



1           The same adaptations were explored in the final outcome regression model, which  
2 again lead to the same conclusions. The observation that these extended models did not have a  
3 significant impact on parameter estimates and inference, shows that there is no evidence  
4 against the adequacy of the models presented in the previous sections.

### 5 **Goodness-of-fit**

6           Finally, as a more formal test of goodness-of-fit, the models presented in the main text  
7 were specified in Stata using the `svy: logit` command and subsequently subjected to a  
8 modified Hosmer–Lemeshow test that incorporates the complex survey design (2). The  
9 traditional Hosmer-Lemeshow test (3) is one of the most popular tests that provides insight  
10 into the general goodness-of-fit of a logistic regression model.

11           While this test is appropriate for simple random sampling designs (where responses  
12 can typically be assumed to be independent and identically distributed, i.i.d), the use of this  
13 test may not be valid with complex survey designs (2,4). To account for this, Archer &  
14 Lemeshow (2) developed a modified Hosmer- Lemeshow test that allows evaluating  
15 goodness-of-fit of a survey-weighted logistic regression. This modified statistic is calculated  
16 by first dividing the sample into deciles of risk based on the observed residuals (i.e. the first  
17 decile consists of subjects with the 10% lowest residuals. Second, a mean residual is  
18 estimated per subgroup by summing the decile’s residuals, weighted by the corresponding  
19 sampling weights and dividing by the total of the sampling weights. Third, the variance of the  
20 estimated vector of mean residuals is estimated via Taylor series linearization. Finally, a Wald  
21 statistic is calculated from the estimated mean residuals and their estimated variances and  
22 subsequently compared to a reference F-distribution (2).

23           We first evaluated model fit in the propensity score model in Stata by using the  
24 `svylogisticgof` command. The modified Hosmer-Lemeshow test was found to be significant at  
25 the .05 significance level, indicating a lack of fit ( $F(9,13804) = 2.85, p = .002$ ). Further

1 investigation of the model revealed that including the continuous variable age was responsible  
2 for the observed lack of fit. Moreover, adding an additional interaction term between age and  
3 education for example produced a much better fit ( $F(9,13804) = 1.12, p = .347$ ). Running the  
4 same regression analysis using estimated propensities based on a model that included  
5 additional interactions between the covariates and age led to the same conclusions as before.

6         Second, we evaluated model fit in the same way for the final outcome regression  
7 model including the propensity score and all covariates. The modified Hosmer-Lemeshow test  
8 was not statistically significant at the  $\alpha = .05$  level ( $F(9,13804) = 1.35, p = .207$ ), suggesting  
9 no lack of fit. Though the propensity score model without any interaction terms was found to  
10 have a problematic fit to the data, inclusion of an interaction term between age and education  
11 remedied this problem. Additionally, our previously presented sensitivity analysis revealed  
12 that both a propensity model with and without interactions yielded similar results. We  
13 therefore decided to continue working with the models as presented in the main text, and  
14 concluded insufficient evidence for lack of fit. Note that, due to the lack of guidelines on the  
15 matter and the large dataset, we did not evaluate potentially influential observations.

**REFERENCES**

1. Millimet DL, & TR. On the specification of propensity scores, with applications to the analysis of trade policies. *Journal of Business & Economic Statistics*. 2009; 27(3): 397-415.
2. Archer KJ, Lemeshow S. Goodness-of-fit Test for a Logistic Regression Model Fitted using Survey Sample Data. *The Stata Journal*. 2006; 6(1): 97-105.
3. Hosmer DW, Lemeshow S. Goodness of fit tests for the multiple logistic regression model. *Communications in statistics-Theory and Methods*. 1980; 9(10): 1043-1069.
4. Archer KJ, Lemeshow S, Hosmer DW. Goodness-of-fit tests for logistic regression models when data are collected using a complex sampling design. *Computational Statistics & Data Analysis*. 2007; 51(9): 4450-4464.