

# Dimension Reduction. Analysis of Cardio training.

## Introduction

In this project we will analyse the information & dataset is a result of the medical examination, based on real patient information.

Also, in this project we will use: PCA & MCA & CA

Dataset contains the following information:

Age   Objective Feature   age   int (days)
Height   Objective Feature   height   int (cm)
Weight   Objective Feature   weight   float (kg)
Gender   Objective Feature   gender   categorical code
Systolic blood pressure   Examination Feature   ap_hi   int
Diastolic blood pressure   Examination Feature   ap_lo   int
Cholesterol   Examination Feature   cholesterol   1: normal, 2: above normal, 3: well above normal
Glucose   Examination Feature   gluc   1: normal, 2: above normal, 3: well above normal
Smoking   Subjective Feature   smoke   binary
Alcohol intake   Subjective Feature   alco   binary
Physical activity   Subjective Feature   active   binary
Presence or absence of cardiovascular disease   Target Variable   cardio   binary

## Data preparation

<pre>install.packages("factoextra") install.packages("gridExtra") install.packages("tidyverse") install.packages("ggplot2")</pre>	<pre>library(factoextra) library(FactoMineR) library(ggplot2) library(dplyr) library(reshape2) library(corrplot) library(gridExtra) library(grid)</pre>
---	---

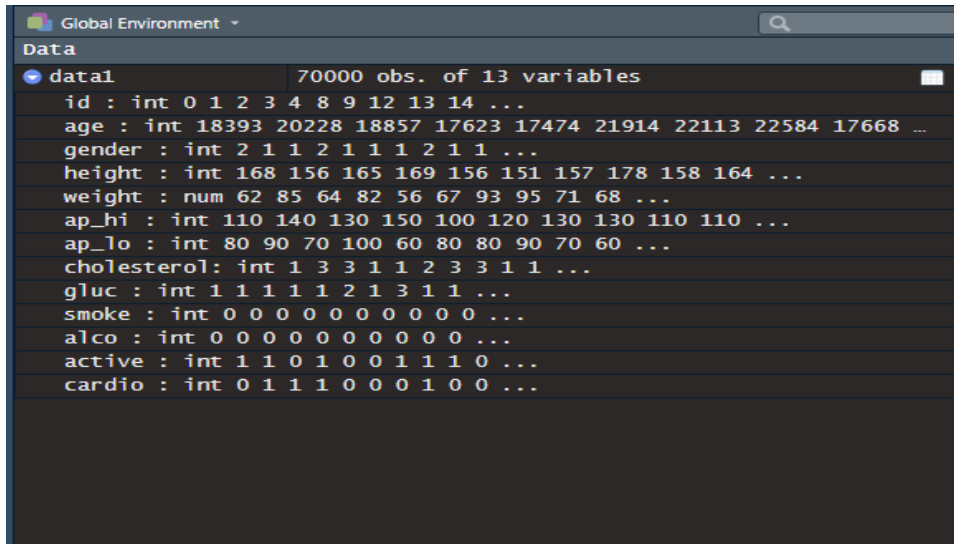
## Loading data

```
setwd("D:\\R and R Studio\\Dimension Reduction\\DimensionReduction")
```

```
getwd()
```

```
data1 <- read.csv("cardio_train.csv", sep = ";")
```

```
View(data1)
```



The screenshot shows the 'Global Environment' window in R Studio. It displays a data object named 'data1' with 70000 observations and 13 variables. The variables and their data types are listed below:

