

# **Analysis on Breast Cancer and Possible Factors**

## **Multivariate Analysis Term Project**

**Ekinsu ÇİÇEK Mehmet Ali ERKAN**

# OUTLINE

- Introduction to the dataset
- Analysis on Research Question 1
- Analysis on Research Question 2
- Analysis on Research Question 3
- Conclusion

# Breast Cancer Data Set

	Age	BMI	Glucose	Insulin	HOMA	Leptin	Adiponectin	Resistin	MCP.1	Classification
1	48	23.50	70	2.70	0.46	8.80	9.70	7.99	417.11	1
2	83	20.69	92	3.11	0.70	8.84	5.42	4.06	468.78	1
3	82	23.12	91	4.49	1.00	17.93	22.43	9.27	554.69	1
4	68	21.36	77	3.22	0.61	9.88	7.16	12.76	928.22	1
5	86	21.11	92	3.5	0.80	6.69	4.81	10.57	773.92	1
6	49	22.85	92	3.2	0.73	6.83	13.67	10.31	530.41	1

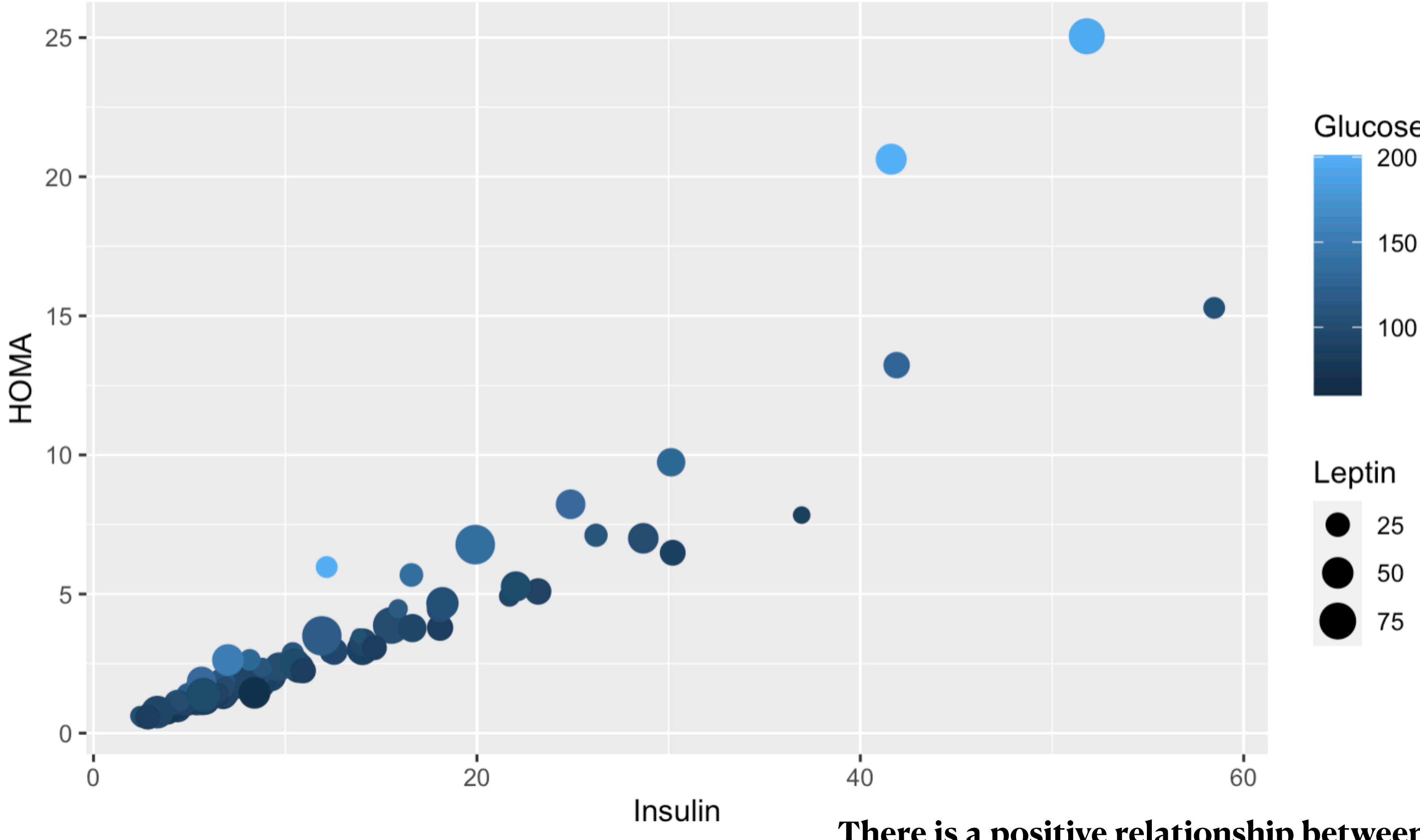
- 116 observation with 10 variables
- 36 outliers in the data.

<b>Variable</b>	<b>Information</b>	<b>Type</b>
<b>Age</b>	years	integer
<b>BMI</b>	kg/m2	numeric
<b>Glucose</b>	mg/dL	integer
<b>Insulin</b>	μU/ml	numeric
<b>HOMA</b>		numeric
<b>Leptin</b>	ng/mL	numeric
<b>Adiponectin</b>	μg/mL	numeric
<b>Resistin</b>	ng/mL	numeric
<b>MCP.1</b>	pg/dL	numeric
<b>Classification</b>	1 = healthy 2 = patients	Factor levels('1','2')

