

¹ A Sparse Additive Model for High-Dimensional
² Interactions with an Exposure Variable

³ Sahir R Bhatnagar^{1,2}, Tianyuan Lu^{3,4}, Amanda Lovato⁵, David L Olds⁶,
⁴ Michael S Kobor⁷, Michael J Meaney⁸, Kieran O'Donnell⁹, Yi Yang¹⁰, and
⁵ Celia MT Greenwood^{1,3,5}

⁶ ¹Department of Epidemiology, Biostatistics and Occupational Health, McGill
⁷ University

⁸ ²Department of Diagnostic Radiology, McGill University
⁹ ³Quantitative Life Sciences, McGill University

¹⁰ ⁴Lady Davis Institute, Jewish General Hospital, Montréal, QC
¹¹ ⁵Statistics Canada, Ottawa, ON

¹² ⁶Department of Pediatrics, University of Colorado School of Medicine, Denver
¹³ ⁷Department of Medical Genetics, University of British Columbia, BC

¹⁴ ⁸Singapore Institute for Clinical Sciences, Singapore; McGill University
¹⁵ ⁹Department of Psychiatry, McGill University

¹⁶ ¹⁰Department of Mathematics and Statistics, McGill University
¹⁷ ¹¹Departments of Oncology and Human Genetics, McGill University

¹⁸ September 8, 2021

```

## /scratch/bhatnagar-lab/sbhatnagar/git_repositories/sail/manuscript/bin
## +-- ADNI.R
## +-- PRS_bootstrap.R
## +-- PRS_eval_functions.R
## +-- PRS_method_functions.R
## +-- PRS_model_functions.R
## +-- PRS_plot_functions.R
## +-- PRS_plots.R
## +-- README_for_CSDA_Figures.Rnw
## +-- README_for_CSDA_Figures.pdf
## +-- README_for_CSDA_Figures.tex
## +-- cache
## | +-- __packages
## | +-- globals_fb504862a537c3ea460cca06e8d1ad90.RData
## | +-- globals_fb504862a537c3ea460cca06e8d1ad90.rdb
## | +-- globals_fb504862a537c3ea460cca06e8d1ad90.rdx
## | +-- packages-simulation-plots_ace50d8bc55d9041b7b41b136d4b9ce3.RData
## | +-- packages-simulation-plots_ace50d8bc55d9041b7b41b136d4b9ce3.rdb
## | +-- packages-simulation-plots_ace50d8bc55d9041b7b41b136d4b9ce3.rdx
## | +-- packages_571ac4aaa2a2dd7f17919dae0269e8e1.RData
## | +-- packages_571ac4aaa2a2dd7f17919dae0269e8e1.rdb
## | \-- packages_571ac4aaa2a2dd7f17919dae0269e8e1.rdx
## +-- figure
## | +-- PRS-intervention-interaction-1.pdf
## | +-- dzclass-interaction-1.pdf
## | +-- plot-mse-sim-1.pdf
## | +-- toy-effects-1.pdf
## | \-- toy-solution-path-1.pdf
## +-- intro.R
## +-- plot_simulation.R
## +-- setup.R
## +-- simulation.R
## +-- support_bootstrap.R
## \-- support_plots.R

