

# 导入包

```
In [1]: import os
import re
import logging
from datetime import datetime, date
import numpy as np
import pandas as pd
from pandas import Series, DataFrame
from sklearn.model_selection import train_test_split
from sklearn.model_selection import StratifiedKFold, KFold
from sklearn.metrics import mean_squared_error
import xgboost as xgb
from xgboost import DMatrix
import lightgbm as lgb
from lightgbm import Dataset, Booster
import matplotlib.pyplot as plt
from astral import LocationInfo
from astral.sun import sunrise, sunset, dawn, noon, dusk

logging.basicConfig(format='%(asctime)s : %(levelname)s : %(message)s', level=logging.INFO)
```

# 读取数据

## 比赛数据

```
In [2]: x_train = pd.read_csv("../data/A榜-训练集_分布式光伏发电预测_气象变量数据.csv", encoding='utf-8')
y_train = pd.read_csv("../data/A榜-训练集_分布式光伏发电预测_实际功率数据.csv", encoding='utf-8')
info_train = pd.read_csv("../data/A榜-训练集_分布式光伏发电预测_基本信息.csv", encoding='utf-8')
x_test = pd.read_csv("../data/A榜-测试集_分布式光伏发电预测_气象变量数据.csv", encoding='utf-8')
y_test = pd.read_csv("../data/submit_example.csv", encoding='utf-8')
info_test = pd.read_csv("../data/B榜-测试集_分布式光伏发电预测_基本信息.csv", encoding='utf-8')
x_test_b = pd.read_csv("../data/B榜-测试集_分布式光伏发电预测_气象变量数据.csv", encoding='utf-8')
y_test_b = pd.read_csv("../data/B_submit_example.csv", encoding='utf-8')
info_test_b = pd.read_csv("../data/B榜-测试集_分布式光伏发电预测_基本信息.csv", encoding='utf-8')
```

```
In [3]: x_train = pd.merge(x_train, info_train[["光伏用户编号", "装机容量(kW)", "经度", "纬度", "时间"]], on="光伏用户编号", utc=False).dt.tz_localize('Asia/Shanghai')
x_test = pd.merge(x_test, info_test[["光伏用户编号", "装机容量(kW)", "经度", "纬度", "时间"]], on="光伏用户编号", utc=False).dt.tz_localize('Asia/Shanghai')
x_test_b = pd.merge(x_test_b, info_test_b[["光伏用户编号", "装机容量(kW)", "经度", "纬度", "时间"]], on="光伏用户编号", utc=False).dt.tz_localize('Asia/Shanghai')
```

```
In [4]: y_train = y_train.set_index(["光伏用户编号", "综合倍率", "时间"]).stack().reset_index(level=3)
y_train["level_3"] = y_train["level_3"].apply(lambda x: int(x[1:]))
y_train["时间"] = pd.to_datetime(y_train["时间"], utc=False).dt.tz_localize('Asia/Shanghai')
y_train["时间"] = y_train["时间"] + (y_train["level_3"] - 1) * 15 * pd.Timedelta(hours=1)
y_train = y_train.drop(columns=["level_3"])

y_test = y_test.set_index(["光伏用户编号", "综合倍率", "时间"]).stack().reset_index(level=3)
y_test["level_3"] = y_test["level_3"].apply(lambda x: int(x[1:]))
y_test["时间"] = pd.to_datetime(y_test["时间"], utc=False).dt.tz_localize('Asia/Shanghai')
```

```
y_test["时间"] = y_test["时间"] + (y_test["level_3"] - 1) * 15 * pd.Timedelta(hours=1)
y_test = y_test.drop(columns=["level_3"])

y_test_b = y_test_b.set_index(["光伏用户编号", "综合倍率", "时间"]).stack().reset_index(level=3)
y_test_b["level_3"] = y_test_b["level_3"].apply(lambda x: int(x[1:]))
y_test_b["时间"] = pd.to_datetime(y_test_b["时间"], utc=False).dt.tz_localize('Asia/Shanghai')
y_test_b["时间"] = y_test_b["时间"] + (y_test_b["level_3"] - 1) * 15 * pd.Timedelta(hours=1)
y_test_b = y_test_b.drop(columns=["level_3"])
```

```
In [5]: df_train = pd.merge(x_train, y_train, on=["光伏用户编号", "时间"], how="left")
df_test = pd.merge(x_test, y_test, on=["光伏用户编号", "时间"], how="left")
df_test_b = pd.merge(x_test_b, y_test_b, on=["光伏用户编号", "时间"], how="left")
df = pd.concat([df_train, df_test, df_test_b], axis=0)
logging.info(df.columns)
```

2024-04-08 08:42:18,352 : INFO : Index(['光伏用户编号', '时间', '气压(Pa)', '相对湿度 (%)', '云量', '10米风速 (10m/s)', '10米风向 (°)', '温度 (K)', '辐照强度 (J/m2)', '降水 (m)', '100m风速 (100m/s)', '100m风向 (°)', '装机容量(kW)', '经度', '纬度', '综合倍率', 'target'], dtype='object')

## 外部数据

```
In [6]: outer_files = os.listdir("../data/")
outer_files = [x for x in outer_files if re.match(r"^\^open-meteo-f\d\.csv$", x)]
outer_dfs = list()
for outer_file in outer_files:
    outer_df = pd.read_csv(os.path.join("../data", outer_file), skiprows=3)
    outer_df["time"] = pd.to_datetime(outer_df["time"], utc=False).dt.tz_localize('Asia/Shanghai')
    outer_df = outer_df.set_index("time", drop=True).asfreq("15T")
    outer_df = outer_df.interpolate(method="time", axis=0)

    outer_df = outer_df.reset_index(drop=False).rename(columns={"time": "时间"})
    outer_df["光伏用户编号"] = re.match(r"^\^open-meteo-f\d\.csv$", outer_file).group(1)
    outer_dfs.append(outer_df)
    outer_data = pd.concat(outer_dfs, axis=0)
    logging.info(outer_data.columns)
```

2024-04-08 08:43:05,956 : INFO : Index(['时间', 'temperature\_2m\_archive\_best\_match (°C)', 'relative\_humidity\_2m\_archive\_best\_match (%)', 'dew\_point\_2m\_archive\_best\_match (°C)', 'apparent\_temperature\_archive\_best\_match (°C)', 'precipitation\_archive\_best\_match (mm)', 'rain\_archive\_best\_match (mm)', 'snowfall\_archive\_best\_match (cm)', 'snow\_depth\_archive\_best\_match (m)', 'weather\_code\_archive\_best\_match (wmo code)', ... 'direct\_normal\_irradiance\_era5\_land (W/m²)', 'global\_tilted\_irradiance\_era5\_land (W/m²)', 'terrestrial\_radiation\_era5\_land (W/m²)', 'shortwave\_radiation\_instant\_era5\_land (W/m²)', 'direct\_radiation\_instant\_era5\_land (W/m²)', 'diffuse\_radiation\_instant\_era5\_land (W/m²)', 'direct\_normal\_irradiance\_instant\_era5\_land (W/m²)', 'global\_tilted\_irradiance\_instant\_era5\_land (W/m²)', 'terrestrial\_radiation\_instant\_era5\_land (W/m²)', '光伏用户编号'], dtype='object', length=222)

```
In [7]:
outer_day_files = os.listdir("../data/")
outer_day_files = [x for x in outer_day_files if re.match(r"^\open-meteo-day-f\d\
outer_day_dfs = list()
for outer_day_file in outer_day_files:
    outer_day_df = pd.read_csv(os.path.join("../data", outer_day_file), skiprows
    outer_day_df["time"] = pd.to_datetime(outer_day_df["time"], utc=False).dt.tz
    outer_day_df = outer_day_df.rename(columns={"time": "时间"})
    outer_day_df["光伏用户编号"] = re.match(r"^\open-meteo-day-(f\d)\.csv$", outer
    outer_dfs.append(outer_day_df)
outer_day_data = pd.concat(outer_day_dfs, axis=0)
logging.info(outer_day_data.columns)
```

```
2024-04-08 08:43:06,090 : INFO : Index(['时间', 'weather_code (wmo code)', 'tempe
rature_2m_max (°C)',
'temperature_2m_min (°C)', 'temperature_2m_mean (°C)',
'apparent_temperature_max (°C)', 'apparent_temperature_min (°C)',
'apparent_temperature_mean (°C)', 'sunrise (iso8661)',
'sunset (iso8661)', 'daylight_duration (s)', 'sunshine_duration (s)',
'precipitation_sum (mm)', 'rain_sum (mm)', 'snowfall_sum (cm)',
'precipitation_hours (h)', 'wind_speed_10m_max (km/h)',
'wind_gusts_10m_max (km/h)', 'wind_direction_10m_dominant (°)',
'shortwave_radiation_sum (MJ/m²)', 'et0_fao_evapotranspiration (mm)',
'光伏用户编号'],
dtype='object')
```

```
In [ ]:
outer_data["日期"] = outer_data["时间"].dt.date
outer_day_data["日期"] = outer_day_data["时间"].dt.date
outer_day_data = outer_day_data.drop(columns=["时间"])
outer_data = pd.merge(outer_data, outer_day_data, on=["光伏用户编号", "日期"], ho
for column in outer_data.columns:
    logging.info(column)
```

整合数据

```
In [9]:
df = pd.merge(df, outer_data[["光伏用户编号",
    "时间",
    "apparent_temperature_archive_best_match (°C)",
    "surface_pressure_archive_best_match (hPa)",
    "is_day_archive_best_match ()",
    "relative_humidity_2m_archive_best_match (%)",
    "cloud_cover_archive_best_match (%)",
    "global_tilted_irradiance_instant_archive_best_match (W/m²)",
    "precipitation_archive_best_match (mm)",
    "wind_speed_10m_archive_best_match (km/h)",
    "wind_speed_100m_archive_best_match (km/h)",
    "wind_direction_10m_archive_best_match (°)",
    "wind_direction_100m_archive_best_match (°)",
    "sunshine_duration_archive_best_match (s)",
    "weather_code_archive_best_match (wmo code)",
    "temperature_2m_max (°C)",
    "temperature_2m_min (°C)",
    "temperature_2m_mean (°C)",
    "precipitation_sum (mm)",
    "rain_sum (mm)",
    "snowfall_sum (cm)",
    "apparent_temperature_max (°C)",
    "apparent_temperature_min (°C)",
```

```
"apparent_temperature_mean (°C)",
"sunrise (iso8661)",
"shortwave_radiation_archive_best_match (W/m²)",
"direct_radiation_archive_best_match (W/m²)",
"diffuse_radiation_archive_best_match (W/m²)",
"direct_normal_irradiance_archive_best_match (W/m²)",
"global_tilted_irradiance_archive_best_match (W/m²)",
"terrestrial_radiation_archive_best_match (W/m²)"
]], how="left", on=["光伏用户编号", "时间"])

df["温度 (K)"] = df["apparent_temperature_archive_best_match (°C)"] + 273.15
df["气压 (Pa)"] = df["surface_pressure_archive_best_match (hPa)"] * 100
df["是白天"] = df["is_day_archive_best_match ()"].copy()
df["相对湿度 (%)"] = df["relative_humidity_2m_archive_best_match (%)"].copy()
df["云量"] = df["cloud_cover_archive_best_match (%)"] / 100
df["日照强度 (J/m2)"] = df["global_tilted_irradiance_instant_archive_best_match
df["降水 (m)"] = df["precipitation_archive_best_match (mm)] / 1000
df["10米风速 (10m/s)"] = df["wind_speed_10m_archive_best_match (km/h)"] / 3600
df["100米风速 (100m/s)"] = df["wind_speed_100m_archive_best_match (km/h)"] / 3600
df["10米风向 (°)"] = df["wind_direction_10m_archive_best_match (°)"].copy()
df["100米风向 (°)"] = df["wind_direction_100m_archive_best_match (°)"].copy()
df["sunrise (iso8661)"] = pd.to_datetime(df["sunrise (iso8661)"], utc=False).dt.
df["sunset (iso8661)"] = pd.to_datetime(df["sunset (iso8661)"], utc=False).dt.
df["时间-日出时间"] = ((df["时间"] - df["sunrise (iso8661)"]).dt.days * 24 * 3600
df["日落时间-时间"] = ((df["sunset (iso8661)"] - df["时间"]).dt.days * 24 * 3600
```

```
df = df.drop(columns=[
    "apparent_temperature_archive_best_match (°C)",
    "surface_pressure_archive_best_match (hPa)",
    "is_day_archive_best_match ()",
    "relative_humidity_2m_archive_best_match (%)",
    "cloud_cover_archive_best_match (%)",
    "global_tilted_irradiance_instant_archive_best_match (W/m²)",
    "precipitation_archive_best_match (mm)",
    "wind_speed_10m_archive_best_match (km/h)",
    "wind_speed_100m_archive_best_match (km/h)",
    "wind_direction_10m_archive_best_match (°)",
    "wind_direction_100m_archive_best_match (°)",
    "sunrise (iso8661)",
    "sunset (iso8661)"
])
```

