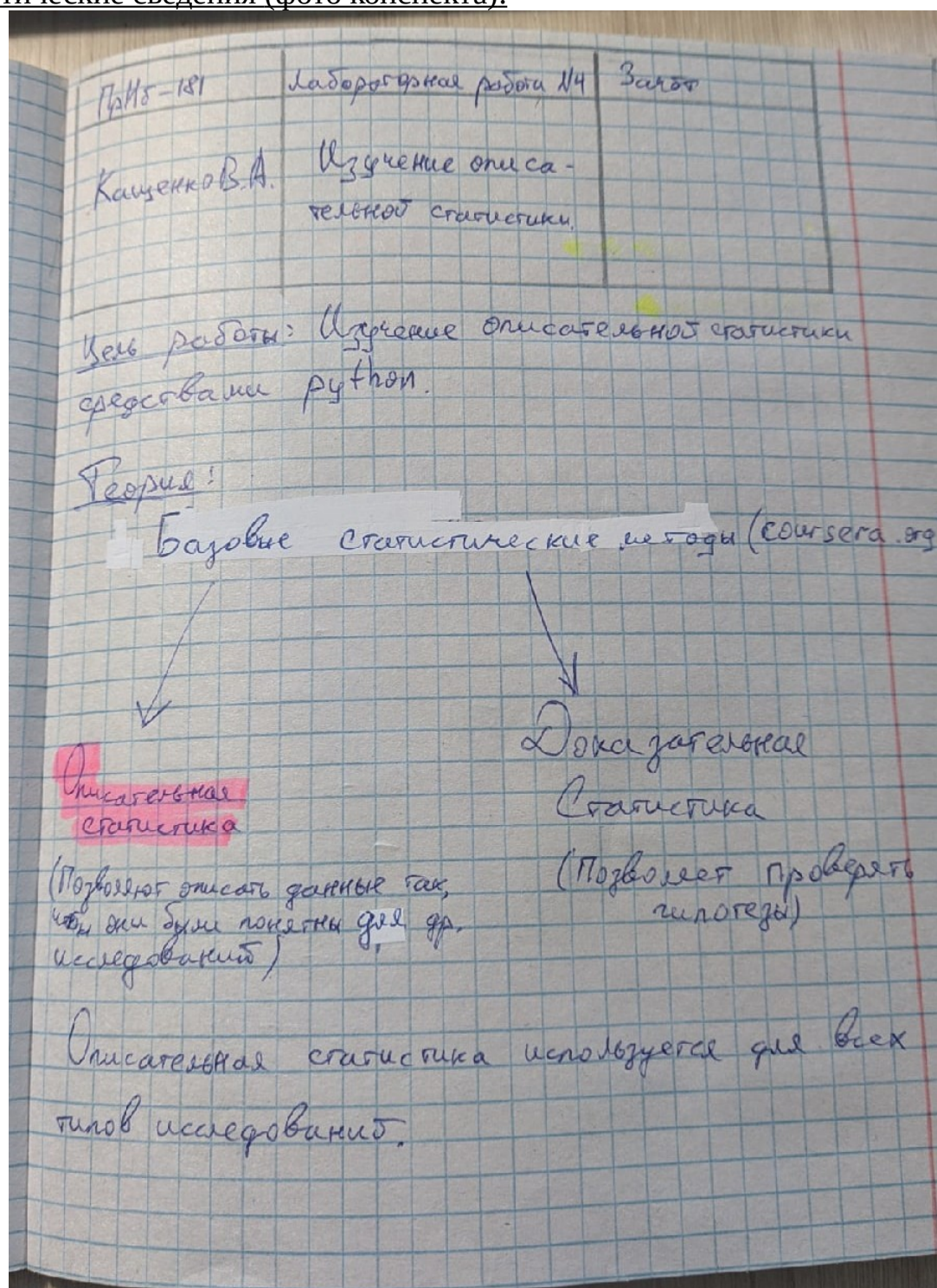


Приб-181	Лабораторная работа №4	Зачёт
Кащенко В. А.	Изучение описательной статистики	

Цель работы: Изучение описательной статистики средствами ЯП Python.

Теоретические сведения (фото конспекта):





# Методы описательной статистики

## Меры центральной тенденции

(для нахождения наиболее типичных значений переменной)

Меры:

- Мода (наиболее часто встречающееся значение)
- Медиана (средина упорядоченного ряда значений)
- Средн. арифметич. (сумма значений / кол-во)

(не самые надежные меры)

(если есть слишком большое значение, выброс,

то он может сместить среднее)

(тогда либо выбрасывают значение, либо пользуются медианой)

## Меры

## разбросности

(на сколько сильно разбросаны значения от типичного)

(чем больше мера разбросности, тем более разноразнобразно будет значение по ряду)

Меры:

- размах (макс. значение - мин. значение)

(выбросы не имеют)

(убирают 25% самых больших значений и 25% самых маленьких)

(это - межквартильный размах)

- стандартное отклонение

(диапазон от среднего нормальных размеров) (~68%) (S)

(если вычесть - нижн. граница нормы)

(если прибавить - верхн. граница)

Меры позволяют кратко описать собранные данные



## Описательная статистика (методика)

Разделы стат. науки, методы описания и представления осн. свойств данных. Позволяет обобщить,

этапы:

- сбор данных
- категоризация
- обобщение
- представление

Пусть  $X_1, X_2 \dots X_n$  - выборка независ. случайных величин. Упорядочим по возрастанию (вариационный ряд)

$$X(1) < X(2) < \dots < X(n),$$

где  $X(1)$  - min,  $X(n)$  - max

Элементы вариационного ряда - порядков. статистики

Величины  $d(i) = X(i+1) - X(i)$  называются разностями или расстояниями между порядковыми статистиками.

Размах выборки

$$R = X(n) - X(1) \quad (\text{max} - \text{min})$$



Выборочное среднее:

$$\bar{X} = (X_1 + X_2 + \dots + X_n) / n$$

### Среднее арифметическое:

Информативная мера "центрального положения" наблюдаемых переменных, если сообщается доверительный интервал (значения вокруг оценки, где с уровнем доверия находится "истинное" (неизвестное) среднее).

Увеличение размера выборки делает оценку среднего более надёжной. (Увеличение разброса надёж. значений уменьшает надёжность оценки)

Вычисление доверит. интервалов основывается на предположении "нормальности" наблюдаемых величин. (Без него оценка может быть ложной, всё для таких выборок.)

Нужно содержательно обобщать данные.

Диаграмма — отвлеч. точка. Можно стать информативнее, используя нужные характеристики данных.



Образ данных можно сформировать:

— зная из чего состоит величина

— зная рассеянность наблюдений.

формула:

$$\bar{X} = (X_1 + X_2 + \dots + X_n) / n$$

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$$

### Медиана:

При упорядочивании (min, max)

Делит ряд упоряд. значений пополам, с равным числом как выше, так и ниже.

При  $n$ -клетках:

$$\frac{(n+1)}{2}$$

$$\left( \text{Для } n=11 \Rightarrow \frac{11+1}{2} = 6 \right)$$

При  $n$ -четн:

средн. медианы чет. (средн. арифметич)

$$n=20$$

$$\frac{20}{2} = 10$$

$$\frac{20+1}{2} = 11$$

(ряд упорядочен)



Мода:

Мода  $M$  — наибольшее частое значение  
(если данные непрерывны — группируем и вычисляем  
модальную группу).

Неск. мод, моды не имеют моды ( $\forall x$  только одна)

Всегда  $> 1$ :

Когда  $\geq 2$  значение встречается одинаков. число  
раз и встречается больше др. значений

Редко используют как общ. характеристику

Среднее арифметическое:

При несиметрич. распредел. не обобщающ.  
показатель

Взвешенное:

$$\frac{w_1 X_1 + w_2 X_2 + \dots + w_n X_n}{w_1 + w_2 + \dots + w_n} = \frac{\sum w_i x_i}{\sum w_i}$$



Размах (интервал измерения):

max - min

При выбросах водит заблуждение.

Квантили:

Предположим, что  $X_0 \rightarrow X_n$

Величина  $X_{0,01}$  квантиль (до нее 1% наблюдений и выше 99%) (первый квантиль)

$X_{0,02}$  - второй квантиль и т.д.

Смысл: левее  $X_p$  лежит  $\sim 100\%$  наблюдений

$X_{0,1}, X_{0,2}, \dots, X_{0,9}$  - делит набор на 10 групп, децили (10-й, 20-й, ..., 90-й квантили)

Величины  $X_{0,25}, X_{0,5}, \dots, X_{0,75}$  (делит на 4 группы) - квантили (25-й, 50-й, 75-й)

50-й квантиль - медиана.



**Дисперсия** - среднее отклонение каждого наблюдения от средн. арифм.

$$D = \sigma^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n}$$

Если не генерован. совокупность, а выборка, то выборочная дисперсия

$$S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}$$

**Вариация** в пределах субъектов и между ними.  
Объект не всегда даёт точные и генетически данные

**Выброс.**

Отклоняющееся значение

**Ассиметрия**

Мера асимметрии распределения вероятностей вещественной случайной величины.



