

# Содержание ¶

- [1 Подготовка данных](#)
- [2 Анализ данных](#)
- [3 Модель](#)
- [4 Чек-лист готовности проекта](#)

## Восстановление золота из руды (учебный проект)

Прототип модели машинного обучения для «Цифры».

Модель предсказывает коэффициент восстановления золота из золотосодержащей руды.

### Этапы:

rougher — флотация

primary\_cleaner — первичная очистка

secondary\_cleaner — вторичная очистка

final — финальные характеристики

### Типы параметров:

input — параметры сырья

output — параметры продукта

state — параметры, характеризующие текущее состояние этапа

calculation — расчётные характеристики

### Технологический процесс

Rougher feed — исходное сырьё

Rougher additions (или reagent additions) — флотационные реагенты: Xanthate, Sulphate, Depressant

Xanthate — ксантогенат (промотер, или активатор флотации);

Sulphate — сульфат (на данном производстве сульфид натрия);

Depressant — депрессант (силикат натрия).

Rougher process (англ. «грубый процесс») — флотация

Rougher tails — отвальные хвосты

Float banks — флотационная установка

Cleaner process — очистка

Rougher Au — черновой концентрат золота

Final Au — финальный концентрат золота

### Параметры этапов

air amount — объём воздуха

fluid levels — уровень жидкости

feed size — размер гранул сырья

feed rate — скорость подачи

# Подготовка данных

In [2]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np
5 import random
6
7 from sklearn.preprocessing import StandardScaler
8
9 from sklearn.metrics import mean_absolute_error, make_scorer
10 from sklearn.model_selection import cross_val_score, cross_val_predict
11 from sklearn.linear_model import LinearRegression
12 from sklearn.ensemble import RandomForestRegressor
13
```

In [3]:

```
1 train = pd.read_csv('/datasets/gold_recovery_train_new.csv')
2 test = pd.read_csv('/datasets/gold_recovery_test_new.csv')
3 full = pd.read_csv('/datasets/gold_recovery_full_new.csv')
4
```

In [4]:

```
1 print(len(train.columns), len(test.columns), len(full.columns)) #длина выборок
```

87 53 87

In [5]:

```
1 train_c = train.columns.to_list()
2 full_c = full.columns.to_list()
3 test_c = test.columns.to_list()
4 if train_c == full_c:
5     print('Столбцы обучающей и полной выборки равны') #столбцы обучающей и полной
6
```

Столбцы обучающей и полной выборки равны

In [6]:

```
1 deleted = [i for i in full_c if not i in test_c]
2 len(deleted) #признаки, отсутствующие в тестовой выборке
```

Out[6]:

34

In [7]:

```
1 nodeleted = [i for i in full_c if not i in deleted]
2 test_c == nodeleted
```

Out[7]:

True

## Тестовая выборка содержит:

Дата

Первичная очистка:

- параметры сырья: сульфаты, депрессант, ксантогенат, размер гранул
- текущее состояние этапа на 8 флотационной установке: объем воздуха и уровень жидкости (4 измерения)

Флотация:

- входящий размер 4 элементов, размер гранул, скорость подачи
- параметры сырья на 10 и 11 флотационной установке (сульфат, ксантогенат)
- текущее состояние этапа на 10 флотационной установке: объем воздуха и уровень жидкости (6 измерений)

Вторичная очистка:

- текущее состояние этапа на 2, 3, 4, 5, 6 установках: объем воздуха и уровень жидкости (2 измерения)

## Отсутствующие столбцы

Финальная:

- параметры продукта: концентрация золота, свинца, серебра, соли
- параметры продукта: эффективность восстановления
- параметры продукта: концентрация 4 элементов в хвостах

Первичная очистка:

- параметры продукта: параметры сырья: сульфаты, депрессант, ксантогенат, размер гранул сырья
- параметры продукта: концентрация 4 элементов после очистки
- параметры продукта: концентрация 4 элементов в хвостах после очистки

Флотация:

- отношение сульфата к концентрации золота
- отношение сульфата к размеру гранул золота, флотационные установки 10 и 11
- соотношение золота и свинца
- параметры продукта: размер гранул серебра, свинца, соли, золота
- параметры продукта: концентрация 4 элементов
- параметры продукта: эффективность восстановления
- параметры продукта: концентрация 4 элементов в хвостах

Вторичная очистка:

- параметры продукта: концентрация 4 элементов в хвостах

## Вывод:

В тестовой выборке отсутствуют столбцы, которые являются результатами работы этапа

**Список недоступных признаков на тесте:**

Финальный концентрат:

параметры продукта: концентрация золота, серебра, свинца, соли

параметры продукта: концентрация 4 элементов в хвостах

параметры продукта: эффективность обогащения

Первичная очистка:

параметры продукта: концентрация 4 элементов

параметры продукта: концентрация 4 элементов в хвостах

Флотация:

расчетные характеристики: концентрация сульфата по отношению к золоту

расчетные характеристики: флотационная установка 10, отношения сульфата к золоту

расчетные характеристики: флотационная установка 11, отношения сульфата к золоту

расчетные характеристики: флотационная установка 10, соотношение свинца

параметры продукта: концентрация 4 элементов

параметры продукта: концентрация 4 элементов в хвостах

параметры продукта: эффективность обогащения

Вторичная очистка:

параметры продукта: концентрация 4 элементов

параметры продукта: концентрация 4 элементов в хвостах

In [8]:

```
1 full.columns[~full.columns.isin(test.columns.to_list())]
```

Out[8]:

```
Index(['final.output.concentrate_ag', 'final.output.concentrate_pb',
      'final.output.concentrate_sol', 'final.output.concentrate_au',
      'final.output.recovery', 'final.output.tail_ag', 'final.output.
tail_pb',
      'final.output.tail_sol', 'final.output.tail_au',
      'primary_cleaner.output.concentrate_ag',
      'primary_cleaner.output.concentrate_pb',
      'primary_cleaner.output.concentrate_sol',
      'primary_cleaner.output.concentrate_au',
      'primary_cleaner.output.tail_ag', 'primary_cleaner.output.tail_
pb',
      'primary_cleaner.output.tail_sol', 'primary_cleaner.output.tail
_au',
      'rougher.calculation.sulfate_to_au_concentrate',
      'rougher.calculation.floatbank10_sulfate_to_au_feed',
      'rougher.calculation.floatbank11_sulfate_to_au_feed',
      'rougher.calculation.au_pb_ratio', 'rougher.output.concentrate_
ag',
      'rougher.output.concentrate_pb', 'rougher.output.concentrate_so
l',
      'rougher.output.concentrate_au', 'rougher.output.recovery',
      'rougher.output.tail_ag', 'rougher.output.tail_pb',
      'rougher.output.tail_sol', 'rougher.output.tail_au',
      'secondary_cleaner.output.tail_ag', 'secondary_cleaner.output.t
ail_pb',
      'secondary_cleaner.output.tail_sol',
      'secondary_cleaner.output.tail_au'],
      dtype='object')
```

In [9]:

```
1 #Recovery = (c*(f-t))/(f*(c-t))*100%<br><br>
2
3 #C – доля золота в концентрате после флотации/очистки;<br>
4 #F – доля золота в сырье/концентрате до флотации/очистки;<br>
5 #T – доля золота в отвальных хвостах после флотации/очистки.m<br>
6
7 c = full['rougher.output.concentrate_au']
8 t = full['rougher.output.tail_au']
9 f = full['rougher.input.feed_au']
10
11 full['rougher.recovery_test'] = (c*(f-t))/(f*(c-t))*100
12
13 #print(full['rougher.recovery_test'][16])
14 #print(full['rougher.output.recovery'][16])
15 print(full[['rougher.recovery_test', 'rougher.output.recovery']])
16
17
```

	rougher.recovery_test	rougher.output.recovery
0	87.107763	87.107763
1	86.843261	86.843261
2	86.842308	86.842308
3	87.226430	87.226430
4	86.688794	86.688794
...	...	...
19434	89.574376	89.574376
19435	87.724007	87.724007
19436	88.890579	88.890579
19437	89.858126	89.858126
19438	89.514960	89.514960

[19439 rows x 2 columns]

In [10]:

```
1
2 c = full['final.output.concentrate_au']
3 t = full['final.output.tail_au']
4 f = full['primary_cleaner.output.concentrate_au']
5
6 full['final.recovery_test'] = (c*(f-t))/(f*(c-t))*100
7
8 print(full[['final.recovery_test', 'final.output.recovery']].sort_values('final.
9
10
```

	final.recovery_test	final.output.recovery
14746	-0.0	0.0
14359	-0.0	0.0
14360	-0.0	0.0
10026	-0.0	0.0
7585	-0.0	0.0
...	...	...
15787	100.0	100.0
6279	100.0	100.0
6163	100.0	100.0
18037	100.0	100.0
11977	100.0	100.0

[19439 rows x 2 columns]

In [11]:

```
1 full[(full['final.recovery_test']<0)&(full['primary_cleaner.output.concentrate_au
2
3
4
5
6
```

Out[11]:

	primary_cleaner.output.concentrate_au	final.output.concentrate_au	final.output.tail_au	final.recovery_test	final.output.recovery
<b>19</b>	0.0	42.509402	2.272460	-inf	0.0
<b>14308</b>	0.0	41.347015	3.060066	-inf	0.0
<b>14113</b>	0.0	35.303158	5.471755	-inf	0.0
<b>14112</b>	0.0	43.186241	4.907862	-inf	0.0
<b>14111</b>	0.0	43.963810	4.374566	-inf	0.0
...	...	...	...	...	...
<b>8263</b>	0.0	42.691757	3.655060	-inf	0.0
<b>8262</b>	0.0	41.862333	4.064143	-inf	0.0
<b>8261</b>	0.0	43.046942	4.322626	-inf	0.0
<b>7443</b>	0.0	44.528214	3.199278	-inf	0.0

In [12]:

```
1
2
3 #Найдите MAE между вашими расчётами и значением признака.
4 mae = mean_absolute_error(full['rougher.recovery_test'], full['rougher.output.recovery'])
5 mae
```

Out[12]:

9.874045668302637e-15

В 'rougher.output.recovery', рассчитанном по формуле различия с исходными данными несущественные. Возможно, они возникли в следствии округления данных в расчетах.

'final.recovery\_test' - значения отличаются сильно и также я заметила, что если на входе была концентрация 0, то формула дает результат -inf.

Замена данных формулой приводит к сильному ухудшению качества модели. Возможно, при расчетах формула отличалась?

In [13]:

```
1 full['rougher.output.recovery'] = full['rougher.recovery_test']
2 #full.loc[full['primary_cleaner.output.concentrate_au']!=0,'final.output.recovery'] = full['rougher.output.recovery']
3
```

In [14]:

```
1 #full[full['final.output.recovery']<0]['primary_cleaner.output.concentrate_au'] = 0
```

In [15]:

```
1 print(full.shape,test.shape,train.shape) #в фул я добавляла одну колонку для теста
2
```

(19439, 89) (5290, 53) (14149, 87)

Дубликатов нет

In [16]:

```
1 for one in [full,test,train]:
2     print(one.duplicated().sum())
```

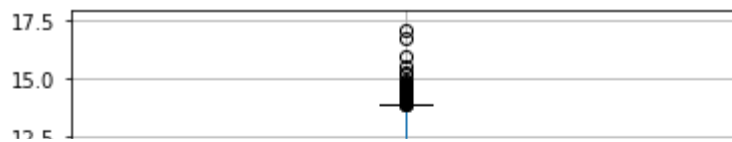
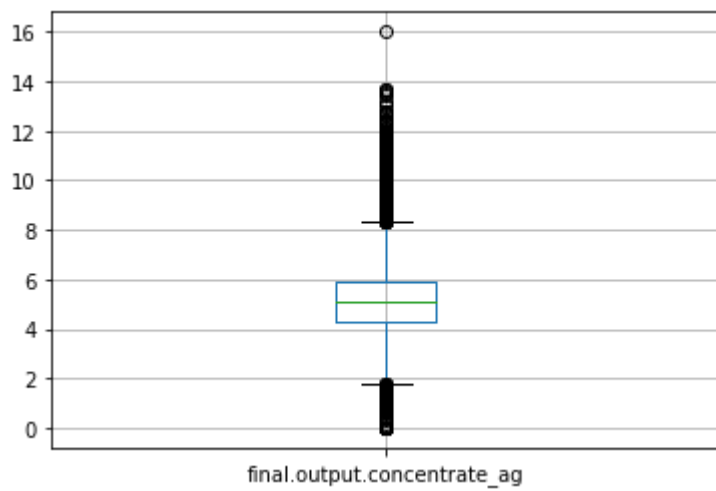
0  
0  
0

Выпады в данных есть, я не могу их удалить без консультации со специалистом



In [17]:

```
1
2 for one in full.columns.drop('date'):
3     full[[one]].boxplot()
4     plt.show()
5
```



Пропуски

параметры продукта заполню ближайшими значениями, параметры сырья - медианами



In [18]:

```
1 def passing(data):
2     #продукт
3     for elem in ['ag','sol','pb','au']:
4         for output in ['final.output.concentrate_',
5                         'final.output.tail_',
6                         'primary_cleaner.output.tail_',
7                         'primary_cleaner.output.concentrate_']:
8             if (output+elem) in data.columns:
9                 data.loc[data[output+elem].isna(), output+elem] = data[output+elem]
10
11     if 'rougher.output.concentrate_sol' in data.columns:
12         data.loc[data['rougher.output.concentrate_sol'].isna(),
13                  'rougher.output.concentrate_sol'] = data['rougher.output.concentrate_sol']
14     if 'rougher.output.tail_ag' in data.columns:
15         data.loc[data['rougher.output.tail_ag'].isna(),
16                  'rougher.output.tail_ag'] = data['rougher.output.tail_ag'].ffill()
17
18     for elem in ['ag','pb','sol']:
19         if ('secondary_cleaner.output.tail_'+elem) in data.columns:
20             data.loc[data['secondary_cleaner.output.tail_'+elem].isna(),
21                      'secondary_cleaner.output.tail_'+elem] = data['secondary_cleaner.output.tail_'+elem].ffill()
22
23     #сырье
24     for reagent in ['sulfate','depressant','xanthate']:
25         if ('primary_cleaner.input.'+reagent) in data.columns:
26             data.loc[data['primary_cleaner.input.'+reagent].isna(),
27                      'primary_cleaner.input.'+reagent] = data['primary_cleaner.input.'+reagent]
28     for reagent in ['sulfate','xanthate']:
29         for inpu in ['rougher.input.floatbank10_', 'rougher.input.floatbank11_']:
30             if (inpu+reagent) in data.columns:
31                 data.loc[data[inpu+reagent].isna(), inpu+reagent] = data[inpu+reagent]
32     for feed in ['feed_pb','feed_rate','feed_size','feed_sol']:
33         if ('rougher.input.'+feed) in data.columns:
34             data.loc[data['rougher.input.'+feed].isna(), 'rougher.input.'+feed] = data['rougher.input.'+feed]
35
36
37     #текущий этап
38     for state in ['primary_cleaner.state.floatbank8_']:
39         for letter in ['a','b','c','d']:
40             for item in ['_air','_level']:
41                 if (state+letter+item) in data.columns:
42                     data.loc[data[state+letter+item].isna(),
43                              state+letter+item] = data[state+letter+item].ffill()
44
45     for letter in ['a','b','c']:
46         for item in ['_air','_level']:
47             if ('rougher.state.floatbank10_'+letter+item) in data.columns:
48                 data.loc[data['rougher.state.floatbank10_'+letter+item].isna(),
49                          'rougher.state.floatbank10_'+letter+item] = data['rougher.state.floatbank10_'+letter+item].ffill()
50
51     if 'rougher.state.floatbank10_e_air' in data.columns:
52         data.loc[data['rougher.state.floatbank10_e_air'].isna(),
53                  'rougher.state.floatbank10_e_air'] = data['rougher.state.floatbank10_e_air'].ffill()
54
55     sec = 'secondary_cleaner.state.floatbank'
56     for fb in range(2,6):
57         fb = str(fb)
58         for letter in ['_a','_b']:
59             for item in ['air','level']:
```

