

# scrapy-cookbook Documentation

ǎŘŚăŸĈ 0.2.2

Xiong Neng

11ǎlJĹ 10, 2017

## Contents

<b>1</b>	<b>ScrapyǎŦŽćĹŦ01- ǎĔĕĕŦĹćřĜ</b>	<b>1</b>
1.1	ǎŦĹ'ĕĕĔscrapy . . . . .	1
1.2	ćŏǎǎ■Ŧćđ'žăĹŦ . . . . .	2
1.3	ScrapyćĹ'žǎǎġăŸǎĔġĹ . . . . .	4
<b>2</b>	<b>ScrapyǎŦŽćĹŦ02- ǎŦŦǎŦŦ'ćđ'žăĹŦ</b>	<b>4</b>
2.1	ǎĹŽǎžžScrapyǎŦĕćĹŦ . . . . .	4
2.2	ǎŦŽǎžĹ'ǎĹŚǎžŋćŽĎItem . . . . .	5
2.3	ćŋŋăŸǎġŸĹSpider . . . . .	5
2.4	ĕĤŘĕǎŦćĹŦŋĕŽŋ . . . . .	6
2.5	ǎđ'ĎćŘĔĕŚĹǎŦŦĕ . . . . .	6
2.6	ǎřĵǎĜžǎĹŚǎŦŦŦǎŦŦǎŦŦŦ . . . . .	7
2.7	ǎĤĹǎ■ŸǎŦŦǎŦŦŦǎĹŦǎŦŦǎŦŦŦǎžŚ . . . . .	7
2.8	ǎŸŦăŸǎŦŦ . . . . .	9
<b>3</b>	<b>ScrapyǎŦŽćĹŦ03- Spiderĕĕĕĕĕĕĕĕ</b>	<b>9</b>
3.1	CrawlSpider . . . . .	9
3.2	XMLFeedSpider . . . . .	10
3.3	CSVFeedSpider . . . . .	11
3.4	SitemapSpider . . . . .	12
<b>4</b>	<b>ScrapyǎŦŽćĹŦ04- Selectorĕĕĕĕĕĕĕĕ</b>	<b>12</b>
4.1	ǎĔŸǎžŦŦĕǎĹ'ǎŦŦǎŦŦ . . . . .	12
4.2	ǎĵĤćŦĹĕǎĹ'ǎŦŦǎŦŦ . . . . .	12
4.3	ǎŦŦǎĕŦŦĕǎĹ'ǎŦŦǎŦŦ . . . . .	14
4.4	ǎĵĤćŦĹǎ■ĕǎĹŽĕǎĹĕĹǎĵĹ . . . . .	15
4.5	XPathćŽŸǎřǎžĕŦŦǎĹĎ . . . . .	16
4.6	XPathǎžžĕŦŦ . . . . .	16
<b>5</b>	<b>ScrapyǎŦŽćĹŦ05- Itemĕĕĕĕĕĕĕĕ</b>	<b>16</b>
5.1	ǎŦŽǎžĹ'Item . . . . .	17
5.2	Item Fields . . . . .	17
5.3	ItemǎĵĤćŦĹĕđ'žăĹŦ . . . . .	17
5.4	Item Loader . . . . .	18

7.3	āĽĒāyČāijRčĽñèŽń . . . . .	27
7.4	éŸšæ■cēcnārAçŽĎç■ŮçTě . . . . .	28
<b>8</b>	<b>ScrapyæTŽćĽŃ08- æŮGäzúäyŎāZčĽĜ</b>	<b>28</b>
8.1	ä;řçŤĽFiles Pipeline . . . . .	28
8.2	ä;řçŤĽImages Pipeline . . . . .	29
8.3	ä;řçŤĽäĽŃā■Ř . . . . .	29
8.4	èĜĽāōŽāzĽ'ātSä;řçōāéAř . . . . .	30
<b>9</b>	<b>ScrapyæTŽćĽŃ09- éČĽçš</b>	<b>30</b>
9.1	éČĽç;šāĽřScrapyd . . . . .	30
9.2	éČĽç;šāĽřScrapy Cloud . . . . .	33
<b>10</b>	<b>ScrapyæTŽćĽŃ10- āĽĽāĀāéĚ■ç;ōçĽñèŽń</b>	<b>33</b>
10.1	èĎŽæĽñèřŘēāŃScrapy . . . . .	33
10.2	āŘŃäyĀèĚŽćĽŃèřŘēāŃād'ŽäyĽspider . . . . .	34
10.3	āōŽāzĽ'èĝĎāĽŽēāĽ . . . . .	35
10.4	āōŽāzĽ'æŮĜçñāItem . . . . .	36
10.5	āōŽāzĽ'ArticleSpider . . . . .	36
10.6	çijŮāĚŽpipelineā■ŸāČĽāĽřæŤřæ■ōāžSäy■ . . . . .	37
10.7	āĚōāĽŤžrun.pyāŘřāĽĽèĎŽæĽñ . . . . .	38
<b>11</b>	<b>ScrapyæTŽćĽŃ11- æĽāæŃřçŽzā;Ť</b>	<b>39</b>
11.1	éĜ■āĚŽstart_requestsæŮzæřŤ . . . . .	41
11.2	ä;řçŤĽFormRequest . . . . .	41
11.3	éĜ■āĚŽ_requests_to_follow . . . . .	42
11.4	éāŤéĽčād'ĎçŘĒæŮzæřŤ . . . . .	43
11.5	āōŃæŤřæžŘçāĽ . . . . .	43
<b>12</b>	<b>ScrapyæTŽćĽŃ12- æĽřāŔŮāĽĽāĀç;řçñŽ</b>	<b>46</b>
12.1	scrapy-splashçōĀāzŃ . . . . .	46
12.2	āōĽ'èçĚdocker . . . . .	46
12.3	āōĽ'èçĚsplash . . . . .	47
12.4	āōĽ'èçĚscrapy-splash . . . . .	47
12.5	éĚ■ç;ōscrapy-splash . . . . .	48
12.6	ä;řçŤĽscrapy-splash . . . . .	48
12.7	ä;řçŤĽāōđäĽŃ . . . . .	50
<b>13</b>	<b>èĀŤçřzæĽŚ</b>	<b>52</b>

# 1 ScrapyæTŻćÍŃ01- ăĚěŮíçŕĜ

ScrapyæYřäyÄäyläyžāzēÇŁñáŔŬç;ŠčñZæȚTræ■ōrijNæRŘăŔŬçzŞæđDæĂgæȚræ■òèĀņcijŪăEŻçŽĐāžTçTlă  
ăLəæAřād'DçREælŬă■YăĆlăŎĖăRšæȚræ■ôç■LăyĂçşzăLŬçŽĐçİNăžRăy■ăĂCăĖüăIJăĂlăYřäyžāzēEęąę  
żşăŔřăžēăžTçTlăIJlėŐăŔŬAPIăLĂêƒTăŽďçŽĐæȚræ■ó(ærTăęĆWeb Ser-  
vices)ælŬëĂÉéĂŽçTlçŽĐç;ŠçzIJçLňěŽnáĂĆ

ScrapyäzšëĈ;äyöä;ääöđĎŔénŸéŸüćŽĎĹñèŽnææEæđüiijŇnærŤæĈĹnăRŮæŮüćŽĎĉ;ŚçñŽèöd'èrAăĂAăEĹ

## 1.1 aŏL'èčĚscrapy

æŁŚçŽĐæŧÑèŕŦçŎřăćČæŸŕcentos6.5

ǎ■ĞçžğpythoňǎŁřǎIJǎŮřčŁ'ŁčŽĐ2.7ijǎŸNǎŸNéíçčŽĐǎŁ'ǎǎIJŁ'ǎ■ééłd' éČǎŁĞǎ■čǎŁřrootčŤǎŁǎŁǎ

çTšāŽŌscrapycŽōāL■āRlēČ;ēfŘēāŇāIJlpython2äyŁiijŇāL'ÄāžēāĒLāŽt'æŮrcentosäyŁēlčČŽDpythonāLřāI.  
Python 2.7.11iijŇāĒŮā;ŠæŮžæšTērŭgoogleäyŇā;Łād'ŽēfŽāăŭčŽDæTŽčlŇāĀČ

åĖŁăŃŁ'ěčĚäÿĂäžZă; İetŰè;rázű

```
yum install python-devel
yum install libffi-devel
yum install openssl-devel
```

çĐúåŘŒåŕL'èčĚpyopensslåžŞ

```
pip install pyopenssl
```

ãóL'èčĚxlm1

```
yum install python-lxml
yum install libxml2-devel
yum install libxslt-devel
```

## Service-identity

```
pip install service-identity
```

ãŁL'èčĚtwisted

```
pip install scrapy
```

# ãŁL'èčĚscrapy

```
pip install scrapy -U
```

ætÑèrTscrapy

scrapy bench

æIJĂçZŁæŁŔăŁşiiJŃăd'łäy■ăózáæYŞăZĘiijA

## 1.2 scrapy

scrapy spider python scrapy stackoverflow.py scrapy

```
import scrapy

class StackOverflowSpider(scrapy.Spider):
    name = 'stackoverflow'
    start_urls = ['http://stackoverflow.com/questions?sort=votes']

    def parse(self, response):
        for href in response.css('.question-summary h3 a::attr(href)'):
            full_url = response.urljoin(href.extract())
            yield scrapy.Request(full_url, callback=self.parse_question)

    def parse_question(self, response):
        yield {
            'title': response.css('h1 a::text').extract()[0],
            'votes': response.css('.question .vote-count-post::text').extract()[0],
            'body': response.css('.question .post-text').extract()[0],
            'tags': response.css('.question .post-tag::text').extract(),
            'link': response.url,
        }
```

scrapy

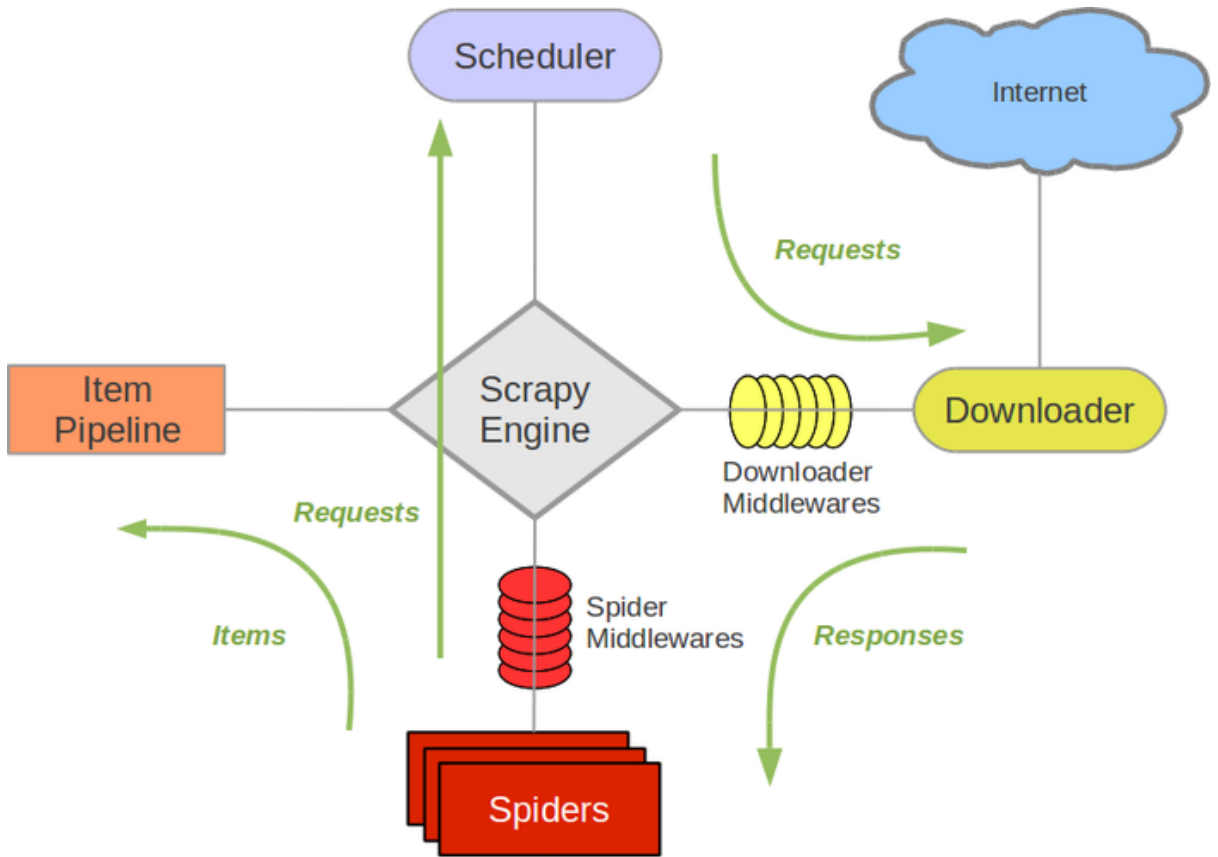
```
scrapy runspider stackoverflow_spider.py -o top-stackoverflow-questions.json
```

scrapy

```
[{
    "body": "... LONG HTML HERE ...",
    "link": "http://stackoverflow.com/questions/11227809/why-is-processing-a-sorted-array-faster-than-an-unsorted-array",
    "tags": ["java", "c++", "performance", "optimization"],
    "title": "Why is processing a sorted array faster than an unsorted array?",
    "votes": "9924"
},
{
    "body": "... LONG HTML HERE ...",
    "link": "http://stackoverflow.com/questions/1260748/how-do-i-remove-a-git-submodule",
    "tags": ["git", "git-submodules"],
}
```

```
    "title": "How do I remove a Git submodule?",
    "votes": "1764"
},
...]
```

Scrapy runspider somefile.py  
start\_urls = ['http://www.example.com/']  
parse(self, response):  
 return scrapy.Request('http://www.example.com/next-page/', callback=self.parse)



Scrapy is a web crawling framework for Python. It is designed to be fast and efficient, and it is easy to use. It is a powerful tool for web scraping and data extraction. It is a framework for writing web crawlers. It is a tool for extracting data from websites. It is a framework for writing web scrapers. It is a tool for extracting data from websites.

### 1.3 Scrapy L'zæĀgäyĀèġĹ

Scrapy is a web crawling framework for Python. It is designed to be fast and efficient, and it is easy to use. It is a powerful tool for web scraping and data extraction. It is a framework for writing web crawlers. It is a tool for extracting data from websites.

1. Scrapy is a web crawling framework for Python. It is designed to be fast and efficient, and it is easy to use. It is a powerful tool for web scraping and data extraction. It is a framework for writing web crawlers. It is a tool for extracting data from websites.
2. Scrapy is a web crawling framework for Python. It is designed to be fast and efficient, and it is easy to use. It is a powerful tool for web scraping and data extraction. It is a framework for writing web crawlers. It is a tool for extracting data from websites.
3. Scrapy is a web crawling framework for Python. It is designed to be fast and efficient, and it is easy to use. It is a powerful tool for web scraping and data extraction. It is a framework for writing web crawlers. It is a tool for extracting data from websites.

4. `ĀĀēāčōçŽĎçijŮčāĀæŤræŇĀāŠŇēĠāĽĽēŕĒāĽĽīijŇçŤĭāžŌād'ĎçŘĒād'ŮæŮĠāĀĀēĭđæāĠāĠĒāŠŇēŤž`
5. `āŖŕæĽĽ'āsŤriijŇāĒĀēōŷā;āā;ĤçŤĭsignals āŠŇāŖŇāē;çŽĎAPI(middlewares, extensions, āŠŇpipelines)æĬēçijŮāĒĒēĠāōŽāžĽ'æŖŠāžŭāĽšēČ;āĀČ`
6. `ād'ġēĠŖçŽĎāĒĒēç;ōæĽĽ'āsŤāŠŇāy■ēŮŕ'āžŭā;Žā;ĤçŤĭriijŽ`
  - cookies and session handling
  - HTTP features like compression, authentication, caching
  - user-agent spoofing
  - robots.txt
  - crawl depth restriction
  - and more
7. `ēĤŸæĬĽ'āĒŮāžŮāē;ād'Žāē;āyĬāyĬriijŇæŕŤæČāŖŕéĠād'■āĽĽ'çŤĭēĬŸēŽŽæĬēçĽŇāŖŮSitemapsāŠŇX āyĀāyĽēūšçĽŇāŖŮāĒČçŕ'āāĒšēāŤçŽĎāĽšā;šçōæĀšæĬ ēĠāĽĽāyŇē;āŽçĽĠriijŇ āyĀāyĽçijšā■ŸDNSēġçæđŖāŽĬ■Ľç■Ľ'āĀČ`

## 2 scrapyĤŽçĭŇ02- āōŇæŤŕ'çđ'žāçŇ

ēĤŽçŕĠæŮĠçŇāæĽšāžŇēĀŽēĤĠāyĀāyĽæŕŤē;ČāōŇæŤŕ'çŽĎāçŇā■ŖæĬæŤŽā;āā;ĤçŤĭScrapyriijŇæĽšēĀĽæĽēŽēĠŇæĽšāžŇāŕĒāōŇæĽŖāēČāyŇāĠāāyĽæ■ēēĽđ'riijŽ

- `āĽŽāžžāyĀāyĽæŮŕçŽĎScrapyāūēçĭŇ`
- `āōŽāžĽ'ā;āæĽ'ĀēĬĀēçĀēçĀæĽ;āŖŮçŽĎItemāržēsā`
- `çijŮāĒĒāyĀāyĽspideræĬēçĽŇāŖŮæšŖāyĽç;šçŇŽāžŭæŖŖāŖŮāĠžæĽ'ĀæĬĽ'çŽĎItemāržēsā`
- `çijŮāĒĒāyĀāyĽItem PipelineæĬē■ŸāČĭæŖŖāŖŮāĠžæĬēçŽĎItemāržēsā`

Scrapyā;ĤçŤĭPythonēŕ■ēĬāçijŮāĒĒriijŇāēČæđĬā;āāržēĤŽēŮĽēŕ■ēĬĀēĤŸāy■çĒēriijŇēŕŭāĒĽāŌžā■ēāžāyŇāš

### 2.1 āĽŽāžžScrapyāūēçĭŇ

āĬĽāžžā;Ťā;āāŮĬæŇççŽĎçŽōā;ŤæĽ'ġēāŇāēČāyŇāš;āžđ'

```
scrapy startproject coolscrapy
```

ārĒāijŽāĽŽāžžcoolscrapyæŮĠāžŭād'žriijŇāĒŮççŽōā;ŤçžšæđĎāēČāyŇriijŽ

```
coolscrapy/
  scrapy.cfg                                # ēČĭç;šēĒ■ç;ōæŮĠāžŭ

coolscrapy/                                #_
→PythonāĬāāĬŮīijŇā;āæĽ'ĀæĬĽ'çŽĎāžççāĀēČ;æŤçēĤŽēĠŇēĬç
  __init__.py
```

```

items.py          # Item
pipelines.py      # pipelines
settings.py       # éĚ; ó
spiders/          #
→ æL'ÄæIJL'çĹñžńspideréČ; æT; èfžäylæŮĜäzúâd'žäyNéÍć
    __init__.py
    ...

```

## 2.2 áŮŽázL'æĹSäzñçŽDItem

æĹSäzñæÄžēĜāĹZāzžäyÄäyĭscrapy.ItemçšzĭijNāzūáŮŽázL'áoČçŽDçšzāđNäyžscrapy.FieldçŽDāšđæÄġĭijN  
æĹSäzñāĜĖāđ'ĜārĖēžŌāŮĖē; ŠæŮřēŮzāĹŮēāĹçŽDāR■ġrāÄAēŠ; æŌēāIJřāĹAāŠNæŠŸēēAçĹnāRŮäyNæĹ

```

import scrapy

class HuxiuItem(scrapy.Item):
    title = scrapy.Field()      # æăĜéćŸ
    link = scrapy.Field()      # éŞ; æŌě
    desc = scrapy.Field()      # çŏÄěĹř
    posttime = scrapy.Field()  # āŔŚāyČæŮúéŮt'

```

ázšēōyā; äēġL'ā; ŮáoŽázL'èfžäyĭItemæIJL'çČzéžzçČēĭijNā; ĖæŸřáoŽázL'áoNāzNāŔŌā; āāŔřāzēā; ŮāĹřēōyāō

## 2.3 çñňäyÄäyĭSpider

èIJŸèŽŽāŕsæŸřā; āáoŽázL'çŽDäyÄäyĭçšzĭijNScrapyā; ĺçŤĹáoČāznæĹēazŌäyÄäyĭdomainĭijĹæĹŮdomainçžĹ  
āIJĹèIJŸèŽŽçšzäy■áoŽázL'āžĖäyÄäyĭāĹĹāġNāNŮçŽDURLäyNē; ĺĹŮēāĹĭijNāzēāŔĹæĀŌæāüēüšēyĹēŞ; æŌē

áoŽázL'äyÄäyĭSpiderĭijNāŔĹēIJÄçžġæĹŸscrapy.SpiderçšzāzūáoŽázŌäyÄäyĭZāšđæÄġĭijŽ

- name: SpiderāŔ■ġřĭijNāŸĖēāzæŸřāŤřäyÄçŽD
- start\_urls: āĹĹāġNāNŮäyNē; ĺēŞ; æŌēURL
- parse(): çŤĹæĹēēġčæđŔäyNē; ĺāŔŌçŽDResponseāržēsāĭijNēřēāržēsāzšæŸřēfžäyĹæŮzæşŤçŽDāŤřäyÄ  
áoČēŧ'şēŧ'çēġčæđŔēŸŤāžđēāŧēĹæŤŕæ■áožūæŔŔāŔŮāĜçŽyāžŤçŽDItemĭijĹēŸŤāžđItemāržēsāĭijĹĭij

æĹSäzñāIJĹcoolscrapy/spidersæŮĜäzúâd'žäyNéÍćæŮřāžzhuxiu\_spider.  
pyĭijNāĖĖáožāçCäyNĭijŽ

```

#!/usr/bin/env python
# -*- encoding: utf-8 -*-
"""
Topic: sample
Desc :

```





```

        item = HuxiuItem()
        item['title'] = sel.xpath('h3/a/text()')[0].extract()
        item['link'] = sel.xpath('h3/a/@href')[0].extract()
        url = response.urljoin(item['link'])
        item['desc'] = sel.xpath('div[@class="mob-sub"]/text()
→') [0].extract()
        # print(item['title'], item['link'], item['desc'])
        yield scrapy.Request(url, callback=self.parse_article)

    def parse_article(self, response):
        detail = response.xpath('//div[@class="article-wrap"]')
        item = HuxiuItem()
        item['title'] = detail.xpath('h1/text()')[0].extract()
        item['link'] = response.url
        item['posttime'] = detail.xpath(
            'div[@class="article-author"]/span[@class="article-time
→"]/text()')[0].extract()
        print(item['title'], item['link'], item['posttime'])
        yield item

```

çŒřaIJłparseaRłæRŕaRŪæĐſaĖř'èúčçŽĐéŞ;æŒčřijNçĐúaRŒřEéŞ;æŒčřaEĖăăóžèğčæđRăžd'çžZăRęăđ'ŪčŽăăRfăžăăſžăžŒŒfŽăyłæđĐăžžæŽř'ăLăăđ'■æİČçŽĐçLněŽńçłNăžRăžEăĂĆ

## 2.6 aříjaĠžæŁſaRŪæŤřæ■ó

æIJăčŒăă■ŤçŽĐăfIă■ŸæŁſaRŪæŤřæ■ŒçŽĐæŪžaijRæŸřajçŤłjsonæaijajRçŽĐæŪĠăžúăfIă■ŸăIJłæIJňăł

```
scrapy crawl huxiu -o items.json
```

ăIJłæijŤçđ'žçŽĐăřRçşçžçſéĠŒéİçèfŽçğ■æŪžaijRèúşăđ'şăžEăĂĆăy■èfĠăçĆăđIJăjăèçAăđĐăžžăđ'■æİČçŽăIJăăčjèĠăăſçijŪăEĖŽItem PipelineăĂĆ

## 2.7 äŕIă■ŸæŤřæ■ŒăLřæŤřæ■ŒăžŞ

äyŁéİçăŁſăžňăžNçž■ăžEăRfăžăăřEăŁſaRŪçŽĐItemaříjaĠžăyžjsonæaijajRçŽĐæŪĠăžúijNăy■èfĠăæIJăă pipelines.pyăŒŽăžL'

```

# -*- coding: utf-8 -*-
import datetime
import redis
import json
import logging
from contextlib import contextmanager

from scrapy import signals
from scrapy.exporters import JsonItemExporter
from scrapy.pipelines.images import ImagesPipeline

```



```
scrapy crawl huxiu
```

éCčázLæL'ĂæIJL'æŮřéŮzçŽDæŮĜčnáěČ;ā■ŸāĆlāLřæŤræ■óāžŠäy■āŌzāžEāĂĆ

## 2.8 äŸNäyĂæ■ě

æIJñčnáāRlæŸřäyēä;ăéCēçŤěāžEçscrapyæIJĀāšžæIJñçŽDāLšèČ;īijNēŸæIJL'ā;Ĺād'ŽénŸçžġL'zæĂġæšqæL

## 3 ScrapyæŤŽćÍŃ03- Spiderèrēēğč

SpideræŸřćĹñèŽnáæAēæđŮçŽDæäyāŸČīijŃćĹñāRŮætAćÍŃāēCāyŃīijŽ

1. āĚĹāĹiāġŃāNŮērŮæsĆURLāĹŮēāīijŃāžŮæŃĠāōŽäyNē;āŖŌād'ĐçRēresponseçŽDāŽdērČāĠ;æŤrā
2. āIJĹparseāŽdērČäy■ēğčæđRresponseāžŮēŸŤāŽđā■ŮāĚŸ,Itemāržèšq,RequestāržèšqæĹŮāōČāžñçŽD
3. āIJĹāŽdērČāĠ;æŤrēĠŃēĹčīijŃā;ăēĂŽēŸĠā;ŸćŤĹéĀL'æŃĹ'āŽīijĹāŖŃæāŮāŖfāžēä;ŸćŤĹBeautifulSoup,L
4. æIJĀāŖŌēŸŤāŽđçŽDēŸŽāžŽIteméĂŽäyŸāijŽèčŃæŃĀžĚāŃŮāĹræŤræ■óāžŠäy■(ä;ŸćŤĹItem Pipeline)æĹŮēĂĚä;ŸćŤĹFeed exportsārĒāĚŮāŸĹā■ŸāĹræŮĠāžŮäy■āĂĆ

ār;çōāēŸŽäyĹætAćÍŃéĂĆāŖĹāžŌæL'ĂæIJL'çŽDēIJŸēŽŽīijŃā;EæŸřScrapyéĠŃēĹcāyžäy■āŖŃçŽDā;ŸćŤĹçŽ

### 3.1 CrawlSpider

éŞ;æŌēçĹñāRŮēIJŸēŽŽīijŃāyŞēŮĹäyžēČčāžŽćĹñāRŮæIJL'çL'zāōŽēġDā;ŃçŽDēŞ;æŌēāĒēāōžēĀŃāĠEāđ  
āēČæđIJā;ăēġL'ā;ŮāōČēŸŸäy■ēŮšāžēēĂĆāŖĹā;ăçŽDēIJĀæsĆīijŃāŖfāžēāĚĹçžġæL'ŸāōČçDŮāŖŌēēEçŽŮçŽ  
āōČēŽd'āžEāžŌscrapy.SpiderçşzçžġæL'ŸçŽDāsđæĂġād'ŮīijNēŸæIJL'äyĂäyĹæŮřçŽDāsđæĂġrules  
äyĂäyĹèrēççEçŽDā;Ńā■ŖīijŽ

```
from coolscrapy.items import HuxiuItem
import scrapy
from scrapy.spiders import CrawlSpider, Rule
from scrapy.linkextractors import LinkExtractor

class LinkSpider(CrawlSpider):
    name = "link"
    allowed_domains = ["huxiu.com"]
    start_urls = [
        "http://www.huxiu.com/index.php"
    ]

    rules = (
        # æŖŖāŖŮāŃžēĚ■æ■čāĹžāijŖ'/group?f=index_group'és;æŌē_
        → (ä;EæŸřäy■èČ;āŃžēĚ■'deny.php')
```

```

#_
→āzūāyTāijŽēĀŠā;ŠčĹñāRŮ(āęĆāđIJæšæIJL'āóŽāžL'callbackiijNěžYēōd'follow=True).
→
    Rule(LinkExtractor(allow=('/group?f=index_group', ), deny=(
→'deny\.php', )),
        # æŘŘāRŮāNžéĚ'/'article/\d+/\d+.html
→'čŽĐéŞ;æŌěiijNāžŭä;£çTlparse_
→itemæIēēğčāđRāōČāžñäYNè;āŘŌčŽĐāĚĚāōžiiijNäy■éĀŠā;Š
        Rule(LinkExtractor(allow=('/article/\d+/\d+\.html', )),_
→callback='parse_item'),
    )

def parse_item(self, response):
    self.logger.info('Hi, this is an item page! %s', response.
→url)

    detail = response.xpath('//div[@class="article-wrap"]')
    item = HuxiuItem()
    item['title'] = detail.xpath('h1/text()')[0].extract()
    item['link'] = response.url
    item['posttime'] = detail.xpath(
        'div[@class="article-author"]/span[@class="article-time
→']/text()')[0].extract()
    print(item['title'], item['link'], item['posttime'])
    yield item

```

## 3.2 XMLFeedSpider

XMLēōcéYĚēIJYēZZiijNčTlælēčĹñāRŮXMLā;ćaijRčŽĐēōcéYĚāĚĚāōžiiijNéĀŽē£Gæ\$ŘäyŁæNĠāōŽčŽĐē  
āŘřā;£çTlinternodes, xml, āŠNhtmläyLčg■ā;ćaijRčŽĐēf■āžčāŽliijNäy■ē£Gā;ŠāĚĚāōžærTē;Čād'ŽčŽĹ  
ézYēōd'āžšæYřāōČiijNāŘřāžēēŁČçIJAāĚĚā■YæŘŘā■GæĀğèČ;iiijNäy■ēIJĀēçAārĚæTt'äyĪDOMāŁæ;ĵāŁřā  
ād'DčŘĚXMLčŽĐæŮūāĀŽæIJĀāē;āĚĹRemoving namespaces

æŌēäyNāIēæŁŚéĀŽē£GčĹñāRŮæŁŚčŽĐā■ŽāōcéōcéYĚXMLæIēāsTčd'žāōČčŽĐä;£çTlæŮzæşTāĀĆ

```

from coolscrapy.items import BlogItem
import scrapy
from scrapy.spiders import XMLFeedSpider

class XMLSpider(XMLFeedSpider):
    name = "xml"
    namespaces = [('atom', 'http://www.w3.org/2005/Atom')]
    allowed_domains = ["github.io"]
    start_urls = [
        "http://www.pycoding.com/atom.xml"
    ]
    iterator = 'xml' #_
→çijžçIJAčŽĐiternodesiijNěšNäiijijāržžāžŌæIJL'namespacečŽĐxmläy■ēāN

```

```

        itertag = 'atom:entry'

    def parse_node(self, response, node):
        # self.logger.info('Hi, this is a <%s> node!', self.itertag)
        item = BlogItem()
        item['title'] = node.xpath('atom:title/text()')[0].extract()
        item['link'] = node.xpath('atom:link/@href')[0].extract()
        item['id'] = node.xpath('atom:id/text()')[0].extract()
        item['published'] = node.xpath('atom:published/text()')[0].
        ↪extract()
        item['updated'] = node.xpath('atom:updated/text()')[0].
        ↪extract()
        self.logger.info(''.join([item['title'], item['link'], item[
        ↪'id'], item['published']]))
        return item

```

### 3.3 CSVFeedSpider

èŁŻäÿlèu\$äyLéíçŽĐXMLFeedSpiderā;ŁčśzäijijNāNžāLnáIJläžŌāóČäijŽäyĀëāNäyĀëāNčŽĐëŁ■äzčrijNē  
æŕRæñàëŁ■äzčëāNčŽĐæŬūāĀŽäijŽërČŦĪparse\_row() æŬžæşŦāĂĆ

```

from coolscrapy.items import BlogItem
from scrapy.spiders import CSVFeedSpider

class CSVSpider(CSVFeedSpider):
    name = "csv"
    allowed_domains = ['example.com']
    start_urls = ['http://www.example.com/feed.csv']
    delimiter = ';'
    quotechar = '"'
    headers = ['id', 'name', 'description']

    def parse_row(self, response, row):
        self.logger.info('Hi, this is a row!: %r', row)
        item = BlogItem()
        item['id'] = row['id']
        item['name'] = row['name']
        return item

```

### 3.4 SitemapSpider

çñŽçĆzāIJŕāŽ;èIJŸëŽŽijNāĒĀèöÿä;ää;ŁçŦĪSitemapsāŖŚçŎŕURLāŖŎçŁňāŔŬæŦŦäÿŁçñŽçĆzāĂĆ  
èŁŸæŦŕæNĀāŦNāëŬçŽĐçñŽçĆzāIJŕāŽ;äzēāŖLäžŎrobots.txtäÿ■āŖŚçŎŕçñŽçĆURL

## 4 Scrapy Selector

Scrapy Selector is a class that represents a part of the HTML document. It is a subclass of `lxml.etree.Element`.

1. `BeautifulSoup` is a Python library that provides a simple and powerful API for parsing HTML and XML documents. It uses `lxml` as the underlying parser.
2. `lxml` is a Python library that provides a fast and powerful API for parsing XML and HTML documents. It uses `libxml2` as the underlying parser.

Scrapy uses `lxml` as the underlying parser. It provides a simple and powerful API for parsing HTML and XML documents. It uses `lxml` as the underlying parser.

Scrapy uses `lxml` as the underlying parser. It provides a simple and powerful API for parsing HTML and XML documents. It uses `lxml` as the underlying parser.

### 4.1 Selector

Scrapy Selector is a class that represents a part of the HTML document. It is a subclass of `lxml.etree.Element`.

- `/html/head/title`: selector for the title element.
- `/html/head/title/text()`: selector for the text of the title element.
- `//td`: selector for the td element.
- `//div[@class="mine"]`: selector for the div element with class="mine".

Scrapy Selector is a class that represents a part of the HTML document. It is a subclass of `lxml.etree.Element`.

- `xpath()`: selector for the xpath expression.
- `css()`: selector for the css selector.
- `extract()`: selector for the extract method.
- `re()`: selector for the re method.

### 4.2 Selector

Scrapy Selector is a class that represents a part of the HTML document. It is a subclass of `lxml.etree.Element`.

```
<html>
<head>
  <base href='http://example.com/' />
  <title>Example website</title>
</head>
```

```
<body>
  <div id='images'>
    <a href='image1.html'>Name: My image 1 <br /><img src='image1_
    thumb.jpg' /></a>
    <a href='image2.html'>Name: My image 2 <br /><img src='image2_
    thumb.jpg' /></a>
    <a href='image3.html'>Name: My image 3 <br /><img src='image3_
    thumb.jpg' /></a>
    <a href='image4.html'>Name: My image 4 <br /><img src='image4_
    thumb.jpg' /></a>
    <a href='image5.html'>Name: My image 5 <br /><img src='image5_
    thumb.jpg' /></a>
  </div>
</body>
</html>
```

Scrapy Shell

```
scrapy shell http://doc.scrapy.org/en/latest/_static/selectors-
sample1.html
```

Selectors

```
>>> response.xpath('//title/text()')
[<Selector (text) xpath=//title/text()>]
>>> response.css('title::text')
[<Selector (text) xpath=//title/text()>]
```

SelectorsList

```
>>> response.css('img').xpath('@src').extract()
[u'image1_thumb.jpg',
 u'image2_thumb.jpg',
 u'image3_thumb.jpg',
 u'image4_thumb.jpg',
 u'image5_thumb.jpg']
```

Selector

```
>>> response.xpath('//div[@id="images"]/a/text()').extract_first()
u'Name: My image 1 '
```

Selector

```
>>> response.xpath('//div[@id="not-exists"]/text()').extract_
first(default='not-found')
'not-found'
```

Selector

```
>>> response.css('title::text').extract()
[u'Example website']
```

äyÑéÍcæYřăĜäăylærTëĭČăĚléÍcžŽDčd'žăĭNijŽ

```
>>> response.xpath('//base/@href').extract()
[u'http://example.com/']

>>> response.css('base::attr(href)').extract()
[u'http://example.com/']

>>> response.xpath('//a[contains(@href, "image")]/@href').extract()
[u'image1.html',
 u'image2.html',
 u'image3.html',
 u'image4.html',
 u'image5.html']

>>> response.css('a[href*=image]::attr(href)').extract()
[u'image1.html',
 u'image2.html',
 u'image3.html',
 u'image4.html',
 u'image5.html']

>>> response.xpath('//a[contains(@href, "image")]/img/@src').
↳extract()
[u'image1_thumb.jpg',
 u'image2_thumb.jpg',
 u'image3_thumb.jpg',
 u'image4_thumb.jpg',
 u'image5_thumb.jpg']

>>> response.css('a[href*=image] img::attr(src)').extract()
[u'image1_thumb.jpg',
 u'image2_thumb.jpg',
 u'image3_thumb.jpg',
 u'image4_thumb.jpg',
 u'image5_thumb.jpg']
```

## 4.3 ātŇăěŮéĂL'æŇřăŽÍ

xpath() āŠŇcss() ēĚTăŽđčŽDæYřéĂL'æŇřăŽÍăĹŮēāliijŇæL'ĂăžěăĭăăŔřăžēčžğcz■ăĭĚčTĭăőČăžŋčŽDæŮ

```
>>> links = response.xpath('//a[contains(@href, "image")]')
>>> links.extract()
[u'<a href="image1.html">Name: My image 1 <br></a>',
 u'<a href="image2.html">Name: My image 2 <br></a>',
```





```
# æĹŮèĀĚäŸNéÍcéŁŻäŸłçŽt'æŎěä;ŁçŤłpăžŖăŖăžě
>>> for p in divs.xpath('p'):
...     print p.extract()
```

## 4.6 XPath

### ä;ŁçŤłtextä;łJäŸžæłäžúæŮŮ

éĀĤăĚ■ä;ŁçŤł.//text(),çŽt'æŎěä;ŁçŤł.

```
>>> sel.xpath("//a[contains(., 'Next Page')]").extract()
[u'<a href="#">Click here to go to the <strong>Next Page</strong></a>']
```

### //node[1]āŠŇ(//node)[1]āŇžāĹń

- //node[1]: éĀĹ'æŇĹ'æĹ'ĀæĹĹ'ä;■äžŎçňňäŸĀäŸł■ŘèŁČçČžä;■ç;ŏçŽĐnodeèŁČçČž
- (//node)[1]: éĀĹ'æŇĹ'æĹ'ĀæĹĹ'çŽĐnodeèŁČçČžīīŇçĐŮāŖŎěĤāŽđçžŖăđĹäŸ■çŽĐçňňäŸĀäŸłnodeè

### éĀŽèŁĠclassæŖæŁ;æŮŮäijŸăĚĹèĀČèŽŖCSS

```
>> from scrapy import Selector
>>> sel = Selector(text='<div class="hero shout"><time datetime=
↳ "2014-07-23 19:00">Special date</time></div>')
>>> sel.css('.shout').xpath('./time/@datetime').extract()
[u'2014-07-23 19:00']
```

## 5 Scrapy

ItemæŸřăĹă■ŸçžŖăđĐæŤřæ■ŏçŽĐăĹĹ'æŮžīīŇScrapyăŖřăžěăŖĚèğçæđŘçžŖăđĹăžěă■ŮăĚŸă;ăāijŘèĤăŽđ  
ItemæŘŖă;ŽăžĚçšžă■ŮăĚŸçŽĐĀĹīījŇăžŮäŸŤăŖřăžěă;ĹăŮžă;ŁçŽĐăčřæŸŎă■ŮăŏŧīījŇă;Ĺăđ'ŽScrapyçžđ

### 5.1 áŎŽăžĹ'Item

áŎŽăžĹ'IteméĹđăŸŸçŏĀă■ŤīījŇăŖĹéĹĀèçĀçžğæĹ'Łscrapy.ItemçšžīījŇăžŮăŖĚăĹ'ĀæĹĹ'ă■ŮăŏŧéČ;áŎŽăž  
FieldçšžăđŇă■ŖăŖ

```
import scrapy

class Product(scrapy.Item):
```

```
name = scrapy.Field()
price = scrapy.Field()
stock = scrapy.Field()
last_updated = scrapy.Field(serializer=str)
```

## 5.2 Item Fields

FieldárzèšqáRřçTlæIěárzærRäyłā■ŮæōļæŇĠăōŽăĚĈæTřæ■őăĂĆăŇăęĆăyŁéÍćlast\_updatedçŽĎžŘ

## 5.3 ItemäĲŁçd'žăĲ

äĲääĲŽçIJŇÁĹItemçŽĎäĲçTlèŮšPythonäy■çŽĎ■ŮăĚyAPIéĎăyçšzäĲĲ

### ăĹŽăžItem

```
>>> product = Product(name='Desktop PC', price=1000)
>>> print product
Product(name='Desktop PC', price=1000)
```

### ěŮăŖŮăĲĲ

```
>>> product['name']
Desktop PC
>>> product.get('name')
Desktop PC

>>> product['price']
1000

>>> product['last_updated']
Traceback (most recent call last):
...
KeyError: 'last_updated'

>>> product.get('last_updated', 'not set')
not set

>>> product['lala'] # getting unknown field
Traceback (most recent call last):
...
KeyError: 'lala'

>>> product.get('lala', 'unknown field')
'unknown field'
```

```

>>> 'name' in product    # is name field populated?
True

>>> 'last_updated' in product    # is last_updated populated?
False

>>> 'last_updated' in product.fields    # is last_updated a declared_
↪field?
True

>>> 'lala' in product.fields    # is lala a declared field?
False

```

## èõççíóåĀij

```

>>> product['last_updated'] = 'today'
>>> product['last_updated']
today

>>> product['lala'] = 'test' # setting unknown field
Traceback (most recent call last):
...
KeyError: 'Product does not support field: lala'

```

## èõĖéŮóæĹ'ÄæĴĴçŽĎåĀij

```

>>> product.keys()
['price', 'name']

>>> product.items()
[('price', 1000), ('name', 'Desktop PC')]

```

## 5.4 Item Loader

Item LoaderäyžæĹŚäznæŘŘäĴžĚçĤŦæĹŘItemçŽĎçŽyåĴŞäĴåĹĴçŽĎæŮzæşŦăĂĆItemäyžæĹŞåŮçŽĎæĴæŁäŮçæĹæĹŚäznéĴdäyŷæŮzäĴçŽĎăŦĚĴŞăĖăăăăăĖĖăĴăőzăŽĴäy■ăĂĆ

äyŊéĴæĹŚäznéĂŽĖĴGäyĂäyĴäĴŊ■ŘăĴăăŦçĎ'žäyĂĖĴăĴçŦĴăŮzæşŦijŽ

```

from scrapy.loader import ItemLoader
from myproject.items import Product

def parse(self, response):
    l = ItemLoader(item=Product(), response=response)
    l.add_xpath('name', '//div[@class="product_name"]')

```

```

l.add_xpath('name', '//div[@class="product_title"]')
l.add_xpath('price', '//p[@id="price"]')
l.add_css('stock', 'p#stock')
l.add_value('last_updated', 'today') # you can also use literal_
↪ values
return l.load_item()

```

æşlæĐRäyŁÉlćŻĐnameå■ŮæŏtæŸřázŎäyđ'äyłxpathèûrâ;Đæûzçt'ráŁääŔŎâ;ŮåŁřāĂĆ

## 5.5 èŁŞåĖĖ/èŁŞåĠžād'ĐçŘĚāZl

æřRäyHtem LoaderåržæřRäyHFieldéČ;æIJL'äyÄäyŁèŁŞåĖĖād'ĐçŘĚāZlāŠŇäyÄäyŁèŁŞåĠžād'ĐçŘĚāZlāĂĆ  
load\_item() æŮŮāĖ■æŁğèāŇèŁŞåĠžād'ĐçŘĚāZlīijŇèŁŤāŽđæIJĂçzŁçzŞæđIJāĂĆ

```

l = ItemLoader(Product(), some_selector)
l.add_xpath('name', xpath1) # (1)
l.add_xpath('name', xpath2) # (2)
l.add_css('name', css) # (3)
l.add_value('name', 'test') # (4)
return l.load_item() # (5)