### Практика:

- 1) Текст + Конспект — ✓
- 2) ~~ЖЗ~~ дисциплины — ✓ анализ, — ✓
- 3) Построение гистограмм — ✓  
График

Вывод: Изучена описательная статистика  
расчетов рутор.



### Постановка задачи:

Выберите не менее 9 дисциплин, которые читаются всему потоку (трем группам на одном курса разом). Каждую дисциплину для данного курса рассматривайте отдельно. Проведите анализ выбранных вами данных с использованием всех средств описательной статистики. Выполните построение соответствующих столбчатых диаграмм, поверх диаграмм должен располагаться график распределения, наиболее близко характеризующий вашу выборку.

### Практическая часть (код программы):

# библиотеки

```
import seaborn as sns
import pandas as pd
import numpy as np
from scipy import stats
```

# функция создания подписей гистограмм описательной статистики

```
def set_up_statistic(marks, ax):
```

```
    # Среднее арифметическое
```

```
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.90,
            'Среднее арифметическое: {0}'.format(np.mean(marks)), fontsize=9)
```

```
    # Медиана
```

```
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.85,
            'Медиана: {0}'.format(np.median(marks)), fontsize=9)
```

```
    # Мода
```

```
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.80,
            'Мода: {0}'.format(stats.mode(marks)[0]), fontsize=9)
```

```
    # Среднее геометрическое
```

```
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.75,
            'Среднее геометрическое: {0}'.format(stats.hmean(marks)), fontsize=9)
```

```
    # Размах
```

```
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.70,
            'Размах: {0}'.format(np.ptp(marks)), fontsize=9)
```

```
    # Межквартильный размах
```

```
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.65,
            'Межквартильный размах: {0}'.format(stats.iqr(marks)), fontsize=9)
```

```
    # Межквантильный диапазон (Интердециальный размах)
```

```
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.60,
            'Интердециальный размах: {0}'.format(stats.iqr(marks, rng=(10, 90))),
            fontsize=9)
```

```
    # Дисперсия
```

```
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.55,
```



```

'Дисперсия: {0}'.format(stats.variation(marks)), fontsize=9)

# Среднеквадратичное отклонение
ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.50,
'Sреднеквадратичное отклонение: {0}'.format(np.std(marks)), fontsize=9)

# Коэффициент асимметрии
ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.45,
'Коэффициент асимметрии : {0}'.format(stats.skew(marks)), fontsize=9)

# Файл csv
data = pd.read_csv('marks_groups.csv')
# Ненужный столбец
data = data.drop('Группа', axis=1)
# Колонки
subjects = [column for column in data.columns][:]

# Создание диаграмм для всех дисциплин
for subject in subjects:
    subject_data = data[[subject]]

    # Построение распределения
    g = sns.displot(subject_data, x=subject, binwidth=4, height=4,
    facet_kws=dict(margin_titles=True), kde=True, color = 'green')
    subject_marks = np.array([mark for mark in subject_data[subject]])

    # Настройка удобного отображения
    set_up_statistic(subject_marks, g.ax)

# Вывод
g.savefig('./diagrams/{0}.svg'.format(subject))

```



## Этапы выполнения работы:

### 1. Подключаем необходимые библиотеки:

```
# библиотеки
import seaborn as sns
import pandas as pd
import numpy as np
from scipy import stats
```

### 2. Определим функцию для подписей графиков:

Это основная функция. Передаются значения интервалов для каждого отдельного вычисления, её подпись, размер шрифта и тип вычисления.

```
# функция создания подписей гистограмм описательной статистики
```

```
def set_up_statistic(marks, ax):
    # Среднее арифметическое
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.90,
            'Среднее арифметическое: {0}'.format(np.mean(marks)), fontsize=9)

    # Медиана
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.85,
            'Медиана: {0}'.format(np.median(marks)), fontsize=9)

    # Мода
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.80,
            'Мода: {0}'.format(stats.mode(marks)[0]), fontsize=9)

    # Среднее геометрическое
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.75,
            'Среднее геометрическое: {0}'.format(stats.hmean(marks)), fontsize=9)

    # Размах
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.70,
            'Размах: {0}'.format(np.ptp(marks)), fontsize=9)

    # Межквартильный размах
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.65,
            'Межквартильный размах: {0}'.format(stats.iqr(marks)), fontsize=9)

    # Межквантильный диапазон (Интердециальный размах)
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.60,
            'Интердециальный размах: {0}'.format(stats.iqr(marks, rng=(10, 90))),
            fontsize=9)

    # Дисперсия
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.55,
            'Дисперсия: {0}'.format(stats.variation(marks)), fontsize=9)

    # Среднеквадратичное отклонение
    ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.50,
            'Среднеквадратичное отклонение: {0}'.format(np.std(marks)), fontsize=9)
```



```
# Коэффициент асимметрии
ax.text(ax.viewLim.intervalx.max(), ax.viewLim.intervaly.max() * 0.45,
'Коэффициент асимметрии : {0}'.format(stats.skew(marks)), fontsize=9)
```

3. Загружаем и подготавливаем графики:

```
# Файл csv
data = pd.read_csv('marks_groups.csv')
# Ненужный столбец
data = data.drop('Группа', axis=1)
# Колонки
subjects = [column for column in data.columns][:]
```

4. Проходим по всем колонкам и строим графики статистики, сохраняем их на диске.

```
# Создание диаграмм для всех дисциплин
for subject in subjects:
    subject_data = data[[subject]]

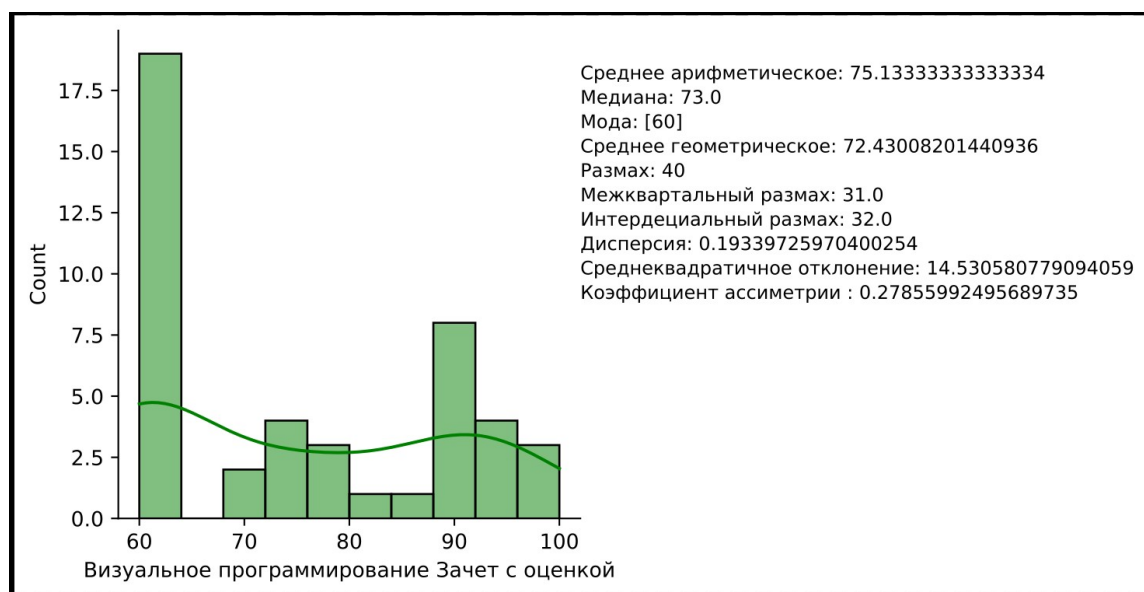
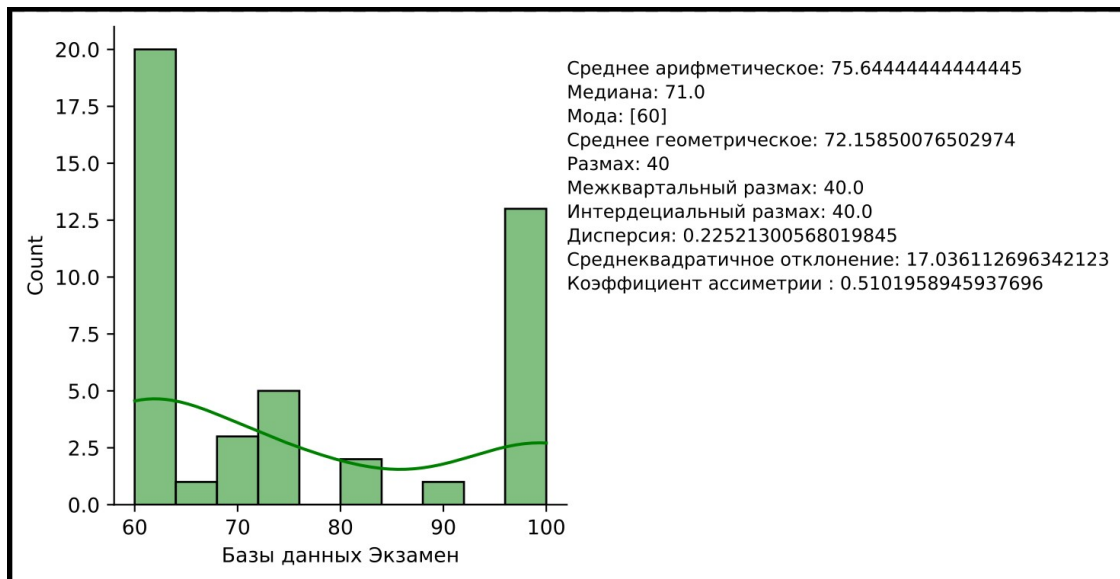
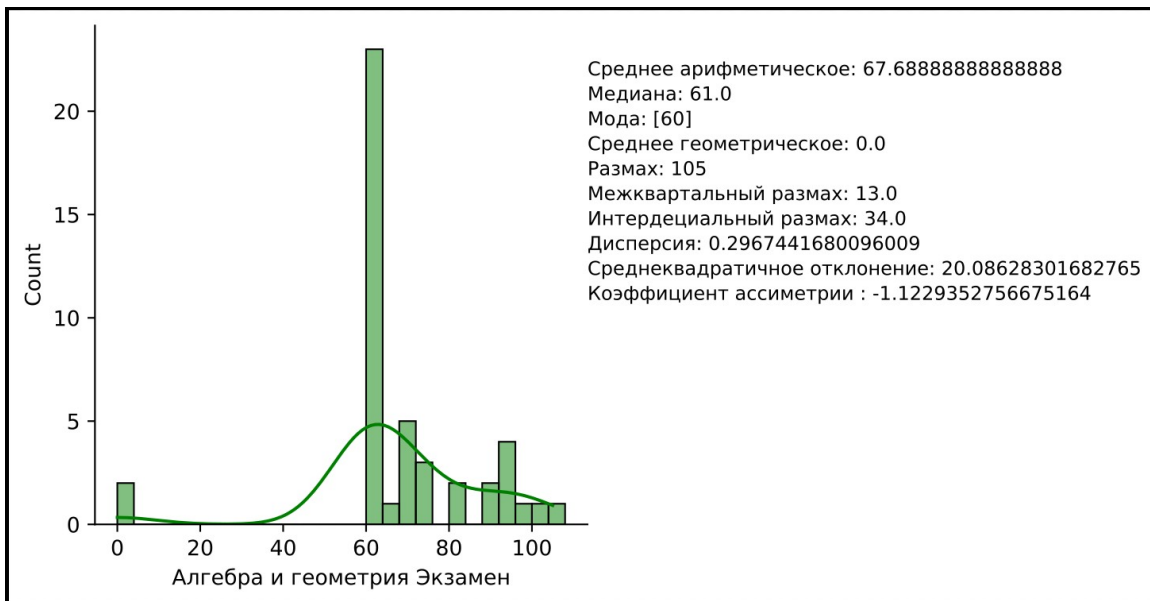
    # Построение распределения
    g = sns.displot(subject_data, x=subject, binwidth=4, height=4,
facet_kws=dict(margin_titles=True), kde=True, color = 'green')
