

第2天两行代码下载网站数据

一、创建

```
1 s1 = 'lenovo'
2 s2 = "QF"
3 s3 = """hello lenovo"""
4 s4 = '''hello 千锋云计算'''
5 s5 = """hello
6 shark
7 """
8 s6 = '''hello
9 world'''
```

二、简单使用

1. \ 转义符

```
1 testimony = 'This shirt doesn\'t fit me'
```

```
1 words = 'hello \nshark'
```

2. + 拼接

```
1 In [1]: file_name= "成功的21个信念"
2
3 In [2]: suffix = '.txt'
4
5 In [3]: file_name = file_name + suffix
6
7 In [4]: file_name
8 Out[4]: '成功的21个信念.txt'
9
```

拼接只能是 字符串和字符串进行操作，不可以用 字符串和 一个非字符串类型的对象相加

```
1 In [5]: '西瓜甜' + 1    ## 这会报错的
```

3. * 复制

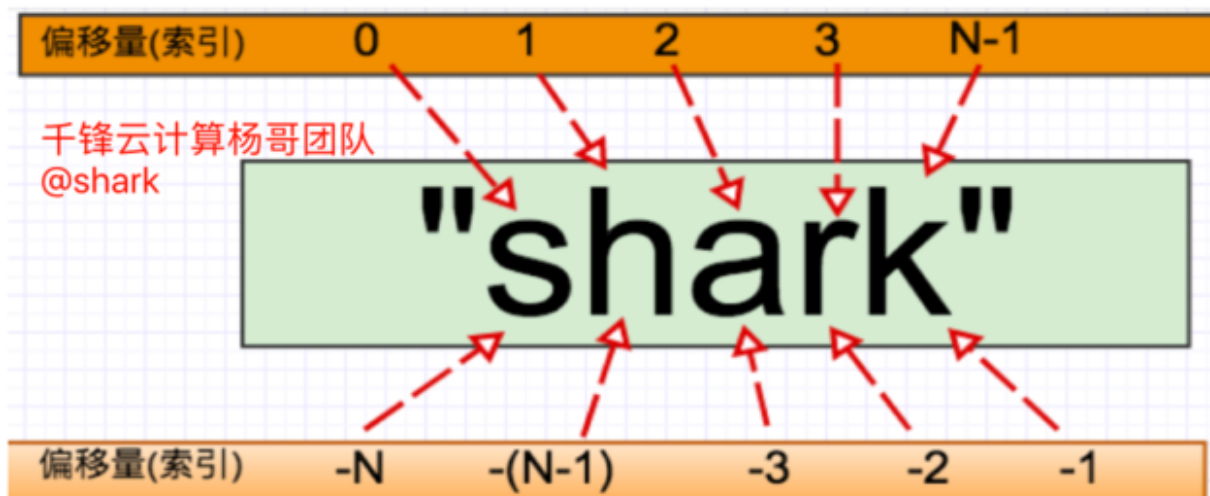
```
1 In [6]: "-" * 10
2 Out[6]: '-----'
3
4 In [7]: print('*' * 10)
5 *****
6
```

三、取值和切片

1. 字符串 是 Python 中的一个 序列类型 的数据结构

- 存放的数据，在其内是有序的。

```
1 | s1 = "shark"
```



序列类型的特点

- 序列里的每个数据被称为序列的一个元素
- 元素在序列里都是有个自己的位置的，这个位置被称为索引或者叫偏移量，也有叫下标的，从 0 开始，从左到右依次递增
- 序列中的每一个元素可以通过这个元素的索引来获取到
- 获取到序列类型数据中的多个元素需要用切片的操作来获取

2. 通过索引取值，获取单个元素

```
1 In [10]: s1 = "shark"
2
3 In [11]: s1[0]
4 Out[11]: 's'
5
6 In [12]: s1[-1]
7 Out[12]: 'k'
8
9 In [13]: s1[3]
10 Out[13]: 'r'
```

3. 切片，获取多个元素

• `[start:end:step]` 分片

千锋云计算杨哥团队
@shark

- * **start** 永远是起始索引号
- * **end** 永远是终止索引号
- * **step** 是可选的步长

分片操作只包含位于起始索引号位置的元素，
不包含位于终止索引号位置的元素；

同时，起始和终止的意义都是针对于从左向右的顺序来定义的

3.1 一般操作

```
1 # 使用切片获取多个元素
2 In [15]: s1[0:3]
3 Out[15]: 'sha'
4
5 # 起始和结尾的索引号可以省略
6 In [16]: s1[:3]
7 Out[16]: 'sha'
8
9 In [17]: s1[1:]
```

```
10 Out[17]: 'hark'
11
12 # 索引可以使用 负数
13 In [18]: s1[2:-1]
14 Out[18]: 'ar'
15
16 In [19]:
```

下面这样的操作，臣妾做不到

```
1 >>> s1[-1:-3]
2 ''
3 >>>
```

因为默认的切片是从左向右开始操作，索引号 `-1` 的右边没有任何索引号了 `-3` 在 `-1` 的左边

3.2 使用步长

- 步长就是每数几个取一个的意思
- 步长是正数时，是从左向右开始操作
- 步长是负数时，是从右向左操作

```
1 In [19]: s2 = 'abcdefg'
2
3 In [20]: s2[::2]
4 Out[20]: 'aceg'
5
6 In [21]: s2[::-1]
7 Out[21]: 'gfedcba'
8
9 In [22]: s2[::-2]
10 Out[22]: 'geca'
11
```

四、字符串方法

1. 统计序列数据的长度

就是获取一个序列数据的元素个数，这个适用于所有的序列类型的数据，比如 字符串、列表、元组。

```
1 # 获取字符串的长度，包含空格和换行符
2 In [25]: s3 = "a \n\t"
3
4 In [26]: len(s3)
5 Out[26]: 4
```

`\n` 是一个换行符 `\t` 是一个 Tab 键

2. in 成员判断

```
1 In [39]: line = 'Size: 8192 MB'
2
3 In [40]: if 'Size' in line:
4     ...:     print(line)
5     ...:
6     Size: 8192 MB
```

注意：空的字符串总是被视为任何其他字符串的子串，因此

`"" in "abc"` 将返回 `True`。

3. strip() 去除字符串两端的空白字符（空格、`\t`、`\n`）

```
1
2 Out[41]: line = '      Size: 8192 MB'
3
4 In [42]: line.strip()
5 Out[42]: 'Size: 8192 MB'
```

4. split() 分割

默认使用空白字符作为分隔符（空格、`\t`、`\n`）和 shell 中的 `awk` 一样道理

```
1 In [47]: line
2 Out[47]: '      Size: 8192 MB'
3
4 In [48]: line.split()
5 Out[48]: ['Size:', '8192', 'MB']
6
7 In [49]: s
8 Out[49]: '\tab\n'
9
10 In [50]: s.split()
11 Out[50]: ['ab']
12
```

可以指定分隔符

```
1 In [51]: line.split(':')
2 Out[51]: ['      Size', ' 8192 MB']
3
4 In [52]: line.split(' : ')
5 Out[52]: ['      Size', '8192 MB']
6
```

5. strip() 移除字符串两端的空白字符**

```
1 In [71]: line = '      Size: 8192 MB'
2
3 In [72]: line.strip()
4 Out[72]: 'Size: 8192 MB'
5
```

`strip()` 返回的是字符串，所以可以连续操作

```
1 In [73]: line.strip().split(': ')
2 Out[73]: ['Size', '8192 MB']
3
4 In [74]: line
5 Out[74]: '      Size: 8192 MB'
6
7 In [75]: k, v = line.strip().split(': ')
8
9 In [76]: k
10 Out[76]: 'Size'
11
12 In [77]: v
13 Out[77]: '8192 MB'
14
```

6. replace() 替换


```

1 In [65]: line = '    <strong>10、命运在自己手里，而不是在
    别人的嘴里</strong></p>'
2
3 In [66]: line.strip()    ## 先去除两端空白字符
4 Out[66]: '<strong>10、命运在自己手里，而不是在别人的嘴里
    </strong></p>'
5
6 In [67]: line.strip().replace('strong>', '')    ##
    将字符串 strong> 替换为空
7 Out[67]: '<10、命运在自己手里，而不是在别人的嘴里</</p>'
8
9 In [68]: line.strip().replace('strong>', '')[1:-6]
10 Out[68]: '10、命运在自己手里，而不是在别人的嘴里'
11

```

7. startswith() 判断字符串以什么为开头

```

1 In [85]: line = 'Locator: DIMM_A2'
2
3 In [86]: line.startswith("Locator:")
4 Out[86]: True

```

8. endswith() 判断字符串以什么为结尾

```

1 In [87]: line = 'Size: 8192 MB'
2
3 In [88]: line.endswith('MB')
4 Out[88]: True
5

```

五、字符串的判断（扩展自修）

```
1 In [1]: s = '123'
2
3 In [2]: s.isdigit()    # 判断是否是纯数字
4 Out[2]: True
5
6 In [3]: s1 = '123adb'
7
8 In [4]: s1.isalnum()   # 判断是否是数字和字母
9 Out[4]: True
10
11 In [5]: s2 = 'adb'
12
13 In [6]: s2.isalpha()  # 判断是否是纯字母
14 Out[6]: True
15
16 In [7]: s2.encode()   # 转换为二进制 bytes 类型
17 Out[7]: b'adb'
18
19 In [8]: s4 = "云计算"
20
21 In [9]: s4.encode()   # 转换为二进制 bytes 类型, 默认编
    码 utf-8
22 Out[9]: b'\xe4\xba\x91\xe8\xae\xa1\xe7\xae\x97'
23
24 In [16]: b = s4.encode()
25
26 In [17]: b.hex()      # bytes 转换成 16 进制
27 Out[17]: 'e4ba91e8aeae7ae97'
28
29 In [18]: b.decode()   # bytes 转换成 str, 默认编码
    utf-8
30 Out[18]: '云计算'
```