特征工程

时间特征

```
In [10]:
df["年"] = df["时间"].dt.year
df["季节"] = df["时间"].dt.quarter
df["月"] = df["时间"].dt.month
df["日"] = df["时间"].dt.day
df["周"] = df["时间"].dt.week
df["分"] = df["时间"].dt.minute // 15 + df["时间"].dt.hour * 4
```

```
C:\Program Files\Python37\lib\site-packages\ipykernel_launcher.py:5: FutureWarning: Series.dt.weekofyear and Series.dt.week have been deprecated. Please use Series.dt.isocalendar().week instead.
"""
```

```
In [11]: # df["分_"] = df["分"].copy()
# df = pd.get_dummies(df, columns=["分_"], prefix_sep="")
df["分"] = df["分"].astype("category")
```

## 根据日出时间日落时间计算时间段

```
In [12]: def solar_time(current_time, dawn_time, sunrise_time, noon_time, sunset_time, dusk_time):
"""
根据太阳判断当前时间段\n
"""
    if dawn_time < current_time < sunrise_time:
        return 1
    elif sunrise_time <= current_time < noon_time:
        return 2
    elif noon_time <= current_time < sunset_time:
        return 3
    elif sunset_time <= current_time < dusk_time:
        return 4
    else:
        return 0
```

```
In [13]: # df["地点"] = df.apply(lambda x: LocationInfo(name=x["光伏用户编号"], region="Ch
# df["黎明时刻"] = df.apply(lambda x: dawn(x["地点"]).observer, date=x["时间"], tz
# df["日出时刻"] = df.apply(lambda x: sunrise(x["地点"]).observer, date=x["时间"], tz
# df["正午时刻"] = df.apply(lambda x: noon(x["地点"]).observer, date=x["时间"], tz
# df["日落时刻"] = df.apply(lambda x: sunset(x["地点"]).observer, date=x["时间"], tz
# df["黄昏时刻"] = df.apply(lambda x: dusk(x["地点"]).observer, date=x["时间"], tz
# df["时间段"] = df.apply(lambda x: solar_time(x["时间"], x["黎明时刻"], x["日出时
# df = df.drop(columns=["地点", "黎明时刻", "日出时刻", "正午时刻", "日落时刻", "黄
```

## 光伏用户编号

```
In [14]: df["光伏用户编号_"] = df["光伏用户编号"].copy()
df = pd.get_dummies(df, columns=["光伏用户编号_"], prefix_sep="")
```

## 气象特征

```
In [15]: df['100m风速 (100m/s)'] = df['100m风速 (100m/s)'] * np.sin(np.pi * df['100m风向
# df['cos_100m风速 (100m/s)'] = df['100m风速 (100m/s)'] * np.cos(np.pi * df['10
df['10米风速 (10m/s)'] = df['10米风速 (10m/s)'] * np.sin(np.pi * df['10米风向 (°
# df['cos_10米风速 (10m/s)'] = df['10米风速 (10m/s)'] * np.cos(np.pi * df['10米风
```

```
In [16]: df["光照/温度"] = df["辐照强度 (J/m2)"] / df["温度 (K)"]
```

## 历史值特征

```
In [17]: dfs = []
for site, df_site in df.groupby("光伏用户编号"):
    df_site = df_site.sort_values("时间")
```

```
df_site["辐照强度 (J/m2) - 1"] = df_site["辐照强度 (J/m2)"].shift(1) - df_si
df_site["辐照强度 (J/m2) - 8"] = df_site["辐照强度 (J/m2)"].shift(8) - df_si
# df_site["辐照强度 (J/m2) - 2"] = df_site["辐照强度 (J/m2)"].shift(2) - df_
dfs.append(df_site)
df = pd.concat(dfs, axis=0)
```

## 处理异常值

```
In [18]: print(df_train["target"].nsmallest(3))
df_train[df_train["target"] < -8]
df.loc[207628:207633, ["光伏用户编号", "时间", "target"]]
```

207630 -8.89000  
37316 -0.0085  
39139 -0.0084

Name: target, dtype: float64

	光伏用户编号	时间	target
207628	f6	2022-08-15 20:30:00+08:00	-0.002
207629	f6	2022-08-15 20:45:00+08:00	-0.002
207630	f6	2022-08-15 21:00:00+08:00	-8.890
207631	f6	2022-08-15 21:15:00+08:00	-0.002
207632	f6	2022-08-15 21:30:00+08:00	NaN
207633	f6	2022-08-15 21:45:00+08:00	NaN

```
In [19]: df.loc[207630, "target"] = -0.002
```

## 光照与当天最强光照的比值

```
In [20]: df["日期"] = df["时间"].dt.date
day_max_values = df[["光伏用户编号", "日期", "辐照强度 (J/m2)"]].groupby(by=["光伏
day_max_values = day_max_values.rename(columns={x: x + "_max" for x in day_max_v
df = pd.merge(df, day_max_values, on=["光伏用户编号", "日期"], how="left").drop(c
df["辐照强度 (J/m2) _max"] = df["辐照强度 (J/m2)"] / df["辐照强度 (J/m2) _max"]
```

## 当天的平均光照

```
In [21]: df["日期"] = df["时间"].dt.date
day_mean_values = df[["光伏用户编号", "日期", "辐照强度 (J/m2)"]].groupb
day_mean_values = day_mean_values.rename(columns={x: x + "_mean" for x in day_me
df = pd.merge(df, day_mean_values, on=["光伏用户编号", "日期", "是白天"], how="lef
```

## 温度与当天最高温最低度的差值

```
In [22]: df["日期"] = df["时间"].dt.date
day_max_values = df[["光伏用户编号", "日期", "温度 (K)"]].groupby(by=["光伏用户编
day_min_values = df[["光伏用户编号", "日期", "温度 (K)"]].groupby(by=["光伏用户编
day_max_values = day_max_values.rename(columns={x: x + "_max" for x in day_max_v
day_min_values = day_min_values.rename(columns={x: x + "_min" for x in day_min_v
df = pd.merge(df, day_max_values, on=["光伏用户编号", "日期"], how="left")
df = pd.merge(df, day_min_values, on=["光伏用户编号", "日期"], how="left").drop(c
```

```
df["温度 (K) _max"] = df["温度 (K) _max"] - df["温度 (K) "]
df["温度 (K) _min"] = df["温度 (K) "] - df["温度 (K) _min"]
df = df.rename(columns={
    "辐射强度 (J/m2) _max": "光照/当天最强光照",
    "温度 (K) _max": "与当天最高温度之差",
    "温度 (K) _min": "与当天最低温度之差"
})
```

## 划分测试集

```
In [23]: df_train = df[(df["时间"] >= df_train["时间"].min()) & (df["时间"] <= df_train["时间"]
df_test = df[(df["时间"] >= df_test["时间"].min()) & (df["时间"] <= df_test["时间"]
df_test_b = df[(df["时间"] >= df_test_b["时间"].min()) & (df["时间"] <= df_test_b
```

## 训练模型

### 评价指标

```
In [24]: def score(y_true, y_pred):
mse = mean_squared_error(y_true, y_pred)
rmse = np.sqrt(mse)
return 1 / (1 + rmse)
```

## lightgbm模型

```
In [25]: params_lgb = {
    "num_boost_round": 1000,
    "learning_rate": 0.02,
    "boosting_type": 'gbdt',
    'objective': 'mse',
    'metric': 'rmse',
    'num_leaves': 127,
    'verbosity': -1,
    'seed': 42,
    'n_jobs': -1,
    'feature_fraction': 0.8,
    'bagging_fraction': 0.9,
    'bagging_freq': 4,
    "early_stopping_round": 100
}
model_lgb = []
```

## xgboost模型

```
In [26]: params_xgb = {
    "num_boost_round": 500,
    "learning_rate": 0.02,
    "booster": "gbtree",
    "objective": "reg:squarederror",
    "eval_metric": "rmse",
    "max_leaves": 127,
    "verbosity": 1,
    "seed": 42,
```

```
"nthread": -1,
"colsample_bytree": 0.6,
"subsample": 0.7,
"early_stopping_rounds": 100
}
model_xgb = []
```

## 交叉验证

```
In [ ]: kfold = KFold(n_splits=5, random_state=42, shuffle=True)

x = df_train.drop(columns=["光伏用户编号", "时间"]).dropna().astype(np.float32)
y = x.pop("target")
mse = 0
for fold, (train_index, val_index) in enumerate(kfold.split(x, y)):
    logging.info(f'##### fold: {fold} #####')
    x_train, x_val, y_train, y_val = x.iloc[train_index], x.iloc[val_index], y.i

    trainset = Dataset(x_train, y_train)
    valset = Dataset(x_val, y_val)
    model = lgb.train(params_lgb, trainset, valid_sets=[trainset, valset], cate
    model.save_model("./models/lgb_%d.txt" % fold)
    model_lgb.append(model)
    lgb_pred = Series(model.predict(x_val, num_iteration=model.best_iteration),

trainset = DMatrix(x_train, y_train, enable_categorical=True, nthread=-1)
valset = DMatrix(x_val, y_val, enable_categorical=True, nthread=-1)
model = xgb.train(params_xgb, trainset, evals=[(trainset, 'train')],(valset,
model.save_model("./models/xgb_%d.json" % fold)
model_xgb.append(model)
xgb_pred = Series(model.predict(x_val, num_iteration=model.best_iteration),

val_pred = (lgb_pred + xgb_pred) / 2
mse += mean_squared_error(y_val.fillna(0), val_pred)
rmse = np.sqrt(mse / kfold.n_splits)
score = 1 / (1 + rmse)
logging.info(f'-----本地分数 {score}-----')
```

## 预测

```
In [28]: x_test_b = df_test_b.drop(columns=["光伏用户编号", "时间"]).astype(np.float32)
y_test_b = x_test_b.pop("target")
y_pred = np.zeros((df_test_b.shape[0], ))
for i in range(0, kfold.n_splits):
    y_pred += model_lgb[i].predict(x_test_b, num_iteration=model_lgb[i].best_ite
    y_pred += model_xgb[i].predict(DMatrix(x_test_b, enable_categorical=True, nt
y_pred = y_pred / 2 / kfold.n_splits
df_test_b["target"] = y_pred
```

C:\Program Files\Python37\lib\site-packages\ipykernel\_launcher.py:8: SettingWithC
opyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stabl
e/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stabl
e/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
In [29]: df_test_b = df_test_b["光伏用户编号", "综合倍率", "年", "月", "日", "分", "target",
df_test_b["时间"] = df_test_b["年"].astype(str) + "-" + df_test_b["月"].astype(st
df_test_b["分"] = "p" + (df_test_b["分"].astype(int) + 1).astype(str)
df_test_b = df_test_b.drop(columns=["年", "月", "日"])
```

```
In [30]: result = pd.pivot(df_test_b, index=["光伏用户编号", "综合倍率", "时间"], columns="
result = result[result["综合倍率"].notnull()]
result["综合倍率"] = result["综合倍率"].astype(int)
```

```
In [31]: result.to_csv("../data/%s.csv" % datetime.now().strftime("%Y%m%d_%H%M%S"), encod
```

In [ ]: