
BeautifulSoup

Marcos Castro e Thomas William

O que é BeautifulSoup ?

Biblioteca Python projetada para facilitar a extração de dados nos documentos html e xml.

Principais Características

- Fornece alguns métodos simples e expressões pythonicas para navegar, pesquisar e modificar uma árvore de análise.
- Converte automaticamente documentos de entrada para Unicode e documentos de saída para UTF-8.
- Utiliza-se bibliotecas html parse lxml, html5lib, permitindo-lhe experimentar estratégias de análise diferentes ou velocidade comercial para a flexibilidade.
- Bastante robusta para trabalhar com html/xml mal formatado.

PARSER

Parser	Modo de uso	Vantagem	Desvantagem	
html.parser	BeautifulSoup(markup, "html.parser")	<ul style="list-style-type: none">• Incluso biblioteca• Velocidade moderada		
lxml	BeautifulSoup(markup, "lxml") BeautifulSoup(markup, "lxml-xml") BeautifulSoup(markup, "xml")	<ul style="list-style-type: none">• Rápida• Suporta xml	<ul style="list-style-type: none">• Dependência externa	C
html5lib	BeautifulSoup(markup, "html5lib")	<ul style="list-style-type: none">• Analisa páginas da mesma forma que um browser• Cria HTML5 válido	<ul style="list-style-type: none">• Lenta• Dependência externa	Python

OBRIGADO !