# 1 Import the Data File (teachfac.csv) from CANVAS

```
> library('GPArotation')
> library('psych')
```

# 3 Define data with suitable name

> df<-teachfac

**4** Inspect the current data frame

> head(df)

This is page 5 of 30. This reveals the data frame to consists of 28 variables:

| V1 | V2 | V3 | V4 | V5 | V6 | V7 | V8 | V9 | V10 | V11 | V12 | V13 | V14 | V15 | V16 | V17 | V18 | V19 | V20 | V21 | V22 | V23 | V24 | V25 | V26 | V27 | V28 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| sector | tsex | wen1 | wen2 | wen3 | wen4 | wen5 | wen6 | wen7 | wen8 | wen9 | wen10 | wen11 | wen12 | wen13 | wen14 | wen15 | wen16 | wen17 | wen18 | wen19 | wen20 | wen21 | strain | watts | energy | order | warmth |
| 3 | 1 | 5 | 5 | 4 | 5 | 4 | 5 | 4 | 5 | 5 | 5 | 4 | 5 | 5 | 5 | 4 | 5 | 5 | 5 | 5 | 5 | 5 | 1.271 | -1 | 6.371 | 6 | 6.384 |
| 1 | 1 | 4 | 3 | 3 | 2 | 2 | 2 | 3 | 2 | 3 | 4 | 2 | 3 | 3 | 4 | 2 | 4 | 2 | 3 | 3 | 4 | 4 | 1.664 | 2.763 | 5 | 3.296 | 5.472 |
| 2 | 1 | 0 | 0 | 0 | 0 | | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | -1 | -1 | -1 | -1 | -1 |
| 2 | 1 | 4 | 4 | 3 | 4 | 2 | 4 | 3 | 4 | 4 | 2 | 2 | 4 | 5 | 4 | 4 | 3 | 2 | 4 | 4 | 4 | 4 | 2.459 | 3.437 | 5.459 | 4.76 | 4.991 |
| 1 | 0 | 5 | 4 | 2 | 4 | 1 | 5 | 3 | 5 | 5 | 5 | 2 | 4 | 5 | 5 | 5 | 4 | 5 | 5 | 4 | 5 | 2.94 | 3.674 | 5.584 | 6 | 6.225 |

# 5 Redefine data frame for variables of interest only (3:23)

> df<-df[ ,3:23]

# Determine number of missing cases

> sum(is.na(df))

Here we can see that it's seven

**7** Check how many missing cases there are in Column 1

> sum(is.na(df[ ,1]))

Here we can see that it's zero

Repeat for all 21 columns to find columns with missing

```
> sum(is.na(df[ ,1]))
> sum(is.na(df[ ,2]))
> sum(is.na(df[ ,3]))
> sum(is.na(df[ ,4]))
> sum(is.na(df[ ,5]))
>           ...
> sum(is.na(df[ ,21]))
```

Here we can stop at column 5 because all seven were found in the same column

**Bonus task (for later): write a code that loops through all columns**

**9** Replace the missing values from column 5 with the median

```
>    temp<-median(df[,5], na.rm=T)
>    temp
>    temp1<-df[,5]
>    temp1[is.na(temp1)]
>    temp1[is.na(temp1)]<-temp
>    df[,5]<-temp1
```

Recheck the number of missing cases (should be zero now)

> sum(is.na(df))

