# Supplementary Information

## Gene-environment dependencies lead to collider bias in models with polygenic scores

Evelina T. Akimova,[1,2,*] Richard Breen, [1,3] David M. Brazel,[2,3] Melinda C. Mills[2,3]

*[1] Department of Sociology, University of Oxford, Oxford, OX1 1JD, United Kingdom*
*[2] Leverhulme Centre for Demographic Science, University of Oxford, Oxford, OX1 1JD, United Kingdom*
*[3] Nuffield College, University of Oxford, Oxford, OX1 1NF, United Kingdom*

*Corresponding author. Email address: evelina.akimova@sociology.ox.ac.uk

**S1. Expressions for the bias in gene-environment models due to interdependency of polygenic scores and environments**

The main text provided a general description of endogenous selection bias when polygenic scores and environments are not independent. Here we further illustrate this issue by deriving exact expressions for the bias under the assumption of linear relationships that can be modelled using regression analysis.

We assume that the data have been generated by the DAG shown in Figure 1A. Here U is an unobserved variable or set of variables that confounds the $E - Y$ relationship (this is equivalent to, but, in our view more transparent than, a depiction that would include correlated error terms for E and Y). We further assume that all the variables have unit standard deviation and that G is exogenous.

*S1.1 Additive model*

When the effects of G and E on Y are additive the true linear models given the data-generating process are:

$$E(E|G,U) = \alpha_0 + \alpha_1 G + \alpha_2 U \qquad \text{(Equation S1)}$$

$$E(Y|G,E,U) = \beta_0 + \beta_1 G + \beta_2 E + \beta_3 U \qquad \text{(Equation S2)}$$

We assume that the parameters $\alpha_1$, $\beta_1$, and $\beta_2$ are all positive. If estimate the models:

$$E(E|G) = a_0 + a_1 G \qquad \text{(Equation S3)}$$

30 $$E(Y|G,E) = b_0 + b_1 G + b_2 E \qquad \text{(Equation S4)}$$

31 the relationships between the true and estimated parameters are:

32 $$b_1 = \beta_1 - \frac{\alpha_1 \alpha_2 \beta_3}{1-\alpha_1^2} \qquad \text{(Equation S5)}$$

33 $$b_2 = \beta_2 + \frac{\alpha_2 \beta_3}{1-\alpha_1^2} \qquad \text{(Equation S6)}$$

34 The proof is as follows. Let $\beta_{YG}$ denote the coefficient of the unconditional

35 regression of Y on G and likewise for $\beta_{YE}$. Then, tracing the paths linking G and

36 Y in the Figure A1a we have:

37 $$\beta_{YG} = \beta_1 + \beta_2 \alpha_1$$

38 and

39 $$\beta_{YE} = \beta_2 + \beta_1 \alpha_1 + \beta_3 \alpha_2$$

40 Given that $\beta_{EG} = \alpha_1$ we then apply the standard formula to derive conditional

41 regression coefficients from unconditional:

42 $$b_1 = \frac{\beta_{YG} - \beta_{YE}\alpha_1}{1 - \alpha_1^2} = \frac{\beta_1 + \beta_2 \alpha_1 - (\beta_2 + \beta_1 \alpha_1 + \beta_3 \alpha_2)\alpha_1}{1 - \alpha_1^2}$$

43 Straightforward algebra yields (S5). $b_2$ is derived similarly. Notice, however,

44 that $a_1 = \alpha_1$ because G and U are unconditionally independent.

45 The bias in both estimates depends on the sign of $\alpha_2 \times \beta_3$: if this is positive

46 the estimate of the partial effect of G on Y, given E, will be downwardly biased and

47 the estimate of the effect of E on Y, given G, will be upwardly biased. If there is

48 no correlation between G and E ($\alpha_1 = 0$) then $b_1$ will be unbiased. If there is no

49 effect of an unmeasured confounder (either $\alpha_2 = 0$ and/or $\beta_3 = 0$) both $b_1$ and $b_2$

50 will be unbiased. The bias in the effect of G on Y has a different sign than the bias

51 in the effect of E on Y: if the bias in the latter is positive, the size of the genetic

52 effect will be underestimated relative to the environmental effect.

