

当我们遇到一个超级大的DataFrame，里面有一列类型为字符串，要将每一行的字符串都用同一方式进行处理，一般会想到遍历整合DataFrame，但是如果直接这样做的话将会耗费很长时间，有时几个小时都处理不完。因此我们将学习pandas快速处理字符串方法。

✧ 一 向量化字符串操作简介

量化操作简化了纯数值的数组操作语法，我们不需要再担心数组的长度或维度，只需要把中心放在操作上面。

而对字符串的向量化需要工具包的支持，如Numpy就没办法直接对字符串进行向量化操作，只能通过繁琐的循环来实现。

Pandas则可以很好的处理这类问题。

✧ 二、str方法的简介

Python会处理字符串起来会很容易，作为工具包的Pandas同样可以简单快速的处理字符串，

几乎把Python内置的字符串方法都给复制过来了，这种方法就是Pandas内置的str方法，

通俗来说就可以将series和index对象中包含字符串的部分简单看作单个字符串处理，达到批量简单快速处理的目的

- `lower()` 将的字符串转换为小写。
- `upper()` 将的字符串转换为大写。
- `len()` 得出字符串的长度。
- `strip()` 去除字符串两边的空格（包含换行符）。
- `split()` 用指定的分割符分割字符串。
- `cat(sep="")` 用给定的分隔符连接字符串元素。
- `contains(pattern)` 如果子字符串包含在元素中，则为每个元素返回一个布尔值 True，否则为 False。

- `replace(a,b)` 将值 a 替换为值 b。
- `count(pattern)` 返回每个字符串元素出现的次数。
- `startswith(pattern)` 如果 Series 中的元素以指定的字符串开头，则返回 True。
- `endswith(pattern)` 如果 Series 中的元素以指定的字符串结尾，则返回 True。
- `findall(pattern)` 以列表的形式返回出现的字符串。
- `find(pattern)` 返回字符串第一次出现的索引位置。

注意：上述所有字符串函数全部适用于 `DataFrame` 对象，同时也可以与 Python 内置的字符串函数一起使用，这些函数在处理 Series/DataFrame 对象的时候会自动忽略缺失值数据（NaN）

```
1 import pandas as pd
2 import numpy as np
```

- `lower()` 将的字符串转换为小写。

```
1 s = pd.Series(['C', 'Python', 'java', 'go', np.nan,
2               '1125', 'javascript'])
3 s.str.lower()
```

```
1 0      c
2 1    python
3 2     java
4 3      go
5 4     NaN
6 5    1125
7 6  javascript
8 dtype: object
```

- `upper()` 将的字符串转换为大写。

```
1 s = pd.Series(['C', 'Python', 'java', 'go', np.nan,
2               '1125', 'javascript'])
3 s.str.upper()
```

```

1 0          C
2 1      PYTHON
3 2          JAVA
4 3          GO
5 4          NaN
6 5          1125
7 6  JAVASCRIPT
8 dtype: object

```

- `len()` 得出字符串的长度。

```

1 s = pd.Series(['C', 'Python', 'java', 'go', np.nan,
2               '1125', 'javascript'])
3 s.str.len()

```

```

1 0      1.0
2 1      6.0
3 2      4.0
4 3      2.0
5 4      NaN
6 5      4.0
7 6     10.0
8 dtype: float64

```

- `strip()` 去除字符串两边的空格（包含换行符）。

```

1 s = pd.Series(['C ', ' Python\t', '   java   ', 'go\t', np.nan,
2               '\t1125 ', '\tjavascript'])
3 s.str.strip()

```

```

1 0          C
2 1      Python
3 2          java
4 3          go
5 4          NaN
6 5          1125
7 6  javascript
8 dtype: object

```

- `split()` 用指定的分割符分割字符串。

```
1 s = pd.Series(['Zhang hua', ' Py thon', 'java', 'go', '11 25',  
2   ', 'javascript'])  
3 print(s.str.split(" "))
```

```
1 0    [Zhang, hua]  
2 1    [, Py, thon]  
3 2          [java]  
4 3          [go]  
5 4    [11, 25, ]  
6 5    [javascript]  
7 dtype: object
```

- `cat(sep="")` 用给定的分隔符连接字符串元素。

```
1 s = pd.Series(['C', 'Python', 'java', 'go', np.nan,  
2   '1125', 'javascript'])  
3 #会自动忽略NaN  
4 s.str.cat(sep="_")
```

```
1 'C_Python_java_go_1125_javascript'
```

```
1 # 先转化列表 ,将列表转化为Series  
2 d_str = s.str.cat(sep="_")  
3 pd.Series(d_str.split('_'))
```

```
1 0      C  
2 1    Python  
3 2     java  
4 3       go  
5 4     1125  
6 5  javascript  
7 dtype: object
```

- `contains(pattern)` 如果子字符串包含在元素中, 则为每个元素返回一个布尔值 True, 否则为 False。

```

1 s = pd.Series(['C ', 'Py thon', 'java', 'go', '1125 ', 'javascript'])
2 s.str.contains(" ")
3

```

```

1 0      True
2 1      True
3 2     False
4 3     False
5 4      True
6 5     False
7 dtype: bool

```

```

1 s[s.str.contains(" ")]

```

```

1 0      C
2 1  Py thon
3 4    1125
4 dtype: object

```

- `replace(a,b)` 将值 a 替换为值 b。

```

1 s = pd.Series(['C ', ' Python', 'java', 'go', '1125 ', 'javascript'])
2 s.str.replace("java", "python")

```

```

1 0      C
2 1    Python
3 2    python
4 3      go
5 4    1125
6 5  pythonscript
7 dtype: object

```

- `count(pattern)` 返回每个字符串元素出现的次数。

```

1 s = pd.Series(['C ', 'Python Python', 'Python', 'go', '1125 ', 'javascript'])
2 s.str.count("Python")

```

```
1 0 0
2 1 2
3 2 1
4 3 0
5 4 0
6 5 0
7 dtype: int64
```

- `startswith(pattern)` 如果 Series 中的元素以指定的字符串开头，则返回 True。
- `endswith(pattern)` 如果 Series 中的元素以指定的字符串结尾，则返回 True。

```
1 s = pd.Series(['C ', ' Python', 'java', 'go', '1125 ', 'javascript'])
2 #若以指定的"j"开头则返回True
3 print(s.str.startswith("j"))
4 s[s.str.startswith("j")]
```

```
1 0 False
2 1 False
3 2 True
4 3 False
5 4 False
6 5 True
7 dtype: bool
```

```
1 2 java
2 5 javascript
3 dtype: object
```

```
1 #若以指定的"a"开头则返回True
2 print(s.str.endswith("a"))
```

```

1 0 False
2 1 False
3 2 True
4 3 False
5 4 False
6 5 False
7 dtype: bool

```

- `repeat(value)` 以指定的次数重复每个元素。

```

1 s = pd.Series(['C ', ' Python', 'java', 'go', '1125 ', 'javascript'])
2 print(s.str.repeat(3))

```

```

1 0 C C C
2 1 Python Python Python
3 2 javajavajava
4 3 gogogo
5 4 1125 1125 1125
6 5 javascriptjavascriptjavascript
7 dtype: object

```

- `find(pattern)` 返回字符串第一次出现的索引位置。

```

1 s = pd.Series(['C ', ' Python', 'java', 'go', '1125 ', 'j1111javascript'])
2 print(s.str.find("a"))
3 # 如果返回 -1 表示该字符串中没有出现指定的字符。

```

```

1 0 -1
2 1 -1
3 2 1
4 3 -1
5 4 -1
6 5 5
7 dtype: int64

```

- `findall(pattern)` 以列表的形式返回出现的字符串。

```

1 s = pd.Series(['C ', ' Python', 'java', 'go', '1125 ', 'javascript'])
2 print(s.str.findall("a"))

```

```

1 0 []
2 1 []
3 2 [a, a]
4 3 []
5 4 []
6 5 [a, a]
7 dtype: object

```

