# Solution for the assignment of the ninth class

Kovacs Marton

9/29/2021

## Importing data

```
processed <- read_tsv("data/boldog_processed.tsv")
```

```
##
## -- Column specification ---------------------------------------------------
## cols(
##   .default = col_double(),
##   neme = col_character(),
##   isk = col_character()
## )
## i Use `spec()` for the full column specifications.
```

## Data exploration

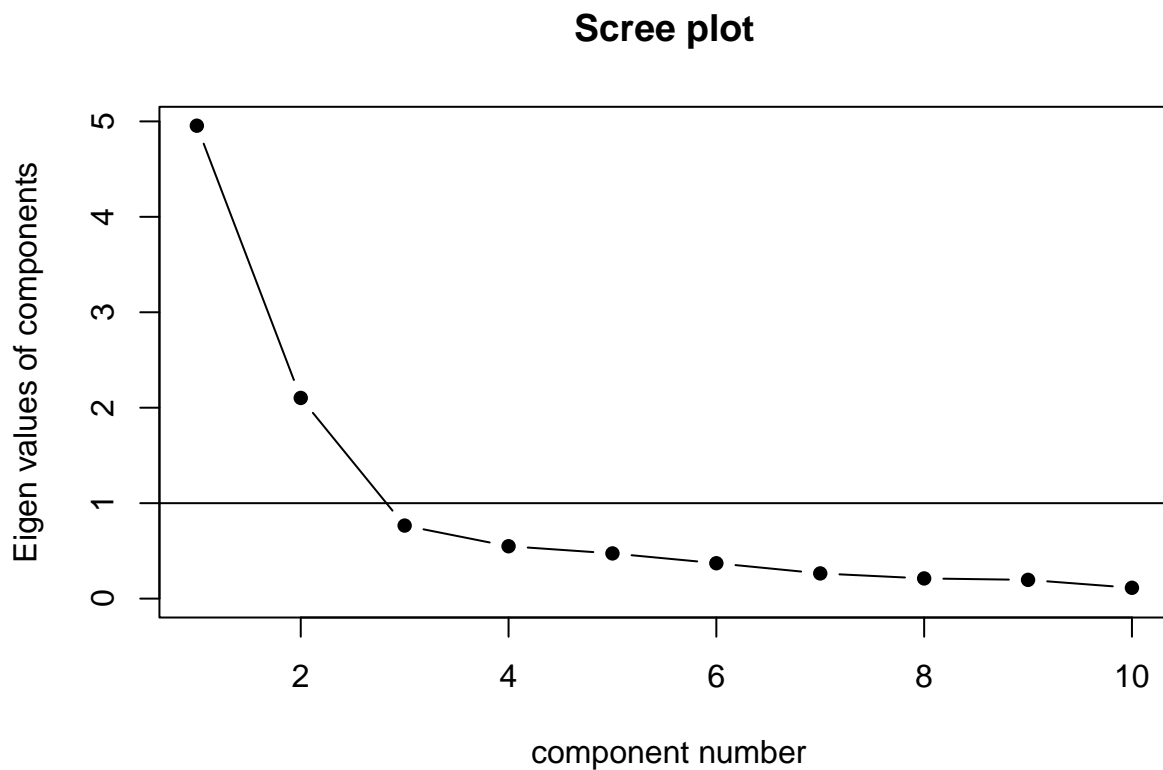```
skimr::skim(processed) %>%
  kable()
```

| skim_type | variable | missing | complete | character.min | character.max | character.empty | character.n_unique | character.whitespace | numeric.mean | numeric.sd | numeric.p0 | numeric.p25 | numeric.p50 | numeric.p75 | numeric.p100 | numeric.hist |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| character | neme | 0 | 1.000 | 2 | 5 | 0 | 2 | 0 | NA | NA | NA | NA | NA | NA | NA | NA |
| character | isk | 0 | 1.000 | 7 | 11 | 0 | 4 | 0 | NA | NA | NA | NA | NA | NA | NA | NA |
| numeric | index | 0 | 1.000 | NA | NA | NA | NA | NA | 972.344 | 486.432 | 1.00 | 598.00 | 1070.00 | 1361.75 | 1900 | <U+2585><U+2585><U |
| numeric | eletkora | 0 | 1.000 | NA | NA | NA | NA | NA | 52.522 | 10.725 | 18.96 | 45.00 | 52.50 | 61.00 | 80 | <U+2581><U+2585><U |
| numeric | gyermeke | 0 | 1.000 | NA | NA | NA | NA | NA | 1.750 | 1.237 | 0.85 | 1.00 | 2.00 | 2.00 | 9 | <U+2585><U+2587><U |
| numeric | anyagi | 0 | 1.000 | NA | NA | NA | NA | NA | 3.100 | 0.564 | 1.96 | 3.00 | 3.00 | 3.00 | 5 | <U+2581><U+2581><U |
| numeric | pic_elmeny_pefont | 0 | 1.000 | NA | NA | NA | NA | NA | 58.460 | 20.418 | 5.01 | 40.00 | 60.00 | 80.00 | 100 | <U+2582><U+2585><U |
| numeric | risti_fi | 3 | 0.994 | NA | NA | NA | NA | NA | 4.164 | 0.957 | 2.09 | 4.00 | 4.00 | 5.00 | 6 | <U+2581><U+2583><U |
| numeric | alt_lelki | 4 | 0.992 | NA | NA | NA | NA | NA | 4.161 | 0.907 | 2.62 | 4.00 | 4.00 | 5.00 | 6 | <U+2582><U+2583><U |
| numeric | alt_eg_a | 7 | 0.986 | NA | NA | NA | NA | NA | 4.121 | 0.704 | 1.61 | 3.00 | 4.00 | 5.00 | 6 | <U+2582><U+2583><U |
| numeric | fizero | 6 | 0.988 | NA | NA | NA | NA | NA | 4.153 | 0.867 | 3.07 | 3.00 | 4.00 | 5.00 | 6 | <U+2582><U+2585><U |
| numeric | aicocska | 0 | 1.000 | NA | NA | NA | NA | NA | 2.570 | 1.435 | 5.89 | 2.00 | 2.00 | 3.00 | 7 | <U+2587><U+2585><U |
| numeric | aggodalo | 0 | 1.000 | NA | NA | NA | NA | NA | 2.772 | 1.432 | 5.23 | 2.00 | 2.50 | 4.00 | 6 | <U+2587><U+2583><U |
| numeric | ideges | 0 | 1.000 | NA | NA | NA | NA | NA | 2.490 | 1.374 | 5.61 | 1.00 | 2.00 | 3.00 | 6 | <U+2587><U+2582><U |
| numeric | feszult | 0 | 1.000 | NA | NA | NA | NA | NA | 2.610 | 1.399 | 1.97 | 1.00 | 2.00 | 3.00 | 6 | <U+2587><U+2583><U |
| numeric | nyugtala | 0 | 1.000 | NA | NA | NA | NA | NA | 2.406 | 1.457 | 8.03 | 1.00 | 2.00 | 3.00 | 6 | <U+2587><U+2582><U |

| skim_type | skim_variable | n_missing | complete_rate | character.min | character.max | character.empty | character.n_unique | character.whitespace | numeric.mean | numeric.sd | numeric.p0 | numeric.p25 | numeric.p50 | numeric.p75 | numeric.p100 | numeric.hist |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| numeric | diener1 | 0 | 1.000 | NA | NA | NA | NA | NA | 5.5140 | 1.3286 | 1 | 5.00 | 6.0000 | 6.0000 | 7 | <U+2581><U+2581><U... |
| numeric | diener2 | 0 | 1.000 | NA | NA | NA | NA | NA | 5.2280 | 1.3219 | 1 | 4.00 | 6.0000 | 6.0000 | 7 | <U+2581><U+2581><U... |
| numeric | diener3 | 0 | 1.000 | NA | NA | NA | NA | NA | 5.5740 | 1.2310 | 1 | 5.00 | 6.0000 | 6.0000 | 7 | <U+2581><U+2581><U... |
| numeric | diener4 | 0 | 1.000 | NA | NA | NA | NA | NA | 5.4180 | 1.2611 | 1 | 5.00 | 6.0000 | 6.0000 | 7 | <U+2581><U+2581><U... |
| numeric | diener5 | 0 | 1.000 | NA | NA | NA | NA | NA | 5.7560 | 1.0727 | 1 | 5.00 | 6.0000 | 7.0000 | 7 | <U+2581><U+2581><U... |
| numeric | diener6 | 0 | 1.000 | NA | NA | NA | NA | NA | 5.7680 | 1.0245 | 1 | 5.00 | 6.0000 | 7.0000 | 7 | <U+2581><U+2581><U... |
