

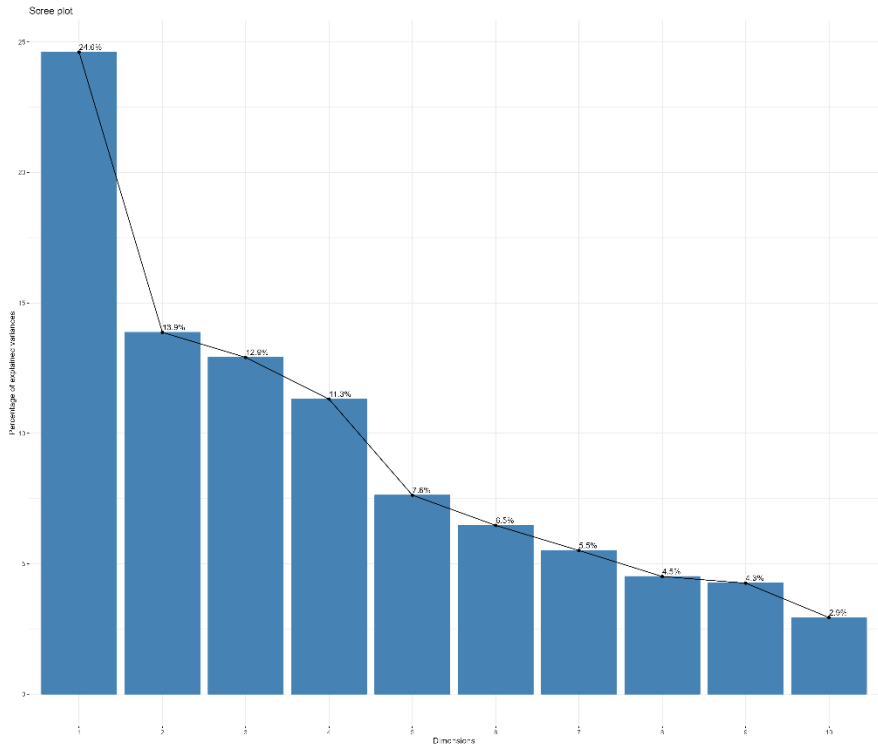


Individuals (the 10 first)

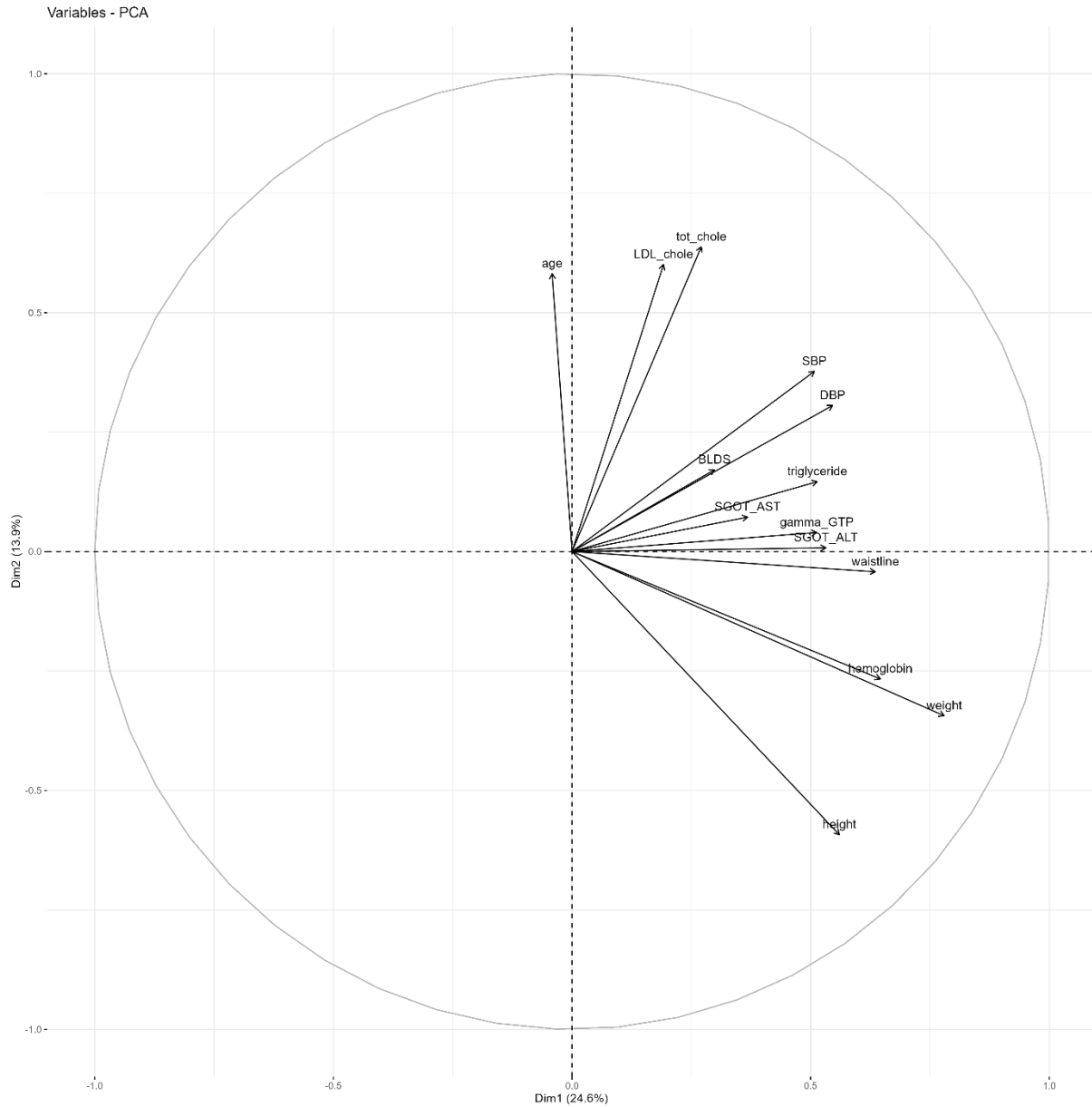
	Dist	Dim.1	ctr	cos2	Dim.2	ctr	cos2	Dim.3	ctr	cos
2										
1	2.619	1.604	0.000	0.375	-1.211	0.000	0.214	-0.993	0.000	0.14
4										
2	3.329	2.295	0.000	0.475	-0.776	0.000	0.054	-1.792	0.000	0.29
0										
3	2.886	0.874	0.000	0.092	-2.149	0.000	0.555	0.993	0.000	0.11
8										
4	3.606	2.642	0.000	0.537	-0.715	0.000	0.039	0.007	0.000	0.00
0										
5	1.516	-0.008	0.000	0.000	0.525	0.000	0.120	0.242	0.000	0.02
6										
6	2.690	0.855	0.000	0.101	1.048	0.000	0.152	0.870	0.000	0.10
5										
7	3.380	-2.961	0.000	0.768	-0.008	0.000	0.000	-0.697	0.000	0.04
3										
8	1.960	0.563	0.000	0.082	-1.038	0.000	0.281	-0.291	0.000	0.02
2										
9	2.577	1.656	0.000	0.413	0.732	0.000	0.081	-0.221	0.000	0.00
7										
10	3.578	1.661	0.000	0.216	-0.200	0.000	0.003	0.126	0.000	0.00
1										

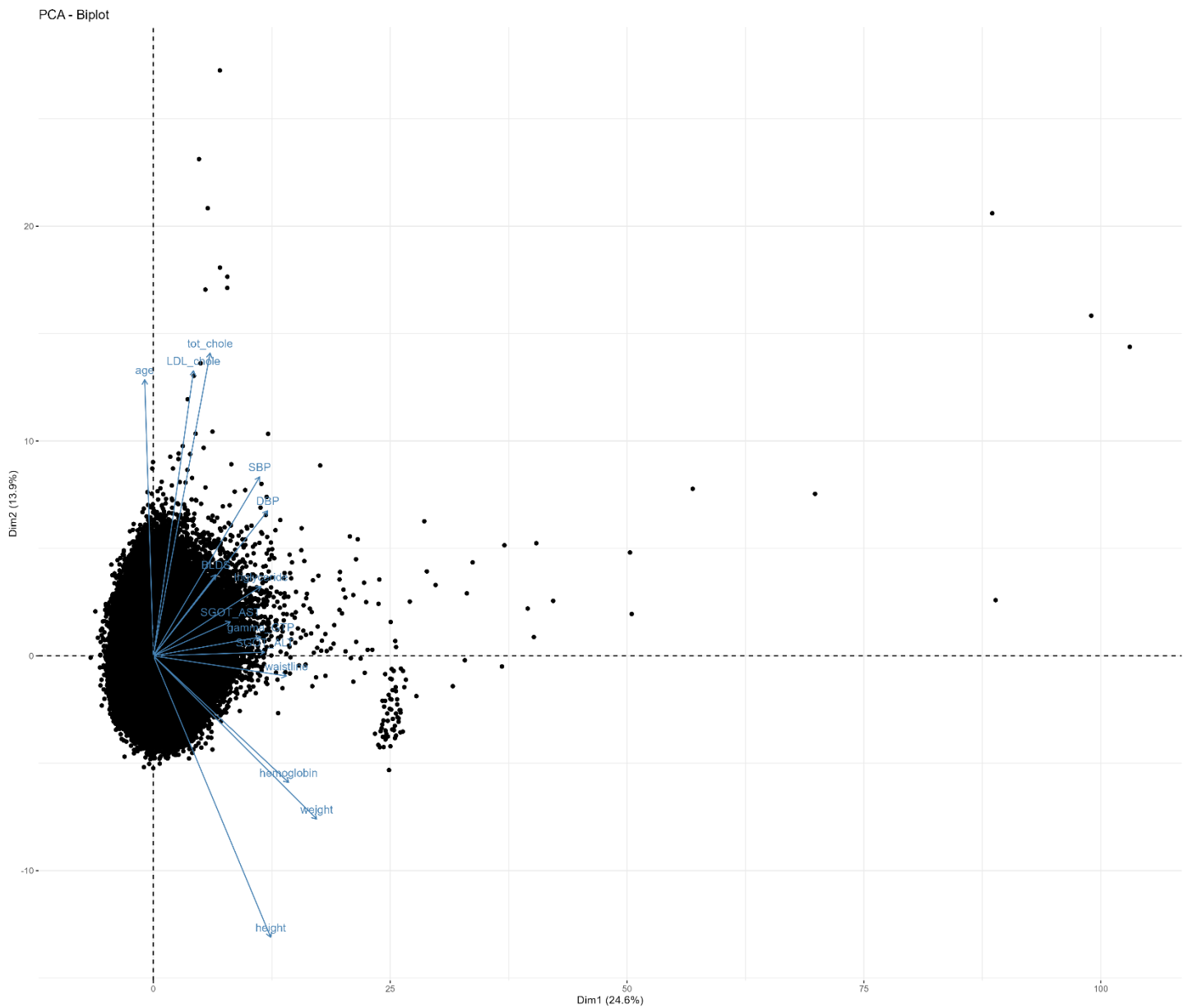
Variables (the 10 first)

	Dim.1	ctr	cos2	Dim.2	ctr	cos2	Dim.3	ctr	cos2
age	-0.042	0.050	0.002	0.582	17.427	0.338	0.464	11.906	0.215
height	0.560	9.107	0.314	-0.593	18.093	0.351	-0.297	4.879	0.088
weight	0.780	17.649	0.608	-0.344	6.086	0.118	-0.146	1.173	0.021
waistline	0.636	11.723	0.404	-0.042	0.092	0.002	0.070	0.270	0.005
DBP	0.546	8.654	0.298	0.306	4.815	0.094	0.277	4.242	0.077
SBP	0.508	7.481	0.258	0.377	7.325	0.142	0.395	8.611	0.156
BLDS	0.299	2.586	0.089	0.171	1.508	0.029	0.306	5.171	0.093
tot_chole	0.271	2.129	0.073	0.638	20.952	0.407	-0.681	25.691	0.464
LDL_chole	0.192	1.067	0.037	0.601	18.586	0.361	-0.719	28.578	0.517
triglyceride	0.514	7.662	0.264	0.146	1.101	0.021	0.013	0.010	0.000
NULL									



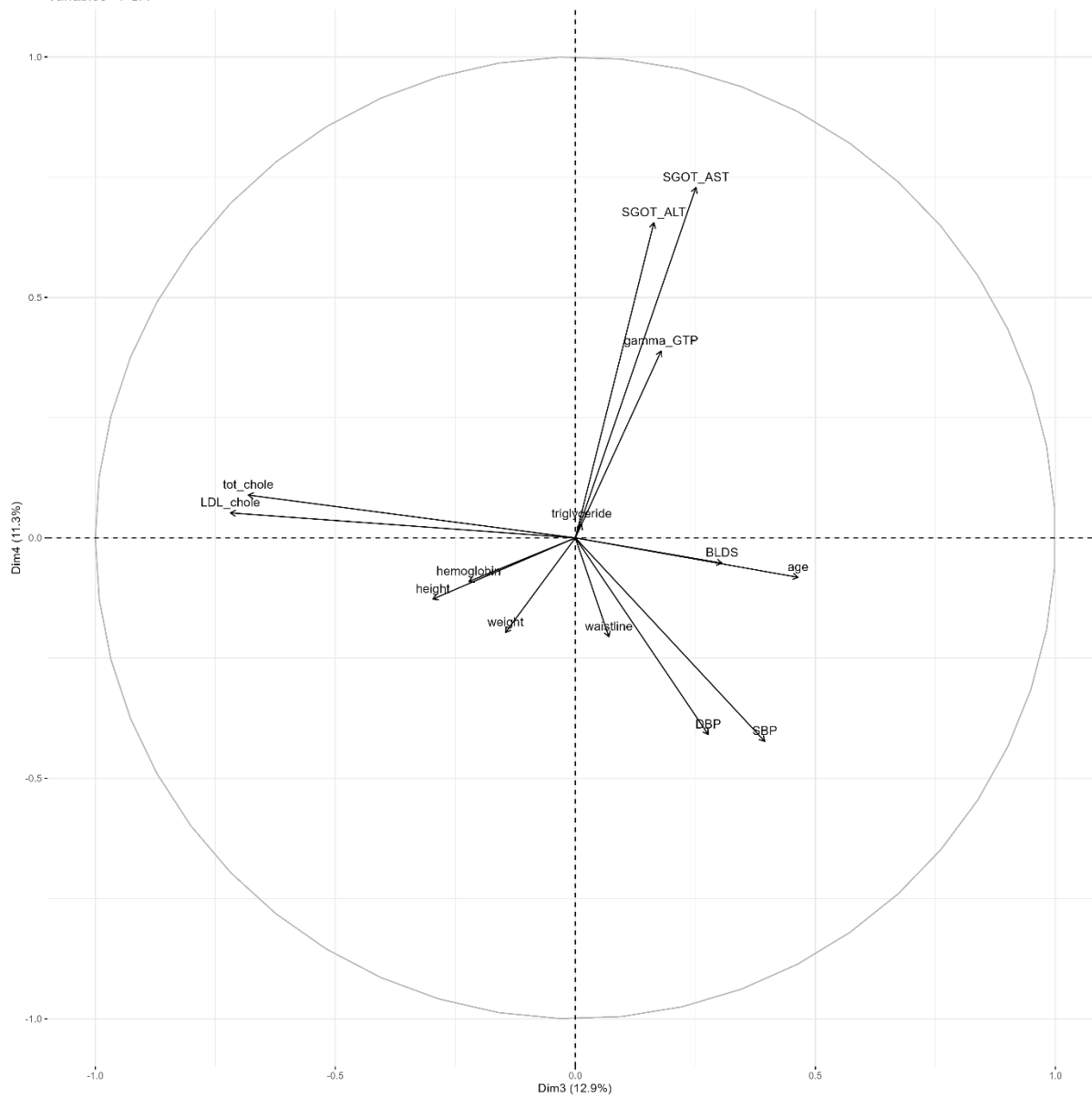
Надалі використовуватимемо перші 5 компонент (оскільки різниця між 5 і 6 компонентами – 1.0% не є значною у порівнянні з різницею між 4 і 5 компонентами).  
РСА значно скоротив кількість регресорів – 5 остаточних проти 14 початкових.

