

Inverse probability tilting for moment condition models  
with missing data, supplemental material:  
Additional proofs, Monte Carlo experiments,  
and further details on the empirical application

This supplemental web appendix contains additional proofs, summarizes the results of a series of Monte Carlo experiments, and provides further details on the empirical application. All notation is as defined in the main text unless explicitly defined otherwise. Equation and Table numbering continues in sequence with that established in the main text. To simplify notation let  $\beta$  denote the true parameter value  $\beta_0$  unless explicitly stated otherwise (similarly the ‘0’ subscript is removed from other objects, such as the propensity score, when doing so does not cause confusion).

**Proof of Proposition 2.1** Wooldridge (2007) proves consistency and asymptotic normality of IPW M-estimators. Here we derive the variance expression given in Proposition 2.1. Let  $\beta = (\gamma', \delta')'$  and define the moment vector and derivative matrix

$$m_i(\beta) = \begin{pmatrix} \frac{D_i}{G_i(\delta)} \psi_i(\gamma) \\ \frac{D_i - G_i(\delta)}{G_i(\delta)(1 - G_i(\delta))} G_{1i}(\delta) t_i \end{pmatrix}, \quad M_i(\beta) = \begin{bmatrix} \frac{D_i}{G_i(\delta)} \frac{\partial \psi_i(\beta)}{\partial \gamma'} & -\frac{D_i G_{1i}(\delta)}{G_i(\delta)^2} \psi_i(\gamma) t_i' \\ 0 & -J_i(\delta) \end{bmatrix}, \quad (52)$$

where  $J_i(\delta)$  is the  $i^{\text{th}}$  unit’s contribution to the Hessian of the log-likelihood for the propensity score parameter  $\delta$ . The solution to  $\sum_{i=1}^N m_i(\hat{\beta})/N = 0$  corresponds to the IPW estimate. The covariance of  $m_i$  is given by

$$\Omega = \begin{pmatrix} \mathbb{E} \left[ \frac{\psi \psi'}{G} \right] & \mathbb{E} \left[ \frac{G_1}{G} \psi t' \right] \\ \mathbb{E} \left[ \frac{G_1}{G} t \psi' \right] & \mathcal{I}(\delta) \end{pmatrix}, \quad (53)$$

while the population mean of  $M_i$  equals

$$M = \begin{pmatrix} \Gamma & -\mathbb{E} \left[ \frac{G_1}{G} \psi t' \right] \\ 0 & -\mathcal{I}(\delta_0) \end{pmatrix}, \quad (54)$$

with  $\mathcal{I}(\delta_0)$  the Fisher information for  $\delta_0$ .

Standard results on GMM imply that  $\sqrt{N}(\hat{\beta} - \beta_0)$  has a limiting sampling vari-

ance of

$$M^{-1}\Omega M^{-1'} = \begin{pmatrix} \Gamma^{-1}\mathbb{E}\left[\frac{\psi\psi'}{G}\right]\Gamma^{-1'} - \Gamma^{-1}\mathbb{E}\left[\frac{G_1}{G}\psi t'\right]\mathcal{I}(\delta_0)^{-1}\mathbb{E}\left[\frac{G_1}{G}t\psi'\right]\Gamma^{-1'} & 0 \\ 0 & \mathcal{I}(\delta_0)^{-1} \end{pmatrix}. \quad (55)$$

An insightful rearrangement of the upper-left-hand block of (55) is (see the detailed calculations supplement)

$$\mathcal{I}(\gamma_0)^{-1} + \Gamma^{-1}\mathbb{E}\left[\left(\left(\frac{D}{G} - 1\right)q - \Pi_S S_\delta\right)\left(\left(\frac{D}{G} - 1\right)q - \Pi_S S_\delta\right)'\right]\Gamma^{-1'}, \quad (56)$$

with  $\Pi_S$  as defined in the statement of Proposition 2.1. Part (ii) of the result follows by inspection.

**The asymptotic variance of three-step AIPW estimators** Here we summarize the first order asymptotic properties of a class of three-step AIPW estimators under Assumptions 1.1 to 1.5. This class includes the estimator proposed by Robins, Rotnitzky and Zhao (1994). As well as the variants proposed by Newey (1994), Hirano and Imbens (2001), and Cao, Tsiatis and Davidian (2009). While the first order properties of AIPW are well-known, we include the results below as they will prove useful for the higher order bias calculations.

We begin by developing some notation. Let  $\beta = (\gamma', \delta)'$  and define the  $K + 2(1 + M) \times 1$  moment vector and derivative matrix

$$m_i(\beta) = \begin{pmatrix} \frac{D_i}{G_i(\delta)}\psi_i(\gamma) \\ \left(\frac{D_i}{G_i(\delta)} - 1\right)t_i \\ \frac{D_i - G_i(\delta)}{G_i(\delta)(1 - G_i(\delta))}G_{1i}(\delta)t_i \end{pmatrix}, \quad M_i(\beta) = \begin{pmatrix} \frac{D_i}{G_i(\delta)}\frac{\partial\psi_i(\beta)}{\partial\gamma'} & -\frac{D_i}{G_i(\delta)}\frac{G_{1i}(\delta)}{G_i(\delta)}\psi_i(\gamma)t_i' \\ 0 & -\frac{D_i}{G_i(\delta)}\frac{G_{1i}(\delta)}{G_i(\delta)}t_i t_i' \\ 0 & J_i(\delta) \end{pmatrix}, \quad (57)$$

with  $J_i(\delta)$  as defined previously. Further define the weight matrix

$$V_i(\beta) = \begin{bmatrix} \frac{D_i}{G_i(\delta)^2}\psi_i(\gamma)\psi_i(\gamma)' & \frac{D_i}{G_i(\delta)}\omega_i(\delta)\psi_i(\gamma)t_i' & 0 \\ \frac{D_i}{G_i(\delta)}\omega_i(\delta)t_i\psi_i(\gamma)' & \nu_i(\delta)\omega_i(\delta)t_i t_i' & 0 \\ 0 & \frac{D_i}{G_i(\delta)}\frac{G_{1i}(\delta)}{G_i(\delta)}t_i t_i' & -J_i(\delta) \end{bmatrix}, \quad (58)$$

where  $\omega_i(\delta) = \omega(X_i, \delta)$  and  $\nu_i(\delta) = \nu(D_i, X_i, \delta)$  are known, scalar-valued, weight functions (the latter with the property that  $\mathbb{E}[\nu_i(\delta_0)|X] = 1$ ).

Finally let

$$\overline{M}(\beta) = \frac{1}{N} \sum_{i=1}^N M_i(\beta), \quad \overline{V}(\beta) = \frac{1}{N} \sum_{i=1}^N V_i(\beta), \quad \overline{m}(\beta) = \frac{1}{N} \sum_{i=1}^N m_i(\beta). \quad (59)$$

Our asymptotic analysis of three-step AIPW estimators exploits their representation as particular iterated GMM estimators.

**Lemma A.1** (ITERATED GMM REPRESENTATION OF AIPW) *The AIPW estimate  $\hat{\gamma}$  which solves (14) in the main text is numerically identical to the iterated GMM estimate  $\hat{\beta} = (\hat{\gamma}', \hat{\delta}')'$  which solves*

$$\overline{M}(\hat{\beta}) \overline{V}(\hat{\beta})^{-1} \overline{m}(\hat{\beta}) = 0. \quad (60)$$

The proof of Lemma A.1 is omitted but may be found online in the ‘detailed calculations’ supplement at <https://files.nyu.edu/bsg1/public/>. Invoking Lemma A.1 we proceed to characterize the large sample properties of (60). The population mean of the derivative of  $m_i(\beta_0)$ , as defined in (57) above, equals

$$M = \begin{bmatrix} \Gamma & -\mathbb{E} \left[ \frac{G_1}{G} \psi t' \right] \\ 0 & -\mathbb{E} \left[ \frac{G_1}{G} t t' \right] \\ 0 & -\mathcal{I}(\delta_0) \end{bmatrix}. \quad (61)$$

The probability limit of the weight matrix (58) is given by

$$V = \begin{bmatrix} \mathbb{E} \left[ \frac{\psi \psi'}{G} \right] & E_\omega & 0 \\ E'_\omega & F_\omega & 0 \\ 0 & \mathbb{E} \left[ \frac{G_1}{G} t t' \right] & \mathcal{I}(\delta_0) \end{bmatrix}, \quad (62)$$

with  $E_\omega = \mathbb{E}[\omega \psi t']$  and  $F_\omega = \mathbb{E}[\omega t t']$ .

The covariance of the moment vector  $m_i = m_i(\beta_0)$  is

$$\Omega = \begin{pmatrix} \mathbb{E} \left[ \frac{\psi \psi'}{G} \right] & E_0 & \mathbb{E} \left[ \frac{G_1}{G} \psi t' \right] \\ E'_0 & F_0 & \mathbb{E} \left[ \frac{G_1}{G} t t' \right] \\ \mathbb{E} \left[ \frac{G_1}{G} t \psi' \right] & \mathbb{E} \left[ \frac{G_1}{G} t t' \right] & \mathcal{I}(\delta_0) \end{pmatrix}, \quad (63)$$

with  $E_0$  and  $F_0$  as defined in the main text.

Using (61), (62) and (63) we get a limiting sampling variance for  $\sqrt{N}(\hat{\beta} - \beta_0)$  of

$$(M'V^{-1}M)^{-1} M'V^{-1}\Omega V^{-1}M (M'V^{-1}M)^{-1'} \quad (64)$$

$$= \begin{pmatrix} \begin{pmatrix} \Gamma^{-1} \left( \mathbb{E} \left[ \frac{\psi\psi'}{G} \right] - E_0 F_0^{-1} E_0' - \Delta_\omega \mathcal{I}(\delta_0)^{-1} \Delta_\omega' \right) \Gamma^{-1'} \\ + \Gamma^{-1} (E_0 F_0^{-1} - E_\omega F_\omega^{-1}) F_0 (E_0 F_0^{-1} - E_\omega F_\omega^{-1})' \Gamma^{-1'} \end{pmatrix} & 0 \\ 0 & \mathcal{I}(\delta_0)^{-1} \end{pmatrix},$$

with

$$\Delta_\omega = \mathbb{E} \left[ \frac{D}{G} \{ \psi - E_\omega F_\omega^{-1} t \} S_\delta' \right].$$

If Assumption 2.1 also holds (64) simplifies to  $diag \{ \mathcal{I}(\gamma_0)^{-1}, \mathcal{I}(\delta_0)^{-1} \}$ .

**Derivation of the higher order bias of AIPW (Theorem 3.1)** Let  $\theta = (\gamma', \delta', \lambda')$  be the  $T = 2K + 3(1 + M)$  vector of parameters of interest. Let  $\hat{\theta}$  be the solution to (35) with

$$b_i(\theta) = - \begin{bmatrix} M_i(\beta)' \lambda \\ m_i(\beta) + V_i(\beta)' \lambda \end{bmatrix}, \quad (65)$$

and  $M_i(\beta)$  and  $V_i(\beta)$  as defined by (57) and (58) above. The AIPW estimate corresponds to the first  $K \times 1$  components of  $\hat{\theta}$  since  $\bar{V}(\hat{\beta})^{-1} \bar{m}(\hat{\beta}) = -\hat{\lambda}$  so that  $\bar{M}(\hat{\beta})' \hat{\lambda} = \bar{M}(\hat{\beta})' \bar{V}(\hat{\beta})^{-1} \bar{m}(\hat{\beta}) = 0$  which, by Lemma A.1, is the first order condition for the AIPW estimator.

Similar to the treatment of two-step and iterated GMM by Newey and Smith (2004, cf., Lemmas A.5 and A.6), we derive the higher order bias properties of the AIPW estimate of  $\gamma_0$  by considering those of a simplified  $O_p(N^{-3/2})$  equivalent estimate.

We have, by Lemma A.4 of Newey and Smith (2004),  $\hat{\theta} - \theta_0 = \tilde{\phi}/\sqrt{N} + O_p(1/N)$  with  $\tilde{\phi} = \sum_{i=1}^N \phi_i/\sqrt{N}$  for  $\phi_i$  as defined in (37) above. This means that for  $q = 1, \dots, K + 1 + M$  we have  $\hat{\beta} - \beta_{q0} = e_q' \tilde{\phi}/\sqrt{N} + O_p(1/N)$ . Now consider the mean

value expansion

$$\begin{aligned}
\bar{V}(\hat{\beta}) &= \bar{V}(\beta_0) + \sum_{q=1}^{K+1+M} \frac{\partial \bar{V}(\bar{\beta})}{\partial \beta_q} (\hat{\beta} - \beta_{q0}) \\
&= \bar{V}(\beta_0) + \sum_{q=1}^{K+1+M} \mathbb{E} \left[ \frac{\partial V_i(\beta_0)}{\partial \beta_q} \right] (\hat{\beta} - \beta_{q0}) + O_p(1/N) \\
&= \bar{V}(\beta_0) + \sum_{q=1}^{K+1+M} \mathbb{E} \left[ \frac{\partial V_i(\beta_0)}{\partial \beta_q} \right] e'_q \tilde{\phi} / \sqrt{N} + O_p(1/N),
\end{aligned}$$

with  $\bar{V}(\beta)$  and  $V_i(\beta)$  as defined in (59) and (58) above;  $\bar{\beta}$  is mean value between  $\beta_0$  and  $\hat{\beta}$ . Let  $\xi_i = V_i - V + \sum_{q=1}^{K+1+M} \mathbb{E} \left[ \frac{\partial V_i(\beta_0)}{\partial \beta_q} \right] e'_q \phi_i$  with  $V$  as given in (62) above. This gives  $\bar{V}(\hat{\beta}) = V + \frac{1}{N} \sum_{i=1}^N \xi_i + O_p(1/N)$ . Now consider the solution to  $\bar{b}^*(\hat{\theta}) = \frac{1}{N} \sum_{i=1}^N b_i^*(\hat{\theta}) = 0$  with

$$b_i^*(\theta) = - \begin{bmatrix} M_i(\beta)' \lambda \\ m_i(\beta) + [V + \xi_i]' \lambda \end{bmatrix}. \quad (66)$$

Using the definitions above we have

$$\begin{aligned}
0 &= \bar{b}^*(\hat{\theta}) \\
&= \bar{b}(\hat{\theta}) + \begin{pmatrix} 0 \\ [\bar{V}(\hat{\beta}) - V - \frac{1}{N} \sum_{i=1}^N \xi_i]' \hat{\lambda} \end{pmatrix} \\
&= \bar{b}(\hat{\theta}) + O_p(N^{-3/2}),
\end{aligned} \quad (67)$$

since  $\bar{V}(\hat{\beta}) - V - \frac{1}{N} \sum_{i=1}^N \xi_i = O_p(1/N)$  and  $\hat{\lambda} = \lambda_0 + O_p(1/\sqrt{N})$ . Appealing to the equivalence implicit in the last line of (67) we henceforth analyze the bias properties of the solution to  $\bar{b}^*(\hat{\theta}) = 0$ . In what follows we redefine  $b_i(\theta)$  to be equal to  $b_i^*(\theta)$  as given in (66).

The terms defined in (37) are given by, recalling that  $\lambda_0 = 0$ ,

$$B = - \begin{pmatrix} 0 & M' \\ M & V \end{pmatrix}, \quad A_i = - \begin{pmatrix} 0 & (M_i - M)' \\ (M_i - M) & \xi_i \end{pmatrix},$$

with  $M$  and  $V$  given by (61) and (62) above.

The partitioned inverse formula gives

$$B^{-1} = - \begin{pmatrix} -\Upsilon & H \\ H' & L \end{pmatrix}, \quad (68)$$

with

$$\Upsilon = (M'V^{-1}M)^{-1}, \quad H = \Upsilon M'V^{-1}, \quad L = V^{-1} - V^{-1}MH. \quad (69)$$

Manipulating we get

$$\Upsilon = \begin{pmatrix} \Gamma^{-1}\Lambda\Gamma^{-1'} & 0 \\ \Pi'_S\Gamma^{-1'} & \mathcal{I}(\delta_0)^{-1} \end{pmatrix} \quad (70)$$

$$H = \begin{pmatrix} \Gamma^{-1} & -\Gamma^{-1}\Pi_0 & 0 \\ 0 & 0 & -\mathcal{I}(\delta_0)^{-1} \end{pmatrix} \quad (71)$$

$$L = \begin{pmatrix} 0 & 0 & -\Lambda^{-1}\Pi_S \\ 0 & F_\omega^{-1} & \Pi'_0\Lambda^{-1}\Pi_S \\ \Pi'_S\Lambda^{-1} & -\mathcal{I}(\delta_0)^{-1} \mathbb{E} \left[ \frac{G_1}{G} tt' \right] \left( F_\omega - E'_\omega \mathbb{E} \left[ \frac{\psi\psi'}{G} \right]^{-1} E_\omega \right)^{-1} & 0 \end{pmatrix}. \quad (72)$$

From (37) and the expressions above we also have  $\phi_i = - \begin{pmatrix} H' & L' \end{pmatrix}' m_i$ , where we evaluate at the population value of  $\theta$ .

The first part of the AIPW bias formula is given by the first  $K$  rows of  $\frac{-B^{-1}}{N} \mathbb{E} [A_i \phi_i]$ . Manipulating, using the expressions given above, we have

$$-B^{-1} \mathbb{E} [A_i \phi_i] = -\mathbb{E} \left[ \begin{pmatrix} -HM_iH + \Upsilon M'_iL - H\xi_iL \\ -LM_iH - H'M'_iL - L\xi_iL \end{pmatrix} m_i \right]. \quad (73)$$

We require expressions for the first  $K$  rows of the matrix  $\mathbb{E} [(HM_iH - \Upsilon M'_iL + H\xi_iL) m_i]$ .

Let  $\{A\}_{1:K,:}$  denote rows 1 to  $K$  of a matrix. Very tedious calculations give

$$\{H\mathbb{E}[M_i H m_i]\}_{1:K,:} = \Gamma^{-1} \mathbb{E} \left[ \frac{1}{G} \frac{\partial \psi(\beta)}{\partial \gamma'} \Gamma^{-1} (\psi - \Pi_0 t) \right] + \Gamma^{-1} \mathbb{E} \left[ \frac{\partial \psi(\beta)}{\partial \gamma'} \Gamma^{-1} \Pi_0 t \right] \quad (74)$$

$$- \{\Upsilon \mathbb{E}[M'_i L m_i]\}_{1:K,:} = \mathcal{I}(\gamma_0)^{-1} \mathbb{E} \left[ \frac{D}{G} \left( \frac{\partial \psi}{\partial \gamma'} \right)' \Lambda^{-1} \Pi_S S_\delta \right] \quad (75)$$

$$\begin{aligned} \{H\mathbb{E}[\xi_i L m_i]\}_{1:K,:} &= \Gamma^{-1} \Pi_0 \mathbb{E} \left[ \omega \left( \frac{D}{G} - \nu \right) \left( \frac{D}{G} - 1 \right) t t' F_0^{-1} t \right] \\ &\quad - \Gamma^{-1} \mathbb{E} \left[ \frac{D}{G} \psi \left( \frac{D}{G} \psi - \omega \Pi_0 t_0 \right)' \Lambda^{-1} \Pi_S S_\delta \right] \\ &\quad + \Gamma^{-1} \mathbb{E} \left[ \omega \left( \frac{D}{G} - \nu \right) \Pi_0 t t' \Pi_0 \Lambda^{-1} \Pi_S S_\delta \right]. \end{aligned} \quad (76)$$

Collecting these terms yields a bias contribution of

$$\begin{aligned} C_L &= \frac{\Gamma^{-1}}{N} \mathbb{E} \left[ \frac{1}{G} \frac{\partial \psi(\beta)}{\partial \gamma'} \Gamma^{-1} (\psi - \Pi_0 t) \right] + \frac{\Gamma^{-1}}{N} \mathbb{E} \left[ \frac{\partial \psi(\beta)}{\partial \gamma'} \Gamma^{-1} \Pi_0 t \right] \\ &\quad + \frac{\mathcal{I}(\gamma_0)^{-1}}{N} \mathbb{E} \left[ \frac{D}{G} \left( \frac{\partial \psi}{\partial \gamma'} \right)' \Lambda^{-1} \Pi_S S_\delta \right] \\ &\quad + \frac{\Gamma^{-1}}{N} \Pi_0 \mathbb{E} \left[ \omega \left( \frac{D}{G} - \nu \right) \left( \frac{D}{G} - 1 \right) t t' F_0^{-1} t \right] \\ &\quad - \frac{\Gamma^{-1}}{N} \mathbb{E} \left[ \frac{D}{G} \psi \left( \frac{D}{G} \psi - \omega \Pi_0 t_0 \right)' \Lambda^{-1} \Pi_S S_\delta \right] \\ &\quad + \frac{\Gamma^{-1}}{N} \mathbb{E} \left[ \omega \left( \frac{D}{G} - \nu \right) \Pi_0 t t' \Pi_0 \Lambda^{-1} \Pi_S S_\delta \right]. \end{aligned} \quad (77)$$

To compute the second component of (36) for  $\widehat{\gamma}_{AIPW}$  we require some additional notation and results. Recall that  $B_q = \mathbb{E}[\partial^2 b_i(\theta) / \partial \theta_q \partial \theta']$ . For  $q = 1, \dots, K + 1 + M$  we have this term equal to

$$B_q = - \begin{pmatrix} 0 & \mathbb{E} \left[ \frac{\partial M'_i}{\partial \beta_q} \right] \\ \mathbb{E} \left[ \frac{\partial M_i}{\partial \beta_q} \right] & 0 \end{pmatrix}, \quad (78)$$

while for  $q = K + 1 + M + 1, \dots, 2K + 3(1 + M)$  is it given by

$$B_q = - \begin{pmatrix} \mathbb{E} \left[ \frac{\partial^2 m_{q-K-1-M}(Z_i, \beta)}{\partial \beta \partial \beta'} \right] & 0 \\ 0 & 0 \end{pmatrix}. \quad (79)$$

We also have

$$\mathbb{E} [\phi_i \phi_i'] = \begin{pmatrix} H\Omega H' & H\Omega L' \\ L\Omega H' & L + L(\Omega - V)L' \end{pmatrix}, \quad (80)$$

with  $\Omega$  as given by (63) above. Note that  $LVL' = L$  so that  $L\Omega L' = LVL' + L(\Omega - V)L' = L + L(\Omega - V)L'$ .

Using these expressions we evaluate the second component of (36) as follows:

$$\begin{aligned} -\frac{B^{-1}}{2N} \mathbb{E} \left[ \sum_{q=1}^T \phi_{q,i} B_q \phi_i \right] &= -\frac{1}{2N} \sum_{q=1}^{K+1+M} \begin{pmatrix} -\Upsilon \mathbb{E} \left[ \frac{\partial M'_i}{\partial \beta_q} \right] L\Omega H' + H \mathbb{E} \left[ \frac{\partial M'_i}{\partial \beta_q} \right] H\Omega H' \\ + H' \mathbb{E} \left[ \frac{\partial M'_i}{\partial \beta_q} \right] L\Omega H' + L \mathbb{E} \left[ \frac{\partial M'_i}{\partial \beta_q} \right] H\Omega H' \end{pmatrix} e_q \\ &\quad (81) \\ &\quad -\frac{1}{2N} \sum_{q=1}^{K+2(1+M)} \begin{pmatrix} -\Upsilon \mathbb{E} \left[ \frac{\partial^2 m_q(Z_i, \beta)}{\partial \beta \partial \beta'} \right] H\Omega L' \\ + H' \mathbb{E} \left[ \frac{\partial^2 m_q(Z_i, \beta)}{\partial \beta \partial \beta'} \right] H\Omega L' \end{pmatrix} e_q. \end{aligned}$$

The first  $K$  rows of (81) contribute to the bias expression for  $\hat{\gamma}_{AIPW}$ . To determine the form of the first  $K$  rows of (81) we only require expressions for the first  $K$  rows of the two matrices in parentheses to the right of the equality in (81).

After tedious calculation we have, for  $q = 1, \dots, K$ ,

$$-\left\{ \Upsilon \mathbb{E} \left[ \frac{\partial M'_i}{\partial \beta_q} \right] L\Omega H' \right\}_{1:K,:} = \left( 0 \quad -\mathcal{I}(\gamma_0)^{-1} \mathbb{E} \left[ \frac{\partial^2 \psi}{\partial \gamma_q \partial \gamma} \right]' \Lambda^{-1} \Pi_S \right) \quad (82)$$

$$\left\{ H \mathbb{E} \left[ \frac{\partial M_i}{\partial \beta_q} \right] H\Omega H' \right\}_{1:K,:} = \left( \Gamma^{-1} \mathbb{E} \left[ \frac{\partial^2 \psi}{\partial \gamma_q \partial \gamma} \right] \mathcal{I}(\gamma_0)^{-1} \quad -\Gamma^{-1} \mathbb{E} \left[ \frac{D}{G} \frac{\partial \psi}{\partial \gamma_q} S'_\delta \right] \mathcal{I}(\delta_0)^{-1} \right) \quad (83)$$

Similarly, for  $q = K + 1, \dots, K + 1 + M$ , we have

$$-\left\{ \Upsilon \mathbb{E} \left[ \frac{\partial M'_i}{\partial \beta_q} \right] L\Omega H' \right\}_{1:K,:} = \left( 0 \quad \mathcal{I}(\gamma_0)^{-1} \mathbb{E} \left[ \frac{D}{G} \frac{\partial \psi}{\partial \gamma'} S'_{\delta, q-K} \right]' \Lambda^{-1} \Pi_S \right) \quad (84)$$

$$\left\{ H \mathbb{E} \left[ \frac{\partial M_i}{\partial \beta_q} \right] H\Omega H' \right\}_{1:K,:} = \left( -\Gamma^{-1} \mathbb{E} \left[ \frac{D}{G} \frac{\partial \psi}{\partial \gamma'} S'_{\delta, q-K} \right] \mathcal{I}(\gamma_0)^{-1} \quad 0 \right) \quad (85)$$



Using (81) and (82) to (85) we get, after some manipulation, a bias contribution of

$$C_{NL1} = -\frac{1}{2N} \sum_{q=1}^K \Gamma^{-1} \mathbb{E} \left[ \frac{\partial^2 \psi}{\partial \gamma_q \partial \gamma} \right] \mathcal{I}(\gamma_0)^{-1} e_q \quad (86)$$

$$- \frac{\mathcal{I}(\gamma_0)^{-1}}{2N} \mathbb{E} \left[ \frac{D}{G} \left( \frac{\partial \psi}{\partial \gamma'} \right)' \Lambda^{-1} \Pi_S S_\delta \right].$$

Now consider the second part of (81). For  $q = 1, \dots, K$  we have

$$- \left\{ \Upsilon \mathbb{E} \left[ \frac{\partial^2 m_q(Z_i, \beta)}{\partial \beta \partial \beta'} \right] H \Omega L' \right\}_{1:K,:} = \left( \mathcal{I}(\gamma_0)^{-1} \mathbb{E} \left[ \frac{D}{G} \frac{\partial \psi_q}{\partial \gamma} S'_\delta \right] \Pi'_S \Lambda^{-1} \quad (\text{NR}) \quad (\text{NR}) \right), \quad (87)$$

while for  $q = K + 1, \dots, K + 1 + M$

$$- \left\{ \Upsilon \mathbb{E} \left[ \frac{\partial^2 m_q(Z_i, \beta)}{\partial \beta \partial \beta'} \right] H \Omega L' \right\}_{1:K,:} = \begin{pmatrix} 0 & 0 & 0 \end{pmatrix}, \quad (88)$$

and finally for  $q = K + 1 + M + 1, \dots, K + 2(1 + M)$

$$- \left\{ \Upsilon \mathbb{E} \left[ \frac{\partial^2 m_q(Z_i, \beta)}{\partial \beta \partial \beta'} \right] H \Omega L' \right\}_{1:K,:} = \begin{pmatrix} 0 & 0 & 0 \end{pmatrix}. \quad (89)$$

Using (81) and (87) to (89) we get, after some manipulation, a bias contribution of

$$C_{NL2} = -\frac{\mathcal{I}(\gamma_0)^{-1}}{2N} \mathbb{E} \left[ \frac{D}{G} \left( \frac{\partial \psi}{\partial \gamma'} \right)' \Lambda^{-1} \Pi_S S_\delta \right]. \quad (90)$$

Equation (23) is given by the sum of (77), (86), and (90).

### Comparison of small sample properties of IPT with leading alternatives

We compare the small sample performance of IPT with that of several alternative missing data estimators. Specifically we consider the parametric inverse probability weighting estimator described by Wooldridge (2007), henceforth IPW; parametric linear imputation as in Rubin (1977), henceforth PI; the non-parametric IPW estimator of Hirano, Imbens and Ridder (2003), henceforth HIR; the nonparametric imputation estimator of Imbens, Newey and Ridder (2005) (with their data-dependent choice of smoothing parameter), henceforth INR; the nonparametric conditional expectation projection estimator of Chen, Hong and Tarozzi (2008), henceforth CHT; and the augmented inverse probability weighting estimator of Robins, Rotnitzky and Zhao

(1994) as described by Tsiatis (2006), henceforth AIPW-RRZ.<sup>24</sup>

We assume that  $Y_1$ , the outcome of interest, is generated according to

$$Y_1 = \alpha_0 + \alpha_1 X + \alpha_2 \Phi\left(\frac{X - a}{b}\right) + U, \quad U|X \sim \mathcal{N}(0, \sigma^2),$$

where  $\Phi(\cdot)$  is the CDF of a standard normal random variable and  $X \sim \mathcal{U}(-1, 1)$ . Observations of  $Y_1$  are missing at random with  $Y_1$  observed if  $D = 1$  where

$$D = \mathbf{1}(\beta_0 + \beta_1 X + \beta_2 \Phi(X/c) - V),$$

with  $V|X, Y_1$  logistic.

We consider four different data generating processes (DGPs). Across all designs we set  $a = 1/2$ ,  $b = 1/5$ ,  $c = 3/20$ , and  $\alpha_1 = -1/4$ . We vary  $\alpha_0$ ,  $\alpha_2$  and  $\beta_0$ ,  $\beta_2$  to induce nonlinearity in, respectively, the CEF of  $Y_1$  given  $X$  and the index of the propensity score model.

Each design is calibrated such that the propensity score ranges from 0.1 to 0.9 with a marginal probability of missingness equal to one half. The target estimand is  $\gamma_0 = \mathbb{E}[Y]$ , which is identically equal to zero for each design. For each experiment the sample size is set equal to  $N = 1,000$  with 5,000 Monte Carlo replications. We vary  $\sigma$  across designs in order to keep the variance bound constant such that, from the perspective of first order efficiency theory, each data generating process is equally ‘difficult’.

The values of the varying parameters are listed in Table 4. In the first design both the outcome CEF and the propensity score are smooth in  $X$ . In the second, the outcome CEF is inhomogenous, while the propensity score remains smooth. In the third, the CEF is smooth, while the propensity score is now inhomogenous. In the fourth design both the CEF and propensity score are inhomogenous.

The IPW, AIPW-RRZ and IPT estimators are based upon a logistic model for the propensity score with  $X$  entering the index linearly. The HIR estimator is based on a series logistic model for the propensity score with  $X$  entering quadratically.<sup>25</sup>

---

<sup>24</sup>Matlab replication files for our Monte Carlo experiments as well as a technical Appendix, describing in detail our implementation of each estimator, is available on the first author’s webpage.

<sup>25</sup>Hirano, Imbens and Ridder’s (2003) theoretical results suggest using a series logit estimator with  $N^{1/9}$  times some constant polynomial terms. For our designs  $1000^{1/9} \approx 2.15$  which we round up to 3, yielding the quadratic specification.

Table 4: Parameter values for the four Monte Carlo experiments.

	Design 1: 'smooth-smooth'	Design 2: 'rough-smooth'	Design 3: 'smooth-rough'	Design 4: 'rough-rough'
$\alpha_0$	0.00000	-0.12510	0.00000	-0.12510
$\alpha_2$	0.00000	0.50000	0.00000	0.50000
$\beta_0$	0.00000	0.00000	2.00000	2.00000
$\beta_1$	2.19722	2.19722	4.19722	4.19722
$\beta_2$	0.00000	0.00000	-4.00000	-4.00000
$\sigma$	0.16183	0.17327	0.17410	0.18640
$\sqrt{\mathcal{I}(\gamma_0)^{-1}/N}$	0.01000	0.01000	0.01000	0.01000

NOTES: The square root of Hahn's (1998) variance bound for each design (divided by  $N^{1/2} = \sqrt{1,000}$ ) is reported in the last row of the table.

The PI model is based on a linear model for  $\mathbb{E}[Y|X]$ . The INR estimator is based on a polynomial series estimate of  $\mathbb{E}[Y|X]$  with the number of series terms chosen to minimize estimated mean square error of  $\hat{\gamma}$  (see Imbens, Newey and Ridder, 2005). The CHT estimator is also based on a polynomial series estimate of  $\mathbb{E}[Y|X]$ , but with the number of series terms fixed at 3 (i.e., a quadratic approximation).

Our designs are chosen to highlight the strengths and weakness of each estimator. In the first design we expect all estimators to perform well. In the second design we expect the IPW and HIR estimators to perform acceptably and the parametric imputation (PI) estimator to perform poorly. In principal the nonparametric imputation estimators, INR and CHT, are consistent and efficient in this design, although we expect them to perform poorly in practice due to the inhomogeneity of  $\mathbb{E}[Y|X]$ . In design 3 we expect a pattern opposite to that of design 2. The AIPW-RRZ and IPT estimators, due to their double-robustness attribute, should perform well in both designs 2 and 3 (as well as design 1). In design 4 we expect all estimators to perform poorly.

The HIR, INR and CHT estimators are consistent and attain Hahn's (1998) bound in all designs. IPW is never efficient but is consistent in designs 1 and 2. PI is consistent and efficient in designs 1 and 3. AIPW-RRZ and IPT are consistent in designs 1, 2 and 3. In design 1 they attain Hahn's (1998) bound. In design 2 their large sample variance lies above Hahn's (1998) bound but below that of parametric

IPW. In design 3 their large sample variance is actually smaller than Hahn’s (1998) bound.<sup>26</sup>

For purposes of comparison we also include an infeasible oracle estimator: the mean of  $Y$  across all missing as well as non-missing observations. This estimator calibrates the cost of missingness in each of our designs.

Table 5 summarizes the results for the first design. Column 1 of the Table reports the asymptotic bias of each estimator and Column 2 the median bias across Monte Carlo replications (both scaled by the relevant estimator’s asymptotic standard error). Column 3 reports each estimator’s asymptotic standard error, Column 4 the median estimated standard error across Monte Carlo replications and Column 5 the standard deviation of the point estimates across Monte Carlo replications. Column 6 reports the coverage of a nominal 95 percent confidence interval. As expected each estimator performs well in design 1, with small sample properties well-approximated by asymptotic distribution theory. An approximate standard error for the Column 2 Monte Carlo scaled bias estimates is  $\sqrt{(\pi/2)/5,000} \approx 0.0177$ ; differences in median bias are not significant across estimators for this design.<sup>27</sup>

Table 6 reports results from the second design. In this design the inhomogeneity of  $\mathbb{E}[Y|X]$  creates problems for PI as well as both the INR and CHT imputation estimators, all of which exhibit significant median bias. Unsurprisingly, given the underlying smoothness of the propensity score, both the IPW and HIR estimators do well. In this design the HIR procedure involves an overfit of the propensity score. Consistent with the theoretical results of their paper, such overfitting results in improved precision (see Graham (2011) for a related discussion). The Monte Carlo sampling standard deviation of the HIR point estimates are about 5 percent lower than the corresponding IPW estimates (cf., Column 5).

Both AIPW-RRZ and IPT estimators also perform well in design 2. When both  $\mathbb{E}[Y|X]$  and  $p(X)$  are correctly modelled the two estimators are first-order equiva-

---

<sup>26</sup>Bang and Robins (2005, p. 966) comment that, when the propensity score is misspecified, but the outcome CEF is not, AIPW is often ‘nearly’ as efficient as parametric imputation. This assessment is based on their Monte Carlo experiments. In the designs whose results are reported in the lower half of Table 2 (p. 966) and the upper and lower portions of Table 4 (p. 970) of their paper, the AIPW estimator with a misspecified propensity score has a Monte Carlo sampling variance that lies below that of the AIPW estimator based on a correct model for the propensity score, but above that of the parametric imputation estimator. Our Design 3 results replicate this ordering.

<sup>27</sup>The standard error of the median bias estimates are approximately  $[(\pi/2)/5000]^{1/2} \simeq 0.01772454$  since for  $N$  and  $B$ , the number of Monte Carlo replications, large enough  $\hat{\gamma}^2/[I(\gamma_0)^{-1}/N]^{1/2} \stackrel{D}{\simeq} N(\text{Bias}_A, 1)$  and  $\sqrt{B}(\text{Bias}_{MC} - \text{Bias}_A) \stackrel{D}{\simeq} N(0, \pi/2)$ .

Table 5: Monte Carlo results for Design 1: both  $p(x)$  and  $E[Y|X=x]$  smooth

	(1)	(2)	(3)	(4)	(5)	(6)
	Asymptotic Bias	Median Bias	Asymptotic Std. Err.	Median Std. Err.	Standard Deviation	Coverage of 95% CI
<b>Oracle</b>	0.0000	-0.0398*	0.0069	0.0069	0.0068	0.951
	Parametric estimators					
<b>IPW</b>	0.0000	-0.0128	0.0107	0.0106	0.0108	0.947
<b>PI</b>	0.0000	-0.0129	0.0096	0.0096	0.0097	0.951
	Double robust parametric estimators					
<b>AIPW</b>	0.0000	-0.0195	0.0100	0.0099 [0.0099]	0.0101	0.948 [0.947]
<b>IPT</b>	0.0000	-0.0205	0.0100	0.0099	0.0101	0.946
	Nonparametric estimators					
<b>HIR</b>	0.0000	-0.0163	0.0100	0.0100	0.0103	0.945
<b>INR</b>	0.0000	-0.0109	0.0100	0.0099	0.0100	0.948
<b>CHT</b>	0.0000	-0.0271	0.0100	0.0099	0.0101	0.948

NOTES: Each row corresponds to a specific estimator as described in the text. Column 1 reports the scaled large sample bias of each estimator (i.e., its probability limit minus the true parameter divided by the square root of its large sample variance,  $(AVar(\hat{\gamma})/N)^{1/2}$ ). Column 2 reports the median Monte Carlo bias of each estimator scaled by its asymptotic standard error (a \* (+) next to these estimates denotes that they are significantly different from zero at the 5 (10) percent level). Column 1 calibrates the scale of inconsistency for each estimator, while a comparison of Columns 1 and 2 allows for an assessment of whether an estimator's actual sampling distribution is centered at its probability limit. Column 3 gives the large sample standard error of each estimator (i.e.,  $(AVar(\hat{\gamma})/N)^{1/2}$ ), Column 4 the median estimated standard error and column 5 the standard deviation of the point estimates across the 5,000 Monte Carlo replications. Column 6 reports the actual coverage of a 95 percent Wald-based confidence interval (the asymptotic variance estimators are described in a technical appendix available online). Two standard errors are provided for the AIPW-RRZ estimator. The first assumes that the working models for both  $E[Y|X]$  and  $p(X)$  are correctly specified. This is the covariance estimator typically used in applied work. The second, given in the square brackets, treats the AIPW-RRZ estimator as a sequential method-of-moments estimator with standard errors calculated accordingly. These standard errors are valid whenever AIPW-RRZ is consistent. The mean number of series terms selected by the INR estimator is 2.26, while the median is 2.

Table 6: Monte Carlo results for Design 2:  $p(x)$  smooth and  $E[Y|X=x]$  rough

	(1)	(2)	(3)	(4)	(5)	(6)
	Asymptotic Bias	Median Bias	Asymptotic Std. Err.	Median Std. Err.	Standard Deviation	Coverage of 95% CI
<b>Oracle</b>	0.0000	-0.0231	0.0063	0.0063	0.0061	0.957
	Parametric estimators					
<b>IPW</b>	0.0000	-0.0136	0.0106	0.0105	0.0107	0.941
<b>PI</b>	-3.0095	-3.0599*	0.0093	0.0114	0.0111	0.278
	Double robust parametric estimators					
<b>AIPW</b>	0.0000	-0.0246	0.0114	0.0124 [0.0113]	0.0116	0.961 [0.943]
<b>IPT</b>	0.0000	-0.0067	0.0104	0.0103	0.0105	0.942
	Nonparametric estimators					
<b>HIR</b>	0.0000	-0.0089	0.0100	0.0099	0.0103	0.941
<b>INR</b>	0.0000	0.1226*	0.0100	0.0104	0.0106	0.945
<b>CHT</b>	0.0000	0.0461*	0.0100	0.0100	0.0102	0.951

NOTES: The mean number of series terms selected by the INR estimator is 3.07, while the median is 5. See notes to Table 5 for additional information.

lent. However under partial misspecification they are no longer first-order equivalent. In design 2, while both are consistent, their asymptotic variances differ, with IPT being more precisely determined (Columns 1 and 3 of Table 6). The small sample properties of the two estimators mirror their large sample ones. Both are approximately median unbiased, but the sampling variability of the IPT estimator is about 10 percent lower than that of AIPW (Columns 2 and 5 of Table 6).

