

ГУАП

КАФЕДРА № 41

ОТЧЕТ  
ЗАЩИЩЕН С ОЦЕНКОЙ  
ПРЕПОДАВАТЕЛЬ

ассистент

должность, уч. степень, звание

подпись, дата

В. В. Боженко

инициалы, фамилия

## ОТЧЕТ О ЛАБОРАТОРНОЙ РАБОТЕ

Анализ зависимостей между признаками в двумерном наборе данных

по курсу: ВВЕДЕНИЕ В АНАЛИЗ ДАННЫХ

РАБОТУ ВЫПОЛНИЛ

СТУДЕНТ ГР. № 4917

подпись, дата

Е.А. Ясиновский

инициалы, фамилия

Санкт-Петербург 2022

**Цель работы:** изучения связи между признаками набора данных.

Вариант 2: Файл liver.csv, в котором предоставлены данные о анализах для диагностирования заболеваний печени у пациентов.

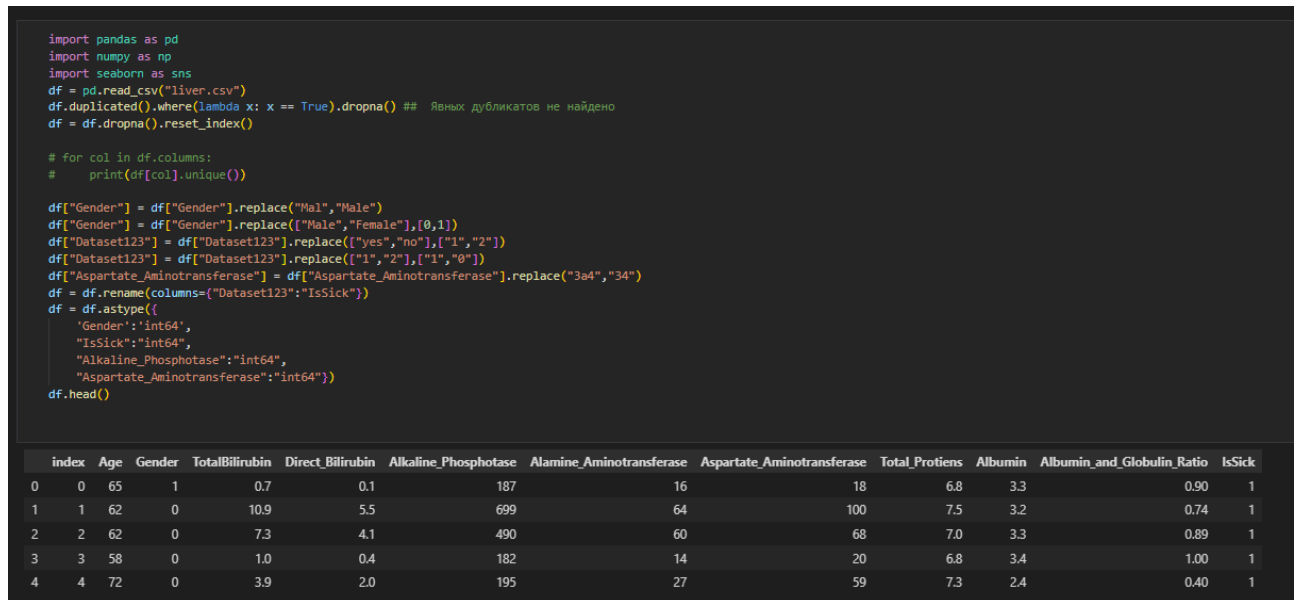


Рисунок 1 — Предварительная подготовка данных, чистка и дубликатов, устранение некорректных строк