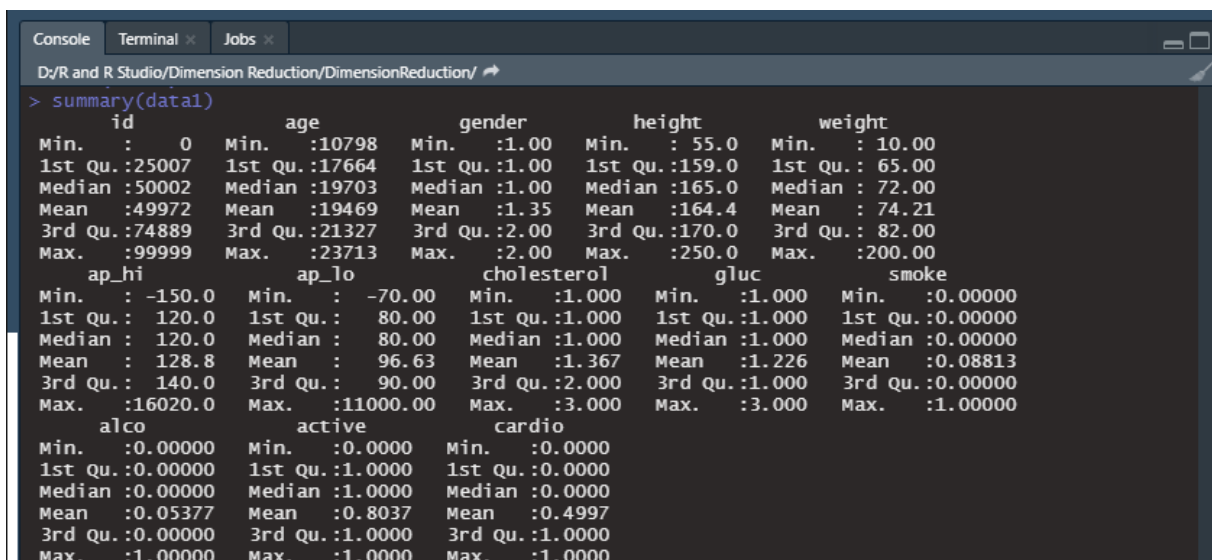
Variable	Data Type
id	int
age	int
gender	int
height	int
weight	num
ap_hi	int
ap_lo	int
cholesterol	int
gluc	int
smoke	int
alco	int
active	int
cardio	int

## Data Analysis

```
head(data1)
```

```
> head(data1)
  id age gender height weight ap_hi ap_lo cholesterol gluc smoke alco active cardio
1  0 18393     2   168    62   110    80           1    1      0     0      1      0
2  1 20228     1   156    85   140    90           3    1      0     0      1      1
3  2 18857     1   165    64   130    70           3    1      0     0      0      1
4  3 17623     2   169    82   150   100           1    1      0     0      1      1
5  4 17474     1   156    56   100    60           1    1      0     0      0      0
6  8 21914     1   151    67   120    80           2    2      0     0      0      0
```

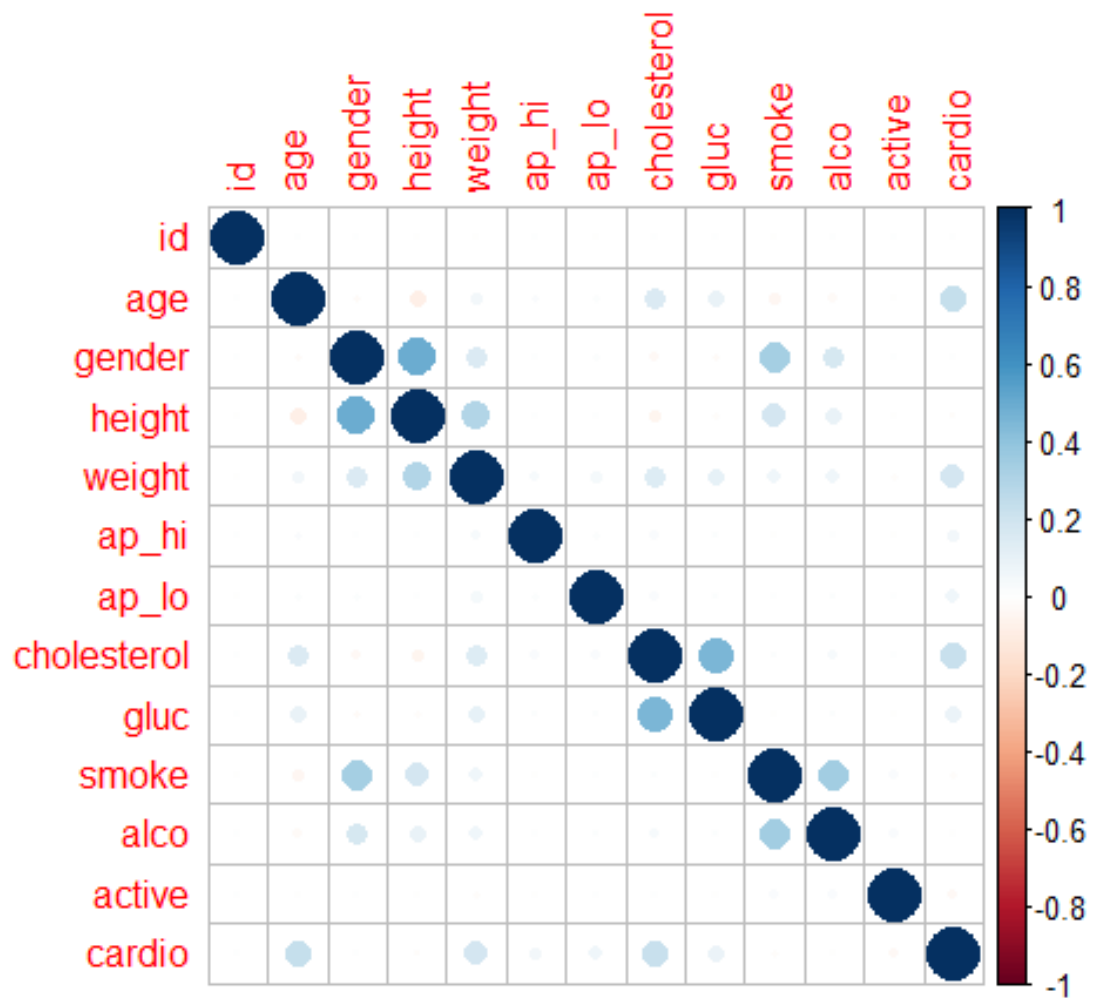
```
summary(data1)
```



The screenshot shows the 'Console' window in R Studio. It displays the output of the `summary(data1)` command, providing a summary of the data for each variable.

Variable	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
id	0	25007	50002	49972	74889	99999
age	10798	17664	19703	19469	21327	23713
gender	1.00	1.00	1.00	1.35	2.00	2.00
height	55.0	159.0	165.0	164.4	170.0	250.0
weight	10.00	65.00	72.00	74.21	82.00	200.00
ap_hi	-150.0	120.0	120.0	128.8	140.0	16020.0
ap_lo	-70.00	80.00	80.00	96.63	90.00	11000.00
cholesterol	1.000	1.000	1.000	1.367	2.000	3.000
gluc	1.000	1.000	1.000	1.226	1.000	3.000
smoke	0.00000	0.00000	0.00000	0.08813	0.00000	1.00000
alco	0.00000	0.00000	0.00000	0.05377	0.00000	1.00000
active	0.0000	1.0000	1.0000	0.8037	1.0000	1.0000
cardio	0.0000	0.0000	0.0000	0.4997	1.0000	1.0000

