

정규표현식

- 문자열을 처리할때 특정 패턴으로 문자열을 처리하는 방법
- 정규표현식 함수
 - match : 가장 앞에서 부터 일치하는 패턴을 찾음
 - search : 문자열에서 일치하는 가장 첫번째 패턴을 찾음
 - findall : 문자열에서 일치하는 모든 패턴을 찾음
 - split : 문자열을 특정 패턴으로 나눕니다.
 - sub : 특정패턴에 맞는 문자열을 대체하기.
- pattern
- 예제
 - 중고나라
 - 0일공이삼790o삼영, 0일공-이삼79-0o삼영
 - 정규표현식을 이용 => 01023790030

1. 정규표현식 함수

```
In [1]: 1 import re
```

```
In [2]: 1 s = "fast campus datascience fighting. datascience fighting. fast campus fighting."
```

```
In [3]: 1 # match
2 result1 = re.match("fast", s)
3 result2 = re.match("campus", s)
4 result1, result2
```

```
Out[3]: (<re.Match object; span=(0, 4), match='fast'>, None)
```

```
In [4]: 1 # search
        2 result1 = re.search("fast", s)
        3 result2 = re.search("campus", s)
        4 result1, result2
```

```
Out[4]: (<re.Match object; span=(0, 4), match='fast'>,
        <re.Match object; span=(5, 11), match='campus'>)
```

```
In [5]: 1 # findall
        2 result1 = re.findall("fast", s)
        3 result2 = re.findall("campus", s)
        4 result1, result2
```

```
Out[5]: (['fast', 'fast'], ['campus', 'campus'])
```

```
In [7]: 1 # split
        2 re.split("campus", s)
```

```
Out[7]: ['fast ', ' datascience fighting. datascience fighting. fast ', ' fighting.']
```

```
In [8]: 1 # sub
        2 re.sub("fast", "slow", s)
```

```
Out[8]: 'slow campus datascience fighting. datascience fighting. slow campus fighting.'
```

2. 패턴

- 문자
 - \d : 숫자
 - \D : 비숫자
 - \w : 숫자, 문자, _
 - \W : 숫자, 문자, _ 제외
 - \s : 공백문자
 - \S : 비공백문자
- 지정자

```
In [9]: 1 import string
        2 pt = string.printable
        3 len(pt), pt
```

```
Out[9]: (100,
        '0123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ!#$%&\'()*+,-./:;<=>?@[\\]^_`{|}~ \t\n\r\x0b\x0c')
```

```
In [11]: 1 result = re.findall("\d", pt)
        2 "".join(result)
```

```
Out[11]: '0123456789'
```

```
In [12]: 1 result = re.findall("\D", pt)
        2 "".join(result)
```

```
Out[12]: 'abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ!#$%&\'()*+,-./:;<=>?@[\\]^_`{|}~ \t\n\r\x0b\x0c'
```

```
In [13]: 1 result = re.findall("\w", pt)
        2 "".join(result)
```

```
Out[13]: '0123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ_'
```

```
In [14]: 1 result = re.findall("\W", pt)
        2 "".join(result)
```

```
Out[14]: '!"#$%&\'()*+,-./:;<=>?@[\\]^_`{|}~ \t\n\r\x0b\x0c'
```

```
In [15]: 1 result = re.findall("\s", pt)
        2 "".join(result)
```

```
Out[15]: ' \t\n\r\x0b\x0c'
```

```
In [16]: 1 result = re.findall("\S", pt)
        2 "".join(result)
```

```
Out[16]: '0123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ!#$%&\'()*+,-./:;<=>?@[\\]^_`{|}~'
```

3. 지정자

- [] : 문자의 묶음
- - : 범위
- . : 하나의 문자
- ? : 0회 또는 1회
- * : 0회 이상
- + : 1회 이상
- {m} : m 회
- {m,n} : m ~ n 회
- () : 그룹핑

```
In [17]: 1 pt
```

```
Out[17]: '0123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ!#$%&\'()*+,-./:;<=>?@[\\]^_`{|}~ \t\n\r\x0b\x0c'
```

```
In [19]: 1 re.findall("[abcd1234]", pt)
```

```
Out[19]: ['1', '2', '3', '4', 'a', 'b', 'c', 'd']
```

```
In [20]: 1 # - 범위
          2 re.findall("[a-d1-4]", pt)
```

```
Out[20]: ['1', '2', '3', '4', 'a', 'b', 'c', 'd']
```

```
In [24]: 1 re.findall("[a-zA-Z]", pt)[:5]
```

```
Out[24]: ['a', 'b', 'c', 'd', 'e']
```

```
In [25]: 1 # . 문자하나
          2 ls = ["123aab123", "a0b", "abc"]
          3 for s in ls:
          4     result = re.findall("a.b", s)
          5     print(s, result)
```

```
123aab123 ['aab']
```

```
a0b ['a0b']
```

```
abc []
```

```
In [26]: 1 # ? : ?앞에 있는 패턴을 0회 또는 1회 반복
2 ls = ["aab", "a0b", "abc", "accb"]
3 for s in ls:
4     # a + 모든문자 0개 또는 1개 + b
5     result = re.findall("a.?b", s)
6     print(s, result)
```

```
aab ['aab']
a0b ['a0b']
abc ['ab']
accb []
```

```
In [28]: 1 # * : 0회 이상 반복
2 ls = ["ac", "abc", "abbbbc", "a3bec"]
3 for s in ls:
4     # a + b가 0개 이상 + c
5     result = re.findall("ab*c", s)
6     print(s, result)
```

```
ac ['ac']
abc ['abc']
abbbbc ['abbbbc']
a3bec []
```

```
In [29]: 1 # + : 1회 이상 반복
2 ls = ["ac", "abc", "abbbbc", "a3bec"]
3 for s in ls:
4     # a + b가 1개 이상 + c
5     result = re.findall("ab+c", s)
6     print(s, result)
```

```
ac []
abc ['abc']
abbbbc ['abbbbc']
a3bec []
```

```
In [32]: 1 # {m}, {m, n}
2 ls = ["ac", "abc", "abbbbc", "abbbbbbbbbc"]
3 for s in ls:
4     result = re.findall("ab{4}c", s)
5     print(s, result)
```

```
ac []
abc []
abbbbc ['abbbbc']
abbbbbbbbbc []
```

```
In [33]: 1 # () 그룹핑
2 ls = ["aaa5.djfi", "abdddc5", "1abbbc", "a3.bec"]
3 for s in ls:
4     result = re.findall("([0-9]+)[.]( [a-z]{2})", s)
5     print(s, result)
```

```
aaa5.djfi [('5', 'dj')]
abdddc5 []
1abbbc []
a3.bec [('3', 'be')]
```

```
In [34]: 1 # 이메일 주소 찾기
2 s = "저의 이메일 주소는 pdj1224@gmail.com 이고 pythontest@daum.net도 사용합니다"
```

```
In [36]: 1 p = "[\w]+@[0-9a-z]+\.[0-9a-z]+"
2 re.findall(p, s)
```

```
Out[36]: ['pdj1224@gmail.com', 'pythontest@daum.net']
```

```
In [42]: 1 # 주민등록번호 : 661231-1987123 -> 661231-*****
2 s = "저의 주민등록 번호는 661231-1987123 입니다."
3 p = "([0-9]{6})[-]?([0-9]{7})"
4 print(re.findall(p, s))
5 re.sub(p, "\g<1>-*****", s)
```

```
[('661231', '1987123')]
```

```
Out[42]: '저의 주민등록 번호는 661231-***** 입니다.'
```

