

0_beautifulsoup.py

```
1 # 10_basic.py
2 # 15_make_soup.py
3 # 20_search.py
4 # 25_navigation.py
5 # 30_edit.py
6 # 40_encoding.py
7 # 50_parse_only_part.py
```

10_basic.py

```
1 # https://www.crummy.com/software/BeautifulSoup/bs4/doc/
2 from bs4 import BeautifulSoup
3 soup = BeautifulSoup(html_doc, 'html.parser')
4
5
6 soup.title # <title>The Dormouse's story</title>
7 soup.title.name # u'title'
8 soup.title.string # u'The Dormouse's story'
9 soup.title.parent.name # u'head'
10
11 #various finder
12 css_soup.select("p.strikeout.body") # css finder
13 soup.p # <p class="title"><b>The Dormouse's story</b></p>
14 soup.p['class'] # u'title'
15 soup.a # <a class="sister" href="http://example.com/elsie" id="link1">Elsie</a>
16 soup.find_all('a') # [<a .., ..]
17 soup.find(id="link3") # <a class="sister" href="http://example.com/tillie" id="link3">Tillie</a>
18 for link in soup.find_all('a'):
19     print(link.get('href')) # http://example.com/elsi, # http://example.com/lacie
```

15_make_soup.py

```
1 from bs4 import BeautifulSoup
2
3 soup = BeautifulSoup(open("index.html"))
4 soup = BeautifulSoup("<html>data</html>")
```

18_output.py

```
1 # HTML
2 soup.prettify() #pretty print
3 str(soup) # non-pretty print
4
5 # String
6 soup.get_text() #all text under the element
```

20_search.py

```
1 search.pyhttps://www.crummy.com/software/BeautifulSoup/bs4/doc/
2 #-----
3 # css selector
4 #-----
5 css_soup.select("p.strikeout.body")
6 soup.select("p nth-of-type(3)") # 3rd child
7 soup.select("head > title")
8 soup.select("p > a:nth-of-type(2)")
9 soup.select("p > #link1") # direct child
10 soup.select("#link1 ~ .sister") # sibling
11 soup.select('a[href]') # existence of an attribute
12 soup.select_one(".sister")
13
14 # attribute value
15 soup.select('a[href="http://example.com/elsie"]') # exact attribute
16 soup.select('a[href^="http://example.com/"]') # negative match
17 soup.select('a[href$="tillie"]') # end match
18 soup.select('a[href*=".com/el"]') # middle match
19
20
21 #-----
```

```

22 # basic
23 #-----
24 soup.find_all('b') # match by tag
25 soup.find_all(re.compile("^b")) # match by tag using regex
26 soup.find_all(["a", "b"]) # match by tag in list
27
28 # function (complex condition)
29 def has_class_but_no_id(tag):
30     return tag.has_attr('class') and not tag.has_attr('id')
31 soup.find_all(has_class_but_no_id)
32
33
34 #-----
35 # find_all_api
36 #-----
37 find_all(name, attrs, recursive, string, limit, **kwargs)
38
39 soup.find_all("title") # tag condition
40 soup.find_all("p", "title") # tag and attr
41 # [<p class="title"><b>The Dormouse's story</b></p>]
42 soup.find_all("a")
43
44 # keyword arguments
45 soup.find_all(id="link2")
46 soup.find_all(href=re.compile("elsie"), id='link1')
47 soup.find(string=re.compile("sisters")) # text contain sisters
48
49 # css class (class is reserved keyword)
50 soup.find_all("a", class_="sister")

```

🔗 25_navigation.py

```

1 #-----
2 # going up/down/side
3 #-----
4 # ---- going down ----
5 soup.head# <head><title>The Dormouse's story</title></head>
6 soup.title# <title>The Dormouse's story</title>
7 soup.body.b # <b>The Dormouse's story</b>
8 soup.a # <a class="sister" href="http://example.com/elsie" id="link1">Elsie</a>
9 soup.find_all('a')
10 # [<a class="sister" href="http://example.com/elsie" id="link1">Elsie</a>,
11 #  <a class="sister" href="http://example.com/lacie" id="link2">Lacie</a>,
12 #  <a class="sister" href="http:
13
14 # children = contents
15 head_tag.contents # [<title>The Dormouse's story</title>]
16 head_tag.children # [<title>The Dormouse's story</title>]
17
18 # descendants (all of a tag's children, recursively)
19 for child in head_tag.descendants:
20     print(child)
21
22 # .string is tricky
23 head_tag.contents # [<title>The Dormouse's story</title>]
24 head_tag.string # u'The Dormouse's story' (because head tag has only one child)
25 print(soup.html.string) # None (because html has many children)
26
27 # whitespace removed strings
28 for string in soup.stripped_strings:
29     print(repr(string))
30
31
32 # ---- going up ----
33 title_tag.parent # <head><title>The Dormouse's story</title></head>
34 # going up recursively
35 link.parents # [ p, body, html, [document], None]
36
37
38 # ---- sideways ----
39 # sibling = include text node
40 sibling_soup.b.next_sibling
41 sibling_soup.c.previous_sibling

```

```

42
43 # multiple
44 sibling_soup.b.next_siblings
45 sibling_soup.c.previous_siblings
46
47 # element = not include text node
48 sibling_soup.b.next_element
49 sibling_soup.c.previous_element
50 sibling_soup.b.next_elements
51 sibling_soup.c.previous_elements

```

30_edit.py

```

1  #-----
2  # change existng tag
3  #-----
4  tag.name = "blockquote" # modify tag name
5  tag['class'] = 'verybold' # modify tag attribute
6  del tag['class'] # delete attribute
7  tag.string= 'not too bold' # modify tag contents string
8  tag.append(" but bolder than usual") # append tag contents
9
10 #-----
11 # insert tag
12 #-----
13 new_tag = soup.new_tag("a", href="http://www.example.com")
14 original_tag.append(new_tag) # create child
15 new_tag.string = "Link text." # can edit element after creating child
16
17 soup.b.string.insert_before(tag)
18 soup.b.i.insert_after(soup.new_string(" ever "))
19
20 #-----
21 # delete tag
22 #-----
23 soup.i.clear() # removes the contents
24 i_tag = soup.i.extract() # completely removes a tag from tree and returns the element
25 soup.i.decompose() # completely removes a tag from tree and discard the tag
26
27 #-----
28 # replace/wrap/unwrap tag
29 #-----
30 a_tag.i.replace_with(soup.new_tag("b"))
31 a_tag.i.replace_with(BeautifulSoup("<b>bold element</b>")) # replace inner html
32 soup.p.string.wrap(soup.new_tag("b"))
33 a_tag.i.unwrap()

```

40_encoding.py

```

1  #output
2  soup.prettify("latin-1")
3  tag.encode("utf-8")
4  tag.encode("latin-1")
5  tag.encode("ascii")

```

50_parse_only_part.py

```

1  # The SoupStrainer class allows you to choose which parts of an
2  # incoming document are parsed
3  from bs4 import SoupStrainer
4
5  # conditions
6  only_a_tags = SoupStrainer("a")
7  only_tags_with_id_link2 = SoupStrainer(id="link2")
8
9  def is_short_string(string):
10     return len(string) < 10
11  only_short_strings = SoupStrainer(string=is_short_string)
12
13  # execute parse
14  BeautifulSoup(html_doc, "html.parser", parse_only=only_a_tags)
15  BeautifulSoup(html_doc, "html.parser", parse_only=only_tags_with_id_link2)
16  BeautifulSoup(html_doc, "html.parser", parse_only=only_short_strings)

```



heemayl commented on Apr 19, 2018

`soup.select('a[href^="http://example.com/"]')` # negative match is *incorrect*; `^` matches the start of the attribute value string (like Regex). So, it should read:

```
soup.select('a[href="http://example.com/"]') # start match
```