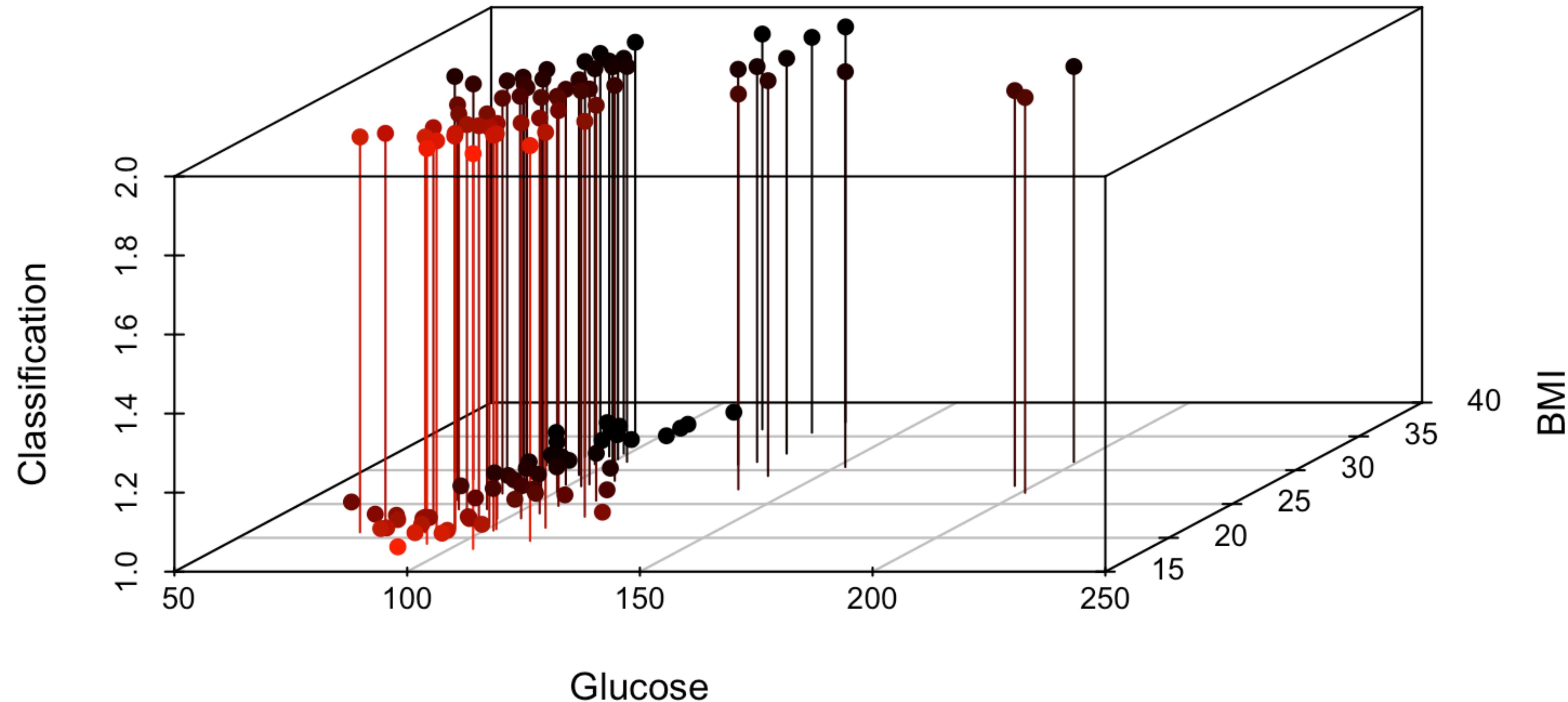
	<i>Age</i>	<i>BMI</i>	<i>Glucose</i>	<i>Insulin</i>	<i>HOMA</i>	<i>Leptin</i>	<i>Adiponectin</i>	<i>Resistin</i>	<i>MCP.1</i>	<i>Classification</i>
<i>Min</i>	24.0	18.37	60.00	2.43	0.46	4.31	1.65	3.21	45.84	1.00
<i>1st Qu.</i>	45.0	22.97	85.75	4.35	0.91	12.31	5.47	6.88	269.98	1.00
<i>Median</i>	56.0	27.66	92.00	5.92	1.38	20.27	8.35	10.82	471.32	2.00
<i>Mean</i>	57.3	27.58	97.79	10.01	2.69	26.61	10.18	14.72	534.65	1.55
<i>3rd Qu.</i>	71.0	31.24	102.00	11.18	2.85	37.37	11.81	17.75	700.09	2.00
<i>Max.</i>	89.0	38.58	201.00	58.46	25.05	90.28	38.04	82.10	1698.44	2.00