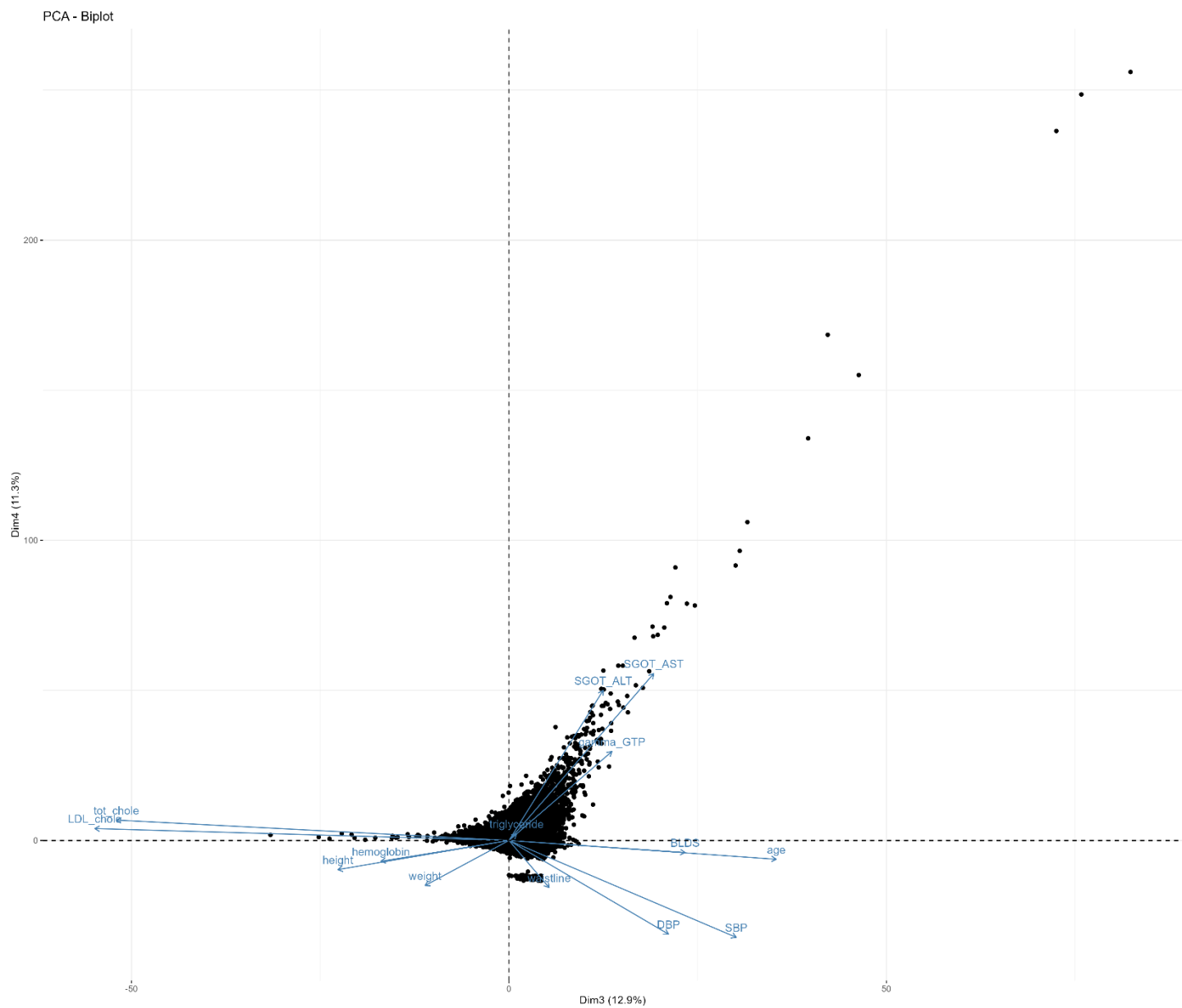




Перша компонента відповідає за ріст гемоглобіну в залежності від комплексного показника тиску, глюкози, тріглицерину, параметрів тіла, SGOT\_AST/ALT, gammga\_GTP. Друга компонента відповідає за ріст гемоглобіну в залежності від холестерину та віку.

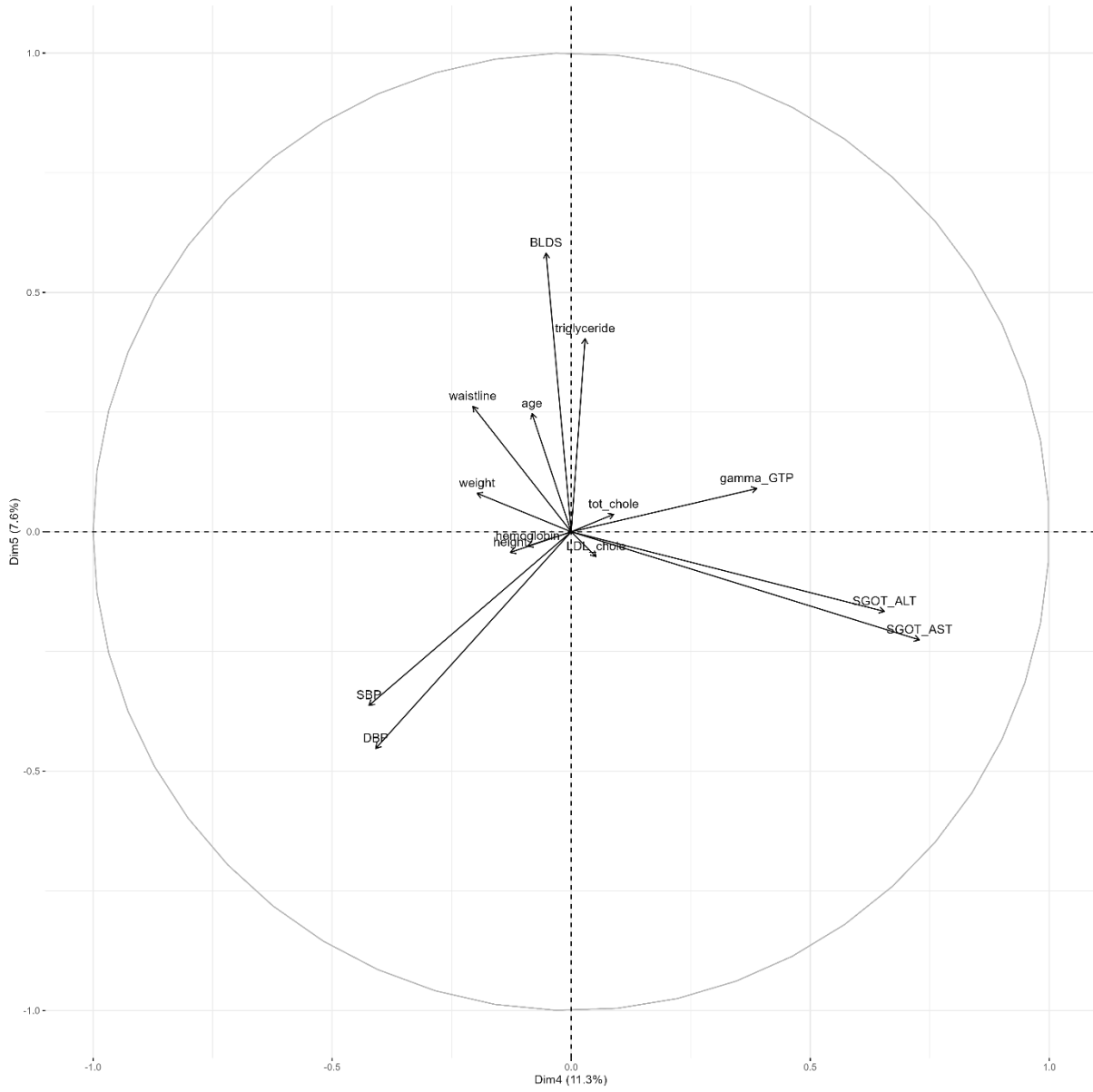
Variables - PCA

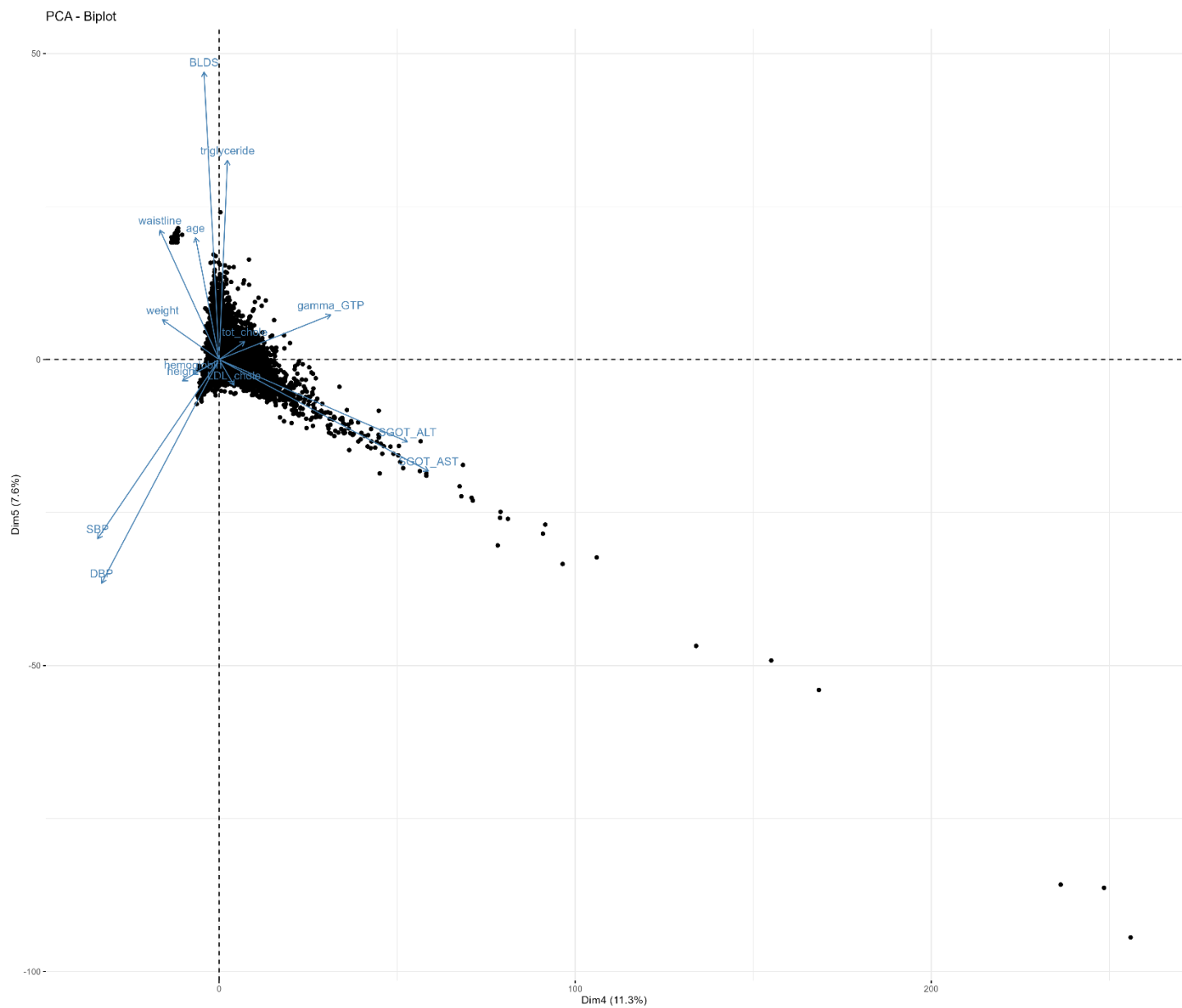




Третя компонента характеризує холестерин та різні види тиску.

Variables - PCA





П'ята компонента характеризує відношення між різними видами тиску і глюкозою.



## Generalized Additive Model

Спробуємо побудувати gam та дослідити її ефективність.

