

Тестовое задание.

Определить абонентов (msisdn), которые являются одной и той же Персоной (человеком).

In [1]:

```
1 import pandas as pd
2 import os
3 from os import walk
4 import datetime
5 pd.set_option('display.max_columns', 500)
6
7 import numpy as np
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10
11 from datetime import datetime, date
12
13 import csv
14
15 import folium
16 from folium import plugins
17 from folium.plugins import HeatMap
18 import branca
19 import branca.colormap as cm
```

Загружаем данные и преобразуем в удобный формат некоторые поля

In [2]:

```
1 current_directory = r'C:\Users\ekaterina.adischeva\Documents\Scripts\BD_task\Техническ
2 current_directory = current_directory.replace('\\', '/')
3 os.chdir(current_directory)
```

In [3]:

```
1 df_data = pd.read_csv('02_Data.csv', sep = ';', dtype = {'imei': str, 'lac' : 'str', 'c
2 facts = pd.read_excel('01_Факты.xlsx', header = None)
3 tac_dict = pd.read_csv('03_устройства.csv', quoting=csv.QUOTE_NONE)
4 event_type = pd.read_excel('04_event_type.xlsx')
```

In [4]:

```
1 tac_dict.columns = tac_dict.columns.str.replace('\\"', '')
2 for column_name in tac_dict.columns:
3     tac_dict[column_name] = tac_dict[column_name].str.replace('\\"', '')
```

Преобразуем в datetime время регистрации события на базовой станции (количество 1/1000 секунды, прошедших с 01.01.1970 до времени регистрации на БС.)

Выделяем tac из поля imei

In [5]:

```

1 df_data['datetime'] = pd.to_timedelta(df_data['tstamp']/1000, unit='S') + pd.Timestamp
2 df_data['datetime'] = pd.to_datetime(df_data['datetime'])
3
4 df_data['tac'] = df_data['imei'].str.slice(0,8)

```

Подгружаем значения из справочников в данные о перемещении абонентов

In [6]:

```

1 df_data = df_data.merge(event_type, left_on = ['event_type'], right_on = 'Homep', how =
2 df_data = df_data.merge(tac_dict, on = ['tac'], how = 'left')

```

Графики и функции

show_chart отображает изменение координат широты и долготы по времени

In [7]:

```

1 def show_chart(fig, data_a, data_b ):
2     """
3     The function draws 2 charts: msisdn lattitude on time and msisdn longitude on time
4
5     """
6
7     data_a_copy = data_a
8     data_b_copy = data_b
9
10    plt.subplot(1, 2, 1)
11    feature = 'long'
12    plt.scatter(data_a_copy['datetime'], data_a_copy[feature], alpha = .5, c = 'blue')
13    plt.plot(data_a_copy['datetime'], data_a_copy[feature+'_max'],'-o',alpha = .3, color = 'blue')
14    plt.plot(data_a_copy['datetime'], data_a_copy[feature+'_min'],'-o',alpha = .3, color = 'blue')
15    plt.scatter(data_b_copy['datetime'], data_b_copy[feature], alpha = .5, c = 'red')
16    plt.plot(data_b_copy['datetime'], data_b_copy[feature+'_max'],'-o',alpha = .3, color = 'red')
17    plt.plot(data_b_copy['datetime'], data_b_copy[feature+'_min'],'-o',alpha = .3, color = 'red')
18    plt.xlim(pd.concat([data_a,data_b])['datetime'].min(), pd.concat([data_a,data_b])['datetime'].max())
19    plt.xticks(rotation=70)
20
21    plt.subplot(1, 2, 2)
22    feature = 'lat'
23    plt.scatter(data_a_copy['datetime'], data_a_copy[feature], alpha = .5, c = 'blue')
24    plt.plot(data_a_copy['datetime'], data_a_copy[feature+'_max'],'-o',alpha = .3, color = 'blue')
25    plt.plot(data_a_copy['datetime'], data_a_copy[feature+'_min'],'-o',alpha = .3, color = 'blue')
26    plt.scatter(data_b_copy['datetime'], data_b_copy[feature], alpha = .5, c = 'red')
27    plt.plot(data_b_copy['datetime'], data_b_copy[feature+'_max'],'-o',alpha = .3, color = 'red')
28    plt.plot(data_b_copy['datetime'], data_b_copy[feature+'_min'],'-o',alpha = .3, color = 'red')
29    plt.xlim(pd.concat([data_a,data_b])['datetime'].min(), pd.concat([data_a,data_b])['datetime'].max())
30    plt.xticks(rotation=70)
31
32    plt.show()
33    return fig

```

show_circles_on_map отображает точки в которых был абонент на карте

In [8]:

```

1 def show_circles_on_map(m, data, latitude_column, longitude_column, color):
2     """
3     The function draws map with circles on it.
4     The center of the map is the mean of coordinates passed in data.
5
6     data: DataFrame that contains columns latitude_column and longitude_column
7     latitude_column: string, the name of column for latitude coordinates
8     longitude_column: string, the name of column for longitude coordinates
9     color: string, the color of circles to be drawn
10    """
11
12    location = (data[latitude_column].mean(), data[longitude_column].mean())
13
14    for _, row in data.iterrows():
15        folium.Circle(
16            radius=100,
17            location=(row[latitude_column], row[longitude_column]),
18            color=color,
19            fill_color=color,
20            fill=True,
21            fill_opacity = .5
22        ).add_to(m)
23
24    return m

```

haversine_array определяет расстояние в метрах по координатам

In [9]:

```

1 def haversine_array(lat1, lng1, lat2, lng2):
2     lat1, lng1, lat2, lng2 = map(np.radians, (lat1, lng1, lat2, lng2))
3     AVG_EARTH_RADIUS = 6371 # in km
4     lat = lat2 - lat1
5     lng = lng2 - lng1
6     d = np.sin(lat * 0.5) ** 2 + np.cos(lat1) * np.cos(lat2) * np.sin(lng * 0.5) ** 2
7     h = 2 * AVG_EARTH_RADIUS * np.arcsin(np.sqrt(d))
8     return h

```

Проверяем данные

Смотрим, где в выборке есть пустые данные

In [10]:

```
1 df_data.count() - df_data.shape[0], df_data.dtypes
```

Out[10]:

```
(lac          0
 cid          0
 msisdn       0
 imei        -11588
 event_type   0
 tstamp       0
 long         0
 lat          0
 max_dist     0
 cell_type    0
 start_angle  0
 end_angle    0
 datetime     0
 tac         -11588
 Номер        0
 Обозначение  0
 Описание     0
 vendor      -23679
 platform    -23679
 type        -23679
 dtype: int64, lac          object
 cid          object
 msisdn       int64
 imei        object
 event_type   int64
 tstamp       int64
 long        float64
 lat         float64
 max_dist     int64
 cell_type    object
 start_angle  float64
 end_angle    float64
 datetime     datetime64[ns]
 tac          object
 Номер        int64
 Обозначение  object
 Описание     object
 vendor       object
 platform     object
 type         object
 dtype: object)
```

Некоторые номера, принадлежащие Персонам, не попали в выборку с детальными данными

In [11]:

```

1 # for 2 persons we don't have second msisdn in our dataset
2 display(facts[~facts[1].isin(df_data['msisdn'])],
3         facts[~facts[0].isin(df_data['msisdn'])])

```

	0	1
1	158510912201	158528852857
17	158528850493	158530004641

0	1
---	---

Смотрим распределение событий по времени. В выборке есть данные, датированные позже 26 мая. Удаляем данные с датой более 26 мая.

In [12]:

```

1 plt.hist(df_data['datetime'].values, bins=100)
2 plt.xticks(rotation=70)
3 plt.show()

```

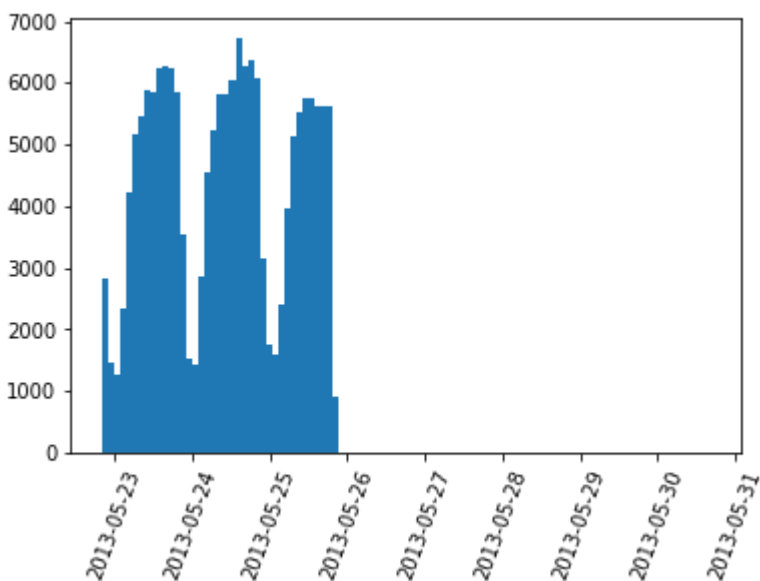
C:\Users\ekaterina.adischeva\AppData\Local\Continuum\anaconda3\lib\site-packages\pandas\plotting_converter.py:129: FutureWarning: Using an implicitly registered datetime converter for a matplotlib plotting method. The converter was registered by pandas on import. Future versions of pandas will require you to explicitly register matplotlib converters.

To register the converters:

```

>>> from pandas.plotting import register_matplotlib_converters
>>> register_matplotlib_converters()
warnings.warn(msg, FutureWarning)

```



In [13]:

```
1 df_data.loc[:, 'date'] = df_data['datetime'].dt.date
2 print(df_data['date'].value_counts())
3 df_data = df_data[(df_data['datetime'] < datetime.strptime('2013-05-26', '%Y-%m-%d'))
4                   # & (df_data['datetime'] >= datetime.strptime('2013-05-23', '%Y-%m-%d'))
5                   ]
```

```
2013-05-24    61642
2013-05-23    59463
2013-05-25    48519
2013-05-22     4450
2013-05-29         24
2013-05-28         18
2013-05-30          7
Name: date, dtype: int64
```

Проверяем базовые станции.

- Расположение базовых станций на карте.
- Расположение ширины угла.
- Расположение зоны действия в метрах.

In [14]:

```
1 CENTER_LAT, CENTER_LONG = 55.7522200, 37.6155600
2 m = folium.Map(location=(CENTER_LAT, CENTER_LONG))
3 show_circles_on_map(m, df_data.sample(500), "lat", "long", "blue")
```

Out[14]:

In [15]:

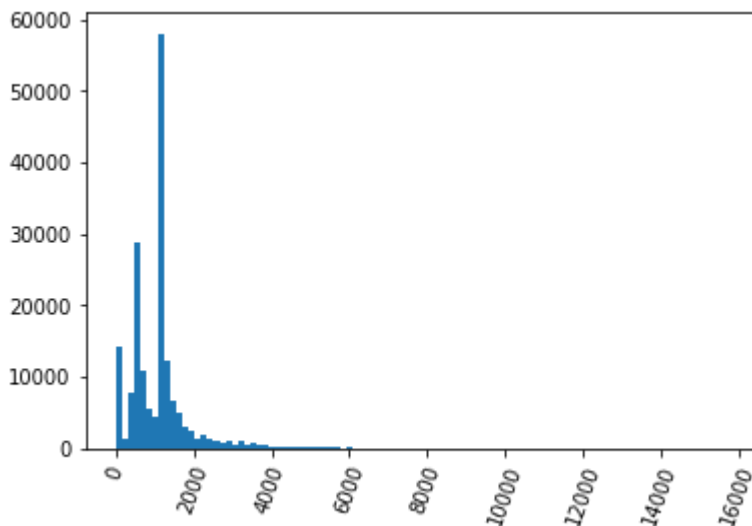
```
1 df_data['angel'] = np.where(df_data['end_angle'] < df_data['start_angle'], df_data['end_angle'] - df_data['start_angle'])
2 df_data['end_angle'] - df_data['start_angle'])
3
4 df_data['angel'].value_counts()
```

Out[15]:

```
60.0    156940
47.0     9719
90.0     3104
360.0    2176
65.0     2135
Name: angel, dtype: int64
```

In [16]:

```
1 plt.hist(df_data['max_dist'].values, bins=100)
2 plt.xticks(rotation=70)
3 plt.show()
```



Указываю координаты зоны действия БС для каждого события

По координатам БС, максимальной дистанции приема и направления и угла покрытия определяем координаты прямоугольника, в которых мог находиться абонент при регистрации на БС

In [17]:

```

1 def BS_rectangle_coordinate(df_):
2     df = df_.copy()
3     lat_length = haversine_array(CENTER_LAT -.5, CENTER_LONG, CENTER_LAT + .5, CENTER_
4     long_length = haversine_array(CENTER_LAT, CENTER_LONG -.5, CENTER_LAT, CENTER_LONG
5     max_dist_mean = df['max_dist'].mean()
6
7     df['max_dist_correct'] = np.where(df['max_dist'] < 50, max_dist_mean, df['max_dist
8     df['max_dist_correct']
9
10    for feature in ['long', 'lat']:
11        if (feature == 'long'):
12            coord_length = long_length
13        else:
14            coord_length = lat_length
15        df[feature + '_start'] = np.where(df['angel'] == 360,
16            df[feature] - df['max_dist_correct']/coord_length/1000,
17            df[feature] + df['max_dist_correct'] * np.where(feature
18
19
20
21
22        df[feature + '_end'] = np.where(df['angel'] == 360,
23            df[feature]+df['max_dist_correct']/coord_length/1000,
24            df[feature] + df['max_dist_correct'] * np.where(feature
25
26
27
28
29        df[feature + '_min'] = df[[feature + '_start', feature + '_end', feature]].min
30        df[feature + '_max'] = df[[feature + '_start', feature + '_end', feature]].max
31
32        df[feature] = (df[feature + '_max'] + df[feature + '_min'])/2
33    return df

```

In [18]:

```

1 df_data_copy = df_data.copy()
2 df_data = BS_rectangle_coordinate(df_data_copy)

```


In [19]:

```
1 df_data.head()
```

Out[19]:

	lac	cid	msisdn	imei	event_type	tstamp	long
0	5029	40798	158529599791	353111050313790	0	1369252800974	37.925474 55.782:
1	7782	56870	158520145943	012929002676510	5	1369252801396	37.422847 55.744:
2	7794	32226	158521798391	351994049226010	5	1369252802013	37.700660 55.580:
3	7758	33528	158537830573	358627016731770	0	1369252802331	37.792199 55.704:
4	5060	17568	158510204039	352458051767250	8	1369252803586	37.190576 55.830:

Создаю bin по времени для создания ряда событий

In [20]:

```

1 from sklearn.preprocessing import KBinsDiscretizer
2 from sklearn.datasets import make_blobs
3
4 def CreateTimeBeans(df_, n_bins_ = 10):
5     df = df_.copy()
6
7     dt_list = np.reshape(df.datetime.values.tolist(), (len(df.datetime.values.tolist()),
8
9     enc = KBinsDiscretizer(n_bins=n_bins_, encode='ordinal', strategy='uniform')
10    dt_array_bins = enc.fit_transform(dt_list)
11
12    df['datetime_bin'] = dt_array_bins[:,0]
13    df['datetime_bin'] = df['datetime_bin'].astype('int')
14
15    return df, df['datetime_bin'].unique()

```

In [21]:

```

1  '''
2  sns.relplot(data=df_data_bin.groupby(['datetime_bin', 'day_of_month'])['lac'].count().reset_index(),
3              y="lac", x="datetime_bin", hue = 'day_of_month')
4  '''

```

Out[21]:

```

'\nsns.relplot(data=df_data_bin.groupby(['datetime_bin', 'day_of_month
\'])['lac'].count().reset_index(),\n              y="lac", x="datetime_bin",
hue = \'day_of_month\')\n'

```

Агрегирую таблицу с событиями для преобразования в ряд.

При необходимости, добавляю технические события в пустые интервалы времени

In [22]:

```

1  # декартовое произведение абонентов и интервалов
2  def CreatePivotBase(df):
3      # декартовое произведение абонентов и часов
4      df_data_msisdn = pd.DataFrame(df['msisdn'].drop_duplicates()).sort_values('msisdn')
5      df_data_datetime = pd.DataFrame(df['datetime_bin'].drop_duplicates()).sort_values('datetime_bin')
6
7      df_data_msisdn['K'] = 1
8      df_data_datetime['K'] = 1
9      df_data_pivot_base = df_data_msisdn.merge(df_data_datetime, on = ['K']).reset_index()
10
11     return df_data_pivot_base

```

In [23]:

```

1  # агрегация данных по интервалам
2  def AggData(df):
3      #Создаем датафрейм в котором записаны часовые срезы расположения абонента и нумерованные события
4
5      df_data_agg = df.groupby(['msisdn', 'datetime_bin'])['long_min', 'long_max', 'lat_min', 'lat_max'].agg(['min', 'max'])
6      df_data_agg['datetime'] = df_data_agg['datetime_bin']
7      df_data_agg = df_data_agg.sort_values(['msisdn', 'datetime_bin'])
8      df_data_agg['state'] = df_data_agg.groupby(['msisdn'])['msisdn'].cumcount()
9      return df_data_agg

```

In [24]:

```

1  # сдвигаю событие на пол часа и объединяю с оригинальными событиями
2  from datetime import timedelta
3  def AddDataDelta(df_):
4      df_delta = df_.copy()
5      df_delta['datetime'] = df_delta['datetime'] - timedelta(seconds = (60*60))
6      df = df_.copy()
7      df = df.append(df_delta)
8      return df

```

In [25]:

```

1  # Промежутки времени, в которых не зафиксировано действие абонента, заполняю линейном
2  # Концевые точки продолжаю до границы интервала наблюдения константами.
3
4  def LinearCombData(df_data_pivot_base, df_data_agg):
5
6      df_data_pivot_base = df_data_pivot_base.merge(df_data_agg[['msisdn', 'datetime_bin',
7      # указываю состояние, в котором произошло предыдущее событие
8      df_data_pivot_base['state_previous'] = 1 - df_data_pivot_base['state'].isna()*1
9      df_data_pivot_base['state_previous'] = df_data_pivot_base.groupby('msisdn')['state']
10
11     # указываю состояние, в котором произойдет следующее событие
12     df_data_pivot_base.drop(columns = ['state'], inplace = True)
13     df_data_pivot_base['state_next'] = df_data_pivot_base['state_previous'] + 1
14     # подгружаю информацию о предыдущем и следующем событии
15     for event in ['previous', 'next']:
16         df_data_pivot_base = df_data_pivot_base.merge(df_data_agg[['msisdn', 'long', 'lat',
17         'lat_max', 'lat_min', 'long_max', 'long_min', 'datetime_bin', 'state']],
18         columns =
19             {
20                 'long': 'long',
21                 'lat': 'lat',
22                 'long_max': 'long_max',
23                 'lat_max': 'lat_max',
24                 'long_min': 'long_min',
25                 'lat_min': 'lat_min',
26                 'datetime': 'datetime_bin',
27                 'state': 'state'
28             },
29             on = ['msisdn', 'datetime_bin'],
30             how = 'left')
31     for feature in ['long', 'long_max', 'long_min', 'lat', 'lat_max', 'lat_min']:
32         # рассчитываю текущее событие для каждого периода как линейную комбинацию окружающих
33         df_data_pivot_base[feature] = np.where(df_data_pivot_base[feature+'_previous'].isna(),
34         df_data_pivot_base[feature+'_next'],
35         np.where(df_data_pivot_base[feature+'_next'].isna(),
36         df_data_pivot_base[feature+'_previous'],
37         df_data_pivot_base[feature+'_previous'] +
38         (df_data_pivot_base['datetime_bin'] - df_data_pivot_base['datetime_bin_previous']) *
39         (df_data_pivot_base[feature+'_next'] - df_data_pivot_base[feature+'_previous']) /
40         (df_data_pivot_base['datetime_bin'] - df_data_pivot_base['datetime_bin_previous']))
41
42     return df_data_pivot_base[['msisdn', 'datetime_bin', 'long', 'long_max', 'long_min', 'lat', 'lat_max', 'lat_min', 'state_previous', 'state_next']]

```

In [26]:

```

1  # Промежутки времени, в которых не зафиксировано действие абонента, заполняю линейном
2  # Концевые точки продолжаю до границы интервала наблюдения константами.
3
4  def CombData(df_data_pivot_base, df_data_agg):
5
6      df_data_pivot_base = df_data_pivot_base.merge(df_data_agg, on = ['msisdn', 'datetime_bin'], how = 'left')
7
8      return df_data_pivot_base[['msisdn', 'datetime_bin', 'long', 'long_max', 'long_min', 'lat', 'lat_max', 'lat_min', 'state_previous', 'state_next']]

```

In [27]:

```

1 # Создает pivot из исходной таблицы
2 def CreatePivot(df_for_pivot):
3     df_pivot = df_for_pivot.pivot(index = 'msisdn', columns = 'datetime_bin', values =
4     df_pivot.columns = df_pivot.columns.astype('str') + '_' + 'long'
5
6     for feature in ['long_max', 'long_min', 'lat', 'lat_max', 'lat_min']:
7         df_agg_feature = df_for_pivot.pivot(index = 'msisdn', columns = 'datetime_bin'
8         df_agg_feature.columns = df_agg_feature.columns.astype('str') + '_' + feature
9         df_pivot = pd.concat([df_pivot, df_agg_feature], axis = 1)
10
11     return df_pivot.reset_index()

```

In [28]:

```

1 # Определяем пары, для которых знаем результаты
2 # * пары - все, которые указаны в таблице facts
3 # * не пары: номер, у которого есть пара с ЛЮБЫМ номером не из своей пары (этот номер
4 # Упрощение, что нет троек номеров.
5
6 def CreateMsisdnPairs(train_pairs, df, frac_):
7     msisdn_pairs = pd.concat([train_pairs[[0, 1]].rename(columns = {0:'msisdn_x', 1:'ms
8     # , facts[[1, 0]].rename(columns = {1:'msisdn_x', 0:'msisdn_y'
9     ], axis = 0)
10     msisdn_pairs_dubli = pd.concat([train_pairs[[0, 1]].rename(columns = {0:'msisdn_x'
11     , train_pairs[[1, 0]].rename(columns = {1:'msisdn_x', 0:'msisdn
12     ], axis = 0)
13     msisdn_pairs['is_pair'] = 1
14     msisdn_pairs['K'] = 1
15     msisdn_pairs_dubli['K'] = 1
16
17     df_data_msisdn = pd.DataFrame(df['msisdn'].drop_duplicates()).sort_values('msisdn'
18     df_data_msisdn['K'] = 1
19
20     # любой номер из пары с другими номерами не создает пары
21     msisdn_pairs_part1 = df_data_msisdn[~df_data_msisdn['msisdn'].isin(msisdn_pairs_dub
22     msisdn_pairs_part1['is_pair'] = 0
23
24     msisdn_pairs_part1 = msisdn_pairs_part1.sample(frac = frac_)
25
26     # объединяю пары и не пары
27     msisdn_pairs = pd.concat([msisdn_pairs, msisdn_pairs_part1
28     # , msisdn_pairs_part2
29     ], axis = 0).drop(columns = ['K'])
30
31     # перемешиваю первый и второй номер
32     msisdn_pairs_part1 = msisdn_pairs.sample(frac = .5)
33     msisdn_pairs_part2 = msisdn_pairs.append(msisdn_pairs_part1).drop_duplicates(keep
34     msisdn_pairs = msisdn_pairs_part1.append(msisdn_pairs_part2)
35     print(msisdn_pairs.shape)
36
37     return msisdn_pairs

```

In [29]:

```

1  #Создаю таблицу для обучения
2
3  def CreateAndProcessPivot(msisdn_pairs, data_agg_pivot, bins_list, identifier,
4                             dist_ind = True,
5                             coord_ind = True,
6                             center_dist_ind = True ):
7
8      df_X = msisdn_pairs.merge(data_agg_pivot.rename(columns = {'msisdn':'msisdn_x'}),
9      # считая расстояние в каждом часе
10     columns_name_dist = list()
11     columns_name_intersect = list()
12     columns_name_lat_x = list()
13     columns_name_long_x = list()
14     columns_name_lat_y = list()
15     columns_name_long_y = list()
16     columns_name_center_dist_x = list()
17     columns_name_center_dist_y = list()
18
19     columns_name = list()
20     columns_name_identifier = list()
21
22     lat_length = haversine_array(CENTER_LAT -.5, CENTER_LONG, CENTER_LAT + .5, CENTER_
23     long_length = haversine_array(CENTER_LAT, CENTER_LONG -.5, CENTER_LAT, CENTER_LONG
24
25     for i in bins_list:
26         str_i = str(int(i))
27         columns_name_dist.append(str_i+"_dist")
28         columns_name_intersect.append(str_i+"_intersect")
29         columns_name_lat_x.append(str_i+"_lat_x")
30         columns_name_long_x.append(str_i+"_long_x")
31         columns_name_lat_y.append(str_i+"_lat_y")
32         columns_name_long_y.append(str_i+"_long_y")
33         columns_name_center_dist_x.append(str_i+"_center_dist_x")
34         columns_name_center_dist_y.append(str_i+"_center_dist_y")
35
36         lat_diff = (df_X[str_i+"_lat_x"].values - df_X[str_i+"_lat_y"].values)*lat_length
37         long_diff = (df_X[str_i+"_long_x"].values - df_X[str_i+"_long_y"].values)*long_length
38         dist = np.sqrt(lat_diff**2 + long_diff**2)
39         df_X[str_i+"_dist"] = dist
40         df_X[str_i+"_intersect"] = ((df_X[str_i+"_lat_max_x"] >= df_X[str_i+"_lat_min_x"] &
41                                     (df_X[str_i+"_lat_max_y"] >= df_X[str_i+"_lat_min_y"] &
42                                     (df_X[str_i+"_long_max_x"] >= df_X[str_i+"_long_min_x"] &
43                                     (df_X[str_i+"_long_max_y"] >= df_X[str_i+"_long_min_y"] &
44         center_x_lat_diff = (df_X[str_i+"_lat_x"].values - CENTER_LAT)*lat_length
45         center_y_lat_diff = (df_X[str_i+"_lat_y"].values - CENTER_LAT)*lat_length
46
47         center_x_long_diff = (df_X[str_i+"_long_x"].values - CENTER_LONG)*long_length
48         center_y_long_diff = (df_X[str_i+"_long_y"].values - CENTER_LONG)*long_length
49
50         df_X[str_i+"_center_dist_x"] = np.sqrt(center_x_lat_diff**2 + center_x_long_diff**2)
51         df_X[str_i+"_center_dist_y"] = np.sqrt(center_y_lat_diff**2 + center_y_long_diff**2)
52
53     # считая среднее расстояние
54     df_X['dist_mean'] = df_X[columns_name_dist].mean(axis = 1)
55     df_X['intersect_mean'] = df_X[columns_name_intersect].mean(axis = 1)
56     df_X['intersect_count'] = df_X[columns_name_dist].count(axis = 1)/len(bins_list)
57     df_X['center_dist_x_mean'] = df_X[columns_name_center_dist_x].mean(axis = 1)
58     df_X['center_dist_y_mean'] = df_X[columns_name_center_dist_y].mean(axis = 1)
59

```

```

60     if(dist_ind):
61         columns_name = [*columns_name, *columns_name_dist]
62
63     if(coord_ind):
64         columns_name = [*columns_name, *columns_name_lat_x, *columns_name_lat_y, *columns_name_lat_z]
65
66     if(center_dist_ind):
67         columns_name = [*columns_name, *columns_name_center_dist_x, *columns_name_center_dist_y, *columns_name_center_dist_z]
68
69     columns_name = [*columns_name,
70                     *['dist_mean', 'intersect_mean', 'intersect_count', 'center_dist_x', 'center_dist_y', 'center_dist_z']]
71
72     columns_name = [*columns_name, *['is_pair', 'msisdn_x', 'msisdn_y']]
73     df_X = df_X[columns_name].set_index(['is_pair', 'msisdn_x', 'msisdn_y'])
74     df_X.columns = df_X.columns + '_' + identifier
75
76     print(df_X.shape)
77
78     return df_X

```

Создаю таблицу

- Линейная комбинация недостающих событий
- Сдвиг по времени событие для создания большего количества общих интервалов

Создаем pivottable

Создаю pivot с координатами в каждый период

Создаем таблицу с данными для обучения

Описание целевого события.

- пары - все, которые указаны в таблице facts
- не пары: номер, у которого есть пара с ЛЮБЫМ номером не из своей пары (этот номер не обязательно из таблицы facts) Упрощение, что нет троек номеров.

Если (a, b) пара, то (b, a) - тоже пара. Если (a, b) - пара, то (a, c) Симметричные пары убираю, для этого выбираю (a > b), а потом перемешиваю.

Пары (a,a) убираем Модифицируем на шум одинаковых абонентов(?) Смешаем время для сравнения, расширяем диапазоны в обе стороны для второго абонента.

Таблица

Создаю таблицу для обучения

In [30]:

```
1 #модель
2 !pip install xgboost
3
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.model_selection import train_test_split
6 from xgboost import XGBClassifier
7
8 from sklearn.metrics import accuracy_score
9 from sklearn.metrics import roc_auc_score
10 from sklearn.metrics import average_precision_score
11
12 from sklearn.metrics import roc_curve
13 from sklearn.metrics import plot_roc_curve
14 from sklearn.metrics import precision_recall_curve
15 from sklearn.metrics import plot_precision_recall_curve
16 from sklearn.metrics import average_precision_score
17 import warnings
18 warnings.filterwarnings('ignore')
```

Requirement already satisfied: xgboost in c:\users\ekaterina.adischeva\appdata\local\continuum\anaconda3\lib\site-packages (1.3.3)

Requirement already satisfied: scipy in c:\users\ekaterina.adischeva\appdata\local\continuum\anaconda3\lib\site-packages (from xgboost) (1.2.1)

Requirement already satisfied: numpy in c:\users\ekaterina.adischeva\appdata\local\continuum\anaconda3\lib\site-packages (from xgboost) (1.16.2)

In [33]:

```

1 Cls_result_agg = pd.DataFrame()
2 BINS_QNT = 90
3 ln_ind = True
4 shift_inf = False
5 cls_ind = 'RF'
6 dist_ind_ln, coord_ind_ln, center_dist_ind_ln = True, False, False
7 dist_ind_shift, coord_ind_shift, center_dist_ind_shift = True, False, False
8 clfs = []
9
10 for bins_q, ln_ind, shift_ind, clf_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln,
11     [30, True, True, 'RF', True, True, False, True, False, False],
12     [30, True, True, 'GB', True, True, False, True, False, False],
13                                     [30, True, True, 'RF',
14                                     [60, True, False, 'GB'
15                                     [60, True, False, 'RF'
16                                     [60, False, True, 'GB'
17                                     [60, False, True, 'RF'
18                                     [60, True, False, 'GB'
19                                     [60, True, False, 'RF'
20                                     [60, False, True, 'GB'
21                                     [60, False, True, 'RF'
22                                     [90, True, False, 'RF'
23                                     [90, True, False, 'GB'
24                                     [90, True, False, 'RF'
25                                     [90, True, False, 'GB'
26                                     [90, False, True, 'RF'
27                                     [90, False, True, 'GB'
28                                     [90, False, True, 'RF'
29                                     [90, False, True, 'GB'
30                                     [90, False, True, 'RF'
31                                     [90, False, True, 'GB'
32                                     [90, False, True, 'RF'
33                                     [90, False, True, 'GB'
34                                     ]:
35
36 df_data_bin_ln, bins_list_ln = CreateTimeBeans(df_data, BINS_QNT)
37 df_data_pivot_base_ln = CreatePivotBase(df_data_bin_ln)
38 df_data_bin_agg_ln = AggData(df_data_bin_ln)
39 df_data_agg_for_pivot_ln = LinearCombData(df_data_pivot_base_ln, df_data_bin_agg_l
40
41 df_data_shift = AddDataDelta(df_data)
42 df_data_bin_shift, bins_list_shift = CreateTimeBeans(df_data_shift, BINS_QNT)
43 df_data_pivot_base_shift = CreatePivotBase(df_data_bin_shift)
44 df_data_bin_agg_shift = AggData(df_data_bin_shift)
45 df_data_agg_for_pivot_shift = CombData(df_data_pivot_base_shift, df_data_bin_agg_s
46
47 df_msisdn_pairs = CreateMsisdnPairs(facts, df_data, 0.25)
48
49 df_data_agg_pivot_ln = CreatePivot(df_data_agg_for_pivot_ln)
50 df_data_agg_pivot_shift = CreatePivot(df_data_agg_for_pivot_shift)
51
52 df_X_ln = CreateAndProcessPivot(df_msisdn_pairs, df_data_agg_pivot_ln, bins_list_l
53 df_X_shift = CreateAndProcessPivot(df_msisdn_pairs, df_data_agg_pivot_shift, bins_
54
55 df_X = pd.DataFrame()
56 if (ln_ind):
57     df_X = pd.concat([df_X_ln, df_X], axis = 1)
58 if (shift_ind):
59     df_X = pd.concat([df_X_shift, df_X], axis = 1)

```



```

60
61 df_X.reset_index(inplace = True)
62
63 X = df_X.drop(columns = ['msisdn_x', 'msisdn_y', 'is_pair']).fillna(0)
64 y = df_X['is_pair']
65
66 X_train, X_test, y_train, y_test = train_test_split(
67     X, y, test_size=0.3, random_state=42)
68 if(clf_ind == 'RF'):
69     clf = RandomForestClassifier(random_state=0)
70 else:
71     clf = XGBClassifier(random_state=42)
72
73 clf.fit(X_train, y_train)
74
75 y_test_predict = clf.predict(X_test)
76 y_train_predict = clf.predict(X_train)
77 y_test_predict_proba = clf.predict_proba(X_test)
78
79 print('[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist
80     [bins_q, ln_ind, shift_inf, cls_ind,
81     dist_ind_ln, coord_ind_ln, center_dist_ind_ln,
82     dist_ind_shift, coord_ind_shift, center_dist_ind_shift])
83
84 print('roc_auc_score:', roc_auc_score(y_train, y_train_predict),
85     'average_precision_score:', average_precision_score(y_train, y_train_predict)
86 print("roc_auc_score:", roc_auc_score(y_test, y_test_predict),
87     'average_precision_score:', average_precision_score(y_test, y_test_predict))
88
89
90 disp = plot_precision_recall_curve(clf, X_test, y_test)
91 average_precision = average_precision_score(y_test, y_test_predict)
92 print(average_precision)
93
94 disp = plot_roc_curve(clf, X_test, y_test)
95 roc_auc = roc_auc_score(y_test, y_test_predict)
96 print(roc_auc)
97
98
99 Cls_result = pd.DataFrame(y_test)
100 Cls_result['proba'] = y_test_predict_proba[:,1]
101 Cls_result['predict'] = y_test_predict
102 Cls_result = Cls_result.join(df_X[['msisdn_x', 'msisdn_y']])
103
104 alpha = .1
105
106 Cls_result['predict_alpha'] = np.where(Cls_result['proba'] > alpha, 1, 0)
107 Cls_result_agg = Cls_result_agg.append(Cls_result.groupby(['is_pair', 'predict_alpha']
108 Cls_result_agg = Cls_result_agg.append(Cls_result.groupby(['is_pair', 'predict']).
109

