

МИНИСТЕРСТВО ОБРАЗОВАНИЯ РОССИЙСКОЙ ФЕДЕРАЦИИ
Федеральное государственное автономное образовательное учреждение
высшего образования
«Национальный исследовательский университет «Высшая школа
экономики»
Факультет компьютерных наук
Образовательная программа «Прикладная математика и информатика»

ОТЧЕТ
О ВЫПОЛНЕНИИ ИТОГОВОЙ РАБОТЫ ПО КУРСУ
CORE CONCEPTS IN DATA ANALYSIS

Выполнили	<i>Студент 3 курса ПМИ ФКН ВШЭ,</i> <i>группы БПМИ174</i>	Никифоров Алексей Владимирович
	<i>Студент 3 курса ПМИ ФКН ВШЭ,</i> <i>группы БПМИ174</i>	Попов Илья Иванович
Руководитель	<i>Профессор департамента анализа</i> <i>данных и искусственного</i> <i>интеллекта, доктор технических</i> <i>наук</i>	Миркин Борис Григорьевич

Москва, 2019

СОДЕРЖАНИЕ

СОДЕРЖАНИЕ	2
ВЫБОР И ОПИСАНИЕ ДАННЫХ	3
КЛАСТЕР-АНАЛИЗ	4
БУТСТРЭП	7
ТАБЛИЦА СОПРЯЖЕННОСТИ	11
МГК/SVD	17
ЛИНЕЙНАЯ РЕГРЕССИЯ	21
СПИСОК ИСПОЛЬЗОВАННЫХ МАТЕРИАЛОВ	25
ПРИЛОЖЕНИЯ	26

ВЫБОР И ОПИСАНИЕ ДАННЫХ

Для выполнения работы мы использовали исторические данные о ноябрьской погоде в Москве в период с 1999 по 2014 год включительно. Данные были взяты с сайта *climate-energy.ru*[1], где автор публикует архивы погоды по Москве, городам Московской области и некоторым городам России. Ниже приводится пример данных из датасета, а также описание признаков.

date	avg_temp	min_temp	max_temp	downfall	pressure	humidity	wind_speed	snow_depth	year
1	1.3	-1.0	2.7	0.0	1020.7	80.0	1.7	0	2014
2	-1.7	-3.0	-0.4	0.0	1026.2	74.0	1.9	0	2014
3	-0.4	-1.0	0.5	5.6	1020.0	89.0	1.0	0	2014
5	3.2	2.0	4.4	0.0	1014.6	92.0	1.0	0	2014
6	5.1	4.0	5.8	0.0	1018.7	95.0	1.0	0	2014
7	5.8	4.0	7.3	4.1	1016.7	97.0	1.0	0	2014
8	7.4	6.0	8.6	0.0	1014.0	88.0	1.1	0	2014
9	4.2	3.0	5.8	5.1	1024.0	86.0	1.0	0	2014
11	4.6	3.0	5.3	0.0	1030.3	92.0	1.2	0	2014
12	5.4	5.0	5.7	0.0	1026.0	82.0	1.0	0	2014

Таблица 1 - Пример данных из датасета

Описание признаков в датасете:

- `date` – целое число – номер дня в месяце;
- `avg_temp` – действительное число – средняя температура в течении дня (в градусах Цельсия);
- `min_temp` – действительное число – минимальная температура в течении дня (в градусах Цельсия);
- `max_temp` – действительное число – максимальная температура в течении дня (в градусах Цельсия);
- `downfall` – действительное число – количество осадков в течении дня (в миллиметрах на один квадратный метр);
- `pressure` – действительное число – среднее атмосферное давление в течении дня (в миллиметрах ртутного столба);
- `humidity` – действительное число – средняя относительная влажность воздуха в течении дня (в процентах);
- `wind_speed` – действительное число – средняя скорость ветра в течении дня (в метрах в секунду);

- `snow_depth` – действительное число – глубина снежного покрова (в сантиметрах);
- `year` – целое число – год сбора данных;

Итого, в нашем датасете имеется 375 записей и 10 признаков. Ссылку на файл в формате CSV со всеми данными можно найти в Списке использованных материалов[2]. Работа выполнялась на языке Python, с использованием среды Google Colab[3] для осуществления совместной работы над проектом, а также с использованием инструментов из библиотек Pandas[4], NumPy[5], Matplotlib[6] и SKLearn[7].

КЛАСТЕР-АНАЛИЗ

Для выполнения задания по реализации кластерного анализа из датасета были выбраны следующие количественные признаки:

- `avg_temp` – средняя температура в течении дня (в градусах Цельсия);
- `pressure` – среднее атмосферное давление в течении дня (в миллиметрах ртутного столба);
- `wind_speed` – средняя скорость ветра в течении дня (в метрах в секунду);

Именно эти признаки были выбраны для проведения анализа потому, что они не коррелируют между собой напрямую, и, как нам кажется, позволяют наилучшим образом определить группы схожих дней.

Для выполнения непосредственно анализа был использован модуль KMeans[8] библиотеки SKLearn[7]. Рассмотрим код реализации анализа К-средних для 5 кластеров:

```
1. from sklearn.cluster import KMeans
2.
3. # Выбираем нужные нам признаки из общего датасета
4. df_features = df[['avg_temp', 'pressure', 'wind_speed']]
5.
6. kmeans_5 = KMeans(
7.     n_clusters = 5,          # 5 кластеров
8.     init = 'random',        # центры кластеров - случайно выбранные точки из датасета
9.     n_init = 10              # количество инициализаций алгоритма
10. )
11. kmeans_5.fit(df_features) # выполнение алгоритма
```

В настройках метода были указаны необходимые параметры работы алгоритма:

- `n_clusters=5` – количество кластеров
- `init='random'` – в качестве центров кластеров при инициализации алгоритма используются случайные объекты из датасета
- `n_init=10` – метод производит 10 случайных инициализаций алгоритма и выбирает то, которое удовлетворяет минимуму критерия метода;

Рассмотрим также визуализацию полученных кластеров, реализованную с помощью методов библиотеки Matplotlib[6] (полный исходный код слишком большой, он приведен ниже в приложении).

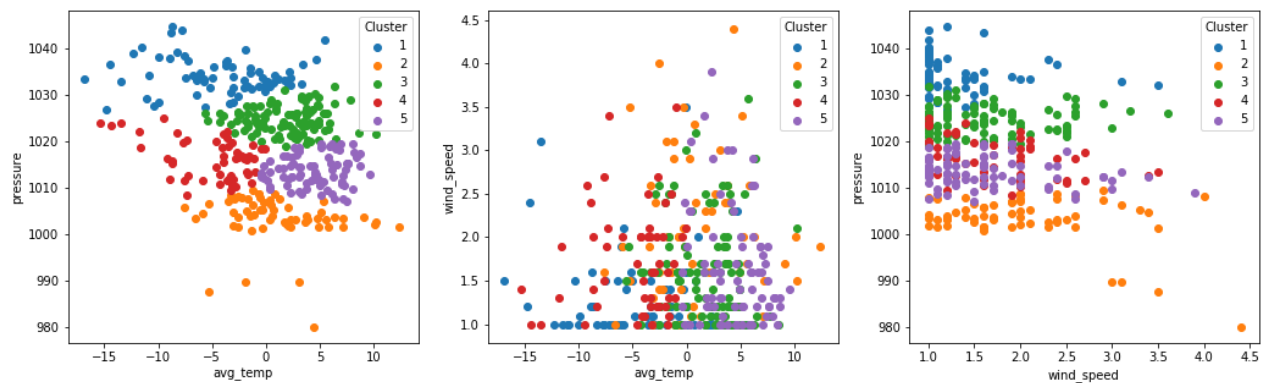


Иллюстрация 1 - Визуализация кластеров, полученных методов KMeans для k=5

Теперь выполним аналогичные действия для k=9:

```
1. from sklearn.cluster import KMeans
2.
3. # Выбираем нужные нам признаки из общего датасета
4. df_features = df[['avg_temp', 'pressure', 'wind_speed']]
5.
6. kmeans_9 = KMeans(
7.     n_clusters = 9,      # 9 кластеров
8.     init = 'random',     # центры кластеров - случайно выбранные точки из датасета
9.     n_init = 10          # количество инициализаций алгоритма
10.)
11. kmeans_9.fit(df_features) # выполнение алгоритма
```

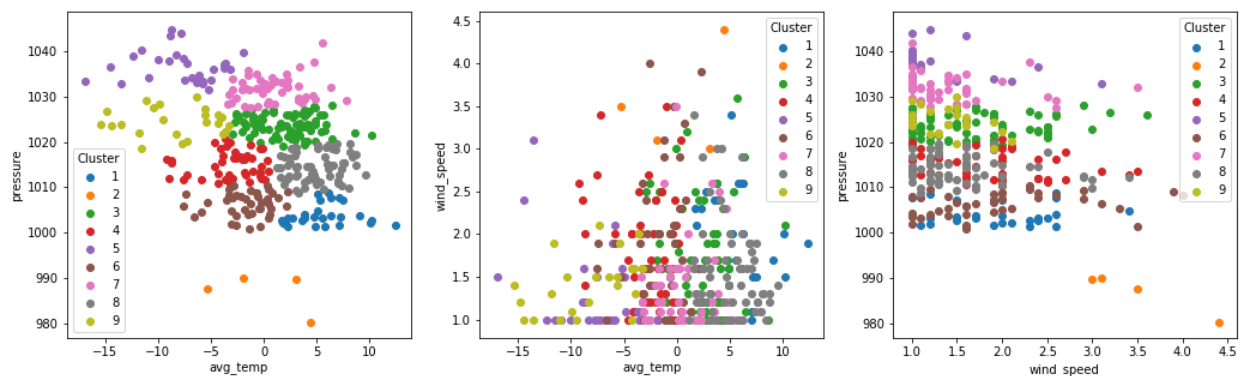


Иллюстрация 2 - Визуализация кластеров, полученных методов KMeans для k=9

Проинтерпретируем оба разбиения с помощью признаков таблицы данных путем сравнения внутрикластерных средних с общими средними. Для этого сделаем следующие визуализации:

	avg_temperature	avg_pressure	avg_wind
cluster			
All data	0.186376	1019.381471	1.645504
Cluster #0	5.559259	1003.914815	1.929630
Cluster #1	0.075000	986.850000	3.500000
Cluster #2	2.156627	1023.771084	1.603614
Cluster #3	-3.089130	1015.043478	1.819565
Cluster #4	-7.672414	1036.734483	1.310345

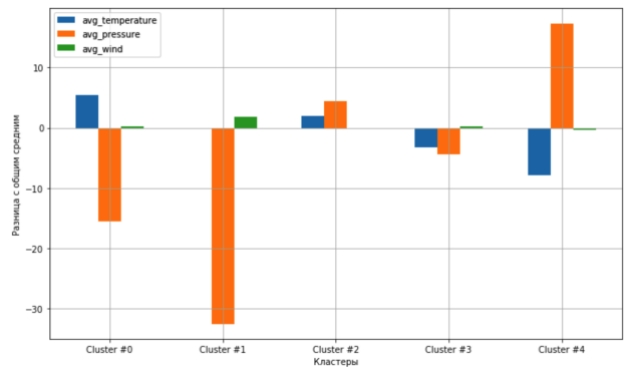


Иллюстрация 3 - Таблица средних и bar chart разницы внутрикластерного и общего средних для разбиения на 5 кластеров

Хорошо видно, что в Кластер №0 попали дни с температурой выше средней и давлением ниже среднего, в Кластер №1 – дни с аномально низким атмосферным давлением, а в Кластер №4 – дни с температурой ниже среднего и атмосферным давлением сильно выше среднего (что в народе называется «мороз и солнце»). Для остальных кластеров внутрикластерные средние значения признаков не сильно отличаются от средних по всему датасету.

	avg_temperature	avg_pressure	avg_wind
cluster			
All data	0.186376	1019.381471	1.645504
Cluster #0	5.559259	1003.914815	1.929630
Cluster #1	0.075000	986.850000	3.500000
Cluster #2	2.156627	1023.771084	1.603614
Cluster #3	-3.089130	1015.043478	1.819565
Cluster #4	-7.672414	1036.734483	1.310345
Cluster #5	-1.665854	1006.353659	1.965854
Cluster #6	0.952000	1031.764000	1.380000
Cluster #7	4.943750	1014.351562	1.575000
Cluster #8	-8.352174	1024.452174	1.417391

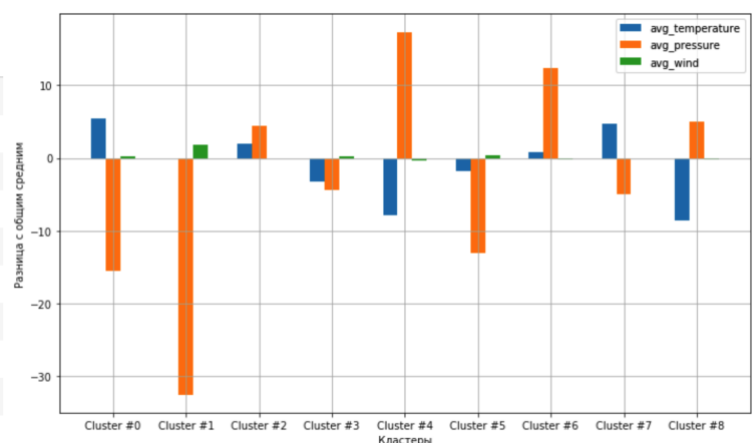


Иллюстрация 4 - Таблица средних и bar chart разницы внутрикластерного и общего средних для разбиения на 9 кластеров

Похожую картину мы наблюдаем и для разбиения на 9 кластеров:

- в Кластер №0 попали дни с температурой выше средней и давлением ниже среднего
- в Кластер №1 попали дни с давлением сильно ниже среднего и скоростью ветра сильно выше среднего
- в Кластер №4 попали дни с температурой ниже средней и аномально высоким давлением («мороз и солнце»)
- в Кластер №5 попали дни с температурой и давлением ниже среднего
- в Кластер №6 попали дни с температурой и давлением выше среднего
- для остальных кластеров можно сказать, что их внутрикластерные средние близки к средним на все датасете

Глядя на оба разбиения, нельзя точно сказать, какое из них лучше с точки зрения интерпретации. Оба разбиения имеют право на жизнь, и они хорошо подойдут для решения разных задач.

