

# Structural Equation Modeling

## P.08 - MIMIC Models and Instrumental Variables

November 15, 2022 (13:39:56)

---

### Lab Description

For this practical you will need the following packages: `lavaan` and `semPlot`. You can install and load these packages using the following code:

```
# Install packages.
install.packages(c("lavaan", "semPlot", "mvtnorm", "GGally"))

# Load the packages.
library(lavaan)
library(semPlot)
```

### Exercise 1

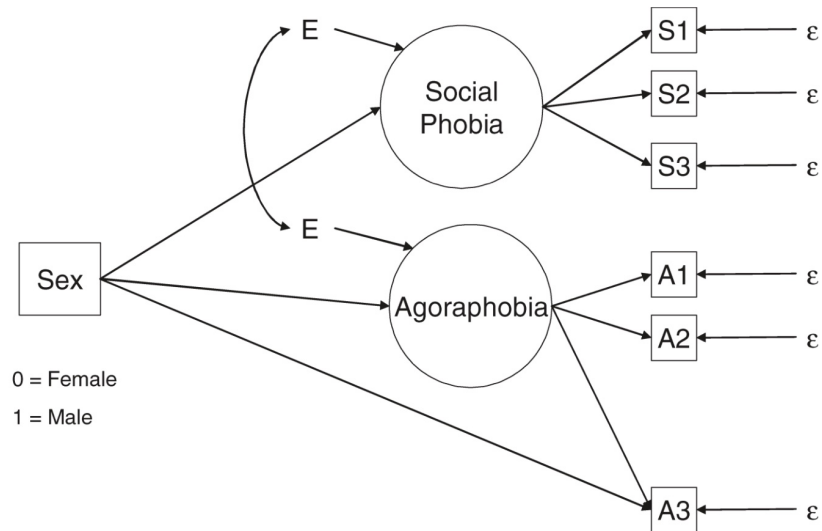
Estimate the model in *Figure 1* in `lavaan` and examine if there is evidence of Differential Item Functioning (DIF) in the measurement instruments. To help you get started, you are provided with the code that contains the correlations and standard deviations corresponding to the model depicted in *Figure 1*.

Standard deviations and correlations.

```
# Standard deviations.
sd <- "2.26 2.73 2.11 2.32 2.61 2.44 0.50"

# Correlations.
cor <- "
  1.000
  0.705 1.000
  0.724 0.646 1.000
  0.213 0.195 0.190 1.000
  0.149 0.142 0.128 0.521 1.000
  0.155 0.162 0.135 0.557 0.479 1.000
  -0.019 -0.024 -0.029 -0.110 -0.074 -0.291 1.000
"

# Get covariances.
cov <- getCov(cor, sds = sd, names = c("S1", "S2", "S3", "A1", "A2", "A3", "sex"))
```



Sample Correlations and Standard Deviations (SDs);  $N = 730$  (365 males, 365 females)

	S1	S2	S3	A1	A2	A3	Sex
S1	1.000						
S2	0.705	1.000					
S3	0.724	0.646	1.000				
A1	0.213	0.195	0.190	1.000			
A2	0.149	0.142	0.128	0.521	1.000		
A3	0.155	0.162	0.135	0.557	0.479	1.000	
Sex	-0.019	-0.024	-0.029	-0.110	-0.074	-0.291	1.000
SD:	2.260	2.730	2.110	2.320	2.610	2.440	0.500

**FIGURE 7.5.** MIMIC model of Social Phobia and Agoraphobia. S1, giving a speech; S2, meeting strangers; S3, talking to people; A1, going long distances from home; A2, entering a crowded mall; A3, walking alone in isolated areas. (All questionnaire items rated on 0–8 scales, where 0 = no fear and 8 = extreme fear.)

Figure 1: Reproduction of Figure 7.5 from Brown (2014, p. 275)

We start by specifying the syntax for the *MIMIC* model.

