

In [65]:

```
# coding=gbk
import os

import numpy as np
from sklearn.metrics import mean_squared_error, mean_absolute_error
from xgboost import plot_importance

import matplotlib.pyplot as plt
import xgboost as xgb

plt.rcParams['font.sans-serif'] = 'SimHei'
plt.rcParams['axes.unicode_minus'] = False

import pandas as pd

pd.set_option('expand_frame_repr', False)
pd.set_option('display.max_rows', 20)
pd.set_option('precision', 2)
```

In [66]:

```
columns = [f'HE{i:02d}' for i in range(1, 25)]

def create_features(df, label=None):
    df['date'] = df.index
    df['hour'] = df['date'].dt.hour
    df['day_of_week'] = df['date'].dt.dayofweek
    df['quarter'] = df['date'].dt.quarter
    df['month'] = df['date'].dt.month
    df['year'] = df['date'].dt.year
    df['day_of_year'] = df['date'].dt.dayofyear
    df['day_of_month'] = df['date'].dt.day
    df['week_of_year'] = df['date'].dt.weekofyear

    X = df[['hour', 'day_of_week', 'quarter', 'month', 'year', 'day_of_year', 'day_of_month', 'week_
    if label:
        y = df[label]
        return X, y
    return X
```

In [67]:

```
def read_data(data_file, sheet_names) -> pd.DataFrame:
    def get_data_per_sheet(sheet_name):
        df = pd.read_excel(data_file, parse_dates=['DATE'], index_col=[0], usecols=['DATE', *columns],
                           sheet_name=sheet_name)

        data = pd.DataFrame()
        for index, row in df.iterrows():
            d = [getattr(row, c) for c in columns]
            t = [pd.to_datetime(index.strftime('%Y-%m-%d') + f' {i - 1:02d}:00:00') for i in range(
                len(d))]
            dd = pd.DataFrame(index=[index]).from_dict({'value': d, 'DATE': t})
            dd.set_index('DATE', inplace=True)
            if len(data):
                data = pd.concat([data, dd])
            else:
                data = dd

        return data

    data = get_data_per_sheet(sheet_names[0])
    for sheet_name in sheet_names[1:]:
        data = pd.concat([data, get_data_per_sheet(sheet_name)])

    return data
```

In [68]:

```
def get_data(data_path) -> pd.DataFrame:
    if not os.path.exists(data_path):
        sheet_names = list(map(int, data_path[:-5].split('_')[1:]))
        data_name = data_path.split('_')[0]

        if data_name == 'data':

            data2014 = read_data('2014PJM数据.xls', sheet_names)
            data2015 = read_data('2015PJM数据.xls', sheet_names)
            data2016 = read_data('2016PJM数据.xls', sheet_names)

            data = pd.concat([data2014, data2015, data2016])
        else:
            data = read_data(f'{data_name}PJM数据.xls', sheet_names)

        data.to_excel(data_path)

    else:
        data = pd.read_excel(data_path, index_col=[0], parse_dates=['DATE'])

    return data
```

In [69]:

```
# data_path = 'data_0.xlsx'
# data = get_data(data_path)

data_path = 'pjm_data.csv'
data = pd.read_csv(data_path, index_col=[0], parse_dates=['DATE'])
```

In [70]:

```
sorted(set(data.index.to_list()))
```

Out[70]:

```
[Timestamp('2002-01-01 01:00:00'),
 Timestamp('2002-01-01 02:00:00'),
 Timestamp('2002-01-01 03:00:00'),
 Timestamp('2002-01-01 04:00:00'),
 Timestamp('2002-01-01 05:00:00'),
 Timestamp('2002-01-01 06:00:00'),
 Timestamp('2002-01-01 07:00:00'),
 Timestamp('2002-01-01 08:00:00'),
 Timestamp('2002-01-01 09:00:00'),
 Timestamp('2002-01-01 10:00:00'),
 Timestamp('2002-01-01 11:00:00'),
 Timestamp('2002-01-01 12:00:00'),
 Timestamp('2002-01-01 13:00:00'),
 Timestamp('2002-01-01 14:00:00'),
 Timestamp('2002-01-01 15:00:00'),
 Timestamp('2002-01-01 16:00:00'),
 Timestamp('2002-01-01 17:00:00'),
 Timestamp('2002-01-01 18:00:00')]
```

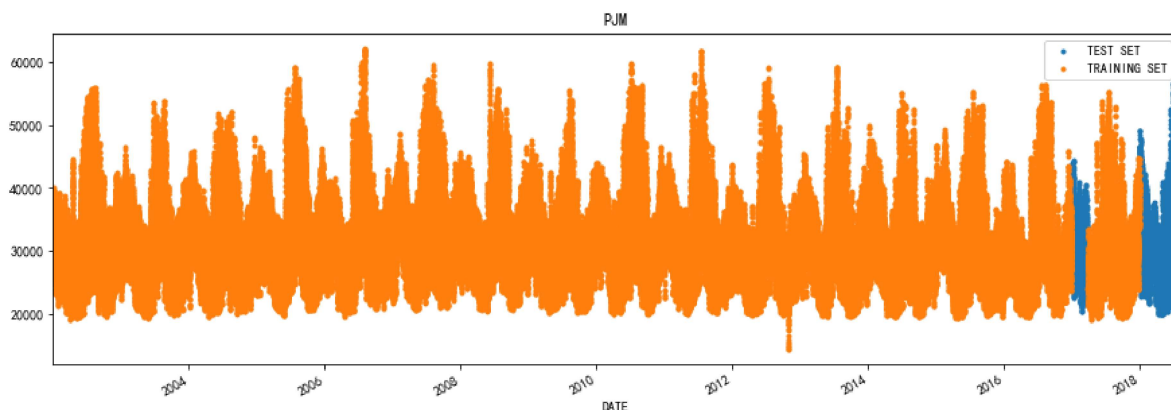
In [72]:

```
split_ratio = 0.95
split_index = int(len(data) * split_ratio)
split_date = data.iloc[split_index].name
train_data = data.iloc[:split_index].copy()
test_data = data.iloc[split_index:].copy()
print(f'split_date={split_date}, len(train_data)={len(train_data)}, len(test_data)={len(test_data)}')
```

```
split_date=2017-03-30 02:00:00, len(train_data)=138097, len(test_data)=7269
```

In [74]:

```
_ = test_data \
    .rename(columns={'value': 'TEST SET'}) \
    .join(train_data.rename(columns={'value': 'TRAINING SET'}), how='outer') \
    .plot(figsize=(15, 5), title='PJM', style='.')
plt.show()
```



In []:

训练集

X_train, y_train = create_features(train_data, label='value')

测试集

X_test, y_test = create_features(test_data, label='value')

模型

In [34]:

模型

reg = xgb.XGBRegressor(n_estimators=10000)

训练

reg.fit(X_train, y_train, eval_set=[(X_train, y_train), (X_test, y_test)], early_stopping_rounds=100)

[0] validation_0-rmse:23080.73438 validation_1-rmse:22288.98438

Multiple eval metrics have been passed: 'validation_1-rmse' will be used for early stopping.

Will train until validation_1-rmse hasn't improved in 100 rounds.

[1] validation_0-rmse:16351.98535 validation_1-rmse:15543.28125

[2] validation_0-rmse:11705.17871 validation_1-rmse:10890.80566

[3] validation_0-rmse:8543.58594 validation_1-rmse:8019.77539

[4] validation_0-rmse:6408.63184 validation_1-rmse:6036.32959

[5] validation_0-rmse:5024.89502 validation_1-rmse:4991.55859

[6] validation_0-rmse:4179.26270 validation_1-rmse:4309.34375

[7] validation_0-rmse:3651.98511 validation_1-rmse:3974.51660

[8] validation_0-rmse:3352.31909 validation_1-rmse:3826.53735

[9] validation_0-rmse:3176.65552 validation_1-rmse:3776.46851

[10] validation_0-rmse:3003.09106 validation_1-rmse:3768.60815

[11] validation_0-rmse:2944.55249 validation_1-rmse:3747.82910

[12] validation_0-rmse:2874.28467 validation_1-rmse:3723.09619

[13] validation_0-rmse:2850.01685 validation_1-rmse:3719.63647

[14] validation_0-rmse:2782.59619 validation_1-rmse:3724.53955

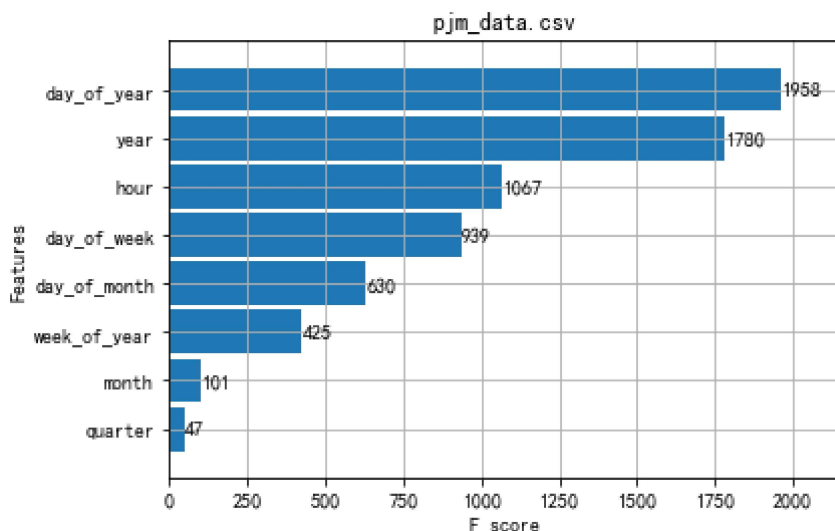
[15] validation_0-rmse:2752.70227 validation_1-rmse:3754.18070

In [43]:

分析重要程度

_ = plot_importance(reg, height=0.9, title=data_path)

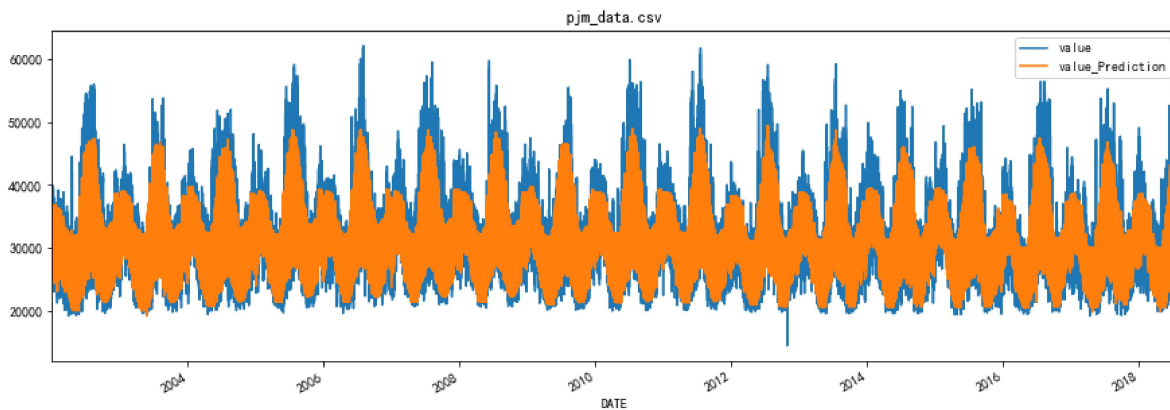
plt.show()



In [44]:

预测结果

```
test_data['value_Prediction'] = reg.predict(X_test)
train_data['value_Prediction'] = reg.predict(X_train)
pjme_all = pd.concat([test_data, train_data], sort=False)
_ = pjme_all[['value', 'value_Prediction']].plot(figsize=(15, 5))
plt.title(data_path)
plt.show()
```



In [48]:

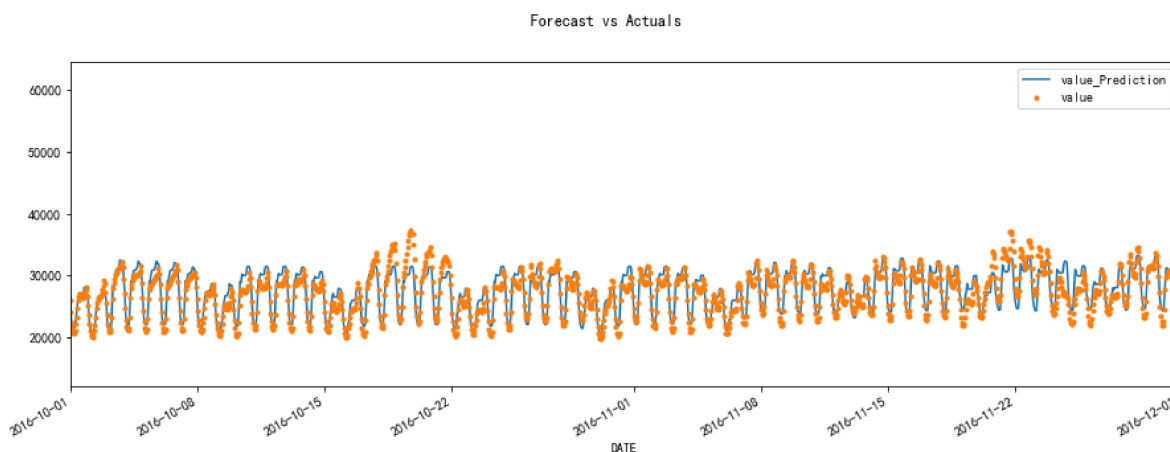
指定日期范围内画出预测结果与真实结果作对比

```
def plot_predict_days(lower, upper):
    f, ax = plt.subplots(1)
    f.set_figheight(5)
    f.set_figwidth(15)
    _ = pjme_all[['value_Prediction', 'value']].plot(ax=ax,
                                                    style=['-', '.'])

    ax.set_xbound(lower=lower, upper=upper)
    # ax.set_ylim(0, 60000)
    plot = plt.suptitle('Forecast vs Actuals')
    plt.show()
```

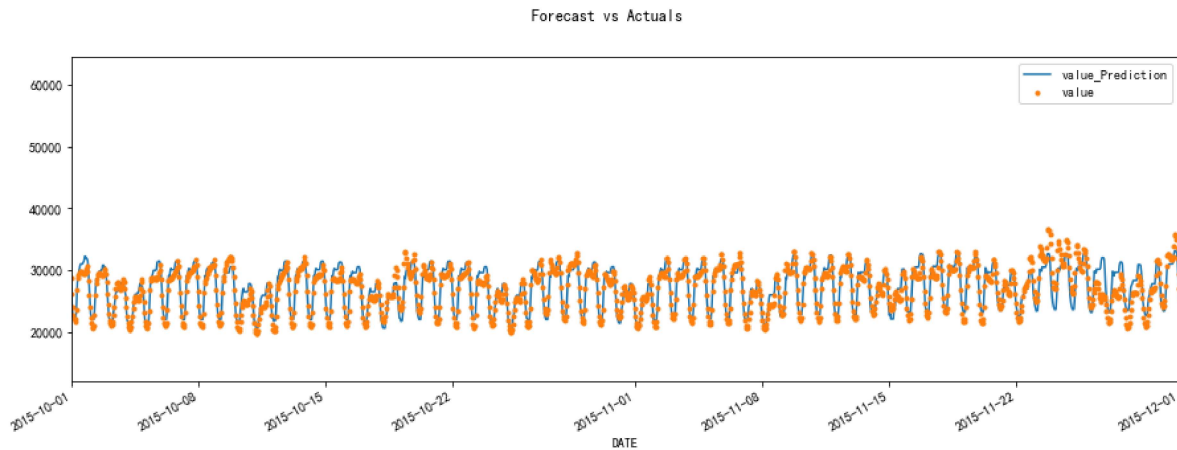
In [49]:

```
plot_predict_days(lower='10-01-2016', upper='12-01-2016')
```



In [50]:

```
plot_predict_days(lower='10-01-2015', upper='12-01-2015')
```



In [52]:

```
# 评估MSE
print(mean_squared_error(y_true=test_data['value'],
                        y_pred=test_data['value_Prediction']))
```

13835695.715198534

In [53]:

```
# 评估MAE
print(mean_absolute_error(y_true=test_data['value'],
                        y_pred=test_data['value_Prediction']))
```

2870.349107994996

In [54]:

```
# 评估MAPE
def mean_absolute_percentage_error(y_true, y_pred):
    y_true, y_pred = np.array(y_true), np.array(y_pred)
    return np.mean(np.abs((y_true - y_pred) / y_true)) * 100

print(mean_absolute_percentage_error(y_true=test_data['value'],
                                    y_pred=test_data['value_Prediction']))
```

8.960519640829324

In [55]:

最好的和最坏的结果

```
test_data['error'] = test_data['value'] - test_data['value_Prediction']
test_data['abs_error'] = test_data['error'].apply(np.abs)
error_by_day = test_data.groupby(['year', 'month', 'day_of_month']) \
    .mean()[['value', 'value_Prediction', 'error', 'abs_error']]
```

In [56]:

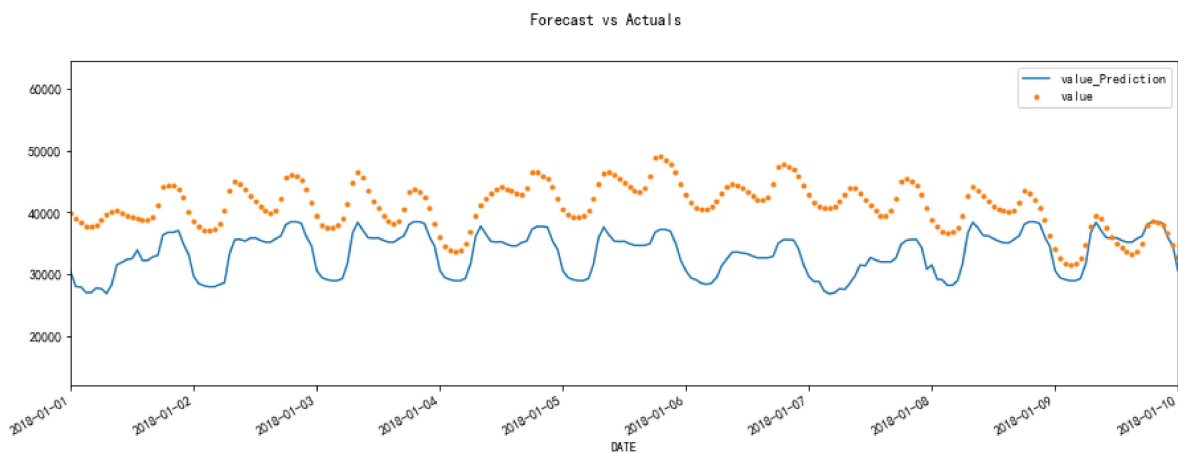
最坏的前十天

```
print(error_by_day.sort_values('abs_error', ascending=False).head(10))
```

			value	value_Prediction	error	abs_error
year	month	day_of_month				
2018	1	6	43565.75	32229.83	11335.92	11335.92
		7	42159.71	31077.75	11081.96	11081.96
		5	44197.79	33921.58	10276.21	10276.21
	7	2	45218.12	35804.73	9413.40	9413.40
		3	45258.04	36020.37	9237.67	9237.67
		1	40584.83	31532.32	9052.51	9052.51
	1	1	40202.48	31700.63	8501.85	8501.85
2017	3	15	36999.58	28991.37	8008.21	8008.21
2018	1	2	41627.33	33879.34	7747.99	7747.99
	5	26	33316.92	25759.93	7556.98	7556.98

In [58]:

```
plot_predict_days(lower='01-01-2018', upper='01-10-2018')
```



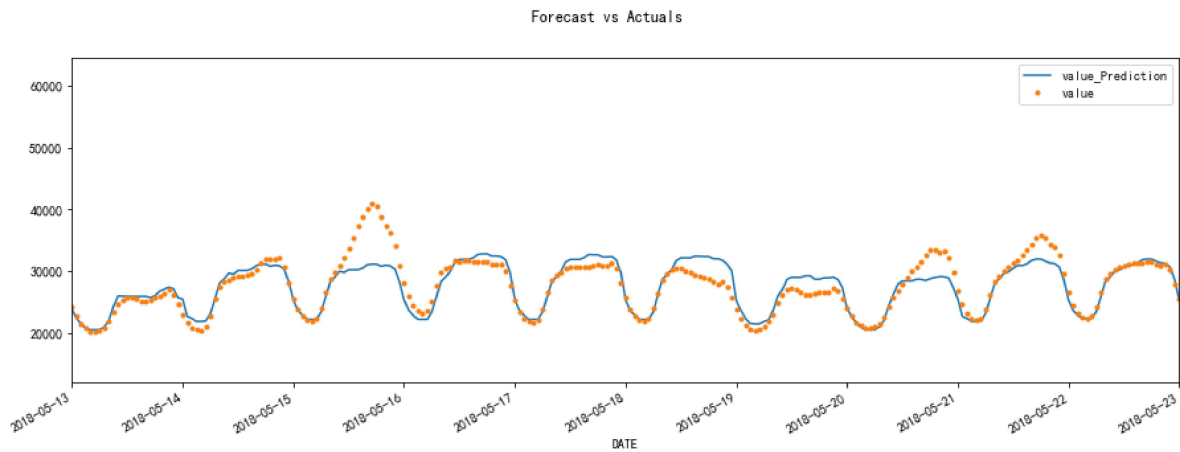
In [57]:

```
# 最好的前十天
print(error_by_day.sort_values('abs_error', ascending=True).head(10))
```

			value	value_Prediction	error	abs_error
year	month	day_of_month				
2018	5	22	28450.00	28329.15	120.85	436.09
		13	24001.50	24418.96	-417.46	485.55
2017	2	3	34424.38	34047.51	376.86	522.64
2018	4	27	26166.88	26500.90	-334.03	526.31
		25	26959.96	26522.16	437.80	540.88
	5	6	23796.08	24297.86	-501.78	634.86
2017	3	29	27230.50	27486.14	-255.64	647.75
2018	4	29	23541.92	23715.10	-173.18	685.14
2017	2	16	33854.96	33628.13	226.83	722.44
2018	2	1	33615.46	33547.46	68.00	802.90

In [59]:

```
plot_predict_days(lower='05-13-2018', upper='05-23-2018')
```



In [62]:

```
plot_predict_days(lower='05-21-2018', upper='05-23-2018')
```

