

Tutorial #1

Table of Contents

Materials & Methods	3
Descriptive analysis	3
Machine learning	3
Results 1/2	4
Preliminary statistics	4
Statistical tests	25
Correlations	30
Individual performances	44
Results 2/2	48
Predictive models	48

Descriptive analysis

The tabular dataset was composed of **32** individuals and **7** variables including **6** biomarkers and a response variable called **Target** (**section Preliminary statistics - Global statistics**). Biomarkers were expressed using many settings such as median (and interquartile range) or mean (and standard deviation). Shapiro-Wilk's test was used to evaluate biomarker normality distribution (**section Preliminary statistics - Statistics by group**). Principal Component Analysis (PCA) was used to evaluate relationship between all continuous variables (**section Variable Contribution (PCA)**). Subjects were divided into **2** groups according to **Target** and compared using non-exhaustive list of tests (**section Statistical tests**). Continuous relationship was assessed between all biomarkers using Spearman rank correlation coefficient. Biserual point correlation was performed to evaluate correlation between continuous and categorical variables. Spearman rank correlation coefficient was measures too. To categorical variables, relationship was measured using chi-squared test (**section Correlations**). Logistic regression was performed to evaluate biomarker individual performances. Youden or Closest topleft thresholds were used to maximize metrics such as AUC or Accuracy (**section Individual performances**).

A p-value < 0.05 was considered statistically significant.

Machine learning

For **Classification** problem, the tabular dataset was split in a training set (train) and a validation set (test) according to a **60:40** ratio. Rows containing missing valyes were removed from the training set. Any Biomarker combinations were computed and evaluated using simple machine learning algorithms such as **logistic regression, naive bayes, k-nearest neighbors, decision three, random forest** and **XGBoost**. Biomarker combinatory performances were indicated by **Accuracy** metric and compared to test set.

The current analysis report was fully generated using Geneseng app.

Preliminary statistics

Global statistics

Table 1: Summary statistics of variables

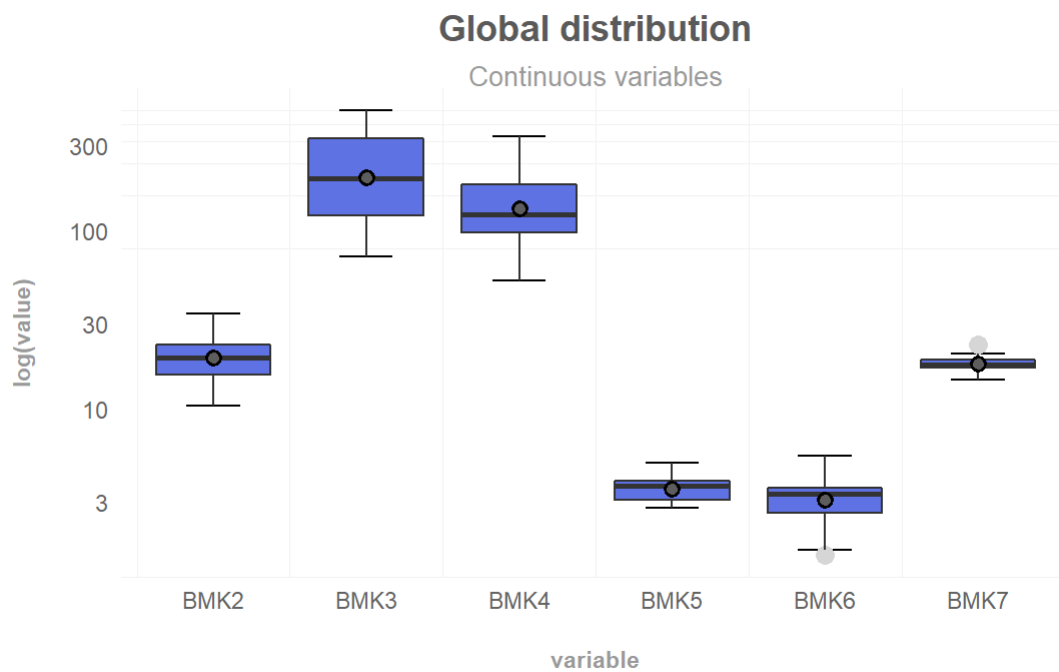
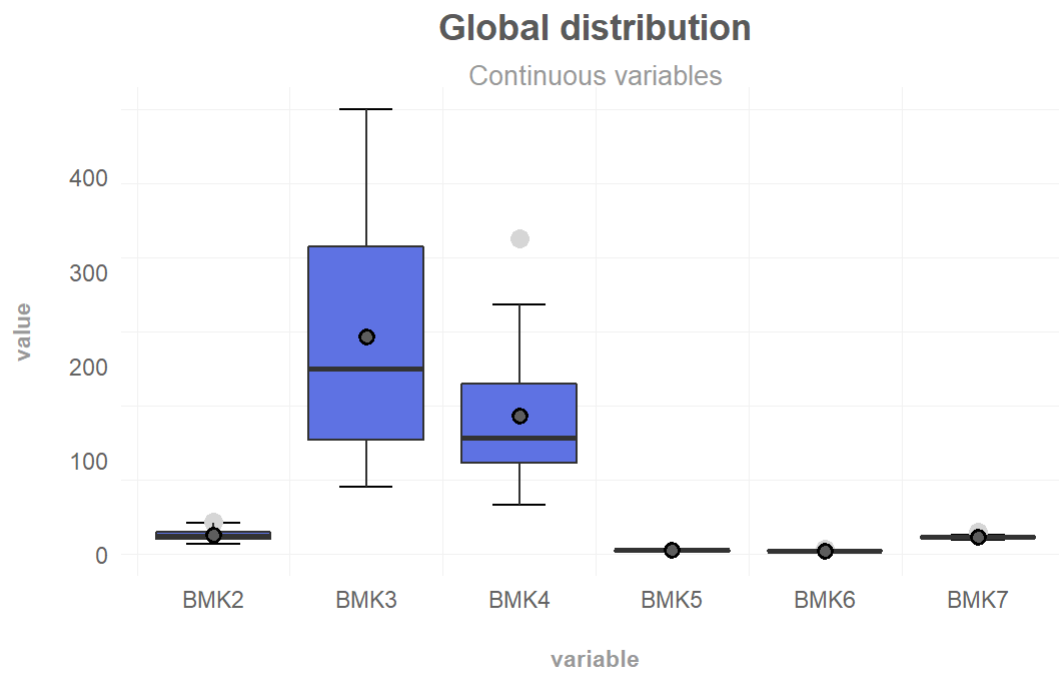
biomarker	n	n distinct	min	median	mean	sd	iqr	max	NA's ¹	Shapiro's test ²	normality
BMK2	32	25	10.40	19.20	20.09	6.03	7.38	33.90	0	1.23e-01	no
BMK3	32	27	71.10	196.30	230.72	123.94	205.18	472.00	0	2.08e-02	no
BMK4	32	22	52.00	123.00	146.69	68.56	83.50	335.00	0	4.88e-02	no
BMK5	32	22	2.76	3.70	3.60	0.53	0.84	4.93	0	1.10e-01	no
BMK6	32	29	1.51	3.33	3.22	0.98	1.03	5.42	0	9.27e-02	no
BMK7	32	30	14.50	17.71	17.85	1.79	2.01	22.90	0	5.94e-01	yes

¹Number of missing values

²Follow normal distribution if p-value > 0.05.

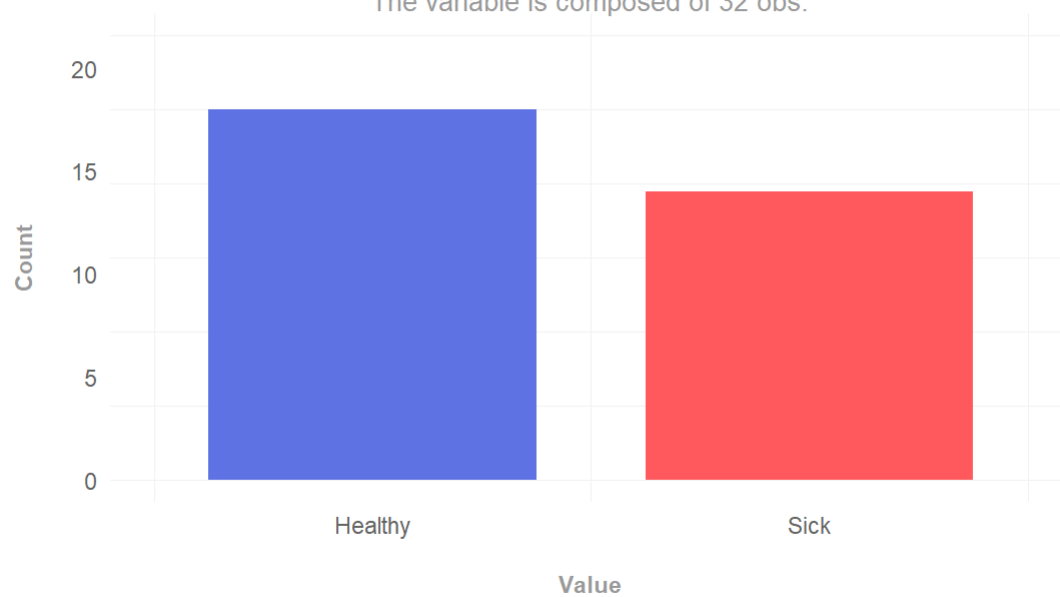
Table 2: Summary statistics of variables

biomarker	value	n	percent
Target	Healthy	18	0.56
Target	Sick	14	0.44



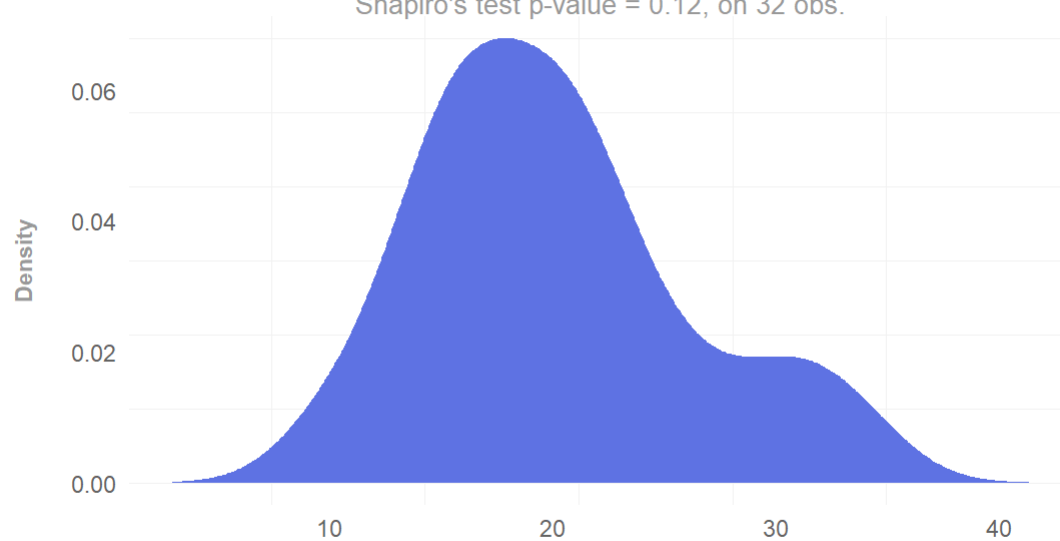
Distribution of Target variable

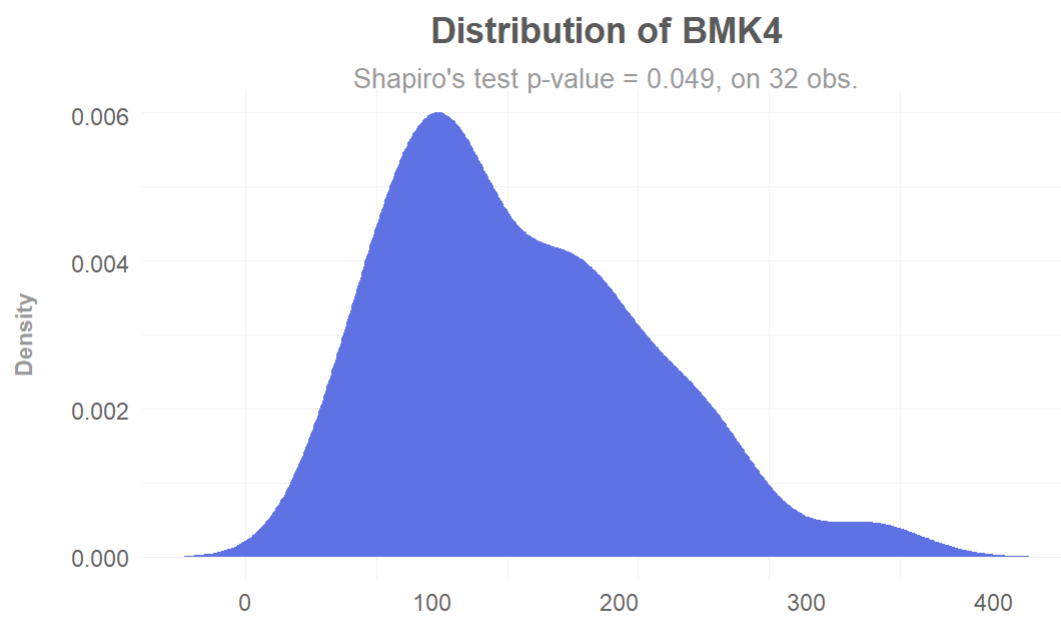
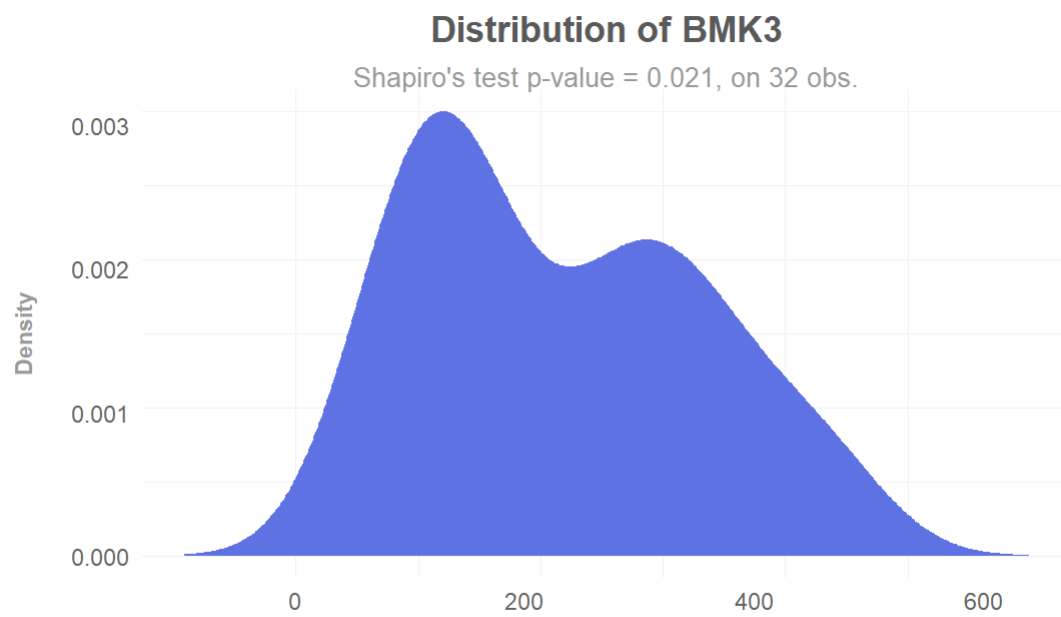
The variable is composed of 32 obs.