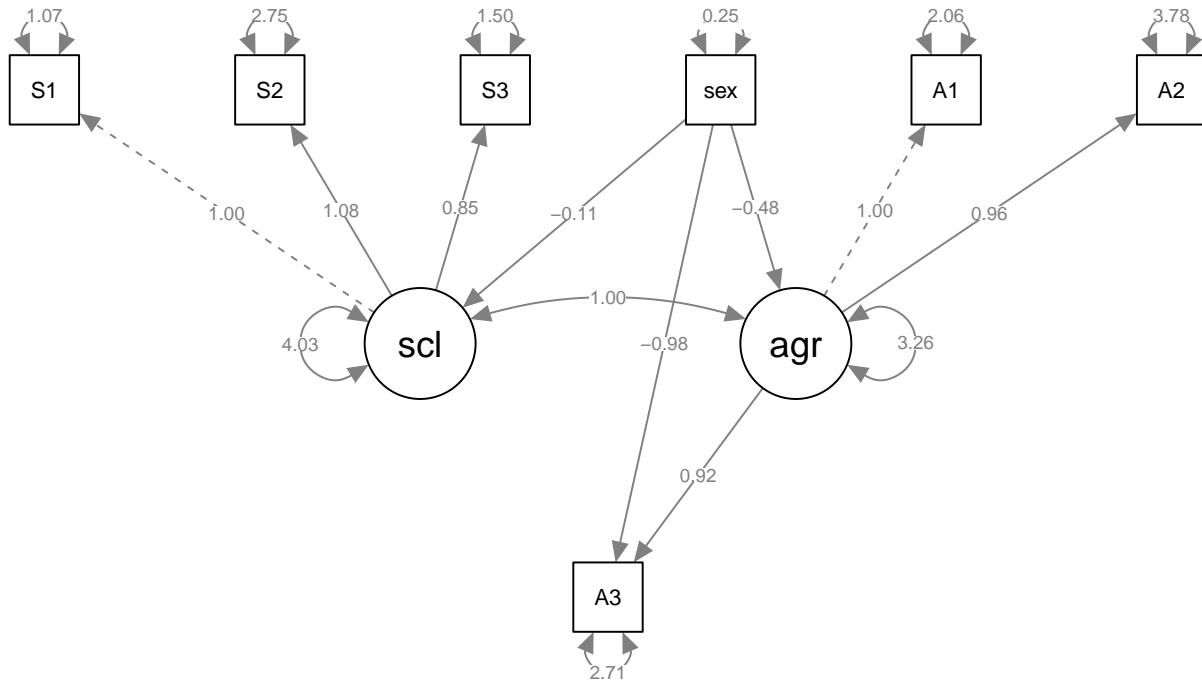
```
# Model syntax.
model_ex_1 <- "
  # Measurement part.
  social =~ S1 + S2 + S3
  agoraph =~ A1 + A2 + A3

  # Regression equations.
  social ~ sex
  agoraph ~ sex
  A3 ~ sex

  # Covariances.
  social ~~ agoraph
"

# Fit the model.
model_ex_1_fit <- cfa(model_ex_1, sample.cov = cov, sample.nobs = 730)

# Visualize the model.
semPaths(model_ex_1_fit, what = "paths", whatLabels = "est")
```



```
# Model summary.
summary(model_ex_1_fit, fit.measures = TRUE, standardized = TRUE, modindices = TRUE)
```

```
## lavaan 0.6-12 ended normally after 52 iterations
```

```
##
```

```
## Estimator ML
```

```
## Optimization method NLMINB
```

```
## Number of model parameters 16
```

```
##
```

```

##      Number of observations                730
##
## Model Test User Model:
##
##      Test statistic                        3.797
##      Degrees of freedom                    11
##      P-value (Chi-square)                 0.975
##
## Model Test Baseline Model:
##
##      Test statistic                        1771.017
##      Degrees of freedom                    21
##      P-value                              0.000
##
## User Model versus Baseline Model:
##
##      Comparative Fit Index (CFI)          1.000
##      Tucker-Lewis Index (TLI)            1.008
##
## Loglikelihood and Information Criteria:
##
##      Loglikelihood user model (H0)        -9167.606
##      Loglikelihood unrestricted model (H1) -9165.707
##
##      Akaike (AIC)                        18367.212
##      Bayesian (BIC)                      18440.701
##      Sample-size adjusted Bayesian (BIC)  18389.896
##
## Root Mean Square Error of Approximation:
##
##      RMSEA                               0.000
##      90 Percent confidence interval - lower 0.000
##      90 Percent confidence interval - upper 0.000
##      P-value RMSEA <= 0.05                1.000
##
## Standardized Root Mean Square Residual:
##
##      SRMR                                0.011
##
## Parameter Estimates:
##
##      Standard errors                      Standard
##      Information                          Expected
##      Information saturated (h1) model      Structured
##
## Latent Variables:
##
##      Estimate Std.Err z-value P(>|z|) Std.lv Std.all
##      social =~
##      S1         1.000          2.007  0.889
##      S2         1.079    0.045 23.967  0.000  2.166  0.794
##      S3         0.855    0.035 24.534  0.000  1.716  0.814

```

```

##   agoraph =~
##       A1           1.000           1.820   0.785
##       A2           0.956   0.066   14.388   0.000   1.739   0.667
##       A3           0.917   0.063   14.495   0.000   1.669   0.684
##
## Regressions:
##           Estimate Std.Err z-value P(>|z|) Std.lv Std.all
##   social ~
##       sex       -0.109   0.158   -0.689   0.491   -0.054   -0.027
##   agoraph ~
##       sex       -0.475   0.160   -2.973   0.003   -0.261   -0.130
##   A3 ~
##       sex       -0.985   0.148   -6.654   0.000   -0.985   -0.202
##
## Covariances:
##           Estimate Std.Err z-value P(>|z|) Std.lv Std.all
##   .social ~~
##   .agoraph      0.999   0.171   5.857   0.000   0.276   0.276
##
## Variances:
##           Estimate Std.Err z-value P(>|z|) Std.lv Std.all
##   .S1           1.072   0.126   8.533   0.000   1.072   0.210
##   .S2           2.750   0.195  14.087   0.000   2.750   0.370
##   .S3           1.501   0.114  13.169   0.000   1.501   0.338
##   .A1           2.062   0.217   9.498   0.000   2.062   0.384
##   .A2           3.777   0.264  14.302   0.000   3.777   0.555
##   .A3           2.705   0.214  12.642   0.000   2.705   0.455
##   .social       4.026   0.284  14.175   0.000   0.999   0.999
##   .agoraph      3.257   0.317  10.269   0.000   0.983   0.983
##
## Modification Indices:
##
##       lhs op      rhs   mi    epc sepc.lv sepc.all sepc.nox
## 1      sex ==    sex 0.000 0.000 0.000 0.000 0.000
## 2    social =~    A1 1.779 0.056 0.113 0.049 0.049
## 3    social =~    A2 0.505 -0.033 -0.067 -0.026 -0.026
## 4  agoraph =~    S1 0.010 -0.004 -0.007 -0.003 -0.003
## 5  agoraph =~    S2 0.461 0.031 0.057 0.021 0.021
## 6  agoraph =~    S3 0.286 -0.019 -0.034 -0.016 -0.016
## 7      S1 ==    S2 0.305 -0.298 -0.298 -0.174 -0.174
## 8      S1 ==    S3 0.459 0.303 0.303 0.239 0.239
## 9      S1 ==    A1 0.322 0.053 0.053 0.036 0.036
## 10     S1 ==    A2 0.018 -0.015 -0.015 -0.007 -0.007
## 11     S1 ==    A3 0.310 -0.054 -0.054 -0.032 -0.032
## 12     S2 ==    S3 0.007 -0.035 -0.035 -0.017 -0.017
## 13     S2 ==    A1 0.025 -0.020 -0.020 -0.008 -0.008
## 14     S2 ==    A2 0.000 -0.002 -0.002 -0.001 -0.001
## 15     S2 ==    A3 0.734 0.110 0.110 0.040 0.040
## 16     S3 ==    A1 0.171 0.039 0.039 0.022 0.022
## 17     S3 ==    A2 0.135 -0.041 -0.041 -0.017 -0.017
## 18     S3 ==    A3 0.531 -0.071 -0.071 -0.035 -0.035

```

```
## 19      A1 ~~      A2 0.599 -0.409 -0.409 -0.147 -0.147
## 20      A1 ~~      A3 0.819 -0.451 -0.451 -0.191 -0.191
## 21      A2 ~~      A3 2.184 0.625 0.625 0.195 0.195
## 22  social ~      A3 0.599 -0.044 -0.022 -0.054 -0.054
## 23 agoraph ~      A3 0.599 0.145 0.080 0.194 0.194
## 24      sex ~ social 0.000 0.990 1.987 3.977 3.977
## 25      sex ~ agoraph 0.000 0.044 0.080 0.160 0.160
## 26      sex ~      A3 0.000 0.018 0.018 0.086 0.086
```

The *MIMIC* model provides a good fit to the data, with a  $\chi^2(11) = 3.80$ ,  $p$ -value = .98,  $RMSEA = 0.00$ , and  $CFI = 1.00$ .