Далее были построены графики рассеяния для интересующих меня параметров

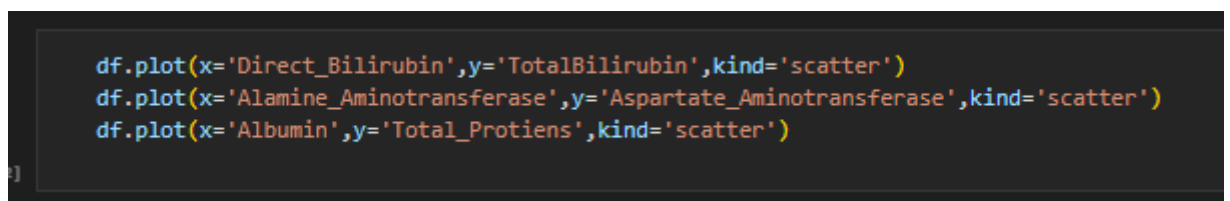


Рисунок 2 — Вызов отрисовки графиков

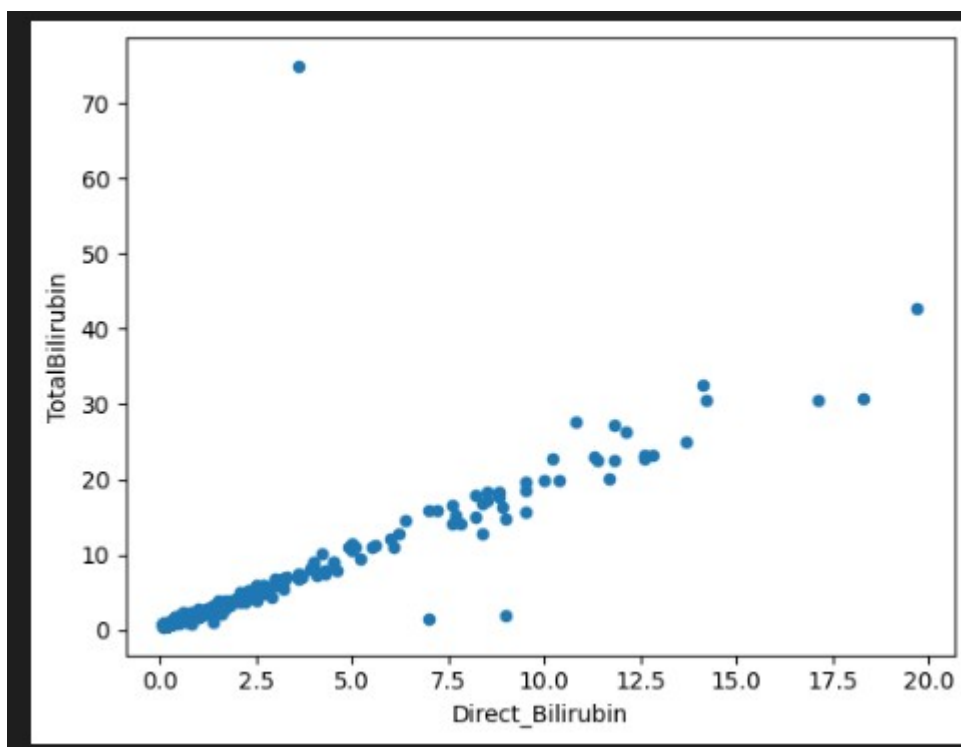


Рисунок 3 — Зависимость общего билирубина от прямого

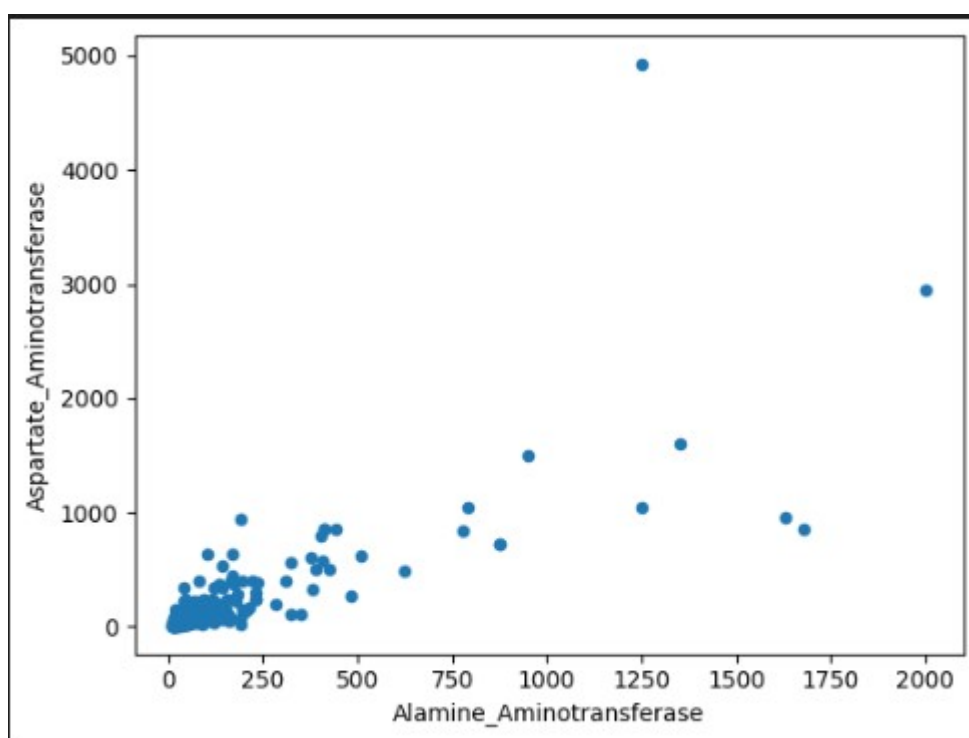


Рисунок 4 — Зависимость аспартатаминотрансфераза от аламиноаминотрансферазы