    subject_marks = np.array([mark for mark in subject_data[subject]])

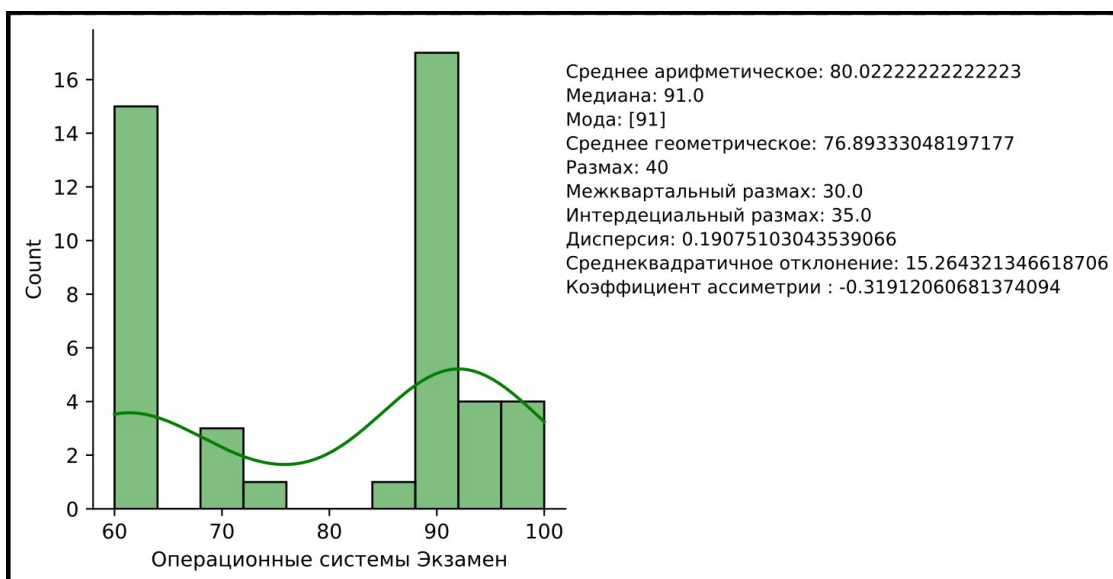
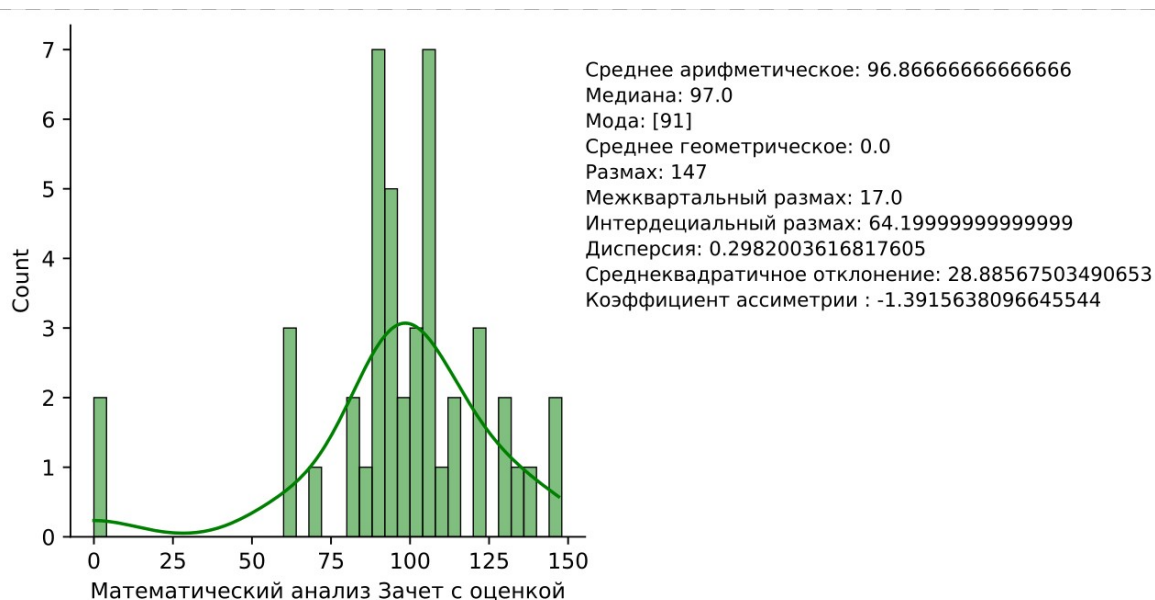
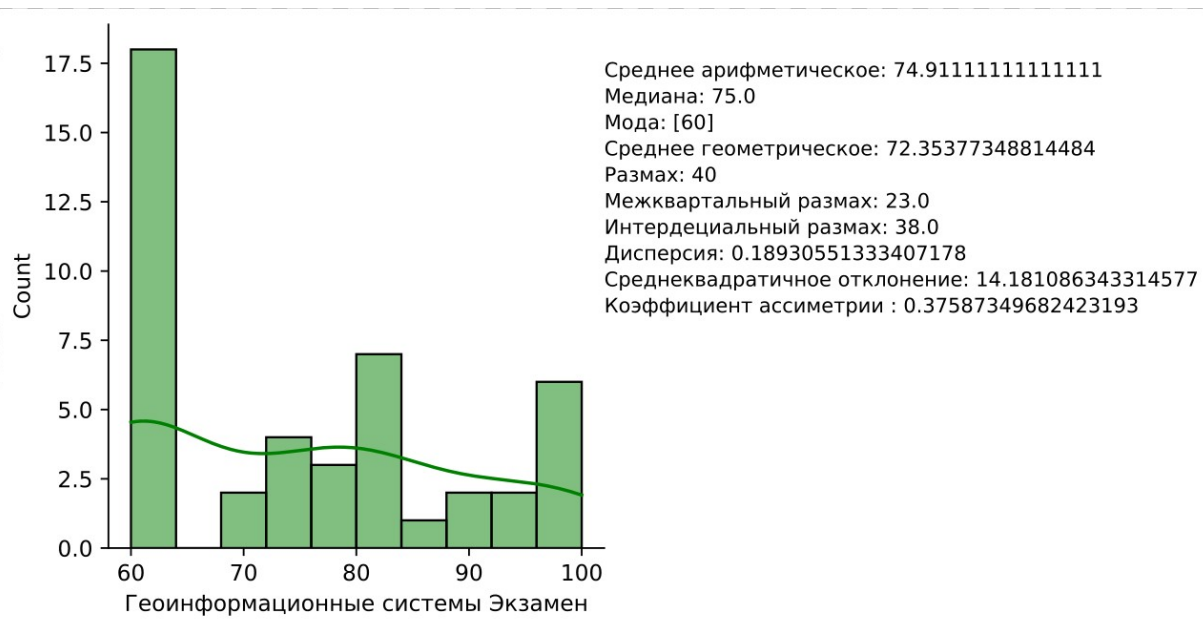
    # Настройка удобного отображения
    set_up_statistic(subject_marks, g.ax)

    # Вывод
    g.savefig('./diagrams/{0}.svg'.format(subject))
```

