

# Link Extractors

A link extractor is an object that extracts links from responses.

The `__init__` method of `:class:`~scrapy.linkextractors.xmlhtml.LxmlLinkExtractor`` takes settings that determine which links may be extracted. `:class:`~scrapy.linkextractors.xmlhtml.LxmlLinkExtractor.extract_links`` returns a list of matching `:class:`~scrapy.link.Link`` objects from a `:class:`~scrapy.http.Response`` object.

**System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs] [topics]link-extractors.rst, line 9); [backlink](#)**

Unknown interpreted text role "class".

**System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs] [topics]link-extractors.rst, line 9); [backlink](#)**

Unknown interpreted text role "class".

**System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs] [topics]link-extractors.rst, line 9); [backlink](#)**

Unknown interpreted text role "class".

**System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs] [topics]link-extractors.rst, line 9); [backlink](#)**

Unknown interpreted text role "class".

Link extractors are used in `:class:`~scrapy.spiders.CrawlSpider`` spiders through a set of `:class:`~scrapy.spiders.Rule`` objects.

**System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs] [topics]link-extractors.rst, line 16); [backlink](#)**

Unknown interpreted text role "class".

**System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs] [topics]link-extractors.rst, line 16); [backlink](#)**

Unknown interpreted text role "class".

You can also use link extractors in regular spiders. For example, you can instantiate `:class:`~scrapy.linkextractors.xmlhtml.LxmlLinkExtractor`` into a class variable in your spider, and use it from your spider callbacks:

**System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs] [topics]link-extractors.rst, line 19); [backlink](#)**

Unknown interpreted text role "class".

```
def parse(self, response):
    for link in self.link_extractor.extract_links(response):
        yield Request(link.url, callback=self.parse)
```

## Link extractor reference

**System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs] [topics]link-extractors.rst, line 32)**

Unknown directive type "module".

```
.. module:: scrapy.linkextractors
   :synopsis: Link extractors classes
```

The link extractor class is `:class:`~scrapy.linkextractors.xmlhtml.LxmlLinkExtractor``. For convenience it can also be imported as `scrapy.linkextractors.LinkExtractor`:

**System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs] [topics]link-extractors.rst, line 35); [backlink](#)**

Unknown interpreted text role "class".

```
from scrapy.linkextractors import LinkExtractor
```

## LxmlLinkExtractor

**System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs] [topics]link-extractors.rst, line 44)**

Unknown directive type "module".

```
.. module:: scrapy.linkextractors.lxmlhtml
   :synopsis: lxml's HTMLParser-based link extractors
```

LxmlLinkExtractor is the recommended link extractor with handy filtering options. It is implemented using lxml's robust HTMLParser.

- param allow:** a single regular expression (or list of regular expressions) that the (absolute) urls must match in order to be extracted. If not given (or empty), it will match all links.
- type allow:** str or list
- param deny:** a single regular expression (or list of regular expressions) that the (absolute) urls must match in order to be excluded (i.e. not extracted). It has precedence over the `allow` parameter. If not given (or empty) it won't exclude any links.
- type deny:** str or list
- param allow\_domains:** a single value or a list of string containing domains which will be considered for extracting the links
- type allow\_domains:** str or list
- param deny\_domains:** a single value or a list of strings containing domains which won't be considered for extracting the links
- type deny\_domains:** str or list
- param deny\_extensions:** a single value or list of strings containing extensions that should be ignored when extracting links. If not given, it will default to `:data:'scrapy.linkextractors.IGNORED_EXTENSIONS'`.

**System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs] [topics]link-extractors.rst, line 72); [backlink](#)**

Unknown interpreted text role "data".