# The Scatter Plot of Insulin, Homa, Leptin and Glucose

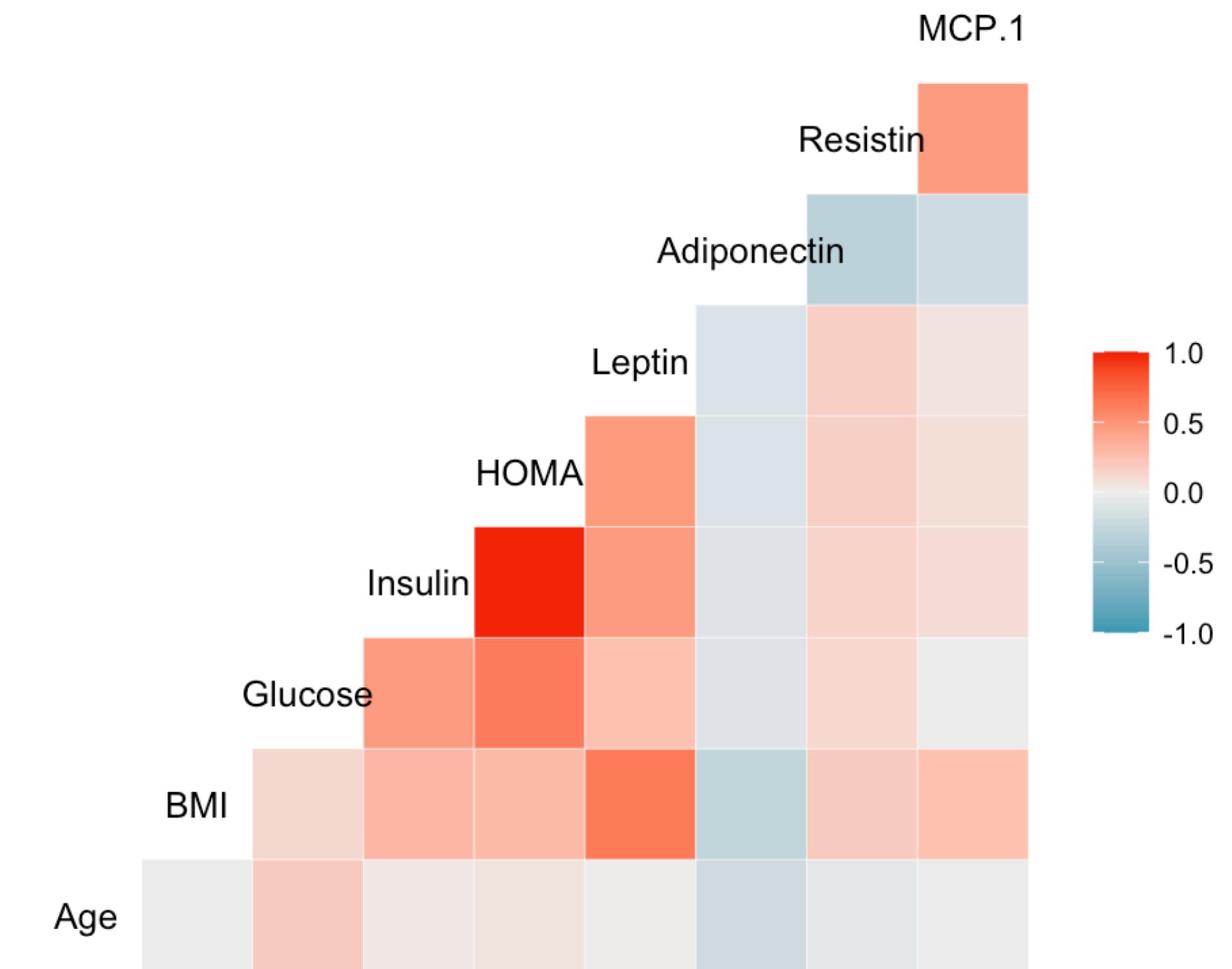
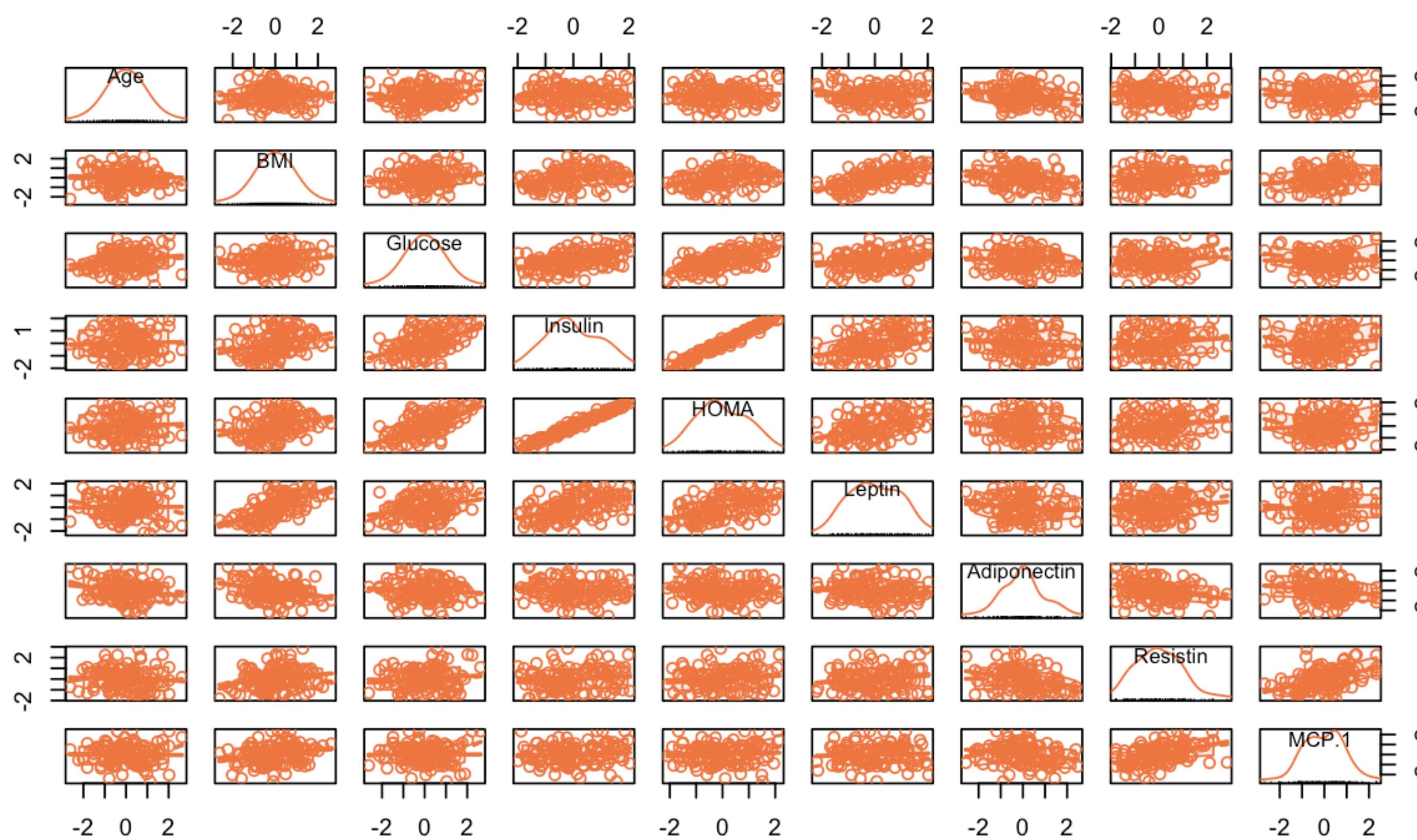
The Scatter Plot of Insulin, Homa, Leptin and Glucose



# Relationship between Glucose - BMI - Classification



# Is there a relationship between Insulin and HOMA level on blood and AGE, BMI and other hormones?



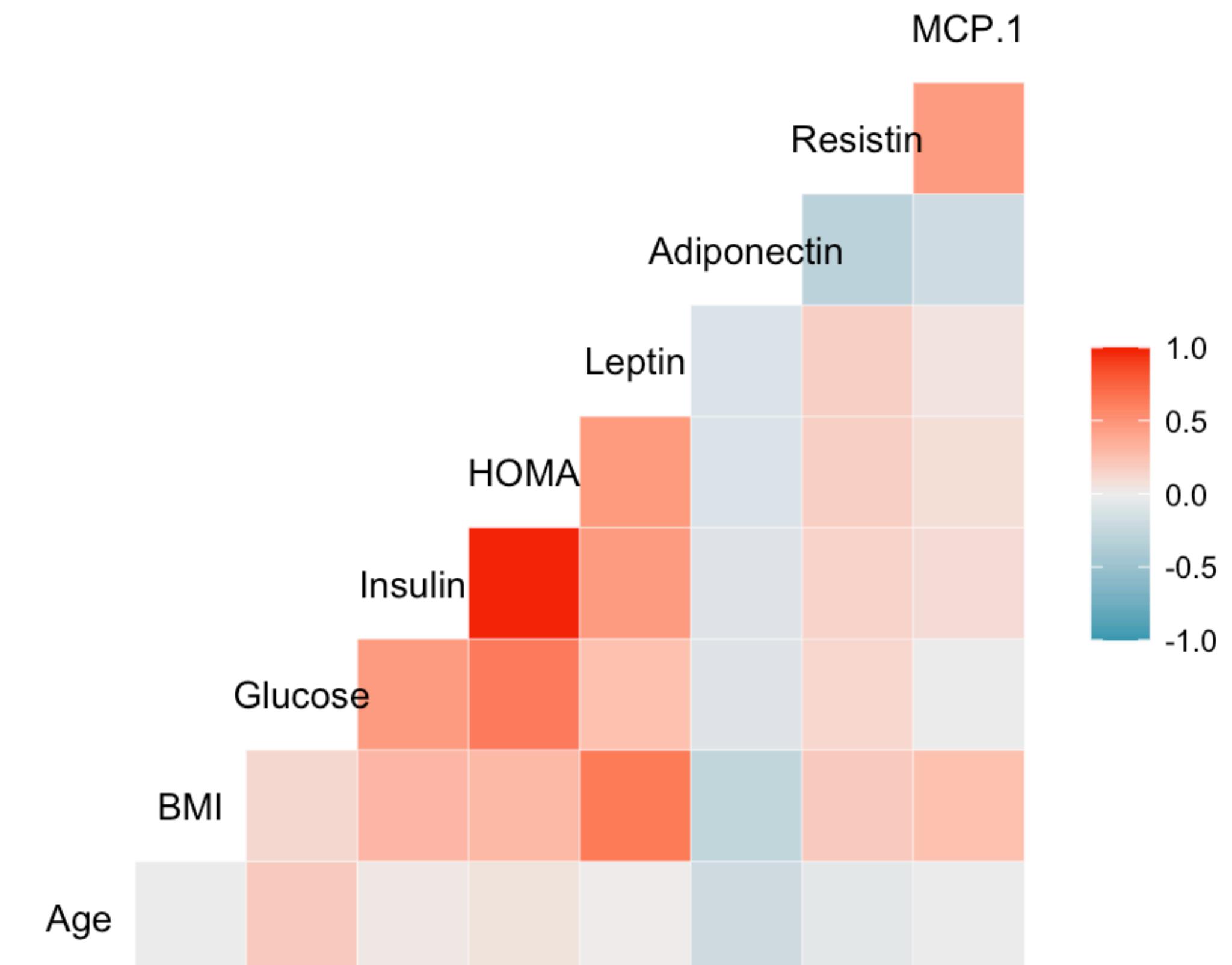
Our dimensions are 116 and 9, we applied principle component analysis before creating a model.

# Scaling

Age	BMI	Glucose	Insulin	HOMA	Leptin	Adiponectin	Resistin	MCP.1
-0.400	-0.532	-2.120	-1.714	-2.094	-1.219	0.2235	-0.514	-0.141
1.594	-1.672	-0.076	-.1378	-1.211	-1.214	-0.728	-1.490	0.034

Since the scale of the variables are different, data was scaled before PCA.

Correlations of the variables seem good for principle component analysis.



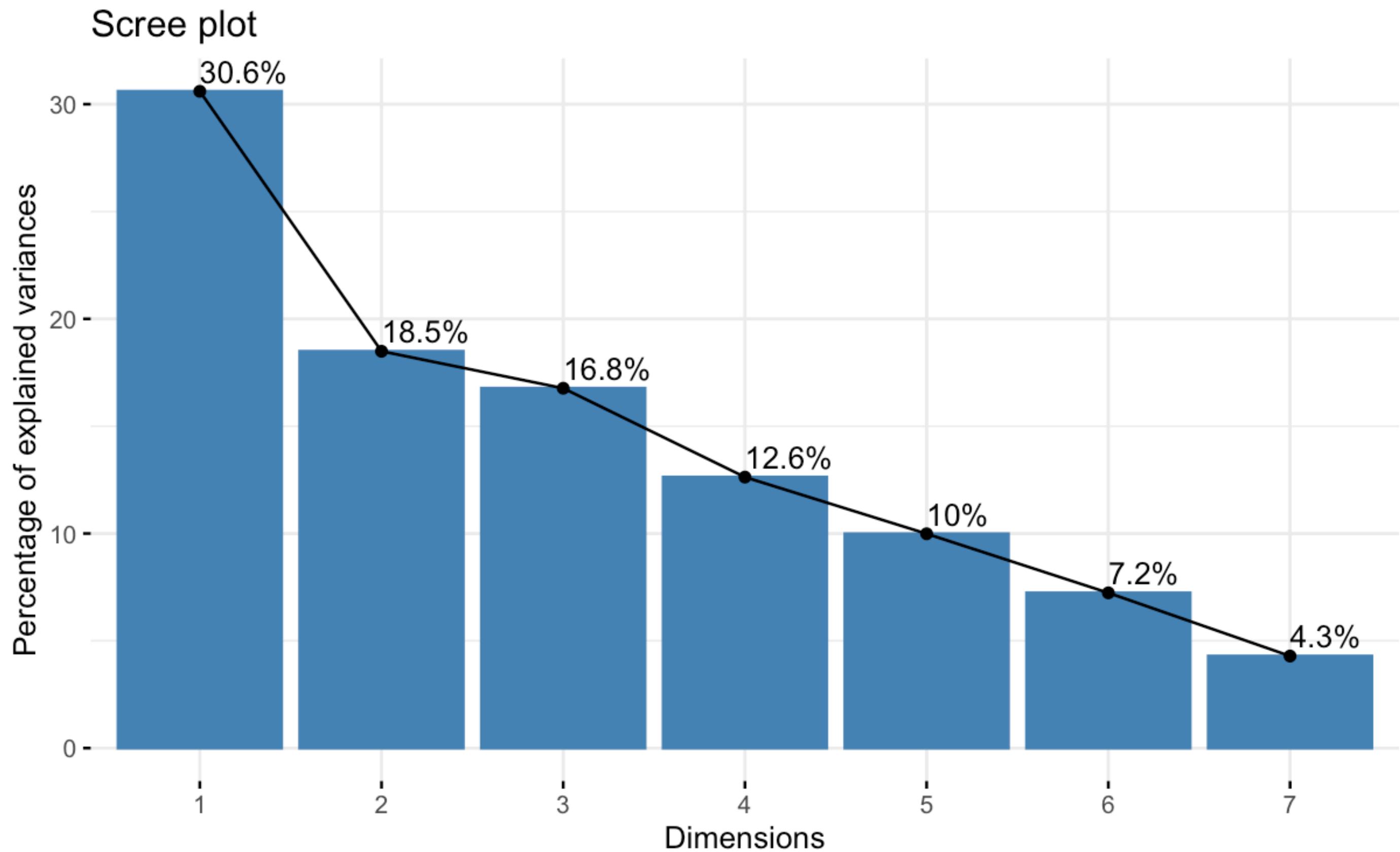
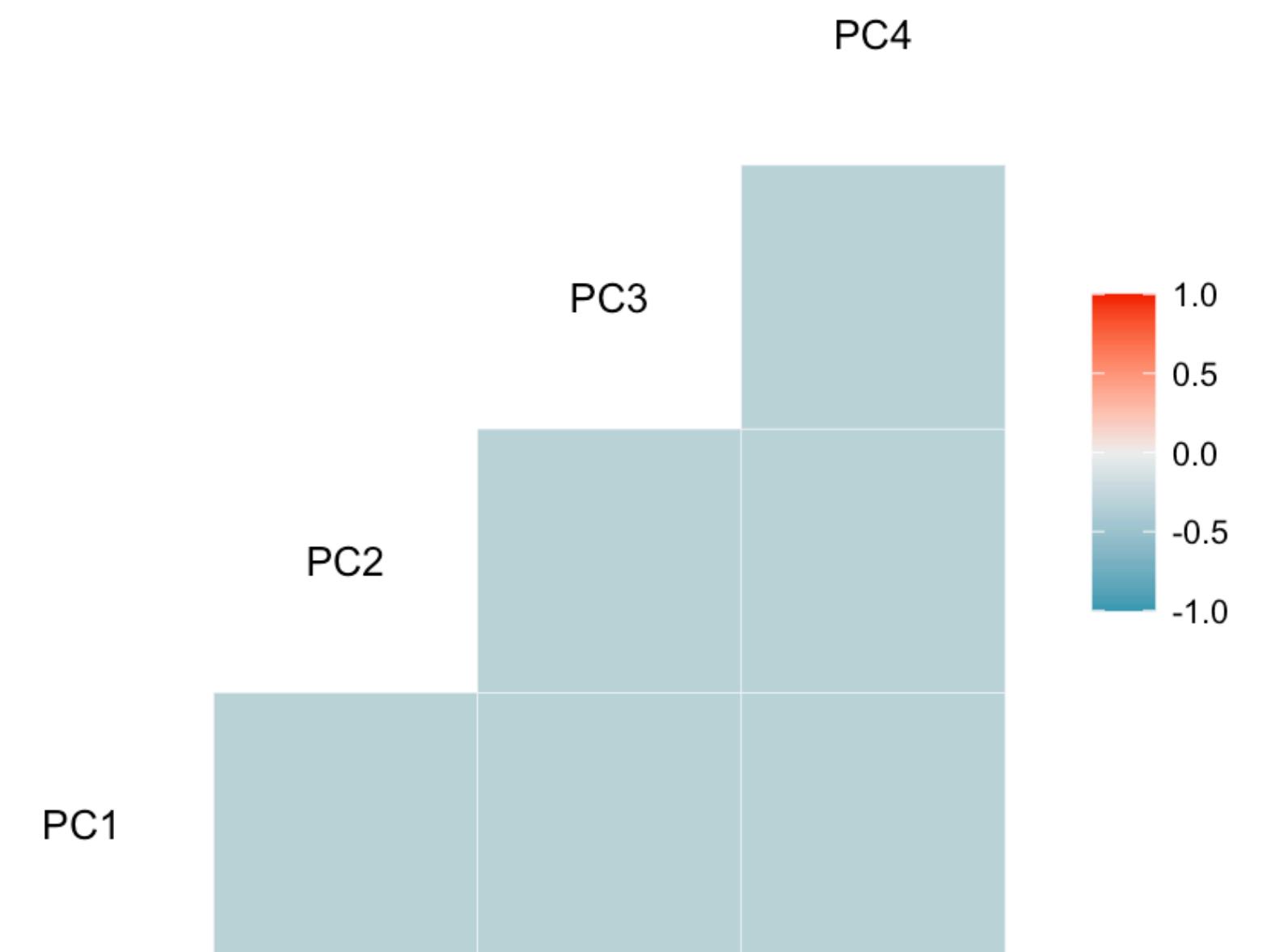
# Principal Components and Eigenvalues

Eigenvalues
1.4635859
1.1376775
1.0833787
0.9401669
0.8361664
0.7116267
0.5480844

The first four components nearly explains 80% variability of the data.

	PC1	PC2	PC3	PC4	PC5	PC6	PC7
Standard Deviation	1.464	1.1377	1.0834	0.9402	0.83617	0.71163	0.54808
Proportion Of Variance	0.306	0.1849	0.1677	0.1263	0.09988	0.07234	0.04291
Cumulative Proportion	0.306	0.4909	0.6586	0.7849	0.88474	0.95709	1.00000

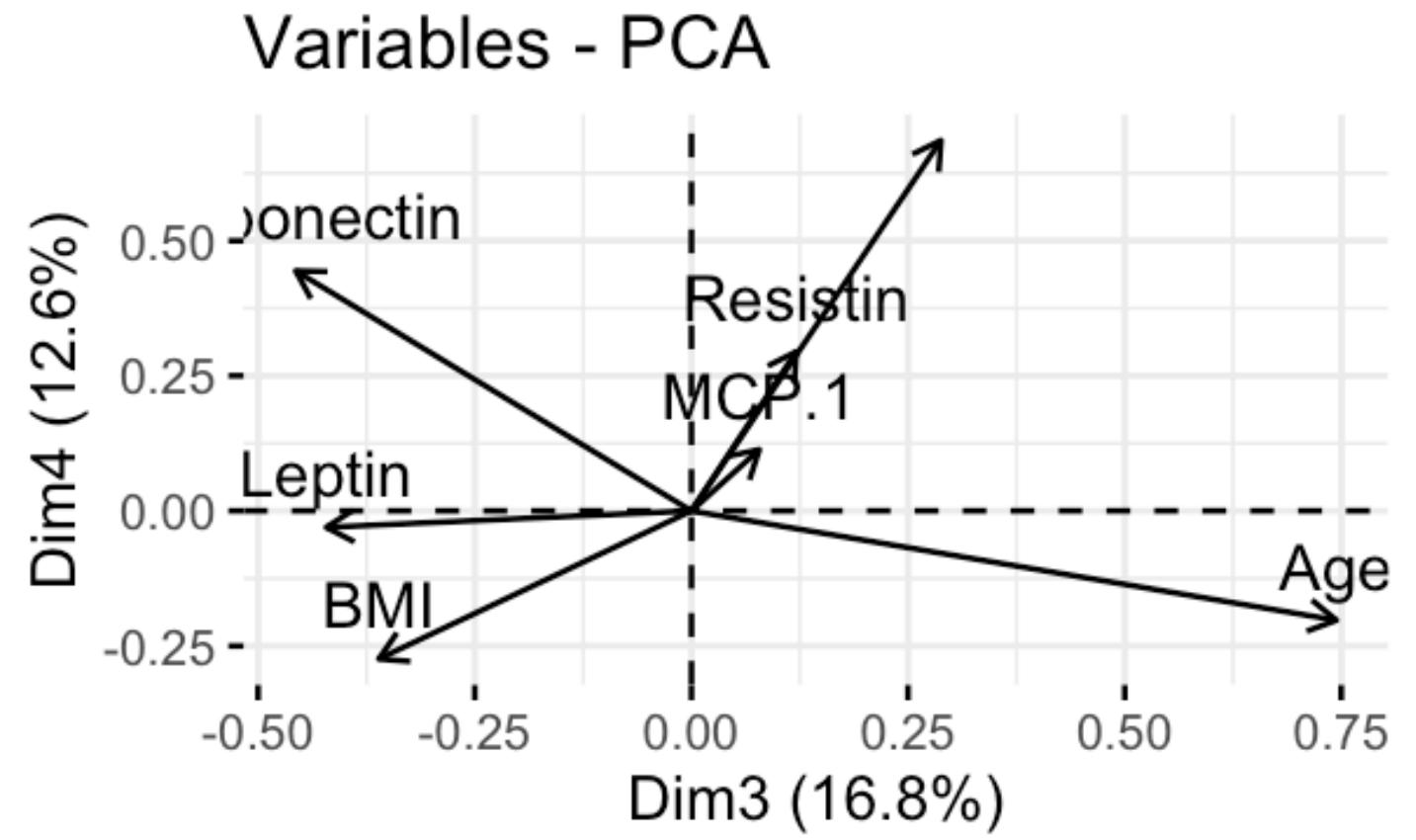
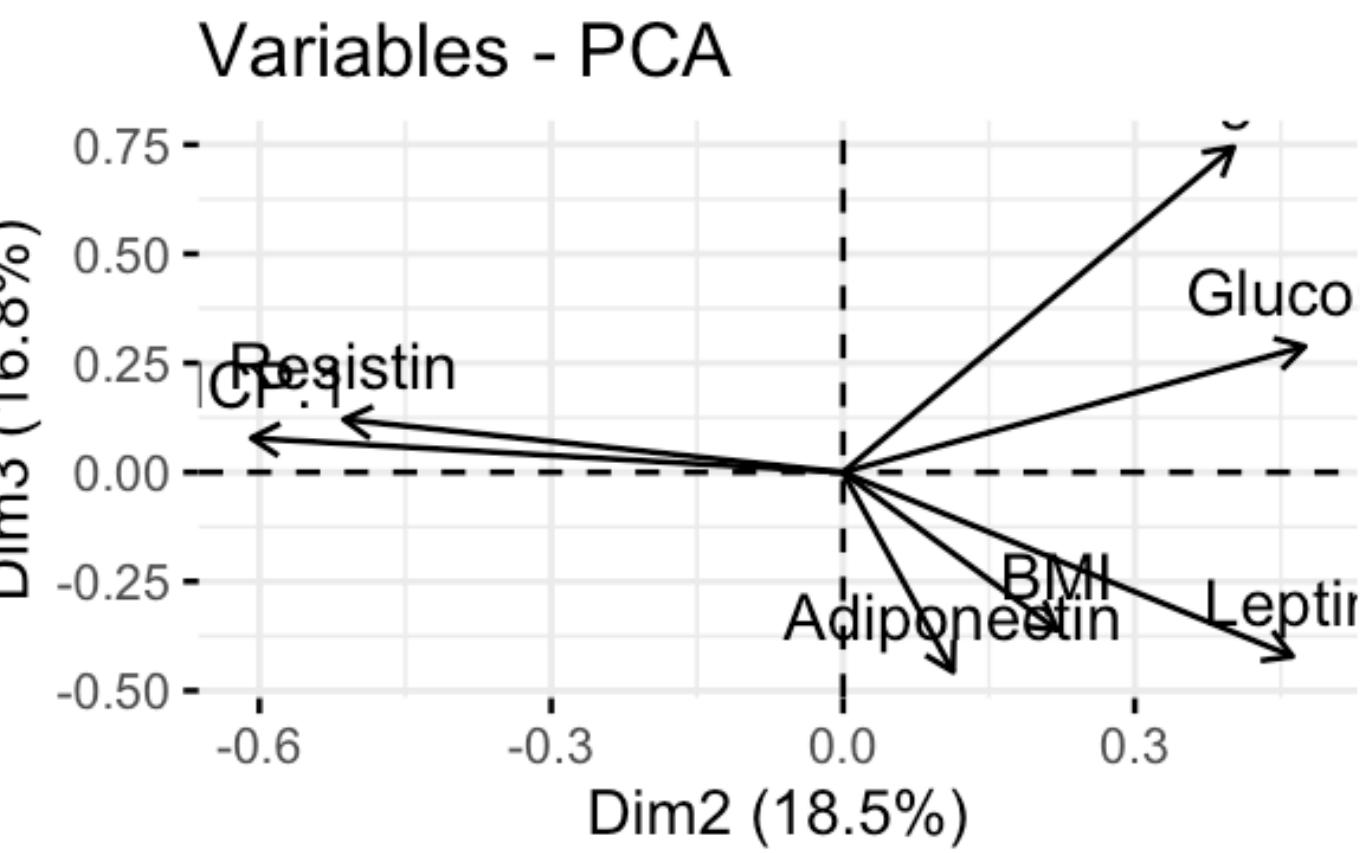
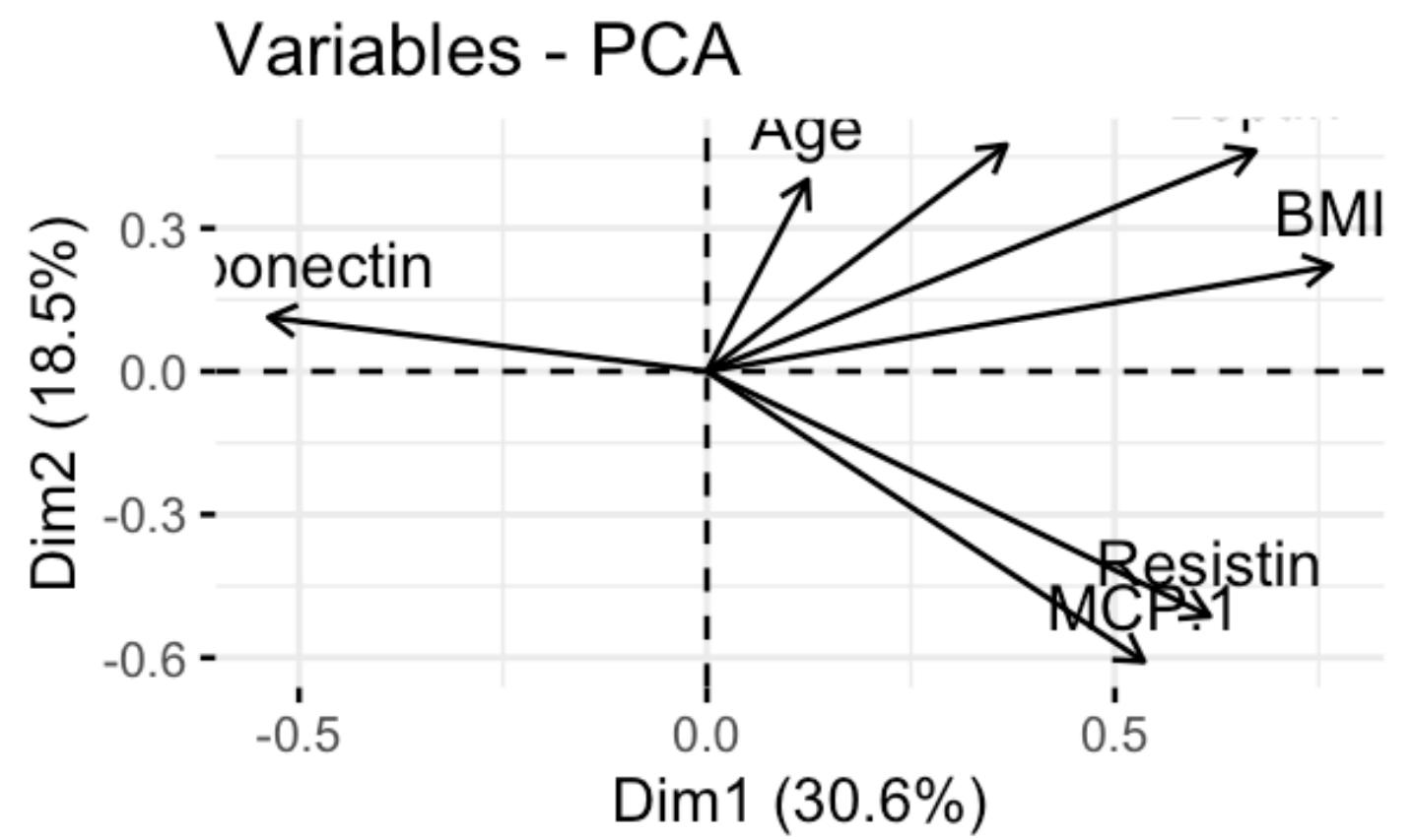
## Correlation Plot of Principal Components



Scree plot shows that first four component are suitable for explaining variability in the data.

