

# **A method for comparing multiple imputation techniques: a case study on the U.S. National COVID Cohort Collaborative**

Elena Casiraghi<sup>1,2,3\*</sup>, Rachel Wong<sup>4,\*</sup>, Margaret Hall<sup>4</sup>, Ben Coleman<sup>5,6</sup>, Marco Notaro<sup>1,2</sup>, Michael D. Evans<sup>7</sup>, Jena S. Tronieri<sup>8</sup>, Hannah Blau<sup>5</sup>, Bryan Laraway<sup>9</sup>, Tiffany J. Callahan<sup>9</sup>, Lauren E. Chan<sup>10</sup>, Carolyn T. Bramante<sup>11</sup>, John B. Buse<sup>12,13</sup>, Richard A. Moffitt<sup>4</sup>, Til Stürmer<sup>14</sup>, Steven G. Johnson<sup>15</sup>, Yu Raymond Shao<sup>16,17</sup>, Justin Reese<sup>3</sup>, Peter N. Robinson<sup>5,6</sup>, Alberto Paccanaro<sup>18,19</sup>, Giorgio Valentini<sup>1,2</sup>, Jared D. Huling<sup>20,\*\*</sup> and Kenneth J. Wilkins<sup>21,\*\*</sup> on behalf of the N3C Consortium

<sup>1</sup> AnacletoLab, Department of Computer Science “Giovanni degli Antoni”, Università degli Studi di Milano, Milan, ITALY

<sup>2</sup> CINI, Infolife National Laboratory, Roma, ITALY

<sup>3</sup> Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA

<sup>4</sup> Department of Biomedical Informatics, Stony Brook University, Stony Brook, NY, USA

<sup>5</sup> The Jackson Laboratory for Genomic Medicine, Farmington, USA

<sup>6</sup> Institute for Systems Genomics, University of Connecticut, Farmington, CT, USA

<sup>7</sup> Biostatistical Design and Analysis Center, Clinical and Translational Science Institute, University of Minnesota, Minneapolis, MN, USA

<sup>8</sup> Department of Psychiatry, Perelman School of Medicine at the University of Pennsylvania, Philadelphia, PA, USA

<sup>9</sup> University of Colorado, Anschutz Medical Campus, Aurora, CO, USA

<sup>10</sup> College of Public Health and Human Sciences, Oregon State University, Corvallis, USA

<sup>11</sup> Division of General Internal Medicine, University of Minnesota, Minneapolis, MN, USA

<sup>12</sup> NC Translational and Clinical Sciences Institute, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA

<sup>13</sup> Division of Endocrinology, Department of Medicine, University of North Carolina School of Medicine, USA

<sup>14</sup> Department of Epidemiology, Gillings School of Global Public Health, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA

<sup>15</sup> Institute for Health Informatics, University of Minnesota, Minneapolis, MN, USA

<sup>16</sup> Harvard-MIT Division of Health Sciences and Technology (HST), 260 Longwood Ave, Boston, USA

<sup>17</sup> Department of Radiation Oncology, UT Southwestern Medical Center, Dallas, USA

<sup>18</sup> School of Applied Mathematics (EMAp), Fundação Getúlio Vargas, Rio de Janeiro, BRAZIL

<sup>19</sup> Department of Computer Science, Royal Holloway, University of London, Egham, UK

<sup>20</sup> Division of Biostatistics, School of Public Health, University of Minnesota, Minneapolis, MN, USA

<sup>21</sup> Biostatistics Program, Office of the Director, National Institute of Diabetes and Digestive and Kidney Diseases, National Institutes of Health, Bethesda, MD, USA

N3C Consortium: Tell Bennet, Christopher Chute, Peter DeWitt, Kenneth Gersing, Andrew Girvin, Melissa Haendel, Jeremy Harper, Janos Hajagos, Stephanie Hong, Emily Pfaff, Jane Reusch, Corneliu Antoniescu, Kimberly Robaski

Corresponding Authors (\*): Elena Casiraghi, Rachel Wong

Co-senior Authorship (\*\*): Jared D Huling and Kenneth J. Wilkins equally contributed to the work

## Supplementary files

### **S1.xlsx:**

This Excel file contains all the results we obtained when using a number of imputation  $m = 42$ .

The file is composed by:

- sheet “mean\_measures\_m42” containing the colored table (Figure 5) showing and detailing the average measures obtained by the tested imputation algorithms across the three outcomes.
- sheets “RB\_mean” (see also [Supplementary Figure S1](#)), “MSE\_mean” (see also [Supplementary Figure S2](#)), “ER\_mean” (see also [Supplementary Figure S3](#)), and “CR\_mean” contain the four win-tie-loss tables (for the RB measure, the MSE measure, the ER measure, the CR measure) obtained by summing the wins, ties, losses obtained by each model over the three outcome variables.

On the right, each of the four sheets contains the mean of the win-tie-loss tables over the three outcomes, where the wins, ties, and losses are computed by comparing the models on the rows to the models on the column by a paired-sided paired rank sign test .

The grid shows numbers in the range  $[-3, +3]$ ; they are computed by representing each win by a +1 value, each tie as a 0 value, each loss as a -1 value.

### **S2.xlsx:**

This Excel file has the same structure of [S1.xlsx](#); it details all the results we obtained when using a number of imputation  $m = 5$ .

### **S3\_MCAR.xlsx:**

This Excel file has the same structure of [S1.xlsx](#) and [S2.xlsx](#); it details all the results we obtained when simulating MCAR missingness in the amputated datasets.

### **S4\_MNAR.xlsx:**

This Excel file has the same structure of [S1.xlsx](#) and [S2.xlsx](#) and [S3\\_MCAR.xlsx](#); it details all the results we obtained when simulating MNAR missingness in the amputated datasets.

## Supplementary Figures:

ML algorithm	univariate imputation method	use outcomes	one-hot encode binned numeric predictors	one-hot encode categorical predictors	univariate imputation order	pmm donors	average of absolute values of RB across outcomes	wins	ties	losses		
amelia		F	T				0.029	1	4	-94		
		T	F	T			0.017	7	16	-50		
		T	F				0.021	8	26	-28		
mice	default	F			monotone		0.012	6	16	-63		
		T	F	F	revmonotone		0.012	5	18	-60		
		F			monotone		0.012	7	14	-63		
		T	F		revmonotone		0.012	7	15	-64		
	logreg	F			monotone		0.013	5	6	-85		
		T	T	T	revmonotone		0.013	5	3	-89		
		F			monotone		0.011	18	21	-34		
		T	F		revmonotone		0.011	22	18	-31		
	norm	F	T			monotone		0.012	7	8	-76	
			T	F		revmonotone		0.012	6	5	-80	
			F			monotone		0.012	10	11	-63	
			T	F		revmonotone		0.013	9	10	-66	
T		F			monotone		0.011	34	16	-23		
		F	T	T	revmonotone		0.011	35	17	-23		
		F			monotone		0.011	40	18	-16		
		F			revmonotone		0.011	41	17	-15		
missRanger	extratrees	F			monotone		0.007	109	4	0		
					monotone	3	0.008	40	8	-38		
					monotone	5	0.009	19	16	-49		
					revmonotone		0.007	107	4	0		
					revmonotone	3	0.008	51	8	-32		
					revmonotone	5	0.009	22	13	-49		
					monotone		0.007	113	4	-2		
					monotone	3	0.010	26	17	-30		
					monotone	5	0.010	12	20	-38		
					revmonotone		0.007	114	4	-2		
		T					monotone		0.010	27	15	-29
							revmonotone	3	0.010	10	18	-45
							monotone		0.006	83	16	0
							monotone	3	0.006	81	10	-9
							monotone	5	0.007	86	14	-14
							revmonotone		0.006	83	16	0
							revmonotone	3	0.006	84	8	-8
							revmonotone	5	0.007	59	14	-15
							monotone		0.006	51	13	-9
							monotone	3	0.008	47	18	-12
F					monotone		0.009	36	18	-21		
					revmonotone		0.006	50	14	-8		
					revmonotone	3	0.008	46	18	-13		
					revmonotone	5	0.008	34	18	-25		
					monotone		0.001	31	24	-6		
					monotone		0.001	32	24	-6		
					monotone		0.001	28	27	-6		
					monotone		0.001	28	27	-6		
IPW	logreg	F	T				0.002	5	11	-48		
		T	F				0.004	3	4	-60		
	RF	F	T				0.013	0	3	-79		
		T	F				0.017	0	2	-82		

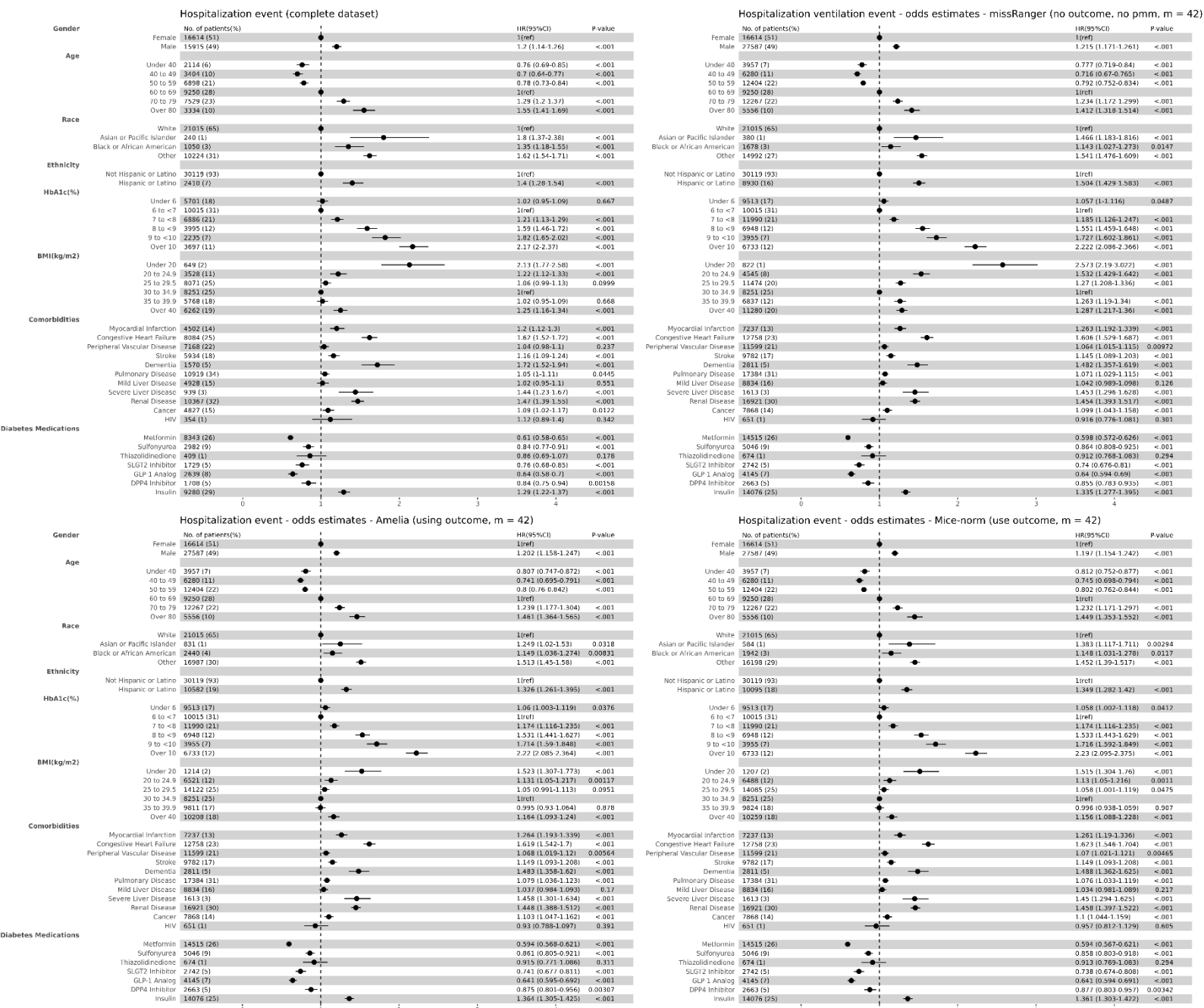
Figure S1: Column “average absolute value of RB across predictors and outcomes” reports the average RB measure across the hospitalization, invasive ventilation, and patients’ survival outcomes (the table is also made available in Supplementary file S1 – sheet “RB\_mean”). Columns “wins”, “ties”, “losses” report the sum of, respectively, wins, ties, and losses computed by comparing the (absolute value of the) RB measures over the three outcomes (the corresponding win-tie-loss grid is shown in the Supplementary material). The comparison between two models over an outcome variable is performed with a sided Wilcoxon signed-rank test comparing the distribution of the (absolute) RB values for all the predictor variables. The winner is the model achieving the lowest RB distribution. All the models but missRanger with no pmm and using the outcome variables in the imputation model are obtaining  $RB \leq 0$ , meaning that the computed estimates are systematically lower than those computed on the complete dataset. missRanger with outcome variable in the imputation model and no pmm is instead bringing to the computation of inflated estimates.

ML algorithm	univariate imputation method	use outcomes	one-hot encode binned numeric predictors	one-hot encode categorical predictors	univariate imputation order	pmm donors	average MSE across outcomes	wins	ties	losses	
amelia		F	T	T	monotone		0.007	18	5	-76	
		F	F		revmonotone		0.004	25	28	-22	
		T	T		monotone		0.005	26	23	-22	
		T	F		revmonotone		0.003	86	19	-1	
mice	default	F	F	F	monotone		0.002	26	21	-32	
		F	F		revmonotone		0.002	25	20	-34	
		T	F		monotone		0.002	30	18	-32	
		T	F		revmonotone		0.002	26	19	-34	
	logreg	F	T	T	monotone		0.003	25	9	-58	
		F	T		revmonotone		0.003	25	7	-60	
		T	T		monotone		0.002	37	26	-15	
		T	T		revmonotone		0.002	39	25	-13	
	norm	F	T	T	monotone		0.002	28	13	-43	
		F	F		revmonotone		0.002	27	9	-48	
		F	F		monotone		0.002	30	16	-33	
		F	F		revmonotone		0.002	28	17	-34	
		T	F	monotone		0.002	54	15	-11		
				revmonotone		0.002	53	16	-9		
				monotone		0.002	57	12	-5		
				revmonotone		0.002	72	14	-2		
missRanger	extratrees	F	T	monotone		0.001	93	2	0		
				monotone	3	0	0.001	38	22	-22	
				monotone	5	0	0.001	28	21	-34	
				revmonotone		0.001	92	3	0		
				revmonotone	3	0	0.001	48	20	-18	
				revmonotone	5	0	0.001	31	21	-29	
				monotone		0.001	88	6	-2		
				monotone	3	0	0.002	34	27	-13	
				monotone	5	0	0.002	30	22	-20	
				revmonotone		0.001	89	6	-2		
			T	F	monotone		0.002	43	25	-8	
					revmonotone	3	0	0.002	43	25	-8
					revmonotone	5	0	0.002	30	23	-24
					monotone		0.001	47	20	-2	
					monotone	3	0	0.001	70	12	-11
					monotone	5	0	0.001	84	13	-12
					revmonotone		0.001	46	21	-2	
					revmonotone	3	0	0.001	74	8	-8
					revmonotone	5	0	0.001	85	10	-17
					monotone		0.002	36	7	-28	
		F	F	monotone		0.002	47	15	-14		
				monotone	5	0	0.002	42	16	-18	
				revmonotone		0.002	36	5	-30		
				revmonotone	3	0	0.002	44	17	-15	
				revmonotone	5	0	0.002	37	18	-20	
				monotone		0.002	37	18	-20		
IPW	logreg	F	T			0.007	13	3	-130		
		F	F			0.007	14	2	-130		
		T	T			0.007	12	2	-132		
	RF	F	F			0.007	12	3	-131		
		F	T			0.008	6	1	-144		
		T	F			0.007	7	1	-143		
						0.007	1	1	-150		
						0.007	0	1	-151		

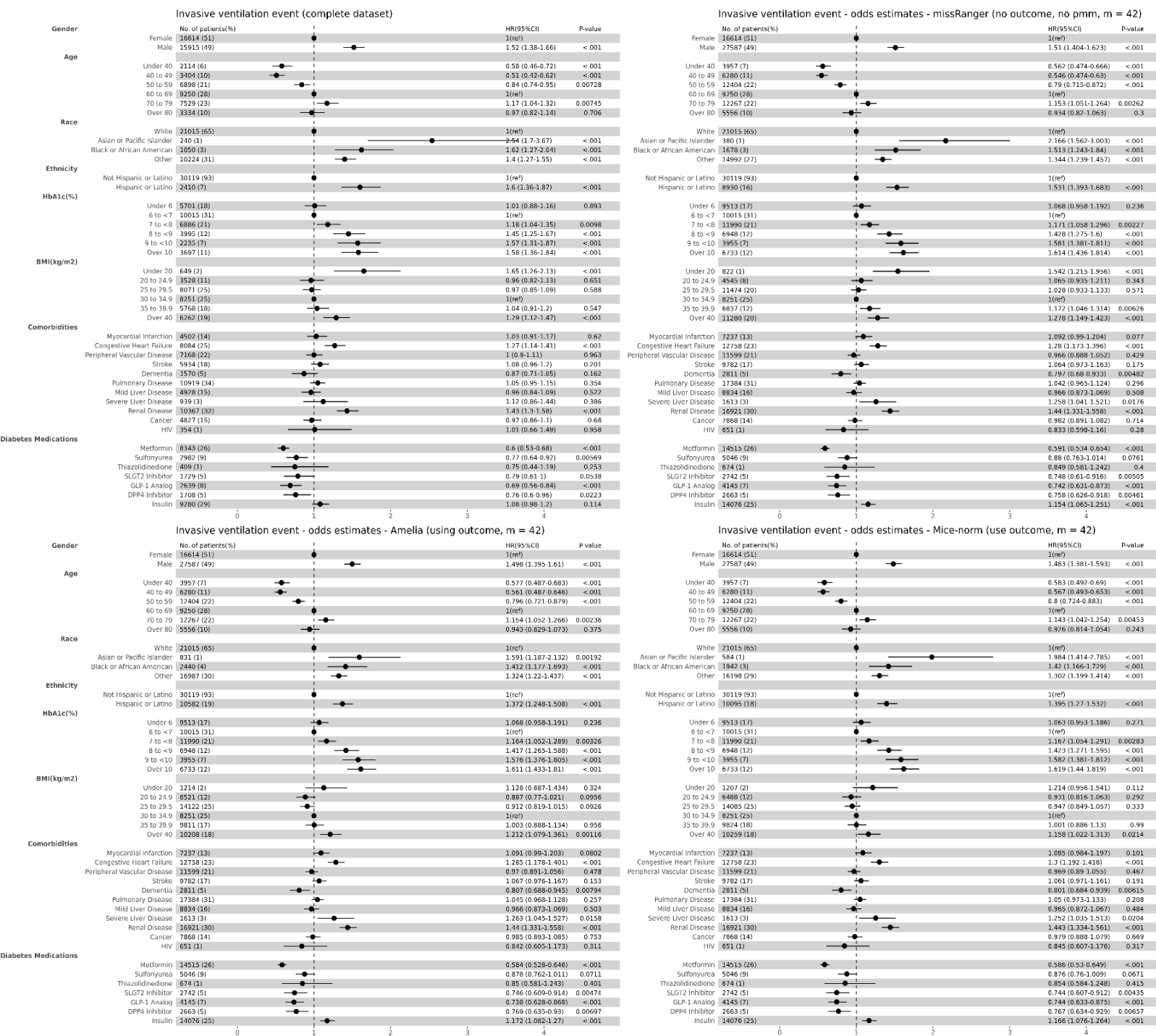
Figure S2: Column “average MSE across predictors and outcomes” reports the average MSE measure across the hospitalization, invasive ventilation, and patients’ survival outcomes (the table is also made available in Supplementary file S1 – sheet “MSE\_mean”). Columns “wins”, “ties”, “losses” report the sum of, respectively, wins, ties, and losses computed by comparing the MSE measures over the three outcomes (the corresponding win-tie-loss grid is shown in the Supplementary material). The comparison between two models over an outcome variable is performed with a sided Wilcoxon signed-rank test comparing the distribution of the MSE values for all the predictor variables. The winner is the model achieving the lowest MSE distribution. The detailed win-tie-loss grids are reported in the Supplementary material.

ML algorithm	univariate imputation method	use outcomes	one-hot encode binned numeric predictors	one-hot encode categorical predictors	univariate imputation order	pmm donors	average ER across outcomes	wins	ties	losses
amelia		F	T	T			0.900	0	5	-84
		F	F				0.949	4	21	-43
		T	T				0.914	6	31	-25
		T	F				0.971	34	23	-13
mice	default	F	F	F	monotone		0.962	6	14	-54
		revmonotone				0.964	6	14	-52	
		monotone				0.963	8	17	-47	
		revmonotone				0.962	9	12	-53	
	logreg	F	T	T	monotone		0.950	5	10	-76
		revmonotone				0.950	5	7	-77	
		monotone				0.952	16	21	-31	
		revmonotone				0.952	17	25	-30	
	norm	F	F	T	monotone		0.952	9	13	-67
		revmonotone				0.950	7	12	-68	
		monotone				0.963	7	16	-53	
		revmonotone				0.961	6	14	-55	
missRanger	extratrees	F	T	T	monotone	3	0.952	110	6	0
					monotone	5	0.958	44	15	-23
					monotone	5	0.956	29	19	-33
					revmonotone	5	0.951	111	6	0
		T	monotone	3	0.958	54	18	-17		
			revmonotone	5	0.955	30	20	-29		
			monotone	3	0.955	103	6	0		
			monotone	5	0.960	24	21	-19		
IPW	RF	F	F	F	monotone	5	0.961	13	24	-24
					revmonotone	3	0.956	103	6	0
					revmonotone	3	0.959	25	21	-19
					revmonotone	5	0.960	11	26	-29
		T	monotone	3	1.032	45	20	0		
			monotone	5	0.971	79	9	-7		
			monotone	5	0.967	85	12	-14		
			revmonotone	3	1.038	42	22	0		
F	revmonotone	3	0.972	82	9	-6				
	revmonotone	5	0.968	80	13	-16				
	monotone	3	1.038	12	27	-15				
	monotone	5	0.979	29	24	-10				
logreg	F	T	F	monotone	3	0.980	24	21	-15	
	revmonotone			3	1.038	14	25	-15		
	monotone			5	0.980	31	21	-12		
	revmonotone			5	0.979	24	21	-19		
amelia		F	T				0.946	18	35	-6
		F	F				0.946	19	34	-6
		T	T				0.950	17	36	-6
		T	F				0.950	19	34	-6
IPW	RF	F	T				0.937	4	8	-52
		F	F				0.946	3	5	-54
		T	T				0.835	0	3	-84
		T	F				0.857	0	3	-89

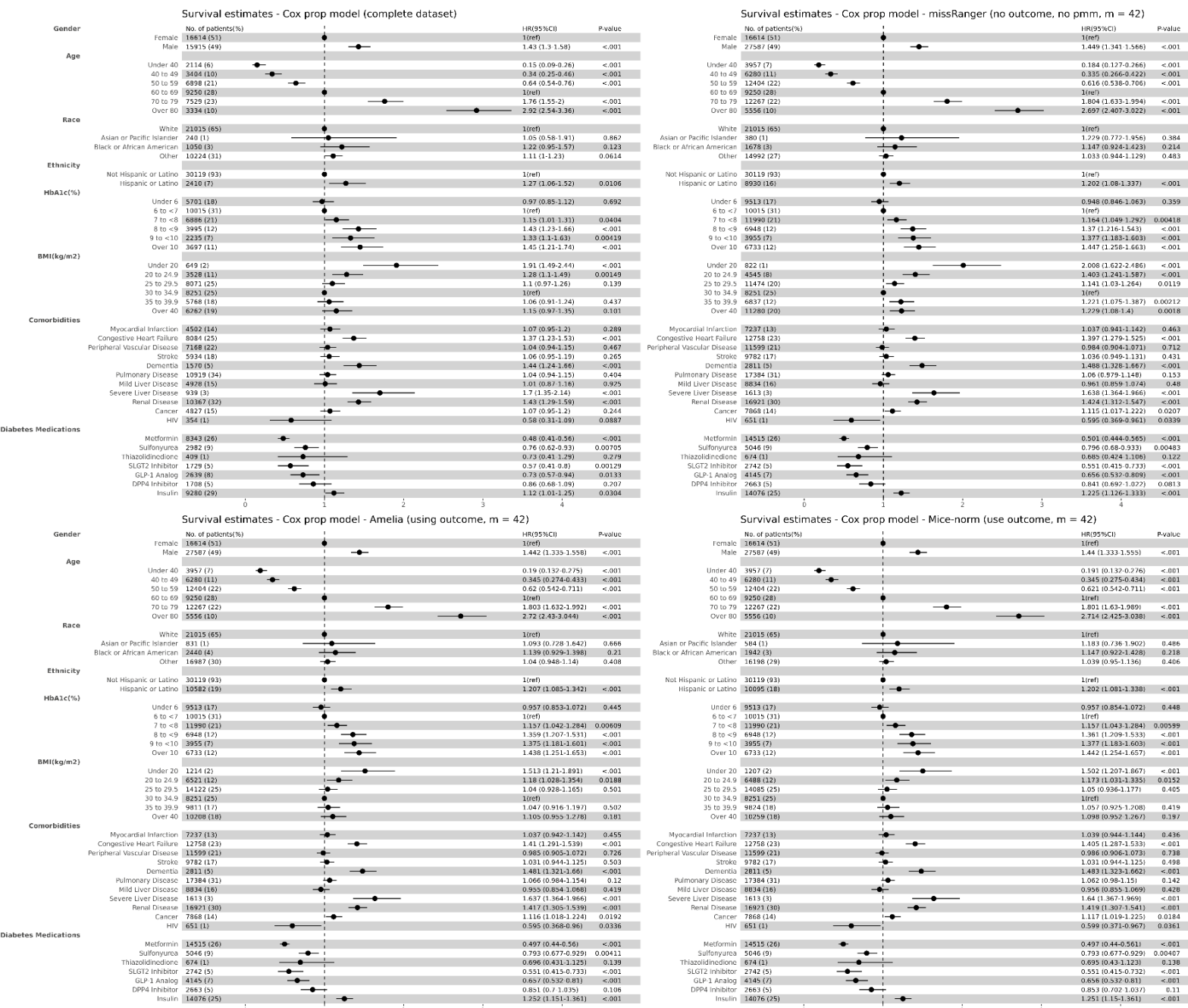
Figure S3: In column “average ER across outcomes” we report the average ER measure across the hospitalization, invasive ventilation, and patients’ survival outcomes (the table is also made available in Supplementary file S1 – sheet “ER\_mean”). Columns “wins”, “ties”, “losses” report the sum of, respectively, wins, ties, and losses computed by comparing the distributions of the ER measures over the three outcomes. Since we would like each  $[ER_i]$  ( $i \in \{1, \dots, d\}$ ) estimate to be as nearest as possible to 1, for the comparison between two models over an outcome variable we used a sided Wilcoxon signed-rank test to compare the following distribution for each model  $f(ER_i) = ||1 - ER_i||$ . The detailed win tie loss grids are reported in the supplementary material.



**Figure S4:** Hospitalization event: estimates obtained on the complete dataset obtained by listwise deletion (top-left) and on the full dataset by the best missRanger (top-right), Amelia (bottom-left), and Mice (bottom-right) models. For missRanger we used no pmm, we did not use the outcome variables in the imputation model, we one-hot encoded categorical predictors and binned numeric predictors (age, BMI, and HbA1c), and we used an univariate imputation order given by the decreasing number of missing values; for Mice-norm we included the outcome variables in the imputation model, we used an univariate imputation order given by the increasing number of missing values; for Amelia we included the outcome variables in the imputation model.



**Figure S5:** Invasive ventilation event: estimates obtained from the complete dataset obtained by listwise deletion (top-left) and on the full dataset by the MI estimation pipelines that using the best missRanger (top-right), Amelia (bottom-left), and Mice (bottom-right) models. For missRanger we used no pmm, we did not use the outcome variables in the imputation model, we one-hot encoded categorical predictors and binned numeric predictors (age, BMI, and HbA1c), and we used an univariate imputation order given by the decreasing number of missing values; for Mice-norm we included the outcome variables in the imputation model, we used an univariate imputation order given by the increasing number of missing values; for Amelia we included the outcome variables in the imputation model.



**Figure S6:** Death event: estimates obtained on the complete dataset obtained by listwise deletion (top-left) and on the full dataset by the MI estimation pipelines that using the best missRanger (top-right), Amelia (bottom-left), and Mice (bottom-right) models. For missRanger we used no pmm, we did not use the outcome variables in the imputation model, we one-hot encoded categorical predictors and binned numeric predictors (age, BMI, and HbA1c), and we used an univariate imputation order given by the decreasing number of missing values; for Mice-norm we included the outcome variables in the imputation model, we used an univariate imputation order given by the increasing number of missing values; for Amelia we included the outcome variables in the imputation model.