Here we can see that it's zero after replacing with the median

# 11 Check the correlation matrix

> correlation<-round(cor(df),2)
> correlation

Do all variables have **at least one** correlation with another variable that is larger than .30?



```
wen21  0.40  0.37 -0.12  0.41 -0.33  0.48 -0.05  0.46  0.63  0.53 -0.15  0.63  0.50  0.56  0.62  0.55  0.00  0.67  0.65  0.52  1.00
> correlation<-round(cor(df),2)
> correlation
       wen1  wen2  wen3  wen4  wen5  wen6  wen7  wen8  wen9 wen10 wen11 wen12 wen13 wen14 wen15 wen16 wen17 wen18 wen19 wen20 wen21
wen1   1.00  0.44 -0.14  0.46 -0.39  0.51 -0.06  0.55  0.46  0.64 -0.09  0.42  0.32  0.65  0.43  0.59 -0.03  0.38  0.45  0.37  0.40
wen2   0.44  1.00  0.03  0.50 -0.16  0.43  0.01  0.41  0.48  0.40 -0.06  0.40  0.36  0.39  0.36  0.37  0.02  0.35  0.52  0.53  0.37
wen3  -0.14  0.03  1.00 -0.01  0.34 -0.14  0.56 -0.14 -0.16 -0.20  0.42 -0.12 -0.06 -0.20 -0.25 -0.21  0.47 -0.19 -0.13 -0.03 -0.12
wen4   0.46  0.50 -0.01  1.00 -0.20  0.42  0.01  0.48  0.48  0.42 -0.07  0.42  0.44  0.40  0.38  0.37  0.03  0.40  0.47  0.46  0.41
wen5  -0.39 -0.16  0.34 -0.20  1.00 -0.35  0.31 -0.31 -0.33 -0.48  0.26 -0.29 -0.18 -0.49 -0.35 -0.44  0.21 -0.32 -0.28 -0.19 -0.33
wen6   0.51  0.43 -0.14  0.42 -0.35  1.00 -0.12  0.48  0.53  0.59 -0.18  0.46  0.41  0.58  0.51  0.52 -0.08  0.43  0.53  0.47  0.48
wen7  -0.06  0.01  0.56  0.01  0.31 -0.12  1.00 -0.05 -0.08 -0.11  0.49 -0.06 -0.04 -0.13 -0.24 -0.14  0.52 -0.15 -0.09  0.02 -0.05
wen8   0.55  0.41 -0.14  0.48 -0.31  0.48 -0.05  1.00  0.52  0.55 -0.14  0.45  0.38  0.60  0.41  0.52 -0.02  0.43  0.51  0.46  0.46
wen9   0.46  0.48 -0.16  0.48 -0.33  0.53 -0.08  0.52  1.00  0.56 -0.20  0.61  0.49  0.61  0.66  0.54 -0.08  0.62  0.73  0.57  0.63
wen10  0.64  0.40 -0.20  0.42 -0.48  0.59 -0.11  0.55  0.56  1.00 -0.12  0.50  0.36  0.72  0.53  0.72 -0.04  0.48  0.52  0.47  0.53
wen11 -0.09 -0.06  0.42 -0.07  0.26 -0.18  0.49 -0.14 -0.20 -0.12  1.00 -0.12 -0.12 -0.19 -0.25 -0.15  0.47 -0.18 -0.22 -0.09 -0.15
wen12  0.42  0.40 -0.12  0.42 -0.29  0.46 -0.06  0.45  0.61  0.50 -0.12  1.00  0.48  0.54  0.61  0.50  0.00  0.67  0.59  0.48  0.63
wen13  0.32  0.36 -0.06  0.44 -0.18  0.41 -0.04  0.38  0.49  0.36 -0.12  0.48  1.00  0.43  0.43  0.38  0.01  0.49  0.57  0.53  0.50
wen14  0.65  0.39 -0.20  0.40 -0.49  0.58 -0.13  0.60  0.61  0.72 -0.19  0.54  0.43  1.00  0.58  0.68 -0.08  0.53  0.60  0.46  0.56
wen15  0.43  0.36 -0.25  0.38 -0.35  0.51 -0.24  0.41  0.66  0.53 -0.25  0.61  0.43  0.58  1.00  0.58 -0.16  0.66  0.62  0.43  0.62
wen16  0.59  0.37 -0.21  0.37 -0.44  0.52 -0.14  0.52  0.54  0.72 -0.15  0.50  0.38  0.68  0.58  1.00 -0.05  0.48  0.53  0.47  0.55
wen17 -0.03  0.02  0.47  0.03  0.21 -0.08  0.52 -0.02 -0.08 -0.04  0.47  0.00  0.01 -0.08 -0.16 -0.05  1.00 -0.06 -0.06  0.03  0.00
wen18  0.38  0.35 -0.19  0.40 -0.32  0.43 -0.15  0.43  0.62  0.48 -0.18  0.67  0.49  0.53  0.66  0.48 -0.06  1.00  0.67  0.52  0.67
wen19  0.45  0.52 -0.13  0.47 -0.28  0.53 -0.09  0.51  0.73  0.52 -0.22  0.59  0.57  0.60  0.62  0.53 -0.06  0.67  1.00  0.68  0.65
wen20  0.37  0.53 -0.03  0.46 -0.19  0.47  0.02  0.46  0.57  0.47 -0.09  0.48  0.53  0.46  0.43  0.47  0.03  0.52  0.68  1.00  0.52
wen21  0.40  0.37 -0.12  0.41 -0.33  0.48 -0.05  0.46  0.63  0.53 -0.15  0.63  0.50  0.56  0.62  0.55  0.00  0.67  0.65  0.52  1.00
>
```