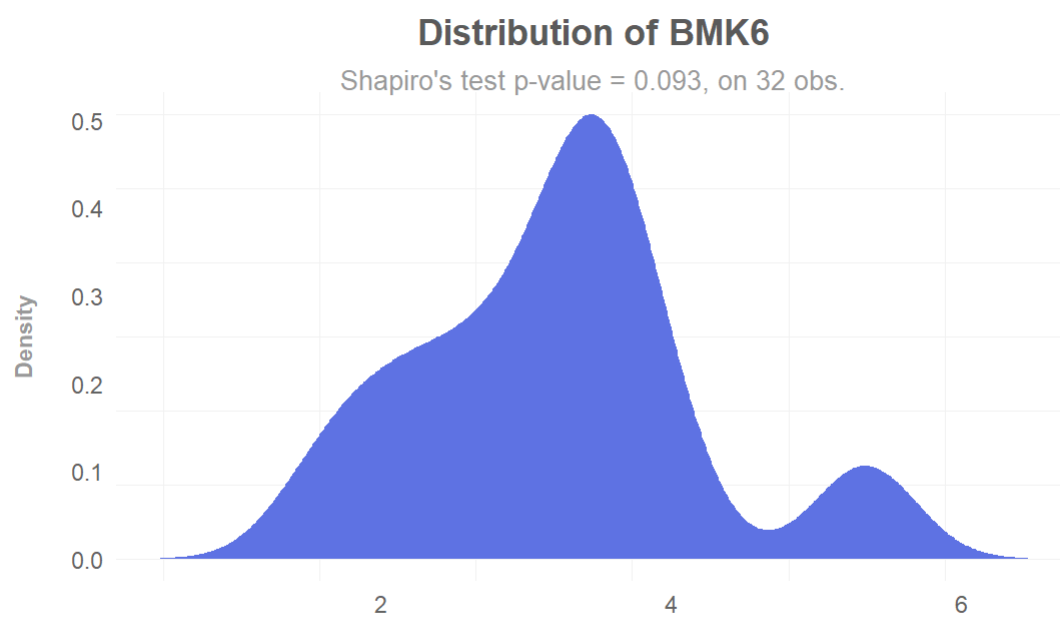
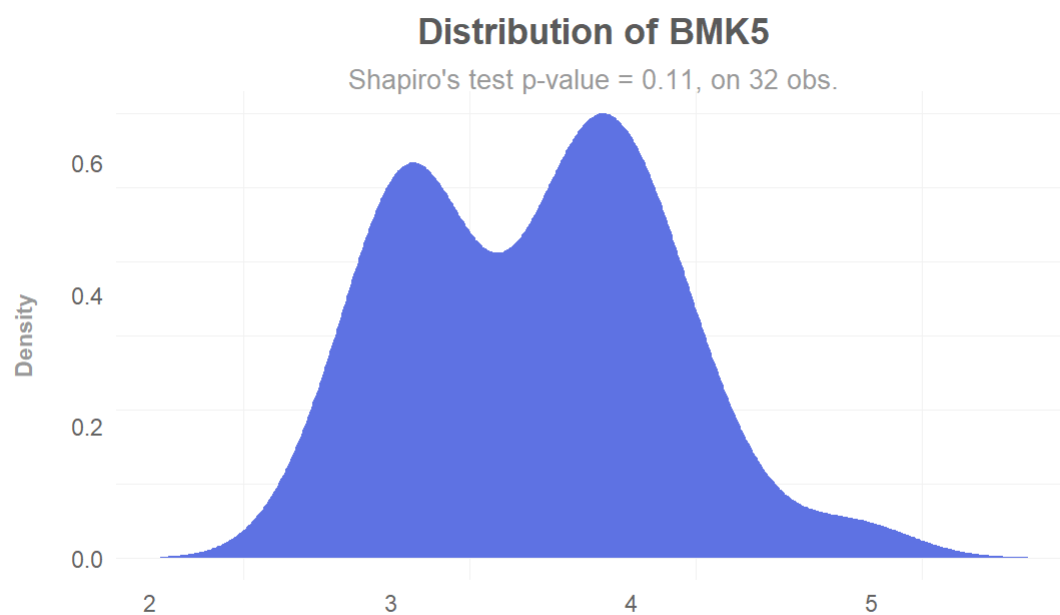


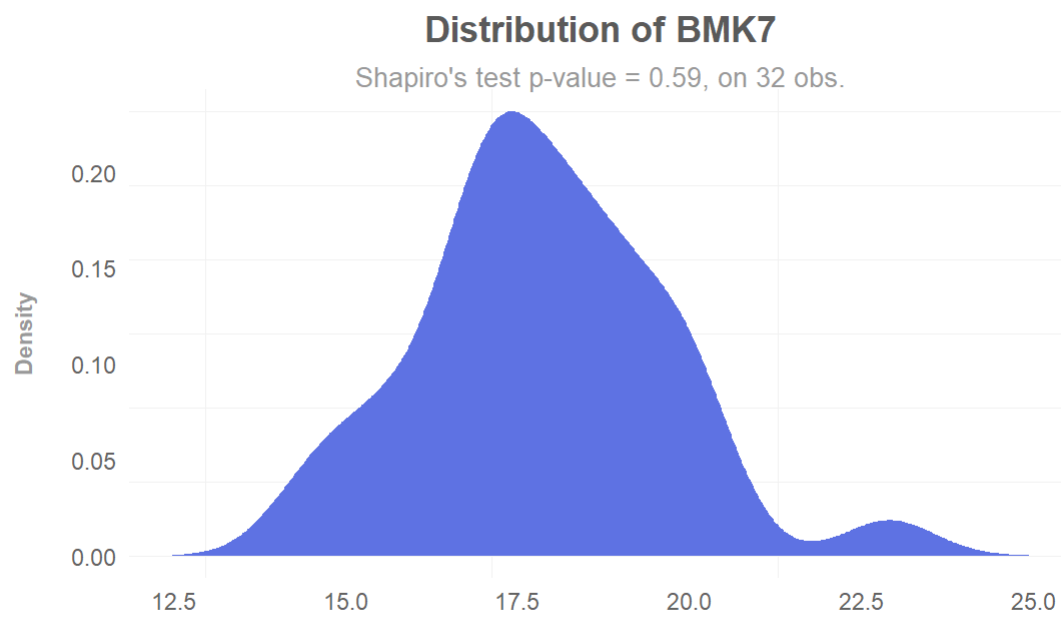
Distribution of BMK2

Shapiro's test p-value = 0.12, on 32 obs.









Statistics by group

Table 3: Summary statistics by group

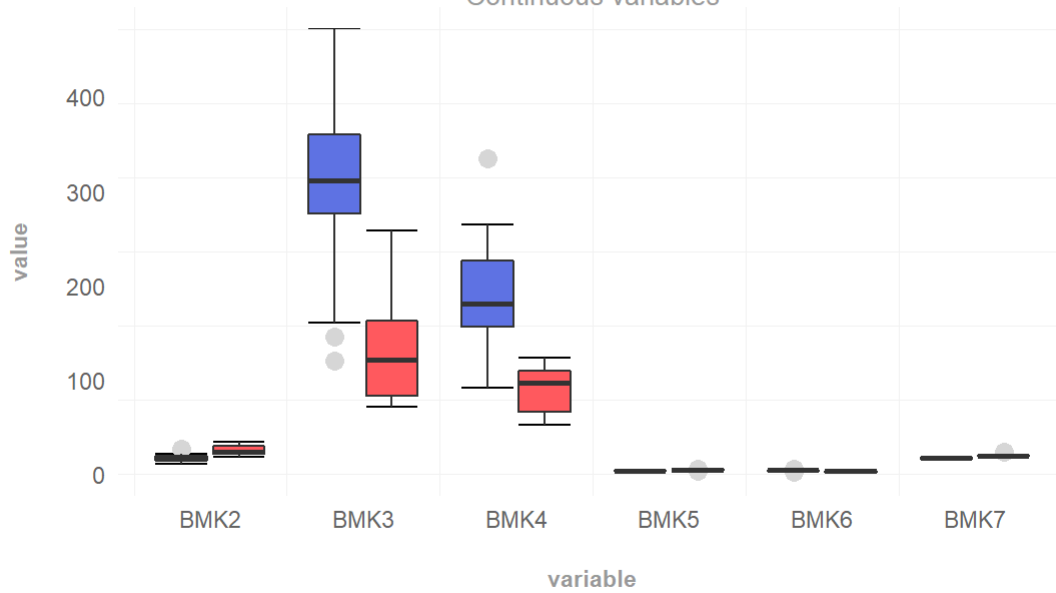
biomarker	group	n	n distinct	min	median	mean	sd	iqr	max	NA's ¹	Shapiro's test ²	normality
BMK2	Healthy	18	15	10.40	15.65	16.62	3.86	4.30	26.00	0	4.49e-01	no
BMK2	Sick	14	11	17.80	22.80	24.56	5.38	8.22	33.90	0	1.67e-01	no
BMK3	Healthy	18	14	120.30	311.00	307.15	106.77	84.20	472.00	0	2.88e-01	no
BMK3	Sick	14	13	71.10	120.55	132.46	56.89	79.35	258.00	0	1.03e-01	no
BMK4	Healthy	18	11	91.00	180.00	189.72	60.28	70.00	335.00	0	5.60e-01	yes
BMK4	Sick	14	12	52.00	96.00	91.36	24.42	43.75	123.00	0	1.10e-01	no
BMK5	Healthy	18	14	2.76	3.18	3.39	0.47	0.63	4.43	0	4.52e-02	no
BMK5	Sick	14	11	2.76	3.92	3.86	0.51	0.36	4.93	0	1.25e-01	no
BMK6	Healthy	18	17	2.14	3.57	3.69	0.90	0.61	5.42	0	9.43e-02	no
BMK6	Sick	14	13	1.51	2.62	2.61	0.72	1.21	3.46	0	1.08e-01	no
BMK7	Healthy	18	17	14.50	17.02	16.69	1.09	1.42	18.00	0	5.93e-02	no
BMK7	Sick	14	13	16.90	19.17	19.33	1.35	1.37	22.90	0	9.63e-02	no

¹Number of missing values

²Follow normal distribution if p-value > 0.05.

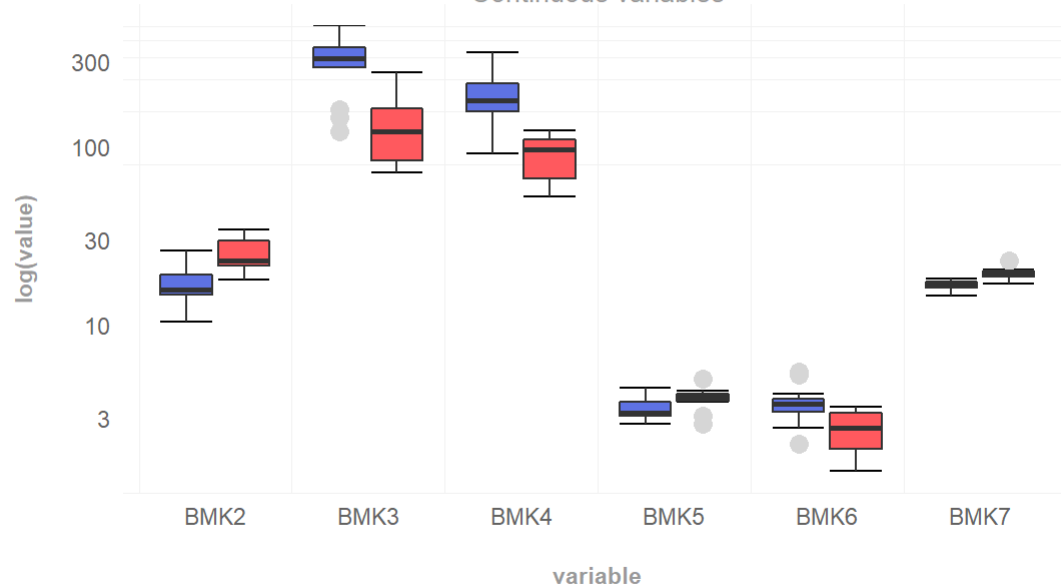
Global distribution

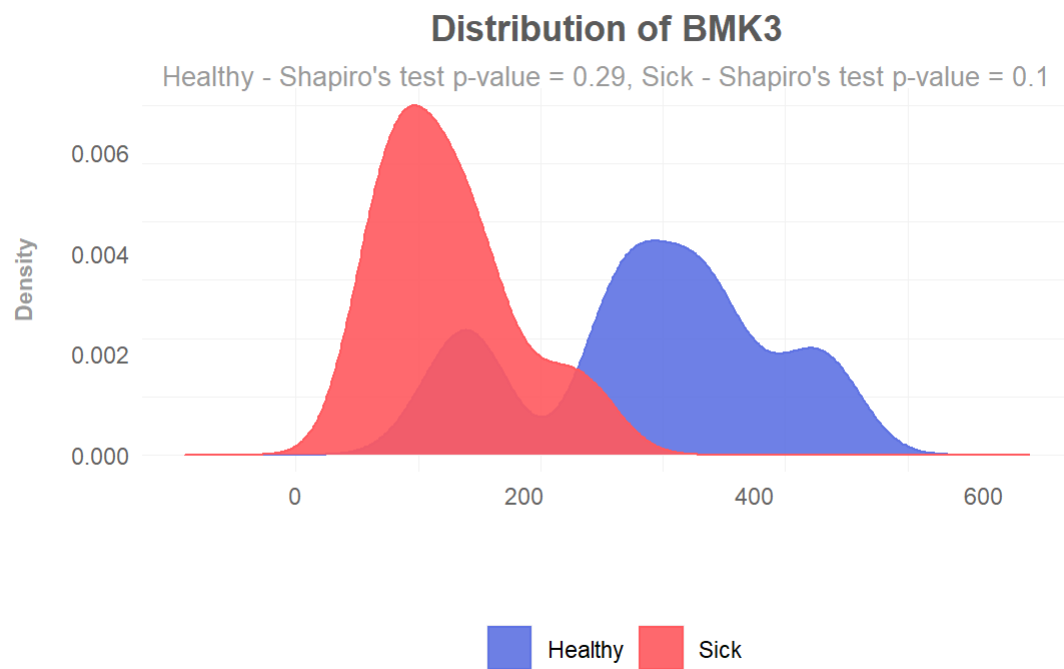
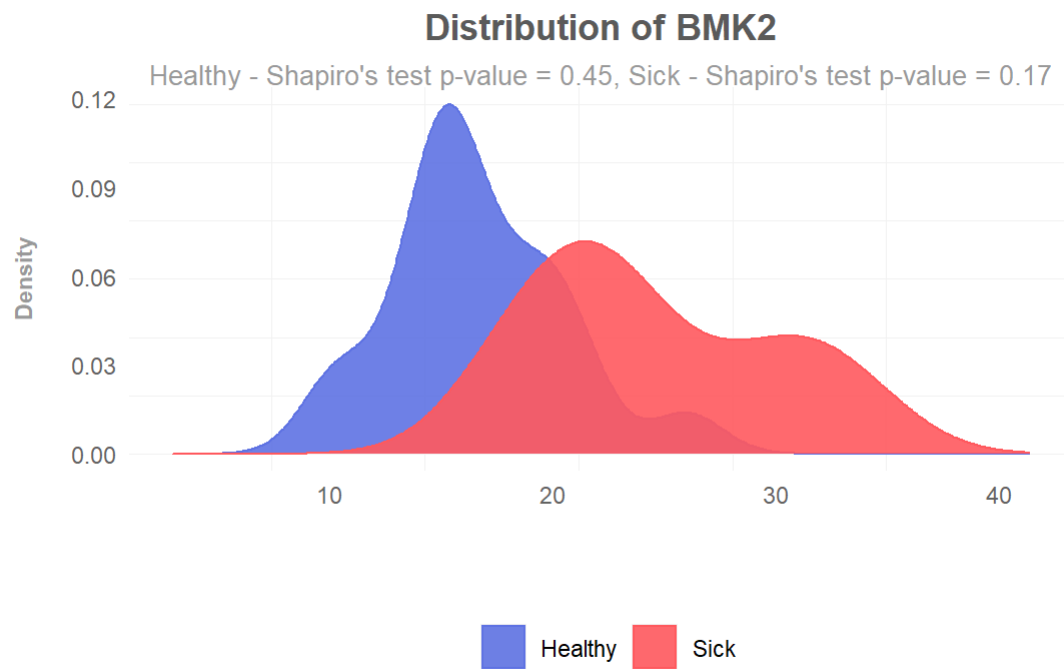
Continuous variables