**System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs] [topics]link-extractors.rst, line 77)**

Unknown directive type "versionchanged".

```
.. versionchanged:: 2.0
   :data:`~scrapy.linkextractors.IGNORED_EXTENSIONS` now includes
   ``7z``, ``7zip``, ``apk``, ``bz2``, ``cdr``, ``dmg``, ``ico``,
   ``iso``, ``tar``, ``tar.gz``, ``webm``, and ``xz``.
```

- type deny\_extensions:** list
- param restrict\_xpaths:** is an XPath (or list of XPath's) which defines regions inside the response where links should be extracted from. If given, only the text selected by those XPath will be scanned for links. See examples below.
- type restrict\_xpaths:** str or list
- param restrict\_css:** a CSS selector (or list of selectors) which defines regions inside the response where links should be extracted from. Has the same behaviour as `restrict_xpaths`.
- type restrict\_css:** str or list
- param restrict\_text:**

a single regular expression (or list of regular expressions) that the link's text must match in order to be extracted. If not given (or empty), it will match all links. If a list of regular expressions is given, the link will be extracted if it matches at least one.

<b>type restrict_text:</b>	str or list
<b>param tags:</b>	a tag or a list of tags to consider when extracting links. Defaults to ('a', 'area').
<b>type tags:</b>	str or list
<b>param attrs:</b>	an attribute or list of attributes which should be considered when looking for links to extract (only for those tags specified in the <code>tags</code> parameter). Defaults to ('href',)
<b>type attrs:</b>	list
<b>param canonicalize:</b>	canonicalize each extracted url (using <code>w3lib.url.canonicalize_url</code> ). Defaults to <code>False</code> . Note that <code>canonicalize_url</code> is meant for duplicate checking; it can change the URL visible at server side, so the response can be different for requests with canonicalized and raw URLs. If you're using <code>LinkExtractor</code> to follow links it is more robust to keep the default <code>canonicalize=False</code> .
<b>type canonicalize:</b>	bool
<b>param unique:</b>	whether duplicate filtering should be applied to extracted links.
<b>type unique:</b>	bool
<b>param process_value:</b>	a function which receives each value extracted from the tag and attributes scanned and can modify the value and return a new one, or return <code>None</code> to ignore the link altogether. If not given, <code>process_value</code> defaults to <code>lambda x: x</code> .

```
System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs]
[topics]link-extractors.rst, line 127)
Unknown directive type "highlight".

.. highlight:: html
```

For example, to extract links from this code:

```
<a href="javascript:goToPage('../other/page.html'); return false">Link text</a>
```

```
System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs]
[topics]link-extractors.rst, line 133)
Unknown directive type "highlight".

.. highlight:: python
```

You can use the following function in `process_value`:

```
def process_value(value):
    m = re.search("javascript:goToPage\('(.*?)'", value)
    if m:
        return m.group(1)
```

<b>type process_value:</b>	<code>collections.abc.Callable</code>
<b>param strip:</b>	whether to strip whitespaces from extracted attributes. According to HTML5 standard, leading and trailing whitespaces must be stripped from <code>href</code> attributes of <code>&lt;a&gt;</code> , <code>&lt;area&gt;</code> and many other elements, <code>src</code> attribute of <code>&lt;img&gt;</code> , <code>&lt;iframe&gt;</code> elements, etc., so <code>LinkExtractor</code> strips space chars by default. Set <code>strip=False</code> to turn it off (e.g. if you're extracting urls from elements or attributes which allow leading/trailing whitespaces).
<b>type strip:</b>	bool

```
System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs]
[topics]link-extractors.rst, line 153)
Unknown directive type "automethod".

.. automethod:: extract_links
```

## Link

**System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs] [topics]link-extractors.rst, line 158)**

Unknown directive type "module".

```
.. module:: scrapy.link
   :synopsis: Link from link extractors
```

**System Message: ERROR/3 (D:\onboarding-resources\sample-onboarding-resources\scrapy-master\docs\topics\[scrapy-master] [docs] [topics]link-extractors.rst, line 161)**

Unknown directive type "autoclass".

```
.. autoclass:: Link
```