

Solution for the assignment of the seventh 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 | single | character | double | integer | boolean | character | numeric | integer | double | integer | double | integer | double | integer | double | 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 | 972.3448604 | 327165 | 98.001070 | 0.00061 | 757000 | <U+2585> | <U+2585> | <U+2585> | |
| numeric | eltekora | 0 | 1.000 | NA | NA | NA | NA | NA | 52.522007 | 2578.96 | 45.00052 | 50061 | 0.000800 | <U+2581> | <U+2585> | <U+2585> | |
| numeric | germeko | 0 | 1.000 | NA | NA | NA | NA | NA | 1.750002 | 3715.85 | 1.00021 | 0002 | 0.000009 | <U+2585> | <U+2587> | <U+2587> | |
| numeric | nyagi | 0 | 1.000 | NA | NA | NA | NA | NA | 3.100005 | 6447.96 | 3.00031 | 0003 | 0.000005 | <U+2581> | <U+2581> | <U+2581> | |
| numeric | pic elmeny | 0 | 1.000 | NA | NA | NA | NA | NA | 58.460204 | 1835.01 | 40.00060 | 00080 | 0.000900 | <U+2582> | <U+2585> | <U+2585> | |
| numeric | testi_fi | 3 | 0.994 | NA | NA | NA | NA | NA | 4.164999 | 9572.09 | 4.00004 | 0005 | 0.000006 | <U+2581> | <U+2583> | <U+2583> | |
| numeric | lelki | 4 | 0.992 | NA | NA | NA | NA | NA | 4.161290 | 7223.62 | 4.00004 | 0005 | 0.000006 | <U+2582> | <U+2583> | <U+2583> | |
| numeric | eg_all | 0 | 0.986 | NA | NA | NA | NA | NA | 4.121704 | 1957.61 | 3.00004 | 0005 | 0.000006 | <U+2582> | <U+2583> | <U+2583> | |
| numeric | fizero | 6 | 0.988 | NA | NA | NA | NA | NA | 4.153846 | 8673.07 | 3.00004 | 0005 | 0.000006 | <U+2582> | <U+2585> | <U+2585> | |
| numeric | cocska | 0 | 1.000 | NA | NA | NA | NA | NA | 2.570001 | 04355.39 | 2.00021 | 0003 | 0.000007 | <U+2587> | <U+2585> | <U+2585> | |
| numeric | igodalo | 0 | 1.000 | NA | NA | NA | NA | NA | 2.772004 | 3252.23 | 2.00021 | 50004 | 0.000006 | <U+2587> | <U+2583> | <U+2583> | |
| numeric | ideges | 0 | 1.000 | NA | NA | NA | NA | NA | 2.490003 | 7465.61 | 1.00021 | 0003 | 0.000006 | <U+2587> | <U+2582> | <U+2582> | |
| numeric | fizult | 0 | 1.000 | NA | NA | NA | NA | NA | 2.610003 | 9921.97 | 1.00021 | 0003 | 0.000006 | <U+2587> | <U+2583> | <U+2583> | |
| numeric | ugtala | 0 | 1.000 | NA | NA | NA | NA | NA | 2.406004 | 5788.03 | 1.00021 | 0003 | 0.000006 | <U+2587> | <U+2582> | <U+2582> | |

[illegible]

1. Look for a variable that's variance is explained by the 8 Diener items by at least 65%.

I will first calculate the total score for the Diener flourishing scale.

```
diener_data <-  
  processed %>%  
  mutate(diener_sum = diener1 + diener2 + diener3 + diener4 + diener5 + diener6 + diener7 + diener8)
```

For this task I am investigating only interval variables. To calculate the Rsquared I calculate the Pearson correlation coefficient and square it.

```
interval_vars <-  
  processed %>%
```

```

select(26:50) %>%
  names()

rsq <- function (data, x, y) cor(data[[x]], data[[y]]) ^ 2

diener_res <-
  tibble::tibble(
    variable = interval_vars,
    rsquared = map_dbl(variable,
      ~ rsq(diener_data, .x, "diener_sum"))
  ) %>%
  mutate(
    rsquared = round(rsquared, 2),
    flag = case_when(rsquared >= .65 ~ TRUE,
      rsquared < .65 ~ FALSE)
  ) %>%
  arrange(desc(rsquared))

```

The total score on the Diener 8 item scale only explains more than 65% of the variance of the *g_jerzpsz* scale.

2. Lets test whether adding aggodalom and ideges variables as main effects will increase the R2.

To do this we have to create a linear regression model.

```

m <- lm(g_jerzpsz ~ diener_sum + aggodalo + ideges, data = diener_data)

summary(m)

##
## Call:
## lm(formula = g_jerzpsz ~ diener_sum + aggodalo + ideges, data = diener_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.31675 -0.35127  0.02655  0.40331  2.07871
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.457202   0.157206   2.908  0.0038 **
## diener_sum    0.090846   0.003008  30.201 <2e-16 ***
## aggodalo     -0.044526   0.028036  -1.588  0.1129
## ideges       -0.060587   0.029656  -2.043  0.0416 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.5348 on 496 degrees of freedom
## Multiple R-squared:  0.6991, Adjusted R-squared:  0.6973
## F-statistic: 384.1 on 3 and 496 DF, p-value: < 2.2e-16

```

Adding those two variables indeed increased the Rsquared to 0.6991.

3. Transforming the kor variable

```
kor_data <-  
  processed %>%  
  dplyr::select(eletkora, index) %>%  
  # bind_cols(kor_data, poly(kor_data$eletkora), degree = 3)  
  mutate(korz = scale(eletkora),  
         korz2 = korz^2,  
         korz3 = korz^3,  
         korz4 = korz^4)
```

Now, lets create a correlation matrix including these variables.

```
Hmisc::rcorr(as.matrix(kor_data), type = "pearson")
```

```
##          eletkora index  korz korz2 korz3 korz4  
## eletkora      1.00  0.00  1.00 -0.09  0.80 -0.14  
## index         0.00  1.00  0.00  0.13 -0.04  0.10  
## korz          1.00  0.00  1.00 -0.09  0.80 -0.14  
## korz2        -0.09  0.13 -0.09  1.00 -0.22  0.88  
## korz3         0.80 -0.04  0.80 -0.22  1.00 -0.34  
## korz4        -0.14  0.10 -0.14  0.88 -0.34  1.00  
##  
## n= 500  
##  
##  
## P  
##          eletkora index  korz  korz2 korz3 korz4  
## eletkora          0.9189 0.0000 0.0501 0.0000 0.0014  
## index    0.9189          0.9189 0.0046 0.4237 0.0211  
## korz      0.0000  0.9189          0.0501 0.0000 0.0014  
## korz2     0.0501  0.0046 0.0501          0.0000 0.0000  
## korz3     0.0000  0.4237 0.0000 0.0000          0.0000  
## korz4     0.0014  0.0211 0.0014 0.0000 0.0000
```

4. Predicting PERMA by korz with polynomial regression. Which power has the largest effect on the outcome variable? Plot the results.

```
polyreg_data <-  
  kor_data %>%  
  left_join(., select(processed, perma, index), by = "index")  
  
m <- lm(perma ~ korz + korz2 + korz3 + korz4, data = polyreg_data)  
  
summary(m)
```

```
##
```

```
## Call:
## lm(formula = perma ~ korz + korz2 + korz3 + korz4, data = polyreg_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -152.046  -15.746    6.024   21.075   58.192
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  174.9017     1.9641  89.051 < 2e-16 ***
## korz          6.1223     2.3521   2.603 0.009521 **
## korz2        -8.0878     2.2799  -3.547 0.000426 ***
## korz3        -0.6376     0.7452  -0.856 0.392650
## korz4         0.9206     0.4145   2.221 0.026791 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 30.57 on 495 degrees of freedom
## Multiple R-squared:  0.0516, Adjusted R-squared:  0.04393
## F-statistic: 6.733 on 4 and 495 DF,  p-value: 2.783e-05
```

The second and fourth degree polynomial variables predicted the perma outcome variable significantly.

5. Look for a variable that has a third degree polynomial relationship with age. Plot the relationship.

To look for the relationship I will look at the correlation between the interval variables and the third degree polynomial of age. I am looking for a significant relationship.

```
third_data <-
  processed %>%
  left_join(., kor_data, by = "index")

third_res <-
  tibble::tibble(
    variable = interval_vars,
    cor_res = map(variable,
      ~ my_cor(
        data = third_data,
        x = .x,
        y = "korz3",
        method = "spearman"
      )),
    cor_r = map_dbl(cor_res, ~ pluck(.x, "estimate", "rho")),
    cor_p = map_dbl(cor_res, ~ pluck(.x, "p.value")),
    flagged = case_when(cor_p <= 0.05 ~ TRUE,
      TRUE ~ FALSE)
  ) %>%
  arrange(desc(cor_r))
```

```
## Warning in cor.test.default(data[[x]], data[[y]], method = method): Cannot
```



```
## compute exact p-value with ties

## Warning in cor.test.default(data[[x]], data[[y]], method = method): Cannot
## compute exact p-value with ties

## Warning in cor.test.default(data[[x]], data[[y]], method = method): Cannot
## compute exact p-value with ties

## Warning in cor.test.default(data[[x]], data[[y]], method = method): Cannot
## compute exact p-value with ties

## Warning in cor.test.default(data[[x]], data[[y]], method = method): Cannot
## compute exact p-value with ties

## Warning in cor.test.default(data[[x]], data[[y]], method = method): Cannot
## compute exact p-value with ties
```

The results show that 10 variables had a significant third degree polynomial relationship with the age variable. The following table is in descending order based on the Spearmans' rho.

```
third_res %>%
  select(-cor_res)
```

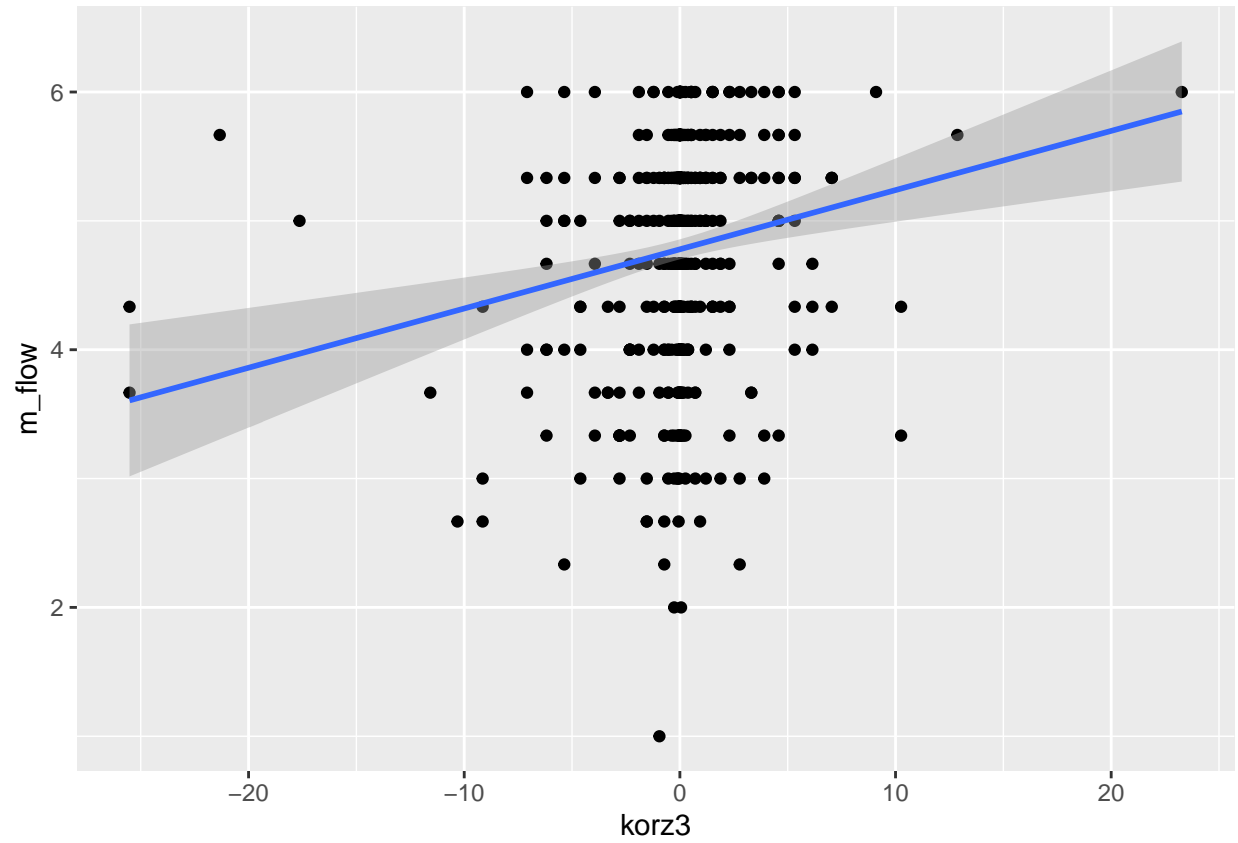
```
## # A tibble: 25 x 4
##   variable cor_r      cor_p flagged
##   <chr>    <dbl>    <dbl> <lgl>
## 1 m_flow  0.207 0.00000319 TRUE
## 2 onreg   0.188 0.0000243 TRUE
## 3 pik_onr 0.176 0.0000771 TRUE
## 4 p_boldog 0.172 0.000113 TRUE
## 5 p_poz_erz 0.168 0.000158 TRUE
## 6 pik_rez  0.165 0.000216 TRUE
## 7 perma    0.145 0.00115 TRUE
## 8 p_elmely 0.125 0.00505 TRUE
## 9 rezil    0.119 0.00776 TRUE
## 10 g_jerz  0.114 0.0111 TRUE
## # ... with 15 more rows
```

Now, I would like to create scatterplot showing the relationship between these variables and age.

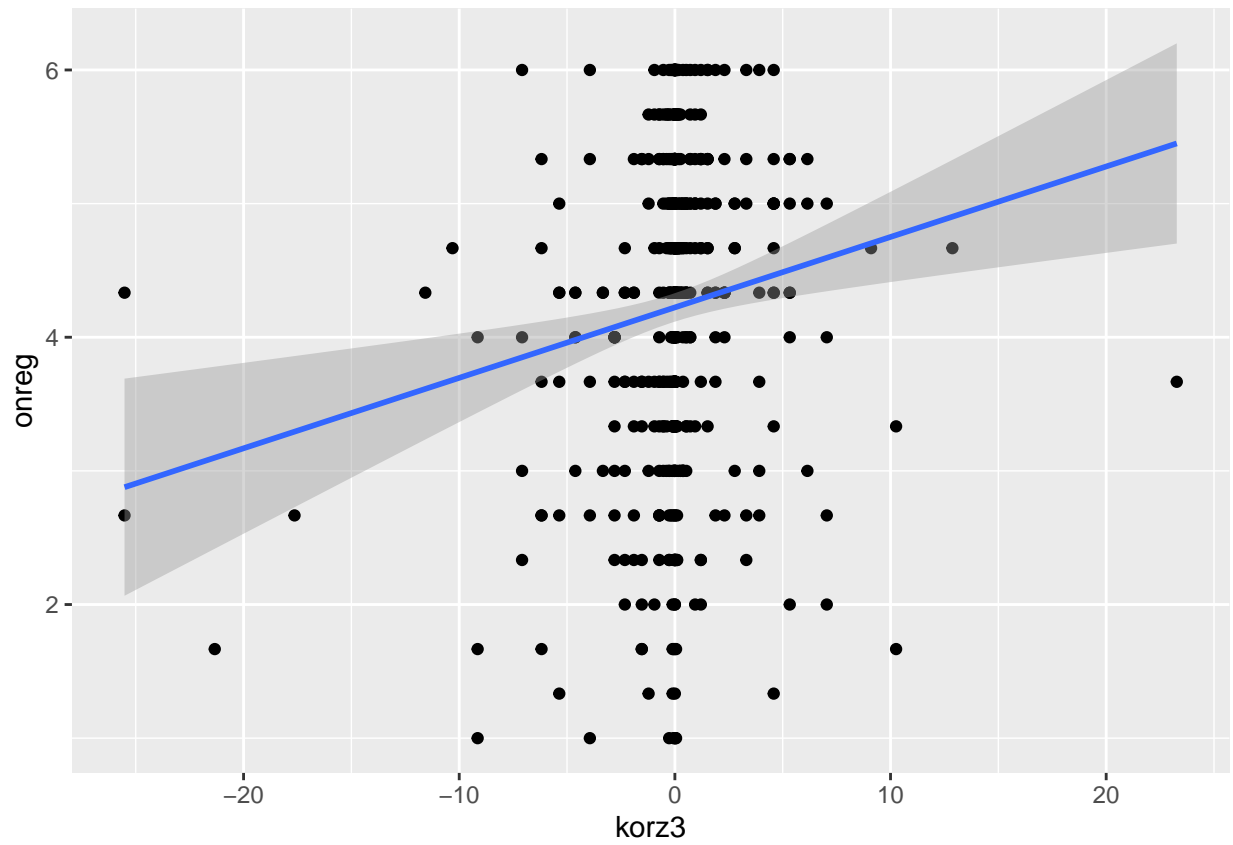
```
third_plot <-
  third_res %>%
  dplyr::filter(flagged) %>%
  mutate(plot = map(variable,
    ~ poly_plot(
      data = third_data,
      x = "korz3",
      y = .x
    )
  ))
```

```
third_plot$plot
```

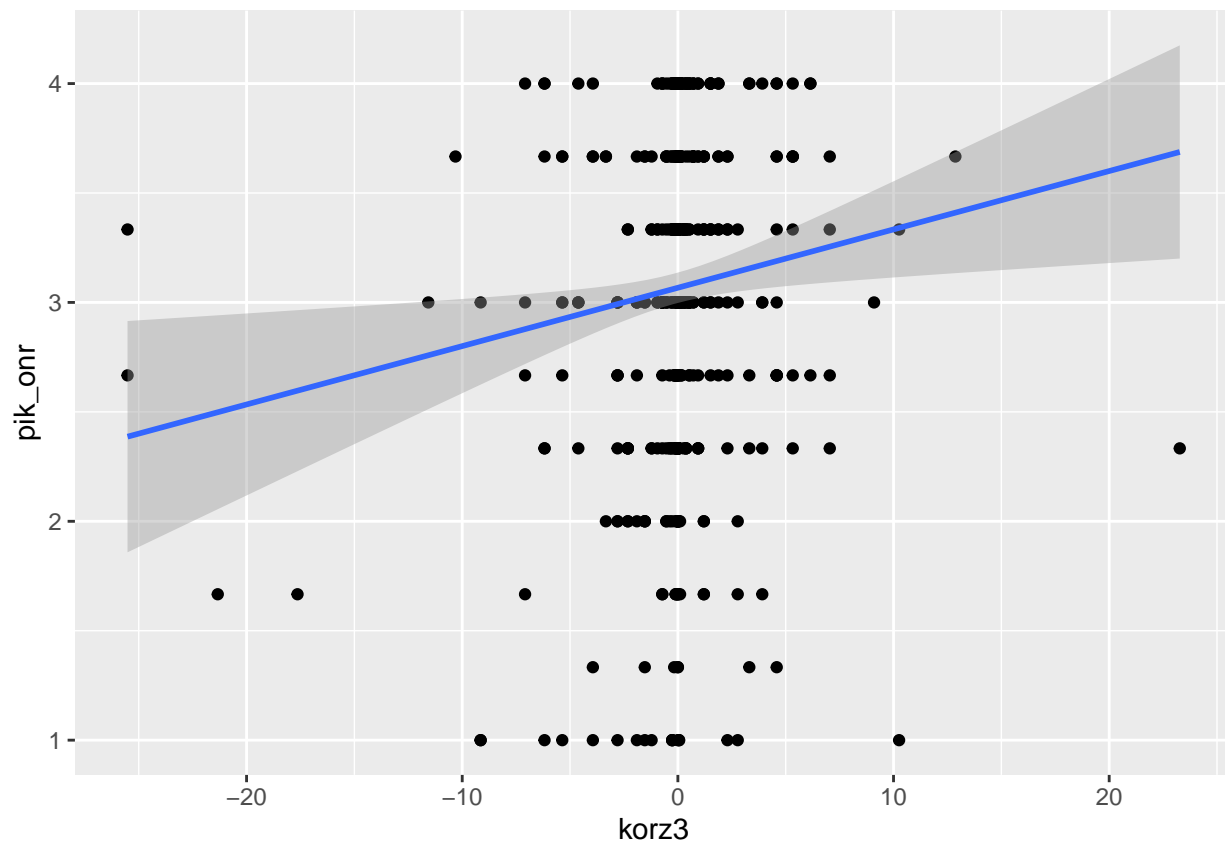
```
## [[1]]
```



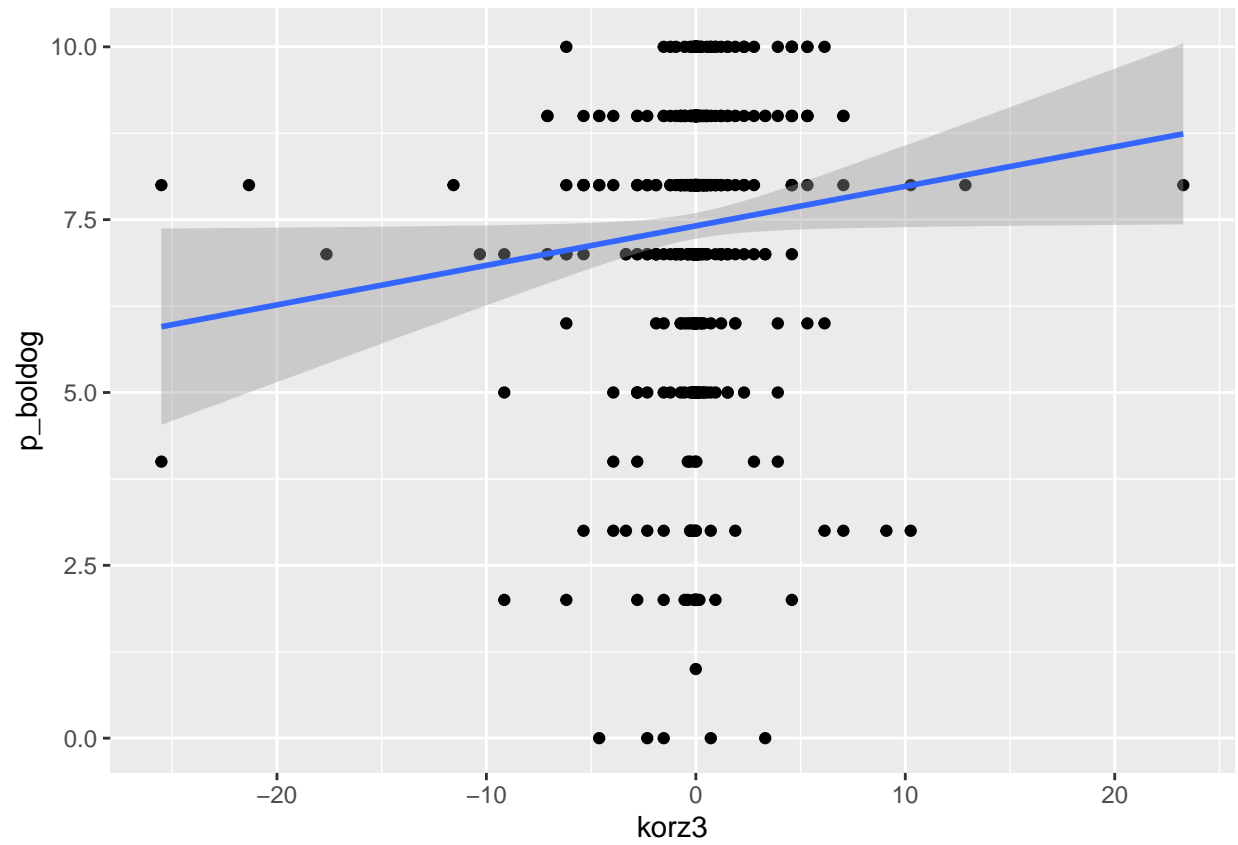
```
##  
## [[2]]
```

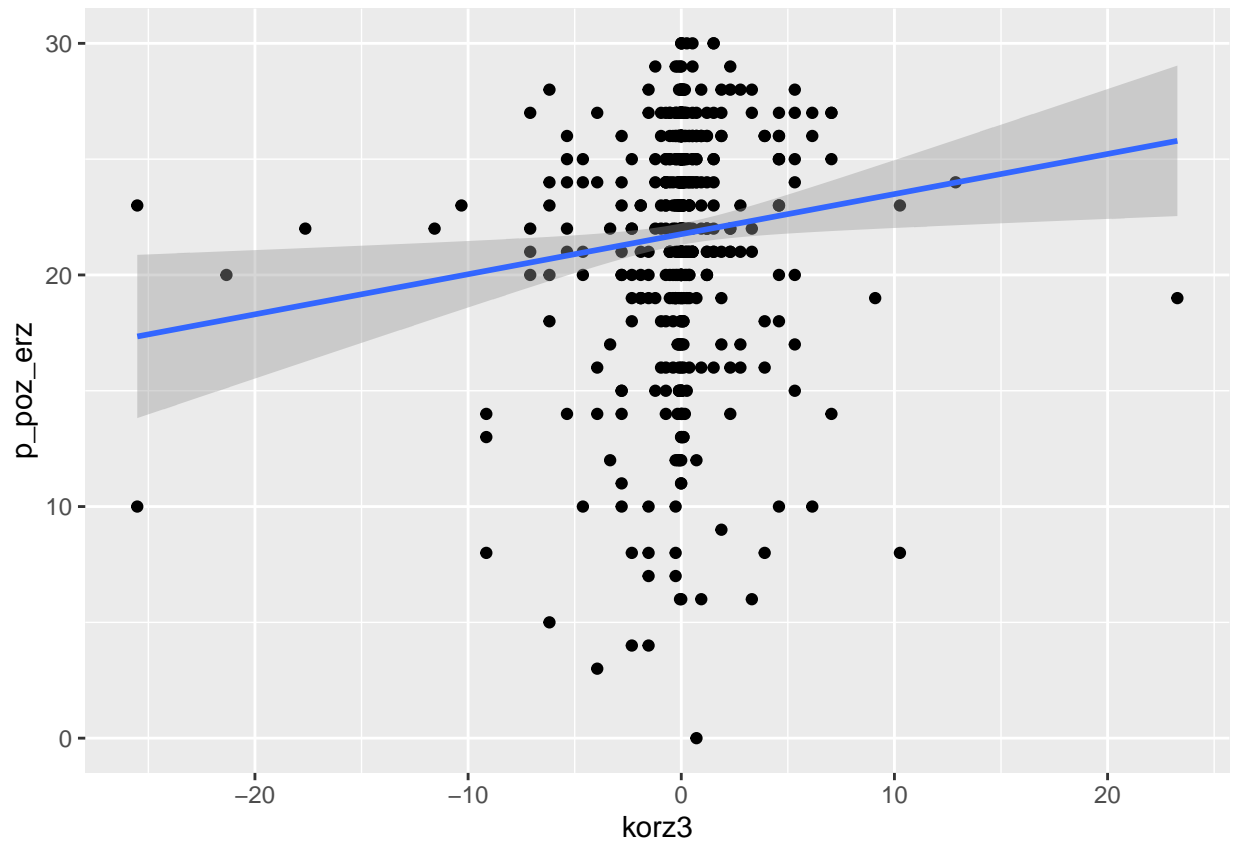
```
##  
## [[3]]
```



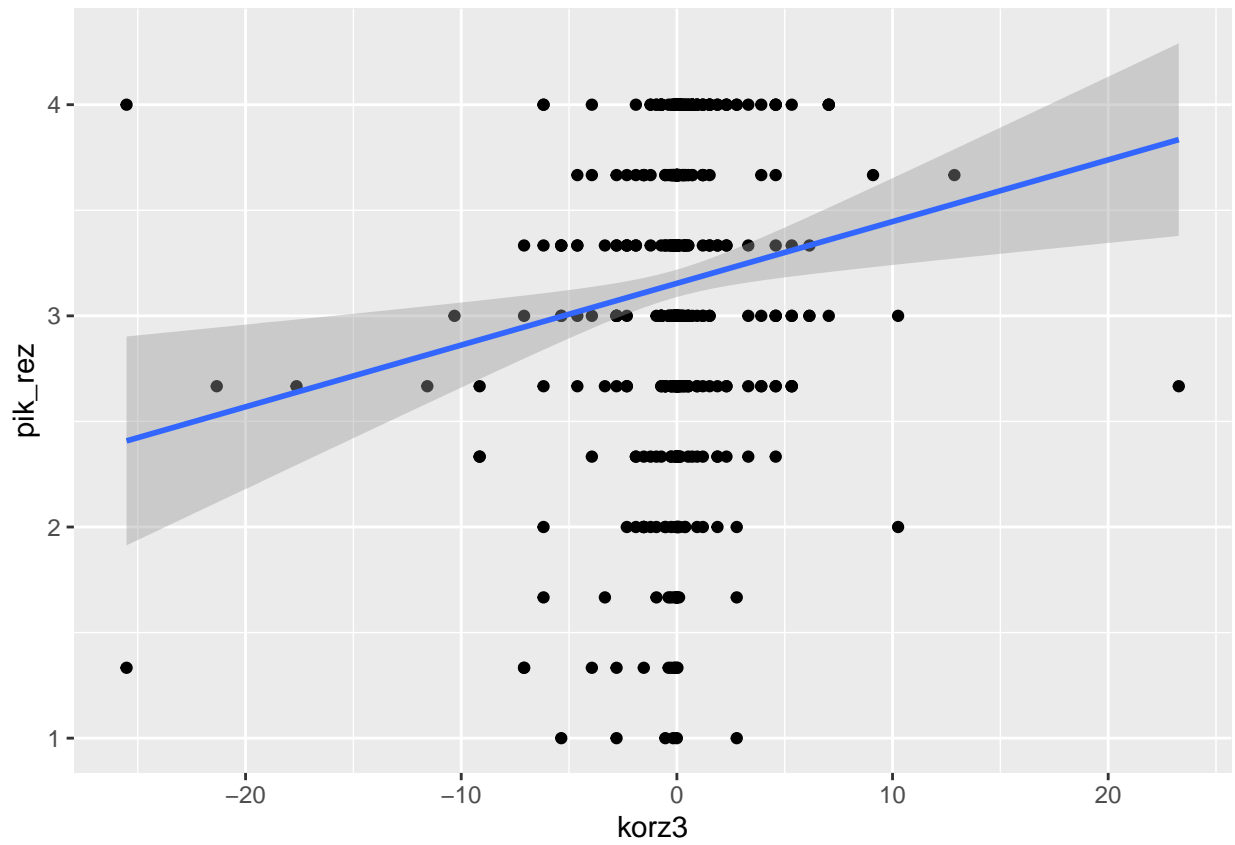
```
##
## [[4]]
```



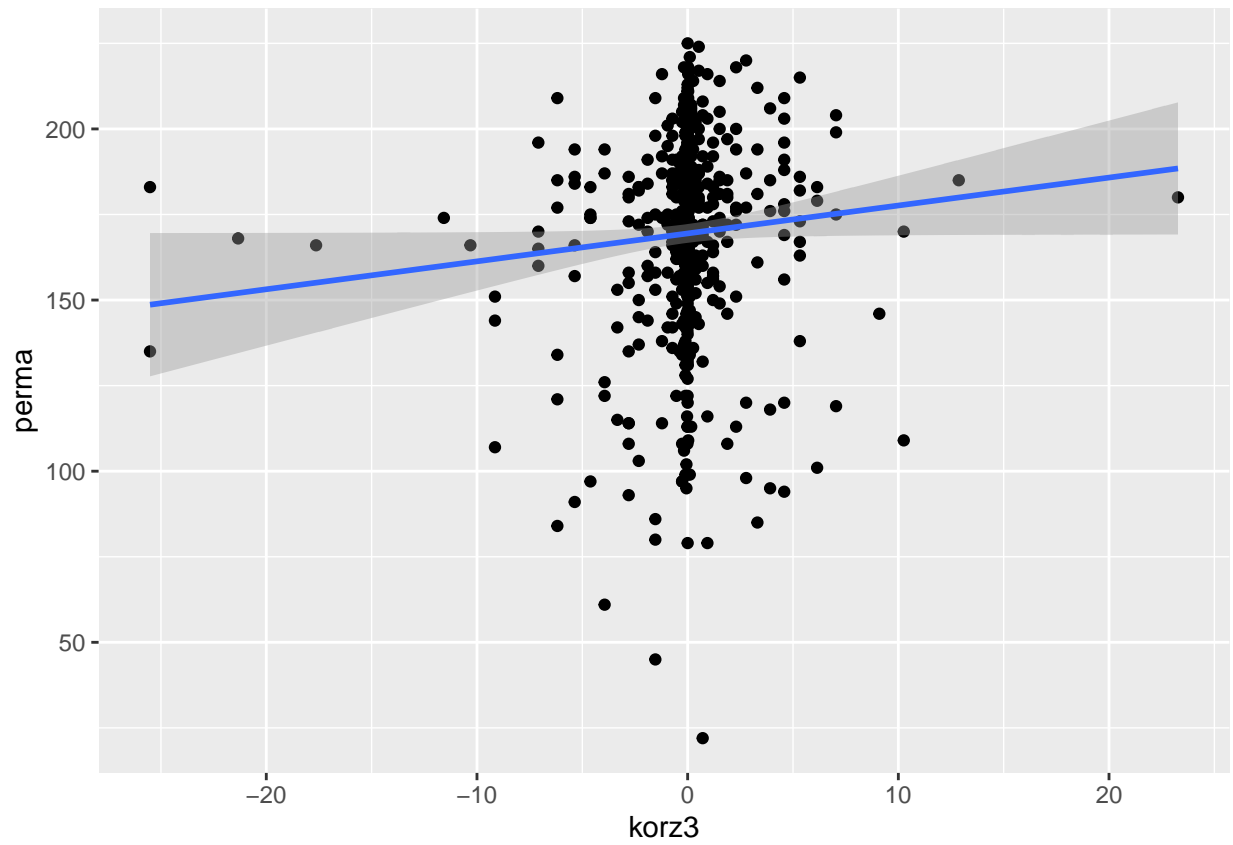
```
##
## [[5]]
```



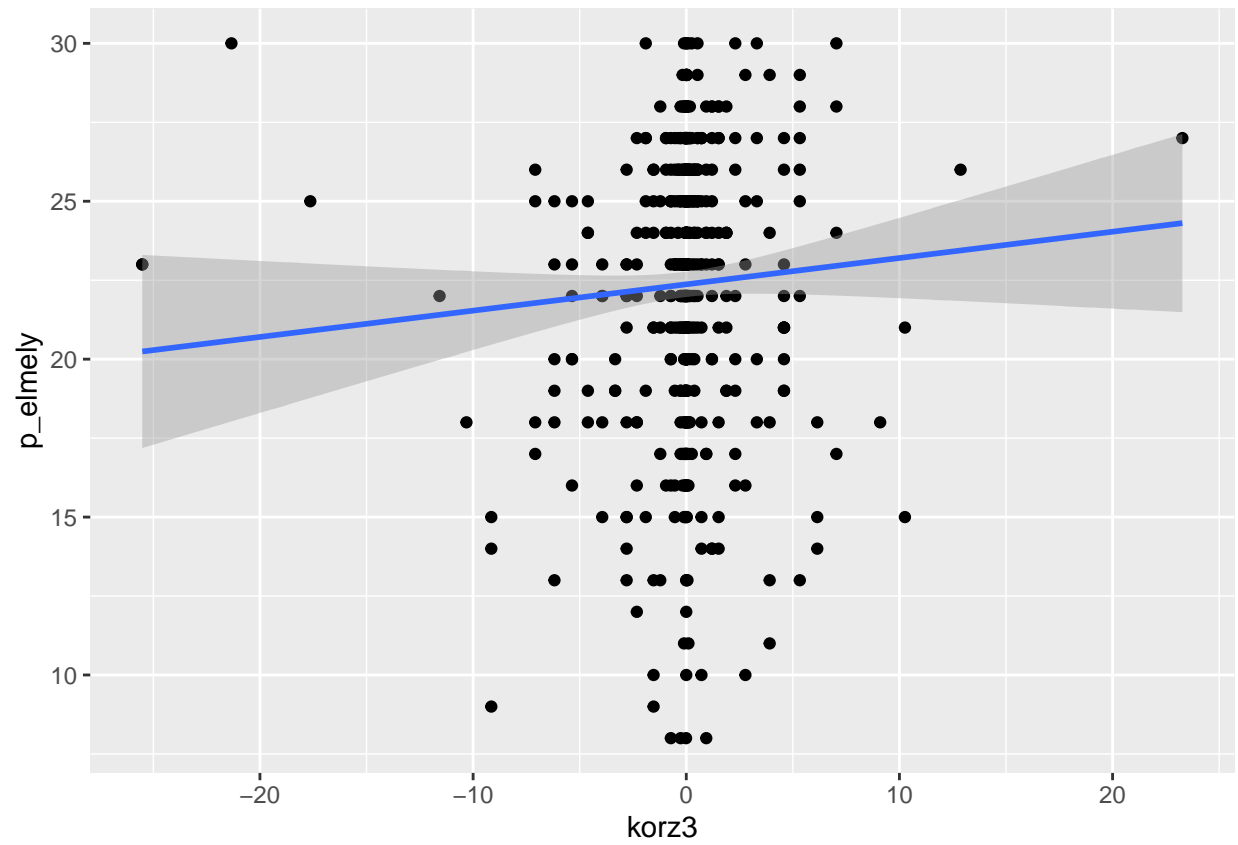
```
##  
## [[6]]
```



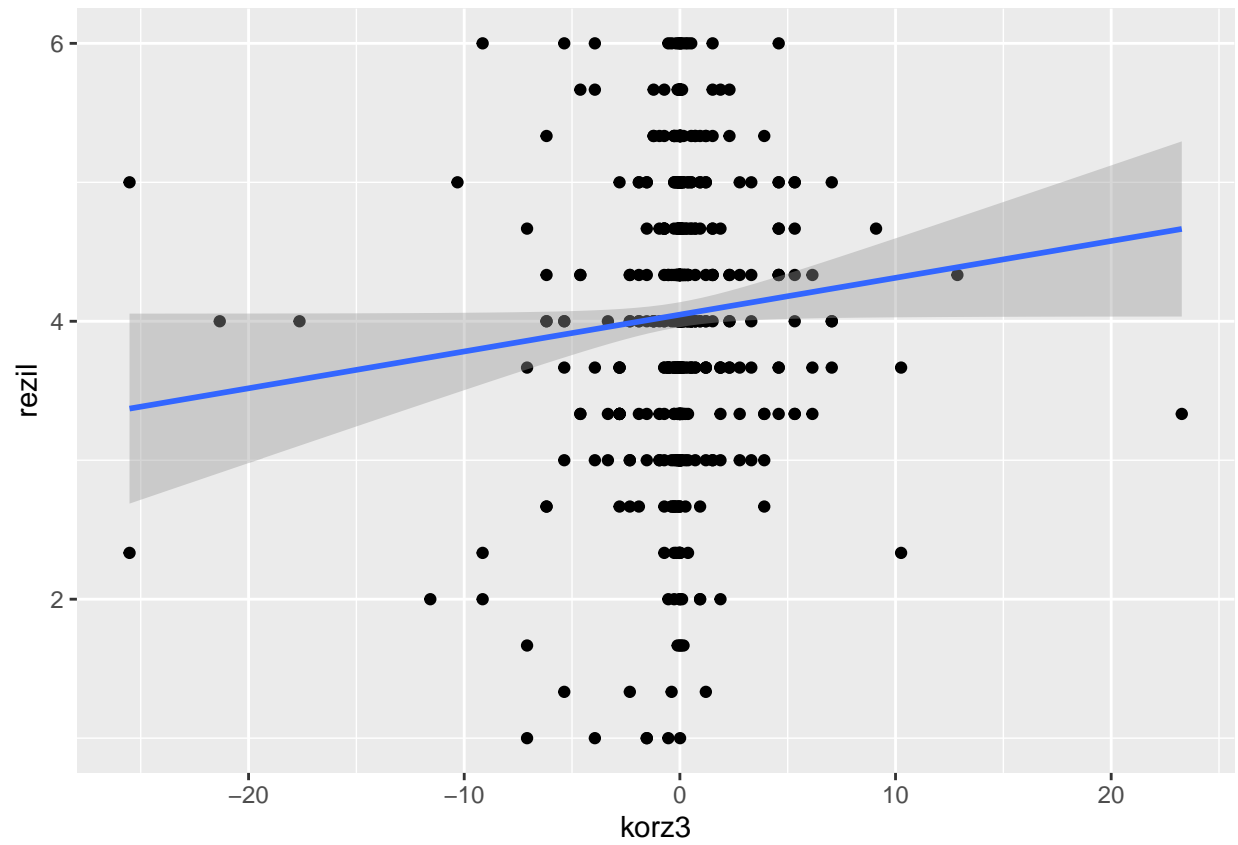
```
##
## [[7]]
```



