

- 1
 - 1.1
 - 1.2
- 2
- 3
- 4
- 5
- 6 -

RMSE

48.

In [1]:

```
#
!pip install lightgbm
```

Requirement already satisfied: lightgbm in e:\anaconda3\envs\ds_practi cum\l i b\si t e-packages (3.3.5)
 Requirement already satisfied: scipy in e:\anaconda3\envs\ds_practi cum\l i b\si t e-p ackages (from lightgbm) (1.8.0)
 Requirement already satisfied: scikit-learn!=0.22.0 in e:\anaconda3\envs\ds_pract i cum\l i b\si t e-packages (from lightgbm) (0.24.1)
 Requirement already satisfied: wheel in e:\anaconda3\envs\ds_practi cum\l i b\si t e-p ackages (from lightgbm) (0.38.4)
 Requirement already satisfied: numpy in e:\anaconda3\envs\ds_practi cum\l i b\si t e-p ackages (from lightgbm) (1.20.1)
 Requirement already satisfied: joblib>=0.11 in e:\anaconda3\envs\ds_practi cum\l i b \si t e-packages (from scikit-learn!=0.22.0->lightgbm) (1.2.0)
 Requirement already satisfied: threadpoolctl>=2.0.0 in e:\anaconda3\envs\ds_pract i cum\l i b\si t e-packages (from scikit-learn!=0.22.0->lightgbm) (3.1.0)

In [2]:

```
#
import os
import time
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from statsmodels.tsa.seasonal import seasonal_decompose
from statsmodels.graphics.tsaplots import plot_acf
from statsmodels.graphics.tsaplots import plot_pacf
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.model_selection import TimeSeriesSplit
```

```

from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import RandomForestRegressor
from lightgbm import LGBMRegressor
from sklearn.metrics import mean_squared_error

```

```

In [3]: pth1 = '/datasets/taxi.csv'
        pth2 = 'https://restricted/datasets/taxi.csv'

        def open_df(pth):
            return pd.read_csv(pth, index_col=[0], parse_dates=[0])

        if os.path.exists(pth1):
            data = open_df(pth1)
        else:
            try:
                data = open_df(pth2)
            except:
                print('Something is wrong, datasets not found!!!')

```

```

In [4]: # df
        def df_info(df):
            display(df.head())
            display(df.info())
            display(df.describe())

```

```

In [5]: df_info(data)

```

| | num_orders |
|---------------------|------------|
| | datetime |
| 2018-03-01 00:00:00 | 9 |
| 2018-03-01 00:10:00 | 14 |
| 2018-03-01 00:20:00 | 28 |
| 2018-03-01 00:30:00 | 20 |
| 2018-03-01 00:40:00 | 32 |

```

<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 26496 entries, 2018-03-01 00:00:00 to 2018-08-31 23:50:00
Data columns (total 1 columns):
#   Column      Non-Null Count  Dtype
---  -
0   num_orders  26496 non-null  int64
dtypes: int64(1)
memory usage: 414.0 KB
None

```

| | num_orders |
|-------|--------------|
| count | 26496.000000 |
| mean | 14.070463 |
| std | 9.211330 |
| min | 0.000000 |
| 25% | 8.000000 |
| 50% | 13.000000 |
| 75% | 19.000000 |
| max | 119.000000 |

num_order s

26496

```
In [6]: if data.index.is_monotonic:
        print(" ")
    else:
        data.sort_index(inplace=True)
        print(" ")
```

```
In [7]: data = data.resample('1H').sum()
```