# Check sampling accuracy with KMO

> kmo(df)

The overall MSA is larger than 0.60

# 13

## Compare correlation matrix to identity matrix

> bartlett.test(df)

Bartlett's test is significant so FA is suitable

> m1<-fa(df, nfactors = 21, rotate="none")



```
RStudio
File  Edit  Code  View  Plots  Session  Build  Debug  Profile  Tools  Help
                                                                                    Project: (None)
teachfac

Filter
      i..sector  tsex  wen1  wen2  wen3  wen4  wen5  wen6  wen7  wen8  wen9  wen10  wen11  wen12  wen13  wen14  wen15  wen16  wen17  w
1          3     1     5     5     4     5     4     5     4     5     5     5      4      5      5      5      4      5      5
           4     3     3     2     2     2     3     2     3     4     2      3      3      4      2      4      2
           0     0     0     0    NA     0     0     0     0     0     0      0      0      0      0      0      0      0
           4     4     3     4     2     4     3     4     4     2     2      4      5      4      4      3      2
```

```
Console  Terminal

                      MR1  MR2  MR3  MR4  MR5  MR6  MR7  MR8  MR9 MR10 MR11 MR12 MR13 MR14 MR15 MR16 MR17 MR18 MR19 MR20 MR21
SS loadings           8.70 2.17 1.05 0.67 0.31 0.26 0.23 0.19 0.15 0.13 0.10 0.10 0.09 0.07 0.06 0.05 0.02 0.01 0.00 0.00 0.00
Proportion Var        0.41 0.10 0.05 0.03 0.01 0.01 0.01 0.01 0.01 0.01 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00
Cumulative Var        0.41 0.52 0.57 0.60 0.61 0.63 0.64 0.65 0.65 0.66 0.66 0.67 0.67 0.68 0.68 0.68 0.68 0.68 0.68 0.68 0.68
Proportion Explained  0.61 0.15 0.07 0.05 0.02 0.02 0.02 0.01 0.01 0.01 0.01 0.01 0.01 0.01 0.00 0.00 0.00 0.00 0.00 0.00 0.00
Cumulative Proportion 0.61 0.76 0.83 0.88 0.90 0.92 0.93 0.95 0.96 0.96 0.97 0.98 0.98 0.99 0.99 1.00 1.00 1.00 1.00 1.00 1.00

Mean item complexity =  1.9
Test of the hypothesis that 21 factors are sufficient.

The degrees of freedom for the null model are   210  and the objective function was   12.26 with Chi Square of  13340.59
The degrees of freedom for the model are -21   and the objective function was  0

The root mean square of the residuals (RMSR) is  0
The df corrected root mean square of the residuals is  NA

The harmonic number of observations is  1097 with the empirical chi square  0  with prob <  NA
The total number of observations was  1097  with Likelihood Chi Square =  0  with prob <  NA

Tucker Lewis Index of factoring reliability =  1.016
Fit based upon off diagonal values = 1
Measures of factor score adequacy
                                             MR1  MR2  MR3  MR4  MR5   MR6   MR7   MR8   MR9  MR10  MR11  MR12  MR13  MR14  MR15
Correlation of (regression) scores with factors  0.99 0.92 0.90 0.83 0.74  0.69  0.70  0.61  0.59  0.54  0.48  0.47  0.50  0.42  0.46
Multiple R square of scores with factors          0.97 0.84 0.82 0.69 0.54  0.47  0.49  0.37  0.35  0.29  0.23  0.22  0.25  0.18  0.22
Minimum correlation of possible factor scores     0.94 0.69 0.63 0.37 0.08 -0.06 -0.02 -0.25 -0.29 -0.41 -0.54 -0.55 -0.50 -0.65 -0.57
                                             MR16  MR17  MR18  MR19  MR20 MR21
Correlation of (regression) scores with factors   0.40  0.28  0.18  0.11  0.03   0
Multiple R square of scores with factors          0.16  0.08  0.03  0.01  0.00   0
Minimum correlation of possible factor scores    -0.69 -0.84 -0.94 -0.97 -1.00  -1
>
```