| numeric | diener7 | 0 | 1.000 | NA | NA | NA | NA | NA | 5.6200 | 1.3738 | 1 | 5.00 | 6.0000 | 7.0000 | 7 | <U+2581><U+2581><U... |
| numeric | diener8 | 0 | 1.000 | NA | NA | NA | NA | NA | 5.5320 | 1.2148 | 1 | 5.00 | 6.0000 | 6.0000 | 7 | <U+2581><U+2581><U... |
| numeric | jollet | 0 | 1.000 | NA | NA | NA | NA | NA | 4.4886 | 1.6610 | 1 | 4.0000 | 4.6667 | 5.3333 | 6 | <U+2581><U+2582><U... |
| numeric | savor | 0 | 1.000 | NA | NA | NA | NA | NA | 4.5146 | 1.7133 | 1 | 4.0000 | 4.6667 | 5.0000 | 6 | <U+2581><U+2582><U... |
| numeric | aic_vhat | 0 | 1.000 | NA | NA | NA | NA | NA | 4.5956 | 0.9393 | 1 | 4.0000 | 4.8000 | 5.2000 | 6 | <U+2581><U+2581><U... |
| numeric | oicreg | 0 | 1.000 | NA | NA | NA | NA | NA | 4.2173 | 1.2217 | 1 | 3.3333 | 4.3333 | 5.0833 | 6 | <U+2582><U+2583><U... |
| numeric | rezil | 0 | 1.000 | NA | NA | NA | NA | NA | 4.0446 | 1.7210 | 1 | 3.3333 | 4.0000 | 4.6667 | 6 | <U+2581><U+2583><U... |
| numeric | nic_flow | 0 | 1.000 | NA | NA | NA | NA | NA | 4.7739 | 0.8909 | 1 | 4.3333 | 5.0000 | 5.3333 | 6 | <U+2581><U+2581><U... |
| numeric | gic_jerz | 0 | 1.000 | NA | NA | NA | NA | NA | 4.1528 | 1.0313 | 1 | 3.4000 | 4.2000 | 5.0000 | 6 | <U+2581><U+2583><U... |
| numeric | gic_jpszi | 0 | 1.000 | NA | NA | NA | NA | NA | 4.2820 | 0.9892 | 1 | 3.7500 | 4.5000 | 5.0000 | 6 | <U+2581><U+2582><U... |
| numeric | gic_jszoc | 0 | 1.000 | NA | NA | NA | NA | NA | 3.9495 | 1.0789 | 1 | 3.0000 | 4.0000 | 5.0000 | 6 | <U+2582><U+2585><U... |
| numeric | gic_jspir | 0 | 1.000 | NA | NA | NA | NA | NA | 4.2455 | 1.2003 | 1 | 3.5000 | 4.2500 | 5.2500 | 6 | <U+2582><U+2583><U... |
| numeric | gic_jerzps | 0 | 1.000 | NA | NA | NA | NA | NA | 4.2174 | 0.9720 | 1 | 3.6187 | 4.3375 | 4.9250 | 6 | <U+2581><U+2582><U... |
| numeric | gic_jspszoc | 0 | 1.000 | NA | NA | NA | NA | NA | 4.0975 | 1.0382 | 1 | 3.3750 | 4.1250 | 5.0000 | 6 | <U+2581><U+2583><U... |
