



Министерство науки и высшего образования Российской Федерации  
Федеральное государственное бюджетное образовательное учреждение  
высшего образования  
«Московский государственный технический университет  
имени Н.Э. Баумана  
(национальный исследовательский университет)»  
(МГТУ им. Н.Э. Баумана)

---

ФАКУЛЬТЕТ

ИНФОРМАТИКА И СИСТЕМЫ УПРАВЛЕНИЯ

КАФЕДРА

СИСТЕМЫ ОБРАБОТКИ ИНФОРМАЦИИ И УПРАВЛЕНИЯ (ИУ5)

# О Т Ч Е Т

## по лабораторной работе

по дисциплине: Технологии машинного обучения

на тему: Разведочный анализ данных. Исследование и визуализация данных

Студент ИУ5-62Б  
(Группа)

(Подпись, дата)

Карягин А.Д.  
(И.О.Фамилия)

Руководитель

Ю.Е.  
Гапанюк

(Подпись, дата)  
(И.О.Фамилия)

2020 г.

# Лабораторная работа №1

## 1) Текстовое описание набора данных

В качестве набора данных мы будем использовать набор данных Diabets dataset <https://scikit-learn.org/stable/datasets/index.html#toy-datasets> Для каждого из  $n = 442$  больных сахарным диабетом были получены десять исходных переменных, возраст, пол, индекс массы тела, среднее артериальное давление и шесть измерений сыворотки крови, а также интересующая нас реакция - количественная мера прогрессирования заболевания через год после исходного уровня.

In

```
import numpy as np import pandas as pd import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline sns.set(style="ticks")

data = pd.read_csv('data/diabetes.tab.txt', sep="\t")
```

[7]:

## 2) Основные характеристики датасета

In [9]:

Out [9]:

	AG E	SE X	B MI	BP	S 1	S2	S3	S 4	S5	S 6	Y
0	59	2	32 .1	101 .0	1 5 7	93. 2	38 .0	4 .0	4.85 98	8 7	15 1
1	48	1	21 .6	87. 0	1 8 3	103 .2	70 .0	3 .0	3.89 18	6 9	75
2	72	2	30 .5	93. 0	1 5 6	93. 6	41 .0	4 .0	4.67 28	8 5	14 1

3	24	1	25	84.	1	131	40	5	4.89	8	20
			.3	0	9	.4	.0	.	03	9	6
					8			0			
4	50	1	23	101	1	125	52	4	4.29	8	13
			.0	.0	9	.4	.0	.	05	0	5
					2			0			

In [10]:

Out[10]: (442, 11)

In [11]:

Всего строк: 442

In [12]:

Out[12]: Index(['AGE', 'SEX', 'BMI', 'BP', 'S1', 'S2', 'S3', 'S4', 'S5', 'S6', 'Y'], dtype='object')

In [13]:

Out[13]: AGE int64  
SEX int64  
BMI float64  
BP float64  
S1 int64  
S2 float64  
S3 float64  
S4 float64  
S5 float64  
S6 int64  
Y int64  
dtype: object

In [14]:

AGE - 0  
SEX - 0  
BMI - 0  
BP - 0  
S1 - 0  
S2 - 0  
S3 - 0  
S4 - 0  
S5 - 0  
S6 - 0  
Y - 0

In [15]:

	AGE	SEX S1	BMI S2	BP S3	
count	442.000000	442.000000	442.000000	442.000000	442.000000
mean	48.518100	1.468326	26.375792	94.647014	189.140271
std	13.109028	0.499561	4.418122	13.831283	34.608052
min	19.000000	1.000000	18.000000	62.000000	97.000000
25%	38.250000	1.000000	23.200000	84.000000	164.250000
50%	50.000000	1.000000	25.700000	93.000000	186.000000
75%	59.000000	2.000000	29.275000	105.000000	209.750000
max	79.000000	2.000000	42.200000	133.000000	301.000000

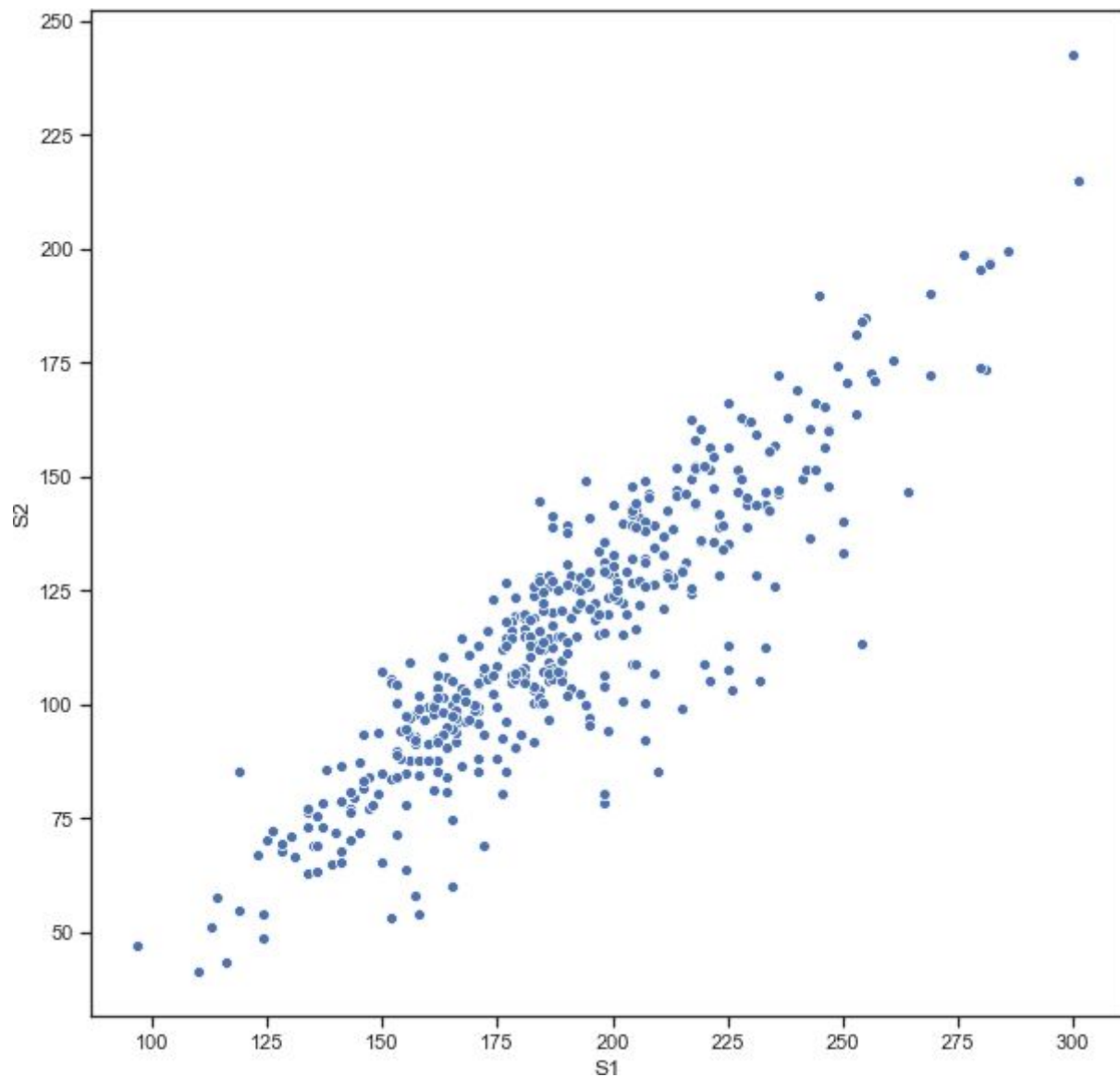
In [21]:

Out [21]: array([2, 1], dtype=int64)

### 3) Визуальное исследование датасета

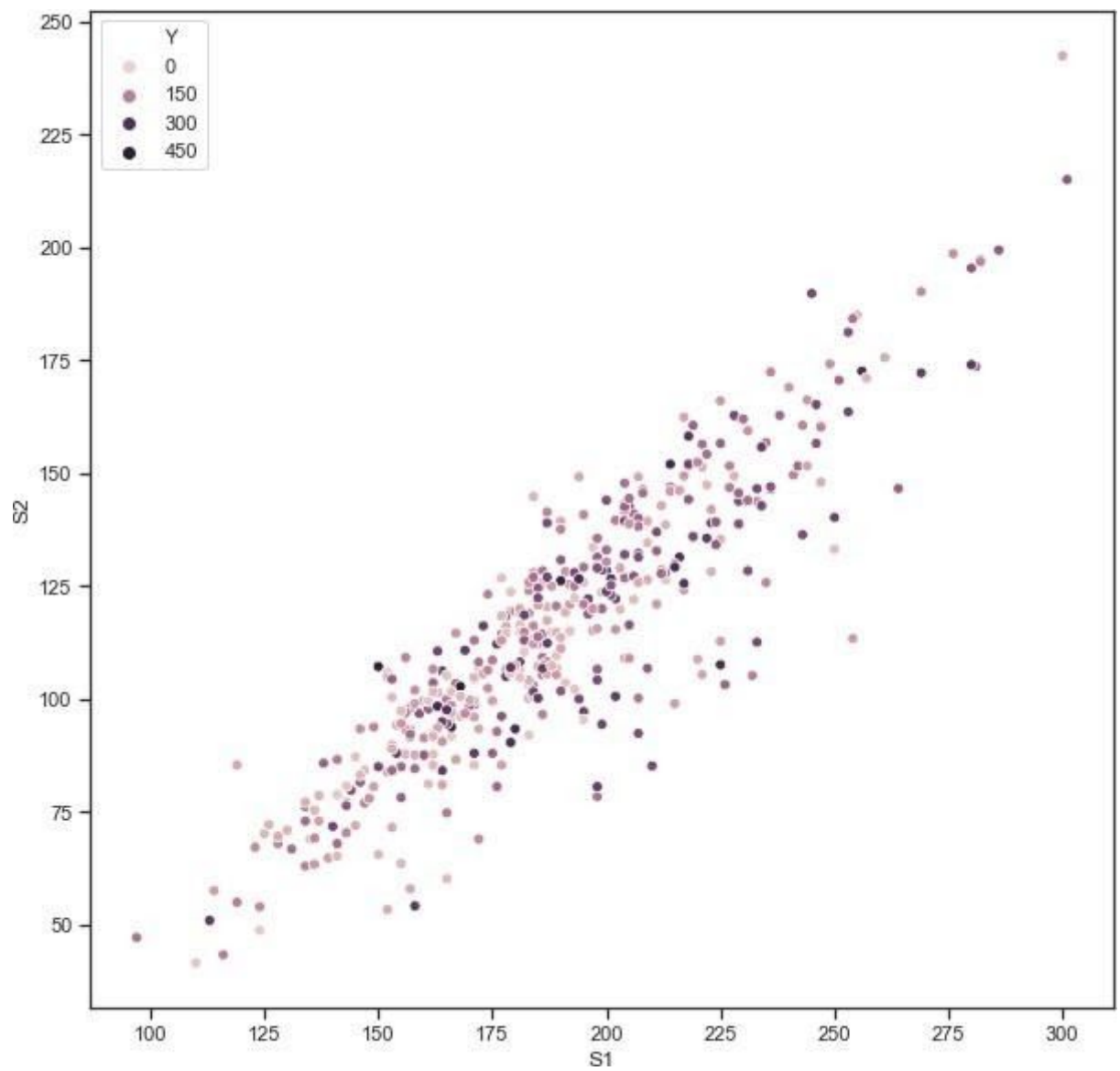
In [38]:

Out [38]: <matplotlib.axes.\_subplots.AxesSubplot at 0xe70c610>



In [42]:

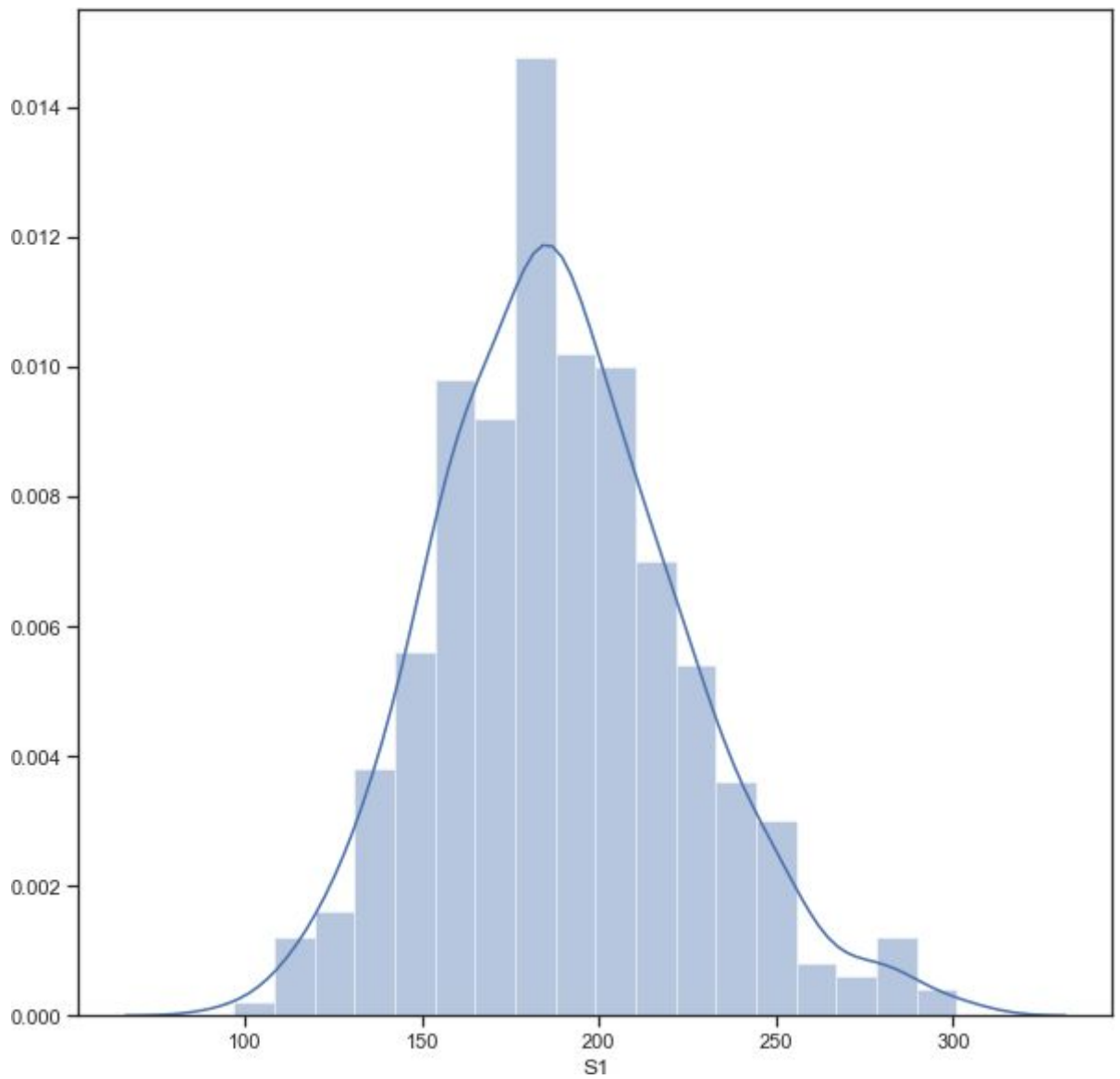
Out [42]: <matplotlib.axes.\_subplots.AxesSubplot at 0xfd81e70>



## Гистограмма

In [43]:

Out[43]: <matplotlib.axes.\_subplots.AxesSubplot at 0xfd816b0>

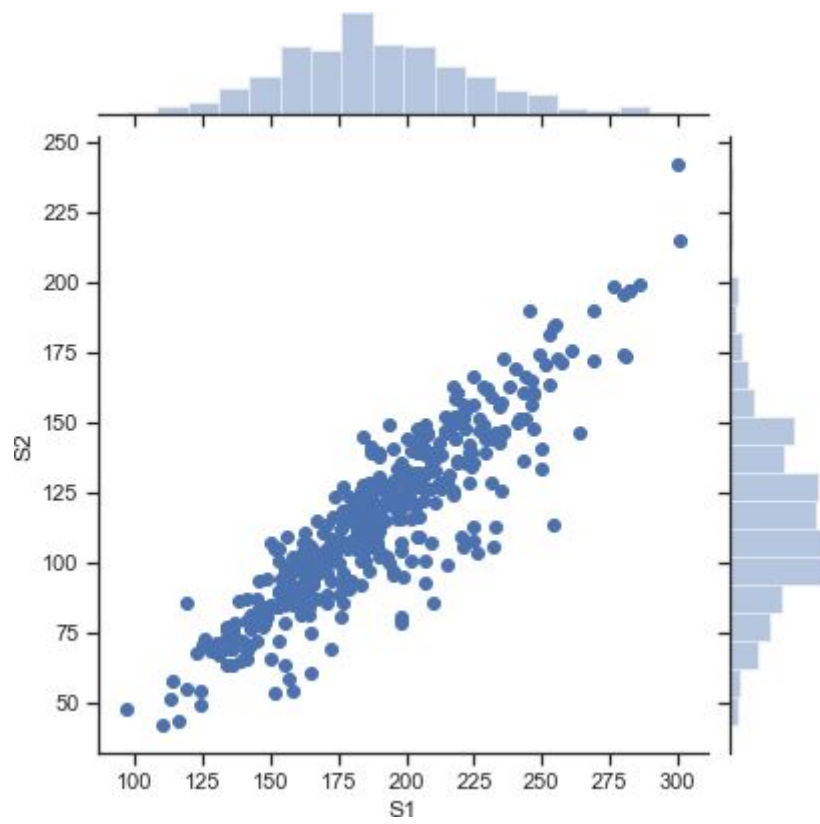


## Jointplot

Комбинация гистограмм и диаграмм рассеивания.

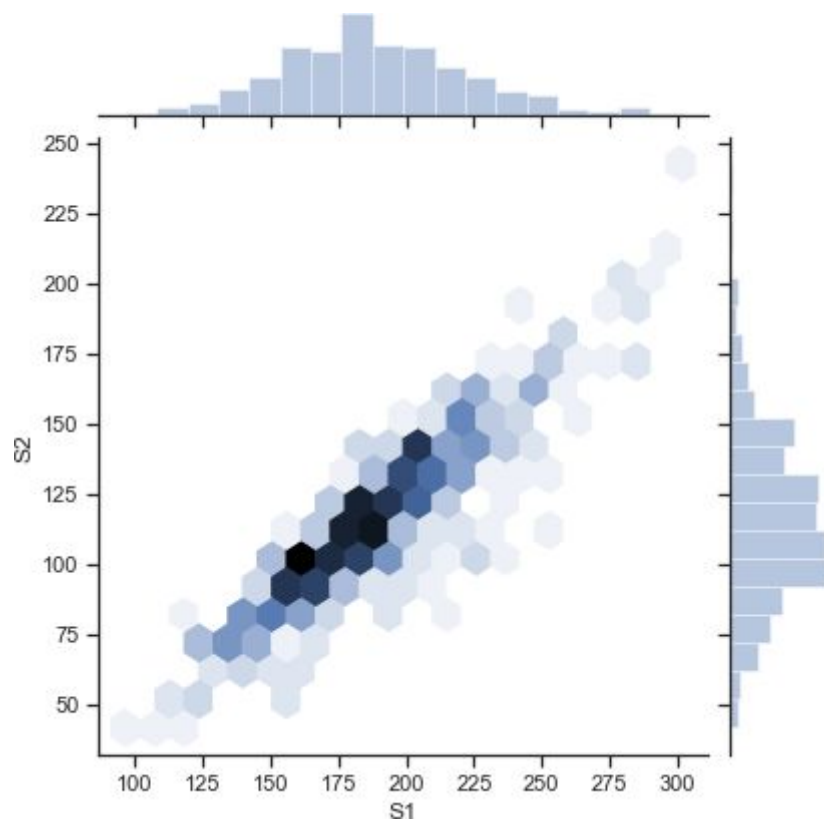
In [44]:

Out[44]: <seaborn.axisgrid.JointGrid at 0xfd663b0>



In [48]:

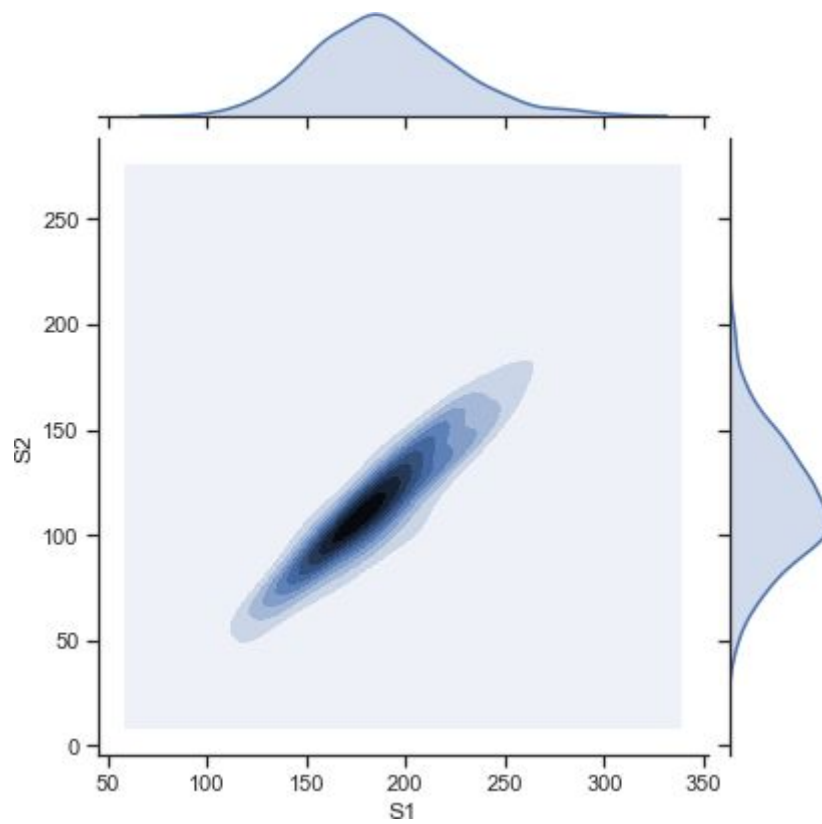
Out[48]: <seaborn.axisgrid.JointGrid at 0x1041bab0>



In [49]:

Out[49]: <seaborn.axisgrid.JointGrid at 0x1079f450>

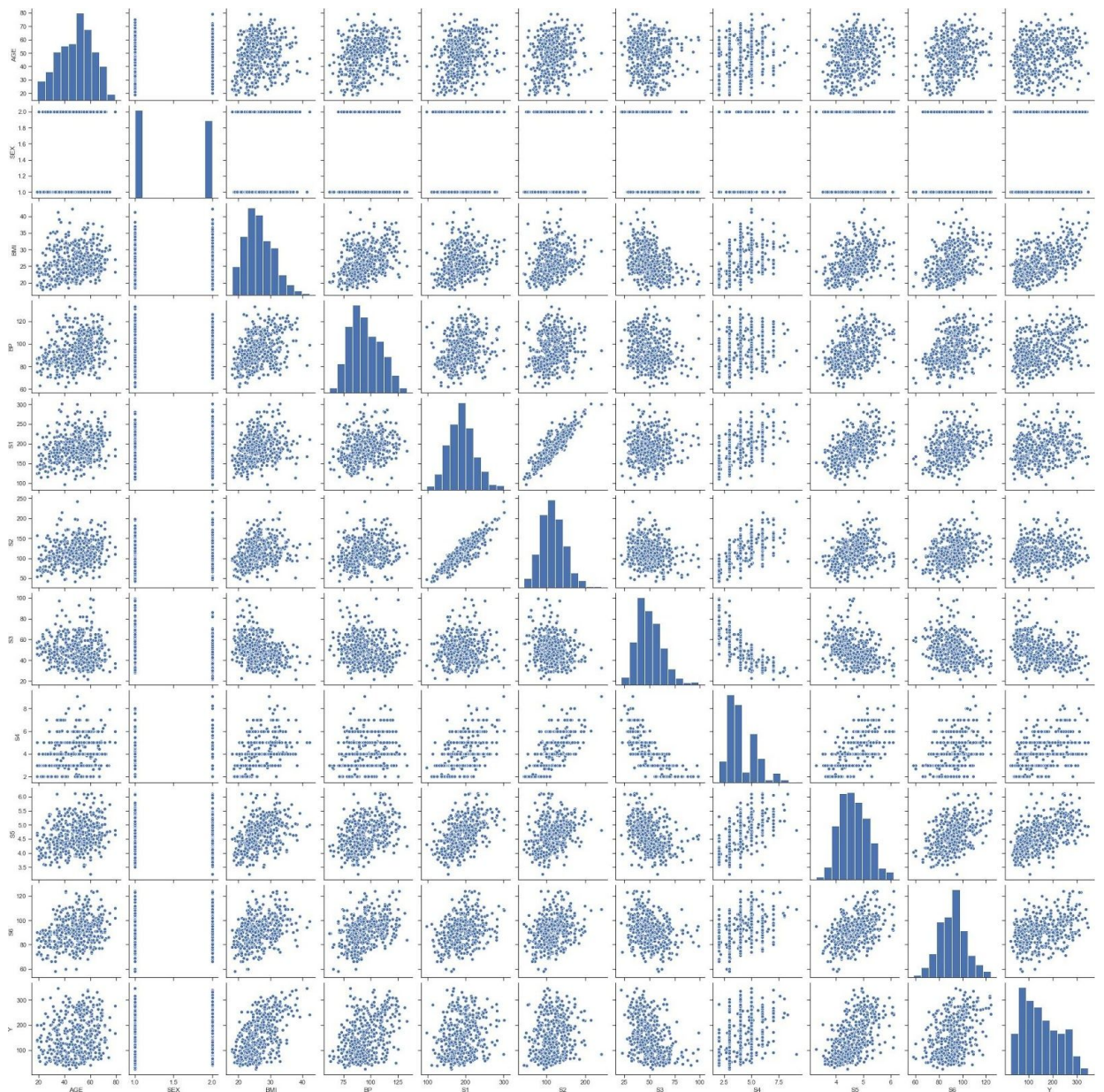




## "Парные диаграммы"

In [54]:

Out [54]: <seaborn.axisgrid.PairGrid at 0x207eb810>



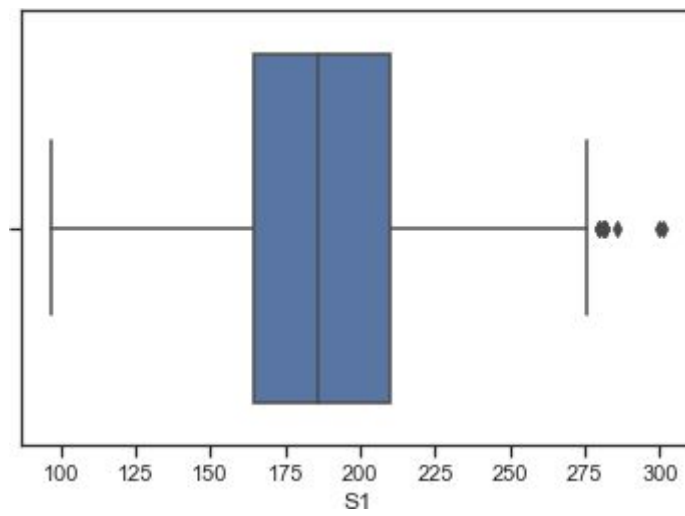
In [ ]:

## Ящик с усами

Отображает одномерное распределение вероятности.

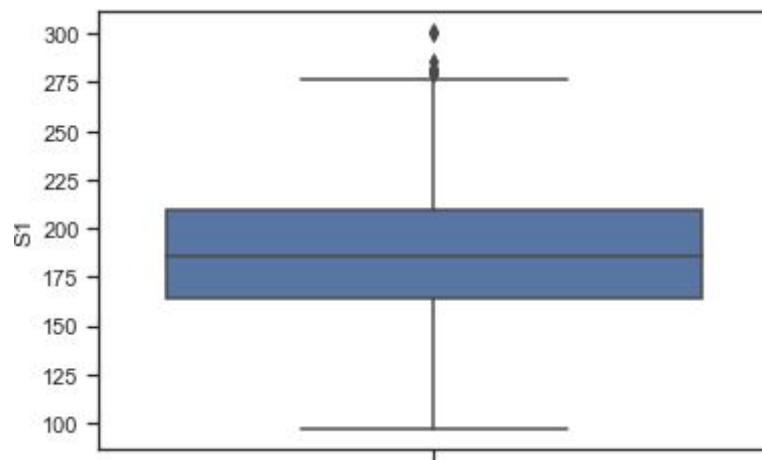
In [56]:

Out[56]: <matplotlib.axes.\_subplots.AxesSubplot at 0x3066f290>



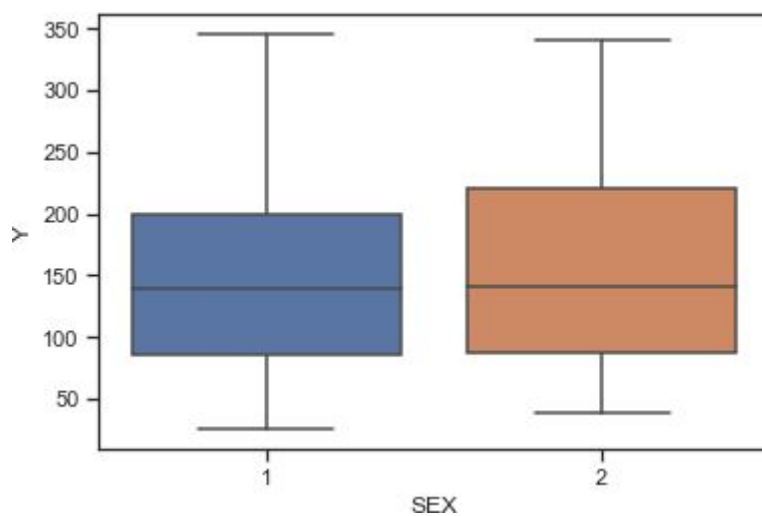
In [57]:

Out[57]: <matplotlib.axes.\_subplots.AxesSubplot at 0x3095a350>



In [61]:

Out[61]: <matplotlib.axes.\_subplots.AxesSubplot at 0x301089d0>

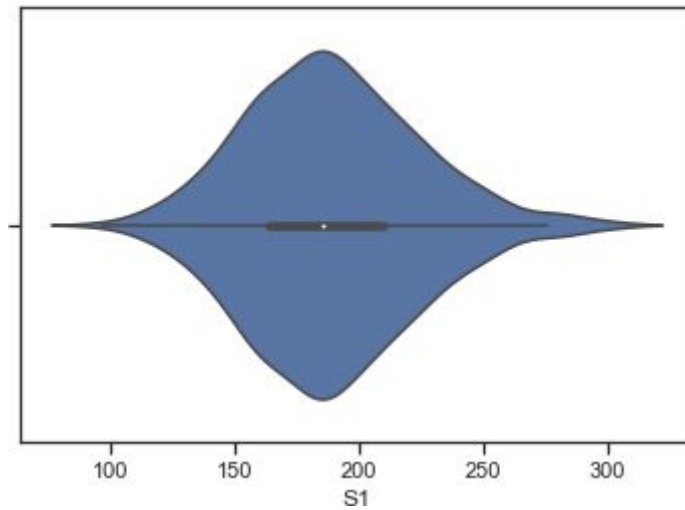


**Violin plot**

Похоже на предыдущую диаграмму, но по краям отображаются распределения плотности

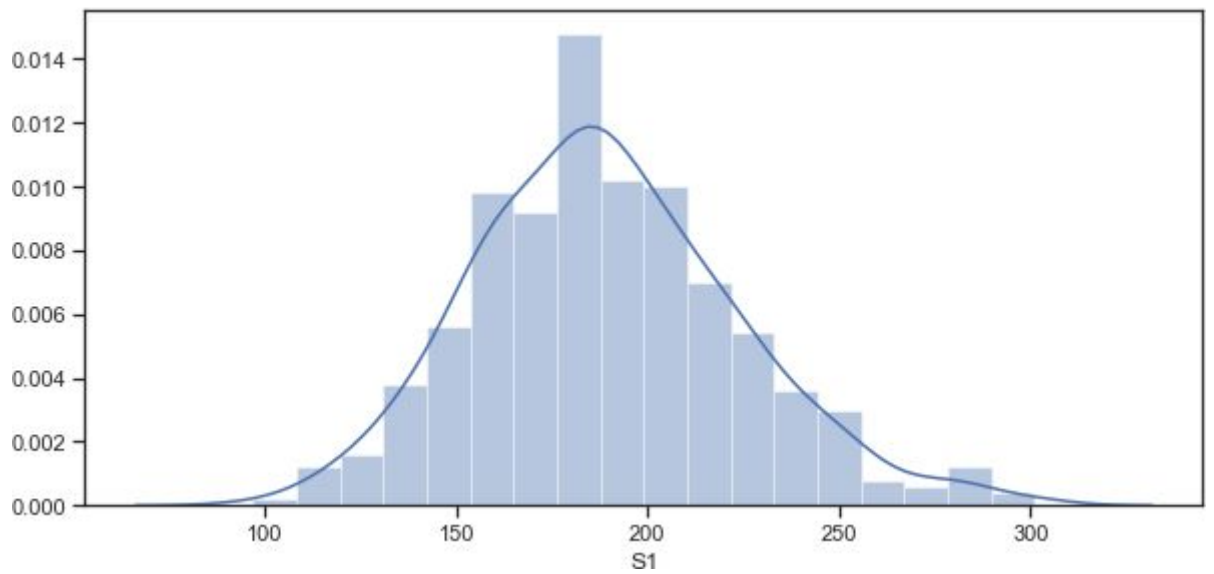
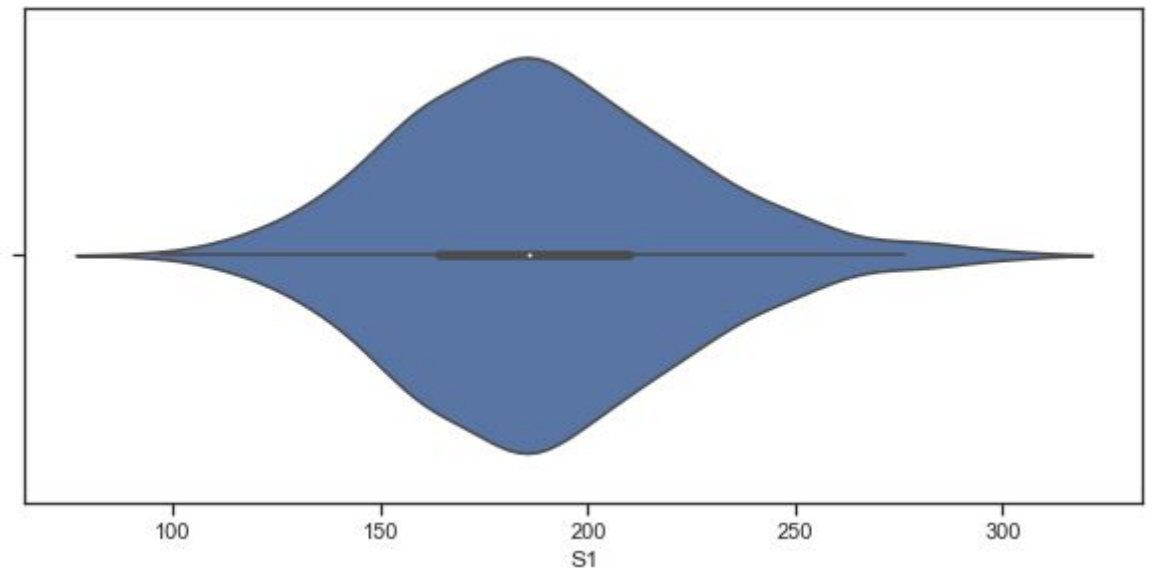
In [59]:

Out[59]: <matplotlib.axes.\_subplots.AxesSubplot at 0x313a38d0>



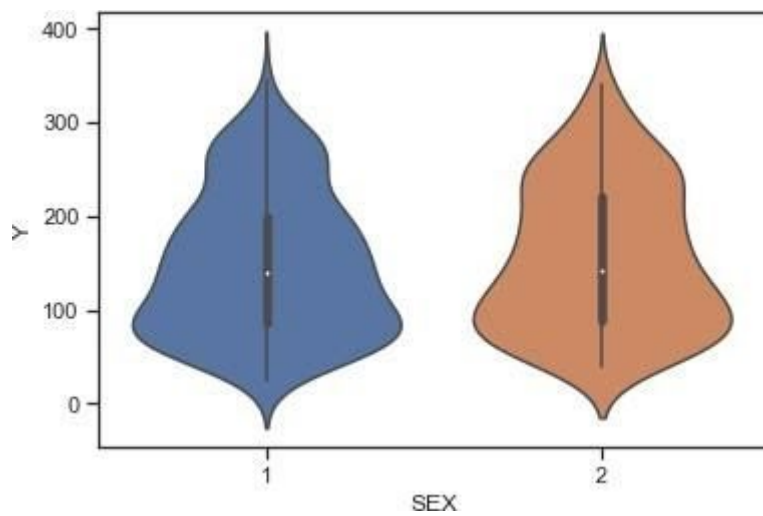
In [62]:

Out[62]: <matplotlib.axes.\_subplots.AxesSubplot at 0x2e5efe10>



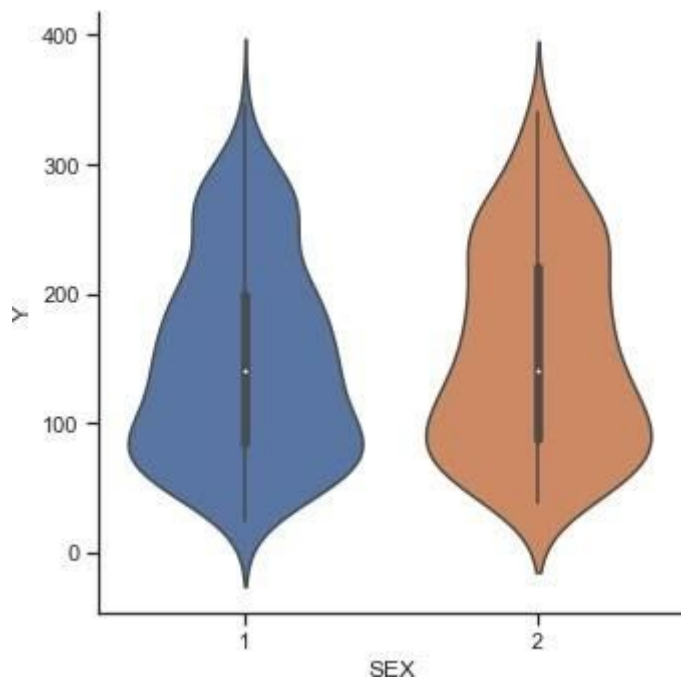
In [65]:

Out[65]: <matplotlib.axes.\_subplots.AxesSubplot at 0x305b53b0>



In [79]:

Out[79]: <seaborn.axisgrid.FacetGrid at 0x340379b0>



## 4) Информация о корреляции признаков

In [68]:

Out[68]:

	AGE	SEX	BMI	BP	S1	S2	S3	S4	
AGE	1.00000 0	0.17373 7	0.18508 5	0.33542 8	0.2600 61	0.21924 3	-0.0751 81	0.20384 1	0.270 7
SEX	0.17373 7	1.00000 0	0.08816 1	0.24101 0	0.0352 77	0.14263 7	-0.3790 90	0.33211 5	0.149 9
BMI	0.18508 5	0.08816 1	1.00000 0	0.39541 1	0.2497 77	0.26117 0	-0.3668 11	0.41380 7	0.446 1
BP	0.33542 8	0.24101 0	0.39541 1	1.00000 0	0.2424 64	0.18554 8	-0.1787 62	0.25765 0	0.393 4
S1	0.26006 1	0.03527 7	0.24977 7	0.24246 4	1.0000 00	0.89666 3	0.05151 9	0.54220 7	0.515 5
S2	0.21924 3	0.14263 7	0.26117 0	0.18554 8	0.8966 63	1.00000 0	-0.1964 55	0.65981 7	0.318 3
S3	-0.0751 81	-0.3790 90	-0.3668 11	-0.1787 62	0.0515 19	-0.1964 55	1.00000 0	-0.7384 93	-0.39 85
S4	0.20384 1	0.33211 5	0.41380 7	0.25765 0	0.5422 07	0.65981 7	-0.7384 93	1.00000 0	0.617 8
S5	0.27077 4	0.14991 6	0.44615 7	0.39348 0	0.5155 03	0.31835 7	-0.3985 77	0.61785 9	1.000 0
S6	0.30173 1	0.20813 3	0.38868 0	0.39043 0	0.3257 17	0.29060 0	-0.2736 97	0.41721 2	0.464 6
Y	0.18788 9	0.04306 2	0.58645 0	0.44148 2	0.2120 22	0.17405 4	-0.3947 89	0.43045 3	0.565 8

In [69]:

Out[69]:

	AGE	SEX	BMI	BP	S1	S2	S3	S4	
AGE	1.00000 0	0.17373 7	0.18508 5	0.33542 8	0.2600 61	0.21924 3	-0.0751 81	0.20384 1	0.270 7

<b>SE</b>	0.17373	1.00000	0.08816	0.24101	0.0352	0.14263	-0.3790	0.33211	0.149
<b>X</b>	7	0	1	0	77	7	90	5	9
<b>BM</b>	0.18508	0.08816	1.00000	0.39541	0.2497	0.26117	-0.3668	0.41380	0.446
<b>I</b>	5	1	0	1	77	0	11	7	1
<b>BP</b>	0.33542	0.24101	0.39541	1.00000	0.2424	0.18554	-0.1787	0.25765	0.393
	8	0	1	0	64	8	62	0	4
<b>S1</b>	0.26006	0.0352	0.24977	0.24246	1.0000	0.89666	0.05151	0.54220	0.515
	1	77	7	4	00	3	9	7	5
<b>S2</b>	0.21924	0.1426	0.26117	0.18554	0.8966	1.00000	-0.19645	0.65981	0.318
	3	37	0	8	63	0	5	7	3
<b>S3</b>	-0.0751	-0.3790	-0.3668	-0.1787	0.0515	-0.1964	1.00000	-0.7384	-0.39
	81	90	11	62	19	55	0	93	85
<b>S4</b>	0.20384	0.3321	0.41380	0.25765	0.5422	0.65981	-0.73849	1.00000	0.617
	1	15	7	0	07	7	3	0	8
<b>S5</b>	0.27077	0.1499	0.44615	0.39348	0.5155	0.31835	-0.39857	0.61785	1.000
	4	16	7	0	03	7	7	9	0
<b>S6</b>	0.30173	0.2081	0.38868	0.39043	0.3257	0.29060	-0.27369	0.41721	0.464
	1	33	0	0	17	0	7	2	6
<b>Y</b>	0.18788	0.0430	0.58645	0.44148	0.2120	0.17405	-0.39478	0.43045	0.565
	9	62	0	2	22	4	9	3	8

In [70]:

Out[70]:

	AGE	SEX	BMI	BP	S1	S2	S3	S4	
<b>AG</b>	1.00000	0.14658	0.13653	0.24211	0.1822	0.15361	-0.0738	0.16089	0.180
<b>E</b>	0	0	5	1	20	2	46	8	5
<b>SE</b>	0.14658	1.00000	0.08042	0.21573	0.0228	0.11020	-0.3261	0.29733	0.143
<b>X</b>	0	0	4	3	09	8	88	5	1
<b>BM</b>	0.13653	0.08042	1.00000	0.28177	0.1941	0.19858	-0.2498	0.33562	0.344
<b>I</b>	5	4	0	0	71	3	31	5	7
<b>BP</b>	0.24211	0.21573	0.28177	1.00000	0.1880	0.14025	-0.1310	0.20594	0.268
	1	3	0	0	67	3	14	8	8
<b>S1</b>	0.18222	0.02280	0.19417	0.18806	1.0000	0.71722	0.01069	0.39336	0.356
	0	9	1	7	00	9	5	7	2
<b>S2</b>	0.15361	0.11020	0.19858	0.14025	0.7172	1.00000	-0.1333	0.50357	0.242
	2	8	3	3	29	0	32	9	2
<b>S3</b>	-0.0738	-0.3261	-0.2498	-0.1310	0.0106	-0.1333	1.00000	-0.6386	-0.31
	46	88	31	14	95	32	0	33	17
<b>S4</b>	0.16089	0.29733	0.33562	0.20594	0.3933	0.50357	-0.6386	1.00000	0.485
	8	5	5	8	67	9	33	0	4
<b>S5</b>	0.18054	0.14317	0.34472	0.26886	0.3562	0.24225	-0.3117	0.48541	1.000
	4	2	0	3	68	0	75	0	0
<b>S6</b>	0.20178	0.16819	0.26637	0.26456	0.2271	0.19408	-0.2005	0.30739	0.316
	4	9	3	6	39	2	45	7	2
<b>Y</b>	0.13070	0.03063	0.39119	0.28935	0.1540	0.12966	-0.2788	0.32473	0.408
	9	0	5	2	16	5	84	4	9

In [71]:

Out[71]:

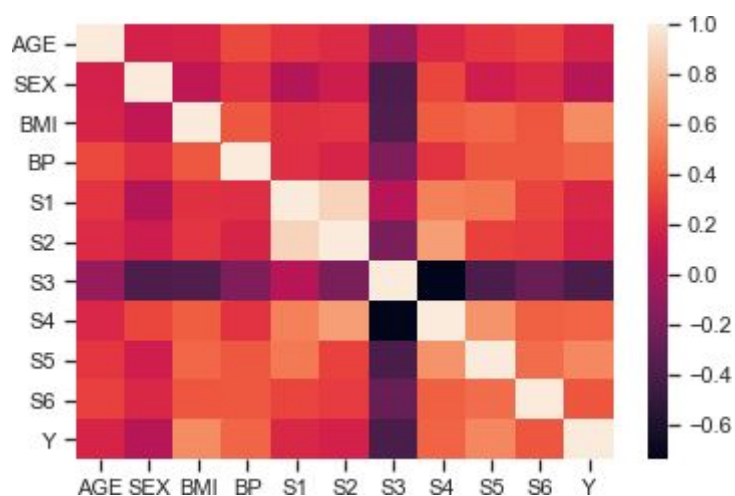
	AGE	SEX	BMI	BP	S1	S2	S3	S4	
<b>AG</b>	1.00000	0.17746	0.20055	0.35085	0.2625	0.22171	-0.1069	0.22101	0.265
<b>E</b>	0	3	4	9	24	1	73	7	1

<b>SE</b>	0.17746	1.00000	0.09807	0.26150	0.0277	0.13469	-0.3945	0.33752	0.174
<b>X</b>	3	0	9	8	90	5	84	4	6
<b>BM</b>	0.20055	0.09807	1.00000	0.39798	0.2878	0.29549	-0.3711	0.45906	0.491
<b>I</b>	4	9	0	5	29	4	72	8	6
<b>BP</b>	0.35085	0.26150	0.39798	1.00000	0.2752	0.20563	-0.1910	0.28079	0.396
	9	8	5	0	24	8	33	9	0
<b>S1</b>	0.26252	0.02779	0.28782	0.27522	1.0000	0.87879	0.01530	0.52067	0.512
	4	0	9	4	00	3	8	4	8
<b>S2</b>	0.22171	0.13469	0.29549	0.20563	0.8787	1.00000	-0.1974	0.65228	0.349
	1	5	4	8	93	0	35	3	9
<b>S3</b>	-0.1069	-0.3945	-0.3711	-0.1910	0.0153	-0.1974	1.00000	-0.7896	-0.45
	73	84	72	33	08	35	0	94	04
<b>S4</b>	0.22101	0.33752	0.45906	0.28079	0.5206	0.65228	-0.7896	1.00000	0.640
	7	4	8	9	74	3	94	0	3
<b>S5</b>	0.26517	0.17462	0.49160	0.39607	0.5128	0.34994	-0.4504	0.64039	1.000
	6	5	9	1	64	7	20	0	0
<b>S6</b>	0.29623	0.20327	0.38466	0.38121	0.3321	0.28648	-0.2908	0.41370	0.453
	5	7	4	9	73	3	63	0	0
<b>Y</b>	0.19782	0.03740	0.56138	0.41624	0.2324	0.19583	-0.4100	0.44893	0.589
	2	1	2	1	29	4	22	1	4



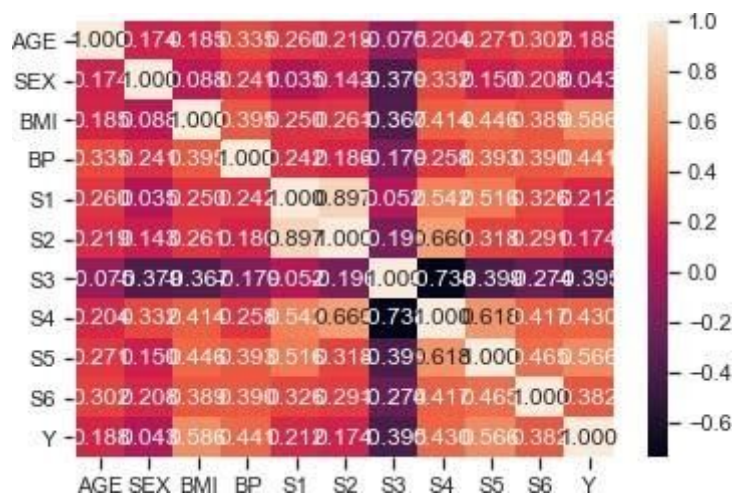
In [72]:

Out[72]: <matplotlib.axes.\_subplots.AxesSubplot at 0x31bc55d0>



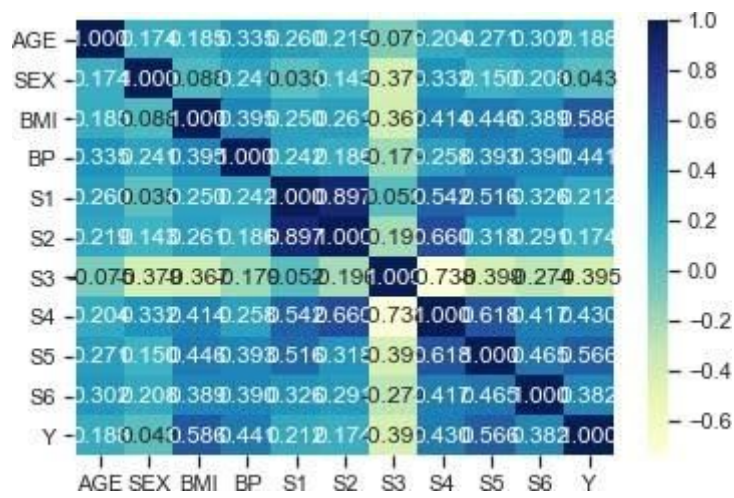
In [74]:

Out[74]: <matplotlib.axes.\_subplots.AxesSubplot at 0x2ccbdf70>



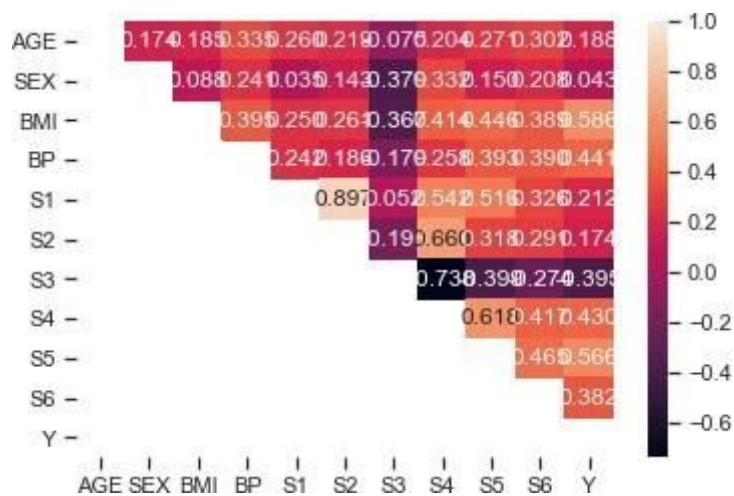
In [75]:

Out[75]: <matplotlib.axes.\_subplots.AxesSubplot at 0x31c7b790>



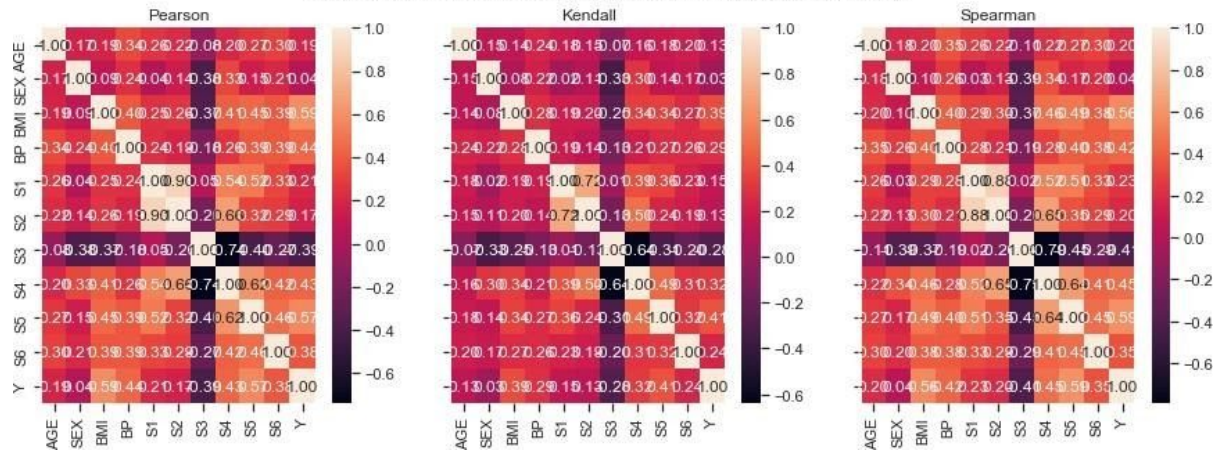
In [76]:

Out[76]: <matplotlib.axes.\_subplots.AxesSubplot at 0x31ae5f10>



In [77]:

Корреляционные матрицы, построенные различными методами



In [ ]: