

Министерство науки и высшего образования Российской Федерации  
Санкт-Петербургский политехнический университет Петра Великого  
Физико-механический институт

Высшая школа теоретической механики и математической физики

Работа допущена к защите

Директор ВШТМиМФ,

д.ф.-м.н., чл.-корр. РАН

\_\_\_\_\_ А. М. Кривцов

«\_\_\_» \_\_\_\_\_ 2023 г.

## **ВЫПУСКНАЯ КВАЛИФИКАЦИОННАЯ РАБОТА БАКАЛАВРА**

**Исследование влияния различных факторов на риск возникновения и  
течение ишемической болезни сердца методами математической  
статистики.**

по направлению подготовки

01.03.03 «Механика и математическое моделирование»

Направленность

01.03.03\_03 Математическое моделирование процессов нефтегазодобычи

Выполнила

Студентка гр. 5030103/90301

Руководитель

Профессор ВШТМиМФ, д.ф.-м.н

А. Р. Курмакаева

В. А. Кузькин

**САНКТ-ПЕТЕРБУРГСКИЙ ПОЛИТЕХНИЧЕСКИЙ  
УНИВЕРСИТЕТ ПЕТРА ВЕЛИКОГО**  
**Физико-механический институт**  
**Высшая школа теоретической механики и математической физики**

УТВЕРЖДАЮ  
Директор ВШТМиМФ  
А. М. Кривцов  
«\_\_» \_\_\_\_\_ 20\_\_ г.

**ЗАДАНИЕ**

**на выполнение выпускной квалификационной работы**

студенту Курмакаевой Алсу Рашитовне, гр. 5030103/90301

1. Тема работы: Исследование влияния различных факторов на риск возникновения и течение ишемической болезни сердца методами математической статистики
2. Срок сдачи студентом законченной работы: 05.06.2023
3. Исходные данные по работе: справочная литература, актуальные научные публикации по теме работы.
4. Содержание работы (перечень подлежащих разработке вопросов): постановка задачи квалификационной работы, анализ подходов к исследованию влияния различных факторов, предварительный анализ и подготовка данных, вычисление эффекта влияния различных факторов, построение предиктивной модели.
5. Перечень графического материала (с указанием обязательных чертежей): не предусмотрено.
6. Консультанты по работе: отсутствуют
7. Дата выдачи задания: 27.02.2023

Руководитель ВКР: \_\_\_\_\_ В. А. Кузькин, д.ф.-м.н., профессор ВШТМиМФ

Задание принял к исполнению: 27.02.2023

Студент: \_\_\_\_\_ А. Р. Курмакаева

## **РЕФЕРАТ**

На 65 с., 16 рисунков, 8 таблицы, 3 приложения

**КЛЮЧЕВЫЕ СЛОВА:** СТАТИСТИЧЕСКИЕ ТЕСТЫ, ИШЕМИЧЕСКАЯ БОЛЕЗНЬ СЕРДЦА, ИБС, МАТЕМАТИЧЕСКАЯ МОДЕЛЬ.

Тема выпускной квалификационной работы: «Исследование влияния различных факторов на риск возникновения и течение ишемической болезни сердца методами математической статистики».

В данной работе проанализированы различные возможные факторы влияния на риск возникновения ишемической болезни сердца. Основной анализ производился с помощью статических тестов на основе критериев Манна-Уитни и ANOVA. Для выявления относительной степени влияния признаков была построена модель логистической регрессии. Для реализации вышеперечисленных источников был использован Python.

## **THE ABSTRACT**

65 pages, 16 pictures, 8 tables, 3 appendices

**KEY WORDS:** STATISTICAL TESTS, ISCHEMIC HEART DISEASE, CHD, MATHEMATICAL MODEL.

The topic of the graduate qualification work: "Research of different factors influence on the risk of coronary heart disease appearance and course by methods of mathematical statistics".

Different possible factors of influence on the risk of coronary heart disease were analyzed in this work. The main analysis was done using Mann-Whitney and ANOVA statistical tests. A logistic regression model was built to reveal the relative degree of influence of the traits. Python was used to implement the above sources.

## СОДЕРЖАНИЕ

ОСНОВНЫЕ ОБОЗНАЧЕНИЯ И СОКРАЩЕНИЯ .....	5
ВВЕДЕНИЕ .....	6
ГЛАВА 1. ОПИСАНИЕ ДАННЫХ.....	7
1.1. Исследование данных на наличие пропусков.....	8
1.2. Выявление корреляций между параметрами исследования.....	11
ГЛАВА 2.ВЫЯВЛЕНИЕ ФАКТОРОВ ВЛИЯНИЯ.....	16
2.1. Описание методов исследования .....	16
2.2. Описание результатов исследования.....	20
ГЛАВА 3. ОТНОСИТЕЛЬНАЯ ОЦЕНКА ФАКТОРОВ ВЛИЯНИЯ .....	27
3.1. Описание регрессионной модели.....	28
3.2. Описание модели решающих деревьев .....	30
3.3. Описание модели catboost.....	32
3.4. Исследование полученных результатов .....	32
ЗАКЛЮЧЕНИЕ.....	41
СПИСОК ИСПОЛЬЗОВАННОЙ ЛИТЕРАТУРЫ.....	43
ПРИЛОЖЕНИЕ .....	45

## **ОСНОВНЫЕ ОБОЗНАЧЕНИЯ И СОКРАЩЕНИЯ**

ТГ – триглицериды

Адсис – Артериальное давление систолическое

Адиас – Артериальное давление диастолическое

ОХС – общий холестерин

ТГ – триглицериды

АСТ – аспартатаминотрансфераза

PS – пульс

ХСЛПНП – холестерин липопротеинов низкой плотности

ХСЛПВП – холестерин липопротеинов высокой плотности

КА – коэффициент атерогенности

ВОЗ – Всемирная организация здравоохранения

## ВВЕДЕНИЕ

Ишемическая болезнь сердца (ИБС) – термин, предложенный Комитетом экспертов ВОЗ в 1962 г. Сегодня ИБС используется для описания острых и хронических заболеваний сердца, полученных в результате недостаточности снабжения миокардом кровью. По статистическим исследованиям Всемирной организации здравоохранения сердечно-сосудистые являются наиболее частой причиной смерти людей к моменту 2022 года, при этом большая часть этих заболеваний приходится на ишемическую болезнь сердца [15]. На данный момент проводится множество различных исследований, позволяющих установить взаимосвязь между различными факторами жизни человека и вероятностью его заболевания [14]. При этом при проведении подобных экспериментов важно учитывать особенности данных, на которых они будут проводиться.

В данной работе исследованы параметры влияния на развитие ишемической болезни сердца. Целью данной работы является выявление с помощью статистических тестов основных факторов, внесших наибольший вклад на возникновение ИБС, а также сравнительный анализ степени влияния каждого. В ходе работы были поставлены следующие задачи:

- 1) Провести статистическую обработку данных.
- 2) Провести статистические тесты и выявить основные факторы влияния.
- 3) Построить предиктивную математическую модель на основе машинного обучения, выявляющую риск возникновения ИБС для оценки степени влияния каждого фактора.
- 4) Сравнить полученные результаты.

## ГЛАВА 1. ОПИСАНИЕ ДАННЫХ

На сегодняшний день статистические исследования проводятся довольно часто для того, чтобы зафиксировать какой-либо эффект. При этом для различных сфер применения существуют разные подходы к исследованиям. Эти подходы отличаются друг от друга требованиями к подготовке теста, стандартами качества, а также предобработкой данных.

В случае исследования различных факторов влияния на риск возникновения ИБС будет использовано когортное (обсервационное) исследование. При проведении подобного анализа лица – объекты исследования распределяются в две группы в соответствии с наличием выбранного фактора. Далее исходя из истории их болезни проводят сравнение частоты развития клинического исхода.

Проведение статистического исследования включает в себя несколько этапов:

- 1) Планирование исследования
- 2) Сбор данных
- 3) Предобработка полученных материалов
- 4) Статистический анализ

В исследуемой выборке 125 наблюдений генеральной совокупности – 125 наблюдаемых пациентов (данные обезличенны).

Предобработка полученных материалов включает в себя следующие задачи:

1. Очистка данных: удаление дубликатов - повторяющихся записей из набора данных, обработка пропущенных значений, обработка выбросов – аномальных значений.
2. Преобразование переменных: преобразования переменных для улучшения их распределений, создание новых переменных на основе имеющихся для введения новых аспектов исследования.
3. Фильтрация данных: удаление ненужных переменных, удаление от нерепрезентативных переменных.

4. Кодирование данных: преобразование категориальных переменных в числовые значения, обработка текстовых данных.
5. Масштабирование данных: приведение значений переменных к определенному диапазону или стандартному формату.

При статистическом анализе необходимо в соответствии условиями сбора генеральной совокупности выбрать критерий сравнения. В данном исследовании будут рассматриваться независимые выборки с учетом многих факторов. Для выполнения описательной статистики используются следующие методы: вычисление медиан и интерквантильных интервалов, пропорций. При сравнении двух независимых групп по одному признаку применяется критерий Манна-Уитни. При одновременном анализе трех и более признаков задействуются такие методы как: регрессионный анализ, логистический регрессионный анализ, анализ древовидных диаграмм [2].

### **1.1 Исследование данных на наличие пропусков**

При проведении статистических исследований, а также построении моделей машинного обучения необходимо учесть пропуски, так как модель на них не сможет обучиться, а также удалить дубликаты. Необходимость удаления дубликатов связана с вычислением общей ошибки в подобных математических моделях, формула которой представлена ниже

$$\Delta_{ml} = \text{bias} + \text{variance} + \text{noise}, \quad (3)$$

где  $\Delta_{ml}$  – ошибка модели,  $\text{bias}$  – смещение,  $\text{variance}$  – разброс,  $\text{noise}$  – шум. Смещение показывает, насколько модель смещена от истинного значения исходной функции. Дубликаты вносят вклад в эту часть ошибки. Разброс в свою очередь отражает вариативность предсказания модели при обучении на различных тренировочных выборках. Шум – это ошибка, связанная с нерегулярностью и непредсказуемостью данных, их чаще всего сокращают с помощью поиска квантильных значений и удалением тех, что выходят за рамки 1 и 3 квантиля или 25 и 75 перцентиля.

Методы работы с пропусками:



- 1) Удаление пропусков
- 2) Заполнение пропусков его статистической характеристикой (модой, средним значением, медианой и тд.)
- 3) Заполнение новым значением, не встречавшимся ранее в выборке
- 4) Восстановление данных с помощью моделей машинного обучения.

Для того чтобы оценить количество пропусков была построена матрица разреженности (рис. 1.1):

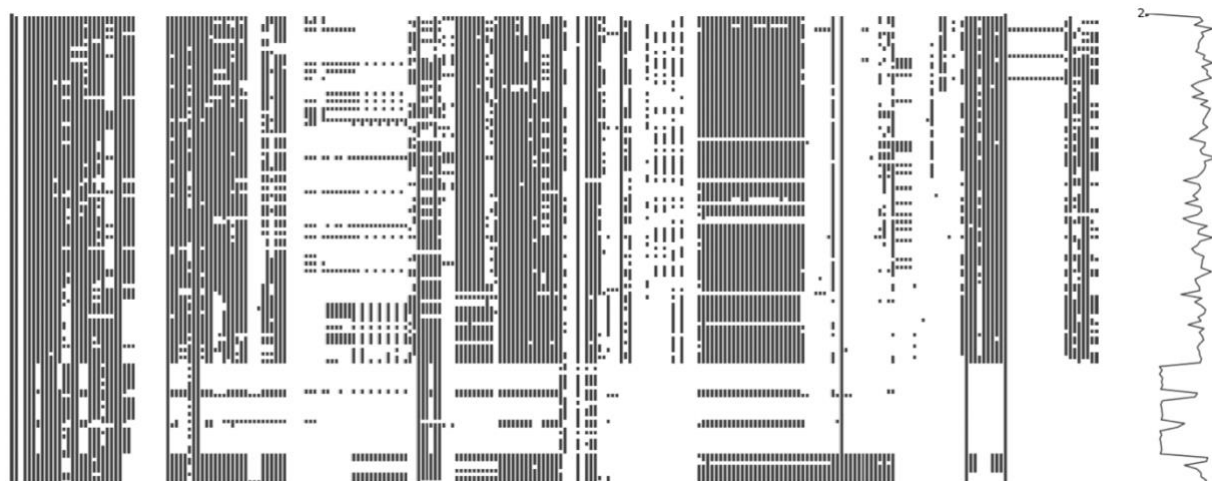


Рис. 1.1. Матрица разреженности исходного датасета.

Датасет состоит из 125 записей – 125 наблюдаемых пациентов и 221 признака – 221 исследуемых характеристик каждого пациента. Исходя из матрицы разреженности можно сделать вывод о том, что не все колонки являются информативными, большинство из них содержат большое число пропусков, также есть те, что и вовсе не заполнены. Для дальнейшего рассмотрения были выбраны следующие факторы: возраст, пол, номер визита, наследственность, адсист, аддиас, курение, рост, вес, ИМТ, менопауза, ОХС, ТГ, глюкоза, креатинин, билирубин, АСТ, артерия стентирования, aortic SP, конечные точки, FABP-4, ng/ml (кровь), FABP-4 ep1 mRNA, УЕЭ (эпикард. жир), FABP-4 мРНК в подкожном жире (УЕЭ), FABP4 rs16909192 (AA-1; AC-2; CC-3), FABP4 rs2290201 (GG-1; GA-2; AA-3), PS, Возраст развития СД, ХСЛПНП, ХСЛПВП, КА, эритроциты, гемоглобин, лейкоциты, лимфоциты, тромбоциты, моноциты.

Для сравнительного анализа из них подошли признаки, заполненность которых ниже 16% (рис. 1.2). Те признаки, заполненность которых была выше 16% были удалены. После описанного способа обработки был получен датасет, размерность которого: 125×18.

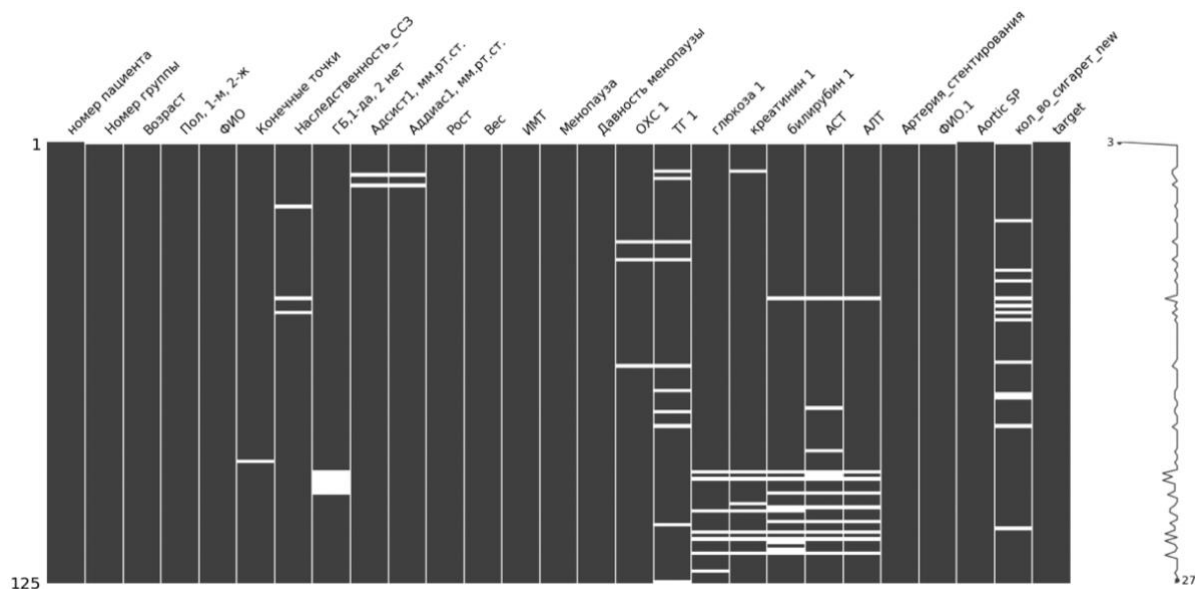


Рис. 1.2. Матрица разреженности нового датасета.

Далее были удалены все записи о пациенте, где в каком-либо описании его характеристик есть пропуски. Размерность полученного датасета 80×19, где 80 – количество записей или объектов исследования, а 19 – количество признаков.

Полученный первичный учетный документ для изучения влияния различных факторов на риск возникновения ИБС представлен в таблице 1.1:

№ п/п	Учетные признаки	Градация признака	Шифр
1	Возраст	28 – 78	-
2	Пол	женский мужской	2 1
3	Наследственность	да нет	1 2
4	Адсис1	100 – 180	-
5	Аддиас	60 – 100	-
6	Курение	да нет	1 2

7	Рост	150 – 186	-
8	Вес	54 – 124	-
9	ИМТ	18.0 – 46.1	-
10	Менопауза	да нет	1 2
11	ОХС	2.2 – 8.58	
12	ТГ	0.35 – 5.23	-
13	Глюкоза	4.3 – 9.0	-
14	Креатинин	0.042 – 0.720	-
15	Билирубин	4.55 – 39.0	-
16	АСТ	12 - 167	-
17	Артерия стентирования	1-13	-
18	Конечные точки	0 >0	0 1
19	ФАВР-4, ng/ml (кровь)	11.49 – 439.06	-
20	ФАВР-4 ері mRNA, УЕЭ (эпикард. жир)	0.078 – 1.908	-
21	ФАВР-4 мРНК в подкожном жире (УЕЭ)	0.38 – 3.24	-
22	ФАВР4 rs16909192 (AA-1; AC-2; CC-3)	GG GA AA	1 2 3
23	ФАВР4 rs2290201 (GG-1; GA-2; AA-3)	AA AC CC	1 2 3
24	PS	58 – 115	-
25	Возраст развития СД	0 – 73	-
26	ХСЛПНП	0.09 – 5.83	-
27	ХСЛПВП	0.49 – 4.09	-
28	КА	0.8 – 9.6	-
29	Эритроциты	3.32 – 6.5	-
30	Гемоглобин	91 – 180	-
31	Лейкоциты	3.4 – 17.8	-
32	Лимфоциты	0.8 – 3.7	-
33	Тромбоциты	108 – 390	-
34	Моноциты	0.2 – 5.0	-

Таблица 1.1. Генеральная совокупность.

## 1.2 Выявление корреляций между параметрами исследования

При проведении исследований будут применены разные методы, в том числе и линейные, поэтому необходимо провести проверку на наличие взаимосвязей между параметрами. Проверка будет произведена с помощью определения корреляции Пирсона, которая показывает степень линейной зависимости между двумя непрерывными значениями. Вычисляется по следующей формуле [1]:

$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}, \quad (4)$$

где  $X_i$  и  $Y_i$  – значения переменных  $X$  и  $Y$  для наблюдения  $i$ ,  $\bar{X}$  и  $\bar{Y}$  – средние значения,  $n$  – количество наблюдений.

Результат проверки представлен на рис. 1.3:

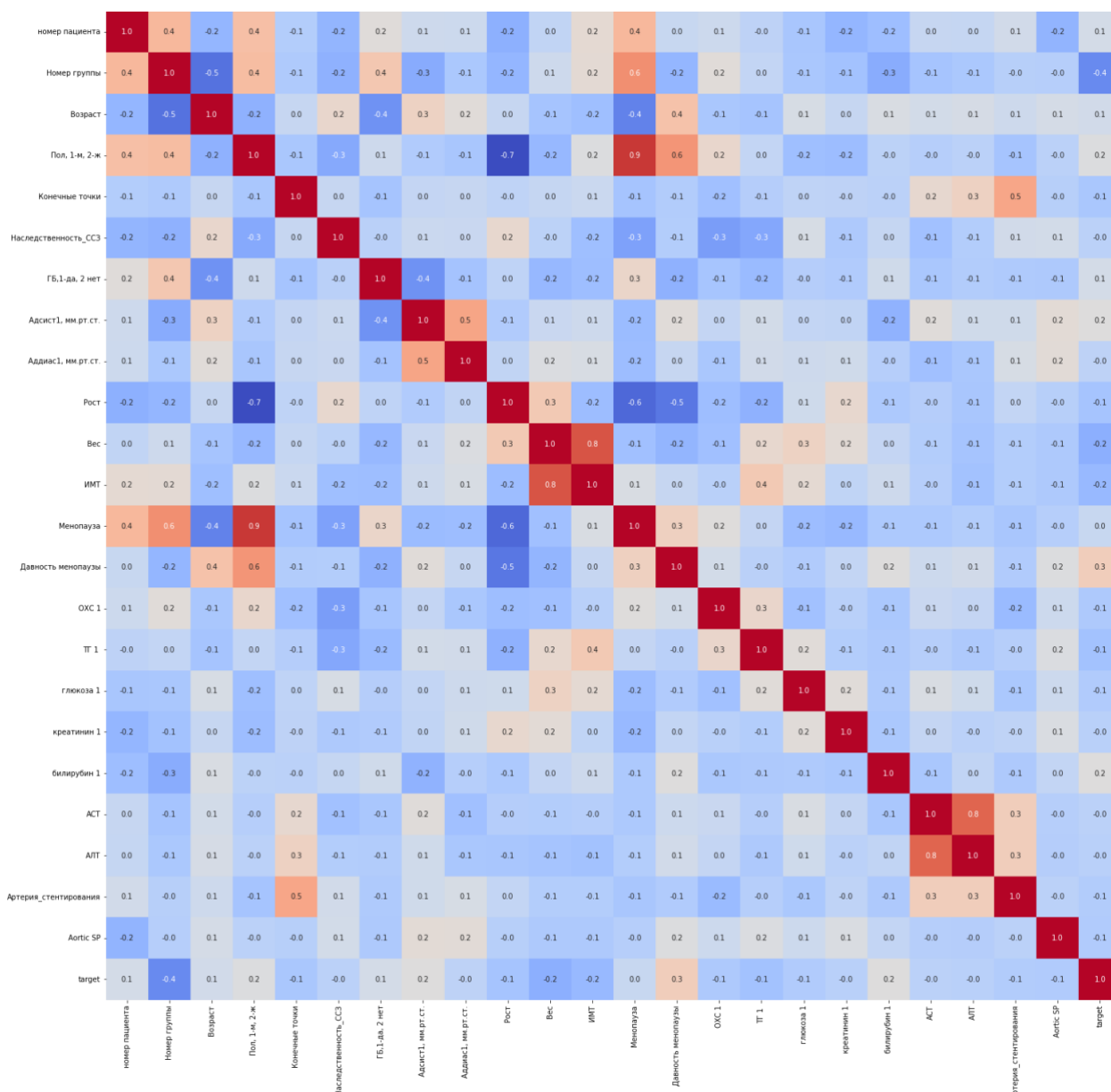


Рис. 1.3. Матрица корреляции исследуемых признаков.