| numeric | pik_mm | 0 | 1.000 | NA | NA | NA | NA | NA | 3.0273 | 1.3426 | 1 | 2.6667 | 3.1667 | 3.5000 | 4 | <U+2581><U+2582><U... |
| numeric | pik_av | 0 | 1.000 | NA | NA | NA | NA | NA | 3.1750 | 1.0623 | 1 | 2.7500 | 3.2500 | 3.7500 | 4 | <U+2581><U+2581><U... |
| numeric | pik_onr | 0 | 1.000 | NA | NA | NA | NA | NA | 3.0640 | 1.0791 | 1 | 2.6667 | 3.3333 | 3.6667 | 4 | <U+2581><U+2582><U... |
| numeric | pik_rez | 0 | 1.000 | NA | NA | NA | NA | NA | 3.1506 | 0.8742 | 1 | 2.6667 | 3.3333 | 3.6667 | 4 | <U+2581><U+2582><U... |
| numeric | pic_poz_erz | 0 | 1.000 | NA | NA | NA | NA | NA | 21.7425 | 2.7513 | 1 | 20.0000 | 23.0000 | 26.0000 | 30 | <U+2581><U+2581><U... |
| numeric | pic_elmely | 0 | 1.000 | NA | NA | NA | NA | NA | 22.3604 | 3.5858 | 1 | 20.0000 | 23.0000 | 26.0000 | 30 | <U+2581><U+2582><U... |
| numeric | pic_poz_kapc | 0 | 1.000 | NA | NA | NA | NA | NA | 22.4425 | 5.6710 | 1 | 19.0000 | 24.0000 | 26.0000 | 30 | <U+2581><U+2581><U... |
| numeric | pic_ert_cel | 0 | 1.000 | NA | NA | NA | NA | NA | 23.3740 | 6.7909 | 1 | 21.0000 | 24.0000 | 27.0000 | 30 | <U+2581><U+2581><U... |
| numeric | pic_telj | 0 | 1.000 | NA | NA | NA | NA | NA | 22.7160 | 3.4200 | 1 | 21.0000 | 24.0000 | 26.0000 | 30 | <U+2581><U+2581><U... |
| numeric | pic_boldog | 0 | 1.000 | NA | NA | NA | NA | NA | 7.4040 | 2.0217 | 1 | 7.0000 | 8.0000 | 9.0000 | 10 | <U+2581><U+2581><U... |
| numeric | pic_egeszs | 0 | 1.000 | NA | NA | NA | NA | NA | 22.2605 | 2.2480 | 1 | 19.0000 | 23.0000 | 26.0000 | 30 | <U+2581><U+2581><U... |
| numeric | pic_neg_erz | 0 | 1.000 | NA | NA | NA | NA | NA | 9.7660 | 5.0708 | 1 | 5.0000 | 9.0000 | 13.0000 | 30 | <U+2587><U+2587><U... |
| numeric | pic_magany | 0 | 1.000 | NA | NA | NA | NA | NA | 3.1620 | 3.0565 | 1 | 0.0000 | 2.0000 | 5.0000 | 10 | <U+2587><U+2582><U... |
| numeric | pierma | 0 | 1.000 | NA | NA | NA | NA | NA | 169.3781 | 0.2672 | 1 | 155.00 | 175.50 | 191.00 | 250 | <U+2581><U+2581><U... |

# 1. Run an EFA with the following variables