> `fa.parallel(df, fa="fa")`

**16** Run a second model with **4 factors** (still un-rotated)

```
> m2<-fa(df, nfactors = 4, rotate="none")
> m2$loadings
```

This is not a simple structure (where each variable loads highly onto one and only one factor) so rotation is required

# Run a third model using a rotation method (e.g. oblimin)

> m3<-fa(df, nfactors = 4, rotate="oblimin")
> m3$loadings

Item 20 is still loading highly across multiple factors (we should remove this item)

**Remove item 20 and run a forth model**

> df<-df[, -20]

> m4<-fa(df, nfactors = 4, rotate="oblimin")

> m4$loadings

Here factor 4 only has two items loading highly onto it. We need to decide whether or not to remove it (for this example we will remove)



```
> df<-df[, -20]
> m4<-fa(df, nfactors = 4, rotate="oblimin")
> m4$loadings

Loadings:
        MR1     MR3     MR2     MR4
wen1   -0.140   0.804           0.202
wen2    0.206   0.178           0.489
wen3           -0.156   0.655   0.118
wen4    0.229   0.195           0.461
wen5           -0.555   0.266   0.137
wen6    0.202   0.448           0.183
wen7                    0.772
wen8    0.116   0.511           0.233
wen9    0.649   0.128           0.158
wen10           0.865
wen11  -0.119           0.645
wen12   0.739
wen13   0.567                   0.208
wen14   0.180   0.723
wen15   0.679   0.130  -0.146
wen16   0.168   0.712
wen17                   0.726
wen18   0.901
wen19   0.725                   0.230
wen21   0.773

                MR1     MR3     MR2     MR4
SS loadings     3.934   3.349   2.082   0.775
Proportion Var  0.197   0.167   0.104   0.039
Cumulative Var  0.197   0.364   0.468   0.507
>
```



Parallel Analysis Scree Plots

Remove item 20 and run a forth model

> m5<-fa(df, nfactors = 3, rotate="oblimin")

> m5$loadings

Much better:



RStudio screenshot showing console output:

```
Cumulative Var 0.197 0.364 0.468 0.507
> m5<-fa(df, nfactors = 3, rotate="oblimin")
> m5$loadings

Loadings:
        MR1    MR3    MR2
wen1           0.848
wen2    0.389  0.245  0.158
wen3          -0.138  0.690
wen4    0.402  0.258  0.150
wen5          -0.537  0.312
wen6    0.268  0.473
wen7                  0.781
wen8    0.200  0.546
wen9    0.726  0.123
wen10          0.861
wen11  -0.171         0.609
wen12   0.741
wen13   0.662
wen14   0.157  0.720
wen15   0.665        -0.177
wen16   0.122  0.690
wen17                 0.684
wen18   0.867 -0.102
wen19   0.833
wen21   0.747

                MR1   MR3   MR2
SS loadings    4.470 3.458 2.124
Proportion Var 0.224 0.173 0.106
Cumulative Var 0.224 0.396 0.503
>
```