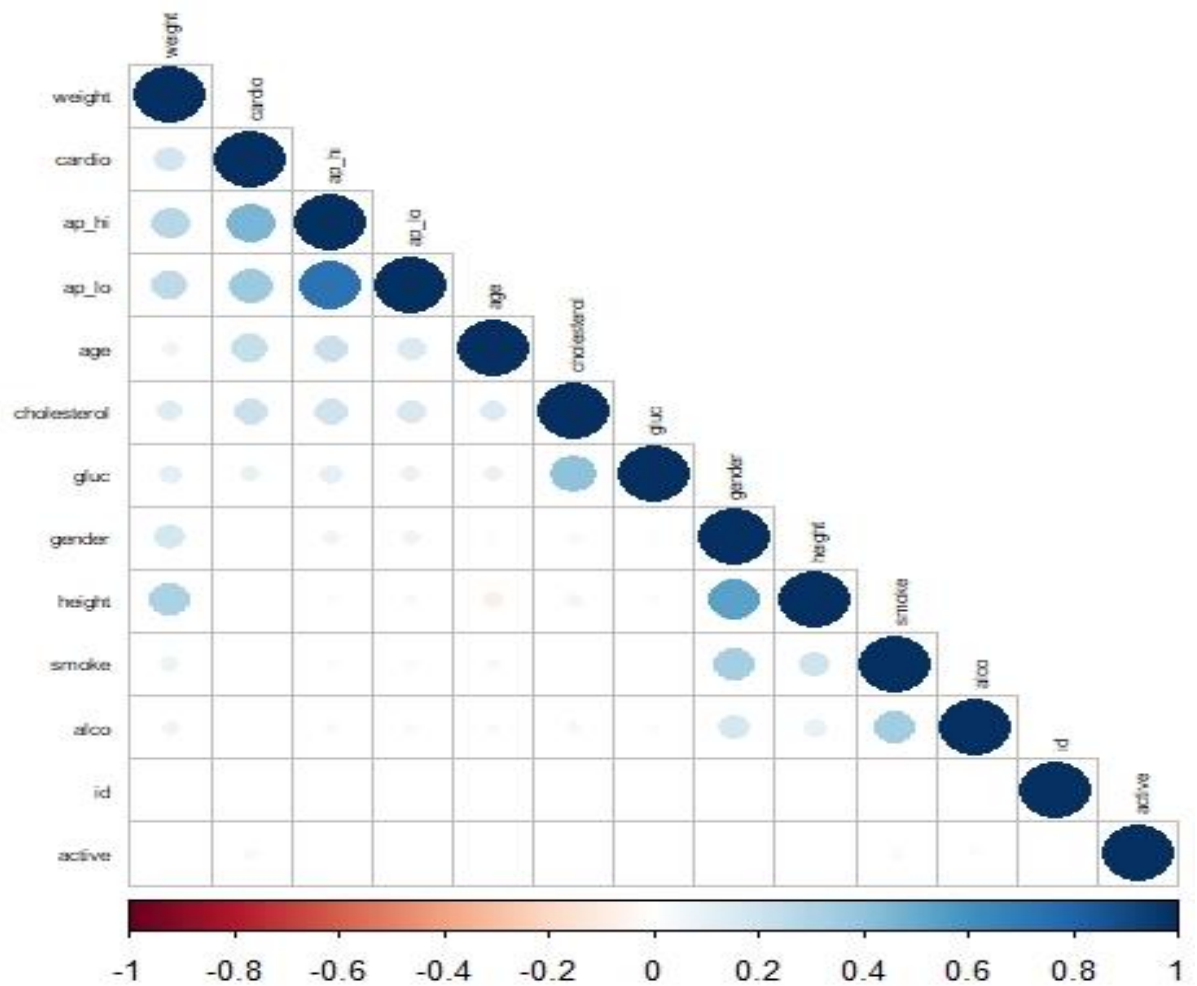
```
corrplot(corr_data)
```



Here, correlation using matrix Cholesterol.

```
cor.matrix <- cor(data1, method = c("spearman"))
```

```
corrplot(cor.matrix, type = "lower", order = "hclust", tl.col = "black", tl.cex = 0.5)
```



## PCA

Choosing number in Component

```
> data.pca <- prcomp(data1, center=TRUE, scale=TRUE)
```

```
> eigen(cor(data1))$values
```

```
[1] 1.9430876 1.7384125 1.1616289 1.0544779 1.0021538 0.9977607 0.9841246 0.9731207  
0.8270677 0.7235359
```

```
[11] 0.6215278 0.5243303 0.4487716
```

