

1 Описание проекта

Допустим, вы работаете в добывающей компании «ГлавРосГосНефть». Нужно решить, где бурить новую скважину.

Вам предоставлены пробы нефти в трёх регионах: в каждом 10 000 месторождений, где измерили качество нефти и объём её запасов. Постройте модель машинного обучения, которая поможет определить регион, где добыча принесёт наибольшую прибыль. Проанализируйте возможную прибыль и риски.

Шаги для выбора локации:

- В избранном регионе ищем месторождения, для каждого определяем значения признаков;
- Строим модель и оцениваем объём запасов;
- Выбираем месторождения с самыми высокими оценками значений. Количество месторождений зависит от бюджета компании и стоимости разработки одной скважины;
- Прибыль равна суммарной прибыли отобранных месторождений

2 Задача

1. Построить модель для определения региона, где добыча принесёт наибольшую прибыль.
2. Проанализировать возможную прибыль и риски.

3 Загрузка и подготовка данных

В [2]:

```
1 # импорт данных
2 import os
3 from urllib.request import urlretrieve
4 from IPython.display import display
5 from pathlib import Path
6 import urllib
7
8
9 # Основные библиотеки для работы с данными
10 import pandas as pd
11 import numpy as np
12
13 # работа с графиками
14 from matplotlib import pyplot as plt
15 import seaborn as sns
16 from IPython.display import display
17 from PIL import Image
18
19 # библиотека для модели
20 from sklearn.linear_model import LinearRegression
21
22 # Настройка моделей и параметров (разбивка датасета, его перемешивание и т.д.)
23 from sklearn.model_selection import train_test_split
24 from sklearn.model_selection import GridSearchCV
25 from tqdm import tqdm
26 from itertools import product
27
28 # Наши основные показатели
29 from sklearn.preprocessing import StandardScaler
30 from sklearn.metrics import mean_squared_error
31
32 from scipy import stats as st
33 from scipy.stats import t
34 from numpy.random import RandomState
35 import urllib
36 import urllib.request
37
```

executed in 1.03s, finished 12:49:19 2021-08-16

В [3]:

```
1 # отключение предупреждений
2 import warnings; warnings.filterwarnings("ignore", category=Warning)
```

executed in 5ms, finished 12:49:19 2021-08-16

В [4]:

```
1 # обновление библиотек Numpy и Pandas для исключения ошибок версий
2 # !pip3 install --upgrade --user --quiet --no-warn-script-location numpy==1.20.1 pandas
```

executed in 12ms, finished 12:49:19 2021-08-16

В [5]:

```
1 # для просмотра всех столбцов таблицы
2 pd.options.display.max_columns = None
```

executed in 11ms, finished 12:49:19 2021-08-16

B [6]:

```
1 # # обновление библиотек
2 # !pip3 install --upgrade --user --quiet matplotlib seaborn
```

executed in 13ms, finished 12:49:19 2021-08-16

B [7]:

```
1 # # обновление библиотек
2 # !pip3 install --upgrade --user --quiet sklearn
```

executed in 8ms, finished 12:49:19 2021-08-16

3.1 Загрузка данных

B [8]:

```
1 def load_from_url(file_url, file_path):
2     """Загрузка файла из сетевого источника и загрузка в file_path"""
3     folder_path = os.path.dirname(file_path)
4     if not os.path.exists(folder_path):
5         print('Создаем папку {folder_path}')
6         os.makedirs(folder_path)
7         print(f'Загружаю файл из "{file_url}" в "{file_path}" ... ', end='')
8         result = urllib.urlretrieve(url=file_url, filename=file_path)
9         print('OK')
10        return result
11
12
```

executed in 9ms, finished 12:49:19 2021-08-16

B [9]:

```
1 def get_data_frame(file_path, *args, **kwargs):
2     """Загружаю файл из локального файла file_path"""
3     print(f'Открываю "{file_path}" с помощью Pandas ...', end='')
4     result = pd.read_csv(file_path, *args, **kwargs)
5     print('OK')
6     return result
7
```

executed in 12ms, finished 12:49:19 2021-08-16

B [10]:

```
1 #
2 def load_dataframe(file_url, file_path, *args, **kwargs):
3     if os.path.isfile(file_path):
4         print(f'Файл "{file_path}" уже был ранее загружен.')
5     else:
6         load_from_url(file_url, file_path)
7     # Создаем датафрейм
8     return get_data_frame(file_path, *args, **kwargs)
9
```

executed in 7ms, finished 12:49:19 2021-08-16

B [11]:

```
1 dataset_url0 = 'https://code.s3.yandex.net/datasets/geo_data_0.csv'
2 dataset_url1 = 'https://code.s3.yandex.net/datasets/geo_data_1.csv'
3 dataset_url2 = 'https://code.s3.yandex.net/datasets/geo_data_2.csv'
4
5 dataset_path0 = 'D:/projects/geo_data_0.csv'
6 dataset_path1 = 'D:/projects/geo_data_1.csv'
7 dataset_path2 = 'D:/projects/geo_data_2.csv'
```

executed in 14ms, finished 12:49:19 2021-08-16

B [12]:

```
1 raw_data0 = load_dataframe(dataset_url0, dataset_path0)
2 raw_data1 = load_dataframe(dataset_url1, dataset_path1)
3 raw_data2 = load_dataframe(dataset_url2, dataset_path2)
```

executed in 536ms, finished 12:49:19 2021-08-16

Файл "D:/projects/geo_data_0.csv" уже был ранее загружен.
Открываю "D:/projects/geo_data_0.csv" с помощью Pandas ...OK
Файл "D:/projects/geo_data_1.csv" уже был ранее загружен.
Открываю "D:/projects/geo_data_1.csv" с помощью Pandas ...OK
Файл "D:/projects/geo_data_2.csv" уже был ранее загружен.
Открываю "D:/projects/geo_data_2.csv" с помощью Pandas ...OK

3.2 Описание данных и подготовка к анализу

- id — уникальный идентификатор скважины;
- f0, f1, f2 — три признака точек;
- product — объём запасов в скважине (тыс. баррелей).

B [13]:

```
1 # Делаем функцию которая будет описывать данные для каждого датафрейма
2 def data_analysis(loc, name):
3     print(name + '\n' + '-' * 50)
4     loc.info()
5     print('-' * 50)
6     print('head()')
7     display(loc.head(3))
8     print('-' * 50)
9     print('nunique()')
10    display(loc.nunique())
11    print('-' * 50)
12    print('Полных дубликатов -', loc.duplicated().sum())
13    print('-' * 50)
14    print('describe()')
15    display(loc.describe().transpose())
16    print('-' * 50)
17    print('corr()')
18    display(loc.corr())
19    g=sns.pairplot(loc)
20    g.fig.suptitle(f'Распределения и зависимости признаков в регионе "{name}"', y=1.08)
21    print()
```

executed in 8ms, finished 12:49:19 2021-08-16

3.2.1 Регион №1

B [14]:

```
1 data_analysis(raw_data0, 'Well1')
```

executed in 6.08s, finished 12:49:25 2021-08-16

Well1

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100000 entries, 0 to 99999
Data columns (total 5 columns):
#   Column      Non-Null Count  Dtype
---  -
0   id           100000 non-null  object
1   f0           100000 non-null  float64
2   f1           100000 non-null  float64
3   f2           100000 non-null  float64
4   product      100000 non-null  float64
dtypes: float64(4), object(1)
memory usage: 3.8+ MB
```

head()

| | id | f0 | f1 | f2 | product |
|---|-------|----------|-----------|----------|------------|
| 0 | txEyH | 0.705745 | -0.497823 | 1.221170 | 105.280062 |
| 1 | 2acmU | 1.334711 | -0.340164 | 4.365080 | 73.037750 |
| 2 | 409Wp | 1.022732 | 0.151990 | 1.419926 | 85.265647 |

nunique()

```
id           99990
f0          100000
f1          100000
f2          100000
product     100000
dtype: int64
```

Полных дубликатов - 0

describe()

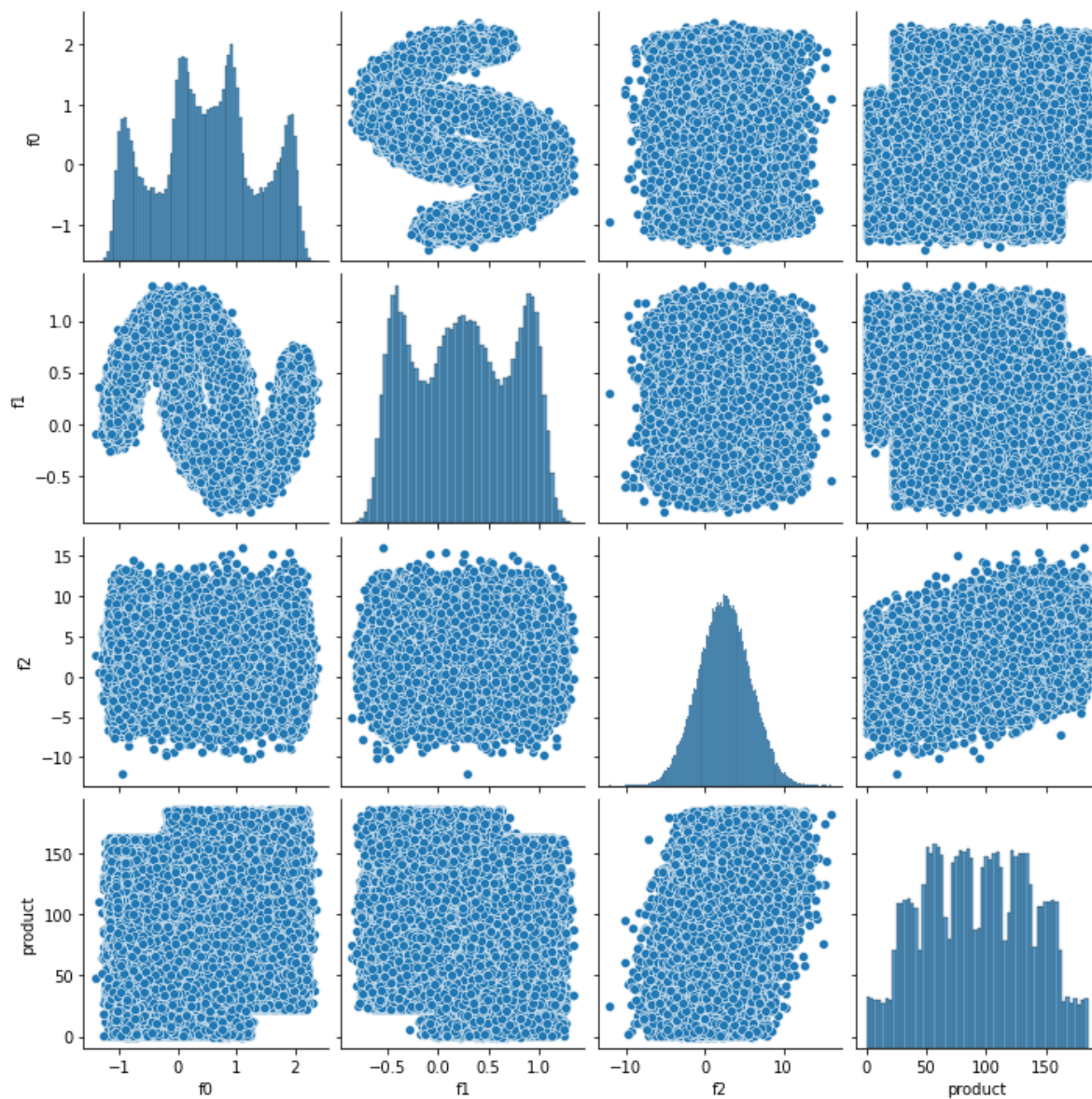
| | count | mean | std | min | 25% | 50% | 75% | n |
|---------|----------|-----------|-----------|------------|-----------|-----------|------------|----------|
| f0 | 100000.0 | 0.500419 | 0.871832 | -1.408605 | -0.072580 | 0.502360 | 1.073581 | 2.3623 |
| f1 | 100000.0 | 0.250143 | 0.504433 | -0.848218 | -0.200881 | 0.250252 | 0.700646 | 1.3437 |
| f2 | 100000.0 | 2.502647 | 3.248248 | -12.088328 | 0.287748 | 2.515969 | 4.715088 | 16.0037 |
| product | 100000.0 | 92.500000 | 44.288691 | 0.000000 | 56.497507 | 91.849972 | 128.564089 | 185.3647 |

corr()

| | f0 | f1 | f2 | product |
|----|----------|-----------|-----------|----------|
| f0 | 1.000000 | -0.440723 | -0.003153 | 0.143536 |

| | f0 | f1 | f2 | product |
|---------|-----------|-----------|----------|-----------|
| f1 | -0.440723 | 1.000000 | 0.001724 | -0.192356 |
| f2 | -0.003153 | 0.001724 | 1.000000 | 0.483663 |
| product | 0.143536 | -0.192356 | 0.483663 | 1.000000 |

Распределения и зависимости признаков в регионе "Well1"



3.2.1.1 Вывод по региону 1

Данные нормальные. Нет пропусков и дубликатов. Также не видно сильных аномалий

3.2.2 Регион №2

B [15]:

```
1 data_analysis(raw_data1, 'Well12')
```

executed in 5.65s, finished 12:49:31 2021-08-16

Well12

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 100000 entries, 0 to 99999

Data columns (total 5 columns):

| # | Column | Non-Null Count | Dtype |
|---|---------|-----------------|---------|
| 0 | id | 100000 non-null | object |
| 1 | f0 | 100000 non-null | float64 |
| 2 | f1 | 100000 non-null | float64 |
| 3 | f2 | 100000 non-null | float64 |
| 4 | product | 100000 non-null | float64 |

dtypes: float64(4), object(1)

memory usage: 3.8+ MB

head()

| | id | f0 | f1 | f2 | product |
|---|-------|------------|-----------|-----------|------------|
| 0 | kBEdx | -15.001348 | -8.276000 | -0.005876 | 3.179103 |
| 1 | 62mP7 | 14.272088 | -3.475083 | 0.999183 | 26.953261 |
| 2 | vyE1P | 6.263187 | -5.948386 | 5.001160 | 134.766305 |

nunique()

| | |
|---------|--------|
| id | 99996 |
| f0 | 100000 |
| f1 | 100000 |
| f2 | 100000 |
| product | 12 |

dtype: int64

Полных дубликатов - 0

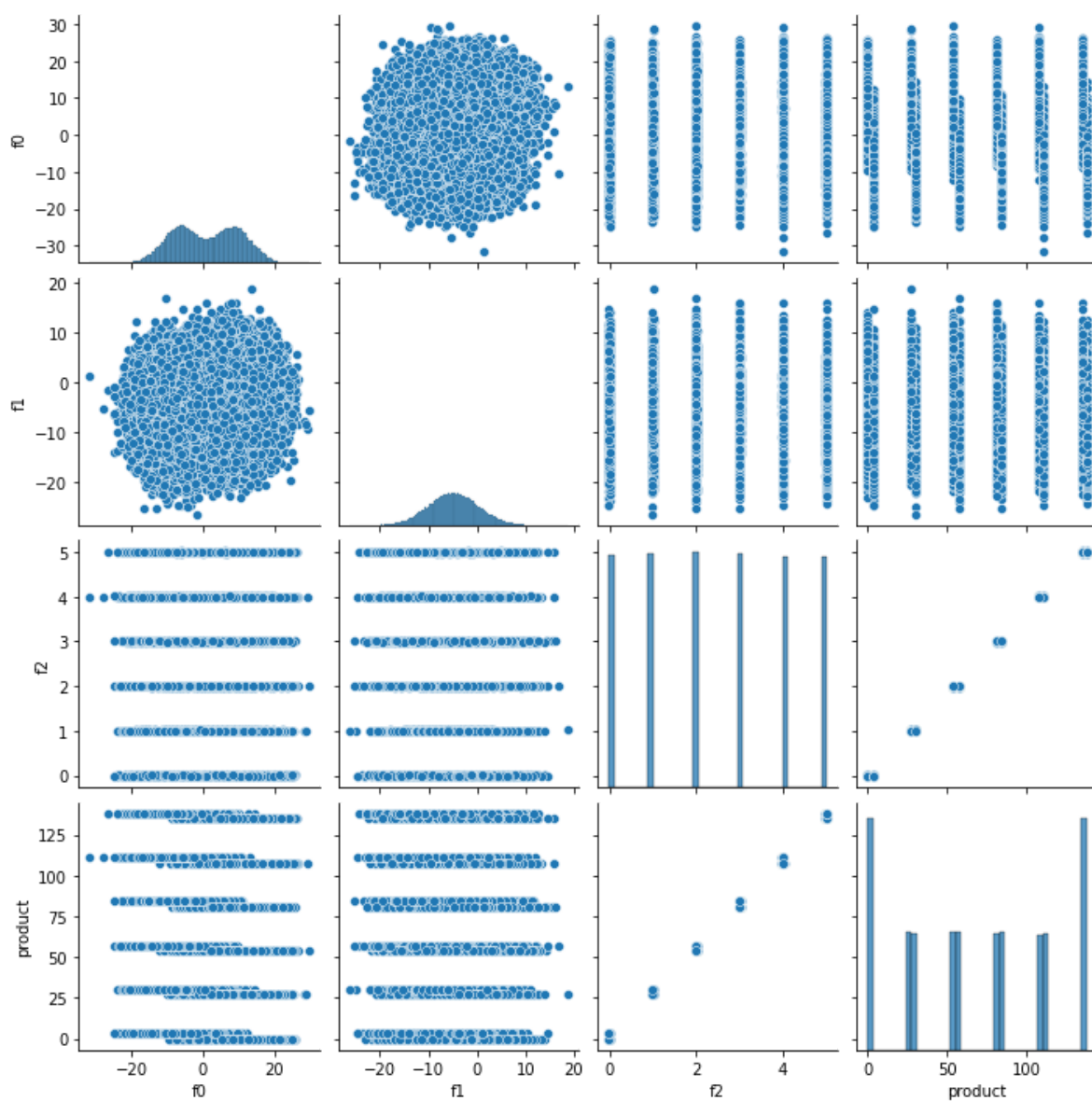
describe()

| | count | mean | std | min | 25% | 50% | 75% | n |
|---------|----------|-----------|-----------|------------|-----------|-----------|------------|----------|
| f0 | 100000.0 | 1.141296 | 8.965932 | -31.609576 | -6.298551 | 1.153055 | 8.621015 | 29.4211 |
| f1 | 100000.0 | -4.796579 | 5.119872 | -26.358598 | -8.267985 | -4.813172 | -1.332816 | 18.7340 |
| f2 | 100000.0 | 2.494541 | 1.703572 | -0.018144 | 1.000021 | 2.011479 | 3.999904 | 5.0191 |
| product | 100000.0 | 68.825000 | 45.944423 | 0.000000 | 26.953261 | 57.085625 | 107.813044 | 137.9454 |

corr()

| | f0 | f1 | f2 | product |
|----------------|-----------|-----------|-----------|-----------|
| f0 | 1.000000 | 0.182287 | -0.001777 | -0.030491 |
| f1 | 0.182287 | 1.000000 | -0.002595 | -0.010155 |
| f2 | -0.001777 | -0.002595 | 1.000000 | 0.999397 |
| product | -0.030491 | -0.010155 | 0.999397 | 1.000000 |

Распределения и зависимости признаков в регионе "Well2"



3.2.2.1 Вывод по региону 2

Данные нормальные. Нет пропусков и дубликатов. Не видно сильных аномалий.

Также есть интересная особенность параметр f_2 очень сильно коррелирует с целым признаком. Далее мы попробуем только на этом признаке обучить модель, посмотрим что из этого получится.

3.2.3 Регион №3

B [16]:

```
1 data_analysis(raw_data2, 'Well13')
```

executed in 6.57s, finished 12:49:38 2021-08-16

Well13

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100000 entries, 0 to 99999
Data columns (total 5 columns):
Column Non-Null Count Dtype
--- ---
0 id 100000 non-null object
1 f0 100000 non-null float64
2 f1 100000 non-null float64
3 f2 100000 non-null float64
4 product 100000 non-null float64
dtypes: float64(4), object(1)
memory usage: 3.8+ MB

head()

| | id | f0 | f1 | f2 | product |
|---|-------|-----------|----------|-----------|-----------|
| 0 | fwXo0 | -1.146987 | 0.963328 | -0.828965 | 27.758673 |
| 1 | WJtFt | 0.262778 | 0.269839 | -2.530187 | 56.069697 |
| 2 | ovLUW | 0.194587 | 0.289035 | -5.586433 | 62.871910 |

nunique()

id 99996
f0 100000
f1 100000
f2 100000
product 100000
dtype: int64

Полных дубликатов - 0

describe()

| | count | mean | std | min | 25% | 50% | 75% | max |
|---------|----------|-----------|-----------|------------|-----------|-----------|------------|------------|
| f0 | 100000.0 | 0.002023 | 1.732045 | -8.760004 | -1.162288 | 0.009424 | 1.158535 | 7.238200 |
| f1 | 100000.0 | -0.002081 | 1.730417 | -7.084020 | -1.174820 | -0.009482 | 1.163678 | 7.844800 |
| f2 | 100000.0 | 2.495128 | 3.473445 | -11.970335 | 0.130359 | 2.484236 | 4.858794 | 16.739400 |
| product | 100000.0 | 95.000000 | 44.749921 | 0.000000 | 59.450441 | 94.925613 | 130.595027 | 190.029800 |

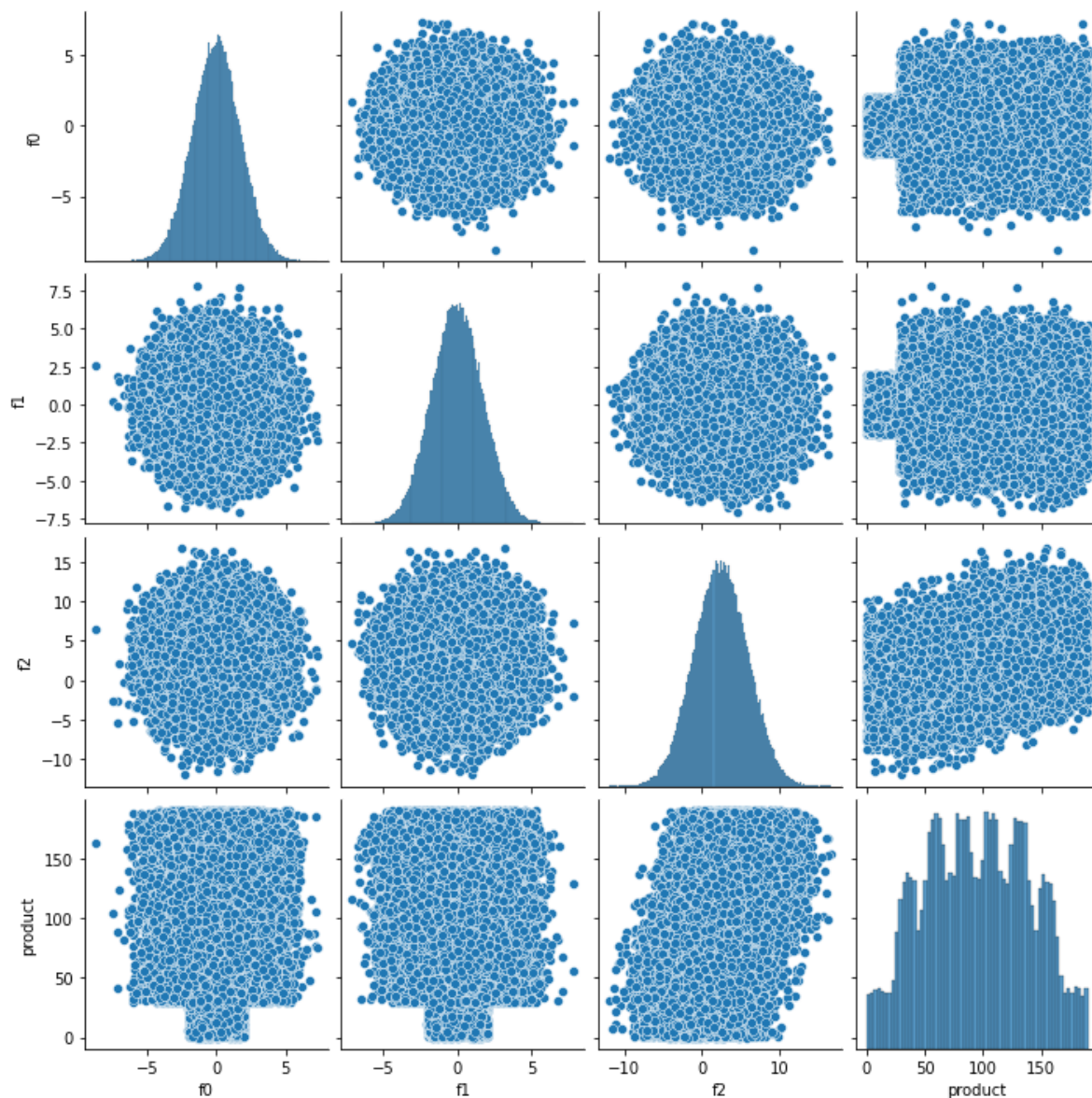
corr()

| | f0 | f1 | f2 | product |
|----|----------|----------|-----------|-----------|
| f0 | 1.000000 | 0.000528 | -0.000448 | -0.001987 |

| | f0 | f1 | f2 | product |
|----|-----------|----------|----------|-----------|
| f1 | 0.000528 | 1.000000 | 0.000779 | -0.001012 |
| f2 | -0.000448 | 0.000779 | 1.000000 | 0.445871 |



Распределения и зависимости признаков в регионе "Well3"



3.2.3.1 Вывод по региону 3

Данные нормальные. Нет пропусков и дубликатов. Не видно сильных аномалий.

4 Обучение и проверка модели

B [17]:

```
1 # Зададим значения randomState для получения одинаковых результатов
2 randomState = 42
3 state = RandomState(randomState)
```

executed in 6ms, finished 12:49:38 2021-08-16

4.1 Регион №1

B [18]:

```
1 # Выделяем отдельно признаки, отдельно целевой признак
2 # Столбец id нам также не нужен для обучения
3 features0 = raw_data0.drop(['product', 'id'], axis=1)
4 target0 = raw_data0['product']
5 # features0
```

executed in 17ms, finished 12:49:38 2021-08-16

B [19]:

```
1 # Делим данные на обучающую и валидационную выборки
2 features_train0, features_valid0, target_train0, target_valid0 = train_test_split(
3     features0, target0, test_size=0.25, random_state=randomState)
```

executed in 28ms, finished 12:49:38 2021-08-16

B [20]:

```
1 # Наша модель которая будет делать все предсказания
2 model = LinearRegression()
```

executed in 4ms, finished 12:49:38 2021-08-16

B [21]:

```
1 # Функция вывода всей нужной нам информации о модели
2 def model_specifications(features_train, features_valid, target_train, target_valid):
3     model.fit(features_train, target_train)
4     predictions_valid = model.predict(features_valid)
5     result = mean_squared_error(target_valid, predictions_valid)**0.5 # RMSE на валидац
6     print("RMSE модели линейной регрессии на валидационной выборке: {:.2f}".format(resu
7
8     raw_material_stock_pred = predictions_valid.mean()
9     raw_material_stock_real = target_valid.mean()
10    print('Средний запас предсказанного сырья - {:.2f}'.format(raw_material_stock_pred)
11    print('Средний запас реального сырья - {:.2f}'.format(raw_material_stock_real))
12    confidence_interval = t.interval(0.95, len(target_valid0), target_valid0.mean(), ta
13    print('95%-ый доверительный интервал:', confidence_interval)
14
```

executed in 13ms, finished 12:49:38 2021-08-16

В [22]:

```
1 model_specifications(features_train0, features_valid0, target_train0, target_valid0)
```

executed in 38ms, finished 12:49:38 2021-08-16

RMSE модели линейной регрессии на валидационной выборке: 37.76

Средний запас предсказанного сырья - 92.40

Средний запас реального сырья - 92.33

95%-ый доверительный интервал: (91.77707015144999, 92.8748425902369)

В [23]:

```
1 # Еще раз разбиваем наш датасет
2 features_train0, features_valid0, target_train0, target_valid0 = train_test_split(
3     features0, target0, test_size=0.25, random_state=randomState)
4
```

executed in 28ms, finished 12:49:38 2021-08-16

Попробуем нормализовать числовые данные, может это позволит уменьшить RMSE

В [24]:

```
1 features = ['f0', 'f1', 'f2']
2 scaler = StandardScaler()
3 scaler.fit(features_train0[features])
4 features_train0[features] = scaler.transform(features_train0[features])
5 # features_train0
```

executed in 32ms, finished 12:49:38 2021-08-16

В [25]:

```
1 # Сбрасываем индексы, это нам понадобится для дальнейшего бутстрэпа
2 target_valid0.reset_index(drop=True, inplace=True)
3
```

executed in 4ms, finished 12:49:38 2021-08-16

В [26]:

```
1 # И повторно запускаем нашу функцию
2 model_specifications(features_train0, features_valid0, target_train0, target_valid0)
```

executed in 33ms, finished 12:49:38 2021-08-16

RMSE модели линейной регрессии на валидационной выборке: 81.12

Средний запас предсказанного сырья - 145.56

Средний запас реального сырья - 92.33

95%-ый доверительный интервал: (91.77707015144999, 92.8748425902369)

Нормализация числовых переменных не дала никакого результата, в дальнейшем делать этого не будем

Для удобства запишем все в таблицу

B [27]:

```
1 results = pd.DataFrame(index=['RMSE на валидационной выборке', 'Средний запас предсказа  
2                               'Средний запас реального сырья'])  
3 results['Сырые данные регион 1'] = ['37.76', '92.40', '92.33']  
4 results
```

executed in 13ms, finished 12:49:38 2021-08-16

Out[27]:

| Сырые данные регион 1 | |
|------------------------------------|-------|
| RMSE на валидационной выборке | 37.76 |
| Средний запас предсказанного сырья | 92.40 |
| Средний запас реального сырья | 92.33 |

B [28]:

```
1 # confidence_interval = t.interval(0.95, len(target_valid0), target_valid0.mean(), targ  
2 # confidence_interval
```

executed in 11ms, finished 12:49:38 2021-08-16

4.2 Регион №2

B [29]:

```
1 features1 = raw_data1.drop(['product', 'id'], axis=1)  
2 target1 = raw_data1['product']  
3 # features1  
4
```

executed in 7ms, finished 12:49:38 2021-08-16

B [30]:

```
1 features_train1, features_valid1, target_train1, target_valid1 = train_test_split(  
2     features1, target1, test_size=0.25, random_state=randomState)  
3  
4 target_valid1.reset_index(drop=True, inplace=True)
```

executed in 21ms, finished 12:49:38 2021-08-16

B [31]:

```
1 model_specifications(features_train1, features_valid1, target_train1, target_valid1)
```

executed in 27ms, finished 12:49:38 2021-08-16

RMSE модели линейной регрессии на валидационной выборке: 0.89

Средний запас предсказанного сырья - 68.71

Средний запас реального сырья - 68.73

95%-ый доверительный интервал: (91.77707015144999, 92.8748425902369)

Также напомним что у нас странным образом коррелировали столбец со значением f2 и "product" — объём запасов в скважине. Попробуем убрать признаки 'f0', 'f1' и обучить модель без них.

B [32]:

```
1 features1_1 = raw_data1.drop(['product', 'id', 'f0', 'f1'], axis=1)
```

executed in 10ms, finished 12:49:38 2021-08-16

B [33]:

```
1 features_train1_1, features_valid1_1, target_train1_1, target_valid1_1 = train_test_split(
2     features1_1, target1, test_size=0.25, random_state=randomState)
3 # features_train1
4
```

executed in 22ms, finished 12:49:38 2021-08-16

B [34]:

```
1 model_specifications(features_train1_1, features_valid1_1, target_train1_1, target_valid1_1)
```

executed in 18ms, finished 12:49:38 2021-08-16

RMSE модели линейной регрессии на валидационной выборке: 1.60

Средний запас предсказанного сырья - 68.72

Средний запас реального сырья - 68.73

95%-ый доверительный интервал: (91.77707015144999, 92.8748425902369)

Результат получился плачевный, отклонение стало только больше.

Вернем признаки в датасет по первому региону т.к. на этих данных модель обучилась лучше

B [35]:

```
1 # features1 = raw_data1.drop(['product', 'id'], axis=1)
```

executed in 4ms, finished 12:49:38 2021-08-16

B [36]:

```
1 results['Сырые данные регион 2'] = ['0.89', '68.71', '68.73']
2 results
```

executed in 13ms, finished 12:49:38 2021-08-16

Out[36]:

| | Сырые данные регион 1 | Сырые данные регион 2 |
|------------------------------------|-----------------------|-----------------------|
| RMSE на валидационной выборке | 37.76 | 0.89 |
| Средний запас предсказанного сырья | 92.40 | 68.71 |
| Средний запас реального сырья | 92.33 | 68.73 |

4.3 Регион №3

B [37]:

```
1 features2 = raw_data2.drop(['product', 'id'], axis=1)
2 target2 = raw_data2['product']
3 # features1
```

executed in 8ms, finished 12:49:38 2021-08-16

B [38]:

```
1 features_train2, features_valid2, target_train2, target_valid2 = train_test_split(
2     features2, target2, test_size=0.25, random_state=randomState)
3
4 target_valid2.reset_index(drop=True, inplace=True)
```

executed in 23ms, finished 12:49:38 2021-08-16

B [39]:

```
1 model_specifications(features_train2, features_valid2, target_train2, target_valid2)
```

executed in 24ms, finished 12:49:38 2021-08-16

RMSE модели линейной регрессии на валидационной выборке: 40.15

Средний запас предсказанного сырья - 94.77

Средний запас реального сырья - 95.15

95%-ый доверительный интервал: (91.77707015144999, 92.8748425902369)

B []:

```
1
```

B [40]:

```
1 results['Сырые данные регион 3'] = ['40.15', '94.77', '95.15']
2 results
```

executed in 12ms, finished 12:49:38 2021-08-16

Out[40]:

| | Сырые данные регион 1 | Сырые данные регион 2 | Сырые данные регион 3 |
|---------------------------------------|--------------------------|--------------------------|--------------------------|
| RMSE на валидационной выборке | 37.76 | 0.89 | 40.15 |
| Средний запас предсказанного сырья | 92.40 | 68.71 | 94.77 |
| Средний запас реального сырья | 92.33 | 68.73 | 95.15 |

4.4 Вывод по блоку

Мы получили 3 основных значения:

- RMSE на валидационной выборке **это корень из среднеквадратичной ошибки наших прогнозов**
- Средний запас предсказанного сырья
- Средний запас реального сырья

Мы видим что во 2-м регионе у нас отклонение от реальных значений самое маленькое, значит тут наша модель сработала лучше всего. Также видно что средний запас предсказанного сырья не сильно отличается от реального запаса, но средние значения нам очень трудно интерпретировать. С учетом того что мы знаем RMSE и видим что в 1-м и 3-м регионе отклонение может достигать почти половины от реальных значений.

Для того чтобы определить наилучший регион для разработок, проведем исследование наших данных техникой Bootstrap.

5 3. Подготовка к расчёту прибыли

5.1 Основные данные

- Бюджет на разработку скважин в регионе — 10 млрд рублей. (budget - Development budget)
- Доход с каждой единицы продукта составляет 450 тыс. рублей (barel - Price per barrel * 1000)

В [41]:

```
1 budget = 10000000 # считаем в тысячах, поэтому получаем 10 миллионов тысяч
2 well = 200 #Количество скважин
3 budget_one = budget / well
4 barel = 450 #это в тысячах (450 000)
5
6 # Well budget - wbudget Бюджет на одну скважину (из расчета что в регионе будет всего 2
7 # quantity of oil - oilq - минимальный объём запасов в скважине для того чтобы проект 6
```

executed in 5ms, finished 12:49:38 2021-08-16

В [42]:

```
1 wbudget = budget / 200 #Бюджет на бурение одной скважины
2 oilq = wbudget / 450
3 print('Минимальный объём запасов в скважине нужен для того чтобы скважина считалась уда
4       '- {:.2f} тысяч баррелей'.format(oilq))
```

executed in 6ms, finished 12:49:38 2021-08-16

Минимальный объём запасов в скважине нужен для того чтобы скважина считалась удачной- 111.11 тысяч баррелей

5.1.1 Регион №1

В [43]:

```
1 model.fit(features_train0, target_train0)
2 predictions_valid0 = model.predict(features_valid0)
3 result0 = mean_squared_error(target_valid0, predictions_valid0)**0.5 # RMSE на валидации
4 # print("RMSE модели линейной регрессии на валидационной выборке: {:.2f}".format(result0))
5
6 raw_material_stock_pred0 = predictions_valid0.mean()
7 raw_material_stock_real0 = target_valid0.mean()
8 # print('Средний запас предсказанного сырья - {:.2f}'.format(raw_material_stock_pred0))
9 # print('Средний запас реального сырья - {:.2f}'.format(raw_material_stock_real0))
10 confidence_interval0 = t.interval(0.95, len(target_valid0), target_valid0.mean(), target_valid0.std())
11 # print('95%-ый доверительный интервал: {:.2f} : {:.2f}'.format(confidence_interval0[0], confidence_interval0[1]))
```

executed in 21ms, finished 12:49:38 2021-08-16

Для удобства заносим все в таблицу

В [44]:

```
1 results2 = pd.DataFrame(index=['RMSE модели линейной регрессии на валидационной выборке',
2                               'Средний запас реального сырья', '95%-ый доверительный интервал'],
3                           columns=['Данные предсказаний регион 1'],
4                           data=[result0, raw_material_stock_pred0, raw_material_stock_real0,
5                                '95%-ый доверительный интервал: {:.2f} : {:.2f}'.format(confidence_interval0[0], confidence_interval0[1])])
6 results2
```

executed in 13ms, finished 12:49:38 2021-08-16

Out[44]:

Данные предсказаний регион 1

| | |
|---|---------------|
| RMSE модели линейной регрессии на валидационной выборке | 81.116629 |
| Средний запас предсказанного сырья | 145.562789 |
| Средний запас реального сырья | 92.325956 |
| 95%-ый доверительный интервал | 91.78 : 92.87 |

5.1.2 Регион №2

В [45]:

```
1 model.fit(features_train1, target_train1)
2 predictions_valid1 = model.predict(features_valid1)
3 result1 = mean_squared_error(target_valid1, predictions_valid1)**0.5 # RMSE на валидации
4 # print("RMSE модели линейной регрессии на валидационной выборке: {:.2f}".format(result1))
5
6 raw_material_stock_pred1 = predictions_valid1.mean()
7 raw_material_stock_real1 = target_valid1.mean()
8 # print('Средний запас предсказанного сырья - {:.2f}'.format(raw_material_stock_pred1))
9 # print('Средний запас реального сырья - {:.2f}'.format(raw_material_stock_real1))
10 confidence_interval1 = t.interval(0.95, len(target_valid1), target_valid1.mean(), target_valid1.std())
11 # print('95%-ый доверительный интервал: ', confidence_interval1)
```

executed in 21ms, finished 12:49:38 2021-08-16

В [46]:

```
1 results2['Данные предсказаний регион 2'] = [result1, raw_material_stock_pred1,  
2 raw_material_stock_real1, '{:.2f} : {:.2f}'  
3  
4 results2
```

executed in 15ms, finished 12:49:38 2021-08-16

Out[46]:

| | Данные предсказаний регион 1 | Данные предсказаний регион 2 |
|--|---------------------------------|---------------------------------|
| RMSE модели линейной регрессии на валидационной выборке | 81.116629 | 0.89028 |
| Средний запас предсказанного сырья | 145.562789 | 68.712878 |
| Средний запас реального сырья | 92.325956 | 68.725381 |
| 95%-ый доверительный интервал | 91.78 : 92.87 | 68.16 : 69.29 |

5.1.3 Регион №3

В [47]:

```
1 model.fit(features_train2, target_train2)  
2 predictions_valid2 = model.predict(features_valid2)  
3 result2 = mean_squared_error(target_valid1, predictions_valid2)**0.5 # RMSE на валидаци  
4 # print("RMSE модели линейной регрессии на валидационной выборке: {:.2f}".format(result2)  
5  
6 raw_material_stock_pred2 = predictions_valid2.mean()  
7 raw_material_stock_real2 = target_valid2.mean()  
8 # print('Средний запас предсказанного сырья - {:.2f}'.format(raw_material_stock_pred2))  
9 # print('Средний запас реального сырья - {:.2f}'.format(raw_material_stock_real2))  
10 confidence_interval2 = t.interval(0.95, len(target_valid2), target_valid2.mean(), target  
11 # print('95%-ый доверительный интервал:', confidence_interval2)
```

executed in 28ms, finished 12:49:38 2021-08-16

В [48]:

```
1 results2['Данные предсказаний регион 3'] = [result2, raw_material_stock_pred2,  
2 raw_material_stock_real2, '{:.2f} : {:.2f}'  
3 results2
```

executed in 15ms, finished 12:49:38 2021-08-16

Out[48]:

| | Данные предсказаний регион 1 | Данные предсказаний регион 2 | Данные предсказаний регион 3 |
|--|------------------------------------|------------------------------------|------------------------------------|
| RMSE модели линейной регрессии на валидационной выборке | 81.116629 | 0.89028 | 56.447249 |
| Средний запас предсказанного сырья | 145.562789 | 68.712878 | 94.771024 |
| Средний запас реального сырья | 92.325956 | 68.725381 | 95.150999 |
| 95%-ый доверительный интервал | 91.78 : 92.87 | 68.16 : 69.29 | 94.60 : 95.71 |

Как мы видим в регионе номер 2 у нас самое маленькое RMSE, это значит что данные у нас там распределены с меньшим разбросом. При этом средний запас сырья там меньше всего.

5.2 Вывод по блоку

Здесь мы подготовили данные, а именно все ключевые значения для расчётов сохранили в отдельных переменных.

Также мы знаем, что минимальный объём запасов в скважине нужный для того чтобы скважина считалась безубыточной должен составлять - 111.11 тысяч бырелей. Здесь же мы видим что ни один из регионов по этому параметру не подходит. Но нужно не забывать что мы рассматриваем массив из 25 000 скважин. Знаит там есть как хорошие, так и плохие.

Наша основная задача теперь состоит в том чтобы через Bootstrap найти эти скважины и посчитать все риски.

6 4. Расчёт прибыли и рисков

6.1 Подготовка данных

В [49]:

```
1 # переводим объекты в Series  
2 predictions_valid0 = pd.Series(predictions_valid0)  
3 predictions_valid1 = pd.Series(predictions_valid1)  
4 predictions_valid2 = pd.Series(predictions_valid2)
```

executed in 6ms, finished 12:49:38 2021-08-16

В [50]:

```
1 # Напишем функцию расчета прибыли
2
3 def profit(pred_valid, target_valid, n):
4     predictions_valid_sorted = pred_valid.sort_values(ascending=False)
5     target_valid_selected = target_valid[predictions_valid_sorted.index][:n]
6     result = target_valid_selected.sum() * barel - budget
7     # print('Прибыль с 200 лучших скважин составила - {:.2f} тысяч рублей'.format(result))
8     return result
9
10 # profit(predictions_valid0, predictions_valid0, 200) # для проверки
```

executed in 6ms, finished 12:49:38 2021-08-16

6.1.1 Регион №1

В [51]:

```
1 # Все значения выручки из нашей выборки будут храниться в переменной values0 (для каждого региона)
2 values0 = []
3
4 # Производим бутстреп
5 for i in range(1000):
6     target_valid_sub0 = target_valid0.sample(n=500, replace=True, random_state=i)
7     pred_valid_subsample0 = predictions_valid0[target_valid_sub0.index]
8     predictions_valid_sorted0 = pred_valid_subsample0.sort_values(ascending=False)
9     values0.append(profit(predictions_valid_sorted0, target_valid_sub0, 200)) #Добавляем результат в список
10
11 values0 = pd.Series(values0)
12
13 mean0 = values0.mean()
14 print("Средняя прибыль: {:.2f}".format(mean0))
15
16 lower0 = values0.quantile(.025)
17 upper0 = values0.quantile(.975)
18
19 # # Определяем 95%-й доверительный интервал
20 # interval0 = t.interval(0.95, len(values0), values0.mean(), values0.sem())
21
22 # print(interval0)
23 print('95% доверительный интервал в регионе (тысяч руб.) - {:.2f} - {:.2f}'.format(lower0, upper0))
```

executed in 1.95s, finished 12:49:40 2021-08-16

Средняя прибыль: 330813.18

95% доверительный интервал в регионе (тысяч руб.) - -221514.71 - 851390.64

В [52]:

```
1 # Рассчитываем риск убытков
2 risk_of_loss0 = values0[values0 < 0].count()/1000
3 print('Риск убытков в регионе 1: {:.2%}'.format(risk_of_loss0))
```

executed in 6ms, finished 12:49:40 2021-08-16

Риск убытков в регионе 1: 12.10%

Для удобства запишем все в таблицу

В [53]:

```
1 results3 = pd.DataFrame(index=['Средняя прибыль в регионе после Bootstrap (тысяч руб.)',
2                               '95% доверительный интервал в регионе (тысяч руб.)',
3                               'Риск убытков в регионе в %'])
4 results3['Данные предсказаний регион 1'] = ['{:.2f}'.format(mean0),
5                                             '{:.2f} : {:.2f}'.format(lower0, upper0),
6 results3
```

executed in 12ms, finished 12:49:40 2021-08-16

Out[53]:

Данные предсказаний регион 1

| | |
|--|------------------------|
| Средняя прибыль в регионе после Bootstrap (тысяч руб.) | 330813.18 |
| 95% доверительный интервал в регионе (тысяч руб.) | -221514.71 : 851390.64 |
| Риск убытков в регионе в % | 12.10% |

6.1.2 Регион №2

В [54]:

```
1 values1 = []
2
3 for i in range(1000):
4     target_valid_sub1 = target_valid1.sample(n=500, replace=True, random_state=state)
5     pred_valid_subsample1 = predictions_valid1[target_valid_sub1.index]
6     predictions_valid_sorted1 = pred_valid_subsample1.sort_values(ascending=False)
7     values1.append(profit(predictions_valid_sorted1, target_valid_sub1, 200))
8
9 values1 = pd.Series(values1)
10
11 mean1 = values1.mean()
12 print("Средняя прибыль: {:.2f}".format(mean1))
13
14 # interval1 = t.interval(0.95, len(values1), values1.mean(), values1.sem())
15
16 # print(interval1)
17 lower1 = values1.quantile(.025)
18 upper1 = values1.quantile(.975)
19
20 print('95% доверительный интервал в регионе (тысяч руб.): {:.2f} - {:.2f}'.format(lower1, upper1))
21
```

executed in 1.98s, finished 12:49:42 2021-08-16

Средняя прибыль: 511530.22

95% доверительный интервал в регионе (тысяч руб.): 91700.56 - 921455.67

В [55]:

```
1 risk_of_loss1 = values0[values1 < 0].count()/1000
2 print('Риск убытков в регионе 1: {:.2%}'.format(risk_of_loss1))
3
```

executed in 11ms, finished 12:49:42 2021-08-16

Риск убытков в регионе 1: 0.60%

В [56]:

```
1 results3['Данные предсказаний регион 2'] = [' {:.2f}'.format(mean1),
2                                             ' {:.2f} : {:.2f}'.format(lower1, upper1),
3 results3
```

executed in 17ms, finished 12:49:42 2021-08-16

Out[56]:

| | Данные предсказаний регион 1 | Данные предсказаний регион 2 |
|---|---------------------------------|---------------------------------|
| Средняя прибыль в регионе после Bootstrap (тысяч руб.) | 330813.18 | 511530.22 |
| 95% доверительный интервал в регионе (тысяч руб.) | -221514.71 : 851390.64 | 91700.56 : 921455.67 |
| Риск убытков в регионе в % | 12.10% | 0.60% |

6.1.3 Регион №3

В [57]:

```
1 values2 = []
2
3 for i in range(1000):
4     target_valid_sub2 = target_valid2.sample(n=500, replace=True, random_state=state)
5     pred_valid_subsample2 = predictions_valid2[target_valid_sub2.index]
6     predictions_valid_sorted2 = pred_valid_subsample2.sort_values(ascending=False)
7     values2.append(profit(predictions_valid_sorted2, target_valid_sub2, 200))
8
9 values2 = pd.Series(values2)
10 lower2 = values2.quantile(0.025)
11
12 mean2 = values2.mean()
13 print("Средняя прибыль: {:.2f}".format(mean2))
14 # interval2 = t.interval(0.95, len(values2), values2.mean(), values2.sem())
15
16 # print(interval2)
17 lower2 = values2.quantile(.025)
18 upper2 = values2.quantile(.975)
19
20 print('95% доверительный интервал в регионе (тысяч руб.): {:.2f} - {:.2f}'.format(lower2, upper2))
```

executed in 1.79s, finished 12:49:44 2021-08-16

Средняя прибыль: 408545.68

95% доверительный интервал в регионе (тысяч руб.): -120624.87 - 960859.44

В [58]:

```
1 risk_of_loss2 = values0[values2 < 0].count()/1000
2 print('Риск убытков в регионе 1: {:.2%}'.format(risk_of_loss2))
```

executed in 7ms, finished 12:49:44 2021-08-16

Риск убытков в регионе 1: 7.50%

В [59]:

```
1 results3['Данные предсказаний регион 3'] = ['{:0.2f}'.format(mean2),
2                                           '{:0.2f} : {:0.2f}'.format(lower2, upper2),
3 results3
```

executed in 16ms, finished 12:49:44 2021-08-16

Out[59]:

| | Данные предсказаний регион 1 | Данные предсказаний регион 2 | Данные предсказаний регион 3 |
|---|------------------------------------|------------------------------------|------------------------------------|
| Средняя прибыль в регионе после Bootstrap (тысяч руб.) | 330813.18 | 511530.22 | 408545.68 |
| 95% доверительный интервал в регионе (тысяч руб.) | -221514.71 : 851390.64 | 91700.56 : 921455.67 | -120624.87 : 960859.44 |
| Риск убытков в регионе в % | 12.10% | 0.60% | 7.50% |

6.2 Вывод

Мы брали по 500 случайных скважин каждого региона, далее сортировали их в порядке убывания по количеству сырья, потом отбирали из них 200 лучших, считали среднее значение и записывали его в переменную. Так мы сделали 1000 раз и получили результаты:

1. Средняя прибыль в регионе номер 2 будет самой большой
2. Мы получили 95% доверительный интервал в этом интервале с 95% вероятностью будет находиться среднее значение прибыли
3. мы также посчитали % убыточных скважин.

Исходя из полученных результатов отвечаем на наши главные вопросы:

1. Построить модель для определения региона, где добыча принесёт наибольшую прибыль. - это регион номер 2
2. Проанализировать возможную прибыль и риски. - наименьшие риски убыточных скважин - так же 2-й регион

Делаем вывод что бурить нужно во 2-м регионе!