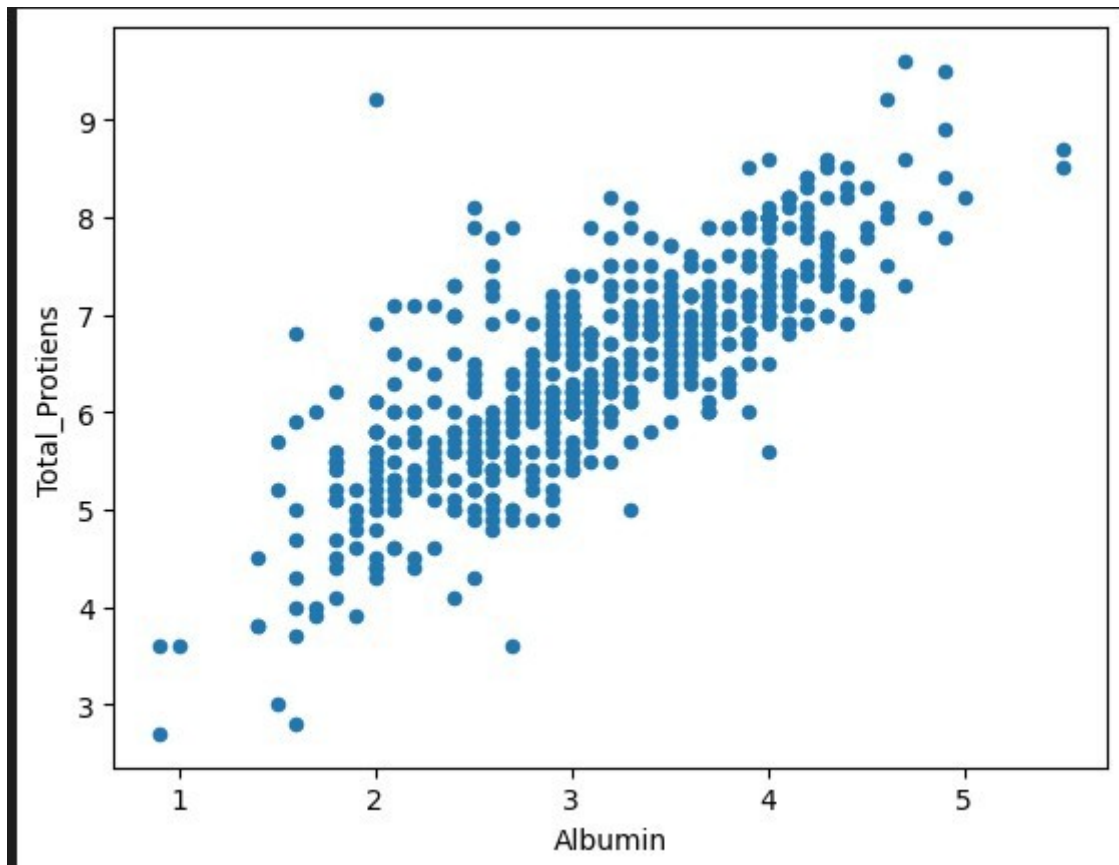


Рисунок 5 — Зависимость общего количества белков от альбумина

Исходя из графиков выше можно увидеть явную прямую зависимость общего количества билирубина от прямого, а также некую размытую прямую зависимость уровня белков от уровня альбумина. В связи с тем, почти все точки графика зависимости аспартатаминотрансфераза от аламиноаминотрансферазы лежат в одной области сложно точно говорить о их прямой зависимости, однако это можно предположить, учитывая то, что крайние правые точки этого графика лежат выше чем его левые точки.

Далее мною была построена матрица рассеяния для всех показателей датафрейма

```
pd.plotting.scatter_matrix(df,figsize=(20,20))
```

