


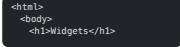


Parse a Document

Turn an HTML string into a navigable tree; pick a parser.

PURPOSE	COMMAND	RESULT
Parse with the built-in parser.	<code>soup = BeautifulSoup(html, "html.parser")</code>	 tree of nodes
Parse fast and lenient (lxml).	<code>BeautifulSoup(html, "lxml")</code>	 pip install lxml
Parse browser-style (html5lib).	<code>BeautifulSoup(html, "html5lib")</code>	 auto-wraps
Parse a file object.	<code>BeautifulSoup(open("page.html"), "html.parser")</code>	
Pretty-print the parsed tree.	<code>print(soup.prettify())</code>	
Parse only XML.	<code>BeautifulSoup(xml, "xml")</code>	 self-closing

name the parser explicitly: html.parser is built in, lxml is fast, html5lib parses exactly like a browser

Read Tags & Attributes

A Tag is a dict of attributes plus a node you can descend into.

PURPOSE	COMMAND	RESULT
Reach the first matching tag.	<code>soup.h1</code>	 first match list
Get a tag's name.	<code>tag.name</code>	 "h1"
Get one attribute (dict-style).	<code>tag["href"]</code>	 href="/p/1"
Safe attribute lookup.	<code>tag.get("href")</code>	 None not an error
All attributes as a dict.	<code>tag.attrs</code>	 {class: [...], data-id: "1"}
Test for an attribute.	<code>tag.has_attr("class")</code>	 True

tag["href"] raises if missing, tag.get("href") returns None; class comes back as a list (multiple classes allowed)

Search: find / find_all


The two workhorses: one match or a list, filtered by name, attrs, text.

PURPOSE	COMMAND	RESULT
First match (or None).	<code>soup.find("a")</code>	
All matches as a list.	<code>soup.find_all("li")</code>	
Filter by CSS class.	<code>soup.find_all("a", class_="name")</code>	 class_ underscore
Filter by any attribute.	<code>soup.find_all("li", attrs={"data-id": "2"})</code>	 others greyed
Match attribute by regex.	<code>soup.find_all("a", href=re.compile(r"^/p/"))</code>	
Cap the number of results.	<code>soup.find_all("li", limit=2)</code>	 stops after 2

find returns the first Tag or None; find_all returns a list; filter by name, class=, attrs={...}, regex, or limit=

Search: CSS Selectors

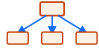


select() / select_one(): the same selectors you'd use in a browser.

PURPOSE	COMMAND	RESULT
All matches for a selector.	<code>soup.select("li.product")</code>	
First selector match.	<code>soup.select_one("a.name")</code>	
Descendant combinator.	<code>soup.select("ul.products .price")</code>	 nested match
Attribute selector.	<code>soup.select('li[data-id="2"]')</code>	
Direct child / grouping.	<code>soup.select("ul > li") ("h1, h2")</code>	 direct child h1, h2
Positional pseudo-class.	<code>soup.select("li:nth-of-type(2)")</code>	 the 2nd li

select() accepts browser-style selectors (Soup Sieve): descendant chains, > children, [attr], grouping, and pseudo-classes

Walk the Tree

Move up to parents, down to children, sideways to siblings.

PURPOSE	COMMAND	RESULT
Direct children (tags + strings).	<code>tag.children</code> <code>tag.contents</code>	 immediate nodes
All descendants (deep).	<code>tag.descendants</code>	 everything below
Go up to the parent.	<code>tag.parent</code> <code>tag.parents</code>	 parents = chain up
Next / previous sibling tag.	<code>tag.find_next_sibling("li")</code>	 skip whitespace
Find an ancestor by filter.	<code>tag.find_parent("li")</code>	 stops at matching li
Find next match anywhere.	<code>tag.find_next("span")</code>	 document order

prefer `find_next_sibling` and `find_parent` over `raw.next_sibling/.parent` when whitespace text nodes get in the way

Pull the Data Out

Get text and attributes, clean and collected.

PURPOSE	COMMAND	RESULT
All visible text, joined.	<code>tag.get_text()</code>	 "Bolts\$3.00"
Text with a separator, trimmed.	<code>tag.get_text(separator=" ", strip=True)</code>	 "Bolt \$3.00" trimmed
One tag's own string.	<code>tag.string</code>	 "Widgets"
Iterate cleaned text pieces.	<code>list(tag.stripped_strings)</code>	 ["Bolt", "\$3.00"]
Collect an attribute across matches.	<code>[a["href"] for a in soup.select("a[href]")]</code>	 List of URLs
Collect text across matches.	<code>[li.get_text(strip=True) for li in soup.select("li")]</code>	 ["Bolt"] clean text list

`get_text()` flattens a subtree; a list comprehension over a search result collects text or an attribute from many nodes

Edit the Tree






Add, change, replace, and remove nodes in place.

PURPOSE	COMMAND	RESULT
Change a tag's text.	<code>tag.string= "Gadgets"</code>	
Set or change an attribute.	<code>tag["href"]= "/changed"</code>	
Build and append a new node.	<code>soup.new_tag("span"); parent.append(t)</code>	
Insert next to a node.	<code>tag.insert_before(n) / insert_after(n)</code>	
Replace one node with another.	<code>tag.replace_with(new)</code>	
Remove a node from the tree.	<code>tag.decompose() tag.extract()</code>	

Beautiful Soup edits in place; `decompose()` destroys a node, `extract()` detaches and returns it for reuse

Serialize Back to a String

Turn the tree (or a piece of it) back into HTML text.

PURPOSE	COMMAND	RESULT
Whole tag back to HTML.	<code>str(soup) str(tag)</code>	
Indented, human-readable HTML.	<code>soup.prettify()</code>	
Inner HTML only (no wrapper).	<code>tag.decode_contents()</code>	
Bytes with an encoding.	<code>soup.encode("utf-8")</code>	
Strip all tags, keep text.	<code>soup.get_text()</code>	
Save to a file.	<code>Path("out.html").write_text(str(soup))</code>	

`str(soup)` renders HTML, `prettify()` indents it, `decode_contents()` gives inner markup, `get_text()` strips to plain text