

A first R analysis

Guido Biele

16 november 2018

Introduction

This document aims to provide a short introduction into statistical analysis with R. The assumption is, that the reader is already familiar with basic concepts in R and the RStudio IDE, and has some basic knowledge about statistics, for example from using SPSS. The goal here is to show how one can

- open an SPSS file
- merge separate data sets
- visually inspect the data
- do a simple data reduction through an exploratory factor analysis
- estimate the association between an exposure and an outcome with a linear regression
- visualize the results of this regression.

Packages

The name R refers at the same time to 3 things:

1. A programming language
2. A software for statistical analyses
3. An ecosystem of *packages* that implement different statistical analyses.

R packages are written (typically by using R, the programming language) by researchers and available for free. These packages are crucial, because the capabilities of the basic R software are limited.

There are thousands of R packages. Therefore, CRAN task view web sites are useful to give an overview about R packages for different application areas. CRAN stands for “Comprehensive R Archive Network”.

For this introduction we will use following packages:

- `haven` allows to read SPSS (or stata) data into R
- `ggplot` facilitates plotting
- `psych` is a package for factor analysis and related things
- `grid` and `cowplot` for general plotting
- `sjPlot` and `ggeffect` for plotting of regression analysis results

Installing packages is easy!

```
install.packages("haven")
install.packages("ggplot2")
install.packages("psych")
install.packages("sjplot")
install.packages("VIM")
```

One can also install multiple packages simultaneously as follows `install.packages(c("haven", "psych"))`.

Loading data

The problem we are looking at is the association between Parenting Style and child mental health symptoms. More specifically, we look at the association between self reported Parenting Style of well educated Norwegian

mothers at child age 5 years and oppositional behavior as well as mood and feeling at age 8. This data can be obtained from moba questionnaires 5 and 8.

If these data are available in SPSS (or stata) format, we can read them with the `haven` package, which we first have to load before we can use it:

```
library(haven)
library(VIM)
```

All packages come with a documentation, and if one is lucky also with one or more *vignettes*, which are tutorial-style documents that show how the functions in the package can be used. This and more information about the packages can be found at <https://cran.rstudio.com/web/packages/haven/>. For other packages, just replace “haven” with the name of the package you are interested in.

OK, now we can load the data:

```
data_directory = "data/"
Q5aar = read_sav(paste0(data_directory,
                         "PDB439_Skjema5aar_v10.sav"))
Q5aar = data.frame(Q5aar)
```

So what have we done here? We first created a new variable `data_directory` and set it to the path to the directory where the data files are. Then, when loading the data into the `data.frame` `Q5aar`, we collated the path to the directory with the name of the data file by using the `paste0` command.

The reason for doing this is that the code becomes a bit easier to read (still not great). For the same reason, I made a line break at the end of `Q5aar = read_sav(paste0(data_directory,,`.

Note that if we would already be in the same directory as the data file, we could just have written `Q5aar = read_sav("PDB439_Skjema5aar_v10.sav")`. You can use the commands `getwd()` and `setwd()` to check and set your working directory (wd).

Next we select the variables we need. This are the pregnancy ID and the child number, as well as the items for measuring parental style (we get SPSS variable names from here).

```
Q5aar = Q5aar[,c("PREG_ID_439", "BARN_NR",
                 paste0("LL", 526:534))]
names(Q5aar)

## [1] "PREG_ID_439"  "BARN_NR"       "LL526"        "LL527"        "LL528"
## [6] "LL529"        "LL530"        "LL531"        "LL532"        "LL533"
## [11] "LL534"
```

Ok, these variable names are not very clear, let's rename them:

```
names(Q5aar)[3:11] = paste0("PS.item.", 1:9)
names(Q5aar)

## [1] "PREG_ID_439"  "BARN_NR"       "PS.item.1"     "PS.item.2"     "PS.item.3"
## [6] "PS.item.4"     "PS.item.5"     "PS.item.6"     "PS.item.7"     "PS.item.8"
## [11] "PS.item.9"
```

It is important to understand the coding of the variables. The `haven` package makes SPSS' *Variable label* and *Value label* information available as attributes of the variables in the `data.frame`. We can access those as follows:

```
attr(Q5aar$PS.item.1, "label")
```

```
## [1] "W_53_1:SKJEMA_5AARB; Du lar barnet ditt forstå når han/hun gjør en god innsats; 53. Hvor ofte h
```

```

## Mer enn ett kryss          Aldri      Nesten aldri      Noen ganger
##          0                  1             2                  3
##          Ofte                Alltid
##          4                  5

```

Here we see one annoying detail of SPSS Moba data: If the mother has ticked more than one box, the data are for our purposes missing, and not zero. We need to fix that.

To do this, we do a quick for loop:

```

item_names = names(Q5aar)[-c(1,2)]
item_quest = c()
for (v in item_names) {
  item_quest = c(item_quest,
                 strsplit(attr(Q5aar[,v],"label"),
                          split = ";")[[1]][2])
  attr(Q5aar[, v],"labels") = attr(Q5aar[, v],"labels")[-1]
  # need to use the "which", otherwise R stumbles over missing values
  value_is_0 = which(Q5aar[,v] == 0)
  Q5aar[value_is_0, v] = NA
}
table(as_factor(Q5aar$PS.item.1))

##
##          Aldri  Nesten aldri  Noen ganger          Ofte      Alltid
##          9          7          359        13795      12110

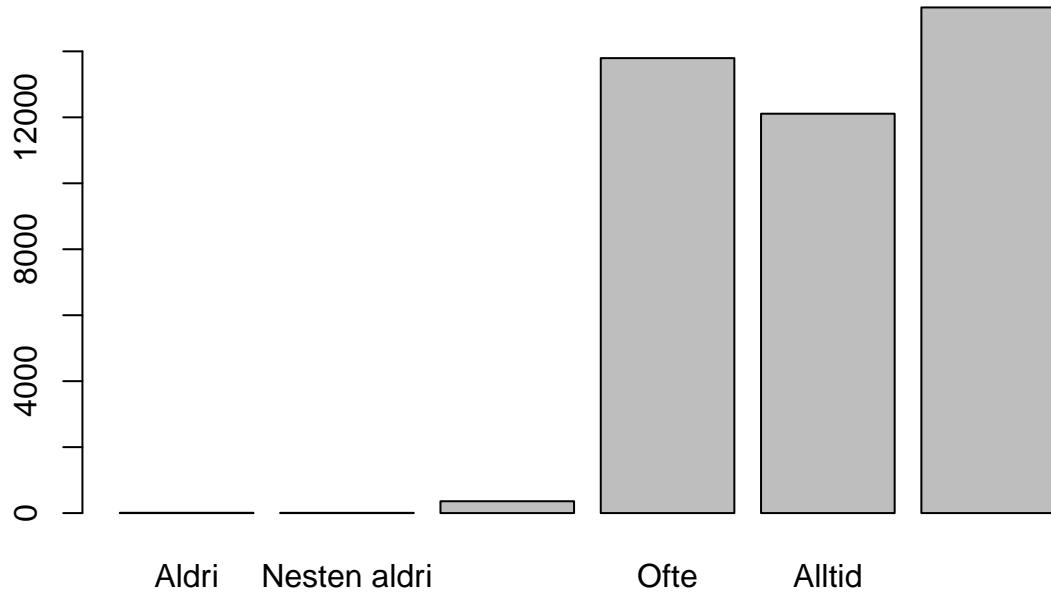
```

Missing data is always a problem, especially with questionnaire data from long lasting studies such as MoBa. Lets look at missing data by showing a table that counts them:

```

table_plus_missing = table(addNA(as_factor(Q5aar$PS.item.1)))
barplot(table_plus_missing)

```



```
table_plus_missing
```

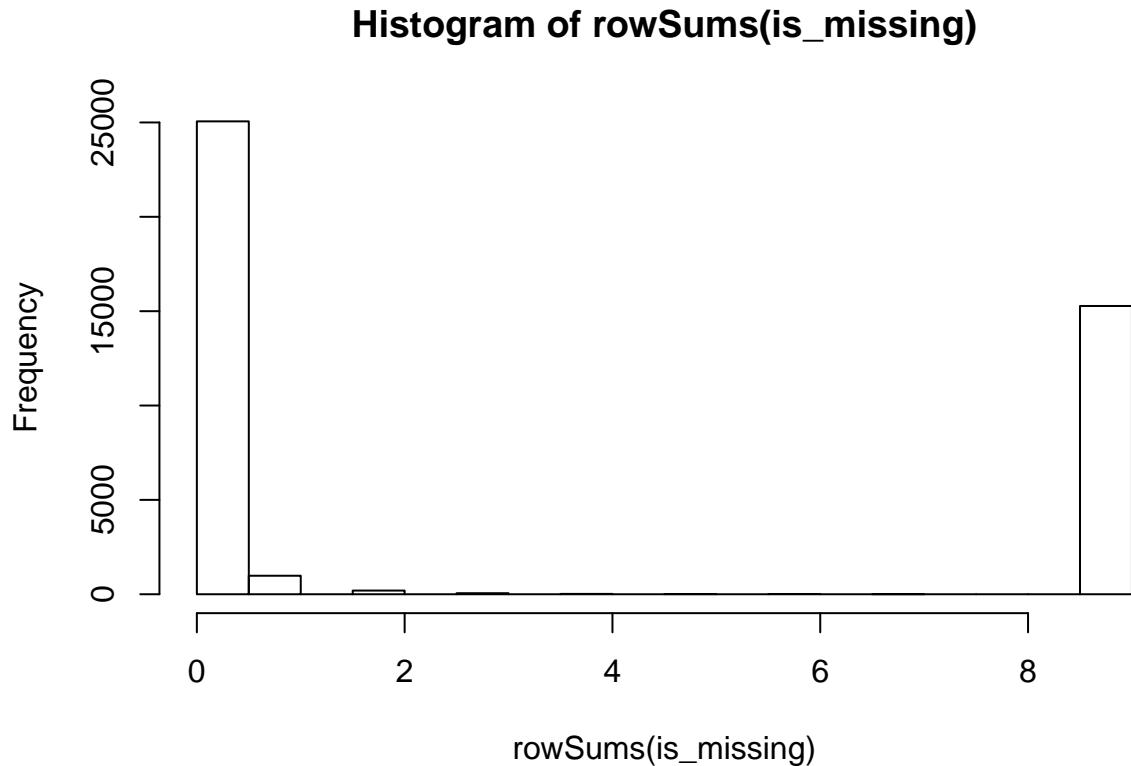
```
##
```

	Aldri	Nesten aldri	Noen ganger	Ofte	Alltid
##	9	7	359	13795	12110
##	<NA>				
##	15332				

There is lots of missing data, probably because MoBa changed the Parenting Style items throughout.

Lets look a little closer at the missing data:

```
is_missing = is.na(Q5aar[,-c(1,2)])
hist(rowSums(is_missing))
```



Indeed there is a substantial number of cases for which als Parenting style items are missing. Let's remove those participants where more than 50% of parenting style items are missing, and impute the missing values for the remaining participants. We use the function `hotdeck` from the package `VIM` to impute values

```
library(VIM)
Q5aar = Q5aar[rowMeans(is_missing[,item_names]) < .5,]
ps_imputations = hotdeck(as.matrix(Q5aar[,item_names]))
Q5aar[,item_names] = ps_imputations[,item_names]
```

A quick exploratory factor analysis of parenting style items in MoBa

Our goal is to reduce the 9 items to a smaller set of variables describing parenting style. One way to do this is to use exploratory factor analysis (EFA) to learn the factor structure of the data. As a first step we'll check how many factors we need, using the `fa.parallel.poly` command from the `psych` package.

As with many things in R, using up to date methods is not hard, as you'll see. The problem is more in knowing that these methods exist. If you do not know which method to use, try googling it. The challenge is in finding out which of the many similar packages one should use. Here are a few tips:

- google R tutorial exploratory factor analysis,
- If a paper about a package is published in the Journal of Statistical Software, this is probably one of the better packages. dOn't worry about the name of the Journal! The articles typically include a tutorial about how to use the key functions.

Now lets move on with the exploratory factor analysis. One thing we need to specify is that polychoric correlations are used, because we have ordinal data.

```

library(psych)
describe(Q5aar[,item_names])

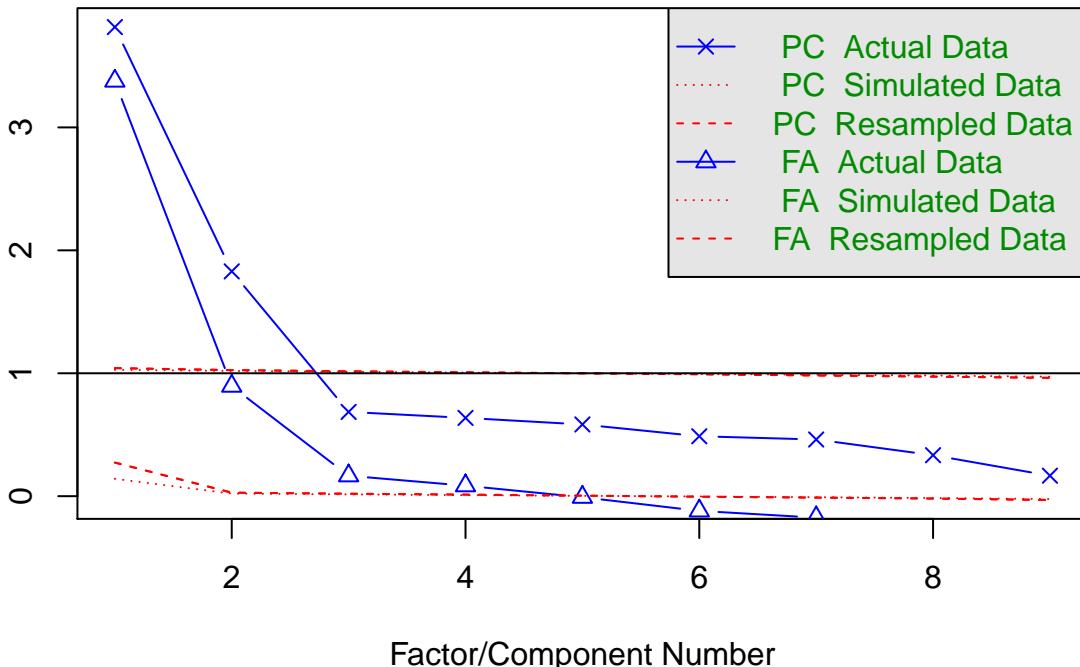
##          vars      n  mean    sd median trimmed  mad min max range skew
## PS.item.1    1 26307 4.45 0.53      4    4.45 0.00  1   5    4 -0.18
## PS.item.2    2 26307 2.13 0.78      2    2.13 1.48  1   5    4  0.19
## PS.item.3    3 26307 4.19 0.45      4    4.14 0.00  1   5    4  0.54
## PS.item.4    4 26307 2.04 0.84      2    2.01 1.48  1   5    4  0.36
## PS.item.5    5 26307 4.77 0.46      5    4.85 0.00  1   5    4 -2.00
## PS.item.6    6 26307 4.65 0.49      5    4.70 0.00  1   5    4 -0.79
## PS.item.7    7 26307 4.47 0.57      5    4.51 0.00  1   5    4 -0.56
## PS.item.8    8 26307 4.22 0.59      4    4.26 0.00  1   5    4 -0.26
## PS.item.9    9 26307 2.33 0.79      2    2.36 1.48  1   5    4 -0.01
##          kurtosis   se
## PS.item.1   -0.68 0.00
## PS.item.2   -0.35 0.00
## PS.item.3    1.52 0.00
## PS.item.4   -0.40 0.01
## PS.item.5    5.15 0.00
## PS.item.6   -0.84 0.00
## PS.item.7   -0.31 0.00
## PS.item.8    0.38 0.00
## PS.item.9   -0.30 0.00

fa.parallel(Q5aar[,item_names], cor = "poly")

```

eigenvalues of principal components and factor analysis

Parallel Analysis Scree Plots



```

## Parallel analysis suggests that the number of factors = 4 and the number of components = 2