```r
vars <- c("p_elmeny_percent", "testi_fi", "alt_lelki", "alt_eg_all", "fizero", "arcocska", "aggodalo", "
```

```r
efa_data <-
  processed %>%
  dplyr::select(all_of(vars))
```

We can calculate the correlation matrix first.

```r
cor_matrix <- cor(efa_data, use = "pairwise.complete.obs")
```

Now we can calculate the Eigenvalues from this.

```
eigenvalues <- eigen(cor_matrix)

eigenvalues$values
```

```
##  [1] 4.9546529 2.1023843 0.7652170 0.5487025 0.4734496 0.3706545 0.2639130
##  [8] 0.2114798 0.1962446 0.1133020
```

Now we can look at a screeplot to decide on the number of factors.

```
scree(cor_matrix, factors = FALSE)
```



**Scree plot**

Based on the screeplot I will go with 2 factors. I will also use oblimin rotation. Using maximum likelihood estimation method.

```
fa_res <- fa(efa_data, nfactors = 2, rotate = "oblimin", fm = "ml")
```

```
## Loading required namespace: GPArotation
```

```
fa_res
```

```
## Factor Analysis using method =  ml
## Call: fa(r = efa_data, nfactors = 2, rotate = "oblimin", fm = "ml")
## Standardized loadings (pattern matrix) based upon correlation matrix
```

```
##                     ML1   ML2   h2   u2  com
## p_elmeny_percent -0.48   0.17 0.31 0.69 1.3
## testi_fi           0.04   0.90 0.79 0.21 1.0
## alt_lelki         -0.34   0.47 0.44 0.56 1.8
## alt_eg_all        -0.01   0.87 0.77 0.23 1.0
## fizero             0.02   0.85 0.71 0.29 1.0
## arcocska           0.48  -0.33 0.44 0.56 1.8
## aggodalo           0.88   0.06 0.75 0.25 1.0
## ideges             0.94   0.02 0.86 0.14 1.0
## feszult            0.95   0.02 0.90 0.10 1.0
## nyugtala           0.65  -0.09 0.47 0.53 1.0
##
##                      ML1   ML2
## SS loadings          3.67 2.77
## Proportion Var       0.37 0.28
## Cumulative Var       0.37 0.64
## Proportion Explained 0.57 0.43
## Cumulative Proportion 0.57 1.00
##
##  With factor correlations of
##       ML1   ML2
## ML1  1.00 -0.31
## ML2 -0.31  1.00
##
## Mean item complexity =  1.2
## Test of the hypothesis that 2 factors are sufficient.
##
## The degrees of freedom for the null model are  45  and the objective function was  6.96 with Chi Squa
## The degrees of freedom for the model are 26  and the objective function was  0.25
##
## The root mean square of the residuals (RMSR) is  0.04
## The df corrected root mean square of the residuals is  0.05
##
## The harmonic number of observations is  496 with the empirical chi square  73.41  with prob <  2.1e-(
## The total number of observations was  500  with Likelihood Chi Square =  124.08  with prob <  9.5e-1!
##
## Tucker Lewis Index of factoring reliability =  0.95
## RMSEA index =  0.087  and the 90 % confidence intervals are  0.072 0.103
## BIC =  -37.5
## Fit based upon off diagonal values = 0.99
## Measures of factor score adequacy
##                                                    ML1  ML2
## Correlation of (regression) scores with factors   0.98 0.96
## Multiple R square of scores with factors          0.95 0.91
## Minimum correlation of possible factor scores     0.91 0.83
```

Which items that has a commonality (h2) lower than 0.25? No, there was not.

The two factors explain 65% of the variance of the data based on Proportion Var values.

## 2. How big is the KMO and Bartlett test?

Check KMO.

4

```
KMO(cor_matrix)
```

```
## Kaiser-Meyer-Olkin factor adequacy
## Call: KMO(r = cor_matrix)
## Overall MSA =  0.88
## MSA for each item =
## p_elmeny_percent         testi_fi         alt_lelki         alt_eg_all
##             0.94             0.81             0.91             0.83
##            fizero         arcocska         aggodalo           ideges
##             0.86             0.91             0.91             0.85
##           feszult         nyugtala
##             0.84             0.96
```

KMO can be used to check whether the data is adequate for factor analysis based on testing for common variance between variables.

It seems like that both the individual and the summarized KMO tests are great.

Check Barlett sphericity test.

```
cortest.bartlett(cor_matrix, n = 500, diag = FALSE)
```

```
## $chisq
## [1] 3442.281
##
## $p.value
## [1] 0
##
## $df
## [1] 45
```

The p value can never be 0 but I suspect that the function rounded the value based on the size of the test statistics. Therefore, we can reject the null hypothesis, so we can perform an EFA.

# 3. Is the model bad based on significance?

```
fa_res
```

```
## Factor Analysis using method =  ml
## Call: fa(r = efa_data, nfactors = 2, rotate = "oblimin", fm = "ml")
## Standardized loadings (pattern matrix) based upon correlation matrix
##                   ML1   ML2   h2   u2 com
## p_elmeny_percent -0.48  0.17 0.31 0.69 1.3
## testi_fi          0.04  0.90 0.79 0.21 1.0
## alt_lelki        -0.34  0.47 0.44 0.56 1.8
## alt_eg_all       -0.01  0.87 0.77 0.23 1.0
## fizero            0.02  0.85 0.71 0.29 1.0
## arcocska          0.48 -0.33 0.44 0.56 1.8
## aggodalo          0.88  0.06 0.75 0.25 1.0
## ideges            0.94  0.02 0.86 0.14 1.0
```

```
## feszult             0.95  0.02 0.90 0.10 1.0
## nyugtala            0.65 -0.09 0.47 0.53 1.0
##
##                      ML1  ML2
## SS loadings         3.67 2.77
## Proportion Var      0.37 0.28
## Cumulative Var      0.37 0.64
## Proportion Explained 0.57 0.43
## Cumulative Proportion 0.57 1.00
##
##  With factor correlations of
##       ML1   ML2
## ML1  1.00 -0.31
## ML2 -0.31  1.00
##
## Mean item complexity =  1.2
## Test of the hypothesis that 2 factors are sufficient.
##
## The degrees of freedom for the null model are  45  and the objective function was  6.96 with Chi Squa
## The degrees of freedom for the model are 26  and the objective function was  0.25
##
## The root mean square of the residuals (RMSR) is  0.04
## The df corrected root mean square of the residuals is  0.05
##
## The harmonic number of observations is  496 with the empirical chi square  73.41  with prob <  2.1e-0
## The total number of observations was  500  with Likelihood Chi Square =  124.08  with prob <  9.5e-1!
##
## Tucker Lewis Index of factoring reliability =  0.95
## RMSEA index =  0.087  and the 90 % confidence intervals are  0.072 0.103
## BIC =  -37.5
## Fit based upon off diagonal values = 0.99
## Measures of factor score adequacy
##                                                 ML1  ML2
## Correlation of (regression) scores with factors  0.98 0.96
## Multiple R square of scores with factors        0.95 0.91
## Minimum correlation of possible factor scores   0.91 0.83
```

Seemingly the psych package only computes confidence interval for RMSEA to check the model fit. The results are RMSEA = 0.087 CI90[0.072, 0.103]. As the confidence interval does not include 0 I conclude that it is a good model fit.

The function also returns the Chi square with a p value. The p value here is significant.

*The total number of observations was 500 with Likelihood Chi Square = 124.08 with prob < 9.5e-15*

## 4. Using a promax rotation

Here I use a promax rotation therefore, I allow for correlation between the factors.

```
fa_promax_res <- fa(efa_data, nfactors = 2, rotate = "promax", fm = "ml")

fa_promax_res
```

```
## Factor Analysis using method =  ml
## Call: fa(r = efa_data, nfactors = 2, rotate = "promax", fm = "ml")
## Standardized loadings (pattern matrix) based upon correlation matrix
##                    ML1   ML2   h2   u2  com
## p_elmeny_percent -0.49  0.12 0.31 0.69 1.1
## testi_fi          0.13  0.95 0.79 0.21 1.0
## alt_lelki        -0.31  0.45 0.44 0.56 1.8
## alt_eg_all        0.07  0.91 0.77 0.23 1.0
## fizero            0.10  0.89 0.71 0.29 1.0
## arcocska          0.48 -0.29 0.44 0.56 1.6
## aggodalo          0.94  0.17 0.75 0.25 1.1
## ideges            0.99  0.13 0.86 0.14 1.0
## feszult           1.00  0.14 0.90 0.10 1.0
## nyugtala          0.68 -0.01 0.47 0.53 1.0
##
##                     ML1  ML2
## SS loadings        3.74 2.70
## Proportion Var     0.37 0.27
## Cumulative Var     0.37 0.64
## Proportion Explained  0.58 0.42
## Cumulative Proportion 0.58 1.00
##
##  With factor correlations of
##       ML1   ML2
## ML1  1.00 -0.48
## ML2 -0.48  1.00
##
## Mean item complexity =  1.2
## Test of the hypothesis that 2 factors are sufficient.
##
## The degrees of freedom for the null model are  45  and the objective function was  6.96 with Chi Squa
## The degrees of freedom for the model are 26  and the objective function was  0.25
##
## The root mean square of the residuals (RMSR) is  0.04
## The df corrected root mean square of the residuals is  0.05
##
## The harmonic number of observations is  496 with the empirical chi square  73.41  with prob <  2.1e-(
## The total number of observations was  500  with Likelihood Chi Square =  124.08  with prob <  9.5e-1!
##
## Tucker Lewis Index of factoring reliability =  0.95
## RMSEA index =  0.087  and the 90 % confidence intervals are  0.072 0.103
## BIC =  -37.5
## Fit based upon off diagonal values = 0.99
## Measures of factor score adequacy
##                                                     ML1  ML2
## Correlation of (regression) scores with factors    0.98 0.96
## Multiple R square of scores with factors           0.96 0.92
## Minimum correlation of possible factor scores      0.91 0.84
```

Plotting the loadings.

```
plot(fa_promax_res$loadings[,1],
     fa_promax_res$loadings[,2],
    xlab = "Factor 1",
```
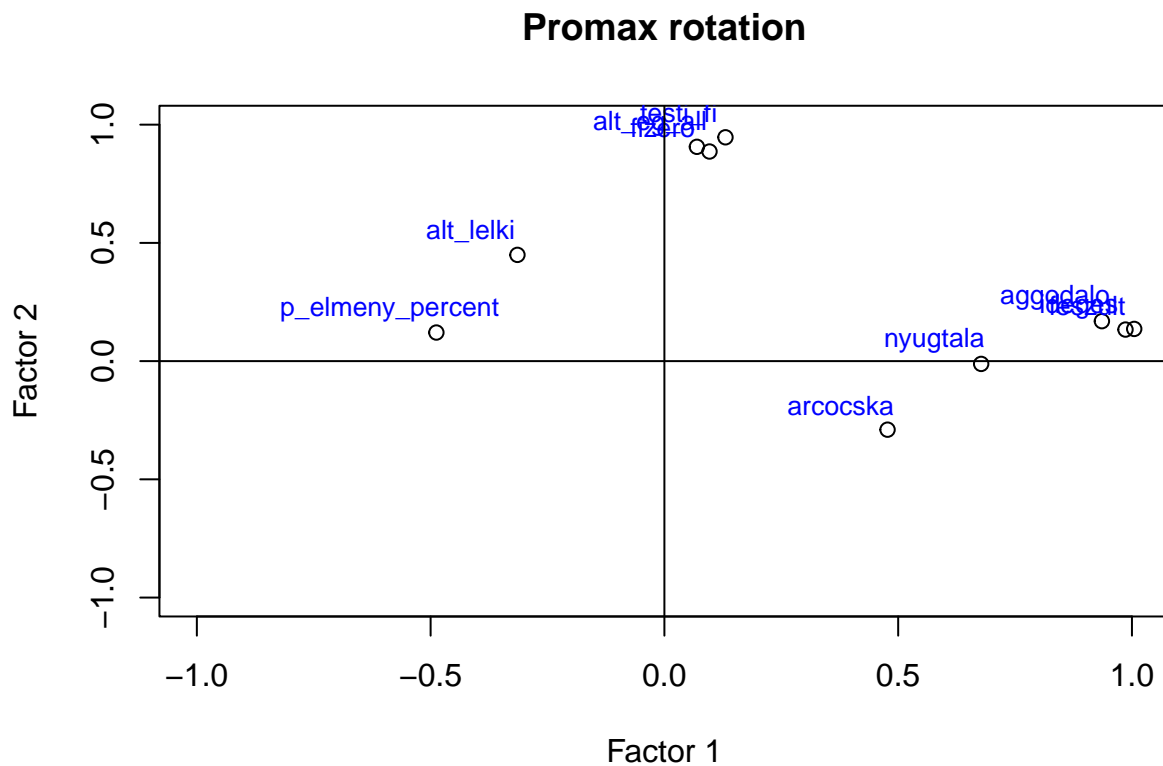
```
    ylab = "Factor 2",
    ylim = c(-1,1),
    xlim = c(-1,1),
    main = "Promax rotation")

