

Pandas 的数据读取与输出

Guangyao Zhao

2022-12-01

Contents

| | |
|-------------|----------|
| 数据读取 | 1 |
| 数据输出 | 4 |

Pandas 提供了丰富的数据类型读取接口，比如 CSV, Excel, JSON 等。在此只介绍前两者。

数据读取

`read_csv()` 的参数如下：

```
1 pandas.read_csv(filepath_or_buffer, # 文件路径
2                 sep=_NoDefault.no_default, # 分隔符, 和 delimiter 一样, 两者选其一即可
3                 delimiter=None,
4                 header='infer', # 表头, 默认为第一行
5                 names=_NoDefault.no_default, # 如果表头设置为 None, 则可自定义表头名称, 即 columns 名称
6                 index_col=None, # 索引列, 默认无
7                 usecols=None, # 选择要使用的列
8                 squeeze=None,
9                 prefix=_NoDefault.no_default, # 给表头添加前缀, 比如原表头名称为 '1,2,3'; 添加前缀 'ORP_'
10                mangle_dupe_cols=True, # 表头名称如果重复, 比如有两列 'ORP', 使用该选项后可将其转化为 'ORP_'
11                dtype=None, # 数据类型
12                engine=None,
13                converters=None,
14                true_values=None,
```

15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

```
false_values=None,  
skipinitialspace=False,  
skiprows=None, # 跳过指定行  
skipfooter=0,  
nrows=None,  
na_values=None, # 空值替换, 若想将 '5, 5.0' 设置为空值, 则可使用 na_values=[5, 5.0]  
keep_default_na=True, # 是否保留 NaN 为空值, 此选项和 na_values 共同决定空值有哪些  
na_filter=True, # 是否检查空值, 如果确定文件中无空值, 则可将其设置为 False, 提高运行速度  
verbose=False,  
skip_blank_lines=True, # 是否跳过空行  
parse_dates=None, # 是否对时间日期进行解析  
infer_datetime_format=False, # 如果开启了 parse_dates 选项, 而且 datetime 字符串的格式都  
keep_date_col=False,  
date_parser=None,  
dayfirst=False,  
cache_dates=True,  
iterator=False,  
chunksize=None,  
compression='infer',  
thousands=None, # 是否设置千分位分隔符  
decimal='.', # 小数点标识符  
lineterminator=None, # 行结束标识符, 比如 'a,b,c~1,2,3~4,5,6', 将 lineterminator='~', 则读  
quotechar='\"',  
quoting=0,  
doublequote=True,  
escapechar=None,  
comment=None, # 注释标识符, 比如 comment='#', 则会被忽视  
encoding=None,  
encoding_errors='strict',  
dialect=None,  
error_bad_lines=None,  
warn_bad_lines=None,  
on_bad_lines=None,  
delim_whitespace=False, # 空格分隔符, 和 delimiter, sep 冲突, 相当于 sep='\s+'  
low_memory=True,  
memory_map=False,
```

```
51 float_precision=None,  
52 storage_options=None)
```

和 `read_csv()` 相比, `read_excel()` 仅仅多了一个 `sheet_name()` 选项:

```
1 pandas.read_excel(io,  
2     sheet_name=0, # 要读取的 sheet 名称, 可以是一个也可以是多个, 如果不指定则默认第一个。int  
3     header=0,  
4     names=None,  
5     index_col=None,  
6     usecols=None,  
7     squeeze=None,  
8     dtype=None,  
9     engine=None,  
10    converters=None,  
11    true_values=None,  
12    false_values=None,  
13    skiprows=None,  
14    nrows=None,  
15    na_values=None,  
16    keep_default_na=True,  
17    na_filter=True,  
18    verbose=False,  
19    parse_dates=False,  
20    date_parser=None,  
21    thousands=None,  
22    decimal='.',  
23    comment=None,  
24    skipfooter=0,  
25    convert_float=None,  
26    mangle_dupe_cols=True,  
27    storage_options=None)
```

可以看到大部分命令和 `read_csv` 基本一致, 但是多了 `sheet_name` 选项。如果 `csv` 文件格式能满足需求, 则尽量避免使用 `excel`, 因为前者更通用, 读取数据快且处理方法更丰富。

数据输出

处理后的表格也可以进行相应的输出。to_csv() 的参数如下：

```
1 DataFrame.to_csv(path_or_buf=None, # 路径
2                 sep=',', # 分隔符
3                 na_rep='', # 空值代替符号
4                 float_format=None,
5                 columns=None, # 表头名称
6                 header=True, # 是否含有表头, 和 columns 结合使用, 也可直接在此传入表头名称
7                 index=True,
8                 index_label=None,
9                 mode='w',
10                encoding=None,
11                compression='infer', # 是否压缩, 不常用
12                quoting=None,
13                quotechar='"',
14                lineterminator=None,
15                chunksize=None,
16                date_format=None,
17                doublequote=True,
18                escapechar=None,
19                decimal='.',
20                errors='strict',
21                storage_options=None)
```

to_excel() 和 to_csv() 相比, 特殊之处在于前者支持多 sheet:

```
1 with pandas.ExcelWriter('xxx.xlsx') as writer:
2     df1.to_excel(writer, sheet_name='name1')
3     df2.to_excel(writer, sheet_name='name2')
```