```

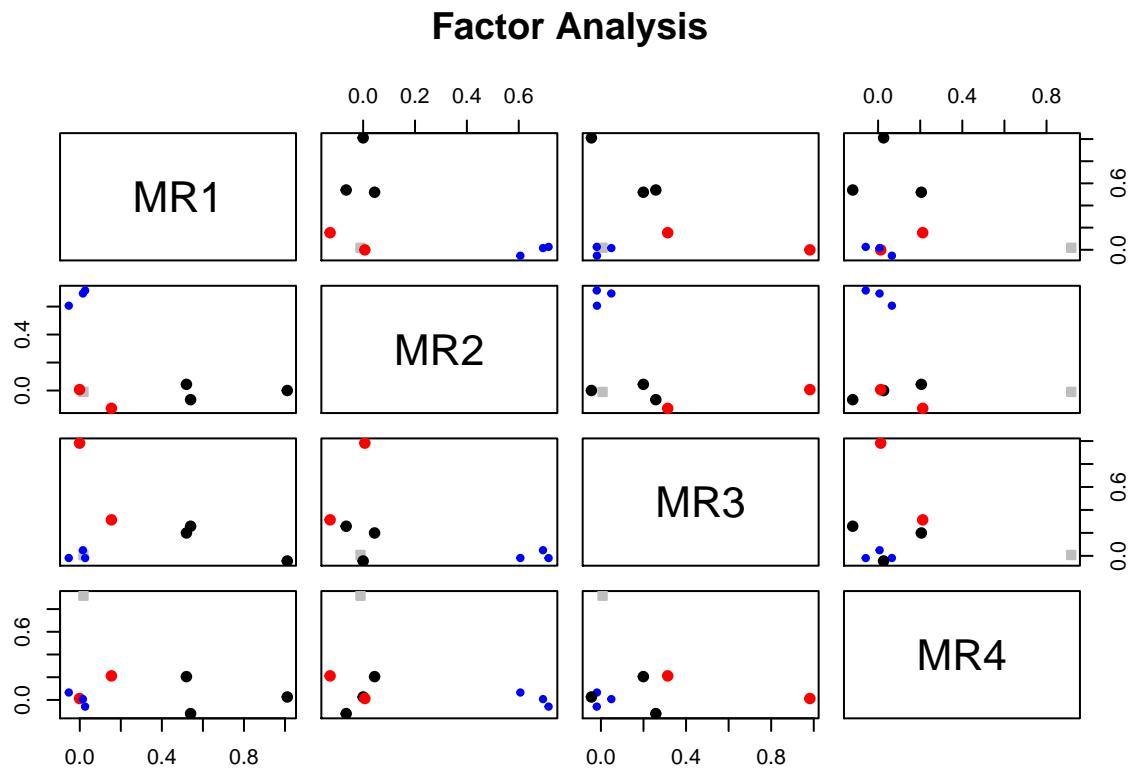
4 factors. I had hoped fewer, but lets continue and do an exploratory factor analysis.

```
fa_ps = fa(Q5aar[,item_names],  
            nfactors = 4,  
            cor = "poly")  
fa_ps  
  
## Factor Analysis using method = minres  
## Call: fa(r = Q5aar[, item_names], nfactors = 4, cor = "poly")  
##  
## Warning: A Heywood case was detected.  
## Standardized loadings (pattern matrix) based upon correlation matrix  
##  
##          MR1   MR2   MR3   MR4   h2      u2 com  
## PS.item.1  0.02 -0.01  0.01  0.92  0.88  0.1243 1.0  
## PS.item.2  0.03  0.72 -0.02 -0.06  0.53  0.4717 1.0  
## PS.item.3  0.15 -0.13  0.31  0.21  0.38  0.6244 2.7  
## PS.item.4 -0.05  0.61 -0.02  0.07  0.37  0.6343 1.0  
## PS.item.5  0.54 -0.07  0.26 -0.12  0.44  0.5624 1.6  
## PS.item.6  1.01  0.00 -0.04  0.03  1.01 -0.0053 1.0  
## PS.item.7  0.52  0.04  0.20  0.21  0.66  0.3431 1.6  
## PS.item.8  0.00  0.01  0.98  0.01  0.97  0.0267 1.0  
## PS.item.9  0.02  0.69  0.05  0.01  0.47  0.5258 1.0  
##  
##          MR1   MR2   MR3   MR4  
## SS loadings      1.85  1.39  1.37  1.08  
## Proportion Var   0.21  0.15  0.15  0.12  
## Cumulative Var   0.21  0.36  0.51  0.63  
## Proportion Explained  0.32  0.24  0.24  0.19  
## Cumulative Proportion 0.32  0.57  0.81  1.00  
##  
## With factor correlations of  
##  
##          MR1   MR2   MR3   MR4  
## MR1  1.00 -0.15  0.61  0.70  
## MR2 -0.15  1.00 -0.08 -0.21  
## MR3  0.61 -0.08  1.00  0.47  
## MR4  0.70 -0.21  0.47  1.00  
##  
## Mean item complexity =  1.3  
## Test of the hypothesis that 4 factors are sufficient.  
##  
## The degrees of freedom for the null model are 36 and the objective function was 3.81 with Chi Squa  
## The degrees of freedom for the model are 6 and the objective function was 0.02  
##  
## The root mean square of the residuals (RMSR) is  0.01  
## The df corrected root mean square of the residuals is  0.02  
##  
## The harmonic number of observations is 26307 with the empirical chi square 169.53 with prob < 5.  
## The total number of observations was 26307 with Likelihood Chi Square = 598.16 with prob < 5.8e-05  
##  
## Tucker Lewis Index of factoring reliability = 0.964  
## RMSEA index = 0.061 and the 90 % confidence intervals are 0.057 0.065  
## BIC = 537.09  
## Fit based upon off diagonal values = 1
```

We can at first look at the RMSEA. The value is 0.06, which is not great but OK.

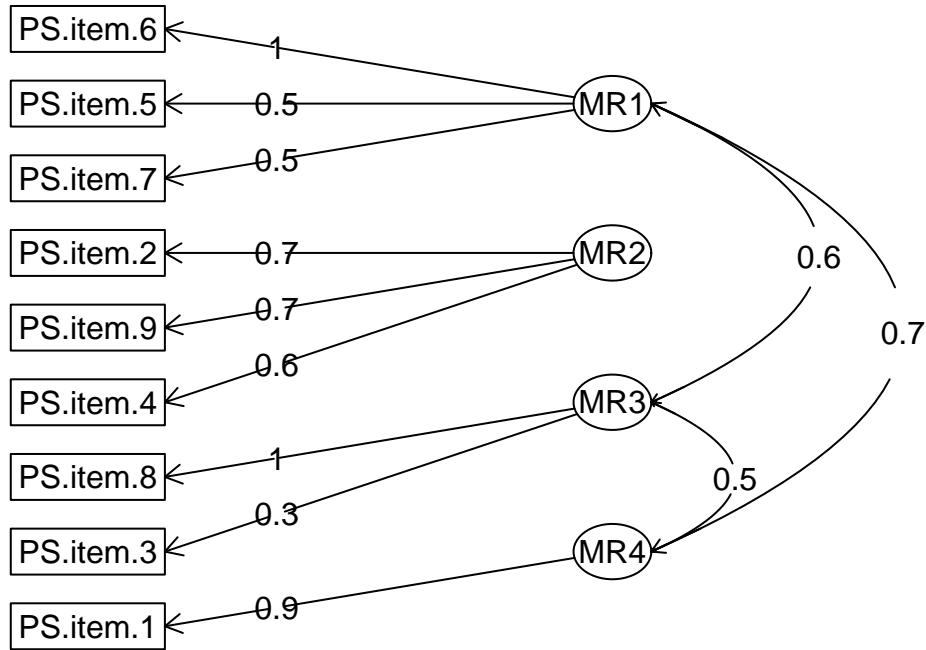
This are the factor loadings:

```
plot(fa_ps)
```



```
fa.diagram(fa_ps)
```

Factor Analysis

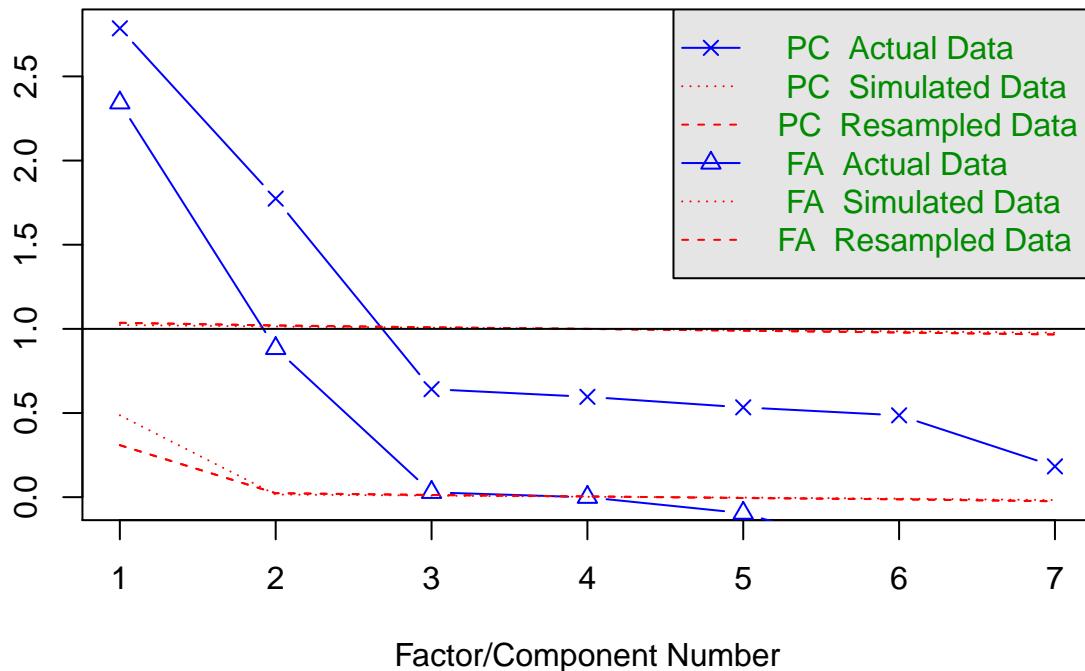


This is not very clear, because a number of items (especially item 3) have high loadings for multiple factors. Maybe even more problematically, items 1 and 8 stand apart. Lets try to redo this without these items.

```
reduced_items = item_names[-c(1,8)]  
fa.parallel(Q5aar[,reduced_items], cor = "poly")
```

eigenvalues of principal components and factor analysis

Parallel Analysis Scree Plots



```
## Parallel analysis suggests that the number of factors = 3 and the number of components = 2
OK, now we only need 3 or maybe only 2 factors. Lets first try with 2.
```

```
fa_ps = fa(Q5aar[,reduced_items],
            nfactors = 2,
            cor = "poly")
fa_ps

## Factor Analysis using method = minres
## Call: fa(r = Q5aar[, reduced_items], nfactors = 2, cor = "poly")
## Standardized loadings (pattern matrix) based upon correlation matrix
##          MR1    MR2    h2    u2 com
## PS.item.2 -0.04  0.72  0.53  0.471 1.0
## PS.item.3  0.53 -0.13  0.32  0.676 1.1
## PS.item.4 -0.02  0.60  0.36  0.639 1.0
## PS.item.5  0.63 -0.04  0.41  0.591 1.0
## PS.item.6  0.97  0.01  0.93  0.065 1.0
## PS.item.7  0.81  0.04  0.65  0.353 1.0
## PS.item.9  0.06  0.69  0.47  0.528 1.0
##
##          MR1    MR2
## SS loadings     2.29 1.39
## Proportion Var  0.33 0.20
## Cumulative Var 0.33 0.53
## Proportion Explained 0.62 0.38
## Cumulative Proportion 0.62 1.00
```

```

## 
## With factor correlations of
##      MR1    MR2
## MR1  1.00 -0.16
## MR2 -0.16  1.00
##
## Mean item complexity =  1
## Test of the hypothesis that 2 factors are sufficient.
##
## The degrees of freedom for the null model are  21  and the objective function was  2.41 with Chi Squa
## The degrees of freedom for the model are 8  and the objective function was  0.03
##
## The root mean square of the residuals (RMSR) is  0.02
## The df corrected root mean square of the residuals is  0.03
##
## The harmonic number of observations is  26307 with the empirical chi square  284.24 with prob <  9.8
## The total number of observations was  26307  with Likelihood Chi Square =  767.5 with prob <  2.1e-1
##
## Tucker Lewis Index of factoring reliability =  0.969
## RMSEA index =  0.06  and the 90 % confidence intervals are  0.057 0.064
## BIC =  686.07
## Fit based upon off diagonal values = 1
## Measures of factor score adequacy
##                                     MR1    MR2
## Correlation of (regression) scores with factors  0.97  0.85
## Multiple R square of scores with factors        0.95  0.72
## Minimum correlation of possible factor scores  0.89  0.45

```

RMSEA and even more the TLI suggest that 2 factors should be OK.

This are the new factor loadings:

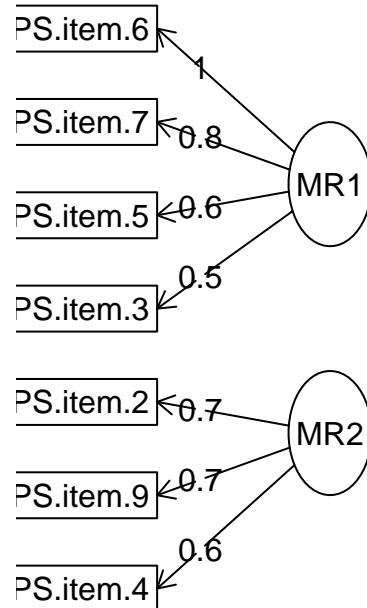
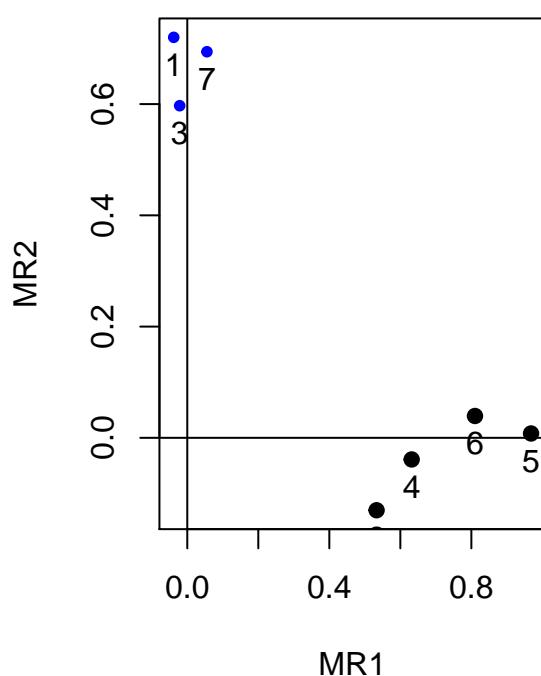
```

par(mfrow = c(1,2))
plot(fa_ps)
fa.diagram(fa_ps)

```

Factor Analysis

Factor Analysis



This looks like a clear factor structure. Lets look at the items to understand what those factor might measure.

```

factor_loadings = loadings(fa_ps)
simple_loadings = apply(factor_loadings,
                       1,
                       function(x)
                         (abs(x) == max(abs(x))))
item_quest_reduced = item_quest[-c(1,8)]

items_factor1 = which(simple_loadings[1,] == T)
items_factor2 = which(simple_loadings[2,] == T)

cat(paste0("Factor1:\n",
           paste(item_quest_reduced[items_factor1],
                 collapse = "\n"),
           "\n"))

## Factor1:
## Du har en hyggelig samtale med barnet ditt
## Du spør barnet ditt om hvordan hans/hennes dag i barnehagen har vært
## Du sier noe pent til barnet ditt når han/hun har gjort noe bra
## Du roser barnet ditt om han/hun oppfører seg pent

cat(paste0("Factor2:\n",
           paste(item_quest_reduced[items_factor1],
                 collapse = "\n"),
           "\n"))
  
```

```
## Factor2:
## Du har en hyggelig samtale med barnet ditt
## Du spør barnet ditt om hvordan hans/hennes dag i barnehagen har vært
## Du sier noe pent til barnet ditt når han/hun har gjort noe bra
## Du roser barnet ditt om han/hun oppfører seg pent
```

