
pytidylib

Release

December 14, 2015

1	Naming conventions	3
2	Installing HTML Tidy	5
3	Installing PyTidyLib	7
4	Small example of use	9
5	Configuration options	11
6	Function reference	13
	Index	15

PyTidyLib is a Python package that wraps the [HTML Tidy](#) library. This allows you, from Python code, to “fix” invalid (X)HTML markup. Some of the library’s many capabilities include:

- Clean up unclosed tags and unescaped characters such as ampersands
- Output HTML 4 or XHTML, strict or transitional, and add missing doctypes
- Convert named entities to numeric entities, which can then be used in XML documents without an HTML doctype.
- Clean up HTML from programs such as Word (to an extent)
- Indent the output, including proper (i.e. no) indenting for `pre` elements, which some (X)HTML indenting code overlooks.

As of the latest PyTidyLib maintenance updates, HTML Tidy itself has currently not been updated since 2008, and it may have trouble with newer HTML. This is just a thin Python wrapper around HTML Tidy, which is a separate project.

As of 0.2.3, both Python 2 and Python 3 are supported with passing tests.

Naming conventions

HTML Tidy is a longstanding open-source library written in C that implements the actual functionality of cleaning up (X)HTML markup. It provides a shared library (`so`, `dll`, or `dylib`) that can variously be called `tidy`, `libtidy`, or `tidylib`, as well as a command-line executable named `tidy`. For clarity, this document will consistently refer to it by the project name, HTML Tidy.

PyTidyLib is the name of the Python package discussed here. As this is the package name, `pip install pytidylib` is correct (they are case-insensitive). The *module* name is `tidylib`, so `import tidylib` is correct in Python code. This document will consistently use the package name, `PyTidyLib`, outside of code examples.

Installing HTML Tidy

You must have both [HTML Tidy](#) and [PyTidyLib](#) installed in order to use the functionality described here. There is no affiliation between the two projects. The following briefly outlines what you must do to install HTML Tidy. See the [HTML Tidy](#) web site for more information.

Linux/BSD or similar: First, try to use your distribution's package management system (`apt-get`, `yum`, etc.) to install HTML Tidy. It might go under the name `libtidy`, `tidylib`, `tidy`, or something similar. Otherwise see *Building from Source*, below.

OS X: You may already have HTML Tidy installed. In the Terminal, run `locate libtidy` and see if you get any results, which should end in `dylib`. Otherwise see *Building from Source*, below.

Windows: (Do not use pre-0.2.0 PyTidyLib.) You may be able to find prebuild DLLs. The DLL sources that were linked to in previous versions of this documentation have since gone 404 without obvious replacements.

Once you have a DLL (which may be named `tidy.dll`, `libtidy.dll`, or `tidylib.dll`), you must place it in a directory on your system path. If you are running Python from the command-line, placing the DLL in the present working directory will work, but this is unreliable otherwise (e.g. for server software).

See the articles [How to set the path in Windows 2000/Windows XP](#) (ComputerHope.com) and [Modify a Users Path in Windows Vista](#) (Question Defense) for more information on your system path.

Building from Source: The HTML Tidy developers have chosen to make the source code downloadable *only* through CVS, and not from the web site. Use the following CVS checkout at the command line:

```
cvs -z3 -d:pserver:anonymous@tidy.cvs.sourceforge.net:/cvsroot/tidy co -P tidy
```

Then see the instructions packaged with the source code or on the [HTML Tidy](#) web site.

Installing PyTidyLib

PyTidyLib is available on the Python Package Index:

```
pip install pytidylib
```

You can also download the latest source distribution from PyPI manually.

Small example of use

The following code cleans up an invalid HTML document and sets an option:

```
from tidylib import tidy_document
document, errors = tidy_document('''<p>f&otilde;o ''',
    options={'numeric-entities':1})
print document
print errors
```

Configuration options

The Python interface allows you to pass options directly to HTML Tidy. For a complete list of options, see the [HTML Tidy Configuration Options Quick Reference](#) or, from the command line, run `tidy -help-config`.

This module sets certain default options, as follows:

```
BASE_OPTIONS = {
    "indent": 1,          # Pretty; not too much of a performance hit
    "tidy-mark": 0,       # No tidy meta tag in output
    "wrap": 0,            # No wrapping
    "alt-text": "",       # Help ensure validation
    "doctype": 'strict',  # Little sense in transitional for tool-generated markup...
    "force-output": 1,    # May not get what you expect but you will get something
}
```

If you do not like these options to be set for you, do the following after importing `tidylib`:

```
tidylib.BASE_OPTIONS = {}
```

Function reference

`tidylib.tidy_document` (*text*, *options=None*, *keep_doc=False*)

Run a string with markup through HTML Tidy; return the corrected one.

text: The markup, which may be anything from an empty string to a complete (X)HTML document. If you pass in a unicode type (py3 str, py2 unicode) you get one back out, and tidy will have some options set that may affect behavior (e.g. named entities converted to plain unicode characters). If you pass in a bytes type (py3 bytes, py2 str) you will get one of those back.

options (dict): Options passed directly to HTML Tidy; see the HTML Tidy docs (<http://tidy.sourceforge.net/docs/quickref.html>) or run tidy -help-config from the command line.

keep_doc (boolean): If True, store 1 document object per thread and re-use it, for a slight performance boost especially when tidying very large numbers of very short documents.

returns (str, str): The tidied markup and unparsed warning/error messages. Warnings and errors are returned just as tidylib returns them.

`tidylib.tidy_fragment` (*text*, *options=None*, *keep_doc=False*)

Tidy a string with markup and return only the <body> contents.

HTML Tidy normally returns a full (X)HTML document; this function returns only the contents of the <body> element and is meant to be used for snippets. Calling tidy_fragment on elements that don't go in the <body>, like <title>, will produce incorrect behavior.

Arguments and return value are the same as tidy_document. Note that HTML Tidy will always complain about the lack of a doctype and <title> element in fragments, and these errors are not stripped out for you.

`tidylib.release_tidy_doc` ()

Release the stored document object in the current thread. Only useful if you have called tidy_document or tidy_fragment with keep_doc=True.

R

`release_tidy_doc()` (in module tidylib), [13](#)

T

`tidy_document()` (in module tidylib), [13](#)

`tidy_fragment()` (in module tidylib), [13](#)