```

¹⁹ **1 Figure 1 - Toy example solution path and effects**

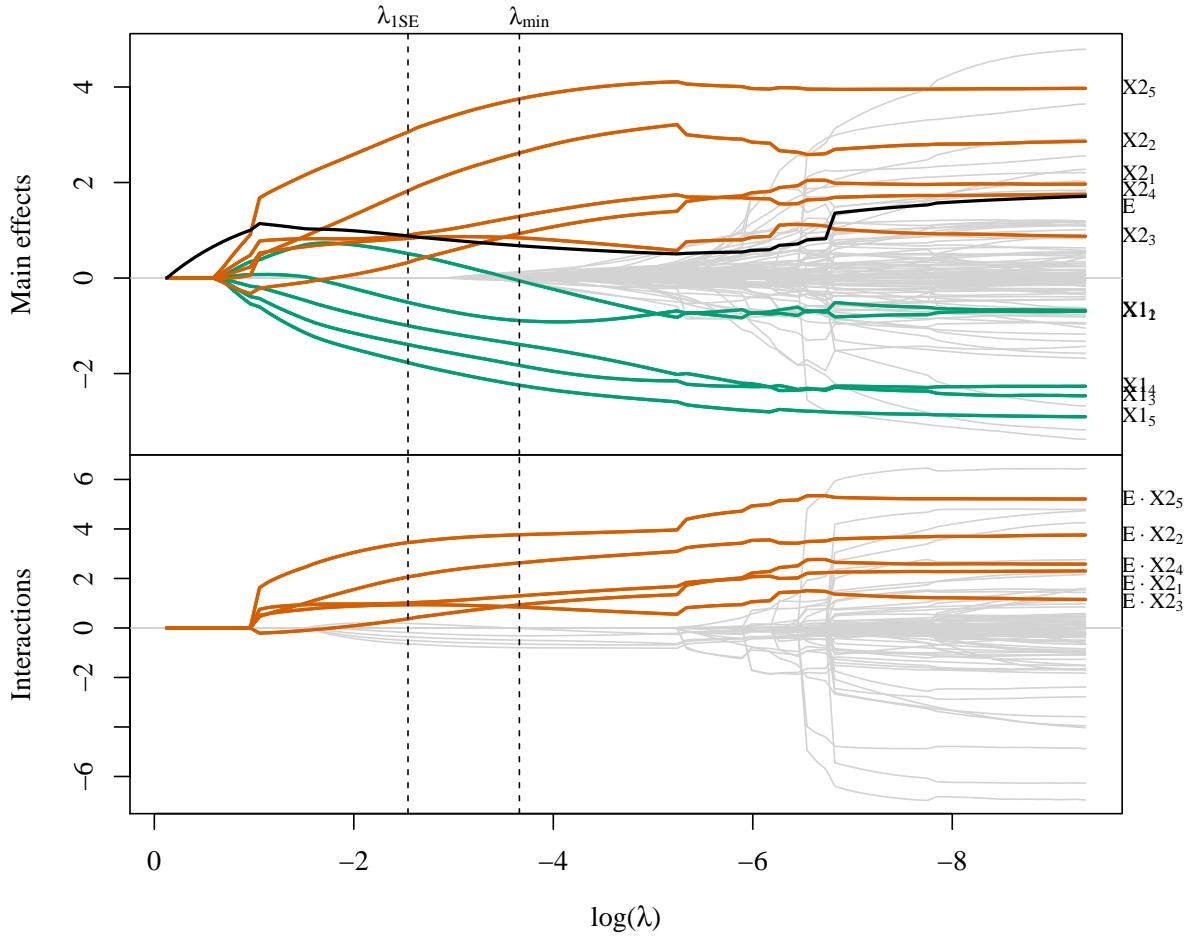


Figure 1: Toy example solution path for main effects (top) and interactions (bottom). $\{X_{11}, X_{12}, X_{13}\}$ and $\{X_{21}, X_{22}, X_{23}\}$ are the three basis coefficients for X_1 and X_2 , respectively. λ_{ISE} is the largest value of penalization for which the CV error is within one standard error of the minimizing value λ_{min} .

20 In Figure 2, we plot the true and estimated component functions $\hat{f}_1(X_1)$ and $E \cdot \hat{f}_2(X_2)$, and
 21 their estimates from this analysis with **sail**. We are able to capture the shape of the correct
 22 functional form, but the means are not well aligned with the data. Lack-of-fit for $f_1(X_1)$
 23 can be partially explained by acknowledging that **sail** is trying to fit a cubic spline to a
 24 linear function. Nevertheless, this example demonstrates that **sail** can still identify trends
 25 reasonably well.

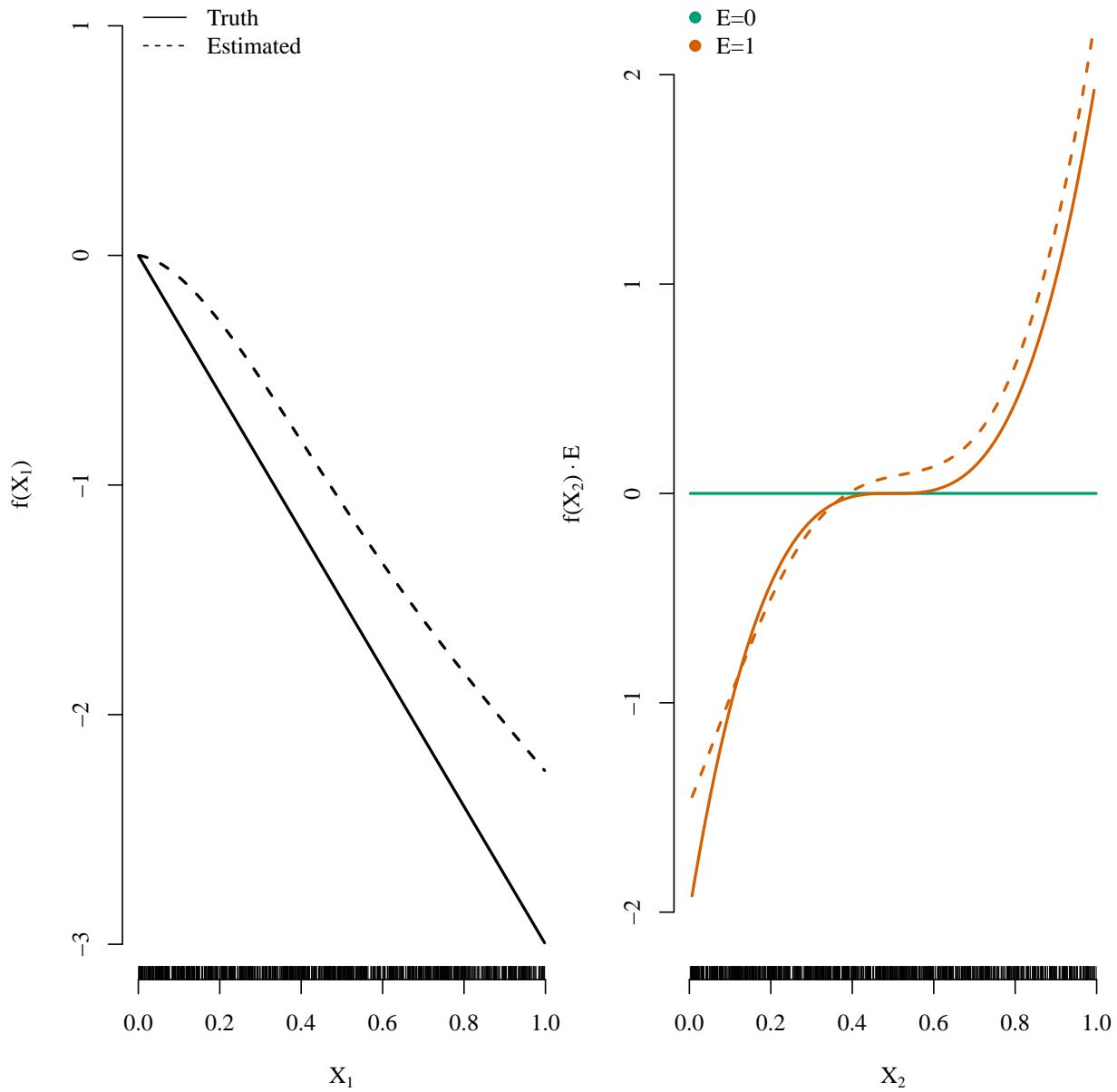


Figure 2: Estimated smooth functions for X_1 and the $X_2 \cdot E$ interaction by the **sail** method based on λ_{min} .

Table 1: Mean (standard deviation) of the number of selected variables ($|\hat{\mathcal{J}}|$), true positive rate (TPR) and false positive rate (FPR) as a percentage from 200 simulations for each of the five scenarios. $|\mathcal{J}|$ is the number of truly associated variables.