text(
  fa_promax_res$loadings[,1]-0.1,
    fa_promax_res$loadings[,2]+0.1,
    colnames(cor_matrix),
    col="blue", cex = .8)

abline(h = 0, v = 0)
```

## Promax rotation



Factor 1

It seems like that only two items have high complexity. p_elmeny_percent, aggoddalo, arcocska, idegesm feszult, and nyugtala belongs to the first factor. p_elmeny has a negative correlation with the factor. testi_fi, alt_lelki, alt_eg_all, fizero belongs to the second factor.

## 5. Is PAF or ML better?

If the variables are normally distributed we can run ML, however, if they are not we should rather run PAF.

Lets check the normality of the variables.

```
map(efa_data, shapiro.test)
```

```
## $p_elmeny_percent
##
##  Shapiro-Wilk normality test
##
## data:  .x[[i]]
## W = 0.94702, p-value = 2.191e-12
##
##
## $testi_fi
##
##  Shapiro-Wilk normality test
##
## data:  .x[[i]]
## W = 0.89491, p-value < 2.2e-16
##
##
## $alt_lelki
##
##  Shapiro-Wilk normality test
##
## data:  .x[[i]]
## W = 0.9095, p-value < 2.2e-16
##
##
## $alt_eg_all
##
##  Shapiro-Wilk normality test
##
## data:  .x[[i]]
## W = 0.91536, p-value = 5.518e-16
##
##
## $fizero
##
##  Shapiro-Wilk normality test
##
## data:  .x[[i]]
## W = 0.92239, p-value = 2.825e-15
##
##
## $arcocska
##
##  Shapiro-Wilk normality test
##
## data:  .x[[i]]
## W = 0.90318, p-value < 2.2e-16
##
##
## $aggodalo
##
##  Shapiro-Wilk normality test
```

```
##
## data:  .x[[i]]
## W = 0.90149, p-value < 2.2e-16
##
##
## $ideges
##
##  Shapiro-Wilk normality test
##
## data:  .x[[i]]
## W = 0.8751, p-value < 2.2e-16
##
##
## $feszult
##
##  Shapiro-Wilk normality test
##
## data:  .x[[i]]
## W = 0.88703, p-value < 2.2e-16
##
##
## $nyugtala
##
##  Shapiro-Wilk normality test
##
## data:  .x[[i]]
## W = 0.84454, p-value < 2.2e-16
```

All of them are significant so we should rather run PAF than ML.