

# | STAT243 Lecture 2.2 Reading data from text files into Python

## | 1 Core Python functions

- 使用 Pandas 读取文本数据：
  - `read_table` 与 `read_csv` 读取分隔符文本；  
关键参数：
    - `sep` 分隔符
    - `header` 是否有表头。
  - 定宽文本使用 `read_fwf()` 读取为 DataFrame。
  - 类型推断：Pandas 会自动推断，但更稳妥做法是通过 `dtype` 明确指定各列类型。
  - `sep` 支持正则：可以用正则表达式按空白或复杂分隔模式拆分。
- 示例 1： `RTADDataSub.csv` 初读入与缺失处理。

```
Python
1 import os, pandas as pd
2
3 # 初次读取：不指定 dtype 时，多数列为 object
4 dat = pd.read_table(os.path.join('.', 'data', 'RTADDataSub.csv'),
5                    sep = ',', header = None)
6 print(dat.dtypes.head())
7 # Output:
8 # 0    object
9 # 1    object
10 # 2    object
11 # 3    object
12 # 4    object
13 # dtype: object
14
15 print(dat.loc[0,1])
16 # Output:
17 # 2336
18
19 print(type(dat.loc[0,1]))
20 # Output:
21 # <class 'str'>
22
23 # 将 'x' 视为缺失
24 dat2 = pd.read_table(os.path.join('.', 'data', 'RTADDataSub.csv'),
25                    sep = ',', header = None, na_values = 'x')
26 print(dat2.dtypes.head())
27 # Output:
28 # 0    object
29 # 1    float64
30 # 2    float64
31 # 3    float64
32 # 4    float64
33 # dtype: object
34
35 print(dat2.loc[:,1].unique())
36 # Output:
37 # array([2336.,   nan, 1450., ...])
```

- 示例 2： `hivSequ.csv` 中通过 `dtype` 精确控制多列类型。



## Python

```
1 dat = pd.read_table(os.path.join '..', 'data', 'hivSequ.csv'),
2                      sep = ',', header = 0,
3                      dtype = {
4                          'PatientID': int,
5                          'Resp': int,
6                          'PR Seq': str,
7                          'RT Seq': str,
8                          'VL-t0': float,
9                          'CD4-t0': int
10                     })
11 print(dat.dtypes)
12 # Output:
13 # PatientID      int64
14 # Resp           int64
15 # PR Seq         object
16 # RT Seq         object
17 # VL-t0         float64
18 # CD4-t0         int64
19 # dtype: object
20
21 print(dat.loc[0, 'PR Seq'])
22 # Output:
23 # CCTCAATCACTCTT...
```

- 只读取部分列：使用 `usecols` 指定需要的列可以减少内存与 I/O。
- 读入前检查：建议先在终端用 `less` 等预览原始文件，以便发现分隔符、缺失标记、列数不齐等问题。
- 若文件列宽不齐（ragged lines）：可以先逐行读为字符串，再按切片处理固定位置字段。



## Python

```
1 file_path = os.path.join '..', 'data', 'precip.txt')
2 with open(file_path, 'r') as file:
3     lines = file.readlines()
4
5 id = [line[3:11] for line in lines]
6 year = [int(line[17:21]) for line in lines]
7 month = [int(line[21:23]) for line in lines]
8 nvalues = [int(line[27:30]) for line in lines]
9 print(year[0:5])
10 # Output:
11 # [2010, 2010, 2010, 2010, 2010]
```

- 上述 `precip.txt` 实际为 fixed-width，使用 `pandas.read_fwf()` 更合适。

### ⚠ Remark ▾

`with` 语句是标准的 Python 文件上下文管理方式，离开代码块会自动关闭文件句柄。

## 2 Connections and streaming

- 除了普通文件，还可以从不同连接读取：压缩文件、归档、子进程输出、网络资源等。



## Python

```
1 import gzip
2 with gzip.open('dat.csv.gz', 'r') as file:
3     lines = file.readlines()
```

```

4 # Output:
5 # [b'...first line...', b'...'] # 示例
6
7 import zipfile
8 with zipfile.ZipFile('dat.zip', 'r') as archive:
9     with archive.open('data.txt', 'r') as file:
10         lines = file.readlines()
11 # Output:
12 # [b'...first line...', b'...'] # 示例
13
14 import subprocess, io
15 command = "ls -al"
16 output = subprocess.check_output(command, shell = True)
17 with io.BytesIO(output) as stream:
18     content = stream.readlines()
19 # Output:
20 # [b'total ...', b'drwxr-xr-x ...', ...] # 示例
21
22 # 直接从 URL 读取 (制表符分隔)
23 df = pd.read_csv("https://download.bls.gov/pub/time.series/cu/cu.item", sep="\t")
24 print(df.shape)
25 # Output:
26 # (N, M) # 示例: 行列尺寸

```

- 流式/分块读取 (online processing / streaming / chunking): 适用于大文件。



#### Python

```

1 file_path = os.path.join('..', 'data', 'RTADDataSub.csv')
2 chunksize = 50
3 with pd.read_csv(file_path, chunksize = chunksize) as reader:
4     for chunk in reader:
5         print(f'Read {len(chunk)} rows.')
6         # Output:
7         # Read 50 rows.

```

- 将字符串当作类文件对象读取, 便于对接只接受 file-like 的 API。



#### Python

```

1 file_path = os.path.join('..', 'data', 'precip.txt')
2 with open(file_path, 'r') as file:
3     text = file.read()
4
5 StringIOtext = io.StringIO(text)
6
7 df = pd.read_fwf(StringIOtext, header = None, widths = [3,8,4,2,4,2])
8 print(df.head().shape)
9 # Output:
10 # (5, 6) # 示例: 前五行的形状

```

- 也可以创建连接用于写出, 但需要先显式打开连接 (写出在下一节)。

## 3 File paths

1. 避免在代码中硬编码绝对路径, 改用项目内的相对路径。



#### Python

```

1 dat = pd.read_csv('../data/cpds.csv')

```

```

2 print(dat.shape)
3 # Output:
4 # (N, M) # 示例

```

2. 使用 UNIX 风格分隔符具有跨平台性；Windows 风格 `\` 在 Linux/Mac 上不可用。

```

Python
1 # good: 跨平台
2 pd.read_csv('../data/cpds.csv')
3 # Output:
4 # (N, M) # 示例
5
6 # bad: 在 Linux/Mac 上不可用
7 pd.read_csv('../\\data\\cpds.csv')
8 # Output:
9 # FileNotFoundError: [Errno 2] No such file or directory: '../\\data\\cpds.csv' # 示例

```

3. 更好的方法是：使用 `os.path.join` 构造不依赖于操作系统的路径。

```

Python
1 pd.read_csv(os.path.join '..', 'data', 'cpds.csv'))

```

## 4 Reading data quickly: Arrow and Polars

- Apache Arrow：列式内存布局，Python 通过 PyArrow 调用。
  - 同列值连续存放，支持高效访问。
  - 可从 Parquet、Arrow 原生格式、文本读取。
  - 按需读取，避免整表常驻内存。
- 其他避免全量入内存方案：`Dask`、`numpy.load(mmap_mode=...)`。
- polars：强调速度的 DataFrame 库，常见对比示例如下。

```

Python
1 import polars, time
2
3 t0 = time.time()
4 dat = pd.read_csv(os.path.join '..', 'data', 'airline.csv'))
5 t1 = time.time()
6
7 dat2 = polars.read_csv(os.path.join '..', 'data', 'airline.csv'), null_values = ['NA'])
8 t2 = time.time()
9 print(f"Timing for Pandas: {t1-t0}.")
10 # Output:
11 # Timing for Pandas: 0.99.
12 print(f"Timing for Polars: {t2-t1}.")
13 # Output:
14 # Timing for Polars: 0.91.

```