Parallel Analysis Scree Plots

- FA Actual Data
- FA Simulated Data
- FA Resampled Data

# Export the residual matrix as a csv file

```
> temp<-round(m5$residual,2)
> write.csv(temp, file="residual.csv")
```

**21** Check which variables have high residuals

=COUNTIF(B2:U2,">.05")

Consider removing item 4

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U | | W | X |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | wen1 | wen2 | wen3 | wen4 | wen5 | wen6 | wen7 | wen8 | wen9 | wen10 | wen11 | wen12 | wen13 | wen14 | wen15 | wen16 | wen17 | wen18 | wen19 | wen21 | | | |
| 2 | wen1 | 0.39 | 0.06 | 0 | 0.06 | 0.02 | 0 | -0.01 | 0.03 | -0.01 | -0.03 | 0 | 0 | -0.01 | 0.01 | -0.01 | -0.02 | -0.03 | 0 | 0 | -0.03 | | 3 | 16% |
| 3 | wen2 | 0.06 | 0.66 | 0.03 | 0.14 | 0.05 | 0.05 | -0.03 | 0.03 | 0.02 | -0.03 | -0.03 | -0.02 | 0.01 | -0.05 | -0.04 | -0.04 | -0.05 | -0.06 | 0.06 | -0.06 | | 4 | 21% |
| 4 | wen3 | 0 | 0.03 | 0.49 | 0.01 | 0.01 | 0.02 | 0.01 | -0.03 | -0.01 | -0.01 | -0.02 | -0.01 | 0 | 0.02 | 0 | 0 | 0 | -0.02 | 0.01 | 0 | | 1 | 5% |
| 5 | wen4 | 0.06 | 0.14 | 0.01 | 0.63 | 0.03 | 0.02 | -0.02 | 0.08 | 0.01 | -0.03 | -0.03 | -0.03 | 0.07 | -0.06 | -0.04 | -0.06 | -0.03 | -0.04 | -0.01 | -0.04 | | 5 | 26% |
| 6 | wen5 | 0.02 | 0.05 | 0.01 | 0.03 | 0.62 | 0.01 | 0 | 0.03 | 0 | -0.02 | -0.01 | -0.02 | 0.01 | -0.03 | 0.01 | -0.01 | -0.03 | -0.03 | 0.02 | -0.04 | | 1 | 5% |
| 7 | wen6 | 0 | 0.05 | 0.02 | 0.02 | 0.01 | 0.51 | -0.01 | 0 | 0 | 0.01 | -0.02 | -0.02 | 0.03 | -0.01 | 0.01 | -0.02 | -0.02 | -0.04 | 0.01 | -0.01 | | 1 | 5% |
| 8 | wen7 | -0.01 | -0.03 | 0.01 | -0.02 | 0 | -0.01 | 0.39 | 0.01 | 0.03 | 0 | 0 | 0 | -0.02 | 0.02 | -0.01 | 0 | 0 | -0.01 | 0 | 0.02 | | 1 | 5% |
| 9 | wen8 | 0.03 | 0.03 | -0.03 | 0.08 | 0.03 | 0 | 0.01 | 0.52 | 0.01 | -0.03 | 0 | -0.01 | 0.01 | 0.02 | -0.05 | -0.02 | -0.01 | -0.01 | 0.01 | -0.02 | | 2 | 11% |
| 10 | wen9 | -0.01 | 0.02 | -0.01 | 0.01 | 0 | 0 | 0.03 | 0.01 | 0.33 | 0 | -0.01 | -0.01 | -0.03 | 0 | 0.03 | -0.01 | -0.02 | -0.03 | 0.04 | -0.01 | | 1 | 5% |
| 11 | wen10 | -0.03 | -0.03 | -0.01 | -0.03 | -0.02 | 0.01 | 0 | -0.03 | 0 | 0.26 | 0.03 | 0.01 | -0.02 | 0 | 0.01 | 0.05 | 0.02 | 0.02 | -0.01 | 0.02 | | 1 | 5% |
| 12 | wen11 | 0 | -0.03 | -0.02 | -0.03 | -0.01 | -0.02 | 0 | -0.03 | -0.01 | 0.03 | 0.59 | 0.03 | -0.02 | -0.01 | 0.02 | 0.02 | 0.05 | 0.04 | -0.03 | 0 | | 1 | 5% |
| 13 | wen12 | 0 | -0.02 | -0.01 | -0.03 | -0.02 | -0.02 | 0 | -0.01 | -0.01 | 0.01 | 0.03 | 0.41 | -0.01 | 0.01 | 0.03 | 0.01 | 0.02 | 0.06 | -0.06 | 0.03 | | 2 | 11% |
| 14 | wen13 | -0.01 | 0.01 | 0 | 0.07 | 0.01 | 0.03 | -0.02 | 0.01 | -0.03 | -0.02 | -0.02 | -0.01 | 0.59 | 0.01 | -0.05 | -0.01 | 0 | -0.02 | 0.03 | 0 | | 2 | 11% |
| 15 | wen14 | 0.01 | -0.05 | 0.02 | -0.06 | -0.03 | -0.01 | 0.02 | 0.02 | 0 | 0 | -0.01 | 0.01 | 0.01 | 0.28 | 0.01 | 0.01 | 0.01 | 0.01 | 0.02 | 0 | | 1 | 5% |
| 16 | wen15 | -0.01 | -0.04 | 0 | -0.04 | 0.01 | 0.01 | -0.01 | -0.05 | 0.03 | 0.01 | 0.02 | 0.03 | -0.05 | 0.01 | 0.38 | 0.05 | 0.01 | 0.03 | -0.03 | 0.02 | | 1 | 5% |
| 17 | wen16 | -0.02 | -0.04 | 0 | -0.06 | -0.01 | -0.02 | 0 | -0.02 | -0.01 | 0.05 | 0.02 | 0.02 | -0.01 | 0.01 | 0.05 | 0.38 | 0.03 | 0.01 | -0.01 | 0.04 | | 1 | 5% |
| 18 | wen17 | -0.03 | -0.05 | 0 | -0.06 | -0.03 | -0.02 | 0 | -0.01 | -0.02 | 0.02 | 0.05 | 0.02 | 0 | 0.01 | 0.01 | 0.03 | 0.54 | 0.03 | -0.01 | 0.03 | | 1 | 5% |
| 19 | wen18 | 0 | -0.06 | -0.02 | -0.04 | -0.03 | -0.04 | -0.01 | -0.01 | -0.03 | 0.02 | 0.04 | 0.06 | -0.02 | 0.01 | 0.03 | 0.01 | 0.03 | 0.34 | -0.01 | 0.04 | | 2 | 11% |
| 20 | wen19 | 0 | 0.06 | 0.01 | -0.01 | 0.02 | 0.01 | 0 | 0.01 | 0.04 | -0.01 | -0.03 | -0.06 | 0.03 | 0.02 | -0.03 | -0.01 | -0.01 | -0.01 | 0.28 | -0.02 | | 2 | 11% |
| 21 | wen21 | -0.03 | -0.06 | 0 | -0.04 | -0.04 | -0.01 | 0.02 | -0.02 | -0.01 | 0.02 | 0 | 0.03 | 0 | 0 | 0.02 | 0.04 | 0.03 | 0.04 | -0.02 | 0.38 | | 1 | 5% |

