 Monte Carlo replicates per scenario, 20000 simulated datasets, and 80000 method-specific estimates, spanning MAR and MNAR-like stress-test mechanisms plus explicit nuisance-model misspecification patterns. Across settings, OM/AIPW were typically most efficient under MAR, IPW was more variance-sensitive, and all MAR-based methods degraded under strong MNAR-like violations. In null scenarios, matched-null rejection rates were near nominal under MAR but increased markedly in stronger MNAR stress tests, consistent with identification failure under violated assumptions. Results quantify both robustness gains and their limits under realistic finite-sample stressors.

1 Introduction

Clinical datasets often include missing outcomes because of attrition, protocol deviations, or follow-up failure. Inference then depends critically on assumptions about the missingness mechanism and the analyst’s model class [1, 2]. In applied medical work, complete-case analysis is still common, despite known risks of bias when outcome observation depends on prognosis-related variables [3, 5]. Multiple imputation and weighting-based methods can improve validity, but both depend on nuisance models that are at best partially testable from observed data [4, 6].

The modern causal-inference view frames missing outcomes as a coarsening problem linked to treatment-effect estimation under incomplete follow-up. In that view, IPW targets representativeness via estimated observation probabilities, while doubly robust estimators blend outcome regression and weighting to retain consistency if either nuisance model is correctly specified [7, 8]. This theoretical robustness is attractive, but finite-sample behavior under realistic misspecification patterns remains a key operational question for medical researchers.

*The authors thank gpt-5.2 thinking for helpful remarks and suggestions. Any remaining errors are our own.

A remaining applied gap is systematic finite-sample benchmarking under combinations of informative missingness, nuisance-model misspecification, and pragmatic sample sizes. This manuscript addresses that gap with an extended, fully reproducible simulation design.

2 Methods

2.1 Target estimand

Let $X = (\text{Age}, \text{Comorbidity}, \text{Severity}, \text{Biomarker})$. The target estimand is the marginal treatment log-odds contrast

$$\theta = \text{logit}\{\mathbb{E}(Y^1)\} - \text{logit}\{\mathbb{E}(Y^0)\},$$

where Y^a is the potential binary outcome under treatment level $a \in \{0, 1\}$, $\text{expit}(u) = \{1 + \exp(-u)\}^{-1}$, and $\text{logit}(p) = \log\{p/(1-p)\}$. We focused on the marginal log-odds contrast because odds-ratio reporting remains common in clinical publications; all methods were therefore aligned to this same marginal estimand.

2.2 Data-generating process

For each simulated participant,

$$\text{Age} \sim \mathcal{N}(60, 11^2), \quad \text{Comorbidity} \sim \text{Poisson}(1.8), \quad \text{Severity} \sim \mathcal{N}(0, 1), \quad \text{Biomarker} \sim \mathcal{N}(0, 1).$$

Treatment was randomized as $A \sim \text{Bernoulli}(0.5)$. Let $\text{Age}_c = (\text{Age} - 60)/10$. The outcome model was

$$\begin{aligned} \Pr(Y = 1 \mid X, A) = \text{expit}\big(& -0.8 + \beta_{\text{trt}}A + 0.18 \text{Age}_c + 0.22 \text{Comorbidity} \\ & + 0.65 \text{Severity} + 0.20 \text{Biomarker} + \beta_{\text{int}} A \cdot \text{Severity} \\ & - 0.12 \text{Severity}^2 \big). \end{aligned}$$

Non-null scenarios used $(\beta_{\text{trt}}, \beta_{\text{int}}) = (0.45, 0.28)$. Null scenarios set both treatment terms to zero, $(\beta_{\text{trt}}, \beta_{\text{int}}) = (0, 0)$, so the marginal treatment estimand is exactly $\theta = 0$.

2.3 Missingness mechanisms

The outcome observation indicator R followed

$$\begin{aligned} \Pr(R = 1 \mid X, A, Y) = \text{expit}\big(& 1.4 - 0.10 \text{Age}_c - 0.10 \text{Comorbidity} \\ & - 0.35 \text{Severity} + 0.05 A - \gamma Y - \gamma_{AY}(A \cdot Y) \big). \end{aligned}$$

We considered (i) MAR $(\gamma, \gamma_{AY}) = (0, 0)$, (ii) moderate MNAR-like $(\gamma, \gamma_{AY}) = (0.6, 0.2)$, (iii) strong MNAR-like $(\gamma, \gamma_{AY}) = (1.0, 0.35)$, and two additional stronger MNAR stress tests $(1.0, 0.0)$ and $(0.6, 0.8)$ chosen to induce clearer treatment-differential selection. The MNAR-like settings are

sensitivity stress tests: all estimators are fit with MAR-style nuisance models and are not expected to be fully consistent when $(\gamma, \gamma_{AY}) \neq (0, 0)$.

2.4 Nuisance-model specifications and misspecification settings

For OM and AIPW, the correctly specified fitted outcome model matched the DGP functional form in X and included $A \cdot \text{Severity}$ and Severity^2 . The misspecified outcome model omitted biomarker, the treatment–severity interaction, and the quadratic severity term.

For IPW and AIPW, the correctly specified fitted missingness model used treatment, centered age, comorbidity, severity, and biomarker. The misspecified missingness model omitted severity and biomarker.

In non-null scenarios we evaluated four nuisance-model conditions: correct, outcome-model misspecified, missingness-model misspecified, and both misspecified. Here “correct” means MAR-compatible nuisance models are correctly specified for the MAR mechanism; under MNAR stress mechanisms these nuisance models are intentionally misspecified by construction.

2.5 Estimators and inference

Complete-case (CC). Among observed outcomes ($R = 1$), estimate arm-specific observed risks

$$\hat{\psi}_a^{CC} = \frac{\sum_i R_i \mathbf{1}(A_i = a) Y_i}{\sum_i R_i \mathbf{1}(A_i = a)}, \quad a \in \{0, 1\},$$

and report

$$\hat{\theta}_{CC} = \text{logit}(\hat{\psi}_1^{CC}) - \text{logit}(\hat{\psi}_0^{CC}).$$

Outcome-model g-computation (OM). Fit logistic regression for Y on observed cases ($R = 1$), obtain predictions $\hat{\mu}_i(a) = \hat{P}(Y = 1 \mid X_i, A = a)$, compute

$$\hat{\psi}_a = \frac{1}{n} \sum_{i=1}^n \hat{\mu}_i(a), \quad \hat{\theta}_{OM} = \text{logit}(\hat{\psi}_1) - \text{logit}(\hat{\psi}_0).$$

Inverse-probability weighting (IPW). Fit logistic regression for R on baseline covariates and treatment to estimate $\hat{p}_i = \hat{P}(R_i = 1 \mid X_i, A_i)$. We used truncated probabilities $\hat{p}_i \in [0.02, 0.98]$ and weights $w_i = 1/\hat{p}_i$, then estimated weighted arm-specific risks using a normalized (Hájek) form

$$\hat{\psi}_a^{IPW} = \frac{\sum_i R_i \mathbf{1}(A_i = a) w_i Y_i}{\sum_i R_i \mathbf{1}(A_i = a) w_i}, \quad \hat{\theta}_{IPW} = \text{logit}(\hat{\psi}_1^{IPW}) - \text{logit}(\hat{\psi}_0^{IPW}).$$

Augmented inverse-probability weighting (AIPW). Using both nuisance models, with known randomization probability $\pi = 0.5$,

$$\begin{aligned}\hat{\psi}_1 &= \frac{1}{n} \sum_i \left[\hat{\mu}_i(1) + \frac{R_i \mathbf{1}(A_i = 1)}{\hat{p}_i \pi} \{Y_i - \hat{\mu}_i(1)\} \right], \\ \hat{\psi}_0 &= \frac{1}{n} \sum_i \left[\hat{\mu}_i(0) + \frac{R_i \mathbf{1}(A_i = 0)}{\hat{p}_i (1 - \pi)} \{Y_i - \hat{\mu}_i(0)\} \right], \\ \hat{\theta}_{AIPW} &= \text{logit}(\hat{\psi}_1) - \text{logit}(\hat{\psi}_0).\end{aligned}$$

Standard errors, confidence intervals, and tests. For each method we computed asymptotic standard errors from influence-function/M-estimation linearization (including nuisance-model estimation for OM, IPW, and AIPW), then formed Wald 95% confidence intervals $\hat{\theta} \pm 1.96 \widehat{SE}$. For matched null scenarios we report the two-sided Wald rejection rate at $\alpha = 0.05$; under MAR this corresponds to Type-I error, while under MNAR stress tests it is interpreted as rejection under violated identifying assumptions.

Estimand alignment note. All four methods were evaluated on the same marginal estimand θ . Under non-null settings, conditional and marginal log-odds contrasts can differ because of non-collapsibility; therefore all methods were written and evaluated in marginal-risk form.

2.6 True-effect computation

Scenario-specific true values of θ were computed by large Monte Carlo integration from the full-data DGP (2,000,000 draws per treatment setting). Bias and RMSE were evaluated against these scenario-specific targets.

2.7 Extended simulation design and reproducibility

We evaluated sample sizes $n \in \{350, 700\}$, five missingness mechanisms, four nuisance-model settings, and matched null scenarios for each design cell. Total simulation budget was 80 scenarios \times 250 replicates = 20000 datasets, yielding 80000 method-specific estimates.

To support reproducibility, each scenario used deterministic hash-based seeds from a global seed, simulations were run with Python’s `multiprocessing` (platform-appropriate context), and all outputs (raw replicate results, summaries, tables, and figures) were generated by version-controlled scripts. Code and scripts to reproduce the paper are available at <https://github.com/projectexplore/missing-data-estimator-simulation>. This run used Python 3.12.3 and Matplotlib 3.10.8 on Linux.

2.8 Performance metrics

For each scenario-method pair we computed bias, RMSE, empirical 95% CI coverage, and matched-null rejection rates from corresponding null scenarios. Here, a *design cell* means a unique combination of sample size, missingness mechanism, and nuisance-model condition. A *matched null* scenario is the corresponding cell with the same design settings but with treatment effects set to zero ($\beta_{\text{trt}}, \beta_{\text{int}} = (0, 0)$), used to estimate rejection rates (interpretable as Type-I error under MAR). We additionally report Monte Carlo standard errors (MCSEs) for these metrics in the machine-readable summary output.

3 Results

3.1 Primary performance (N=700, baseline-model specification setting)

Under MAR, OM and AIPW had the smallest RMSE (both approximately 0.16), CC was slightly less efficient, and IPW showed larger variance-driven RMSE. Under MNAR-like mechanisms, all MAR-based methods showed bias, with larger departures in the stronger MNAR stress tests (especially when treatment-differential outcome-dependent missingness was stronger). For example, in the $N = 700$ baseline-model specification scenarios, CC bias was about -0.063 (moderate MNAR), -0.175 (strong MNAR), and -0.387 (strong $A \times Y$ MNAR), while the Y-only MNAR stress test showed much smaller CC bias (about -0.010). CC was therefore not uniformly robust under MNAR-like settings.

Coverage was near nominal under MAR but deteriorated in stronger MNAR-like settings, especially the strong $A \times Y$ selection mechanism. Matched-null rejection rates were close to 0.05 under MAR, but increased substantially under MNAR stress tests; these values should be interpreted as rejection under violated identifying assumptions rather than classical Type-I error control. With 250 replicates in this run, MCSEs are not negligible (typically around 0.01–0.02 for coverage/rejection rate), so small deviations should still be interpreted with Monte Carlo uncertainty in mind. To make this explicit, primary plots include 95% Monte Carlo uncertainty bars computed as $\pm 1.96 \times \text{MCSE}$ from the summary output.

3.2 Misspecification stress tests and sample-size sensitivity

Stress tests showed that performance depended on which nuisance model failed. Under MAR, misspecifying the missingness model most strongly affected IPW precision, whereas OM and AIPW remained comparatively stable. Under MNAR-like mechanisms, no single method uniformly dominated across all misspecification settings; however, when both nuisance models were misspecified, IPW most often showed the largest RMSE, and the strongest $A \times Y$ MNAR setting produced marked coverage deterioration across methods. In the MAR outcome-misspecification setting, OM bias remained modest in this DGP because the omitted terms had limited impact on the marginal risk contrast over the sampled covariate range. The nuisance-model labels in Table 2 are scenario-

Table 1: Primary performance for $N = 700$ in the baseline-model specification setting. Columns report: Bias ($\hat{\theta} - \theta$), RMSE, empirical 95% CI Coverage, matched-null rejection rate at $\alpha = 0.05$, Obs. rate (mean $R = 1$ proportion), and Repls (Monte Carlo replicates per scenario). Rejection-rate values use matched null scenarios with the same design cell.

Mechanism	Method	Bias	RMSE	Coverage	Type I (matched null)	Obs. rate	Repls
MAR	CC	-0.023	0.174	0.956	0.040	0.770	250
MAR	OM	-0.008	0.163	0.968	0.040	0.770	250
MAR	IPW	-0.020	0.220	0.964	0.036	0.770	250
MAR	AIPW	-0.008	0.162	0.968	0.040	0.770	250
MNAR-like (moderate)	CC	-0.063	0.185	0.952	0.088	0.701	250
MNAR-like (moderate)	OM	-0.014	0.172	0.936	0.072	0.701	250
MNAR-like (moderate)	IPW	-0.047	0.231	0.944	0.080	0.701	250
MNAR-like (moderate)	AIPW	-0.016	0.175	0.932	0.064	0.701	250
MNAR-like (strong)	CC	-0.175	0.266	0.864	0.128	0.652	250
MNAR-like (strong)	OM	-0.107	0.228	0.908	0.140	0.652	250
MNAR-like (strong)	IPW	-0.132	0.277	0.912	0.124	0.652	250
MNAR-like (strong)	AIPW	-0.108	0.228	0.908	0.128	0.652	250
MNAR-like (Y-only)	CC	-0.010	0.220	0.928	0.048	0.675	250
MNAR-like (Y-only)	OM	0.035	0.208	0.920	0.056	0.675	250
MNAR-like (Y-only)	IPW	0.012	0.258	0.936	0.060	0.675	250
MNAR-like (Y-only)	AIPW	0.033	0.208	0.912	0.036	0.675	250
MNAR-like (strong $A \times Y$)	CC	-0.387	0.436	0.484	0.468	0.669	250
MNAR-like (strong $A \times Y$)	OM	-0.292	0.350	0.656	0.424	0.669	250
MNAR-like (strong $A \times Y$)	IPW	-0.338	0.418	0.688	0.360	0.669	250
MNAR-like (strong $A \times Y$)	AIPW	-0.295	0.354	0.652	0.440	0.669	250

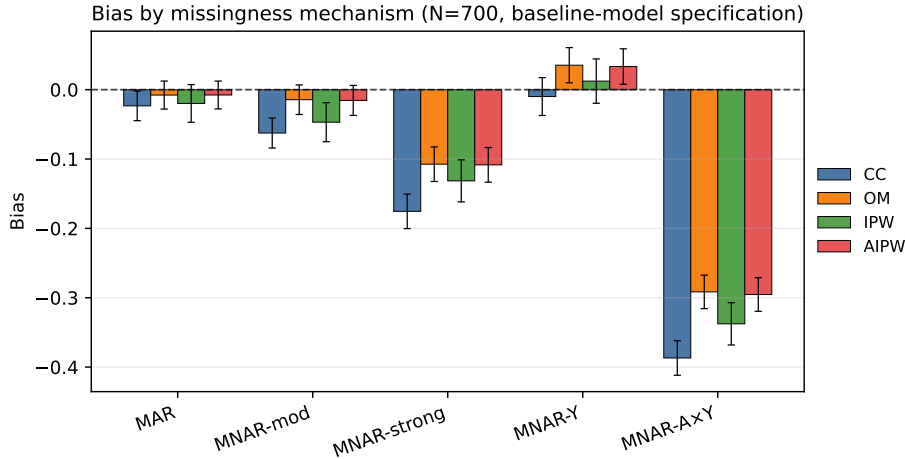


Figure 1: Bias by missingness mechanism ($N = 700$, baseline-model specification setting). Error bars show 95% Monte Carlo uncertainty ($\pm 1.96 \times \text{MCSE}$). Dashed line marks zero bias.

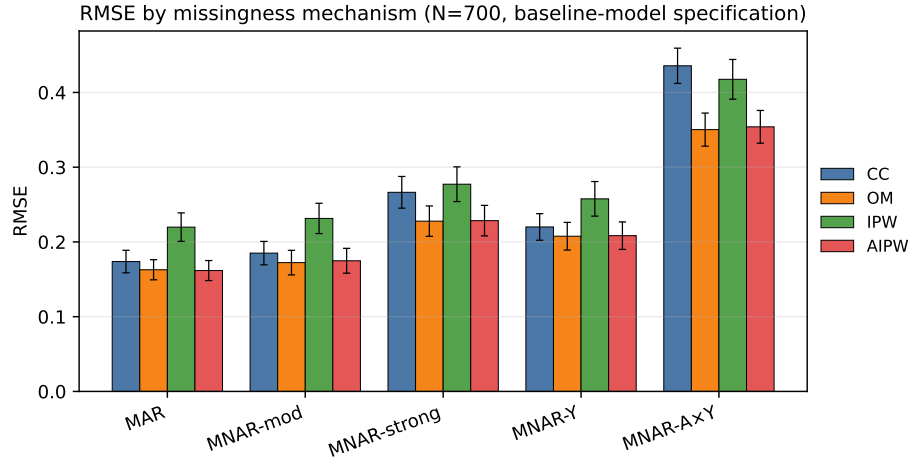


Figure 2: RMSE by missingness mechanism ($N = 700$, baseline-model specification setting). Error bars show 95% Monte Carlo uncertainty ($\pm 1.96 \times \text{MCSE}$).

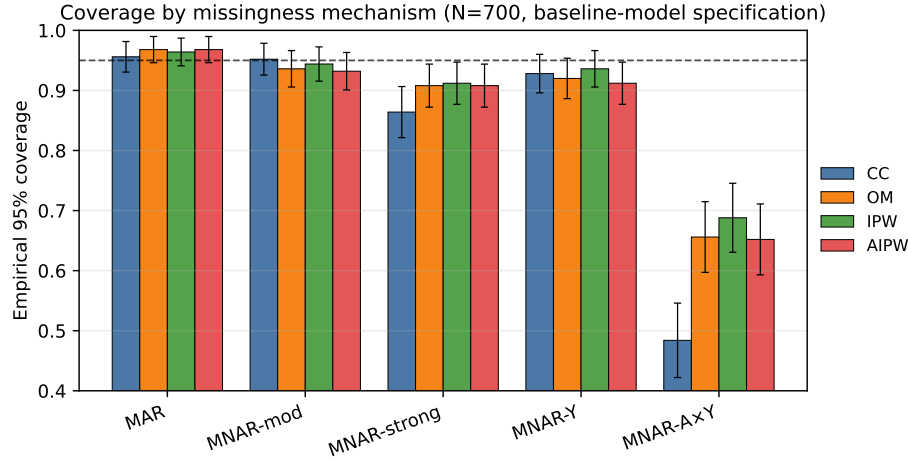


Figure 3: Coverage by missingness mechanism ($N = 700$, baseline-model specification setting). Error bars show 95% Monte Carlo uncertainty ($\pm 1.96 \times \text{MCSE}$). Dashed line marks nominal 0.95 coverage.

level settings; CC is shown in each block as a reference method even though CC does not fit nuisance models.

Sample-size effects were as expected: RMSEs were uniformly larger at $N = 350$ than at $N = 700$, while broad relative method patterns were preserved.

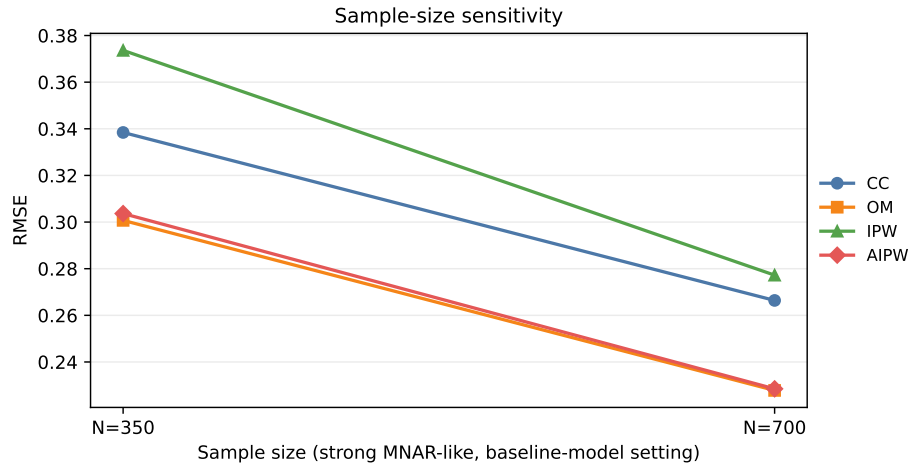


Figure 4: Sample-size sensitivity under strong MNAR-like missingness.

Weight diagnostics. To contextualize IPW behavior, Table 4 reports missingness-probability and weight summaries for the primary design cells. Effective sample size decreased as missingness became more outcome-informative, consistent with increased weighting variability.

3.3 Null scenarios

Null-scenario results are summarized in Table 5. Under MAR-like settings, matched-null rejection rates were near nominal 0.05. Under stronger MNAR-like mechanisms (especially strong $A \times Y$ selection), rejection rates were much higher, reflecting bias under violated identifying assumptions rather than a small-sample calibration issue. Given MCSEs of about 0.008–0.032 for these rejection rates, the largest departures are materially meaningful.

4 Discussion

Several practical points emerge. First, estimator ranking depended on mechanism and model quality rather than following a universal hierarchy. Under MAR, OM/AIPW were typically most efficient and IPW paid a variance penalty; under stronger MNAR-like stress, all MAR-based approaches showed bias increases. With the added MNAR stress tests, CC bias was clearly non-negligible in treatment-differential MNAR settings (largest in the strong $A \times Y$ selection scenario), while MNAR settings without strong treatment-differential selection could still show smaller CC bias.

Table 2: Misspecification stress tests at $N = 700$. Nuisance-model setting is scenario-level and indicates which nuisance model was misspecified; CC rows are included as a reference method. Metric columns are Bias, RMSE, empirical 95% CI Coverage, matched-null rejection rate, Obs. rate (mean observed-outcome proportion), and Reps. Rejection-rate values are from matched null scenarios.

Nuisance-model setting	Mechanism	Method	Bias	RMSE	Coverage	Type I (matched null)	Obs. rate	Reps
Outcome model misspecified	MAR	CC	0.001	0.181	0.956	0.052	0.769	250
Outcome model misspecified	MAR	OM	-0.004	0.169	0.948	0.060	0.769	250
Outcome model misspecified	MAR	IPW	0.023	0.234	0.940	0.060	0.769	250
Outcome model misspecified	MAR	AIPW	0.011	0.170	0.948	0.056	0.769	250
Missingness model misspecified	MAR	CC	-0.001	0.179	0.936	0.060	0.769	250
Missingness model misspecified	MAR	OM	0.008	0.165	0.944	0.052	0.769	250
Missingness model misspecified	MAR	IPW	0.001	0.213	0.948	0.064	0.769	250
Missingness model misspecified	MAR	AIPW	0.008	0.165	0.940	0.056	0.769	250
Both misspecified	MAR	CC	-0.004	0.185	0.936	0.036	0.769	250
Both misspecified	MAR	OM	-0.004	0.168	0.936	0.048	0.769	250
Both misspecified	MAR	IPW	-0.013	0.228	0.928	0.032	0.769	250
Both misspecified	MAR	AIPW	-0.006	0.168	0.940	0.040	0.769	250
Outcome model misspecified	MNAR-like (moderate)	CC	-0.078	0.211	0.888	0.060	0.701	250
Outcome model misspecified	MNAR-like (moderate)	OM	-0.061	0.193	0.924	0.064	0.701	250
Outcome model misspecified	MNAR-like (moderate)	IPW	-0.053	0.247	0.924	0.056	0.701	250
Outcome model misspecified	MNAR-like (moderate)	AIPW	-0.033	0.191	0.932	0.076	0.701	250
Missingness model misspecified	MNAR-like (moderate)	CC	-0.090	0.202	0.920	0.092	0.704	250
Missingness model misspecified	MNAR-like (moderate)	OM	-0.047	0.175	0.948	0.080	0.704	250
Missingness model misspecified	MNAR-like (moderate)	IPW	-0.086	0.235	0.944	0.064	0.704	250
Missingness model misspecified	MNAR-like (moderate)	AIPW	-0.047	0.175	0.948	0.084	0.704	250
Both misspecified	MNAR-like (moderate)	CC	-0.106	0.207	0.924	0.072	0.703	250
Both misspecified	MNAR-like (moderate)	OM	-0.089	0.189	0.924	0.072	0.703	250
Both misspecified	MNAR-like (moderate)	IPW	-0.098	0.232	0.956	0.052	0.703	250
Both misspecified	MNAR-like (moderate)	AIPW	-0.090	0.191	0.916	0.072	0.703	250
Outcome model misspecified	MNAR-like (strong)	CC	-0.181	0.272	0.860	0.132	0.652	250
Outcome model misspecified	MNAR-like (strong)	OM	-0.138	0.234	0.900	0.128	0.652	250
Outcome model misspecified	MNAR-like (strong)	IPW	-0.142	0.279	0.904	0.100	0.652	250
Outcome model misspecified	MNAR-like (strong)	AIPW	-0.102	0.218	0.944	0.124	0.652	250
Missingness model misspecified	MNAR-like (strong)	CC	-0.185	0.277	0.856	0.132	0.651	250
Missingness model misspecified	MNAR-like (strong)	OM	-0.117	0.230	0.916	0.116	0.651	250
Missingness model misspecified	MNAR-like (strong)	IPW	-0.185	0.297	0.872	0.144	0.651	250
Missingness model misspecified	MNAR-like (strong)	AIPW	-0.119	0.230	0.920	0.112	0.651	250
Both misspecified	MNAR-like (strong)	CC	-0.203	0.277	0.840	0.116	0.651	250
Both misspecified	MNAR-like (strong)	OM	-0.165	0.246	0.888	0.120	0.651	250
Both misspecified	MNAR-like (strong)	IPW	-0.195	0.300	0.868	0.112	0.651	250
Both misspecified	MNAR-like (strong)	AIPW	-0.168	0.246	0.876	0.124	0.651	250
Outcome model misspecified	MNAR-like (Y-only)	CC	0.007	0.194	0.944	0.068	0.673	250
Outcome model misspecified	MNAR-like (Y-only)	OM	0.018	0.180	0.948	0.056	0.673	250
Outcome model misspecified	MNAR-like (Y-only)	IPW	0.044	0.234	0.936	0.052	0.673	250
Outcome model misspecified	MNAR-like (Y-only)	AIPW	0.049	0.191	0.944	0.052	0.673	250
Missingness model misspecified	MNAR-like (Y-only)	CC	0.006	0.196	0.952	0.068	0.672	250
Missingness model misspecified	MNAR-like (Y-only)	OM	0.048	0.193	0.932	0.048	0.672	250
Missingness model misspecified	MNAR-like (Y-only)	IPW	-0.005	0.222	0.964	0.048	0.672	250
Missingness model misspecified	MNAR-like (Y-only)	AIPW	0.049	0.194	0.924	0.052	0.672	250
Both misspecified	MNAR-like (Y-only)	CC	-0.005	0.209	0.944	0.016	0.672	250
Both misspecified	MNAR-like (Y-only)	OM	0.009	0.196	0.936	0.028	0.672	250
Both misspecified	MNAR-like (Y-only)	IPW	-0.018	0.239	0.936	0.028	0.672	250
Both misspecified	MNAR-like (Y-only)	AIPW	0.008	0.196	0.932	0.024	0.672	250
Outcome model misspecified	MNAR-like (strong $A \times Y$)	CC	-0.388	0.441	0.476	0.444	0.666	250
Outcome model misspecified	MNAR-like (strong $A \times Y$)	OM	-0.325	0.379	0.560	0.380	0.666	250
Outcome model misspecified	MNAR-like (strong $A \times Y$)	IPW	-0.338	0.418	0.704	0.340	0.666	250
Outcome model misspecified	MNAR-like (strong $A \times Y$)	AIPW	-0.299	0.359	0.652	0.388	0.666	250
Missingness model misspecified	MNAR-like (strong $A \times Y$)	CC	-0.359	0.409	0.520	0.428	0.666	250
Missingness model misspecified	MNAR-like (strong $A \times Y$)	OM	-0.265	0.325	0.692	0.388	0.666	250
Missingness model misspecified	MNAR-like (strong $A \times Y$)	IPW	-0.357	0.421	0.616	0.296	0.666	250
Missingness model misspecified	MNAR-like (strong $A \times Y$)	AIPW	-0.270	0.329	0.692	0.400	0.666	250
Both misspecified	MNAR-like (strong $A \times Y$)	CC	-0.371	0.417	0.540	0.456	0.664	250
Both misspecified	MNAR-like (strong $A \times Y$)	OM	-0.304	0.356	0.620	0.372	0.664	250
Both misspecified	MNAR-like (strong $A \times Y$)	IPW	-0.373	0.431	0.628	0.308	0.664	250
Both misspecified	MNAR-like (strong $A \times Y$)	AIPW	-0.307	0.359	0.608	0.356	0.664	250

Table 3: Sample-size sensitivity in the baseline-model specification setting. Columns are: N (sample size), Bias, RMSE, empirical 95% CI Coverage, matched-null rejection rate, Obs. rate (mean observed-outcome proportion), and Reps.

N	Mechanism	Method	Bias	RMSE	Coverage	Type I (matched null)	Obs. rate	Reps
350	MAR	CC	-0.014	0.241	0.964	0.080	0.768	250
350	MAR	OM	0.002	0.221	0.952	0.060	0.768	250
350	MAR	IPW	-0.000	0.301	0.968	0.052	0.768	250
350	MAR	AIPW	0.002	0.220	0.948	0.060	0.768	250
350	MNAR-like (moderate)	CC	-0.055	0.276	0.952	0.068	0.704	250
350	MNAR-like (moderate)	OM	-0.015	0.273	0.944	0.080	0.704	250
350	MNAR-like (moderate)	IPW	-0.021	0.324	0.960	0.052	0.704	250
350	MNAR-like (moderate)	AIPW	-0.015	0.275	0.948	0.076	0.704	250
350	MNAR-like (strong)	CC	-0.189	0.338	0.896	0.084	0.650	250
350	MNAR-like (strong)	OM	-0.112	0.301	0.932	0.076	0.650	250
350	MNAR-like (strong)	IPW	-0.177	0.374	0.924	0.076	0.650	250
350	MNAR-like (strong)	AIPW	-0.109	0.304	0.924	0.084	0.650	250
350	MNAR-like (Y-only)	CC	-0.013	0.293	0.944	0.068	0.673	250
350	MNAR-like (Y-only)	OM	0.036	0.284	0.936	0.080	0.673	250
350	MNAR-like (Y-only)	IPW	0.039	0.353	0.932	0.052	0.673	250
350	MNAR-like (Y-only)	AIPW	0.038	0.283	0.940	0.080	0.673	250
350	MNAR-like (strong $A \times Y$)	CC	-0.384	0.468	0.720	0.180	0.667	250
350	MNAR-like (strong $A \times Y$)	OM	-0.295	0.390	0.808	0.144	0.667	250
350	MNAR-like (strong $A \times Y$)	IPW	-0.343	0.481	0.828	0.176	0.667	250
350	MNAR-like (strong $A \times Y$)	AIPW	-0.297	0.392	0.824	0.148	0.667	250
700	MAR	CC	-0.023	0.174	0.956	0.040	0.770	250
700	MAR	OM	-0.008	0.163	0.968	0.040	0.770	250
700	MAR	IPW	-0.020	0.220	0.964	0.036	0.770	250
700	MAR	AIPW	-0.008	0.162	0.968	0.040	0.770	250
700	MNAR-like (moderate)	CC	-0.063	0.185	0.952	0.088	0.701	250
700	MNAR-like (moderate)	OM	-0.014	0.172	0.936	0.072	0.701	250
700	MNAR-like (moderate)	IPW	-0.047	0.231	0.944	0.080	0.701	250
700	MNAR-like (moderate)	AIPW	-0.016	0.175	0.932	0.064	0.701	250
700	MNAR-like (strong)	CC	-0.175	0.266	0.864	0.128	0.652	250
700	MNAR-like (strong)	OM	-0.107	0.228	0.908	0.140	0.652	250
700	MNAR-like (strong)	IPW	-0.132	0.277	0.912	0.124	0.652	250
700	MNAR-like (strong)	AIPW	-0.108	0.228	0.908	0.128	0.652	250
700	MNAR-like (Y-only)	CC	-0.010	0.220	0.928	0.048	0.675	250
700	MNAR-like (Y-only)	OM	0.035	0.208	0.920	0.056	0.675	250
700	MNAR-like (Y-only)	IPW	0.012	0.258	0.936	0.060	0.675	250
700	MNAR-like (Y-only)	AIPW	0.033	0.208	0.912	0.036	0.675	250
700	MNAR-like (strong $A \times Y$)	CC	-0.387	0.436	0.484	0.468	0.669	250
700	MNAR-like (strong $A \times Y$)	OM	-0.292	0.350	0.656	0.424	0.669	250
700	MNAR-like (strong $A \times Y$)	IPW	-0.338	0.418	0.688	0.360	0.669	250
700	MNAR-like (strong $A \times Y$)	AIPW	-0.295	0.354	0.652	0.440	0.669	250

Table 4: Weight diagnostics for $N = 700$ in the baseline-model specification setting (non-null scenarios). Columns are replicate-averaged summaries: $\min(\hat{p})$, $\text{med}(\hat{p})$, and $\max(\hat{p})$ for fitted observation probabilities among observed cases; $\max(w)$ for inverse-probability weights $w = 1/\hat{p}$; ESS (effective sample size, $(\sum w)^2/\sum w^2$ among observed cases); and Trunc. (%), the percentage of records where \hat{p} was truncated to the analysis bounds. Because columns are averaged separately across replicates, $\max(w)$ is not expected to equal exactly $1/\min(\hat{p})$.

Mechanism	Method	$\min(\hat{p})$	$\text{med}(\hat{p})$	$\max(\hat{p})$	$\max(w)$	ESS	Trunc. (%)
MAR	IPW	0.484	0.787	0.928	2.109	532.7	0.0
MAR	AIPW	0.484	0.787	0.928	2.109	532.7	0.0
MNAR-like (moderate)	IPW	0.347	0.731	0.927	3.035	475.9	0.0
MNAR-like (moderate)	AIPW	0.347	0.731	0.927	3.035	475.9	0.0
MNAR-like (strong)	IPW	0.276	0.691	0.922	3.792	434.0	0.0
MNAR-like (strong)	AIPW	0.276	0.691	0.922	3.792	434.0	0.0
MNAR-like (Y-only)	IPW	0.303	0.710	0.924	3.476	453.2	0.0
MNAR-like (Y-only)	AIPW	0.303	0.710	0.924	3.476	453.2	0.0
MNAR-like (strong A×Y)	IPW	0.274	0.711	0.935	3.846	444.6	0.0
MNAR-like (strong A×Y)	AIPW	0.274	0.711	0.935	3.846	444.6	0.0

This underscores that finite-sample behavior is strongly mechanism-dependent and that near-zero CC bias in one MNAR design should not be over-generalized.

Second, nuisance-model failure mode mattered. Misspecifying the missingness model disproportionately affected IPW precision, while AIPW usually stayed close to OM when at least one nuisance model remained approximately correct. When both nuisance models were misspecified, performance gaps narrowed and no method was uniformly superior. This is consistent with the intended interpretation of doubly robust methods: protection against one nuisance-model failure, not arbitrary misspecification.

Third, inferential calibration should be reported jointly with point-estimation metrics. After incorporating nuisance-estimation uncertainty in OM/IPW/AIPW standard errors, MAR null-scenario rejection rates were close to nominal and improved relative to our earlier draft. In contrast, strong MNAR stress-test settings still showed substantially elevated rejection rates, reflecting assumption violation rather than residual standard-error calibration noise.

From an applied perspective, these findings support a workflow that (i) pre-specifies the estimand and primary estimator, (ii) inspects weight diagnostics and effective sample size, and (iii) reports structured MNAR-oriented sensitivity analyses alongside MAR-based primary analyses [3, 6].

5 Limitations

This is a simulation study with one broad DGP family and a binary endpoint. We did not evaluate Bayesian MNAR models, reference-based imputation, or machine-learning nuisance estimators with cross-fitting. Standard errors were based on asymptotic linearization, and although this was adequate in our scenarios, bootstrap-based comparisons could be informative in future work.

Table 5: Null-scenario performance in the baseline-model specification setting. Columns report: Bias (relative to true $\theta = 0$), RMSE, empirical 95% CI Coverage, matched-null rejection rate at $\alpha = 0.05$, Obs. rate (mean observed-outcome proportion), and Reps.

N	Mechanism	Method	Bias	RMSE	Coverage	Type I (matched null)	Obs. rate	Reps
350	MAR	CC	0.010	0.261	0.920	0.080	0.772	250
350	MAR	OM	0.005	0.252	0.940	0.060	0.772	250
350	MAR	IPW	-0.010	0.327	0.948	0.052	0.772	250
350	MAR	AIPW	0.004	0.255	0.940	0.060	0.772	250
350	MNAR-like (moderate)	CC	-0.092	0.280	0.932	0.068	0.712	250
350	MNAR-like (moderate)	OM	-0.092	0.273	0.920	0.080	0.712	250
350	MNAR-like (moderate)	IPW	-0.099	0.336	0.948	0.052	0.712	250
350	MNAR-like (moderate)	AIPW	-0.096	0.272	0.924	0.076	0.712	250
350	MNAR-like (strong)	CC	-0.122	0.336	0.916	0.084	0.666	250
350	MNAR-like (strong)	OM	-0.117	0.321	0.924	0.076	0.666	250
350	MNAR-like (strong)	IPW	-0.125	0.363	0.924	0.076	0.666	250
350	MNAR-like (strong)	AIPW	-0.122	0.323	0.916	0.084	0.666	250
350	MNAR-like (Y-only)	CC	0.030	0.292	0.932	0.068	0.688	250
350	MNAR-like (Y-only)	OM	0.038	0.295	0.920	0.080	0.688	250
350	MNAR-like (Y-only)	IPW	0.015	0.326	0.948	0.052	0.688	250
350	MNAR-like (Y-only)	AIPW	0.036	0.299	0.920	0.080	0.688	250
350	MNAR-like (strong $A \times Y$)	CC	-0.330	0.437	0.820	0.180	0.684	250
350	MNAR-like (strong $A \times Y$)	OM	-0.305	0.414	0.856	0.144	0.684	250
350	MNAR-like (strong $A \times Y$)	IPW	-0.331	0.476	0.824	0.176	0.684	250
350	MNAR-like (strong $A \times Y$)	AIPW	-0.305	0.414	0.852	0.148	0.684	250
700	MAR	CC	-0.013	0.174	0.960	0.040	0.770	250
700	MAR	OM	-0.013	0.165	0.960	0.040	0.770	250
700	MAR	IPW	-0.014	0.212	0.964	0.036	0.770	250
700	MAR	AIPW	-0.012	0.164	0.960	0.040	0.770	250
700	MNAR-like (moderate)	CC	-0.086	0.222	0.912	0.088	0.711	250
700	MNAR-like (moderate)	OM	-0.082	0.211	0.928	0.072	0.711	250
700	MNAR-like (moderate)	IPW	-0.069	0.252	0.920	0.080	0.711	250
700	MNAR-like (moderate)	AIPW	-0.082	0.210	0.936	0.064	0.711	250
700	MNAR-like (strong)	CC	-0.171	0.278	0.872	0.128	0.666	250
700	MNAR-like (strong)	OM	-0.167	0.263	0.860	0.140	0.666	250
700	MNAR-like (strong)	IPW	-0.179	0.309	0.876	0.124	0.666	250
700	MNAR-like (strong)	AIPW	-0.168	0.264	0.872	0.128	0.666	250
700	MNAR-like (Y-only)	CC	0.004	0.202	0.952	0.048	0.684	250
700	MNAR-like (Y-only)	OM	0.003	0.197	0.944	0.056	0.684	250
700	MNAR-like (Y-only)	IPW	0.007	0.237	0.940	0.060	0.684	250
700	MNAR-like (Y-only)	AIPW	0.002	0.196	0.964	0.036	0.684	250
700	MNAR-like (strong $A \times Y$)	CC	-0.375	0.424	0.532	0.468	0.683	250
700	MNAR-like (strong $A \times Y$)	OM	-0.348	0.399	0.576	0.424	0.683	250
700	MNAR-like (strong $A \times Y$)	IPW	-0.374	0.447	0.640	0.360	0.683	250
700	MNAR-like (strong $A \times Y$)	AIPW	-0.351	0.401	0.560	0.440	0.683	250

6 Conclusion

For trial-like binary outcomes with missingness, no single estimator was uniformly best across all stress conditions. OM and AIPW were typically most efficient under MAR-like settings, IPW was more variance-sensitive, and all MAR-based methods degraded under strong MNAR-like violations. Robust applied analysis should therefore combine principled primary estimation with explicit sensitivity analysis and transparent diagnostic reporting.

References

- [1] Rubin DB. Inference and missing data. *Biometrika*. 1976;63(3):581–592. DOI: 10.1093/biomet/63.3.581
- [2] Little RJA, Rubin DB. *Statistical Analysis with Missing Data*. 3rd ed. Hoboken, NJ: Wiley; 2019.
- [3] Sterne JAC, White IR, Carlin JB, et al. Multiple imputation for missing data in epidemiological and clinical research: potential and pitfalls. *BMJ*. 2009;338:b2393. DOI: 10.1136/bmj.b2393
- [4] White IR, Royston P, Wood AM. Multiple imputation using chained equations: Issues and guidance for practice. *Stat Med*. 2011;30(4):377–399. DOI: 10.1002/sim.4067
- [5] Seaman SR, White IR. Review of inverse probability weighting for dealing with missing data. *Stat Methods Med Res*. 2013;22(3):278–295. DOI: 10.1177/0962280210395740
- [6] Carpenter JR, Kenward MG. *Multiple Imputation and its Application*. Chichester, UK: Wiley; 2013.
- [7] Robins JM, Rotnitzky A, Zhao LP. Estimation of regression coefficients when some regressors are not always observed. *J Am Stat Assoc*. 1994;89(427):846–866. DOI: 10.1080/01621459.1994.10476818
- [8] Bang H, Robins JM. Doubly Robust Estimation in Missing Data and Causal Inference Models. *Biometrics*. 2005;61(4):962–973. DOI: 10.1111/j.1541-0420.2005.00377.x