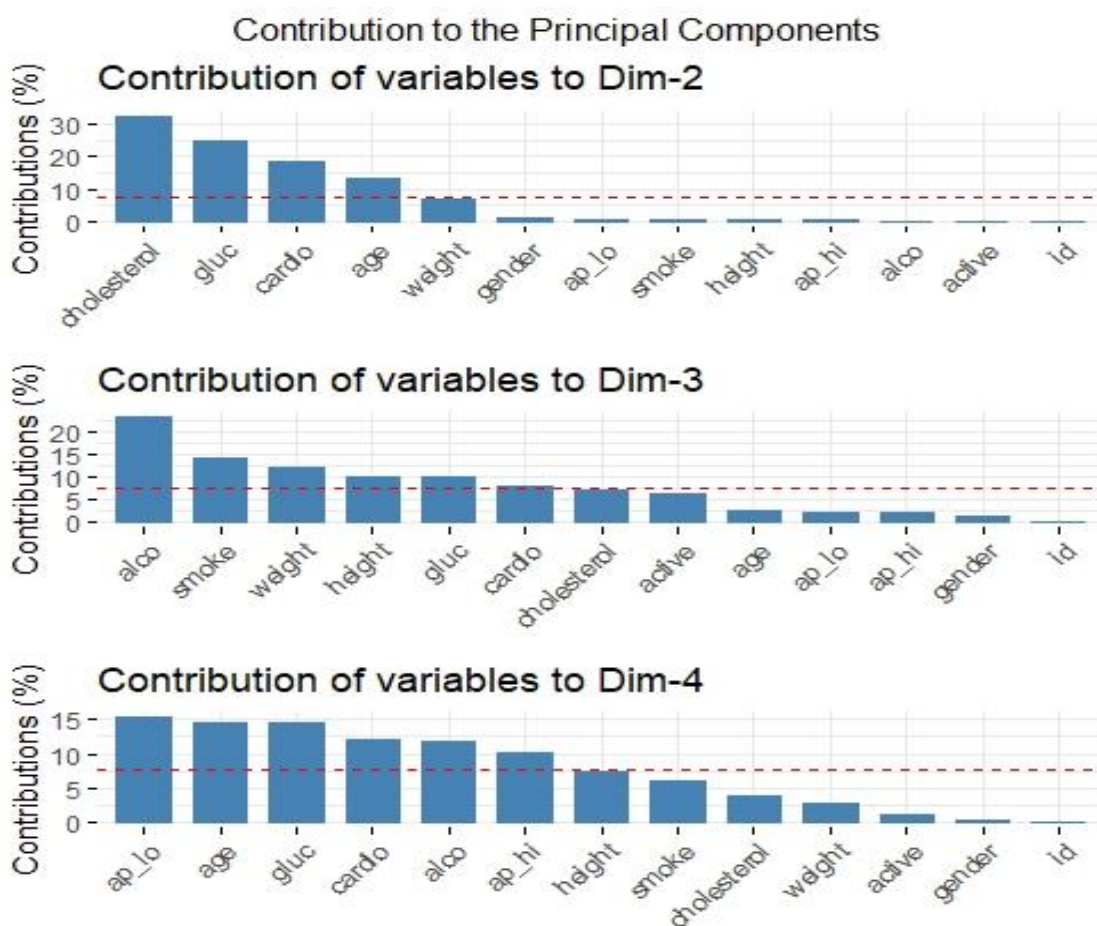
```
var <- get_pca_var(data.pca)
```

```
a<-fviz_contrib(data.pca, "var",axes = 2)
```

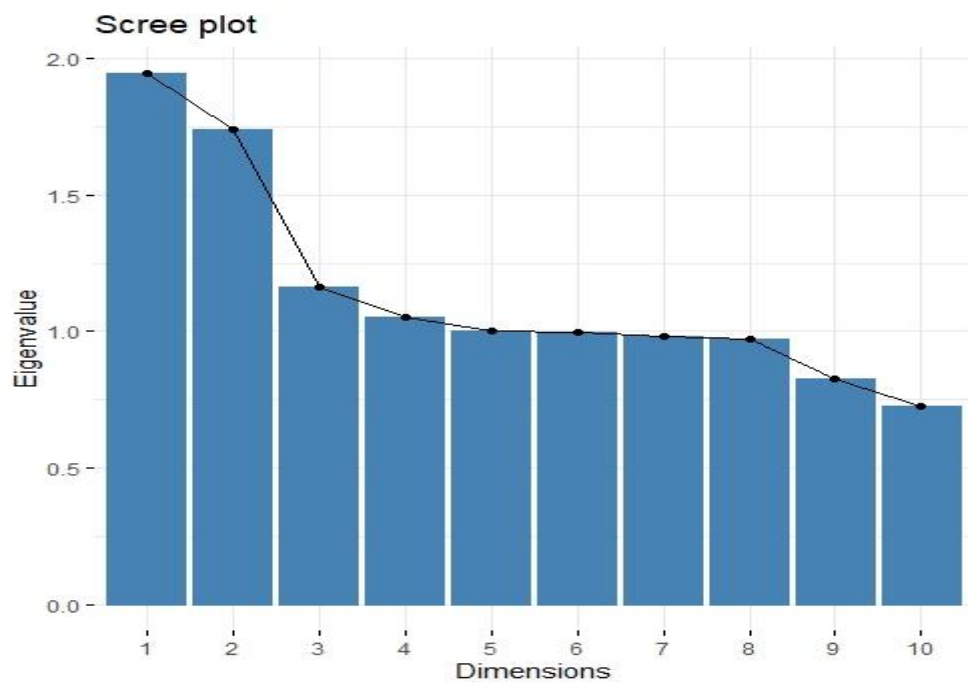
```
b<-fviz_contrib(data.pca, "var",axes = 3)
```

```
c<-fviz_contrib(data.pca, "var",axes = 4)
```

```
grid.arrange(a,b,c,top='Contribution to the Principal Components')
```

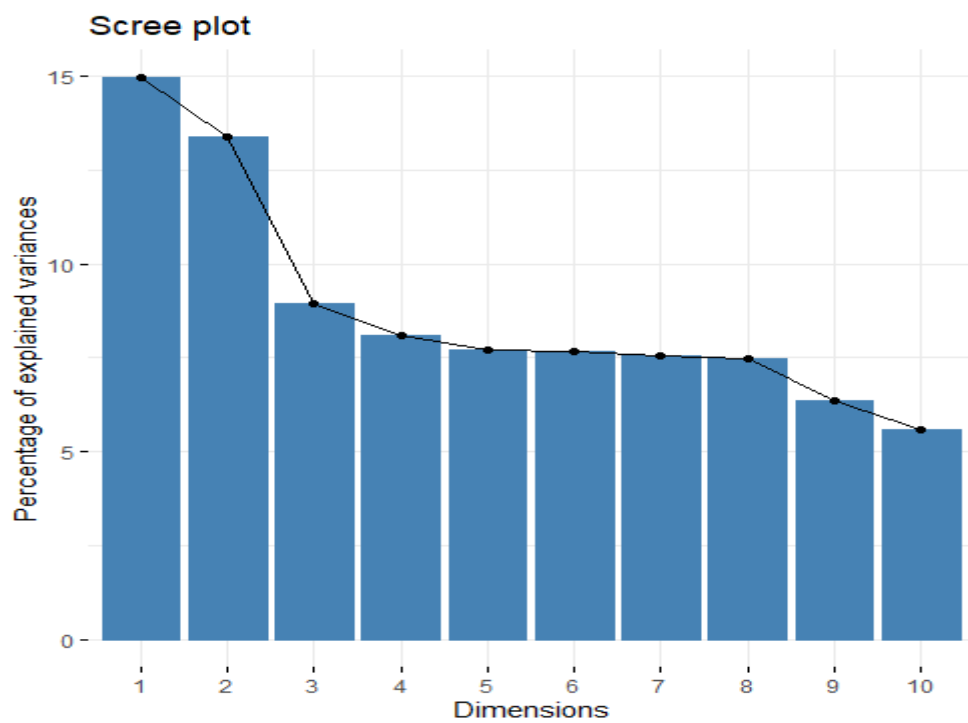


```
fviz_eig(data.pca, choice='eigenvalue')
```



Two different contributions to see dimensionality space result.

```
fviz_eig(data.pca)
```



```
summary(data.pca)
```

