

# Exercise 7

Tobias Raidl, 11717659

2023-12-12

## Contents

1	1
2	3
3	5
4	29
5	30
6	30

```
cars = read.csv("cardata.csv")
cars = na.omit(cars)
df = cars[, -c(1:9,15,16,18)]
```

## 1

Use the function `pfa()` from the package `StatDA` for principal factor analysis, with the argument `scores="regression"` to also obtain scores. Inspect the biplot and try to interpret the factors.

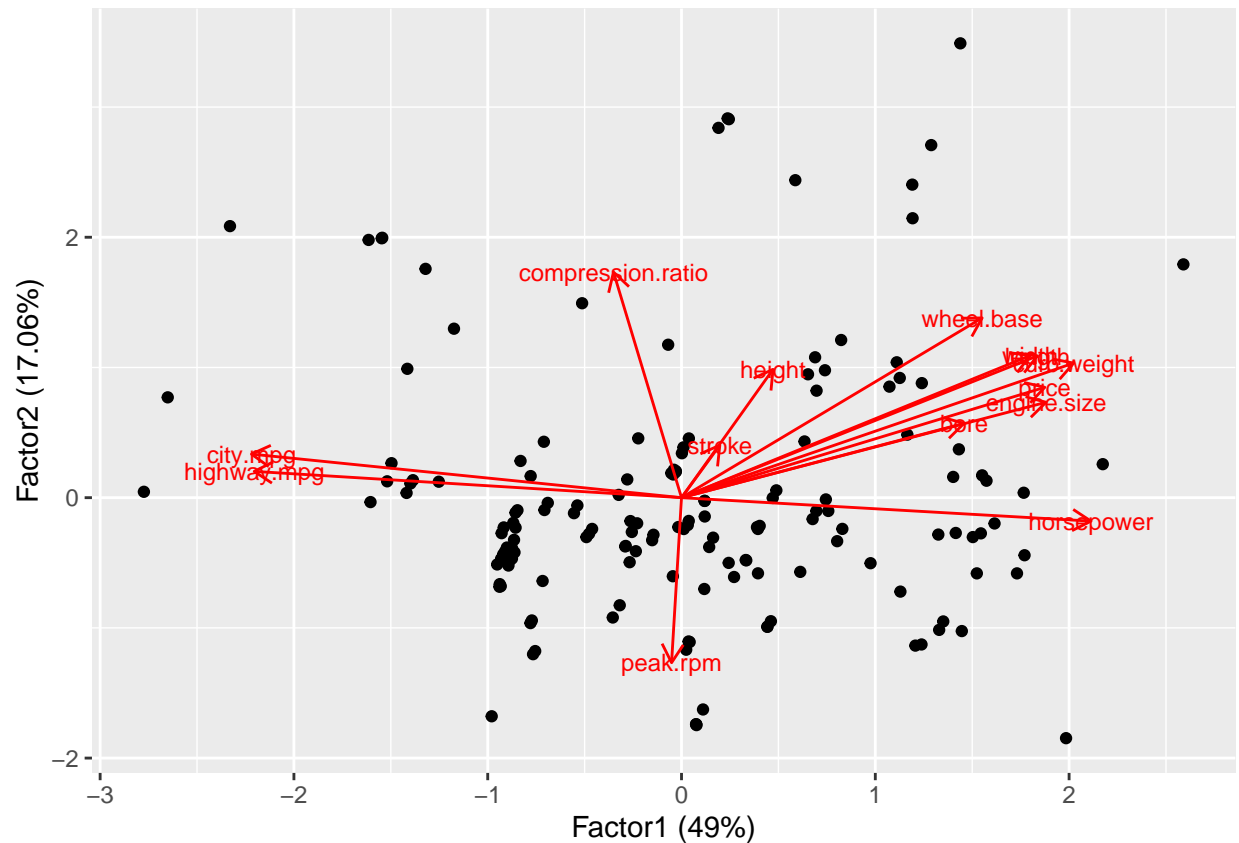
```
library(StatDA)
```

```
## Warning: package 'StatDA' was built under R version 4.3.2
## Loading required package: sgeostat
## Registered S3 method overwritten by 'geoR':
##   method      from
##   plot.variogram sgeostat
```

```
library(ggfortify)
```

```
## Warning: package 'ggfortify' was built under R version 4.3.2
## Loading required package: ggplot2
## Warning: package 'ggplot2' was built under R version 4.3.2
```

```
df_scaled = scale(df)
pfa = pfa(df_scaled, factors=2, scores="regression")
autoplot(pfa, data=df_scaled, loadings = TRUE, loadings.colour = 'red',
         loadings.label = TRUE, loadings.label.size = 3)
```



```
pfa
```

```
##
## Call:
## pfa(x = df_scaled, factors = 2, scores = "regression")
##
## Uniquenesses:
##      wheel.base      length      width      height
##      0.199         0.153      0.177      0.778
##      curb.weight    engine.size    bore      stroke
##      0.041         0.243      0.544      0.963
##      compression.ratio horsepower    peak.rpm    city.mpg
##      0.420         0.165      0.701      0.064
##      highway.mpg    price
##      0.089         0.214
##
## Loadings:
##      Factor1 Factor2
## wheel.base  0.669  0.595
## length      0.792  0.469
## width       0.776  0.470
## height      0.203  0.425
## curb.weight 0.872  0.445
## engine.size 0.811  0.316
## bore       0.628  0.246
## stroke      0.174
```

```
## compression.ratio -0.151  0.747
## horsepower        0.911
## peak.rpm          -0.546
## city.mpg          -0.956  0.144
## highway.mpg       -0.950
## price             0.808  0.365
##
##               Factor1 Factor2
## SS loadings    6.859  2.388
## Proportion Var 0.490  0.171
## Cumulative Var 0.490  0.661
##
## Test of the hypothesis that 2 factors are sufficient.
## The chi square statistic is NA on 64 degrees of freedom.
## The p-value is NA
```

The first factor covers variables like horsepower, size, weight, highway.mpg and city.mpg (horizontal), and the second factor covers variables such as peak.rpm, height and compression rate. Wheel base for example is similarly covered by both factors.

## 2

Estimate robustly mean and covariance using the MCD estimator. Do the same as above, but provide the results from the MCD to the covmat argument to obtain a robust factor analysis solution. Note that scaling should then also be done robustly. Does the interpretation of the factors change?

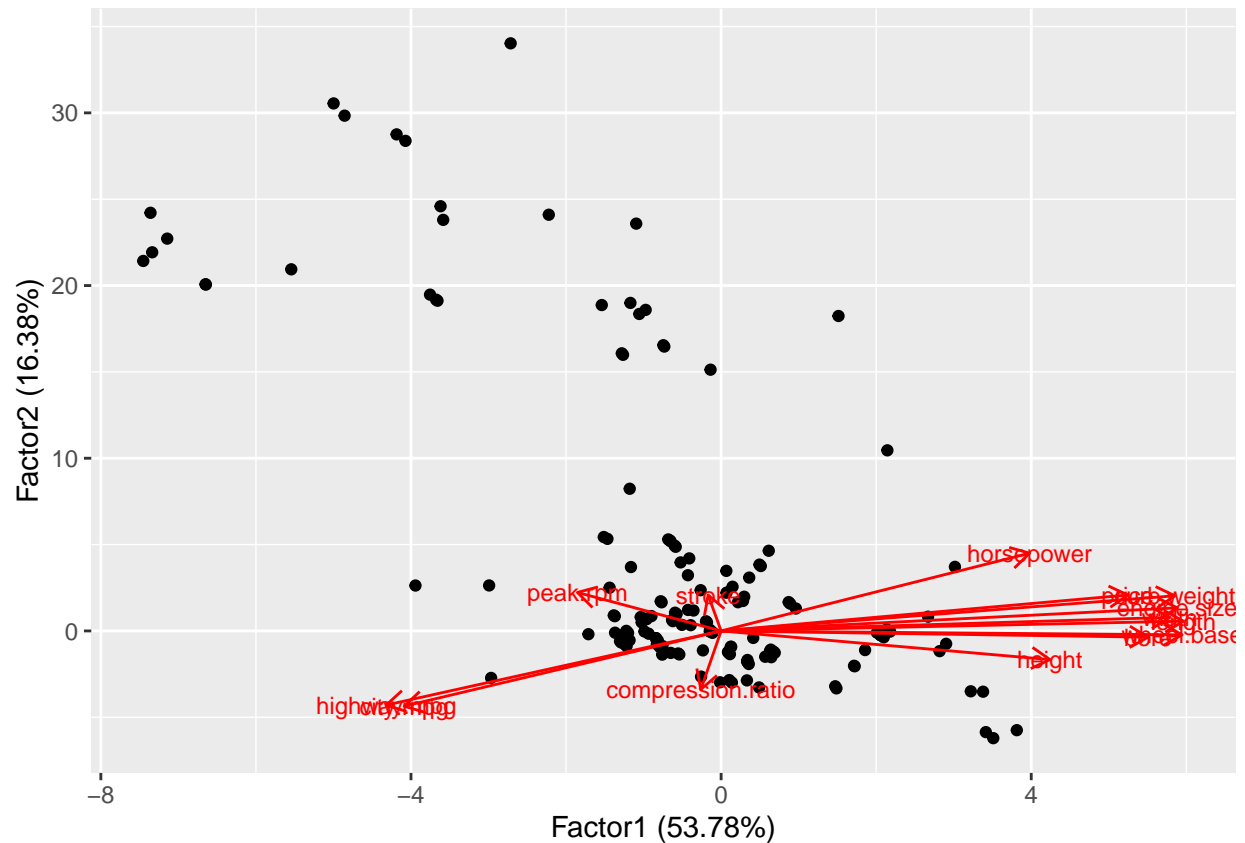
```
library(robustbase)
```

```
## Warning: package 'robustbase' was built under R version 4.3.2
```

```
library(DescTools)
```

```
## Warning: package 'DescTools' was built under R version 4.3.2
```

```
covmat = covMcd(df)$cov
df_robtscaled = RobScale(df)
set.seed(69)
robust_pfa = pfa(df_robtscaled, factors=2, covmat=covmat, scores="regression")
autoplot(robust_pfa, data=df_robtscaled, loadings = TRUE, loadings.colour = 'red',
         loadings.label = TRUE, loadings.label.size = 3)
```



```
robust_pfa
```

```
##
## Call:
## pfa(x = df_rob_scaled, factors = 2, covmat = covmat, scores = "regression")
##
## Uniquenesses:
##      wheel.base      length      width      height
##      0.123          0.111      0.154      0.486
##      curb.weight    engine.size    bore      stroke
##      0.049          0.098      0.245      0.892
##      compression.ratio horsepower    peak.rpm    city.mpg
##      0.714          0.106      0.792      0.115
##      highway.mpg    price
##      0.081          0.211
##
## Loadings:
##      Factor1 Factor2
## wheel.base  0.936
## length      0.939
## width        0.911  0.124
## height       0.668 -0.261
## curb.weight  0.921  0.320
## engine.size  0.926  0.211
## bore         0.867
## stroke       0.327
```

```
## compression.ratio      -0.533
## horsepower             0.626   0.708
## peak.rpm               -0.292   0.351
## city.mpg               -0.644  -0.686
## highway.mpg            -0.680  -0.676
## price                  0.826   0.327
##
##               Factor1 Factor2
## SS loadings      7.530   2.293
## Proportion Var   0.538   0.164
## Cumulative Var   0.538   0.702
##
## The degrees of freedom for the model is 64 and the fit was NA
```

