



**Нижегородский государственный университет им. Н.И. Лобачевского**  
**Институт информационных технологий, математики и механики**

# ***Наглядный вероятностно-статистический анализ данных***

## **Практическое задание 4**

**«Числовые характеристики многомерных случайных величин.  
Основы корреляционного анализа»**

Пройдакова Екатерина Вадимовна,  
доцент кафедры ТВиАД ИИТММ

# Содержание

---

- ❑ Анализ многомерных данных: выделение зависимости
  - Пример: данные о результатах экзаменов
  - Числовые характеристики многомерных данных
  - Диаграммы рассеивания
  - Пример: данные о рейтингах покупателей ТЦ
- ❑ Практическое задание

# 1. АНАЛИЗ МНОГОМЕРНЫХ ДАННЫХ: ВЫЯВЛЕНИЕ ЗАВИСИМОСТИ

# 1. Анализ многомерных данных: выделение зависимости

## 1.1 Пример: данные о результатах экзаменов

Рассмотрим данные об итоговых экзаменах студентов США (Источник данных <https://www.kaggle.com/spscientist/students-performance-in-exams>) по столбцам:

**A** - пол (бинарные данные),

**B** - расовая принадлежность (одна из 5 этнических групп – номинальные данные),

**C** - уровень образования родителей (ординальные данные, можно упорядочить),

**D** - наличие льгот (бинарные данные),

**E** - прохождение подготовительного курса (бинарные данные),

**F, G, H** - оценки за экзамены по математике, чтению, письму (0-100, дискретные данные).

**Цель:**  
определить степень зависимости между характеристиками

	A	B	C	D	E	F	G	H
1	gender	race/ethnicity	parental level of education	lunch	test preparation course	math score	reading score	writing score
2	female	group B	bachelor's degree	standard	none	72	72	74
3	female	group C	some college	standard	completed	69	90	88
4	female	group B	master's degree	standard	none	90	95	93
5	male	group A	associate's degree	free/reduced	none	47	57	44
6	male	group C	some college	standard	none	76	78	75
7	female	group B	associate's degree	standard	none	71	83	78
8	female	group B	some college	standard	completed	88	95	92

# 1. Анализ многомерных данных: выделение зависимости

## 1.1 Пример: данные о результатах экзаменов

- Подготовим данные к анализу: импортируем для работы в Python

```
1 import xlrd
2
3 def importData(bookName, sheetInd):
4     res = []
5     #указывается полный путь до файла
6     book = xlrd.open_workbook(bookName)
7     #извлекаем лист с данными по индексу
8     sheet = book.sheet_by_index(sheetInd)
9     # загружаем данные построчно в список
10    for i in range(1, sheet.nrows): # с единицы, так как не загружаем заголовки
11        res.append(sheet.row_values(i, 0, sheet.ncols))
12    return res
13
14 # ===== Данные о тестировании учеников =====
15
16 examsData = importData('04_Экзамены.xlsx', 0)
```

# 1. Анализ многомерных данных: выделение зависимости

## 1.1 Пример: данные о результатах экзаменов

- ❑ Преобразуем данные в более удобный для нас формат, введем порядок для ординальных данных

```
18 types = {0:"Пол", 1:"Расовая группа", 2:"Образование родителей",
19          3:"Наличие льгот", 4:"Подготовительный курс",
20          5:"Оценка по математике", 6:"Оценка по чтению", 7:"Оценка по письму"}
21
22 genderTypes = {"male": "М", "female": "Ж"}
23 raceTypes = {"group A": 0, "group B": 1, "group C": 2, "group D": 3, "group E": 4}
24 edTypes = {"some high school": 0, "high school": 1, "some college": 2,
25            "associate's degree": 3, "bachelor's degree": 4, "master's degree": 5}
26 lunchTypes = {"standard": "нет", "free/reduced": "да"}
27 prepTypes = {"none": "нет", "completed": "да"}
28
29 for sample in examsData:
30     sample[0] = genderTypes[sample[0]]
31     sample[1] = raceTypes[sample[1]]
32     sample[2] = edTypes[sample[2]]
33     sample[3] = lunchTypes[sample[3]]
34     sample[4] = prepTypes[sample[4]]
35
36 # функция, возвращающая наблюдения для конкретного показателя
37 def dataByType(data, type):
38     res = []
39     for rec in data:
40         res.append(rec[type])
41     return res
42
43 # Индексы числовых данных
44 numData = [5, 6, 7]
```