Рисунок 6 — Построение матрицы рассеяния

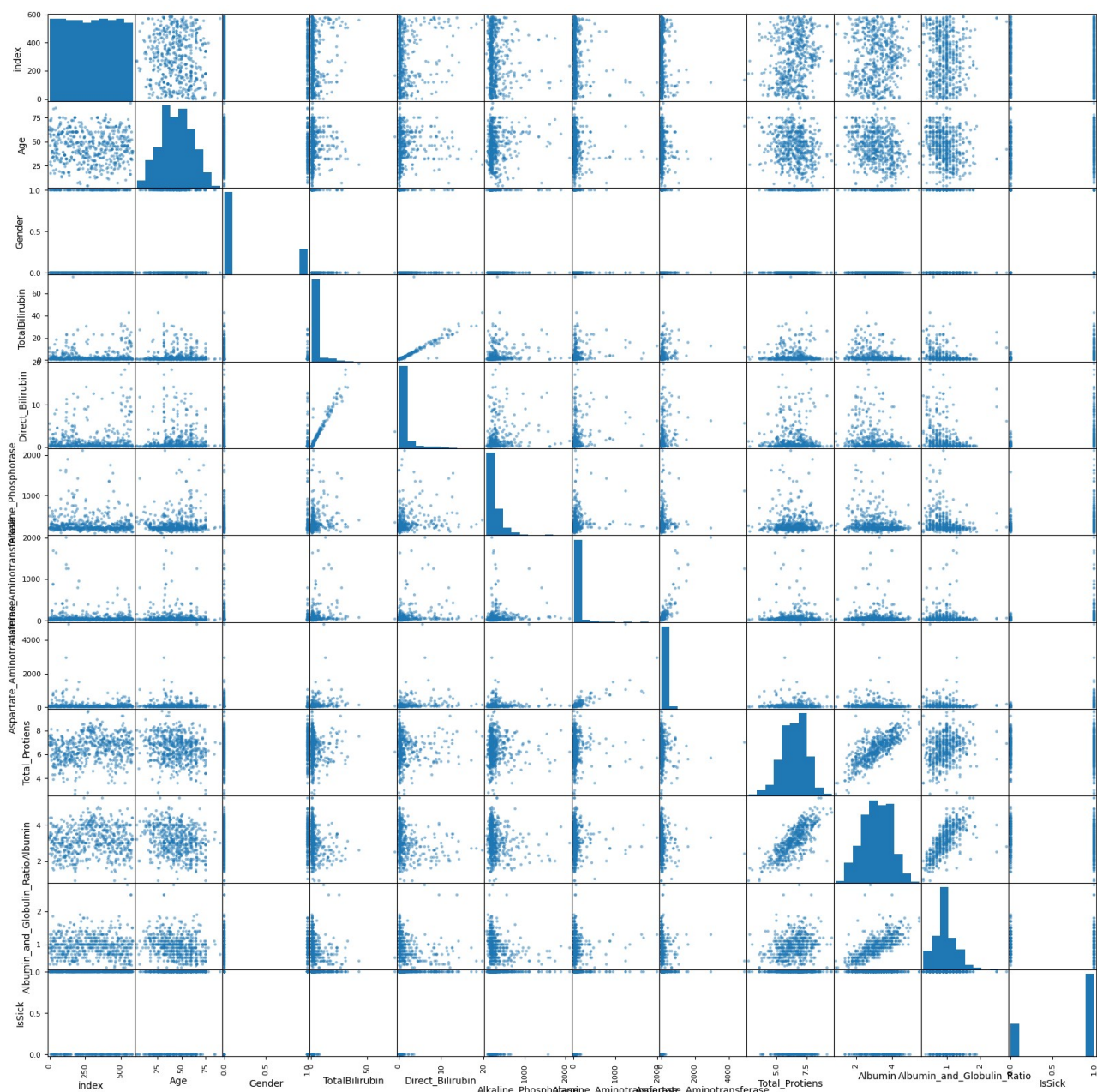


Рисунок 7 — Матрица рассеяния

Уже исходя из данных полученных в матрице рассеяния можно было сделать вывод зависимости ключевого показателя от некоторых других параметров, а именно: билирубина, щелочной фосфатазы, аламиноаминотрансферазы и ампарататаминотрансферазы, и почти не зависит от уровня белков, уровня альбумина и отношения уровня альбумина к уровню глобулина.

Проверить эти предположения мы можем посредством нахождения коэффициента Пирсона для заданных показателей.

```

BilirubCorr = df['TotalBilirubin'].corr(df["Direct_Bilirubin"])
print(BilirubCorr)
amf = df["Alamine_Aminotransferase"].corr(df["Aspartate_Aminotransferase"])
print(amf)

```

0.8745632082572773  
0.7919116749914477

Рисунок 8 — Поиск коэффициента Пирсона

Рассчитанные коэффициенты подтверждают предположение о корреляции выбранных показателей, так как коэффициента Пирсона не положителен, и находится на достаточном удалении от нуля, чтобы не считать это погрешностью.

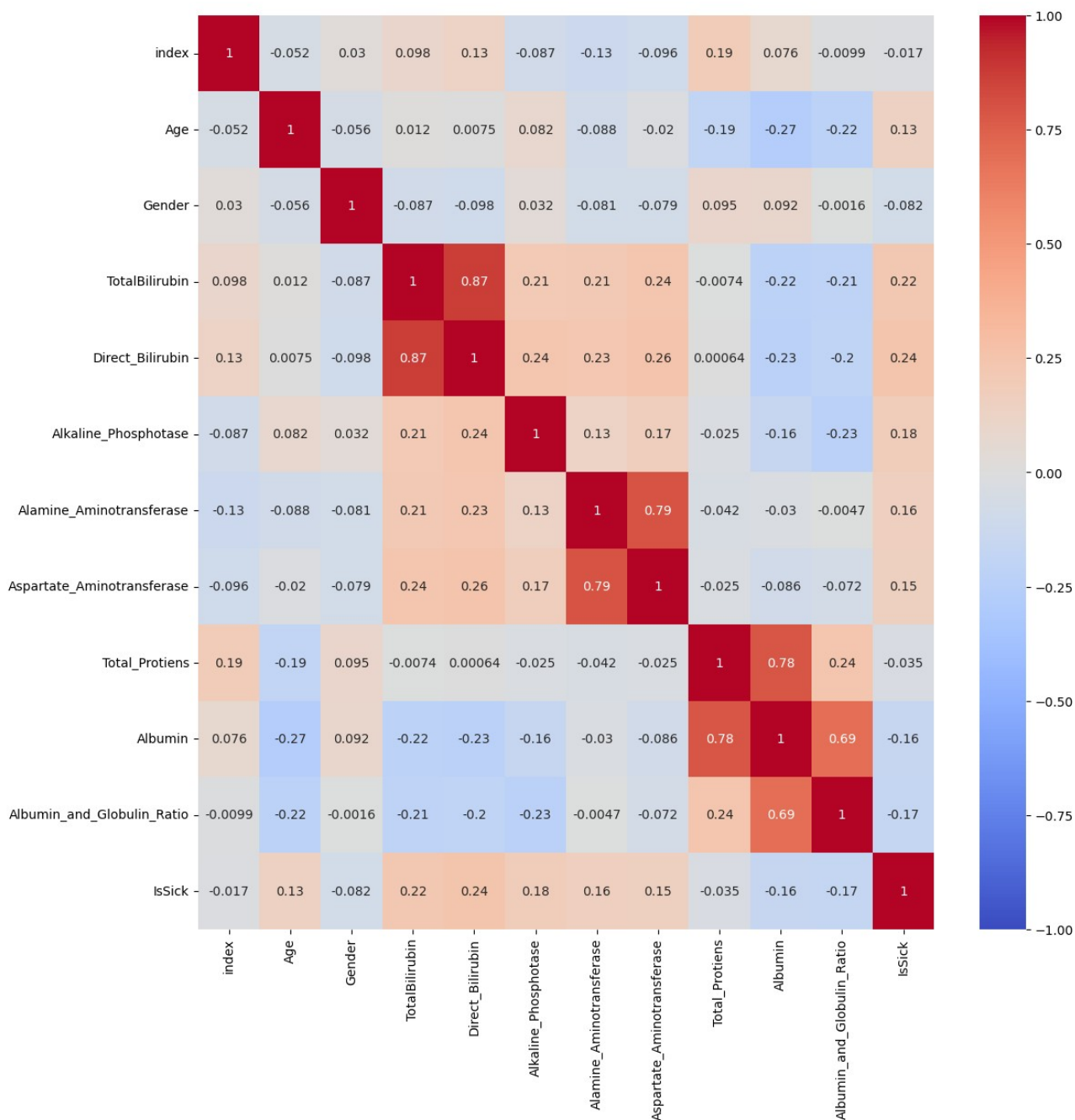
Также для наглядного визуального поиска корреляции и ковариации можно построить тепловую карту корреляции.

```

import matplotlib.pyplot as plt
fig, ax = plt.subplots(figsize=(13,13))
sns.heatmap(df.corr(), vmin=-1, vmax=1, annot=True, cmap="coolwarm", ax=ax)

```

Рисунок 9 — Построение тепловой карты корреляции



Рисуно 10 — Тепловая карта

При рассмотрении данной тепловой карты зависимости и корреляция\ковариация различных показателей становятся более наглядными, однако скорее всего стоит учесть, что если показатель корреляции близок к нулю(слишком мал), то стоит относиться к такому показателю не как с слабой корреляции, а скорее как к погрешности.

Исходный файл Jupyter Notebook находится на сервере GitHub по адресу <https://github.com/EgorYasinovskiy/Data-Analys/blob/master/JIP2/main.ipynb>

**Вывод:** в ходе выполнения данной лабораторной работы научился строить графики зависимости параметров в датафрейме а также матрицы рассеяния среди всех или указанных параметров в датафрейме. С помощью этого произвел анализ данных о пациентах и нашел зависимости уровня некоторых белков и гормонов от диагноза пациента. В ходе данной лабораторной работы выяснилось, что у пациентов с больной печенью обычно всегда выше

уровень следующих показателей: билирубин, щелочная фосфатаза, аламиноаминотранфераза и ампартатаминотрансфераза