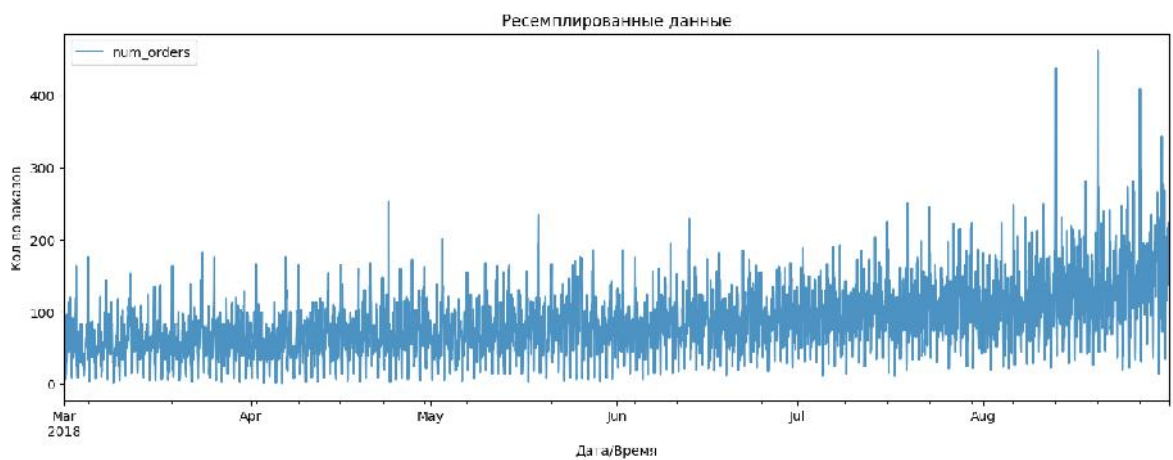
```

•
•
•
•
num_order s
354369;
;
```

```
In [8]: print(' ', data.index.min())
        print(' ', data.index.max())
        print(' ', data.index.max() - data.index.min())

: 2018-03-01 00:00:00
: 2018-08-31 23:00:00
: 183 days 23:00:00
```

```
In [9]: data.plot(figsize=(15, 5), alpha=0.8)
plt.title('')
plt.xlabel(' / ')
plt.ylabel(' - ')
plt.show()
```



```
In [10]: decomposed = seasonal_decompose(data['num_orders'])

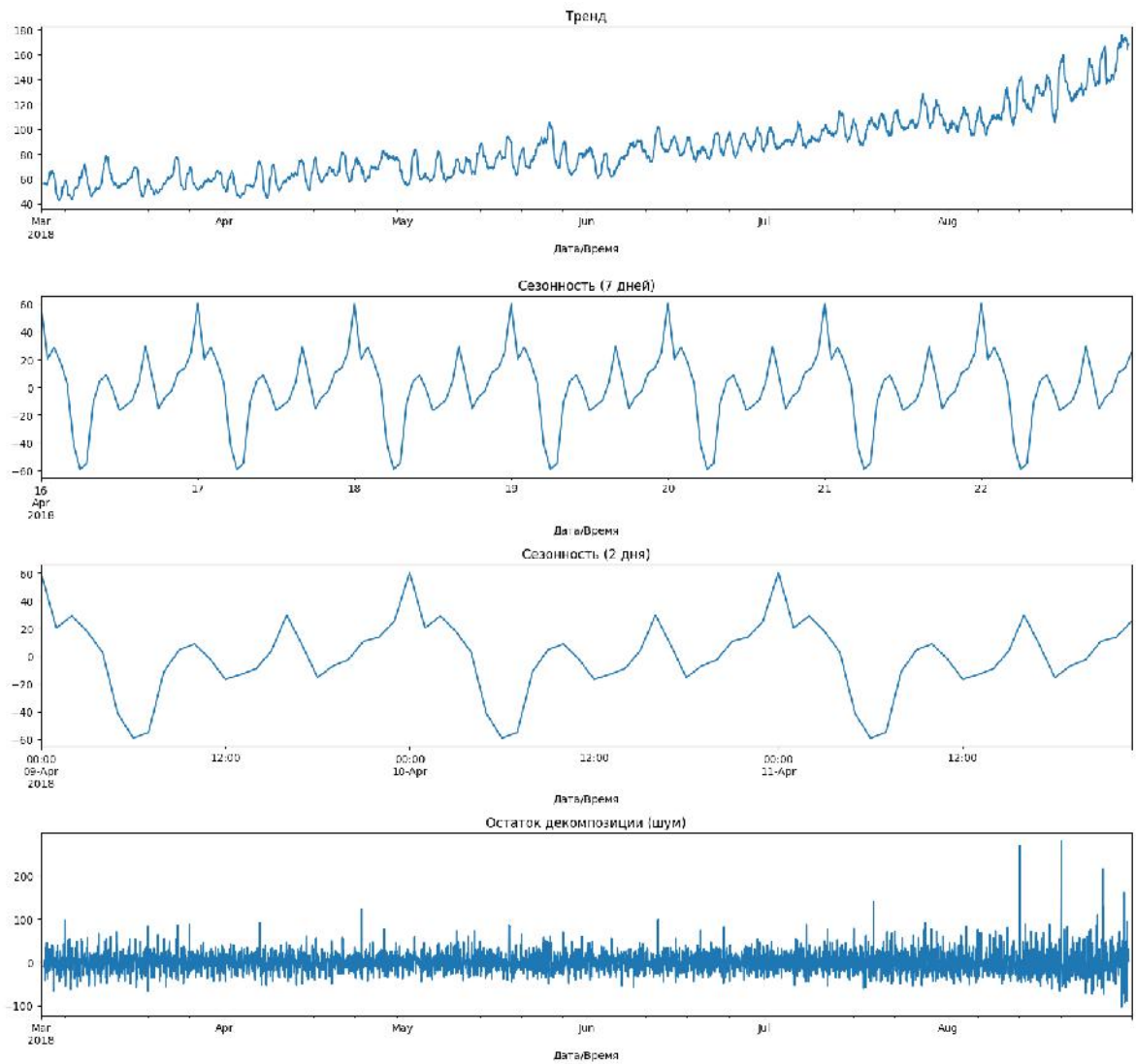
plt.figure(figsize=(15, 14))

#
plt.subplot(411)
decomposed.trend.plot(ax=plt.gca())
plt.title('')
plt.xlabel(' / ')

#
plt.subplot(412)
decomposed.seasonal['2018-04-16': '2018-04-22'].plot(ax=plt.gca())
plt.title('(7)')
plt.xlabel(' / ')

#
plt.subplot(413)
decomposed.seasonal['2018-04-09': '2018-04-11'].plot(ax=plt.gca())
plt.title('(2)')
plt.xlabel(' / ')

plt.subplot(414)
decomposed.resid.plot(ax=plt.gca())
plt.title('( )')
plt.xlabel(' / ')
plt.tight_layout()
```

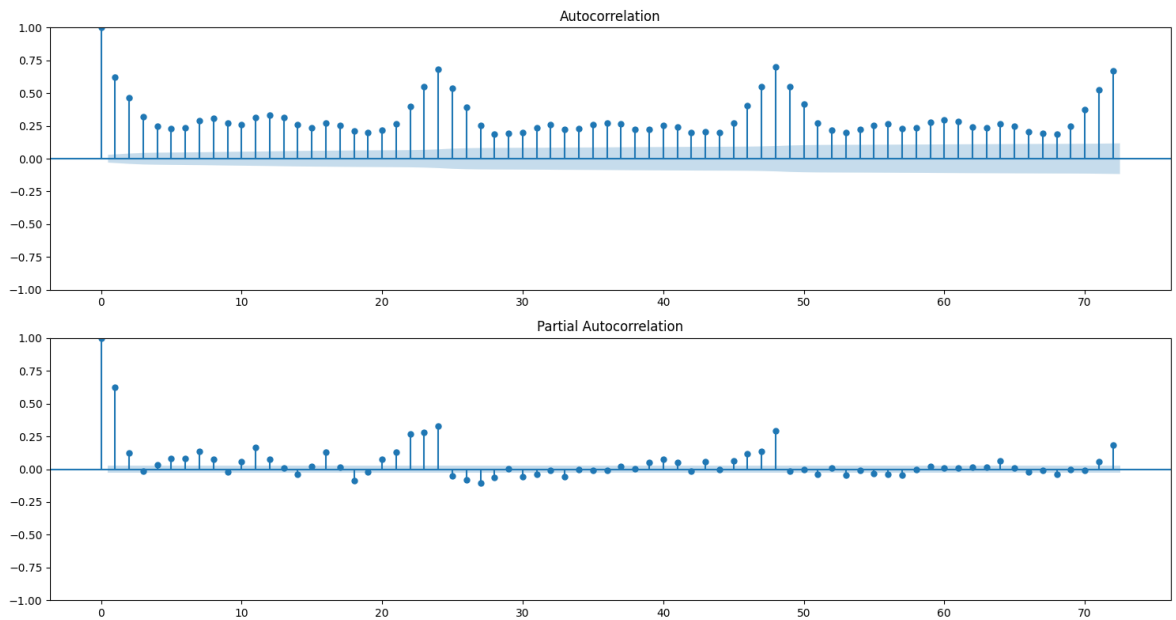


00-00 6-00 11-00, 15-00 00-00

ACF PCF

« »

```
In [11]: plt.figure(figsize=(15, 8))
plt.subplot(211)
plt.acf(data, ax=plt.gca(), lags=24*3)
plt.subplot(212)
plt.pacf(data, method='yw', ax=plt.gca(), lags=24*3)
plt.tight_layout()
```



ACF

24

, PCF

7, 11, 16,

24.

:

•

•

■

■

•

ACF PCF,

«

»

90:10

```
In [12]: def make_features(df, max_lag, rolling_mean_size):
df = df.copy()
df['month'] = df.index.month
df['day'] = df.index.day
df['dayofweek'] = df.index.dayofweek
df['is_weekend'] = df.index.weekday.isin([5, 6]) * 1
df['hour'] = df.index.hour

for lag in range(1, max_lag + 1):
    df['lag_{}'.format(lag)] = df['num_orders'].shift(lag)

df['rolling_mean'] = df['num_orders'].shift().rolling(rolling_mean_size).mean

return df
```

```
In [13]: lag = 24
rolling = 7 #7, 11, 16, 24
```

```
#
data_tng = make_features(data, lag, rolling)
data_tng = data_tng.dropna()
features = data_tng.drop(['num_orders'], axis=1)
target = data_tng['num_orders']

#
(features_train, features_test,
 target_train, target_test) = train_test_split(features, target,
                                                shuffle=False,
                                                test_size=0.1)
print(" ", features_train.shape)
print(" ", features_test.shape)
```

```
(3952, 30)
(440, 30)
```

```
: LinearRegression, RandomForestRegressor,
LGBMRegressor
```

```
In [14]: # C - TimeSeriesSplit
tscv = TimeSeriesSplit(n_splits=3)
```

```
In [15]: %%time
model_lrg = LinearRegression(n_jobs=-1)

scores = cross_val_score(model_lrg, features_train, target_train,
                          scoring='neg_root_mean_squared_error', cv=tscv)

cv_lrg_best_score = round(-np.mean(scores), 6)

cv_lrg_results = ['LinearRegression',
                  cv_lrg_best_score]

print("score", cv_lrg_best_score)
```

```
score 27.507456
CPU times: total: 15.6 ms
Wall time: 44 ms
```

```
In [16]: %%time
model_rfr = RandomForestRegressor(random_state=12345)

param_grid_rfr = {
    'n_estimators': range(10, 210, 50),
    'max_depth': [None] + [i for i in range(2, 11)]
}

#
cv_rfr = GridSearchCV(estimator=model_rfr,
                      param_grid=param_grid_rfr,
                      cv=tscv,
                      n_jobs=-1,
                      scoring='neg_root_mean_squared_error',
```

```

        verbose=10
    )
cv_rfr.fit(features_train, target_train)

cv_rfr_best_params = cv_rfr.best_params_
cv_rfr_best_score = round(-cv_rfr.best_score_, 6)

cv_rfr_results = ['RandomForest',
                  cv_rfr_best_score]

print("best params", cv_rfr_best_params, "score", cv_rfr_best_score)

```

Fitting 3 folds for each of 40 candidates, totalling 120 fits
 best params {'max_depth': None, 'n_estimators': 160} score 26.818797
 CPU times: total: 8.08 s
 Wall time: 34.5 s

LightGBM

In [17]:

```

%%time

model_lgb = LGBMRegressor()

param_grid_lgb = {
    'max_depth': [25, 50],
    'learning_rate': [0.01, 0.03],
    'n_estimators': range(10, 800, 50),
}

cv_lgb = GridSearchCV(estimator=model_lgb,
                      param_grid=param_grid_lgb,
                      cv=5,
                      n_jobs=-1,
                      scoring='neg_root_mean_squared_error',
                      verbose=10)

cv_lgb.fit(features_train, target_train)

cv_lgb_best_params = cv_lgb.best_params_
cv_lgb_best_score = round(-cv_lgb.best_score_, 6)

cv_lgb_results = ['LightGBM',
                  cv_lgb_best_score]

print("best params", cv_lgb_best_params, "score", cv_lgb_best_score)

```

Fitting 3 folds for each of 64 candidates, totalling 192 fits
 best params {'learning_rate': 0.01, 'max_depth': 25, 'n_estimators': 710} score 26.526482
 CPU times: total: 3.66 s
 Wall time: 23.5 s

In [18]:

```

model_analytics = pd.DataFrame([cv_lrg_results, cv_rfr_results, cv_lgb_results],
                               columns=['', '', '(RMSE)'])

model_analytics

```


Out [18]:

| | | (RMSE) |
|---|------------------|-----------|
| 0 | LinearRegression | 27.507456 |
| 1 | RandomForest | 26.818797 |
| 2 | LightGBM | 26.526482 |

```
                                : LinearRegression ,
RandomForestRegressor , LGBMRegressor .
                                GridSearchCV
```

RMSE.

```
LGBMRegressor .
```

In [19]:

```
#
def model_analysis (features_train, target_train, features_test, target_test, model_name):
    start = time.time()
    model.fit(features_train, target_train)
    end = time.time()
    fit_time = end - start

    start = time.time()
    model_pred = model.predict(features_test)
    end = time.time()
    pred_time = end - start

    rmse = round(mean_squared_error(target_test, model_pred, squared = False), 6)

    return [model_name, rmse, fit_time, pred_time, model_pred]
```

In [20]:

```
model_lgb = LGBMRegressor(**cv_lgb_best_params)
model = model_analysis(features_train, target_train, features_test, target_test, model_lgb)

print('      :', model[0],
      '\n      (RMSE):', model[1],
      '\n      :', model[2], ' .',
      '\n      :', model[3], ' .',
      )
```

```
      : LightGBM
      (RMSE): 39.765486
      : 0.5669946670532227
      : 0.004998922348022461
```

In [21]:

```
predict_test = pd.Series(model[4], index=target_test.index)
plt.figure(figsize=(16, 5))
plt.title('')
```

```
plt.xlabel(' / ')
plt.ylabel(' ')
plt.plot(target_test, alpha=0.5, label=' ')
plt.plot(predict_test, label=' ')
plt.legend()
plt.show()
```



LGBMRegressor

{'learning_rate': 0.01, 'max_depth': 25, 'n_estimators': 710}

«

»

:

-

:

• - 10%

;

• RMSE

48.

:

1. :

•

:

■

■

num_order s

354369

;

■

•

:

■

;

■

2. :

•

;

•

:

:

- -
- ACF PCF,
 - « »
- 3.
 -
 -
 -
 - 90:10
 - cross_val _score GridSearchCV
 - :
 - LinearRegression;
 - RandomForestRegressor;
 - LGBMRegressor.
 - RMSE,
 - LGBMRegressor
- 4.
 - LightGBM:
 - (RMSE): 39.765486;
 - : 0.5669946670532227 ;
 - : 0.004998922348022461 ;
 - : {'learning_rate': 0.01, 'max_depth': 25, 'n_estimators': 710}.
 -