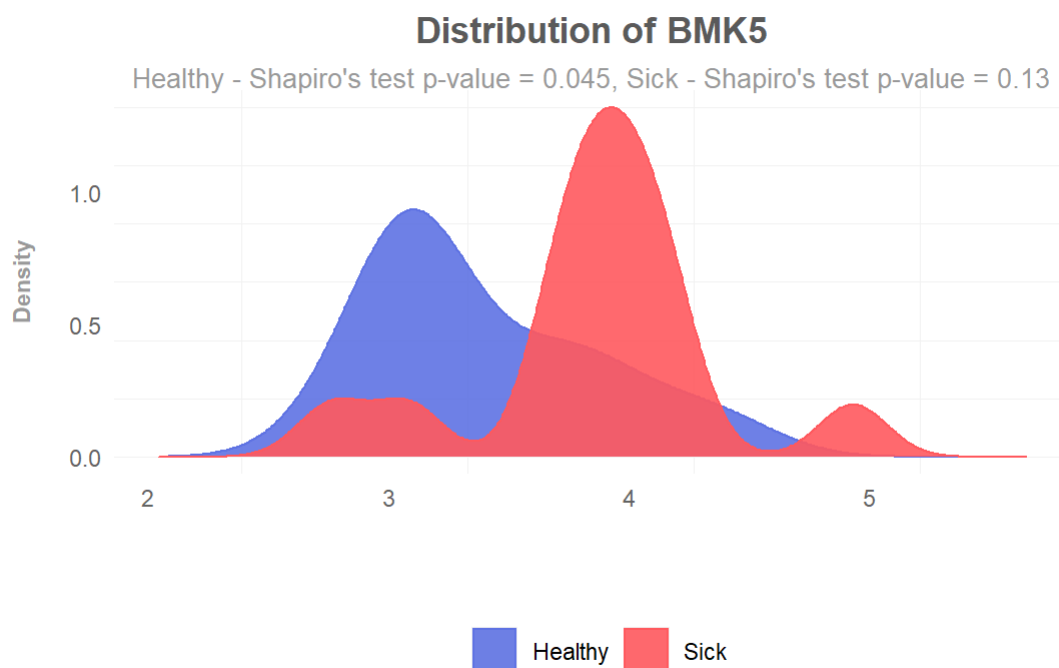
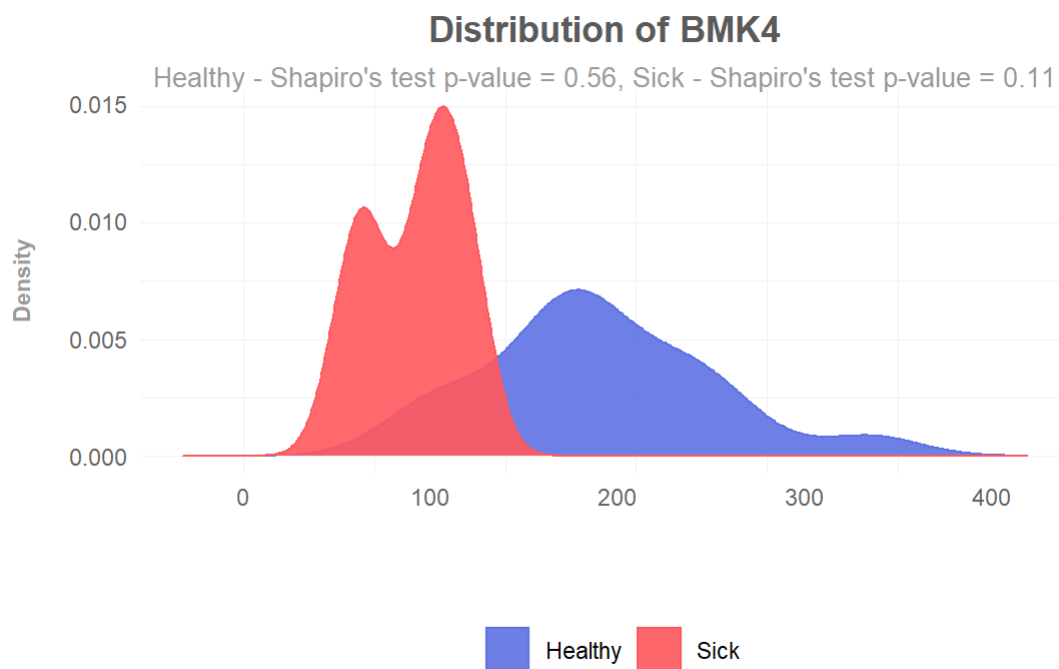


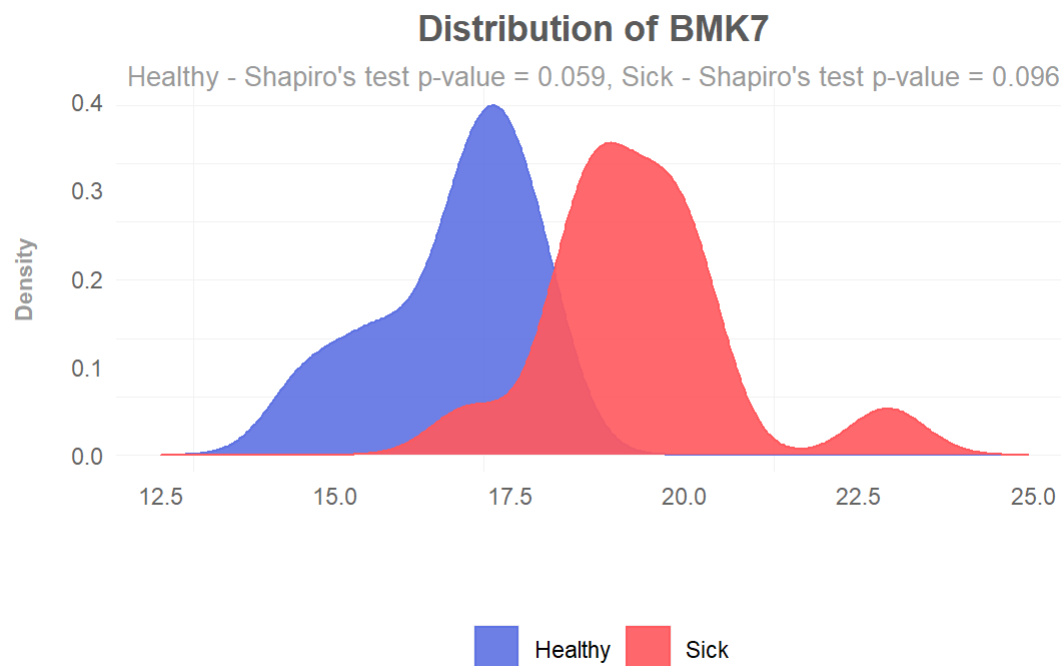
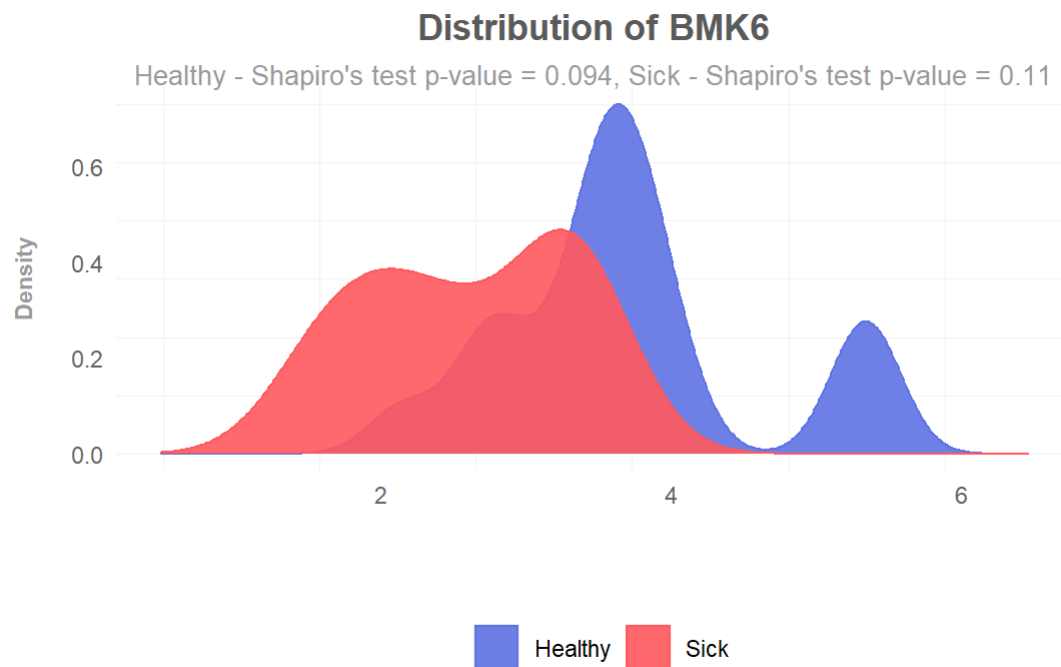
Global distribution

Continuous variables









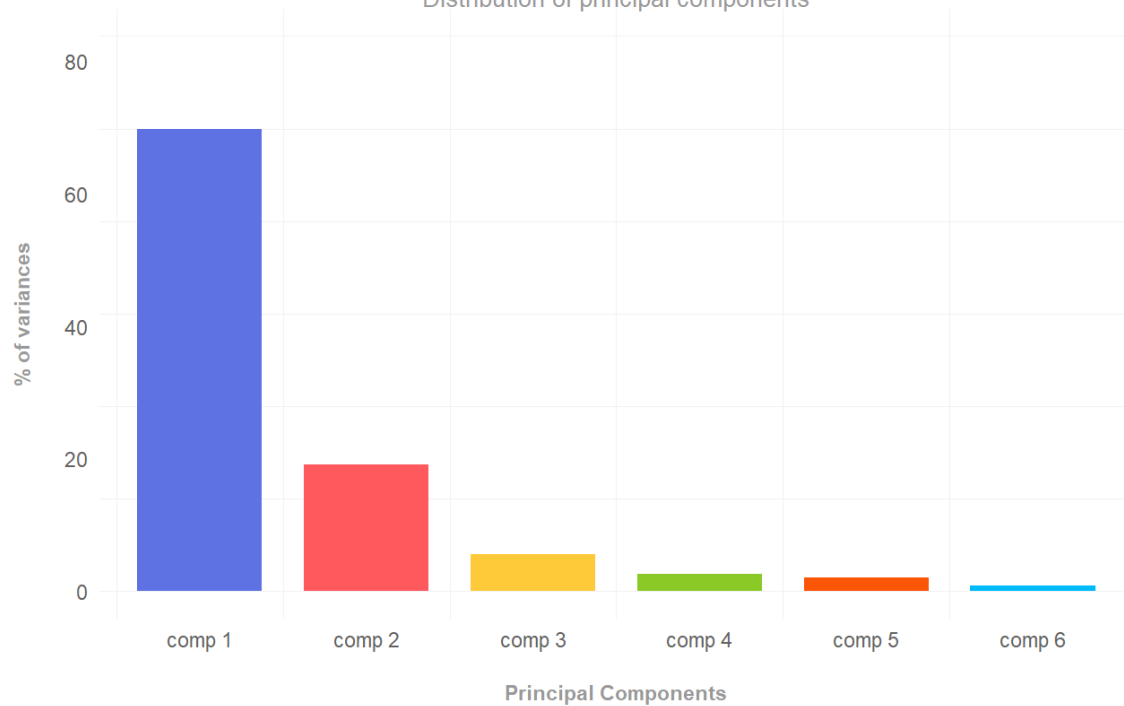
Variable Contribution (PCA)

Table 4: Principal Component Analysis - Variable contribution (%)

	Dim.1	Dim.2	Dim.3	Dim.4	Dim.5
BMK2	21.04	0.34	3.79	61.16	1.24
BMK3	21.72	0.37	0.94	36.00	8.68
BMK4	18.14	13.07	2.14	1.51	64.92
BMK5	13.48	19.06	64.08	0.05	2.07
BMK6	19.24	8.97	17.45	1.09	5.30
BMK7	6.39	58.19	11.60	0.18	17.80

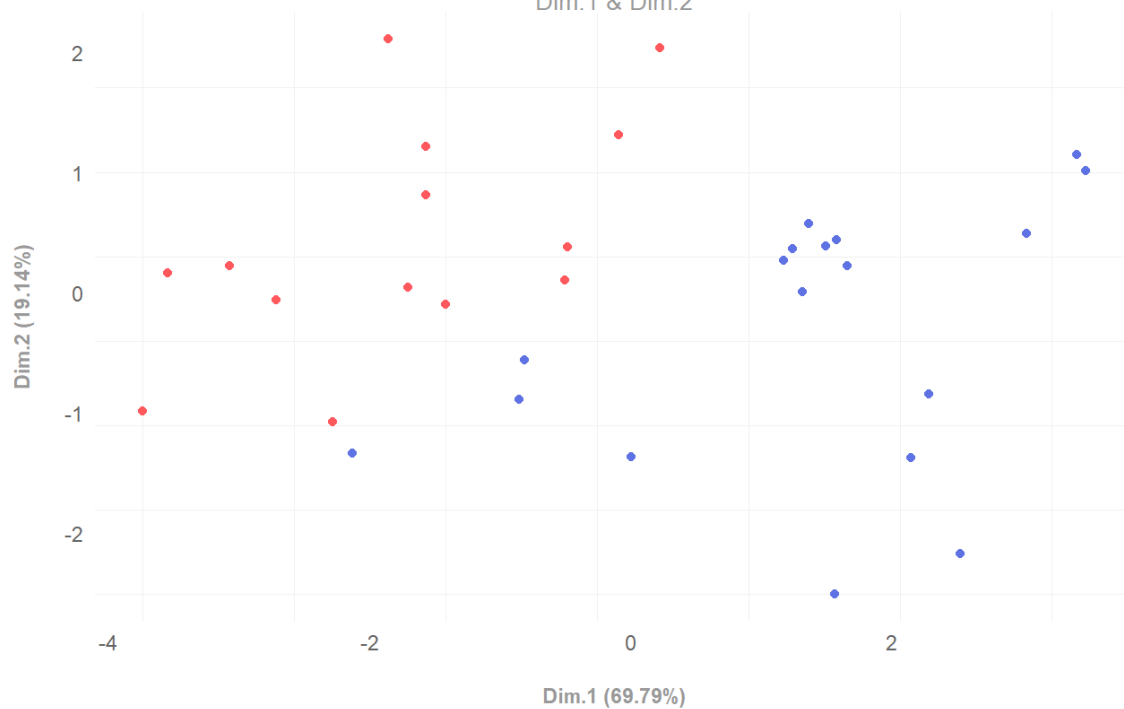
Explained Variances

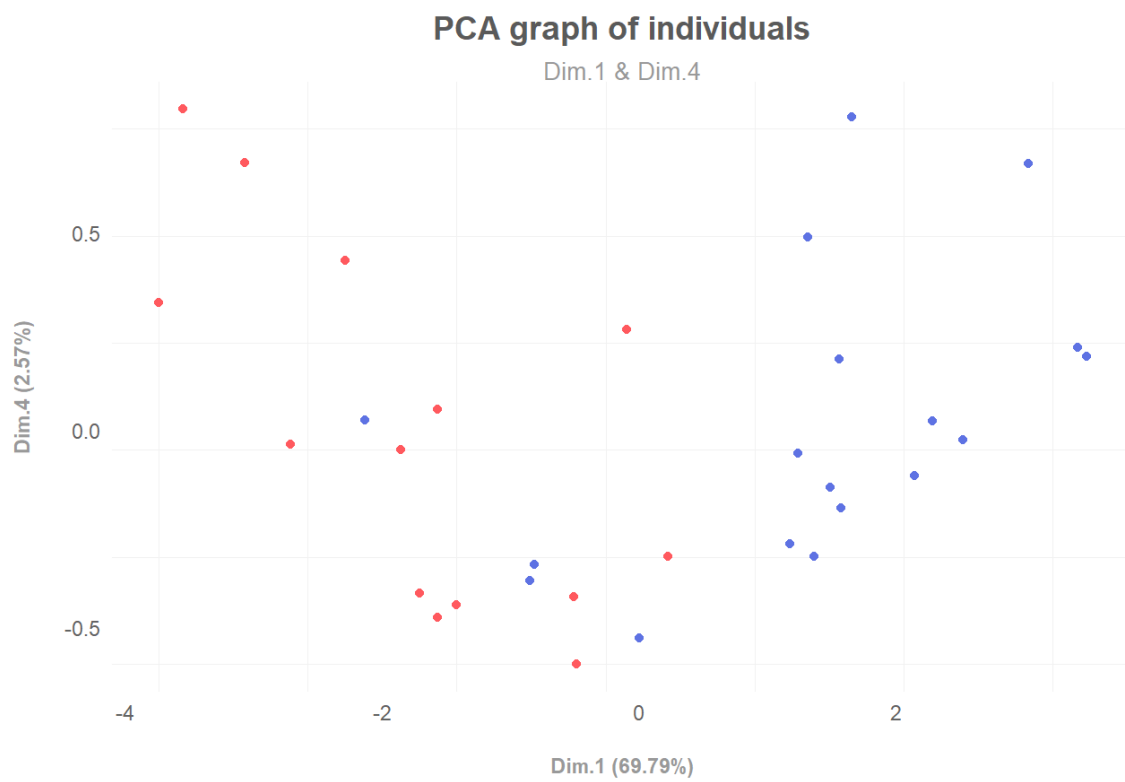
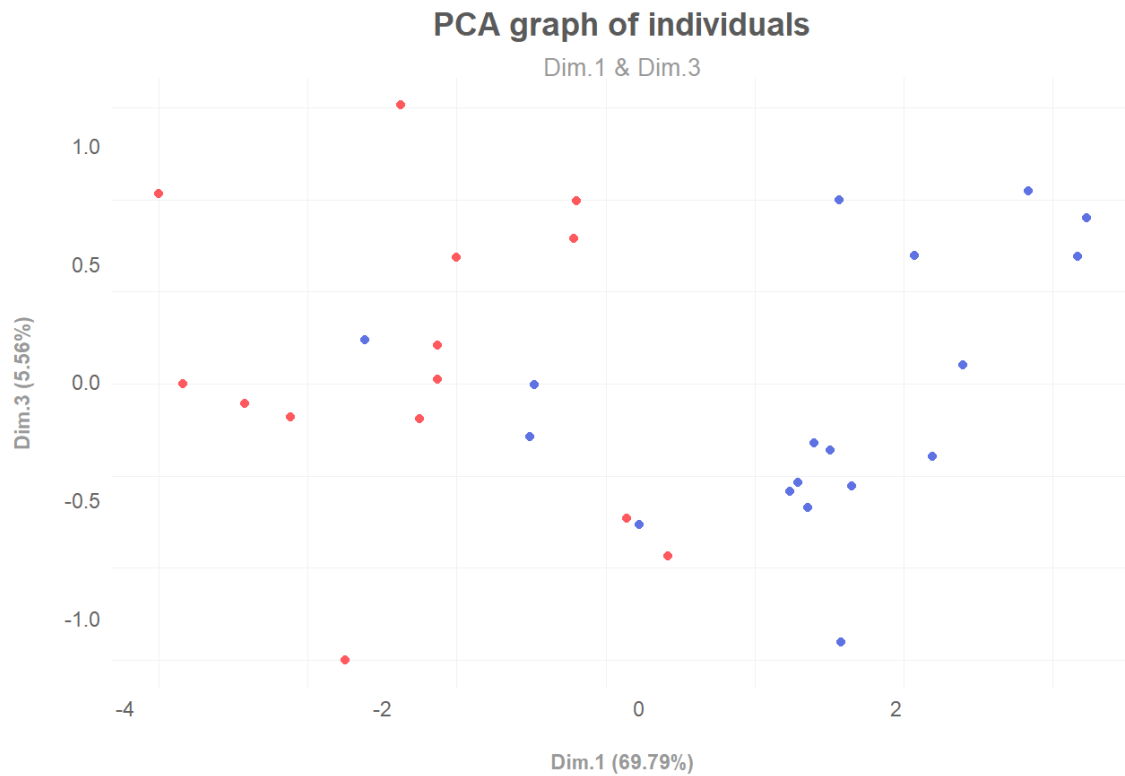
Distribution of principal components