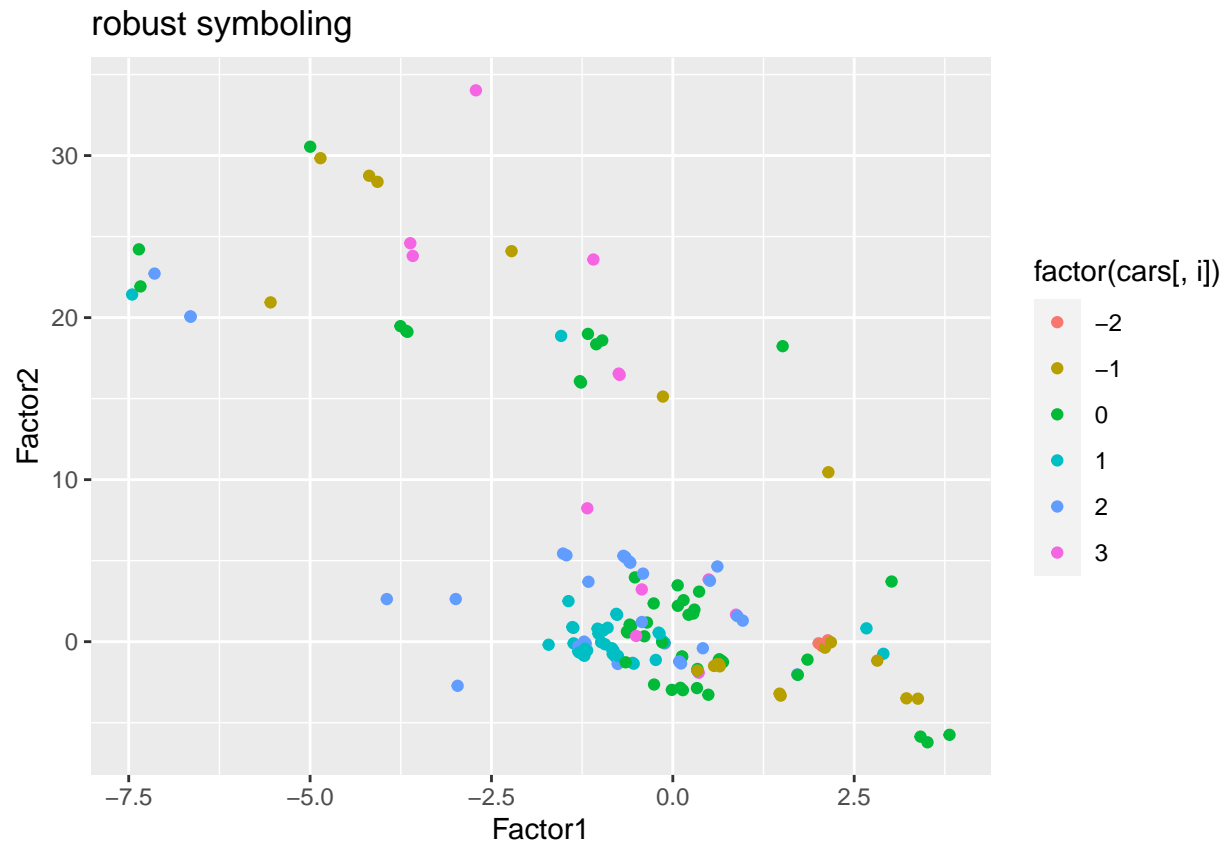
The explained variance by the first 2 factors is now ~5% higher than in the non-robust approach.

### 3

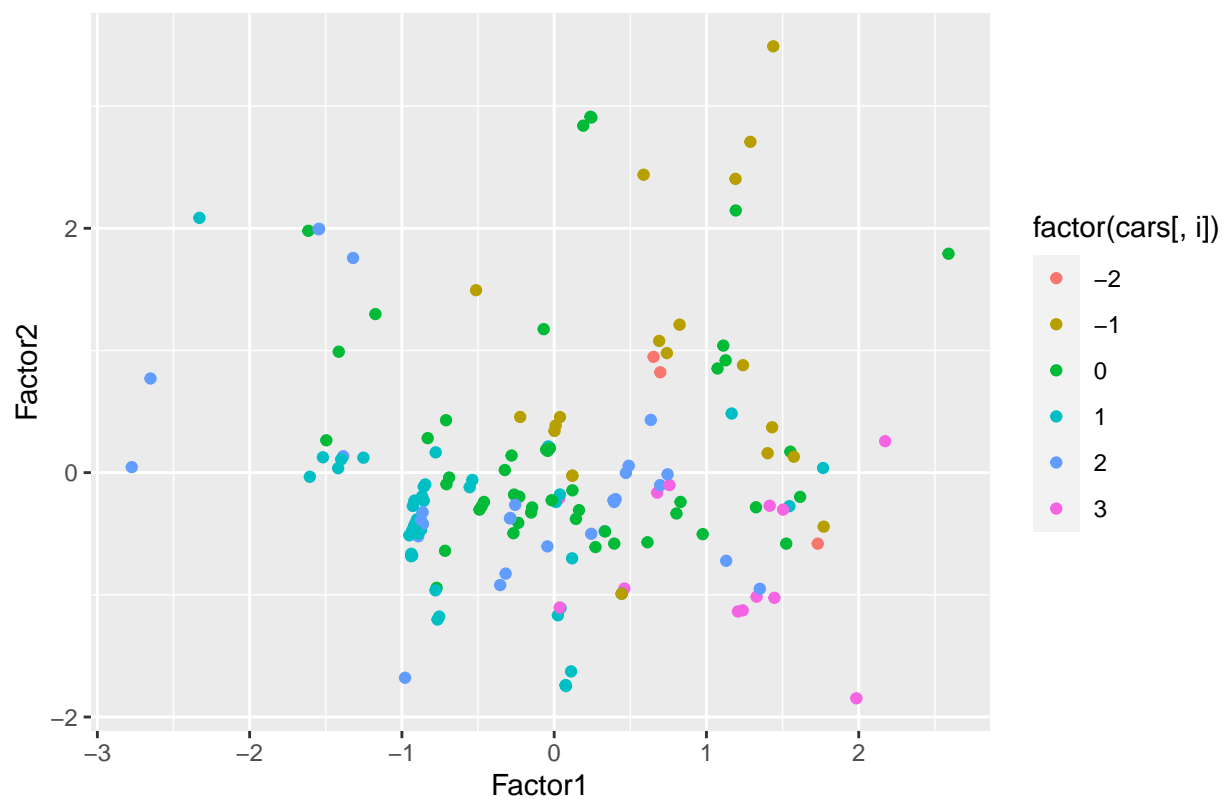
Look at the scores from robust factor analysis and try to identify a variable from the data set which is explaining outliers in the scores. So, essentially, plot the scores, with color according to other (categorical) variables). Do the same for the non-robust scores. Does the same variable also lead to an explanation of the outlyingness?

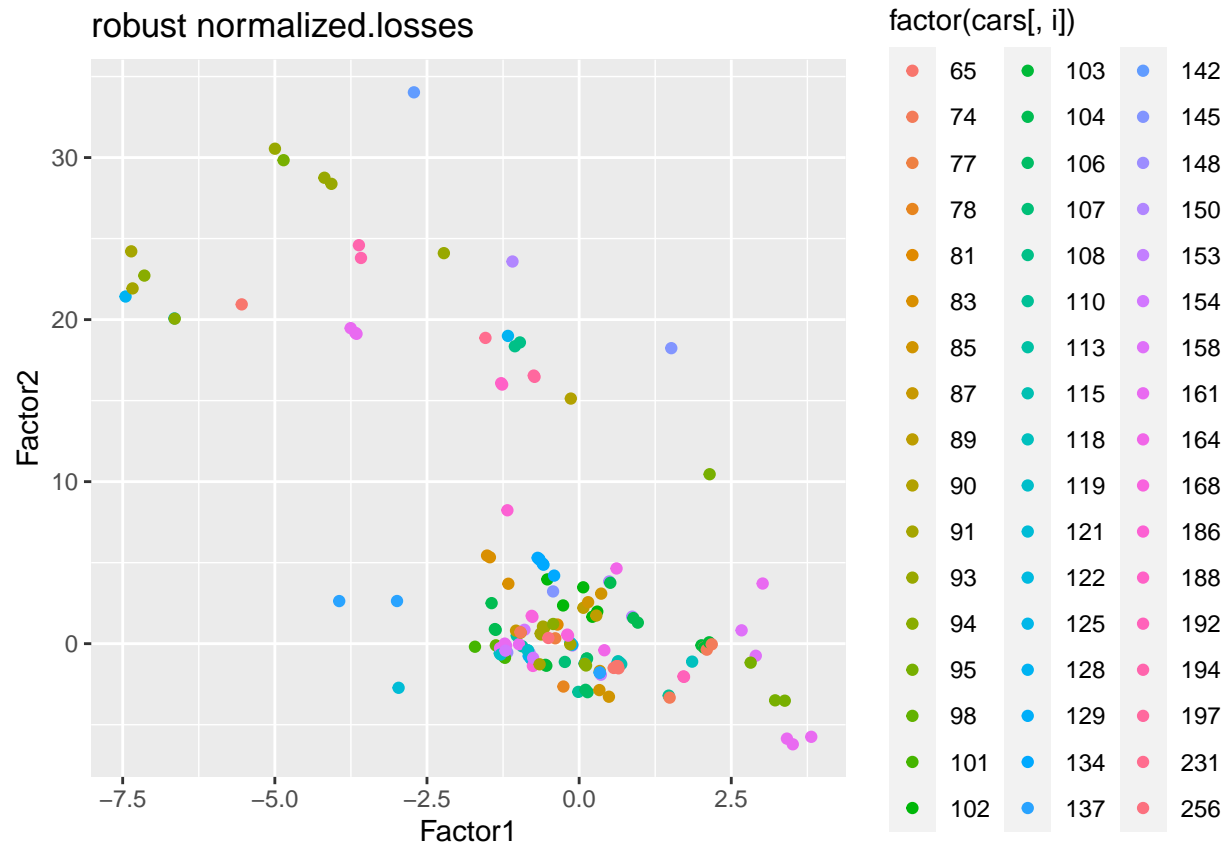
```
library(ggplot2)

facvars = c(1:9,15,16,18)
for(i in facvars) {
  print(
    ggplot(robust_pfa$scores, aes(x=Factor1, y=Factor2, col=factor(cars[,i]))) +
    geom_point() +
    ggtitle(paste("robust", names(cars)[i]))
  )
  print(
    ggplot(pfa$scores, aes(x=Factor1, y=Factor2, col=factor(cars[,i]))) +
    geom_point() +
    ggtitle(paste("non-robust", names(cars)[i]))
  )
}
```

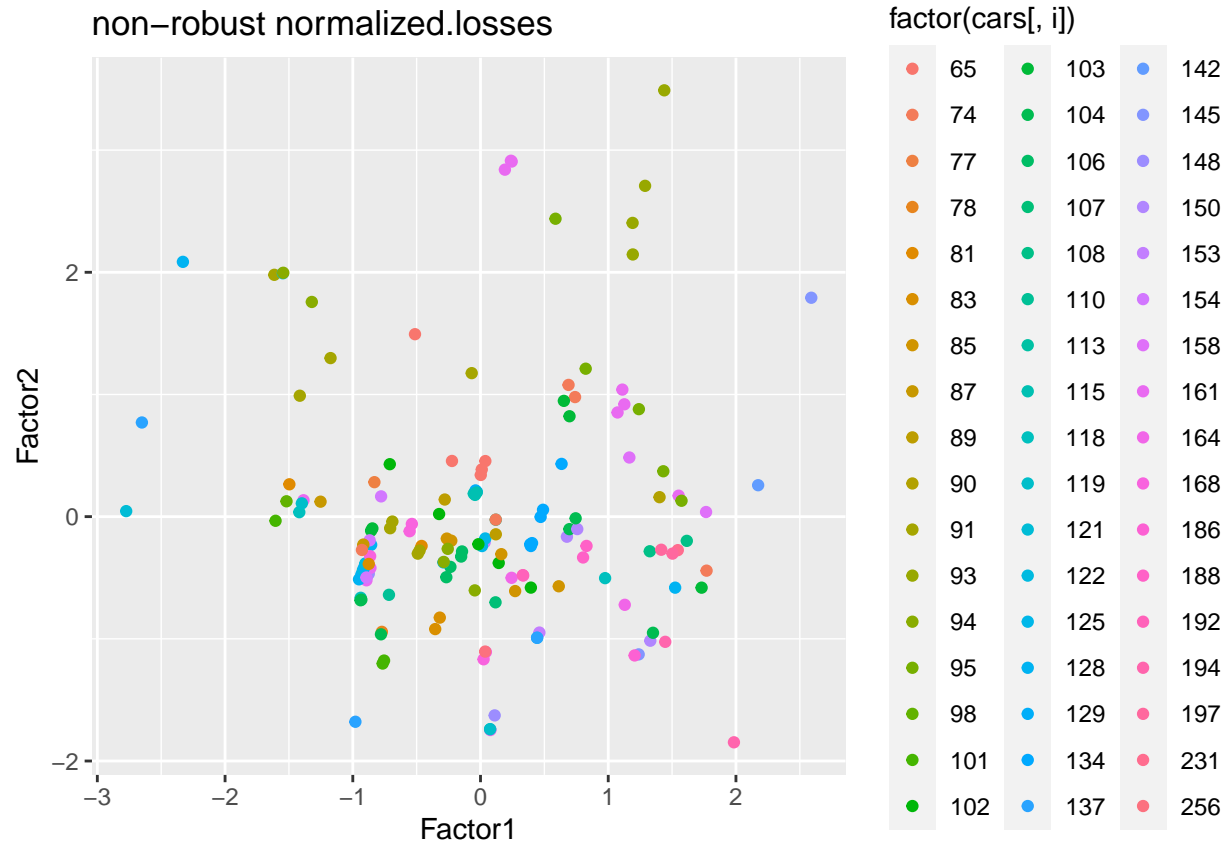


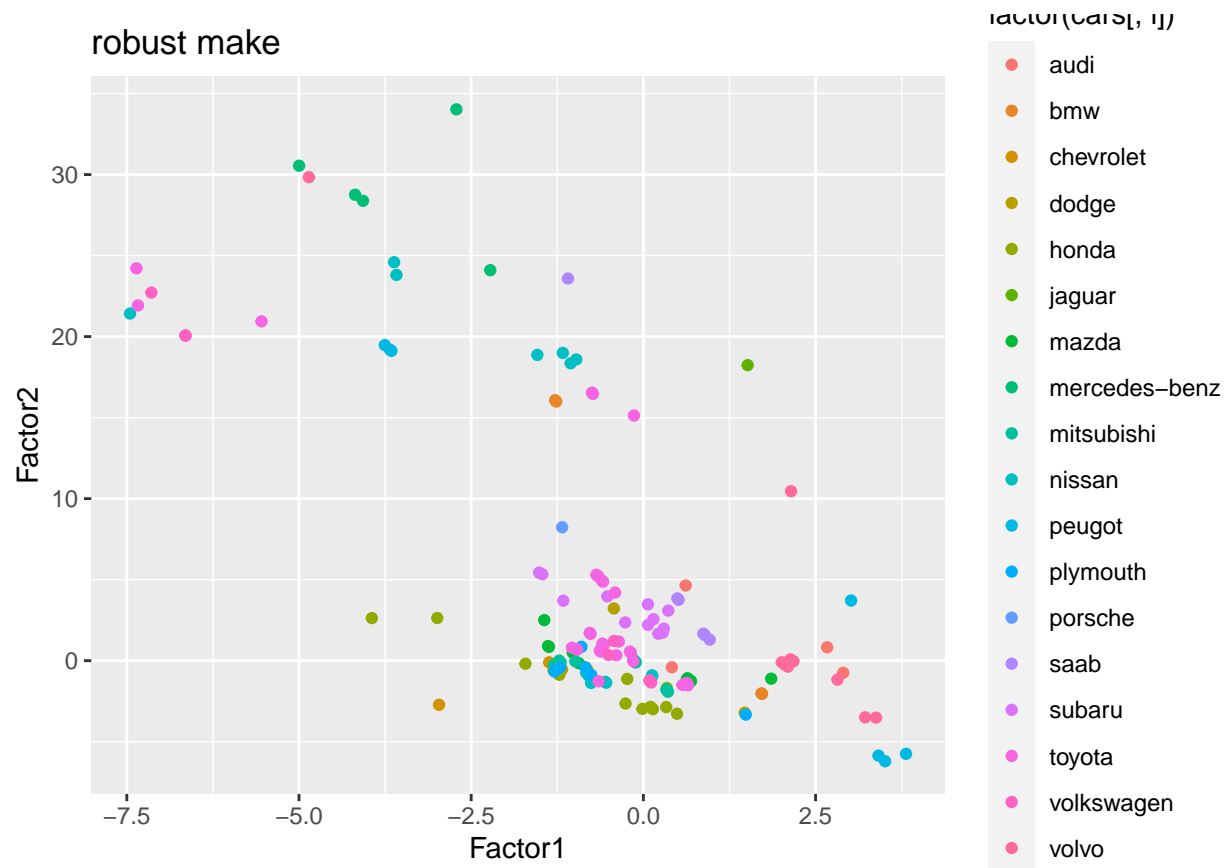
non-robust symboling

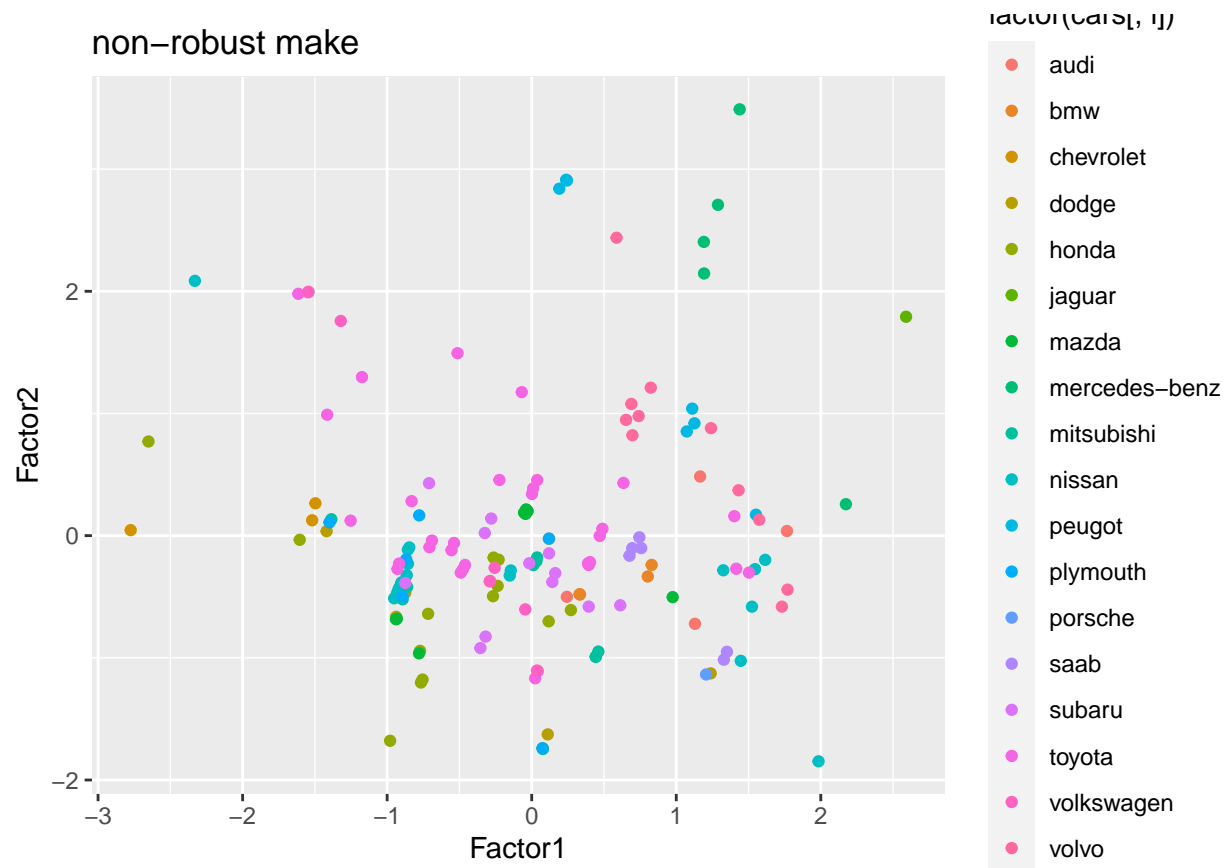


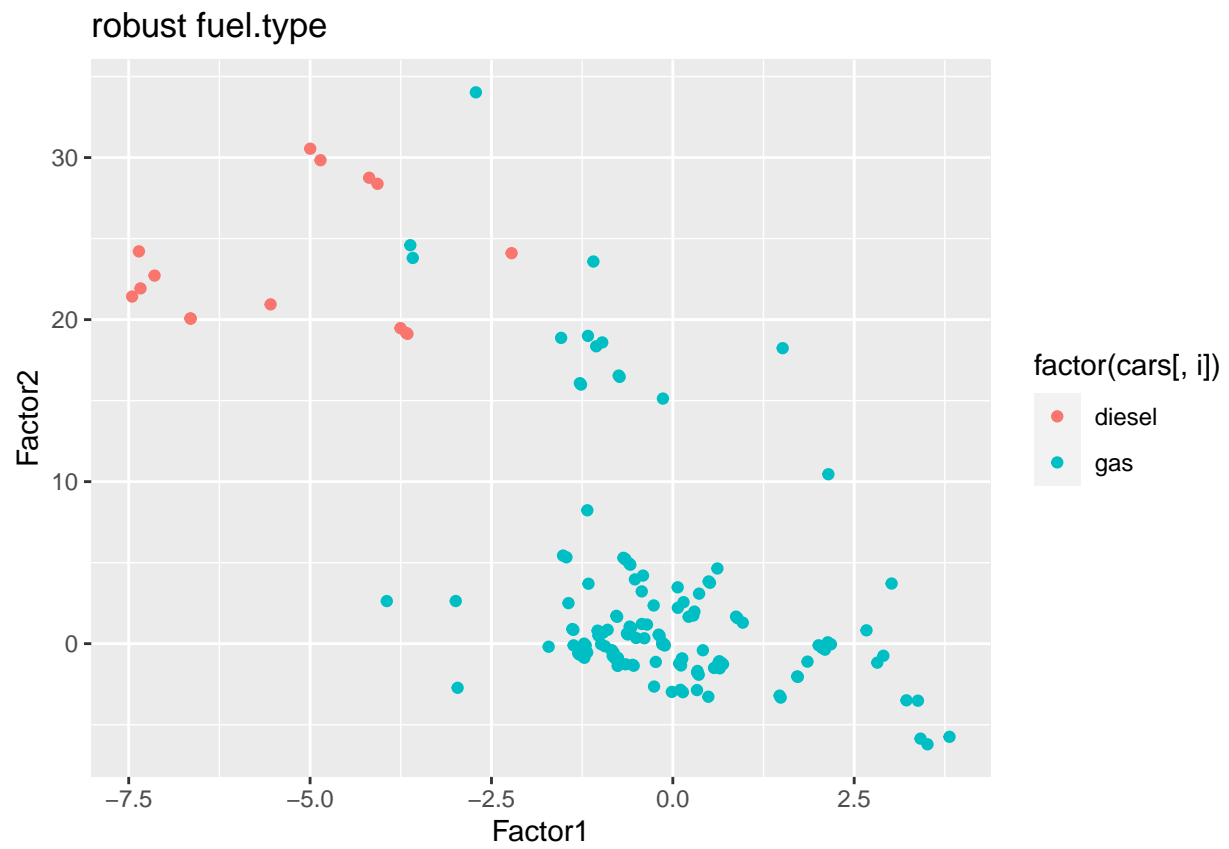


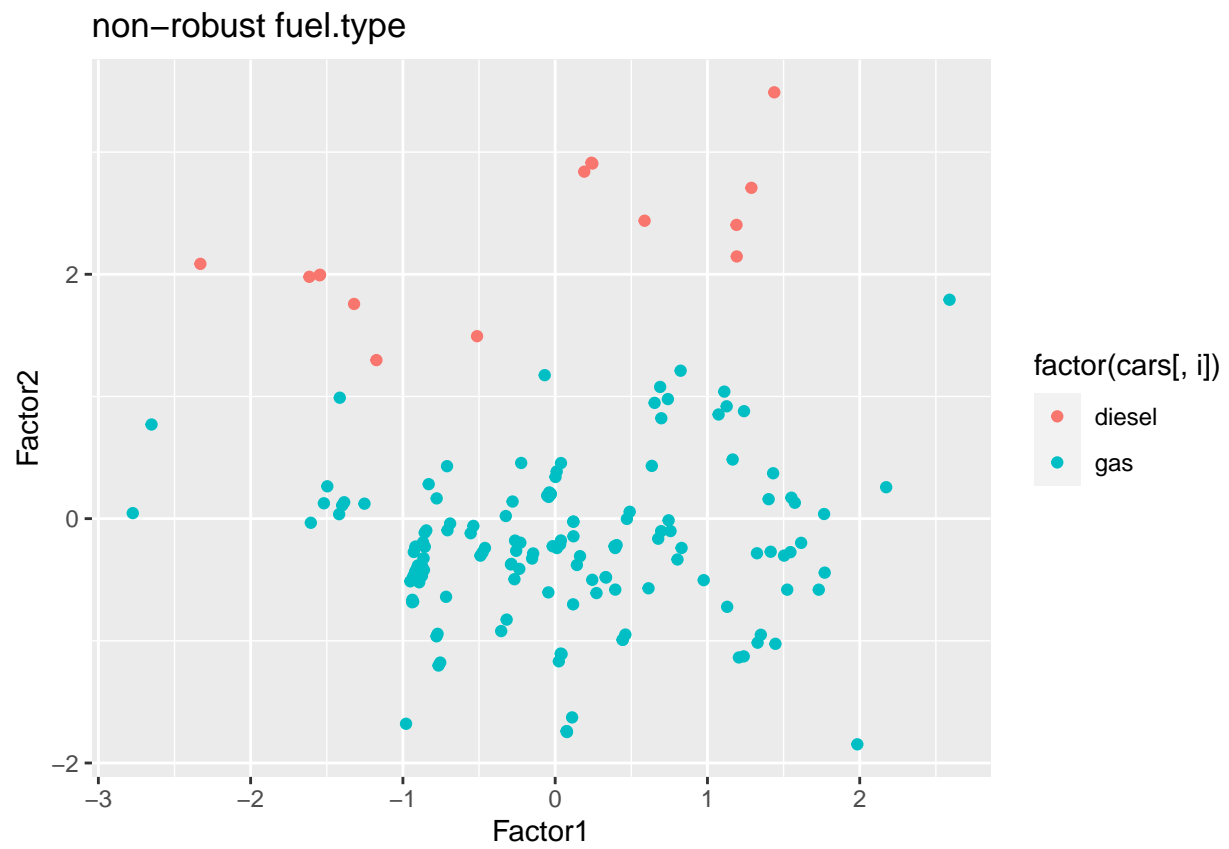


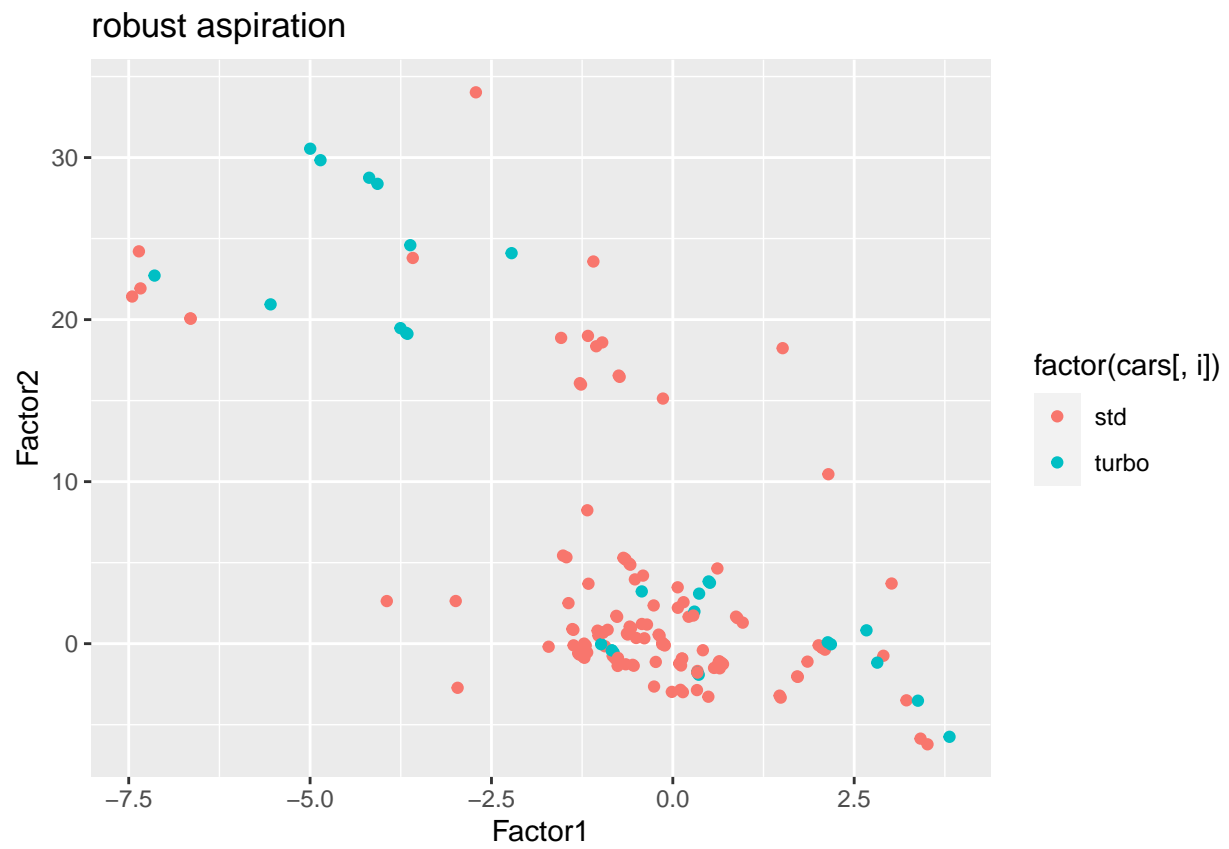




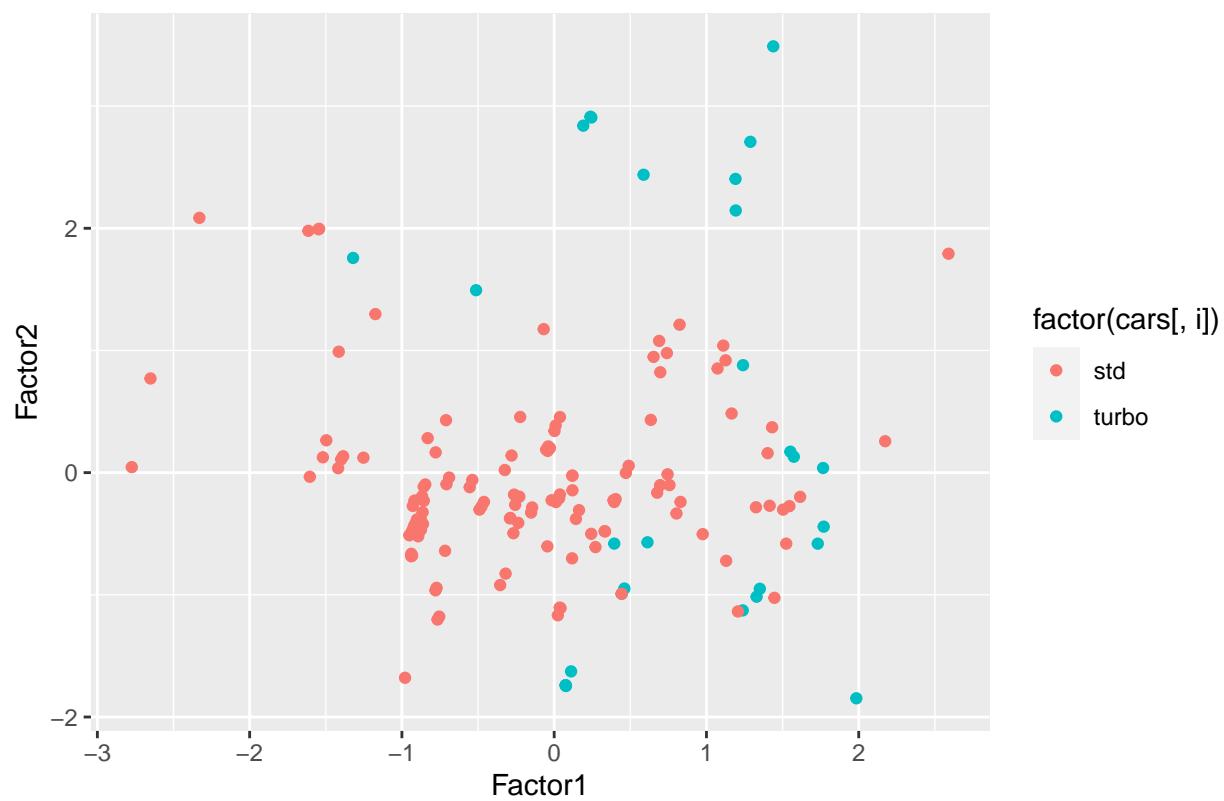


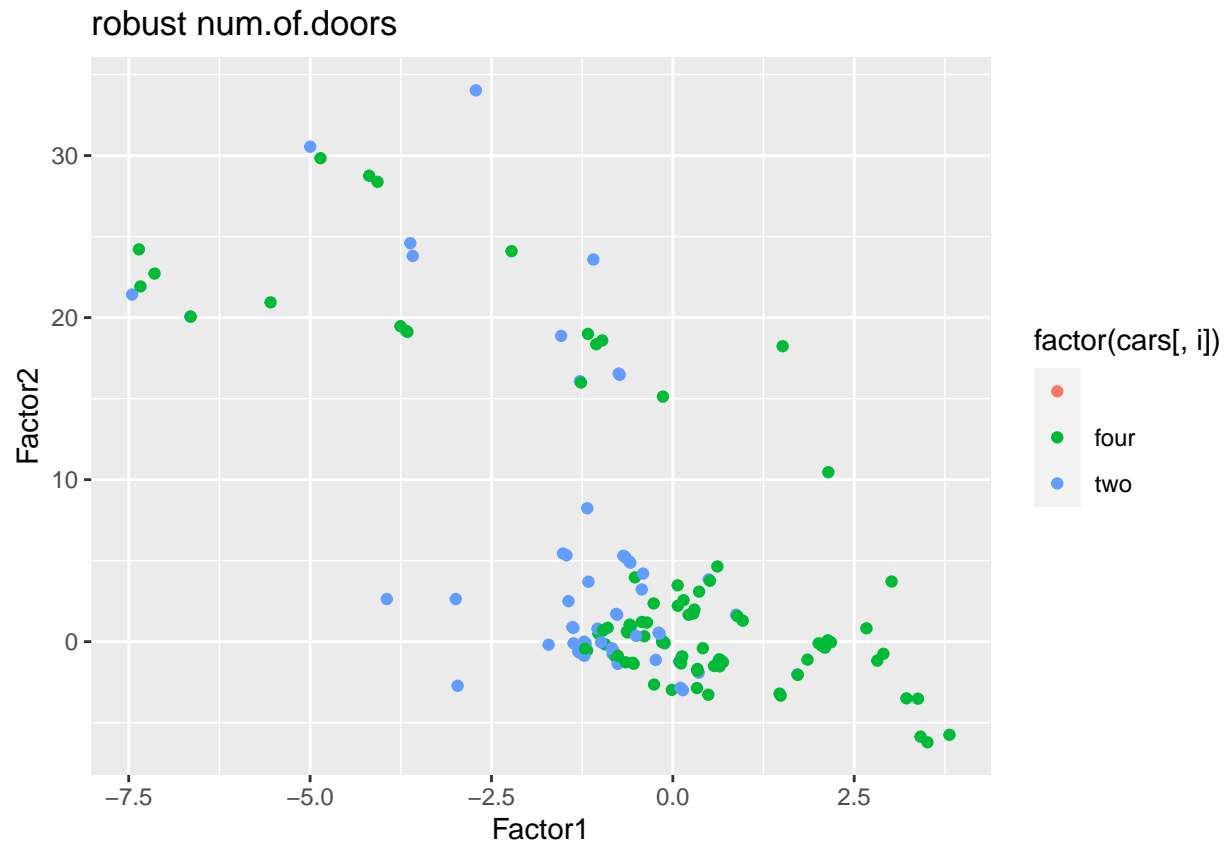




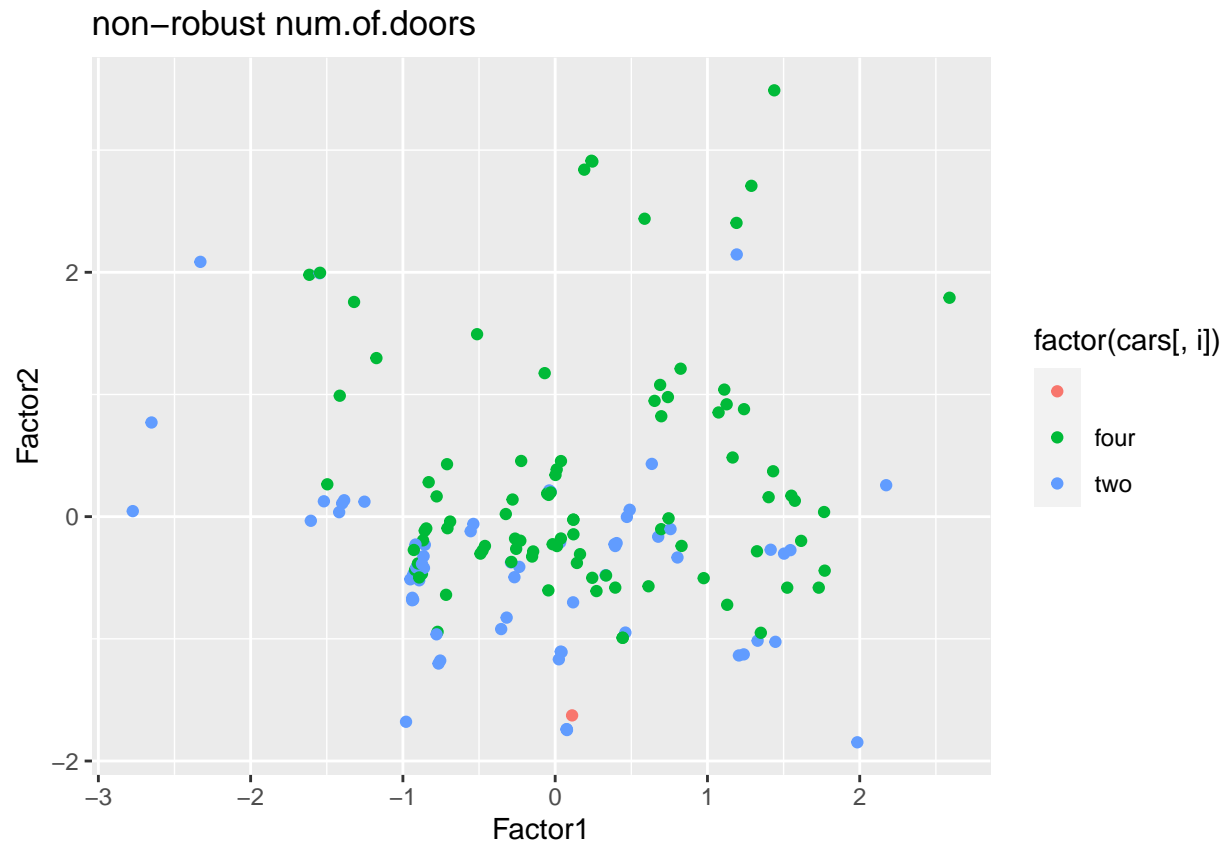


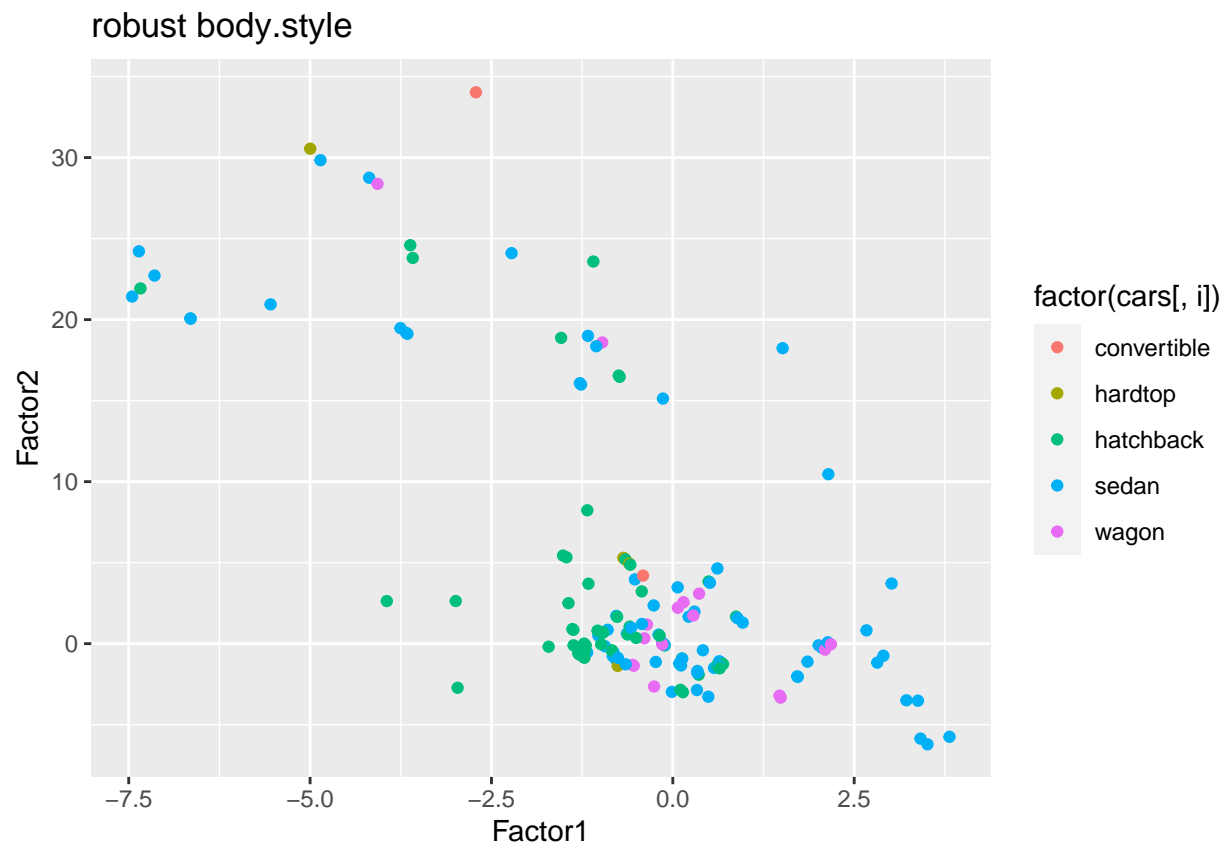
non-robust aspiration

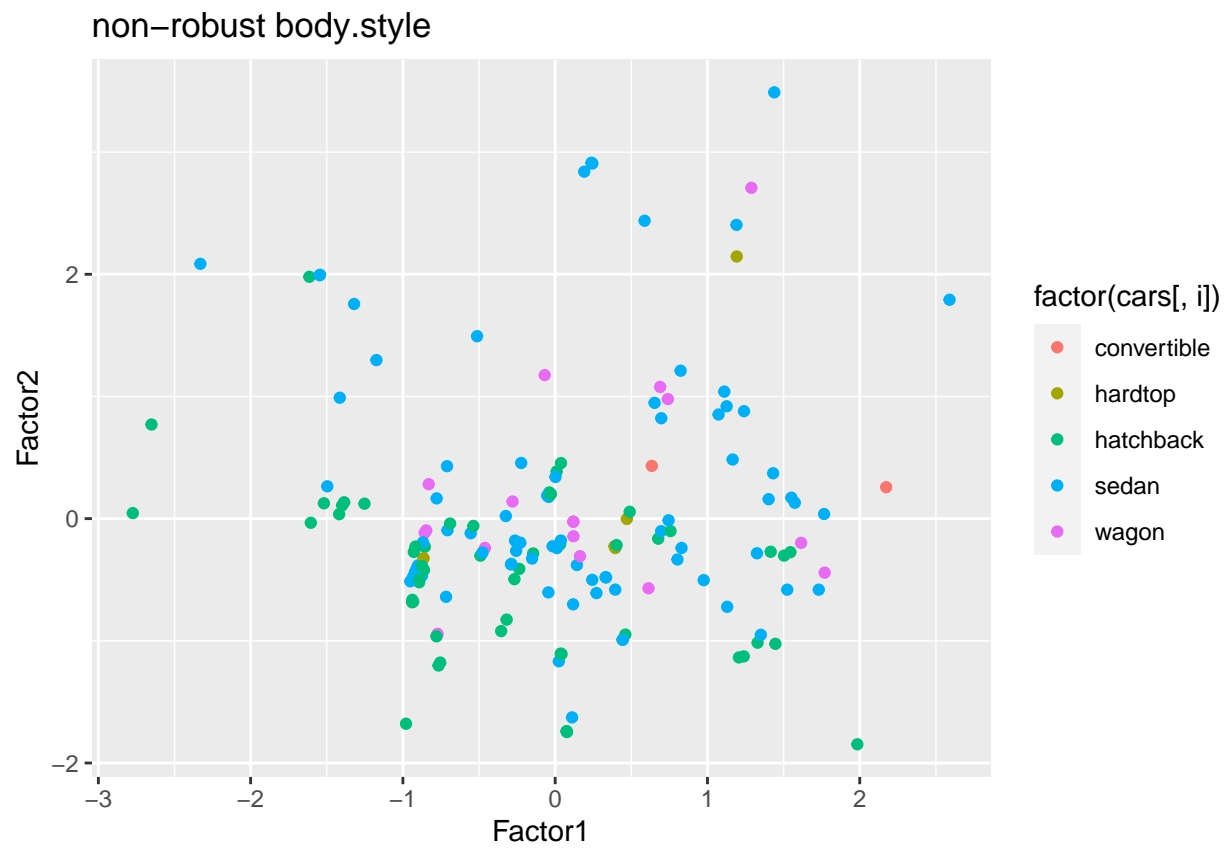


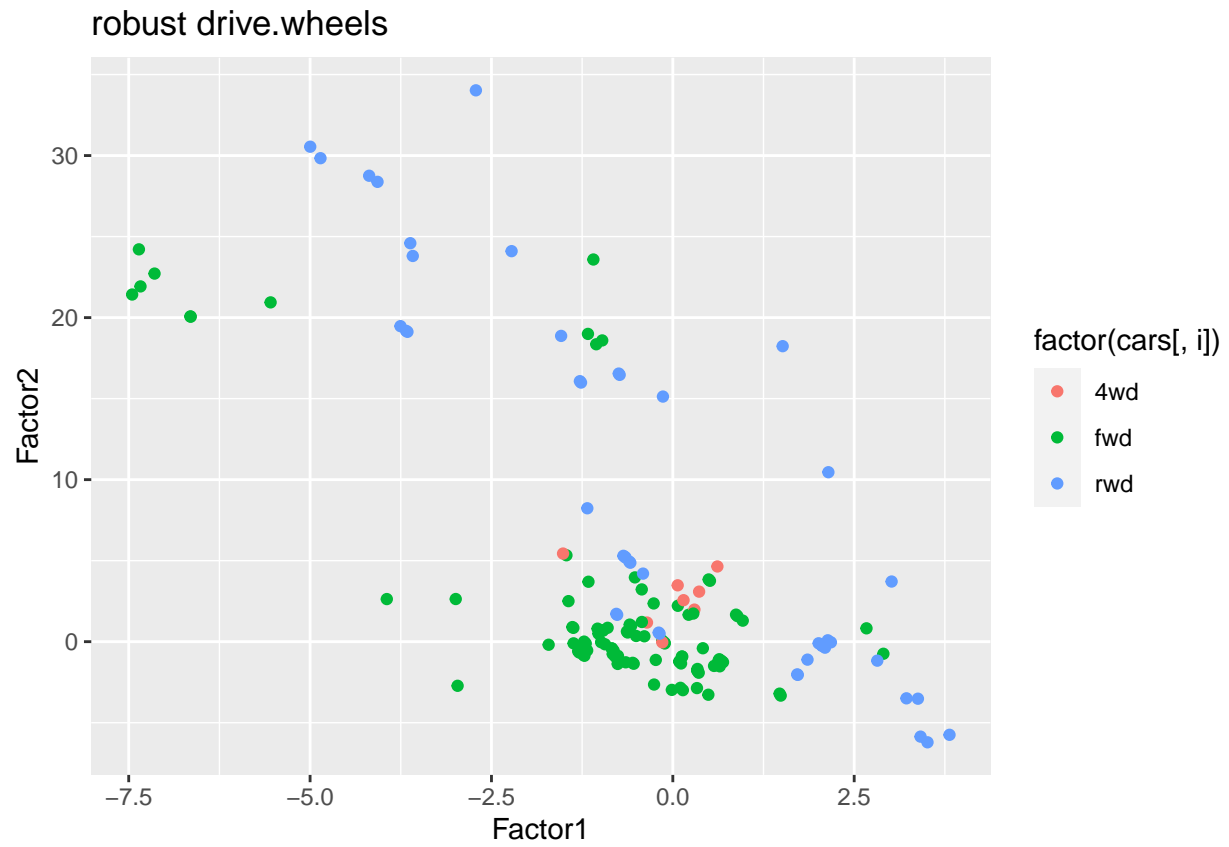




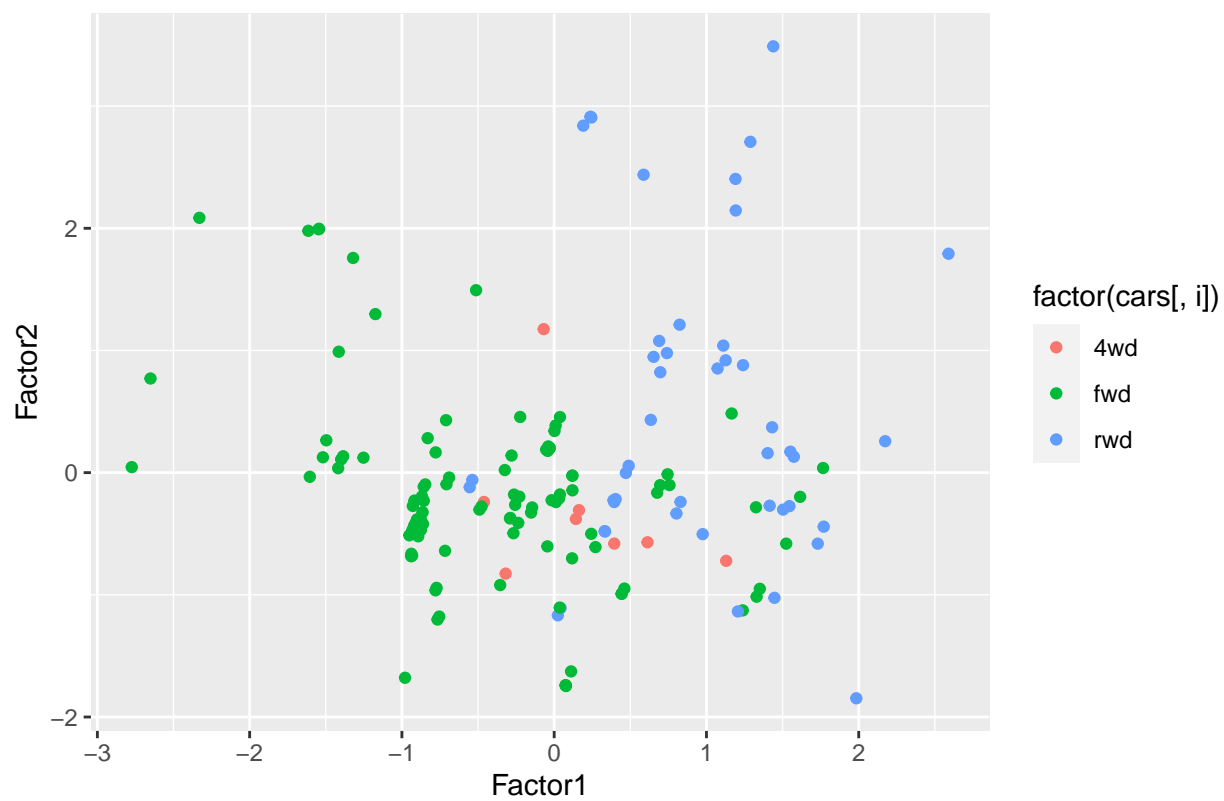


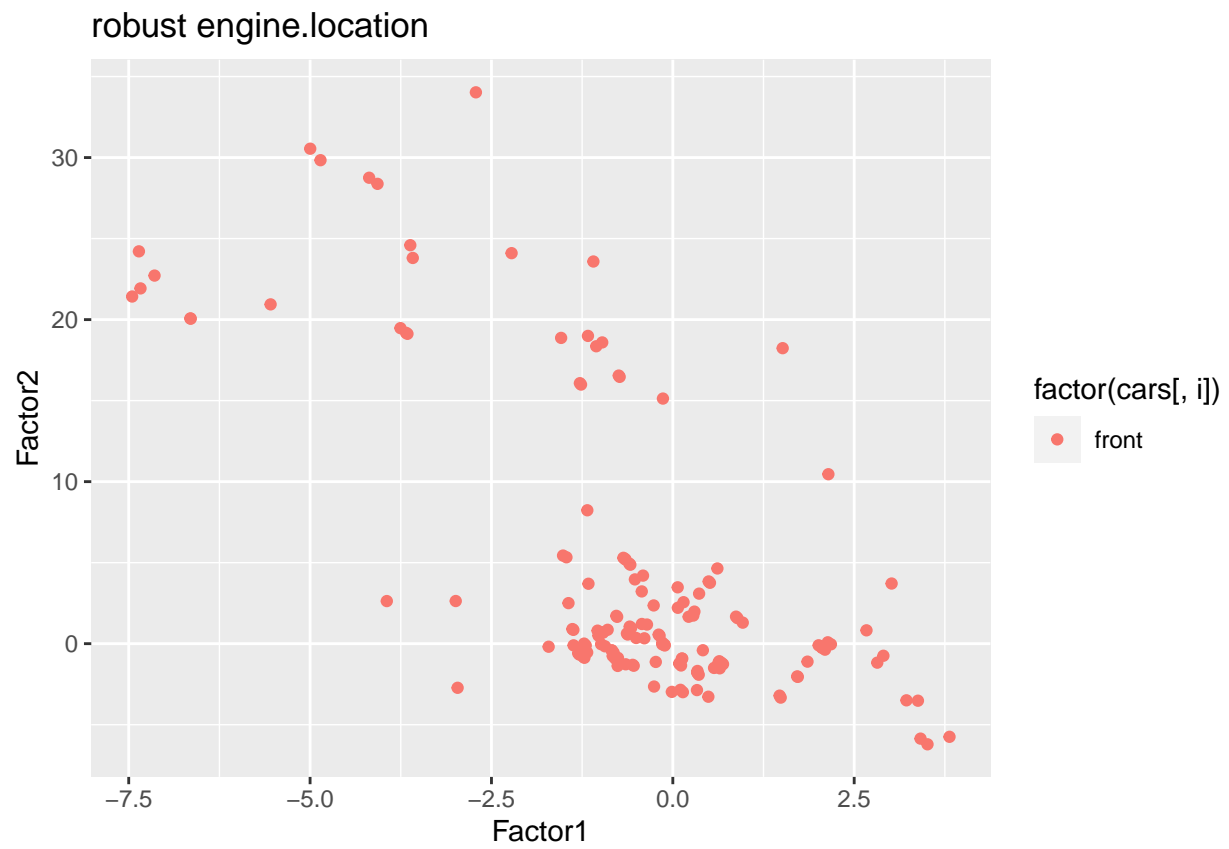


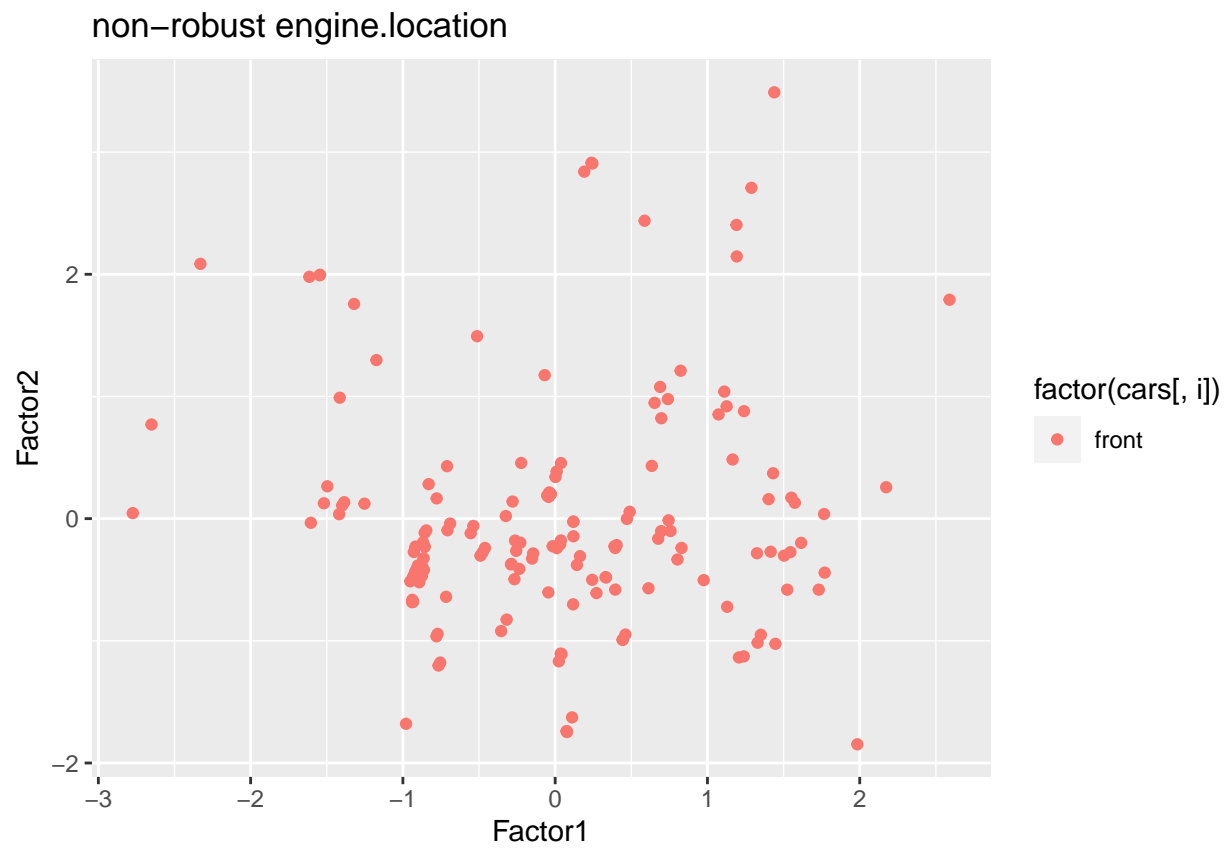


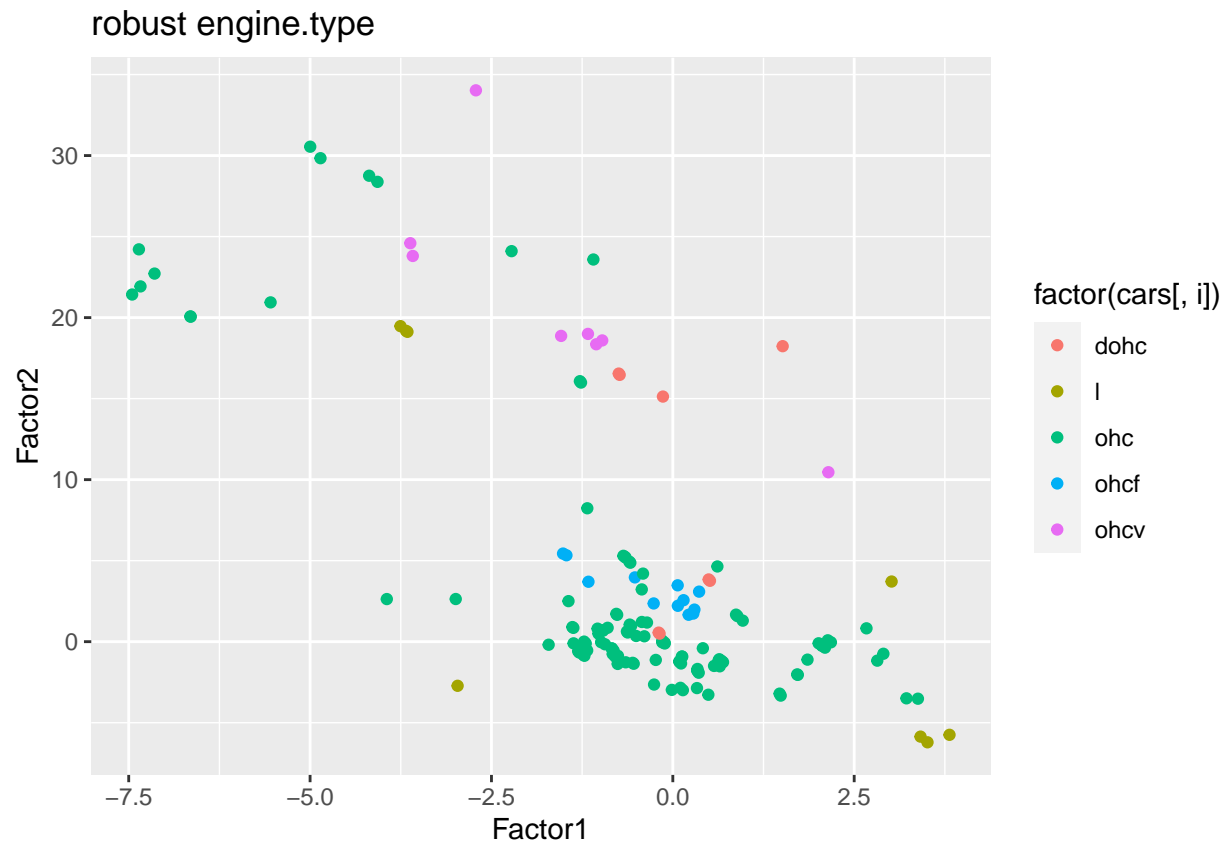


non-robust drive.wheels

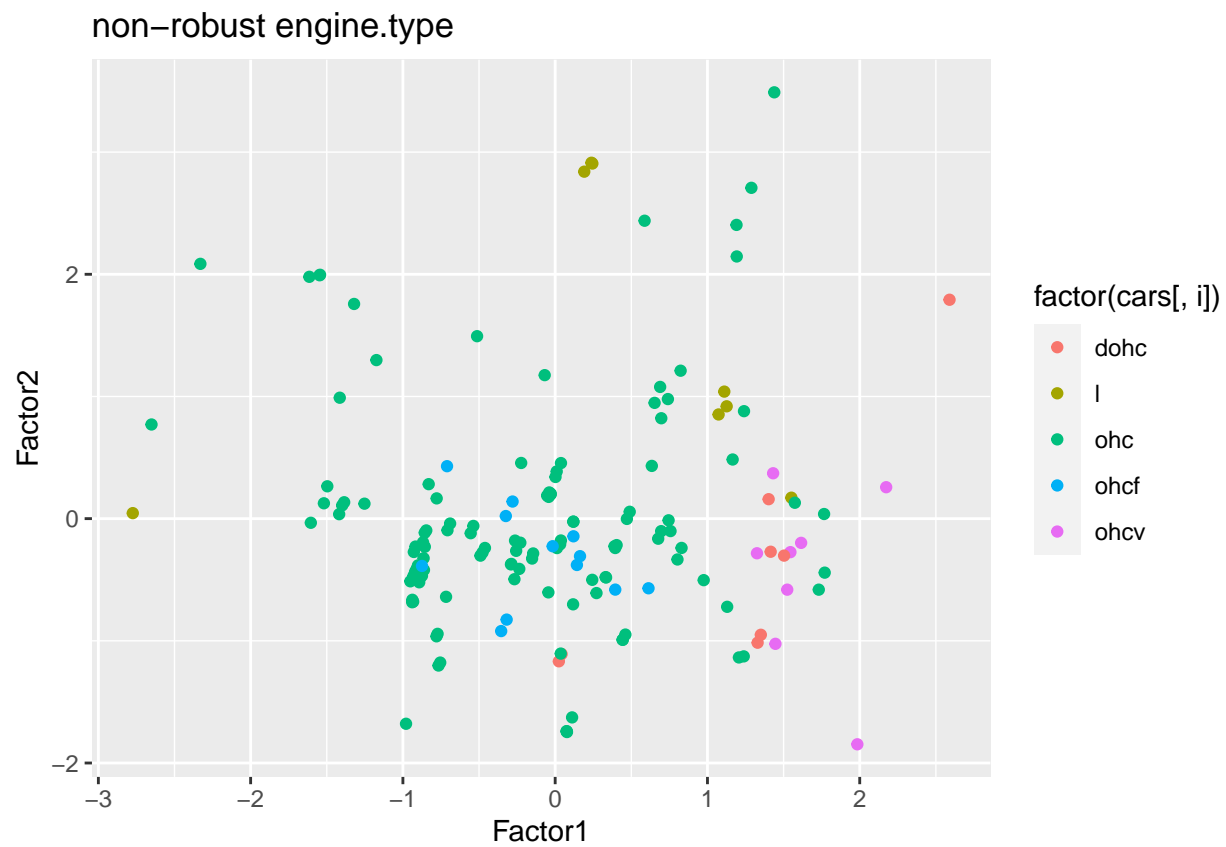


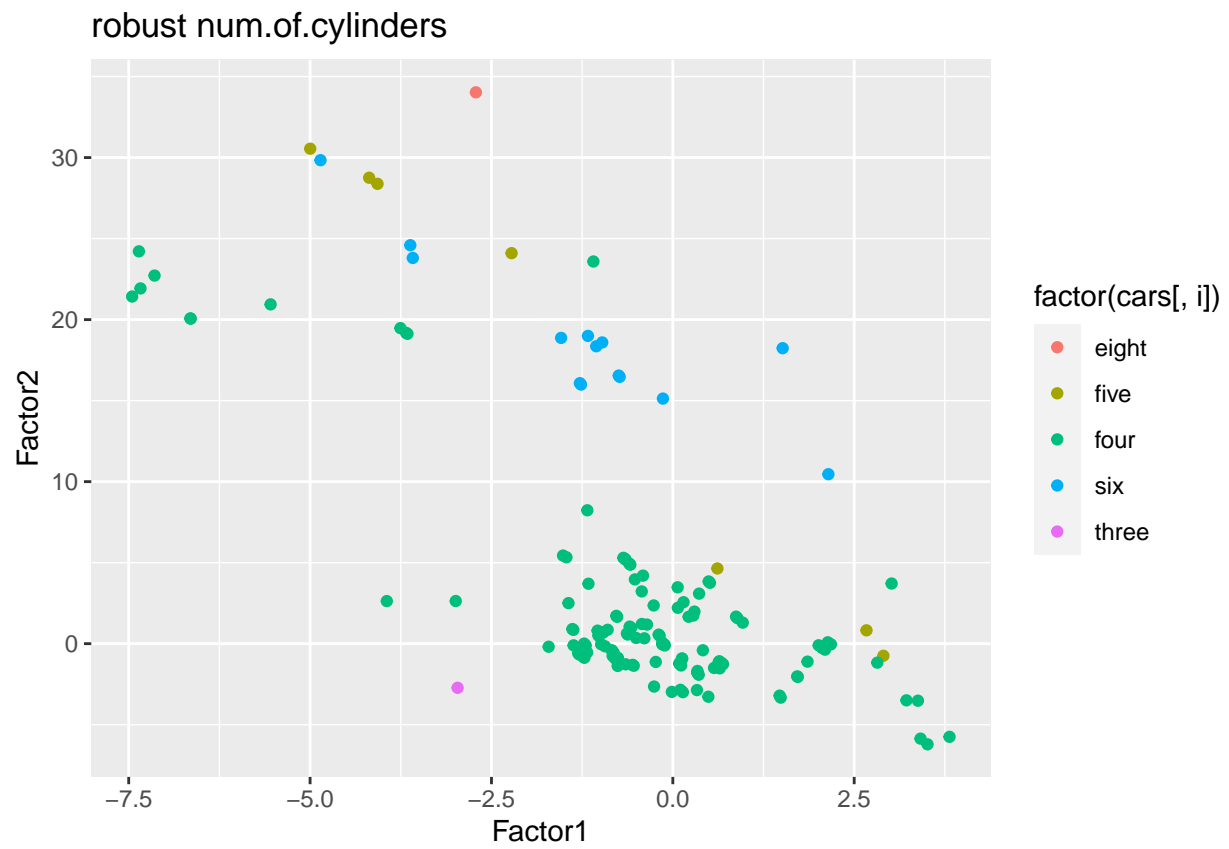


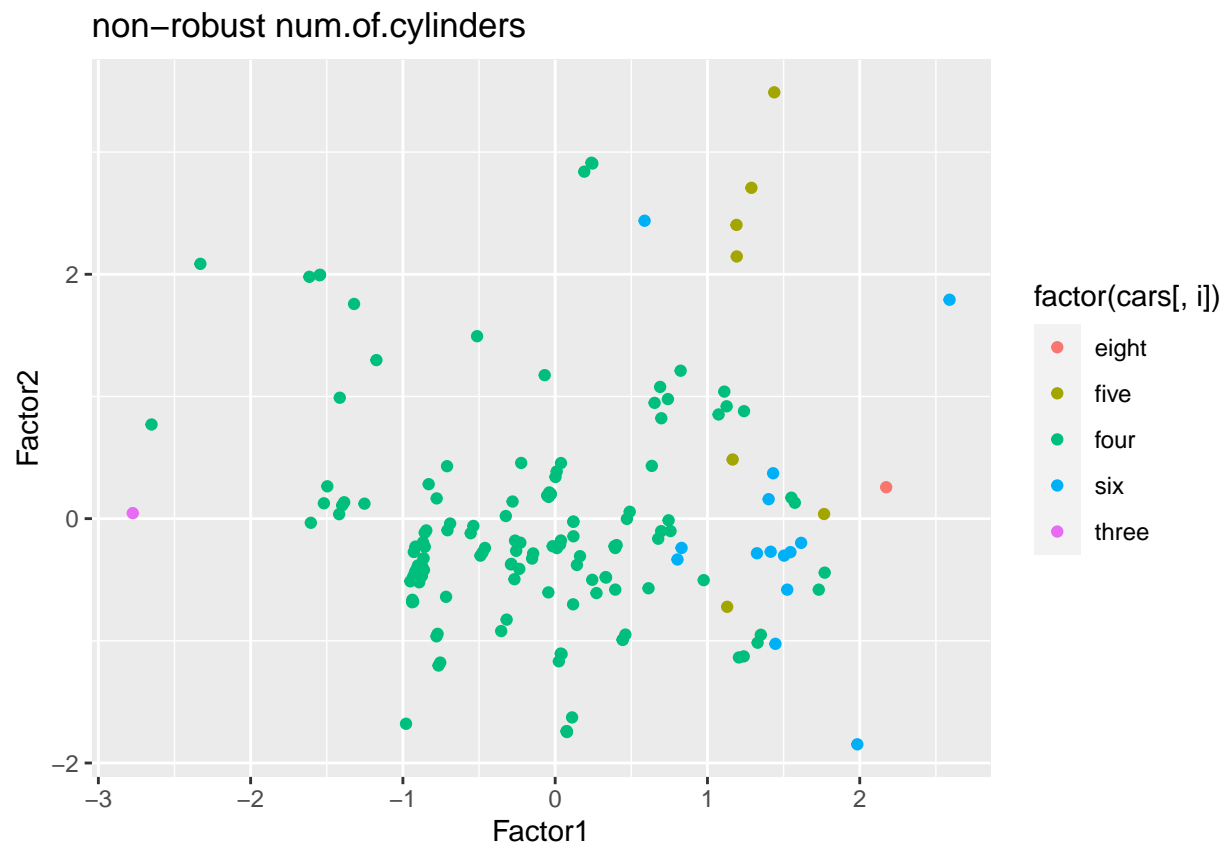


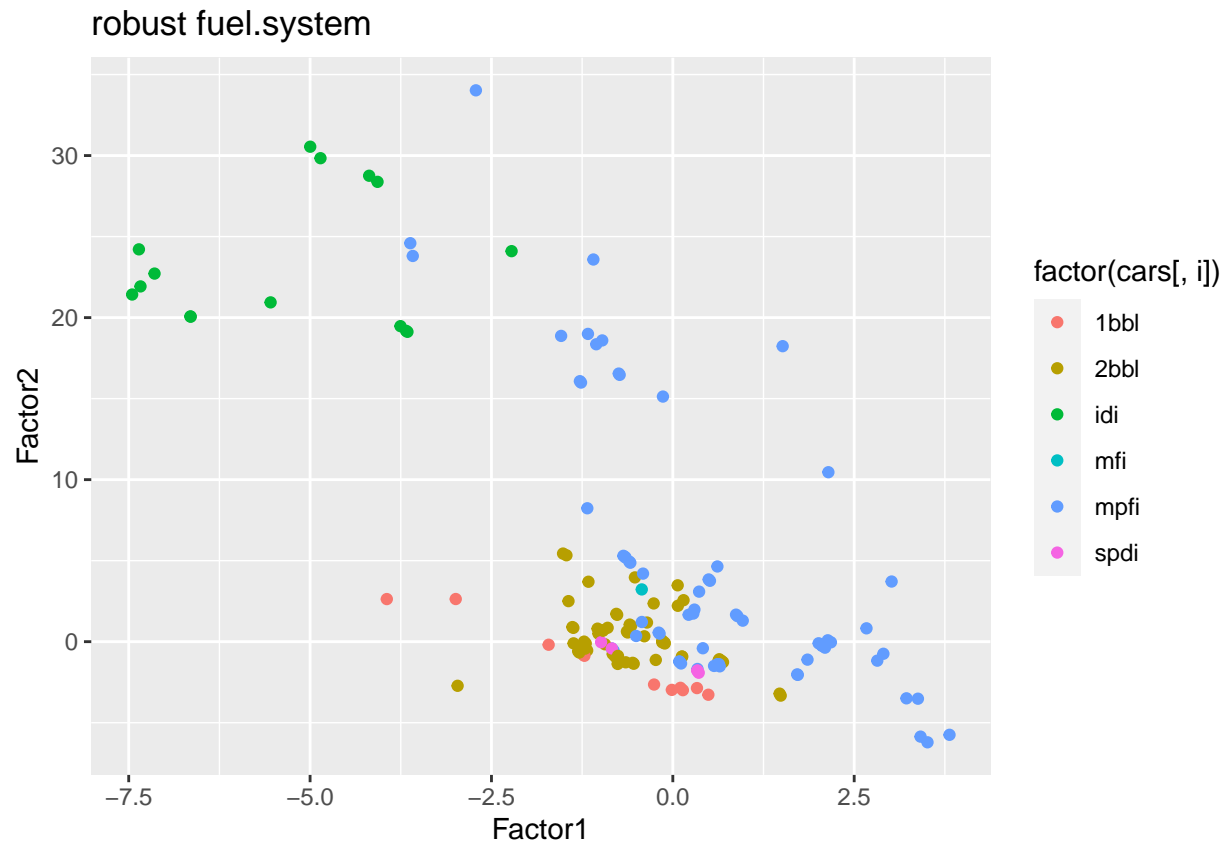


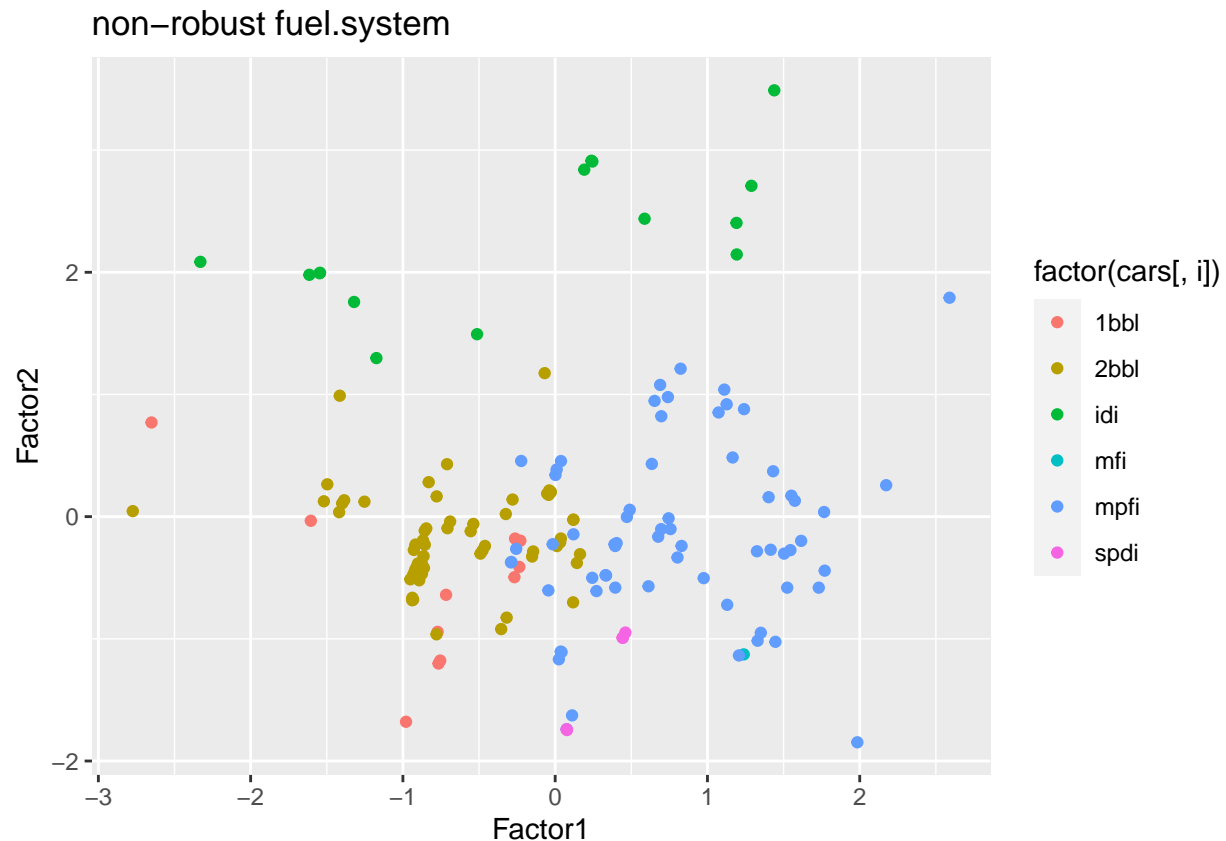












The one outlier on the top left is the only one where the variable *body.style* is “convertible”. Also The *num.of.cylinders* is eight for this observation only. In the non- robust method this is not detectable.

#### 4

With `print()` applied to the factor analysis output object you can see the variance proportions of the factors. How are these values computed?

```
print(robust_pfa)
```

```
##
## Call:
## pfa(x = df_rob_scaled, factors = 2, covmat = covmat, scores = "regression")
##
## Uniquenesses:
##      wheel.base      length      width      height
##      0.123          0.111      0.154      0.486
##      curb.weight    engine.size    bore      stroke
##      0.049          0.098      0.245      0.892
##      compression.ratio horsepower    peak.rpm    city.mpg
##      0.714          0.106      0.792      0.115
##      highway.mpg    price
##      0.081          0.211
##
## Loadings:
##      Factor1 Factor2
## wheel.base  0.936
```

```
## length          0.939
## width           0.911   0.124
## height          0.668  -0.261
## curb.weight     0.921   0.320
## engine.size     0.926   0.211
## bore            0.867
## stroke          0.327
## compression.ratio -0.533
## horsepower      0.626   0.708
## peak.rpm        -0.292   0.351
## city.mpg        -0.644  -0.686
## highway.mpg     -0.680  -0.676
## price           0.826   0.327
##
##               Factor1 Factor2
## SS loadings      7.530   2.293
## Proportion Var   0.538   0.164
## Cumulative Var   0.538   0.702
##
## The degrees of freedom for the model is 64 and the fit was NA
```

## 5

Compute the robust principal components based on the MCD estimator, and focus on the first two components. Rotate the components according to the varimax criterion (which is also the default for `pfa()`). This can be done by using the function `varimax()` from the package `GPArotation`.

```
library(GPArotation)
```

```
## Warning: package 'GPArotation' was built under R version 4.3.2
```

```
robust_pca = prcomp(~., data=data.frame(df_rob_scaled), covmat=covmat)
```

```
## Warning: In prcomp.default(x, ...) :
## extra argument 'covmat' will be disregarded
```

```
summary(robust_pca)
```

```
## Importance of components:
##               PC1      PC2      PC3      PC4      PC5      PC6      PC7
## Standard deviation  6.635 3.3721 1.46392 1.1215 0.99795 0.70272 0.61711
## Proportion of Variance 0.714 0.1844 0.03476 0.0204 0.01615 0.00801 0.00618
## Cumulative Proportion 0.714 0.8984 0.93317 0.9536 0.96972 0.97773 0.98390
##               PC8      PC9      PC10     PC11     PC12     PC13     PC14
## Standard deviation  0.52822 0.45348 0.42848 0.36615 0.33932 0.22296 0.15896
## Proportion of Variance 0.00453 0.00334 0.00298 0.00217 0.00187 0.00081 0.00041
## Cumulative Proportion 0.98843 0.99176 0.99474 0.99692 0.99878 0.99959 1.00000
```

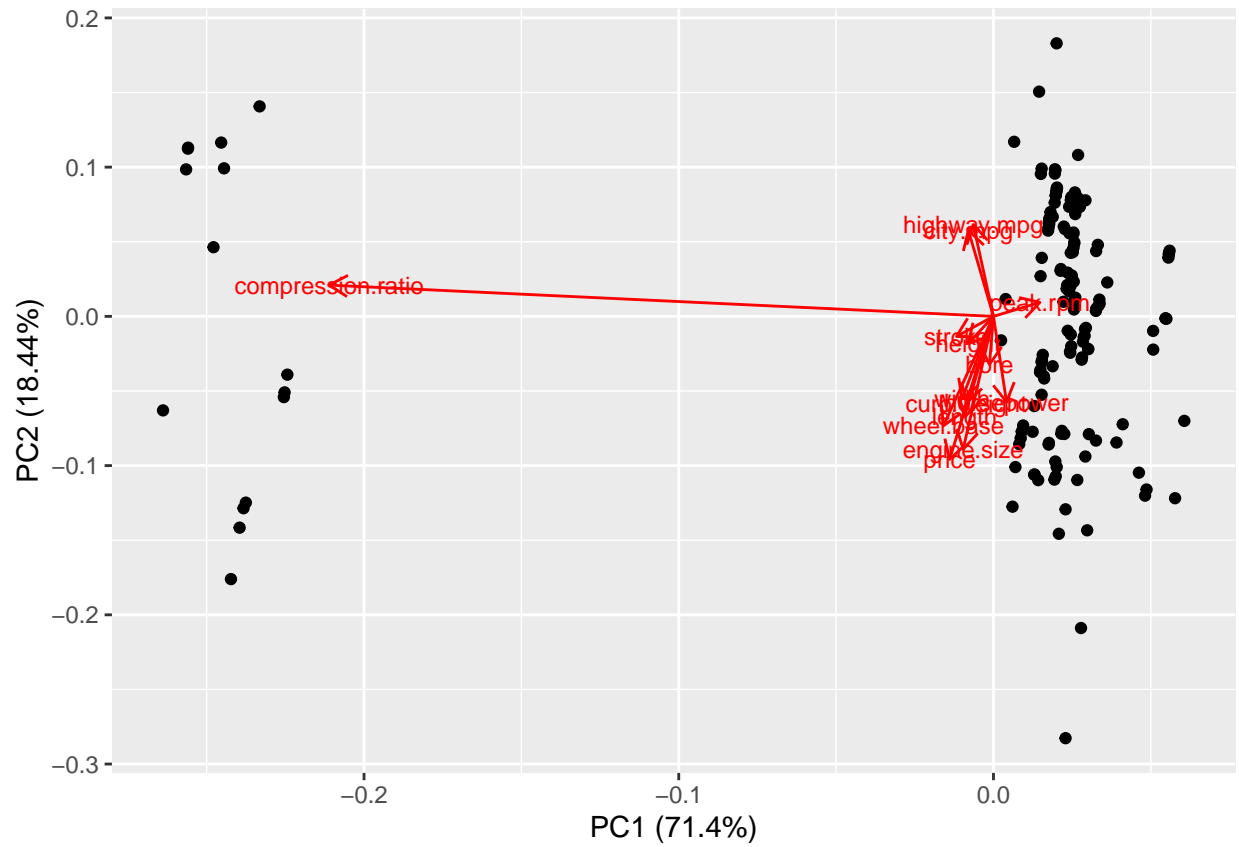
## 6

Compute also the scores to the rotated principal components, and present loadings and scores in a biplot. What are major differences to the robust factor analysis solution?

```
varimax_loadings = varimax(robust_pca$rotation)$loadings
varimax_loadings
```

```
##
## Loadings:
##          PC1 PC2 PC3 PC4 PC5 PC6 PC7 PC8 PC9 PC10 PC11 PC12 PC13 PC14
## wheel.base          -1
## length              1
## width              -1
## height             -1
## curb.weight              1
## engine.size        -1
## bore                1
## stroke              1
## compression.ratio -1
## horsepower              1
## peak.rpm             -1
## city.mpg              1
## highway.mpg          -1
## price                -1
##
##          PC1  PC2  PC3  PC4  PC5  PC6  PC7  PC8  PC9  PC10
## SS loadings  1.000 1.000 1.000 1.000 1.000 1.000 1.000 1.000 1.000 1.000
## Proportion Var 0.071 0.071 0.071 0.071 0.071 0.071 0.071 0.071 0.071 0.071
## Cumulative Var 0.071 0.143 0.214 0.286 0.357 0.429 0.500 0.571 0.643 0.714
##          PC11 PC12 PC13 PC14
## SS loadings  1.000 1.000 1.000 1.000
## Proportion Var 0.071 0.071 0.071 0.071
## Cumulative Var 0.786 0.857 0.929 1.000
```

```
scores = df_rob_scaled %*% t(solve(varimax_loadings))
autoplot(robust_pca, data=df_rob_scaled, loadings = TRUE, loadings.colour = 'red',
         loadings.label = TRUE, loadings.label.size = 3)
```



Here compression ratio makes up most of the first component.