

Это **вспомогательный ноутбук**, необходимый лишь для обработки всех данных и формирования общего файла с предсказаниями.

Никаких новых функций в нём нет.

```
In [3]: 1 import pandas as pd
        2 import numpy as np
        3 from tqdm import tqdm_notebook
        4 import matplotlib.pyplot as plt
        5 import seaborn as sns
        6 import warnings
        7 from sklearn.decomposition import PCA
        8 from sklearn.manifold import TSNE
        9 from umap.umap_ import UMAP
       10
       11 from sklearn.preprocessing import StandardScaler
       12
       13 from sklearn.cluster import \
       14     KMeans, \
       15     AgglomerativeClustering, \
       16     DBSCAN, \
       17     SpectralClustering
       18
       19 warnings.filterwarnings('ignore')
       20 sns.set_style('whitegrid')
       21 sns.set(font_scale=1.5)
```

Файл со всеми ~115000 сессиями.

```
In [4]: 1 all_data = pd.read_csv('data/processed_data.csv')
```

Генерация новых признаков.

```

In [5]: 1 def add_features(data):
2
3     new_data = data.copy()
4
5     new_data['delay'] = new_data['ts']
6     new_data.loc[1:, 'delay'] = (new_data['delay'].values[1:] -
7                                 new_data['delay'].values[:-1])
8
9     dist = np.sqrt((new_data['x'].values[1:] - new_data['x'].values[:-1])**2 +
10                  (new_data['y'].values[1:] - new_data['y'].values[:-1])**2)
11
12     new_data['dist'] = 0
13     new_data.loc[1:, 'dist'] = dist
14
15     new_data.loc[new_data['begin'] == 1, 'delay'] = 0
16     new_data.loc[new_data['begin'] == 1, 'av_speed'] = 0
17
18     new_data['av_speed'] = new_data['dist'] / new_data['delay']
19
20     new_data['x_diff'] = new_data.x
21     new_data.loc[1:, 'x_diff'] = new_data.x.values[1:] - new_data.x.values[:-1]
22
23     new_data['y_diff'] = new_data.y
24     new_data.loc[1:, 'y_diff'] = new_data.y.values[1:] - new_data.y.values[:-1]
25
26     new_data.loc[new_data['begin'] == 1, 'x_diff'] = 0
27     new_data.loc[new_data['begin'] == 1, 'y_diff'] = 0
28
29     angles = [np.arctan2(y_diff, x_diff) for y_diff, x_diff in zip(new_data.y_diff, new_data.x_diff)]
30
31     new_data['angle'] = angles
32     new_data.loc[new_data['angle'] == 0, 'angle'] = np.nan
33     new_data.loc[new_data['begin'] == 1, 'angle'] = 0.
34
35     # когда такси стоит, оно сохраняет направление
36     new_data.loc[:, 'angle'].fillna(method='pad', inplace=True)
37
38     return new_data

```

```

In [6]: 1 new_all_data = add_features(all_data)

```

Полученные данные.

```
In [7]: 1 new_all_data.head()
```

Out[7]:

	status	y	ts	x	begin	session	delay	dist	av_speed	x_diff	y_diff	angle
0	0.0	0.000000	0.0	0.000000	1	0	0.0	0.000000	NaN	0.000000	0.000000	0.000000
1	0.0	0.291129	9.0	-0.644802	0	0	9.0	0.707478	0.078609	-0.644802	0.291129	2.717491
2	0.0	1.229173	17.0	-0.243663	0	0	8.0	1.020216	0.127527	0.401139	0.938044	1.166697
3	0.0	13.056778	25.0	2.447144	0	0	8.0	12.129826	1.516228	2.690807	11.827605	1.347101
4	0.0	21.475057	33.0	-3.184499	0	0	8.0	10.128318	1.266040	-5.631643	8.418279	2.160397

Предсказания с использованием методов кластеризации.

```

In [10]: 1 def clustering_prediction(fit_predict_methods, data):
2         session_values = np.unique(data.session.values)
3
4         session_predictions = []
5         cluster_distances = []
6         for i, session in tqdm_notebook(enumerate(session_values), total=len(session_values)):
7
8             # Данные рассматриваемой сессии
9             session_data = data[data.session == session]
10            session_data = session_data[session_data.ts != 0]
11
12            columns = ['x', 'y', 'ts', 'av_speed', 'angle', 'dist']
13
14            # Выделяем из данных сессии данные с нужными статусами
15            status_1_data = session_data[session_data.status == 1].loc[:, columns].copy()
16            status_0_data = session_data[session_data.status == 0].loc[:, columns].copy()
17            status_2_data = session_data[session_data.status == 2].loc[:, columns].copy()
18
19            # Соединяем всё вместе
20            status_data_w_dist = pd.concat([status_0_data.iloc[-2:, :],
21                                           status_1_data,
22                                           status_2_data.iloc[:2, :]],
23                                           axis=0)
24
25            if status_data_w_dist.shape[0] < 3:
26                status_data_w_dist = status_data_w_dist.append(status_data_w_dist.iloc[0, :])
27
28            status_data = status_data_w_dist.loc[:, columns[:-1]].copy()
29
30            if status_data.shape[0] < 3:
31                status_data = status_data.append(status_data.iloc[0, :])
32
33
34            #Стандартизуем данные
35            scaler = StandardScaler(with_mean=False)
36            status_data.iloc[:, :] = scaler.fit_transform(status_data.copy())
37
38
39            x_preds = []
40            y_preds = []
41            dist_preds = []

```

```

42     for method in fit_predict_methods:
43         pred = method(status_data)
44
45         # Найдём точки, где сменяется кластер
46         candidates = []
47         for j in range(len(pred) - 1):
48             if pred[j] != pred[j+1]:
49                 candidates.append(j)
50
51         # Если кластера всего два, то возьмём последнюю рассматриваемую
52         # точку в качестве правой границы
53         if len(candidates) == 1:
54             candidates.append(len(pred) - 1)
55
56         # Предсказания и среднее расстояние в кластере
57         x_preds.append(status_data_w_dist['x'].values[candidates[-1]])
58         y_preds.append(status_data_w_dist['y'].values[candidates[-1]])
59         dist_preds.append(np.mean(status_data_w_dist['dist'].values[candidates[0]:candidates[-1]+1]))
60
61
62         x_pred = np.mean(x_preds)
63         y_pred = np.mean(y_preds)
64
65         session_predictions.append([i, x_pred, y_pred])
66         cluster_distances.append(np.mean(dist_preds))
67
68     return np.array(session_predictions), np.array(cluster_distances)
69
70 fit_predict_methods = [KMeans(n_clusters=3,
71                               random_state=42,
72                               n_init=100,
73                               max_iter=1000,
74                               n_jobs=-1).fit_predict,
75                        AgglomerativeClustering(n_clusters=3,
76                                                linkage='ward').fit_predict,
77                        SpectralClustering(n_clusters=3, n_jobs=-1).fit_predict]

```

```
In [11]: 1 # Очень долго выполняющийся код
        2 # all_predictions, all_cluster_distances = clustering_prediction(fit_predict_methods, new_all_data)
```

100%

115204/115204 [4:46:09<00:00, 6.71it/s]

Предсказание с помощью поиска первой точки со статусом 2.

```
In [12]: 1 def first_status_2(data):
        2     d = {}
        3
        4     session_values = np.unique(data.session.values)
        5     for session in tqdm_notebook(session_values):
        6         session_data = data[data.session == session]
        7         pickup_row = session_data[session_data.status == 2].head(1)
        8         d[session] = (pickup_row.x.values[0], pickup_row.y.values[0])
        9
        10    return d
```

Формирование словаря всех предсказаний.

```

In [16]: 1 def final_preds(predictions, cluster_distances, data):
          2     # пороговое значение
          3     threshold = np.quantile(cluster_distances, 0.1)
          4     print(f'Пороговое значение межкластерного расстояния: {threshold}')
          5
          6     valid_predictions = predictions[cluster_distances < threshold]
          7
          8     preds = {}
          9     for pred in valid_predictions:
         10         preds[int(pred[0])] = (pred[1], pred[2], 0)
         11
         12     first2 = first_status_2(data)
         13
         14     for key in first2:
         15         if key not in preds:
         16             preds[key] = (first2[key][0], first2[key][1], 1)
         17
         18     return preds

```

```

In [17]: 1 preds = final_preds(all_predictions, all_cluster_distances, all_data)

```

Пороговое значение межкластерного расстояния: 1.1775877242421593

100%

115204/115204 [30:07<00:00, 63.73it/s]

```

In [18]: 1 preds[1], preds[100]

```

```

Out[18]: ((467.8795440394896, 1009.9926489507196, 1),
          (373.0633040684311, -1113.3754465019024, 1))

```

Сохранение данных.

```
In [19]: 1 preds_data = pd.DataFrame(preds).T
2 preds_data.columns = ['x', 'y', 'method']
3 preds_data = preds_data.sort_index()
4
5 preds_data.head()
```

Out[19]:

	x	y	method
0	537.323154	73.614187	1.0
1	467.879544	1009.992649	1.0
2	-698.044056	-39.040766	1.0
3	-140.035784	2103.199406	1.0
4	-1033.518299	-107.762341	1.0

```
In [20]: 1 preds_data.to_csv('data/all_predictions.csv')
```

```
In [ ]: 1
```