Remove item 4 and run a sixth model

> df<-df[ , -4]

> m6<-fa(df, nfactors = 3, rotate="oblimin")

> m6$loadings

Factor loadings are still good:



```
> df<-df[ , -4]
> m6<-fa(df, nfactors = 3, rotate="oblimin")
> m6$loadings

Loadings:
        MR1     MR3     MR2
wen1            0.844
wen2    0.372   0.246   0.143
wen3           -0.135   0.687
wen5           -0.537   0.306
wen6    0.263   0.476
wen7                    0.787
wen8    0.191   0.545
wen9    0.721   0.128
wen10           0.868
wen11  -0.165           0.615
wen12   0.743
wen13   0.648
wen14   0.156   0.729
wen15   0.668          -0.168
wen16   0.122   0.698
wen17                   0.692
wen18   0.873  -0.101
wen19   0.830
wen21   0.753


                    MR1     MR3     MR2
SS loadings         4.282   3.425   2.110
Proportion Var      0.225   0.180   0.111
Cumulative Var      0.225   0.406   0.517
>
```

Parallel Analysis Scree Plots

# 23 Export the residual matrix as a csv file for the new model

=COUNTIF(B2:U2,">.05")

Much better:

> temp<-round(m6$residual,2)

> write.csv(temp, file="residual2.csv")

Obtain a summary factor diagram to see which items load where

> fa.diagram(m6)

```
> factor1<-df[ , c(2,8,11,12,14,17,18,19)]
> factor2<-df[ , c(1,4,5,7,9,13,15)]
> factor3<-df[ , c(3, 6, 10, 16)]
```



RStudio

File  Edit  Code  View  Plots  Session  Build  Debug  Profile  Tools  Help

Go to file/function        Addins ▾            Project: (None) ▾

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|
| wen1 | wen2 | wen3 | wen5 | wen6 | wen7 | wen8 | wen9 | wen10 | wen11 | wen12 | wen13 | wen14 | wen15 | wen16 | wen17 | wen18 | wen19 | wen21 |
| 5 | 5 | 4 | 4 | 5 | 4 | 5 | 5 | 5 | 4 | 5 | 5 | 5 | 4 | 5 | 5 | 5 | 5 | 5 |
| 4 | 3 | 3 | 2 | 2 | 3 | 2 | 3 | 4 | 2 | 3 | 3 | 4 | 2 | 4 | 2 | 3 | 3 | 4 |
| 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 4 | 3 | 4 | 4 | 2 | 2 | 4 | 5 | 4 | 4 | 3 | 2 | 4 | 4 | 4 |
| 1 | 5 | 3 | 5 | 5 | 5 | 2 | 4 | 5 | 5 | 5 | 5 | 4 | 5 | 5 | 5 |
| 1 | 4 | 3 | 4 | 3 | 5 | 3 | 5 | 5 | 5 | 3 | 5 | 3 | 3 | 3 | 5 |

```
wen8    -0.0055910205  -0.0109811818   0.018977957  -0.0170791128
wen9    -0.0253105364  -0.0359215569   0.037087635  -0.0177528029
wen10    0.0133567073   0.0157788116  -0.011662946   0.0198706776
wen11    0.0441446138   0.0343384712  -0.033213983  -0.0004973291
wen12    0.0178133094   0.0560385897  -0.059669048   0.0188425009
wen13    0.0034911709  -0.0153449022   0.035410966   0.0009042886
wen14    0.0007226721   0.0040101430   0.011426186  -0.0034468308
wen15    0.0047501167   0.0279620798  -0.032026475   0.0118456443
wen16    0.0237185028  -0.0002357627  -0.015609289   0.0293600091
wen17    0.5346715355   0.0261360468  -0.016940365   0.0229904098
wen18    0.0261360468   0.3345794835  -0.012970528   0.0359288395
wen19   -0.0169403652  -0.0129705285   0.282819500  -0.0202929218
wen21    0.0229904098   0.0359288395  -0.020292922   0.3721188611

$dof
[1] 117

$chi
[1] 220.4391

$nh
[1] 1097

$rms
[1] 0.02423974

$EPVAL
[1] 2.430039e-08

> View(df)
> |
```



Factor Analysis

# Assess the reliability for factor 1

> Alpha(factor1)

---

RStudio

File  Edit  Code  View  Plots  Session  Build  Debug  Profile  Tools  Help

Addins

Project: (None)

teachfac  df

Filter

| | wen1 | wen2 | wen3 | wen5 | wen6 | wen7 | wen8 | wen9 | wen10 | wen11 | wen12 | wen13 | wen14 | we |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 5 | 5 | 4 | 4 | 5 | 4 | 5 | 5 | 5 | 4 | 5 | 5 | 5 | |
| 2 | 4 | 3 | 3 | 2 | 2 | 3 | 2 | 3 | 4 | 2 | 3 | 3 | 4 | |
| 3 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 4 | 4 | 4 | 3 | 2 | 4 | 3 | 4 | 4 | 2 | 2 | 4 | 5 | 4 | |
| 5 | 5 | 4 | 2 | 1 | 5 | 3 | 5 | 5 | 5 | 2 | 4 | 5 | 5 | |

Showing 1 to 6 of 1,097 entries

Environment  History  Connections

Import Dataset    List

Global Environment

| | | |
|---|---|---|
| fac1 | 1097 obs. of 8 variables | |
| factor1 | 1097 obs. of 8 variables | |
| factor2 | 1097 obs. of 7 variables | |
| factor3 | 1097 obs. of 4 variables | |
| m3 | List of 51 | |
| m4 | List of 51 | |
| m5 | List of 51 | |
| m6 | List of 51 | |
| teachfac | 1097 obs. of 28 variables | |
| temp | num [1:19, 1:19] 0.4 0.07 0 0... | |

Values

Console  Terminal

```
[1] 2.430039e-08

> alpha(factor1)

Reliability analysis
Call: alpha(x = factor1)

  raw_alpha std.alpha G6(smc) average_r S/N    ase mean  sd median_r
       0.91      0.91    0.91      0.55 9.8 0.0041  3.5 0.8      0.6

 lower alpha upper     95% confidence boundaries
0.9 0.91 0.92

 Reliability if an item is dropped:
       raw_alpha std.alpha G6(smc) average_r  S/N alpha se  var.r med.r
wen2        0.91      0.91    0.91      0.60 10.5   0.0041 0.0062  0.62
wen9        0.89      0.89    0.88      0.53  8.0   0.0051 0.0135  0.57
wen12       0.89      0.89    0.89      0.54  8.4   0.0048 0.0144  0.57
wen13       0.91      0.91    0.90      0.58  9.6   0.0043 0.0131  0.62
wen15       0.89      0.89    0.89      0.55  8.5   0.0048 0.0129  0.57
wen18       0.89      0.89    0.88      0.54  8.2   0.0049 0.0119  0.57
wen19       0.89      0.89    0.88      0.53  7.8   0.0051 0.0135  0.50
wen21       0.89      0.89    0.89      0.54  8.3   0.0049 0.0134  0.57

 Item statistics
         n raw.r std.r r.cor r.drop mean   sd
wen2  1097  0.61  0.62  0.53   0.50  4.0 0.93
wen9  1097  0.85  0.84  0.82   0.78  3.4 1.18
wen12 1097  0.80  0.80  0.77   0.73  3.5 1.02
wen13 1097  0.69  0.69  0.62   0.59  3.4 1.05
wen15 1097  0.80  0.79  0.76   0.72  3.0 1.16
```
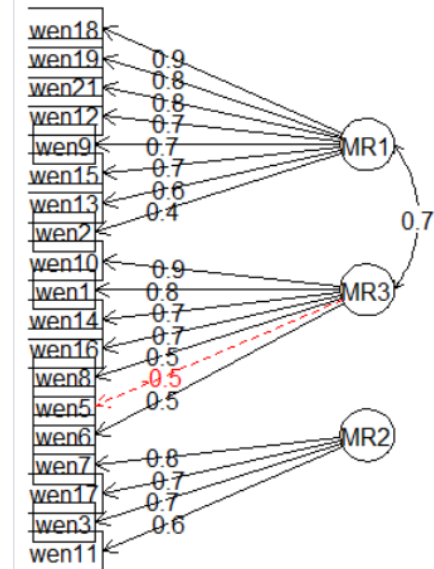
Files  Plots  Packages  Help  Viewer

Zoom  Export

**Factor Analysis**

wen18, wen19, wen21, wen12, wen9, wen15, wen13, wen2, wen10, wen1, wen14, wen16, wen8, wen5, wen6, wen7, wen17, wen3, wen11 → MR1, MR3, MR2 (0.9, 0.8, 0.8, 0.7, 0.7, 0.7, 0.6, 0.4, 0.9, 0.8, 0.7, 0.7, 0.5, -0.5, 0.5, 0.8, 0.7, 0.7, 0.6, 0.7)