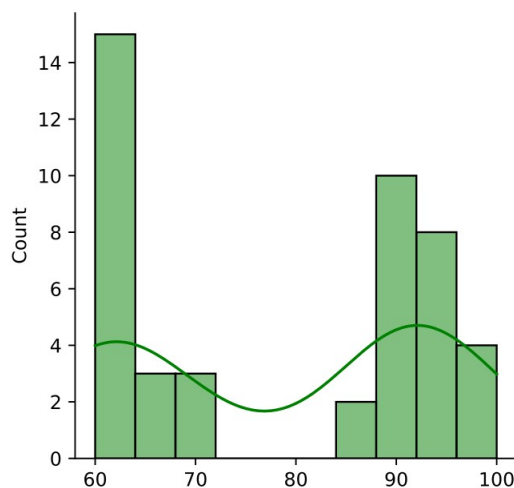


## Практическая часть (результат подведения статистики):



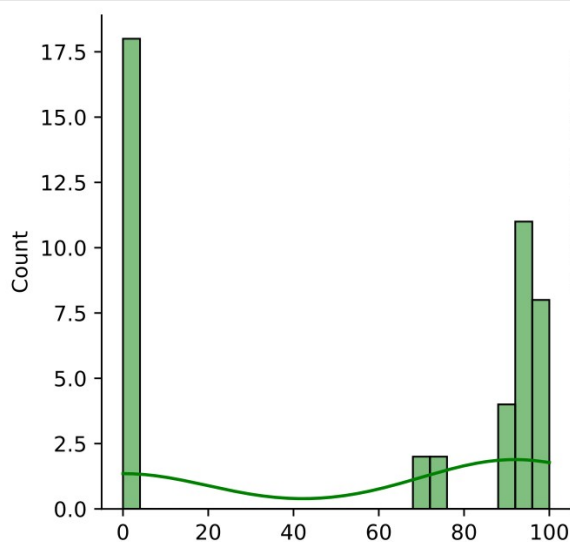






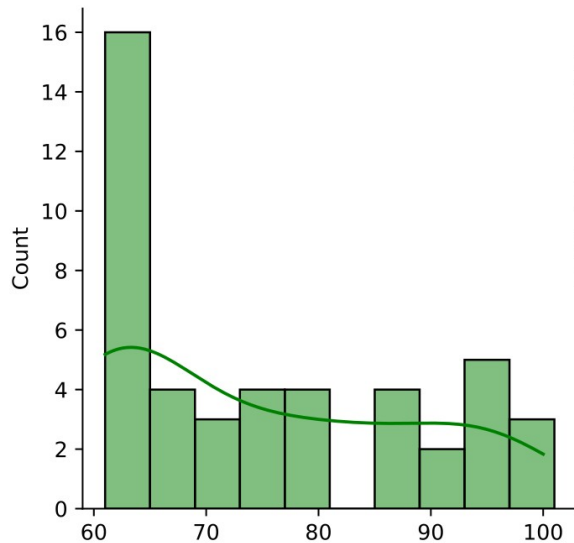
Среднее арифметическое: 78.62222222222222  
 Медиана: 85.0  
 Мода: [60]  
 Среднее геометрическое: 75.54729836263625  
 Размах: 40  
 Межквартильный размах: 30.0  
 Интерквартильный размах: 33.0  
 Дисперсия: 0.19493857693235464  
 Среднеквадратичное отклонение: 15.326504115259349  
 Коэффициент асимметрии : -0.1032811881076914

Производственная практика, научно-исследовательская работа Зачет с оценкой\_x



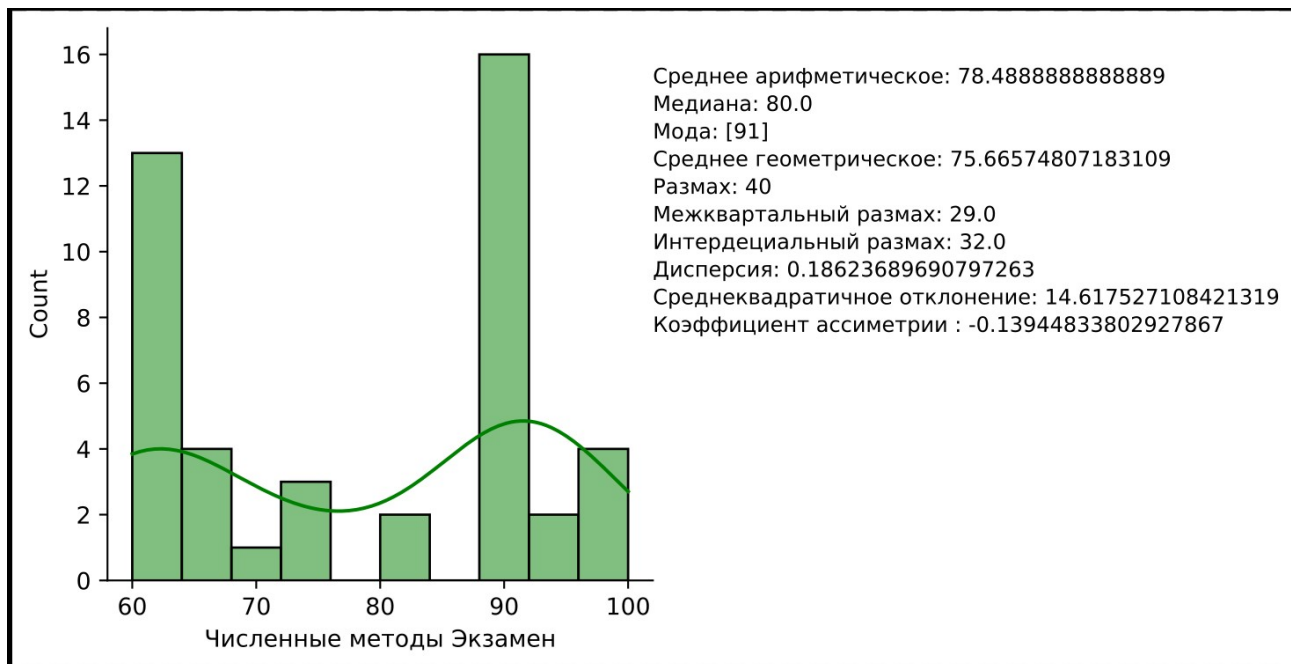
Среднее арифметическое: 54.577777777777776  
 Медиана: 91.0  
 Мода: [0]  
 Среднее геометрическое: 0.0  
 Размах: 100  
 Межквартильный размах: 93.0  
 Интерквартильный размах: 97.2  
 Дисперсия: 0.8252709936546764  
 Среднеквадратичное отклонение: 45.04145689813078  
 Коэффициент асимметрии : -0.34987421784185907

Производственная практика, научно-исследовательская работа Зачет с оценкой\_y



Среднее арифметическое: 74.66666666666667  
 Медиана: 71.0  
 Мода: [61]  
 Среднее геометрическое: 72.43030268554797  
 Размах: 39  
 Межквартильный размах: 24.0  
 Интерквартильный размах: 34.0  
 Дисперсия: 0.1790173004978877  
 Среднеквадратичное отклонение: 13.366625103842281  
 Коэффициент асимметрии : 0.4779546133027651

Теория вероятностей и математическая статистика Зачет с оценкой



### Практическая часть (подведение итогов):

Были проанализированы результаты сессий по 10 предметам. Данные визуализированы с помощью гистограмм и библиотек python для дальнейшей работы с ними. (описательная статистика)

Вывод: Изучена описательная статистика средствами Python.