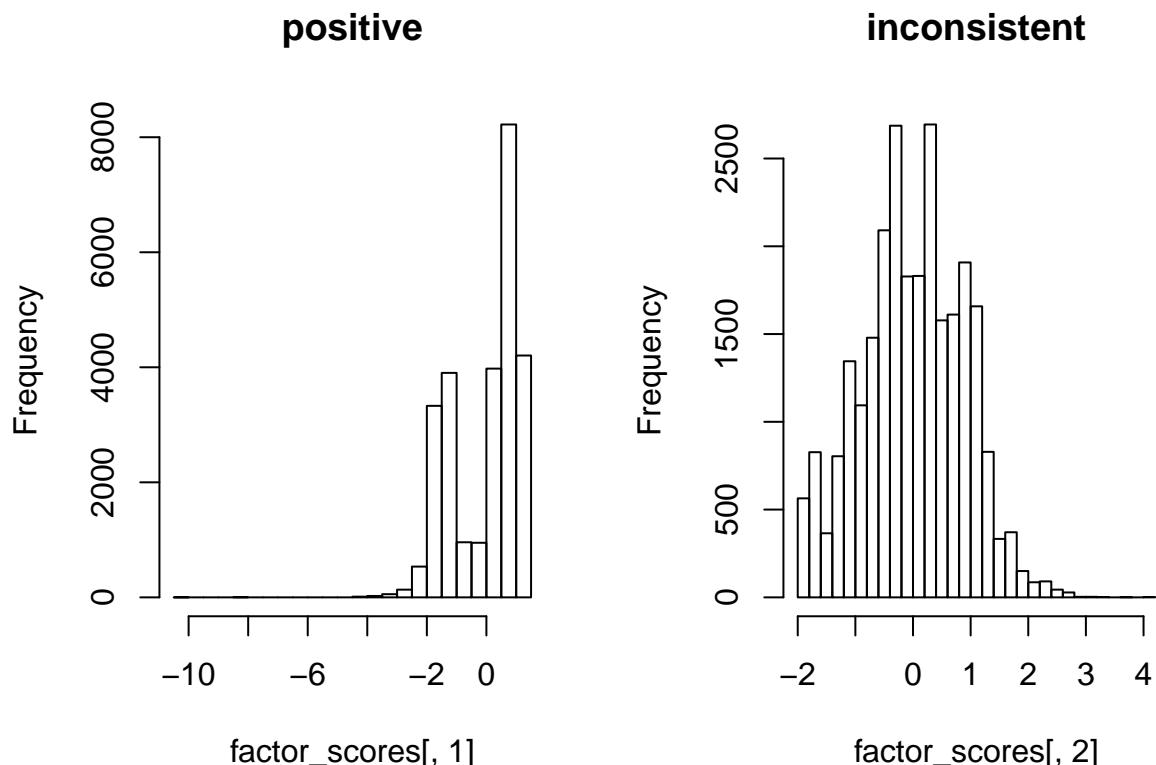
This looks as if there is one factor around positive reinforcement and communication with the child. Lets call this *positive parenting*. The other factor appears to be about the inability to follow through with threats of punishment, lets call this *inconsistent parenting* ;-). If we scroll back and look at correlation of the two factors, it looks as if *positive parenting* has a weak negative correlation with *inconsistent parenting*.

Now, if we want to use the parenting dimensions as predictors, we need to extract the factor scores. Factor scores are likely stored in the `fa_ps` object, which has the results of our exploratory factor analysis. If we want to know what this object contains, we can simply click on it in the “Environment” tab at the right hand side of the RStudio IDE.

```
factor_scores = fa_ps$scores
colnames(factor_scores) = c("positive", "inconsistent")
```

Lets quickly look at the factor scores, before we add them to the data.frame with the 5 year data.:

```
par(mfrow = c(1,2))
hist(factor_scores[,1],
     breaks = 25,
     main = "positive")
hist(factor_scores[,2],
     breaks = 25,
     main = "inconsistent")
```



```
Q5aar = cbind(Q5aar,factor_scores)
```

If we had a great parenting scale, the histograms should look normally distributed. Especially the histogram for positive parenting deviates from that. It shows that (a) there are a few extreme outliers to the left and (b) there seems to be something like a ceiling effect. That is, a good number of mothers report to be **VERY** positive in their parenting. This is another instance where it would be great to have a social desirability scale in MoBa.

Some more data wrangling, including calculation of sum scores

Anyhow, let's continue and put together outcome data. For this, we load the MoBa 8 years data file.

```
data_directory = "data/"
Q8aar = read_sav(paste0(data_directory,
                         "PDB439_Skjema8aar_v10.sav"))
Q8aar = Q8aar[,c("PREG_ID_439","BARN_NR",
                 paste("NN",c(68:80),      # SMFQ
                       211:226,      # CCC-2 sort
                       227:233, 374, # Språkk 20
                       111:118,      # CD
                       119:136,      # ADHD
                       137:144),      # OD
                 sep = ""))]
```

Q8aar = data.frame(Q8aar)

Now we are coming to a part, which always takes more time than we would wish: cleaning the data. We want to calculate sum-scores, but this is complicated by two facts: MoBa data come with

- 0 instead of NA for data with unambiguous responses
- all scales start at 1
- some of the data are missing.

We will address these problems in that order, starting with renaming some variables. To make this easier, we'll create simple functions where this is useful.

First we rename all variables:

```
rename_variables = function(df,old_names,new_names) {
  names(df)[names(df) %in% old_names] = new_names
  return(df)
}

Q8aar = rename_variables(Q8aar,
                         paste0("NN",68:80),
                         paste0("SMFQ.i",1:13,"_8y"))
Q8aar = rename_variables(Q8aar,
                         paste0("NN",c(227:233,374)),
                         paste0("SPRAAK20.i",1:8,"_8y"))
Q8aar = rename_variables(Q8aar,
                         paste0("NN",211:226),
                         paste0("CCCs.i",1:16,"_8y"))
for (item in paste0("CCCs.i",c(10:16),"_8y"))
  Q8aar[,item] = abs(Q8aar[,item]-5)
Q8aar = rename_variables(Q8aar,
                         paste0("NN",111:118),
```

```

            paste0("CD.i",1:8,"_8y"))
Q8aar = rename_variables(Q8aar,
                         paste0("NN",119:136),
                         paste0("ADHD.i",1:18,"_8y"))
Q8aar = rename_variables(Q8aar,
                         paste0("NN",137:144),
                         paste0("OD.i",1:8,"_8y"))
Q8aar = rename_variables(Q8aar,
                         paste0("NN",145:149),
                         paste0("SCARED.i",1:5,"_8y"))

```

Next we replace all 0 values with NA and subtract 1 from all items. We use the `grep` command, which makes it easy for us to find all variables that match a particular pattern. We also temporarily create a new `data.frame` `tmp_items`, in which we keep the item data.

```

all_items = grep("\\.i[0-9]",names(Q8aar), value = T)
tmp_items = Q8aar[,all_items]
tmp_items[tmp_items == 0] = NA
tmp_items = tmp_items-1

```

We can use imputation to replace missing data. I prefer to do this only for cases, for which we have at least 50% of the data. So we remove participants with more than 50% missing data.

```

proportion_missing = rowMeans(is.na(tmp_items))
Q8aar = Q8aar[proportion_missing < .5,]
tmp_items = tmp_items[proportion_missing < .5,]

```

Now we use the function `hotdeck` from the package `VIM` to impute values. (`hotdeck` is fast). Other approaches like nearest neigbour imputation of multiple chained imputation are more accurate, by they are to slow when analysing bigger data-sets in a turorial.

```

tmp_items_imputed = hotdeck(as.matrix(tmp_items))
tmp_items_imputed = tmp_items_imputed[,1:ncol(tmp_items)]

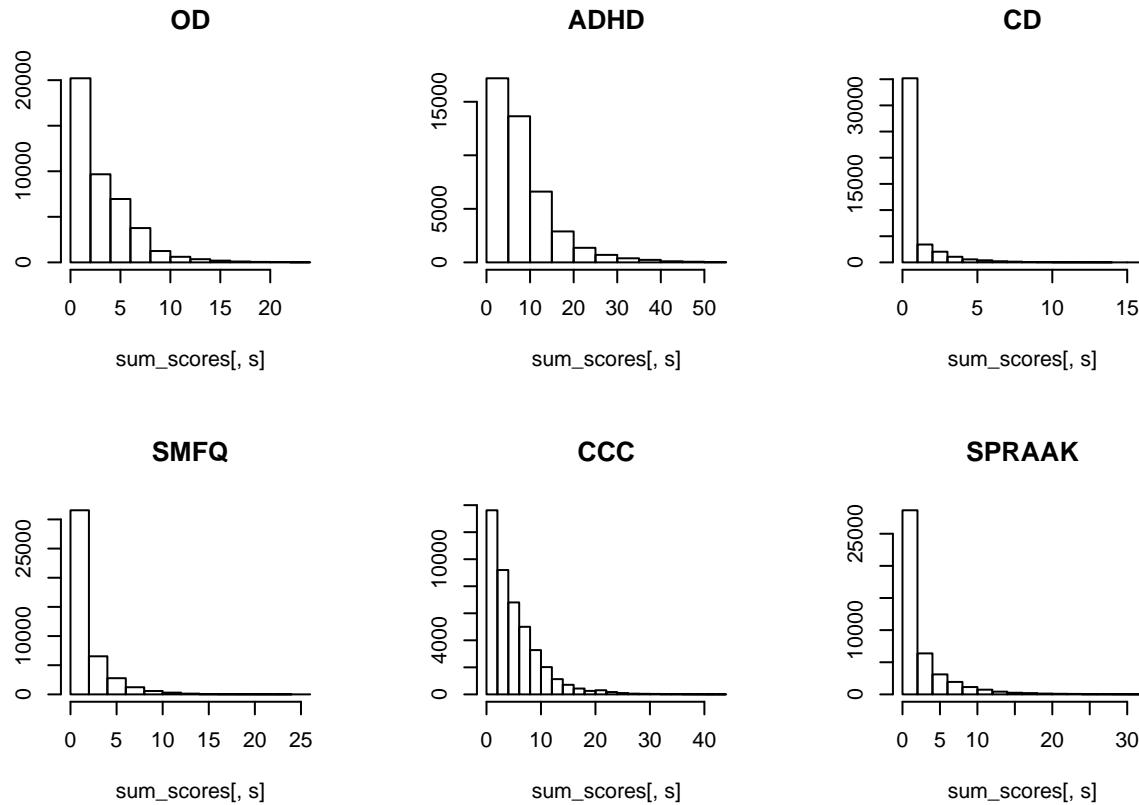
```

Finally, we can calculate sum scores.

```

scales = c("OD","ADHD","CD","SMFQ","CCC","SPRAAK")
sum_scores = matrix(NA,
                    nrow = nrow(Q8aar),
                    ncol = length(scales))
colnames(sum_scores) = scales
par(mfrow = c(2,3))
for (s in scales) {
  items = grep(s,names(tmp_items_imputed), value = T)
  sum_scores[,s] = rowSums(tmp_items_imputed[,items])
  hist(sum_scores[,s],
        main = s,
        ylab = "")
}

```



```
head(sum_scores)
```

```
##      OD ADHD CD SMFQ CCC SPRAAK
## [1,]  2   5   0   1   4   0
## [2,]  1  11   0   2   2   1
## [3,]  2   7   0   0   3   0
## [4,]  0   0   0   1   1   0
## [5,] 12  16   2   5  18  13
## [6,]  6   9   0   0   5   0
```

This looks more or less as expected. Now we put the Q8aar data together, and merge 5 and 8 years data.

```
Q8aar = cbind(Q8aar, sum_scores)
my_data = merge(Q5aar[,c("PREG_ID_439","BARN_NR",colnames(factor_scores))],
                Q8aar[,c("PREG_ID_439","BARN_NR",scales)],
                by = c("PREG_ID_439","BARN_NR"))
head(my_data)
```

```
##   PREG_ID_439 BARN_NR   positive inconsistent OD ADHD CD SMFQ CCC SPRAAK
## 1 100002      1 -0.9721485  -0.6293652  3   4   0   1   9   4
## 2 100003      1  0.7377795  -0.5939807  4   3   0   1   1   0
## 3 100004      1 -1.5494872  -0.6626380  2   5   0   1  10   1
## 4 100006      1  0.7377795  -0.5939807  3   6   0   1   1   0
## 5 100007      1  1.2419478 -1.1441941  1   0   1   0   1   0
## 6 100009      1 -1.5600513  0.3932045  3   8   1   3  11   7
```

```
my_data$positive = as.numeric(scale(my_data$positive))
my_data$inconsistent = as.numeric(scale(my_data$inconsistent))
```

```
my_data = my_data[my_data$positive > -4, ]
```

Linear (and other) regressions in R

For the next section, you can run your own regression by changing the outcome. Use the command `names(my_data)` to see, which variables are available.

Finally, we can do simple linear regression models:

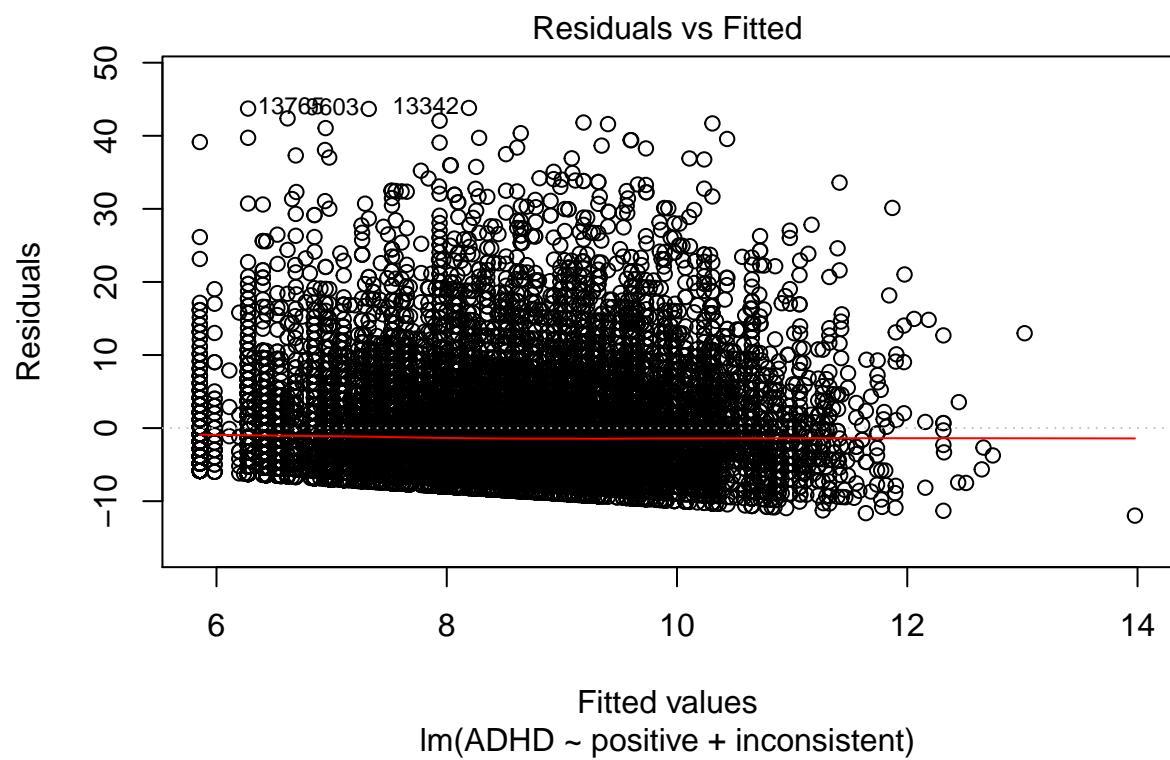
```
ADHD_model = lm(ADHD ~ positive + inconsistent, my_data)
summary(ADHD_model)
```