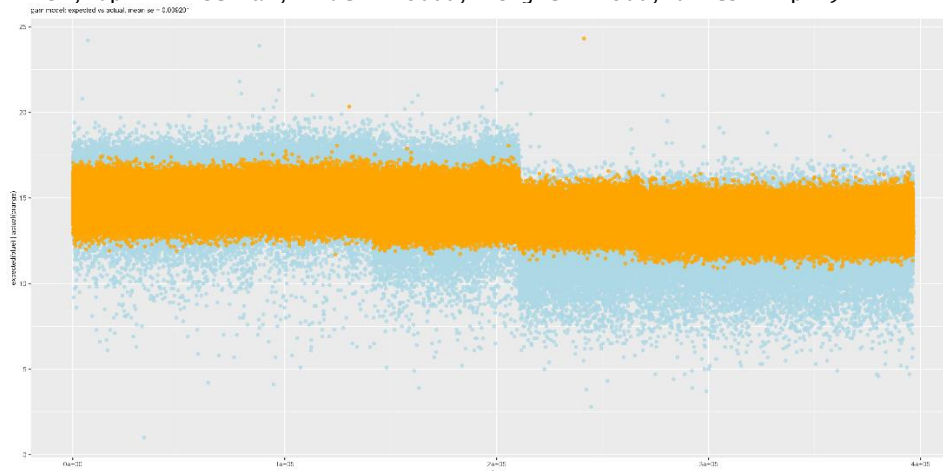
В якості незалежних параметрів модель прийматиме всі 5 головних компонент, виявлених під час PCA.

Розділимо дані на навчальну та тестову вибірки, перетворимо початкові характеристики у PC:

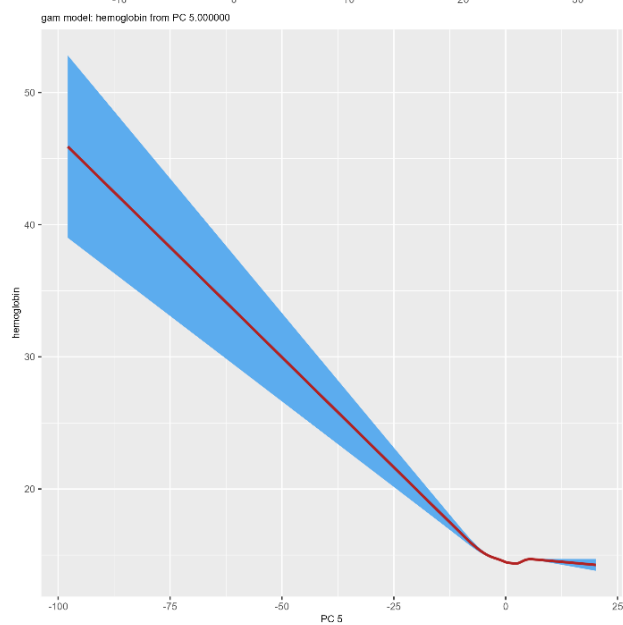
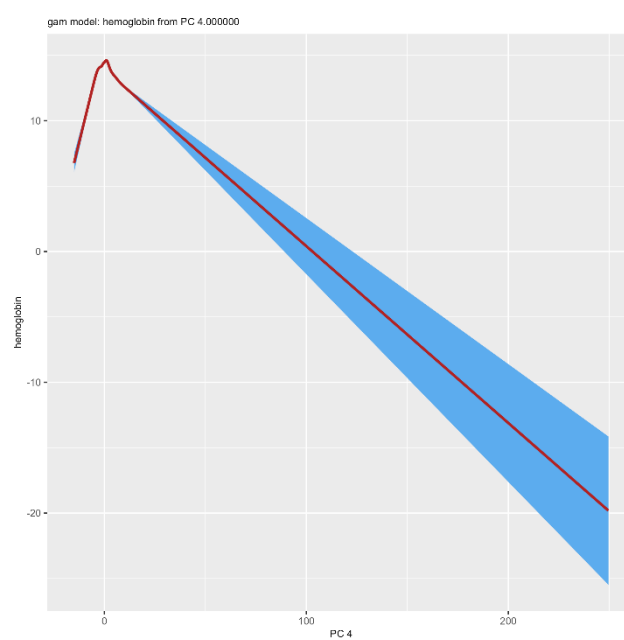
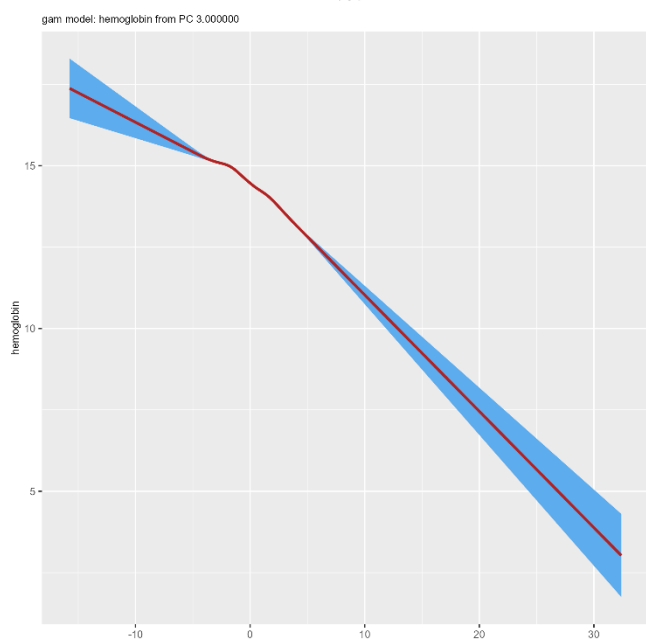
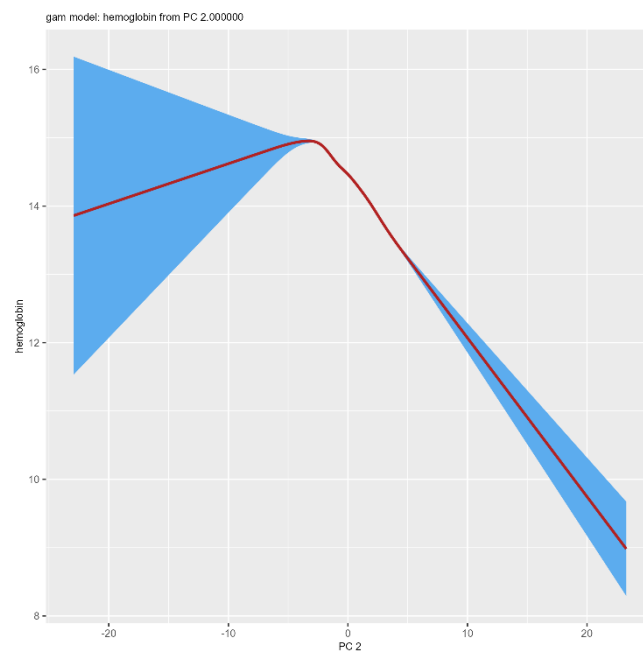
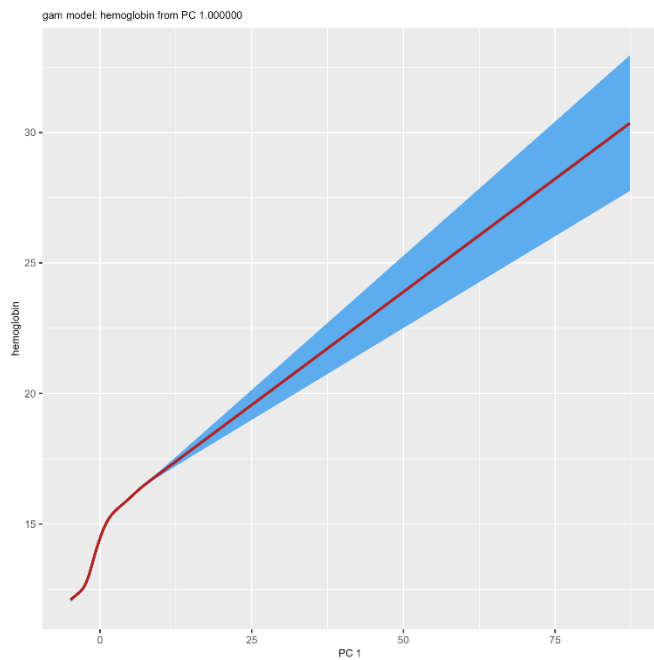
```
data_train_raw <- data %>% sample_n(nrow(data) * 0.6)
data_test_raw <- setdiff(data, data_train_raw)
data_train_raw <- data_train_raw %>% select(age, height, weight, waistline, DBP, SBP, BLDS,
                                             tot_chole, LDL_chole, triglyceride,
                                             SGOT_AST, SGOT_ALT, gamma_GTP, hemoglobin)
data_train_raw_nh <- data_train_raw %>% select(-hemoglobin)
data_test_raw <- data_test_raw %>% select(age, height, weight, waistline, DBP, SBP, BLDS,
                                             tot_chole, LDL_chole, triglyceride,
                                             SGOT_AST, SGOT_ALT, gamma_GTP, hemoglobin)
data_test_raw_nh <- data_test_raw %>% select(-hemoglobin)
data_train_pca <- PCA(data_train_raw_nh, graph = FALSE)
data_test_pca <- PCA(data_test_raw_nh, graph = FALSE)
data_train_pca_15 <- as_tibble(data_train_pca$ind$coord[, 1:5])
data_test_pca_15 <- as_tibble(data_test_pca$ind$coord[, 1:5])
data_train_pca_15 <- data_train_pca_15 %>% mutate(hemoglobin = data_train_raw$hemoglobin)
data_test_pca_15 <- data_test_pca_15 %>% mutate(hemoglobin = data_test_raw$hemoglobin)
```

