

07时间日期

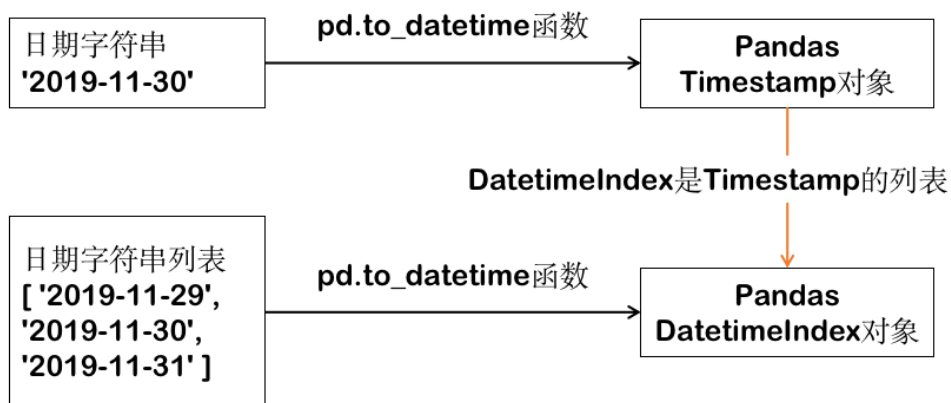
Pandas日期处理的作用：将2018-01-01、1/1/2018等多种日期格式映射成统一的格式对象，在该对象上提供强大的功能支持

几个概念：

1. `pd.to_datetime`：pandas的一个函数，能将字符串、列表、series变成日期形式
2. `Timestamp`：pandas表示日期的对象形式
3. `DatetimeIndex`：pandas表示日期的对象列表形式

其中：

- `DatetimeIndex` 是 `Timestamp` 的列表形式
- `pd.to_datetime` 对单个日期字符串处理会得到 `Timestamp`
- `pd.to_datetime` 对日期字符串列表处理会得到 `DatetimeIndex`



一、问题：怎样统计每周、每月、每季度的最高温度？

1、读取天气数据到dataframe

```
1 fpath = "res/beijing_tianqi_2018.csv"
2 df = pd.read_csv(fpath)
3 # 替换掉温度的后缀℃
4 df.loc[:, "bwendu"] = df["bwendu"].str.replace("℃",
5     "").astype('int32')
6 df.loc[:, "ywendu"] = df["ywendu"].str.replace("℃",
7     "").astype('int32')
8 df.head()
```

2、将日期列转换成pandas的日期

```
1 df.set_index(pd.to_datetime(df["ymd"]), inplace=True)
2 df.head()
3 # DatetimeIndex是Timestamp的列表形式
4 df.index[0]
```

3、方便的对DatetimeIndex进行查询

```
1 # 筛选固定的某一天
2 df.loc['2018-01-05']
3 # 日期区间
4 df.loc['2018-01-05':'2018-01-10']
5 # 按月份前缀筛选
6 df.loc['2018-03']
7 # 按月份前缀筛选
8 df.loc["2018-07":"2018-09"]
9 # 按年份前缀筛选
10 df.loc["2018"]
```

4、方便的获取周、月、季度

`Timestamp`、`DatetimeIndex`支持大量的属性可以获取日期分量：

https://pandas.pydata.org/pandas-docs/stable/user_guide/timeseries.html#time-date-components

```
1 # 周数字列表
2 df.index.week
3 # 月数字列表
4 df.index.month
5 # 季度数字列表
6 df.index.quarter
```

5、统计每周、每月、每个季度的最高温度

```
1 # 统计每周的数据
2 df.groupby(df.index.week)["bwendu"].max()
3 # 统计每个月的数据
4 df.groupby(df.index.month)["bwendu"].max()
5 # 统计每个季度的数据
6 df.groupby(df.index.quarter)["bwendu"].max()
```

二、Pandas怎么处理日期索引的缺失？

可以用两种方法实现：

- 1、`DataFrame.reindex`，调整dataframe的索引以适应新的索引
- 2、`DataFrame.resample`，可以对时间序列重采样，支持补充缺失值

0.构造数据

```
1 df = pd.DataFrame({
2     "pdate": ["2019-12-01", "2019-12-02", "2019-12-04", "2019-
3     12-05"],
4     "pv": [100, 200, 400, 500],
5     "uv": [10, 20, 40, 50],
6 })
7 df # 缺失了2019-12-03的数据
8 df.set_index("pdate").plot()
9 plt.show()
```

方法1：使用 `pandas.reindex` 方法

```
1 df_date = df.set_index("pdate")
2 df_date
3 df_date.index
4 # 1.将df的索引设置为日期索引
5 df_date = df_date.set_index(pd.to_datetime(df_date.index))
6 df_date
7 df_date.index
8 # 2.使用pandas.reindex填充缺失的索引
9 # 生成完整的日期序列
10 pdates = pd.date_range(start="2019-12-01", end="2019-12-05")
11 pdates
12 # 填充0
13 df_date_new = df_date.reindex(pdates, fill_value=0)
14 df_date_new
15 df_date_new.plot() # 画图
16 plt.show()
```

方法2：使用 `pandas.resample` 方法

`resample` 的含义：

改变数据的时间频率，比如把天数据变成月份，或者把小时数据变成分钟级别

`resample` 的语法：

```
(DataFrame or Series).resample(arguments).(aggregate function)
```

`resample` 的采样规则参数：

https://pandas.pydata.org/pandas-docs/stable/user_guide/timeseries.html#offset-aliases

```
1 # 1.先将索引变成日期索引
2 df_new2 =
  df.set_index(pd.to_datetime(df["pdate"])).drop("pdate", axis=1)
3 df_new2
4 df_new2.index
5 # 2.使用resample的方法按照天重采样
6 # 由于采样会让区间变成一个值，所以需要指定mean等采样值的设定方法
7 df_new2 = df_new2.resample("D").mean().fillna(0)
8 df_new2
9 # resample的使用方式
10 df_new2.resample("2D").mean()
```