		Linear			Non-linear			Non-linear			
		Main Effects		Interactions		Main Effects		Interactions		Interactions	
	lasso	adaptive lasso	lassoBT	GLinternet	HierBasis	SPAM	gamsel	sail	adaptive sail	sail	weak
1a) Strong heredity ($ \mathcal{J} = 7$)											
$ \hat{\mathcal{J}} $	28 (15)	8 (4)	35 (18)	40 (20)	133 (48)	42 (19)	46 (21)	37 (15)	8 (3)	21 (3)	
TPR	53.9 (8.4)	49.3 (10.1)	61.7 (11.5)	66.4 (14.0)	65.2 (8.1)	60.9 (8.5)	56.9 (7.7)	89.5 (8.2)	81.4 (13.0)	82.1 (10.9)	
FPR	1.2 (0.7)	0.2 (0.2)	1.5 (0.9)	1.8 (1.0)	6.5 (2.4)	1.9 (0.9)	2.1 (1.1)	1.5 (0.7)	0.1 (0.1)	0.8 (0.1)	
1b) Weak heredity ($ \mathcal{J} = 5$)											
$ \hat{\mathcal{J}} $	19 (12)	4 (2)	20 (13)	38 (23)	24 (23)	28 (16)	21 (15)	24 (19)	5 (3)	14 (10)	
TPR	40.7 (3.6)	40.1 (1.4)	40.8 (3.8)	64.1 (14.9)	42.2 (6.3)	53.9 (9.4)	42.7 (6.8)	52.4 (11.4)	46.4 (10.1)	55.0 (13.7)	
FPR	0.9 (0.6)	0.1 (0.1)	0.9 (0.7)	1.7 (1.1)	1.1 (1.1)	1.2 (0.8)	1.0 (0.7)	1.0 (0.9)	0.2 (0.1)	0.6 (0.5)	
1c) Interactions Only ($ \mathcal{J} = 2$)											
$ \hat{\mathcal{J}} $	12 (12)	3 (2)	14 (13)	38 (21)	12 (13)	13 (12)	12 (12)	10 (18)	2 (2)	26 (30)	
TPR	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	81.4 (27.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	1.0 (6.9)	0.0 (0.0)	22.9 (36.9)	
FPR	0.6 (0.6)	0.6 (6.9)	0.7 (0.7)	1.8 (1.0)	0.6 (0.7)	0.7 (0.6)	0.6 (0.6)	0.5 (0.9)	0.1 (0.1)	1.3 (1.5)	
2) Linear Effects ($ \mathcal{J} = 7$)											
$ \hat{\mathcal{J}} $	37 (17)	8 (3)	48 (19)	51 (23)	37 (19)	42 (19)	37 (16)	34 (18)	11 (4)	20 (4)	
TPR	70.4 (3.7)	67.2 (6.7)	72.3 (6.3)	93.4 (8.5)	70.3 (3.8)	65.0 (8.1)	70.4 (3.7)	93.9 (9.9)	86.0 (18.5)	68.1 (14.9)	
FPR	1.6 (0.8)	0.2 (0.2)	2.2 (1.0)	2.2 (1.2)	1.6 (0.9)	1.9 (0.9)	1.6 (0.8)	1.4 (0.9)	0.2 (0.2)	0.7 (0.2)	
3) Main Effects Only ($ \mathcal{J} = 5$)											
$ \hat{\mathcal{J}} $	29 (14)	7 (4)	31 (15)	34 (18)	154 (17)	46 (21)	56 (20)	44 (19)	9 (2)	22 (2)	
TPR	75.9 (10.9)	66.5 (15.3)	76.0 (10.9)	77.0 (9.5)	97.5 (6.6)	93.1 (10.7)	81.3 (9.5)	91.5 (10.3)	84.1 (9.2)	85.2 (12.1)	
FPR	1.3 (0.7)	0.2 (0.2)	1.3 (0.8)	1.5 (0.9)	7.5 (0.9)	2.1 (1.0)	2.6 (1.0)	2.0 (0.9)	0.2 (0.1)	0.9 (0.1)	

₂₆ 2 Figure 2 - Test set MSE

₂₇ 3 Table 1 - TPR, FPR, Nactive

₂₈ 4 Simulation scenario 1 Plots for sail

₂₉ 4.1 Main effects

```
## Error: Could not load simulation. Check that 'dir' is a directory with 'files/sim-aug_14_2021.Rdata' in it.

## Error in outputs_or_evals(sim, sim@output_refs, TRUE, subset, index, methods, : trying to get slot "output_refs"
from an object of a basic class ("function") with no slots

## Error in lapply(out@out, function(i) f.hat.fit(object = i, xvar = xvar)[["fX"]]): object 'out' not found
```

