

# Solution for the assignment of the eight 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	1.750000	23715.85	1.00021	00002	0.000009	<U+2585>	<U+2587>	<U+2587>	
numeric	nyagi	0	1.000	NA	NA	NA	NA	NA	3.100000	56447.96	3.00031	00003	0.000005	<U+2581>	<U+2581>	<U+2581>	
numeric	pic elmeny	0	1.000	NA	NA	NA	NA	NA	58.460204	11835.01	40.00060	000080	0.000900	<U+2582>	<U+2585>	<U+2585>	
numeric	testi_fi	3	0.994	NA	NA	NA	NA	NA	4.164999	95722.09	4.00004	00005	0.000006	<U+2581>	<U+2583>	<U+2583>	
numeric	lelki	4	0.992	NA	NA	NA	NA	NA	4.161290	7223.62	4.00004	00005	0.000006	<U+2582>	<U+2583>	<U+2583>	
numeric	eg_all	0	0.986	NA	NA	NA	NA	NA	4.121704	11957.61	3.00004	00005	0.000006	<U+2582>	<U+2583>	<U+2583>	
numeric	fizero	6	0.988	NA	NA	NA	NA	NA	4.153846	18673.07	3.00004	00005	0.000006	<U+2582>	<U+2585>	<U+2585>	
numeric	cocska	0	1.000	NA	NA	NA	NA	NA	2.570000	14355.39	2.00021	00003	0.000007	<U+2587>	<U+2585>	<U+2585>	
numeric	godalo	0	1.000	NA	NA	NA	NA	NA	2.772000	43252.23	2.00021	50004	0.000006	<U+2587>	<U+2583>	<U+2583>	
numeric	ides	0	1.000	NA	NA	NA	NA	NA	2.490000	37465.61	1.00021	00003	0.000006	<U+2587>	<U+2582>	<U+2582>	
numeric	fuzult	0	1.000	NA	NA	NA	NA	NA	2.610000	39921.97	1.00021	00003	0.000006	<U+2587>	<U+2583>	<U+2583>	
numeric	ugtala	0	1.000	NA	NA	NA	NA	NA	2.406000	45788.03	1.00021	00003	0.000006	<U+2587>	<U+2582>	<U+2582>	

[illegible]

### 1. PCA on Diener flourishing scale

First, we select the variables of interest.

```
diener_data <- dplyr::select(processed, contains("diener"))
```

In R, we can run the PCA with just one function. We can set the `scale` argument to `true` to scale our variables to unit variance and center the variables.

```
pca <-  
  diener_data %>%  
  prcomp(scale = TRUE, center = TRUE)  
  
summary(pca)
```

```
## Importance of components:
##           PC1      PC2      PC3      PC4      PC5      PC6      PC7
## Standard deviation  2.3837 0.70037 0.67176 0.6203 0.55088 0.54801 0.46048
## Proportion of Variance 0.7102 0.06131 0.05641 0.0481 0.03793 0.03754 0.02651
## Cumulative Proportion 0.7102 0.77154 0.82795 0.8760 0.91399 0.95153 0.97803
##           PC8
## Standard deviation   0.41922
## Proportion of Variance 0.02197
## Cumulative Proportion 1.00000
```

As the data are standardized the standard deviation values for each principal component are the Eigenvalues.

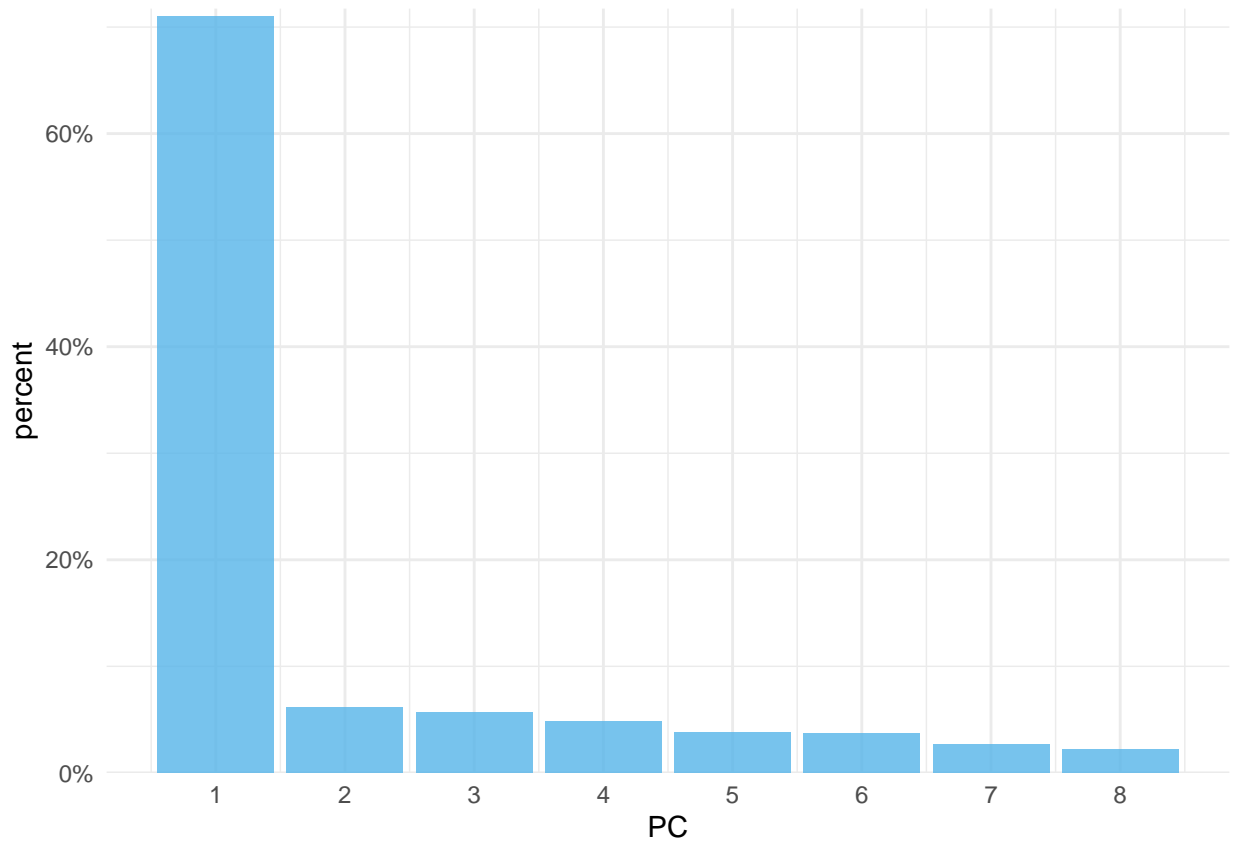
As we can see only the first PC has an Eigenvalue greater than 1.

The first PC accounts for 71.02% of the variance, while the second PC accounts for the 6.13% of the variance.

Together they explain 77.15% of the variance.

Lets plot the results.

```
pca %>%
  tidy(matrix = "eigenvalues") %>%
  ggplot(aes(PC, percent)) +
  geom_col(fill = "#56B4E9", alpha = 0.8) +
  scale_x_continuous(breaks = 1:8) +
  scale_y_continuous(
    labels = scales::percent_format(),
    expand = expansion(mult = c(0, 0.01))
  ) +
  theme_minimal()
```



## 2. PCA of MET

First, we select the variables of interest.

```
met_subscales <- c("jollet", "savor", "a_vhat", "onreg", "rezil")
met_data <- dplyr::select(processed, all_of(met_subscales))
```

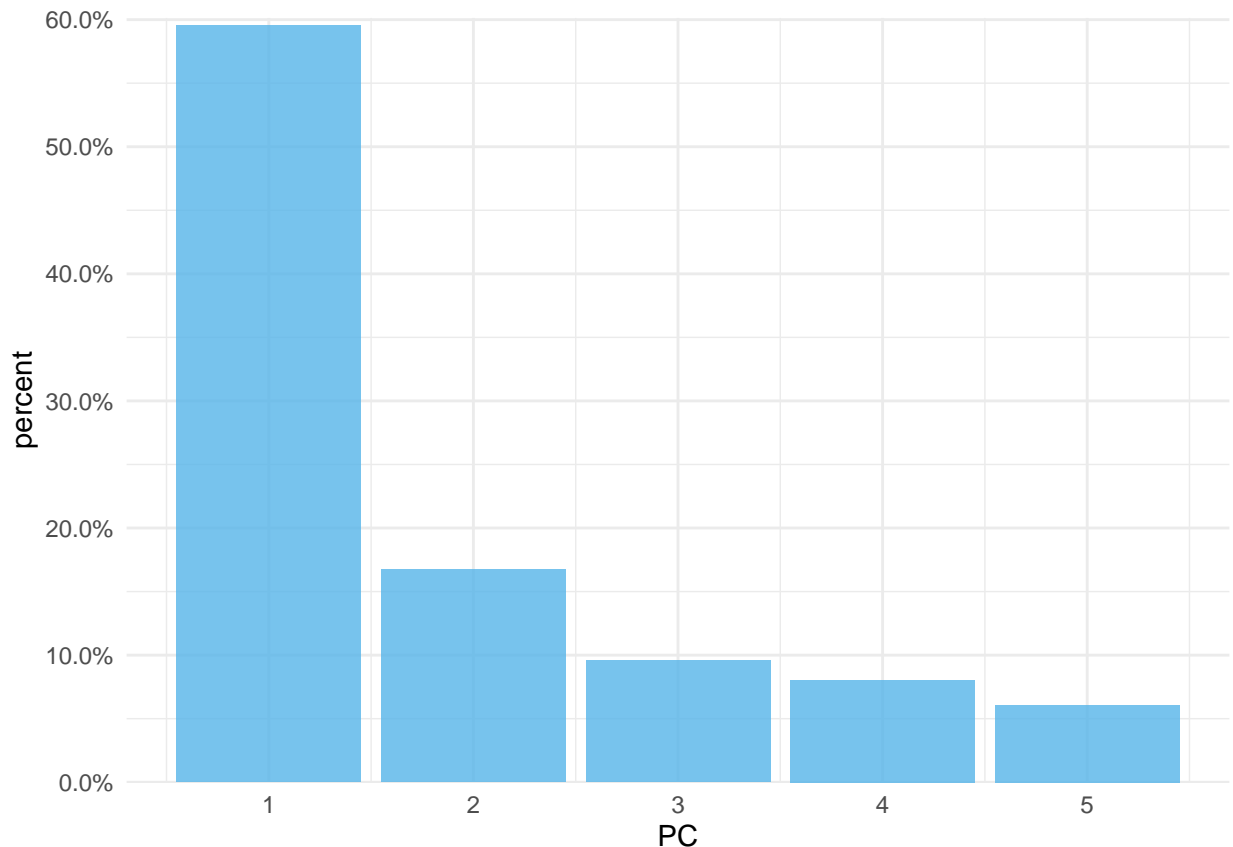
We run the PCA while standardizing the variables.

```
pca_met <-
  met_data %>%
  prcomp(scale = TRUE, center = TRUE)
summary(pca_met)
```

```
## Importance of components:
##              PC1      PC2      PC3      PC4      PC5
## Standard deviation  1.7254  0.9153  0.69228  0.63458  0.55085
## Proportion of Variance 0.5954  0.1676  0.09585  0.08054  0.06069
## Cumulative Proportion 0.5954  0.7629  0.85878  0.93931  1.00000
```

As we can see only the first PC has an Eigenvalue greater than 1, and it accounts for 59.54% of the variance.

```
pca_met %>%
  tidy(matrix = "eigenvalues") %>%
  ggplot(aes(PC, percent)) +
  geom_col(fill = "#56B4E9", alpha = 0.8) +
  scale_x_continuous(breaks = 1:5) +
  scale_y_continuous(
    labels = scales::percent_format(),
    expand = expansion(mult = c(0, 0.01))
  ) +
  theme_minimal()
```



### 3. Use varimax and promax rotation on the first three principal components of the second task.

I could use the `prcomp` function as before, but we would have to use a cutting point for the Eigenvalues to only include 3 PC. Thus, I will use the `psych::principal` function, where I can set the number of PC in the function call.

Thus, I run a new PCA with three components and varimax rotation. This function returns the standardized scores by default. The function calculates the correlation matrix by pairwise deletion of the missing values.

```
pca_met_varimax <- psych::principal(met_data, rotate = "varimax", nfactors = 3, scores = TRUE)
```

Lets take a look at the Eigenvalues.

```
pca_met_varimax$value
```

```
## [1] 2.9768986 0.8377353 0.4792455 0.4026888 0.3034318
```

The rest of the results.

```
pca_met_varimax
```

```
## Principal Components Analysis
## Call: psych::principal(r = met_data, nfactors = 3, rotate = "varimax",
##       scores = TRUE)
## Standardized loadings (pattern matrix) based upon correlation matrix
##          RC1  RC3  RC2  h2    u2 com
## jollet 0.67 0.55 0.14 0.77 0.230 2.0
## savor  0.91 0.10 0.12 0.86 0.144 1.1
## a_vhat 0.74 0.45 0.12 0.76 0.240 1.7
## onreg  0.14 0.22 0.97 1.00 0.002 1.1
## rezil  0.24 0.89 0.25 0.91 0.090 1.3
##
##
##          RC1  RC3  RC2
## SS loadings      1.90 1.36 1.04
## Proportion Var    0.38 0.27 0.21
## Cumulative Var    0.38 0.65 0.86
## Proportion Explained 0.44 0.32 0.24
## Cumulative Proportion 0.44 0.76 1.00
##
## Mean item complexity = 1.5
## Test of the hypothesis that 3 components are sufficient.
##
## The root mean square of the residuals (RMSR) is 0.08
## with the empirical chi square 57.46 with prob < NA
##
## Fit based upon off diagonal values = 0.98
```

The factor loadings are quite small for *onreg* and *rezil* on RC1. Based on the communalities (h2) it seems that most of the variance explained by these two items can be explained by other items.

The complexity scores suggest that *jollet* and maybe *a\_vhat* the items measure more than one latent construct.

The proportion var suggest that the RC1 explains 38% of the variance in the responses, while RC3 27% and RC2 21%.

The fit based upon off diagonal values support a good fit.

Now we should run the same analysis but allow factors to correlate, therefore we use an oblique rotation called promax.

```
pca_met_promax <- psych::principal(met_data, rotate = "promax", nfactors = 3, scores = TRUE)
```

Lets see the results of this solution.

```
pca_met_promax
```

```
## Principal Components Analysis
## Call: psych::principal(r = met_data, nfactors = 3, rotate = "promax",
##   scores = TRUE)
## Standardized loadings (pattern matrix) based upon correlation matrix
##          RC1   RC3   RC2   h2   u2 com
## jollet  0.54  0.43 -0.03 0.77 0.230 1.9
## savor   1.07 -0.27  0.04 0.86 0.144 1.1
## a_vhat  0.68  0.27 -0.03 0.76 0.240 1.3
## onreg   0.02  0.04  0.97 1.00 0.002 1.0
## rezil  -0.16  1.04  0.04 0.91 0.090 1.0
##
##
##          RC1   RC3   RC2
## SS loadings      1.92 1.40 0.98
## Proportion Var    0.38 0.28 0.20
## Cumulative Var    0.38 0.66 0.86
## Proportion Explained 0.45 0.33 0.23
## Cumulative Proportion 0.45 0.77 1.00
##
## With component correlations of
##          RC1   RC3   RC2
## RC1 1.00 0.66 0.29
## RC3 0.66 1.00 0.42
## RC2 0.29 0.42 1.00
##
## Mean item complexity = 1.3
## Test of the hypothesis that 3 components are sufficient.
##
## The root mean square of the residuals (RMSR) is 0.08
## with the empirical chi square 57.46 with prob < NA
##
## Fit based upon off diagonal values = 0.98
```

Allowing for correlation reduced the complexity and we have a much cleaner latent structure than before with varimax. RC1 highly correlates (0.66) with RC3 and only medium with RC2 (0.29), there is a medium to high correlation (0.42) between RC3 and RC2.