```

```

(129303, 3)
(127227, 435)
(127227, 91)
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist
t_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [30, Tru
e, False, 'RF', True, True, False, True, False, False]
roc_auc_score: 1.0 average_precision_score: 1.0
roc_auc_score: 0.5483346556393673 average_precision_score: 0.0422082339843

```

```
6712
0.04220823398436712
0.5483346556393673
(129302, 3)
(127154, 435)
(127154, 91)
[04:11:19] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release
_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluat
ion metric used with the objective 'binary:logistic' was changed from 'err
or' to 'logloss'. Explicitly set eval_metric if you'd like to restore the
old behavior.
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dis
t_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [30, Tru
e, False, 'RF', True, True, False, True, False, False]
roc_auc_score: 1.0 average_precision_score: 1.0
roc_auc_score: 0.634523571703866 average_precision_score: 0.13511345785836
57
0.1351134578583657
0.634523571703866
(129314, 3)
(127190, 263)
(127190, 91)
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dis
t_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [30, Tru
e, False, 'RF', True, False, True, True, False, False]
roc_auc_score: 0.9939759036144578 average_precision_score: 0.9879630390193
753
roc_auc_score: 0.5693395154027089 average_precision_score: 0.0542312362620
563
0.0542312362620563
0.5693395154027089
(129305, 3)
(127256, 435)
(127256, 91)
[04:15:37] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release
_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluat
ion metric used with the objective 'binary:logistic' was changed from 'err
or' to 'logloss'. Explicitly set eval_metric if you'd like to restore the
old behavior.
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dis
t_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [60, Tru
e, False, 'RF', True, True, False, True, False, False]
roc_auc_score: 1.0 average_precision_score: 1.0
roc_auc_score: 0.5776991036982644 average_precision_score: 0.0847560474613
4122
0.08475604746134122
0.5776991036982644
(129303, 3)
(127142, 435)
(127142, 91)
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dis
t_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [60, Tru
e, False, 'RF', True, True, False, True, False, False]
roc_auc_score: 1.0 average_precision_score: 1.0
roc_auc_score: 0.560566700857963 average_precision_score: 0.07002436604198
395
0.07002436604198395
0.560566700857963
(129300, 3)
(127196, 263)
(127196, 435)
```

```
[04:19:51] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
```

```
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [60, False, False, 'RF', True, False, True, True, True, False]
```

```
roc_auc_score: 1.0 average_precision_score: 1.0
```

```
roc_auc_score: 0.513797080794038 average_precision_score: 0.004389437033931125
```

```
0.004389437033931125
```

```
0.513797080794038
```

```
(129317, 3)
```

```
(127211, 263)
```

```
(127211, 435)
```

```
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [60, False, False, 'RF', True, False, True, True, True, False]
```

```
roc_auc_score: 1.0 average_precision_score: 1.0
```

```
roc_auc_score: 0.5 average_precision_score: 0.0009957027565244733
```

```
0.0009957027565244733
```

```
0.5
```

```
(129286, 3)
```

```
(127182, 263)
```

```
(127182, 91)
```

```
[04:22:08] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
```

```
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [60, True, False, 'RF', True, False, True, True, False, False]
```

```
roc_auc_score: 0.99375 average_precision_score: 0.9875112325474295
```

```
roc_auc_score: 0.5768574876692202 average_precision_score: 0.08478097711487831
```

```
0.08478097711487831
```

```
0.5768574876692202
```

```
(129317, 3)
```

```
(127224, 263)
```

```
(127224, 91)
```

```
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [60, True, False, 'RF', True, False, True, True, False, False]
```

```
roc_auc_score: 1.0 average_precision_score: 1.0
```

```
roc_auc_score: 0.5713367705577745 average_precision_score: 0.06687726442494835
```

```
0.06687726442494835
```

```
0.5713367705577745
```

```
(129304, 3)
```

```
(127170, 263)
```

```
(127170, 263)
```

```
[04:25:58] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
```

```
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [60, False, False, 'RF', True, False, True, True, False, True]
```

```

roc_auc_score: 0.9943820224719101 average_precision_score: 0.9887752785007
013
roc_auc_score: 0.5166273182760159 average_precision_score: 0.0090934706822
888
0.0090934706822888
0.5166273182760159
(129311, 3)
(127101, 263)
(127101, 263)
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dis
t_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [60, Fal
se, False, 'RF', True, False, True, True, False, True]
roc_auc_score: 1.0 average_precision_score: 1.0
roc_auc_score: 0.5 average_precision_score: 0.0011276913797172905
0.0011276913797172905
0.5
(129285, 3)
(127183, 91)
(127183, 91)
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dis
t_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [90, Tru
e, False, 'RF', True, False, False, True, False, False]
roc_auc_score: 1.0 average_precision_score: 1.0
roc_auc_score: 0.62156915297164 average_precision_score: 0.169133017173903
06
0.16913301717390306
0.62156915297164
(129297, 3)
(127155, 91)
(127155, 91)
[04:29:51] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release
_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluat
ion metric used with the objective 'binary:logistic' was changed from 'err
or' to 'logloss'. Explicitly set eval_metric if you'd like to restore the
old behavior.
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dis
t_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [90, Tru
e, False, 'RF', True, False, False, True, False, False]
roc_auc_score: 0.9941860465116279 average_precision_score: 0.9883833279684
293
roc_auc_score: 0.5605535864495826 average_precision_score: 0.0613662776611
37024
0.061366277661137024
0.5605535864495826
(129307, 3)
(127065, 349)
(127065, 91)
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dis
t_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [90, Tru
e, False, 'RF', False, True, False, True, False, False]
roc_auc_score: 0.99375 average_precision_score: 0.9875112429029176
roc_auc_score: 0.49998687009269716 average_precision_score: 0.001023084994
7534103
0.0010230849947534103
0.49998687009269716
(129297, 3)
(127164, 349)
(127164, 91)
[04:32:14] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release
_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluat
ion metric used with the objective 'binary:logistic' was changed from 'err

```

or' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

```
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [90, True, False, 'RF', False, True, False, True, False, False]
```

roc_auc_score: 1.0 average_precision_score: 1.0

roc_auc_score: 0.5254835907796995 average_precision_score: 0.008295863158248479

0.008295863158248479

0.5254835907796995

(129296, 3)

(127213, 349)

(127213, 91)

```
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [90, False, False, 'RF', False, True, False, True, False, False]
```

roc_auc_score: 1.0 average_precision_score: 1.0

roc_auc_score: 0.5142594875291773 average_precision_score: 0.010414701463857735

0.010414701463857735

0.5142594875291773

(129307, 3)

(127176, 349)

(127176, 91)

[04:34:15] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

```
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [90, False, False, 'RF', False, True, False, True, False, False]
```

roc_auc_score: 0.9235294117647059 average_precision_score: 0.8472048532071355

roc_auc_score: 0.5440782966036699 average_precision_score: 0.04493016507837638

0.04493016507837638

0.5440782966036699

(129315, 3)

(127183, 263)

(127183, 349)

```
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [90, False, False, 'RF', True, False, True, False, True, False]
```

roc_auc_score: 1.0 average_precision_score: 1.0

roc_auc_score: 0.5 average_precision_score: 0.0013104442405975625

0.0013104442405975625

0.5

(129309, 3)

(127172, 263)

(127172, 349)

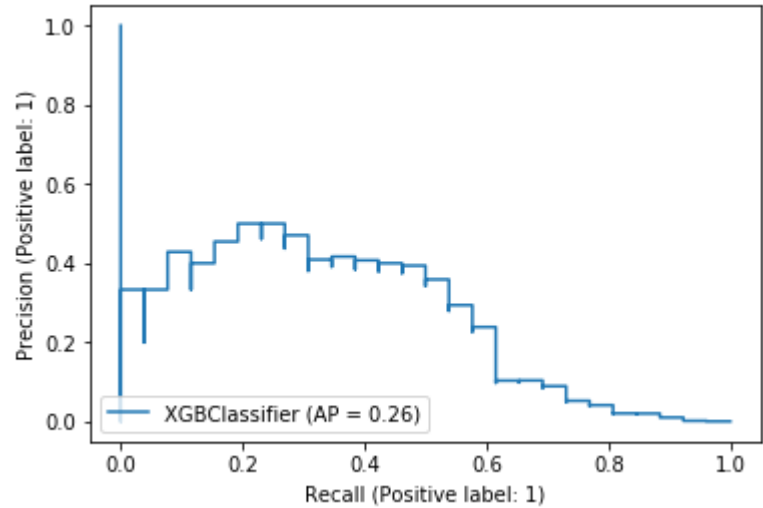
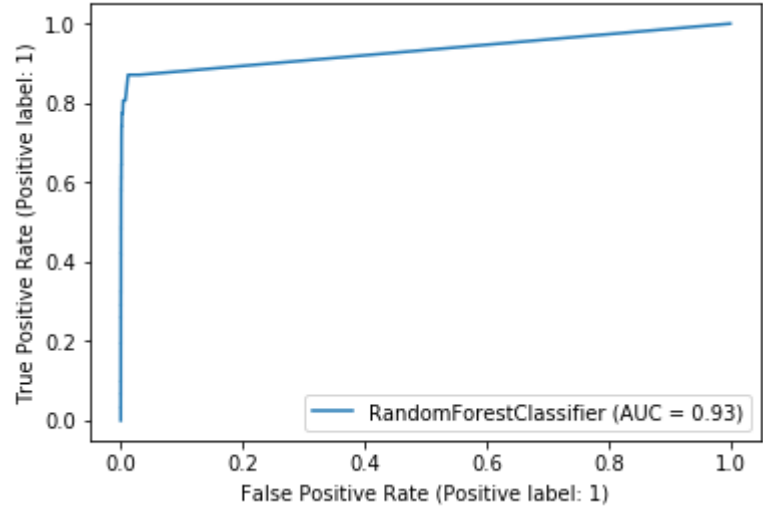
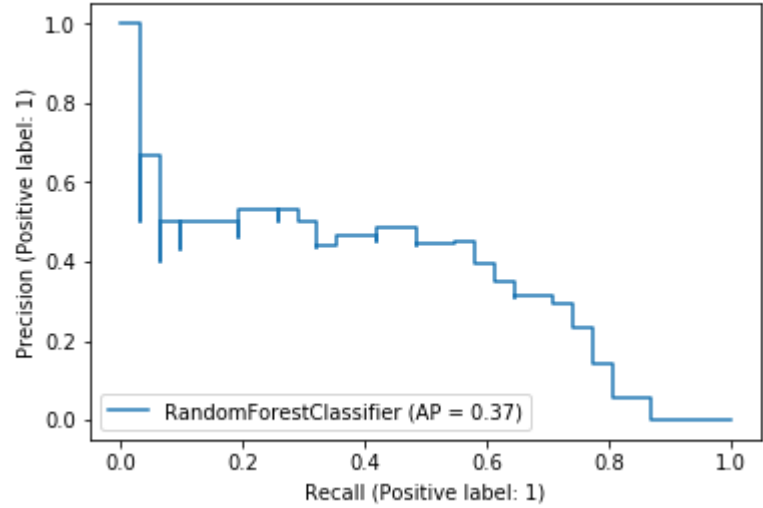
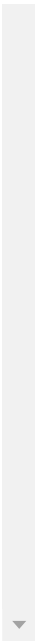
[04:35:59] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

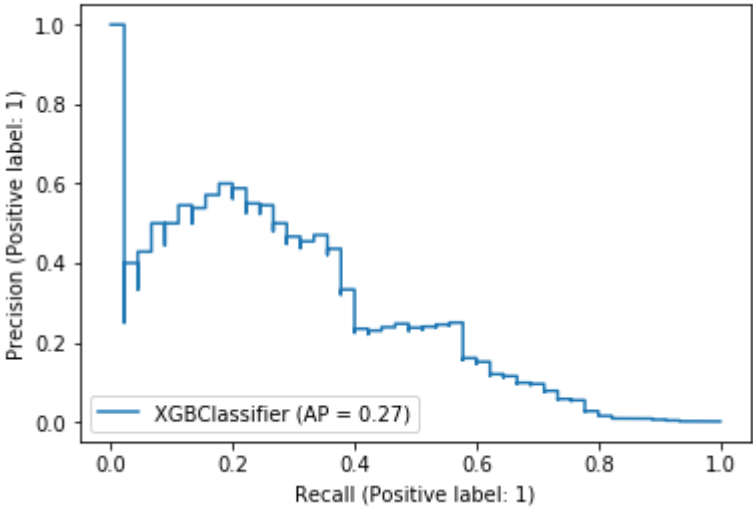
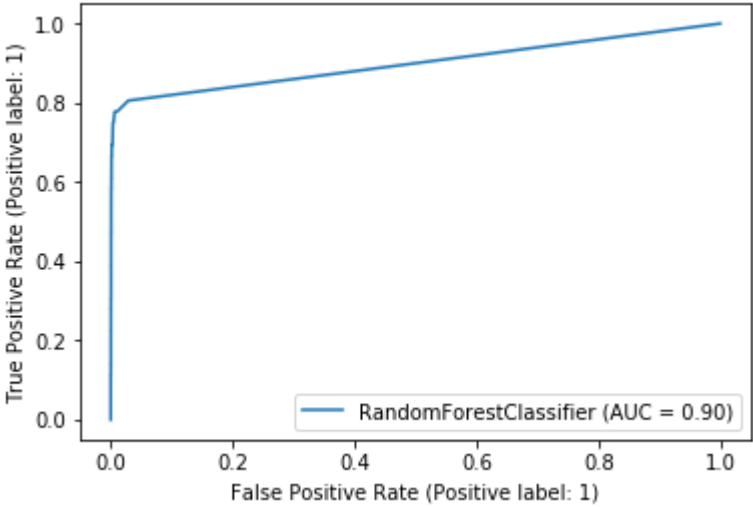
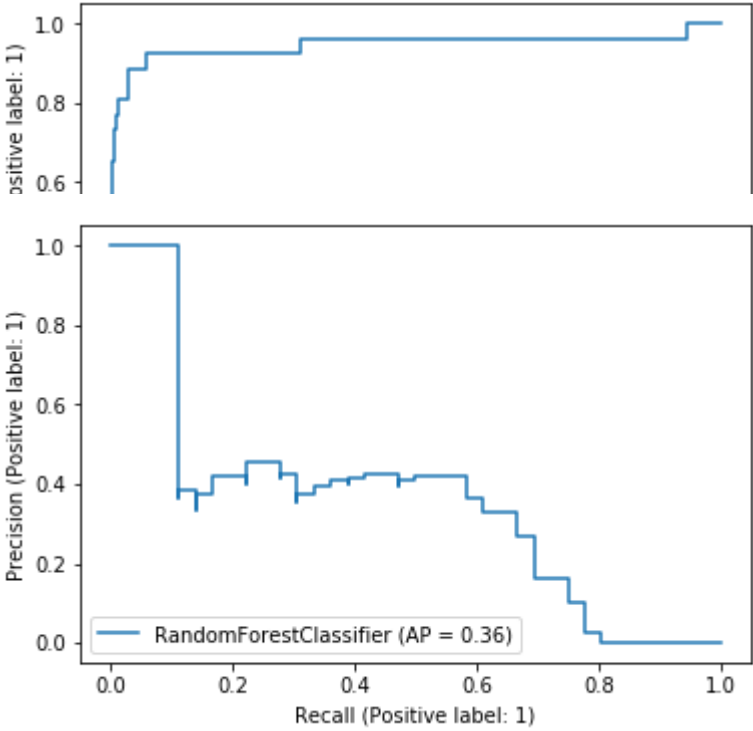
```
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [90, False, False, 'RF', True, False, True, False, True, False]
```

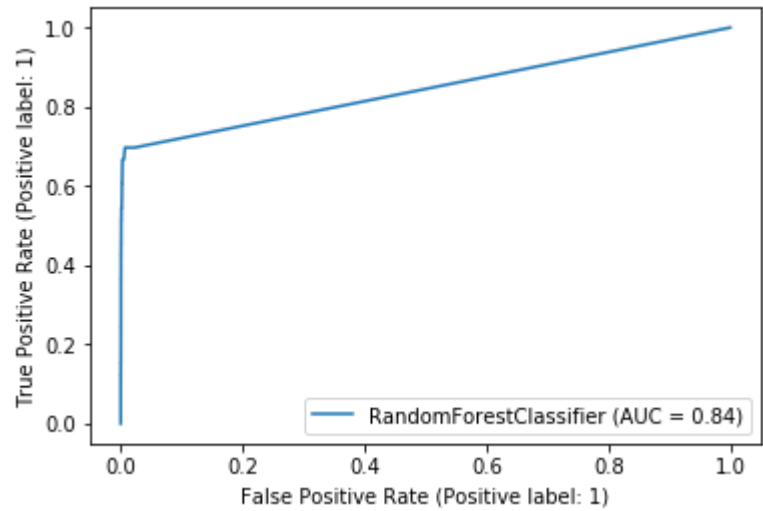
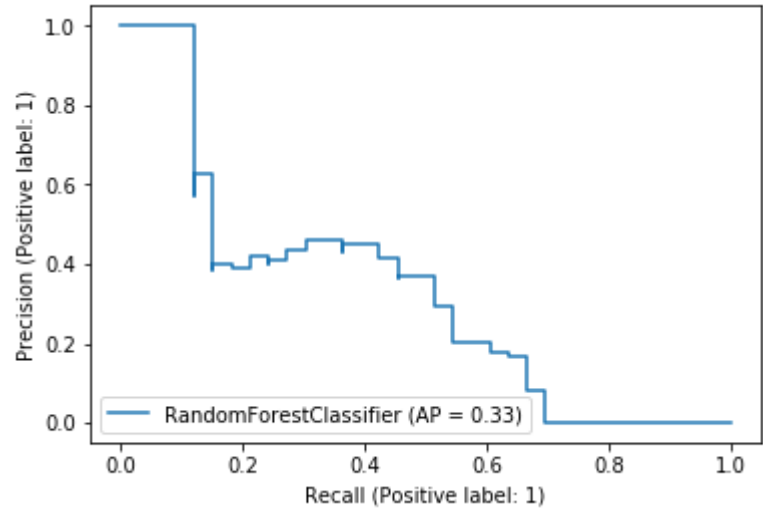
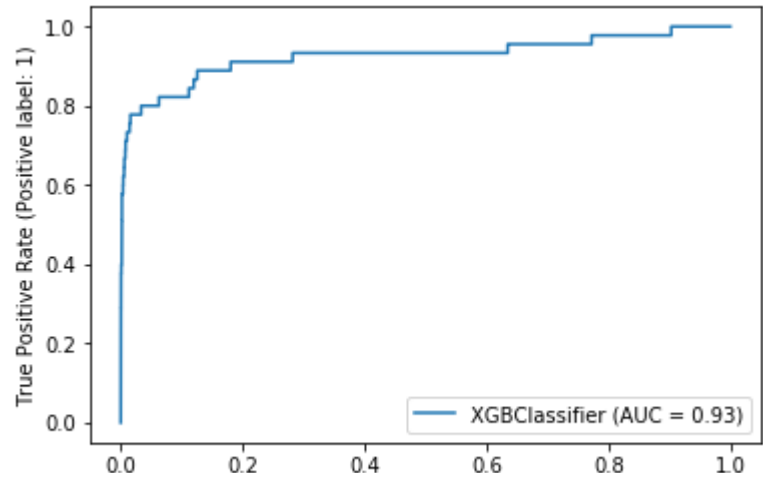
roc_auc_score: 1.0 average_precision_score: 1.0

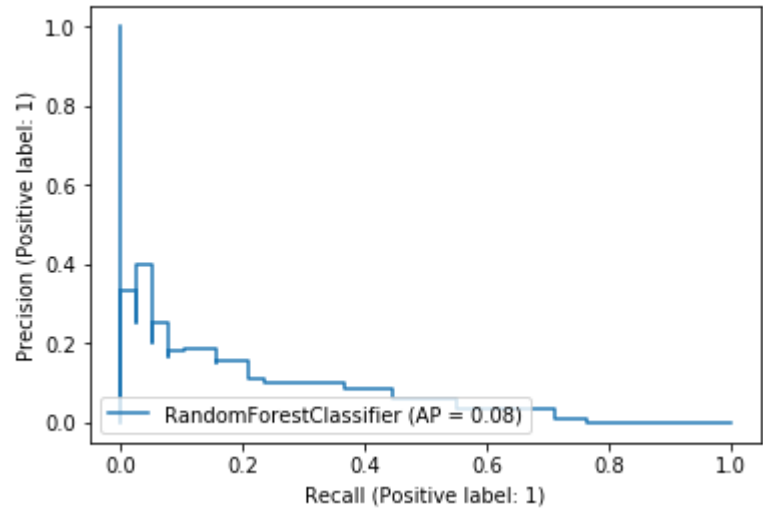
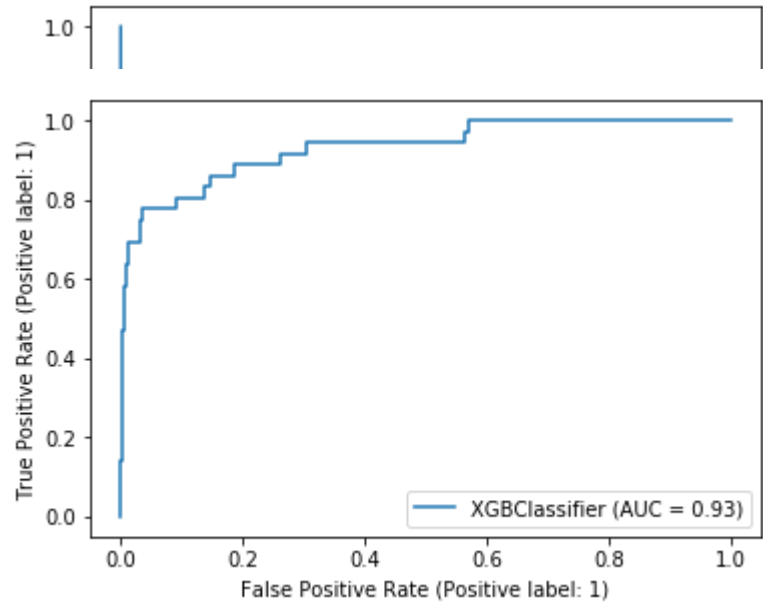
```
roc_auc_score: 0.5127549184039095 average_precision_score: 0.005269520209759
253
0.005269520209759253
0.5127549184039095
(129309, 3)
(127176, 521)
(127176, 263)
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_
ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [90, False,
False, 'RF', False, True, True, True, False, True]
roc_auc_score: 1.0 average_precision_score: 1.0
roc_auc_score: 0.5 average_precision_score: 0.0007600975021623464
0.0007600975021623464
0.5
(129311, 3)
(127195, 521)
(127195, 263)
[04:38:04] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_
1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation
metric used with the objective 'binary:logistic' was changed from 'error' to
'logloss'. Explicitly set eval_metric if you'd like to restore the old behav
ior.
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_
ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [90, False,
False, 'RF', False, True, True, True, False, True]
roc_auc_score: 1.0 average_precision_score: 1.0
roc_auc_score: 0.5166404399101296 average_precision_score: 0.011871089097955
63
0.01187108909795563
0.5166404399101296
(129315, 3)
(127176, 263)
(127176, 521)
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_
ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [90, False,
False, 'RF', True, False, True, False, True, True]
roc_auc_score: 1.0 average_precision_score: 1.0
roc_auc_score: 0.5 average_precision_score: 0.001048410347810133
0.001048410347810133
0.5
(129314, 3)
(127130, 263)
(127130, 521)
[04:40:18] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_
1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation
metric used with the objective 'binary:logistic' was changed from 'error' to
'logloss'. Explicitly set eval_metric if you'd like to restore the old behav
ior.
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_
ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [90, False,
False, 'RF', True, False, True, False, True, True]
roc_auc_score: 1.0 average_precision_score: 1.0
roc_auc_score: 0.4999343745898412 average_precision_score: 0.001153674716169
8
0.0011536747161698
0.4999343745898412
```

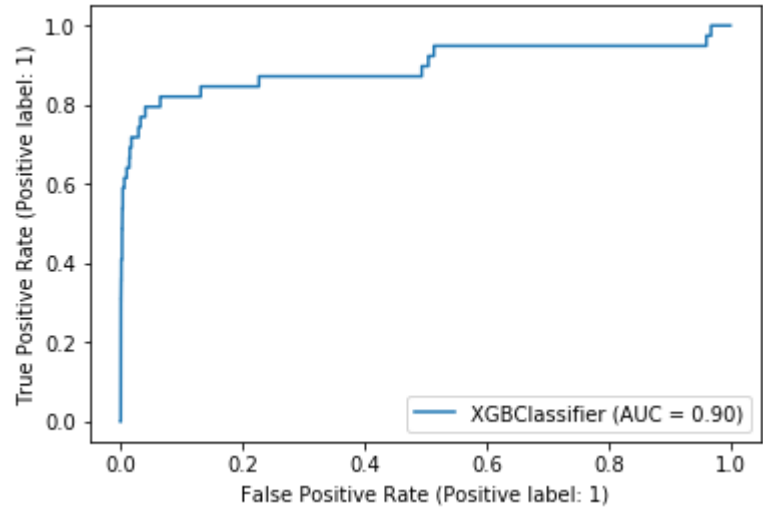
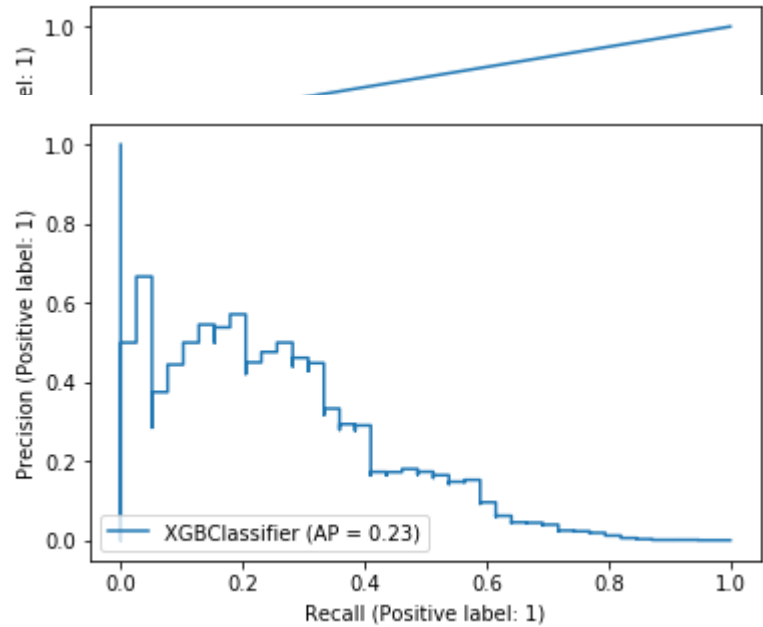
▲

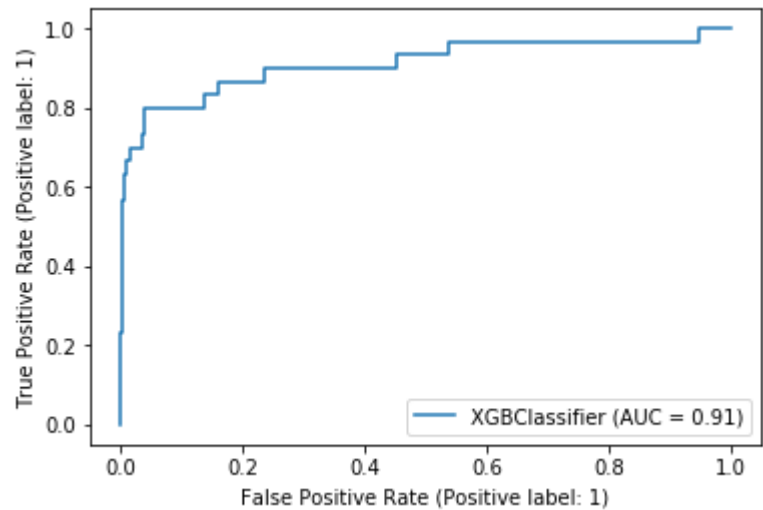
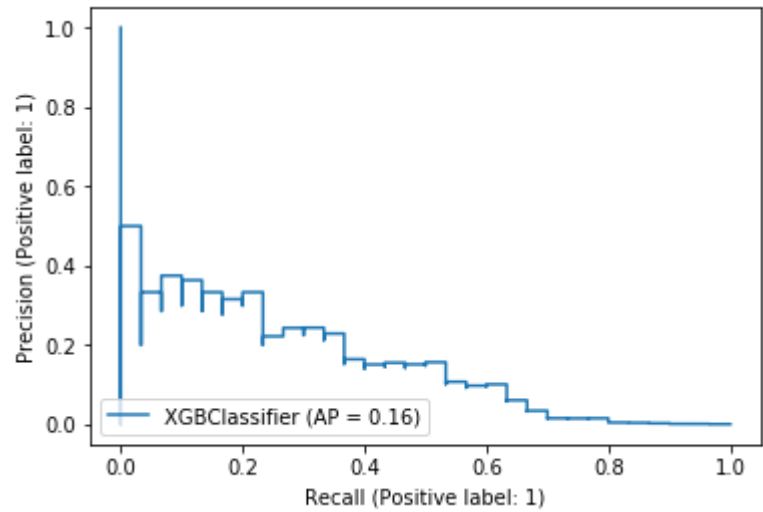
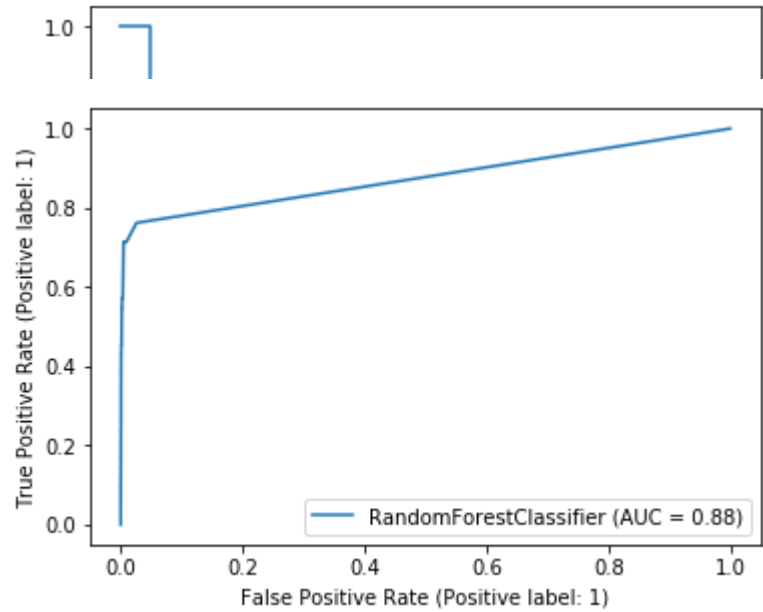


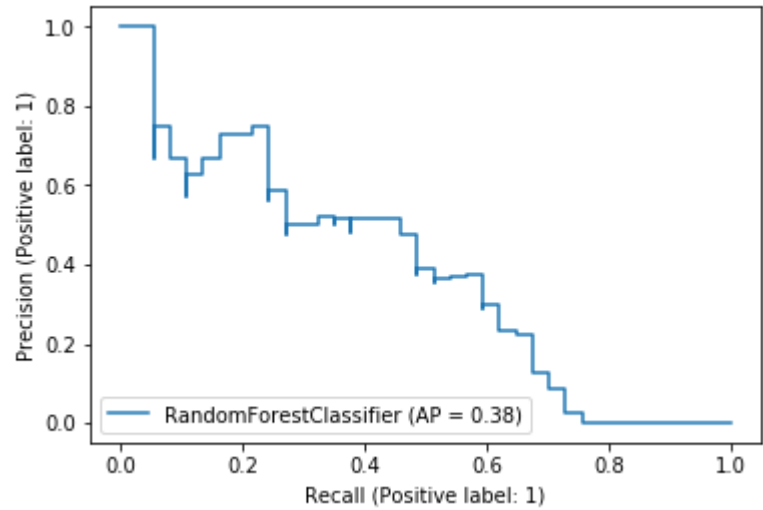
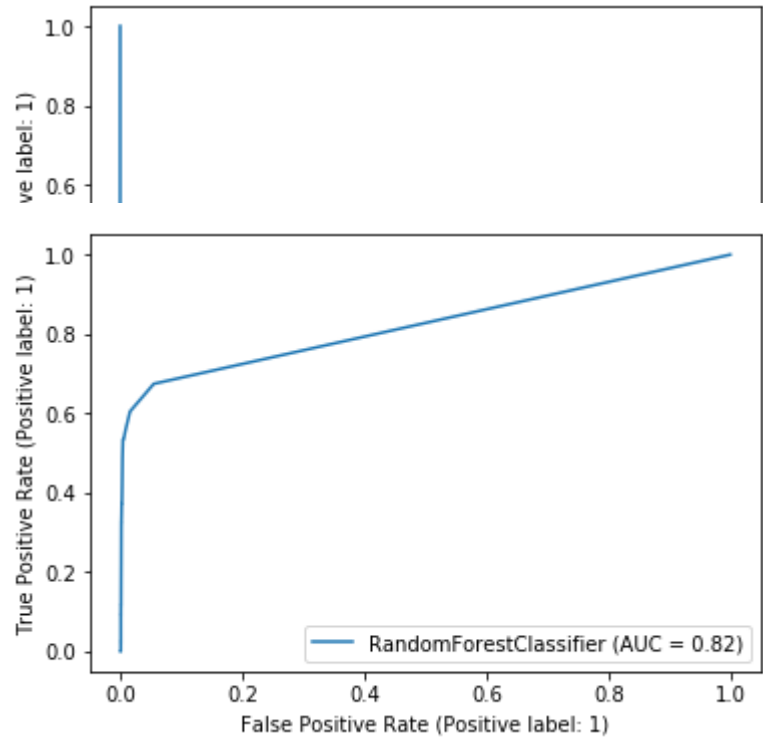


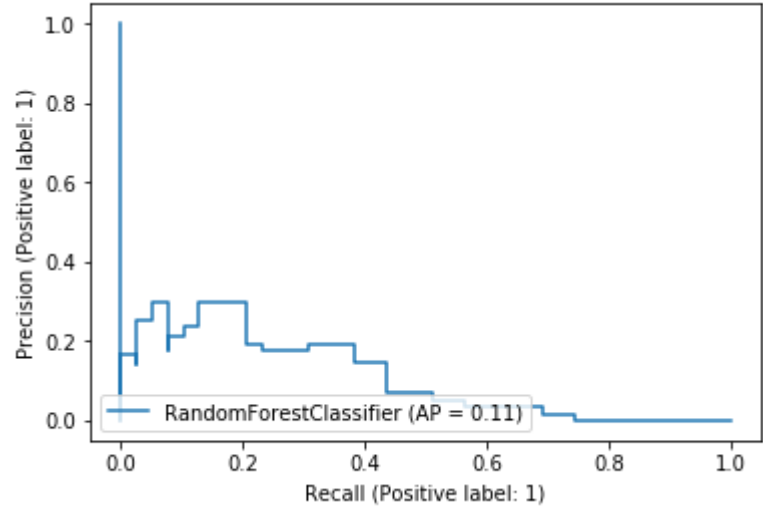
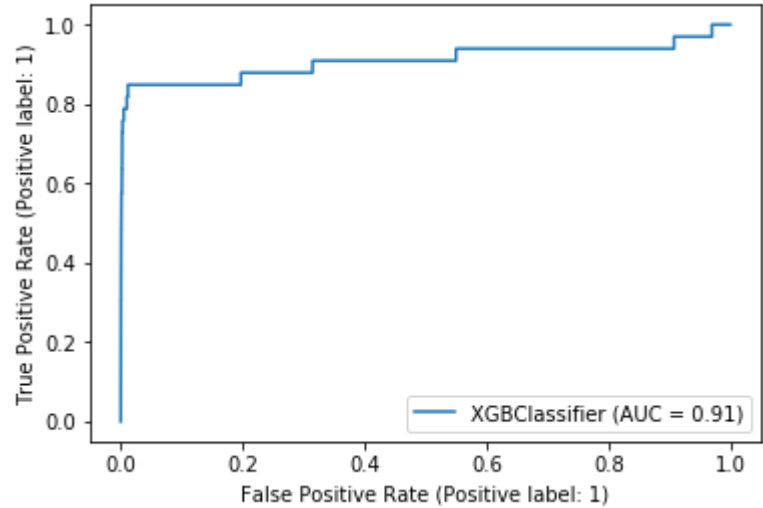
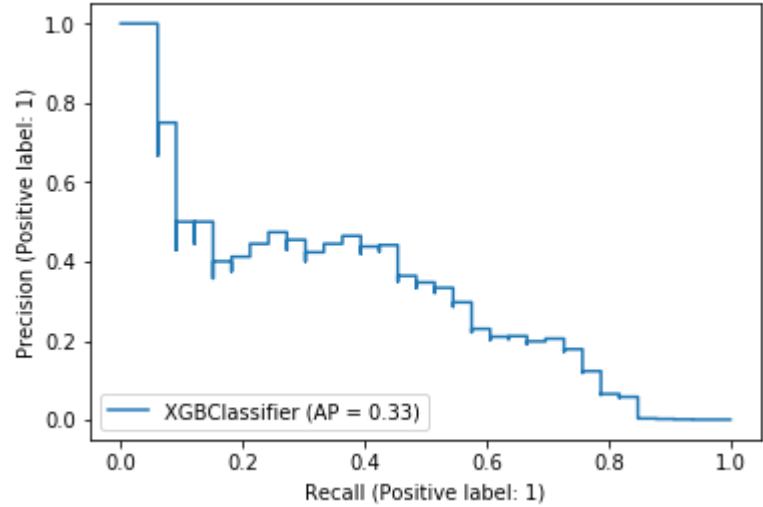
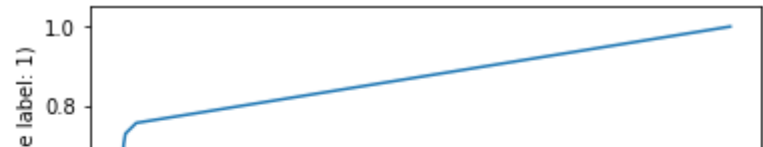


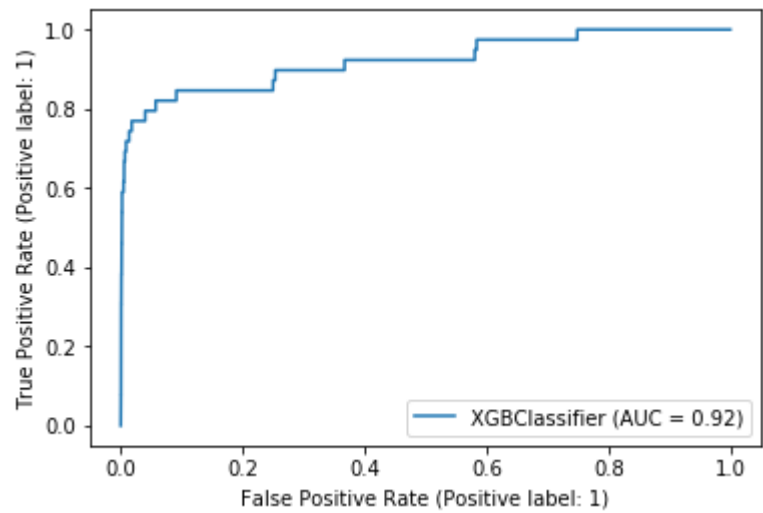
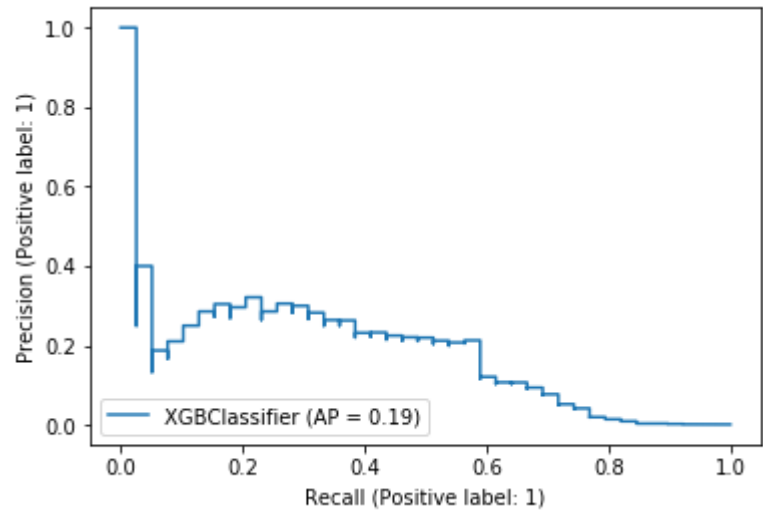
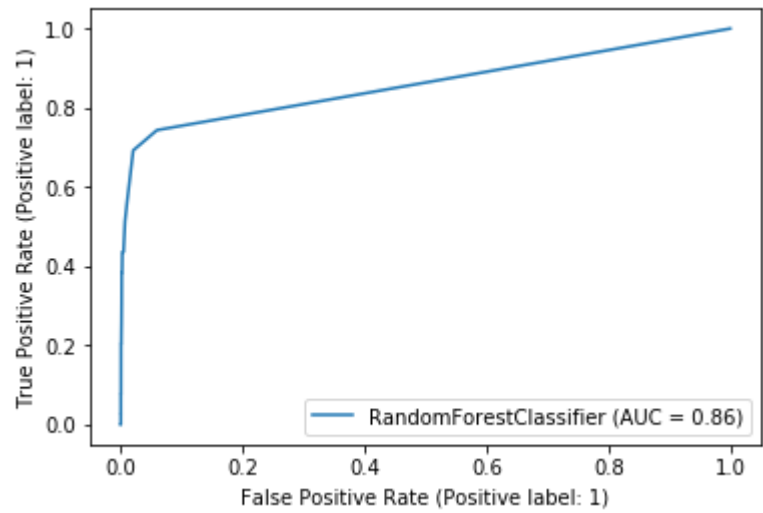


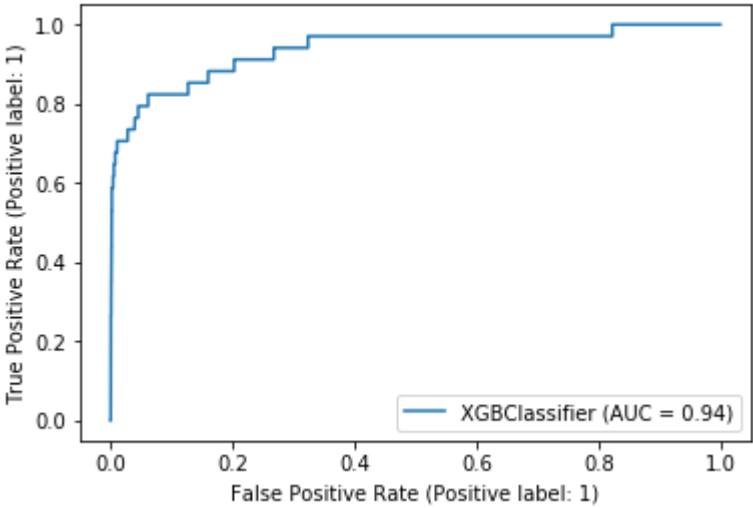
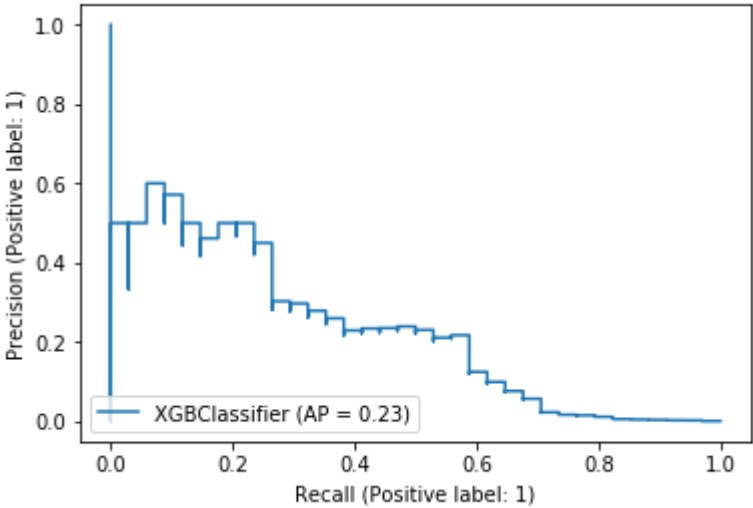
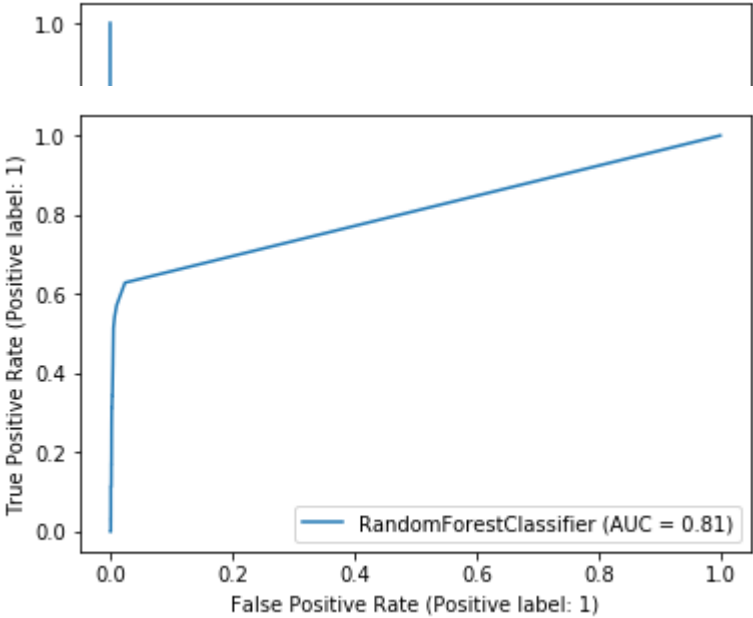


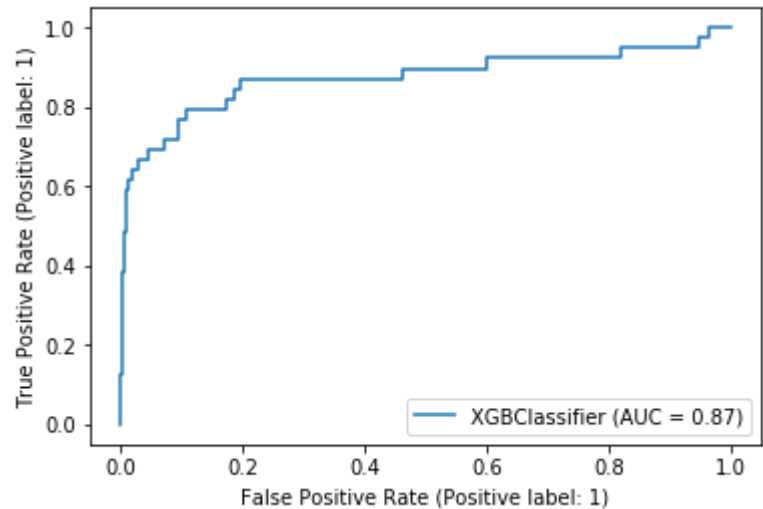
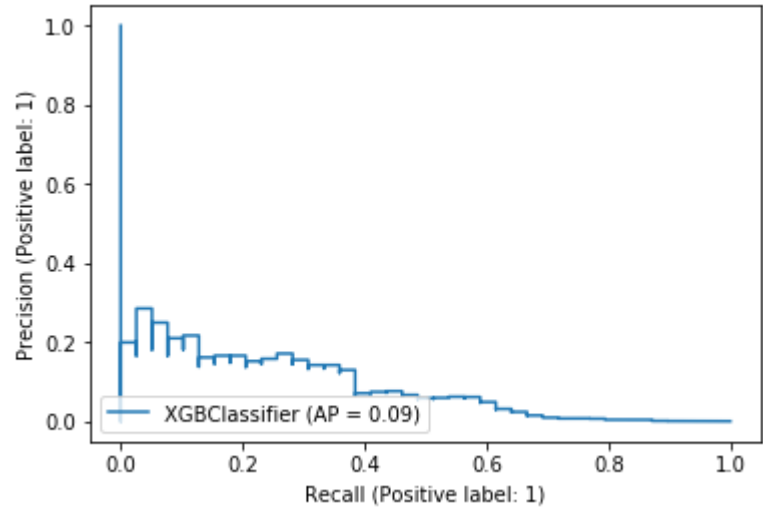
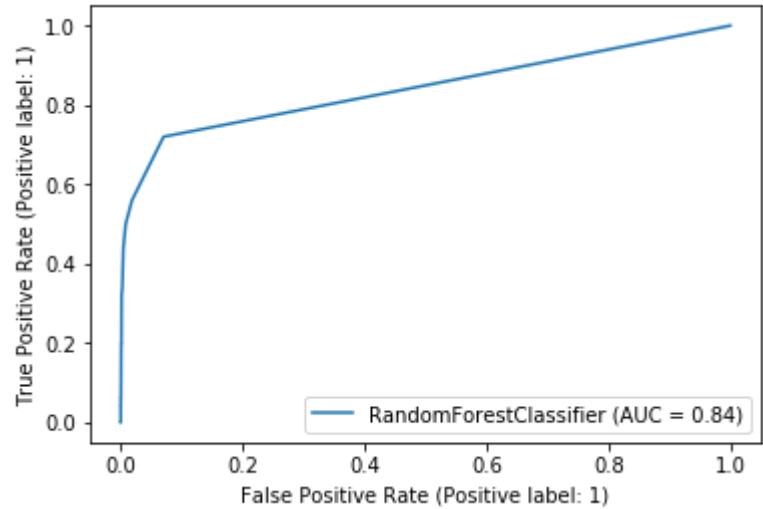


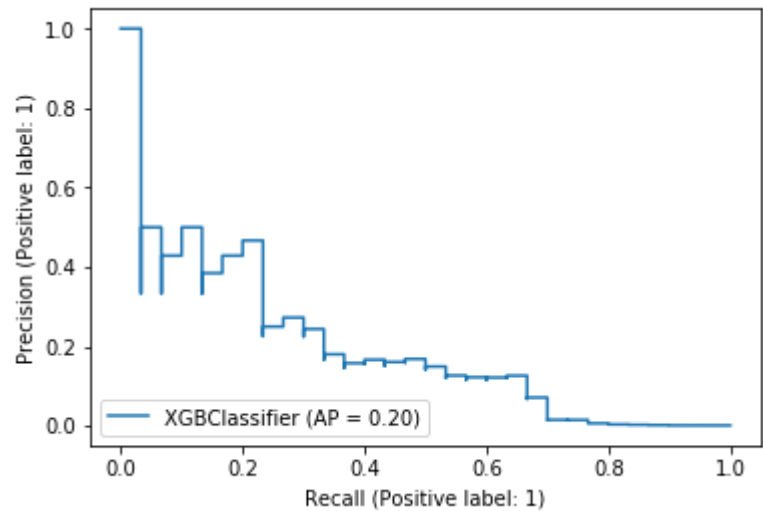
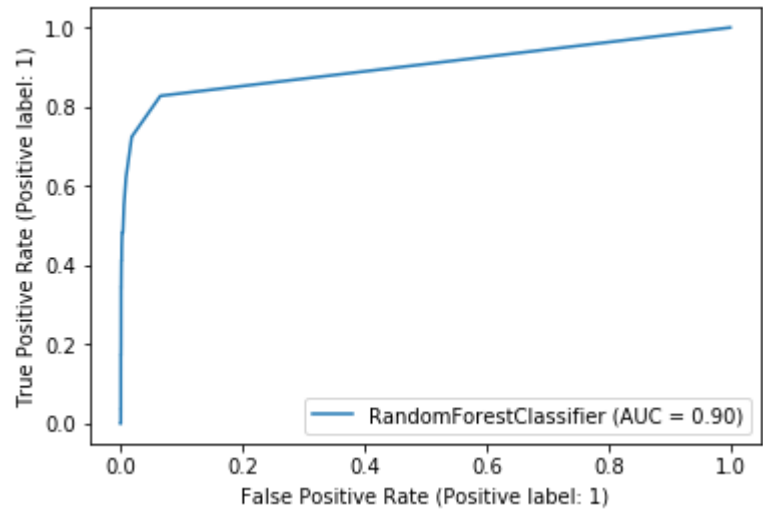
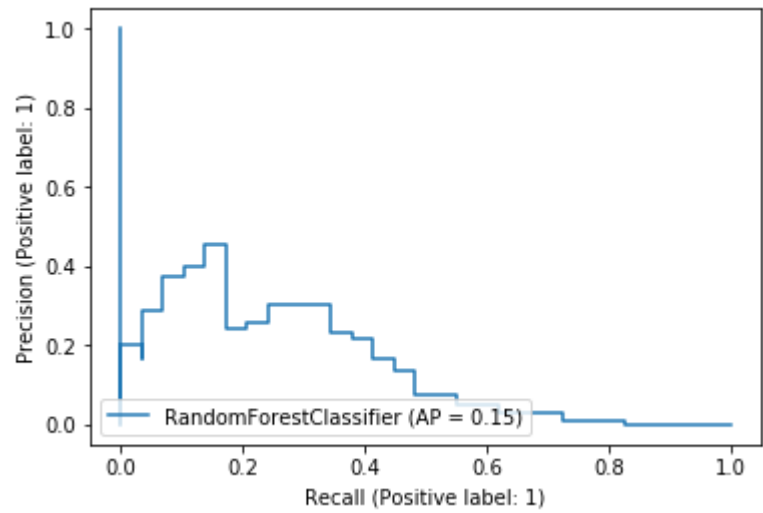


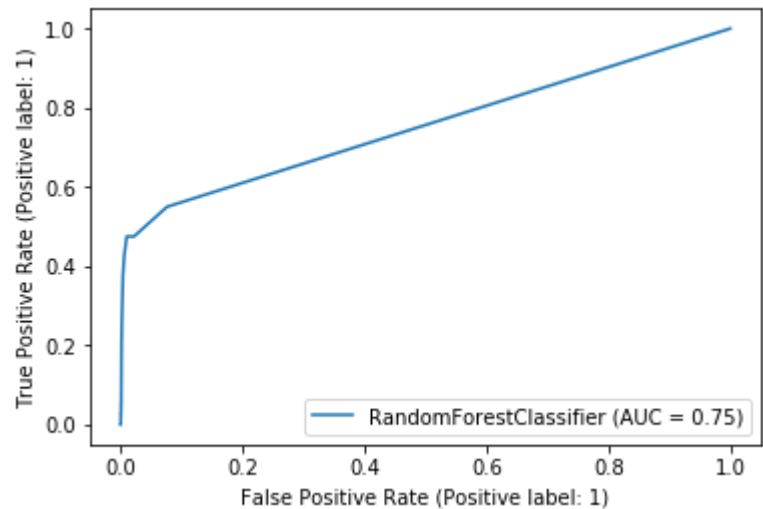
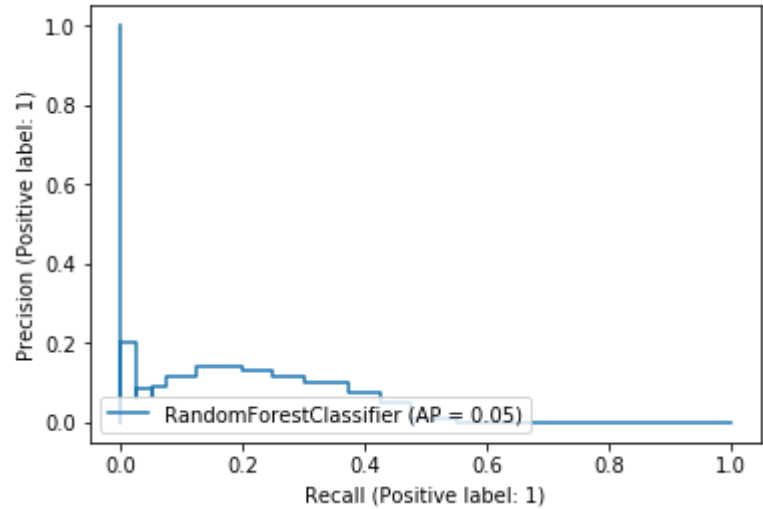
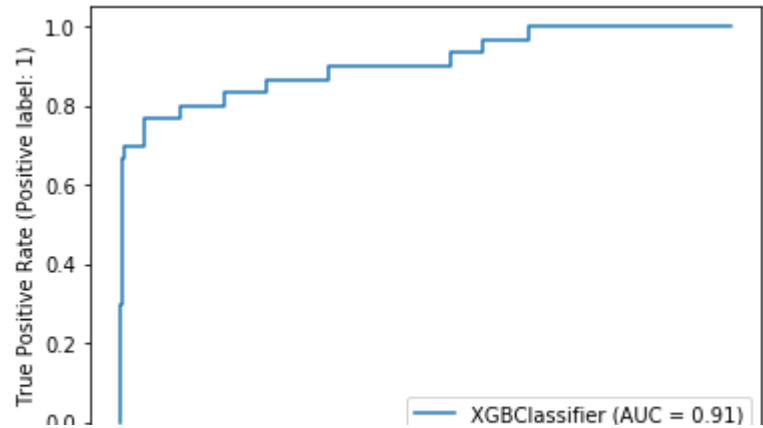


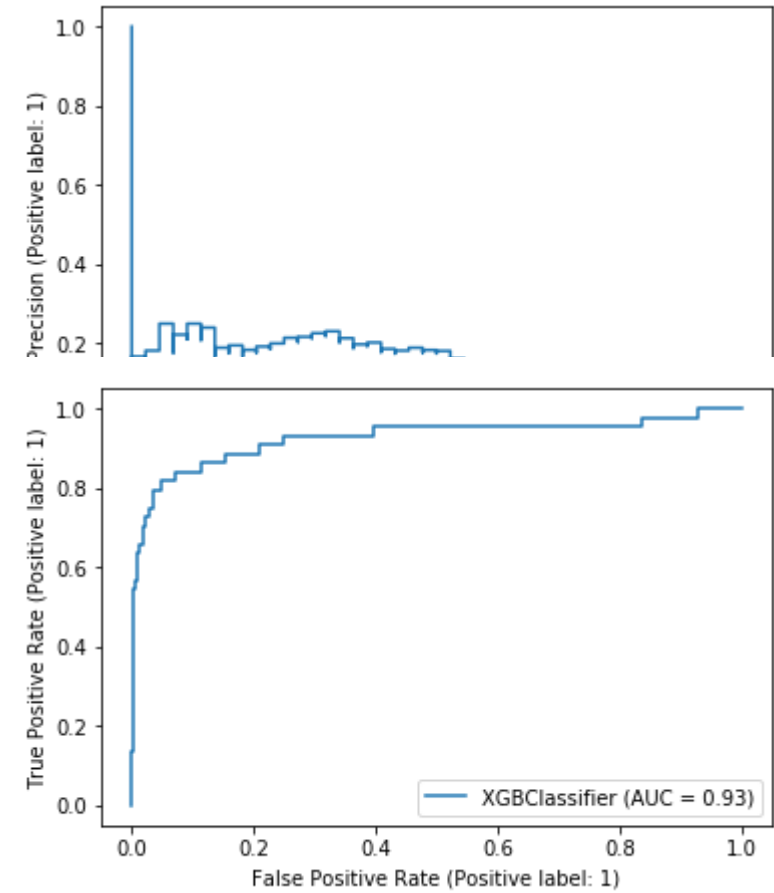












1	# Модель
---	----------

In []:

```
1  #модель
2  #!pip install xgboost
3
4  from sklearn.ensemble import RandomForestClassifier
5  from sklearn.model_selection import train_test_split
6  from xgboost import XGBClassifier
7
8  X = df_X.drop(columns = ['msisdn_x', 'msisdn_y', 'is_pair']).fillna(0)
9  y = df_X['is_pair']
10
11 X_train, X_test, y_train, y_test = train_test_split(
12     X, y, test_size=0.3, random_state=42)
13
14 #clf = RandomForestClassifier(random_state=0)
15 clf = XGBClassifier(random_state=42)
16
17 clf.fit(X_train, y_train)
18
19 y_test_predict = clf.predict(X_test)
20 y_train_predict = clf.predict(X_train)
21 y_test_predict_proba = clf.predict_proba(X_test)
22
23 from sklearn.metrics import accuracy_score
24 from sklearn.metrics import roc_auc_score
25 from sklearn.metrics import average_precision_score
26
27 print('roc_auc_score:', roc_auc_score(y_train, y_train_predict),
28       'average_precision_score:', average_precision_score(y_train, y_train_predict))
29 print("roc_auc_score:", roc_auc_score(y_test, y_test_predict),
30       'average_precision_score:', average_precision_score(y_test, y_test_predict))
31
32 from sklearn.metrics import roc_curve
33 from sklearn.metrics import plot_roc_curve
34 from sklearn.metrics import precision_recall_curve
35 from sklearn.metrics import plot_precision_recall_curve
36
37 disp = plot_precision_recall_curve(clf, X_test, y_test)
38 from sklearn.metrics import average_precision_score
39 average_precision = average_precision_score(y_test, y_test_predict)
40 print(average_precision)
41
42 disp = plot_roc_curve(clf, X_test, y_test)
43 roc_auc = roc_auc_score(y_test, y_test_predict)
44 print(roc_auc)
```

Смотрим, где неправильно предсказала модель

In []:

```

1 Cls_result = pd.DataFrame(y_test)
2 Cls_result['proba'] = y_test_predict_proba[:,1]
3 Cls_result['predict'] = y_test_predict
4 Cls_result = Cls_result.join(df_X[['msisdn_x', 'msisdn_y']])
5
6 alpha = .1
7
8 Cls_result['predict_alpha'] = np.where(Cls_result['proba'] > alpha, 1, 0)
9 display(
10 Cls_result.groupby(['is_pair', 'predict_alpha']).agg({'msisdn_x': 'count', 'proba': 'mean'})
11 Cls_result.groupby(['is_pair', 'predict']).agg({'msisdn_x': 'count', 'proba': 'mean'}))

```

In []:

```

1 a, b, is_pair, proba = Cls_result[(Cls_result['is_pair'] == 1) &
2                               (Cls_result['predict_alpha'] == 0)].sample(n = 1).iloc[0][[
3 #a, b, is_pair, proba = 158503476701, 158503536671, 0, 0
4 print('msisdn_1, msisdn_2', a, b, 'is_pair: ', is_pair, 'proba:', proba )
5 fig = plt.figure(figsize=(15, 4))
6 show_chart(fig, df_data[df_data['msisdn'] == a], df_data[df_data['msisdn'] == b])
7 plt.show()
8
9
10 m = folium.Map(location=(55.7522200, 37.6155600))
11 show_circles_on_map(m, df_data[df_data['msisdn'] == a ], "lat", "long", "blue")
12 show_circles_on_map(m, df_data[df_data['msisdn'] == b ], "lat", "long", "red")
13

```

In [38]:

```
1 Cls_result_agg_copy = Cls_result_agg.copy()
```

In [43]:

```
1 Cls_result_agg_2_test[(Cls_result_agg_2_test['is_pair'] == 1) & (Cls_result_agg_2_test
```

Out[43]:

	is_pair	msisdn_x	predict	predict_alpha	proba
2	1	15	NaN	0.0	0.027333
2	1	16	NaN	0.0	0.028750
2	1	10	NaN	0.0	0.041000
2	1	14	NaN	0.0	0.027857
2	1	18	NaN	0.0	0.023889

In [39]:

```

1 Cls_result_agg = pd.DataFrame()
2 BINS_QNT = 90
3 ln_ind = True
4 shift_inf = False
5 cls_ind = 'RF'
6 dist_ind_ln, coord_ind_ln, center_dist_ind_ln = True, False, False
7 dist_ind_shift, coord_ind_shift, center_dist_ind_shift = True, False, False
8 clfs = []
9
10 for bins_q, ln_ind, shift_ind, clf_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln,
11     [30, True, False, 'RF', True, False, False, True, False, False],
12     [60, True, False, 'RF', True, False, False, True, False, False],
13     [90, True, False, 'RF', True, False, False, True, False, False],
14     [120, True, False, 'RF', True, False, False, True, False, False],
15     [150, True, False, 'RF', True, False, False, True, False, False],
16     ]:
17
18     df_data_bin_ln, bins_list_ln = CreateTimeBeans(df_data, BINS_QNT)
19     df_data_pivot_base_ln = CreatePivotBase(df_data_bin_ln)
20     df_data_bin_agg_ln = AggData(df_data_bin_ln)
21     df_data_agg_for_pivot_ln = LinearCombData(df_data_pivot_base_ln, df_data_bin_agg_ln)
22
23     df_data_shift = AddDataDelta(df_data)
24     df_data_bin_shift, bins_list_shift = CreateTimeBeans(df_data_shift, BINS_QNT)
25     df_data_pivot_base_shift = CreatePivotBase(df_data_bin_shift)
26     df_data_bin_agg_shift = AggData(df_data_bin_shift)
27     df_data_agg_for_pivot_shift = CombData(df_data_pivot_base_shift, df_data_bin_agg_shift)
28
29     df_msisdn_pairs = CreateMsisdnPairs(facts, df_data, 0.25)
30
31     df_data_agg_pivot_ln = CreatePivot(df_data_agg_for_pivot_ln)
32     df_data_agg_pivot_shift = CreatePivot(df_data_agg_for_pivot_shift)
33
34     df_X_ln = CreateAndProcessPivot(df_msisdn_pairs, df_data_agg_pivot_ln, bins_list_ln)
35     df_X_shift = CreateAndProcessPivot(df_msisdn_pairs, df_data_agg_pivot_shift, bins_list_shift)
36
37     df_X = pd.DataFrame()
38     if (ln_ind):
39         df_X = pd.concat([df_X_ln, df_X], axis = 1)
40     if (shift_ind):
41         df_X = pd.concat([df_X_shift, df_X], axis = 1)
42
43     df_X.reset_index(inplace = True)
44
45     X = df_X.drop(columns = ['msisdn_x', 'msisdn_y', 'is_pair']).fillna(0)
46     y = df_X['is_pair']
47
48     X_train, X_test, y_train, y_test = train_test_split(
49         X, y, test_size=0.3, random_state=42)
50     if(clf_ind == 'RF'):
51         clf = RandomForestClassifier(random_state=0)
52     else:
53         clf = XGBClassifier(random_state=42)
54
55     clf.fit(X_train, y_train)
56
57     y_test_predict = clf.predict(X_test)
58     y_train_predict = clf.predict(X_train)
59     y_test_predict_proba = clf.predict_proba(X_test)

```

```

60
61 print('[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_
62       [bins_q, ln_ind, shift_inf, cls_ind,
63       dist_ind_ln, coord_ind_ln, center_dist_ind_ln,
64       dist_ind_shift, coord_ind_shift, center_dist_ind_shift])
65
66 print('roc_auc_score:', roc_auc_score(y_train, y_train_predict),
67       'average_precision_score:', average_precision_score(y_train, y_train_predict)
68 print("roc_auc_score:", roc_auc_score(y_test, y_test_predict),
69       'average_precision_score:', average_precision_score(y_test, y_test_predict))
70
71
72 disp = plot_precision_recall_curve(clf, X_test, y_test)
73 average_precision = average_precision_score(y_test, y_test_predict)
74 print(average_precision)
75
76 disp = plot_roc_curve(clf, X_test, y_test)
77 roc_auc = roc_auc_score(y_test, y_test_predict)
78 print(roc_auc)
79
80
81 Cls_result = pd.DataFrame(y_test)
82 Cls_result['proba'] = y_test_predict_proba[:,1]
83 Cls_result['predict'] = y_test_predict
84 Cls_result = Cls_result.join(df_X[['msisd_n_x', 'msisd_n_y']])
85
86 alpha = .1
87
88 Cls_result['predict_alpha'] = np.where(Cls_result['proba'] > alpha, 1, 0)
89 Cls_result_agg = Cls_result_agg.append(Cls_result.groupby(['is_pair', 'predict_alpha']).agg('max'))
90 Cls_result_agg = Cls_result_agg.append(Cls_result.groupby(['is_pair', 'predict']).agg('max'))
91

```

```

(129324, 3)
(127199, 91)
(127199, 91)
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_
ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [30, True, F
alse, 'RF', True, False, False, True, False, False]
roc_auc_score: 1.0 average_precision_score: 1.0
roc_auc_score: 0.6153714992543985 average_precision_score: 0.208478471214320
28
0.20847847121432028
0.6153714992543985
(129326, 3)
(127214, 91)
(127214, 91)
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_
ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [60, True, F
alse, 'RF', True, False, False, True, False, False]
roc_auc_score: 1.0 average_precision_score: 1.0
roc_auc_score: 0.6361800635114387 average_precision_score: 0.107348216010261
25
0.10734821601026125
0.6361800635114387
(129301, 3)
(127169, 91)
(127169, 91)
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_

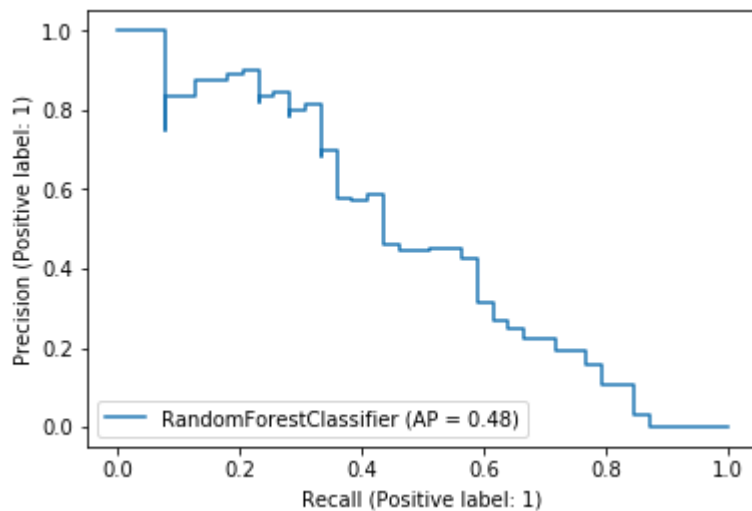
```

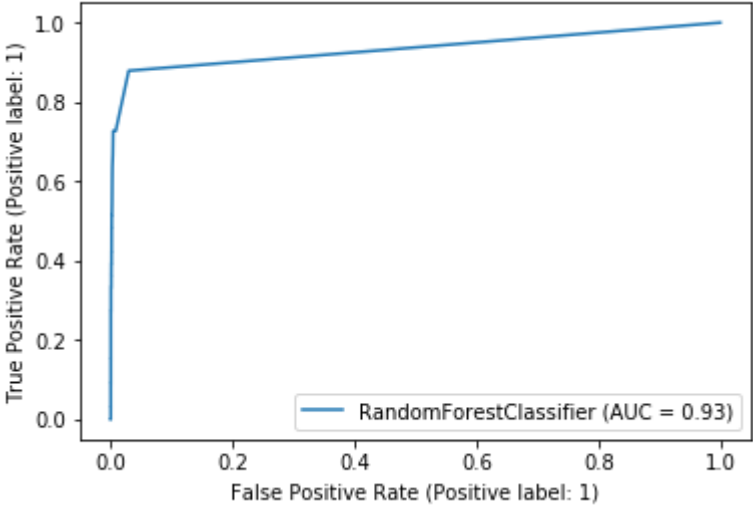
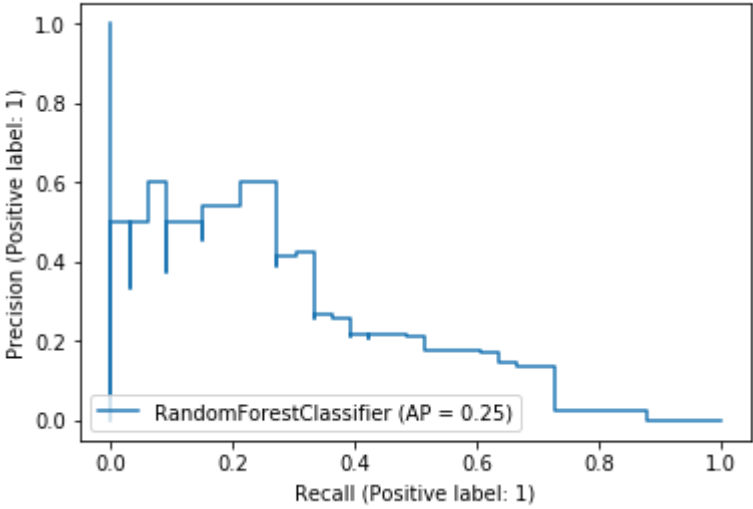
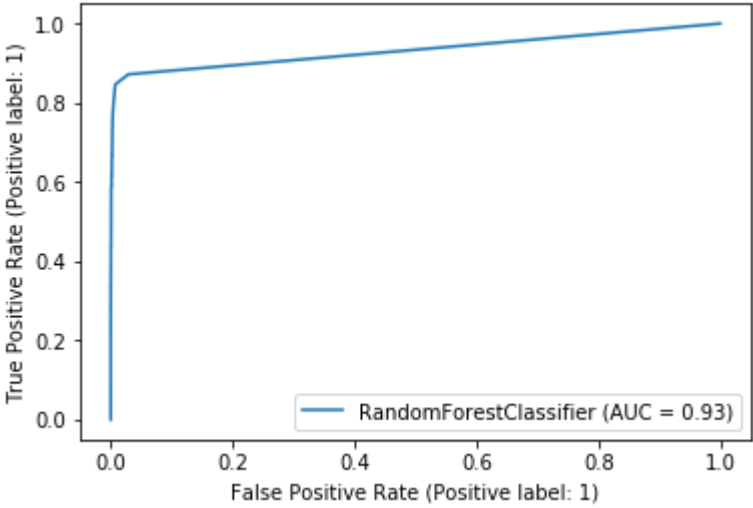


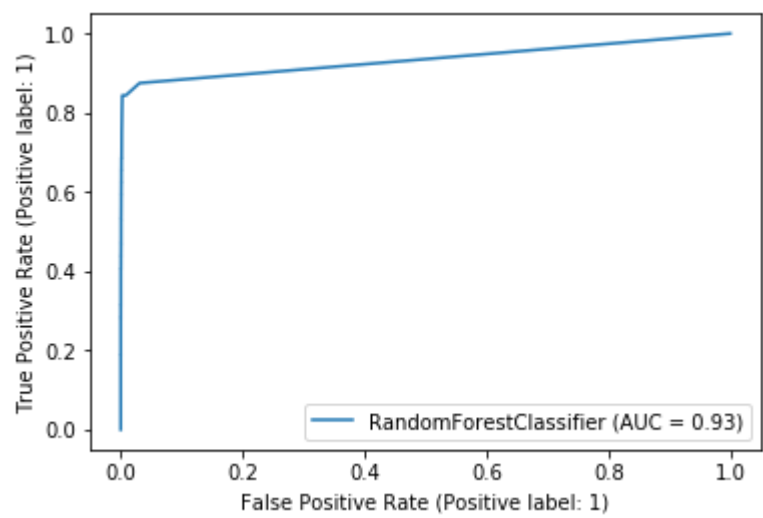
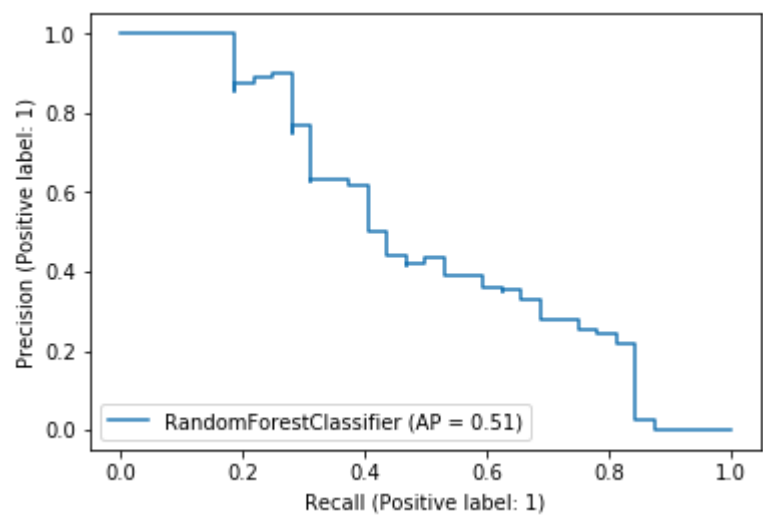
```

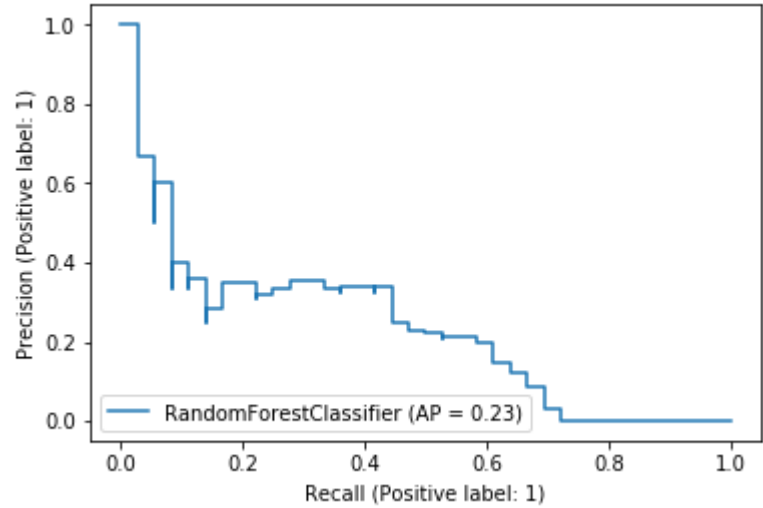
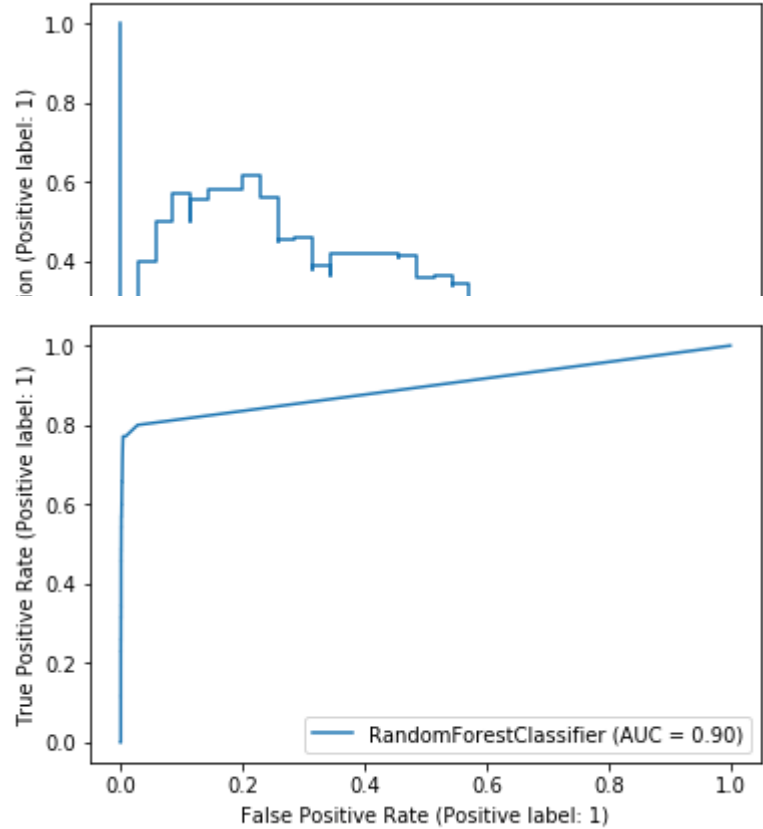
ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [90, True, F
alse, 'RF', True, False, False, True, False, False]
roc_auc_score: 1.0 average_precision_score: 1.0
roc_auc_score: 0.6249868831816154 average_precision_score: 0.222851301407564
68
0.22285130140756468
0.6249868831816154
(129289, 3)
(127148, 91)
(127148, 91)
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_
ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [120, True,
False, 'RF', True, False, False, True, False, False]
roc_auc_score: 1.0 average_precision_score: 1.0
roc_auc_score: 0.6142069947895191 average_precision_score: 0.131320070301026
4
0.1313200703010264
0.6142069947895191
(129307, 3)
(127165, 91)
(127165, 91)
[bins_q, ln_ind, shift_inf, cls_ind, dist_ind_ln, coord_ind_ln, center_dist_
ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [150, True,
False, 'RF', True, False, False, True, False, False]
roc_auc_score: 1.0 average_precision_score: 1.0
roc_auc_score: 0.5415879554319497 average_precision_score: 0.028642784330857
724
0.028642784330857724
0.5415879554319497

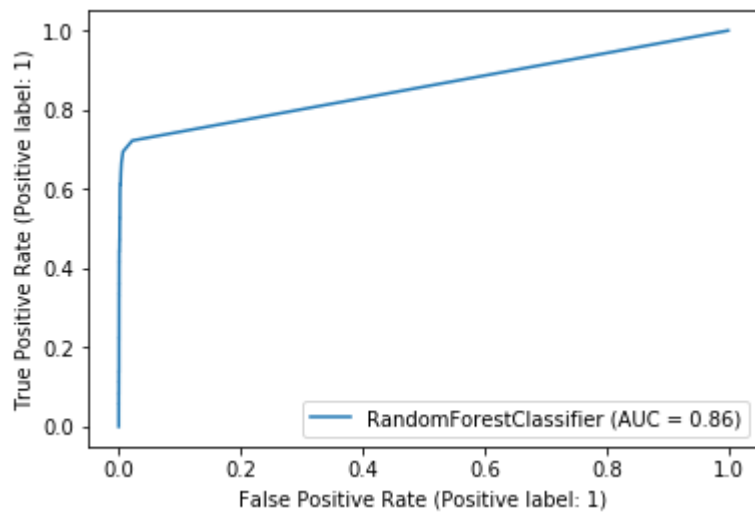
```











In [40]:

```
1 Cls_result_agg_2_test = Cls_result_agg.copy()
```

Использование модели

In [41]:

```

1 Cls_result_agg = pd.DataFrame()
2 BINS_QNT = 90
3 ln_ind = True
4 shift_inf = False
5 cls_ind = 'RF'
6 dist_ind_ln, coord_ind_ln, center_dist_ind_ln = True, False, False
7 dist_ind_shift, coord_ind_shift, center_dist_ind_shift = True, False, False
8 clfs = []
9
10 for bins_q, ln_ind, shift_ind, clf_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln,
11     [30, True, False, 'GF', True, False, False, True, False, False],
12     [60, True, False, 'GF', True, False, False, True, False, False],
13     [90, True, False, 'GF', True, False, False, True, False, False],
14     [120, True, False, 'GF', True, False, False, True, False, False],
15     [150, True, False, 'GF', True, False, False, True, False, False],
16     ]:
17
18     df_data_bin_ln, bins_list_ln = CreateTimeBeans(df_data, BINS_QNT)
19     df_data_pivot_base_ln = CreatePivotBase(df_data_bin_ln)
20     df_data_bin_agg_ln = AggData(df_data_bin_ln)
21     df_data_agg_for_pivot_ln = LinearCombData(df_data_pivot_base_ln, df_data_bin_agg_ln)
22
23     df_data_shift = AddDataDelta(df_data)
24     df_data_bin_shift, bins_list_shift = CreateTimeBeans(df_data_shift, BINS_QNT)
25     df_data_pivot_base_shift = CreatePivotBase(df_data_bin_shift)
26     df_data_bin_agg_shift = AggData(df_data_bin_shift)
27     df_data_agg_for_pivot_shift = CombData(df_data_pivot_base_shift, df_data_bin_agg_shift)
28
29     df_msisdn_pairs = CreateMsisdnPairs(facts, df_data, 0.25)
30
31     df_data_agg_pivot_ln = CreatePivot(df_data_agg_for_pivot_ln)
32     df_data_agg_pivot_shift = CreatePivot(df_data_agg_for_pivot_shift)
33
34     df_X_ln = CreateAndProcessPivot(df_msisdn_pairs, df_data_agg_pivot_ln, bins_list_ln)
35     df_X_shift = CreateAndProcessPivot(df_msisdn_pairs, df_data_agg_pivot_shift, bins_list_shift)
36
37     df_X = pd.DataFrame()
38     if (ln_ind):
39         df_X = pd.concat([df_X_ln, df_X], axis = 1)
40     if (shift_ind):
41         df_X = pd.concat([df_X_shift, df_X], axis = 1)
42
43     df_X.reset_index(inplace = True)
44
45     X = df_X.drop(columns = ['msisdn_x', 'msisdn_y', 'is_pair']).fillna(0)
46     y = df_X['is_pair']
47
48     X_train, X_test, y_train, y_test = train_test_split(
49         X, y, test_size=0.3, random_state=42)
50     if(clf_ind == 'RF'):
51         clf = RandomForestClassifier(random_state=0)
52     else:
53         clf = XGBClassifier(random_state=42)
54
55     clf.fit(X_train, y_train)
56
57     y_test_predict = clf.predict(X_test)
58     y_train_predict = clf.predict(X_train)
59     y_test_predict_proba = clf.predict_proba(X_test)

```

```

60
61 print('[bins_q, ln_ind, shift_inf, clf_ind, dist_ind_ln, coord_ind_ln, center_dist_
62       [bins_q, ln_ind, shift_inf, cls_ind,
63       dist_ind_ln, coord_ind_ln, center_dist_ind_ln,
64       dist_ind_shift, coord_ind_shift, center_dist_ind_shift])
65
66 print('roc_auc_score:', roc_auc_score(y_train, y_train_predict),
67       'average_precision_score:', average_precision_score(y_train, y_train_predict)
68 print("roc_auc_score:", roc_auc_score(y_test, y_test_predict),
69       'average_precision_score:', average_precision_score(y_test, y_test_predict))
70
71
72 disp = plot_precision_recall_curve(clf, X_test, y_test)
73 average_precision = average_precision_score(y_test, y_test_predict)
74 print(average_precision)
75
76 disp = plot_roc_curve(clf, X_test, y_test)
77 roc_auc = roc_auc_score(y_test, y_test_predict)
78 print(roc_auc)
79
80
81 Cls_result = pd.DataFrame(y_test)
82 Cls_result['proba'] = y_test_predict_proba[:,1]
83 Cls_result['predict'] = y_test_predict
84 Cls_result = Cls_result.join(df_X[['msisdn_x', 'msisdn_y']])
85
86 alpha = .1
87
88 Cls_result['predict_alpha'] = np.where(Cls_result['proba'] > alpha, 1, 0)
89 Cls_result_agg = Cls_result_agg.append(Cls_result.groupby(['is_pair', 'predict_alpha']).agg('count'))
90 Cls_result_agg = Cls_result_agg.append(Cls_result.groupby(['is_pair', 'predict']).agg('count'))
91

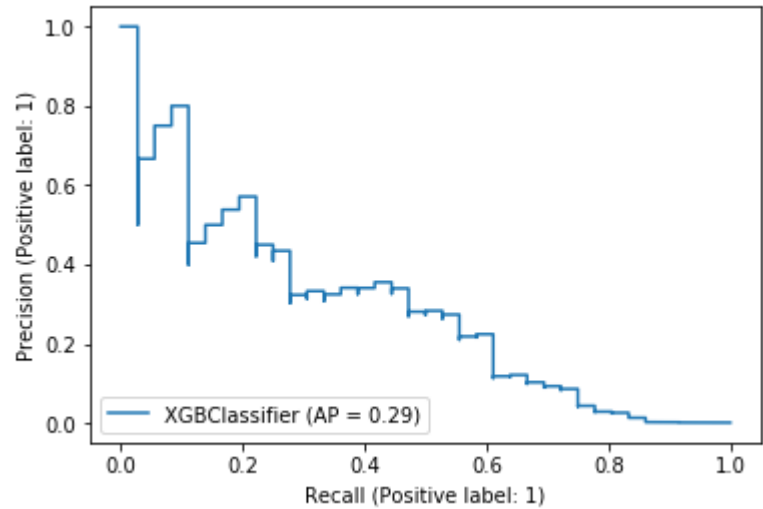
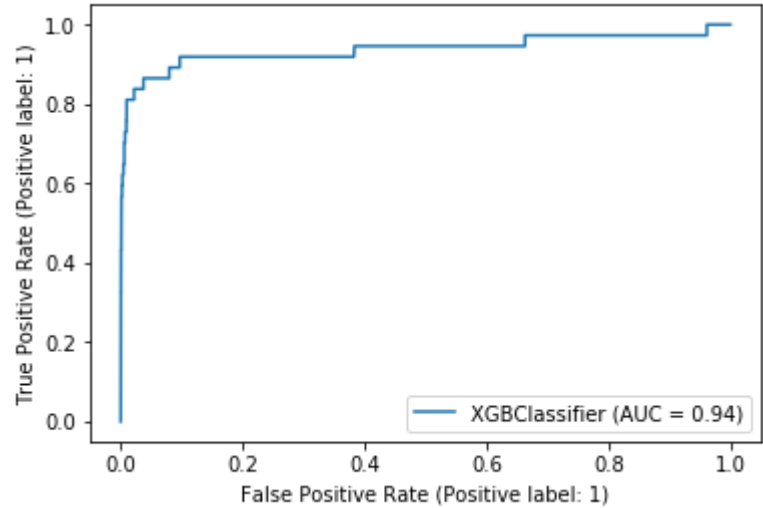
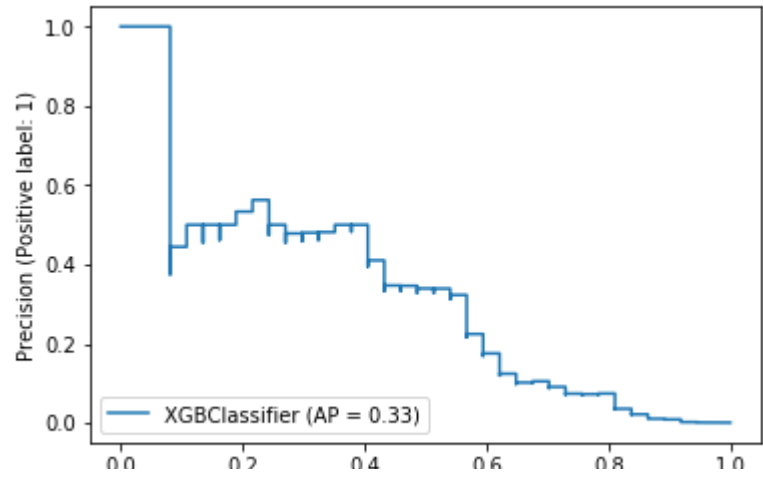
```

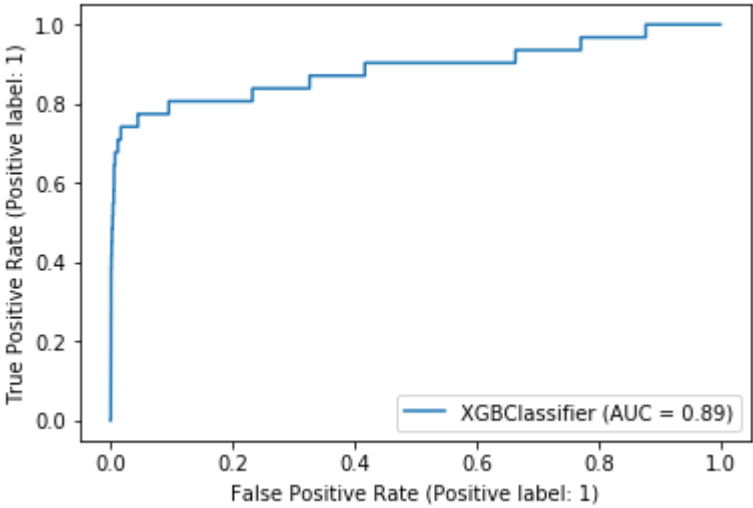
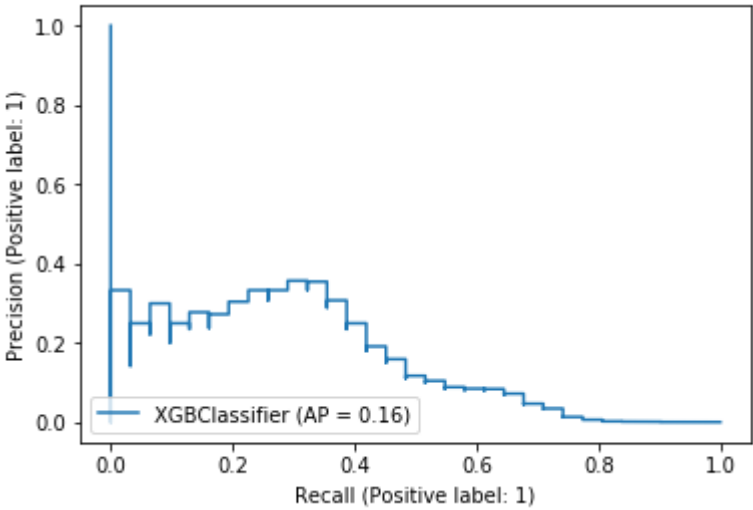
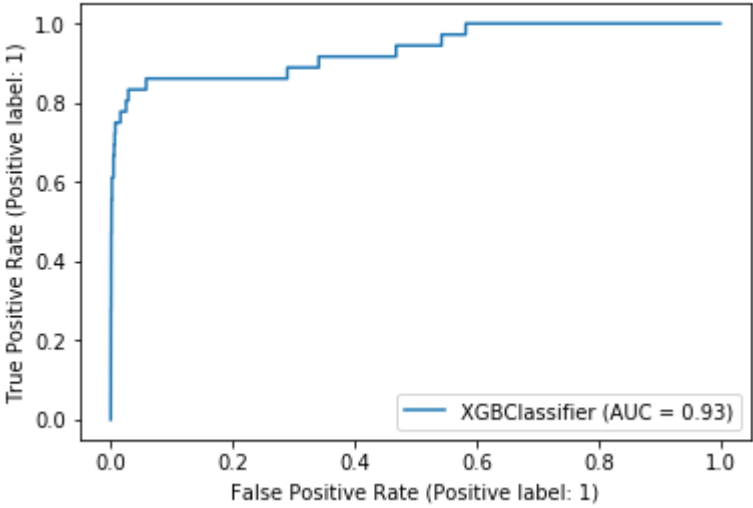
```

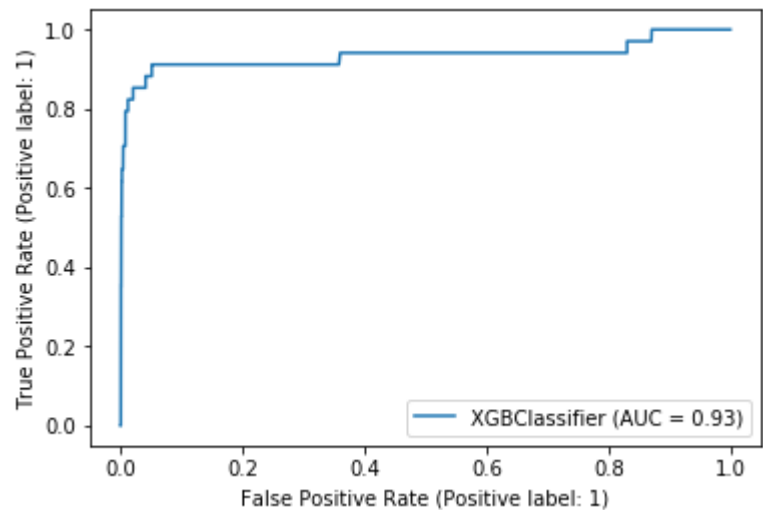
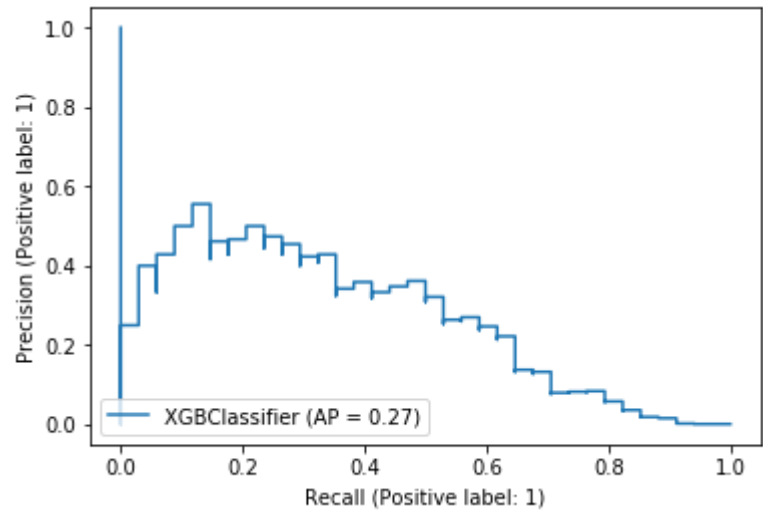
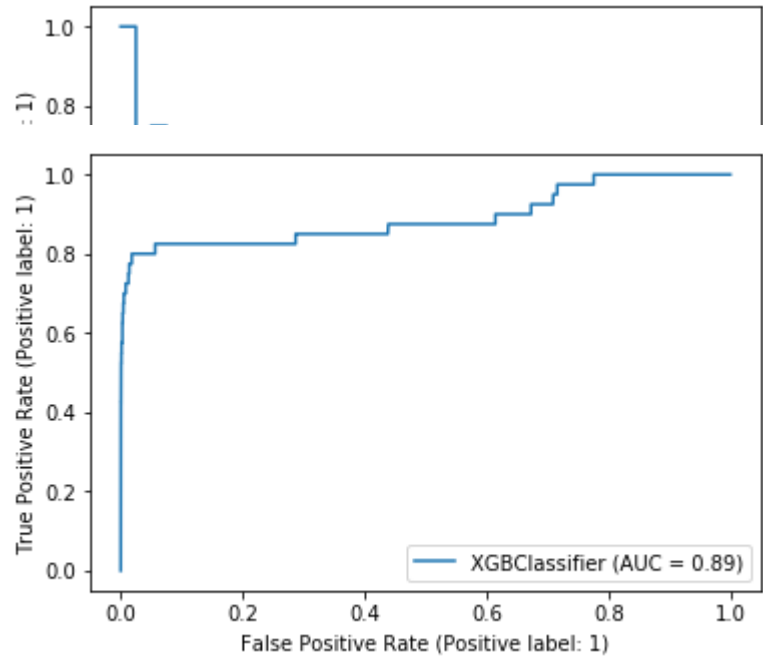
(129308, 3)
(127100, 91)
(127100, 91)
[12:26:08] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release
_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluat
ion metric used with the objective 'binary:logistic' was changed from 'err
or' to 'logloss'. Explicitly set eval_metric if you'd like to restore the
old behavior.
[bins_q, ln_ind, shift_inf, clf_ind, dist_ind_ln, coord_ind_ln, center_dis
t_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [30, Tru
e, False, 'RF', True, False, False, True, False, False]
roc_auc_score: 0.9878048780487805 average_precision_score: 0.9756322355850
287
roc_auc_score: 0.6215297412236482 average_precision_score: 0.1375586542482
6874
0.13755865424826874
0.6215297412236482
(129315, 3)
(127199, 91)
(127199, 91)
[12:27:12] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release
_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluat
ion metric used with the objective 'binary:logistic' was changed from 'err
or' to 'logloss'. Explicitly set eval_metric if you'd like to restore the
old behavior.
[bins_q, ln_ind, shift_inf, clf_ind, dist_ind_ln, coord_ind_ln, center_dis

```

```
t_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [60, True, False, 'RF', True, False, False, True, False, False]
roc_auc_score: 0.9879518072289157 average_precision_score: 0.9759260765250154
roc_auc_score: 0.5554899800650509 average_precision_score: 0.0502212904728628
0.0502212904728628
0.5554899800650509
(129312, 3)
(127089, 91)
(127089, 91)
[12:28:19] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[bins_q, ln_ind, shift_inf, clf_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [90, True, False, 'RF', True, False, False, True, False, False]
roc_auc_score: 0.9886363636363636 average_precision_score: 0.9772952087816862
roc_auc_score: 0.5482558493991411 average_precision_score: 0.023066893907725525
0.023066893907725525
0.5482558493991411
(129314, 3)
(127161, 91)
(127161, 91)
[12:29:28] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[bins_q, ln_ind, shift_inf, clf_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [120, True, False, 'RF', True, False, False, True, False, False]
roc_auc_score: 0.9936708860759493 average_precision_score: 0.9873530065921989
roc_auc_score: 0.6123950379175523 average_precision_score: 0.11993025027253819
0.11993025027253819
0.6123950379175523
(129322, 3)
(127138, 91)
(127138, 91)
[12:30:33] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.3.0/src/learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
[bins_q, ln_ind, shift_inf, clf_ind, dist_ind_ln, coord_ind_ln, center_dist_ind_ln, dist_ind_shift, coord_ind_shift, center_dist_ind_shift] [150, True, False, 'RF', True, False, False, True, False, False]
roc_auc_score: 0.9941176470588236 average_precision_score: 0.9882465305777127
roc_auc_score: 0.6028362116338085 average_precision_score: 0.09678631244852824
0.09678631244852824
0.6028362116338085
```





Type *Markdown* and LaTeX: α^2

