

In [1]:

```
import pandas as pd
import numpy as np
from datetime import datetime
from tqdm import tqdm
import matplotlib.pyplot as plt
```

In [2]:

```
data = pd.read_csv('training.csv')
data.head()
```

Out[2]:

	customer_id	visit_date	visit_spend
0	2	2010-04-01	5.97
1	2	2010-04-06	12.71
2	2	2010-04-07	34.52
3	2	2010-04-12	7.89
4	2	2010-04-14	17.17

In []:

In []:

In [7]:

```
def TopHat(x):
    return 0.5*(abs(x)<10)
def KDE(kernel, w, width, observable_points, x):
    summ = 0
    for i in range(len(observable_points)):
        summ += w[i]*TopHat((x - observable_points[i])/width)
    return summ
```

In [8]:

```
start_train = datetime.strptime('2010-04-05', '%Y-%m-%d').date()
end_train = datetime.strptime('2011-05-29', '%Y-%m-%d').date()
```

In [13]:

```
customers = sorted(np.unique(data['customer_id']))
```

In [14]:

```
data_train = data[(data['visit_date'] < '2011-05-30') & (data['visit_date'] >= '2010-04-05')  
data_test = data[data['visit_date'] >= '2011-05-30']
```

За y_{true} взята сумма покупки в предполагаемый первый день визита из исследования Павла Лукьяненко.

In [16]:

```
y_true = []
y1_true = []
for client in tqdm(customers):
    cur_data = data_test[(data_test['customer_id'] == client)]
    visit = np.array(cur_data[['visit_date', 'visit_spend']]).T
    if visit.shape[1] > 0:
        minind = visit[0].argmin()
        first_visit = visit[0][minind]
        if (datetime.strptime(first_visit, '%Y-%m-%d').date() - end_train).days < 7:
            y_true.append(datetime.strptime(first_visit, '%Y-%m-%d').date().weekday())
            y1_true.append(visit[1][minind])
        else:
            y_true.append(-1)
            y1_true.append(-1)
    else:
        y_true.append(-1)
        y1_true.append(-1)
```

```
100%|███████████| 
| 100000/100000 [03:46<00:00, 441.56it/s]
```

Для каждого клиента создается два элемента:

1. Массив всех его сумм покупок за все время в train (далее отсортированные).
2. Суммы покупок в предсказываемый день недели.

В первом из лекции - это s_1, \dots, s_m , второе - $s'_1, \dots, s'_{m'}$

In [18]:

```
customers_spend_concat = []
for client in tqdm(range(len(customers))):
    cur_data = data_train[(data_train['customer_id'] == customers[client])]
    visit = np.array(cur_data[['visit_date', 'visit_spend']]).T
    m_array = []
    for date in range(len(visit[0])):
        if (datetime.strptime(visit[0][date], '%Y-%m-%d').date() - start_train).days % 7 == 0:
            m_array.append(visit[1][date])
    customers_spend_concat.append([visit[1], m_array])
```

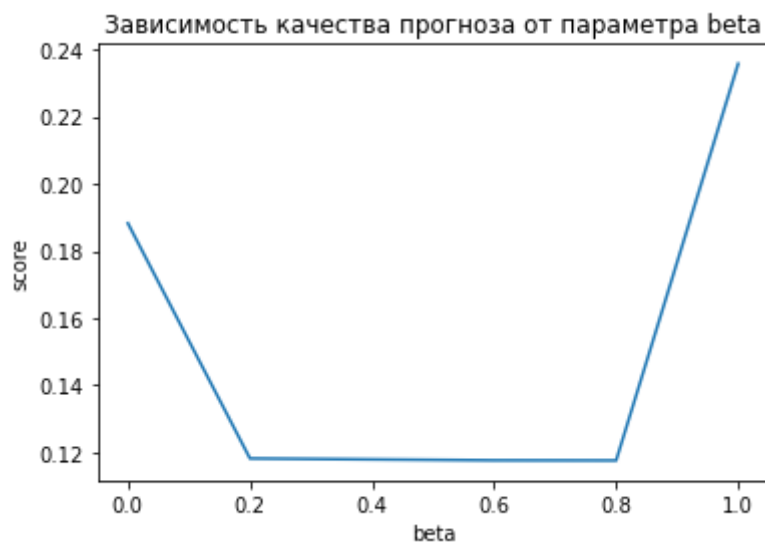
```
100%|███████████████████████████████████████████████████████████████████████████████  
██████████ | 100000/100000 [32:31<00:00, 51.25it/s]
```




1 0.23579

In [24]:

```
plt.plot(betas, res_beta)
plt.xlabel('beta')
plt.ylabel('score')
plt.title('Зависимость качества прогноза от параметра beta')
plt.show()
```



In [25]:

```
res_beta
```

Out[25]:

```
[0.18823, 0.11816, 0.1176, 0.11756, 0.23579]
```

Как видно из графика, в моем случае хорошо работает при β , близких к нулю и к единице, то есть при использовании только данных всех сумм покупок за весь период или только сумм покупок, совершенных в угадываемый день недели.

Могу предположить, что это происходит из-за выбора predict. Было бы неплохо увеличить тестовые данные и совместить результаты на по нескольким суммам по каждому клиенту, но выполненный код очень неоптимален и демонстрирует медленную работу, что не дает выполнить дополнительные исследования

In [24]:

In [25]:

In []:

In [26]: