

Погодные условия.

In [115]:

```
1 import scipy.stats as sps
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import pandas as pd
6 from tqdm.notebook import tqdm
7
8 import datetime
9
10 import plotly.graph_objects as go
11 import plotly.express as px
12 import plotly.offline
13
14 from sklearn.linear_model import Lasso
15 from sklearn.metrics import mean_squared_error
16 from sklearn.model_selection import train_test_split
17 from sklearn.preprocessing import LabelEncoder
18
19 from statsmodels.stats.multitest import multipletests
20
21 from sklearn.preprocessing import StandardScaler
22
23 sns.set(font_scale=1.5)
```

Загрузим сначала датасет с данными о погоде в каждый день.

In [34]:

```
1 weather_data = pd.read_csv('cycle-share-dataset/weather.csv',
2                             parse_dates=[0], na_values='-')
3
4 trips = pd.read_csv('cycle-share-dataset/trip.csv', error_bad_lines=False,
5                     parse_dates=[1, 2])
6
7
8 weather_data.head()
```

b'Skipping line 50794: expected 12 fields, saw 20\n'

Out[34]:

	Date	Max_Temperature_F	Mean_Temperature_F	Min_TemperatureF	Max_Dew_Point_F	Mean
0	2014-10-13	71	62.0	54	55	
1	2014-10-14	63	59.0	55	52	
2	2014-10-15	62	58.0	54	53	
3	2014-10-16	71	61.0	52	49	
4	2014-10-17	64	60.0	57	55	

5 rows × 21 columns

Посмотрим на колонки в датасете о погоде.

In [33]:

```
1 weather_data.columns
```

Out[33]:

```
Index(['Date', 'Max_Temperature_F', 'Mean_Temperature_F', 'Min_TemperatureF',
      'Max_Dew_Point_F', 'MeanDew_Point_F', 'Min_Dewpoint_F', 'Max_Humidity',
      'Mean_Humidity', 'Min_Humidity', 'Max_Sea_Level_Pressure_In',
      'Mean_Sea_Level_Pressure_In', 'Min_Sea_Level_Pressure_In',
      'Max_Visibility_Miles', 'Mean_Visibility_Miles', 'Min_Visibility_Miles',
      'Max_Wind_Speed_MPH', 'Mean_Wind_Speed_MPH', 'Max_Gust_Speed_MPH',
      'Precipitation_In', 'Events'],
      dtype='object')
```

Исследуем сначала зависимость дневного числа поездок от различных погодных факторов.

Добавим признаки числа, месяца и дня недели для дальнейшего удобства.

In [47]:

```
1 time_data = trips[['trip_id', 'starttime', 'stoptime']].copy()
2
3 time_data['start_year'] = time_data['starttime'].dt.year
4 time_data['start_month'] = time_data['starttime'].dt.month
5 time_data['start_day'] = time_data['starttime'].dt.day
6 time_data['start_date'] = time_data['starttime'].dt.date
7 time_data['start_time'] = time_data['starttime'].dt.time
8 time_data['start_weekday'] = time_data['starttime'].dt.dayofweek
9
10 time_data['stop_year'] = time_data['stoptime'].dt.year
11 time_data['stop_month'] = time_data['stoptime'].dt.month
12 time_data['stop_day'] = time_data['stoptime'].dt.day
13 time_data['stop_date'] = time_data['stoptime'].dt.date
14 time_data['stop_time'] = time_data['stoptime'].dt.time
15 time_data['stop_weekday'] = time_data['stoptime'].dt.dayofweek
16
17 time_data.drop(columns=['starttime', 'stoptime'], inplace=True)
```

Количество поездок в зависимости от дня.

In [48]:

```
1 trips_by_date = time_data.groupby('start_date').count()['trip_id']
2
3 trips_by_date = pd.DataFrame(trips_by_date).reset_index()
4
5 trips_by_date.columns = ['Date', 'trip_num']
6
7 trips_by_date.loc[:, 'Date'] = pd.to_datetime(trips_by_date['Date'])
8
9 trips_by_date.head()
```

Out[48]:

	Date	trip_num
0	2014-10-13	818
1	2014-10-14	982
2	2014-10-15	626
3	2014-10-16	790
4	2014-10-17	588

Соединим две полученные таблицы.

In [49]:

```
1 trips_weather = weather_data.merge(trips_by_date, how='left', on='Date')
2
3 trips_weather.head()
```

Out[49]:

	Date	Max_Temperature_F	Mean_Temperature_F	Min_TemperatureF	Max_Dew_Point_F	Mean
0	2014-10-13	71	62.0	54	55	
1	2014-10-14	63	59.0	55	52	
2	2014-10-15	62	58.0	54	53	
3	2014-10-16	71	61.0	52	49	
4	2014-10-17	64	60.0	57	55	

5 rows × 22 columns

Можно видеть, что признак `Events` задаётся в текстовом формате.

В дальнейшем, при обучении моделей и проверке гипотез будем проводить label encoding для этого признака.

In [133]:

```
1 events = trips_weather.Events.dropna()
2
3 le = LabelEncoder()
4 le.fit(events)
```

Out[133]:

LabelEncoder()

Попробуем сначала визуально оценить зависимость числа поездок от указанных погодных параметров.

In [50]:

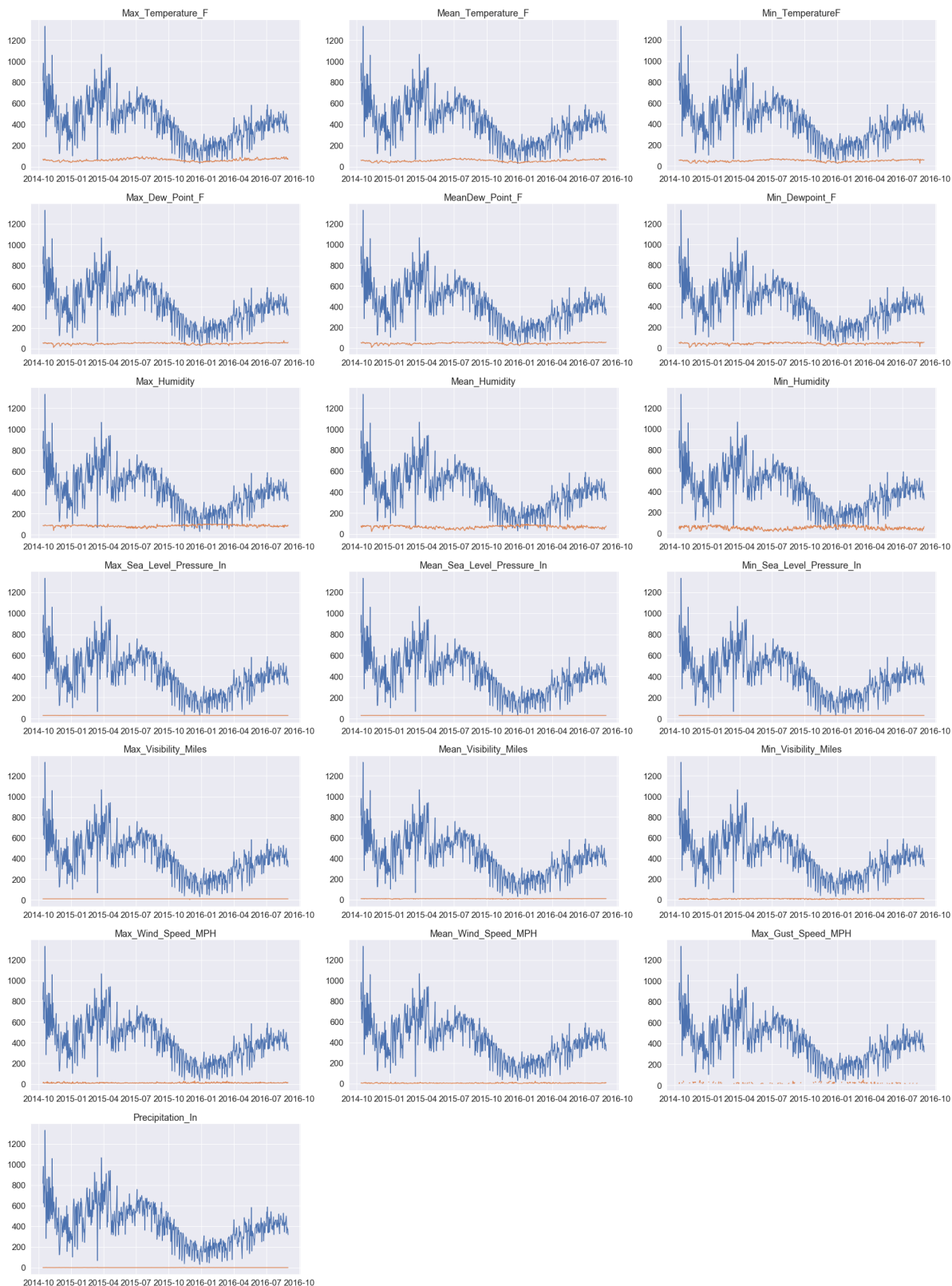
```
1 trips_weather.columns[1:-2].shape[0]
```

Out[50]:

19

In [136]:

```
1 plt.figure(figsize=(26, 35))
2
3 for i, column in enumerate(trips_weather.columns[1:-2]):
4     plt.subplot(7, 3, i+1)
5     plt.title(column)
6
7     plt.plot(trips_weather.Date, trips_weather.trip_num, label='trip_num')
8     plt.plot(trips_weather.Date, trips_weather[column], label='признак')
9
10 plt.tight_layout()
```



По графикам можно предположить, что количество поездок зависит от температуры, влажности, точки росы.

Также, можно видеть, что в колонке `Max_gust_Speed_MPH` много пропущенных значений. Уберём его пока из рассмотрения.

Вспомним предыдущий семестр. Будем использовать регрессию для отбора признаков. Наиболее подходящей для отбора признаков является Лассо-регрессия. Для отбора признаков разделим сначала выборку на `train` и `test`. Затем, будем проходиться по всем подмножествам признаков и считать качество обученной модели.

In [130]:

```
1 class BestFeaturesSelection:
2     def __init__(self, estimator, scoring, parameters=dict(),
3                   test_size=0.3, random_state=17, minimize=True):
4         """
5         Отбор наилучших признаков
6         estimator: конструктор класса, например, LinearRegression
7         paramters: параметры, передаваемые конструктору estimator,
8                     например dict(fit_intercept=False)
9         scoring: функция риска, например, mean_squared_error
10        minimize: минимизировать ли функционал качества
11                    (иначе - максимизировать)
12        """
13
14        self.estimator = estimator
15        self.parameters = parameters
16        self.scoring = scoring
17        self.test_size = test_size
18        self.random_state = random_state
19        self.minimize=minimize
20
21    def fit(self, X, y):
22        """
23        Подбор лучшего подмножества признаков
24        и обучение модели на нём
25        """
26
27        # разделение выборки на test и train.
28        X_train, X_test, y_train, y_test = train_test_split(X, y,
29                                                            test_size=0.3,
30                                                            random_state=17)
31
32        # стандартизирование признаков
33        scaler = StandardScaler()
34
35        self.results_ = [] # список пар (вектор использованных признаков,
36                             # значение функции потерь)
37        features_count = X.shape[1]
38
39        for bitmask in tqdm(range(1, 2 ** features_count)):
40            subset = [i == "1" for i in \
41                     np.binary_repr(bitmask, width=features_count)]
42            # binary_repr возвращает строку длины width с двоичным
43            # представлением числа и ведущими нулями
44
45            x_train = scaler.fit_transform(X_train[:, subset])
46            x_test = scaler.transform(X_test[:, subset])
47
48            reg = self.estimator(**self.parameters)
49            reg.fit(x_train, y_train)
50
51            score = self.scoring(y_test, reg.predict(x_test[:, subset]))
52            # вычисление качества модели
53
54            self.results_.append((subset, score))
55
56        self.results_.sort(key = lambda pair: pair[1],
57                           reverse=not self.minimize)
58        # сортируем по второму элементу в нужном порядке
59
60        self._best_subset = self.results_[0][0]
```

```

60         self._best_estimator = self.estimator(**self.parameters)
61         self._best_estimator.fit(X_train[:, self._best_subset], y_train)
62
63         return self._best_estimator
64
▼ 65     def predict(self, X):
66         """
67         Предсказание модели,
68         обученной на наилучшем подмножестве признаков.
69         """
70
71         return self._best_estimator.predict(X[:, self._best_subset]);

```

In [131]:

```

1  valid_data = trips_weather.iloc[:, 1:-3].copy()
2  valid_data['trip_num'] = trips_weather.trip_num
3  valid_data.dropna(inplace=True)
4
5  X = valid_data.values[:, :-1]
6  y = valid_data.values[:, -1]
7
▼ 8  bfs = BestFeaturesSelection(Lasso,
9                               mean_squared_error)

```


In [132]:

```
1 bfs.fit(X, y)
```

100%

262143/262143 [05:42<00:00, 766.13it/s]

```
/Users/gregorypolyakov/opt/anaconda3/envs/stats/lib/python3.7/site-packages/sklearn/linear_model/_coordinate_descent.py:476: ConvergenceWarning:
```

```
Objective did not converge. You might want to increase the number of iterations. Duality gap: 794.6356730181724, tolerance: 793.7842871794873
```

```
/Users/gregorypolyakov/opt/anaconda3/envs/stats/lib/python3.7/site-packages/sklearn/linear_model/_coordinate_descent.py:476: ConvergenceWarning:
```

```
Objective did not converge. You might want to increase the number of iterations. Duality gap: 794.6356730181724, tolerance: 793.7842871794873
```

```
/Users/gregorypolyakov/opt/anaconda3/envs/stats/lib/python3.7/site-packages/sklearn/linear_model/_coordinate_descent.py:476: ConvergenceWarning:
```

```
Objective did not converge. You might want to increase the number of iterations. Duality gap: 794.9687269581482, tolerance: 793.7842871794873
```

```
/Users/gregorypolyakov/opt/anaconda3/envs/stats/lib/python3.7/site-packages/sklearn/linear_model/_coordinate_descent.py:476: ConvergenceWarning:
```

```
Objective did not converge. You might want to increase the number of iterations. Duality gap: 794.9687269581482, tolerance: 793.7842871794873
```

```
/Users/gregorypolyakov/opt/anaconda3/envs/stats/lib/python3.7/site-packages/sklearn/linear_model/_coordinate_descent.py:476: ConvergenceWarning:
```

```
Objective did not converge. You might want to increase the number of iterations. Duality gap: 794.9392988774925, tolerance: 793.7842871794873
```

```
/Users/gregorypolyakov/opt/anaconda3/envs/stats/lib/python3.7/site-packages/sklearn/linear_model/_coordinate_descent.py:476: ConvergenceWarning:
```

```
Objective did not converge. You might want to increase the number of iterations. Duality gap: 794.9392988774925, tolerance: 793.7842871794873
```

```
/Users/gregorypolyakov/opt/anaconda3/envs/stats/lib/python3.7/site-packages/sklearn/linear_model/_coordinate_descent.py:476: ConvergenceWarning:
```

```
Objective did not converge. You might want to increase the number of iterations. Duality gap: 794.9732963014394, tolerance: 793.7842871794873
```

```
/Users/gregorypolyakov/opt/anaconda3/envs/stats/lib/python3.7/site-packages/sklearn/linear_model/_coordinate_descent.py:476: ConvergenceWarning:
```

```
Objective did not converge. You might want to increase the number of iterations. Duality gap: 794.9732963014394, tolerance: 793.7842871794873
```

Out[132]:

```
Lasso(alpha=1.0, copy_X=True, fit_intercept=True, max_iter=1000,
      normalize=False, positive=False, precompute=False, random_state=None,
      selection='cyclic', tol=0.0001, warm_start=False)
```

In [134]:

```
1 values = np.array(bfs._best_subset).reshape(-1, 1)
2
3 answer = pd.DataFrame(values, index = valid_data.columns[:-1],
4                       columns=['Значимый ли'])
5 answer
```

Out[134]:

	Значимый ли
Max_Temperature_F	False
Mean_Temperature_F	True
Min_TemperatureF	True
Max_Dew_Point_F	False
MeanDew_Point_F	True
Min_Dewpoint_F	False
Max_Humidity	True
Mean_Humidity	True
Min_Humidity	True
Max_Sea_Level_Pressure_In	True
Mean_Sea_Level_Pressure_In	False
Min_Sea_Level_Pressure_In	False
Max_Visibility_Miles	False
Mean_Visibility_Miles	False
Min_Visibility_Miles	False
Max_Wind_Speed_MPH	False
Mean_Wind_Speed_MPH	True
Max_Gust_Speed_MPH	False

Можем видеть, что почти все наши предположения подтвердились. Значимым также оказался достаточно неожиданный признак `Max_Sea_Level_Pressure_In`.

Проведём аналогичное исследование с помощью коэффициентов корреляции. Для каждого признака будем проводить тесты (вычисления коэффициентов корреляции) Пирсона, Спирмена и Кендалла. Затем сделаем поправку на множественную проверку гипотез.

Если один из трёх критериев после МПГ для какого либо признака отвергает гипотезу о независимости, то отвергаем независимость для данного признака.

In [108]:

```
1 p_values = []
2
3 for column in trips_weather.columns[1:-2]:
4
5     data = trips_weather.loc[:, ['trip_num', column]].dropna()
6     p_values.append(sps.spearmanr(data.trip_num, data[column])[1])
7     p_values.append(sps.kendalltau(data.trip_num, data[column])[1])
8     p_values.append(sps.pearsonr(data.trip_num, data[column])[1])
```

In [114]:

```
1 rejected = multipletests(p_values)[0]
2
3 spear_rejected = rejected[0::3]
4 ken_rejected = rejected[1::3]
5 pears_rejected = rejected[2::3]
6
7 rejected = np.logical_or(spear_rejected, ken_rejected)
8 rejected = np.logical_or(rejected, pears_rejected)
```

In [113]:

```
1 answer = pd.DataFrame(rejected, index = trips_weather.columns[1:-2],
2                        columns=[ 'Отверглась ли гипотеза о независимости' ])
3 answer
```

Out[113]:

Отверглась ли гипотеза о независимости	
Max_Temperature_F	True
Mean_Temperature_F	True
Min_TemperatureF	True
Max_Dew_Point_F	True
MeanDew_Point_F	True
Min_Dewpoint_F	True
Max_Humidity	True
Mean_Humidity	True
Min_Humidity	True
Max_Sea_Level_Pressure_In	False
Mean_Sea_Level_Pressure_In	True
Min_Sea_Level_Pressure_In	True
Max_Visibility_Miles	False
Mean_Visibility_Miles	True
Min_Visibility_Miles	True
Max_Wind_Speed_MPH	True
Mean_Wind_Speed_MPH	True
Max_Gust_Speed_MPH	True
Precipitation_In	True

Вывод:

После корреляционного анализа мы получили, что для почти всех признаков отверглась гипотеза о независимости от количества поездок. Примечательно, что количество поездок зависит даже от атмосферного давления.

Для подбора признаков при обучении каких-либо моделей стоит скорее использовать результат, полученный с помощью лассо-регрессии.

To Do:

Для датасетов о перебросах и групповых поездок провести аналогичное исследование.

In []:

1	
---	--