

The Kyoto College of
Graduate Studies
for Informatics

kcg.edu

コンピュータプログラミング概論

秋期第10回 eラーニング資料

安 平勲

h_an@kcg.ac.jp

一般的な文字列操作

- ◆ 文字列を結合・分割できる
- ◆ 英字の大文字と小文字をスワップできる
- ◆ 部分文字列を探索・入れ替えできる
- ◆ 文字列のフォーマットもできる
- ◆ . . .

([公式ドキュメントをクリック](#))

正規表現とは？

- **正規表現** (Regular Expression) とは、「検索」や「置換」で指定する**文字列をパターン表現する方法**で、プログラミング言語やテキストエディタなどで利用できる
- 正規表現には、「パターンを表現するための記号 = **メタ文字**」が多数用意されており、それらを組み合わせることで、「aから始まる英単語」「3桁の数字」「行頭の2文字」といった柔軟な文字列を指定することができる

※例えば、任意のEXCELファイル名を示す『*.xlsx』の*がメタ文字

ただし、正規表現にはプログラム言語やテキストエディタなどにより「方言」があり、メタ文字や記述方法が異なるので注意が必要

Pythonでのメタ文字

- 下記は1例。その他のメタ文字や詳しい使い方は公式ドキュメントを参考 <https://docs.python.jp/3/library/re.html>

| 特殊文字 | 説明 | 例 | matchする例 | matchしない例 |
|----------|-------------------|---------|----------|-----------|
| . | 改行以外の任意の一文字 | a.c | abc | ac abbc |
| ^ | 文字列の先頭 | ^ab | abc | zab |
| \$ | 文字列の末尾 | ab\$ | zab | abc |
| * | 直前の文字の0回以上の繰り返し | ab* | a ab abb | aa ac |
| + | 直前の文字の1回以上の繰り返し | ab+ | ab abb | a |
| ? | 直前の文字の0回または1回 | ab? | a ab | abb |
| {m} | 直前の文字のm回の繰り返し | a{2} | aa | a aaa |
| {m,n} | 直前の文字のm回からn回の繰り返し | a{1,2} | a aa | aaa |
| \ | 特殊文字のエスケープ | \. | . | a |
| [文字の集合] | 集合の中の1文字 | [a-c] | a b c | d |
| [^文字の集合] | 集合に含まれない1文字 | [^a-c] | d | a b c |
| | いずれか | a b | a b | c |
| () | グループ化 | (ab cd) | ab cd | abc ac |

メタ文字を使った検索パターン

| 探したい文字列 | メタ文字によるパターン | 当てはまる例 | 当てはまらない例 |
|-------------|-------------|-------------------------------|-----------------|
| 「京都」で始まる文字列 | ^京都.* | 京都府 京都タワー 京都好き | 滋賀県 今日も京都は晴れ |
| 「ね。」で終わる文字列 | .*ね。\$ | いいね。 ねねね。 だよね。 | いいわ。 いいのね！ |
| 「〇〇cm」 | [0-9]+cm | 30cm 1,000,000cm 5 0 cm | 30m 50メートル |

正規表現をいつ使うか？

■ 文字列の書式チェック

- － 例えば、メールアドレスでは使える文字に制限がある。アットマーク(@)が必須であるといった条件に加え、その範囲内であれば任意の文字列を利用することができる
- － このような場合「メールアドレスとして正しいか」という書式のチェックに正規表現を利用する

■ SNSでの記事検索

- － 例えば、ウェブ・スクレイピング技術を使い、自社製品の評価をしているブロガーの記事を特定パターン（WEBのタイトルから抽出など）で収集・分析する

■ 自社サイトへのNGワードの投稿検索

Pythonでの正規表現操作

- **標準ライブラリのモジュール `re`** の提供メソッドを使い、正規表現された **パターン** と **マッチ** する **文字列** を検索

| re のメソッド | 目的 |
|----------------------------------|---|
| <code>match(パターン, 文字列)</code> | 文字列の先頭で正規表現とマッチするか判定する |
| <code>search(パターン, 文字列)</code> | 文字列を操作して、正規表現がどこにマッチするか調べる |
| <code>findall(パターン, 文字列)</code> | 正規表現にマッチする部分文字列を全て探しだしリストとして返す |
| <code>finditer(パターン, 文字列)</code> | 正規表現にマッチする部分文字列を全て探しだし iterator として返す |

`match()` と `search()` はマッチするものが見つからなければ `None` を返す。
成功すればそれらは **Match オブジェクト・インスタンス** を返します。このオブジェクトにはマッチした情報が含まれる: マッチの開始と終了位置, マッチした部分文字列, など
(次頁参照)