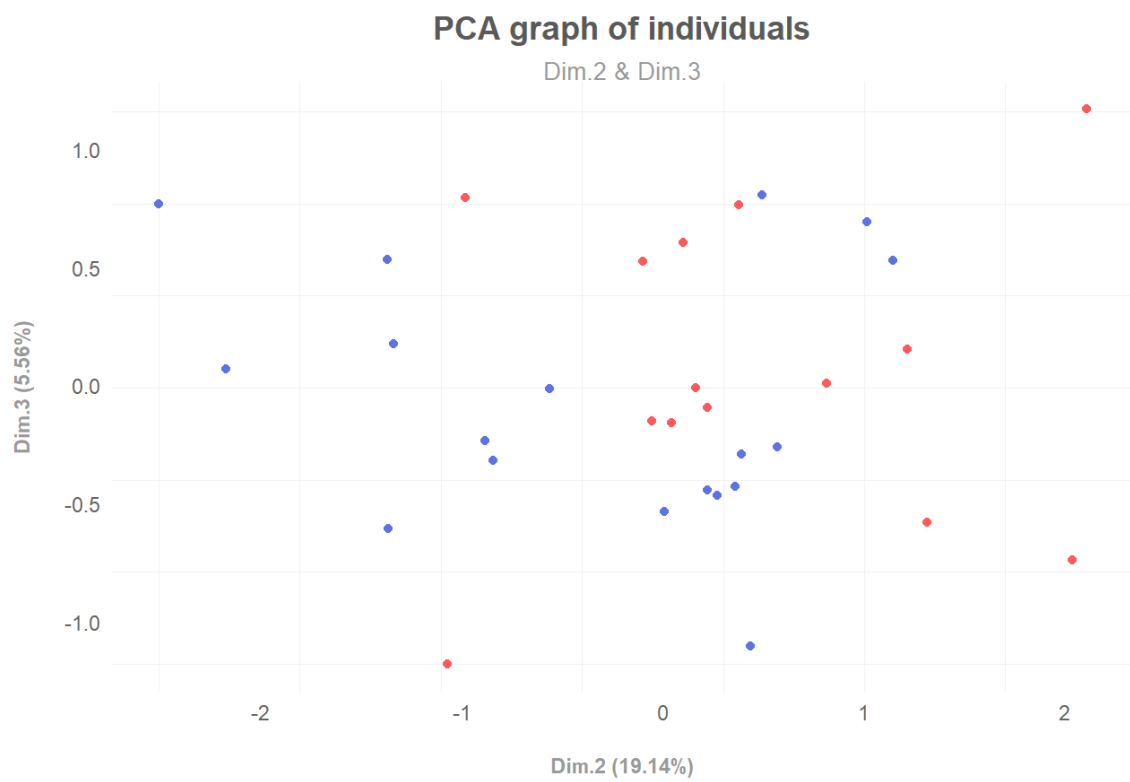
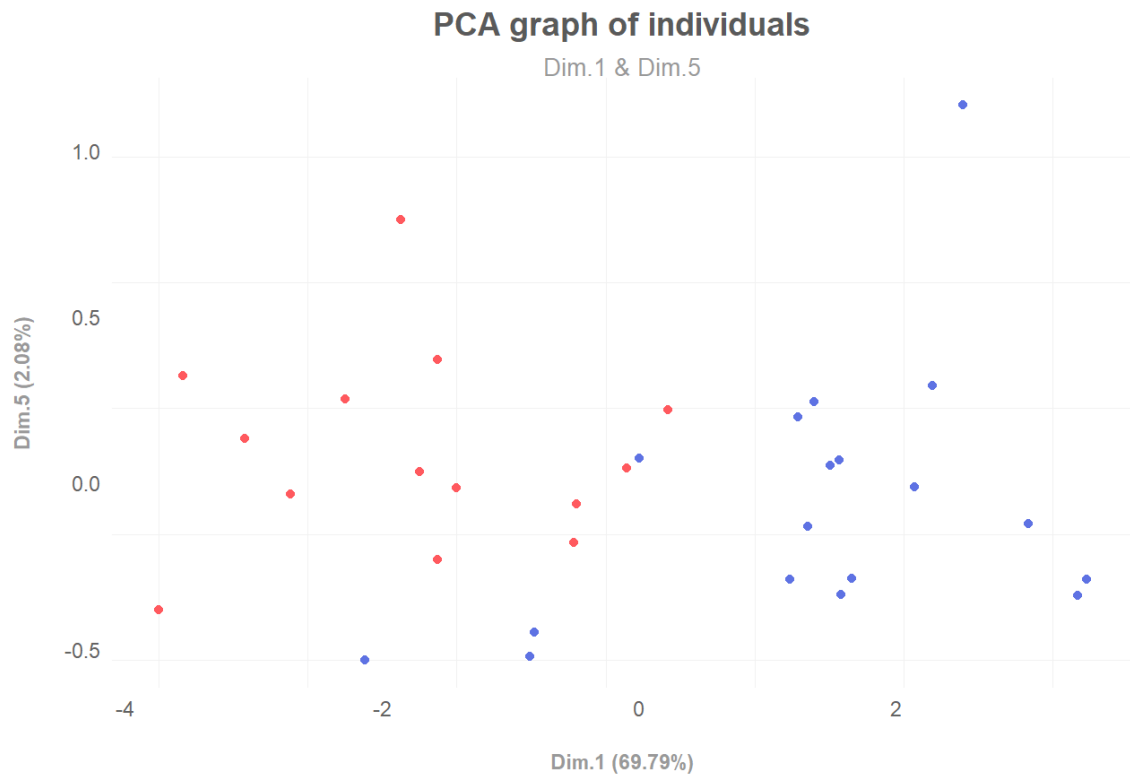


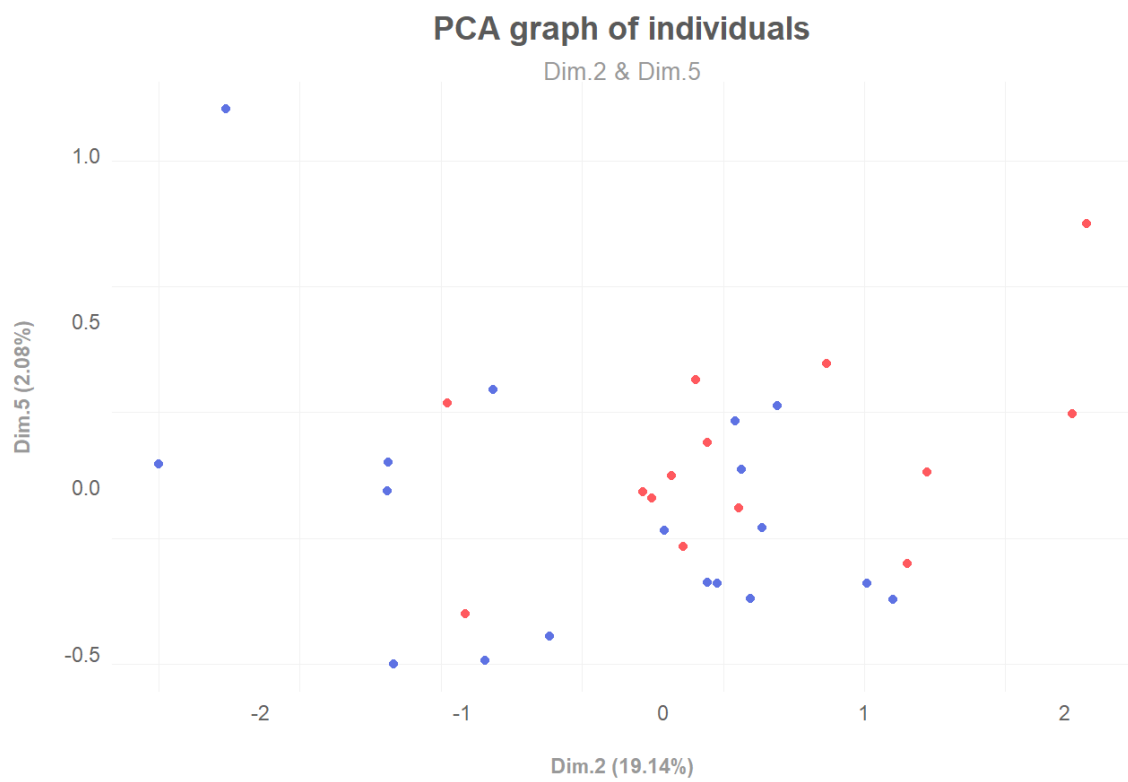
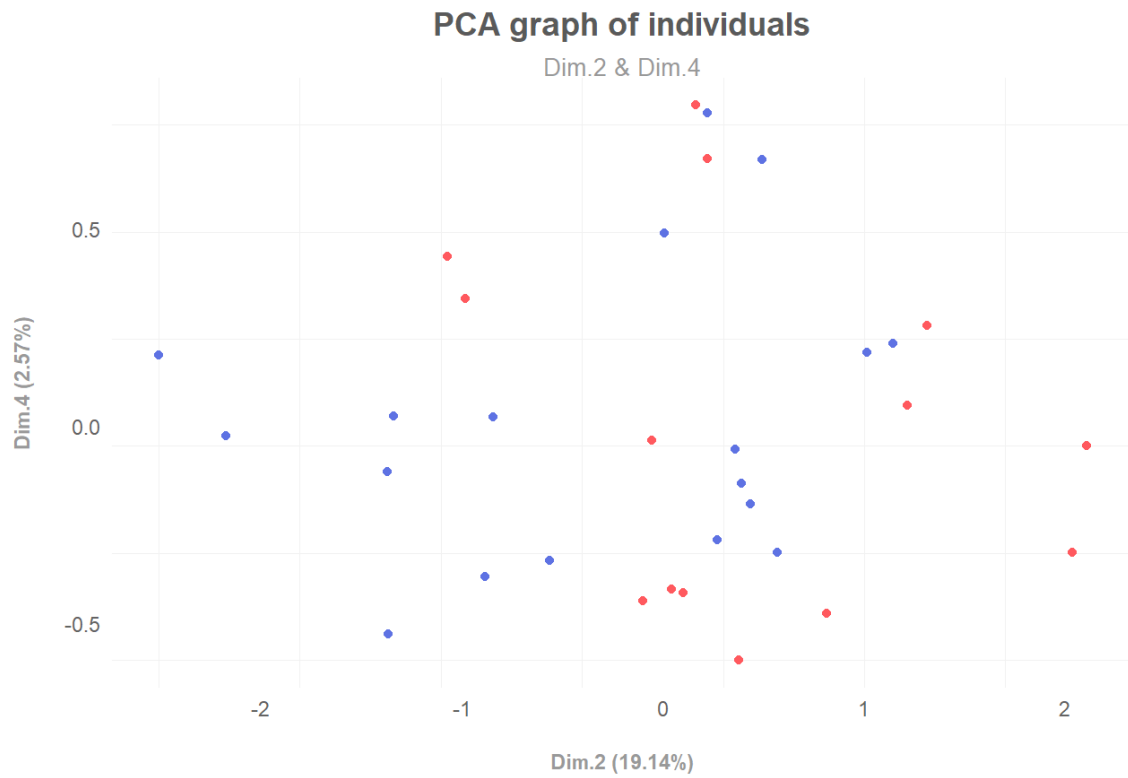
PCA graph of individuals

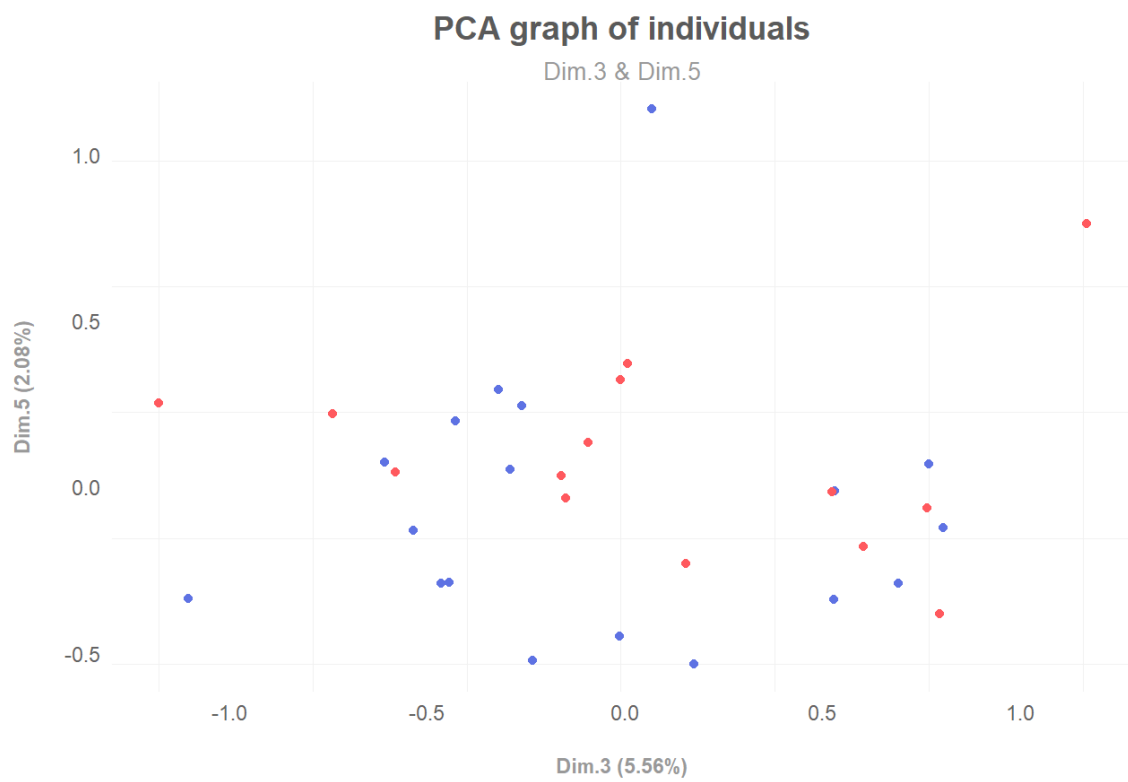
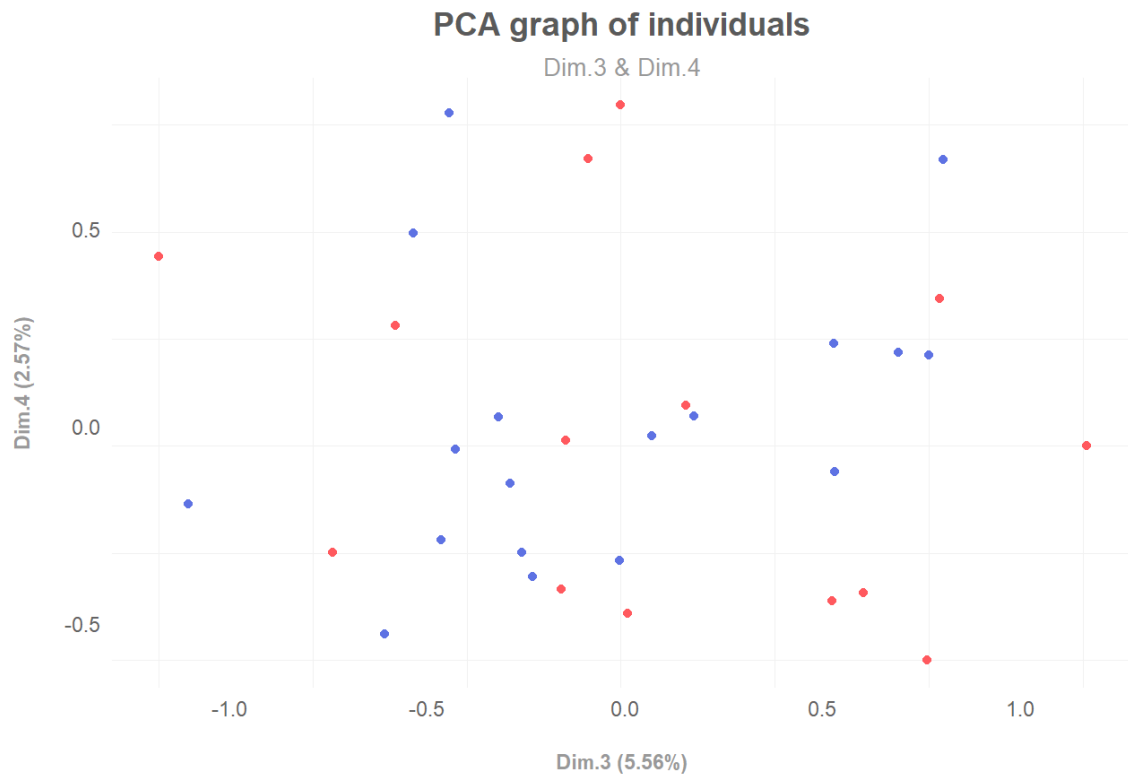
Dim.1 & Dim.2

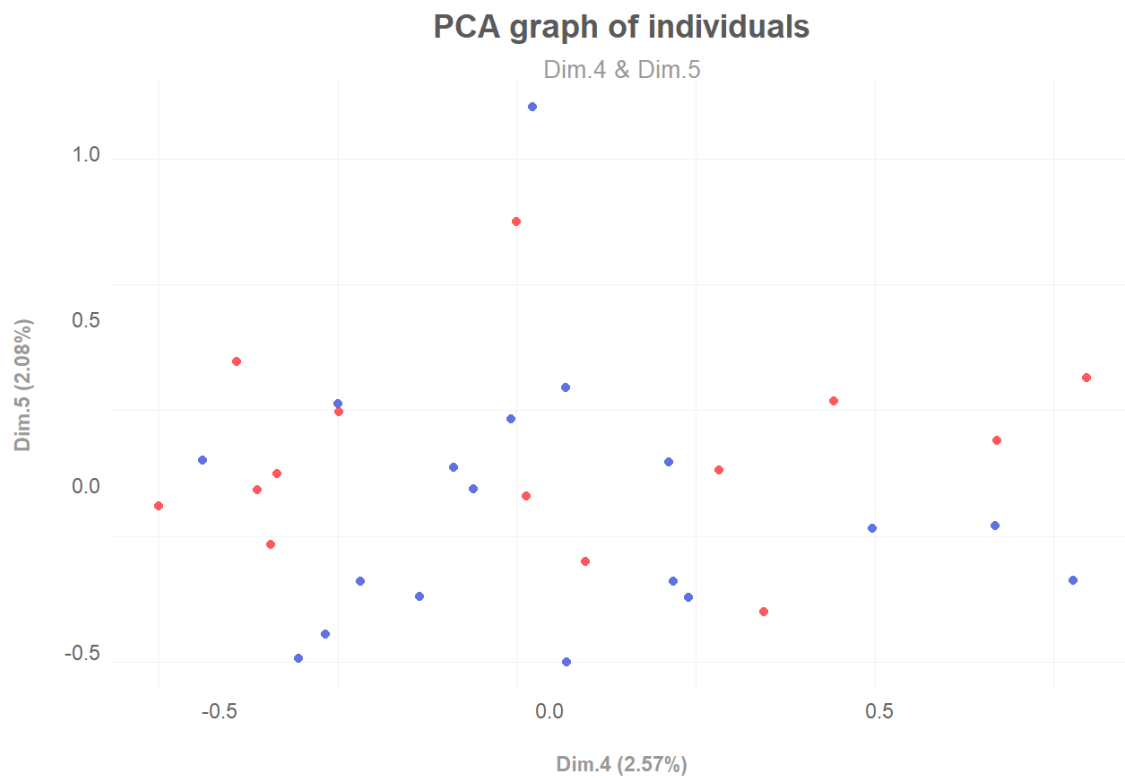










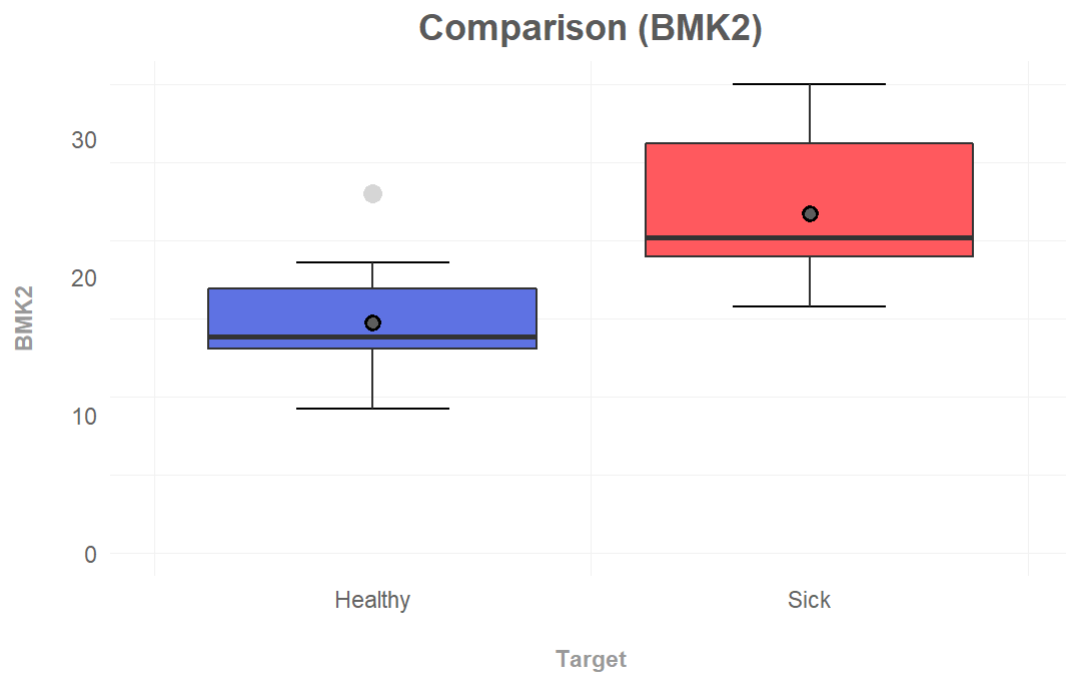


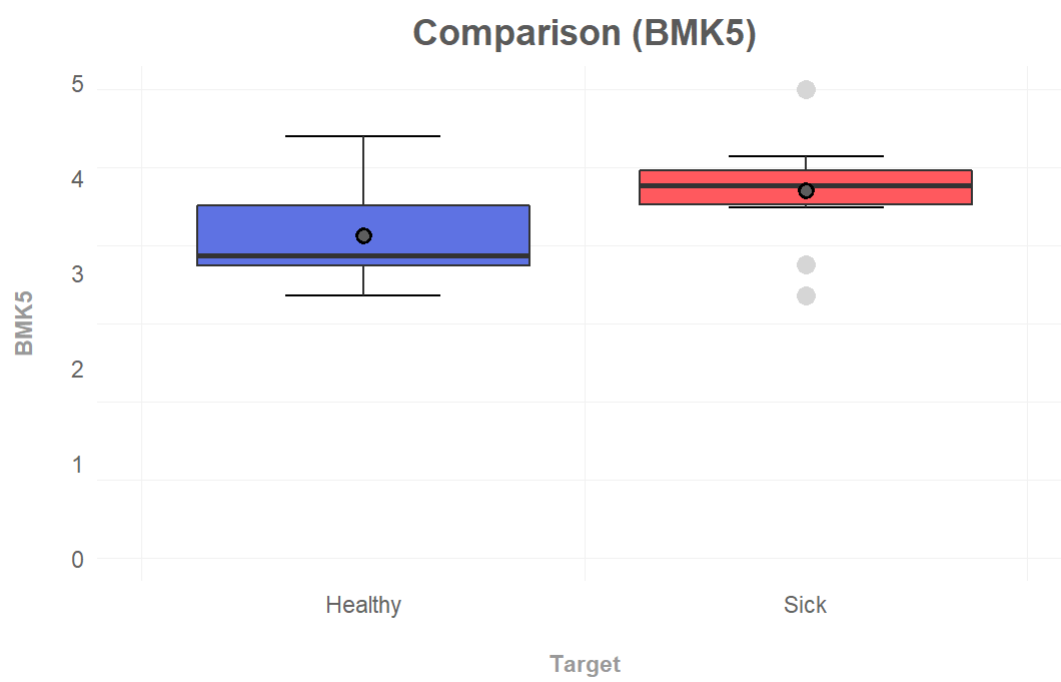
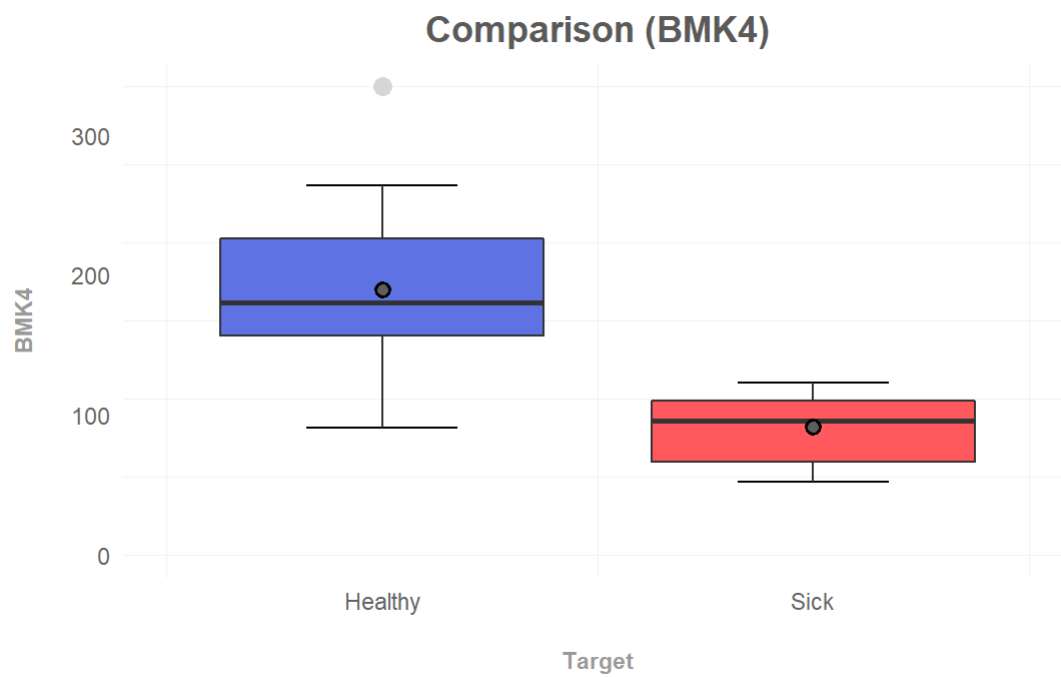
Statistical tests

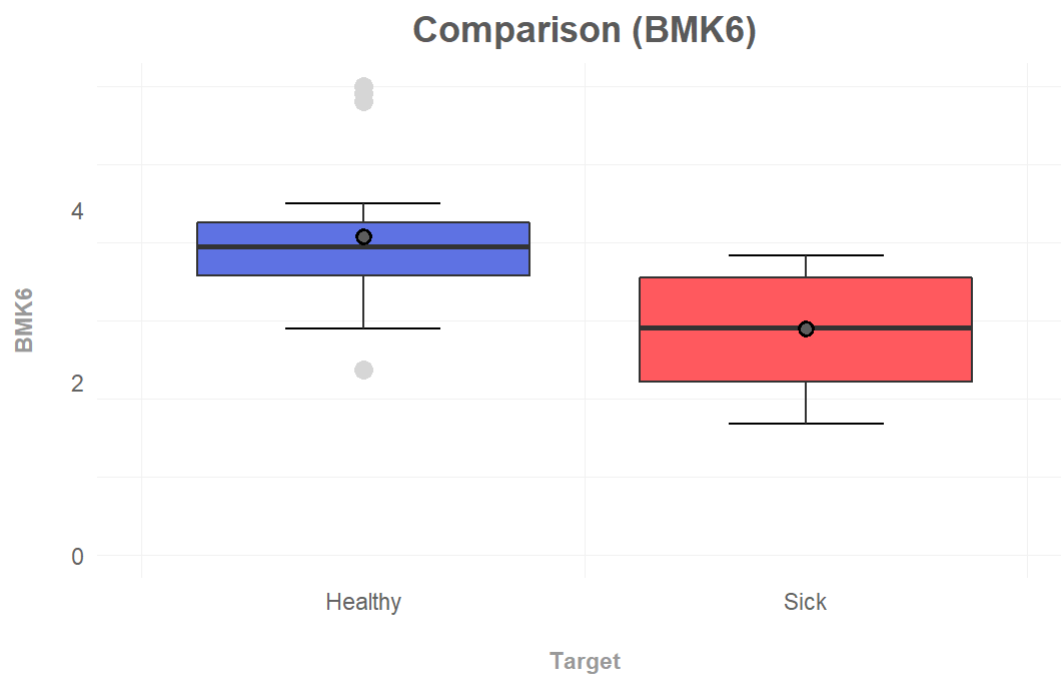
Continuous variables

Table 5: Overview of statistical tests for continuous variables

biomarker	group	Mann-Whitney's test	F-test	T-test	Welch's t-test	Paired t-test
BMK2	Target	8.35e-05	2.00e-01	3.42e-05	1.10e-04	3.42e-05
BMK3	Target	5.60e-05	2.61e-02	5.24e-06	2.48e-06	5.24e-06
BMK4	Target	2.86e-05	2.00e-03	2.94e-06	1.82e-06	2.94e-06
BMK5	Target	1.27e-02	7.88e-01	1.17e-02	1.29e-02	1.17e-02
BMK6	Target	1.08e-03	3.96e-01	9.80e-04	7.28e-04	9.80e-04
BMK7	Target	1.05e-05	4.00e-01	1.03e-06	3.52e-06	1.03e-06







Categorical variables

```
## [1] "No data available."
```

Correlations

Continuous relationships

Table 6: Summary of continuous correlations

	BMK2	BMK3	BMK4	BMK5	BMK6	BMK7
BMK2	1.000	-0.909	-0.895	0.651	-0.886	0.467
BMK3	-0.909	1.000	0.851	-0.684	0.898	-0.460
BMK4	-0.895	0.851	1.000	-0.520	0.775	-0.667
BMK5	0.651	-0.684	-0.520	1.000	-0.750	0.092
BMK6	-0.886	0.898	0.775	-0.750	1.000	-0.225
BMK7	0.467	-0.460	-0.667	0.092	-0.225	1.000

