

## Appendix 1. Data generation process.

We generated a seven-step ordinal outcome variable ( $y$ ), on which a dichotomous exposure variable ( $x_e$ ) had no effect, whereas confounding variables ( $x_1, \dots, x_5$ ) did. We did not model a relationship or correlation between the confounding variables. In order to be able to create  $y$ , we first created a latent variable ( $y^*$ ), which reflects the presence of confounding:(12)

$$y^* = \sum_{j=1}^5 \beta_j x_{ji} + \varepsilon_i, \text{ with } \varepsilon_i \sim N(0,1)$$

In which  $x_{1i}, \dots, x_{5i}$  are the confounding variables for patient  $i$  and  $\beta_1, \dots, \beta_5$  are the regression coefficients of these variables. The term  $\varepsilon_i$  denotes the error term, which varies per subject. Thus, observations with more confounding present (i.e. more exposed to dichotomous confounding present or with higher values for continuous confounding variables) generally have higher overall risk values.

Out of  $y^*$ , we created six variables ( $y^*_1 - y^*_6$ ), that represented the six cut-off points between the seven outcome categories:

$$y^*_{1-6} = \text{std}(y^*) + \text{invnormal}(\alpha_{1-6})$$

We standardized  $y^*$  so that we could shift the placement of the latent variable distributions accurately, by adding the inverse normal (invnormal) of values between 0 and 1 ( $\alpha_{1-6}$ ). For  $\alpha_{1-6}$  we used  $1/7, \dots, 6/7$ , as we aimed to generate a uniform distribution of  $y$ . In other words,  $\alpha_{1-6}$  are the fractions of observations we aim to keep below each cut-off point of the mRS. Thus, every observation had six latent variables ( $y^*_1 - y^*_6$ ), which only differed in the addition of the normal sextiles. Finally, we generated the ordinal outcome  $y$  based on the six latent variables created from  $y^*$

$$\text{if } y^*_1 > 0 \text{ then } y = 6, \text{ if } y^*_2 > 0 \text{ then } y = 5, \dots, \text{ if } y^*_6 > 0 \text{ then } y = 1, \text{ else } = 0$$

This states, that if condition  $y^*_1 > 0$  for an observation holds true, then  $y = 6$ . If  $y^*_1 > 0$  is false, it checks whether  $y^*_2 > 0$  holds true. If  $y^*_2 > 0$  is true, then  $y = 5$ , et cetera. If the last condition ( $y^*_6 > 0$ ) is false, then  $y = 0$ . So, observations with  $y = 6$ , have a positive  $y^*_1$  despite adding a negative term to the equation. At the other end of the spectrum, observations with  $y = 0$  have a  $y^*_6$  smaller than 0, despite adding a positive term. So, in short, observations with more

confounding present are more likely to have higher  $y^*$ , and are therefore more likely to end up in higher categories of  $y$ .

To complete the requirement for confounding, the exposure also needed to be dependent on the confounding variables. Therefore, we constructed the binary exposure ( $X_e$ ) based on whether the confounding variable terms ( $\beta_{c1-ci} X_{c1-ci} - \beta_{c1-ci} X_{c1-ci}(\text{mean})$ ) and a base risk of 0.5 are higher than a randomly generated number between 0 and 1. Here, we subtracted the mean of each confounding variable by its regression coefficient to balance the equation around 0.5, and thus a similar number of cases exposed in each run and scenario.