Исходя из графика можно сделать следующие выводы: общая матрица корреляций показывает отсутствие линейной зависимости между большей частью рассматриваемых факторов влияния, однако есть также сильно скоррелированные между собой факторы:

- 1) Менопауза, пол, давность менопаузы (при попарном рассмотрении)
- 2) Стаж курения, курение
- 3) Вес, ИМТ
- 4) Рост, пол

При проведении статистических тестов будут рассмотрены все факторы, но при проведении исследования с помощью линейной логистической модели, из

каждой выделенной группы высокой корреляции будет оставлен лишь один признак (пол, курение, ИМТ).

Избавление от высоко скоррелированных признаков помогает избежать проблемы мультиколлинеарности, которая в свою очередь усложняет модель, также удаление этих признаков поможет сделать модель более интерпретируемой и лучшей с точки зрения производительности.

Для дополнительного изучения для группы, с наличием ИБС, были также построены матрицы корреляций зависимости различных факторов с FABP4 [10] в крови и в эпидуральном жире (рис. 1.4, рис. 1.5):

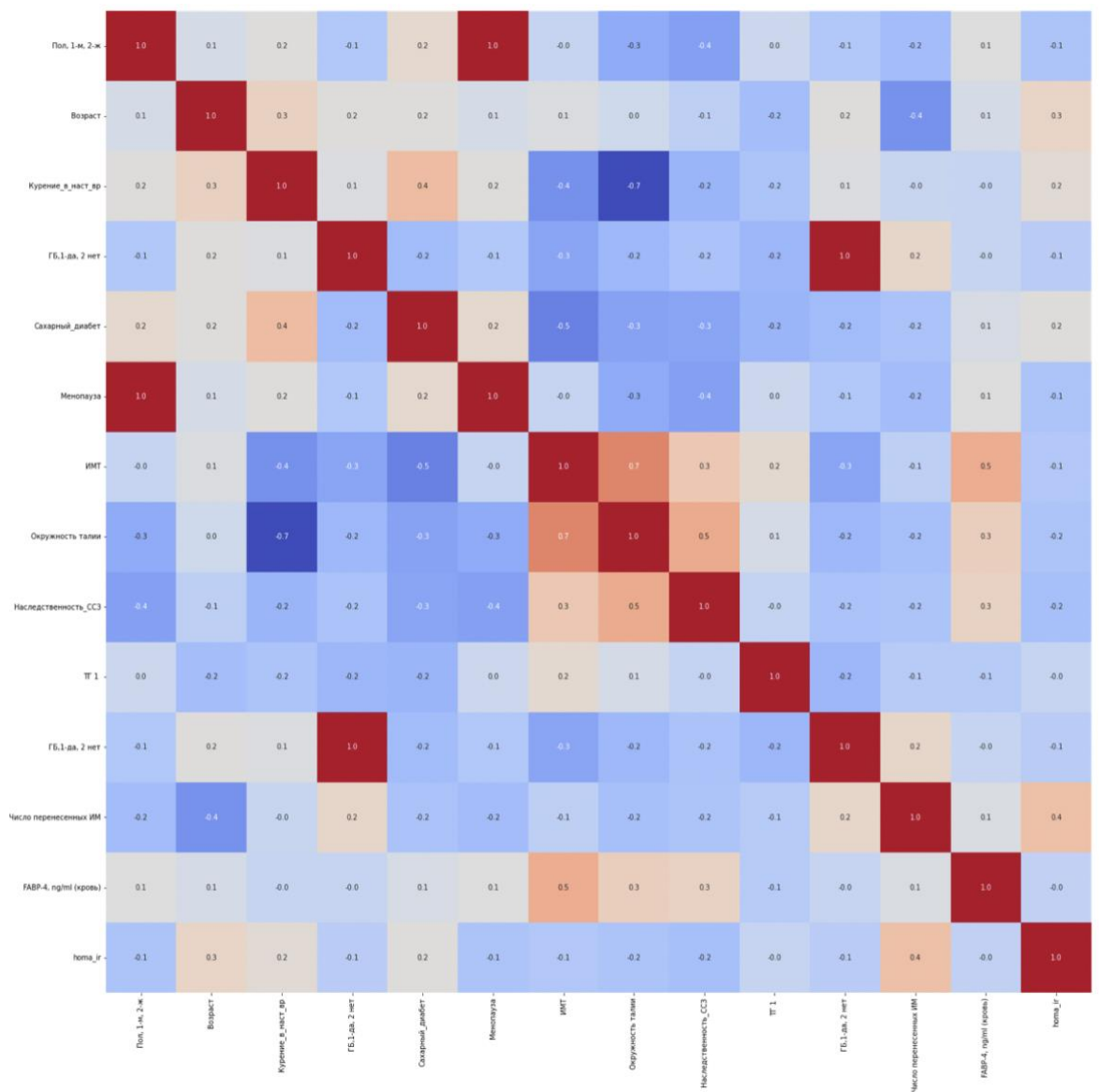


Рис. 1.4. Матрица корреляции для FABP4 в крови.

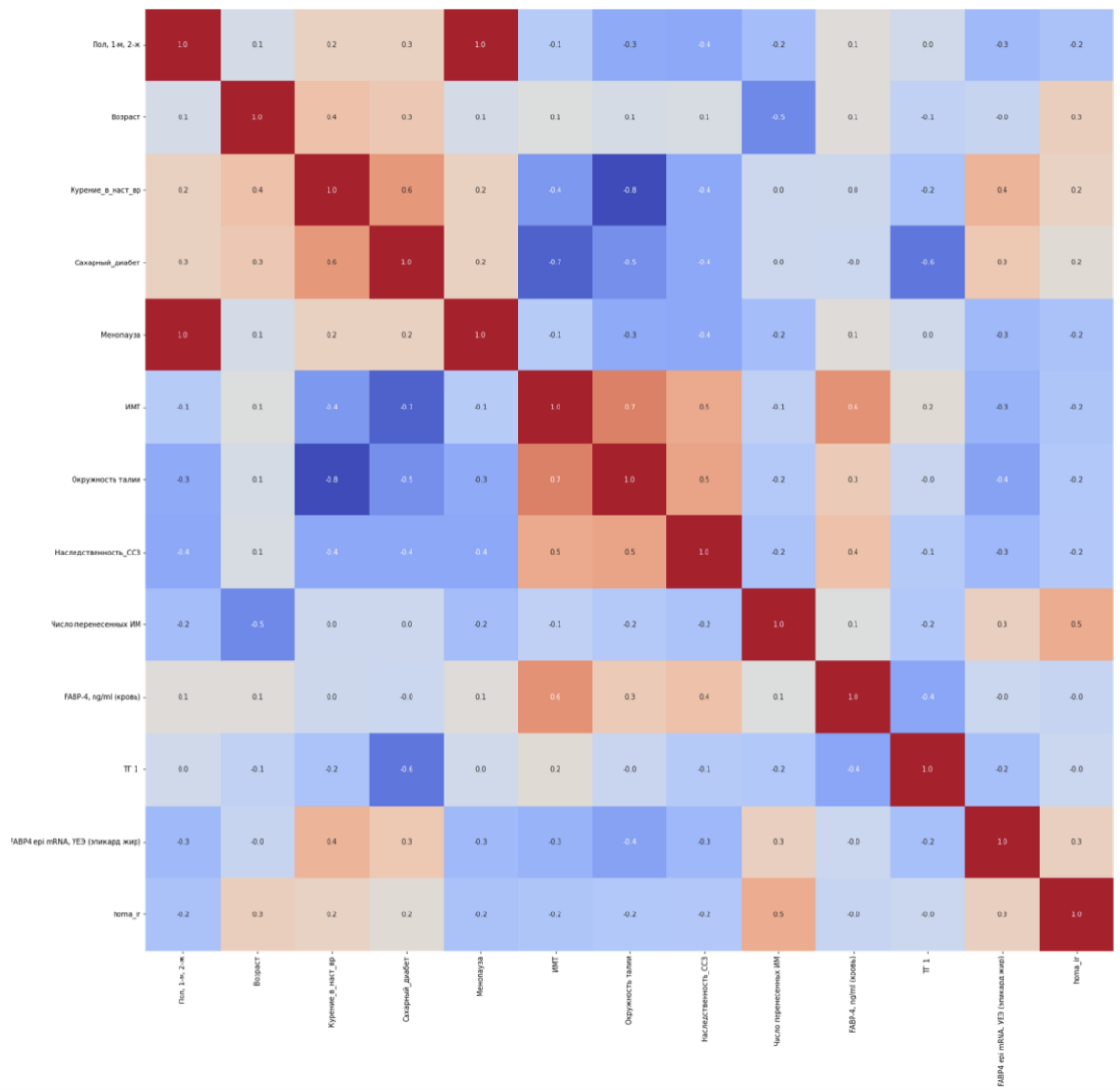


Рис. 1.5. Матрица корреляции для FABP4 в эпидуральном жире.

## ГЛАВА 2. ВЫЯВЛЕНИЕ ФАКТОРОВ ВЛИЯНИЯ

### 2.1 Описание методов исследования

На этапе планирования исследования необходимо поставить с целью и задачи эксперимента. После определиться с генеральной совокупностью, а также с выборочной статистической совокупностью с учетом ее репрезентативности. Для определения генеральной совокупности используется расчет необходимого числа наблюдений, допустимой ошибки (p-value).

Формула вычисления количества необходимых наблюдений:

$$n = \frac{Nt^2pq}{N\Delta^2 + t^2pq}, \quad (1)$$

где  $\Delta$  – предельная ошибка показателя,  $p$  – величина исследуемого показателя,  $q = (100 - p)$  или  $(1 - p)$ ,  $N$  – число наблюдений в генеральной совокупности,  $t$  – коэффициент, показывающий вероятность достоверности полученного результат (Обычно берется равным 2).

Формула вычисления предельной ошибки показателя:

$$\Delta = t \sqrt{\frac{pq}{n}}, \quad (2)$$

Расчет размера подвыборки нужен при моделировании метода Bootstrap [5]. Использование этого метода помогает не учитывать ограничение по распределению исследуемого параметра, которое необходимо брать во внимание при выборе статистического теста. При проверке данных на нормальность с помощью теста Шапиро-Уилка [13] было выявлено, что не все данные проходят эту проверку (таблица 2.1). Поэтому использование метода Bootstrap – необходимость для соблюдения ограничений для использования выбранных статистических критериев.



№ п/п	Учетные признаки	Нормальность распределения признаки
1	Возраст	ненормальное
2	Пол	-
4	Наследственность	-
5	Адсис	ненормальное
6	Аддиас	нормальное
7	Курение	-
8	Рост	нормальное
9	Вес	нормальное
10	ИМТ	нормальное
11	Менопауза	-
12	ОХС	нормальное
13	ТГ	ненормальное
14	Глюкоза	ненормальное
15	Креатинин	ненормальное
16	Билирубин	ненормальное
17	АСТ	ненормальное
18	Артерия стентирования	-
19	Конечные точки	-
20	FABP-4, ng/ml (кровь)	ненормальное
21	FABP-4 epi mRNA, УЕЭ (эпикард. жир)	ненормальное
22	FABP-4 мРНК в подкожном жире (УЕЭ)	ненормальное
23	FABP4 rs16909192 (AA-1; AC-2; CC-3)	-
24	FABP4 rs2290201 (GG-1; GA-2; AA-3)	-
25	PS	ненормальное
26	Возраст развития СД	ненормальное
27	ХСЛПНП	ненормальное
28	ХСЛПВП	ненормальное
29	КА	ненормальное
30	Эритроциты	ненормальное
31	Гемоглобин	ненормальное
32	Лейкоциты	ненормальное
33	Лимфоциты	ненормальное
34	Тромбоциты	ненормальное
35	Моноциты	ненормальное

Таблица 2.1. Результат проверки нормальности распределени.

Основные шаги при моделировании метода Bootstrap описаны ниже [3]:

1. Исходные данные.

На этом этапе формируется генеральная совокупность наблюдений о факторах риска ИБС и соответствующих результатов.

2. Ресемплирование.

При ресемплировании из исходной выборки выбирается  $n$  случайных наблюдений с возвращением. По полученной выборке находится среднее значение. Этот процесс повторяется от 1000 до 10000 раз. При таком моделировании процесса верна центральная предельная теорема, по условиям которой полученное распределение выборочных средних является нормальным.

3. Анализ распределения.

На этапе анализа производится оценка статистик с помощью различных критериев.

При рассмотрении количественной переменной в качестве генеральной совокупности принимаются сами значения фактора, а при анализе категориальной переменной используется вероятность наблюдения ИБС у объекта исследования.

Критерий для статистического исследования численных факторов – критерий Манна-Уитни [4]. Необходимые условия применимости критерия – независимость выборки, а также нормальное распределение исследуемого параметра, выполнены. Основные шаги использования данного критерия:

1. Определение гипотез.

Нулевая гипотеза – распределение риска ишемической болезни сердца одинаково для обеих групп.

$$H_0: F_1(x) = F_2(x), \quad (5)$$

где  $F_1(x)$  и  $F_2(x)$  – функции распределения риска для исследуемого признака соответственно.

Альтернативная гипотеза – распределение риска ишемической болезни сердца различно для двух групп.

$$H_1: F_1(x) \neq F_2(x), \quad (6)$$

## 2. Ранжирование.

Проведение ранжирования по риску ишемической болезни сердца в каждой группе. Присвоения рангов наблюдения в обеих группах ( $U_1$  и  $U_2$ ).

## 3. Распределение статистики $U$ .

Расчет суммы рангов для каждой группы ( $U_1$  и  $U_2$ ) [9]. Сравнение значений для  $U_1$  и  $U_2$ .

## 4. Распределение статистики $U$ .

При достаточно большой выборке, статистика  $U$  приближается к нормальному распределению. При малых выборках используется таблица критических значений для проверки статистической значимости.

## 5. Проверка статистического критерия.

Сравниваем значения статистики  $U$  с критическим значением для заданного уровня значимости –  $p_{\text{value}} = 0.01$ . Если значения статистики  $U$  меньше  $p_{\text{value}}$ , отвергаем нулевую гипотезу и считаем различия в распределении риска между группами статистически значимыми.

Критерий для исследования категориальных факторов – ANOVA [6]. Этот критерий подходит для описания и сравнения двух и более групп. Необходимые условия применимости критерия проверены. Принцип работы ANOVA заключается в выполнении следующих шагов:

### 1. Определение гипотез.

Этап аналогичен этапу 1 при описании предыдущего критерия. Основное различие в объектах выборки, если ранее рассматривались средние значения подвыборок из самих объектов начальной выборки, то на этом этапе рассматривается вероятность возникновения ИБС.

### 2. Разбиение данных на группы.

Исходные данные разбиваются на две или более группы в зависимости от исследуемого фактора.

### 3. Вычисление статистик.

По полученной выборки из вероятностей для каждой группы вычисляется оценки дисперсии от общего среднего значения и от средних внутри группы:

$$MS_B = \sum \frac{(\bar{X}_j - \bar{X}_G)^2}{k - 1} n, \quad (7)$$

$$MS_W = \frac{s_1^2 + s_2^2 + \dots + s_k^2}{k}, \quad (8)$$

где  $MS_B$  – оценка общей дисперсии по разбросу между группами,  $\bar{X}_j$  – среднее группы  $j$ ,  $\bar{X}_G$  – общее среднее,  $k$  – число групп,  $n$  – размер группы,  $MS_W$  – оценка общей дисперсии по разбросу внутри группы,  $s_1^2 + s_2^2 + \dots + s_k^2$  – сумма квадратов стандартных отклонений внутри группы.

Далее идет расчет F-статистики:

$$F = \frac{MS_B}{MS_W}, \quad (9)$$

### 4. Проверка статистического критерия.

Проводится аналогично предыдущему тесту.

## 2.2 Описание результатов исследования

Для описания результатов исследования были построены графики «Ящички с усами» для группы объектов с ИБС – группа 1 и для группы объектов без ИБС – группа 0. Пример построенных графиков представлен на рисунке 2.1, остальные графике представлены в приложении:

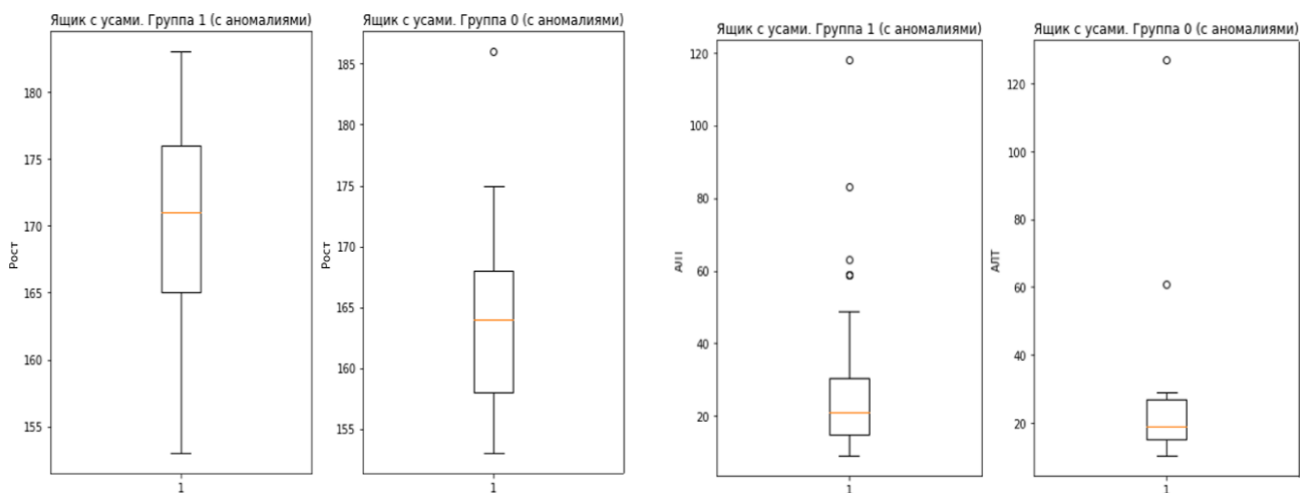


Рис. 2.1. Ящики с усами для количественных признаков.

Полученные графики помогают сравнить относительное расположение медиан (отмечены оранжевой линией на графике) исследуемого признака, а также указывают на 1 и 3 квантиль (границы «усов» графика), помимо этого на графике показаны аномальные значение – выпуклые точки.

При детальном рассмотрении графика видно, что есть факторы, медианы распределения которых значительно отличаются друг от друга, а также те, значения которых не сильно различимы.

В группу сильных различий попали такие факторы как: возраст, Адсист1, рост, вес, пол, давность менопаузы, ОХС, креатинин, билирубин, АСТ. Именно эти факторы являются первыми кандидатами на роль тех факторов, что сильно влияют на риск возникновения ишемической болезни сердца.

Результаты оценки влияния каждого признака с помощью статистического теста с заданной статистической значимостью  $p\text{-value} = 0.01$  представлена в таблице 2.2, где alpha – полученная статистическая значимость, которая в последствии сравнивалась с  $p\text{-value}$ :

Учетные признаки	alpha
Возраст	0.000000
Пол	0.000000
Номер визита	0.000000
Наследственность	0.000000
Адсист	0.000000

Аддиас	0.186201
Курение	0.000000
Рост	0.000000
Вес	0.000000
ИМТ	0.000000
Менопауза	0.000000
ОХС	0.000020
ТГ	0.000010
Глюкоза	0.740582
Креатинин	0.001923
Билирубин	0.000010
АСТ	0.037024
Артерия стентирования	0.000000
Aortic SP	0.000000
Конечные точки	0.000000
FABP-4, ng/ml (кровь)	0.003354
FABP-4 epi mRNA, УЕЭ (эпикард. жир)	0.000000
FABP-4 мРНК в подкожном жире (УЕЭ)	0.000000
FABP4 rs16909192 (AA-1; AC-2; CC-3)	0.000000
FABP4 rs2290201 (GG-1; GA-2; AA-3)	0.000000
PS	0.000001
Возраст развития СД	0.000000
ХСЛПНП	0.853627
ХСЛПВП	0.190666
КА	0.000000
Эритроциты	0.001039
Гемоглобин	0.000000
Лейкоциты	0.001332
Лимфоциты	0.297920
Тромбоциты	0.069957
Моноциты	0.000000

Таблица 2.2. Значения p-value.

Исходя из полученных результатов в группу влияния вошли признаки, статистическая значимость alpha которых меньше значения 0.01: возраст, пол, наследственность, адсист, курение, рост, вес, ИМТ, менопауза, ОХС, ТГ, креатинин, билирубин, артерия стенирования, aortic SP, FABP-4, ng/ml (кровь),

FABP-4 epi mRNA, УЕЭ (эпикард. жир), FABP-4 мРНК в подкожном жире (УЕЭ), FABP4 rs16909192 (AA-1; AC-2; CC-3), FABP4 rs2290201 (GG-1; GA-2; AA-3), PS, Возраст развития СД, КА, эритроциты, гемоглобин, лейкоциты, моноциты.

Для факторов, значения p\_value которых больше 0.01 нет достаточных оснований для отклонения нулевой гипотезы.

Аналогично для группы людей с ИБС были получены таблицы влияния различных факторов на значения FABP4 в крови и эпидуральном жире:

	Alpha - FABP-4, ng/ml (кровь)	Alpha - FABP4 epi mRNA, УЕЭ (эпикард жир)
Пол	0.000000	0.081698
Возраст	0.000000	0.000599
Курение	0.000000	0.000000
Количество инфарктов миокарда	0.000000	0.000000
Артериальная гипертензия	0.000000	0.024372
Сахарный диабет	0.000000	0.003124
Индекс НОМА-IR	0.000000	0.000000
Менопауза	0.000000	0.060416
ИМТ_1	0.000000	0.000000
ИМТ_2	0.000000	0.000000
Обхват талии для мужчин	0.000000	0.021807
Наследственность	0.000000	0.000000
Триглицериды	0.163332	0.000000

Таблица 2.3. Значения p-value для FABP4.

Помимо этого, для более детального исследования и рассмотрения как могут повлиять на риск возникновения ишемической болезни сердца два

фактора, прошедших проверку по критическому значению, были построены дополнительные графики исследования зависимости относительно пола.

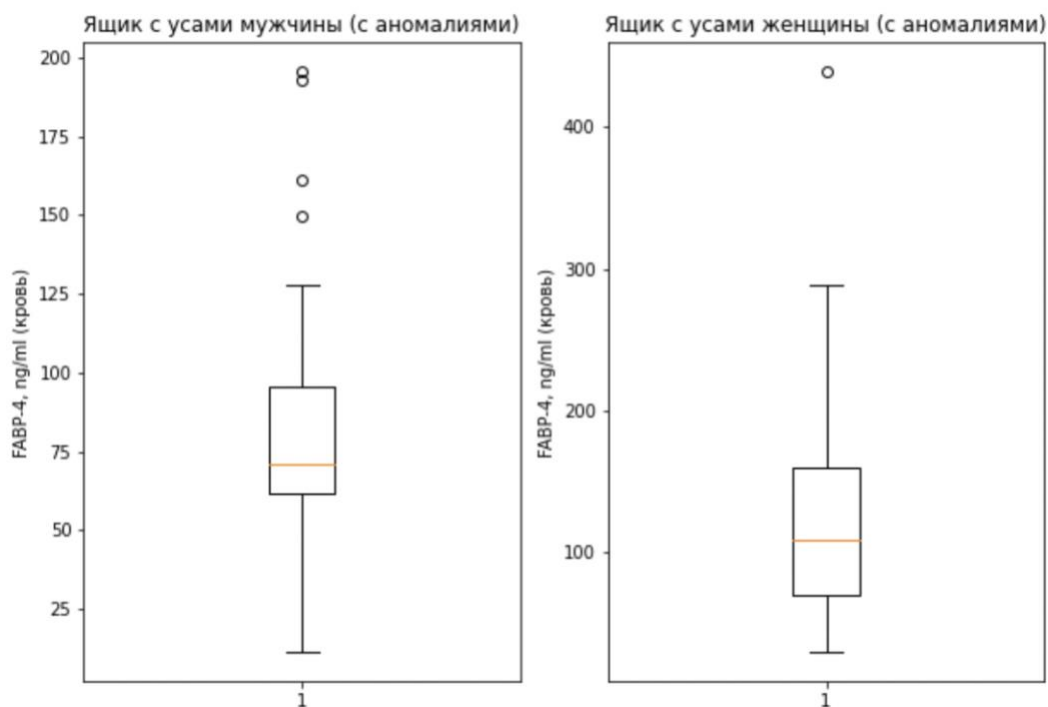


Рис. 2.2. Зависимость распределения фактора «FABP4-кровь» от пола.

Первый график (рис. 2.2) показывает наличие расхождений по параметру «FABP4-кровь» в двух группах – мужчины и женщины. Это также видно по разным значениям медиан (оранжевая линия) и относительным смещениям значений 1 и 3 квантилей.

Далее группа мужчин и женщин была поделена еще две подгруппы по наличию и отсутствию ИБС соответственно. Фактор «FABP4-кровь» был рассмотрен повторно для каждой подгруппы.



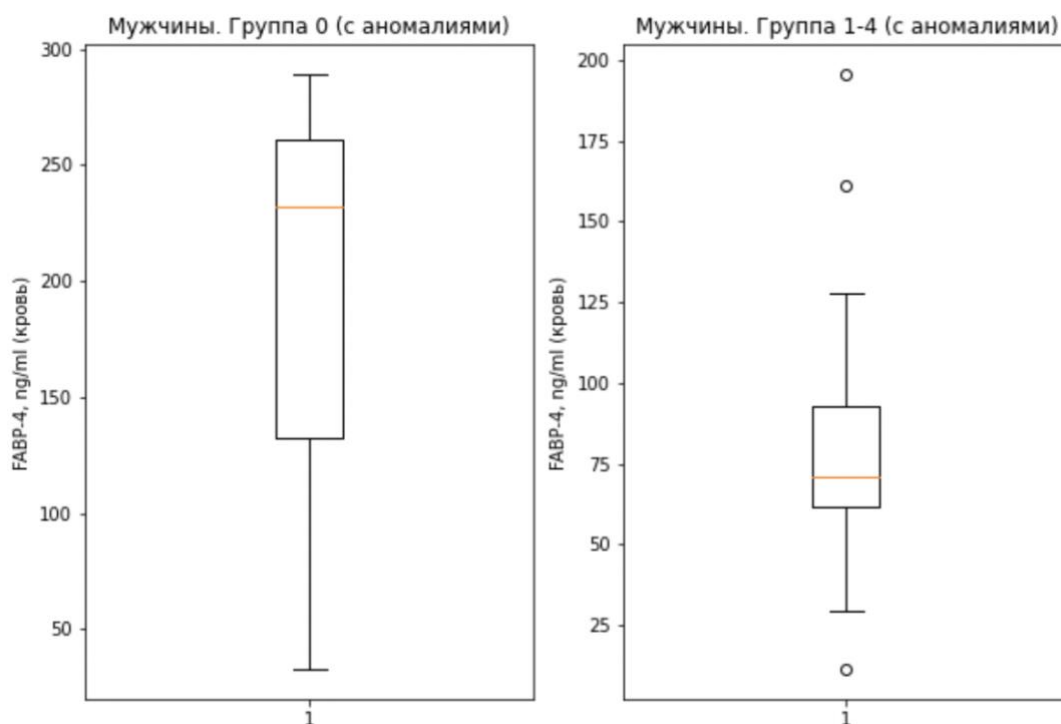


Рис. 2.3. Зависимость распределения фактора FABP4-кровь для мужчин с ИБС и без.

Полученный график для группы мужчин говорит об еще большем увеличении разрыва в статистически значимых величинах. Аналогичный график построен для группы женщин (рис. 2.4.):

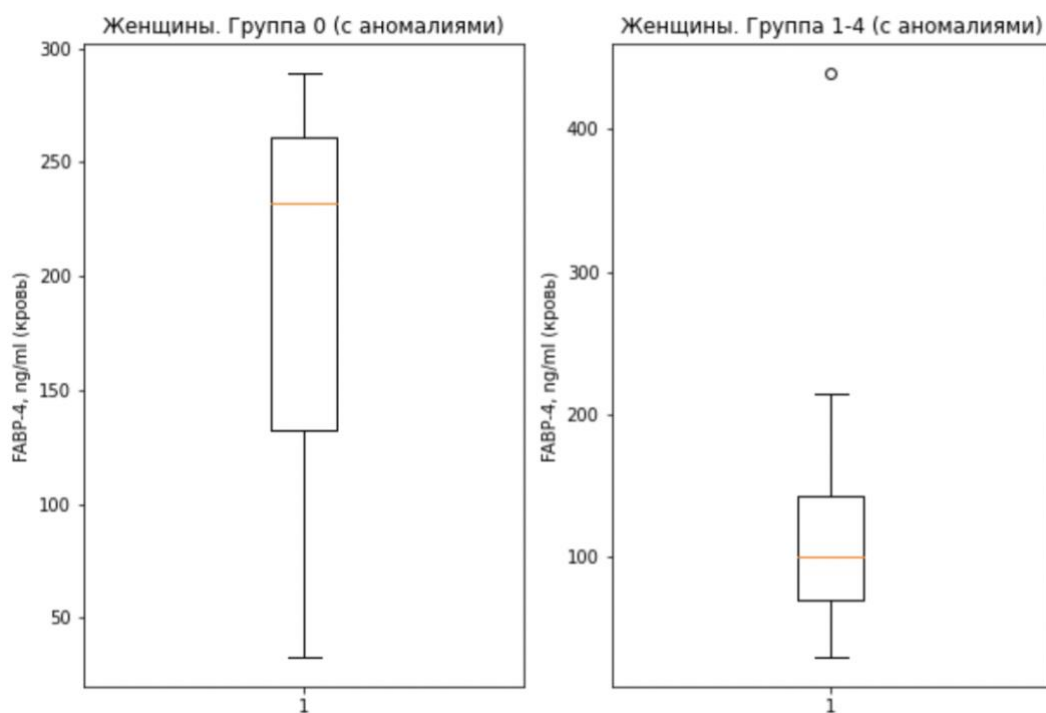


Рис. 2.4. Зависимость распределения фактора «FABP4-кровь» для женщин с ИБС и без.

На рисунке 2.4 также видно достоверное отличие в медианном значении.

Статистическая значимость отличия мужчин от женщин в группе людей с наличием ИБС по значению FABP4 в крови равна 0.000000, по значению FABP4 в эпидуральном жире равна 0.081698, в группе людей с отсутствием ИБС по значению FABP4 в крови равна 0.000000, по значению FABP4 в эпидуральном жире равна 0.000000.

Получившиеся результаты свидетельствуют о том, что наличие нескольких статистически значимых факторов влияния могут усиливать вероятность наличия или отсутствия ишемической болезни сердца. Детальный анализ относительной и совокупной зависимости факторов будет произведен с помощью предиктивной математической модели.

### ГЛАВА 3. ОТНОСИТЕЛЬНАЯ ОЦЕНКА ФАКТОРОВ ВЛИЯНИЯ

Для относительной оценки степени влияния каждого фактора будут использованы 3 математические модели на основе машинного обучения: модель логистической регрессии, модель решающих деревьев, catboost. Основная задача моделей – провести классификацию объектов по наличию ИБС. Далее будет выбрана модель с наибольшей предиктивной способностью, которая будет оценена по таким метрикам как: ROC – AUC, Precision, Recall, F1, также будет выведена метрика accuracy – доля верно определенных объектов. Формулы для вычисления:

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad (7)$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}, \quad (8)$$

$$\text{F1} = 2 \frac{\text{Recall} * \text{Precision}}{\text{Recall} + \text{Precision}}, \quad (9)$$

$$\text{ROC – AUC} = \text{площадь под графиком в осях TPR, FPR}, \quad (10)$$

где TP – верно угаданный 0 класс (без ИБС), FP – неверно предсказанный 0 класс, FN – неверно предсказанный 1 класс (с ИБС), TPR – доля верно определенных положительных объектов, FPR – доля верно определенных отрицательных объектов.

Выбор метрик совершался исходя из дисбаланса классов в обучающей выборке модели, количество объектов, болеющих ИБС в рассматриваемой выборке гораздо больше (рис. 3.1):

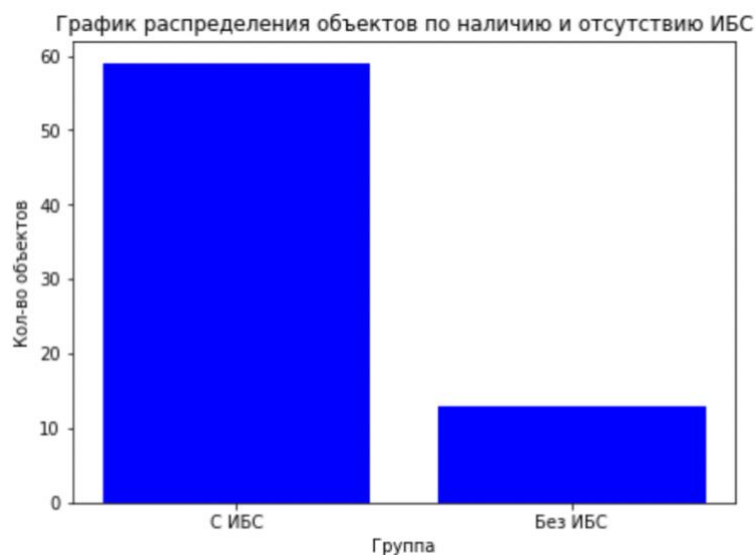


Рис. 3.1. Количество объектов в исследуемых группах.

После выбора наилучшей модели будут выведены коэффициенты – «веса» факторов, которые покажут относительную степень влияния каждого фактора.

### 3.1 Описание регрессионной модели

Исследование с помощью логистической регрессии состоит из следующих основных шагов:

1. Логистическая функция (сигмоидальная функция):

$$\sigma(z) = \frac{1}{1 + e^{-z}}, \quad (11)$$

где  $z$  – линейная комбинация входных признаков и их весов

2. Линейная комбинация входных признаков и их весов [12]:

$$z = w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n, \quad (12)$$

где  $w_0, w_1, w_2, \dots, w_n$  – веса признаков,  $x_1, x_2, \dots, x_n$  – значения входных признаков.

3. Функция отклика – вероятность отнесения к первому и второму классу, соответствующие наличию и отсутствию ИБС:

$$P(y = 1|x) = \sigma(z), \quad (13)$$

$$P(y = 0|x) = 1 - \sigma(z), \quad (14)$$

где  $P(y = 1|x)$  – вероятность отнесения объекта с данными признаками к классу, соответствующему отсутствию ИБС при условии  $x$ ,  $P(y = 0|x)$  – вероятность отнесения объекта с данными признаками к классу, соответствующему наличию ИБС при условии  $x$ .

4. Функция правдоподобия (likelihood):

$$L(w) = \prod_{i=1}^N P(y_i|x_i), \quad (15)$$

где  $L(w)$  – функция правдоподобия,  $N$  – количество наблюдений,  $y_i$  – истинные метки классов для наблюдения  $i$ ,  $x_i$  – значения признаков для наблюдения  $i$ .

5. Функция потерь (logloss):

$$J(w) = -\frac{1}{N} \sum_{i=1}^N [y_i \log(P(y = 1|x_i)) + (1 - y_i) \log(P(y = 0|x_i))], \quad (16)$$

где  $y_i$  – истинные метки классов для наблюдения  $i$ ,  $P(y = 1|x_i)$  – вероятность отнесения наблюдения  $i$  к классу, соответствующему отсутствию ИБС при условии  $x$ ,  $P(y = 0|x_i)$  – вероятность отнесения объекта с данными признаками к классу, соответствующему наличию ИБС при условии  $x$ .

6. Обучение модели:

Нахождение оптимальных весов  $w_i$  путем минимизации функции потерь  $J(w)$  с использованием метода градиентного спуска.

7. Регуляризация модели:

Для улучшения результатов разделяющей функции будут также использованы L1 (лассо) и L2 (гребневая) регуляризация. После будет выбран оптимальный вариант. Основная суть регуляризации заключается в добавлении штрафного члена в функцию потерь для ограничения величины коэффициентов модели.

Функция потерь при L1 регуляризации:

$$J(\theta) = J(w) + \lambda \sum_{j=1}^n |\theta_j|, \quad (17)$$

Функция потерь при L2 регуляризации:

$$J(\theta) = J(w) + \lambda \sum_{j=1}^n \theta_j^2, \quad (18)$$

где  $J(\theta)$  – функция потерь,  $\theta$  – вектор параметров модели (весов),  $\lambda$  – параметр регуляризации,  $J_\theta(w)$  – вычисленная функция потерь (10).

Результатом работы данной функции является вероятность отношения объекта к классу.

Данная модель является относительно простой, что говорит о ее высокой степени интерпретируемости. Также она имеет широкое применение в таких сферах как медицина, маркетинг, реклама, финансы и так далее.

### 3.2 Описание модели решающих деревьев

Решающие деревья относятся к древовидным моделям машинного обучения, которые в свою очередь применяются в статистических исследованиях при одновременном анализе трех и более признаков. Также виды древовидных моделей, такие как случайный лес, градиентный бустинг и так далее, обладают рядом преимуществ, делающие их полезным инструментом для анализа медицинских данных. Нами будет рассмотрено две древовидные модели: модель решающих деревьев и catboost для прогнозирования болезни.

Основные этапы реализации решения с помощью деревьев решений очень похожи на человеческий процесс принятия решений. Основной принцип работы алгоритма – нахождение последовательности простых решающих правил. Для определения качества разбиения выборки по выбранным правилам решающие деревья используют критерии оценки однородности. Алгоритм работы:

1. Определение узла.

Выбор признака и трешхолда - порога по нему для. Разбиение набора данных по полученному признаку и его трешхолду – определение данных в «листы» [7].

2. Оценка однородности данных, может быть осуществлена по следующим формулам:

Индекс Джинни:

$$\text{Gini}(p) = 1 - \sum_{i=1}^C p_i^2 \quad (19)$$

Энтропия:

$$\text{Entropy}(p) = - \sum_{i=1}^C p_i \log_2(p_i), \quad (20)$$

где  $\text{Gini}(p)$  – значение индекса Джинни для узла,  $\text{Entropy}(p)$  – значение энтропии для узла,  $C$  – количество классов,  $p_i$  – доля примеров класса  $i$  в узле.

3. Расчет прироста информации для оптимального разделения признака в узле.

Information Gain:

$$\text{IG}(D, A) = H(D) - \sum_{v=1}^V \frac{|D_v|}{|D|} H(D_v), \quad (21)$$

где  $\text{IG}(D, A)$  – прирост информации при разделении по признаку  $A$  в узле  $D$ ,  $H(D_v)$  – энтропия или индекс Джинни для узла  $D$ ,  $V$  – количество различных значений признака  $A$

4. Поиск следующего узла для оптимального разбиения, далее возвращаемся на этап 1. Процесс повторяется до тех пор, пока мы не достигнем заданной точности в оценке однородности.

При этом модели решающих деревьев склонны к переобучению, для избежания этой проблемы применена регуляризация, основанную на

следующих моментах: ограничение по максимальной глубине дерева, ограничение на минимальное количество объектов в листе, ограничение на максимальное количество листьев в дереве, требование улучшения качества при делении текущей подвыборки на две не менее чем на заданное число процентов. Выбор наиболее подходящего метода производился экспериментально.

### **3.3 Описание модели catboost**

Модель catboost основана на градиентном бустинге [8]. Она предназначена для решения задач машинного обучения. Ее главной особенностью является метод работы с категориальными признаками, которые встречаются и в исследуемой нами выборке. В отличие от других моделей catboost не требует дополнительной предварительной обработки категориальных признаков [11].

Алгоритм работы модели:

1. Построение базовых моделей.

В качестве базовых моделей принимаются деревья решений. На этом этапе происходит реализация стандартного алгоритма дерева решений для.

2. Вычисление функции потерь.

Для решения задачи бинарной классификации берется стандартная функция потерь  $\log \text{loss}$ , описанная ранее (16).

3. Обучение.

Обучение модели является итеративным процессом, в котором происходит последовательное построение деревьев с учетом градиента функции потерь предыдущего дерева.

При необходимости также используется регуляризация модели.

### **3.4 Исследование полученных результатов**

В результате было рассмотрено 4 математические модели на основе машинного обучения: модель логистической регрессии, модель



логистической регрессии с регуляризацией, модель решающих деревьев, catboost.

Все модели склонны к переобучению из-за небольшого количества объектов в обучающей выборке. Поэтому для достоверной оценки был применен метод кросс-валидации. При его использовании весь датасет делился на 5 равных частей. После проводилось поочередное построение 5 разных моделей каждого метода. При этом для тестовой выборки бралась 1 часть из полученных пяти, остальные 4 шли на обучающую выборку, после тестовая часть менялась, обучалась новая модели и так далее. Для каждой модели вычислялись заранее выбранные метрики. Также при окончании эксперимента было вычислено стандартное отклонение метрик, которое также необходимо для учета корректности работы моделей. Чем ниже этот показатель, тем лучше модель отвечает поставленной задаче.

Полученные результаты моделей были представлены в виде графика функции ROC-AUC, рис. 3.2:

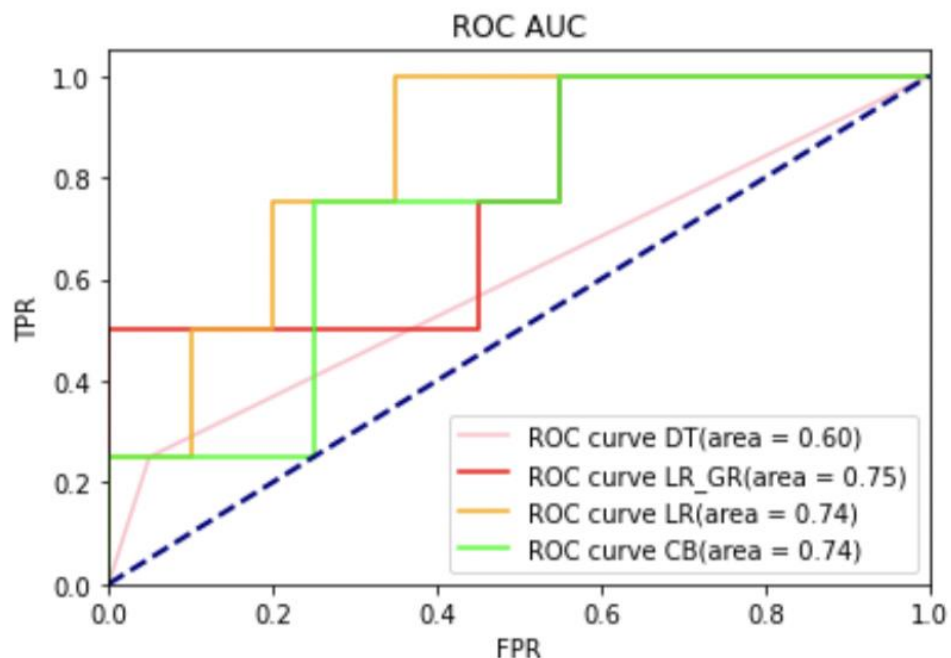


Рис. 3.2. Распределение количества объектов в исследуемых группах.

Наилучший результат показывает модель логистической регрессии, при регуляризационном коэффициенте  $\lambda = 0$ .

Помимо этого, была выведена таблица метрик, рис. 3.3:

	CatBoost	Logreg	Logreg_GR	Decision Tree
<b>accuracy</b>	0.833333	0.750000	0.750000	0.833333
<b>precision</b>	0.000000	0.796296	0.796296	0.803030
<b>recall</b>	0.000000	0.750000	0.750000	0.833333
<b>roc_auc</b>	0.680586	0.666117	0.645604	0.527106
<b>f1</b>	0.000000	0.621053	0.621053	0.619048

Рис. 3.3. Метрики качества моделей.

Исходя из таблицы, наилучший результат работы по метрике ROC-AUC имеет модель catboost, однако ее Precision и Recall равны нулю, что говорит о том, что она вовсе не предсказывает нулевой класс, то есть модель считает все объекты объектами первого класса (с ИБС), что говорит о ее плохой предиктивной способности. При этом наилучший результат также показывает логистическая регрессия. Её значение точности и полноты для определения каждого класса представлены на рисунке 3.4:

	precision	recall	f1-score	support
0	0.89	0.80	0.84	20
1	0.33	0.50	0.40	4
accuracy			0.75	24
macro avg	0.61	0.65	0.62	24
weighted avg	0.80	0.75	0.77	24

Рис. 3.4. Метрики качества модели логистической регрессии.

Значение метрики recall, равное 0.50, говорит, что модель верно нашла 2 объекта нулевого класса из 4 необходимых. При кросс-валидации стандартное отклонение метрик получилось равным 0.15, что говорит о дополнительных возможностях для улучшения модели, которое можно реализовать с помощью создания синтетической выборки или добавлением новых объектов в рассмотрение.

Для сравнительного анализа были выведены коэффициенты логистической регрессии, таблица 3.1:

Менопауза	-4.2444	x <sub>1</sub>
ГБ,1-да, 2 нет	-4.0702	x <sub>2</sub>
Конечные точки	-2.8613	x <sub>3</sub>
Адсист1, мм.рт.ст.	-1.4117	x <sub>4</sub>
билирубин 1	-1.1107	x <sub>5</sub>
Возраст	-1.0734	x <sub>6</sub>
АСТ	-1.0413	x <sub>7</sub>
Рост	-0.5423	x <sub>8</sub>
ОХС 1	-0.5386	x <sub>9</sub>
глюкоза 1	-0.5030	x <sub>10</sub>
Адиас1, мм.рт.ст.	-0.2021	x <sub>11</sub>
кол_во_сигарет_new	-0.1183	x <sub>12</sub>
креатинин 1	-0.0915	x <sub>13</sub>
ИМТ	1.2755	x <sub>14</sub>
ТГ 1	2.2534	x <sub>15</sub>
Наследственность_ССЗ	2.9108	x <sub>16</sub>
Aortic SP	3.4545	x <sub>17</sub>
Артерия_стентирования	4.4683	x <sub>18</sub>

Таблица 3.1. Значения коэффициентов.

Исходя из полученных коэффициентов можно восстановить функцию логистической регрессии. Вид общей функции при полученном исследовании:

$$\sigma(z) = \frac{1}{1 + e^{-z}}, \quad (22)$$

$$z = -4.2444x_1 - 4.0702x_2 - 2.8613x_3 - 1.4117x_4 - 1.1107x_5 - 1.0734x_6 - 1.0413x_7 - 0.5423x_8 - 0.5386x_9 - 0.5030x_{10} - 0.2021x_{11} - 0.1183x_{12} - 0.0915x_{13} + 1.2755x_{14} + 2.2534x_{15} + 2.9108x_{16} + 3.4545x_{17} + 4.4683x_{18} \quad (23)$$

где  $x_i$  – соответствующие значения переменных.

Также по коэффициентам можно оценить относительную важность исследуемых факторов. Степень влияния признака оценивается по абсолютному значению коэффициента, таблица 3.2:

<b>креатинин 1</b>	0.0915
<b>кол_во_сигарет_new</b>	0.1183
<b>Аддиас1, мм.рт.ст.</b>	0.2021
<b>глюкоза 1</b>	0.5030
<b>ОХС 1</b>	0.5386
<b>Рост</b>	0.5422
<b>АСТ</b>	1.0413
<b>Возраст</b>	1.0734
<b>билирубин 1</b>	1.1107
<b>ИМТ</b>	1.2755
<b>Адсист1, мм.рт.ст.</b>	1.4117
<b>ТГ 1</b>	2.2534
<b>Конечные точки</b>	2.8613
<b>Наследственность_ССЗ</b>	2.9108
<b>Aortic SP</b>	3.4545
<b>ГБ,1-да, 2 нет</b>	4.0702
<b>Менопауза</b>	4.2444
<b>Артерия_стентирования</b>	4.4683

Таблица 3.2. Абсолютные значения коэффициентов.

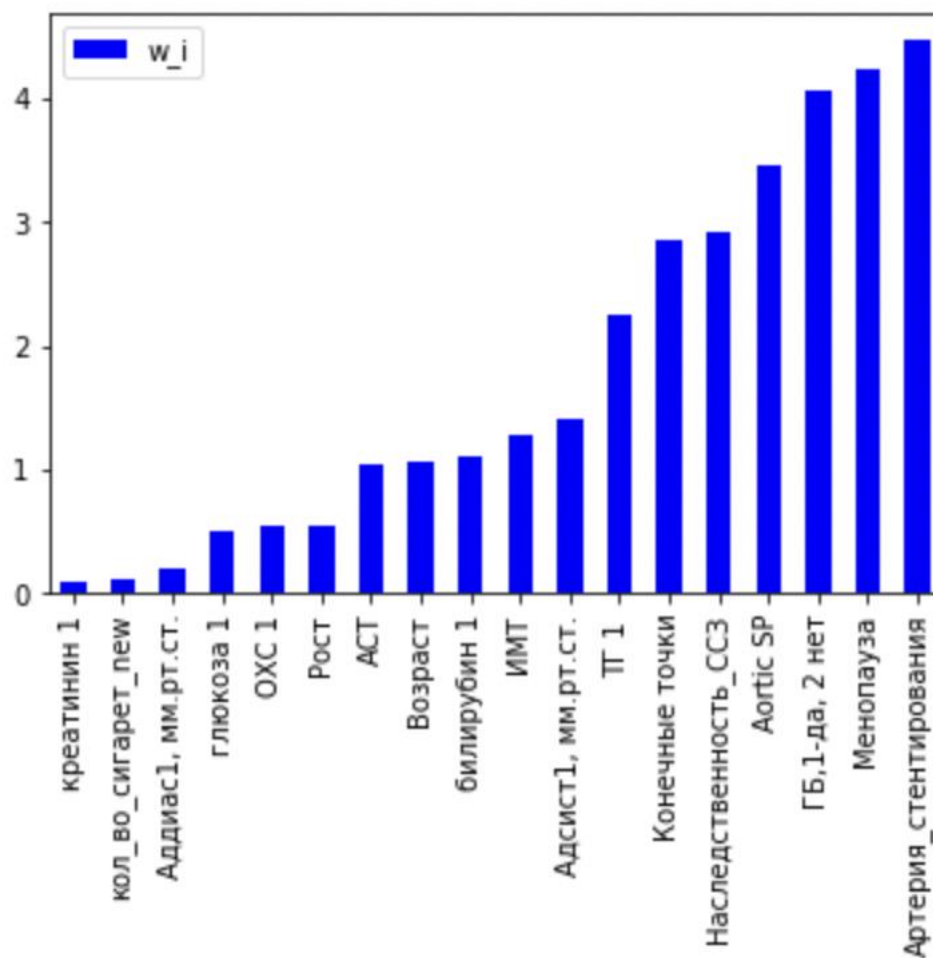


Рис. 3.5. График степени влияния факторов на риск возникновения ИБС.

Чем больше значение коэффициента, тем сильнее он влияет на риск возникновения ИБС. Наибольший вес в принятии решения моделью об отсутствии или наличии ишемической болезни сердца внесли признаки: артерия стентирования, наличие менопаузы, которая в свою очередь пропорциональна полу с точки зрения линейной корреляции, что также говорит о сильном влиянии пола в данном решении математической модели линейной регрессии. Также есть признаки, которые повлияли на результат логистической регрессии несильно: креатинин, количество сигарет, глюкоза, аддиас. Из них два признака: глюкоза, аддиас не имеют достаточных обоснований предполагать, что эти факторы влияют на риск возникновения ИБС со статистической значимостью, равной 0.01.

Стоит также отметить, что при исследовании были найдены группы высокой корреляции, из каждой группы был оставлен лишь один фактор для дальнейшего обучения модели. Поэтому при рассмотрении какого-либо фактора и его коэффициента, в случае если этот фактор входит в одну из групп высокой корреляции, его коэффициент и степень влияния распространяется на все факторы группы.

Аналогично были построены регрессионные модели для предсказания значений FАВР4 в крови и эпидуральном жире в группе людей с ИБС для проведение многофакторного анализа, а также относительных степеней влияния:

<b>ТГ 1</b>	-28.3547
<b>Курение_в_наст_вр</b>	-26.4444
<b>ИМТ</b>	-23.9488
<b>homa_ir</b>	-0.8356
<b>Пол, 1-м, 2-ж</b>	16.2227
<b>Менопауза</b>	16.2227
<b>Число перенесенных ИМ</b>	20.0496
<b>Наследственность_ССЗ</b>	27.6155
<b>Возраст</b>	56.7007
<b>Сахарный_диабет</b>	61.0786
<b>Окружность талии</b>	160.8378

Таблица 3.3. Значения коэффициентов влияния факторов на значение FАВР4 в крови.

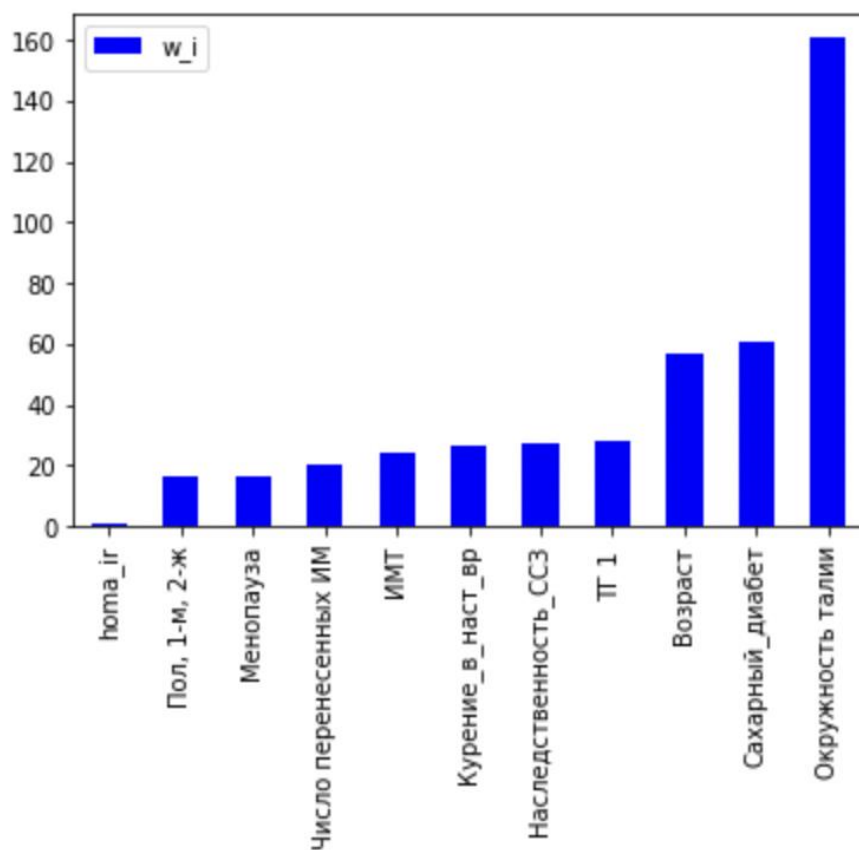


Рис. 3.6. График степени влияния факторов на значение FABP4 в крови.

<b>Наследственность_ССЗ</b>	-0.6130
<b>Возраст</b>	-0.3369
<b>Число перенесенных ИМ</b>	-0.2547
<b>Пол, 1-м, 2-ж</b>	-0.2476
<b>Менопауза</b>	-0.2476
<b>homa_ir</b>	-0.1826
<b>Окружность талии</b>	-0.1281
<b>ИМТ</b>	-0.0185
<b>ТГ 1</b>	0.0272
<b>Сахарный_диабет</b>	0.0386
<b>ФАВР-4, ng/ml (кровь)</b>	0.0869
<b>Курение_в_наст_вр</b>	0.1602

Таблица 3.4. Значения коэффициентов влияния факторов на значение FABP4 в эпидуральном жире.

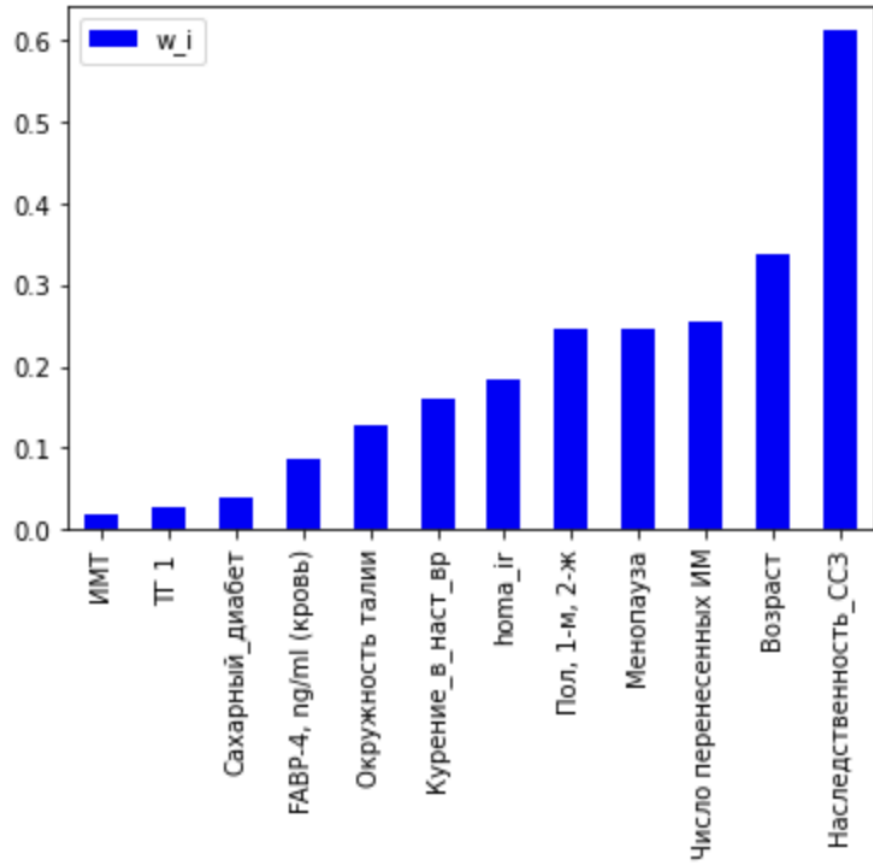


Рис. 3.7. График степени влияния факторов на значение FABP4 в эпидуральном жире.



## ЗАКЛЮЧЕНИЕ

В ходе исследовательской работы был проведен анализ 221 различных факторов влияния на риск возникновения ишемической болезни сердца, при этом количество пациентов, вошедших в исследуемую группу, составило 125 человек. Данные были исследованы в обезличенном виде. Из общего набора факторов на основе анализа разреженности данных была выявлена группа наиболее репрезентативных признаков – 34 различных факторов.

Были сформулированы гипотезы о возможном влиянии каждого из этих факторов на возникновение болезни. Из них 29 факторов показали статистически значимое влияние на риск возникновения ИБС с заданной точностью 0.01.

Для проверки гипотез были выбраны следующие статистические критерии: Манна-Уитни для количественных факторов и ANOVA для категориальных. Критерий Манна-Уитни позволяет выявлять различия в значении параметра между малыми выборками. Дисперсионный анализ ANalysis Of Variance (ANOVA) дает возможность проводить множественные сравнения. Условия применимости обоих критериев были соблюдены – независимость наблюдений в выборке была учтена при сборе данных, нормальность распределений была реализована при помощи метода bootstrap и проверена с помощью критерия Шапиро-Уилка.

Также с помощью математических моделей, основанных на алгоритмах машинного обучения, была произведена оценка влияния одновременного наличия 18 различных факторов. Помимо этого, было найдено относительное влияние каждого признака. Для анализа были рассмотрены 3 модели: логистическая регрессия, решающие деревья и catboost. Все три модели были обучены решать задачу бинарной классификации – предсказания наличия и отсутствия ИБС у пациента по имеющимся факторам. Далее была выбрана модель с наилучшими предиктивными способностями с точки зрения метрик

машинного обучения, модель логистической регрессии. На основе коэффициентов полученной модели создан график степени влияния каждого фактора.

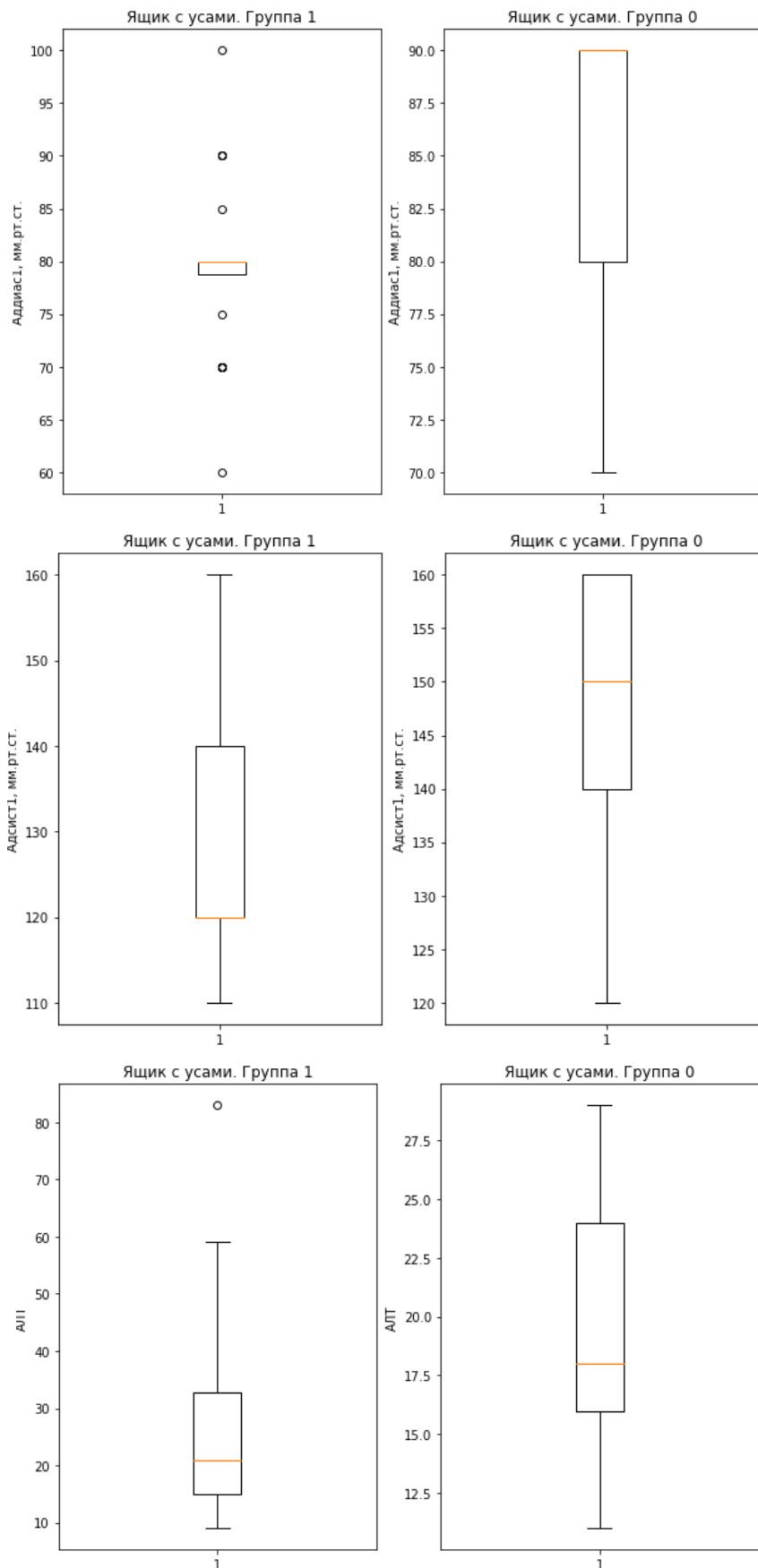
## СПИСОК ИСПОЛЬЗОВАННОЙ ЛИТЕРАТУРЫ

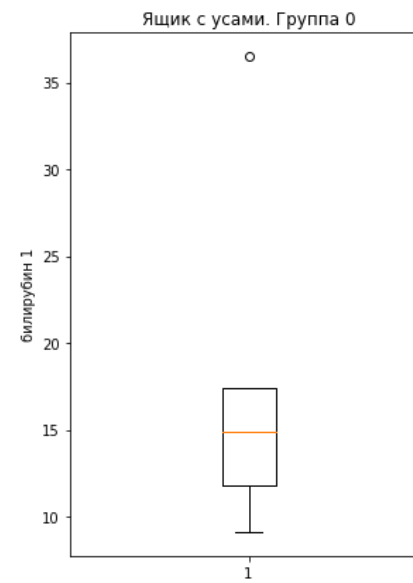
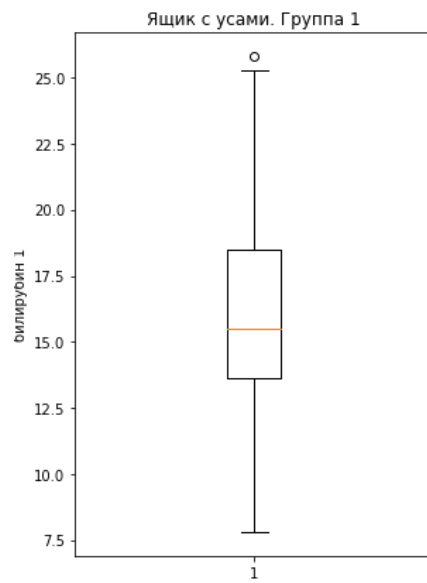
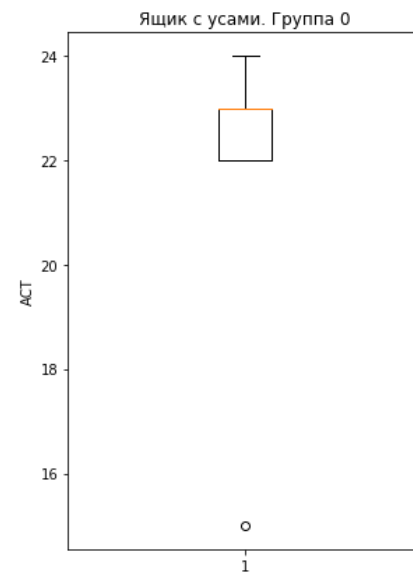
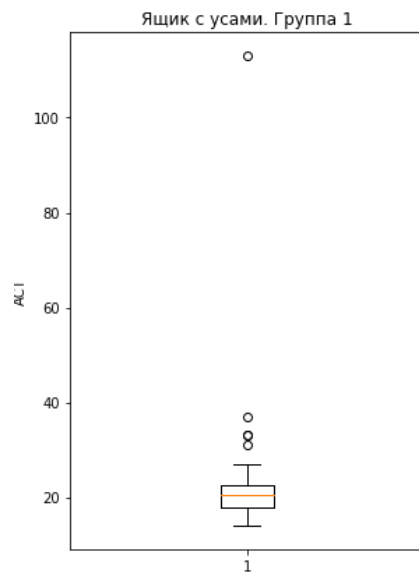
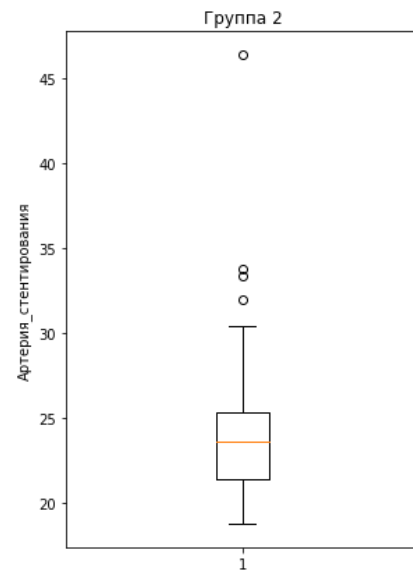
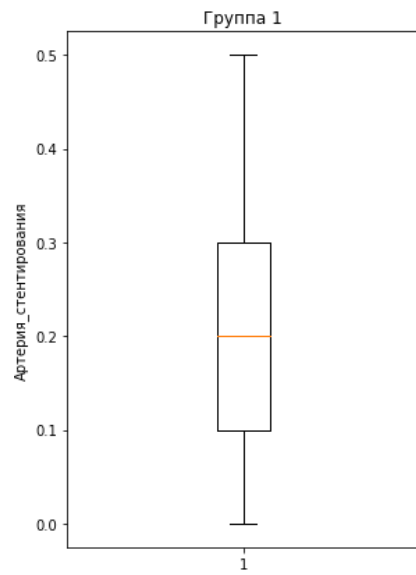
1. Сидняев Н.И. Теория вероятностей и математическая статистика: учебник для вузов, Москва: издательство Юрайт, 2023. 219 с.
2. Шарафутдинова Н.Х., Киреева Э.Ф., Николаева И.Е., Павлова М.Ю., Халфин Р.М., Шарафутдинов М.А., Борисова М.В., Латыпов А.Б., Галикеева А.Ш. – Статистические методы в медицине и здравоохранении, учеб. пособие, Уфа: ФГБОУ ВО БГМУ Минздрава России, 2018. – 131 с.
3. Chernick M. R. Bootstrap methods: A guide for practitioners and researchers. – John Wiley & Sons, 2011.
4. Conover W. J. Practical nonparametric statistics. – John Wiley & Sons, 1999. – Т. 350.
5. Efron B. Tibshirani RJ. An introduction to the bootstrap //Monographs on Statistics and Applied Probability. – 1993. – Т. 57. – С. 1-177.
6. Hinkelmann K., Kempthorne O. Design and analysis of experiments, volume 1: Introduction to experimental design. – John Wiley & Sons, 2007. – Т. 1.
7. Breiman L. Random forests //Machine learning. – 2001. – Т. 45. – С. 5-32.
8. Friedman J. H. Greedy function approximation: a gradient boosting machine //Annals of statistics. – 2001. – С. 1189-1232.
9. Mann H. B., Whitney D. R. On a test of whether one of two random variables is stochastically larger than the other //The annals of mathematical statistics. – 1947. – С. 50-60.
10. Miroshnikova V. V. et al. FABP4 and omentin-1 gene expression in epicardial adipose tissue from coronary artery disease patients //Genetics and molecular biology. – 2021. – Т. 44.
11. Prokhorenkova L. et al. CatBoost: unbiased boosting with categorical features //Advances in neural information processing systems. – 2018. – Т. 31.
12. Seber G. A. F., Lee A. J. Linear regression analysis. – John Wiley & Sons, 2003. – Т. 330.

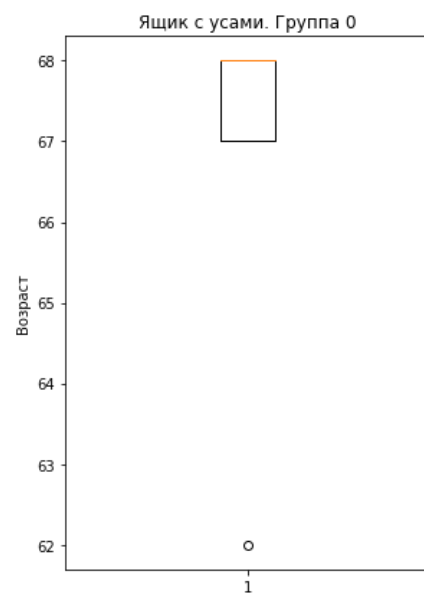
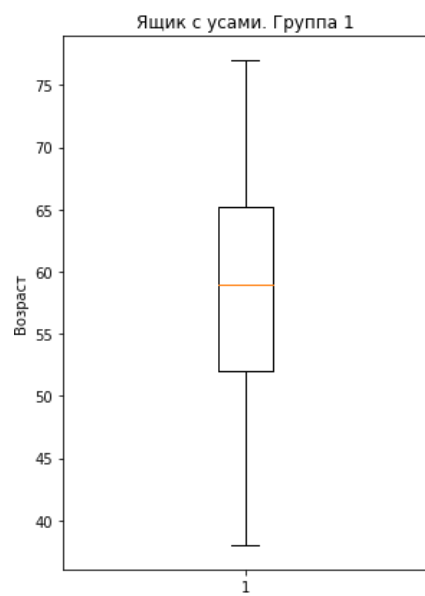
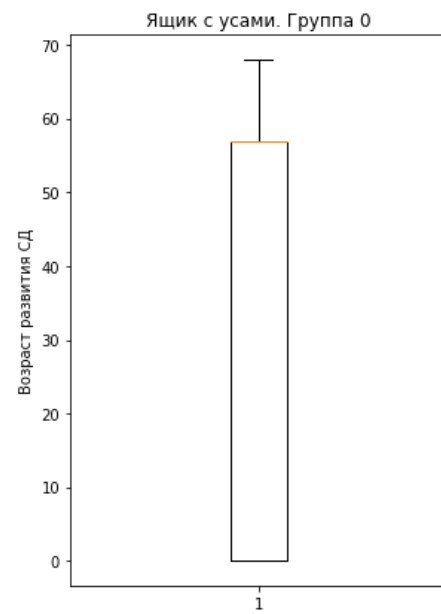
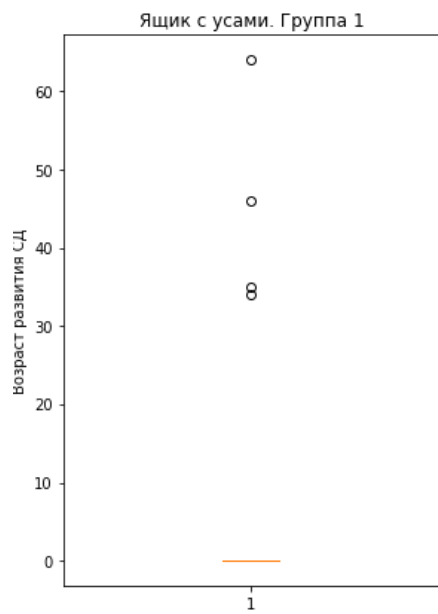
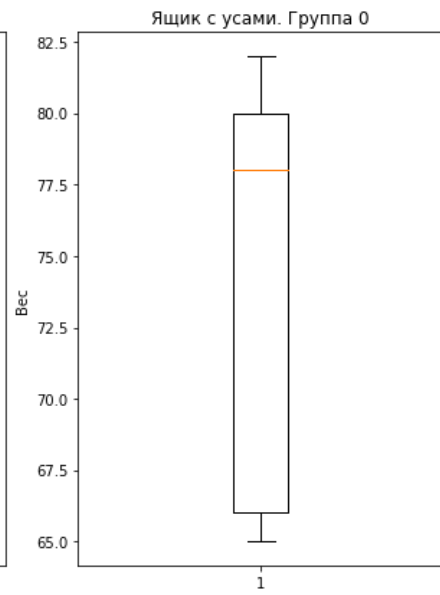
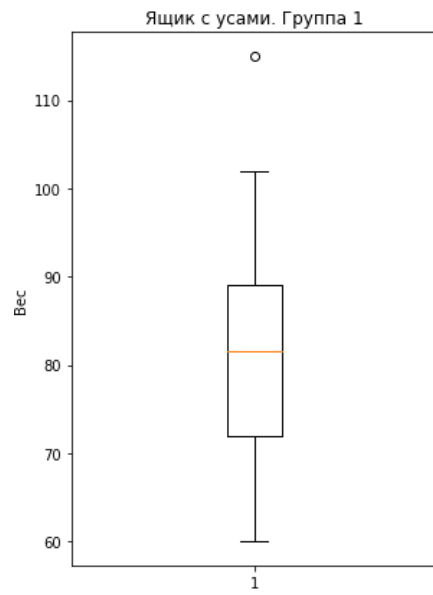
13. Shapiro S. S., Wilk M. B. An analysis of variance test for normality (complete samples) //Biometrika. – 1965. – Т. 52. – №. 3/4. – С. 591-611.
14. Yusuf S. et al. Global burden of cardiovascular diseases: part I: general considerations, the epidemiologic transition, risk factors, and impact of urbanization //Circulation. – 2001. – Т. 104. – №. 22. – С. 2746-2753.
15. Cardiovascular diseases (CVDs) // World Health Organization official site 2020. – URL: [https://www.who.int/en/news-room/fact-sheets/detail/cardiovascular-diseases-\(cvds\)](https://www.who.int/en/news-room/fact-sheets/detail/cardiovascular-diseases-(cvds)) — (дата обращения: 18.03.2023).

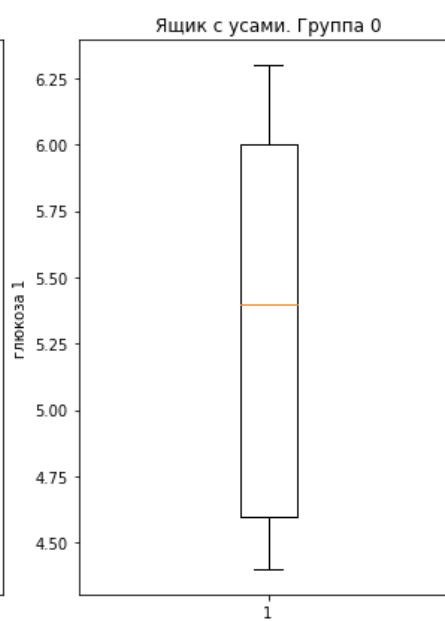
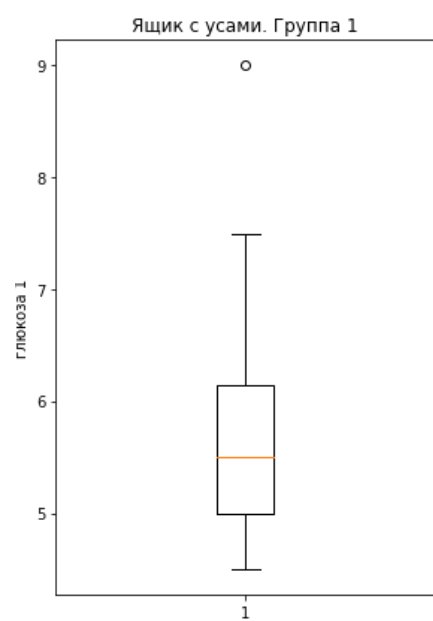
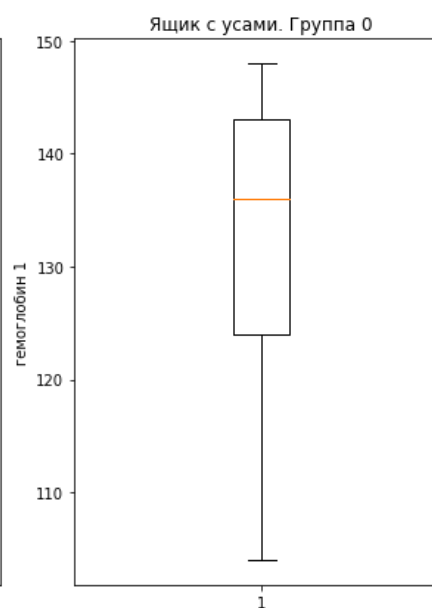
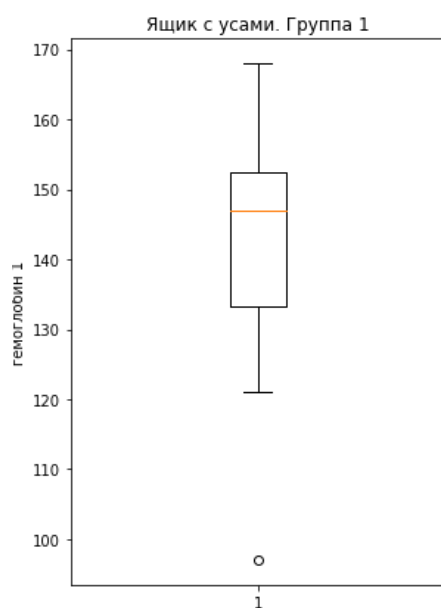
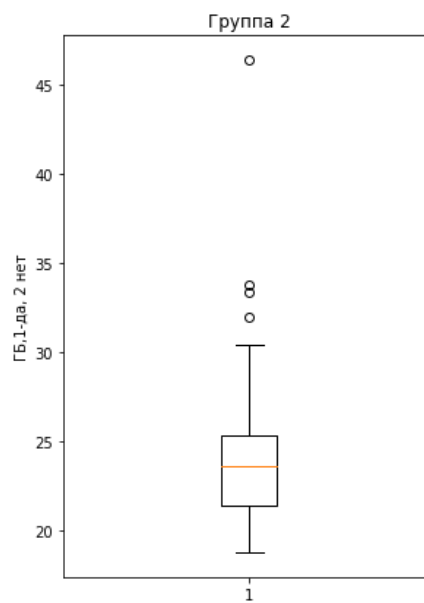
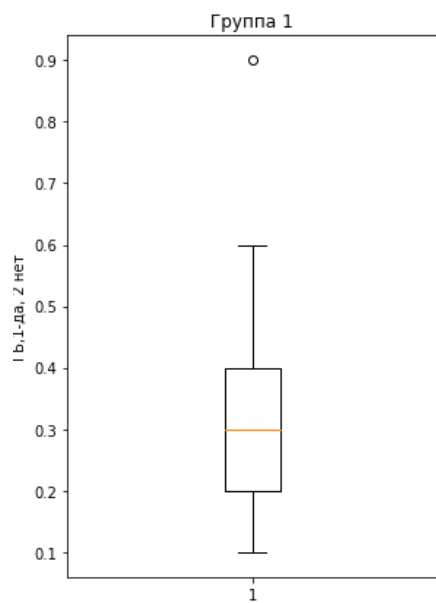
## ПРИЛОЖЕНИЕ

Графики ящиков с усами для репрезентативных факторов влияния на риск возникновения ИБС, группа 1 – люди с ИБС, группа 0 – люди без ИБС:

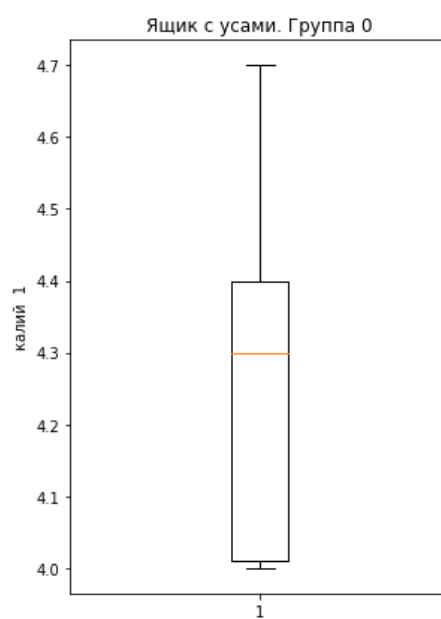
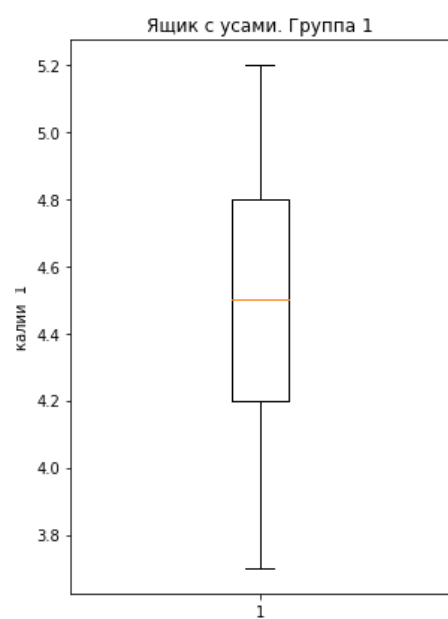
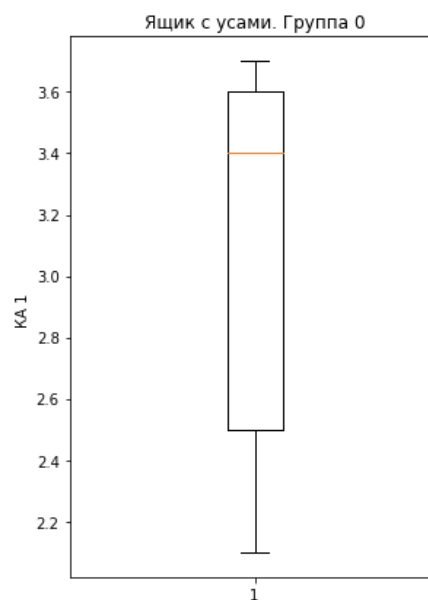
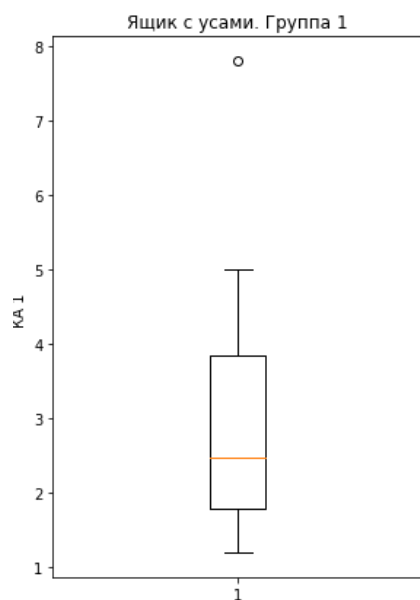
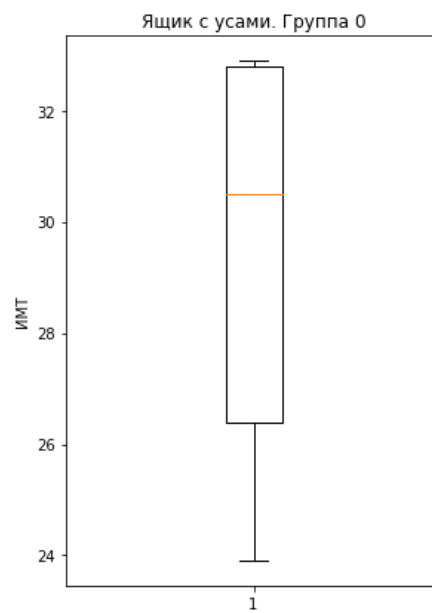
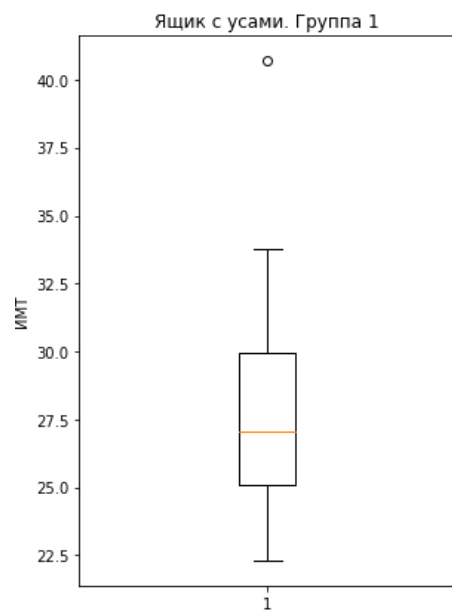


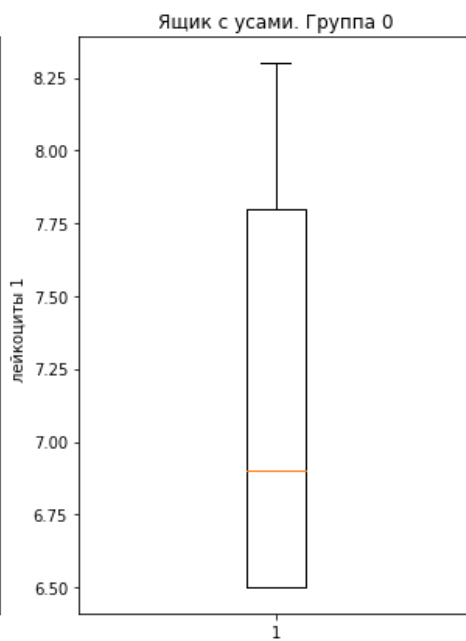
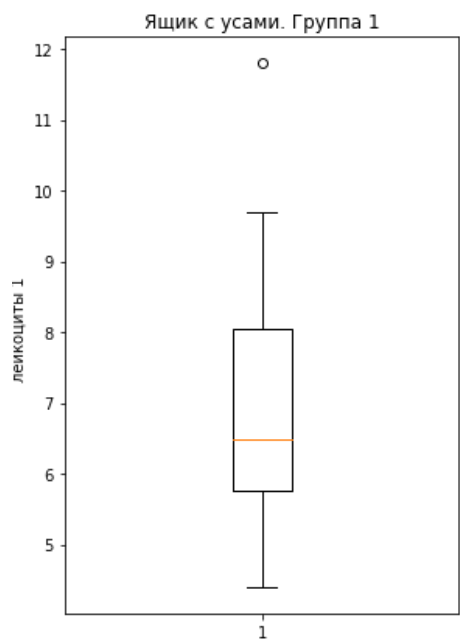
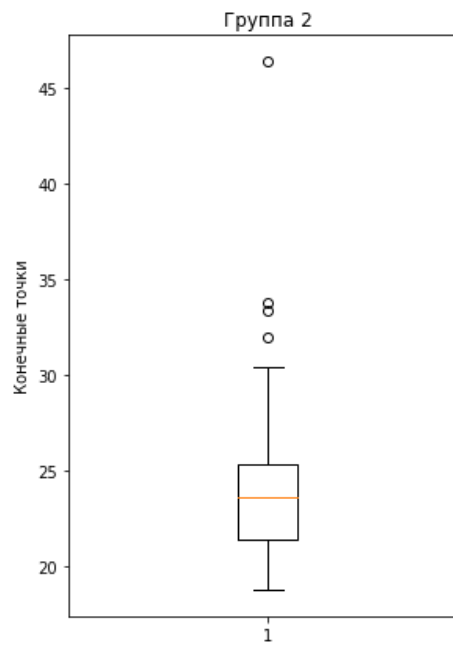
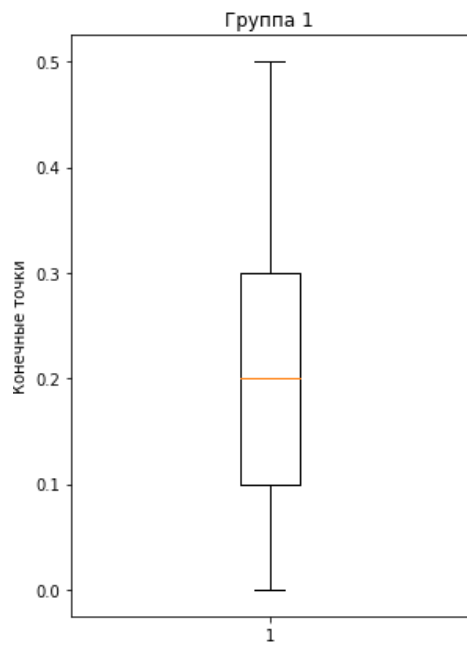
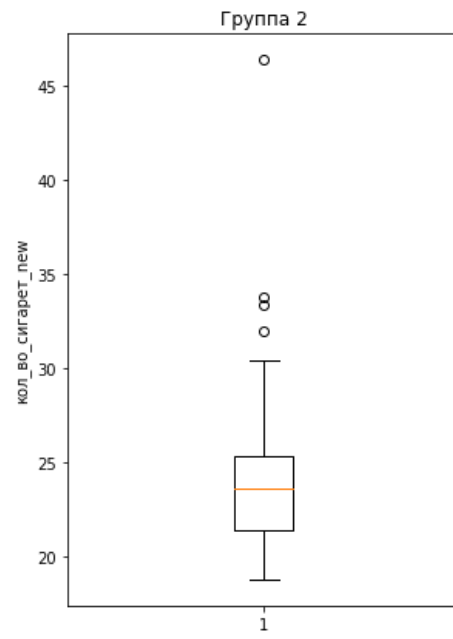
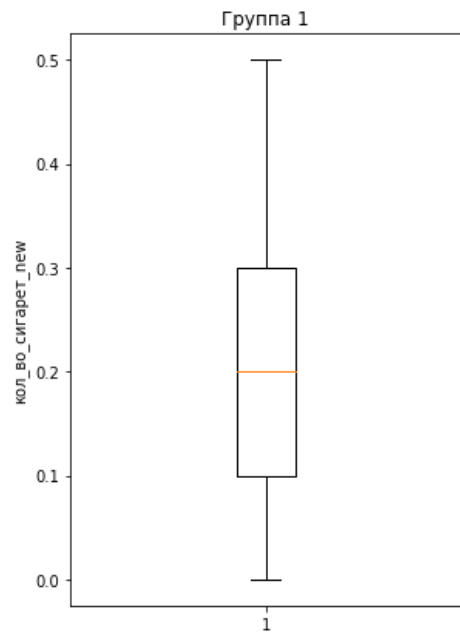


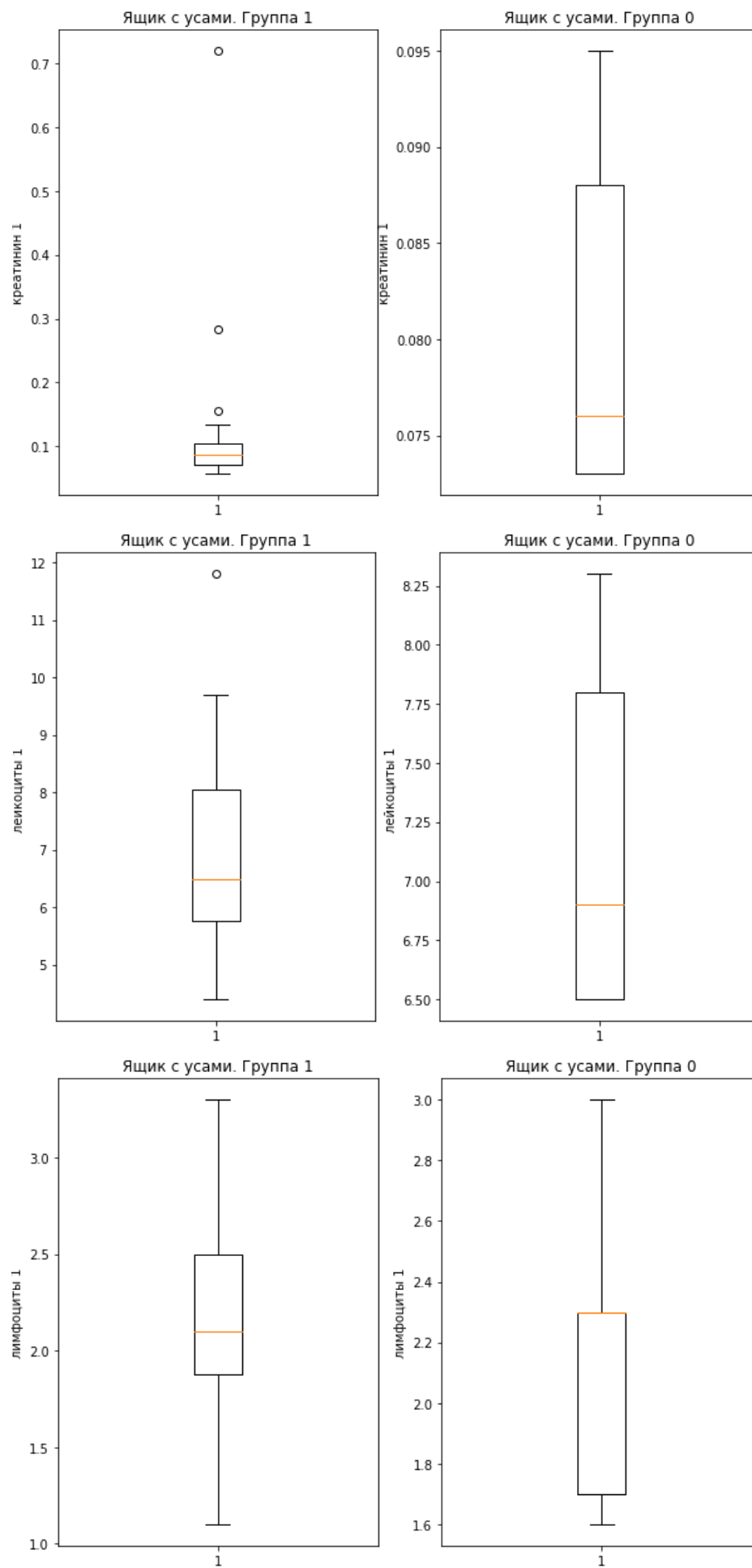


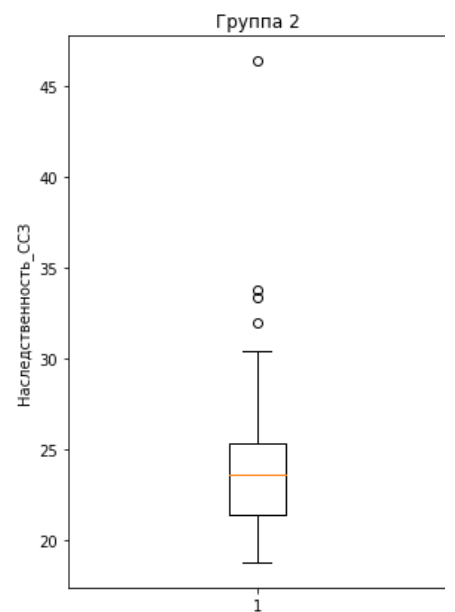
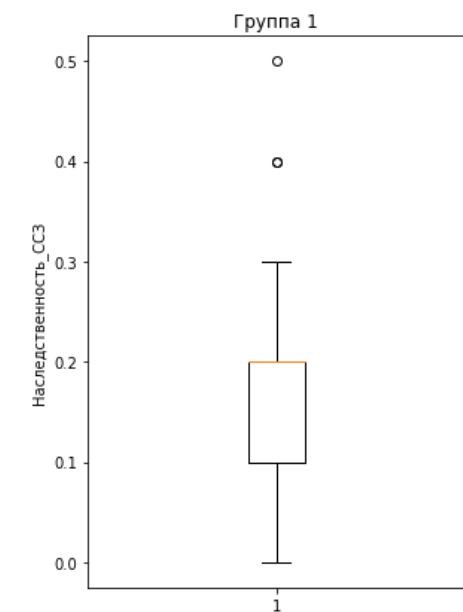
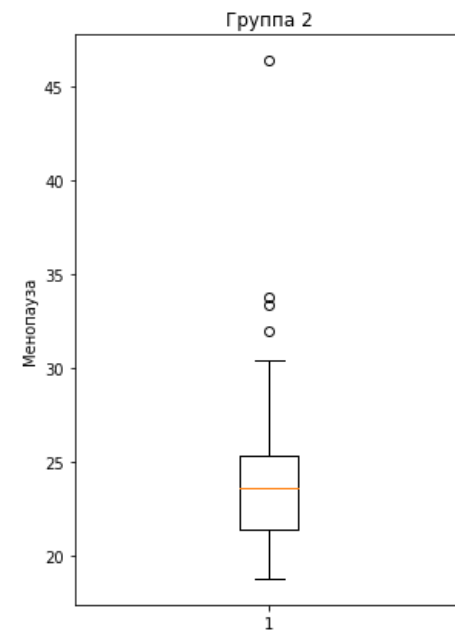
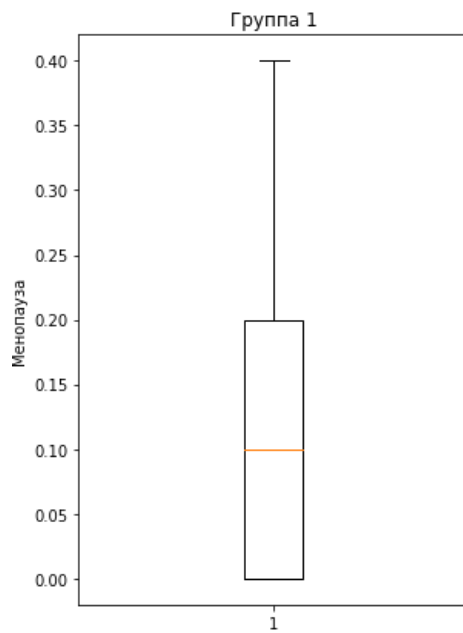
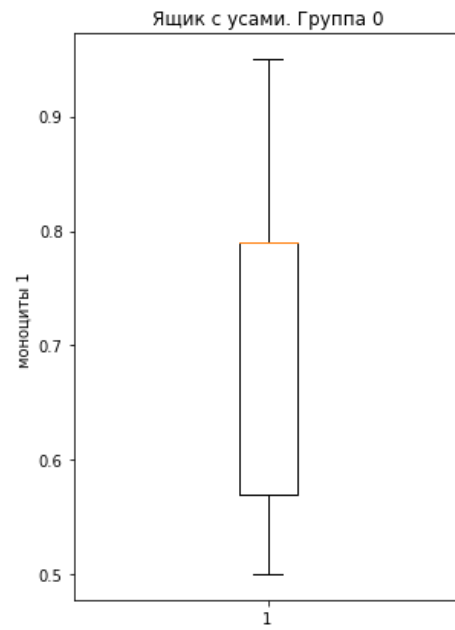
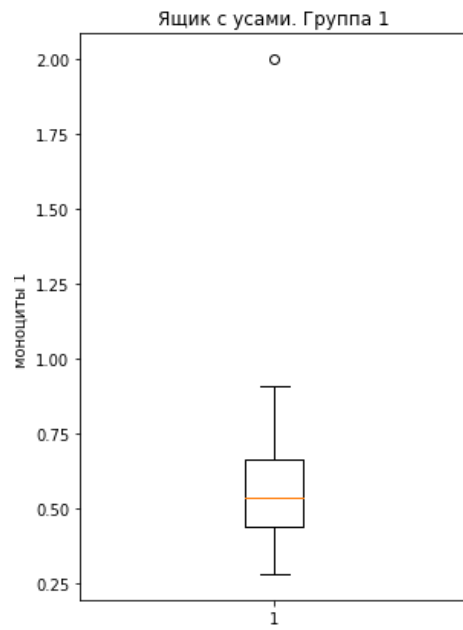


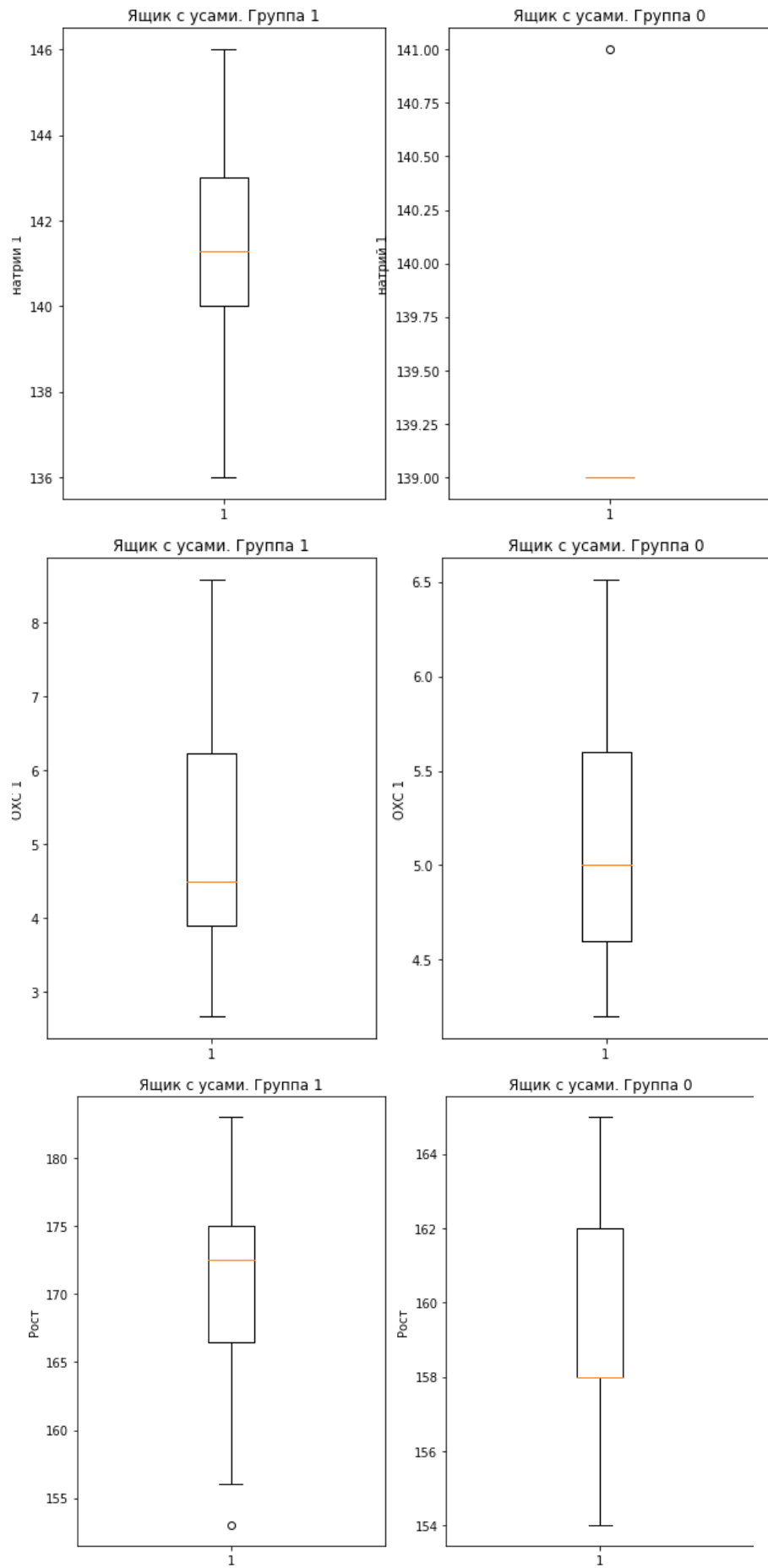


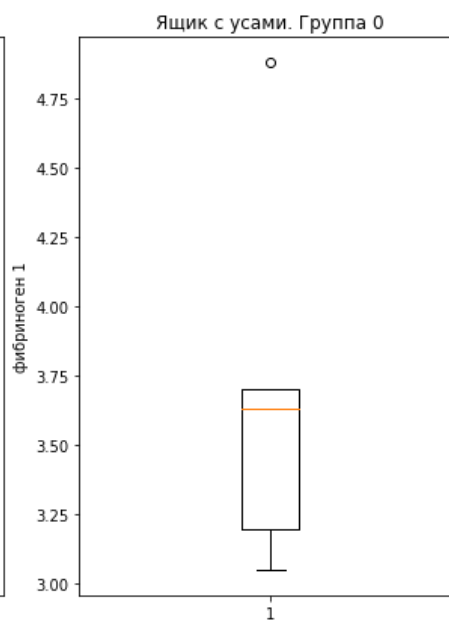
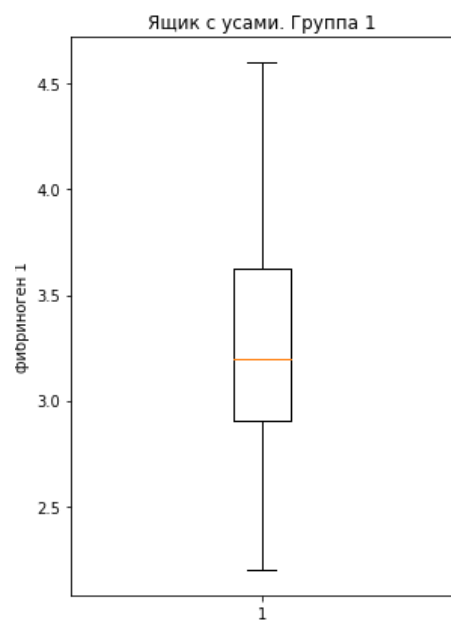
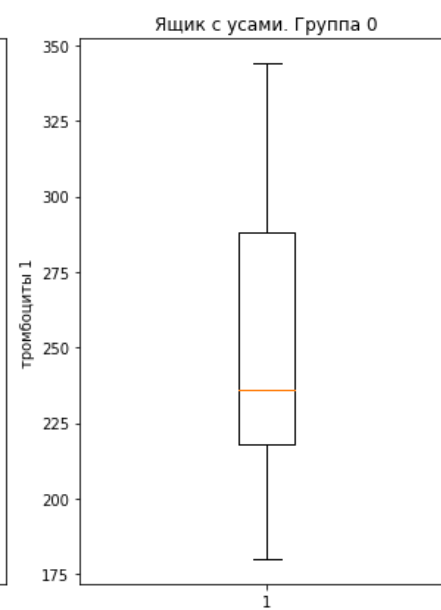
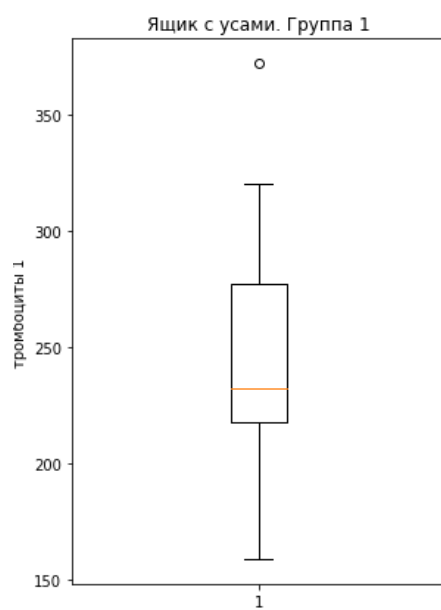
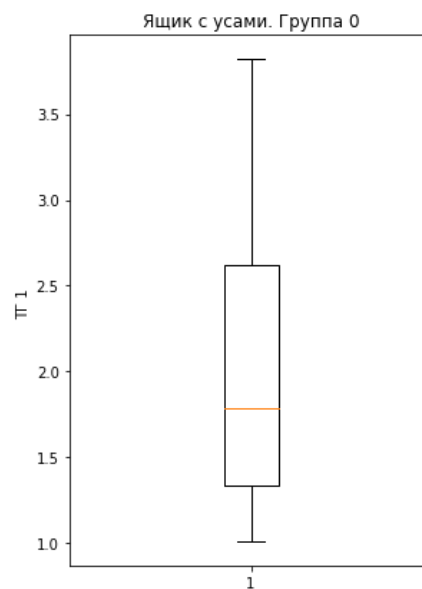
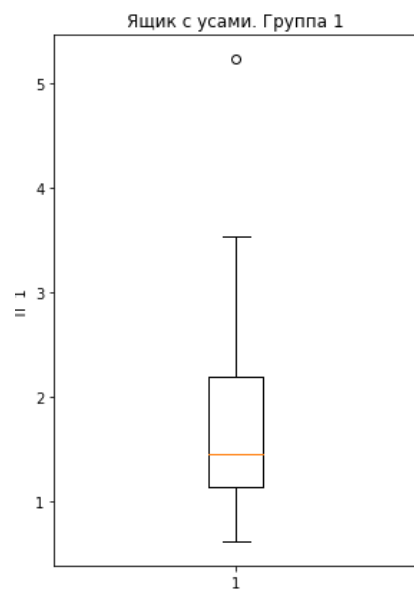


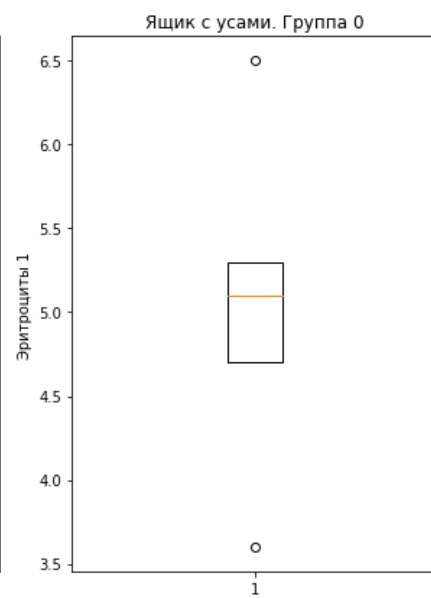
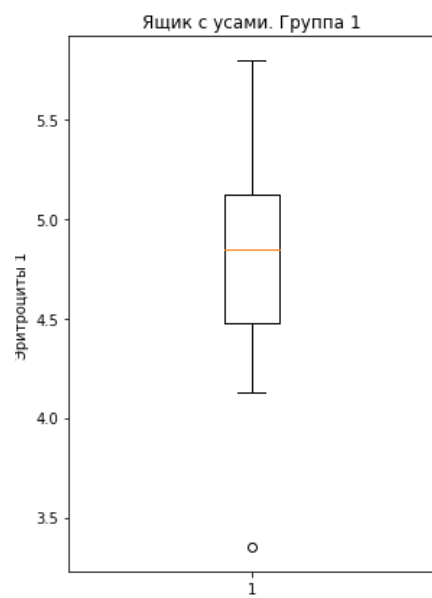
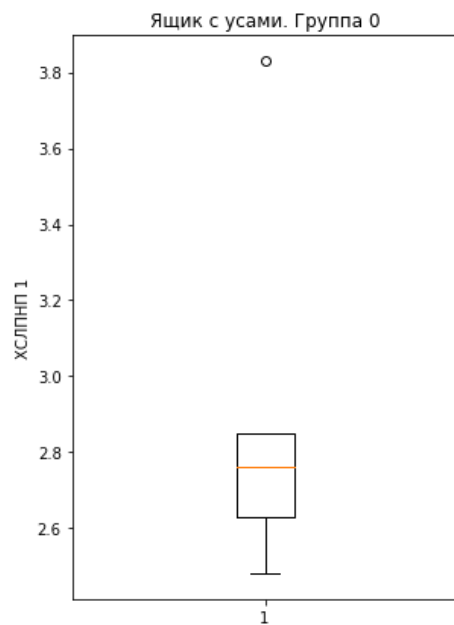
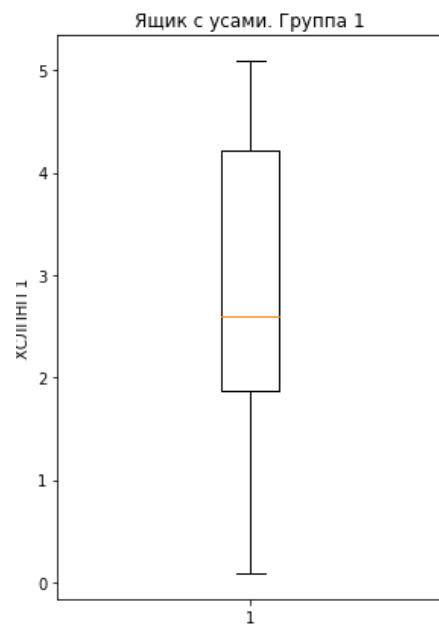
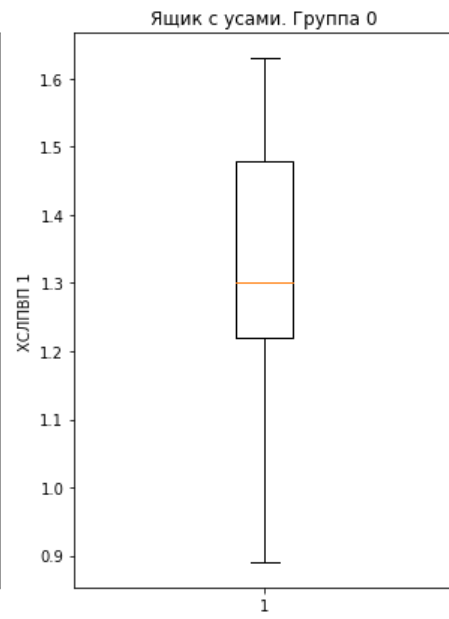
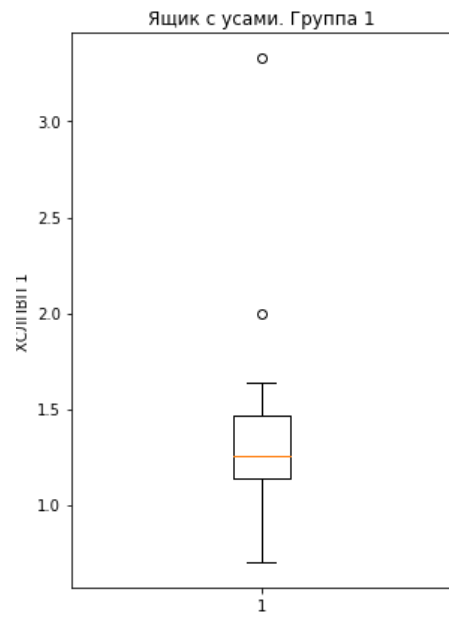


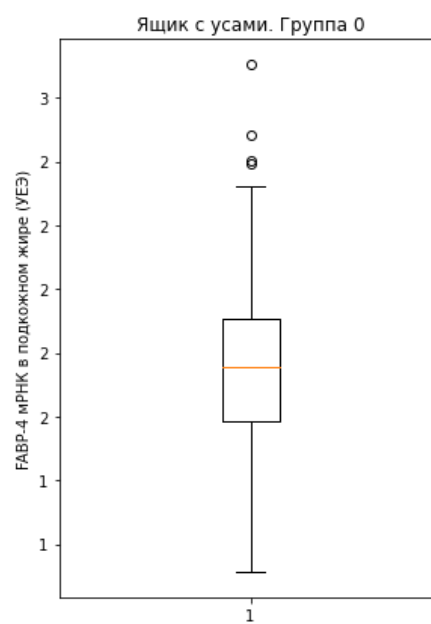
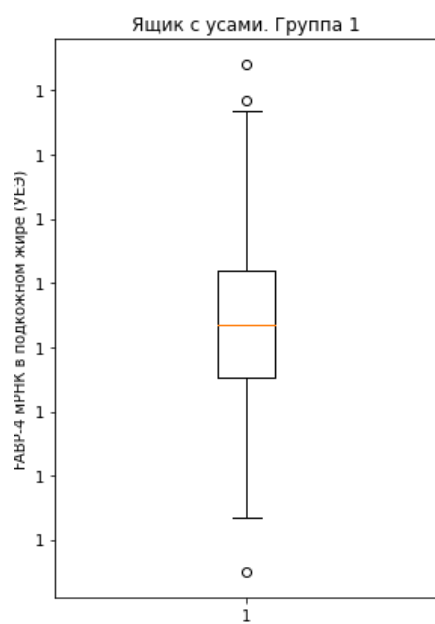
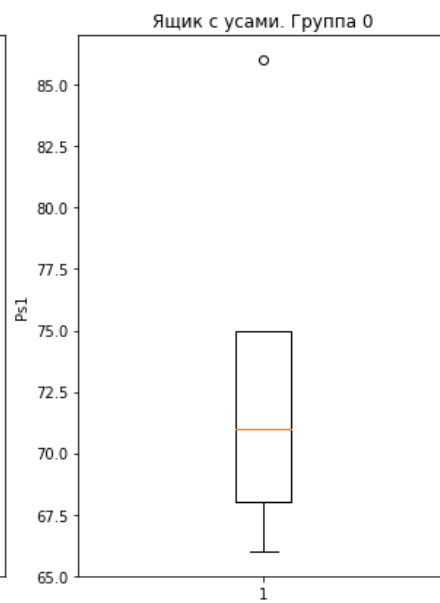
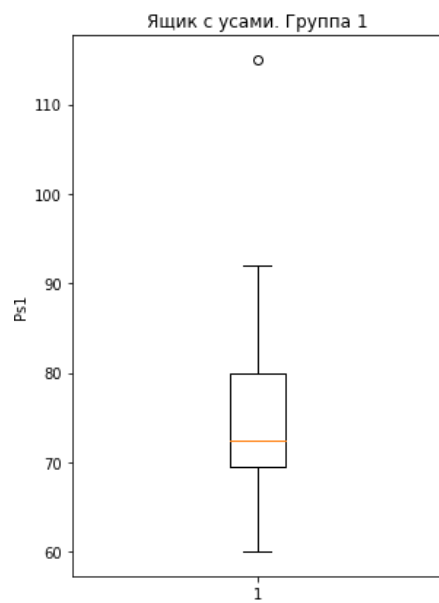




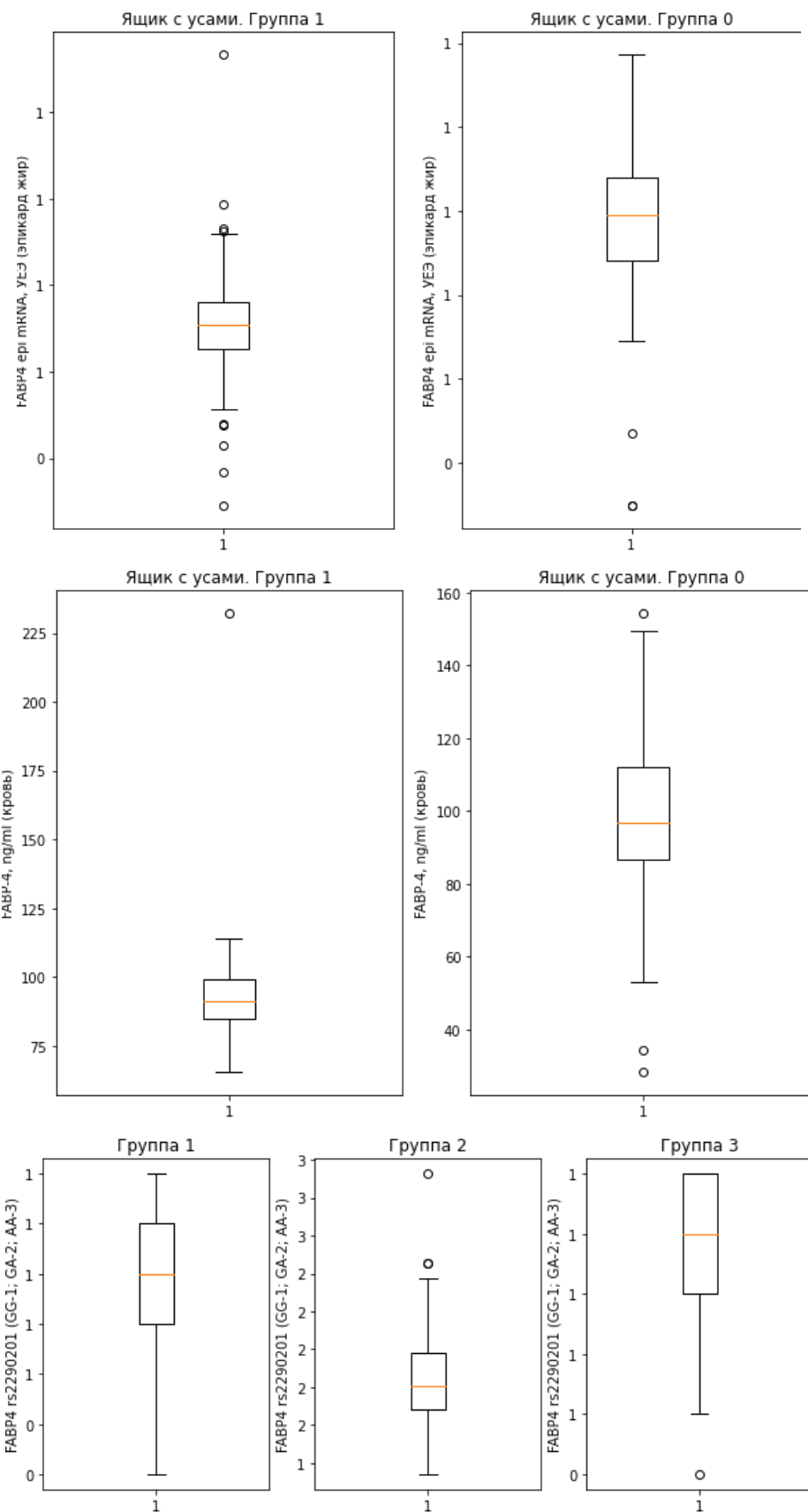


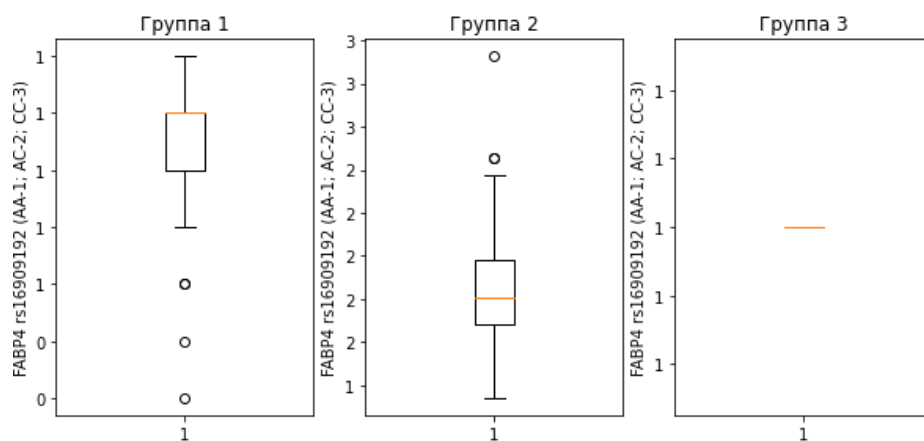






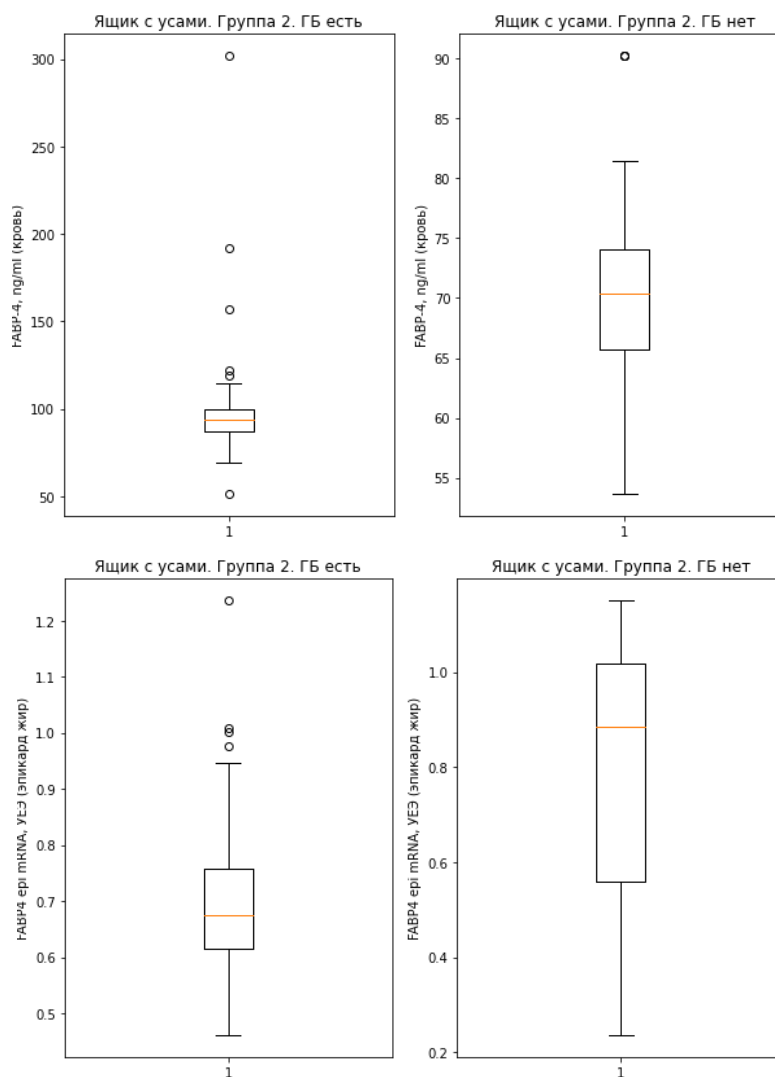


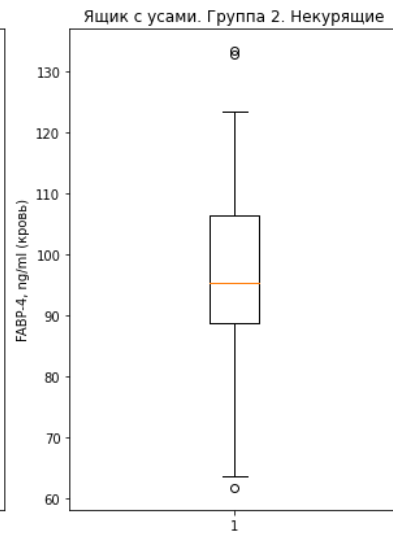
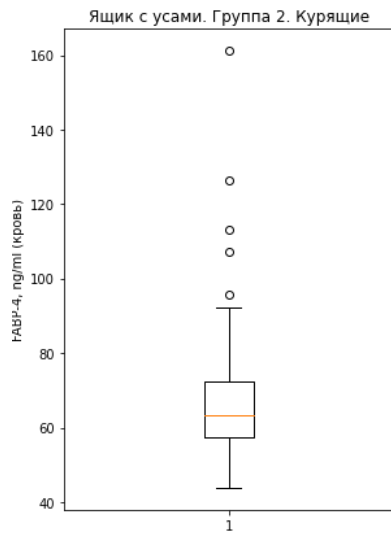
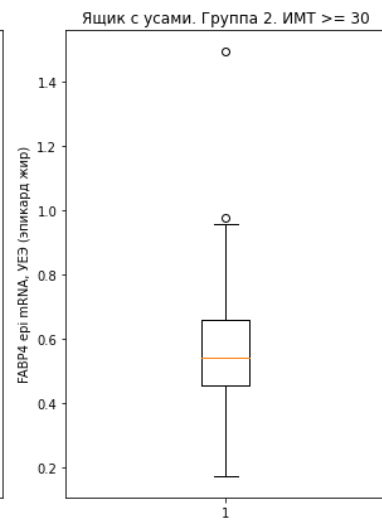
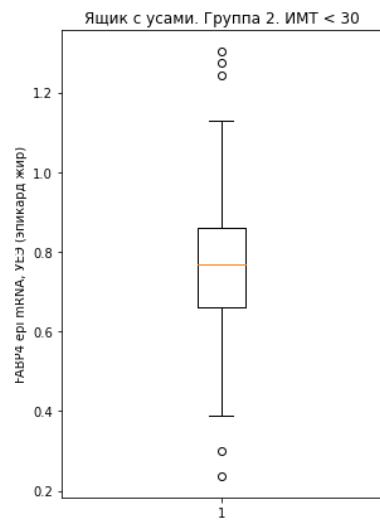
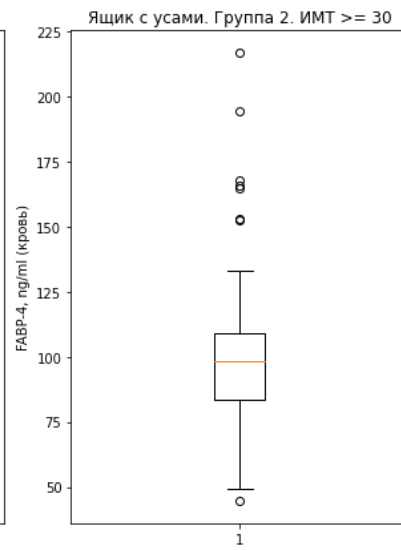
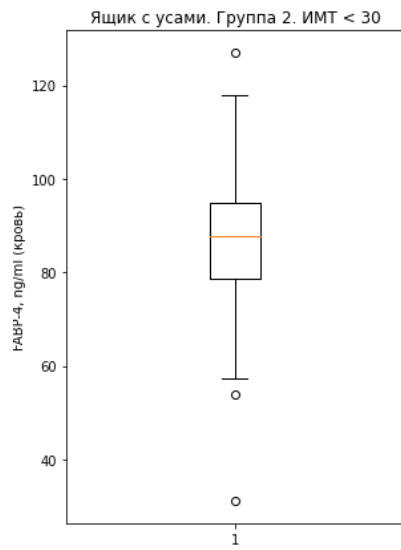


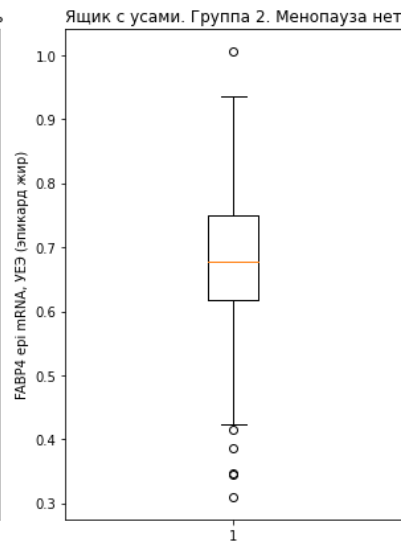
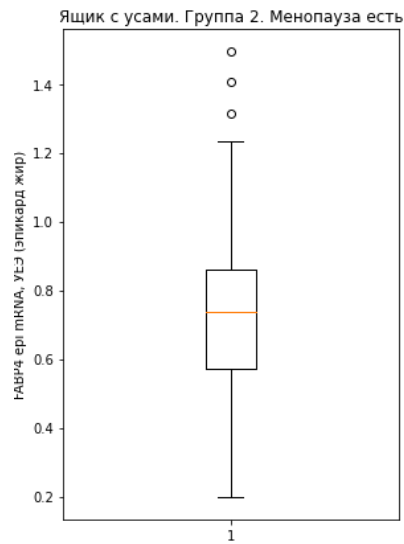
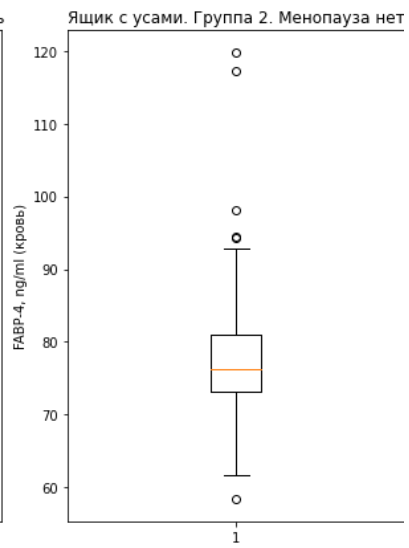
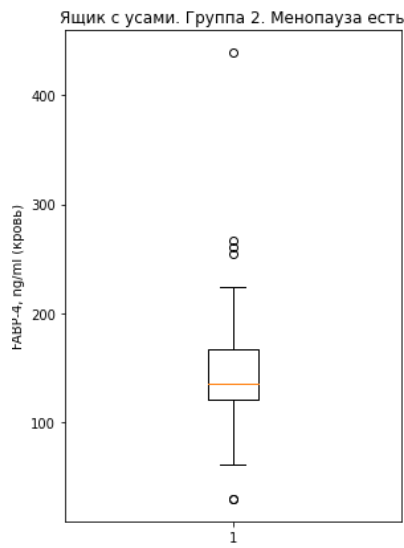
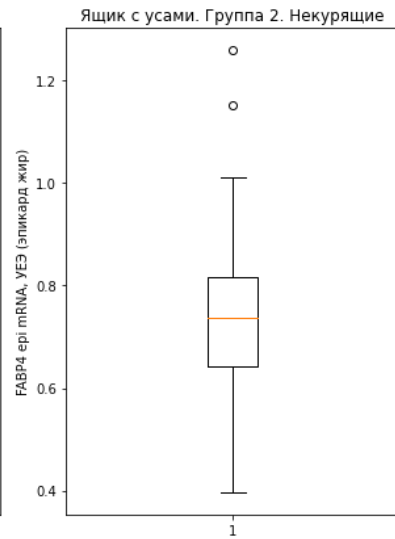
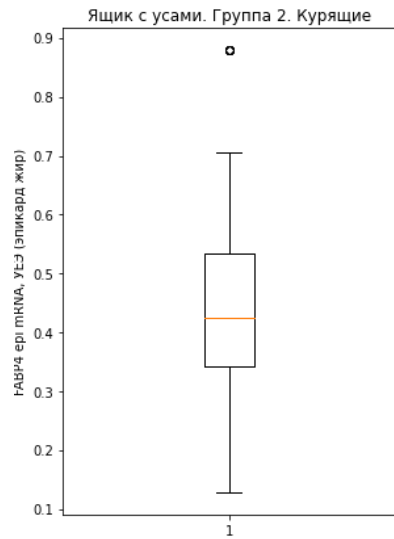


Приложение 1.

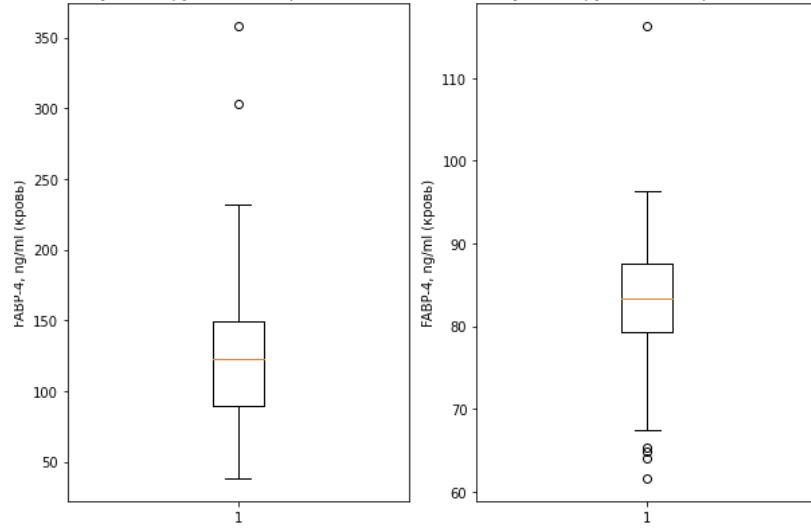
Графики ящиков с усами для репрезентативных факторов влияния на уровень FABP4 в крови и эпидуральном жире для группы людей с ИБС:



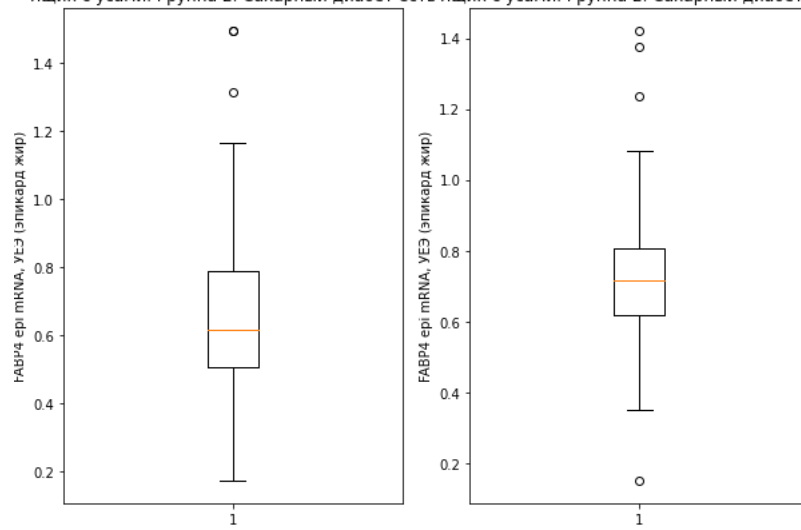




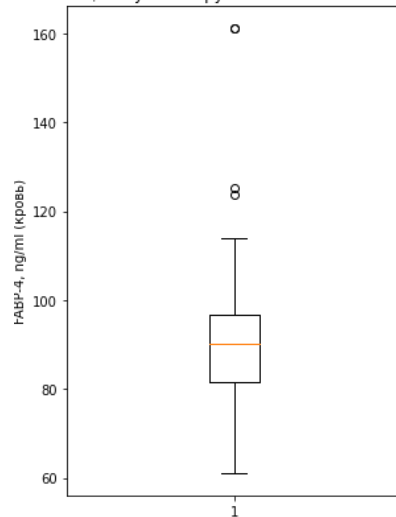
Ящик с усами. Группа 2. Сахарный диабет есть Ящик с усами. Группа 2. Сахарный диабет нет



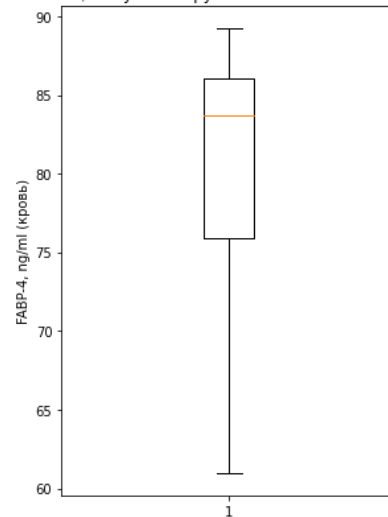
Ящик с усами. Группа 2. Сахарный диабет есть Ящик с усами. Группа 2. Сахарный диабет нет

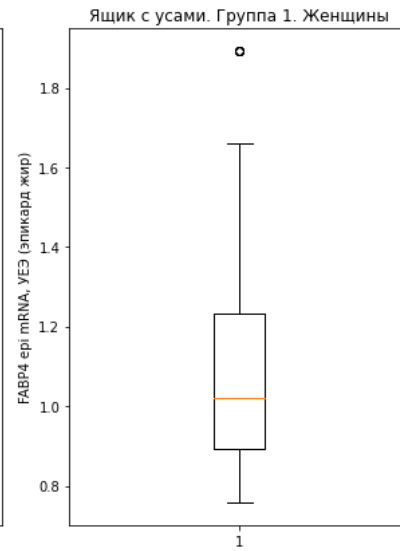
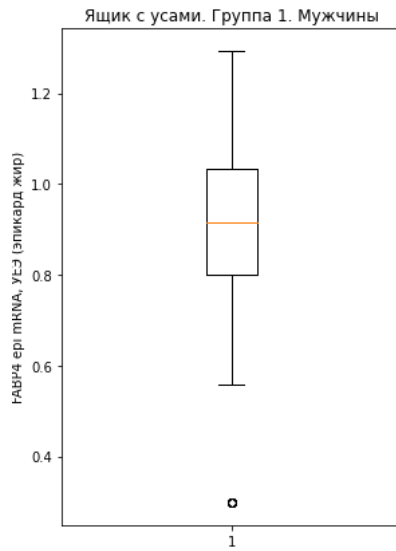
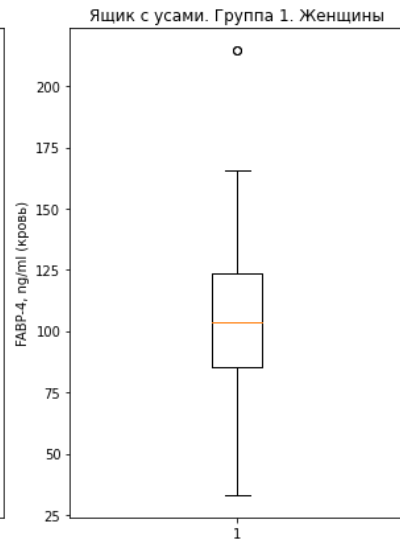
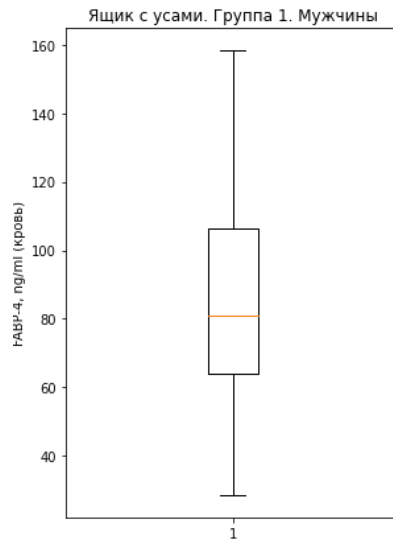
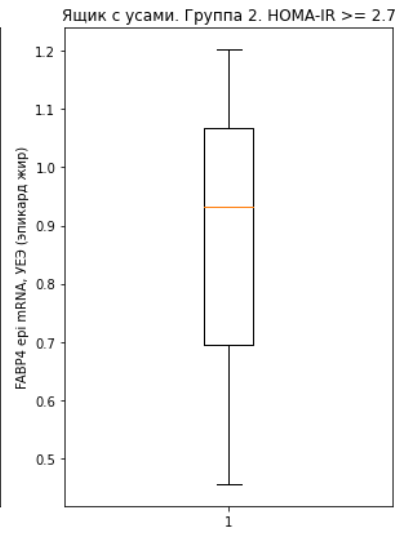
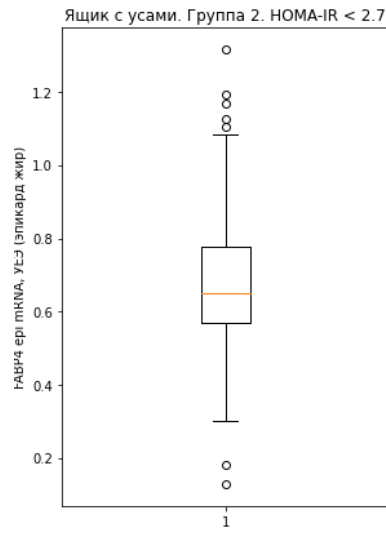


Ящик с усами. Группа 2. НОМА-IR < 2.7

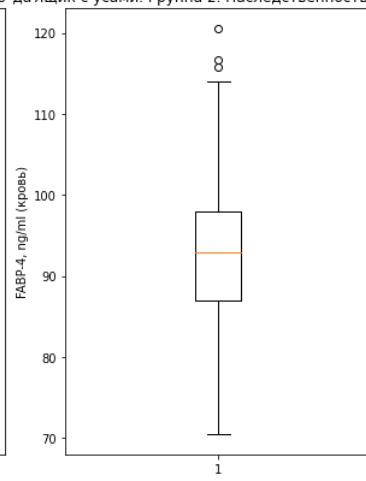
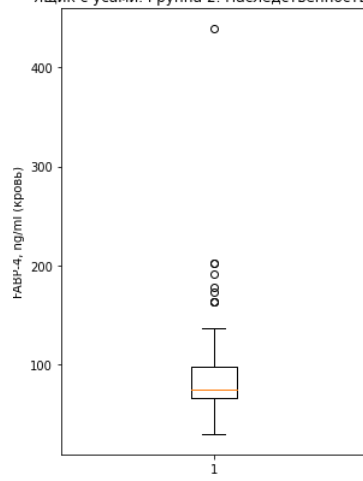


Ящик с усами. Группа 2. НОМА-IR >= 2.7

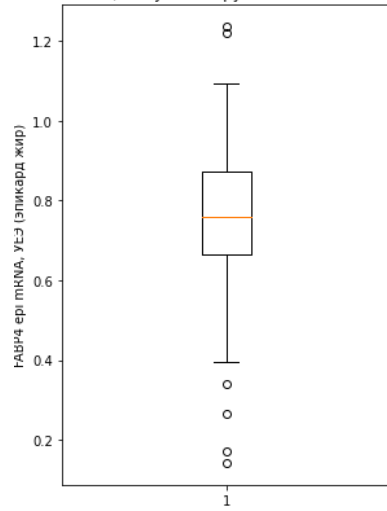




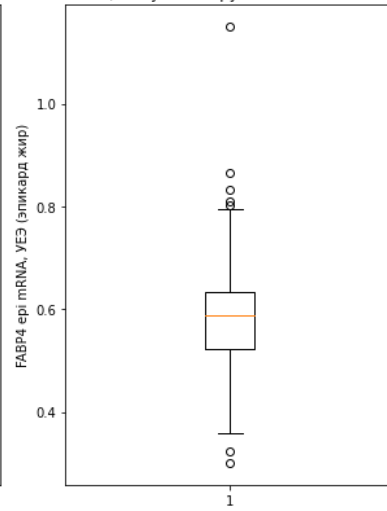
Ящик с усами. Группа 2. Наследственность-да Ящик с усами. Группа 2. Наследственность-нет



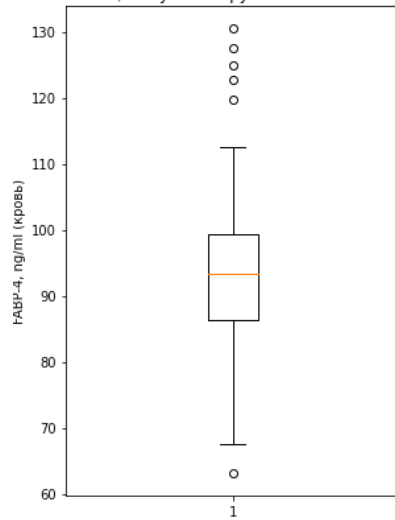
Ящик с усами. Группа 2. ТГ < 1.7



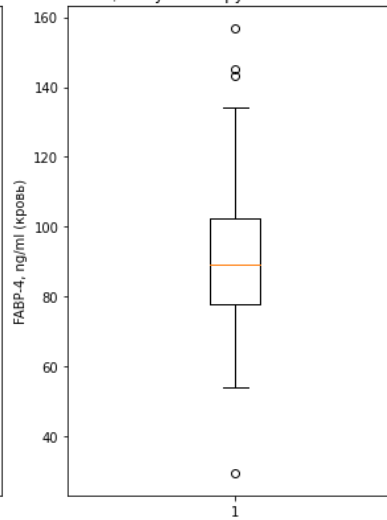
Ящик с усами. Группа 2. ТГ >= 1.7

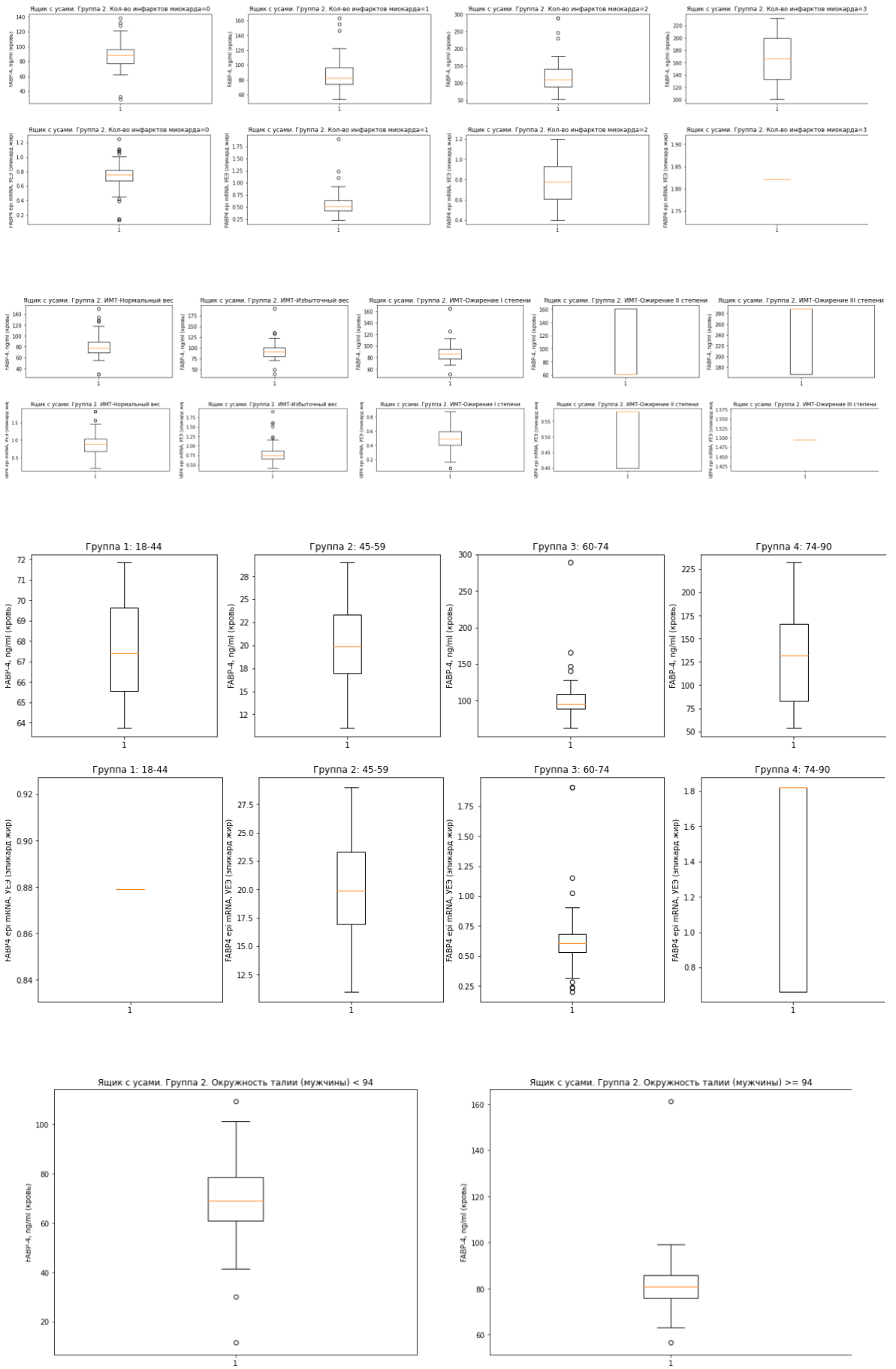


Ящик с усами. Группа 2. ТГ < 1.7

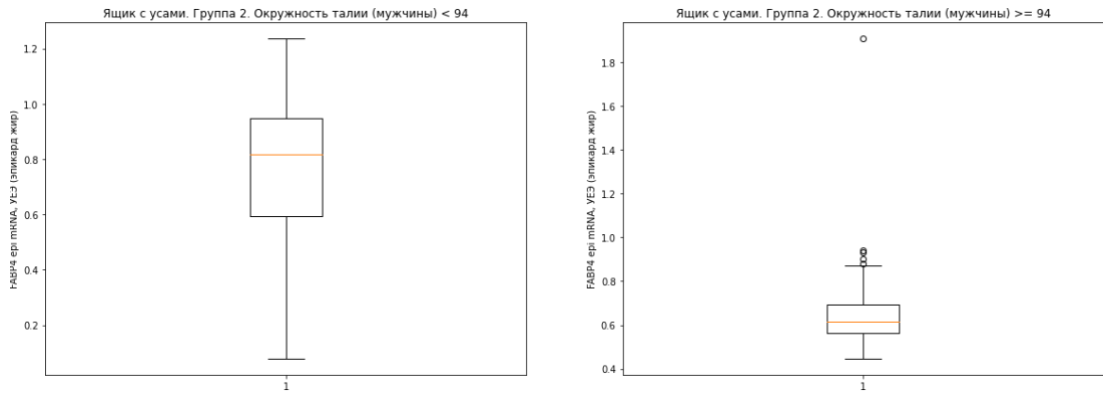


Ящик с усами. Группа 2. ТГ >= 1.7



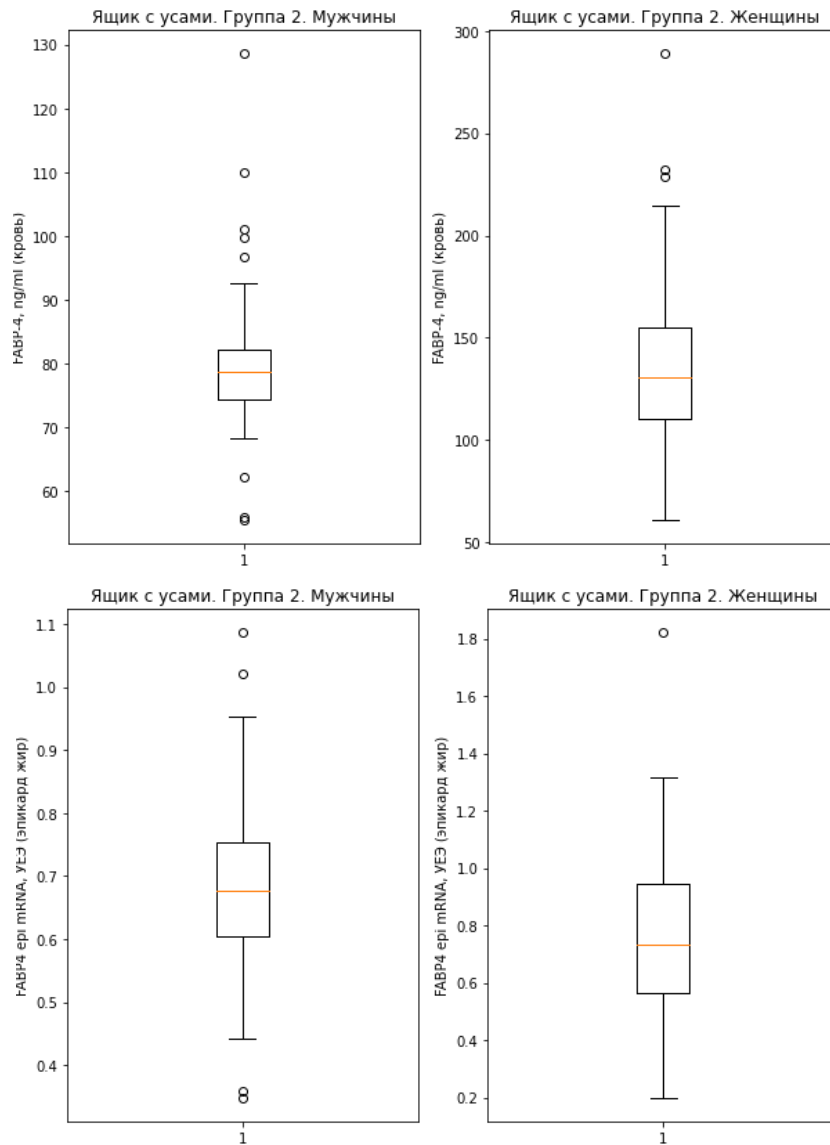






Приложение 2.

Графики ящиков с усами для репрезентативных факторов влияния на уровень FABP4 в крови и эпидуральном жире для группы людей без ИБС:



Приложение 3.