Regarding the evidence for DIF, the following paragraph from Brown (2014, p. 280) is relevant:

Consistent with the researcher's predictions, the results of the *MIMIC* model show that the A3 indicator is not invariant for males and females (akin to intercept non-invariance in multiple-groups CFA). This is reflected by the significant direct effect of **sex** on the A3 indicator ( $z = 6.65$ ,  $p < .001$ ) that is not mediated by **agoraphobia**. In other words, when the latent variable of **agoraphobia** is held constant, there is a significant direct effect of **sex** on the A3 indicator. Thus, at any given value of the factor, women score significantly higher on the A3 indicator than men (by .985 units, or nearly a full point on the 0–8 scale). This is evidence of *differential item functioning*; that is, the item behaves differently as an indicator of **agoraphobia** in men and women.

## Exercise 2

Open the dataset `card.csv` available on Canvas in the folder corresponding to the current practical. This dataset contains several variables used by David Card (1995) to estimate the causal effect of education on wages using proximity to college as an instrumental variable. You can find more information about this dataset at [this link](#).

Set the working directory to the location where your data file has been downloaded and load the data.

```
# For example.
setwd("/Users/mihai/Downloads")

# Load data.
data_ex_2 <- read.csv("card.csv")

# Inspect the data.
View(data_ex_2)
```

- Estimate a model in which you only regress `lwage` on `educ` (i.e., without including an instrumental variable). What do you conclude from this regression?

```
# Model syntax.
model_ex_2_a <- "lwage ~ educ"

# Fit the model.
model_ex_2_a_fit <- sem(model_ex_2_a, data_ex_2)
```

```
# Visualize the model.
semPaths(model_ex_2_a_fit, what = "paths", whatLabels = "est")
```



```
# Model summary.
summary(model_ex_2_a_fit, standardized = TRUE, rsquare = TRUE)
```

```
## lavaan 0.6-12 ended normally after 1 iterations
##
## Estimator ML
## Optimization method NLMINB
## Number of model parameters 2
##
## Number of observations 3010
##
## Model Test User Model:
##
## Test statistic 0.000
## Degrees of freedom 0
##
## Parameter Estimates:
##
## Standard errors Standard
## Information Expected
## Information saturated (h1) model Structured
##
## Regressions:
## Estimate Std.Err z-value P(>|z|) Std.lv Std.all
## lwage ~
## educ 0.052 0.003 18.159 0.000 0.052 0.314
##
## Variances:
## Estimate Std.Err z-value P(>|z|) Std.lv Std.all
## .lwage 0.177 0.005 38.794 0.000 0.177 0.901
##
## R-Square:
```

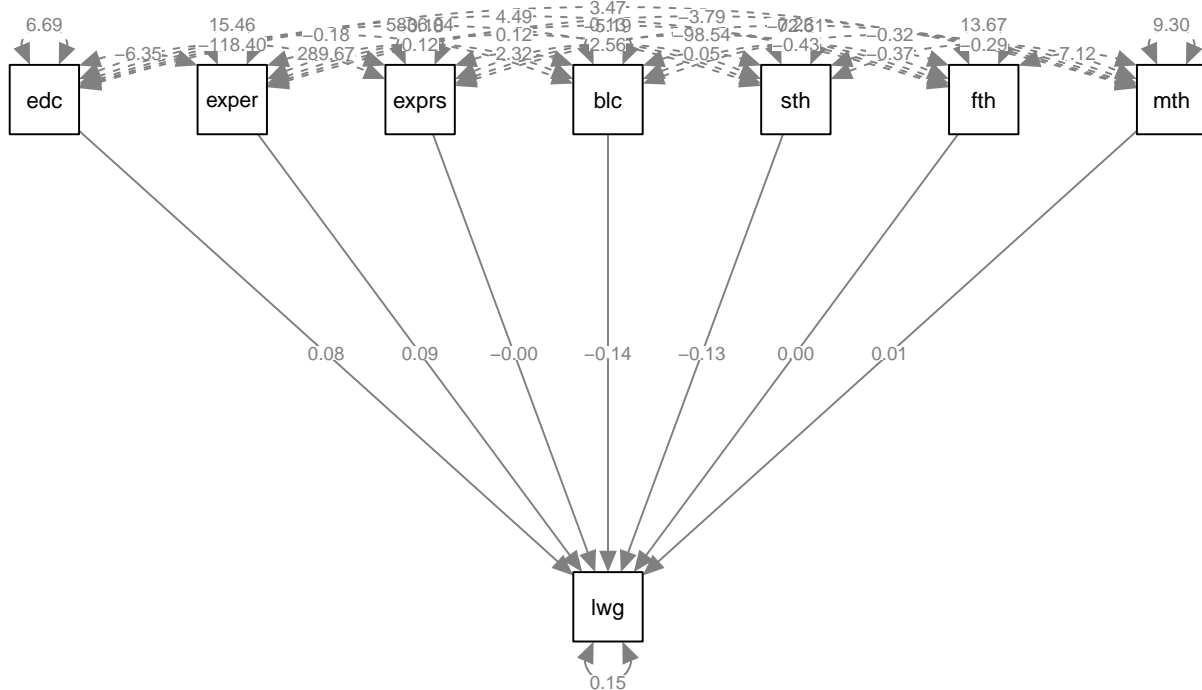
```
##           Estimate
##    lwage      0.099
```

- b. Re-estimate the model at point (a), but this time with the following control variables added: `exper`, `expersq`, `black`, `south`, `fatheduc`, and `motheduc`.

```
model_ex_2_b <- "
    lwage ~ educ + exper + expersq + black + south + fatheduc + motheduc
"

# Fit the model.
model_ex_2_b_fit <- sem(model_ex_2_b, data_ex_2)

# Visualize the model.
semPaths(model_ex_2_b_fit, what = "paths", whatLabels = "est")
```



```
# Model summary.
summary(model_ex_2_b_fit, standardized = TRUE, rsquare = TRUE)
```

```
## lavaan 0.6-12 ended normally after 1 iterations
```

```
##
## Estimator ML
## Optimization method NLMINB
## Number of model parameters 8
##
## Used Total
## Number of observations 2220 3010
##
## Model Test User Model:
##
## Test statistic 0.000
```



```
## Degrees of freedom          0
##
## Parameter Estimates:
##
## Standard errors            Standard
## Information                Expected
## Information saturated (h1) model Structured
##
## Regressions:
##      Estimate Std.Err z-value P(>|z|) Std.lv Std.all
## lwage ~
## educ      0.076   0.004  17.563   0.000   0.076   0.450
## exper      0.089   0.008  11.132   0.000   0.089   0.800
## expersq    -0.002   0.000   -6.001   0.000  -0.002  -0.421
## black     -0.145   0.024   -5.919   0.000  -0.145  -0.120
## south     -0.131   0.018   -7.353   0.000  -0.131  -0.144
## fatheduc    0.002   0.003    0.658   0.511   0.002   0.017
## motheduc    0.007   0.004    2.039   0.041   0.007   0.050
##
## Variances:
##      Estimate Std.Err z-value P(>|z|) Std.lv Std.all
## .lwage      0.147   0.004  33.317   0.000   0.147   0.759
##
## R-Square:
##      Estimate
## lwage      0.241
```

The problem with treating the direct association between `educ` and `lwage` as a causal effect is that there are likely many omitted variables that affect both education and wages. We could control for those variables by measuring them and including them in the model (i.e., as we did at point *b*). But there is no way we can control for all possible confounding variables, especially because some variables are difficult to measure (e.g., ability). It is therefore likely that education is correlated with the error term in the regression (i.e., a form of endogeneity), and that our regression coefficient is, in turn, biased to an unknown degree. David Card proposed to solve this problem by introducing proximity to college as an instrumental variable. Specifically, `nearc4` was a dummy indicator variable for whether or not the person was raised in a local labor market that included a four-year college.

c. Re-estimate the model at point *(b)* with the following additions:

- add `nearc4` as an instrumental variable for `educ`, while controlling for `fatheduc` and `motheduc`
- add a covariance between the error terms of `educ` and `lwage`

Does this model provide evidence of endogeneity of `educ`? Why (not)?

```
model_ex_2_c <- "
  lwage ~ educ + exper + expersq + black + south + fatheduc + motheduc
  educ ~ nearc4 + fatheduc + motheduc
  lwage ~~ educ
"

# Fit the model.
```

```

model_ex_2_c_fit <- sem(model_ex_2_c, data_ex_2)

# Model summary.
summary(model_ex_2_c_fit, standardized = TRUE, rsquare = TRUE)

```

```

## lavaan 0.6-12 ended normally after 40 iterations
##
##      Estimator                      ML
##      Optimization method          NLMINB
##      Number of model parameters      13
##
##                                     Used      Total
##      Number of observations          2220      3010
##
## Model Test User Model:
##
##      Test statistic                  790.934
##      Degrees of freedom                4
##      P-value (Chi-square)              0.000
##
## Parameter Estimates:
##
##      Standard errors                Standard
##      Information                    Expected
##      Information saturated (h1) model  Structured
##
## Regressions:
##
##      Estimate  Std.Err  z-value  P(>|z|)  Std.lv  Std.all
##      lwage ~
##      educ      0.226    0.065    3.454    0.001    0.226    1.241
##      exper      0.089    0.008   11.321    0.000    0.089    0.741
##      expersq    -0.002    0.000   -5.959    0.000   -0.002   -0.389
##      black     -0.147    0.024   -6.036    0.000   -0.147   -0.114
##      south     -0.121    0.018   -6.666    0.000   -0.121   -0.124
##      fatheduc   -0.031    0.015   -2.079    0.038   -0.031   -0.245
##      motheduc   -0.023    0.014   -1.641    0.101   -0.023   -0.148
##      educ ~
##      nearc4      0.364    0.103    3.526    0.000    0.364    0.065
##      fatheduc     0.216    0.017   13.030    0.000    0.216    0.309
##      motheduc     0.203    0.020   10.154    0.000    0.203    0.239
##
## Covariances:
##
##      Estimate  Std.Err  z-value  P(>|z|)  Std.lv  Std.all
##      .lwage ~~
##      .educ      -0.747    0.326   -2.291    0.022   -0.747   -0.660
##
## Variances:
##
##      Estimate  Std.Err  z-value  P(>|z|)  Std.lv  Std.all
##      .lwage      0.259    0.098    2.636    0.008    0.259    1.165
##      .educ       4.963    0.149   33.317    0.000    4.963    0.742
##

```

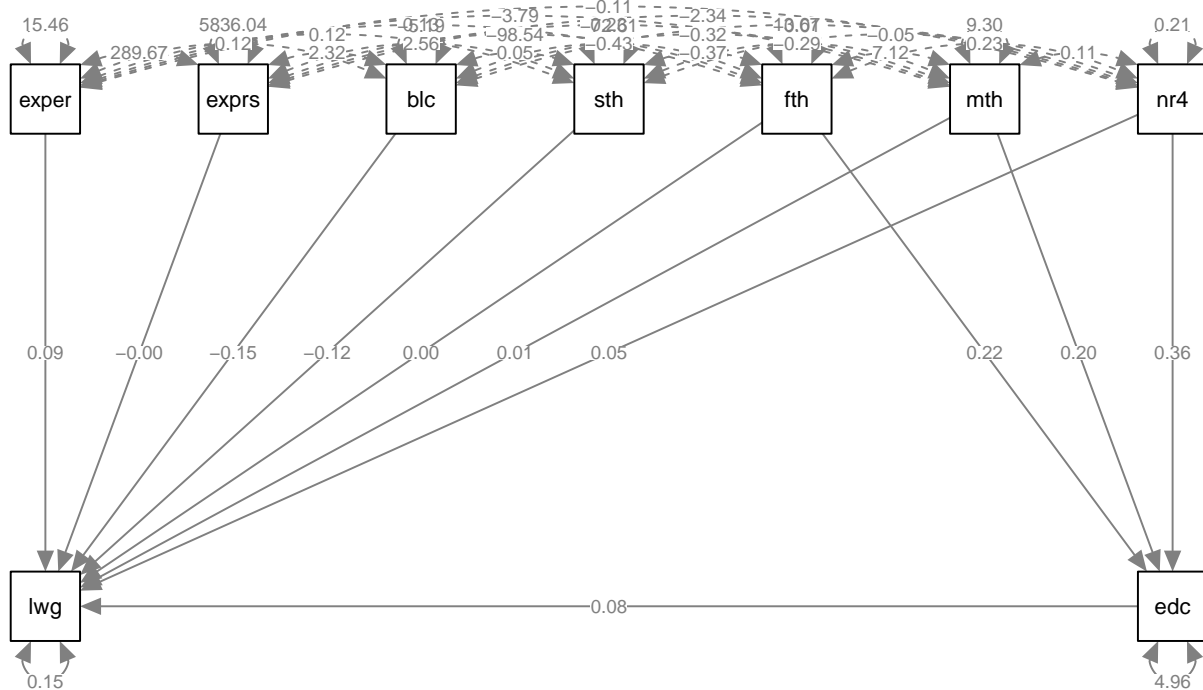
```
## R-Square:
##           Estimate
##    lwage    -0.165
##    educ     0.258
```

- d. Evaluate whether `nearc4` is a weak or strong instrument for dealing with the endogeneity of the variable `educ`. Specifically, consider the criteria that a strong instrument must meet in order to adequately correct for endogeneity.

```
model_ex_2_d <- "
  lwage ~ educ + exper + expersq + black + south + fatheduc + motheduc + nearc4
  educ ~ nearc4 + fatheduc + motheduc
"

# Fit the model.
model_ex_2_d_fit <- sem(model_ex_2_d, data_ex_2)

# Visualize the model.
semPaths(model_ex_2_d_fit, what = "paths", whatLabels = "est")
```



```
# Model summary.
summary(model_ex_2_d_fit, standardized = TRUE, rsquare = TRUE)
```

```
## lavaan 0.6-12 ended normally after 1 iterations
##
## Estimator ML
## Optimization method NLMINB
## Number of model parameters 13
##
## Used Total
## Number of observations 2220 3010
```

```

##
## Model Test User Model:
##
##   Test statistic              790.934
##   Degrees of freedom          4
##   P-value (Chi-square)        0.000
##
## Parameter Estimates:
##
##   Standard errors              Standard
##   Information                  Expected
##   Information saturated (h1) model Structured
##
## Regressions:
##           Estimate Std.Err z-value P(>|z|) Std.lv Std.all
##   lwage ~
##     educ      0.075   0.004  20.720   0.000   0.075   0.414
##     exper      0.089   0.008  11.321   0.000   0.089   0.741
##     expersq    -0.002   0.000  -5.959   0.000  -0.002  -0.389
##     black     -0.147   0.024  -6.036   0.000  -0.147  -0.114
##     south     -0.121   0.018  -6.666   0.000  -0.121  -0.124
##     fatheduc    0.001   0.003   0.416   0.678   0.001   0.010
##     motheduc    0.008   0.004   2.140   0.032   0.008   0.049
##     nearc4     0.055   0.018   3.027   0.002   0.055   0.054
##   educ ~
##     nearc4     0.364   0.103   3.526   0.000   0.364   0.065
##     fatheduc    0.216   0.017  13.030   0.000   0.216   0.309
##     motheduc    0.203   0.020  10.154   0.000   0.203   0.239
##
## Variances:
##           Estimate Std.Err z-value P(>|z|) Std.lv Std.all
##   .lwage      0.146   0.004  33.317   0.000   0.146   0.658
##   .educ       4.963   0.149  33.317   0.000   4.963   0.742
##
## R-Square:
##           Estimate
##   lwage      0.342
##   educ       0.258

```

Overall, `nearc4` does not appear to be a particularly strong instrumental variable!

## References

Brown, T. A. (2014). *Confirmatory factor analysis for applied research*. Guilford Publications.