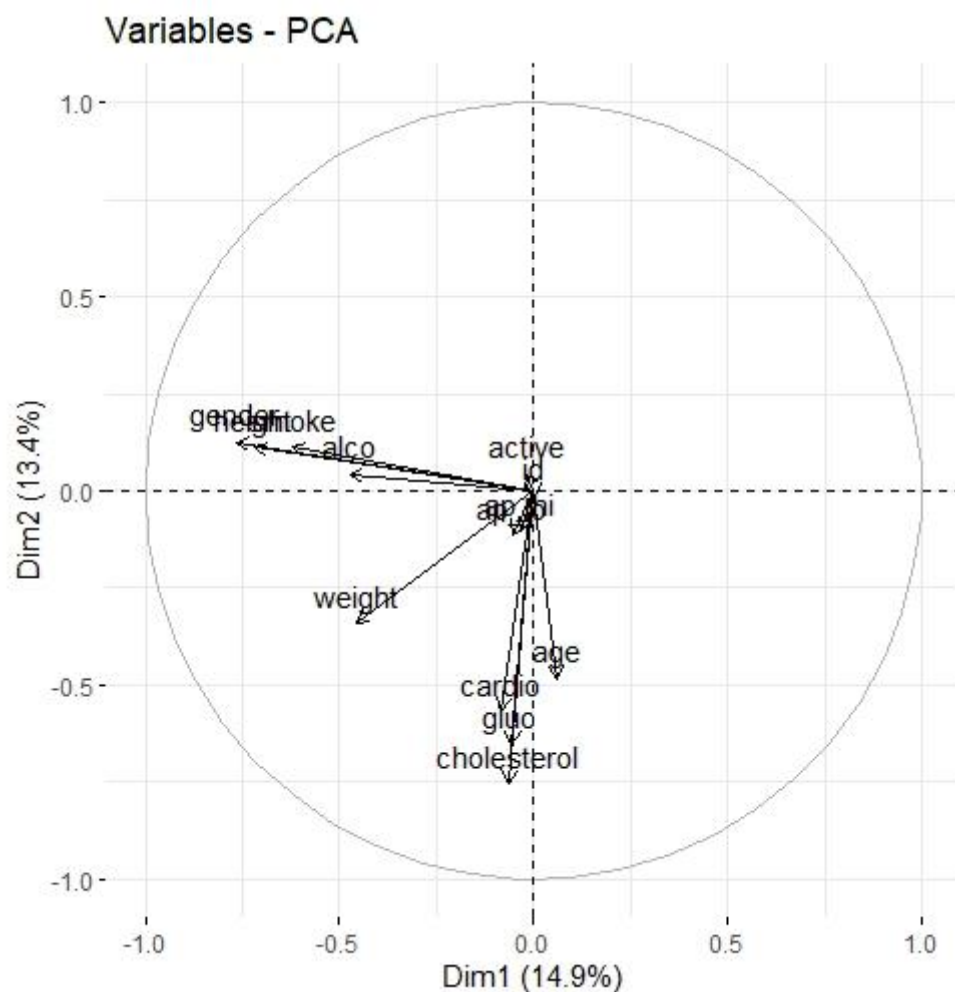
```
> summary(data.pca)
Importance of components:
      PC1      PC2      PC3      PC4      PC5      PC6      PC7      PC8      PC9      PC10     PC11
Standard deviation  1.3939  1.3185  1.07779  1.02688  1.00108  0.99888  0.9920  0.98647  0.90943  0.85061  0.78837
Proportion of Variance 0.1495 0.1337 0.08936 0.08111 0.07709 0.07675 0.0757 0.07486 0.06362 0.05566 0.04781
Cumulative Proportion 0.1495 0.2832 0.37255 0.45366 0.53075 0.60750 0.6832 0.75806 0.82168 0.87734 0.92515
      PC12     PC13
Standard deviation  0.72411 0.66990
Proportion of Variance 0.04033 0.03452
Cumulative Proportion 0.96548 1.00000
```

If we look at the plot of components and variance that they explain. In the analysis, the all component will be taken into consideration. (Mainly, 3 component)

### Component Analysis

Clare visoin of issues

```
fviz_pca_var(data.pca, col.ind = "Age")
```



## Conclusion

To conclude the project, analysing current data, it gives us more understanding about most cases in cardiovascular disease.

## References

Source information; /kaggle/input/cardio\_train.csv

Movite based on: [RPubs - Dimension Reduction for nominal data and qualitative data](#)