

# [учебная задача]

## Выбор локации для скважины

Предоставлены пробы нефти в трёх регионах: в каждом 10 000 месторождений, где измерили качество нефти и объём её запасов. Нужно построить модель машинного обучения, которая поможет определить регион, где добыча принесёт наибольшую прибыль.

### 1 Загрузка и подготовка данных

Ввод [1]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 from sklearn.linear_model import LinearRegression
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import mean_squared_error
7 from sklearn.metrics import r2_score
8
9 from numpy.random import RandomState
10 from scipy import stats as st
```

Ввод [2]:

```
1 df0 = pd.read_csv('geo_data_0.csv')
2 df1 = pd.read_csv('geo_data_1.csv')
3 df2 = pd.read_csv('geo_data_2.csv')
```

#### 1.1 Описание данных

id — уникальный идентификатор скважины;

f0, f1, f2 — три признака точек (неважно, что они означают, но сами признаки значимы);

product — объём запасов в скважине (тыс. баррелей).

Ввод [3]:

```
1 # Порядок действий (функции see()):
2 # Уточним размер таблицы данных
3 # Уточним количество уникальных id скважин
4 # Уточним названия столбцов
5 # Уточним типы данных
6 # Посмотрим на данные
7 # Для более полного описания "разлета" значений воспользуемся методом describe()
8
9 def see(df):
10     print(df.shape)
11     print()
12     print(pd.DataFrame(list(df['id'].unique()))[0].count(), '- уникальных ID')
13     print()
14     print(df.columns)
15     print()
16     print(df.info())
17     display(df.head())
18     display(df.describe())
19     print('-----')
20
21 def see_list(A):
22     for x in A:
23         see(x) # вызов функции rmse
```

Ввод [4]:

```
1 df = [df0, df1, df2]
2 see_list(df)
3
4 product 100000 non-null float64
```

dtypes: float64(4), object(1)  
memory usage: 3.8+ MB  
None

	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
3	KcrkZ	-13.081196	-11.506057	4.999415	137.945408
4	AHL4O	12.702195	-8.147433	5.004363	134.766305

Ввод [5]:

```
1 print('Объем запасов в источнике geo_data_0 - {:.0f}'.format(df0['product'].sum()), ' (т')
2 print('Объем запасов в источнике geo_data_1 - {:.0f}'.format(df1['product'].sum()), ' (т')
3 print('Объем запасов в источнике geo_data_2 - {:.0f}'.format(df2['product'].sum()), ' (т')
```

Объем запасов в источнике geo\_data\_0 - 9250000 (тыс. баррелей)

Объем запасов в источнике geo\_data\_1 - 6882500 (тыс. баррелей)

Объем запасов в источнике geo\_data\_2 - 9500000 (тыс. баррелей)

## Выводы:

Типы данных не требуется менять.

Переименовывать столбцы не требуется

Пропуски в данных отсутствуют

Прямое кодирование и масштабирование признаков не требуется

Категориальный признак id не несет в себе полезных значений для обучения модели, поэтому исключим id.

product = target

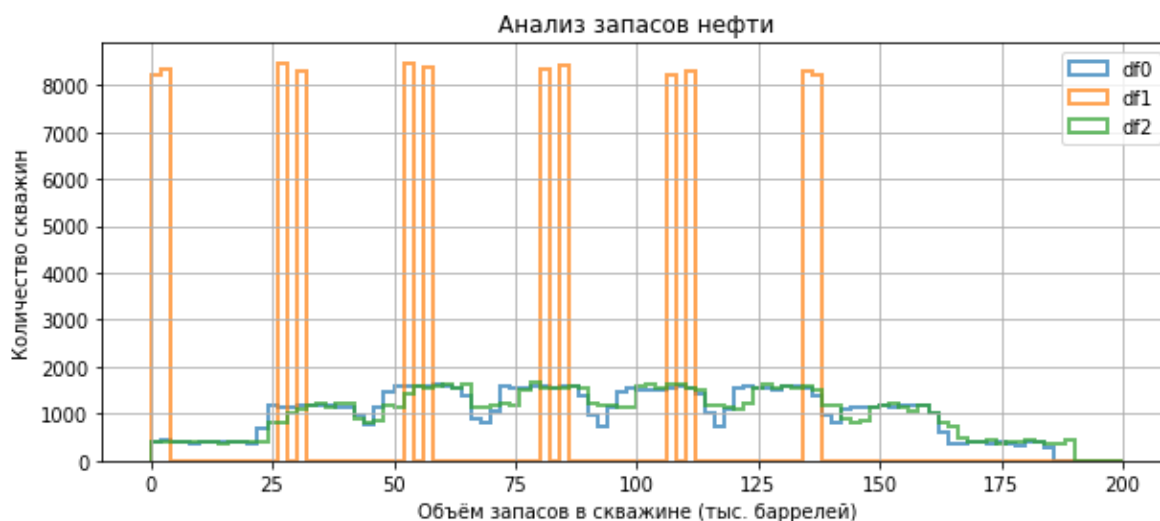
f0, f1, f2 = features

---

Данных много, поэтому нужен графический анализ данных

Ввод [6]:

```
1 ax = df0.sort_values(by='product', ascending=False).plot(
2     kind='hist',
3     y='product',
4     histtype='step',
5     range=(0, 200),
6     bins=100,
7     linewidth=2,
8     alpha=0.7,
9     label='df0',
10 )
11 df1.sort_values(by='product', ascending=False).plot(
12     kind='hist',
13     y='product',
14     histtype='step',
15     range=(0, 200),
16     bins=100,
17     linewidth=2,
18     alpha=0.7,
19     label='df1',
20     ax=ax,
21 )
22
23 df2.sort_values(by='product', ascending=False).plot(
24     kind='hist',
25     y='product',
26     histtype='step',
27     range=(0, 200),
28     bins=100,
29     linewidth=2,
30     alpha=0.7,
31     label='df2',
32     ax=ax,
33     grid=True,
34     legend=True,
35     figsize = (10,4)
36 );
37 plt.title('Анализ запасов нефти')
38 plt.xlabel('Объём запасов в скважине (тыс. баррелей)')
39 plt.ylabel('Количество скважин');
```

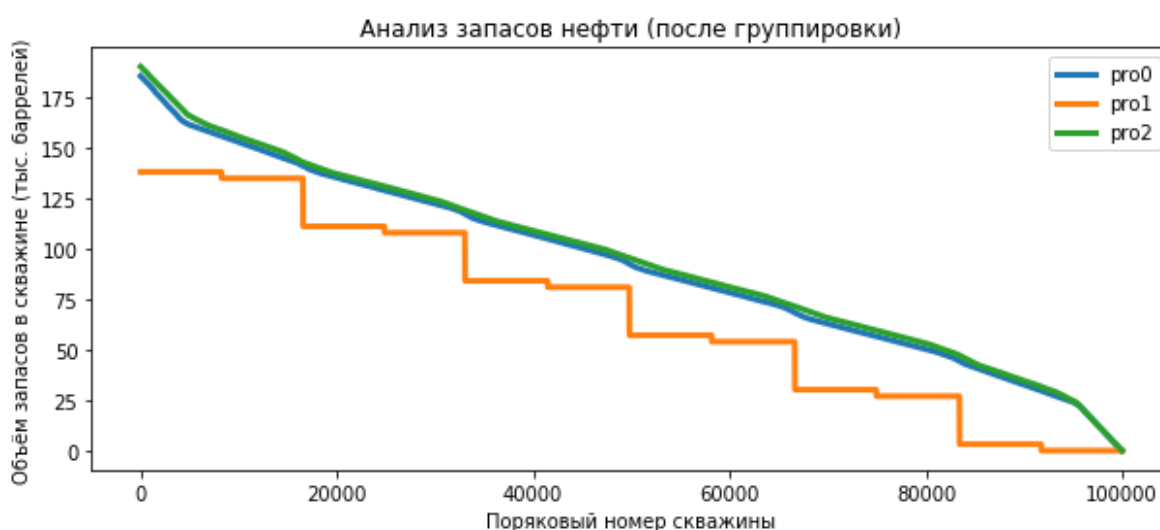


Судя по распределению значений df1 далёк от нормального распределения.

Попробуем сгруппировать значений "по убыванию" и построить график еще раз.

Ввод [7]:

```
1 pro_df = pd.DataFrame()
2 pro_df['pro0'] = df0.sort_values(by='product', ascending=False)['product'].reset_index()
3 pro_df['pro1'] = df1.sort_values(by='product', ascending=False)['product'].reset_index()
4 pro_df['pro2'] = df2.sort_values(by='product', ascending=False)['product'].reset_index()
5
6 pro_df.reset_index().plot(x='index', y=['pro0', 'pro1', 'pro2'], linewidth=3, figsize =
7 plt.title('Анализ запасов нефти (после группировки)')
8 plt.xlabel('Порядковый номер скважины')
9 plt.ylabel('Объём запасов в скважине (тыс. баррелей)');
10 plt.show()
```



Вывод:

Про качество нефти здесь нет возможности говорить, но по объему запасов лидирует источник geo\_data\_2.csv.

Источник geo\_data\_1.csv содержит данные о сильно похожих друг на друга скважинах, как-будто они искусственно созданы человеком (по ГОСТу), есть такая тенденция добывать нефть на воде и закачивать в искусственные скважины на суше (это в США, а у нас эти скважины называются - резервуары и расположены они на нефтебазах).

Высока вероятность того, что источник geo\_data\_1.csv - собственные запасы предприятия или очень "искусственно синтезированные данные".

Ввод [8]:

```
1 # функция создания данных для обучения
2 def data(df):
3     features = df.drop(['id', 'product'], axis=1)
4     target = df['product']
5     features_train, features_valid, target_train, target_valid = train_test_split(
6     features, target, test_size=0.25, random_state=12345) # 25% данные для валидации
7     return features_train, features_valid, target_train, target_valid
```

Проверим как производится разбивка данных:

Ввод [9]:

```
1 features_train, features_valid, target_train, target_valid = data(df0)
2 print('Размер features_train', features_train.shape)
3 print('Размер target_train', target_train.shape)
4 print()
5 print('Размер features_valid', features_valid.shape)
6 print('Размер target_valid', target_valid.shape)
```

Размер features\_train (75000, 3)

Размер target\_train (75000,)

Размер features\_valid (25000, 3)

Размер target\_valid (25000,)

## 2 Обучение и проверка модели

Принятая модель - линейная регрессия

Ввод [10]:

```
1 model = LinearRegression()
```

Ввод [13]:

```
1 def rmse(df): # тоже функция, но с оценкой метрик
2     f_train, f_valid, t_train, t_valid = data(df) # вызов функции data
3     model.fit(f_train, t_train) # обучение модели на тренировочной выборке
4     predictions_valid = model.predict(f_valid) # предсказания модели на валидационной выборке
5     rmse = mean_squared_error(t_valid, predictions_valid)**0.5 # RMSE на валидационной выборке
6     print('RMSE модели линейной регрессии на валидационной выборке: {:.3f}'.format(rmse))
7     print('R2 =', r2_score(t_valid, predictions_valid))
8     print('Средний запас предсказанного сырья: {:.2f}(тыс. баррелей)'.format(predictions_valid.mean()))
9     print('-----')
10
11 def rmse_list(A):
12     for x in A:
13         rmse(x) # вызов функции rmse
```

Ввод [14]:

```
1 rmse_list(df)
```

RMSE модели линейной регрессии на валидационной выборке: 37.579

R2 = 0.27994321524487786

Средний запас предсказанного сырья: 92.59(тыс. баррелей)

-----

RMSE модели линейной регрессии на валидационной выборке: 0.893

R2 = 0.9996233978805127

Средний запас предсказанного сырья: 68.73(тыс. баррелей)

-----

RMSE модели линейной регрессии на валидационной выборке: 40.030

R2 = 0.20524758386040443

Средний запас предсказанного сырья: 94.97(тыс. баррелей)

-----

Модели обучены, результаты оставляют желать лучшего (ошибка в 40 тыс.баррелей - это очень много).

### 3 Подготовка к расчёту прибыли

Перейдем к настройке алгоритма решения бизнес-задачи:

1. Синтезируем 1000 раз выборку по 500 скважин
2. Выполним сортировку по убыванию, выберем 200 лучших скважин (сверху) с большим дебетом.

Ввод [15]:

```
1 state = RandomState(12345) # настроим случайность для Bootstrap
```

### 4 Расчёт прибыли и рисков

Ввод [16]:

```
1 BUDGET = 10_000_000_000 # рублей бюджет на разработку скважин в регионе
2 ONE_POINT = 450_000 # рублей/тыс. баррелей
3 HOLE_COUNT = 200 # лучших скважин для разработки в регионе
```

Среднее значение дебета скважины для реализации безубыточного сценария (точка безубыточности):

Ввод [17]:

```
1 BREAK_EVEN_POINT = BUDGET / ONE_POINT / HOLE_COUNT # точка безубыточности
2 print('Желаемый средний объём запасов в скважине {:.1f} (тыс. баррелей)'.format(BREAK_EVEN_POINT))
3 ONE_HOLE_BUDGET = BUDGET / HOLE_COUNT # бюджет на разработку одной скважины
4 print('Бюджет на разработку одной скважины {:.0f} (рублей)'.format(ONE_HOLE_BUDGET))
```

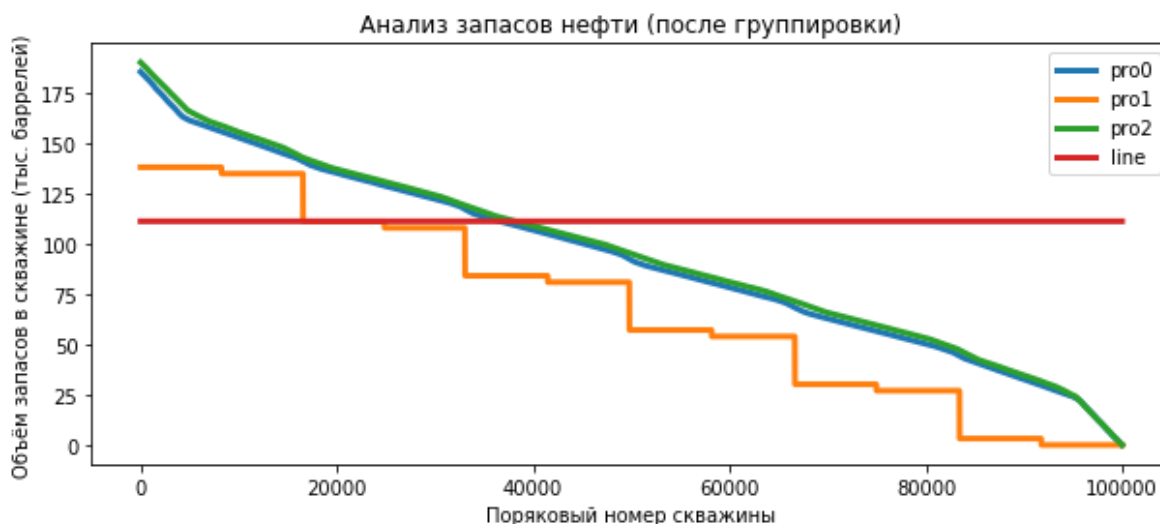
Желаемый средний объём запасов в скважине 111.1 (тыс. баррелей)

Бюджет на разработку одной скважины 50000000 (рублей)

Добавим к нашим графикам запасов по регионам границу среднего запаса (тыс. баррелей).

Ввод [18]:

```
1 pro_df['line'] = BREAKEVEN_POINT
2
3 pro_df.reset_index().plot(x='index', y=['pro0', 'pro1', 'pro2', 'line'], linewidth=3,
4 plt.title('Анализ запасов нефти (после группировки)')
5 plt.xlabel('Порядковый номер скважины')
6 plt.ylabel('Объём запасов в скважине (тыс. баррелей)');
7 plt.show()
```



Рассмотрим 200 независимых исходов событий.

Если в каждом таком независимом исходе будет полученная скважина с значением запаса 111.1 (тыс. баррелей) или выше, то следствием будет прибыль, а не убыток.

Пока оттолкнемся от этой посылки попробуем оценить вероятность прибыли (количество хороших исходов), а риск убытков это оставшее вероятностное пространство, один минус (вероятность получить прибыльный проект).

Ввод [19]:

```
1 def test_1000(df): # функция расчета прибыли 1000 раз
2     f_train, f_valid, t_train, t_valid = data(df) # вызов функции data
3     model.fit(f_train, t_train) # обучение модели на тренировочной выборке
4     predictions_valid = model.predict(f_valid)
5
6     values = [] # значение прибыли\убытка
7
8     for _ in range(1000): # создаем тестовую выборку из 200 лучших скважин
9         # создаем 500 случайных скважин
10        df_predict_500 = pd.Series(predictions_valid).sample(500, random_state=state)
11        # вытягиваем 200 индексов лучших значений
12        df_best_predict_200_index = df_predict_500.sort_values(ascending=False).head(200)
13        # сразу 200 лучших вызову по индексу (здесь не нужно 500 натуральных)
14        df_best_valid_200 = t_valid.reset_index(drop=True)[df_best_predict_200_index]
15        values.append(df_best_valid_200.sum() * 450_000 - 10_000_000_000)
16    return pd.DataFrame(values)
```

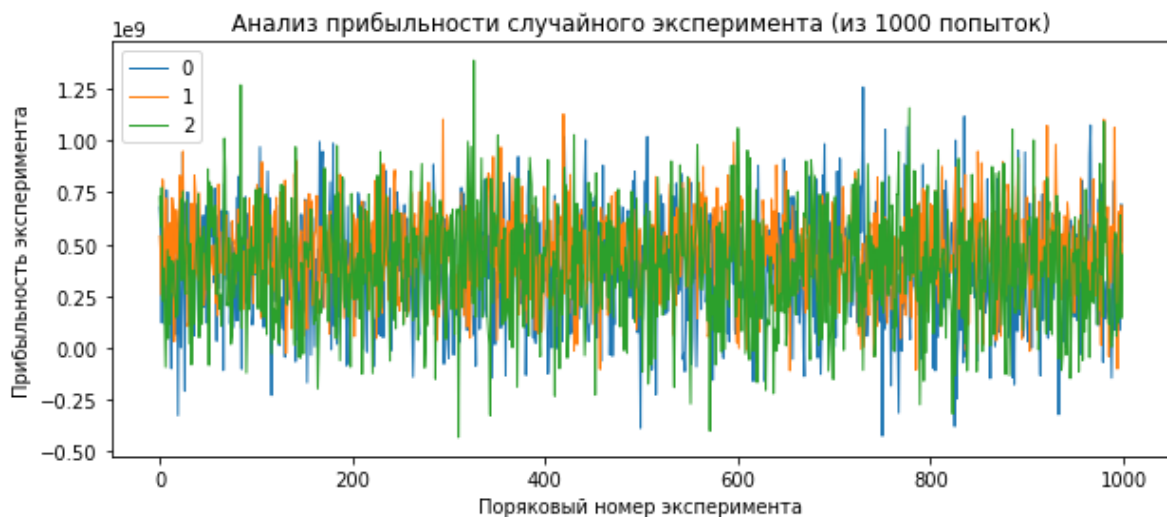


Ввод [20]:

```
1 revenue_1000_0 = test_1000(df0)
2 revenue_1000_1 = test_1000(df1)
3 revenue_1000_2 = test_1000(df2)
4
5 revenue = pd.DataFrame()
6 revenue['0'] = revenue_1000_0
7 revenue['1'] = revenue_1000_1
8 revenue['2'] = revenue_1000_2
```

Ввод [21]:

```
1 revenue.reset_index().plot(x='index', y=['0', '1', '2'], linewidth=1, figsize = (10,4))
2 plt.title('Анализ прибыльности случайного эксперимента (из 1000 попыток)')
3 plt.xlabel('Порядковый номер эксперимента')
4 plt.ylabel('Прибыльность эксперимента');
5 plt.show()
```



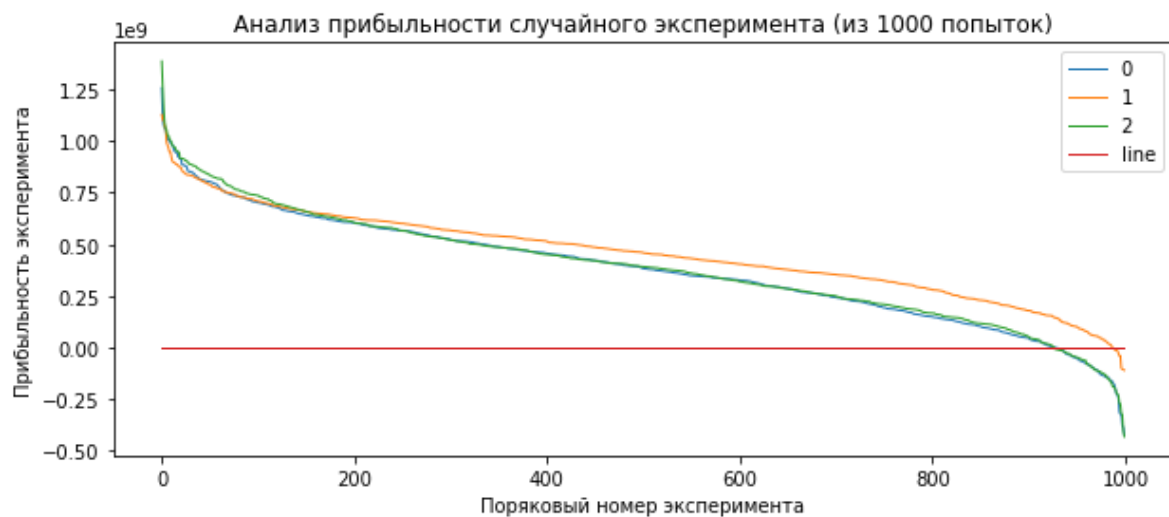
Попробуем выполнить сортировку и показать сколько было прибыльных случаев в регионе 1.

Ввод [22]:

```
1 revenue_sort = pd.DataFrame()
2 revenue_sort['0'] = revenue_1000_0.sort_values(by=0, ascending=False).reset_index(drop=
3 revenue_sort['1'] = revenue_1000_1.sort_values(by=0, ascending=False).reset_index(drop=
4 revenue_sort['2'] = revenue_1000_2.sort_values(by=0, ascending=False).reset_index(drop=
```

Ввод [23]:

```
1 revenue_sort['line'] = 0
2 revenue_sort.reset_index().plot(x='index', y=['0', '1', '2', 'line'], linewidth=1, figs
3 plt.title('Анализ прибыльности случайного эксперимента (из 1000 попыток)')
4 plt.xlabel('Порядковый номер эксперимента')
5 plt.ylabel('Прибыльность эксперимента');
6 plt.show()
```



Ввод [24]:

```
1 def interval(df):
2     revenue = df[0].mean()
3
4     if revenue > 0:
5         print('Прибыльный регион')
6         print('Средняя прибыль {:.0f} руб.'.format(revenue))
7     else:
8         print('Убыточный регион')
9         print('Убыток {:.0f} руб.'.format(revenue))
10
11     print('95%-ый доверительный интервал полученного значения: ',
12           df.sort_values(by=0, ascending=False).quantile([0.025, 0.975]).values)
13
14     p1 = df[df[0] >= 0][0].count()/df.shape[0]
15     print('Вероятность безубыточного проекта в регионе: {:.2f}'.format(p1))
16     print('Вероятность (риск) убыточного проекта в регионе: {:.2f}'.format(1 - p1))
17     print()
18     print('_____')
19     print()
20
21 def interval_list(A):
22     i = 0
23     for x in A:
24         print('Регион: ', i)
25         i += 1
26         interval(x) # вызов функции interval
```

Ввод [25]:

```
1 revenue_1000 = [revenue_1000_0, revenue_1000_1, revenue_1000_2]
2 interval_list(revenue_1000)
```

Прибыльный регион  
Средняя прибыль 454785435 руб.  
95%-ый доверительный интервал полученного значения: [[4.67300848e+07]  
[8.40213356e+08]]  
Вероятность безубыточного проекта в регионе: 0.99  
Вероятность (риск) убыточного проекта в регионе: 0.01

---

Регион: 2  
Прибыльный регион  
Средняя прибыль 389217074 руб.  
95%-ый доверительный интервал полученного значения: [[-1.15609566e+08]  
[ 9.06512590e+08]]  
Вероятность безубыточного проекта в регионе: 0.93  
Вероятность (риск) убыточного проекта в регионе: 0.07

---

**Выводы:**

Обучены 3 модели на данных из трех разных регионов.

Каждая модель "наугадывала" себе 1000 раз (исходя из значений признаков  $f_0/f_1/f_2$ ) по 200 лучших скважин методом bootstrap.

Графическим методом визуализированы прибыльные эксперименты.

Регион 1 имеет самый низкий риск убыточного проекта. К разработке месторождения рекомендован регион 1.