```

60         if (sec+fb+letter+item) in data.columns:
61             data.loc[data[sec+fb+letter+item].isna(),
62                     sec+fb+letter+item] = data[sec+fb+letter+item].ffill
63     for item in ['air', 'level']:
64         if ('secondary_cleaner.state.floatbank6_a_'+item) in data.columns:
65             data.loc[data['secondary_cleaner.state.floatbank6_a_'+item].isna(),
66                     'secondary_cleaner.state.floatbank6_a_'+item] = data['secondary_cleaner.state.floatbank6_a_'+item].ffill
67
68
69
70     #калькуляция
71
72     for conc in ['sulfate_to_au_concentrate',
73                 'floatbank10_sulfate_to_au_feed',
74                 'floatbank11_sulfate_to_au_feed']:
75         if ('rougher.calculation.'+conc) in data.columns:
76             data.loc[data['rougher.calculation.'+conc].isna(),
77                     'rougher.calculation.'+conc] = data['rougher.calculation.'+conc].ffill
78
79     return data
80 full = passing(full)
81 train = passing(train)
82 test = passing(test)
83

```

In [19]:

```

1  #test.isna().sum()

```

## Анализ данных

Как меняется концентрация металлов (Au, Ag, Pb) на различных этапах очистки

In [20]:

```
1 concentrate_au = pd.DataFrame()
2 concentrate_ag = pd.DataFrame()
3 concentrate_pb = pd.DataFrame()
4
5 concentrate_au['input_au'] = full['rougher.input.feed_au']
6 concentrate_ag['input_ag'] = full['rougher.input.feed_ag']
7 concentrate_pb['input_pb'] = full['rougher.input.feed_pb']
8 for one in ['rougher', 'primary_cleaner', 'final']:
9     lev = '.output.concentrate_'
10
11     concentrate_au[one+'_au'] = full[one+lev+'au']
12     concentrate_ag[one+'_ag'] = full[one+lev+'ag']
13     concentrate_pb[one+'_pb'] = full[one+lev+'pb']
14
15 concentrate_au
```

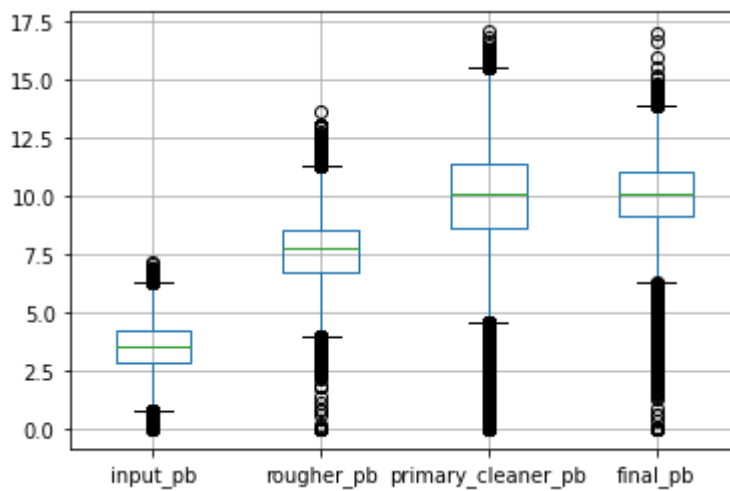
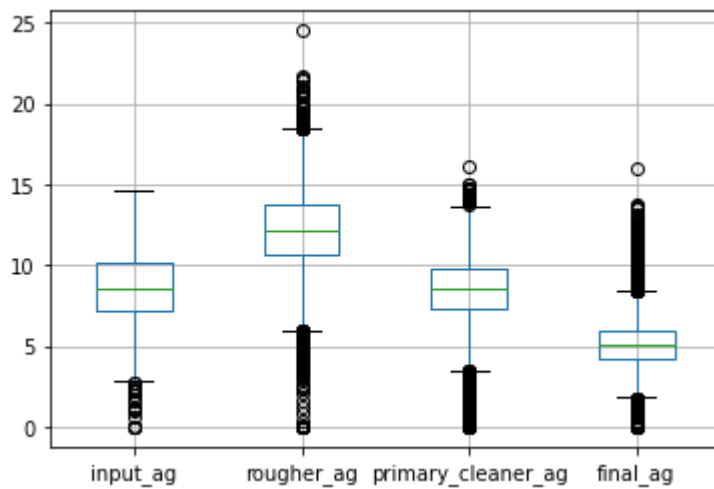
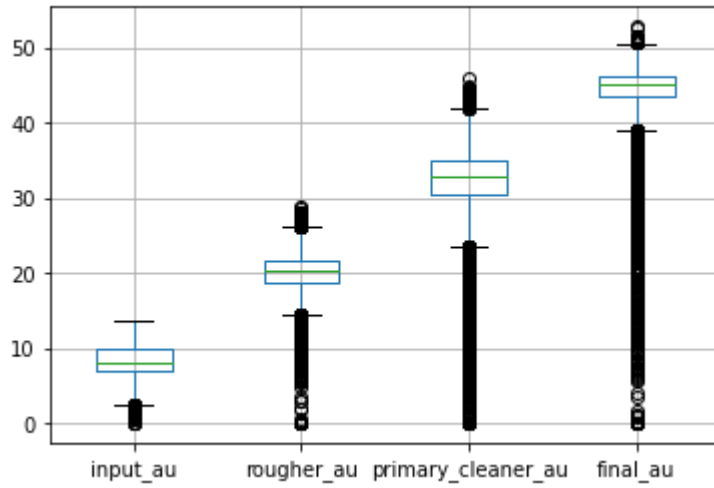
Out[20]:

	input_au	rougher_au	primary_cleaner_au	final_au
0	6.486150	19.793808	34.174427	42.192020
1	6.478583	20.050975	34.118526	42.701629
2	6.362222	19.737170	33.969464	42.657501
3	6.118189	19.320810	28.260743	42.689819
4	5.663707	19.216101	33.044932	42.774141
...	...	...	...	...
19434	5.335862	18.603550	32.940215	46.713954
19435	4.838619	18.441436	32.925325	46.866780
19436	4.525061	15.111231	31.856742	46.795691
19437	4.362781	17.834772	30.770892	46.408188
19438	4.365491	17.804134	30.356618	46.299438

19439 rows × 4 columns

In [21]:

```
1 concentrate_au.boxplot()
2 plt.show()
3 concentrate_ag.boxplot()
4 plt.show()
5 concentrate_pb.boxplot()
6 plt.show()
```

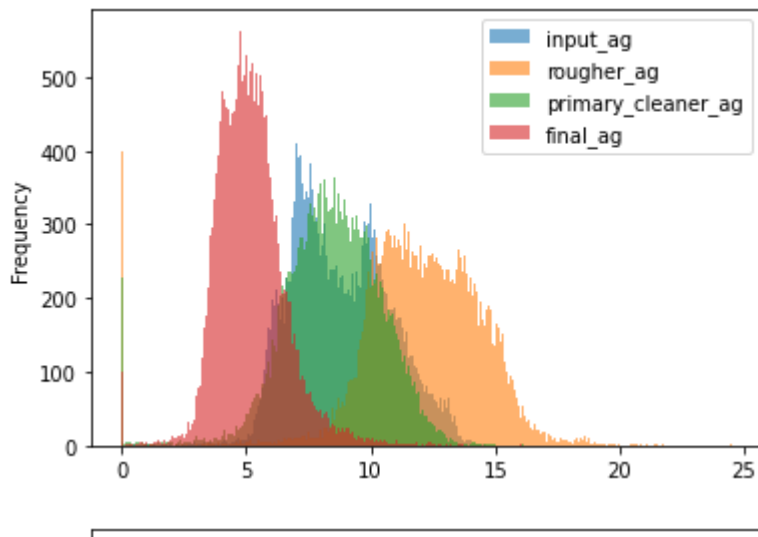


In [22]:

```
1 concentrate_ag.plot(kind = 'hist',alpha=0.6, bins=300)
2 concentrate_au.plot(kind = 'hist',alpha=0.6, bins=300)
3 concentrate_pb.plot(kind = 'hist',alpha=0.6, bins=300)
```

Out[22]:

<AxesSubplot:ylabel='Frequency'>

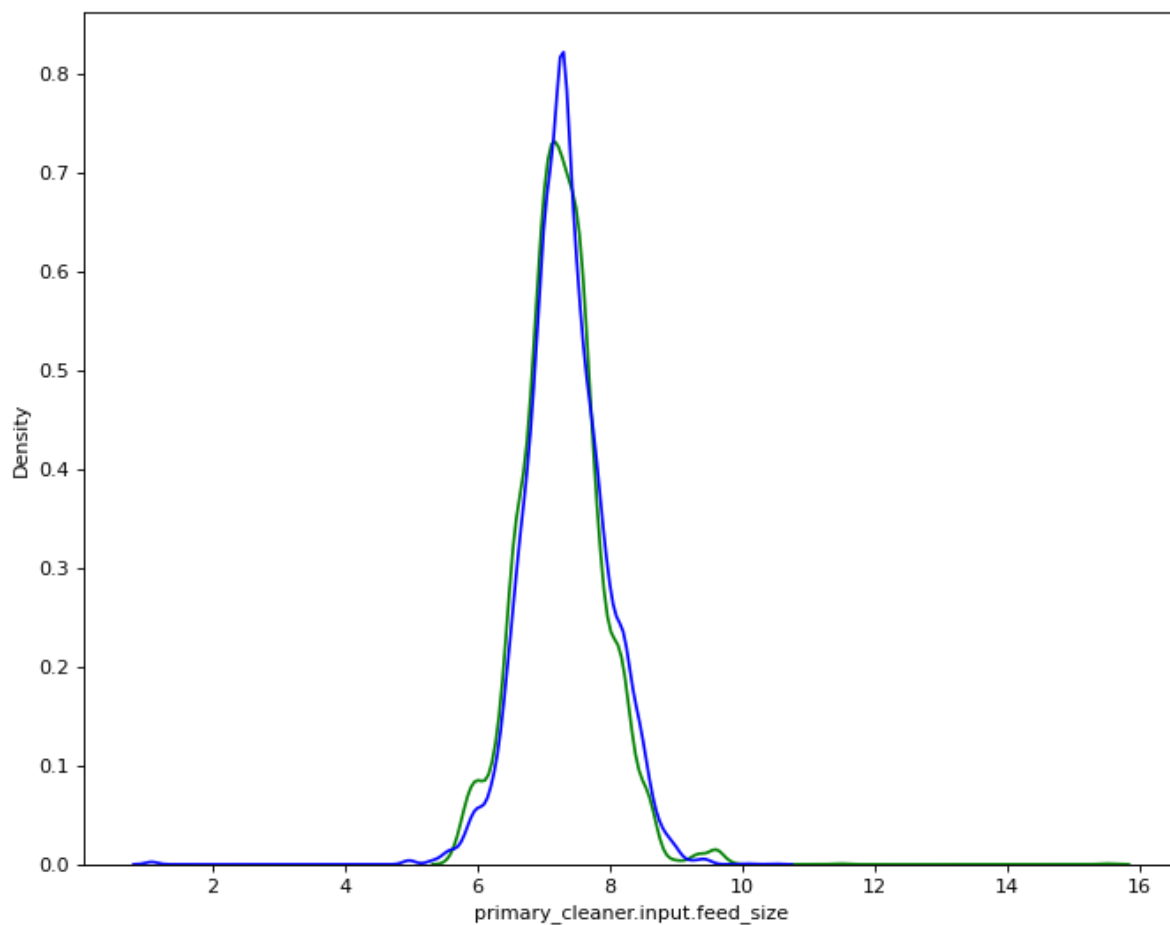


**Вывод:** концентрация золота увеличивается на каждом этапе, концентрация серебра увеличивается на флотации и уменьшается на каждом следующем этапе, концентрация свинца растет на флотации и на первом этапе очистки и не увеличивается во время второй очистки

Сравниваю распределение размеров гранул сырья на обучающей и тестовой выборках.

In [23]:

```
1 plt.figure(figsize=(10,8), dpi= 80)
2 sns.kdeplot(test['primary_cleaner.input.feed_size'],color="g", alpha=.7)
3 sns.kdeplot(train['primary_cleaner.input.feed_size'], color="b", alpha=.7)
4 plt.show()
5
6 plt.figure(figsize=(10,8), dpi= 80)
7 sns.kdeplot(test['rougher.input.feed_size'],color="g", alpha=.7)
8 sns.kdeplot(train['rougher.input.feed_size'], color="b", alpha=.7)
9 plt.show()
```





Вывод:

Суммар

концен

Density

In [ ]

1

In [24]

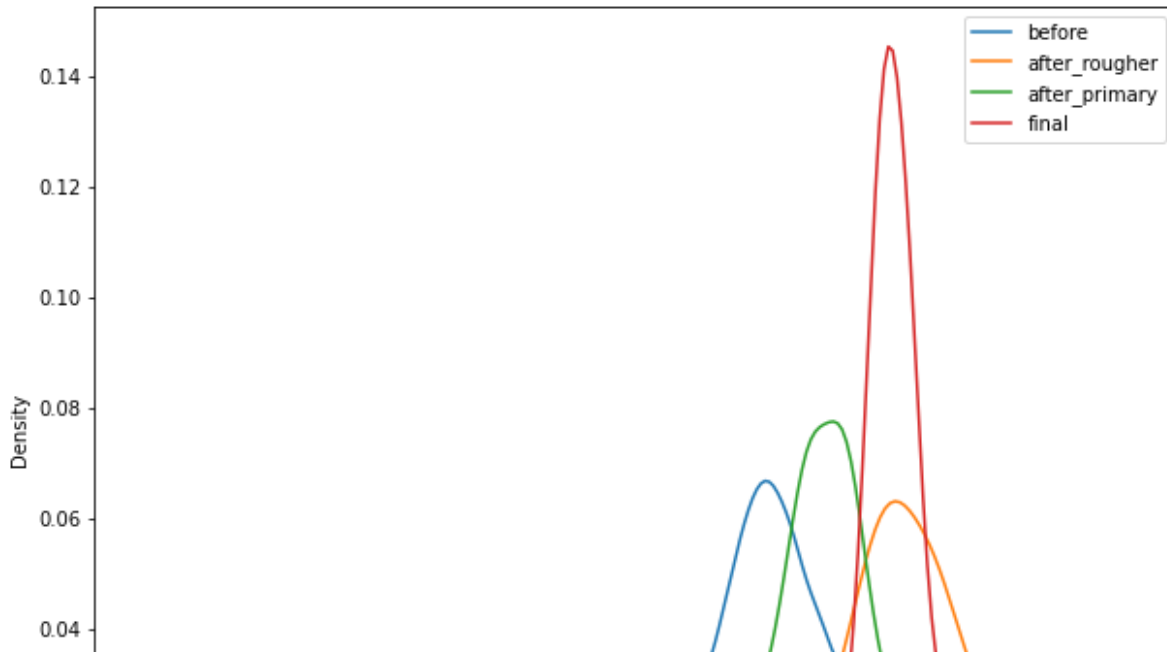
```
1 #full.columns
2 sum_conc = pd.DataFrame()
3 for one in ['rougher.input.feed_',
4             'rougher.output.concentrate_',
5             'primary_cleaner.output.concentrate_',
6             'final.output.concentrate_']:
7     sum_conc[one] = full[one+'ag'] + full[one+'au'] + full[one+'pb'] + full[one+'sol']
8
9
10 sum_conc.columns = ['rougher.input', 'rougher.output', 'primary.output', 'final.out
11 sum_conc
```

Out[24]:

	rougher.input	rougher.output	primary.output	final.output
0	51.680034	66.424950	72.640924	63.644396
1	50.659114	67.012710	72.543485	63.957723
2	50.609929	66.103793	72.095042	64.311180
3	51.061546	65.752751	59.957723	63.573449
4	47.859163	65.908382	71.321010	64.004667
...	...	...	...	...
19434	53.415050	70.781325	59.001692	68.098589
19435	53.696482	70.539603	59.703912	68.274362
19436	54.589604	55.376330	59.561096	68.226068
19437	54.027355	69.201689	57.216686	68.200449

In [25]:

```
1 plt.figure(figsize=(10,8))
2 legend = True
3 sns.kdeplot(sum_conc['rougher.input'],alpha=.7, label = 'before')
4 sns.kdeplot(sum_conc['rougher.output'], alpha=.7, label = 'after_rougher')
5 sns.kdeplot(sum_conc['primary.output'], alpha=.7, label = 'after_primary')
6 sns.kdeplot(sum_conc['final.output'], alpha=.7, label='final')
7 plt.legend()
8 plt.show()
```



In [ ]:

```
1
```

**Выводы:** Есть небольшое количество данных с практически нулевой концентрацией веществ на всех этапах. Суммарная концентрация меняется в зависимости от этапа обработки. Сначала снижается после флотации, затем возрастает после первого этапа очистки. Стремительно возрастает после второго этапа очистки

Заменяем  $X$  на  $Z$  и вычислим, чему будет равно предсказание и вектор весов.

$$\begin{aligned} a_1 &= Zw_1 \\ qquad(2.2) \\ w_1 &= (Z^T Z)^{-1} Z^T y \end{aligned}$$

Подставим уравнение 2.2 правую часть уравнения 2.3 и получим следующее:

$$a_1 = Z(Z^T Z)^{-1} Z^T y$$

Заменяем все  $Z$  правой частью уравнения 2.1:

$$\begin{aligned} a_1 &= XP((XP)^T(XP))^{-1}(XP)^T y \\ qquad(2.5) \end{aligned}$$

Для следующего шага понадобится следующее свойство обратной матрицы:

$$(AB)^{-1} = B^{-1}A^{-1}$$

Раскроем  $((XP)^T(XP))^{-1}$  в два шага:

$$a_1 = XP((XP)^T(XP))^{-1}(XP)^T y = XP(XP)^{-1}((XP)^T)^{-1}(XP)^T y = XPP^{-1}X^{-1}((XP)^T)^{-1}(XP)^T y$$

*qqquad*(2.6)

Умножение приведет к  $PP^{-1} = E$ . Для следующего этапа воспользуемся свойством транспонированной матрицы:

$$(AB)^T = B^T A^T$$

Умножение на единичную матрицу ничего не меняет. Раскроем  $((XP)^T)^{-1}(XP)^T$  в три шага:

$$a_1 = XEX^{-1}((XP)^T)^{-1}(XP)^T y = XX^{-1}(P^T X^T)^{-1}P^T X^T y = XX^{-1}(X^T)^{-1}(P^T)^{-1}P^T X^T y$$

*qqquad*(2.7)

Умножение приведет к  $(P^T)^{-1}P^T = E$ . Посмотрим, что осталось от уравнения 2.7:

$$a_1 = XX^{-1}(X^T)^{-1}EX^T y = X(X^T X)^{-1}X^T y = Xw = a$$

*qqquad*(2.8)

"Как видно, значение предсказания  $a$  не меняется, если умножать матрицу признаков на обратимую матрицу."

In [26]:

```
1 for one in ['rougher.input.feed_',
2             'rougher.output.concentrate_',
3             'primary_cleaner.output.concentrate_',
4             'final.output.concentrate_']:
5     train = train[train[one+'au'] + train[one+'ag'] + train[one+'pb'] + train[one+'sol']
6     train[train[one+'au'] + train[one+'ag'] + train[one+'pb'] + train[one+'sol'] < 1]
```

Out[26]:

date	final.output.concentrate_ag	final.output.concentrate_pb	final.output.concentrate_sol	final.o
------	-----------------------------	-----------------------------	------------------------------	---------

0 rows × 87 columns

## Модель

Метрики качества

In [27]:

```
1 #перенесла отсюда функцию поближе к коду
```

In [28]:

```
1 #train.columns
2
```

In [29]:

```
1 train.index = train.date
2 test.index = test.date
3 full.index = full.date
4
5 test_target = pd.DataFrame()
6 test_target['date'] = test['date']
7 test_target.index = test_target.date
8 test_target[['rougher.output.recovery', 'final.output.recovery']] = full[['rougher.output.recovery', 'final.output.recovery']]
9
10 train[['rougher.output.recovery', 'final.output.recovery']] = full[['rougher.output.recovery', 'final.output.recovery']]
11 train_target = train[['rougher.output.recovery', 'final.output.recovery']]
12 train = train[test.columns]
13
14
15 train.shape, train_target.shape, test.shape, test_target.shape
```

Out[29]:

```
((13599, 53), (13599, 2), (5290, 53), (5290, 3))
```

In [30]:

```
1 test_target[test_target['final.output.recovery'].isna()]=0
2 test_target[test_target['final.output.recovery'].isna()]
```

Out[30]:

```
      date  rougher.output.recovery  final.output.recovery
date
```

---

удалю из обучающей выборки нули в целевом признаке (это ухудшило качество модели)

In [31]:

```
1 #train_target = train_target[(train_target['rougher.output.recovery']!=0)&(train_target['final.output.recovery']!=0)]
2 #train_target = train_target.dropna()
3 #
4 #train = train[train.index.isin(train_target.index)]
5 #train.shape, train_target.shape
```

In [32]:

```
1 train_target[train_target['final.output.recovery'].isna()]
```

Out[32]:

```
      rougher.output.recovery  final.output.recovery
date
```

---

Посмотрю на корреляцию признаков

In [33]:

```
1
2 #pd.set_option('display.max_rows',None)
3 #pd.set_option('display.max_columns',None)
4 full.corr()
5 #primary_cleaner.output.concentrate_ag
6
7
```

Out [33]:

	final.output.concentrate_ag	final.output.concentrate_pb
final.output.concentrate_ag	1.000000	0.063618
final.output.concentrate_pb	0.063618	1.000000
final.output.concentrate_sol	0.345467	-0.041863
final.output.concentrate_au	-0.038185	0.374439
final.output.recovery	0.185964	0.151998
...	...	...
secondary_cleaner.state.floatbank5_b_level	0.156545	-0.059971
secondary_cleaner.state.floatbank6_a_air	0.162637	-0.042138
secondary_cleaner.state.floatbank6_a_level	0.059583	0.067421
rougher.recovery_test	0.084139	0.039727
final.recovery_test	-0.000850	0.037343

88 rows × 88 columns

Мне кажется, нет признаков, которые бы сильно влияли на целевые, а те, которые могли бы влиять, уже удалены в тестовой выборке

In [34]:

```
1 train = train.drop('date', axis=1)
2 test = test.drop('date', axis=1)
3 test_target = test_target.drop('date', axis=1)
4
```

3.2. Обучите разные модели и оцените их качество кросс-валидацией. Выберите лучшую модель и проверьте её на тестовой выборке. Опишите выводы.

In [35]:

```
1 def calc_smape(target, prediction):
2     #функция принимает на вход целевой признак и предсказание и вычисляет smape
3     #симметричное среднее абсолютное процентное отклонение
4     smape = 0
5     for i in range(len(target)):
6         smape += (1/len(target)) * (abs(target[i]-prediction[i])) / ((abs(target
7     return smape
8
9 def calc_itog_smape(target, prediction):
10    target= np.array(target)
11    #функция вычисляет итоговую метрику, принимая на вход smape двух величин
12    smape_one = calc_smape(target[:,0],prediction[:,0])
13    smape_two = calc_smape(target[:,1], prediction[:,1])
14    smape_itog = 0.25*smape_one+0.75*smape_two
15
16    return smape_itog
17
18
```

## Константная модель

In [36]:

```
1 constant_predict = pd.DataFrame(index = train_target.index, columns = ['rougher.
2 constant_predict['rougher.output.recovery'] = constant_predict['rougher.output.r
3 constant_predict['final.output.recovery'] = constant_predict['final.output.recov
4 constant_predict = constant_predict.to_numpy()
5 constant_smape = calc_itog_smape(train_target, constant_predict)
6 constant_smape
```

Out[36]:

9.339089690445489

## Случайный лес

In [ ]:

```
1 В
2
```

In [37]:

```
1 best_est = 0
2 best_depth = 0
3 best_cross = -10000
4 best_split = 0
5 best_leaf = 0
6 prediction = pd.DataFrame()
7 for est in range(25,26): #подобрала лучшие гиперпарамы
8     for depth in range(3,4):
9         for split in range(1,2):
10             for leaf in range(136,137):
11                 model_f = RandomForestRegressor(n_estimators=est,
12                                                  max_depth=depth,
13                                                  random_state=12345,
14                                                  min_samples_split=split/10,
15                                                  min_samples_leaf=leaf)
16
17                 scor = make_scorer(calc_itog_smape, greater_is_better=False)
18                 cross = cross_val_score(model_f, train, train_target, cv=5, scor
19                 if cross.mean()>best_cross:
20                     best_est = est
21                     best_depth = depth
22                     best_cross = cross.mean()
23                     best_split = split
24                     best_leaf = leaf
25                 #print(cross.mean(), est, depth, split, leaf)
26 print(best_cross*(-1), 'Лучшие гиперпараметры:', best_est, best_depth, best_sp1
27
28
29
```

8.662527213987541 Лучшие гиперпараметры: 25 3 1 136

Стандартизация - не влияет на результат

In [38]:

```
1 pd.options.mode.chained_assignment = None
2 scaler = StandardScaler()
3 scaler.fit(train)
4 train_scaled = scaler.transform(train)
5 test_scaled = scaler.transform(test)
6
```

In [46]:

```
1 prediction_1 = pd.DataFrame()
2 model_1 = LinearRegression()
3 model_1.fit(train, train_target)
4 scor = make_scorer(calc_itog_smape, greater_is_better=False)
5 cross_1 = cross_val_score(model_1, train, train_target, cv=5, scoring = scor)
6
7 print(cross_1.mean(), cross_1)
8
```

-10.014241894204414 [-12.03005649 -8.86985447 -9.30815727 -7.843783  
54 -12.01935769]

In [45]:

```
1 model_f.fit(train,train_target)
2 predict_test_f = model_f.predict(test)
3 smape_test_f = calc_itog_smape(test_target, predict_test_f)
4 smape_test_f
```

Out[45]:

9.32012722153264

Лучшая модель - случайный лес. SMAPE на тестовой выборке - 9.32