# 1. Анализ многомерных данных: выделение зависимости

## 1.2 Числовые характеристики многомерных данных

- ❑ Необходимые библиотеки: `numpy`, `matplotlib`
- ❑ Рассмотрим 3 случайные величины  $\xi_1$ ,  $\xi_2$ ,  $\xi_3$  – оценки по математике, чтению, письму. Для них проводим разведывательный анализ.

```
45
46 import matplotlib.pyplot as plt
47 import numpy as np
48
49 # Вычисление числовых характеристик
50
51 scores = [] # массив выборочных значений для оценок
52 Mscores = [] # математическое ожидание оценок
53 for i in numData:
54     scores.append(dataByType(examsData, i))
55     Mscores.append(np.mean(dataByType(examsData, i)))
56
57 scoresCov = np.cov(scores) # ковариационная матрица
58 scoresCorr = np.corrcoef(scores) # матрица коэффициентов корреляции
59 print("Математическое ожидание оценок по соответствующим предметам = " + str(Mscores))
60 print("Ковариационная матрица = \n" + str(scoresCov))
61 print("Матрица коэффициентов корреляции = \n" + str(scoresCorr))
62
63 print("Смещенная оценка дисперсии оценок по математике = " +
64       str(np.var(dataByType(examsData, 5))))
65 print("Несмещенная оценка дисперсии оценок по математике = " +
66       str(np.var(dataByType(examsData, 5), ddof = 1)))
67
```

# 1. Анализ многомерных данных: выделение зависимости

## 1.2 Числовые характеристики многомерных данных

Математическое ожидание оценок по соответствующим предметам = [66.089, 69.169, 68.054]

Ковариационная матрица =

```
[[229.918998 180.99895796 184.93913313]
 [180.99895796 213.1656046 211.78666066]
 [184.93913313 211.78666066 230.90799199]]
```

Матрица коэффициентов корреляции =

```
[[1. 0.81757966 0.80264205]
 [0.81757966 1. 0.95459808]
 [0.80264205 0.95459808 1.]]
```

Смещенная оценка дисперсии оценок по математике = 229.68907899999996

Несмещенная оценка дисперсии оценок по математике = 229.91899799799796

- ❑ Математическое ожидание многомерной случайной величины является вектором из математических ожиданий одномерных величин.
- ❑ На главной диагонали ковариационной матрицы располагаются дисперсии одномерных величин.



# 1. Анализ многомерных данных: выделение зависимости

## 1.2 Числовые характеристики многомерных данных

- ❑ По умолчанию функция `numpy.var()` считает смещенную оценку дисперсии. Использование параметра `ddof = 1` (delta degree of freedom) приводит к вычислению несмещенной оценки дисперсии.
- ❑ Стандартизированный вариант характеристики рассеивания – коэффициент корреляции. Наибольшее значение коэффициента корреляции между  $\xi_2$  и  $\xi_3$  – **оценки по чтению и письму** (значительный вклад вносит линейная зависимость).
- ❑ Недиagonальные элементы ковариационной матрицы и матрицы коэффициентов корреляции положительны, то есть **между оценками по любым двум предметам имеется положительная связь**: высокие значения оценки по одному предмету сопровождаются высокими значениями по другому предмету.

# 1. Анализ многомерных данных: выделение зависимости

## 1.3 Диаграммы рассеивания (scatterplot)

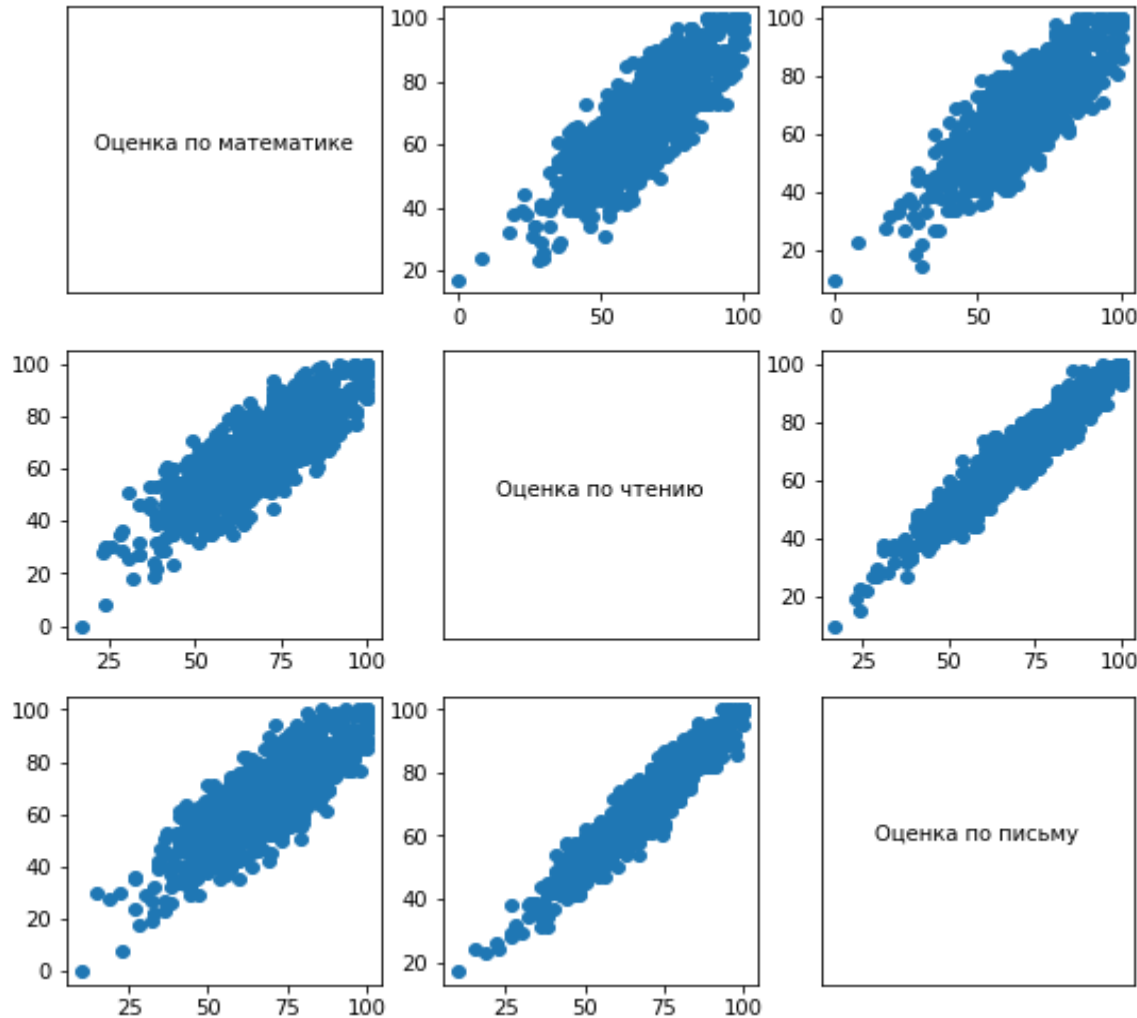
Визуализация многомерных данных – *диаграмма рассеивания* (**scatter**, **scatterplot**)

Каждая точка на диаграмме отвечает одному значению наблюдения за случайными векторами

```
68 # Диаграмма рассеивания для оценок по различным предметам
69 def scatterMatrix(data, ind, names):
70     fig, axes = plt.subplots(len(ind), len(ind), figsize = (3*len(ind), 3*len(ind)))
71     for i in range(len(ind)):
72         for j in range(len(ind)):
73             if (i == j):
74                 axes[i][j].annotate(types[ind[i]], (0.5, 0.5), ha = "center")
75                 axes[i][j].xaxis.set_visible(False)
76                 axes[i][j].yaxis.set_visible(False)
77             else:
78                 axes[i][j].scatter(dataByType(data, ind[i]), dataByType(data, ind[j]))
79     plt.show()
80
81 scatterMatrix(examsData, numData, types)
82
```

# 1. Анализ многомерных данных: выделение зависимости

## 1.3 Диаграммы рассеивания (scatterplot)



- Довольно сильная положительная связь между оценками по двум любым предметам
- Линейная зависимость между случайными величинами  $\xi_2$  и  $\xi_3$  (оценки по чтению и письму) наиболее явна

