

Regular Expression Basics: Takeaways

by Dataquest Labs, Inc. - All rights reserved © 2021

Syntax

REGULAR EXPRESSION MODULE

- Importing the regular expression module:

```
import re
```

- Searching a string for a regex pattern:

```
re.search(r"blue", "Rhythm and blues")
```

PANDAS REGEX METHODS

- Return a boolean mask if a regex pattern is found in a series:

```
s.str.contains(pattern)
```

- Extract a regex capture group from a series:

```
s.str.extract(pattern_with_capture_group)
```

ESCAPING CHARACTERS

- Treating special characters as ordinary text using backslashes:

```
\[pdf\]
```

Concepts

- Regular expressions, often referred to as regex, are a set of syntax components used for matching sequences of characters in strings.
- A pattern is described as a regular expression that we've written. We say regular expression has matched if it finds the pattern exists in the string.
- Character classes allow us to match certain classes of characters.
- A set contains two or more characters that can match in a single character's position.
- Quantifiers specify how many of the previous characters the pattern requires.
- Capture groups allow us to specify one or more groups within our match that we can access separately.
- Negative character classes are character classes that match every character except a character class.
- An anchor matches something that isn't a character, as opposed to character classes which match specific characters.
- A word boundary matches the space between a word character and a non-word character, or a word character and the start/end of a string
- Common character classes:

Character Class	Pattern	Explanation
Set	<code>[fud]</code>	Either <code>f</code> , <code>u</code> , or <code>d</code>
Range	<code>[a-e]</code>	Any of the characters <code>a</code> , <code>b</code> , <code>c</code> , <code>d</code> , or <code>e</code>
Range	<code>[0-3]</code>	Any of the characters <code>0</code> , <code>1</code> , <code>2</code> , or <code>3</code>
Range	<code>[A-Z]</code>	Any uppercase letter
Set + Range	<code>[A-Za-z]</code>	Any uppercase or lowercase character
Digit	<code>\d</code>	Any digit character (equivalent to <code>[0-9]</code>)
Word	<code>\w</code>	Any digit, uppercase, or lowercase character (equivalent to <code>[A-Za-z0-9_]</code>)
Whitespace	<code>\s</code>	Any space, tab or linebreak character
Dot	<code>.</code>	Any character except newline

- Common quantifiers:

Quantifier	Pattern	Explanation
Zero or more	<code>a*</code>	The character <code>a</code> zero or more times
One or more	<code>a+</code>	The character <code>a</code> one or more times
Optional	<code>a?</code>	The character <code>a</code> zero or one times
Numeric	<code>a{3}</code>	The character <code>a</code> three times
Numeric	<code>a{3,5}</code>	The character <code>a</code> three, four, or five times
Numeric	<code>a{,3}</code>	The character <code>a</code> one, two, or three times
Numeric	<code>a{8,}</code>	The character <code>a</code> eight or more times

- Common negative character classes:

Character Class	Pattern	Explanation
Negative Set	<code>[^fud]</code>	Any character except <code>f</code> , <code>u</code> , or <code>d</code>
Negative Set	<code>[^1-3Z\s]</code>	Any characters except <code>1</code> , <code>2</code> , <code>3</code> , <code>Z</code> , or whitespace characters
Negative Digit	<code>\D</code>	Any character except digit characters
Negative Word	<code>\W</code>	Any character except word characters
Negative Whitespace	<code>\S</code>	Any character except whitespace characters

- Common anchors:

Anchor	Pattern	Explanation
Beginning	<code>^abc</code>	Matches <code>abc</code> only at the start of a string
End	<code>abc\$</code>	Matches <code>abc</code> only at the end of a string
Word boundary	<code>s\b</code>	Matches <code>s</code> only when it's followed by a word boundary
Word boundary	<code>s\B</code>	Matches <code>s</code> only when it's not followed by a word boundary

Resources

- [re module](#)
- [Building regular expressions](#)

Advanced Regular Expressions: Takeaways



by Dataquest Labs, Inc. - All rights reserved © 2022

Syntax

CAPTURE GROUPS

- Extracting text using a capture group:

```
s.str.extract(pattern_with_capture_group)
```

- Extracting text using multiple capture groups:

```
s.str.extract(pattern_with_multiple_capture_groups)
```

SUBSTITUTION

- Substituting a regex match:

```
s.str.replace(pattern, replacement_text)
```

Concepts

- Capture groups allow us to specify one or more groups within our match that we can access separately.

Pattern	Explanation
---------	-------------

<code>(yes)no</code>	Matches <code>yesno</code> , capturing <code>yes</code> in a single capture group.
----------------------	--

<code>(yes)(no)</code>	Matches <code>yesno</code> , capturing <code>yes</code> and <code>no</code> in two capture groups.
------------------------	--

- Backreferences allow us to repeat a capture group within our regex pattern by referring to them with an integer in the order they are captured.

Pattern	Explanation
---------	-------------

<code>(yes)no\1</code>	Matches <code>yesnoyes</code>
------------------------	-------------------------------

<code>(yes)(no)\2\1</code>	Matches <code>yesnonoyes</code>
----------------------------	---------------------------------

- Lookarounds let us define a positive or negative match before or after our string.

Pattern	Explanation
---------	-------------

<code>zzz(?=abc)</code>	Matches <code>zzz</code> only when it is followed by <code>abc</code>
-------------------------	---

<code>zzz(?!abc)</code>	Matches <code>zzz</code> only when it is not followed by <code>abc</code>
-------------------------	---

<code>(?<=abc)zzz</code>	Matches <code>zzz</code> only when it is preceded by <code>abc</code>
-----------------------------	---

<code>(?<!abc)zzz</code>	Matches <code>zzz</code> only when it is not preceded by <code>abc</code>
-----------------------------	---

Resources

- [re module](#)
- [RegExr Regular Expression Builder](#)

List Comprehensions and Lambda Functions: Takeaways



by Dataquest Labs, Inc. - All rights reserved © 2021

Syntax

WORKING WITH JSON FILES

- Open a JSON data set from a file to Python objects:

```
f = open('filename.json')
json.load(f)
```

- Convert JSON data from a string to Python objects:

```
json.loads(json_string)
```

- Convert JSON data stored in Python objects to string form:

```
json.dumps(json_obj)
```

LIST COMPREHENSIONS

- Converting a for loop to a list comprehension:

- Using a for loop:

```
letters=['a', 'b', 'c', 'd']
caps=[]
for l in letters:
    caps.append(l.upper())
```

- Using a List comprehension:

```
caps = [l.upper() for l in letters]
```

- Common list comprehension patterns:

- Transforming a list

```
ints = [25, 14, 13, 84, 43, 6, 77, 56]
doubled_ints = [i * 2 for i in ints]
```

- Creating test data

```
tenths = [i/10 for i in range(5)]
```

- Reducing a list

```
big_ints = [i for i in ints if i >= 50]
```

LAMBDA FUNCTIONS

- Converting a definition to a lambda function:

- Defining a function:

```
def double(x):  
    return x * 2
```

- Defining a lambda function:

```
run_function(function=lambda x: x * 2)
```

THE TERNARY OPERATOR

- Create a one-line version of an if/else statement:

```
"val_1 is bigger" if val_1 > val_2 else "val_1 is not bigger"
```

Concepts

- JSON is a language independent format for storing structured data.
 - In Python, it can be represented by a series of nested lists, dictionaries, strings, and numeric objects.
- A list comprehension provides a concise way of creating lists using a single line of code, where:
 - You start with an iterable object
 - Optionally Transform the items in the iterable object
 - Optionally reduce the items in the iterable object using an if statement
 - Create a new list
- Lambda functions can be defined in a single line, which lets you define a function at the time you need it.
- The ternary operator can be used to replace an if/else statement with a single line.

Resources

- [Official JSON specification](#)
- [Python Documentation: JSON Module](#)
- [Python Documentation: List Comprehensions](#)
- [Python Documentation: Lambda Functions](#)

Working with Missing Data: Takeaways

by Dataquest Labs, Inc. - All rights reserved © 2021

Syntax

- Replacing matching values with a single value:

```
s.mask(s == var, value_to_replace)
```

- Replacing matching values with corresponding values from a series:

```
sl.mask(s == var, series_to_replace)
```

- A function to create a null matrix

```
def plot_null_matrix(df, figsize=(18,15)):  
    # initiate the figure  
    plt.figure(figsize=figsize)  
    # create a boolean dataframe based on whether values are null  
    df_null = df.isnull()  
    # create a heatmap of the boolean dataframe  
    sns.heatmap(~df_null, cbar=False, yticklabels=False)  
    plt.show()
```

- A function to create a null correlation heatmap

```
def plot_null_correlations(df):  
    # create a correlation matrix only for columns with at least  
    # one missing value  
    cols_with_missing_vals = df.columns[df.isnull().sum() > 0]  
    missing_corr = df[cols_with_missing_vals].isnull().corr()  
    # create a triangular mask to avoid repeated values and make  
    # the plot easier to read  
    missing_corr = missing_corr.iloc[1:, :-1]  
    mask = np.triu(np.ones_like(missing_corr), k=1)  
    # plot a heatmap of the values  
    plt.figure(figsize=(20,12))  
    ax = sns.heatmap(missing_corr, vmin=-1, vmax=1,  
                     cmap='RdBu', mask=mask, annot=True)  
    # round the labels and hide labels for values near zero  
    for text in ax.texts:  
        t = float(text.get_text())  
        if -0.05 < t < 0.01:  
            text.set_text('')  
        else:  
            text.set_text(round(t, 2))  
    plt.show()
```

Concepts

- Imputation is the process of replacing missing values with other values.

- Imputing can be a better option than simply dropping values because you retain more of your original data.
- You might find values for imputation by:
 - Deriving the value from related columns.
 - Using the most common non-null value from a column.
 - Using an placeholder for missing values.
 - Augmenting factual data (e.g. location data) using an external resource.
- Using plots can help identify patterns in missing values which can help with imputation.

Resources

- [pandas documentation](#)