#### 4. Save the factors from the promax rotation and correlate the with the PERMA test's items. Which items has the highest correlation with the first factor.

Calculating the factor scores.

```
scores <- factor.scores(met_data, pca_met_promax)
scores <- as.data.frame(scores$scores)
```

Binding calculated scores to the original data.

```
perma_vars <- c("p_poz_erz", "p_elmely", "p_poz_kapc", "p_ert_cel", "p_telj", "p_boldog", "p_egeszs", "p_neg_erz", "p_magany")

processed_scores <-
  processed %>%
  bind_cols(., scores) %>%
  select(RC1, RC2, RC3, all_of(perma_vars))
```

Running the correlations.

```
corr_matrix <- rcorr(as.matrix(processed_scores))

corr_matrix$r
```

##		RC1	RC2	RC3	p_poz_erz	p_elmely	p_poz_kapc
##	RC1	1.0000000	0.2864299	0.6621899	0.6913069	0.4993597	0.5658908
##	RC2	0.2864299	1.0000000	0.4246157	0.4356578	0.2492797	0.2877081
##	RC3	0.6621899	0.4246157	1.0000000	0.7040195	0.4838113	0.5060462
##	p_poz_erz	0.6913069	0.4356578	0.7040195	1.0000000	0.5775692	0.6774156
##	p_elmely	0.4993597	0.2492797	0.4838113	0.5775692	1.0000000	0.4796039
##	p_poz_kapc	0.5658908	0.2877081	0.5060462	0.6774156	0.4796039	1.0000000
##	p_ert_cel	0.6604940	0.2450133	0.6050290	0.7531407	0.6305894	0.6495659
##	p_telj	0.6053890	0.2457362	0.5759880	0.7229889	0.5485736	0.5369546
##	p_boldog	0.6577439	0.3347994	0.6490360	0.8331380	0.5095568	0.7279601
##	p_egeszs	0.4361927	0.2566204	0.4511676	0.5274025	0.3379240	0.4356007
##	p_neg_erz	-0.5526917	-0.4805750	-0.6201278	-0.6866772	-0.4282924	-0.4972611
##	p_magany	-0.3836024	-0.2836810	-0.4253238	-0.4525517	-0.2901264	-0.4671761
##		p_ert_cel	p_telj	p_boldog	p_egeszs	p_neg_erz	p_magany
##	RC1	0.6604940	0.6053890	0.6577439	0.4361927	-0.5526917	-0.3836024
##	RC2	0.2450133	0.2457362	0.3347994	0.2566204	-0.4805750	-0.2836810
##	RC3	0.6050290	0.5759880	0.6490360	0.4511676	-0.6201278	-0.4253238
##	p_poz_erz	0.7531407	0.7229889	0.8331380	0.5274025	-0.6866772	-0.4525517
##	p_elmely	0.6305894	0.5485736	0.5095568	0.3379240	-0.4282924	-0.2901264
##	p_poz_kapc	0.6495659	0.5369546	0.7279601	0.4356007	-0.4972611	-0.4671761
##	p_ert_cel	1.0000000	0.7500587	0.7304123	0.4943210	-0.5041553	-0.3912630
##	p_telj	0.7500587	1.0000000	0.6500552	0.5273559	-0.5219503	-0.3331066
##	p_boldog	0.7304123	0.6500552	1.0000000	0.4527424	-0.6081244	-0.5236918
##	p_egeszs	0.4943210	0.5273559	0.4527424	1.0000000	-0.4147676	-0.2980372
##	p_neg_erz	-0.5041553	-0.5219503	-0.6081244	-0.4147676	1.0000000	0.5487226
##	p_magany	-0.3912630	-0.3331066	-0.5236918	-0.2980372	0.5487226	1.0000000

Interestingly there is a strong correlation between almost all of the variables. The strongest correlation is between RC1 and the PERMA items is with *p\_poz\_erz* ( $r = 0.69$ )