Pythonでの正規表現操作

- **Match オブジェクト・インスタンス** はいくつかのメソッドと属性（インスタンス変数）を持っており、重要なのは以下

| メソッド/属性 | 目的 |
|---------|-------------------------------|
| group() | 正規表現にマッチした文字列を返す |
| start() | マッチの開始位置を返す |
| end() | マッチの終了位置を返す |
| span() | マッチの位置 (start, end) を含むタプルを返す |

Pythonでの正規表現操作：matchメソッド

```
import re
pattern = "https?://[^\s/]+"
URL = "https://www.amazon.co.jp/dd/A09TCPZX"

result = re.match(pattern, URL)
print(result)
print(result.group())
```

```
<re.Match object; span=(0, 25), match='https://www.amazon.co.jp/'>
https://www.amazon.co.jp/
```

‘http://’または‘https://’で始まり， ‘/’までを取り出すプログラム

➤ **正規表現**のパターン；`https?://[^\s/]+`

`https?`；httpまたはhttps。 `?`は前の文字(s)が0または1つを意味する

`[^\s/]+`；スラッシュ（/）以外`^`の文字`[]`の繰返し`+`で/が出るまで

➤ **match**メソッドは文字列の先頭からパターンを検索

- ・ match オブジェクトを変数resultで参照
- ・ groupメソッドは一致したmatch オブジェクトの文字列を返す

Pythonでの正規表現操作：searchメソッド

```
import re
pattern = "https?://[^\s/]+/"
URL = """アマゾンのサイト：https://www.amazon.co.jp/dd/A09TCPZX
      KCGIのサイト：https://www.kcg.edu/curriculum"""

result = re.search(pattern, URL)
print(result)
print(result.group())
```

```
<re.Match object; span=(9, 34), match='https://www.amazon.co.jp/'>
https://www.amazon.co.jp/
```

‘http://’または‘https://’で始まり， ‘/’までを取り出すプログラム 2

- 正規表現のパターン；`https?://[^\s/]+/` は同じ
- 検索対象の文字列が違う（urlは文中）
- **search**メソッドはパターンが文字列のどこにマッチするか検索

Pythonでの正規表現操作：findallメソッド

```
import re
pattern = "https?://[^\s/]+"
URL = """アマゾンのサイト：https://www.amazon.co.jp/dd/A09TCPZX
      KCGIのサイト：https://www.kcg.edu/curriculum"""

result = re.findall(pattern, URL)
print(result)
print(result[1], len(result))

['https://www.amazon.co.jp/', 'https://www.kcg.edu/']
https://www.kcg.edu/ 2
```

‘http://’または‘https://’で始まり，‘/’までを取り出すプログラム 3

- 正規表現のパターン；`https?://[^\s/]+`のパターンも検索対象の文字列も同じ
- **findall**メソッドはマッチする部分文字列を全て探しだしリストとして返す

Pythonでの正規表現操作

全角文字の検索

```
import re
```

```
text = """京都情報大学院大学に通学しています。  
京都市に住んでいません、  
大阪市に住んでいます。"""
```

```
result = re.match('京都', text) # matchメソッド  
print(result)  
print(result.group())  
print(result.start())  
print(result.end())  
print(result.span())
```

```
<re.Match object; span=(0, 2), match='京都'>  
京都  
0  
2  
(0, 2)
```

- 全角(2バイトコード)文字の検索
- matchメソッドで先頭から検索

Pythonでの正規表現操作

```
import re
text = """京都情報大学院大学に通学しています。
京都市に住んでいません、
大阪市に住んでいます。"""
```

```
result = re.match('大学', text)
if result: # 見つからればTrueが返る
    print(result.group())
else: # 見つからないとNoneが返る
    print(result)
```

None

✓ '大学'では先頭からは見つからない

```
import re
text = """京都情報大学院大学に通学しています。
京都市に住んでいません、
大阪市に住んでいます。"""
```

```
# '大学'を後ろに含む文字列パターンで検索
result = re.match('.*大学', text)
print(result)
print(result.group())
```

<re.Match object; span=(0, 9), match='京都情報大
京都情報大学院大学

➤ .* (ワイルドカード)を使い, matchメソッドで検索

Pythonでの正規表現操作

```
import re
text = """京都情報大学院大学に通学しています。
京都市に住んでいません、
大阪市に住んでいます。"""

# searchメソッド
result = re.search('大学', text)
print(result)

<re.Match object; span=(4, 6), match='大学'>
```

- searchメソッドなら文中の‘大学’を見つけられる

```
import re
text = """京都情報大学院大学に通学しています。
京都市に住んでいません、
大阪市に住んでいます。"""

# matchメソッドと等価
result = re.search('.*大学', text)
print(result.group())
print(result.start())
print(result.end())
print(result)

京都情報大学院大学
0
9
<re.Match object; span=(0, 9), match='京都情報大学院大
```

- ‘.* 大学’(ワイルドカード)を使った検索ではmatchメソッドとsearchメソッドは等価

Pythonでの正規表現操作

```
import re
text = '''京都情報大学院大学に通学しています。
        京都市に住んでいません、
        大阪市に住んでいます。'''
result = re.findall('京都.', text) # findallメソッド
print(result)                     # listを返す
result = re.findall('.*市', text) # '市'を後ろに持つ文字列
print(result)
```

```
['京都情', '京都市']
['¥u3000¥u3000¥u3000¥u3000¥u3000京都市', '¥u3000¥u3000¥u3000¥u3000¥u3000大阪市']
```

```
result = re.split('。|、', text) # splitメソッド 正規表現!
print(result)                   # パターンで分割。リストを返す
```

```
['京都情報大学院大学に通学しています', '¥n¥u3000¥u3000¥u3000¥u3000¥u3000京都市に住んでいません', '¥n¥u3000¥u3000¥u3000¥u3000¥u3000大阪市に住んでいます', '']
```

```
result = re.sub('市', '府', text) # subメソッド
print(result)                   # パターンで置き換え。文字列を返す
```

```
京都情報大学院大学に通学しています。
京都府に住んでいません、
大阪府に住んでいます。
```

- findallメソッドによる検索。結果すべてをリストで返す
- ✓ '全角空白'も文字列 ※改行(¥n)が区切り

- reモジュールのsplitメソッドの使用例
- ✓ '|' (メタ文字のor) でゆらぎを検索

- ✓ reモジュールのsubメソッドの使用例

Pythonでのその他の特殊文字

| 特殊文字 | 説明 | 同義のパターン |
|------|----------------|-------------------------|
| ¥d | 任意の数字 | ASCIIの場合, [0-9] |
| ¥D | ¥d以外 | |
| ¥s | 任意の空白文字 | ASCIIの場合, [¥t¥n¥r¥f¥v] |
| ¥S | ¥s以外 | |
| ¥w | 任意のUnicode単語文字 | ASCIIの場合, [a-zA-Z0-9_] |
| ¥W | ¥w以外 | |
| ¥A | 文字列の先頭 | ^ |
| ¥Z | 文字列の末尾 | \$ |
| ¥b | 単語の境界(¥wと¥Wの間) | |
| ¥B | 単語の境界以外の文字間 | |



いくつかの特殊文字がエスケープシーケンスとバッティング

- ・ ' ¥b 'はエスケープシーケンスでは'バックスペース'

Pythonでの正規表現操作

```
import re
strings = """家の郵便番号は612-1010です。
             携帯電話は080(1000)2000です"""
#郵便番号をだけを抽出したい
pattern1 = r"¥d{3}-¥d{4}"

result = re.findall(pattern1, strings)
print(result)

['612-1010']
```

- 特殊文字¥dは[0-9]と等価
- 特殊文字を使うパターンには先頭に' r 'を付加
⇒pythonは' r 'で¥を特殊文字と認識

Pythonでの正規表現操作

```
import re
content = 'hellow Python, 123, end'

# 先頭が数字の文字列を探索
pattern1 = '¥b(¥d+)' # ¥bはエスケープ (バックスペース)
pattern2 = r'¥b(¥d+)' # ¥bは特殊文字 (単語以外の境界)
result1 = re.findall(pattern1, content)
result2 = re.findall(pattern2, content)

print(result1)
print(result2)
```

[
['123']

- 先頭に'r'を付加しないと、'¥b'をバックスペースと認識
- 'r'を付加すると、¥b'を単語以外の境界と認識

Pythonでの正規表現操作

```
import re
content = 'hellow Python, 123, end'

# ()で探したい任意の数字を指定
pattern = r'.*?(¥d+).*'
result = re.match(pattern, content)
if result:
    # group()で全文字
    print(result.group(0))
    # group(1)で探したい数字
    print(result.group(1))

# メタ文字'?'を省くと
pattern1 = r'.*(¥d+).*'
result1 = re.match(pattern1, content)
print(result1.group(0))
print(result1.group(1))
```

```
hellow Python, 123, end
123
hellow Python, 123, end
3
```

特殊文字を使うパターンの検索

- 任意の数字列を探す2つの例
(数字列の前後にワイルドカード)
- ? (直前の文字の0回または1回)の有無で結果が異なる
- ✓ groupメソッドの引数に注目