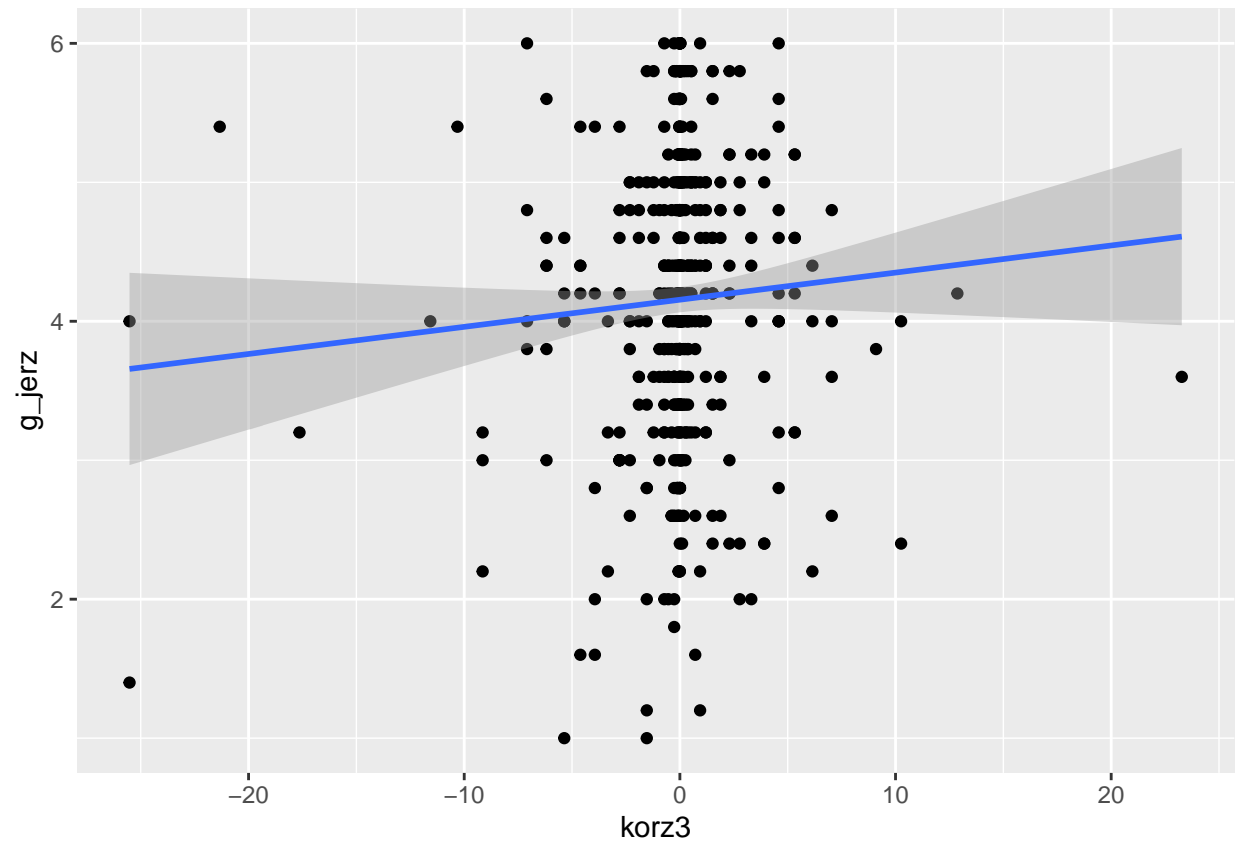
```
##  
## [[8]]
```



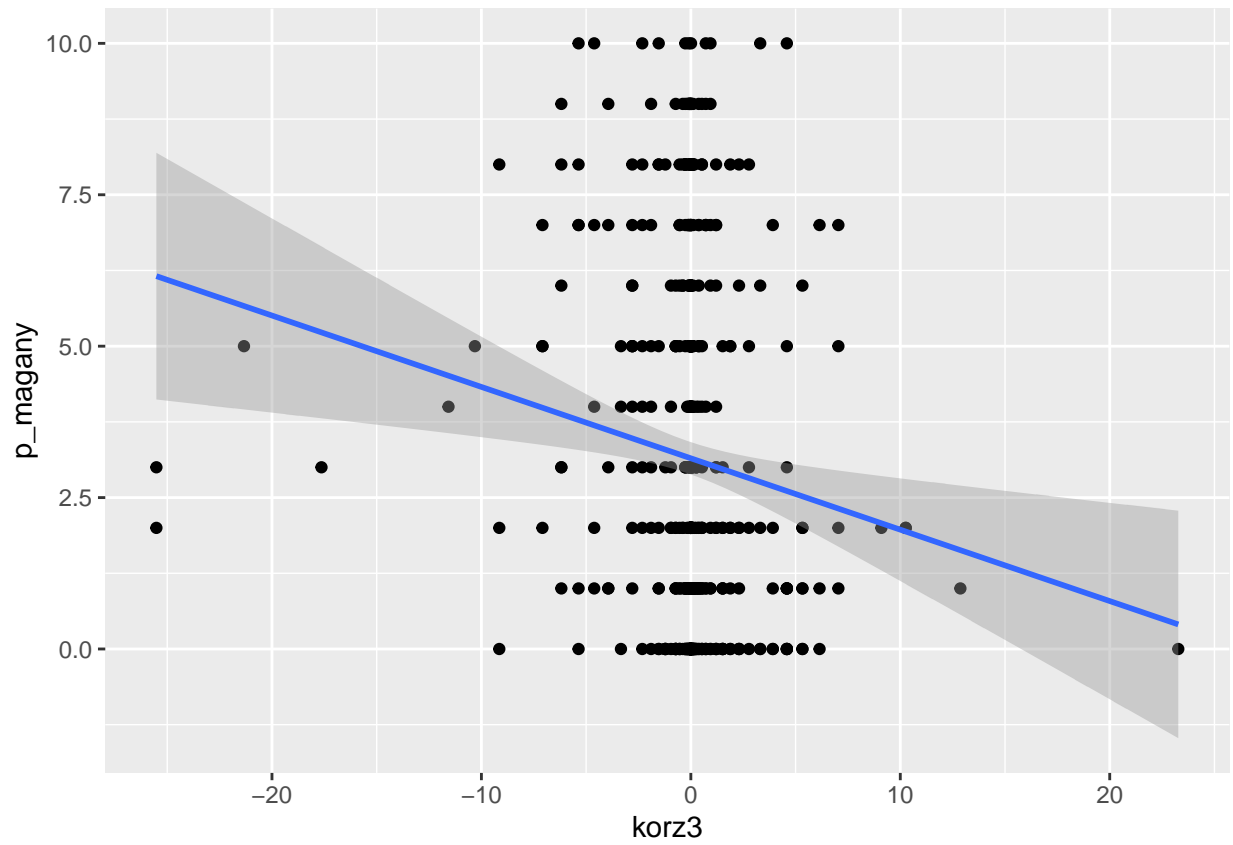
```
##
## [[9]]
```



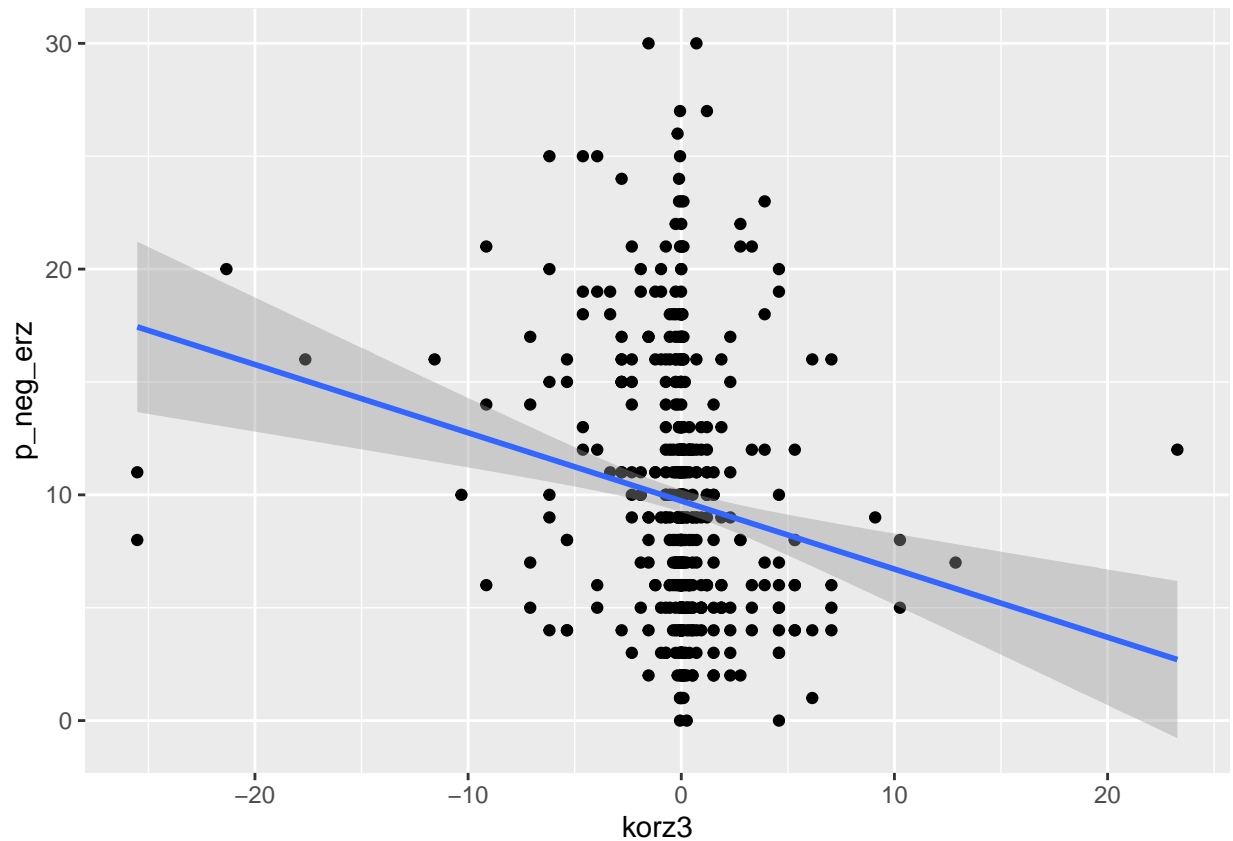
```
##  
## [[10]]
```

```
##  
## [[11]]
```



```
##
## [[12]]
```



6. Multiple linear regression where pelmeny is the outcome variable and Aggodalom, Ideges, Feszült, Nyugtalan are the predictor variables

```
m <- lm(p_elmeny_percent ~ aggodalo + ideges + feszult + nyugtala, data = processed)
summary(m)
```

```
##
## Call:
## lm(formula = p_elmeny_percent ~ aggodalo + ideges + feszult +
##      nyugtala, data = processed)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -61.397 -12.963   2.806  13.565  46.543
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   79.8519    1.8996  42.036  <2e-16 ***
## aggodalo      -0.3307    1.0691  -0.309   0.7572
## ideges        -3.6558    1.3685  -2.671   0.0078 **
```

```
## fészult      -3.0396      1.4024  -2.167   0.0307 *
## nyugtala     -1.4293      0.7555  -1.892   0.0591 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 18.52 on 495 degrees of freedom
## Multiple R-squared:  0.2585, Adjusted R-squared:  0.2525
## F-statistic: 43.15 on 4 and 495 DF,  p-value: < 2.2e-16
```

The rsquared suggest that the predictor variables explain 25.85% of the variance in the outcome variable. For the predictor variables the VIF is the highest for the *fészult* variable.

```
vif <- VIF(m)

vif
```

```
## aggodalo    ideges  fészult nyugtala
## 3.413217 5.149895 5.603635 1.765334
```

And the tolerances are:

```
1 / vif
```

```
## aggodalo    ideges  fészult nyugtala
## 0.2929787 0.1941787 0.1784556 0.5664649
```

Plotting the standardized residuals.

```
to_plot <- broom::augment(m)

ggplot(to_plot,
       aes(x = .fitted, y = .resid)) +
  geom_point() +
  papaja::theme_apapa()
```