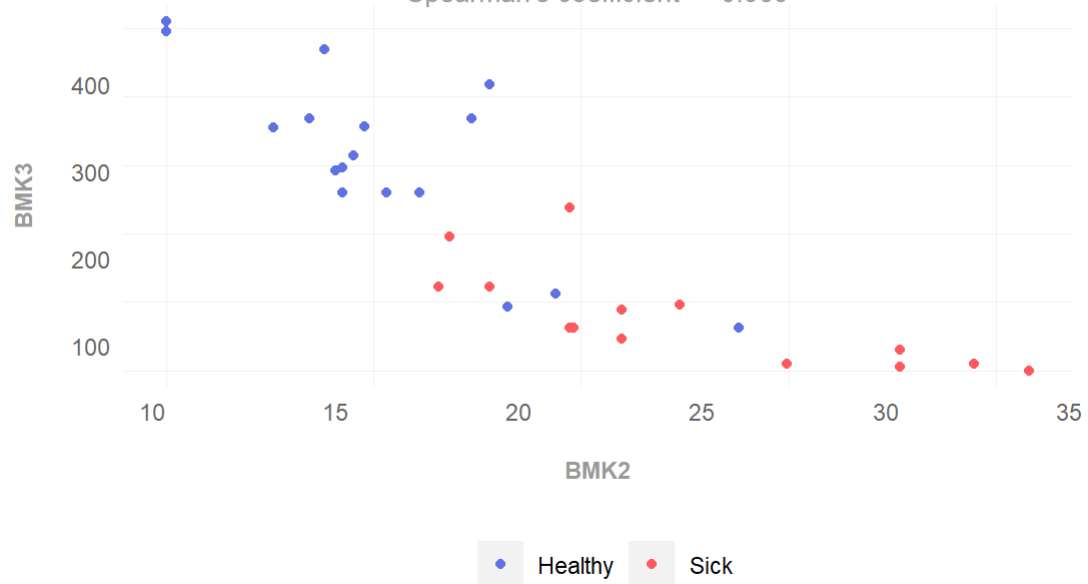
Correlation summary

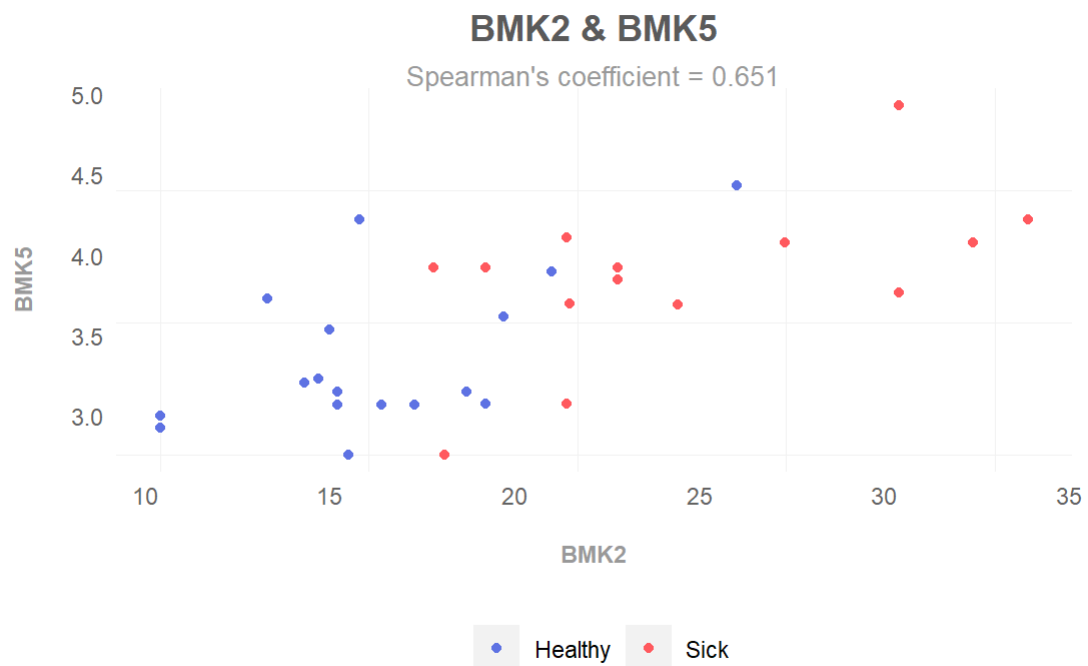
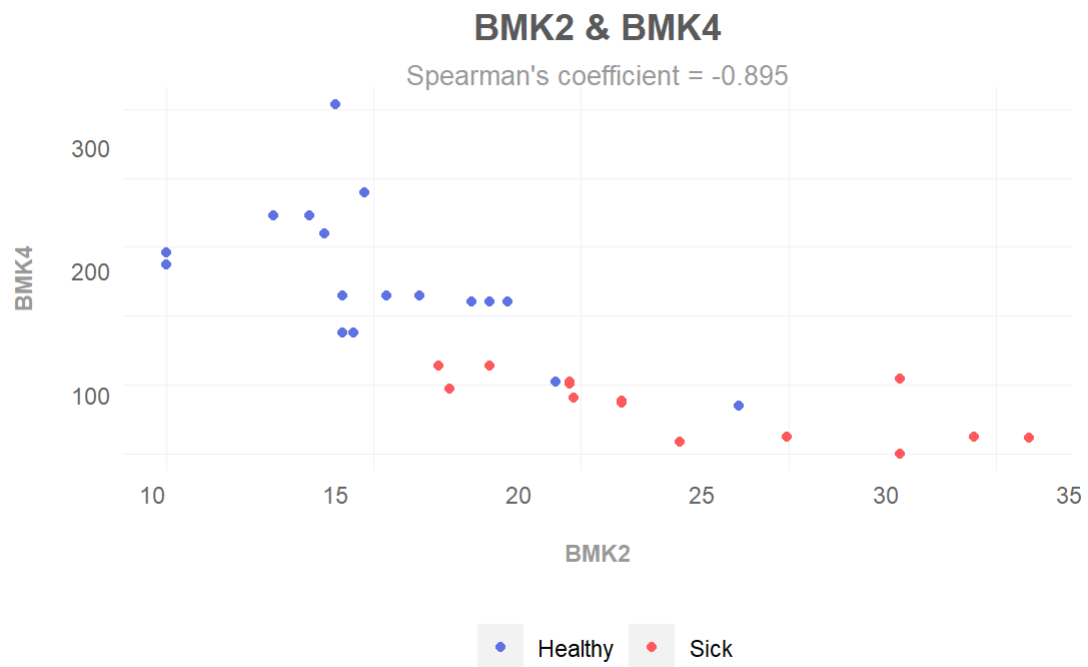
Spearman's coefficient was used to establish relationships

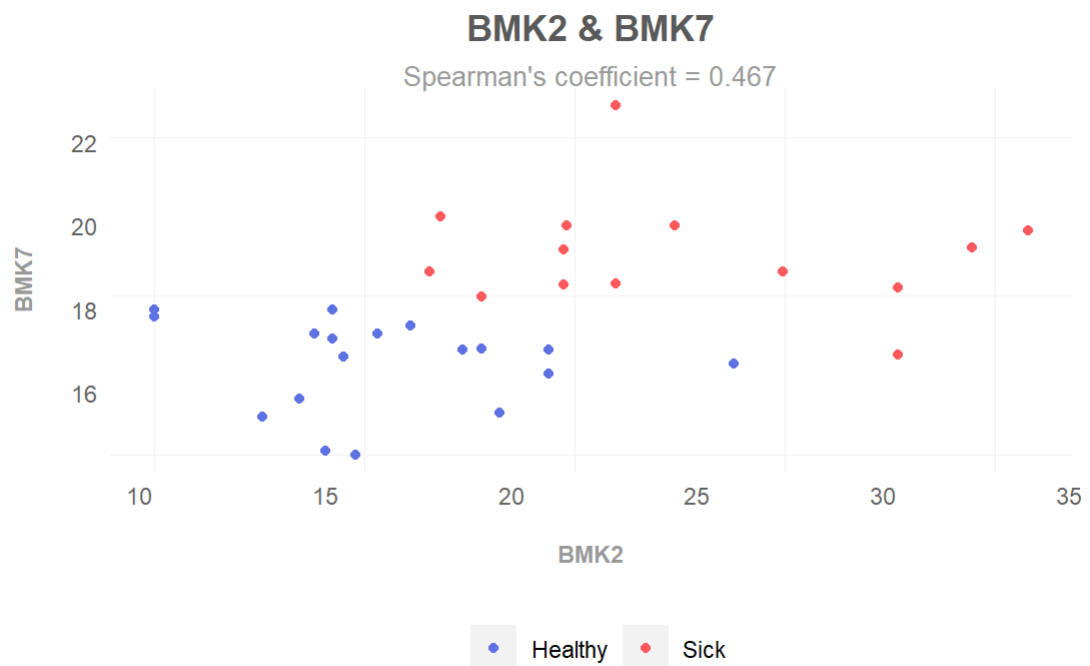
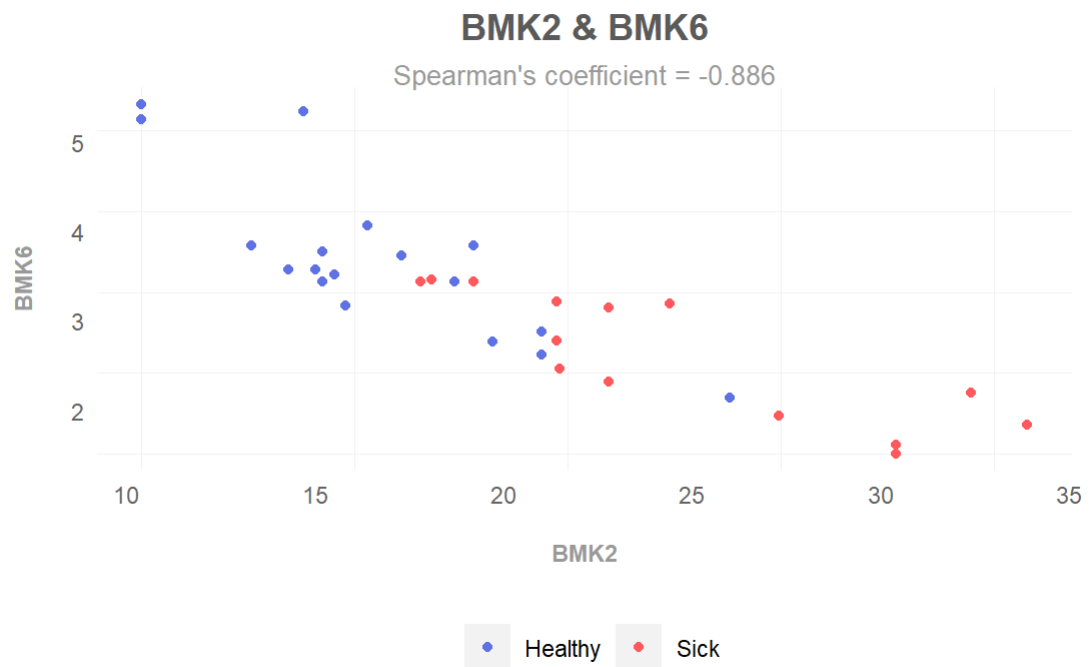
BMK7						1
BMK6					1	-0.225
BMK5				1	-0.75	0.092
BMK4			1	-0.52	0.775	-0.667
BMK3		1	0.851	-0.684	0.898	-0.46
BMK2	1	-0.909	-0.895	0.651	-0.886	0.467
	BMK2	BMK3	BMK4	BMK5	BMK6	BMK7

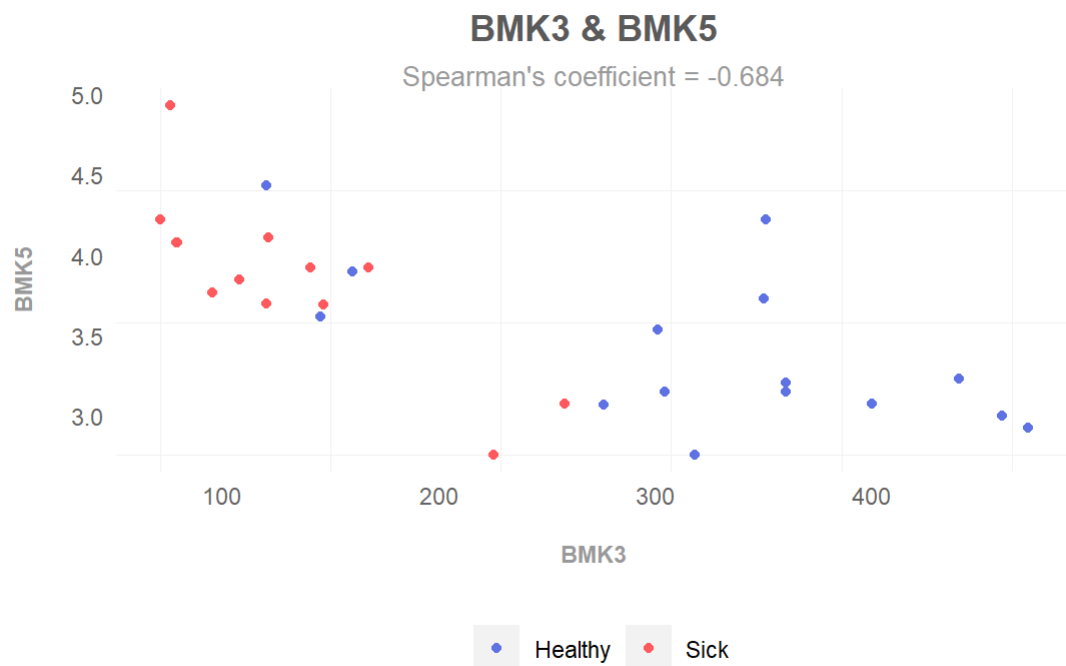
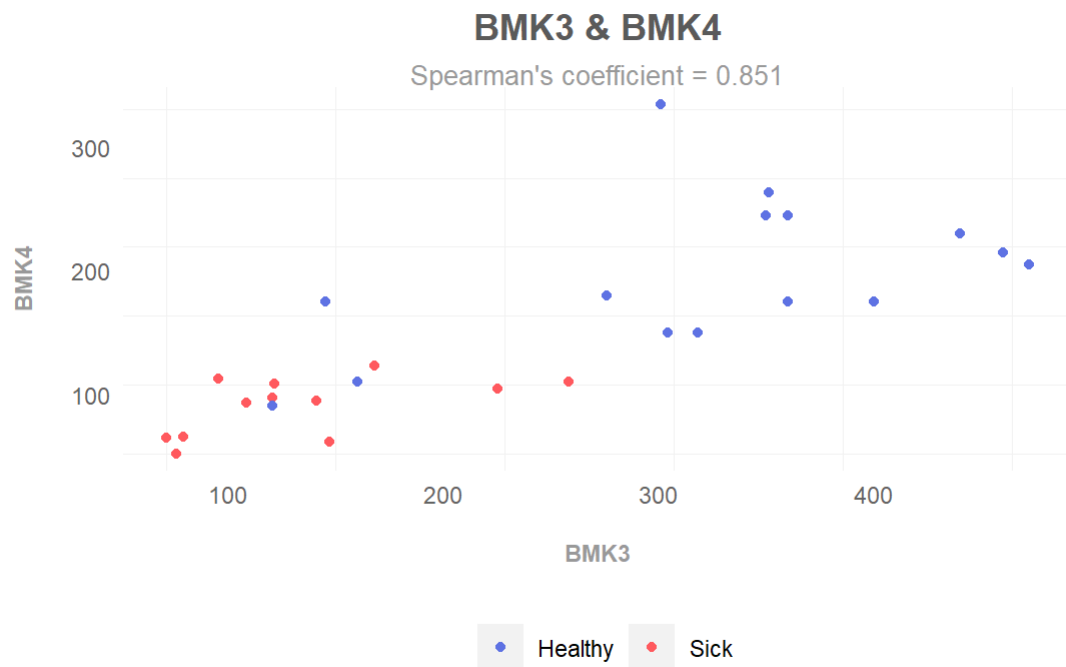
BMK2 & BMK3

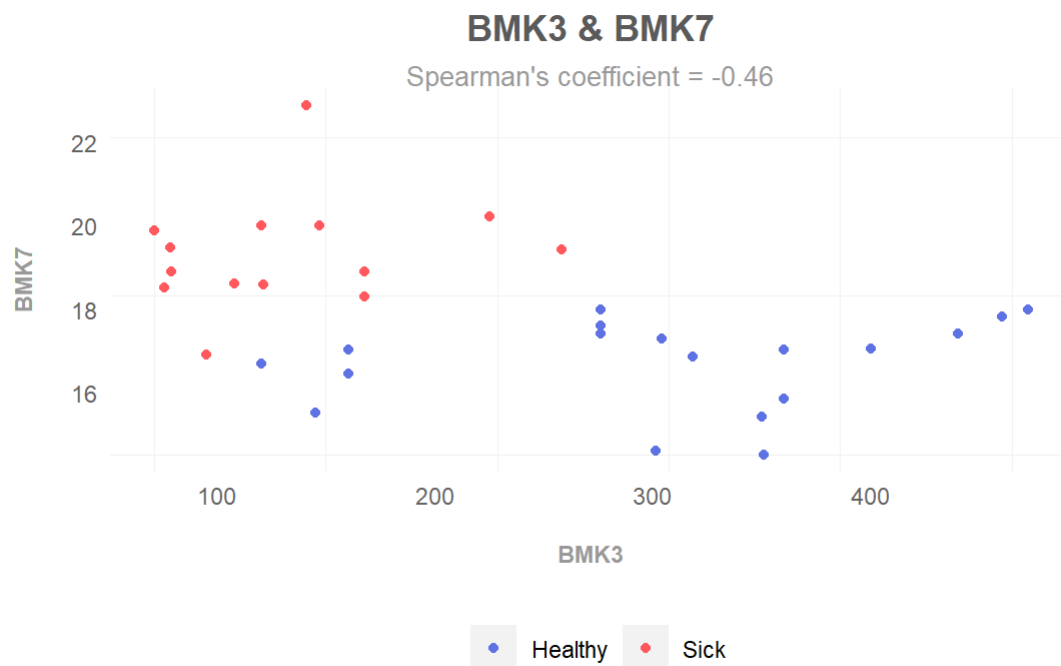
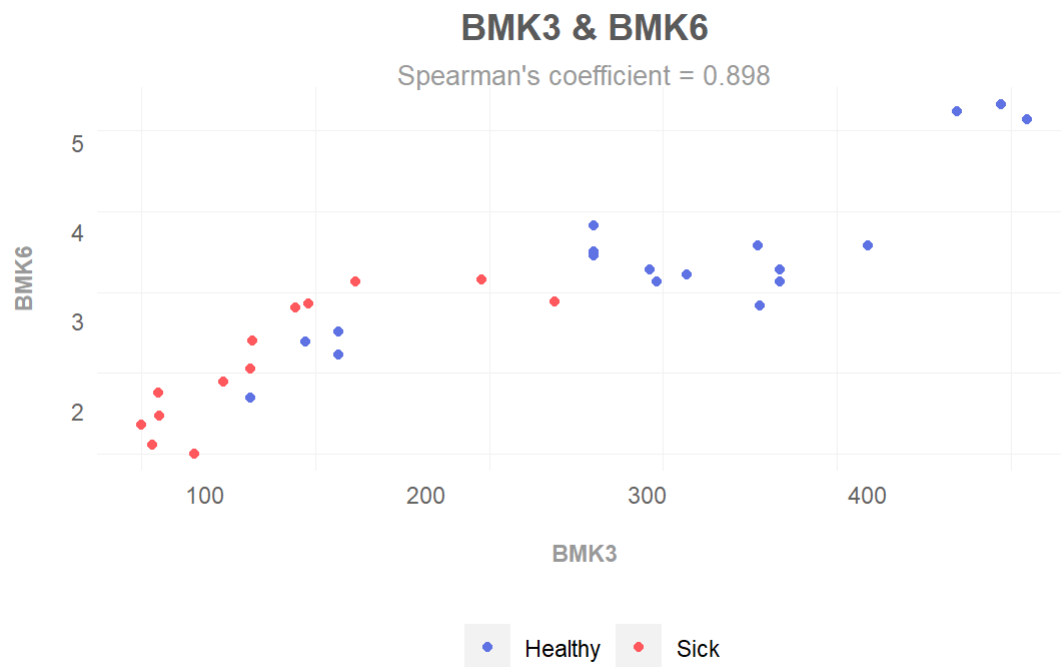
Spearman's coefficient = -0.909

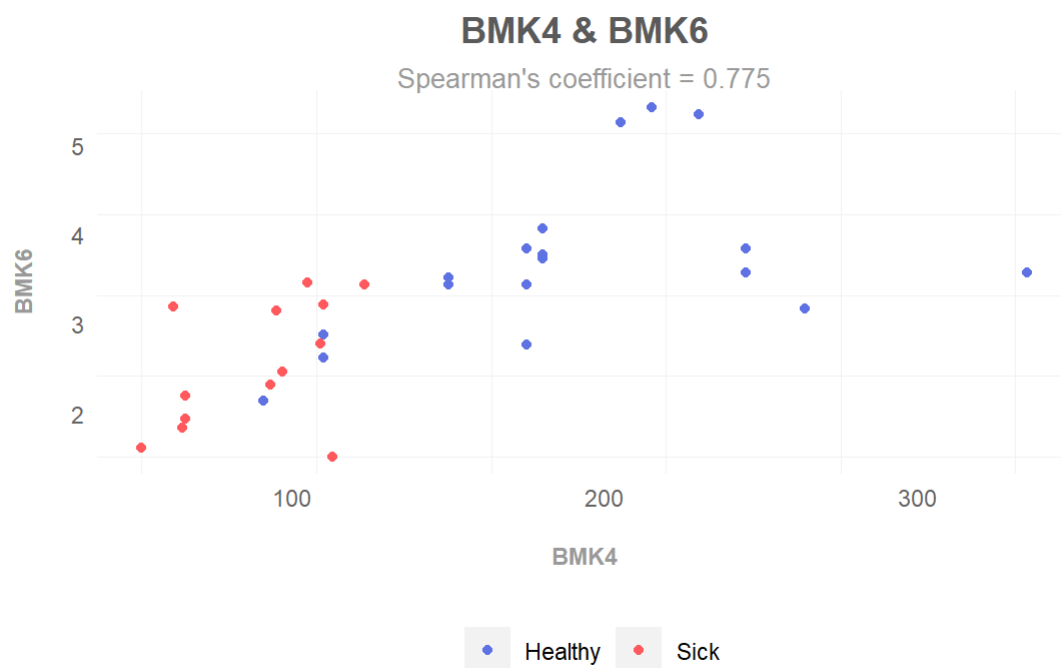
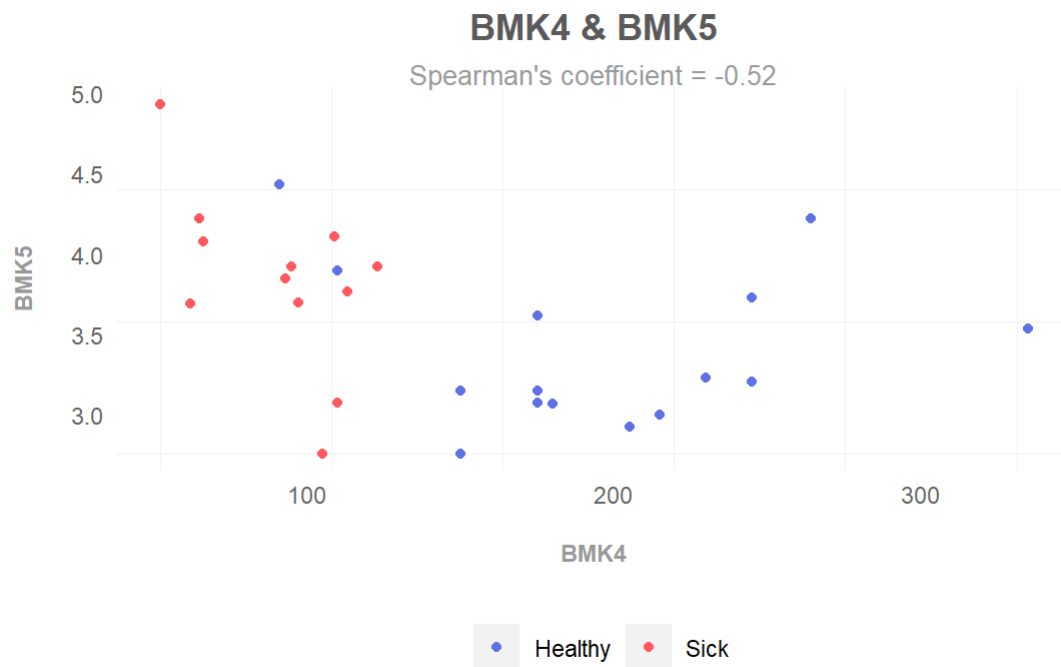


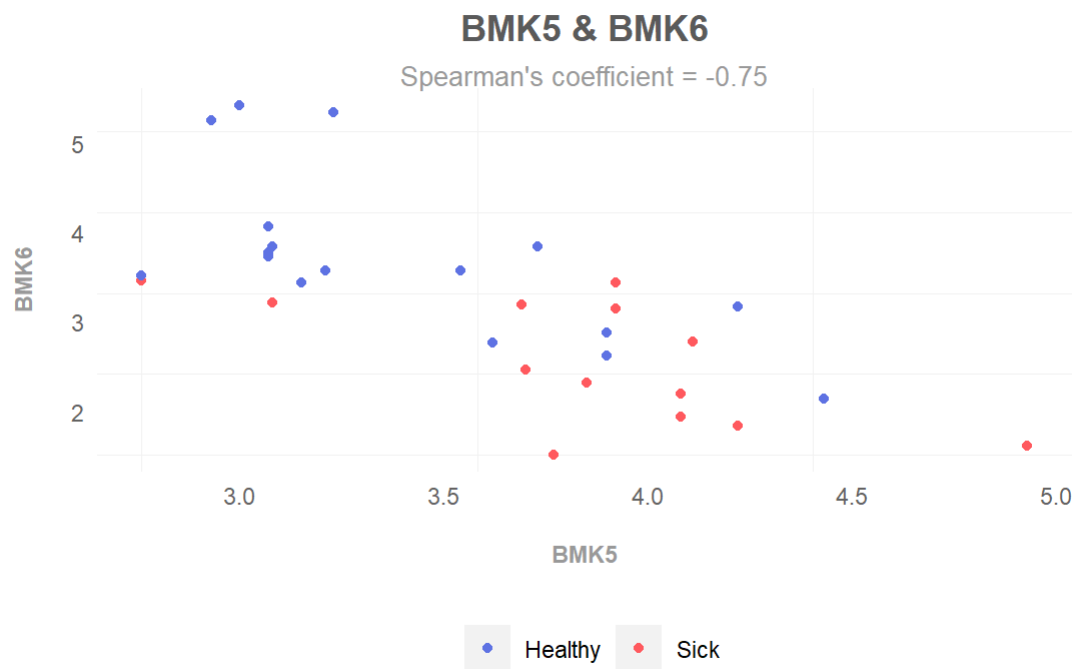
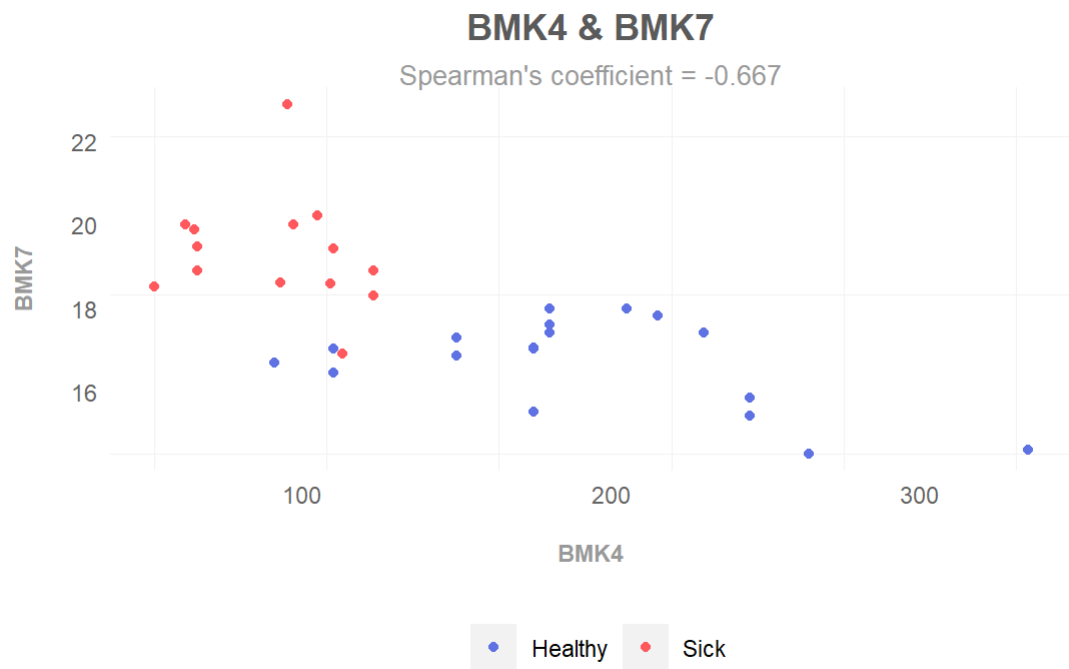


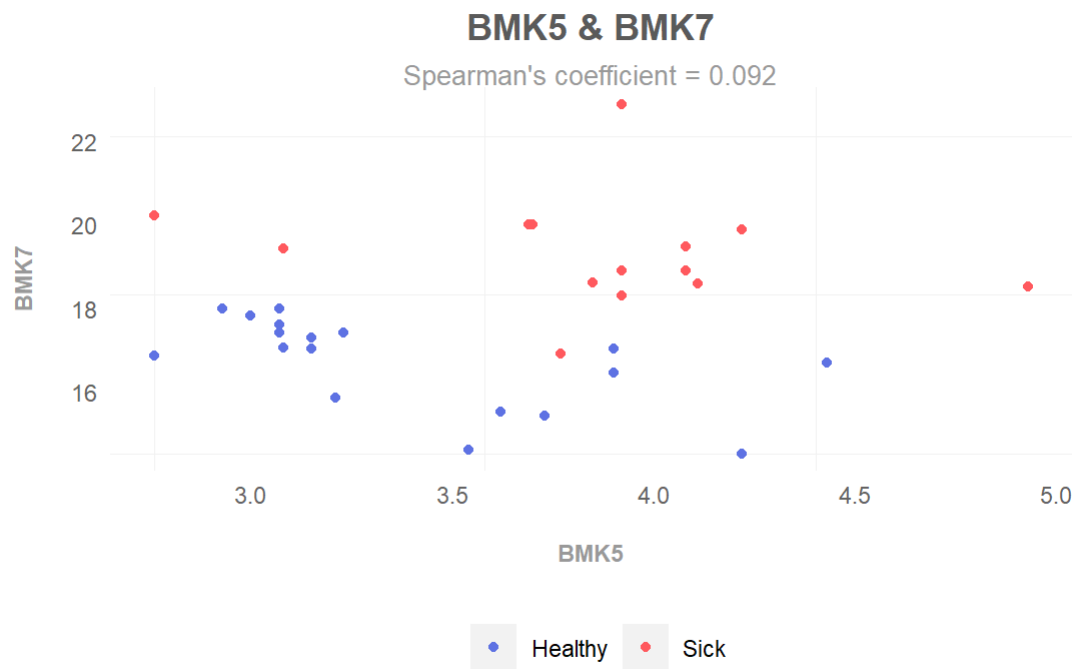












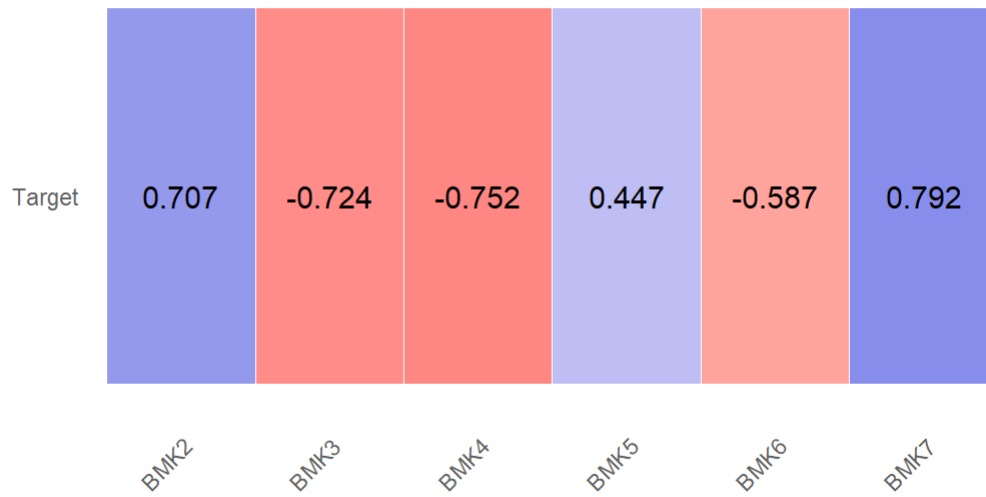
Continuous & Categorical relationships

Table 7: Summary of continuous and categorical correlations

	BMK2	BMK3	BMK4	BMK5	BMK6	BMK7
Target	0.707	-0.724	-0.752	0.447	-0.587	0.792

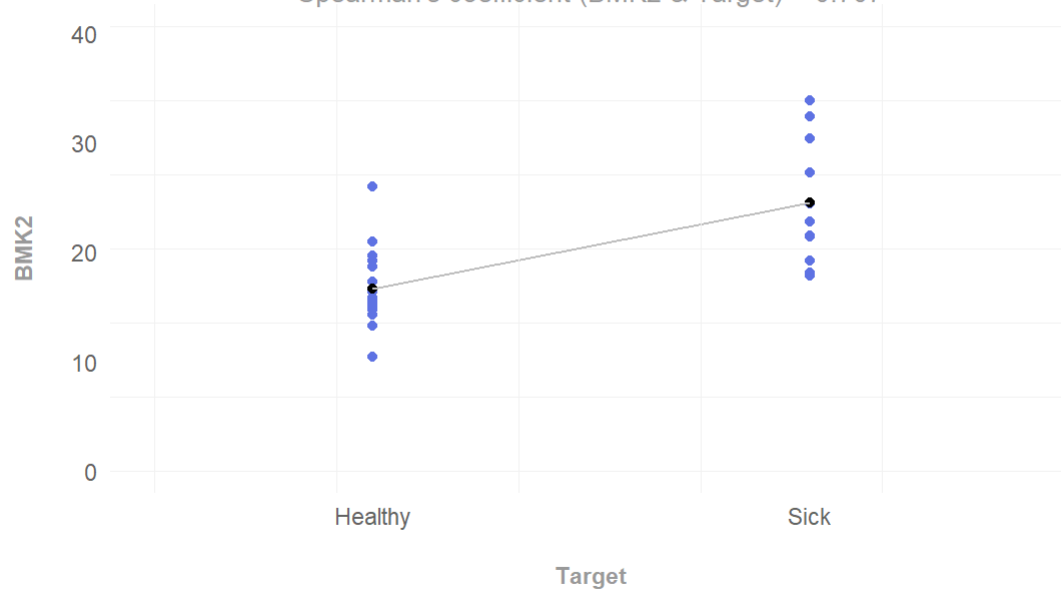
Correlation summary

Spearman's coefficient was used to establish relationships



Point biserial correlation

Spearman's coefficient (BMK2 & Target) = 0.707



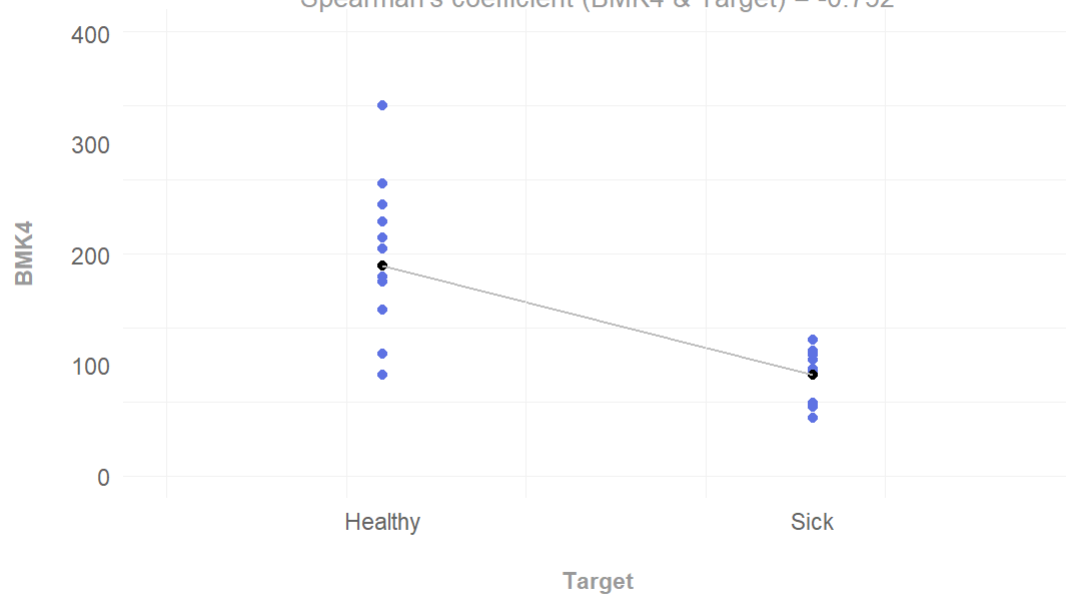
Point biserial correlation

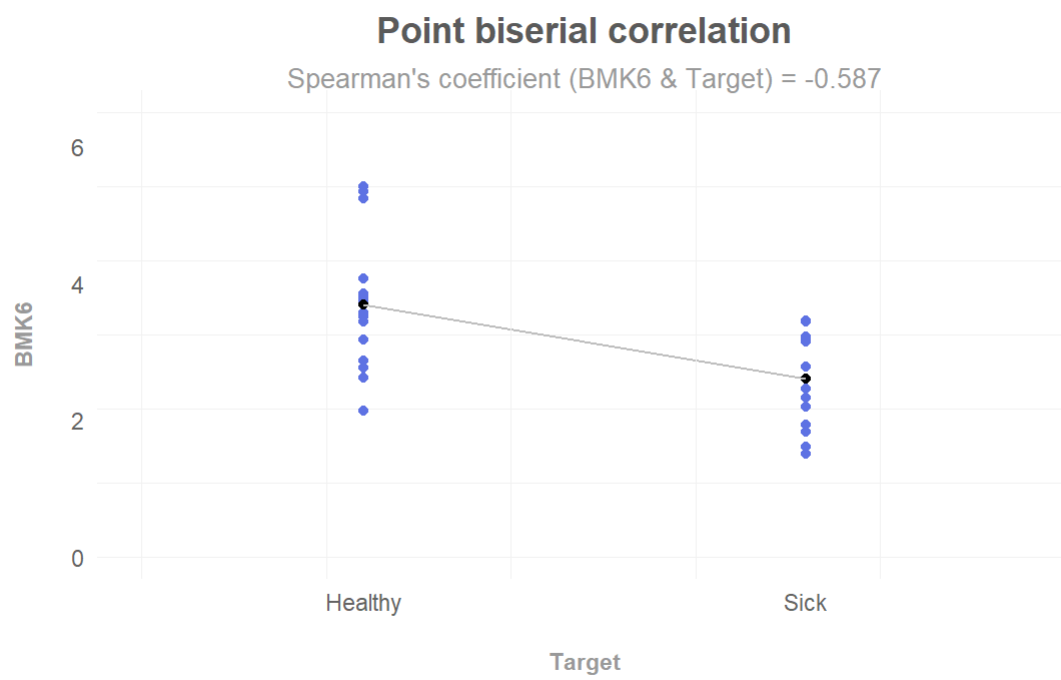
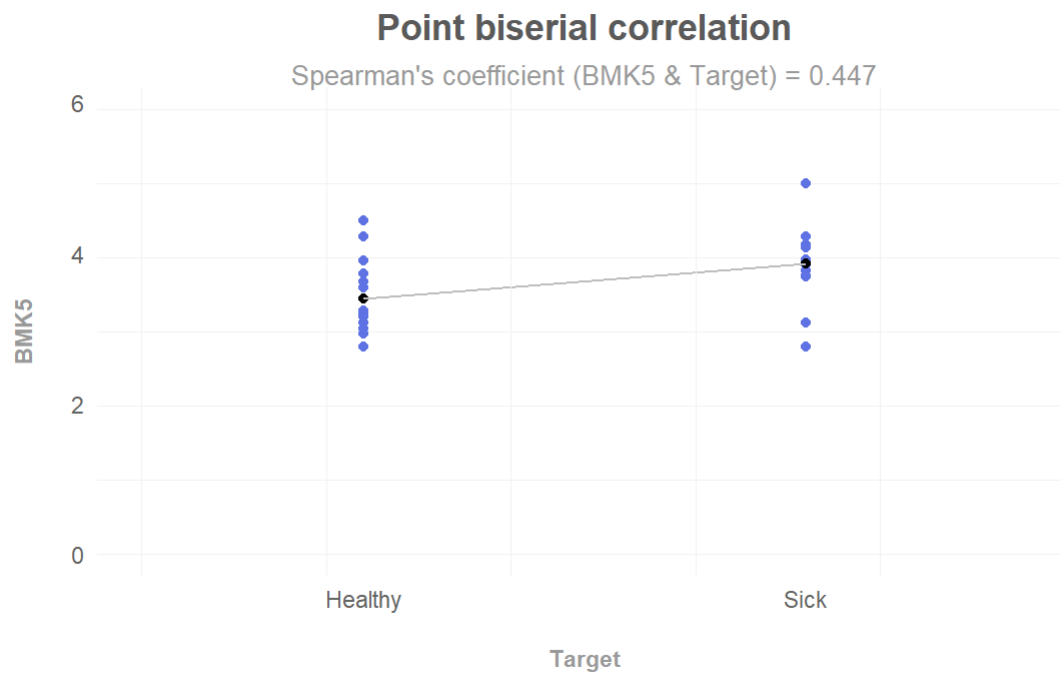
Spearman's coefficient (BMK3 & Target) = -0.724



Point biserial correlation

Spearman's coefficient (BMK4 & Target) = -0.752





Categorical relationships

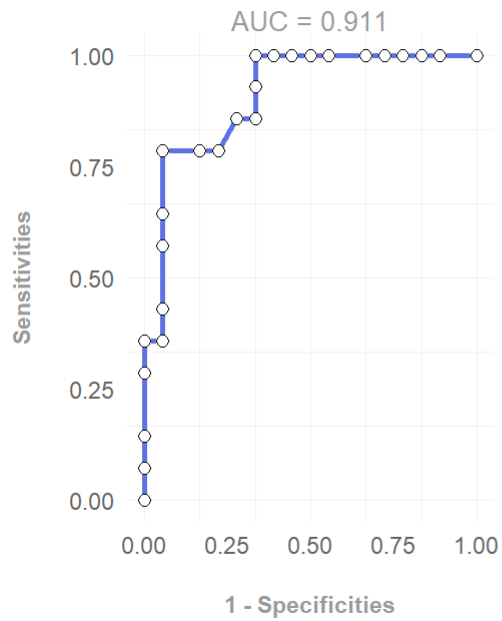
```
## [1] "No data available."
```

Individual performances

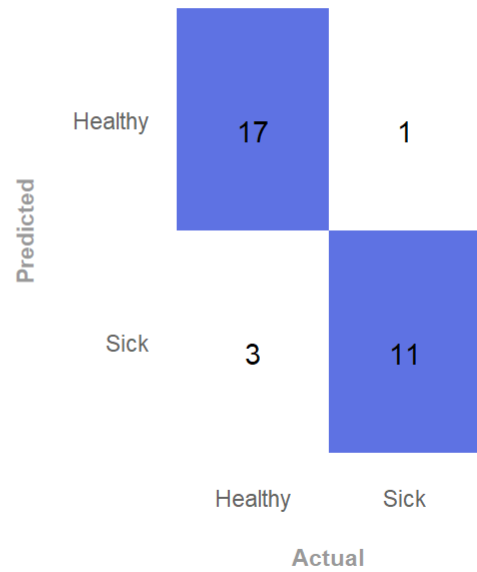
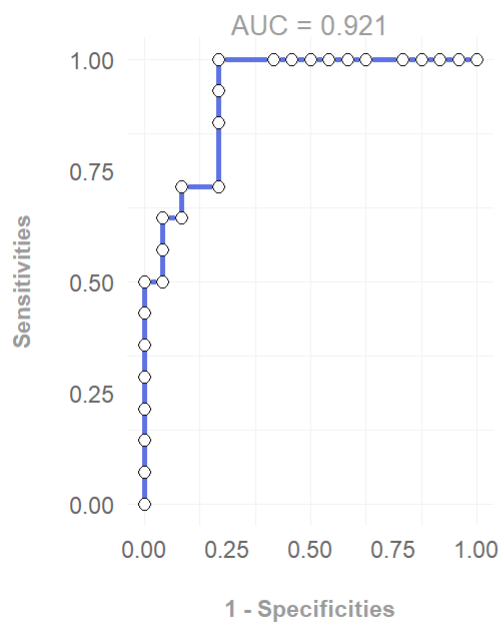
Classification

Table 8: Summary of individual performances

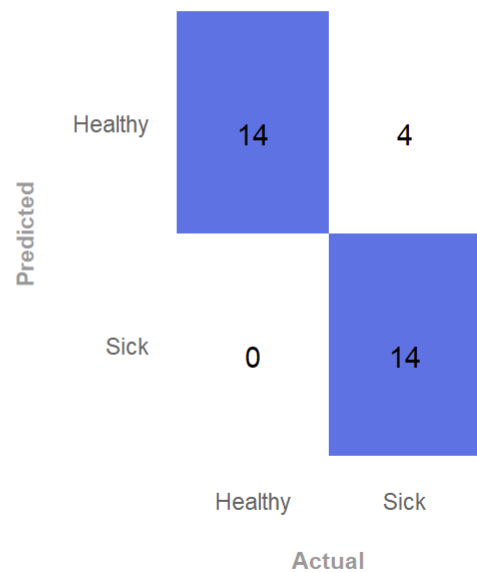
biomarker	class1	class2	logLoss	best.method	threshold	auc	sens	spe	PPV	NPV	accuracy	no.info.rate	balanced.accuracy	precision	f1	TP	FP	TN	FN
BMK2	Healthy	Sick	0.3989584	youden	21.200	0.911	0.850	0.917	0.944	0.786	0.875	0.625	0.883	0.944	0.895	17	1	11	3
BMK3	Healthy	Sick	0.3546224	youden	266.900	0.921	1.000	0.778	0.778	1.000	0.875	0.562	0.889	0.778	0.875	14	4	14	0
BMK4	Healthy	Sick	0.2630872	youden	136.500	0.937	1.000	0.824	0.833	1.000	0.906	0.531	0.912	0.833	0.909	15	3	14	0
BMK5	Healthy	Sick	0.5806063	youden	3.655	0.760	0.867	0.706	0.722	0.857	0.781	0.531	0.786	0.722	0.788	13	5	12	2
BMK6	Healthy	Sick	0.4901052	youden	3.490	0.841	1.000	0.667	0.611	1.000	0.781	0.656	0.833	0.611	0.759	11	7	14	0
BMK7	Healthy	Sick	0.2199398	youden	18.150	0.960	0.947	1.000	1.000	0.929	0.969	0.594	0.974	1.000	0.973	18	0	13	1

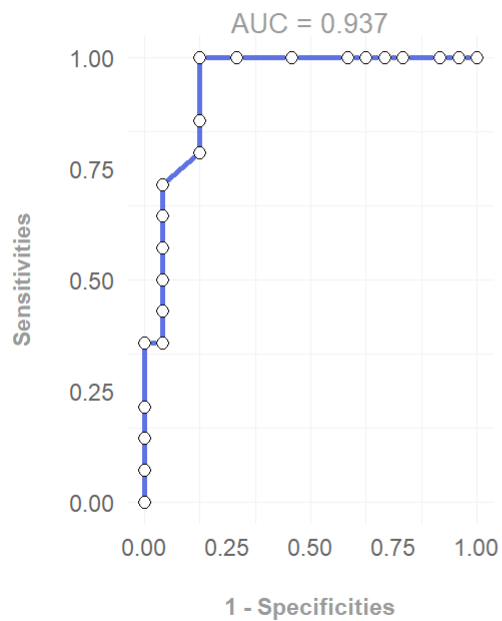
ROC Curve (BMK2)**Confusion Matrix**

Healthy is used as positive case

**ROC Curve (BMK3)****Confusion Matrix**

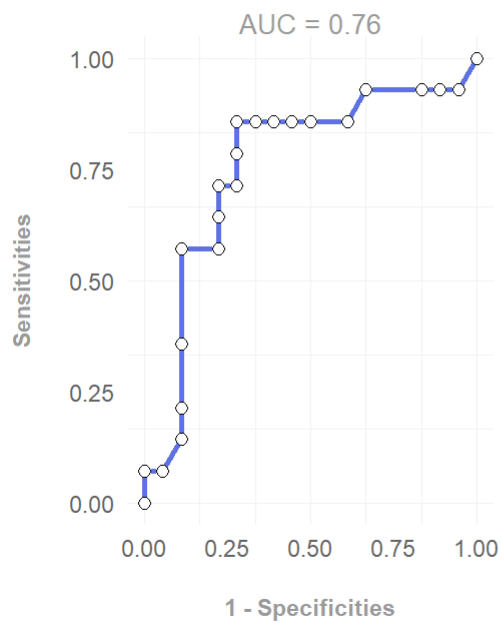
Healthy is used as positive case



ROC Curve (BMK4)**Confusion Matrix**

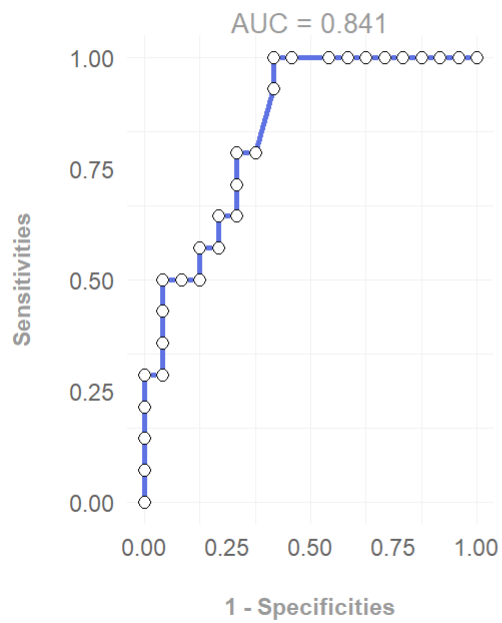
Healthy is used as positive case

Predicted	Actual	
	Healthy	Sick
Healthy	15	3
Sick	0	14

ROC Curve (BMK5)**Confusion Matrix**

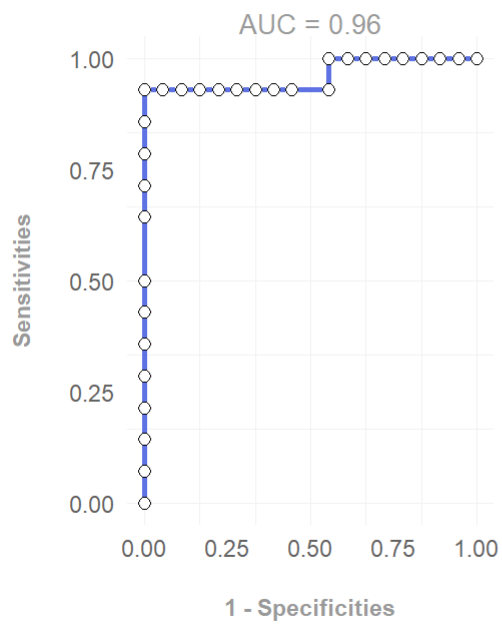
Healthy is used as positive case

Predicted	Actual	
	Healthy	Sick
Healthy	13	5
Sick	2	12

ROC Curve (BMK6)**Confusion Matrix**

Healthy is used as positive case

Predicted	Actual	
	Healthy	Sick
Healthy	11	7
Sick	0	14

ROC Curve (BMK7)**Confusion Matrix**

Healthy is used as positive case

Predicted	Actual	
	Healthy	Sick
Healthy	18	0
Sick	1	13

Predictive models

```
## [1] "No yet available."
```




GENESEN