# 1. Анализ многомерных данных: выделение зависимости

## 1.3 Диаграммы рассеивания (scatterplot)

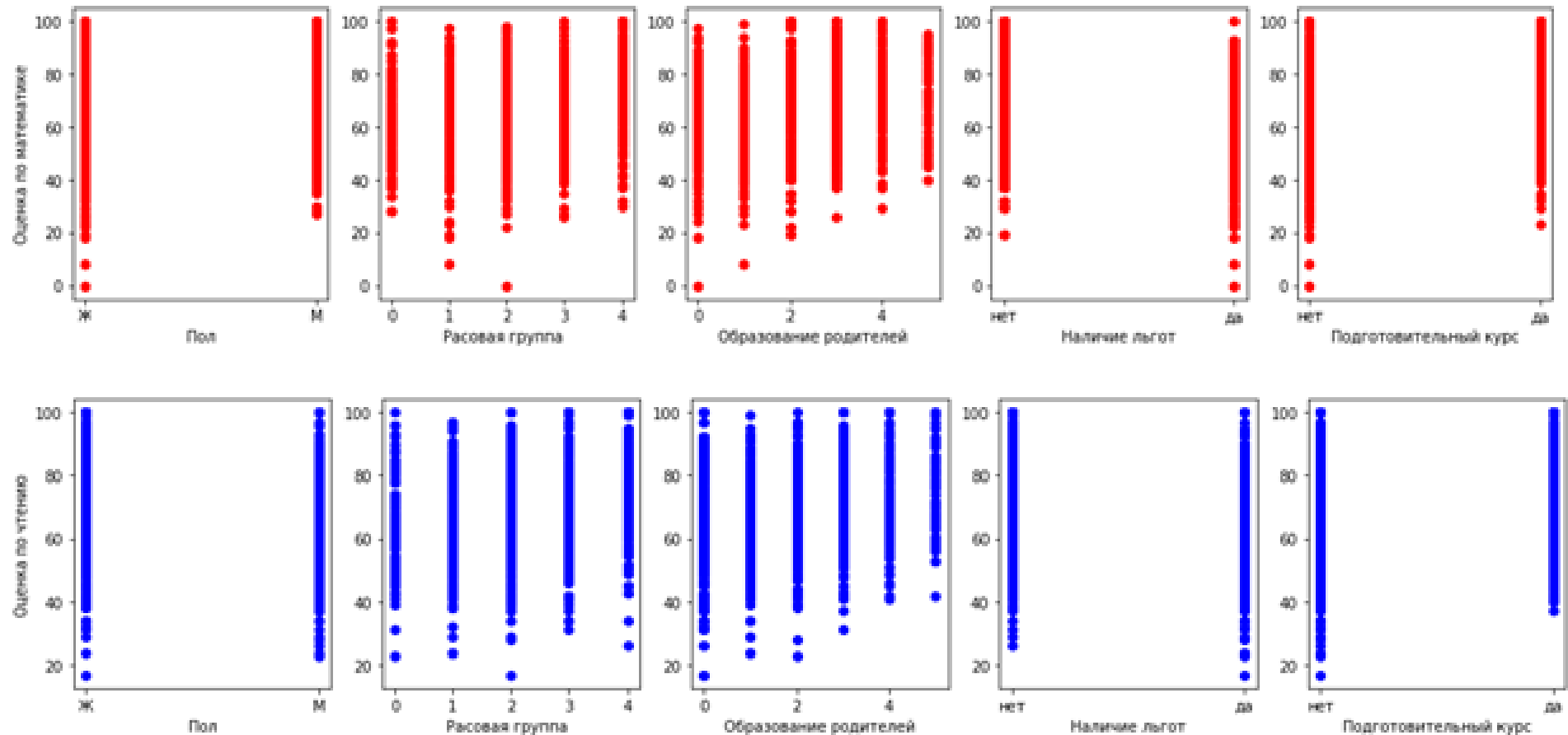
Диаграммы рассеивания можно строить не только для количественных, но и **для качественных данных**.

```
83 # Диаграммы рассеивания для качественных данных
84 colors = ['red', 'blue']
85 fig, axes = plt.subplots(2, 5, figsize = (18, 8))
86 for i in range(5):
87     for j in range(2):
88         axes[j][i].scatter(dataByType(examsData, i), dataByType(examsData, 5+j),
89                             c = colors[j])
90         axes[j][i].set_xlabel(types[i])
91         if (i == 0):
92             axes[j][i].set_ylabel(types[5+j])
93 plt.show()
94
```

# 1. Анализ многомерных данных: выделение зависимости

## 1.3 Диаграммы рассеивания (scatterplot)

Результаты построения:



# 1. Анализ многомерных данных: выделение зависимости

## 1.4 Пример: данные о рейтингах покупателей ТЦ

- ❑ Имеются данные о клиентах некоторого торгового центра (ТЦ): пол, возраст, годовой доход.
- ❑ На основе информации, собранной через членские карточки клиентов, каждому клиенту выставляется рейтинговая оценка (spending score).
- ❑ Задача: формирование представлений о целевой аудитории ТЦ. Для этого необходимо первоначально провести простой корреляционный анализ для определения влияния трех наблюдаемых факторов на рейтинговую оценку покупателя.

	A	B	C	D
	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
1				
2	Male	19	15	39
3	Male	21	15	81
4	Female	20	16	6
5	Female	23	16	77
6	Female	31	17	40
7	Female	22	17	76
8	Female	35	18	6

- ❑ Источник данных: <https://www.kaggle.com/vjchoudhary7/customer-segmentation-tutorial-in-python>

# 1. Анализ многомерных данных: выделение зависимости

## 1.4 Пример: данные о рейтингах покупателей ТЦ

Рассматриваем три случайных величины:  $\eta_1$ ,  $\eta_2$ ,  $\eta_3$  – возраст, годовой доход и рейтинг покупателя. Найдем для них числовые характеристики

```
95 # ===== Данные о рейтинге покупателей =====
96
97 mallData = importData('04_КлиентыТЦ.xlsx', 0)
98
99 types = {0:"Пол", 1:"Возраст", 2:"Годовой доход", 3:"Рейтинг"}
100
101 numData = [1, 2, 3]
102
103 mallNum = []
104 Mmalls = []
105 for i in numData:
106     mallNum.append(dataByType(mallData, i))
107     Mmalls.append(np.mean(dataByType(mallData, i)))
108
109 mallCov = np.cov(mallNum) # ковариационная матрица
110 mallCorr = np.corrcoef(mallNum) # матрица коэффициентов корреляции
111 print("Математическое ожидание = " + str(Mmalls))
112 print("Ковариационная матрица = \n" + str(mallCov))
113 print("Матрица коэффициентов корреляции = \n" + str(mallCorr))
114
115 scatterMatrix(mallData, [0, 1, 2, 3], types)
116
```

# 1. Анализ многомерных данных: выделение зависимости

## 1.4 Пример: данные о рейтингах покупателей ТЦ

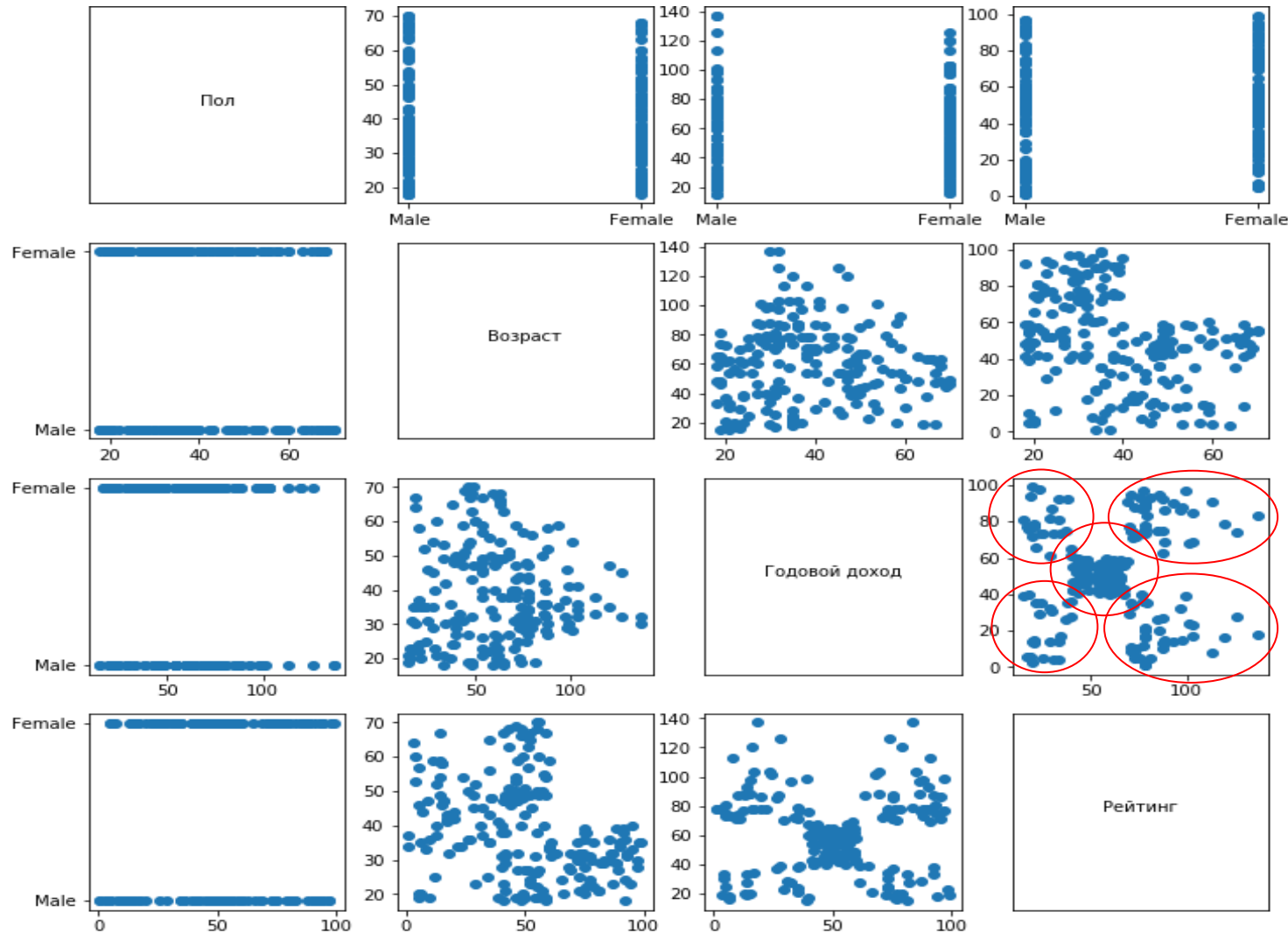
```
Математическое ожидание = [38.85, 60.56, 50.2]
Ковариационная матрица =
[[ 195.13316583   -4.54874372 -118.04020101]
 [  -4.54874372  689.83557789   6.71658291]
 [-118.04020101   6.71658291  666.85427136]]
Матрица коэффициентов корреляции =
[[ 1.           -0.01239804 -0.32722685]
 [-0.01239804   1.           0.00990285]
 [-0.32722685  0.00990285   1.           ]]
```

- ❑ Возраст ( $\eta_1$ ) имеет отрицательную умеренную корреляцию с рейтингом ( $\eta_3$ ), а годовой доход ( $\eta_2$ ) – положительную корреляцию с рейтингом ( $\eta_3$ ), хотя и очень слабую
- ❑ Далее строим диаграммы рассеивания.



# 1. Анализ многомерных данных: выделение зависимости

## 1.4 Пример: данные о рейтингах покупателей ТЦ



# 1. Анализ многомерных данных: выделение зависимости

## 1.4 Пример: данные о рейтингах покупателей ТЦ

---

### Первичные статистические выводы:

- ❑ Покупательский рейтинг несколько больше у женщин.
- ❑ С возрастом рейтинг покупателей падает.
- ❑ Торговый центр заинтересован, скорее, в покупателях, с большим годовым доходом.
- ❑ Вид диаграмм рассеивания рейтинга и годового дохода свидетельствует о существовании определенных целевых групп покупателей

## 2. ПРАКТИЧЕСКОЕ ЗАДАНИЕ

## 2. Практическое задание

**Данные** – из практического задания № 2: файл *02\_Автомобили.xls*

**1. Написать функции для подсчета** следующих выборочных числовых характеристик многомерных случайных величин: математическое ожидание, ковариация, ковариационная матрица, коэффициент корреляции, матрица коэффициентов корреляции.

Для подсчета каждой характеристики необходима отдельная функция. При написании **не использовать библиотечные функции** подсчета числовых характеристик.

**2. Проверить правильность работы** функций, сравнив их выходы с выходами функций библиотеки `numpy` на примере выборочных данных, полученных по результатам наблюдений за величинами  $\chi$  – степень серьезности аварии,  $\xi$  – видимость дороги,  $\zeta$  – скорость ветра,  $\varsigma$  – влажность,  $\gamma$  – температура в момент совершения аварии.

## 2. Практическое задание

3. Построить диаграммы рассеивания для всевозможных пар случайных величин  $\chi, \xi, \zeta, \varsigma, \gamma$ .

**Сделать выводы** о влиянии величин друг на друга.

4. Построить диаграммы рассеивания для определения зависимости величины  $\chi$  – **степень серьезности аварии от факторов, выраженных в качественных данных**: отметки о наличии вблизи места аварии лежачего полицейского, перекрестка, знака «Уступи дорогу», транспортной развязки, знака «нет выхода», железнодорожных путей, кругового движения, остановки общественного транспорта (автобусов, поездов и т.п.), знака «стоп», знаков или других мер успокоения движения, светофоров, поворотной петли.

**Сделать выводы** о влиянии указанных факторов на степень серьезности аварии