



EMPLOYING MACHINE LEARNING TO HANDLE MISSING DATA

PRESENTED BY: ANINDRO BHATTACHARYA / MARCH 23, 2022

AGENDA

- X Mechanisms of Missingness
- X NeuMiss Networks
- X Doubly Robust Estimators
- X Next Steps





MECHANISMS OF MISSINGNESS

NOTATION

- X V : set of observable variables
- X V_o : set of variables that are observed in all records in the population
- X V_m : set of variables missing in at least one record
- X U : set of unobserved variables (latent variables)
- X R : missingness pattern



NOTATIONS IN ACTION

U

F_1	F_2	F_3	F_4
1	2	3	NA
423	2	NA	10
34	32	42	NA

V_o

V_m

V



MISSING COMPLETELY AT RANDOM (MCAR)

- X Probability that V_m is missing is independent of V_m or any other variable in the study
- X e.g., deciding to reveal income levels based on coin flips



MISSING AT RANDOM (MAR)

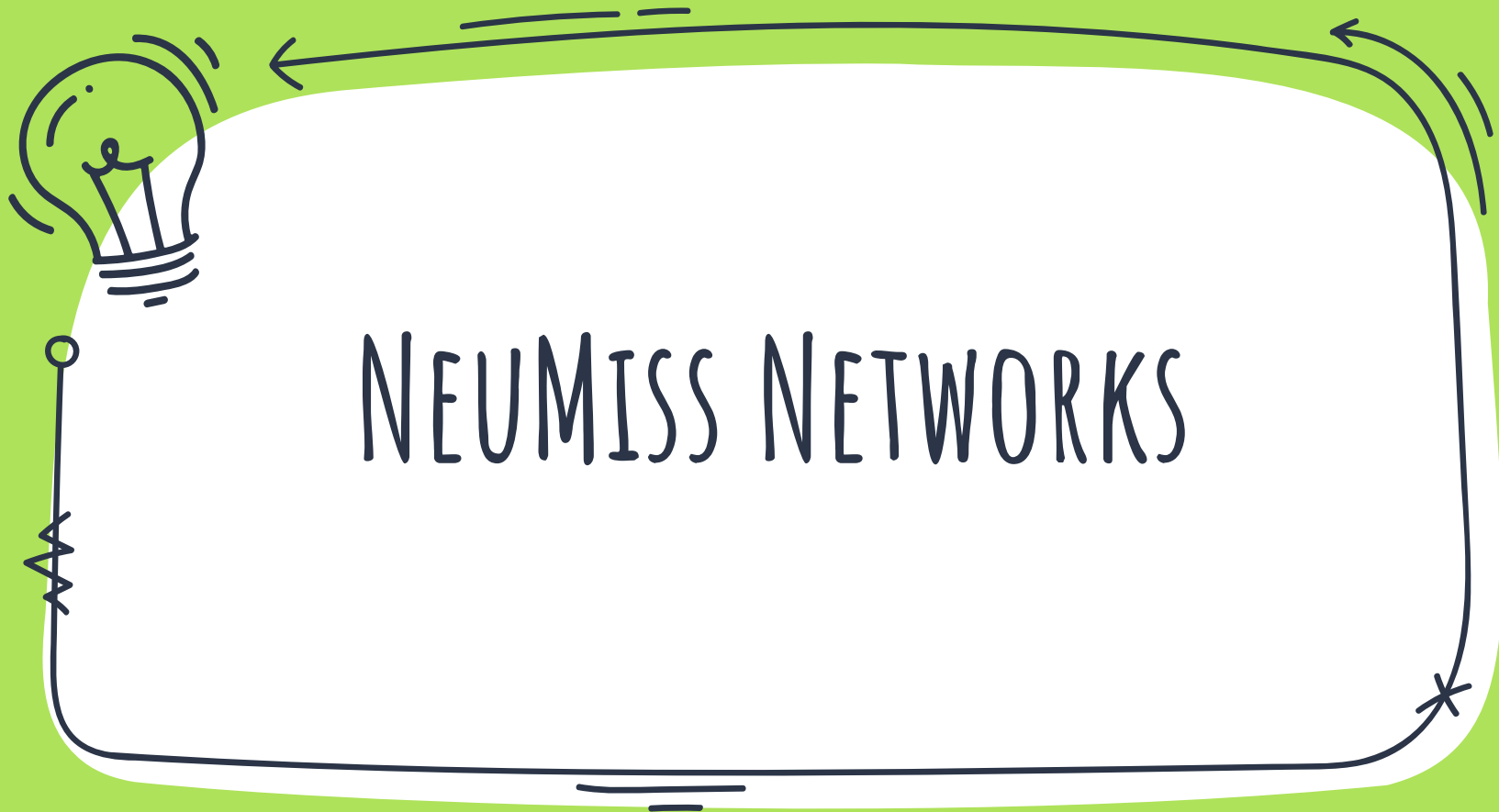
- X For all data cases Y , $P(R|Y_{obs}, Y_{mis}) = P(R|Y_{obs})$
 - X Y_{obs} : observed component of Y
 - X Y_{mis} : observed component of Y
- X e.g., patients under the age of 65 are less likely to have bone density tests taken



MISSING NOT AT RANDOM (MNAR)

- X Neither MCAR or MAR
- X Missingness is related to factors not measured in the study
 - X Probability that V_m is missing is dependent on some $u \in U$
- X e.g., smoking status is not recorded in patients admitted as an emergency





NOTATION

X $X \in \mathbb{R}^d$: set of features

X $Y \in \mathbb{R}$: outcomes

X $M \in \{0,1\}^d$

X $\forall 1 \leq j \leq d, M_j = 1 \Leftrightarrow X_j$ is not observed

X For realizations m of M ,

X $obs(m)$: indices of zero entries of m

X $mis(m)$: indices of non-zero entries of m



NOTATION - EXAMPLE

X Suppose:

X $x = (1.1, 2.3, -3.1, 8, 5.27)$

X $m = (0, 1, 0, 0, 1)$

X then:

X $\tilde{x} = (1.1, \text{NA}, -3.1, 8, \text{NA})$

X $\text{obs}(m) = \{0, 2, 3\}$

X $\text{mis}(m) = \{1, 4\}$

X $x_{\text{obs}(m)} = (1.1, -3.1, 8)$

X $x_{\text{mis}(m)} = (2.3, 5.27)$



BAYES PREDICTOR

$$f^*(X_{obs(M)}, M) = \mathbb{E} [Y | X_{obs(M)}, M]$$

Assumption 1 (Gaussian data). *The distribution of X is Gaussian, that is, $X \sim \mathcal{N}(\mu, \Sigma)$.*

Assumption 2 (MCAR mechanism). *For all $m \in \{0, 1\}^d$, $P(M = m | X) = P(M = m)$.*

Assumption 3 (MAR mechanism). *For all $m \in \{0, 1\}^d$, $P(M = m | X) = P(M = m | X_{obs(m)})$.*

Proposition 2.1 (MAR Bayes predictor). *Assume that the data are generated via the linear model defined in equation (1) and satisfy Assumption 1. Additionally, assume that either Assumption 2 or Assumption 3 holds. Then the Bayes predictor f^* takes the form*

$$f^*(X_{obs}, M) = \beta_0^* + \langle \beta_{obs}^*, X_{obs} \rangle + \langle \beta_{mis}^*, \mu_{mis} + \Sigma_{mis, obs} (\Sigma_{obs})^{-1} (X_{obs} - \mu_{obs}) \rangle, \quad (4)$$

where we use obs (resp. mis) instead of $obs(M)$ (resp. $mis(M)$) for lighter notations.



DERIVING THE BAYES PREDICTOR

$$\begin{aligned} f_{\tilde{X}}^*(\tilde{X}) &= \mathbb{E}[Y|\tilde{X}] \\ &= \mathbb{E}[\beta_0^* + \langle \beta^*, X \rangle \mid M, X_{obs(M)}], \text{ by linear model} \\ &= \beta_0^* + \langle \beta_{obs(M)}^*, X_{obs(M)} \rangle + \langle \beta_{mis(M)}^*, \mathbb{E}[X_{mis(M)} \mid M, X_{obs(M)}] \rangle. \end{aligned}$$

If missingness is MAR or MCAR,

$$\mathbb{E}[X_{mis(M)} \mid M, X_{obs(M)}] = \mathbb{E}[X_{mis(M)} \mid X_{obs(M)}].$$

Since $X \sim N(\mu, \Sigma)$,

$$\mathbb{E}[X_{mis(m)} \mid X_{obs(m)}] = \mu_{mis(m)} + \Sigma_{mis(m), obs(m)} (\Sigma_{obs(m)})^{-1} (X_{obs(m)} - \mu_{obs(m)})$$

Therefore,

$$f^*(X_{obs}, M) = \beta_0^* + \langle \beta_{obs}^*, X_{obs} \rangle + \langle \beta_{mis}^*, \mu_{mis} + \Sigma_{mis, obs} (\Sigma_{obs})^{-1} (X_{obs} - \mu_{obs}) \rangle$$



THE NEUMANN SERIES

$$(I - S)^{-1} = \sum_{j=0}^{\infty} S^j = I + S + S^2 + \dots$$



EXAMPLE WITH NEUMANN SERIES

Estimate inverse of $T = \begin{bmatrix} 1 & 0.5 \\ 0.25 & 1 \end{bmatrix}$. Actual value of $T^{-1} = \begin{bmatrix} 1.14285714 & -0.57142857 \\ -0.28571429 & 1.14285714 \end{bmatrix}$.

1st order approximation of $T^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$.

3rd order approximation of $T^{-1} = \begin{bmatrix} 1.125 & -0.5 \\ -0.25 & 1.125 \end{bmatrix}$.

10th order approximation of $T^{-1} = \begin{bmatrix} 1.14282227 & -0.57141113 \\ -0.28570557 & 1.14282227 \end{bmatrix}$.

20th order approximation of $T^{-1} = \begin{bmatrix} 1.14285714 & -0.57142857 \\ -0.28571429 & 1.14285714 \end{bmatrix}$.



APPROXIMATING INVERSE COVARIANCES WITH NEUMANN SERIES

- X $S^{(0)}$: starting point for approximation ($d \times d$)
- X $S_{obs(m)}^{(0)}$: submatrix of $S^{(0)}$ obtained by selecting components for which $m = 0$
- X Order-0 approximation of $(\Sigma_{obs(m)})^{-1}$



APPROXIMATING INVERSE COVARIANCES WITH NEUMANN SERIES

X For all $m \in \{0, 1\}^d$, define the order- l approximation $S_{obs(m)}^{(l)}$ of $(\Sigma_{obs(m)})^{-1}$ via:

$$S_{obs(m)}^{(\ell)} = (Id - \Sigma_{obs(m)}) S_{obs(m)}^{(\ell-1)} + Id$$



ORDER- L APPROXIMATION OF BAYES PREDICTOR

$$f^*(X_{obs}, M) = \beta_0^* + \langle \beta_{obs}^*, X_{obs} \rangle + \langle \beta_{mis}^*, \mu_{mis} + \Sigma_{mis, obs} (\Sigma_{obs})^{-1} (X_{obs} - \mu_{obs}) \rangle$$

\approx

$$f_\ell^*(X_{obs}, M) = \langle \beta_{obs}^*, X_{obs} \rangle + \langle \beta_{mis}^*, \mu_{mis} + \Sigma_{mis, obs} S_{obs(m)}^{(\ell)} (X_{obs} - \mu_{obs}) \rangle$$



NEUMISS ARCHITECTURE

- X Approximates the Bayes predictor
- X Inverses $(\Sigma_{obs(m)})^{-1}$ are computed using an unrolled version of the iterative algorithm



NEUMISS ARCHITECTURE

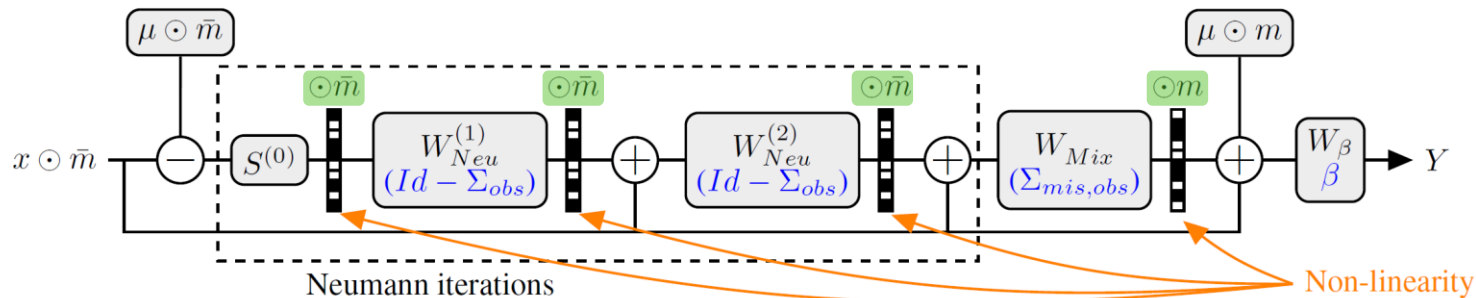


Figure 1: **NeuMiss network architecture with a depth of 4** — $\bar{m} = 1 - m$. Each weight matrix $W^{(k)}$ corresponds to a simple transformation of the covariance matrix indicated in blue.

RESULTS FROM NEUMISS

- X Generated synthetic data according to multivariate Gaussian distribution
 - X $\Sigma = UU^T + \text{diag}(\epsilon)$
 - X $U \in \mathbb{R}^{d \times d/2}$ and entries of U drawn from $N(0, 1)$
 - X ϵ : vector of entries drawn uniformly in $[0.01, 0.1]$



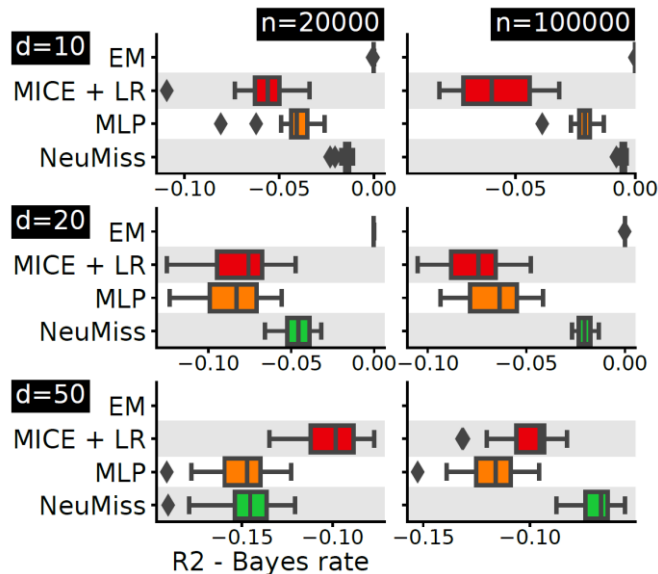
RESULTS FROM NEUMISS

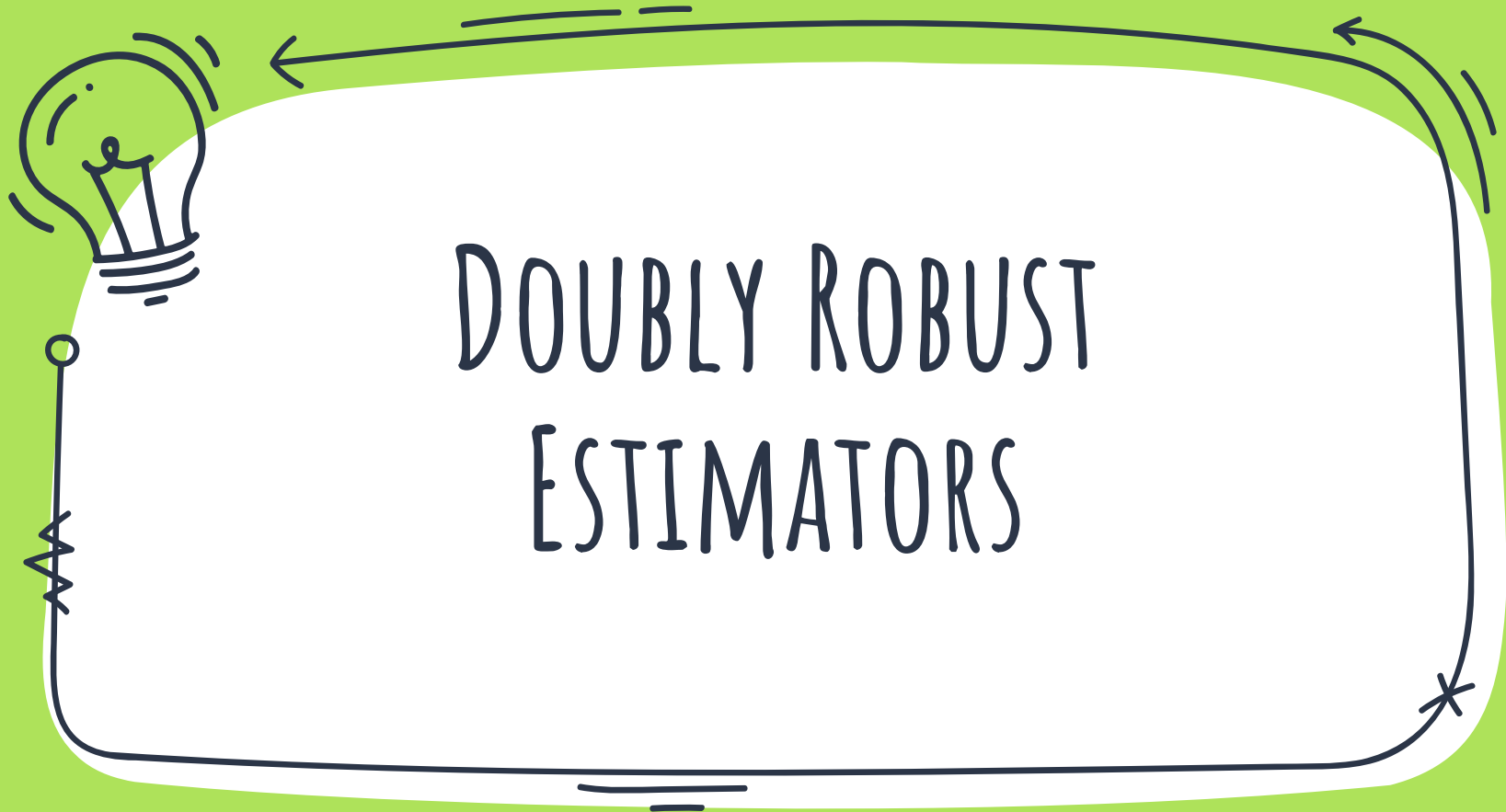
- X Compared performance of NeuMiss to 3 other algorithms
 - X **EM:** Expectation-Maximization algorithm run to estimate parameters of $P(X, Y)$ with missing values
 - Predictions made using $E[X|Y]$
 - X **MICE+LR:** data is imputed using conditional imputation (scikit-learn IterativeImputer), LR fit on imputed data
 - X **MLP:** multilayer perceptron with 1 hidden layer followed by ReLU nonlinearity; data imputed with 0



RESULTS FROM NEUMISS

- X EM gives best results when it can be run
 - X Cannot be run when $d \geq 50$
- X NeuMiss is second best
 - X Except when $d=50$, $n=20\ 000$
- X NeuMiss works best for high $\frac{\text{\# of samples}}{\text{\# of parameters}}$ ratio





DOUBLY ROBUST ESTIMATORS

DOUBLY ROBUST (DR) ESTIMATORS

- X Remain consistent when either:
 - X A model for the treatment assignment mechanism is correctly specified
 - X A model for counterfactual data is correctly specified



NOTATION

X **Full data:** $\mathbf{L} = (\mathbf{V}', \mathbf{Y})'$

X \mathbf{V} : always observed baseline variables

X \mathbf{Y} : scalar outcome which is missing on some subjects

X **Observed data:** $\mathbf{O} = (\Delta, \mathbf{L}_{\text{obs}})$

X $\mathbf{L}_{\text{obs}} = \mathbf{L}$ when $\Delta = 1$

X $\mathbf{L}_{\text{obs}} = \mathbf{V}$ when $\Delta = 0$



GOAL

X Estimate unconditional mean μ of Y based on n i.i.d. copies of \mathbf{O}_i , where $i = 1, \dots, n$

X Assumptions:

X Y is MAR

X $P(\Delta = 1 \mid Y, \mathbf{V}) = P(\Delta = 1 \mid \mathbf{V}) \equiv \pi(\mathbf{V}) > 0$

X $\mu = E(Y) = E\{E(Y \mid \mathbf{V})\}$

Propensity Score



MEAN OF Y IN TERMS OF OBSERVED DATA DISTRIBUTION

$$\begin{aligned}\mu &= E(Y) \\ &= E\{E(Y \mid V)\} \\ &= E\{E(Y \mid \Delta = 1, V)\} \star \\ &= E\left(\frac{\Delta Y}{\pi(V)}\right) \star\end{aligned}$$



METHOD 1



1. Fit model for PS $\pi(\mathbf{V})$ based on parametric model $\pi(\mathbf{V}; \boldsymbol{\alpha})$
2. Estimate μ with Horvitz-Thompson (HT) estimator

$$\hat{\mu}_{HT} = n^{-1} \sum_i \frac{\Delta_i Y_i}{\pi(\mathbf{V}_i; \hat{\boldsymbol{\alpha}})}, \text{ where } \hat{\boldsymbol{\alpha}} \text{ is the MLE of } \boldsymbol{\alpha}$$



METHOD 2



1. Fit model for $\Psi\{s(\mathbf{V}; \boldsymbol{\beta})\}$ for $E(Y \mid \Delta = 1, \mathbf{V})$
 - ✗ Ψ^{-1} : known link function
 - ✗ $s(\mathbf{V}; \boldsymbol{\beta})$: known regression function ($\boldsymbol{\beta}$ is an unknown finite-dimensional parameter)
2. Estimate μ by OR estimator

$$\hat{\mu}_{OR} = n^{-1} \sum_i \Psi\{s(\mathbf{V}_i; \tilde{\boldsymbol{\beta}})\}$$



DR ESTIMATOR (A COMBINATION OF METHODS 1 AND 2)

X Model $E(Y | \Delta = 1, V)$ as $e(V; \beta, \phi) = \Psi\{s(V; \beta) + \phi\pi^{-1}(V; \hat{\alpha})\}$

X $\hat{\mu}_{dr} = n^{-1} \sum_i \Psi\{s(V_i; \hat{\beta}) + \phi\pi^{-1}(V_i; \hat{\alpha})\}$

X $\phi = \Delta_i(Y_i - s(V_i; \hat{\beta}))$



SIMULATION STUDY RESULTS

The naively calculated mean on the complete data is 0.03116.

Estimator	$E(\hat{Y})$	Bias
Horvitz-Thompson	0.07259	-0.04143
Outcome Regression	0.03214	-0.00098
Doubly Robust	0.03043	0.00072

Table 5: Estimates of $E(Y)$ with correctly specified π and s models

Models are incorrectly specified by fitting them on the original dataset masked with MCAR missingness and missingness rate of 0.8.

Estimator	$E(\hat{Y})$	Bias
Horvitz-Thompson	0.05721	-0.02605
Doubly Robust	0.03043	0.00072

Table 6: Estimates of $E(Y)$ with incorrectly specified π and correctly specified s models

Estimator	$E(\hat{Y})$	Bias
Outcome Regression	0.00970	0.02145
Doubly Robust	0.06804	-0.01688

Table 7: Estimates of $E(Y)$ with correctly specified π and incorrectly specified s models

Estimator	$E(\hat{Y})$	Bias
Doubly Robust	0.06804	-0.01688

Table 8: Estimates of $E(Y)$ with incorrectly specified π and incorrectly specified s models



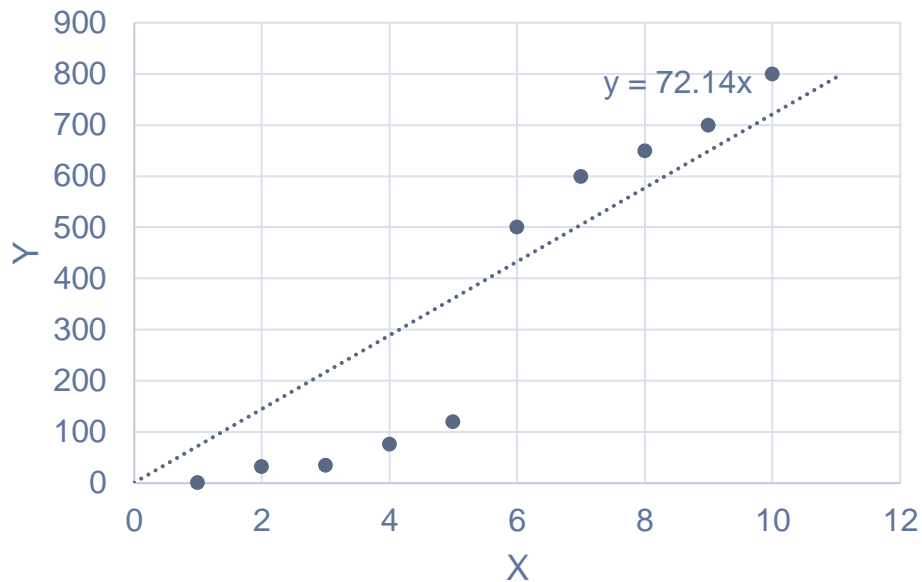


COMBINING NEUMISS AND DR ESTIMATORS

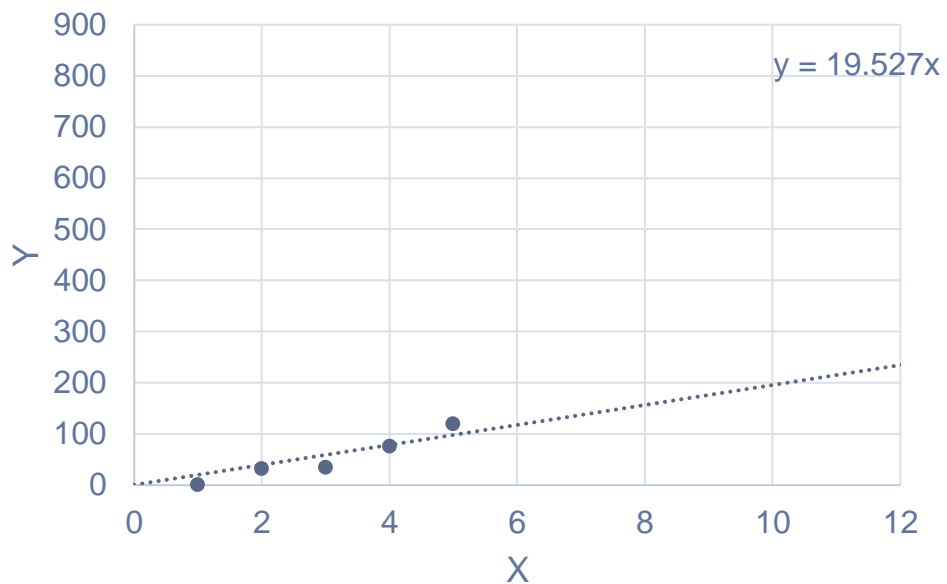
- X Use propensity score model to decide how likely a subject's outcome is to be missing
- X Use propensity score model to modify Bayes' predictor and over-compensate for subjects that are more likely to be missing
- X Use updated Bayes' predictor to impute data



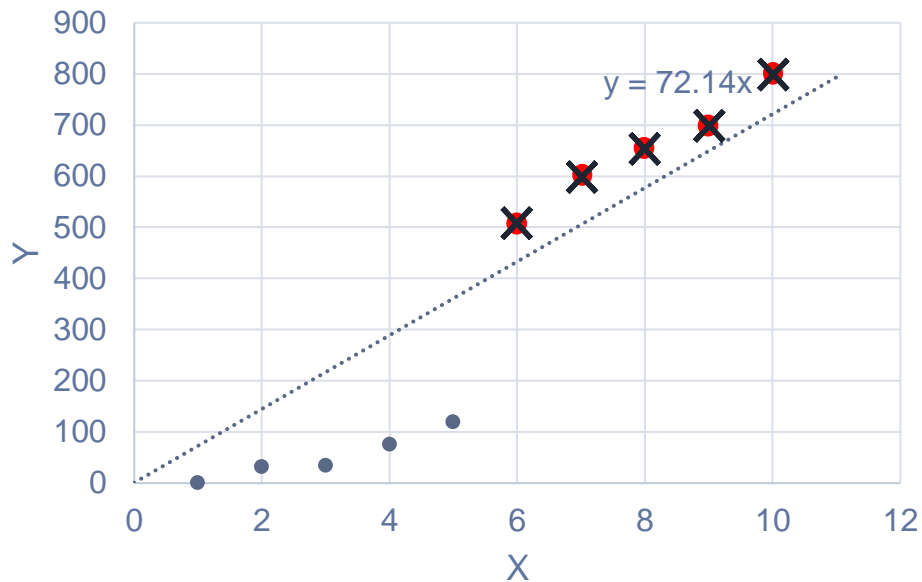
VISUAL EXAMPLE



VISUAL EXAMPLE



VISUAL EXAMPLE





ANY SUGGESTIONS?