Побудуємо модель на тренувальній вибірці, протестуємо на тестовій вибірці, порівняємо expected та actual значення:

```
gam_pca <- gam(hemoglobin ~ s(Dim.1) + s(Dim.2) + s(Dim.3) + s(Dim.4) + s(Dim.5),
               data = data_train_pca_15)
prediction_result <- predict.gam(gam_pca, newdata = data_test_pca_15, se.fit = TRUE)
compare_prediction_true <- tibble(
  index = seq(1, nrow(data_test_pca_15)),
  expected = data_test_pca_15$hemoglobin,
  actual = prediction_result$fit
)
score = mean(prediction_result$se.fit)
hist_ea <- ggplot(compare_prediction_true, aes(x = index))
hist_ea <- hist_ea +
  labs(x = "index", y = "expected(blue) | actual(orange)",
       title = sprintf("gam model: expected vs actual, mean se = %f", score)) +
  theme(axis.text = element_text(size = 8),
        axis.title = element_text(size = 8),
        plot.title = element_text(size = 8))
hist_ea <- hist_ea +
  geom_point(aes(y = expected), color = "lightblue", alpha = 0.75) +
  geom_point(aes(y = actual), color = "orange", alpha = 0.75)
ggsave("./output/Lab4/gam_expected_actual.png", plot = hist_ea,
        limitsize = FALSE, dpi = "retina", width = 5000, height = 2500, units = "px")
```



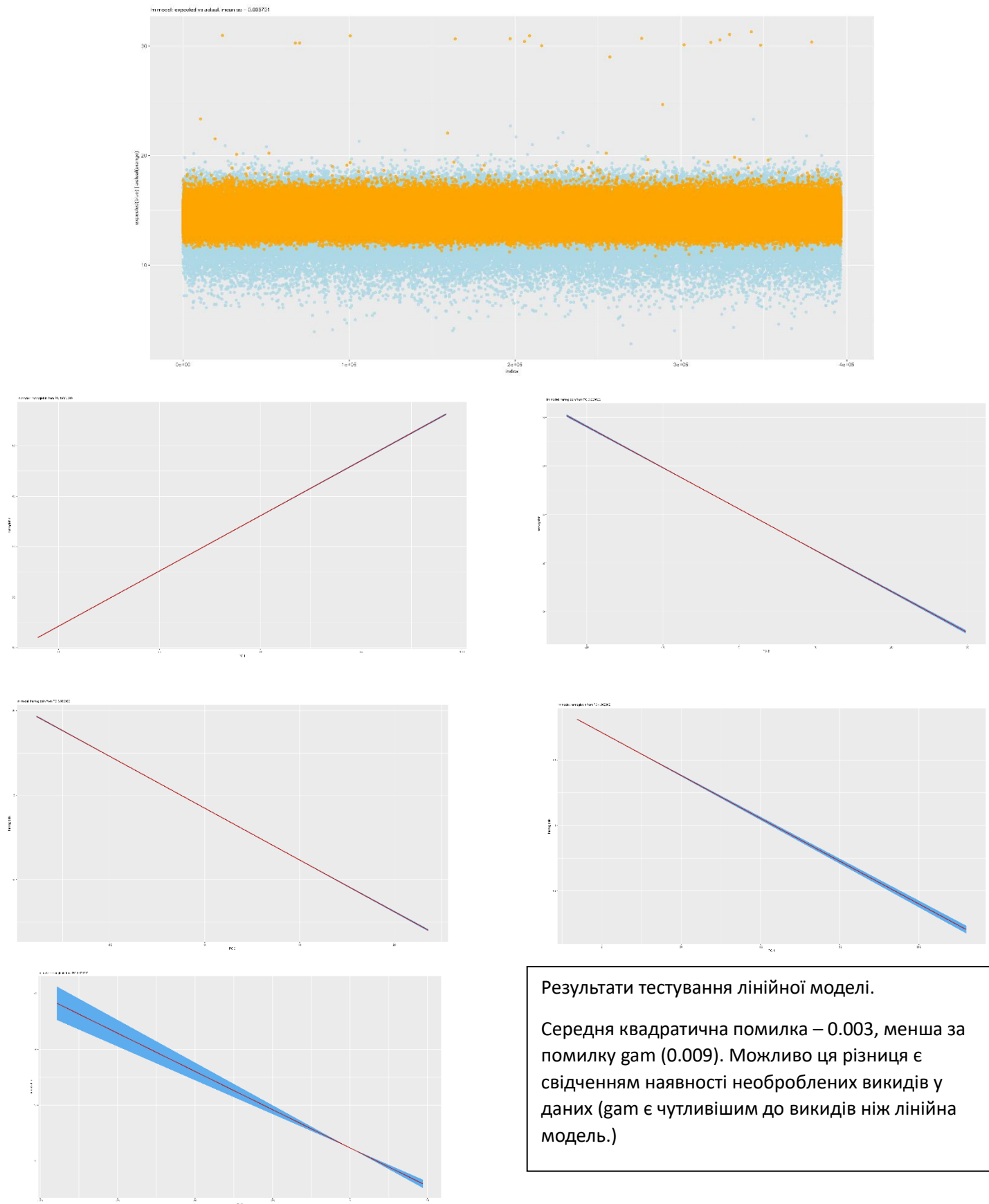
Середня квадратична помилка – 0.009



Залежності гемоглобіну від кожної  
ГОЛОВНОЇ КОМПОНЕНТИ згідно gam

## Linear model

Побудуємо лінійну модель по всім головним компонентам і протестуємо її подібно до тестування gam:



Результати тестування лінійної моделі.

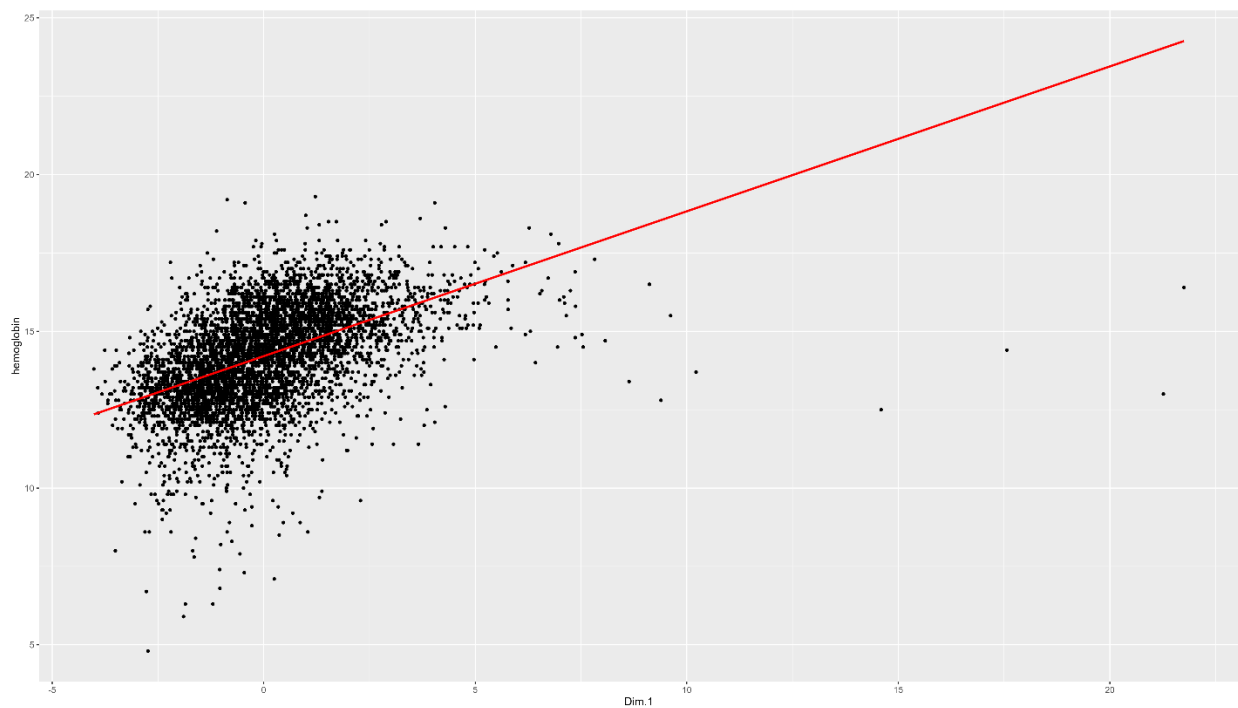
Середня квадратична помилка – 0.003, менша за помилку gam (0.009). Можливо ця різниця є свідченням наявності необроблених викидів у даних (gam є чутливішим до викидів ніж лінійна модель.)

	hemoglobin
Dim.1	0.438*** (0.003)
Dim.2	-0.170*** (0.002)
Dim.3	-0.306*** (0.002)
Dim.4	-0.131*** (0.007)
Dim.5	-0.027*** (0.005)
Constant	14.230*** (0.002)
Observations	594,770
Adjusted R <sup>2</sup>	0.336
<i>Note:</i> * p<0.1;   ** p<0.05;   *** p<0.01	

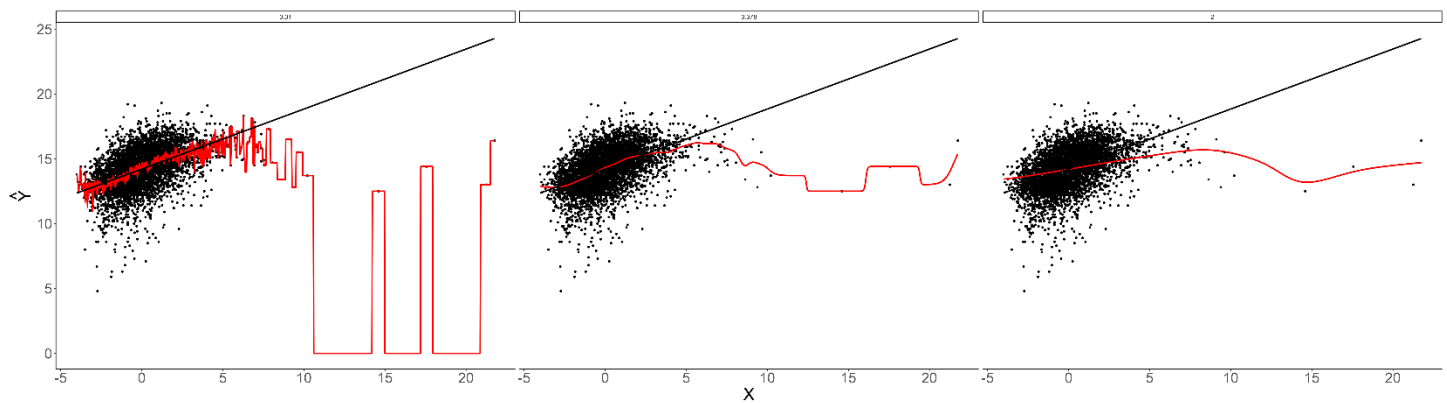
Всі 5 головних компонент лінійної моделі є статистично значущими.

## Надараї-Вотсон і локальна регресія

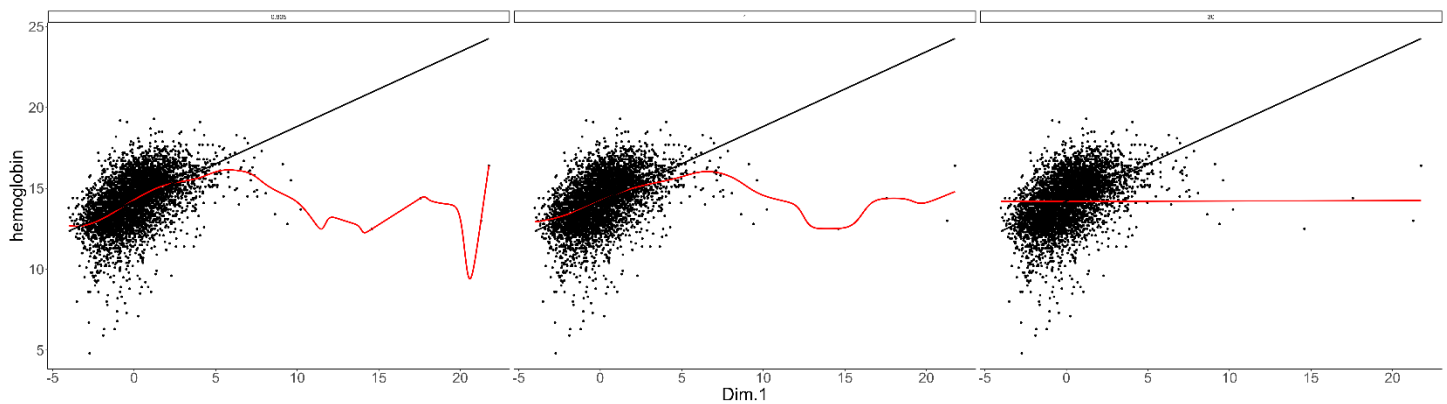
Оцінимо залежність гемоглобіну від першої головної компоненти.



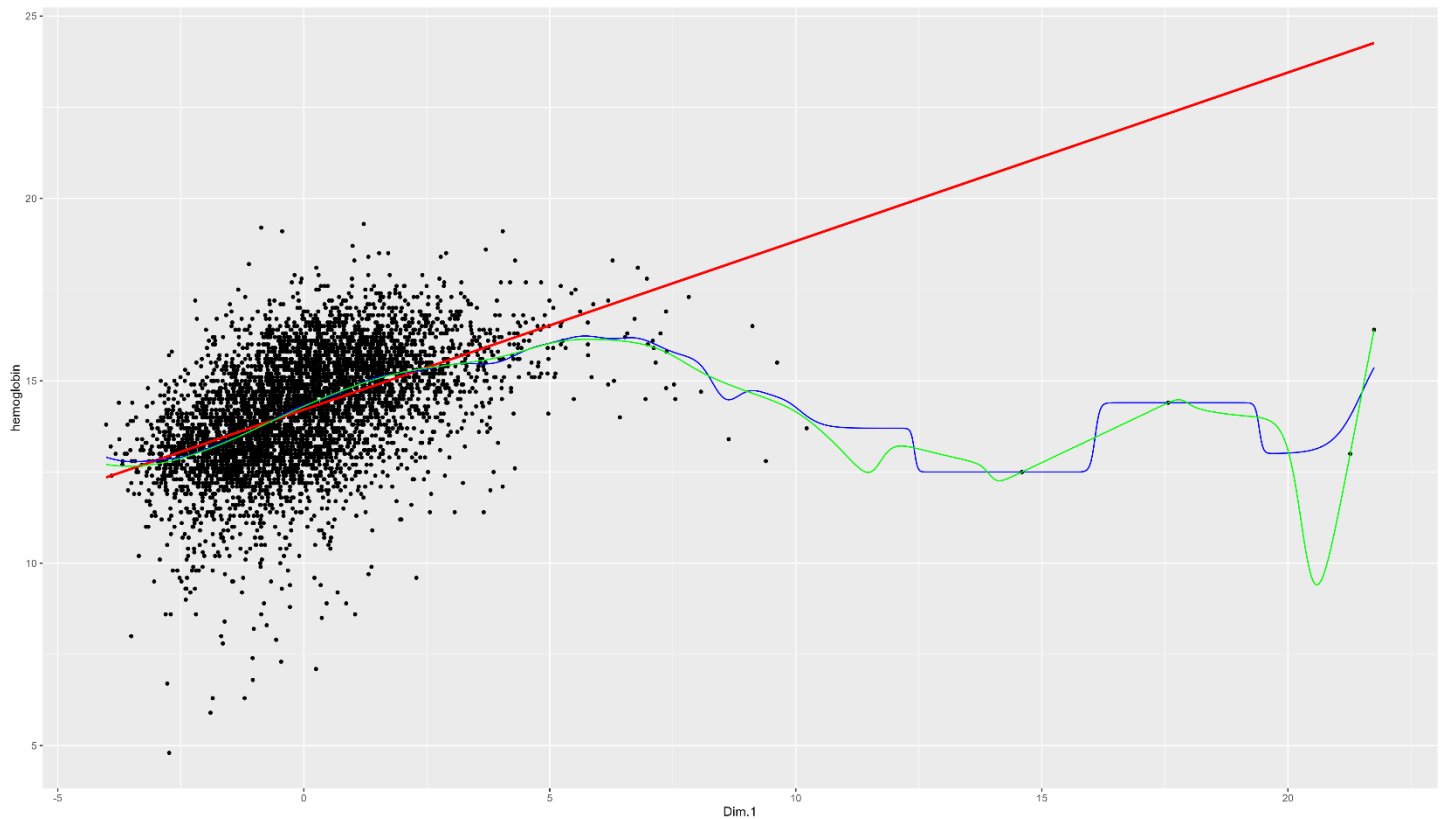
Лінійна модель Dim.1 -> hemoglobin



Надараї-Ватсон

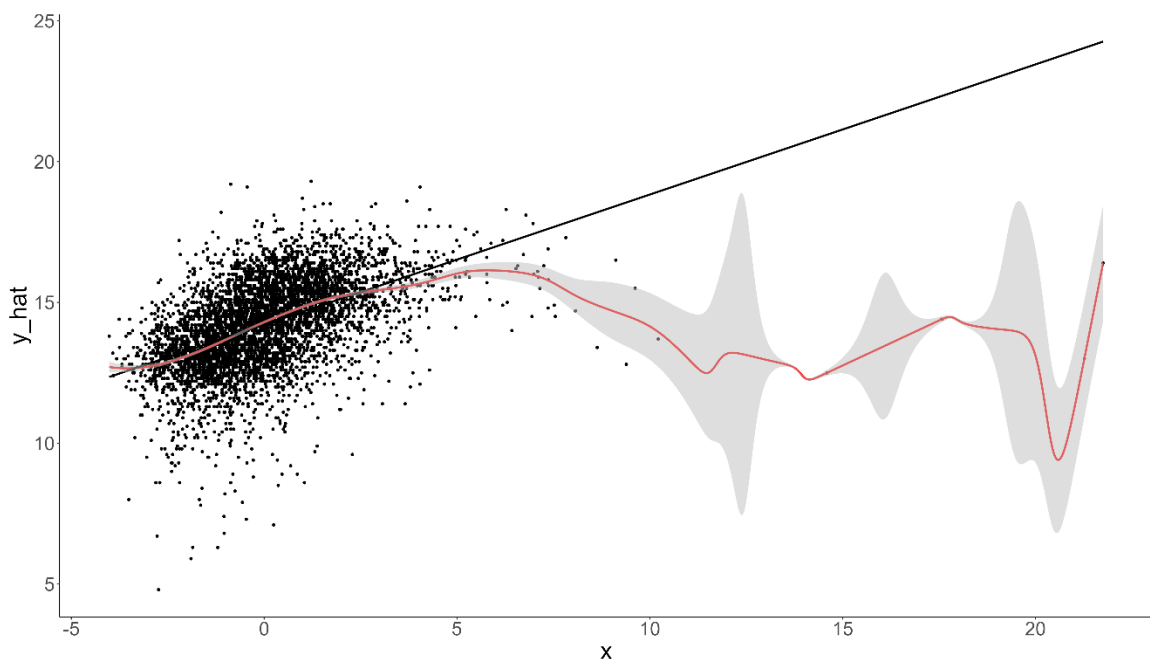


Локальна регресія



Порівняння лінійної регресії, Надарай-Ватсона, локальної регресії (синій – Надарай-Ватсон, зелений – локальна регресія)

Локальна регресія демонструє більшу гладкість, що потенційно може свідчити про більшу стійкість локальної регресії до перенавчання. При цьому моделі явно демонструють, чому варто більш уважно проводити дослідження викидів у даних :)



Локальна регресія і похибки

Похибка майже відсутня у місцях великого скупчення точок, збільшується з віддаленням від головного скупчення, досягає найбільших значень у викидах.