

# Web Scraping

Extracción de datos de Internet

# Fontes de datos

Administracións públicas

Portais de emprego

Redes sociais

Tendas en liña

Xornais dixitais

# Casos de uso

Aprendizaxe automática

Estudios de mercado

Vixilancia

# Requisitos mínimos

Cliente de HTTP  
Analizador de HTML

# Problemas

Dificultade

Rendemento

Seguridade

Mantemento

Heteroxeneidade

Detección



Scrapy

# Facilidade

Comportamento intuitivo

CSS e XPath

Exportación

Caché

Intérprete

# Rendimento

Asincronía

Control da simultaneidade



# Seguridad

Límites de dominios

Límites de tamaño

Límites de tempo

Protección contra vulnerabilidades

# Mantenimento

Elementos

Parametrización

Arañas base

Contratos

Xestión de ficheiros

Spidermon

Scrapyd

# Flexibilidad

Extensiones  
Python

# Flexibilidade

dateparser

extract

js2xml

parsel

price-parser

w3lib

# Camuflaxe

Comportamento

Proxy

scrapy-splash

scrapy-selenium

Titorial

[i.gal/scrappy](https://i.gal/scrappy)