БУТСТРЭП

Найдем 95% доверительный интервал для среднего значения признака `avg_temp` на всем множестве объектов, используя бутстрэп:

```
1. from random import choices
2. from statistics import mean
3.
4. # Список средних
5. means = []
6.
7. # Количество экспериментов
8. bootstrap_iters = 5000
9. # Количество выбираемых значений для эксперимента
10. bootstrap_choices = len(df_features)
11.
12. for i in range(bootstrap_iters):
13.     # Добавляем в список
14.     means.append(
15.         # Среднее значение
16.         mean(
17.             choices(
18.                 # Выбираем случайный элемент из значений признака
19.                 df_features_5['avg_temp'].tolist(),
20.                 # k-раз
21.                 k=bootstrap_choices
22.             )
23.         )
24.     )
25.
26. # Визуализация
27. fig, graph = plt.subplots(figsize=(10, 6))
28. plt.hist(means, bins=50, color='#00b5ff', edgecolor='black')
```

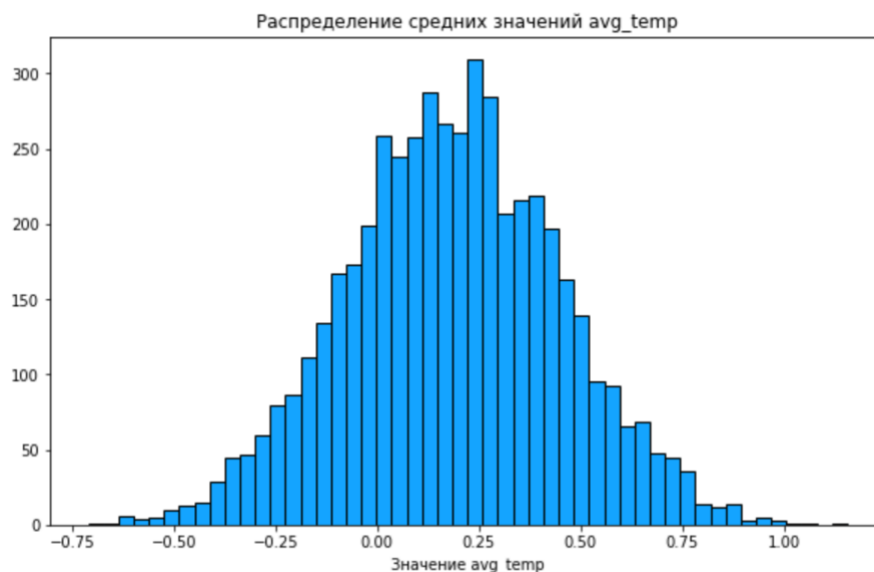


Иллюстрация 5 - Распределение средних значений `avg_temp`

Найдем 95% доверительный интервал способом "без опоры":

```
1. # Сортируем список со средними
2. means_ordered = sorted(means)
3. # Отсекаем 2.5% значений снизу и берем наименьший оставшийся
4. left = means_ordered[int(len(means_ordered)*0.05/2 - 1)]
5. # Отсекаем 2.5% значений сверху и берем наибольший оставшийся
6. right = means_ordered[int(len(means_ordered) - (len(means_ordered)*0.05/2) - 1)]
7.
8. print(f'Границы 95% доверительного интервала:
9. \tЛевая:\t{left}
10. \tПравая:\t{right}
11. ')
12.
13.
14. >>>Границы 95% доверительного интервала:
15. >>> Левая: -0.34005449591280656
16. >>> Правая: 0.6956403269754768
```

Найдем 95% доверительный интервал способом "с опорой":

```
1. from statistics import mean, stdev
2.
3. print(f'Границы 95% доверительного интервала:
4. \tЛевая:\t{mean(means) - 1.96*stdev(means)}
5. \tПравая:\t{mean(means) + 1.96*stdev(means)}
6. ')
7.
8.
9. >>>Границы 95% доверительного интервала:
10. >>> Левая: -0.33938573245791703
11. >>> Правая: 0.7020028441745383
```

Сравним средние по признаку avg_temp в кластерах №3 и №4 (из разбиения на 5 кластеров), используя бутстрэп. Для этого немного модифицируем приведенный выше код:

```
1. means = []
2. bootstrap_iters = 5000
3.
4. for i in range(bootstrap_iters):
5.     means.append(
6.         # Из среднего значения признака в эксперименте на 3 кластере
7.         mean(
8.             choices(
9.                 df_features_5[df_features_5.cluster == 3]['avg_temp'].tolist(),
10.                 k=len(df_features_5[df_features_5.cluster == 3])
11.             )
12.         )
13.         # вычитаем
14.         -
15.         # среднее значение признака в эксперименте на 4 кластере
16.         mean(
17.             choices(
18.                 df_features_5[df_features_5.cluster == 4]['avg_temp'].tolist(),
19.                 k=len(df_features_5[df_features_5.cluster == 4])
20.             )
21.         )
22.     )
23.
24. # Визуализация
25. fig, graph = plt.subplots(figsize=(10, 6))
26. plt.hist(means, bins=50, color='#00b5ff', edgecolor='black')
```

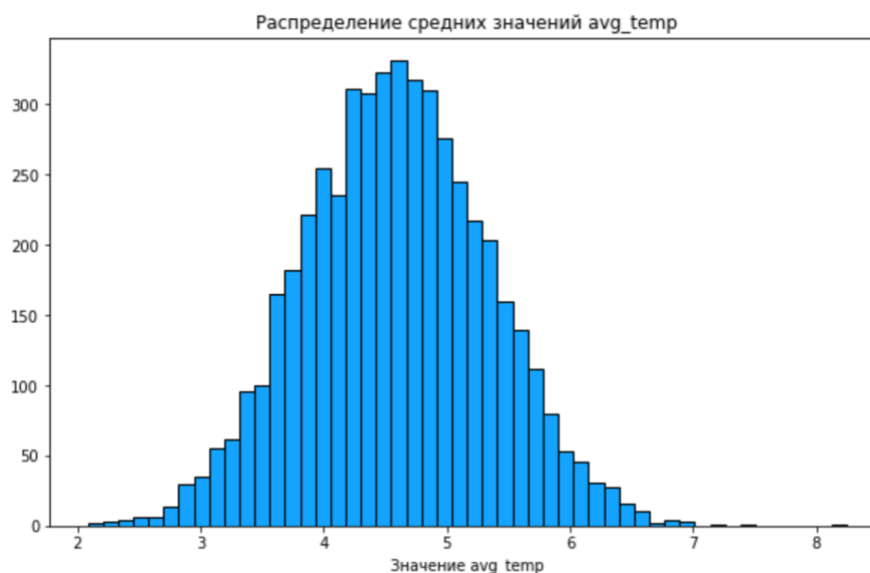



Иллюстрация 6 - Распределение разницы средних внутрикластерных значений признака *avg_temp* в кластерах №3 и №4

Ноль не попадает в это распределение, из чего можно сделать вывод, что средние значения признака *avg_temp* в кластерах №3 и №4 сильно различаются.

Рассмотрим распределения по-отдельности и рассчитаем 95% доверительный интервал для полученного выше распределения.

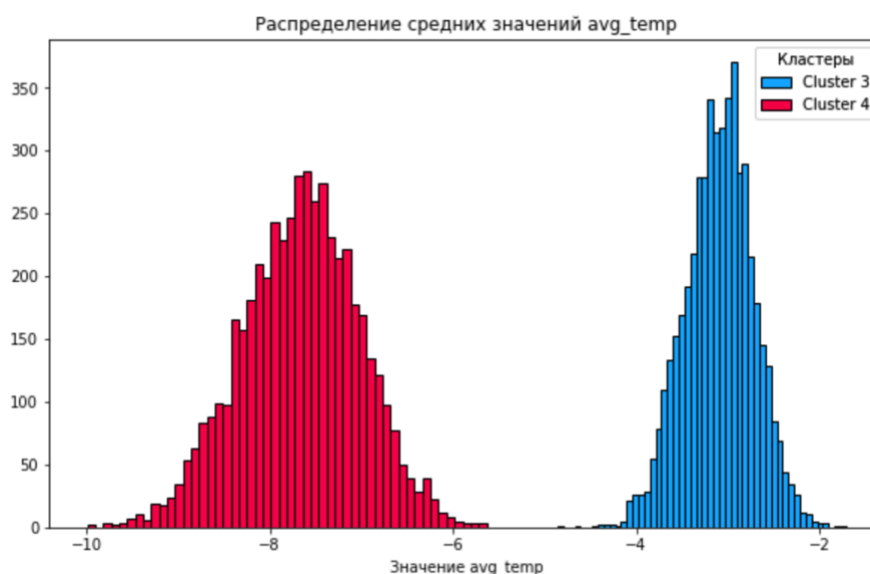


Иллюстрация 7 - Распределение средних внутрикластерных значений признака *avg_temp* в кластерах №3 и №4

Методика расчета доверительных интервалов абсолютно аналогична приведенной выше, поэтому ниже будут приведены только конечные результаты:

```
1. >>>Границы 95% доверительного интервала (без опоры):
2. >>>   Левая: 3.112893553223388
3. >>>   Правая: 6.068665667166417
4. >>>
5. >>>Границы 95% доверительного интервала (с опорой):
6. >>>   Левая: 3.1109550134891464
7. >>>   Правая: 6.043845406300958
```

Для кластера №1 (из разбиения на 5 кластеров) сравним среднее на всем множестве для признака `avg_temp` с его средним внутри кластера, используя бутстрэп:

```
1. means = []
2.
3. bootstrap_iters = 5000
4.
5. for i in range(bootstrap_iters):
6.     means.append(
7.         mean(
8.             choices(
9.                 df_features_5[df_features_5.cluster == 1]['avg_temp'].tolist(),
10.                 k=len(df_features_5[df_features_5.cluster == 1])
11.             )
12.         )
13.     -
14.     mean(
15.         choices(
16.             df_features_5['avg_temp'].tolist(),
17.             k=len(df_features_5)
18.         )
19.     )
20. )
21.
22. fig, graph = plt.subplots(figsize=(10, 6))
23. plt.hist(means, bins=50, color='#00b5ff', edgecolor='black')
```

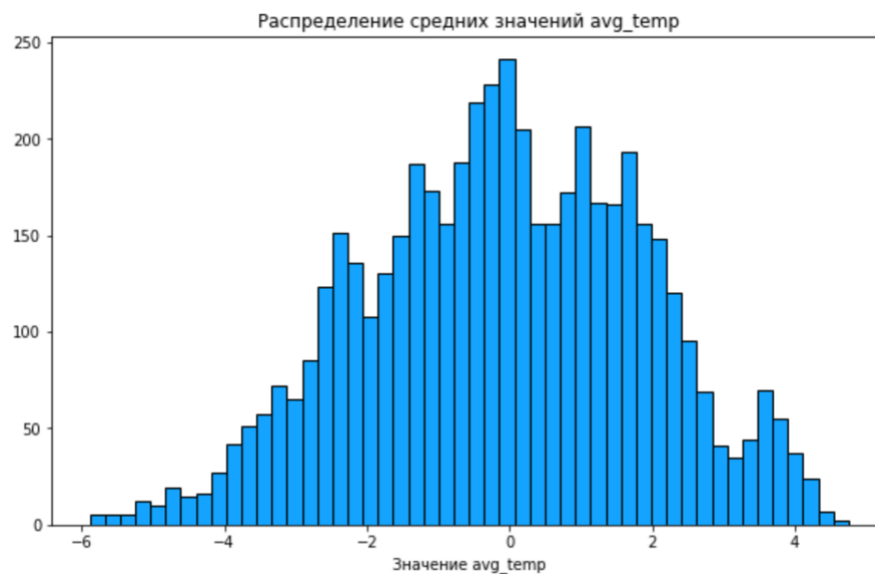


Иллюстрация 8 - Распределение разницы среднего внутрикластерного значения признака `avg_temp` в кластере №1 с средним значением на всей выборке

Ноль попадает в это распределение, из чего можно сделать вывод, что средние значения признака `avg_temp` в кластере №1 и на всей выборке не сильно различаются.

Рассмотрим распределения по-отдельности и рассчитаем 95% доверительный интервал для полученного выше распределения.

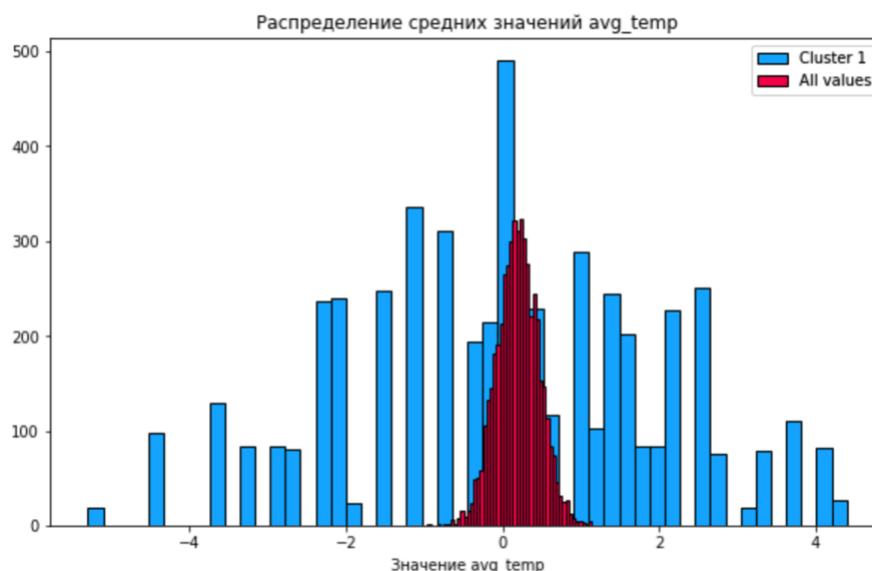


Иллюстрация 9 - Распределение среднего внутрикластерного значения признака *avg_temp* в кластере №1 и на всей выборке

Методика расчета доверительных интервалов абсолютно аналогична приведенной выше, поэтому ниже будут приведены только конечные результаты:

```
1. >>>Границы 95% доверительного интервала:
2. >>>   Левая: -3.8915531335149858
3. >>>   Правая: 3.684264305177112
4. >>>
5. >>>Границы 95% доверительного интервала:
6. >>>   Левая: -3.9516607523396794
7. >>>   Правая: 3.747105384492268
```

ТАБЛИЦА СОПРЯЖЕННОСТИ

Сформируем три номинальных признака *x1*, *x2* и *x3*. В качестве первого номинального признака возьмем разбиение диапазона значений признака *avg_temp* на 4 примерно равных по количеству объектов интервала.

```
1. # Разбиение
2. temp_sorted = df['avg_temp'].sort_values()
3. average_temp_borders = [
4.     temp_sorted.iloc[0],
5.     temp_sorted.iloc[len(temp_sorted.index) // 4],
6.     temp_sorted.iloc[(len(temp_sorted.index) // 4) * 2],
7.     temp_sorted.iloc[(len(temp_sorted.index) // 4) * 3],
8.     temp_sorted.iloc[len(temp_sorted.index)-1] + 0.1
9. ]
10. # Визуализация
11. fig, graph = plt.subplots(figsize=(10, 6))
12. plt.hist(df['avg_temp'], bins=15, color='#00b5ff', label='avg_temp', edgecolor='black')
13. plt.axvline(x=average_temp_borders[1], color='#f91155')
14. plt.axvline(x=average_temp_borders[2], color='#f91155')
15. plt.axvline(x=average_temp_borders[3], color='#f91155')
16. # Вывод
17. print(f'')
18. Границы интервала: [{average_temp_borders[0]}; {average_temp_borders[1]}}; \tЗначений в
    интервале: {len(df[(df.avg_temp >= average_temp_borders[0]) & (df.avg_temp < average_temp_borders[1])])}
```

```

19.     Границы интервала: [{average_temp_borders[1]}; {average_temp_borders[2]}};\tЗначений в
        интервале: {len(df[(df.avg_temp >= average_temp_borders[1]) & (df.avg_temp < average_temp
        _borders[2])])}
20.     Границы интервала: [{average_temp_borders[2]}; {average_temp_borders[3]}};\tЗначений в
        интервале: {len(df[(df.avg_temp >= average_temp_borders[2]) & (df.avg_temp < average_temp
        _borders[3])])}
21.     Границы интервала: [{average_temp_borders[3]}; {average_temp_borders[4]}};\tЗначений в
        интервале: {len(df[(df.avg_temp >= average_temp_borders[3]) & (df.avg_temp < average_temp
        _borders[4])])}
22.     '''

```

```

Границы интервала: [-16.9; -2.9);    Значений в интервале: 93
Границы интервала: [-2.9; 0.3);      Значений в интервале: 89
Границы интервала: [0.3; 3.8);       Значений в интервале: 94
Границы интервала: [3.8; 12.5);      Значений в интервале: 99

```

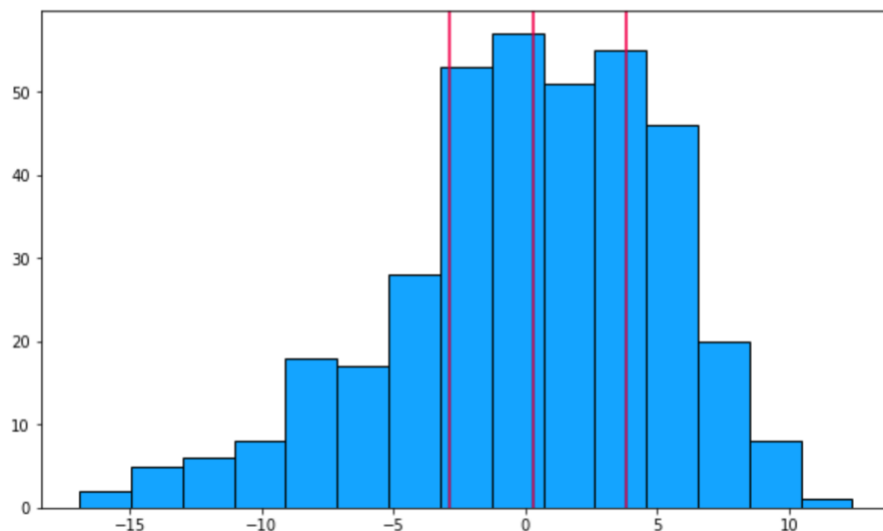


Иллюстрация 10 - Разбиение объектов по признаку `avg_temp` на 4 приблизительно равных по количеству объектов интервала

Для выделения второго номинального признака рассмотрим распределение признака `snow_depth`:

```

1. fig, graph = plt.subplots(figsize=(10, 6))
2. plt.hist(df['snow_depth'], bins=30, color='#00b5ff', label='snow_depth', edgecolor='black'
3. )
4. snow_depth_borders = [0, 0.1, 5, 13, 27]
5. print(f'''Рассмотрим распределение, и выберем для него границы интервалов
6.     Границы интервала: [{snow_depth_borders[0]}; {snow_depth_borders[1]}};\tЗначений в инт
        ервале: {len(df[(df.snow_depth >= snow_depth_borders[0]) & (df.snow_depth < snow_depth_bor
        ders[1])])}
7.     Границы интервала: [{snow_depth_borders[1]}; {snow_depth_borders[2]}};\tЗначений в инт
        ервале: {len(df[(df.snow_depth >= snow_depth_borders[1]) & (df.snow_depth < snow_depth_bor
        ders[2])])}
8.     Границы интервала: [{snow_depth_borders[2]}; {snow_depth_borders[3]}};\tЗначений в и
        нтервале: {len(df[(df.snow_depth >= snow_depth_borders[2]) & (df.snow_depth < snow_depth_b
        orders[3])])}
9.     Границы интервала: [{snow_depth_borders[3]}; {snow_depth_borders[4]}};\tЗначений в инт
        ервале: {len(df[(df.snow_depth >= snow_depth_borders[3]) & (df.snow_depth < snow_depth_bor
        ders[4])])}
10.    '''
11.
12. plt.axvline(x=snow_depth_borders[1], color='#f91155')
13. plt.axvline(x=snow_depth_borders[2], color='#f91155')
14. plt.axvline(x=snow_depth_borders[3], color='#f91155')

```

Рассмотрим распределение, и выберем для него границы интервалов

Границы интервала: [0; 0.1);	Значений в интервале: 249
Границы интервала: [0.1; 5);	Значений в интервале: 46
Границы интервала: [5; 13);	Значений в интервале: 56
Границы интервала: [13; 27);	Значений в интервале: 24

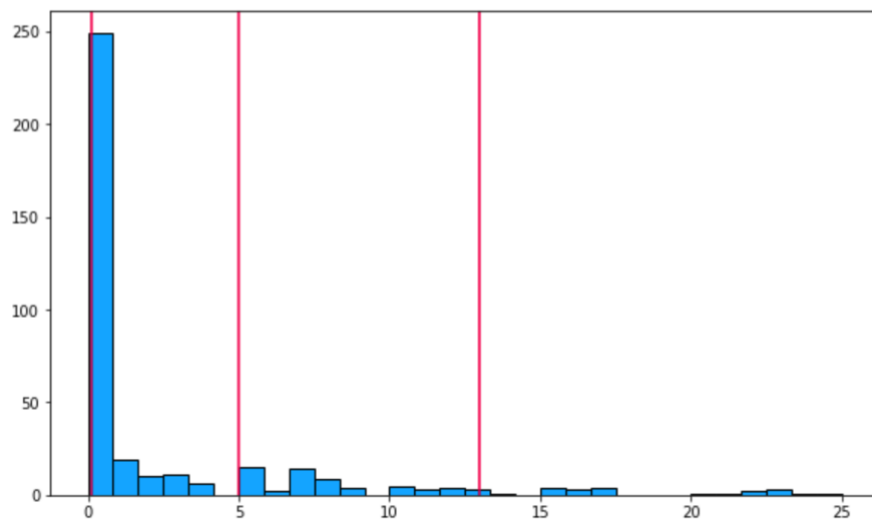


Иллюстрация 11 - Разбиение объектов по признаку snow_depth на 4 интервала

Для выделения третьего номинального признака рассмотрим распределение признака downfall:

```

1. fig, graph = plt.subplots(figsize=(10, 6))
2. plt.hist(df['downfall'], bins=30, color='#00b5ff', label='downfall', edgecolor='black')
3.
4. downfall_borders = [0, 0.1, 4, 7, 28]
5. print(f'''Рассмотрим распределение, и выберем для него границы интервалов
6.     Границы интервала: [{downfall_borders[0]}; {downfall_borders[1]}); \tЗначений в интерва
7.     ле: {len(df[(df.downfall >= downfall_borders[0]) & (df.downfall < downfall_borders[1])])}
8.     Границы интервала: [{downfall_borders[1]}; {downfall_borders[2]}); \tЗначений в интерва
9.     ле: {len(df[(df.downfall >= downfall_borders[1]) & (df.downfall < downfall_borders[2])])}
10.    Границы интервала: [{downfall_borders[2]}; {downfall_borders[3]}); \tЗначений в интер
11.    вале: {len(df[(df.downfall >= downfall_borders[2]) & (df.downfall < downfall_borders[3])])}
12.    Границы интервала: [{downfall_borders[3]}; {downfall_borders[4]}); \tЗначений в интер
13.    вале: {len(df[(df.downfall >= downfall_borders[3]) & (df.downfall < downfall_borders[4])])}
14.    ''')
```

Рассмотрим распределение, и выберем для него границы интервалов

Границы интервала: [0; 0.1);	Значений в интервале: 245
Границы интервала: [0.1; 4);	Значений в интервале: 91
Границы интервала: [4; 7);	Значений в интервале: 28
Границы интервала: [7; 28);	Значений в интервале: 11

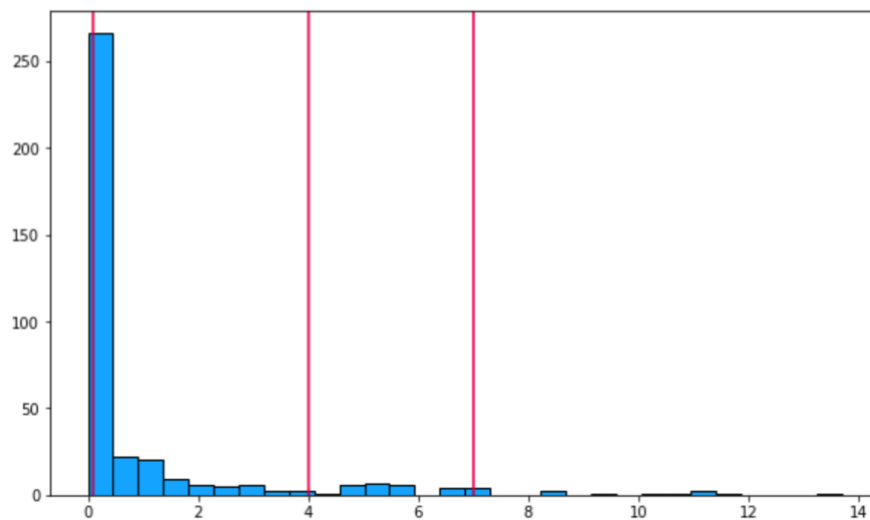


Иллюстрация 12 - Разбиение объектов по признаку downfall на 4 интервала

Сформируем две таблицы сопряженности, x1 и x2, x1 и x3:

```

1. # Таблица сопряженности x1(avg_temp) и x2(snow_depth)
2. t_s = np.arange(16).reshape(4, 4)
3. # Таблица сопряженности x1(avg_temp) и x3(downfall)
4. t_d = np.arange(16).reshape(4, 4)
5.
6. for i in range(4):
7.     for j in range(4):
8.         t_s[j][i] = len(
9.             df[
10.                 (snow_depth_borders[i] <= df['snow_depth']) &
11.                 (df['snow_depth'] < snow_depth_borders[i+1]) &
12.                 (average_temp_borders[j] <= df['avg_temp']) &
13.                 (df['avg_temp'] < average_temp_borders[j+1])
14.             ].index
15.         )
16.         t_d[j][i] = len(
17.             df[
18.                 (downfall_borders[i] <= df['downfall']) &
19.                 (df['downfall'] < downfall_borders[i+1]) &
20.                 (average_temp_borders[j] <= df['avg_temp']) &
21.                 (df['avg_temp'] < average_temp_borders[j+1])
22.             ].index
23.         )

```

		downfall				Total
		[0, 0.1)	[0.1, 4)	[4, 7)	[7, 28)	
avg_temp	[-16.9, 2.2)	69	20	3	1	93
	[-2.2, 1.0)	54	21	10	4	89
	[1.0, 3.6)	59	25	6	4	94
	[3.6, 14.0)	63	25	9	2	99
	Total	245	91	28	11	375

Таблица 2 - Таблица сопряженности x1(avg_temp) и x3(downfall)

		snow_depth				
		[0, 0.1)	[0.1, 5)	[5, 13)	[13, 27)	Total
avg_temp	[-16.9, -2.2)	32	12	31	18	93
	[-2.2, 1.0)	45	21	18	5	89
	[1.0, 3.6)	73	13	7	1	94
	[3.6, 14.0)	99	0	0	0	99
	Total	249	46	56	24	375

Таблица 3 - Таблица сопряженности $x_1(\text{avg_temp})$ и $x_2(\text{snow_depth})$

Построим таблицы условных вероятностей (x_1 от x_2 и x_1 от x_3):

```

1. conditional_probability_t_s = np.zeros((4, 4))
2. conditional_probability_t_d = np.zeros((4, 4))
3.
4. for i in range(4):
5.     for j in range(4):
6.         conditional_probability_t_s[j,i] = round(t_s[j,i] / sum(t_s[:,i]), 2)
7.         conditional_probability_t_d[j,i] = round(t_d[j,i] / sum(t_d[:,i]), 2)

```

		snow_depth			
		[0, 0.1)	[0.1, 5)	[5, 13)	[13, 27)
avg_temp	[-16.9, -2.2)	0,1285	0,2609	0,5536	0,7500
	[-2.2, 1.0)	0,1807	0,4565	0,3214	0,2083
	[1.0, 3.6)	0,2932	0,2826	0,1250	0,0417
	[3.6, 14.0)	0,3976	0,0000	0,0000	0,0000

Таблица 4 - Таблица условных вероятностей $x_1(\text{avg_temp})$ и $x_2(\text{snow_depth})$

		downfall			
		[0, 0.1)	[0.1, 4)	[4, 7)	[7, 28)
avg_temp	[-16.9, 2.2)	0,2816	0,2198	0,1071	0,0909
	[-2.2, 1.0)	0,2204	0,2308	0,3571	0,3636
	[1.0, 3.6)	0,2408	0,2747	0,2143	0,3636
	[3.6, 14.0)	0,2571	0,2747	0,3214	0,1818

Таблица 5 - Таблица условных вероятностей $x_1(\text{avg_temp})$ и $x_3(\text{downfall})$

Построим таблицы коэффициентов Кеттле:

```

1. probability_t_s = np.zeros((4, 4))
2. probability_t_d = np.zeros((4, 4))
3.
4. for i in range(4):
5.     for j in range(4):
6.         probability_t_s[j,i] = round(t_s[j,i] / t_s.sum(), 2)
7.         probability_t_d[j,i] = round(t_d[j,i] / t_d.sum(), 2)
8.
9. kettle_t_s = np.zeros((4, 4))
10. kettle_t_d = np.zeros((4, 4))
11.
12. for i in range(4):
13.     for j in range(4):
14.         kettle_t_s[j,i] = round((probability_t_s[j,i] / (sum(probability_t_s[:,i]) * sum(probability_t_s[j,:]) - 1) * 100, 2)
15.         kettle_t_d[j,i] = round((probability_t_d[j,i] / (sum(probability_t_d[:,i]) * sum(probability_t_d[j,:]) - 1) * 100, 2)

```

		snow_depth			
		[0, 0.1)	[0.1, 5)	[5, 13)	[13, 27)
avg_temp	[-16.9, -2.2)	-48,1798	5,1893	123,2143	202,4194
	[-2.2, 1.0)	-23,8527	92,3547	35,4334	-12,2191
	[1.0, 3.6)	16,9572	12,7428	-50,1330	-83,3777
	[3.6, 14.0)	50,6024	-100,0000	-100,0000	-100,0000

Таблица 6 - Таблица коэффициентов Кетле $x1(avg_temp)$ и $x2(snow_depth)$

		downfall			
		[0, 0.1)	[0.1, 4)	[4, 7)	[7, 28)
avg_temp	[-16.9, 2.2)	13,5616	-11,3789	-56,7972	-63,3431
	[-2.2, 1.0)	-7,1314	-2,7658	50,4815	53,2176
	[1.0, 3.6)	-3,9297	9,5978	-14,5137	45,0677
	[3.6, 14.0)	-2,5974	4,0626	21,7532	-31,1295

Таблица 7 - Таблица коэффициентов Кетле $x1(avg_temp)$ и $x3(downfall)$

Как видно из таблиц, средняя глубина снега растёт с уменьшением средней температуры, а осадки наиболее вероятны при температуре около нуля.

Вычислим и визуализируем среднее значение индекса Кетле на построенных таблицах сопряженности:

```

1. ind_kettle_t_s = np.zeros((4, 4))
2. ind_kettle_t_d = np.zeros((4, 4))
3.
4. for i in range(4):
5.     for j in range(4):
6.         ind_kettle_t_s[j,i] = probability_t_s[j,i]*kettle_t_s[j,i]/100
7.         ind_kettle_t_d[j,i] = probability_t_d[j,i]*kettle_t_s[j,i]/100
8.
9. print(f'
10. Среднее значение индекса Кетле  $x1(avg\_temp)$  и  $x2(snow\_depth)$ : {ind_kettle_t_s.sum()}
11. Среднее значение индекса Кетле  $x1(avg\_temp)$  и  $x3(downfall)$ : {ind_kettle_t_d.sum()}
12. ')

```

		snow_depth				
		[0, 0.1)	[0.1, 5)	[5, 13)	[13, 27)	Total
avg_temp	[-16.9, -2.2)	0,0017	0,1019	0,0972	0,2306	0,4313
	[-2.2, 1.0)	0,0517	0,0170	-0,0016	0,2112	0,2783
	[1.0, 3.6)	0,0044	-0,0094	-0,0022	0,2356	0,2285
	[3.6, 14.0)	0,0000	0,0000	0,0000	0,2614	0,2614
	Total	0,0578	0,1095	0,0933	0,9389	1,1995

Таблица 8 - Таблица индексов Кетле $x1(avg_temp)$ и $x2(snow_depth)$ и среднее значение индекса Кетле (выделено желтым)

		downfall				
		[0, 0.1)	[0.1, 4)	[4, 7)	[7, 28)	Total
avg_temp	[-16.9, 2.2)	0,0250	-0,0061	-0,0045	-0,0017	0,0127
	[-2.2, 1.0)	-0,0103	-0,0015	0,0135	0,0057	0,0073
	[1.0, 3.6)	-0,0062	0,0064	-0,0023	0,0048	0,0027
	[3.6, 14.0)	-0,0044	0,0027	0,0052	-0,0017	0,0019
	Total	0,0041	0,0015	0,0118	0,0071	0,0246

Таблица 9 - Таблица индексов Кетле $x1(avg_temp)$ и $x3(downfall)$ и среднее значение индекса Кетле (выделено желтым)

Прокомментируем смысл значений индекса Кетле на нескольких примерах:

Для температуры больше 3.5 градусов и глубины снега больше 0 значение индекса Кетле -100, т.е. такого не бывает, что вполне логично, так как снег при такой температуре просто растает. Для низкой температуры и глубокого снега значение индекса Кетле очень высокое, т.е. эти события почти всегда происходят одновременно.

Для большого количества осадков и низкой температуры индекс Кетле очень низкий, а для температуры около 0 и осадков индекс Кетле высокий, т.е. осадки наиболее вероятны при температуре у 0, и крайне маловероятны при морозах.

МГК/SVD

Из наших данных более всего под понятие «более или менее относящиеся к одному и тому же аспекту данных», являются признаки:

- `avg_temp` – действительное число – средняя температура в течении дня (в градусах Цельсия);
- `min_temp` – действительное число – минимальная температура в течении дня (в градусах Цельсия);
- `max_temp` – действительное число – максимальная температура в течении дня (в градусах Цельсия);

Подготовим данные для дальнейшего использования:

```

1. av_t = df[['avg_temp', 'min_temp', 'max_temp']].to_numpy()
2. amplitude_scoring = pd.DataFrame((av_t - av_t.mean(axis = 0))/(av_t.max(axis = 0) - av_t.min(axis = 0)), columns = ['avg_temp', 'min_temp', 'max_temp'])
3. amplitude_scoring['snow_depth'] = df['snow_depth'].to_numpy()
4. z_scoring = pd.DataFrame((av_t - av_t.mean(axis = 0))/np.sqrt(((av_t - av_t.mean(axis = 0))*(av_t - av_t.mean(axis = 0))).mean(axis = 0)), columns = ['avg_temp', 'min_temp', 'max_temp'])
5. z_scoring['snow_depth'] = df['snow_depth'].to_numpy()

```

Для визуализации данных с использованием МГК и стандартизации на размах будем использовать модуль `PCA`^[9] библиотеки `SKLearn`^[7]. Визуализируем наши данные с

использованием МГК и стандартизации на размах, при этом цветом отметим наличие/отсутствие снега:

```
1. from sklearn.decomposition import PCA
2.
3. pca_a = PCA(n_components = 2)
4. XPCAreduced = pca_a.fit_transform(amplitude_scoring[['avg_temp', 'min_temp', 'max_temp']]).
    to_numpy())
5. a_s = pd.DataFrame(XPCAreduced, columns = ['col1', 'col2'])
6. a_s['snow_depth'] = df['snow_depth'].to_numpy()
7.
8. fig, graph = plt.subplots(figsize=(10, 6))
9. plt.scatter(
10.     a_s[a_s['snow_depth'] == 0]['col1'],
11.     a_s[a_s['snow_depth'] == 0]['col2'],
12.     label='Нет снега',
13.     cmap='rainbow'
14. )
15. plt.scatter(
16.     a_s[a_s['snow_depth'] > 0]['col1'],
17.     a_s[a_s['snow_depth'] > 0]['col2'],
18.     label='Есть снег',
19.     cmap='rainbow'
20. )
21. plt.title('Визуализация данных с использованием МГК и стандартизации на размах')
22. plt.xlabel('dim_1')
23. plt.ylabel('dim_2')
24. plt.grid()
25. plt.legend(title='Снег')
```

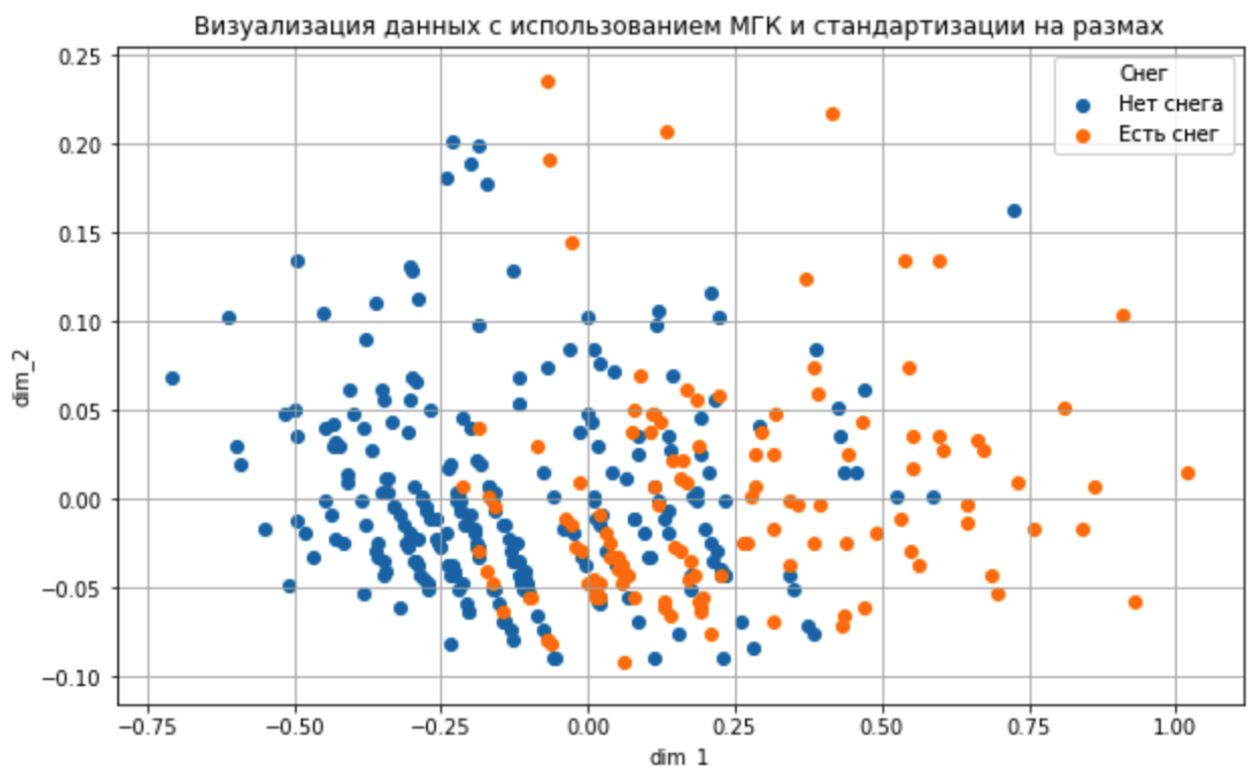


Иллюстрация 13 - Визуализация данных с использованием МГК и стандартизации на размах

Теперь визуализируем наши данные с использованием МГК и z-scoring, при этом цветом отметим наличие/отсутствие снега:

```

1. from sklearn.decomposition import PCA
2.
3. pca_z = PCA(n_components = 2)
4. XPCAreduced = pca_z.fit_transform(z_scoring[['avg_temp', 'min_temp', 'max_temp']].to_numpy())
5. z_s = pd.DataFrame(XPCAreduced, columns = ['col1', 'col2'])
6. z_s['snow_depth'] = df['snow_depth'].to_numpy()
7.
8. fig, graph = plt.subplots(figsize=(10, 6))
9. plt.scatter(
10.     z_s[z_s['snow_depth'] == 0]['col1'],
11.     z_s[z_s['snow_depth'] == 0]['col2'],
12.     label='Нет снега',
13.     cmap='rainbow'
14. )
15. plt.scatter(
16.     z_s[z_s['snow_depth'] > 0]['col1'],
17.     z_s[z_s['snow_depth'] > 0]['col2'],
18.     label='Снег',
19.     cmap='rainbow'
20. )
21.
22. plt.title('Визуализация данных с использованием МГК и z-scoring')
23. plt.xlabel('dim_1')
24. plt.ylabel('dim_2')
25. plt.grid()
26. plt.legend(title='Снег')

```

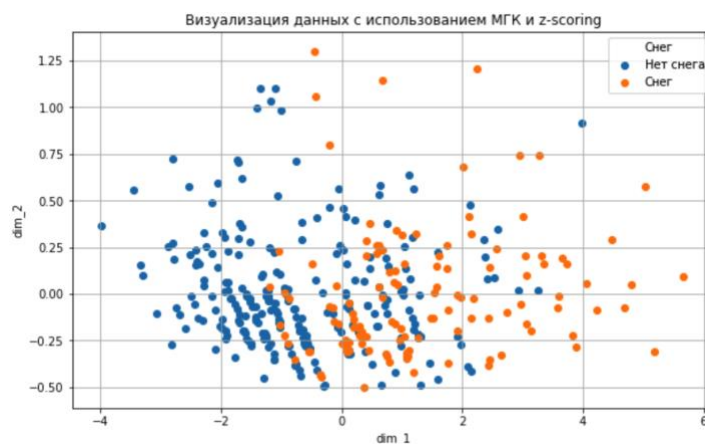


Иллюстрация 14 - Визуализация данных с использованием МГК и z-scoring

А сейчас применим традиционный метод МГК для визуализации этих же данных и сравним с результатом с использованными выше библиотечными методами:

```

1. X = np.transpose(df[['avg_temp', 'min_temp', 'max_temp']].to_numpy())
2. Xcentered = (X[0] - X[0].mean(), X[1] - X[1].mean(), X[2] - X[2].mean())
3. m = (X[0].mean(), X[1].mean(), X[2].mean())
4. covmat = np.cov(Xcentered)
5. _, vecs = np.linalg.eig(covmat)
6. v = [vecs[:,0], vecs[:,2]]
7. Xnew = pd.DataFrame(np.transpose(np.dot(v, Xcentered)), columns = ['col1', 'col2'])
8. Xnew['snow_depth'] = df['snow_depth'].to_numpy()
9.
10. fig, graph = plt.subplots(figsize=(10, 6))
11. plt.scatter(
12.     Xnew[Xnew['snow_depth'] == 0]['col1'],
13.     Xnew[Xnew['snow_depth'] == 0]['col2'],
14.     label='Нет снега',
15.     cmap='rainbow'
16. )
17. plt.scatter(
18.     Xnew[Xnew['snow_depth'] > 0]['col1'],

```

```

19.     Xnew[Xnew['snow_depth'] > 0]['col2'],
20.     label='Есть снег',
21.     cmap='rainbow'
22. )
23.
24. plt.title('Визуализация данных с использованием традиционного МГК, реализованного вручную')
25. plt.legend(title='Снег')
26. plt.xlabel('dim_1')
27. plt.ylabel('dim_2')
28. plt.grid()

```

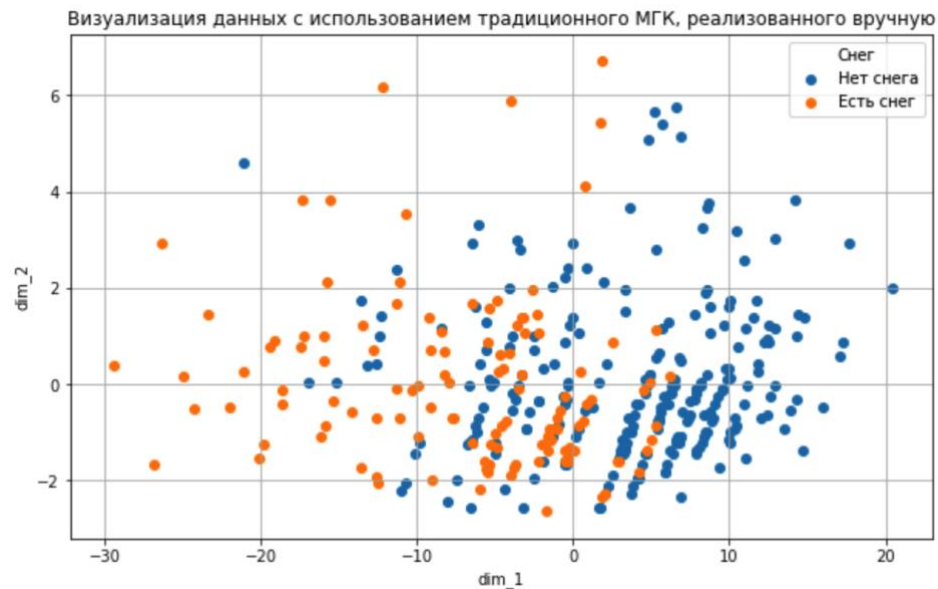


Иллюстрация 15 - Визуализация данных с использованием традиционного МГК, реализованного вручную

С помощью традиционного МКГ, реализованного вручную, мы получили такой же результат, как и при использовании стандартизации на размах и z-scoring, только отзеркаленный по одной из осей.

Рассмотрим главные компоненты:

```

1. print ("Первая главная компонента: ", v[0], "\nВторая главная компонента: ", v[1])
2.
3. >>>Первая главная компонента: [0.57942185 0.60603756 0.54496678]
4. >>>Вторая главная компонента: [ 0.01553545 -0.67673907  0.73605902]

```

В первой главной компоненте коэффициенты почти равны, таким образом первая компонента — это некий общий показатель температуры, на основе всех трёх температурных признаков. Во второй главной компоненте коэффициент перед среднесуточной температурой практически равен 0, а коэффициенты перед минимальной и максимальной температурой почти равны по модулю и противоположны по знаку, значит это разность минимальной и максимальной температуры (умноженная на какой-то коэффициент), т.е. амплитуда колебания температуры в течении суток.

ЛИНЕЙНАЯ РЕГРЕССИЯ

Попробуем найти в наших данных два признака с более-менее линейным полем рассеяния. Для этого воспользуемся функцией `pairplot`^[10] библиотеки `seaborn`^[11].

```
1. import seaborn as sns
2.
3. sns.set()
4. cols = ['date', 'avg_temp', 'min_temp', 'max_temp', 'downfall', 'pressure', 'humidity', 'wind_speed', 'snow_depth', 'year']
5. sns.pairplot(df[cols])
```

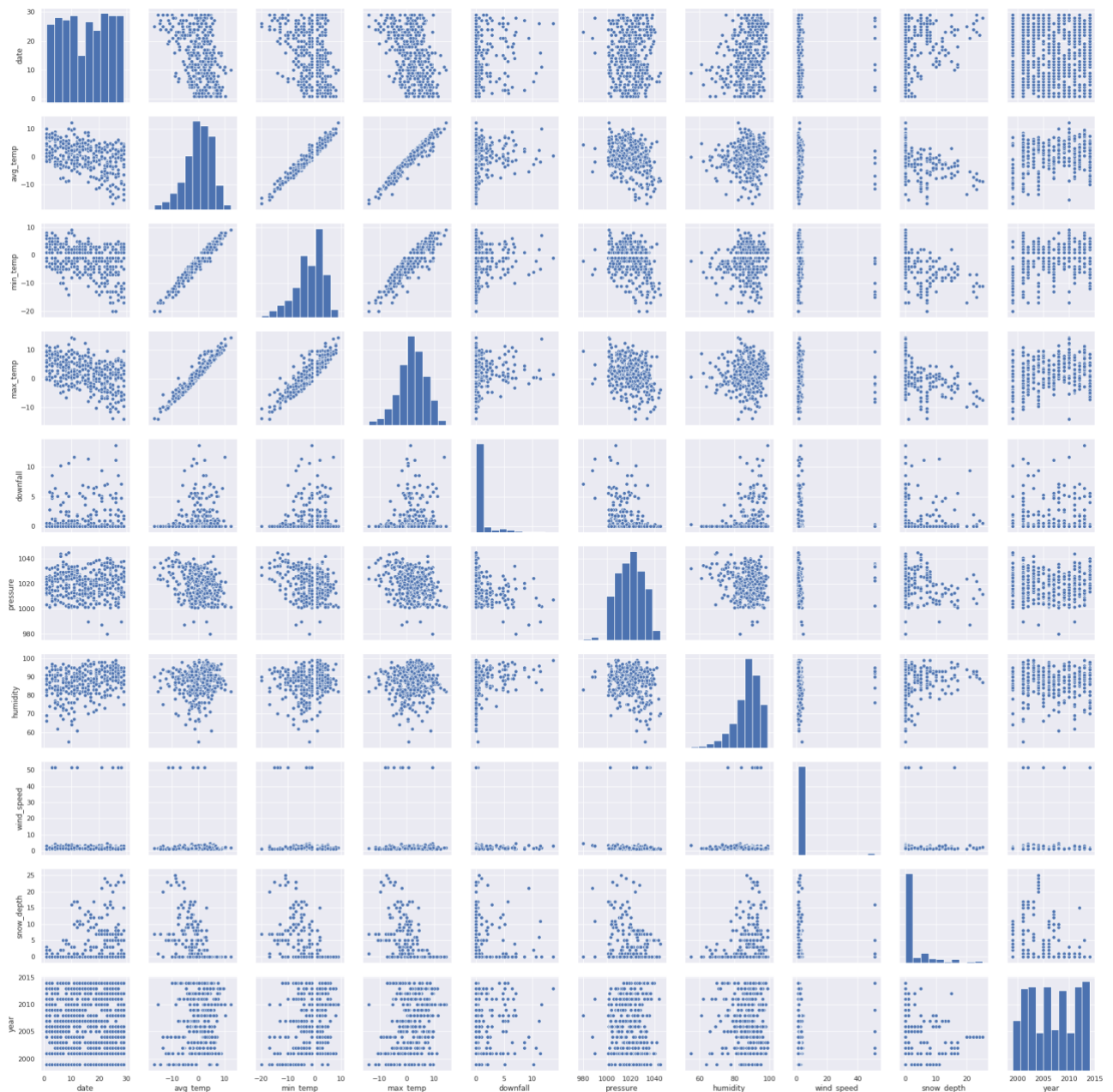


Иллюстрация 16 - Взаимосвязи между всеми признаками в наших данных

Если не считать взаимосвязи между `avg_temp`, `min_temp` и `max_temp` (логично, что они все описывают температуру), то больше остальных на линейное поле рассеяния похожа взаимосвязь между `date` и `avg_temp` (логично, ибо чем ближе к зиме – тем холоднее). Далее будем рассматривать именно эту пару признаков.

Рассмотрим это поле рассеивания более подробно:

```

1. fig, graph = plt.subplots(figsize=(13, 9))
2. for i in range(1, 31):
3.     plt.scatter(
4.         df[df.date == i].date,
5.         df[df.date == i].avg_temp,
6.         label=f'{i} ноября',
7.         cmap='rainbow'
8.     )
9. plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))

```

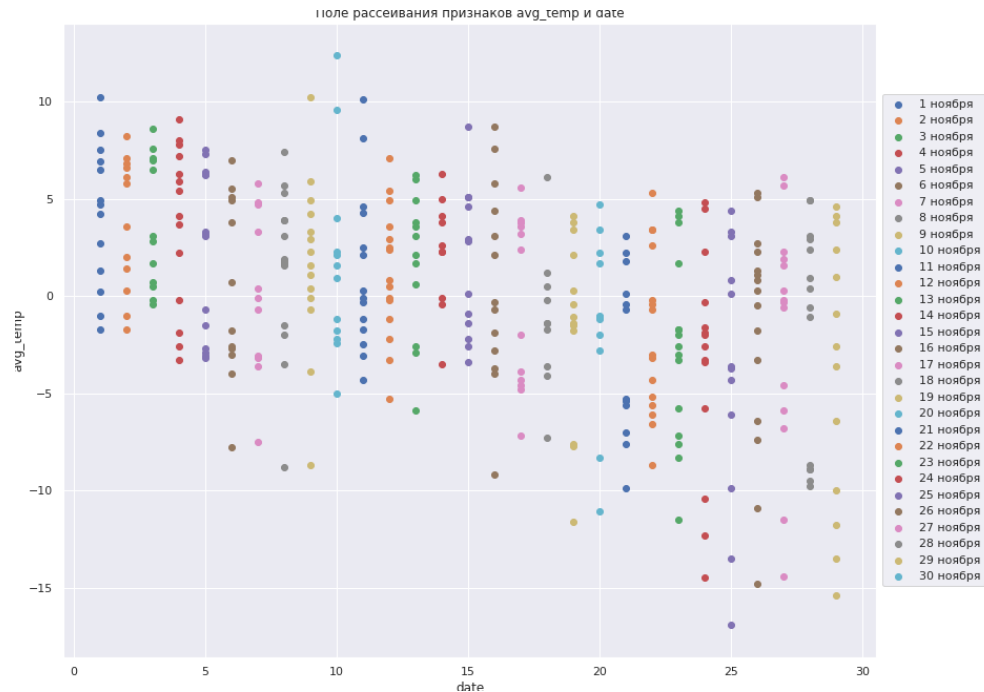


Иллюстрация 17 - Поле рассеивания признаков avg_temp и date

Обучим линейную регрессию:

- date - признак
- avg_temp - целевое значение

Для реализации линейной регрессии будем использовать проверенный временем модуль `LinearRegression`[12] библиотеки `SKLearn`[7].

```

1. from sklearn.linear_model import LinearRegression
2.
3. model = LinearRegression()
4. model.fit(np.array(df['date'])[:,np.newaxis], np.array(df['avg_temp'])[:,np.newaxis])
5.
6. print(f'Уравнение линейной регрессии имеет вид:\ty = {model.coef_.item()} * x + {model.int
  ercept_.item()}')
7.
8. >>>Уравнение линейной регрессии имеет вид: y = -
  0.2550634686459896 * x + 3.9509900360879273

```

Визуализируем результат:

```

1. fig, graph = plt.subplots(figsize=(13, 9))
2. for i in range(1, 31):
3.     plt.scatter(
4.         df[df.date == i].date,
5.         df[df.date == i].avg_temp,

```

```

6.         label=f'{i} ноября',
7.         cmap='rainbow'
8.     )
9.
10. x = np.linspace(0,30)
11. y = model.coef_.item()*x + model.intercept_.item()
12. plt.plot(x, y, label='Уравнение линейной регрессии')
13. plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))

```



Иллюстрация 18 - Поле рассеивания признаков avg_temp и date, уравнение линейной регрессии

Найдем значения коэффициентов корреляции и детерминации. Для первого воспользуемся функцией `pearsonr`[13] библиотеки `SKLearn`[7], а для коэффициента детерминации (также известен в машинном обучении как R^2) – функцией `r2_score`[14] библиотеки `SKLearn`[7].

```

1. from scipy.stats import pearsonr
2. from sklearn.metrics import r2_score
3.
4. corr = pearsonr(model.coef_.item()*df['date'] + model.intercept_.item(), df['avg_temp'])[0]
5. r2 = r2_score(df['avg_temp'], model.coef_.item()*df['date'] + model.intercept_.item())
6.
7. print(f'Значение коэффициента корреляции Пирсона:\t{corr}')
8. print(f'Значение коэффициента детерминации:\t{r2}')
9.
10. >>>Значение коэффициента корреляции Пирсона:    0.4204443991275562
11. >>>Значение коэффициента детерминации:        0.1767734927577317

```

Известно, что чем значение коэффициента детерминации к единице, тем точнее предсказание. В нашем же случае он ближе к нулю, чем к единице, и значит предсказание не самое точное. Сделаем интересное наблюдение, а именно - проверим, какую температуру предскажет модель для нескольких чисел ноября, и сравним её с действительностью за окном.

```

1. print(f'')
2. 1 ноября 2019\tСредняя дневная температура (по Яндекс.Погода): -
  0.5\tПредсказание:\t{model.predict(np.array([1])[0,:np.newaxis]).item()}'

```



```

3. 8 ноября 2019\tСредняя дневная температура (по Яндекс.Погода): +3'\tПредсказание:\t{model.
   predict(np.array([8])[:,np.newaxis]).item()}'
4. 15 ноября 2019\tСредняя дневная температура (по Яндекс.Погода): +3.5'\tПредсказание:\t{mod
   el.predict(np.array([15])[:,np.newaxis]).item()}'
5. 22 ноября 2019\tСредняя дневная температура (по Яндекс.Погода): -
   3.5'\tПредсказание:\t{model.predict(np.array([22])[:,np.newaxis]).item()}'
6. '''
7.
8. 1 ноября 2019    Средняя дневная температура (по Яндекс.Погода): -
   0.5'    Предсказание:    3.695926567441938'
9. 8 ноября 2019    Средняя дневная температура (по Яндекс.Погода): +3'    Предсказание:    1.9104
   822869200104'
10. 15 ноября 2019   Средняя дневная температура (по Яндекс.Погода): +3.5'    Предсказание:    0.
   12503800639808293'
11. 22 ноября 2019   Средняя дневная температура (по Яндекс.Погода): -3.5'    Предсказание:    -
   1.6604062741238446'

```

Предсказание не то, чтобы очень точное, но это и очень грубое приближение в конце концов. А по данным модели, в целом, тренд прослеживается, и примерно понятно, когда пора надеть свитер.

Рассчитаем среднюю относительную ошибку регрессионного уравнения на всех объектах нашей таблицы данных:

```

1. tmp = df[['date', 'avg_temp']]
2. tmp['avg_temp_predict'] = tmp.apply(lambda x: x.date * model.coef_.item() + model.intercep
   t_.item(), axis=1)
3. tmp['diff_percentage'] = tmp.apply(lambda x: abs((x.avg_temp - x.avg_temp_predict) * 100.0
   / x.avg_temp), axis=1)
4. print(f'Значение средней относительной ошибки: {tmp.diff_percentage.mean()}%')
5.
6. >>>Значение средней относительной ошибки: 175.4859768883784%

```

Собственно, как и было видно ранее из значения коэффициента детерминации, линейная регрессия от одного признака дает очень грубое приближение. Мало того, что погода – вещь крайне непредсказуемая (особенно в Москве), так мы еще и пытаемся приблизить её линейной функцией от одного аргумента.

СПИСОК ИСПОЛЬЗОВАННЫХ МАТЕРИАЛОВ

- [1] Источник датасета с погодой в Москве - https://climate-energy.ru/weather/archive_weather_276120.php
- [2] Датасет, использованный в ходе работы
https://github.com/mgcrp/hse_ccda_2019/blob/master/CLEAR_dataset_november_weather.csv
- [3] Среда для разработки Google Colab - <https://colab.research.google.com>
- [4] Pandas – Python-библиотека для анализа данных - <https://pandas.pydata.org>
- [5] NumPy – Python-библиотека для матричных и векторных операций <https://numpy.org>
- [6] Matplotlib – Python-библиотека для построения графиков и визуализации -
<https://matplotlib.org>
- [7] SKLearn – Python-библиотека для машинного обучения и анализа данных <https://scikit-learn.org/stable/>
- [8] Модуль Kmeans библиотеки SKLearn – <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html>
- [9] Модуль PCA библиотеки SKLearn – <https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA.html>
- [10] <https://seaborn.pydata.org/generated/seaborn.pairplot.html>
- [11] Seaborn – Python-библиотека для продвинутого построения графиков и визуализации - <https://seaborn.pydata.org/index.html>
- [12] Модуль LinearRegression библиотеки SKLearn - https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html
- [13] Функция pearsonr библиотеки SKLearn -
<https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.pearsonr.html>
- [14] Функция r2_score библиотеки SKLearn - https://scikit-learn.org/stable/modules/generated/sklearn.metrics.r2_score.html

ПРИЛОЖЕНИЯ

Приложение 1 – Репозиторий GitHub с материалами

https://github.com/mgcrp/hse_ccda_2019

Приложение 2 – Jupyter Notebook с исходным кодом

https://github.com/mgcrp/hse_ccda_2019/blob/master/homework_final.ipynb

Приложение 3 – Полный исходный код домашней работы

```
1. # -*- coding: utf-8 -*-
2. """HW - Миркин.ipynb
3.
4. Automatically generated by Colaboratory.
5.
6. Original file is located at
7.     https://colab.research.google.com/drive/18owHBnh00KAUPuZUXtP5w5tWwHhsYxXb
8.
9. # Курсовой проект - Миркин
10.
11. __Группа__ - Никифоров, Попов
12. """
13.
14. # Подключаю свой Google Drive, там лежит датасет
15. from google.colab import drive
16. drive.mount('/content/gdrive')
17.
18. # Импортирую датасет
19. import pandas as pd
20. df = pd.read_csv('gdrive/My Drive/Colab Notebooks/november_weather.csv', delimiter='\t')
21.
22. # Загрузить датасет (способ для пидорасов)
23. import io
24. from google.colab import files
25.
26. uploaded = files.upload()
27.
28. df = pd.read_csv(io.BytesIO(uploaded['november_weather.csv']), delimiter='\t')
29.
30. """# Задание 2, Кластер-анализ
31. 1. Выберите не менее трех количественных признаков, объясните выбор и примените метод К-
    средних:
32.     * для K=5
33.     * для K=9
34.     * в обоих случаях сделайте порядка 10 случайных инициализаций и выберите то, которое дос
    тавляет минимум критерию метода
35. 2. Проинтерпретируйте оба разбиения с помощью признаков таблицы данных путем сравнения вну
    трикластерных средних с общими средними. Объясните, какое из разбиений лучше с точки зрени
    я интерпретации.
36.
37. ---
38.
39. __Выбираем признаки__:
40. * `avg_temp` - средняя температура в течении дня;
41. * `pressure` - среднее атмосферное давление в течении дня;
42. * `wind_speed` - сркдная скорость ветра;
43. """
44.
```

```

45. df_features = df[['avg_temp', 'pressure', 'wind_speed']].drop(df[df.pressure == '-
'].index).drop(df[df.wind_speed > 45].index)
46. df_features.pressure = df_features.apply(lambda x: float(x.pressure), axis=1)
47.
48. # Commented out IPython magic to ensure Python compatibility.
49. import numpy as np
50. from random import choices
51. import matplotlib.pyplot as plt
52. from sklearn.cluster import KMeans
53. from statistics import mean, stdev
54.
55. # %matplotlib inline
56.
57. """__1.1 - kmeans для 5 кластеров__"""
58.
59. kmeans_5 = KMeans(
60.     n_clusters = 5,          # 5 кластеров
61.     init = 'random',         # центры кластеров - случайно выбранные точки
62.     n_init = 10              # количество инициализаций алгоритма
63. )
64. kmeans_5.fit(df_features) # выполнение алгоритма
65.
66. # Добавлю к исходному датасету созданные кластеры
67. df_features_5 = df_features
68. df_features_5['cluster'] = kmeans_5.labels_
69.
70. # Визуализация
71. fig, graphs = plt.subplots(1, 3, figsize=(17.5, 5))
72.
73. plt.axes(graphs[0])
74. for i in range(5):
75.     plt.scatter(
76.         df_features_5[df_features_5.cluster == i].avg_temp,
77.         df_features_5[df_features_5.cluster == i].pressure,
78.         label=i+1,
79.         cmap='rainbow'
80.     )
81. plt.xlabel('avg_temp')
82. plt.ylabel('pressure')
83. plt.legend(title='Cluster')
84.
85. plt.axes(graphs[1])
86. for i in range(5):
87.     plt.scatter(
88.         df_features_5[df_features_5.cluster == i].avg_temp,
89.         df_features_5[df_features_5.cluster == i].wind_speed,
90.         label=i+1,
91.         cmap='rainbow'
92.     )
93. plt.xlabel('avg_temp')
94. plt.ylabel('wind_speed')
95. plt.legend(title='Cluster')
96.
97. plt.axes(graphs[2])
98. for i in range(5):
99.     plt.scatter(
100.         df_features_5[df_features_5.cluster == i].wind_speed,
101.         df_features_5[df_features_5.cluster == i].pressure,
102.         label=i+1,
103.         cmap='rainbow'
104.     )
105.     plt.xlabel('wind_speed')
106.     plt.ylabel('pressure')
107.     plt.legend(title='Cluster')
108.
109. fig.savefig('gdrive/My Drive/Colab Notebooks/task2_five_clusters.png')
110. print()
111.
112. """__1.2 - kmeans для 9 кластеров__"""

```

```

113.
114.     kmeans_9 = KMeans(
115.         n_clusters = 9,           # 9 кластеров
116.         init = 'random',         # центры кластеров - случайно выбранные точки
117.         n_init = 10              # количество инициализаций алгоритма
118.     )
119.     kmeans_9.fit(df_features) # выполнение алгоритма
120.
121.     # Добавлю к исходному датасету созданные кластеры
122.     df_features_9 = df_features
123.     df_features_9['cluster'] = kmeans_9.labels_
124.
125.     # Визуализация
126.     fig, graphs = plt.subplots(1, 3, figsize=(17.5, 5))
127.
128.     plt.axes(graphs[0])
129.     for i in range(9):
130.         plt.scatter(
131.             df_features_9[df_features_9.cluster == i].avg_temp,
132.             df_features_9[df_features_9.cluster == i].pressure,
133.             label=i+1,
134.             cmap='rainbow'
135.         )
136.     plt.xlabel('avg_temp')
137.     plt.ylabel('pressure')
138.     plt.legend(title='Cluster')
139.
140.     plt.axes(graphs[1])
141.     for i in range(9):
142.         plt.scatter(
143.             df_features_9[df_features_9.cluster == i].avg_temp,
144.             df_features_9[df_features_9.cluster == i].wind_speed,
145.             label=i+1,
146.             cmap='rainbow'
147.         )
148.     plt.xlabel('avg_temp')
149.     plt.ylabel('wind_speed')
150.     plt.legend(title='Cluster')
151.
152.     plt.axes(graphs[2])
153.     for i in range(9):
154.         plt.scatter(
155.             df_features_9[df_features_9.cluster == i].wind_speed,
156.             df_features_9[df_features_9.cluster == i].pressure,
157.             label=i+1,
158.             cmap='rainbow'
159.         )
160.     plt.xlabel('wind_speed')
161.     plt.ylabel('pressure')
162.     plt.legend(title='Cluster')
163.
164.     fig.savefig('gdrive/My Drive/Colab Notebooks/task2_nine_clusters.png')
165.     print()
166.
167.     """__2 - Проинтерпретируйте оба разбиения с помощью признаков таблицы данных путем
    сравнения внутрикластерных средних с общими средними. Объясните, какое из разбиений лучше
    с точки зрения интерпретации__"""
168.
169.     metrics_5cluster = pd.DataFrame(columns=['cluster', 'avg_temperature', 'avg_pressur
    e', 'avg_wind'])
170.     metrics_5cluster = metrics_5cluster.append(
171.         pd.Series(
172.             {
173.                 'cluster': 'All data',
174.                 'avg_temperature': mean(df_features_5.avg_temp),
175.                 'avg_pressure': mean(df_features_5.pressure),
176.                 'avg_wind': mean(df_features_5.wind_speed)
177.             }, ignore_index=True)
178.

```

```

179.     for i in range(5):
180.         metrics_5cluster = metrics_5cluster.append(
181.             pd.Series(
182.                 {
183.                     'cluster': f'Cluster #{i}',
184.                     'avg_temperature': mean(df_features_5[df_features_5.cluster == i].avg_t
emp),
185.                     'avg_pressure': mean(df_features_5[df_features_5.cluster == i].pressure
),
186.                     'avg_wind': mean(df_features_5[df_features_5.cluster == i].wind_speed)
187.                 }), ignore_index=True)
188.
189.     metrics_5cluster = metrics_5cluster.set_index('cluster')
190.     metrics_5cluster
191.
192.     # data to plot
193.     n_groups = 5
194.
195.     # create plot
196.     fig, ax = plt.subplots(figsize=(10, 6))
197.     index = np.arange(n_groups)
198.     bar_width = 0.2
199.     opacity = 0.8
200.
201.     metrics_5cluster_viz = metrics_5cluster[metrics_5cluster.index != 'All data']
202.
203.     plt.bar(index, metrics_5cluster_viz.avg_temperature - [metrics_5cluster.avg_tempera
ture['All data']]*5, bar_width, label='avg_temperature')
204.     plt.bar(index + bar_width, metrics_5cluster_viz.avg_pressure - [metrics_5cluster.av
g_pressure['All data']]*5, bar_width, label='avg_pressure')
205.     plt.bar(index + 2*bar_width, metrics_5cluster_viz.avg_wind - [metrics_5cluster.avg_
wind['All data']]*5, bar_width, label='avg_wind')
206.
207.     plt.xlabel('Кластеры')
208.     plt.ylabel('Разница с общим средним')
209.     plt.xticks(index + bar_width, metrics_5cluster_viz.index)
210.     plt.legend()
211.     plt.grid()
212.
213.     plt.tight_layout()
214.     plt.show()
215.
216.     metrics_9cluster = pd.DataFrame(columns=['cluster', 'avg_temperature', 'avg_pressur
e', 'avg_wind'])
217.     metrics_9cluster = metrics_9cluster.append(
218.         pd.Series(
219.             {
220.                 'cluster': 'All data',
221.                 'avg_temperature': mean(df_features_9.avg_temp),
222.                 'avg_pressure': mean(df_features_9.pressure),
223.                 'avg_wind': mean(df_features_9.wind_speed)
224.             }), ignore_index=True)
225.
226.     for i in range(9):
227.         metrics_9cluster = metrics_9cluster.append(
228.             pd.Series(
229.                 {
230.                     'cluster': f'Cluster #{i}',
231.                     'avg_temperature': mean(df_features_9[df_features_9.cluster == i].avg_t
emp),
232.                     'avg_pressure': mean(df_features_9[df_features_9.cluster == i].pressure
),
233.                     'avg_wind': mean(df_features_9[df_features_9.cluster == i].wind_speed)
234.                 }), ignore_index=True)
235.
236.     metrics_9cluster = metrics_9cluster.set_index('cluster')
237.     metrics_9cluster

```

```

238.
239.     # data to plot
240.     n_groups = 9
241.
242.     # create plot
243.     fig, ax = plt.subplots(figsize=(10, 6))
244.     index = np.arange(n_groups)
245.     bar_width = 0.2
246.     opacity = 0.8
247.
248.     metrics_9cluster_viz = metrics_9cluster[metrics_9cluster.index != 'All data']
249.
250.     plt.bar(index, metrics_9cluster_viz.avg_temperature - [metrics_9cluster.avg_tempera
251.         ture['All data']]*9, bar_width, label='avg_temperature')
252.     plt.bar(index + bar_width, metrics_9cluster_viz.avg_pressure - [metrics_9cluster.av
253.         g_pressure['All data']]*9, bar_width, label='avg_pressure')
254.     plt.bar(index + 2*bar_width, metrics_9cluster_viz.avg_wind - [metrics_9cluster.avg_
255.         wind['All data']]*9, bar_width, label='avg_wind')
256.
257.     plt.xlabel('Кластеры')
258.     plt.ylabel('Разница с общим средним')
259.     plt.xticks(index + bar_width, metrics_9cluster_viz.index)
260.     plt.legend()
261.     plt.grid()
262.
263.     plt.tight_layout()
264.     plt.show()
265.
266.     """__Вывод__: *yet to do*
267.
268.     # Задание 3 - Бутстреп
269.
270.     Для одного из полученных в д.з. 2 разбиений
271.     * Найдите 95% доверительный интервал для среднего значения какого-
272.     либо признака на всем множестве объектов, используя бутстрэп
273.     * Сравните средние по какому-либо признаку в двух кластерах, используя бутстрэп
274.     * Для одного из кластеров сравните среднее на всем множестве для какого-
275.     либо признака с его средним внутри кластера, используя бутстрэп
276.
277.     **Примечание**: каждое применение бутстрэпа должно быть обоими методами, с опорой и
278.     без
279.
280.     __1 - Найдите 95% доверительный интервал для среднего значения какого-
281.     либо признака на всем множестве объектов, используя бутстрэп__
282.
283.     Анализировать будем признак `avg_temp`
284.
285.     """
286.
287.     from random import choices
288.     from statistics import mean
289.
290.     # Список средних
291.     means = []
292.
293.     # Количество экспериментов
294.     bootstrap_iters = 5000
295.     # Количество выбираемых значений для эксперимента
296.     bootstrap_choices = len(df_features)
297.
298.     for i in range(bootstrap_iters):
299.         # Добавляем в список
300.         means.append(
301.             # Среднее значение
302.             mean(
303.                 choices(
304.                     # Выбираем случайный элемент из значений признака
305.
306.                     df_features_5['avg_temp'].tolist(),
307.                     # k-раз

```

```

299.             k=bootstrap_choices
300.         )
301.     )
302. )
303.
304. # Визуализация
305. fig, graph = plt.subplots(figsize=(10, 6))
306. plt.hist(means, bins=50, color='#00b5ff', edgecolor='black')
307. plt.xlabel('Значение avg_temp')
308. plt.title('Распределение средних значений avg_temp')
309.
310. print()
311.
312. """Найдем 95% доверительный интервал способом "без опоры"""
313.
314. # Сортируем список со средними
315. means_ordered = sorted(means)
316. # Отсекаем 2.5% значений снизу и берем наименьший оставшийся
317. left = means_ordered[int(len(means_ordered)*0.05/2 - 1)]
318. # Отсекаем 2.5% значений сверху и берем наибольший оставшийся
319. right = means_ordered[int(len(means_ordered) - (len(means_ordered)*0.05/2) - 1)]
320.
321. print(f'"""
322.     Границы 95% доверительного интервала:
323.     \tЛевая:\t{left}
324.     \tПравая:\t{right}
325.     ""')
326.
327. """Найдем 95% доверительный интервал способом "с опорой"""
328.
329. from statistics import mean, stdev
330.
331. print(f'"""
332.     Границы 95% доверительного интервала:
333.     \tЛевая:\t{mean(means) - 1.96*stdev(means)}
334.     \tПравая:\t{mean(means) + 1.96*stdev(means)}
335.     ""')
336.
337. """__2 - Сравните средние по какому-
либо признаку в двух кластерах, используя бутстрэп__
338.
339. Будем использовать полученное ранее разбиение на 5 кластеров
340.
341. Пусть сравниваемыми будут кластеры `3` и `4`
342. """
343.
344. means = []
345. bootstrap_iters = 5000
346.
347. for i in range(bootstrap_iters):
348.     means.append(
349.         # Из среднего значения признака в эксперименте на 3 кластере
350.         mean(
351.             choices(
352.                 df_features_5[df_features_5.cluster == 3]['avg_temp'].tolist(),
353.                 k=len(df_features_5[df_features_5.cluster == 3])
354.             )
355.         )
356.         # вычитаем
357.         -
358.         # среднее значение признака в эксперименте на 4 кластере
359.         mean(
360.             choices(
361.                 df_features_5[df_features_5.cluster == 4]['avg_temp'].tolist(),
362.                 k=len(df_features_5[df_features_5.cluster == 4])
363.             )
364.         )
365.     )
366.

```

```

367.     # Визуализация
368.     fig, graph = plt.subplots(figsize=(10, 6))
369.     plt.hist(means, bins=50, color='#00b5ff', edgecolor='black')
370.     plt.xlabel('Значение avg_temp')
371.     plt.title('Распределение средних значений avg_temp')
372.
373.     print()
374.
375.     """0 не попадает в это распределение - средние значения признака avg_temp в кластер
    ах `3` и `4` сильно различаются.
376.
377.     Рассмотрим распределения по-
    отдельности и рассчитаем 95% доверительный интервал для полученного выше распределения.
378.     """
379.
380.     means_cluster3 = []
381.     means_cluster4 = []
382.
383.     bootstrap_iters = 5000
384.
385.     for i in range(bootstrap_iters):
386.         means_cluster3.append(
387.             mean(
388.                 choices(
389.                     df_features_5[df_features_5.cluster == 3]['avg_temp'].tolist(),
390.                     k=len(df_features_5[df_features_5.cluster == 3])
391.                 )
392.             )
393.         )
394.         means_cluster4.append(
395.             mean(
396.                 choices(
397.                     df_features_5[df_features_5.cluster == 4]['avg_temp'].tolist(),
398.                     k=len(df_features_5[df_features_5.cluster == 4])
399.                 )
400.             )
401.         )
402.
403.     fig, graph = plt.subplots(figsize=(10, 6))
404.     plt.hist(means_cluster3, bins=50, color='#00b5ff', label='Cluster 3', edgecolor='black')
405.     plt.hist(means_cluster4, bins=50, color='#f91155', label='Cluster 4', edgecolor='black')
406.     plt.legend(title='Кластеры')
407.     plt.xlabel('Значение avg_temp')
408.     plt.title('Распределение средних значений avg_temp')
409.
410.     print()
411.
412.     """Найдем 95% доверительный интервал способом "без опоры"""
413.
414.     means_ordered = sorted(means) # C
    оптимизируем массив со средними
415.     left = means_ordered[int(len(means_ordered)*0.05/2 - 1)] # 0
    отсекаем 2.5% значений снизу и берем наименьший оставшийся
416.     right = means_ordered[int(len(means_ordered) - (len(means_ordered)*0.05/2) - 1)] # 0
    отсекаем 2.5% значений сверху и берем наибольший оставшийся
417.     print(f'
418.         Границы 95% доверительного интервала:
419.         \tЛевая:\t{left}
420.         \tПравая:\t{right}
421.     ')
422.
423.     """Найдем 95% доверительный интервал способом "с опорой"""
424.
425.     print(f'
426.         Границы 95% доверительного интервала:
427.         \tЛевая:\t{mean(means) - 1.96*stdev(means)}
428.         \tПравая:\t{mean(means) + 1.96*stdev(means)}

```



```

429.     '')
430.
431.     """__3 - Для одного из кластеров сравните среднее на всем множестве для какого-
        либо признака с его средним внутри кластера, используя бутстрэп__
432.
433.     Будем использовать полученное ранее разбиение на 5 кластеров
434.
435.     Сравнивать с выборкой будем значения признака на кластере `1`
436.     """
437.
438.     means = []
439.
440.     bootstrap_iters = 5000
441.
442.     for i in range(bootstrap_iters):
443.         means.append(
444.             mean(
445.                 choices(
446.                     df_features_5[df_features_5.cluster == 1]['avg_temp'].tolist(),
447.                     k=len(df_features_5[df_features_5.cluster == 1])
448.                 )
449.             )
450.             -
451.             mean(
452.                 choices(
453.                     df_features_5['avg_temp'].tolist(),
454.                     k=len(df_features_5)
455.                 )
456.             )
457.         )
458.
459.     fig, graph = plt.subplots(figsize=(10, 6))
460.     plt.hist(means, bins=50, color='#00b5ff', edgecolor='black')
461.     plt.xlabel('Значение avg_temp')
462.     plt.title('Распределение средних значений avg_temp')
463.
464.     print()
465.
466.     """0 попадает в это распределение - средние значения признака в кластере `1` и на в
        сей выборке похожи.
467.
468.     Рассмотрим распределения по-
        отдельности и рассчитаем 95% доверительный интервал для полученного выше распределения.
469.     """
470.
471.     means_cluster1 = [] # П
        принцип аналогичен используемому выше
472.     means_allvalues = []
473.
474.     bootstrap_iters = 5000
475.
476.     for i in range(bootstrap_iters):
477.         means_cluster1.append(
478.             mean(
479.                 choices(
480.                     df_features_5[df_features_5.cluster == 1]['avg_temp'].tolist(),
481.                     k=len(df_features_5[df_features_5.cluster == 1])
482.                 )
483.             )
484.         )
485.         means_allvalues.append(
486.             mean(
487.                 choices(
488.                     df_features_5['avg_temp'].tolist(),
489.                     k=len(df_features_5)
490.                 )
491.             )
492.         )
493.

```

```

494.     fig, graph = plt.subplots(figsize=(10, 6))
495.     plt.hist(means_cluster1, bins=50, color='#00b5ff', label='Cluster 1', edgecolor='black')
496.     plt.hist(means_allvalues, bins=50, color='#f91155', label='All values', edgecolor='black')
497.     plt.legend()
498.     plt.xlabel('Значение avg_temp')
499.     plt.title('Распределение средних значений avg_temp')
500.
501.     print()
502.
503.     """Найдем 95% доверительный интервал способом "без опоры"""
504.
505.     means_ordered = sorted(means) # C
506.     left = means_ordered[int(len(means_ordered)*0.05/2 - 1)] # 0
507.     right = means_ordered[int(len(means_ordered) - (len(means_ordered)*0.05/2) - 1)] # 0
508.     print(f'Границы 95% доверительного интервала:
509.           \tЛевая:\t{left}
510.           \tПравая:\t{right}
511.           ')
512.
513.
514.     """Найдем 95% доверительный интервал способом "с опорой"""
515.
516.     print(f'Границы 95% доверительного интервала:
517.           \tЛевая:\t{mean(means) - 1.96*stdev(means)}
518.           \tПравая:\t{mean(means) + 1.96*stdev(means)}
519.           ')
520.
521.
522.     """# Задание 4 - Таблица сопряженности
523.
524.     1. Сформируйте три номинальных признака x1, x2 и x3 (один из них, но не больше, может быть взят из Вашей таблицы данных)
525.     2. Сформируйте две таблицы сопряженности, x1 и x2, x1 и x3. Постройте также матрицы условных вероятностей (x1 от x2 и x1 от x3), а также таблицы коэффициентов Кетле. Прокомментируйте связи между категориями x1 и категориями x2 и x3.
526.     3. Вычислите и визуализируйте среднее значение индекса Кетле на построенных таблицах сопряженности.
527.     4. Прокомментируйте смысл значений индекса Кетле на двух-трех примерах.
528.
529.     __1 - Сформируйте три номинальных признака x1, x2 и x3 (один из них, но не больше, может быть взят из Вашей таблицы данных)__
530.
531.     В качестве первого номинального признака возьмем разбиение диапазона значений признака `avg_temp` на 4 равных по количеству объектов интервала
532.     """
533.
534.     # Разбиение
535.     temp_sorted = df['avg_temp'].sort_values()
536.     average_temp_borders = [
537.         temp_sorted.iloc[0],
538.         temp_sorted.iloc[len(temp_sorted.index) // 4],
539.         temp_sorted.iloc[(len(temp_sorted.index) // 4) * 2],
540.         temp_sorted.iloc[(len(temp_sorted.index) // 4) * 3],
541.         temp_sorted.iloc[len(temp_sorted.index)-1] + 0.1
542.     ]
543.     # Визуализация
544.     fig, graph = plt.subplots(figsize=(10, 6))
545.     plt.hist(df['avg_temp'], bins=15, color='#00b5ff', label='avg_temp', edgecolor='black')
546.     plt.axvline(x=average_temp_borders[1], color='#f91155')
547.     plt.axvline(x=average_temp_borders[2], color='#f91155')
548.     plt.axvline(x=average_temp_borders[3], color='#f91155')
549.     # Вывод
550.     print(f'

```

```

551.         Границы интервала: [{average_temp_borders[0]}; {average_temp_borders[1]}};\t3на
          чений в интервале: {len(df[(df.avg_temp >= average_temp_borders[0]) & (df.avg_temp < avera
ge_temp_borders[1])])}}
552.         Границы интервала: [{average_temp_borders[1]}; {average_temp_borders[2]}};\t3на
          чений в интервале: {len(df[(df.avg_temp >= average_temp_borders[1]) & (df.avg_temp < avera
ge_temp_borders[2])])}}
553.         Границы интервала: [{average_temp_borders[2]}; {average_temp_borders[3]}};\t3на
          чений в интервале: {len(df[(df.avg_temp >= average_temp_borders[2]) & (df.avg_temp < avera
ge_temp_borders[3])])}}
554.         Границы интервала: [{average_temp_borders[3]}; {average_temp_borders[4]}};\t3на
          чений в интервале: {len(df[(df.avg_temp >= average_temp_borders[3]) & (df.avg_temp < avera
ge_temp_borders[4])])}}
555.         '')
556.
557.         ""Для выделения второго номинального признака рассмотрим распределение признака `s
now_depth`""
558.
559.         fig, graph = plt.subplots(figsize=(10, 6))
560.         plt.hist(df['snow_depth'], bins=30, color='#00b5ff', label='snow_depth', edgecolor=
'black')
561.
562.         snow_depth_borders = [0, 0.1, 5, 13, 27]
563.         print(f''''Рассмотрим распределение, и выберем для него границы интервалов
564.         Границы интервала: [{snow_depth_borders[0]}; {snow_depth_borders[1]}};\t3значени
          й в интервале: {len(df[(df.snow_depth >= snow_depth_borders[0]) & (df.snow_depth < snow_de
pth_borders[1])])}}
565.         Границы интервала: [{snow_depth_borders[1]}; {snow_depth_borders[2]}};\t3значени
          й в интервале: {len(df[(df.snow_depth >= snow_depth_borders[1]) & (df.snow_depth < snow_de
pth_borders[2])])}}
566.         Границы интервала: [{snow_depth_borders[2]}; {snow_depth_borders[3]}};\t\t3значе
          ний в интервале: {len(df[(df.snow_depth >= snow_depth_borders[2]) & (df.snow_depth < snow_
depth_borders[3])])}}
567.         Границы интервала: [{snow_depth_borders[3]}; {snow_depth_borders[4]}};\t3значени
          й в интервале: {len(df[(df.snow_depth >= snow_depth_borders[3]) & (df.snow_depth < snow_de
pth_borders[4])])}}
568.         '')
569.
570.         plt.axvline(x=snow_depth_borders[1], color='#f91155')
571.         plt.axvline(x=snow_depth_borders[2], color='#f91155')
572.         plt.axvline(x=snow_depth_borders[3], color='#f91155')
573.
574.         print()
575.
576.         ""Для выделения третьего номинального признака рассмотрим распределение признака `
downfall`""
577.
578.         fig, graph = plt.subplots(figsize=(10, 6))
579.         plt.hist(df['downfall'], bins=30, color='#00b5ff', label='downfall', edgecolor='bla
ck')
580.
581.         downfall_borders = [0, 0.1, 4, 7, 28]
582.         print(f''''Рассмотрим распределение, и выберем для него границы интервалов
583.         Границы интервала: [{downfall_borders[0]}; {downfall_borders[1]}};\t3значений в
          интервале: {len(df[(df.downfall >= downfall_borders[0]) & (df.downfall < downfall_borders[
1])])}}
584.         Границы интервала: [{downfall_borders[1]}; {downfall_borders[2]}};\t3значений в
          интервале: {len(df[(df.downfall >= downfall_borders[1]) & (df.downfall < downfall_borders[
2])])}}
585.         Границы интервала: [{downfall_borders[2]}; {downfall_borders[3]}};\t\t3значений
          в интервале: {len(df[(df.downfall >= downfall_borders[2]) & (df.downfall < downfall_border
s[3])])}}
586.         Границы интервала: [{downfall_borders[3]}; {downfall_borders[4]}};\t\t3значений
          в интервале: {len(df[(df.downfall >= downfall_borders[3]) & (df.downfall < downfall_border
s[4])])}}
587.         '')
588.
589.         plt.axvline(x=downfall_borders[1], color='#f91155')
590.         plt.axvline(x=downfall_borders[2], color='#f91155')
591.         plt.axvline(x=downfall_borders[3], color='#f91155')

```

```

592.
593.     print()
594.
595.     """__2 - Сформируйте две таблицы сопряженности, x1 и x2, x1 и x3. Постройте также м
атрицы условных вероятностей (x1 от x2 и x1 от x3), а также таблицы коэффициентов Кетле. П
рокомментируйте связи между категориями x1 и категориями x2 и x3__"""
596.
597.     # Таблица сопряженности x1(avg_temp) и x2(snow_depth)
598.     t_s = np.arange(16).reshape(4, 4)
599.     # Таблица сопряженности x1(avg_temp) и x3(downfall)
600.     t_d = np.arange(16).reshape(4, 4)
601.
602.     for i in range(4):
603.         for j in range(4):
604.             t_s[j][i] = len(
605.                 df[
606.                     (snow_depth_borders[i] <= df['snow_depth']) &
607.                     (df['snow_depth'] < snow_depth_borders[i+1]) &
608.                     (average_temp_borders[j] <= df['avg_temp']) &
609.                     (df['avg_temp'] < average_temp_borders[j+1])
610.                 ].index
611.             )
612.             t_d[j][i] = len(
613.                 df[
614.                     (downfall_borders[i] <= df['downfall']) &
615.                     (df['downfall'] < downfall_borders[i+1]) &
616.                     (average_temp_borders[j] <= df['avg_temp']) &
617.                     (df['avg_temp'] < average_temp_borders[j+1])
618.                 ].index
619.             )
620.
621.     print(f'"""
622.         Таблица сопряженности x1(avg_temp) и x2(snow_depth):
623.         {t_s}
624.
625.         Таблица сопряженности x1(avg_temp) и x3(downfall):
626.         {t_d}
627.         ""')
628.
629.     conditional_probability_t_s = np.zeros((4, 4))
630.     conditional_probability_t_d = np.zeros((4, 4))
631.
632.     for i in range(4):
633.         for j in range(4):
634.             conditional_probability_t_s[j,i] = round(t_s[j,i] / sum(t_s[:,i]), 2)
635.             conditional_probability_t_d[j,i] = round(t_d[j,i] / sum(t_d[:,i]), 2)
636.
637.     print(f'"""
638.         Таблица условных вероятностей x1(avg_temp) и x2(snow_depth):
639.         {conditional_probability_t_s}
640.
641.         Таблица условных вероятностей x1(avg_temp) и x3(downfall):
642.         {conditional_probability_t_d}
643.         ""')
644.
645.     probability_t_s = np.zeros((4, 4))
646.     probability_t_d = np.zeros((4, 4))
647.
648.     for i in range(4):
649.         for j in range(4):
650.             probability_t_s[j,i] = round(t_s[j,i] / t_s.sum(), 2)
651.             probability_t_d[j,i] = round(t_d[j,i] / t_d.sum(), 2)
652.
653.     ketle_t_s = np.zeros((4, 4))
654.     ketle_t_d = np.zeros((4, 4))
655.
656.     for i in range(4):
657.         for j in range(4):

```

```

658.         ketle_t_s[j,i] = round((probability_t_s[j,i] / ( sum(probability_t_s[:,i]) *
sum(probability_t_s[j,:]) )) - 1)*100, 2)
659.         ketle_t_d[j,i] = round((probability_t_d[j,i] / ( sum(probability_t_d[:,i]) *
sum(probability_t_d[j,:]) )) - 1)*100, 2)
660.
661.         print(f'''
662.             Таблица коэффициентов Кетле x1(avg_temp) и x2(snow_depth):
663.             {ketle_t_s}
664.
665.             Таблица коэффициентов Кетле x1(avg_temp) и x3(downfall):
666.             {ketle_t_d}
667.             ''')
668.
669.         """Как видно из таблиц, средняя глубина снега растёт с уменьшением средней
температуры, а осадки наиболее вероятны при температуре около 0.
670.
671.         __3 - Вычислите и визуализируйте среднее значение индекса Кетле на построенных табл
ицах сопряженности__
672.         """
673.
674.         ind_ketle_t_s = np.zeros((4, 4))
675.         ind_ketle_t_d = np.zeros((4, 4))
676.
677.         for i in range(4):
678.             for j in range(4):
679.                 ind_ketle_t_s[j,i] = probability_t_s[j,i]*ketle_t_s[j,i]/100
680.                 ind_ketle_t_d[j,i] = probability_t_d[j,i]*ketle_t_s[j,i]/100
681.
682.         print(f'''
683.             avg ketle t-s: {ind_ketle_t_s.sum()}
684.             avg ketle t-d: {ind_ketle_t_d.sum()}
685.             ''')
686.
687.         print(f'''
688.             {avg_ketle_t_s}
689.
690.             {avg_ketle_t_d}
691.             ''')
692.
693.         """__4 - Прокомментируйте смысл значений индекса Кетле на двух-трех примерах__
694.
695.         Для температуры больше 3.5 градусов и глубины снега больше 0 значение индекса Кетле
-
100, т.е. такого не бывает, что вполне логично, так как снег при такой температуре просто
растает. Для низкой температуры и глубокого снега значение индекса Кетле очень высокое, т.
е. эти события почти всегда происходят одновременно.
696.
697.         Для большого количества осадков и низкой температуры индекс Кетле очень низкий, а д
ля температуры около 0 и осадков индекс Кетле высокий, т.е. осадки наиболее вероятны при т
емпературе у 0, и крайне маловероятны при морозах.
698.
699.         # Задание 5: МГК/SVD
700.
701.         1. Выберите в Ваших данных 3-
6 признаков, более или менее относящихся к одному и тому же аспекту данных; откомментируйте
702.
703.         2. Визуализируйте Ваши данные дважды, один раз с использованием стандартизации на
размах, второй – путем z-scoring. Выберите какую-
либо группу объектов и выделите ее на визуализациях цветом или формой «точек»
704.
705.         3. Примените традиционный метод МГК для визуализации (для какого-
либо способа нормализации) и убедитесь, что получено то же самое, что и выше. Если нет – о
бьясните, в чем дело
706.
707.         4. Постарайтесь проинтерпретировать полученные компоненты
708.
709.         __1 - Выберите в Ваших данных 3-
6 признаков, более или менее относящихся к одному и тому же аспекту данных; откомментируйте__
710.         """

```

```

709.     av_t = df[['avg_temp', 'min_temp', 'max_temp']].to_numpy()
710.     amplitude_scoring = pd.DataFrame((av_t - av_t.mean(axis = 0))/(av_t.max(axis = 0) -
av_t.min(axis = 0)), columns = ['avg_temp', 'min_temp', 'max_temp'])
711.     amplitude_scoring['snow_depth'] = df['snow_depth'].to_numpy()
712.     z_scoring = pd.DataFrame((av_t - av_t.mean(axis = 0))/np.sqrt( ((av_t - av_t.mean(a
xis = 0))*(av_t - av_t.mean(axis = 0))).mean(axis = 0) ), columns = ['avg_temp', 'min_temp
', 'max_temp'])
713.     z_scoring['snow_depth'] = df['snow_depth'].to_numpy()
714.
715.     """__2 - Визуализируйте Ваши данные дважды, один раз с использованием стандартизаци
и на размах, второй – путем z-scoring. Выберите какую-
либо группу объектов и выделите ее на визуализациях цветом или формой «точек»__"""
716.
717.     print('Визуализируем данные с использованием стандартизации на размах:\n')
718.     fig, graphs = plt.subplots(1, 3, figsize=(17.5, 5))
719.     plt.axes(graphs[0])
720.     plt.scatter(
721.         amplitude_scoring[amplitude_scoring['snow_depth'] == 0]['avg_temp'],
722.         amplitude_scoring[amplitude_scoring['snow_depth'] == 0]['min_temp'],
723.         label='no_snow',
724.         cmap='rainbow'
725.     )
726.     plt.scatter(
727.         amplitude_scoring[amplitude_scoring['snow_depth'] > 0]['avg_temp'],
728.         amplitude_scoring[amplitude_scoring['snow_depth'] > 0]['min_temp'],
729.         label='snow',
730.         cmap='rainbow'
731.     )
732.     plt.xlabel('avg_temp')
733.     plt.ylabel('min_temp')
734.     plt.legend(title='is_snow')
735.
736.     plt.axes(graphs[1])
737.     plt.scatter(
738.         amplitude_scoring[amplitude_scoring['snow_depth'] == 0]['avg_temp'],
739.         amplitude_scoring[amplitude_scoring['snow_depth'] == 0]['max_temp'],
740.         label='no_snow',
741.         cmap='rainbow'
742.     )
743.     plt.scatter(
744.         amplitude_scoring[amplitude_scoring['snow_depth'] > 0]['avg_temp'],
745.         amplitude_scoring[amplitude_scoring['snow_depth'] > 0]['max_temp'],
746.         label='snow',
747.         cmap='rainbow'
748.     )
749.     plt.xlabel('avg_temp')
750.     plt.ylabel('max_temp')
751.     plt.legend(title='is_snow')
752.
753.     plt.axes(graphs[2])
754.     plt.scatter(
755.         amplitude_scoring[amplitude_scoring['snow_depth'] == 0]['min_temp'],
756.         amplitude_scoring[amplitude_scoring['snow_depth'] == 0]['max_temp'],
757.         label='no_snow',
758.         cmap='rainbow'
759.     )
760.     plt.scatter(
761.         amplitude_scoring[amplitude_scoring['snow_depth'] > 0]['min_temp'],
762.         amplitude_scoring[amplitude_scoring['snow_depth'] > 0]['max_temp'],
763.         label='no_snow',
764.         cmap='rainbow'
765.     )
766.     plt.xlabel('min_temp')
767.     plt.ylabel('max_temp')
768.     plt.legend(title='is_snow')
769.
770.     from sklearn.decomposition import PCA
771.
772.     pca_a = PCA(n_components = 2)

```

```

773.     XPCAreduced = pca_a.fit_transform(amplitude_scoring[['avg_temp', 'min_temp', 'max_t
emp']]).to_numpy()
774.     a_s = pd.DataFrame(XPCAreduced, columns = ['col1', 'col2'])
775.     a_s['snow_depth'] = df['snow_depth'].to_numpy()
776.
777.     fig, graph = plt.subplots(figsize=(10, 6))
778.     plt.scatter(
779.         a_s[a_s['snow_depth'] == 0]['col1'],
780.         a_s[a_s['snow_depth'] == 0]['col2'],
781.         label='Нет снега',
782.         cmap='rainbow'
783.     )
784.     plt.scatter(
785.         a_s[a_s['snow_depth'] > 0]['col1'],
786.         a_s[a_s['snow_depth'] > 0]['col2'],
787.         label='Есть снег',
788.         cmap='rainbow'
789.     )
790.     plt.title('Визуализация данных с использованием МГК и стандартизации на размах')
791.     plt.xlabel('dim_1')
792.     plt.ylabel('dim_2')
793.     plt.grid()
794.     plt.legend(title='Снег')
795.
796.     print()
797.
798.     fig, graphs = plt.subplots(1, 3, figsize=(17.5, 5))
799.     fig.suptitle('Визуализация данных с использованием z-scoring', fontsize=16)
800.
801.     plt.axes(graphs[0])
802.     plt.scatter(
803.         z_scoring[z_scoring['snow_depth'] == 0]['avg_temp'],
804.         z_scoring[z_scoring['snow_depth'] == 0]['min_temp'],
805.         label='Нет снега',
806.         cmap='rainbow'
807.     )
808.     plt.scatter(
809.         z_scoring[z_scoring['snow_depth'] > 0]['avg_temp'],
810.         z_scoring[z_scoring['snow_depth'] > 0]['min_temp'],
811.         label='Есть снег',
812.         cmap='rainbow'
813.     )
814.     plt.xlabel('avg_temp')
815.     plt.ylabel('min_temp')
816.     plt.legend(title='Снег')
817.
818.     plt.axes(graphs[1])
819.     plt.scatter(
820.         z_scoring[z_scoring['snow_depth'] == 0]['avg_temp'],
821.         z_scoring[z_scoring['snow_depth'] == 0]['max_temp'],
822.         label='Нет снега',
823.         cmap='rainbow'
824.     )
825.     plt.scatter(
826.         z_scoring[z_scoring['snow_depth'] > 0]['avg_temp'],
827.         z_scoring[z_scoring['snow_depth'] > 0]['max_temp'],
828.         label='Есть снег',
829.         cmap='rainbow'
830.     )
831.     plt.xlabel('avg_temp')
832.     plt.ylabel('max_temp')
833.     plt.legend(title='Снег')
834.
835.     plt.axes(graphs[2])
836.     plt.scatter(
837.         z_scoring[z_scoring['snow_depth'] == 0]['min_temp'],
838.         z_scoring[z_scoring['snow_depth'] == 0]['max_temp'],
839.         label='Нет снега',
840.         cmap='rainbow'

```

```

841.     )
842.     plt.scatter(
843.         z_scoring[z_scoring['snow_depth'] > 0]['min_temp'],
844.         z_scoring[z_scoring['snow_depth'] > 0]['max_temp'],
845.         label='Есть снег',
846.         cmap='rainbow'
847.     )
848.     plt.xlabel('min_temp')
849.     plt.ylabel('max_temp')
850.     plt.legend(title='Снег')
851.     fig.savefig('gdrive/My Drive/Colab Notebooks/task5_zscoring.png')
852.
853.     print()
854.
855.     from sklearn.decomposition import PCA
856.
857.     pca_z = PCA(n_components = 2)
858.     XPCAreduced = pca_z.fit_transform(z_scoring[['avg_temp', 'min_temp', 'max_temp']].to_
o_numpy())
859.     z_s = pd.DataFrame(XPCAreduced, columns = ['col1', 'col2'])
860.     z_s['snow_depth'] = df['snow_depth'].to_numpy()
861.
862.     fig, graph = plt.subplots(figsize=(10, 6))
863.     plt.scatter(
864.         z_s[z_s['snow_depth'] == 0]['col1'],
865.         z_s[z_s['snow_depth'] == 0]['col2'],
866.         label='Нет снега',
867.         cmap='rainbow'
868.     )
869.     plt.scatter(
870.         z_s[z_s['snow_depth'] > 0]['col1'],
871.         z_s[z_s['snow_depth'] > 0]['col2'],
872.         label='Снег',
873.         cmap='rainbow'
874.     )
875.
876.     plt.title('Визуализация данных с использованием МГК и z-scoring')
877.     plt.xlabel('dim_1')
878.     plt.ylabel('dim_2')
879.     plt.grid()
880.     plt.legend(title='Снег')
881.
882.     print()
883.
884.     X = np.transpose(df[['avg_temp', 'min_temp', 'max_temp']].to_numpy())
885.     Xcentered = (X[0] - X[0].mean(), X[1] - X[1].mean(), X[2] - X[2].mean())
886.     m = (X[0].mean(), X[1].mean(), X[2].mean())
887.     covmat = np.cov(Xcentered)
888.     _, vecs = np.linalg.eig(covmat)
889.     v = [vecs[:,0], vecs[:,2]]
890.     Xnew = pd.DataFrame(np.transpose(np.dot(v,Xcentered)), columns = ['col1', 'col2'])
891.
892.     Xnew['snow_depth'] = df['snow_depth'].to_numpy()
893.
894.     fig, graph = plt.subplots(figsize=(10, 6))
895.     plt.scatter(
896.         Xnew[Xnew['snow_depth'] == 0]['col1'],
897.         Xnew[Xnew['snow_depth'] == 0]['col2'],
898.         label='Нет снега',
899.         cmap='rainbow'
900.     )
901.     plt.scatter(
902.         Xnew[Xnew['snow_depth'] > 0]['col1'],
903.         Xnew[Xnew['snow_depth'] > 0]['col2'],
904.         label='Есть снег',
905.         cmap='rainbow'
906.     )

```



```

907. plt.title('Визуализация данных с использованием традиционного МГК, реализованного в
ручную')
908. plt.legend(title='Снег')
909. plt.xlabel('dim_1')
910. plt.ylabel('dim_2')
911. plt.grid()
912.
913. print()
914.
915. print ("Первая главная компонента: ", v[0], "\nВторая главная компонента: ", v[1])
916.
917. """В первой главной компоненте коэффициенты почти равны, таким образом первая компо
нента это некий общий показатель температуры, на основе всех трёх температурных признаков.
918.
919. Во второй главной компоненте коэффициент перед среднесуточной температурой практиче
ски равен 0, а коэффициенты перед минимальной и максимальной температурой почти равны по м
одулю и противоположны по знаку, значит это разность минимальной и максимальной температур
ы (домноженная на какой-
то коэффициент), т.е. амплитуда колебания температуры в течении суток.
920.
921. # Задание 6 - Линейная регрессия
922.
923. * По возможности, найдите два признака в Ваших данных с более или менее «линейным»
полем рассеяния; представьте его на графике.
924. * Постройте уравнение линейной регрессии одного из этих признаков через другой при
знак. Сделайте комментарий о смысле величины регрессионного коэффициента при этом другом п
ризнаке.
925. * Найдите значения коэффициентов корреляции и детерминации; сделайте комментарий о
смысле величины последнего.
926. * Сделайте предсказание величины целевого признака на двух-
трех объектах; прокомментируйте результат.
927. * Рассчитайте среднюю относительную ошибку регрессионного уравнения на всех объект
ах Вашей таблицы данных и
928. сравните ее с величиной коэффициента детерминации.
929. """
930.
931. import seaborn as sns
932.
933. sns.set()
934. cols = ['date', 'avg_temp', 'min_temp', 'max_temp', 'downfall', 'pressure', 'humidity', 'w
ind_speed', 'snow_depth', 'year']
935. sns.pairplot(df[cols])
936.
937. fig = plt.gcf()
938. fig.savefig('gdrive/My Drive/Colab Notebooks/task6_pairplot.png')
939. print()
940.
941. """Более остальных на линейное поле рассеивания похоже поле рассеивания признаков `
date` и `avg_temp`"""
942.
943. fig, graph = plt.subplots(figsize=(13, 9))
944. for i in range(1, 31):
945.     plt.scatter(
946.         df[df.date == i].date,
947.         df[df.date == i].avg_temp,
948.         label=f'{i} ноября',
949.         cmap='rainbow'
950.     )
951. plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
952. plt.tight_layout()
953. plt.title('Поле рассеивания признаков avg_temp и date')
954. plt.xlabel('date')
955. plt.ylabel('avg_temp')
956.
957. fig.savefig('gdrive/My Drive/Colab Notebooks/task6_scatter_plot.png')
958.
959. print()

```

```

960.
961.     """Обучим линейную регрессию:
962.         * `date` - признак
963.         * `avg_temp` - целевое значение
964.     """
965.
966.     from sklearn.linear_model import LinearRegression
967.
968.     model = LinearRegression()
969.     model.fit(np.array(df['date'])[:,np.newaxis], np.array(df['avg_temp'])[:,np.newaxis
970. ])
971.     print(f'Уравнение линейной регрессии имеет вид:\ty = {model.coef_.item()} * x + {mo
972. del.intercept_.item()}')
973.
974.     fig, graph = plt.subplots(figsize=(13, 9))
975.     for i in range(1, 31):
976.         plt.scatter(
977.             df[df.date == i].date,
978.             df[df.date == i].avg_temp,
979.             label=f'{i} ноября',
980.             cmap='rainbow'
981.         )
982.
983.     x = np.linspace(0,30)
984.     y = model.coef_.item()*x + model.intercept_.item()
985.     plt.plot(x, y, label='Уравнение линейной регрессии')
986.     plt.legend(loc='center left', bbox_to_anchor=(1, 0.5))
987.     plt.tight_layout()
988.     plt.title('Поле рассеивания признаков avg_temp и date')
989.     plt.xlabel('date')
990.     plt.ylabel('avg_temp')
991.
992.     fig.savefig('gdrive/My Drive/Colab Notebooks/task6_scatter_plot.png')
993.
994.     print()
995.
996.     from scipy.stats import pearsonr
997.     from sklearn.metrics import r2_score
998.
999.     corr = pearsonr(model.coef_.item()*df['date'] + model.intercept_.item(), df['avg_te
1000. mp'])[0]
1001.     r2 = r2_score(df['avg_temp'], model.coef_.item()*df['date'] + model.intercept_.item
1002. ())
1003.
1004.     print(f'Значение коэффициента корреляции Пирсона:\t{corr}')
1005.     print(f'Значение коэффициента детерминации:\t\t{r2}')
1006.
1007.     """Сделаем интересное наблюдение, а именно - проверим, какую погоду температуру мод
1008. ель на нескольких числах ноября и сравним её с действительной за окном;"""
1009.
1010.     print(f''''
1011.         1 ноября 2019\tСредняя дневная температура (по Яндекс.Погода): -
1012.         0.5\tПредсказание:\t{model.predict(np.array([1])[:,np.newaxis]).item()}
1013.         8 ноября 2019\tСредняя дневная температура (по Яндекс.Погода): +3\tПредсказание:\t
1014.         {model.predict(np.array([8])[:,np.newaxis]).item()}
1015.         15 ноября 2019\tСредняя дневная температура (по Яндекс.Погода): +3.5\tПредсказание
1016.         :\t{model.predict(np.array([15])[:,np.newaxis]).item()}
1017.         22 ноября 2019\tСредняя дневная температура (по Яндекс.Погода): -
1018.         3.5\tПредсказание:\t{model.predict(np.array([22])[:,np.newaxis]).item()}
1019.         ''')
1020.
1021.     """Не то, чтобы прямо точно, но мы и не синоптики в конце концов. А по данным модел
1022. и, в целом, тренд прослеживается, и примерно понятно, когда стоит надеть свитер :)"
1023.
1024.     tmp = df[['date', 'avg_temp']]
1025.     tmp['avg_temp_predict'] = tmp.apply(lambda x: x.date * model.coef_.item() + model.i
1026. ntercept_.item(), axis=1)

```

```
1017.     tmp['diff_percentage'] = tmp.apply(lambda x: abs((x.avg_temp - x.avg_temp_predict)
      * 100.0 / x.avg_temp), axis=1)
1018.     print(f'Значение средней относительной ошибки: {tmp.diff_percentage.mean()}%')
```