53 The example of coefficient deflation from Papageorge and Thom[1] can be

54 demonstrated following Equation S5. For instance, considering the case on

55 nonroutine interactive job tasks as the dependent variable, we see that the

56 baseline coefficient of the educational attainment polygenic score is 0.185, which

57 reflects a model without any environmental and phenotypic covariates (Table 6 in

58 Papageorge and Thom[1]). In the model with educational controls (respondent's

59 years of schooling and parental education), the polygenic score coefficient drops to

60 0.055 reflecting a 70% negative change. Since the dependent variable is

61 standardised, we can assess the relative importance of collider bias which is $\frac{\alpha_1 \alpha_2 \beta_3}{1-\alpha_1^2}$

62 from Equation S5 under additional assumptions. If we allow the coefficient of the

63 correlation between educational attainment polygenic score and respondents

64 years of schooling $\alpha_1$=0.300, and the presence of unobserved confounder U,

65 positively correlated with both years of schooling and job task (for example, living

66 in advantaged neighbourhood as a child), we have $\alpha_2 = 0.250$ and $\beta_3 = 0.250$.

67 These are all plausible and rather modest suggestions following correlation matrix

68 from Table 6 in Papageorge and Thom[1], leading the inflation bias to be:

69 $$\frac{\alpha_1 \alpha_2 \beta_3}{1-\alpha_1^2} = \frac{0.300 \times 0.250 \times 0.250}{1-0.300^2} = 0.021$$

70 which explains 16% downward change of polygenic score coefficient.

71 *S1.2 G×E interaction model*

4

72    The DAG in Figure 2A shows the case in which the effect of G on Y varies with E.

73    In this case, the true linear models given the data-generating process are:

74    $$E(E|G,U) = \alpha_0 + \alpha_1 G + \alpha_2 U \qquad \text{(Equation S7)}$$

75    $$E(Y|G,E,U) = \beta_0 + \beta_1 G + \beta_2 E + \beta_3 U + \beta_4 (GE) \quad \text{(Equation S8)}$$

76    We estimate:

77    $$E(E|G) = a_0 + a_1 G \qquad \text{(Equation S9)}$$

78    $$E(Y|G,E) = b_0 + b_1 G + b_2 E + b_4 GE \qquad \text{(Equation S10)}$$

79        In this case, $b_4$ is an unbiased estimate of $\beta_4$ because the backdoor path from

80    G-E to Y is blocked by E. The bias in $b_1$ and $b_2$ will be the same as above. In the

81    case in which E is a binary variable, coded 0 and 1, $b_4$ will be an unbiased estimate

82    of the difference in the effect of G at $E = 1$ and $E = 0$, but the estimate of the

83    baseline effect of G on Y when $E = 0$ will be biased.

84    *S1.3 Bias in $R^2$*

85    The $R^2$ for models S4 and S10 will be biased. In the additive case, for example,

86    the true $R^2$ attributable to G and E is:

87    $$\frac{\beta_1^2 var(G) + \beta_2^2 var(E) + 2\beta_1\beta_2 cov(G,E)}{var(Y)} = \beta_1^2 + \beta_2^2 + 2\beta_1\beta_2\alpha_1 \quad \text{(Equation S11)}$$

88    (using the assumption that all the variables have unit standard deviation). But

89    the reported $R^2$ from model S4 is:

90 $$b_1^2 + b_2^2 + 2b_1b_2\alpha_1 \qquad\qquad \text{(Equation S12)}$$

91 Substituting S5 and S6 into S12 we calculate the inflation of $R^2$ due to

92 confounding and collider bias. This is:

93 $$R^2 bias = \alpha_2\beta_3 \left[\frac{\alpha_2\beta_3}{1-\alpha_1^2} + 2\beta_2\right] \qquad\qquad \text{(Equation S13)}$$

94 Confounder bias arises from $\alpha_2\beta_3$. The derivative of S13 with respect to this

95 is positive provided that $1 - \alpha_1^2 > 0$. The derivative of S13 with respect to $\alpha_1$

96 (which captures the association between G and E) is:

97 $$\frac{2\alpha_1\alpha_2^2\beta_3^2}{(1-\alpha_1^2)^2}$$

98 The sign of this depends on the sign of the numerator. When it is positive

99 both confounding and collider bias will inflate the reported $R^2$. As an example,

100 consider a case in which $\beta_1 = 0.465, \beta_2 = 0.505, \beta_3 = 0.231, \alpha_1 = 0.209, \alpha_2 = 0.693$.

101 Then the observed $R^2 = 0.758$, whereas the true share of the variance in Y

102 explained by G and E is 0.569. The inflation bias here is:

103 $$0.693 \times 0.231 \left[\frac{0.693 \times 0.231}{1 - 0.209^2} + 2 \times 0.505\right] = 0.188$$

104 If the correlation between G and E had been larger and/or if the confounding of E

105 had been greater, the reported $R^2$ would have been larger because of the greater

106 bias.

107 **S2. Simulation analyses**

6

108   The code for the simulations and figures is available on Zenodo (DOI:

109   10.5281/zenodo.4184673) and GitHub (https://github.com/eva-akimova/collider-

110   simulations.git). For the figures presented in the main text, we simulated

111   scenarios of OLS regressions where G – E association varies between 0 and .5; G

112   – Y and E – Y coefficients are both positive and .6; uncontrolled confounder, U, is

113   positively and modestly, moderately or strongly correlated with covariate, E, and

114   outcome, Y, ($r$ = .12, $r$ = .25, and $r$ = .38 respectively for the three scenarios). For

115   the gene-environment interaction models we simulated the same settings and

116   added GxE coefficient at a fixed value of 0.1 for all scenarios.

117   Here, we expand our simulation analyses and further illustrate the bias in

118   gene-environment interaction models where unobserved confounder, U, interacts

119   with covariate, E, at a fixed value of 0.2 for all scenarios. Simulation results

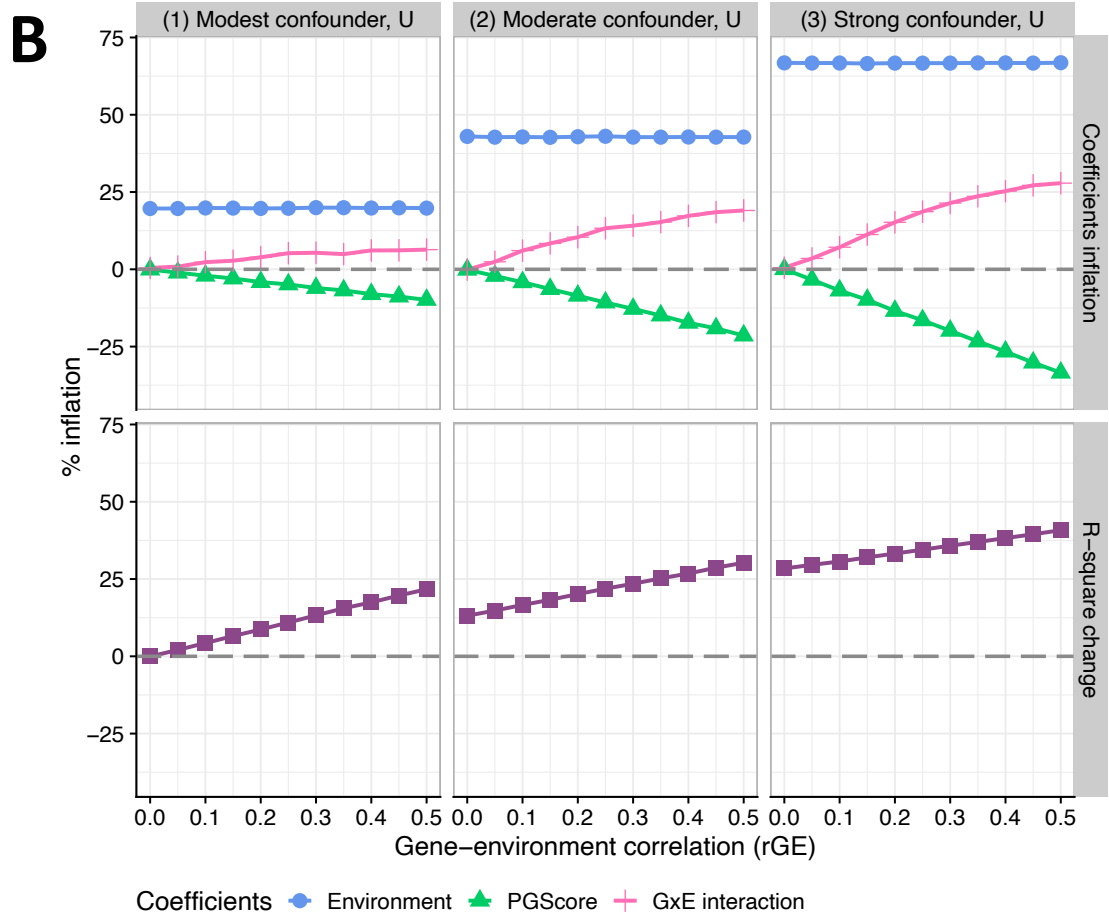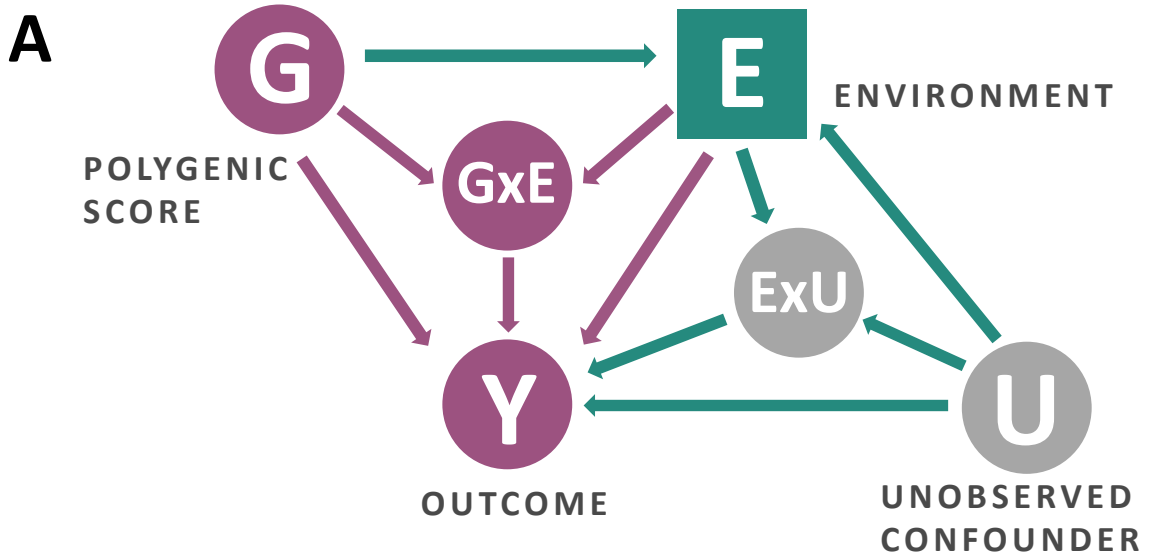120   presented below in Figure S1 along with a DAG to illustrate the bias.

121

**Figure S1. Collider bias in polygenic gene-environment interaction models**

124    Figure S1. Collider bias in polygenic gene-environment interaction models. Panel

125    A. Schematic diagram of the collider bias which occurs between polygenic score,

126    environment, and outcome in cases of gene-environment interdependence. Dark

127    purple circles represent variables, unobserved confounders and it interaction term

128    with E are shown in grey circles, collider variables are indicated in squares. By

129    adding E into the model with the polygenic score G, we make E a collider. A collider

130    that is not conditioned on, blocks the path between its sources (G and U); once a

131    collider is controlled for, the path is opened as indicated by green nodes. ExU

132    interaction term is also on the bias path once E is conditioned on. Panel B (top).

133    Spurious regression estimates for the polygenic score and environment along with

134    inflated interaction terms from the series of OLS simulations reflecting a range of

135    gene-environment interdependence and the presence of modest, moderate, or

136    strong confounder, U. Collider bias due to positive values of gene-environment

137    correlation and the presence of an uncontrolled confounder, which is positively

138    correlated with covariate and outcome, results in deflation of polygenic score

139    estimates. Deflation is greater the higher the gene-environment correlation;

140    greater confounding also results in greater bias. The interaction term is affected

141    proportionally to the strength of rGE and unobserved confounder, U. Panel B

142    (bottom). R-squared inflation plot from the series of OLS simulations; collider bias

143    results in inflated values of explained variance statistics. R-squared statistics for

144    the model with endogenous covariate and polygenic score includes not only the true

145    share of the variance in Y explained by G and E (baseline estimate indicated by 0),

146    but also the elements of variance that are due to gene-environment correlation and

147    confounder(s), U.

## 148  **References**

149  1. Papageorge, N.W., and Thom, K. (2019). Genes, education, and labor market outcomes:

150      Evidence from the Health and Retirement Study. Journal of the European Economic

151      Association jvz072.

152