Correlation of four components shows that they are linear independent.

	<b>PC1</b>	<b>PC2</b>	<b>PC3</b>	<b>PC4</b>
<b>Age</b>	0.1234575	0.4013896	0.74425187	-0.20239975
<b>BMI</b>	0.7653600	0.2202011	-0.36066473	-0.27393188
<b>Glucose</b>	0.3668419	0.4746081	0.28793469	0.68563259
<b>Leptin</b>	0.6719581	0.4620263	-0.42135976	-0.03041939
<b>Adiponectin</b>	-0.5371351	0.1131741	-0.45666194	0.44421208
<b>Resistin</b>	0.6161955	-0.5131530	0.12084092	0.29446678
<b>MCP.1</b>	0.5354946	-0.6081549	0.07827094	0.11338845



# The Relationship Between Components and Variables

# Principle Component Regression

$\text{Insulin} = -5.157\text{e-}17 + 2.894\text{e-}01 \text{PC1} + 2.546\text{e-}01 \text{PC2} - 6.256\text{e-}02 \text{PC3} + 2.698\text{e-}01 \text{PC4}$	
<b>Residual Standard Error : 0.8318</b>	<b>MSE : 0.6735252</b>
<b>Multiple R- Squared : 0.3322</b>	<b>Adjusted R- Squared : 0.3082</b>
<b>F – Statistics : 13.81</b>	<b>P – Value: 5.588e-09</b>

$\text{HOMA} = 3.032\text{e-}16 + 3.128\text{e-}01 \text{PC1} + 3.113\text{e-}01 \text{PC2} - 5.496\text{e-}03 \text{PC3} + 3.962\text{e-}01 \text{PC4}$	
<b>Residual Standard Error : 0.8318</b>	<b>MSE : 0.5332795</b>
<b>Multiple R- Squared : 0.3322</b>	<b>Adjusted R- Squared : 0.3082</b>
<b>F – Statistics : 13.81</b>	<b>P – Value: 5.588e-09</b>

# Multivariate Linear Regression

$\text{Insulin} = -0.0001605 - 0.0386830 \text{Age} + 0.0125844 \text{BMI} + 0.3735138 \text{Glucose} + 0.3591940 \text{Leptin} - 0.0088957 \text{Adiponectin} + 0.0071172 \text{Resistin} + 0.0588322 \text{MCP.1}$	
<b>Residual Standard Error : 0.8329</b>	<b>MSE : 0.6597791</b>
<b>Multiple R- Squared : 0.3484</b>	<b>Adjusted R- Squared : 0.3062</b>
<b>F – Statistics : 8.251</b>	<b>P – Value: 5.894e-08</b>

$\text{HOMA} = -0.0002302 - 0.0450229 \text{Age} + 0.0196093 \text{BMI} + 0.5482072 \text{Glucose} + 0.3134629 \text{Leptin} - 0.0094384 \text{Adiponectin} + 0.0047811 \text{Resistin} + 0.0542865 \text{MCP.1}$	
<b>Residual Standard Error : 0.7323</b>	<b>MSE : 0.5131732</b>
<b>Multiple R- Squared : 0.4963</b>	<b>Adjusted R- Squared : 0.4637</b>
<b>F – Statistics : 15.2</b>	<b>P – Value: 1.046e-13</b>

All models are significant but R-Squared values of them very low. Also, mean standard error of the models are high.

# Glucose - BMI one sample T2-test

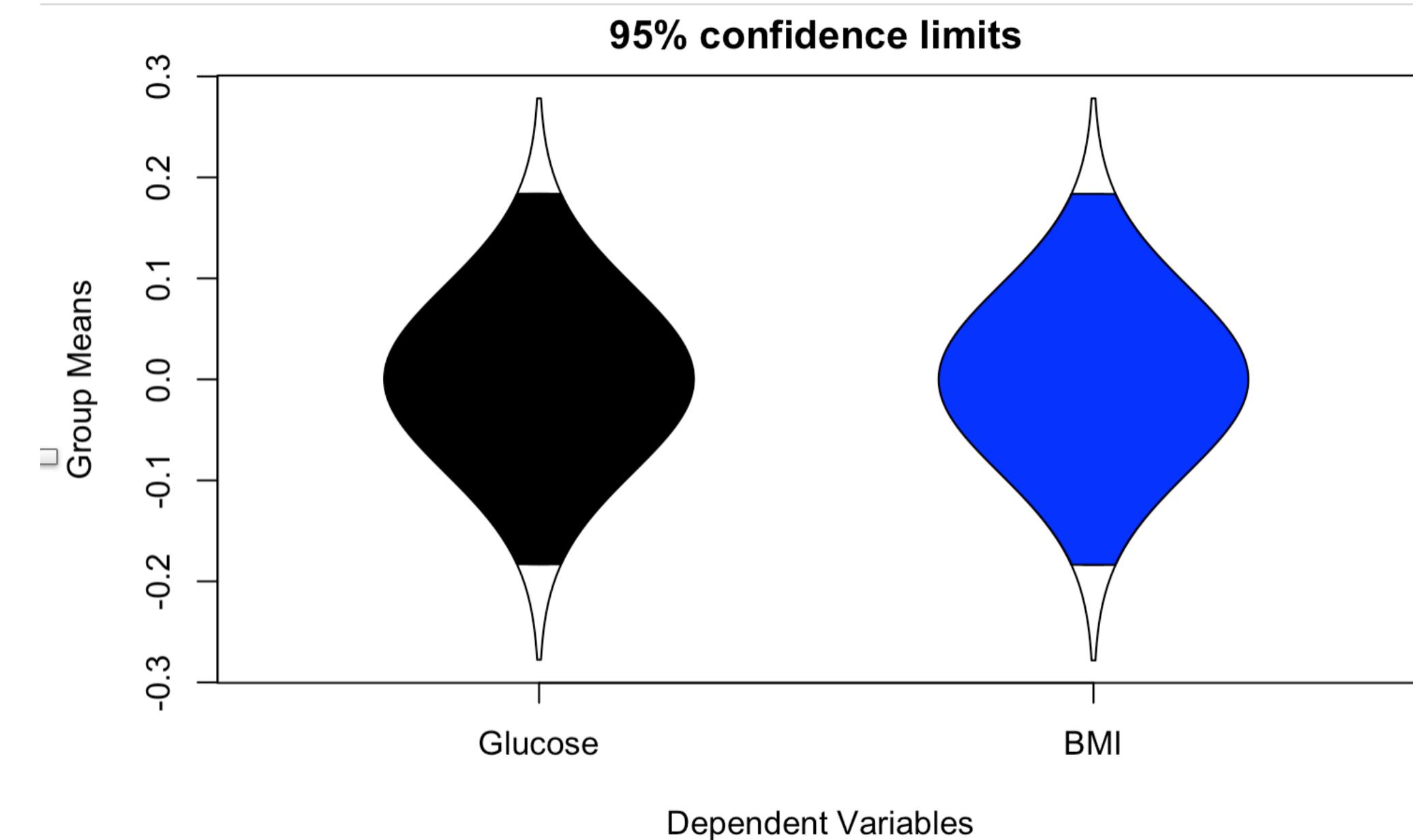
## Normality Check

Test	statistics	statistics	result
Mardia	0.59	0.96	yes
Skewness			
Mardia	-0.96	0.33	yes
Skewness			

Hotelling's one sample T2-test

data	Glucose-BMI
Value	T.2=4915.5, df1=2, df=114, p <u>&lt;=</u> 2e-16

## Response Matrix



# Does the classification have a significant effect on Glucose and BMI variables?

Factor 1	Factor 2
52	64

Box's M-test for Homogeneity of Covariance Matrices

data	Glucose-BMI
Value	Chi-Sq= 2,18, df=3, p = 0.53

*variance-covariance matrices are equal*

Classification	variable	statistics	p
1	BMI	0.99	0.96
1	Glucose	0.99	0.98
2	BMI	0.99	0.96
2	Glucose	0.99	0.94

*data follows univariate normality*

Hotelling's two sample T2-test

data	Glucose-BMI by Classification
Value	T.2= 17,308, df=2, p = 2. <u>77</u> e-07

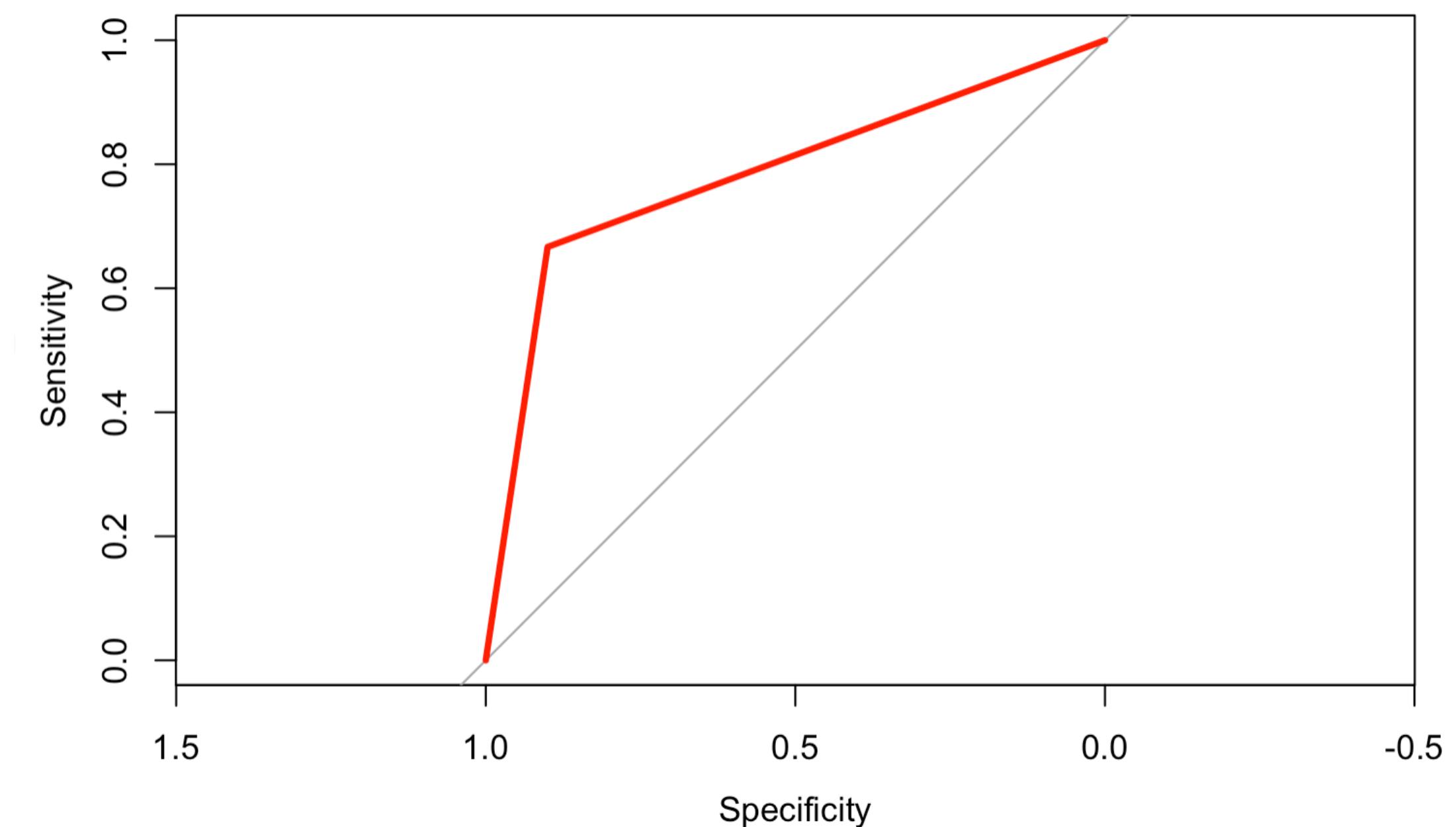
# How accurate is the classification for detecting the disease?

<b>Accuracy</b>	<b>0.77</b>
<b>95%CI</b>	<b>(0.54, 0.92)</b>
<b>No Information Rate</b>	<b>(0.54)</b>
<b>P value [Acc&gt;Nir]</b>	<b>(0.02)</b>
<b>Kappa</b>	<b>0.55</b>
<b>Mcnemar's Test P-Value</b>	<b>0.37</b>
<b>Sensitivity</b>	<b>0.90</b>
<b>Specificity</b>	<b>0.66</b>
<b>Pos Pred Value</b>	<b>0.69</b>
<b>Neg Pred Value</b>	<b>0.88</b>
<b>Prevalence</b>	<b>0.45</b>
<b>Detection Rate</b>	<b>0.40</b>
<b>Detection Prevalance</b>	<b>0.59</b>
<b>Balanced Accuracy</b>	<b>0.78</b>

## Confusion Matrix

		Reference	
		1	2
Prediction	1	9	4
	2	1	8

## ROC CURVE



# Conclusion