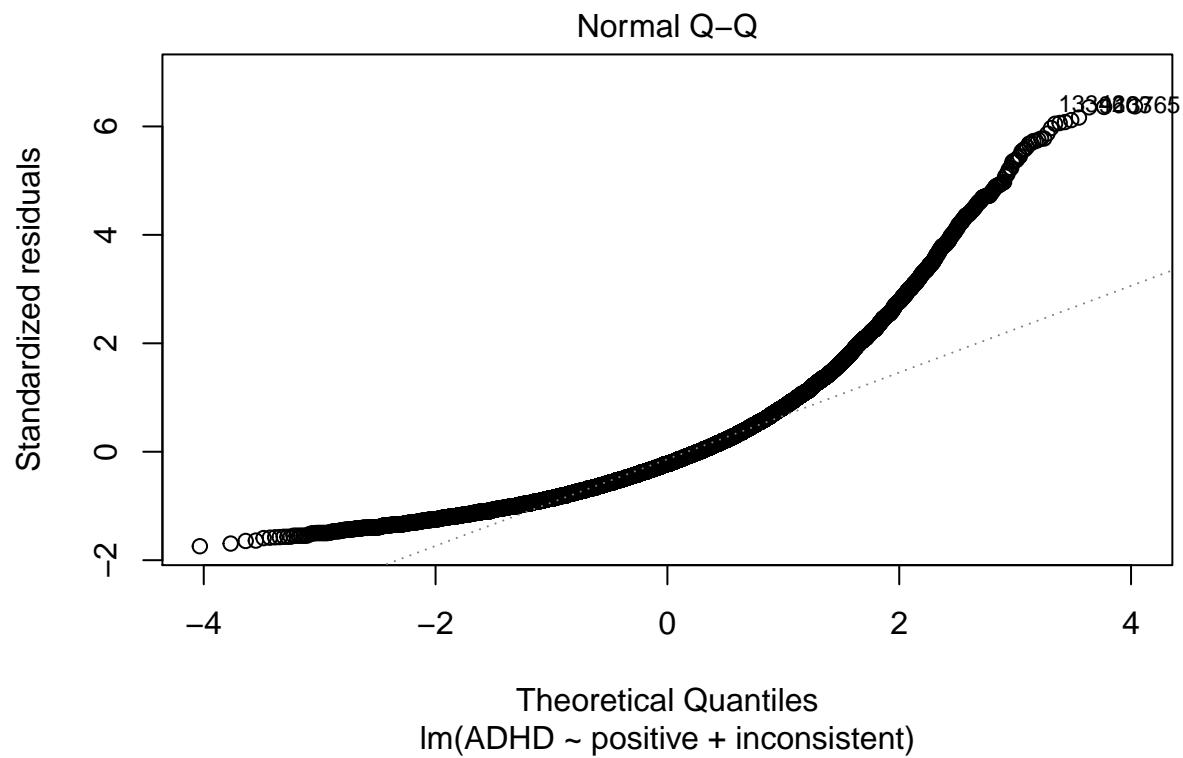
```
##
## Call:
## lm(formula = ADHD ~ positive + inconsistent, data = my_data)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -11.978  -4.652  -1.529   2.764  43.808
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 8.48847   0.05097 166.534 < 2e-16 ***
## positive   -0.40120   0.05290  -7.584 3.51e-14 ***
## inconsistent 1.04253   0.05249  19.862 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.877 on 18200 degrees of freedom
## Multiple R-squared:  0.02958,    Adjusted R-squared:  0.02947
## F-statistic: 277.4 on 2 and 18200 DF,  p-value: < 2.2e-16
```

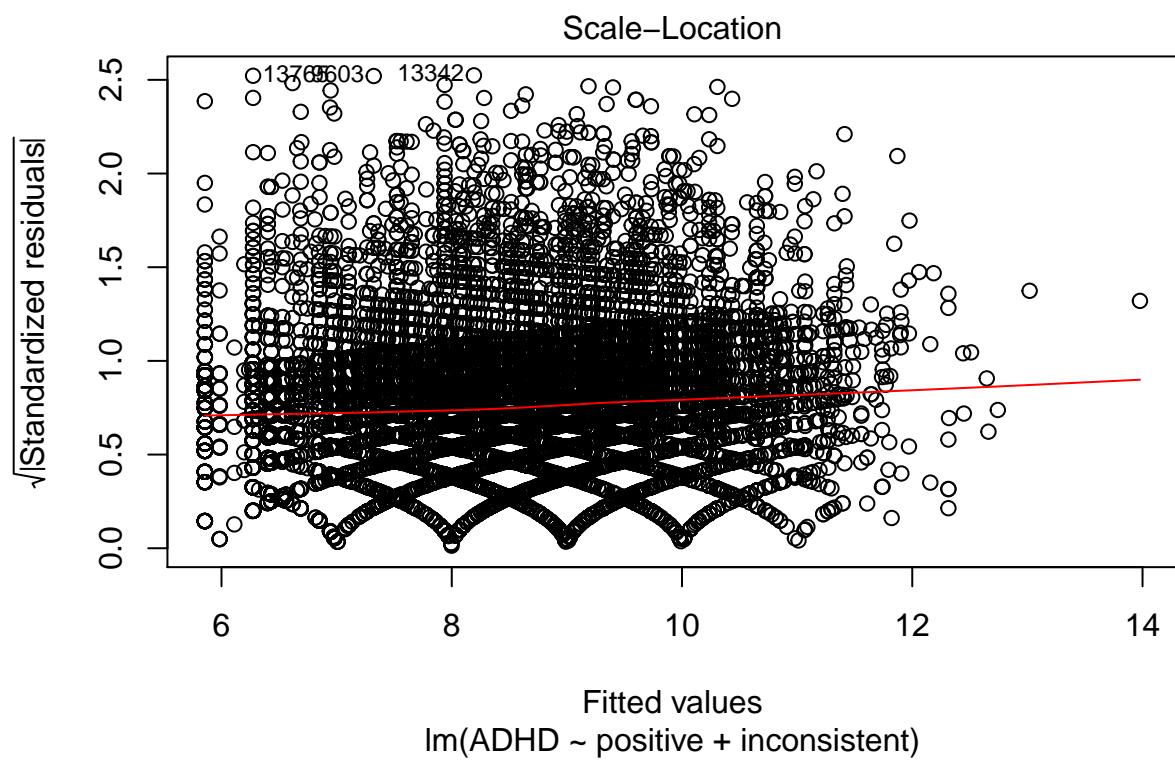
This is exciting! Parenting style at age 5 predicts ADHD symptoms at age 8!

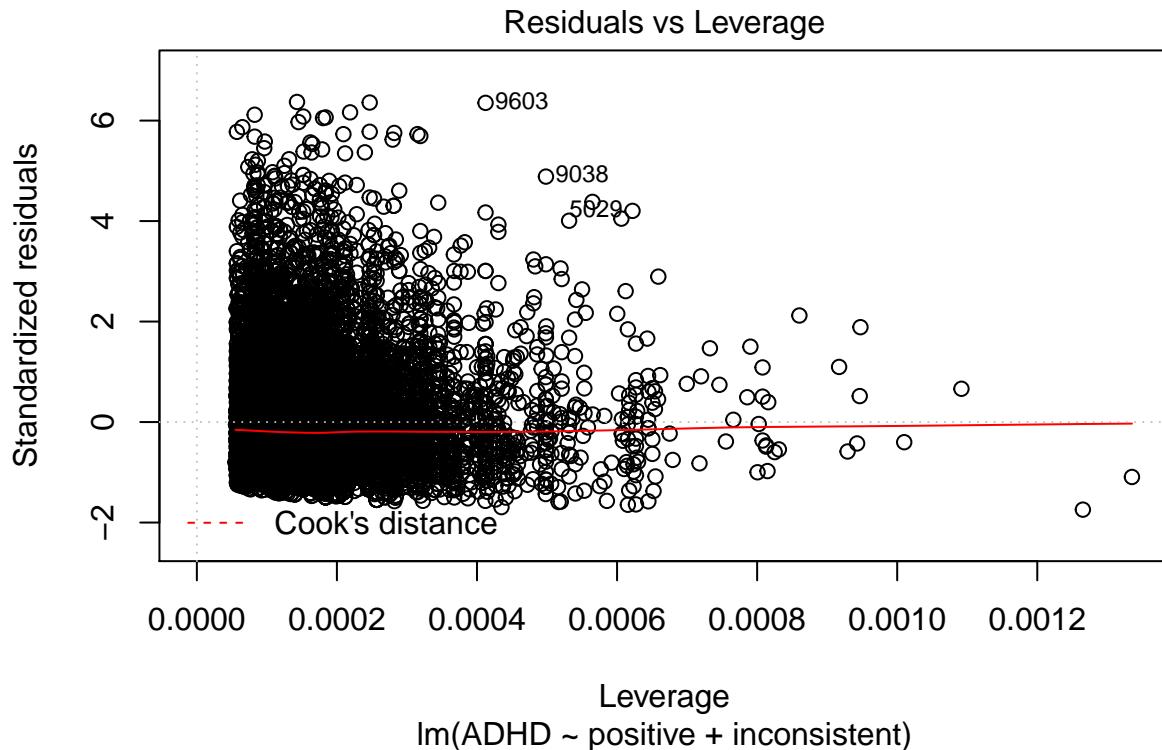
Before getting carried away, it makes sense to look at a few diagnostic plots:

```
par(ask=F)
plot(ADHD_model)
```









Here we can go through an explanation of the diagnostic plots.

As expected, this shows that a simple linear regression is probably not appropriate here.

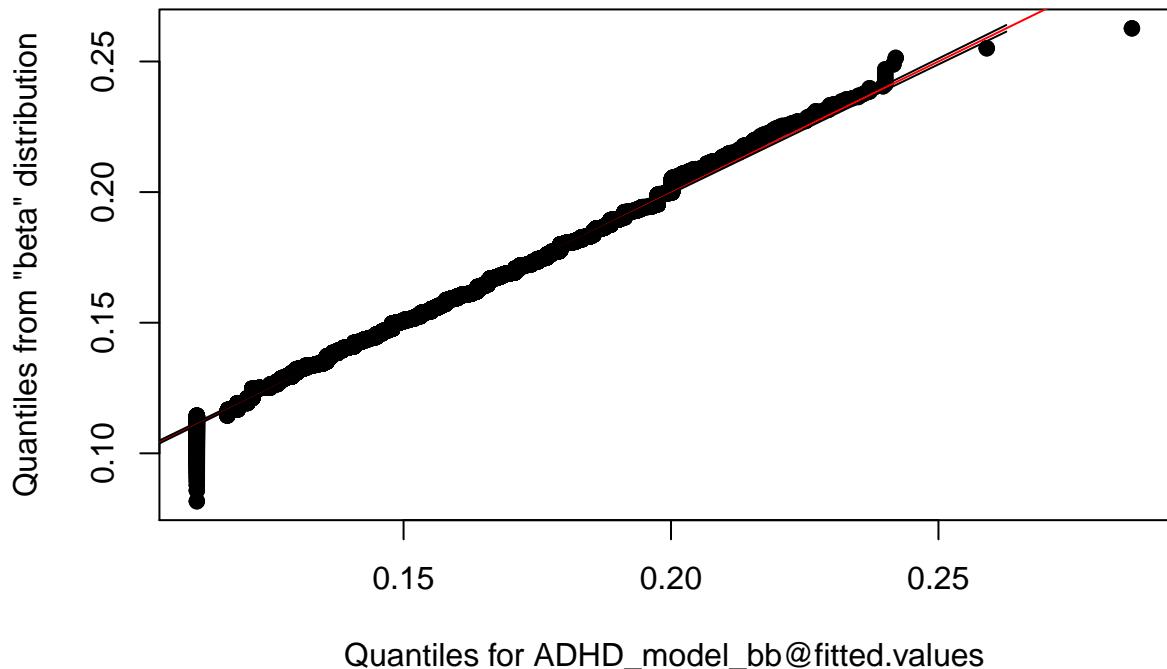
A beta-binomial regression

A beta-binomial regression would be more appropriate here, because the outcome variable can only be between 0 and an upper bound. There are several possibilities to do such regressions in R. On mature package that can do this is **VGAM**. Alternatives include **brms** and **rstanarm**.

```
library(VGAM)
library(qualityTools)
ADHD_model_bb = vglm(cbind(ADHD, 54-ADHD) ~ poly(positive,2) + poly(inconsistent,2),
                      betabinomial,
                      data = my_data,
                      trace = TRUE)

## VGLM    linear loop  1 :  loglikelihood = -56950.6521
## VGLM    linear loop  2 :  loglikelihood = -56942.2978
## VGLM    linear loop  3 :  loglikelihood = -56942.2948
## VGLM    linear loop  4 :  loglikelihood = -56942.2948
qqPlot(ADHD_model_bb@fitted.values, "beta", start = list(shape1 = 10, shape2 = 2))
```

Q–Q Plot for "beta" distribution



```
summary(ADHD_model_bb)
```

```
##
## Call:
## vglm(formula = cbind(ADHD, 54 - ADHD) ~ poly(positive, 2) + poly(inconsistent,
##      2), family = betabinomial, data = my_data, trace = TRUE)
##
##
## Pearson residuals:
##          Min      1Q  Median      3Q     Max
## logit(mu) -2.495 -0.6252 -0.0487  0.5370  6.682
## logit(rho) -0.812 -0.7132 -0.4786  0.1698 20.519
##
## Coefficients:
##                               Estimate Std. Error z value Pr(>|z|)
## (Intercept):1           -1.67518   0.00680 -246.351 < 2e-16 ***
## (Intercept):2           -2.25709   0.01311 -172.154 < 2e-16 ***
## poly(positive, 2)1      -8.07727   0.87687  -9.211 < 2e-16 ***
## poly(positive, 2)2      -3.96380   0.86140  -4.602 4.19e-06 ***
## poly(inconsistent, 2)1 19.16213   0.88902  21.554 < 2e-16 ***
## poly(inconsistent, 2)2  0.18766   0.84915   0.221    0.825
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Number of linear predictors:  2
##
```

```

## Names of linear predictors: logit(mu), logit(rho)
##
## Log-likelihood: -56942.29 on 36400 degrees of freedom
##
## Number of iterations: 4
##
## No Hauck-Donner effect found in any of the estimates

```

Even though the QQ plot shows that beta-binomial model is more appropriate, the tutorial continues with the result of the regression model. This is because some of the topics we are going to cover are easy to implement for a linear regression model, but take a bit more time for alternative models.

Processing regression results: Tables and plots

With the package `apaTables` we can format results tables from regressions as described in the APA Manual 6.

```

library(apaTables)
apa.reg.table(ADHD_model, filename = "Table2 APA.doc", table.number = 2)

##
##
## Table 2
##
## Regression results using ADHD as the criterion
##
##
## Predictor      b      b_95%_CI  beta    beta_95%_CI sr2 sr2_95%_CI
## (Intercept)  8.49**  [8.39, 8.59]
## positive    -0.40** [-0.50, -0.30] -0.06 [-0.07, -0.04] .00 [.00, .00]
## inconsistent 1.04**  [0.94, 1.15]  0.15  [0.13, 0.16] .02 [.02, .03]
##
## r            Fit
##
## -.09**
## .16**
##          R2 = .030**
##          95% CI [.02, .03]
##
##
## Note. A significant b-weight indicates the beta-weight and semi-partial correlation are also significant.
## b represents unstandardized regression weights. beta indicates the standardized regression weights.
## sr2 represents the semi-partial correlation squared. r represents the zero-order correlation.
## Square brackets are used to enclose the lower and upper limits of a confidence interval.
## * indicates p < .05. ** indicates p < .01.
##
```

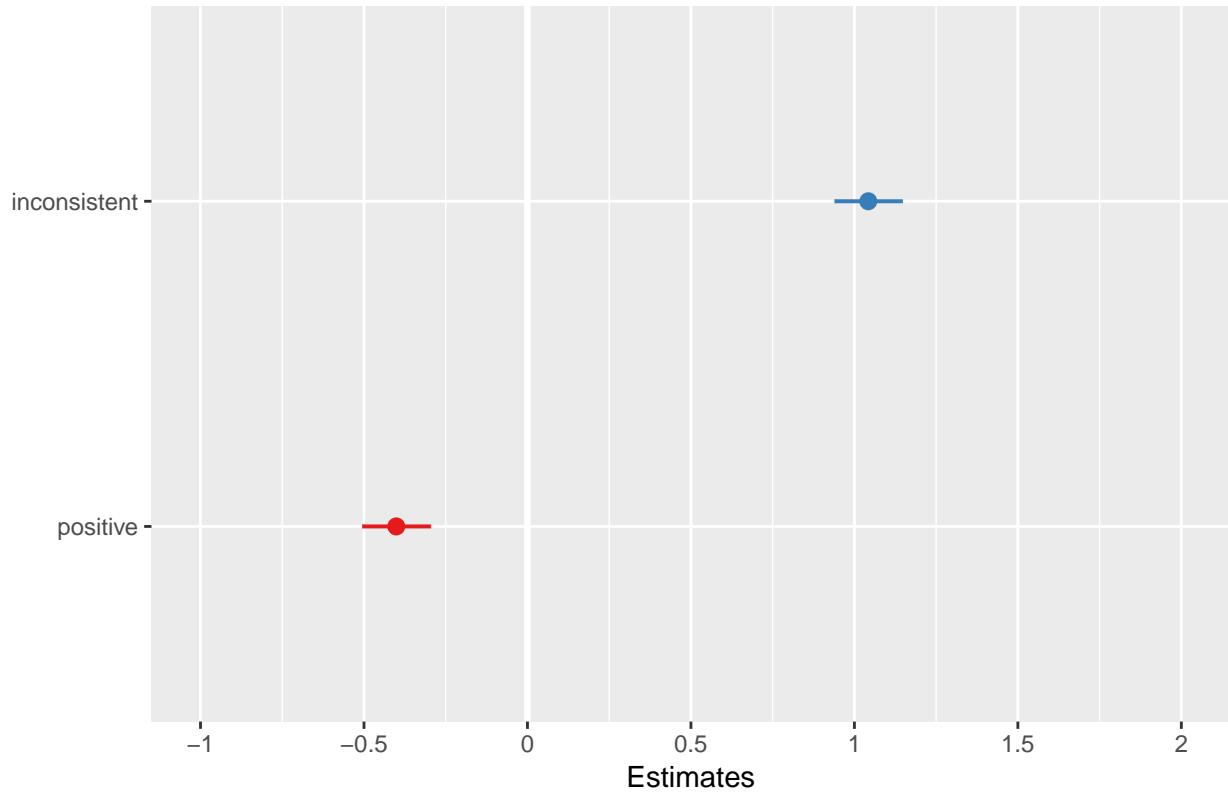
It is also relatively easy to plot model results. First, we simply plot the regression coefficients and confidence intervals.

```

library(sjPlot)
library(ggeffects)
library(ggplot2)
plot_model(ADHD_model)

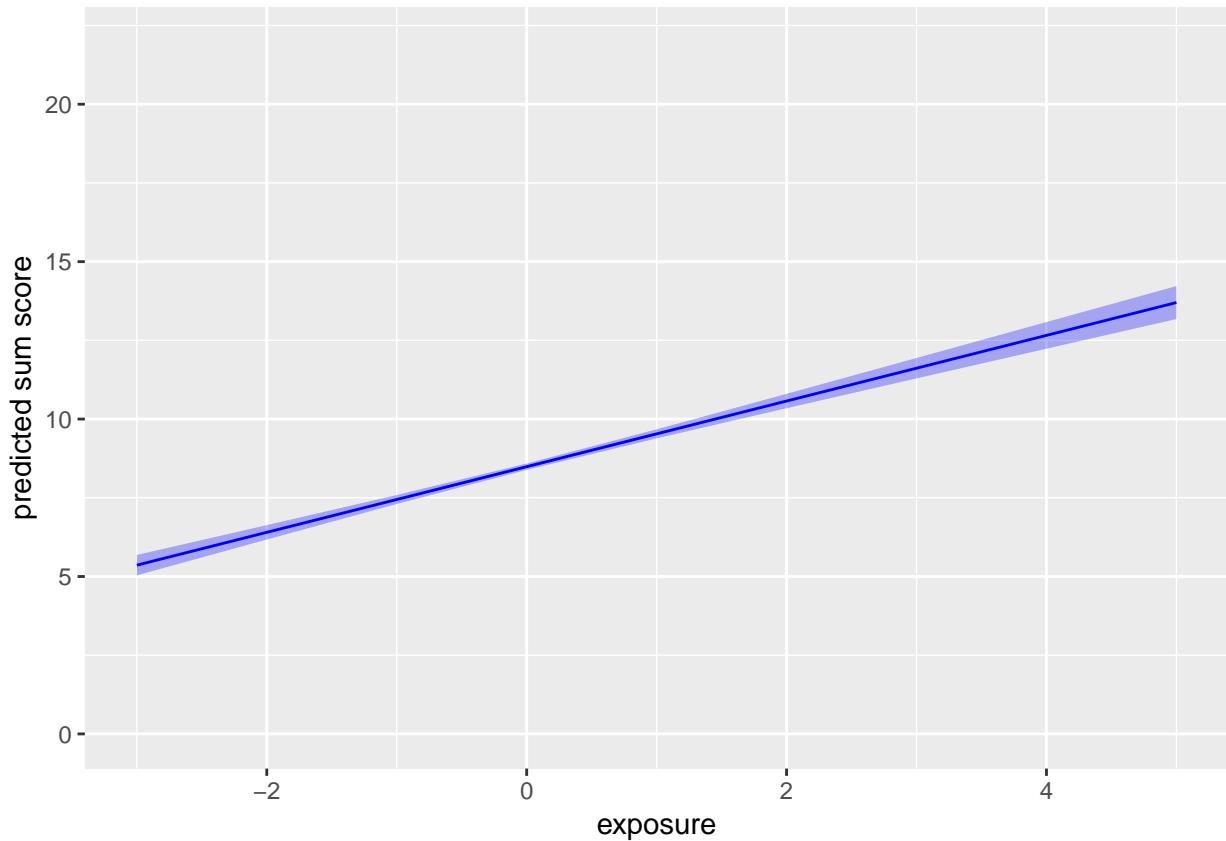
```

ADHD



```
effect_inconsistent = ggpredict(ADHD_model, terms = "inconsistent", typical = "mean")
effect_positive = ggpredict(ADHD_model, terms = "positive", typical = "mean")

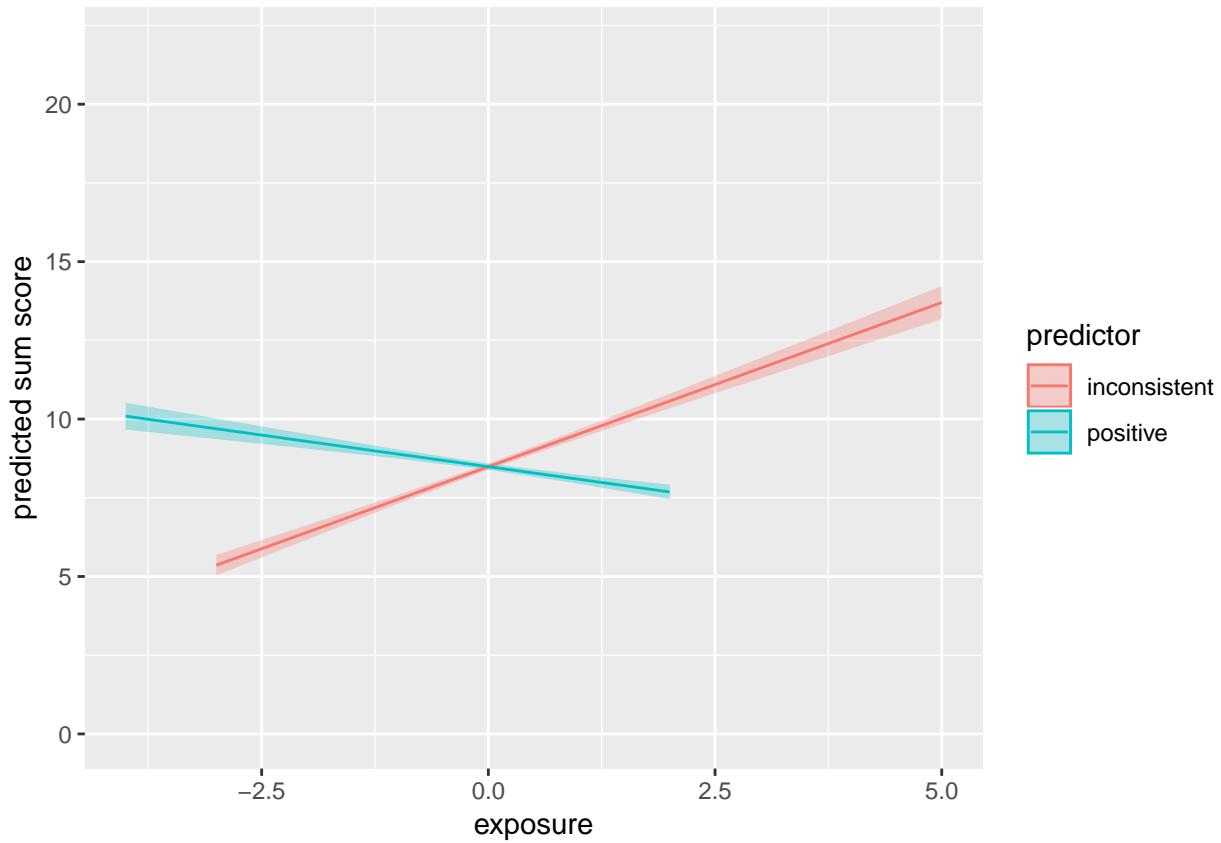
ggplot(effect_inconsistent, aes(x = x, y = predicted)) +
  geom_ribbon(aes(ymin=conf.low, ymax=conf.high, x=x, fill = "band"), alpha = 0.3, fill = "blue") +
  geom_line(col = "blue") +
  ylab("predicted sum score") +
  xlab("exposure") +
  coord_cartesian(ylim = c(0,quantile(my_data$ADHD,c(.95))))
```



We can also plot the effects of teh two predictors together:

```
my_effects = rbind(cbind(effect_inconsistent,predictor = "inconsistent"),
                    cbind(effect_positive,predictor = "positive"))

ggplot(my_effects, aes(x = x, y = predicted, col = predictor, group = predictor)) +
  geom_ribbon(aes(ymin=conf.low, ymax=conf.high, x=x, fill = predictor, col = NULL), alpha = 0.3) +
  geom_line() +
  ylab("predicted sum score") +
  xlab("exposure") +
  coord_cartesian(ylim = c(0,quantile(my_data$ADHD,c(.95))))
```



Adding and plotting non-linear effects

This shows a linear effect of parenting style. Is it possible that there are non-linear effects? We can easily test this by adding squared terms to the model.

```
ADHD_model2 = lm(ADHD ~ poly(positive, 2, raw = T) + poly(inconsistent, 2, raw = T), my_data)
summary(ADHD_model2)
```

```
##
## Call:
## lm(formula = ADHD ~ poly(positive, 2, raw = T) + poly(inconsistent,
##     2, raw = T), data = my_data)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -14.559  -4.692  -1.528   2.782  44.130
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)
## (Intercept)                 8.53420   0.08468 100.777 < 2e-16 ***
## poly(positive, 2, raw = T)1  -0.50734   0.06242 -8.127 4.67e-16 ***
## poly(positive, 2, raw = T)2  -0.17881   0.05824 -3.070 0.00214 **
## poly(inconsistent, 2, raw = T)1 1.03376   0.05254 19.677 < 2e-16 ***
## poly(inconsistent, 2, raw = T)2  0.13054   0.04055  3.219 0.00129 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## Residual standard error: 6.874 on 18198 degrees of freedom
## Multiple R-squared:  0.03056,    Adjusted R-squared:  0.03034
## F-statistic: 143.4 on 4 and 18198 DF,  p-value: < 2.2e-16
```

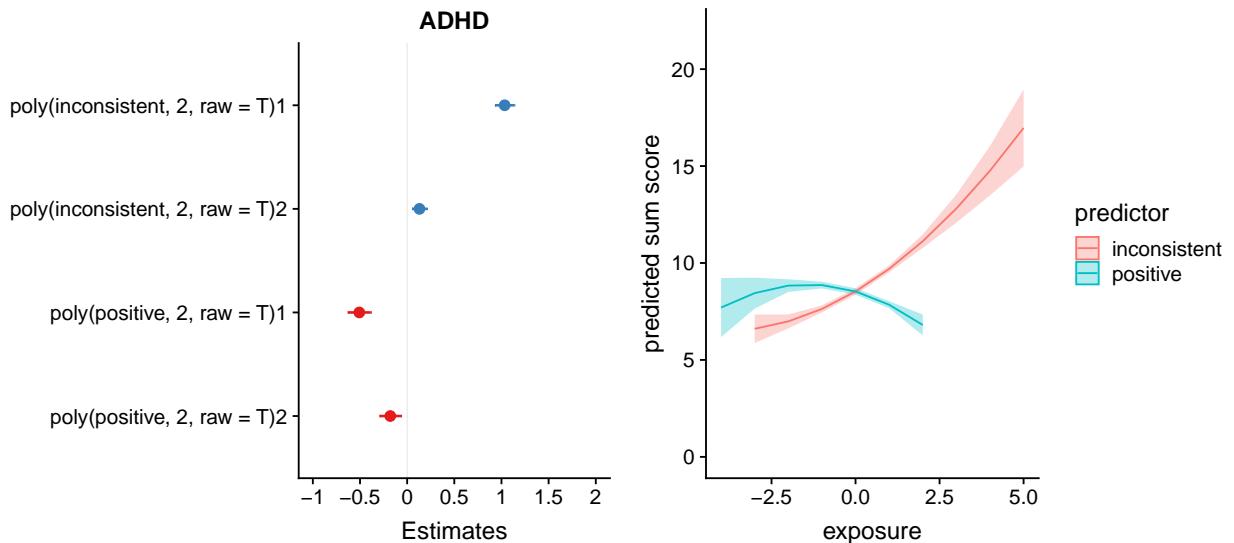
We'll use a new library, `cowplot` to put multiple ggplot figures together.

```
library(cowplot)
coef_plot = plot_model(ADHD_model2)
effect_inconsistent = ggpredict(ADHD_model2, terms = "inconsistent", typical = "mean")
effect_positive = ggpredict(ADHD_model2, terms = "positive", typical = "mean")

my_effects = rbind(cbind(effect_inconsistent,predictor = "inconsistent"),
                    cbind(effect_positive,predictor = "positive"))

effect_plot = ggplot(my_effects, aes(x = x, y = predicted, col = predictor, group = predictor)) +
  geom_ribbon(aes(ymin=conf.low, ymax=conf.high, x=x, fill = predictor, col = NULL), alpha = 0.3) +
  geom_line() +
  ylab("predicted sum score") +
  xlab("exposure") +
  coord_cartesian(ylim = c(0,quantile(my_data$ADHD,c(.95)))))

cowplot:::plot_grid(coef_plot,
                     effect_plot, nrow = 1)
```



Adding interaction effects

Finally, we can also add an interaction effect:

```
ADHD_model2 = lm(ADHD ~ poly(positive,2,raw = T) +
                  poly(inconsistent,2, raw = T) +
                  positive:inconsistent,
                  my_data)
summary(ADHD_model2)
```

```
##
```

```

## Call:
## lm(formula = ADHD ~ poly(positive, 2, raw = T) + poly(inconsistent,
##      2, raw = T) + positive:inconsistent, data = my_data)
##
## Residuals:
##      Min      1Q  Median      3Q      Max
## -14.558  -4.692  -1.528   2.782  44.130
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)
## (Intercept)                 8.5341008  0.0858618 99.393 < 2e-16 ***
## poly(positive, 2, raw = T)1  -0.5073202  0.0624897 -8.118 5.02e-16 ***
## poly(positive, 2, raw = T)2  -0.1787037  0.0603514 -2.961 0.00307 **
## poly(inconsistent, 2, raw = T)1 1.0337384  0.0526070 19.650 < 2e-16 ***
## poly(inconsistent, 2, raw = T)2  0.1306154  0.0422000  3.095 0.00197 **
## positive:inconsistent      0.0003715  0.0556015  0.007 0.99467
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.874 on 18197 degrees of freedom
## Multiple R-squared:  0.03056,   Adjusted R-squared:  0.03029
## F-statistic: 114.7 on 5 and 18197 DF,  p-value: < 2.2e-16

```

This does not seem to make a big difference.

This analysis suggests that mothers who report inconsistent parenting at child age 3 (they do not follow through with a threatened punishment) later report more ADHD problems with their children. There are several reasons that we should not immediately endorse a causal interpretation. This could in fact be reverse causation, if children with problems at age 8 also had problems at age 5, and it was these problems at age five that caused inconsistent parenting. Some would certainly argue that this is all genetic: the same genes that are responsible for mental health problems also cause inconsistent parenting (e.g. parents with mental health problems might be more prone to be inconsistent.). Genes are an example of confounders: common causes of exposure and outcomes. We do not have genetic information here, but we can easily think of other potential common causes like parental education and age, parity. In addition, we might want to adjust for other covariates like the child's gender.

Investigating potential confounders through plotting

The next lines of R code load these variables and adds them to our data.

```

data_directory = "data/"
Q1 = read_sav(paste0(data_directory,
                      "PDB439_Skjema1_v10.sav"))
Q1 = data.frame(Q1[,c("PREG_ID_439",
                      paste("AA",c(11,1123:1127,1300:1303),
                      sep = ""))])
names(Q1)[-1] = c("Q1_YEAR", "CivilStatus", "mEDUcomp", "mEDUcurr",
                  "fEDUcomp", "fEDUcurr", "n_people_19p_Q1",
                  "n_children_12_18_Q1", "n_children_6_11_Q1",
                  "n_children_0_5_Q1"))

Q1$CivilStatus = factor(Q1$CivilStatus,
                        levels = 1:6,
                        labels = c("married", "separated", "cohabitating",

```

```

        "widow", "single", "Other"))

EDUlabels = c("basic", "1-2 high school", "vocational high", "general high", "Bachelor", "Master" )

Q1$mEDU = Q1$mEDUcomp
idx = is.na(Q1$mEDUcomp) & !is.na(Q1$mEDUcurr) & Q1$mEDUcurr != 1
Q1$mEDU[idx] = Q1$mEDUcurr[idx] - 1
Q1$mEDU = ordered(Q1$mEDU, labels = EDUlabels)
Q1$fEDU = Q1$fEDUcomp
idx = is.na(Q1$fEDUcomp) & !is.na(Q1$fEDUcurr) & Q1$fEDUcurr != 1
Q1$fEDU[idx] = Q1$fEDUcurr[idx] - 1
Q1$fEDU = ordered(Q1$fEDU, labels = EDUlabels)

rm(EDUlabels, idx)

data_directory = "data/"
MFR = read_sav(paste0(data_directory,
                      "PDB439_MFR_520_v10.sav"))
MFR = data.frame(MFR[, c("PREG_ID_439", "BARN_NR", "KJONN", "MORS_ALDER",
                         "FAAR", "LEVENDEFODTE_5", "PARITET_5", "FLERFODSEL",
                         "VEKT", "SVLEN_DG", "APGAR1", "APGAR5", "FMND")])
names(MFR)[-c(1, 2)] = c("gender", "mAge", "birth_year",
                         "life_births", "parity", "multiple_births",
                         "birthweight", "pregnancy_duration",
                         "APGAR1", "APGAR5", "birthmonth")
MFR$gender = factor(MFR$gender, levels = 1:2, labels = c("boy", "girl"))
MFR[is.na(MFR$life_births), life_births := parity]
MFR[, birthmonth := as.numeric(birthmonth)]

my_data = merge(my_data,
                 MFR,
                 by = c("PREG_ID_439", "BARN_NR"),
                 all.x = T,
                 all.y = F)
my_data = merge(my_data,
                 Q1,
                 by = c("PREG_ID_439"),
                 all.x = T,
                 all.y = F)

```

Now we are ready to look at a few tables and plots. First, we use the `group_by` function of the `dplyr` package to show the mean exposures grouped by education.

```

library(dplyr)
options(digits = 2)
my_stats = my_data %>%
  group_by(mEDU) %>%
  summarise(m = mean(inconsistent),
            sd = sd(inconsistent),
            N = n())
my_stats

```

```
## # A tibble: 7 x 4
```

```

##   mEDU           m      sd      N
##   <ord>         <dbl> <dbl> <int>
## 1 basic          0.209  1.08   158
## 2 1-2 high school 0.0309  1.05   405
## 3 vocational high 0.00317 1.04   1512
## 4 general high   -0.0471 1.03   1974
## 5 Bachelor        -0.0245 0.980  7979
## 6 Master          0.0383  0.992  5724
## 7 <NA>            0.0187  1.05   451

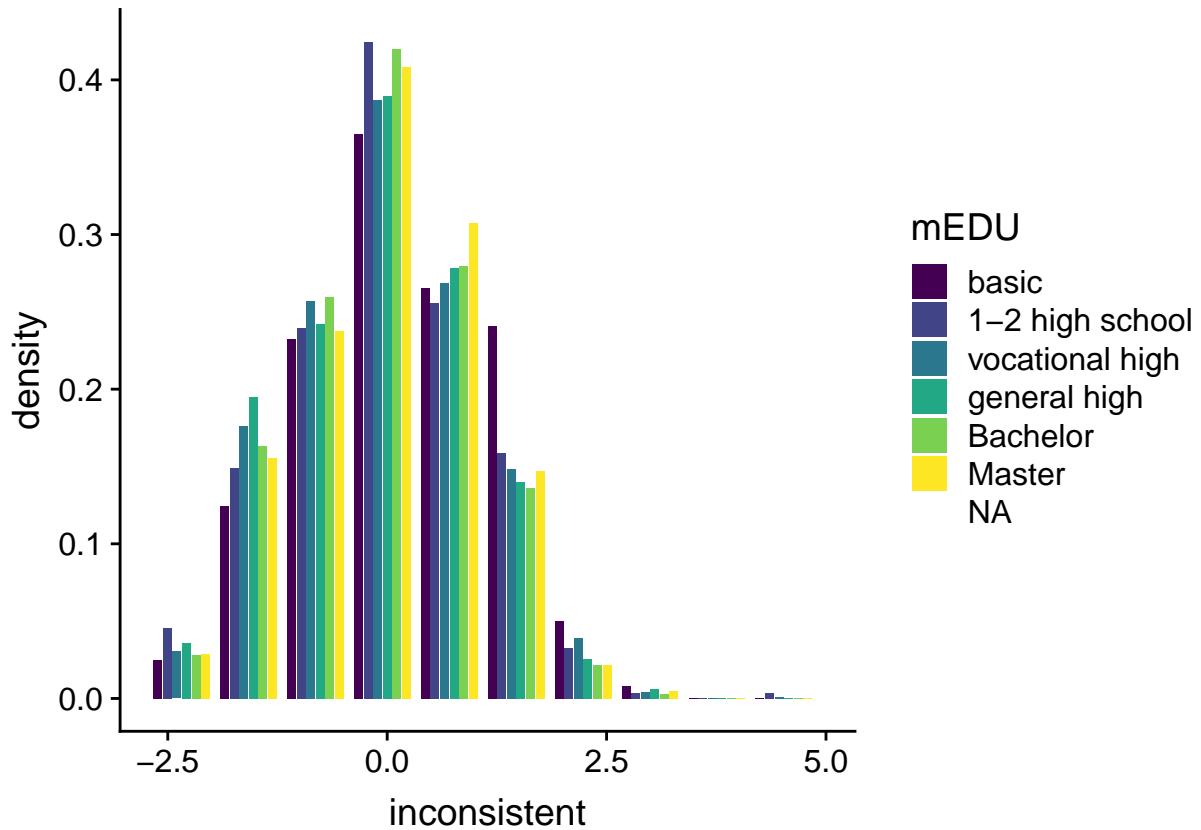
```

This does not look very clear. We can also plot the raw data using the nice `ggplot` package. Lets start with a histogram.

```

library(ggplot2)
ggplot(my_data, aes(x = inconsistent, fill = mEDU)) +
  geom_histogram(aes(y = ..density..), position="dodge2", bins = 10)

```



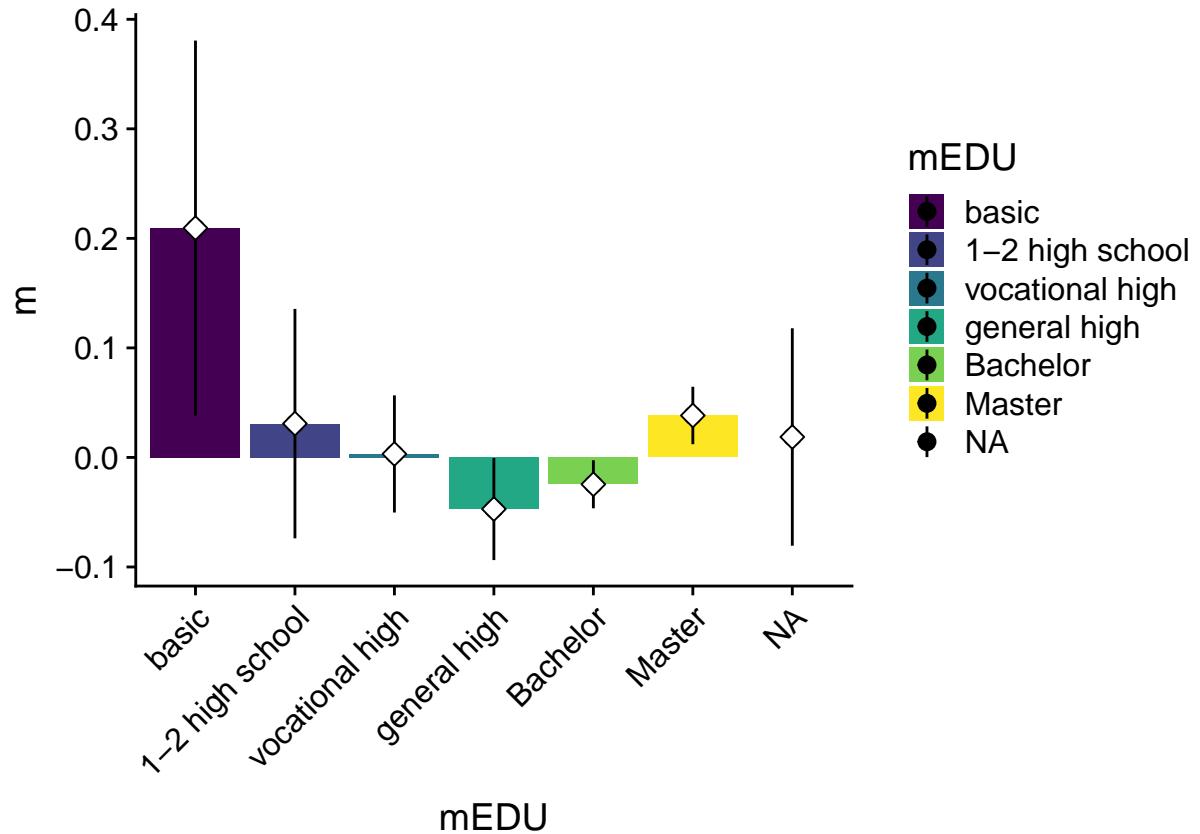
We can simply plot the `my_stats` table we made above, after adding confidence intervals:

```

my_stats$CIlower = my_stats$m - my_stats$sd/sqrt(my_stats$N)*2
my_stats$CIupper = my_stats$m + my_stats$sd/sqrt(my_stats$N)*2

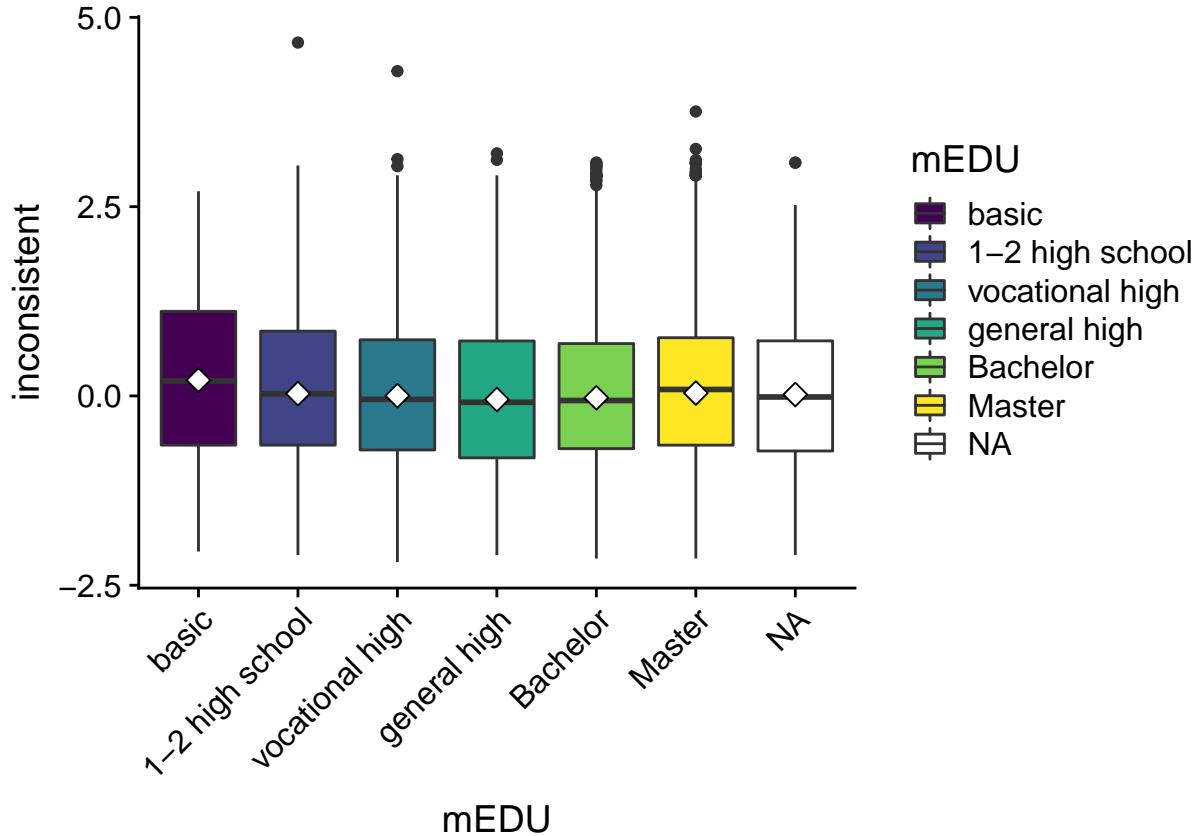
ggplot(my_stats, aes(x = mEDU, y = m, fill = mEDU)) + geom_bar(stat = "identity") +
  geom_pointrange(aes(x = mEDU, ymin = CIlower, ymax = CIupper)) +
  stat_summary(fun.y="median", geom="point", shape=23, size=3, fill="white") +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))

```



Now it looks as if the differences between groups are substantial. However, remember that the exposure variables has an sd of 1, so the differences are maybe not so important. To get a better view of the variations within and between groups, we can plot boxplots:

```
ggplot(my_data, aes(x = mEDU, y = inconsistent, fill = mEDU)) +
  geom_boxplot() +
  stat_summary(fun.y="mean", geom="point", shape=23, size=3, fill="white") +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



This shows that even though there are reliable differences between groups, the variation within groups is much larger than the variation between groups.

Now let's look at the association with all potential measured confounders (this looks like a lot of code, but it is really only repetition of the same code!):

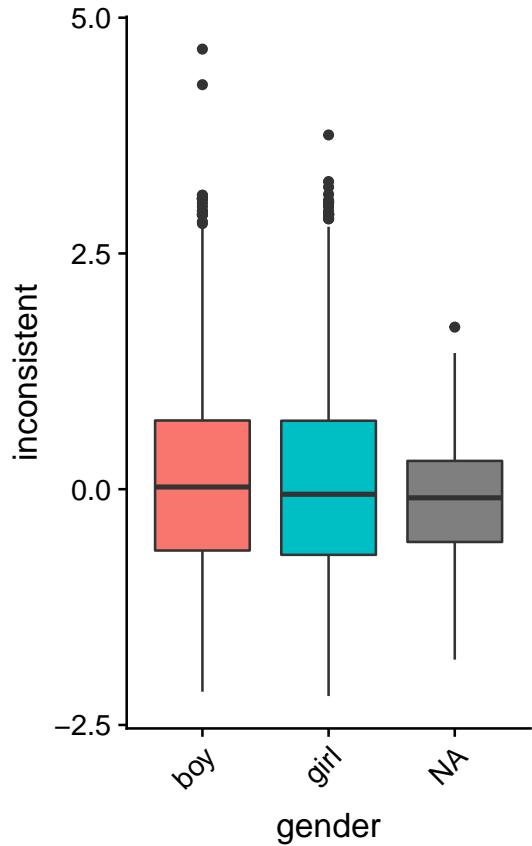
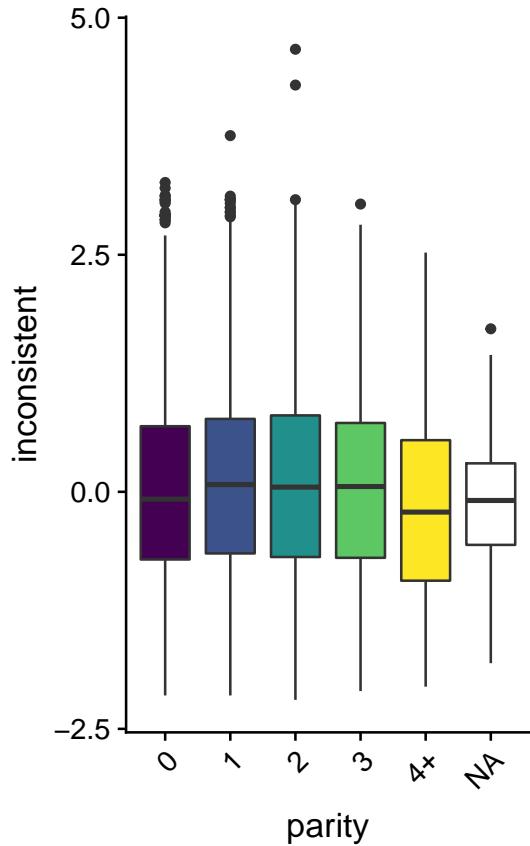
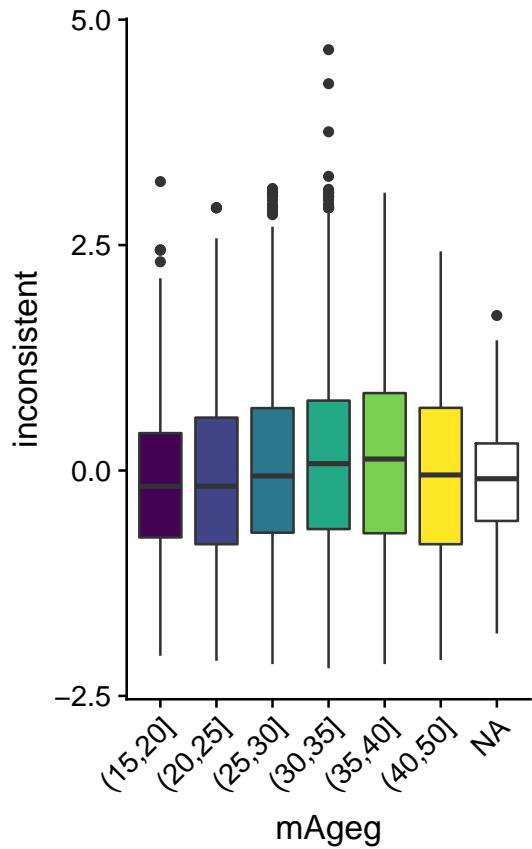
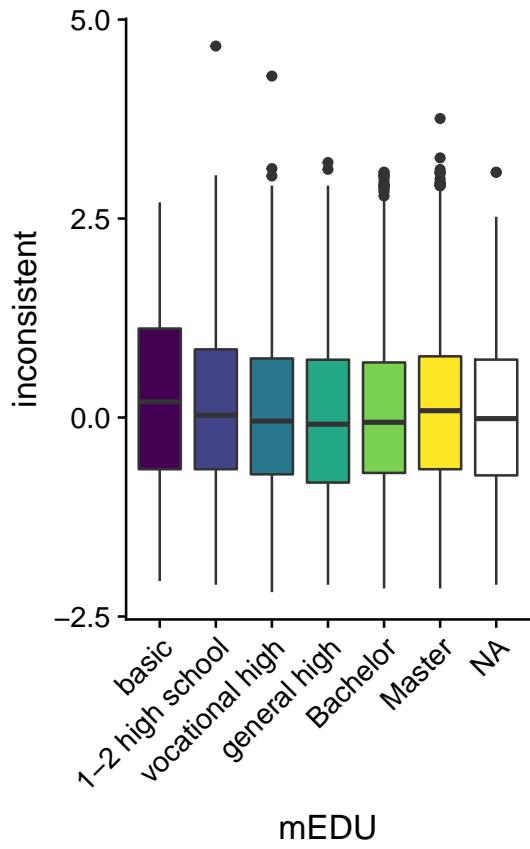
```

breaks = c(seq(15,40, by = 5),50)
my_data$mAgeg = cut(my_data$mAge, breaks = breaks, ordered_result = T)
my_data$parity = as_factor(my_data$parity, ordered = T)
my_data$birthmonth = ordered(my_data$birthmonth)
my_data$birth_year = ordered(my_data$birth_year)

by_edu = ggplot(my_data, aes(x = mEDU, y = inconsistent, fill = mEDU)) +
  geom_boxplot(show.legend = F) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
by_age = ggplot(my_data, aes(x = mAgeg, y = inconsistent, fill = mAgeg)) +
  geom_boxplot(show.legend = F) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
by_parity = ggplot(my_data, aes(x = parity, y = inconsistent, fill = parity)) +
  geom_boxplot(show.legend = F) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
by_gender = ggplot(my_data, aes(x = gender, y = inconsistent, fill = gender)) +
  geom_boxplot(show.legend = F) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
cowplot::plot_grid(by_edu,
                   by_age,
                   by_parity,

```

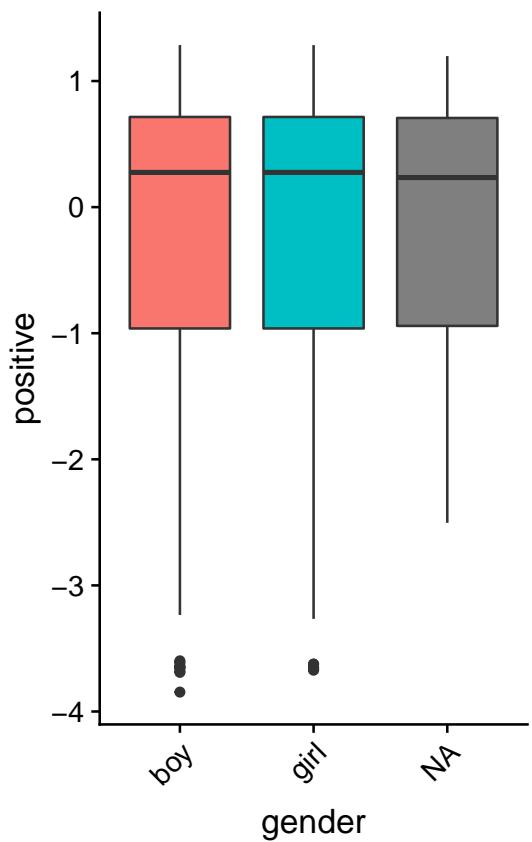
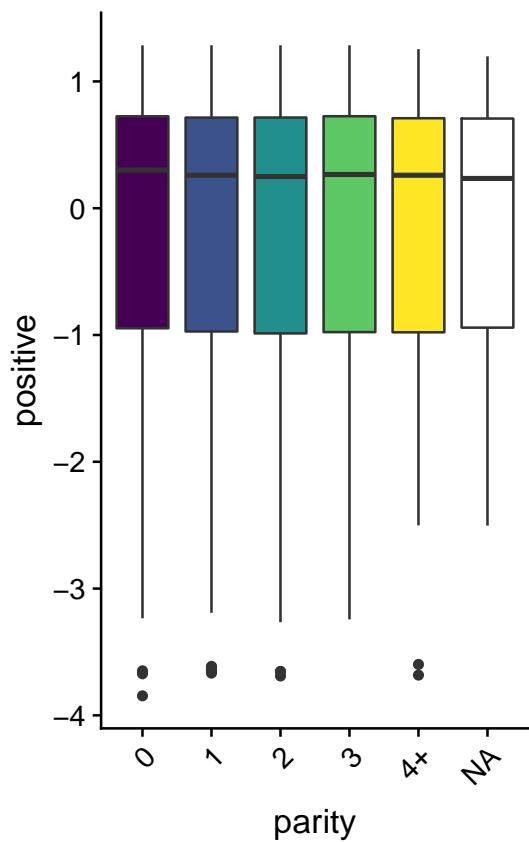
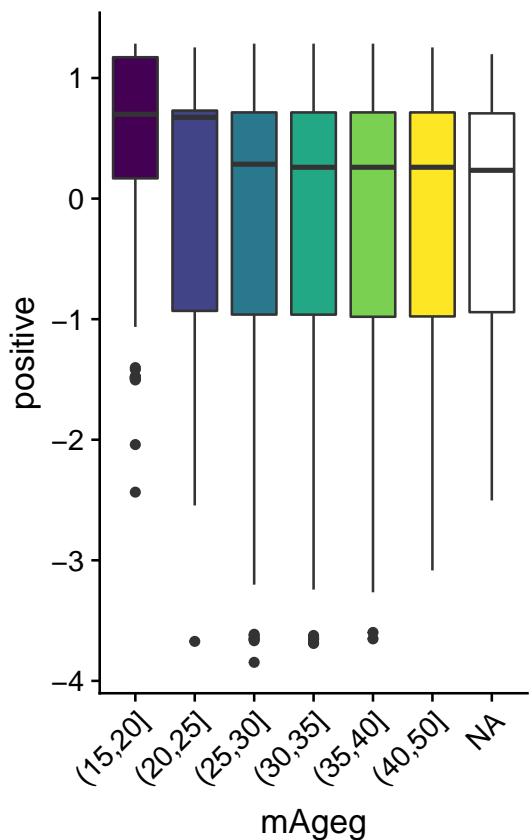
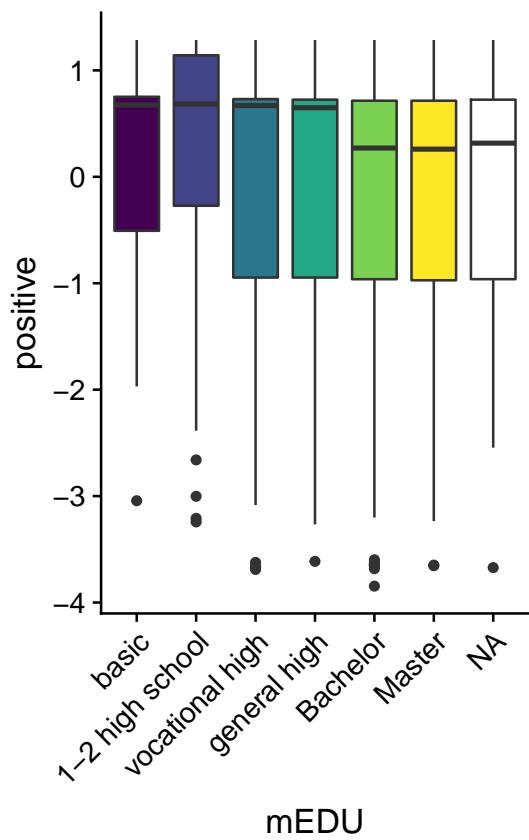
by_gender)



Try to do the same plot but stratify by birthmonth.

Here is the same plot for positive parenting:

```
by_edu = ggplot(my_data, aes(x = mEDU, y = positive, fill = mEDU)) +
  geom_boxplot(show.legend = F) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
by_age = ggplot(my_data, aes(x = mAgeg, y = positive, fill = mAgeg)) +
  geom_boxplot(show.legend = F) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
by_parity = ggplot(my_data, aes(x = parity, y = positive, fill = parity)) +
  geom_boxplot(show.legend = F) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
by_gender = ggplot(my_data, aes(x = gender, y = positive, fill = gender)) +
  geom_boxplot(show.legend = F) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
cowplot::plot_grid(by_edu,
                    by_age,
                    by_parity,
                    by_gender)
```



These associations also do not seem particularly strong.

We conclude this plotting tour with a look at the association between confounders and ADHD symptoms.

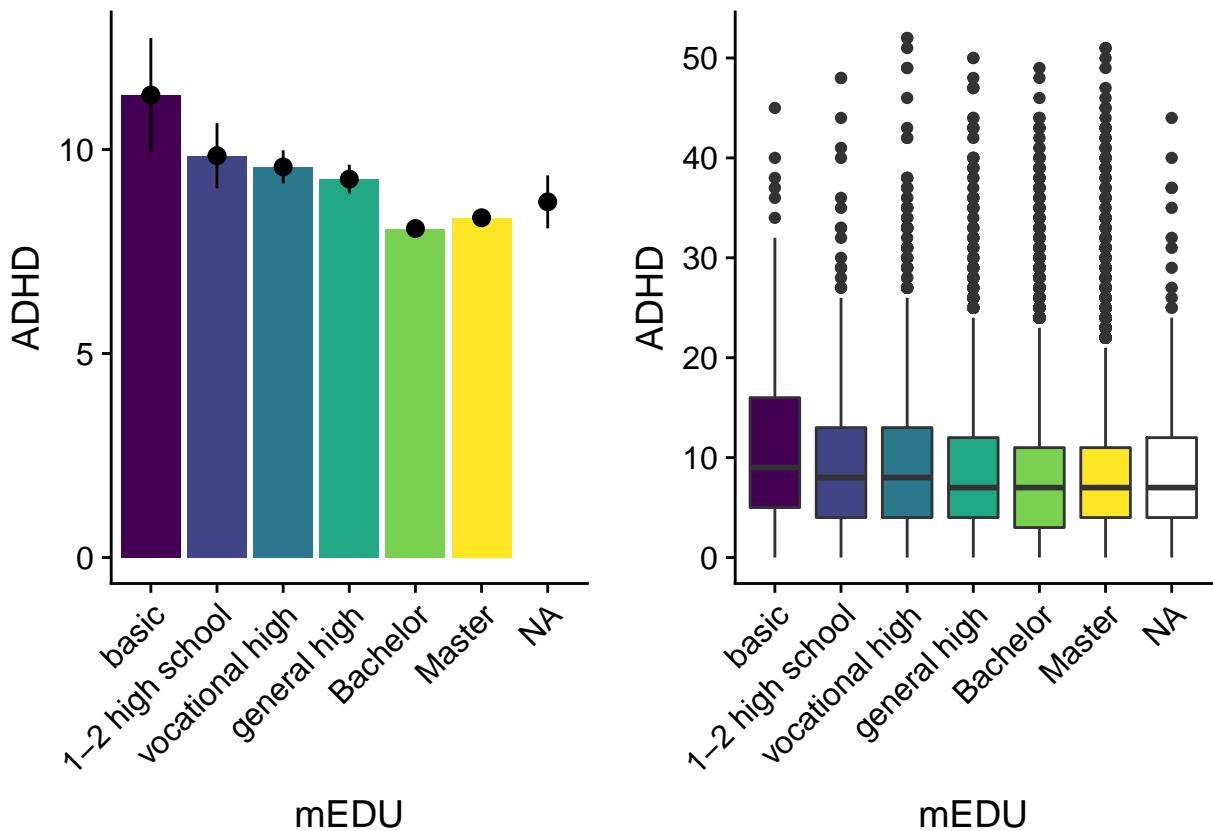
```
my_stats = my_data %>%
  group_by(mEDU) %>%
  summarise(m = mean(ADHD),
            sd = sd(ADHD),
            N = n()
  )

my_stats$CIlower = my_stats$m - my_stats$sd/sqrt(my_stats$N)*2
my_stats$CIupper = my_stats$m + my_stats$sd/sqrt(my_stats$N)*2

mean_ci_plot = ggplot(my_stats,aes(x = mEDU, y = m, fill = mEDU)) +
  geom_bar(stat = "identity", show.legend = F) +
  geom_pointrange(aes(x = mEDU, ymin = CIlower, ymax = CIupper), show.legend = F) +
  ylab("ADHD") +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))

box_plot = ggplot(my_data, aes(x = mEDU, y = ADHD, fill = mEDU)) +
  geom_boxplot(show.legend = F) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))

cowplot:::plot_grid(mean_ci_plot,
                     box_plot)
```



OK, now that we have seen that exposure and outcome are (weakly) associated with a likely confounder, it makes sense to include them into the model.

A final regression model: Ordinal predictors

```

ADHD_model4 = lm(ADHD ~
  poly(positive, 2, raw = T) +
  poly(inconsistent, 2, raw = T) +
  mEDU +
  poly(mAge, 2, raw = T) +
  parity +
  gender,
  data = my_data)
summary(ADHD_model4)

##
## Call:
## lm(formula = ADHD ~ poly(positive, 2, raw = T) + poly(inconsistent,
##      2, raw = T) + mEDU + poly(mAge, 2, raw = T) + parity + gender,
##      data = my_data)
##
## Residuals:
##    Min      1Q  Median      3Q     Max
## -15.73  -4.43  -1.46   2.71  44.78
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)    
## (Intercept)                11.73203   1.94379   6.04  1.6e-09 ***
## poly(positive, 2, raw = T)1 -0.60990   0.06207  -9.83 < 2e-16 ***
## poly(positive, 2, raw = T)2 -0.18476   0.05780  -3.20  0.0014 **  
## poly(inconsistent, 2, raw = T)1  1.01558   0.05207 19.51 < 2e-16 ***
## poly(inconsistent, 2, raw = T)2  0.10674   0.04022  2.65  0.0080 **  
## mEDU.L                   -2.60559   0.35622  -7.31  2.7e-13 ***
## mEDU.Q                   0.33185   0.31432   1.06  0.2911    
## mEDU.C                   0.08137   0.27785   0.29  0.7696    
## mEDU^4                   0.45981   0.23616   1.95  0.0515 .    
## mEDU^5                   0.27128   0.18336   1.48  0.1390    
## poly(mAge, 2, raw = T)1  -0.12652   0.12565  -1.01  0.3140    
## poly(mAge, 2, raw = T)2  0.00218   0.00199   1.09  0.2743    
## parity.L                 -1.61496   0.37906  -4.26  2.1e-05 ***
## parity.Q                 0.97659   0.32255   3.03  0.0025 **  
## parity.C                 0.08687   0.27508   0.32  0.7521    
## parity^4                 -0.04882   0.19909  -0.25  0.8063    
## gendergirl                -2.40471   0.10073 -23.87 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.7 on 17700 degrees of freedom
##   (486 observations deleted due to missingness)
## Multiple R-squared:  0.0807, Adjusted R-squared:  0.0799
## F-statistic: 97.2 on 16 and 17700 DF,  p-value: <2e-16

```

This looks a bit wild. The reason is that if we include ordinal variables like education and parity into the

model, R will by default implement polynomials to a degree equal to `number_levels - 1`. However, we can correct this by setting the contrast to for example a polynomial of order 2.

```

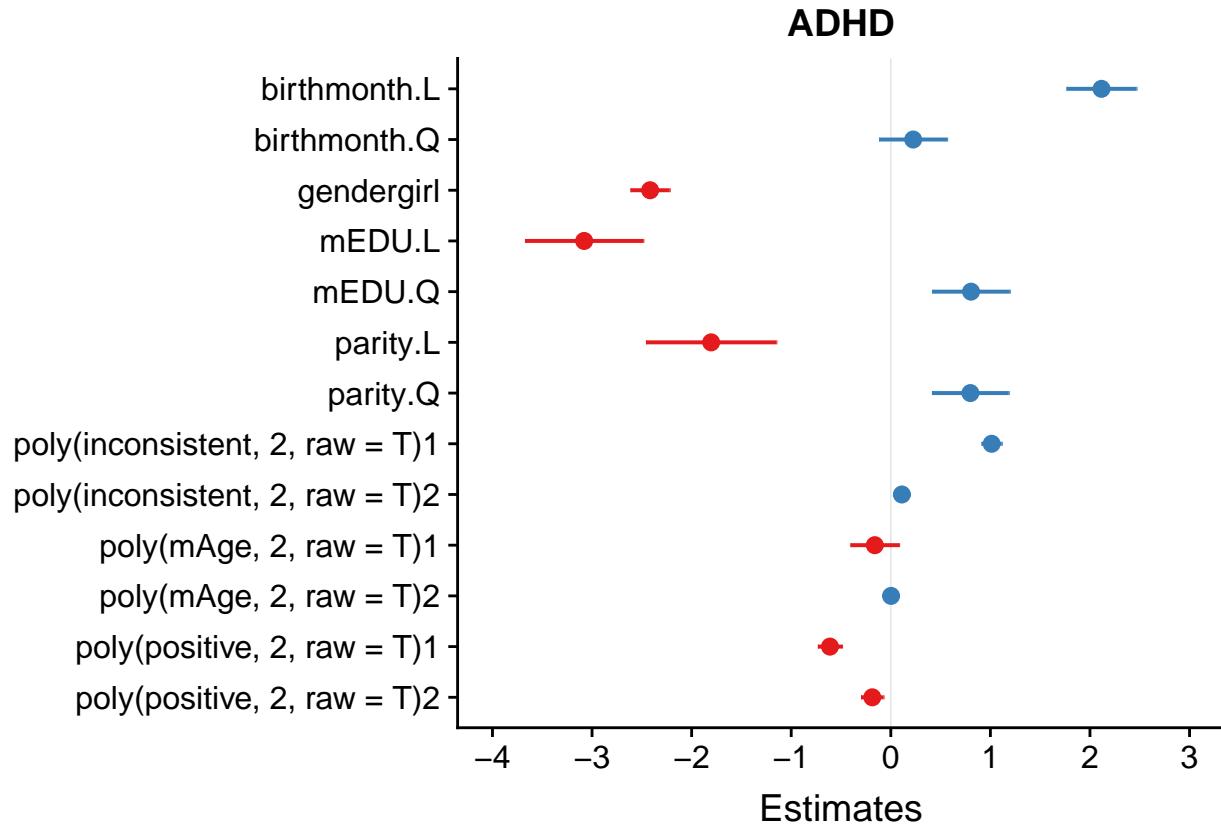
n_levels_Edu = length(levels(my_data$mEDU))
n_levels_parity = length(levels(my_data$parity))
n_levels_bm = length(levels(my_data$birthmonth))
contrasts(my_data$mEDU, 2) <- contr.poly(n_levels_Edu)
contrasts(my_data$parity, 2) <- contr.poly(n_levels_parity)
contrasts(my_data$birthmonth, 2) <- contr.poly(n_levels_bm)

ADHD_model4 = lm(ADHD ~
                  poly(positive, 2, raw = T) +
                  poly(inconsistent, 2, raw = T) +
                  mEDU +
                  poly(mAge, 2, raw = T) +
                  parity +
                  gender +
                  birthmonth,
                  data = my_data)
summary(ADHD_model4)

##
## Call:
## lm(formula = ADHD ~ poly(positive, 2, raw = T) + poly(inconsistent,
##   2, raw = T) + mEDU + poly(mAge, 2, raw = T) + parity + gender +
##   birthmonth, data = my_data)
##
## Residuals:
##    Min      1Q Median      3Q     Max
## -15.52  -4.42  -1.46   2.67  44.44
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)    
## (Intercept)                12.40140  1.92513   6.44  1.2e-10 ***
## poly(positive, 2, raw = T)1  -0.61092  0.06185  -9.88  < 2e-16 ***
## poly(positive, 2, raw = T)2  -0.18578  0.05760  -3.23  0.0013 **  
## poly(inconsistent, 2, raw = T)1  1.01291  0.05189 19.52  < 2e-16 ***
## poly(inconsistent, 2, raw = T)2  0.11124  0.04008  2.78  0.0055 **  
## mEDU.L                   -3.07929  0.30323 -10.16  < 2e-16 ***
## mEDU.Q                   0.80574  0.19853   4.06  5.0e-05 ***
## poly(mAge, 2, raw = T)1   -0.16066  0.12453  -1.29  0.1970  
## poly(mAge, 2, raw = T)2   0.00267  0.00198   1.35  0.1767  
## parity.L                  -1.80388  0.33399  -5.40  6.7e-08 ***
## parity.Q                  0.79911  0.19683   4.06  4.9e-05 *** 
## gendergirl                -2.41677  0.10040 -24.07  < 2e-16 ***
## birthmonth.L                2.11601  0.17936  11.80  < 2e-16 *** 
## birthmonth.Q                0.22470  0.17422   1.29  0.1971  
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.7 on 17703 degrees of freedom
##   (486 observations deleted due to missingness)
## Multiple R-squared:  0.0866, Adjusted R-squared:  0.0859
## F-statistic: 129 on 13 and 17703 DF,  p-value: <2e-16

```

```
plot_model(ADHD_model4)
```



This table suggests that parenting style still has an effect, after we controlled for a number of potential confounders. It is, however, difficult to understand the importance (effect sizes) for the different predictors, because they are measured on different scales (have different means and standard deviations). To get a better idea of effect sizes we can again plot marginal effects.

We use a small loop, so that we don't have to write the same lines of code repeatedly for the different predictors.

```
ylim = c(0,quantile(my_data$ADHD,c(.95)))

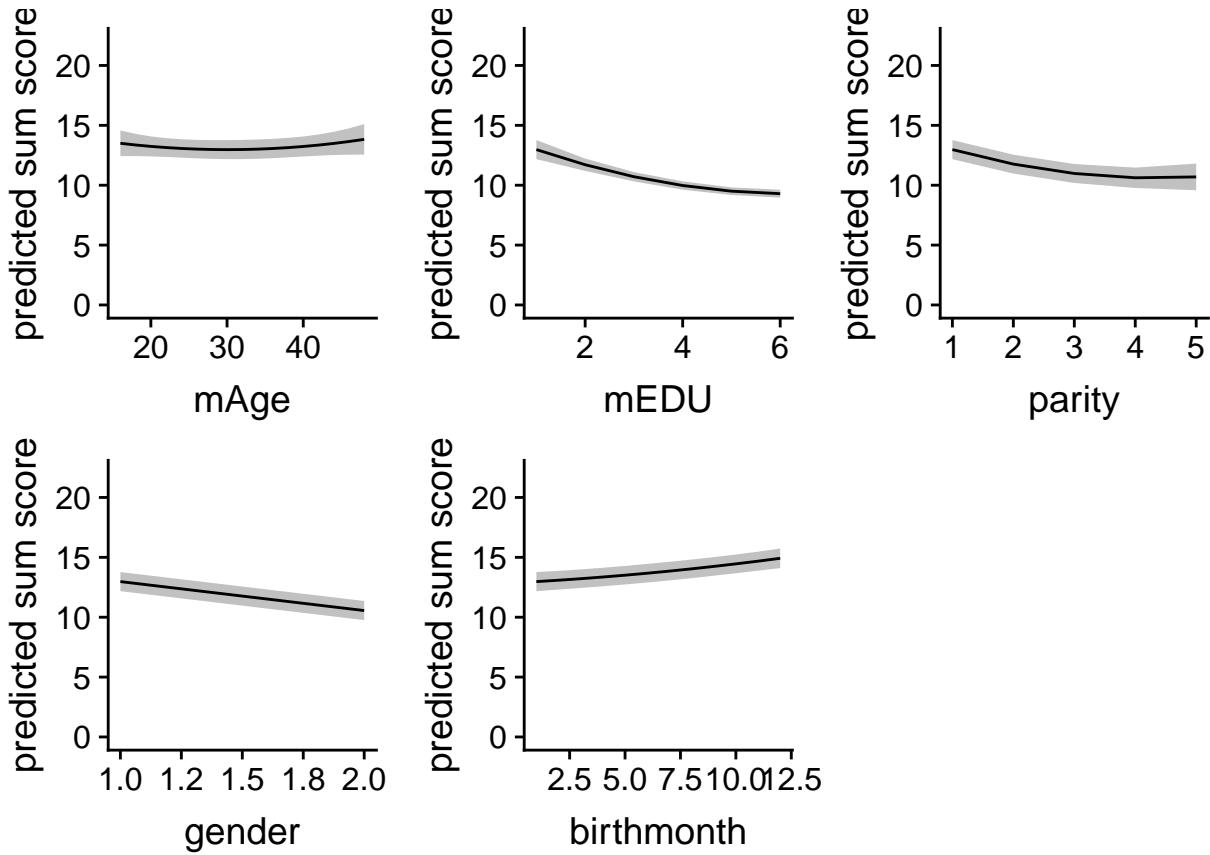
predictors = c("positive","inconsistent","mAge","mEDU","parity","gender", "birthmonth")
effect_plots = vector(mode = "list", length = length(predictors))
names(effect_plots) = predictors

for ( p in 1:length(predictors)) {
  efct = ggpredict(ADHD_model4, terms = predictors[p])
  effect_plots[[p]] = ggplot(efct, aes(x = x, y = predicted)) +
    geom_ribbon(aes(ymin=conf.low,
                     ymax=conf.high,
                     x=x,
                     col = NULL),
                alpha = 0.3) +
    geom_line() +
    ylab("predicted sum score") +
    xlab(predictors[p]) +
```

```
        coord_cartesian(ylim = ylim)
    }
```

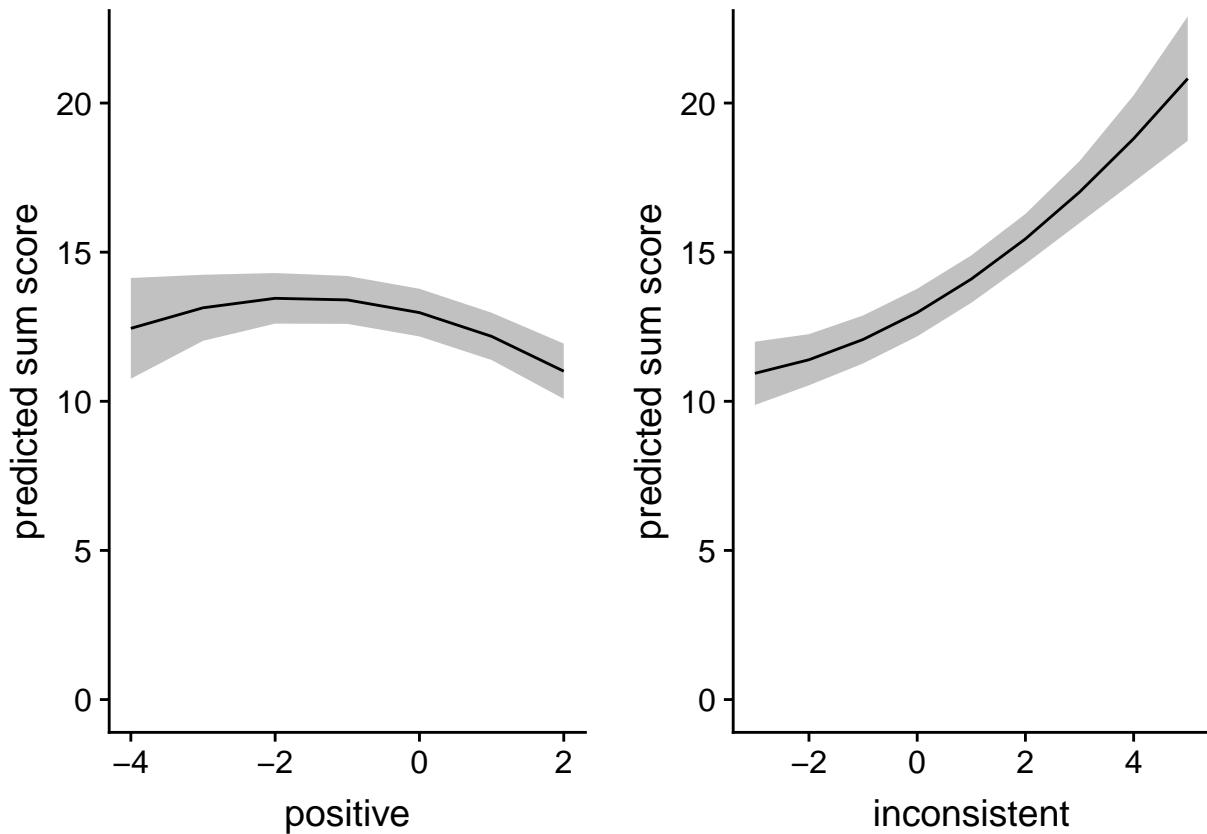
First we look at the effect of confounders:

```
cowplot::plot_grid(plotlist = effect_plots[3:length(predictors)])
```



And finally on the effect of our exposures:

```
cowplot::plot_grid(plotlist = effect_plots[1:2])
```



Interestingly, this suggests that the association between inconsistent parenting style at child age 5 and ADHD sum score at child age 8 increases after we control for some potential confounders.