Table 7 reports results from the third design. In this case the IPW and HIR perform extremely poorly, reflecting their fragility vis-a-vis misspecification of the propensity score. The remaining estimators all perform well. Importantly, the AIPW-RRZ and IPT estimators, while based on misspecified models for the propensity score, are nevertheless consistent. However, as in design 2, their asymptotic variances differ with IPT's again being smaller, albeit negligibly so. These large sample properties are reflected in small samples. In design 3 both AIPW-RRZ and IPT are approximately median unbiased with similar sampling variances (Columns 2 and 5 of Table 7).

Table 8 reports results from the our fourth design. In this design the AIPW-RRZ and IPT estimators are inconsistent as are IPW and PI. While HIR, INR and CHT are consistent and efficient in principle, in practice they exhibit significant median

Table 7: Monte Carlo results for Design 3:  $p(x)$  rough and  $E[Y|X=x]$  smooth

	(1)	(2)	(3)	(4)	(5)	(6)
	Asymptotic Bias	Median Bias	Asymptotic Std. Err.	Median Std. Err.	Standard Deviation	Coverage of 95% CI
<b>Oracle</b>	0.0000	0.0178	0.0072	0.0072	0.0072	0.949
Parametric estimators						
<b>IPW</b>	-0.3204	-0.3270*	0.0093	0.0092	0.0094	0.937
<b>PI</b>	0.0000	-0.0078	0.0093	0.0093	0.0095	0.944
Double robust parametric estimators						
<b>AIPW</b>	0.0000	-0.0106	0.0094	0.0093 [0.0093]	0.0095	0.943 [0.944]
<b>IPT</b>	0.0000	-0.0115	0.0094	0.0093	0.0095	0.944
Nonparametric estimators						
<b>HIR</b>	0.0000	-0.3028*	0.0100	0.0092	0.0094	0.937
<b>INR</b>	0.0000	-0.0002	0.0100	0.0093	0.0099	0.934
<b>CHT</b>	0.0000	-0.0016	0.0100	0.0093	0.0096	0.942

NOTES: The mean number of series terms selected by the INR estimator is 3.08, while the median is 2. See notes to Table 5 for additional information.

Table 8: Monte Carlo results for Design 4:  $p(x)$  rough and  $E[Y|X=x]$  rough

	(1)	(2)	(3)	(4)	(5)	(6)
	Asymptotic Bias	Median Bias	Asymptotic Std. Err.	Median Std. Err.	Standard Deviation	Coverage of 95% CI
<b>Oracle</b>	0.0000	-0.0081	0.0066	0.0066	0.0067	0.947
Parametric estimators						
<b>IPW</b>	0.0614	0.0271	0.0096	0.0095	0.0096	0.945
<b>PI</b>	-0.7273	-0.7635*	0.0089	0.0100	0.0099	0.9006
Double robust parametric estimators						
<b>AIPW</b>	-0.0487	-0.0859*	0.0098	0.0098 [0.0097]	0.0095	0.946 [0.946]
<b>IPT</b>	0.0499	0.0224	0.0096	0.0096	0.0097	0.946
Nonparametric estimators						
<b>HIR</b>	0.0000	0.0324 <sup>+</sup>	0.0100	0.0092	0.0092	0.947
<b>INR</b>	0.0000	0.1821*	0.0100	0.0091	0.0102	0.919
<b>CHT</b>	0.0000	0.6957*	0.0100	0.0092	0.0094	0.873

NOTES: The mean number of series terms selected by the INR estimator is 5, while the median is 5. See notes to Table 5 for additional information.

bias.

The performance of IPT across designs 2 and 3 highlights the practical value of using a ‘doubly robust’ estimator. If either the propensity score or outcome CEF is well approximated by the implicit parametric models used, then IPT will perform well. While the HIR, INR and CHT estimators have attractive large sample properties irrespective of the form of the propensity score and outcome CEF, in practice their small performance is highly sensitive to the smoothness of one of these two functions. The HIR estimator performs well when the propensity score is smooth, but very poorly when it is not. The INR and CHT imputation estimators, in contrast, perform best when the outcome response is smooth.

The Monte Carlo experiments also highlight differences between the AIPW-RRZ estimator of Robins, Rotnitzky and Zhao (1994) and our IPT procedure. Under misspecification of either  $\mathbb{E}[Y|X]$  or  $p(X)$ , but not both simultaneously, they remain consistent but have different large sample variances. In the designs considered here IPT has a smaller asymptotic sampling variance than AIPW-RRZ in such cases.

**Monte Carlo ‘verification’ of Higher Order bias expressions in Theorem 3.1** To assess the small sample relevance of the higher-order bias expressions given in Theorem 3.1 we conducted two additional Monte Carlo experiments. Consider the DGP

$$Y_1 = (\gamma_0 - \sigma_X^2) + \Pi h(X) + U, \quad U|X \sim \mathcal{N}(0, \sigma_U^2), \quad (91)$$

where  $\Pi = (0, 1)$  and  $h(X) = (X, X^2)'$  with  $X \sim \mathcal{N}(0, \sigma_X^2)$ . It is straightforward to show that the variance bound for  $\gamma_0 = \mathbb{E}[Y] = 0$  when  $Y_1$  is MCAR with probability  $Q_0$  is equal to

$$\mathcal{I}(\gamma_0)^{-1} = \frac{\sigma_U^2}{Q_0} + 2\sigma_X^4.$$

We also have  $\mathbb{V}(h(X)) = \text{diag}\{\sigma_X^2, 2\sigma_X^4\}$  so that evaluating the bias expressions given in Theorem 3.1 yields

$$\text{Bias}(\hat{\gamma}_{IPT}) = 0, \quad \text{Bias}(\hat{\gamma}_{AIPW-NEWHEY}) = -\frac{1}{N} \frac{1 - Q_0}{Q_0} 6\sigma_X^2.$$

When the propensity score and the conditional variance of  $Y_1$  are both constant in  $X$ , the biases of the other AIPW estimators listed in Table 1 coincide with that of IPT. We therefore focus on a  $\hat{\gamma}_{IPT}$  and  $\hat{\gamma}_{AIPW-NEWHEY}$  comparison in what follows.



Table 9: Higher order bias comparisons of IPT and AIPW-NEWWEY with data MCAR: Design 1

Q	IPT			AIPW-NEWWEY		
	Bias <sub>A</sub>	Bias <sub>MC</sub>	Coverage of 95% CI	Bias <sub>A</sub>	Bias <sub>MC</sub>	Coverage of 95% CI
0.8	0.0000	-0.0100	0.9520	-0.0300	-0.0390	0.7802
0.6	0.0000	-0.0080	0.9482	-0.0800	-0.0879	0.9010
0.4	0.0000	0.0045	0.9440	-0.1800	-0.1750	0.9676
0.2	0.0000	-0.0445	0.9460	-0.4800	-0.4667	0.9858

**Notes:** The table reports Monte Carlo estimates of the scaled median bias ( $Bias_{MC}$ ) of the IPT and AIPW-NEWWEY estimates of  $\gamma_0$  (i.e., the median bias of  $\hat{\gamma}$  divided by its asymptotic standard error,  $[I(\gamma_0)^{-1}/N]^{1/2}$ ). The DGP is as specified in equation (91) and the text immediately following. Reported statistics are based on 5,000 Monte Carlo draws. The column labelled  $Bias_A$  gives the asymptotic bias of  $\hat{\gamma}$  – again divided by  $[I(\gamma_0)^{-1}/N]^{1/2}$  – as calculated analytically using the expressions given in Theorem 3.1 of the main text.

In moderately sized samples,  $\hat{\gamma}_{AIPW-NEWWEY}$  may exhibit non-negligible bias when  $\sigma_X^2$  is large relative to  $\sigma_U^2$  and/or  $Q_0$  is small. This bias, which arises despite the fact that  $h(X)$  is low dimensional, comes from skewness in the distribution of  $X^2$ .

Table 9 reports the asymptotic bias, and Monte Carlo estimates of the median bias, of  $\hat{\gamma}_{IPT}$  and  $\hat{\gamma}_{AIPW-NEWWEY}$  (scaled by their asymptotic standard error  $\sqrt{\mathcal{I}(\gamma_0)^{-1}/N}$ ) for the above DGP with  $\sigma_X^2 = 4$ ,  $Q = 0.8, 0.6, 0.4$  and  $0.2$ , and  $\sigma_U^2$  varied across designs in order to keep the variance bound fixed at 40 (this implies an asymptotic standard error for  $\hat{\gamma}$  of  $1/5$ ). We set  $N = 1,000$  and perform 5,000 Monte Carlo replications. Each of these designs is first-order equivalent, as are  $\hat{\gamma}_{IPT}$  and  $\hat{\gamma}_{AIPW-NEWWEY}$ .

Across each of these designs IPT is median unbiased (up to simulation error), consistent with Theorem 3.1. In contrast AIPW-NEWWEY is very biased, particularly for the designs with a high probability of missingness. Across each design Theorem 3.1 does a good job of predicting the small sample median bias of AIPW-NEWWEY.

A second example illustrates how the asymptotic bias of AIPW-NEWWEY may grow with the dimension of set of conditioning variables. Assume that

$$Y_1 = \gamma_0 + \Pi X + U, \quad U|X \sim \mathcal{N}(0, \sigma_U^2), \quad (92)$$

where  $X$  is an  $M \times 1$  vector of independent standardized  $\chi_1^2$  random variables (i.e.,  $X = \sqrt{1/2}(W - 1)$  with  $W_m \sim \chi_1^2$  for  $m = 1, \dots, M$ ) and  $Y_1$  is missing completely at random with a probability of observation equal to  $Q_0 = 1/2$ . The target parameter is  $\gamma_0 = \mathbb{E}[Y_1] = 0$ . Let  $\Pi$  be a column vector of ones such that  $\Pi\Pi' = M$ , it is straightforward to show that the variance bound for  $\gamma_0$  is equal to

$$\mathcal{I}(\gamma_0)^{-1} = \frac{\sigma_U^2}{Q_0} + M.$$

Under this DGP evaluating the bias expressions given in Theorem 3.1 yields

$$\text{Bias}(\hat{\gamma}_{IPT}) = 0, \quad \text{Bias}(\hat{\gamma}_{AIPW-NEWHEY}) = -\frac{M}{N}\sqrt{8},$$

so that  $\hat{\gamma}_{IPT}$  is asymptotically unbiased while  $\hat{\gamma}_{AIPW-NEWHEY}$  has a higher order bias which increases linearly with the dimension of  $X$ .

We consider six designs corresponding to  $M = 1, 5, 10, 15, 20$  and  $25$ . In all designs  $N = 1,000$ . Again, we vary  $\sigma_U^2$  across designs in order to keep the variance bound fixed at 40.

Asymptotic bias and the median bias of the Monte Carlo estimates (again based on 5,000 replications), both scaled by the asymptotic standard error, for IPT and AIPW-NEWHEY are reported in Table 10. As theory would predict the median bias of  $\hat{\gamma}_{AIPW-NEWHEY}$  is increasing in the dimension of  $X$ , while  $\hat{\gamma}_{IPT}$  is median unbiased across all designs. In all cases our Monte Carlo median bias estimates are within simulation error of their theoretical counterparts. For the  $M = 25$  design the median bias of  $\hat{\gamma}_{AIPW-NEWHEY}$  is about one third of its asymptotic standard error.

**Additional details of empirical application** Our initial goal was to reconstruct the exact NLSY79 extract used by Johnson and Neal (1998). A preliminary inspection of the data, however, revealed that a preadolescence test score was available for only a handful of Hispanic respondents. We therefore decided to exclude Hispanics from our analysis. We targeted all non-Hispanic male respondents in the cross-sectional sample, as well as in the supplemental Black sample, born during or after 1962 ( $N = 1,612$ ). These individuals were aged 16 to 18 when they took the NLSY79's administration of the Armed Services Vocational Aptitude Battery (ASVAB) from which the AFQT score is constructed.

Table 10: Higher order bias comparisons of IPT and AIPW-NEWWEY with data MCAR: Design 2

M	IPT			AIPW-NEWWEY		
	Bias <sub>A</sub>	Bias <sub>MC</sub>	Coverage of 95% CI	Bias <sub>A</sub>	Bias <sub>MC</sub>	Coverage of 95% CI
1	0.0000	0.0126	0.9518	-0.0141	-0.0015	0.9512
5	0.0000	-0.0126	0.9512	-0.0707	-0.0863	0.9480
10	0.0000	0.0026	0.9484	-0.1414	-0.1317	0.9430
15	0.0000	-0.0092	0.9528	-0.2121	-0.2118	0.9446
20	0.0000	-0.0208	0.9460	-0.2828	-0.3061	0.9292
25	0.0000	0.0228	0.9546	-0.3536	-0.3259	0.9276

**Notes:** The table reports Monte Carlo estimates of the scaled median bias ( $Bias_{MC}$ ) of the IPT and AIPW-NEWWEY estimates of  $\gamma_0$  (i.e., the median bias of  $\hat{\gamma}$  divided by its asymptotic standard error,  $[I(\gamma_0)^{-1}/N]^{1/2}$ ). The DGP is as specified in equation (92) and the text immediately following. Reported statistics are based on 5,000 Monte Carlo draws. The column labelled  $Bias_A$  gives the asymptotic bias of  $\hat{\gamma}$  – again divided by  $[I(\gamma_0)^{-1}/N]^{1/2}$  – as calculated analytically using the expressions given in Theorem 3.1 of the main text.

Of the 1,612 targeted individuals 159 were missing 1990 to 1993 average hourly wage data, 58 a valid AFQT score, and 24 both these items.<sup>28</sup> Our base sample therefore consists of the 1,371 individuals with valid wage and AFQT data. As in Johnson and Neal (1998) we conditioned on this non-response.

AFQT scores are reported in terms of normed percentiles (i.e., relative to population of American youths aged 18 to 23).<sup>29</sup> We transformed these scores onto the real line using the inverse normal CDF. The residual associated with the least squares fit of the transformed scores onto a vector of birth year dummies is our AFQT measure.<sup>30</sup> The average hourly wage for Whites in the base sample is \$12.32 in 1993

<sup>28</sup>Hourly wages are equal to their average over the 1990 to 1993 survey waves. In cases where an individual was not employed or interviewed in a given year, the average is over those years with wage information. Reported wage values less than \$1 and greater than \$75 per hour are discarded. Wages are measured in 1993 dollars. As in Johnson and Neal (1998) we exclude any AFQT score where the testing protocol was non-standard.

<sup>29</sup>We used the 1989 scoring of the AFQT.

<sup>30</sup>Neal and Johnson (1996) appear to have used the residual associated with the least squares fit of AFQT *percentiles* onto a set of birth year dummies (see their Figure 2, p. 886). Even after standardizing this variable to have mean zero and unit variance its distribution is non-normal. Our approach results in an AFQT score distribution that is indistinguishable from a normal one.

prices. Blacks earn on average \$2.87 less per hour. The mean AFQT score for Whites is 0.160, while that for Blacks is 1.031 standard deviations lower.

For 11 percent of the respondents ( $N = 144$ ) in our base sample an IQ test score from between the ages of 7 and 12 is available. This score is recorded as a (nationally normed) percentile. As with the AFQT scores we transformed these scores onto the real line using the inverse normal CDF. The residual associated with the least squares fit of these transformed scores onto a vector of birth year and age when tested (in years) dummies is our EARLYTEST measure. We view this score as measure of acquired cognitive skills not of innate ability (cf., Fryer and Levitt, forthcoming).<sup>31</sup>

Columns 1, 2 and 3 of Table 11 respectively the report average values of our main variables across the full sample, the complete case subsample ( $D = 1$ ), and the subsample with missing early test scores ( $D = 0$ ). Column 4 gives the Column 2 versus Column 3 difference. While the full and complete case samples are similar in terms of wages and age, they are rather different in terms of racial composition and AFQT.

Table 12 provides a synopsis of the results given in the main paper. Conditioning on age alone the Black-White gap in wages is 28 percent (Column 1). Conditioning on AFQT, a measure of cognitive skills acquired by age 18, this gap falls to 11 percent. Conditioning on our early test measure, a measure of cognitive skills acquired by age 12, the gap is 18 percent. We interpret this result as implying that approximately two thirds of the pre-market effect found by Neal and Johnson (1996) reflects skill differences already present by age 12. This interpretation is justified by two additional pieces of evidence. First, when we include both AFQT and EARLYTEST simultaneously in our model we find that the coefficient on the latter is insignificantly different from zero, while that on the former is insignificantly different from the estimate which conditions on AFQT alone (Column 4 of Table 12). This suggests that AFQT and EARLYTEST measure similar types of skills, with the later measure being relevant for labor market outcomes. A similar interpretation is suggested by Table 13 which shows that roughly two thirds of the AFQT gap across Blacks and Whites can be accounted for by skill differences present by age 12 (i.e., our EARLYTEST variable).

Table 14 reports the CMLE and IPT estimates of the propensity score parameter.

---

<sup>31</sup>If a respondent's record includes multiple test scores from the age 7 to 12 period, we use the average percentile score across all available tests. A *STATA* dictionary file of the NLSY79 extract used in our analysis as well as a do file which replicates the data manipulations described here is available online at <https://files.nyu.edu/bsg1/public/>.

Table 11: Comparison of the full sample with the complete case subsample

	(1) Full Sample	(2) $D = 1$	(3) $D = 0$	(4) Difference
LOGWAGE		6.999	6.982	0.0167 (0.0433)
YEAROFBIRTH		62.91	62.99	-0.0842 (0.0824)
BLACK		0.1113	0.1603	-0.0490 (0.0218)*
AFQT		0.1960	-0.0246	0.2205 (0.0953)*
$N$	1,371	144	1,227	

NOTES: Samples are as described in text. The 1979 baseline sampling weights are used when computing all summary statistics. A ‘\*\*\*’, ‘\*’ and ‘+’ denotes that the Column 4 difference is significantly different from zero at the 1, 5 and 10 percent levels. Standard errors (in parentheses) allow for arbitrary patterns of heteroscedasticity and dependence across units residing in the same household at baseline.

## References

- [1] Fryer, Roland G. and Steven D. Levitt. (forthcoming). “Testing for racial differences in the mental ability of young children,” *American Economic Review*.
- [2] Imbens, Guido W., Whitney K. Newey and Geert Ridder (2005). “Mean-square-error calculations for average treatment effects,” *IEPR Working Paper 05.34*.
- [3] Rubin, Donald B. (1977). “Assignment to treatment group on the basis of a covariate,” *Journal of Educational Statistics* 2 (1): 1 - 26.

Table 12: Pre-adolescent skills and adult wages: a synopsis

	(1)	(2)	(3)	(4)
	<i>OLS</i>	<i>OLS</i>	<i>IPT</i>	<i>IPT</i>
YEAROFBIRTH	-0.0458 (0.0151)**	-0.0466 (0.0147)**	-0.0537 (0.0162)**	-0.0443 (0.0147)**
BLACK	-0.2776 (0.0261)**	-0.1079 (0.0284)**	-0.1837 (0.0356)**	-0.1132 (0.0293)**
AFQT	—	0.1645 (0.0146)**	—	0.1866 (0.0284)**
EARLYTEST	—	—	0.1112 (0.0296)**	-0.0332 (0.0374)

NOTES: Samples are as described in the main text. The 1979 baseline sampling weights are used when computing all estimates. A ‘\*\*’, ‘\*’ and ‘+’ denotes that a coefficient is significantly different from zero at the 1, 5 and 10 percent levels. Standard errors (in parentheses) allow for arbitrary patterns of heteroscedasticity and dependence across units residing in the same household at baseline.

Table 13: Understanding the AFQT gap at age 18

	(1)	(2)
	<i>OLS</i>	<i>IPT</i>
YEAROFBIRTH	0.0050 (0.0330)	-0.0502 (0.0493)
BLACK	-1.0314 (0.0519)**	-0.3780 (0.1205)**
EARLYTEST	—	0.7739 (0.0790)**

NOTES: Samples are as described in the main text. The 1979 baseline sampling weights are used when computing all estimates. A ‘\*\*’, ‘\*’ and ‘+’ denotes that a coefficient is significantly different from zero at the 1, 5 and 10 percent levels. Standard errors (in parentheses) allow for arbitrary patterns of heteroscedasticity and dependence across units residing in the same household at baseline.

Table 14: Propensity score estimates and balance assessment (Part 1 of 2)

	Panel A: P-Score		Panel B: Balance		
	(1) CMLE	(2) IPT	(1) FS	(2) CC-IPW	(3) CC-IPT
BLACK	-51.737 (62.573)	108.174 (90.367)	0.155	0.171	0.155
BORNIN1963	-67.995 (96.750)	-292.099 (83.740)**	0.325	0.314	0.325
BORNIN1964	-118.245 (67.851)	-370.900 (96.265)**	0.330	0.326	0.330
BLACK×BORNIN1963	0.717 (0.602)	1.029 (0.596) <sup>+</sup>	0.052	0.057	0.052
BLACK×BORNIN1964	-0.566 (0.689)	-0.721 (0.590)	0.052	0.065	0.052
AFQT	1.012 (2.615)	-1.650 (0.945) <sup>+</sup>	0.000	-0.029	0.000
LOGWAGE	-7.109 (6.882)	-3.414 (5.392)	6.984	6.845	6.984
LOGWAGE×AFQT	-0.095 (0.376)	0.297 (0.145)*	0.175	0.004	0.175
AFQT <sup>2</sup>	-0.040 (0.128)	-0.094 (0.074)	0.975	0.926	0.975
LOGWAGE <sup>2</sup>	0.484 (0.490)	0.202 (0.381)	48.965	47.829	48.965
BLACK×AFQT	-7.836 (6.337)	10.380 (12.650)	-0.135	-0.149	-0.135
BLACK×LOGWAGE	13.680 (17.777)	-31.672 (25.935)	1.045	1.139	1.045
BLACK×LOGWAGE×AFQT	1.046 (0.880)	-1.570 (1.771)	-0.892	-0.979	-0.892
BLACK×AFQT <sup>2</sup>	-0.052 (0.284)	0.117 (0.387)	0.215	0.211	0.215
BLACK×LOGWAGE <sup>2</sup>	-0.908 (1.260)	2.305 (1.863)	7.082	7.641	7.082

NOTES: See notes to Part 2 of the table.

Table 15: Propensity score estimates and balance assessment (Part 2 of 2)

	Panel A: P-Score		Panel B: Balance		
	(1) CMLE	(2) IPT	(1) FS	(2) CC-IPW	(3) CC-IPT
BORNIN1963×AFQT	-5.723 (6.238)	-52.878 (17.479)**	0.000	-0.006	0.000
BORNIN1963×LOGWAGE	18.208 (28.007)	78.757 (22.800)**	2.276	2.206	2.276
BORNIN1963×LOGWAGE×AFQT	0.815 (0.898)	7.467 (2.483)**	0.061	0.045	0.061
BORNIN1963×AFQT <sup>2</sup>	-0.111 (0.210)	-1.045 (0.358)**	0.332	0.329	0.332
BORNIN1963×LOGWAGE <sup>2</sup>	-1.221 (2.023)	-5.295 (1.545)**	15.991	15.554	15.991
BORNIN1964×AFQT	-1.716 (5.026)	-13.540 (6.657)*	0.000	-0.027	0.000
BORNIN1964×LOGWAGE	35.051 (19.618) <sup>+</sup>	109.870 (28.233)**	2.286	2.223	2.286
BORNIN1964×LOGWAGE×AFQT	0.209 (0.734)	1.904 (0.953)*	0.049	-0.135	0.049
BORNIN1964×AFQT <sup>2</sup>	-0.052 (0.218)	-0.185 (0.130)	0.289	0.271	0.289
BORNIN1964×LOGWAGE <sup>2</sup>	-2.592 (1.417)	-8.115 (2.064)**	15.893	15.220	15.893
CONSTANT/SUM OF WEIGHTS	24.079 (24.145)	12.077 (19.099)	1.000	0.983	1.000
<i>N</i>	1,371	1,371	1,371	144	144

NOTES: Panel A of the table reports conditional maximum likelihood and IPT propensity score parameter estimates. The 1979 baseline sampling weights are used when computing both estimates. A ‘\*\*\*’, ‘\*’ and ‘+’ denotes that a point estimate is significantly different from zero at the 1, 5 and 10 percent levels. Standard errors (in parentheses) allow for arbitrary patterns of dependence across units residing in the same household at baseline. Column 1 of Panel B gives the full sample unweighted mean (i.e., only sampling weights used) of each of the propensity score regressors. Columns 2 and 3 give the corresponding complete case inverse probability weighted means using, respectively, the CML and IPT propensity score estimates.