₃₀ 4.2 Interaction effects

```
## Error in lapply(X = X, FUN = FUN, ...): object 'out' not found

## Error in quantile(resX3, probs = c(0.25, 0.5, 0.75), type = 1): object 'resX3' not found

## Error in sail:::plotInter(object = out@out[[indX3[1]]]$fit, xvar = xv, : object 'out' not found

## pdf
## 2

## Error in sail:::plotInter(object = out@out[[indX3[1]]]$fit, xvar = xv, : object 'out' not found

## pdf
## 2

## Error in sail:::plotInter(object = out@out[[indX3[2]]]$fit, xvar = xv, : object 'out' not found

## pdf
## 2

## Error in sail:::plotInter(object = out@out[[indX3[3]]]$fit, xvar = xv, : object 'out' not found

## pdf
## 2

## Error in lapply(X = X, FUN = FUN, ...): object 'out' not found
```

```
## Error in quantile(resX4, probs = c(0.25, 0.5, 0.75), type = 1): object 'resX4' not found

## Error in sail:::plotInter(object = out@out[[indX4[1]]]$fit, xvar = xv, : object 'out' not found

## pdf
## 2

## Error in sail:::plotInter(object = out@out[[indX4[1]]]$fit, xvar = xv, : object 'out' not found

## pdf
## 2

## Error in sail:::plotInter(object = out@out[[indX4[2]]]$fit, xvar = xv, : object 'out' not found

## pdf
## 2

## Error in sail:::plotInter(object = out@out[[indX4[3]]]$fit, xvar = xv, : object 'out' not found

## pdf
## 2
```

³¹ **5 Real Data Analysis**

³² **5.1 Gene-environment interactions in the Nurse Family Partnership program**

³⁴ **5.2 Study to Understand Prognoses Preferences Outcomes and Risks of Treatment**

Table 2: Comparison of analytic methods for selecting interactions using the Nurse Family Partnership program and the SUPPORT datasets. Averages (standard deviations in parentheses) are based on 200 bootstrap samples.

Method	Nurse Family Partnership		SUPPORT	
	Mean Squared Error	$ \hat{\mathcal{J}} $	AUC	$ \hat{\mathcal{H}} $
<code>sail</code>	199.7 (38.1)	5 (5)	0.66 (0.01)	25 (3)
<code>lassoBT</code>	194.01 (35.379)	7 (9)	0.65 (0.009)	49 (14)
<code>GLinternet</code> ^a	—	—	0.65 (0.009)	58 (7)

^a `GLinternet` results not reported for NFP data since the algorithm did not converge in many of the bootstrap samples.

^b $|\hat{\mathcal{J}}|$ is the number of variables selected by the method.

Test Set MSE

Based on 200 simulations

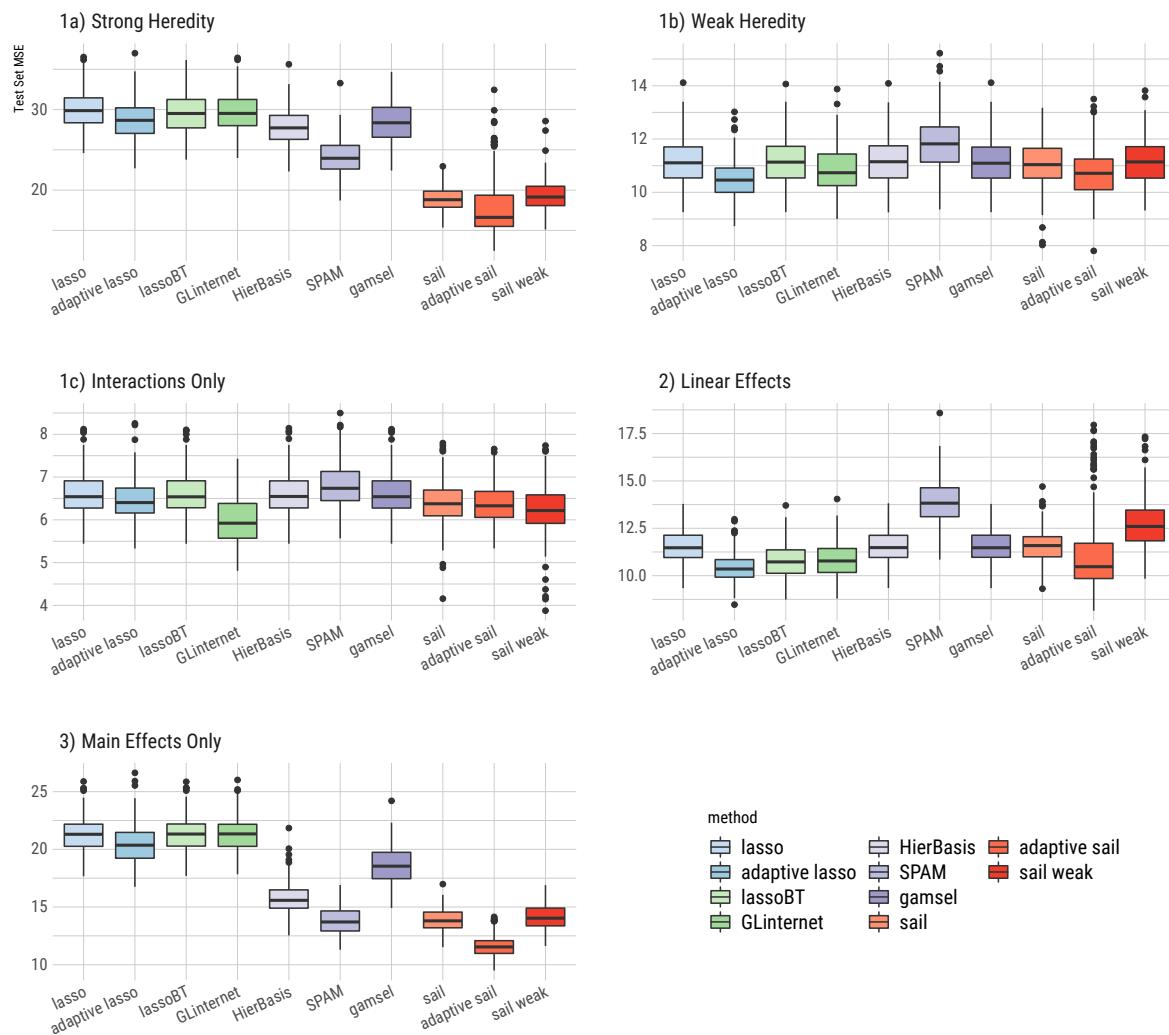


Figure 3: Boxplots of the test set mean squared error from 200 simulations for each of the five simulation scenarios.

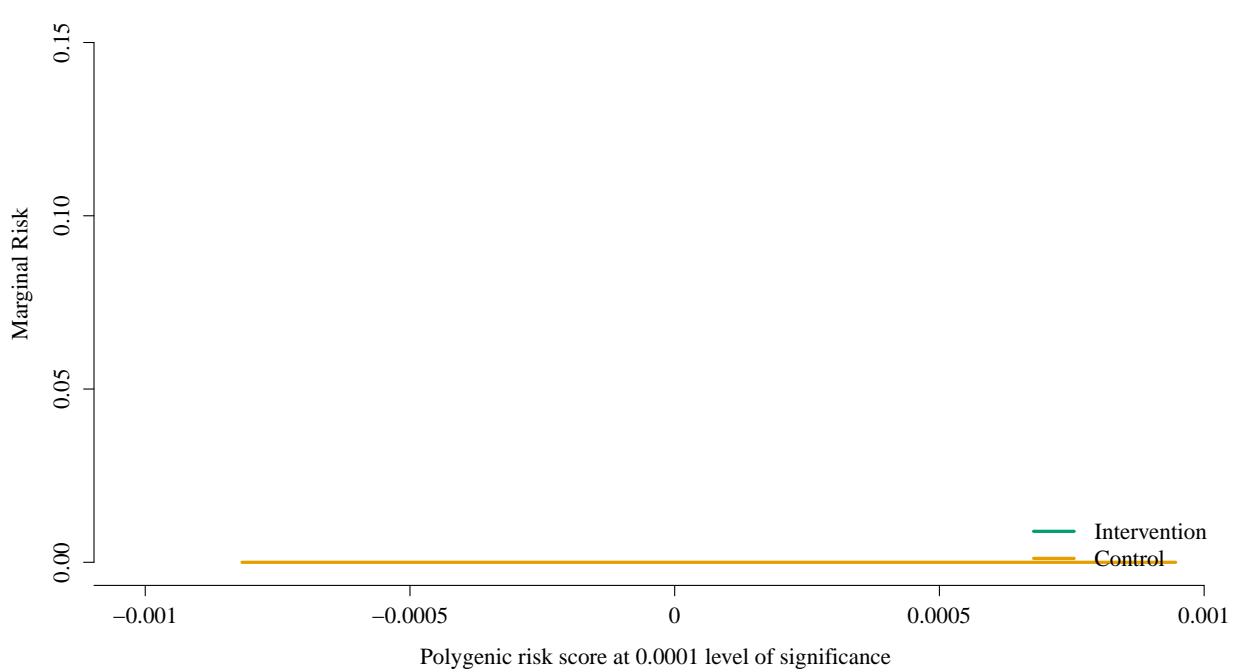


Figure 4: Estimated interaction effect identified by the weak heredity `sail` using cubic B-splines and $\alpha = 0.1$ for the Nurse Family Partnership data. The selected model, chosen via 10-fold cross-validation, contained three variables: the main effects for the intervention and the PRS for educational attainment using genetic variants significant at the 0.0001 level, as well as their interaction.

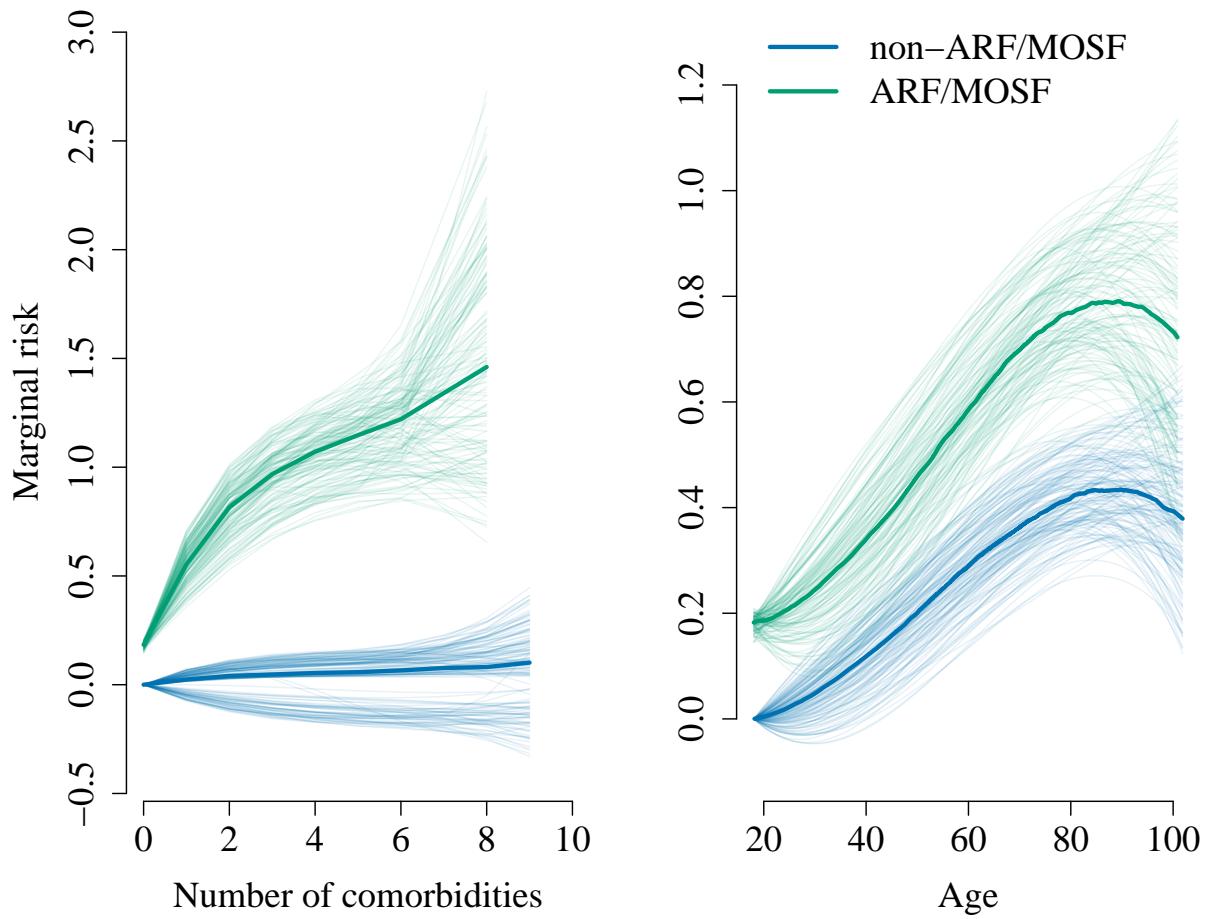


Figure 5: Illustration of estimated interaction effects identified by `sail` for the SUPPORT data. Median prediction curves in dark colors based on 200 train/validate/test splits represent the estimated marginal interaction effects. Coefficients estimated in each of the 200 train/validate/test splits were used to generate prediction curves representing a 90% confidence interval colored in corresponding light colors.