```

æŁğèāŇætĴĂçlŇæŸřèŁŻæăŭçŽĐriijŽ

1. xpath1äy■çŽĐæŤřæ■ŏèćŇæŔŔârŮāĠžæİēriijŇçĐŮāŔŎäijæèŁŞāŁřnameå■ŮæŏtçŽĐèŁŞåĖĖād'ĐçŘĚāZlāĂĆ  
LoaderéĠŇÉlć(èŁŻæŮŮāĂŽæşşæIJL'ètŇŇāĴijçzŽitem)
2. xpath2æŤřæ■ŏèćŇæŔŔârŮāĠžæİēriijŇçĐŮāŔŎäijæèŁŞçzŽ(1)äy■āŔŇæăŭçŽĐèŁŞåĖĖād'ĐçŘĚāZlīijŇæŁçzŞæđIJāĂĆ
3. èŮş2äyĂæăŭ
4. èŮş3äyĂæăŭriijŇäy■èŁĠèŁŻæŇqæŸřçŽt'æŎèçŽĐā■Ůéİcā■ŮçŇęäyşāĴijriijŇāĖŁèĵŇæ■ćæŁŔäyÄäyŁā■Ť
5. äyŁÉlć4æ■èçŽĐæŤřæ■ŏèćŇäijæèŁŞçzŽnameçŽĐèŁŞåĠžād'ĐçŘĚāZlīijŇārĖæIJĂçzŁçzŽĐçzŞæđIJètŇŇā

## 5.6 èĠlāŏŽžāZL'Item Loader

äĵŁçŤlćşzāŏŽžāZL'èŇ■æşŤriijŇäyŇÉlćæŸřäyÄäyŁä;Ňā■Ŕ

```

from scrapy.loader import ItemLoader
from scrapy.loader.processors import TakeFirst, MapCompose, Join

class ProductLoader(ItemLoader):

    default_output_processor = TakeFirst()

    name_in = MapCompose(unicode.title)
    name_out = Join()

    price_in = MapCompose(unicode.strip)

```

```
# ...
```

éĀŽēĠ\_ināŠŇ\_outāŘŮčijĀæĭēāōŽāzL'è;ŠāĚēāŠŇè;ŠāĠžād'ĎčŘĚāŽĭijŇāzūāyŤēĤŸāŘfāzēāōŽāzL'ézŸ  
default\_input\_processorāŠŇItemLoader.default\_input\_processor.

## 5.7 āĬĬFieldāōŽāzL'äy■āčřæŸŎè;ŠāĚē/è;ŠāĠžād'ĎčŘĚāŽĭ

èĤŸæĬĬL'äyĭāĬĬræŮzāŘfāzēēĭđāyŸæŮzā;ĤčŽĎæūzāĤāè;ŠāĚē/è;ŠāĠžād'ĎčŘĚāŽĭijŇéĈčāřsæŸřčŽt'æŎēāĬĬ

```
import scrapy
from scrapy.loader.processors import Join, MapCompose, TakeFirst
from w3lib.html import remove_tags

def filter_price(value):
    if value.isdigit():
        return value

class Product(scrapy.Item):
    name = scrapy.Field(
        input_processor=MapCompose(remove_tags),
        output_processor=Join(),
    )
    price = scrapy.Field(
        input_processor=MapCompose(remove_tags, filter_price),
        output_processor=TakeFirst(),
    )
```

äijŸāĚĤčžgĭijŽ

1. āĬĬItem Loaderäy■āōŽāzL'čŽĎfield\_ināŠŇfield\_out
2. FiledāĚĈæŤřæ■ō(input\_processorāŠŇoutput\_processorāĚšéŤōā■Ů)
3. Item Loaderäy■čŽĎžŸèōd'čŽĎ

TipsĭijŽāyĀēĤāēĭēēōšĭijŇārĚē;ŠāĚēād'ĎčŘĚāŽĭāōŽāzL'āĬĬItem Load-  
erčŽĎāōŽāzL'äy■field\_inĭijŇčŽĎāŘŮārĚē;ŠāĠžād'ĎčŘĚāŽĭāōŽāzL'āĬĬFieldāĚĈæŤřæ■ōäy■

## 5.8 Item LoaderäyĤäyŇæŮĠ

Item LoaderäyĤäyŇæŮĠcēcnāL'ĀæĬĬL'è;ŠāĚē/è;ŠāĠžād'ĎčŘĚāŽĭāĚsāžŇĭijŇārŤāēĈā;āæĬĬL'äyĀäyĭēgčādĤ

```
def parse_length(text, loader_context):
    unit = loader_context.get('unit', 'm')
    # ... length parsing code goes here ...
    return parsed_length
```

āĤĭāgŇāŇŮāŠŇāĤōæŤžāyĤäyŇæŮĠčŽĎāĬij

```
loader = ItemLoader(product)
loader.context['unit'] = 'cm'

loader = ItemLoader(product, unit='cm')

class ProductLoader(ItemLoader):
    length_out = MapCompose(parse_length, unit='cm')
```

## 5.9 aEĚçjőçŽDăd'DçŘĚăŽÍ

1. Identity  $\text{aTēazšāy} \blacksquare \text{āĀŽ}$
2. TakeFirst  $\text{ēfTāZdčnñāyÄäyléIdçl'žāĀijijNéĀŽāyyçTīā;IJè;ŠāGžad'DçŘEāZÍ}$
3. Join  $\text{ārEçzŠædIJèfdètũælěrijNézYēōd'ā;fçTīçl'žæāij}'$  ‘
4. Compose  $\text{ārEāG;æTřéS;æŌëetũæIēā;cæLŘçōaeAŞætĀijjNāžgçTşæIJĀāRŌçŽDè;ŠāGž}$
5. MapCompose  $\text{èùšāyLeİcçŽDComposęszāijijijjNāNžāLnāIJlāžŌāEĒēCīçzŠædIJāIJlāG;æTřāy}\blacksquare\text{ŽD:āōČçŽDè;ŠāEēāĀijæYrāRřēf}\blacksquare\text{āzčçŽDriijNéeŪāÉLārEçnñāyÄäylāG;æTřā;Iæñqā;IJçTīlāžŌæL'ĀæIJLā}$
6. SelectJmes  $\text{ā;fçTīljsonēũrfā;DæIēæšēērcāĀijāžūēfTāZdçzŠædIJ}$

## 6 ScrapyæTŻćÍŒ06- Item Pipeline

äĭŞayÄäyĭtemècñèIJYèZZçĹñàRŪăĹrăzNăRŌăijŽècñăRŚéĂAçzZĭtem  
 PipelineĭjŇçDŭăRŌăd'ŽăyĭtçzDăzŭăNLçĖĖgăzăžăRăd'DçRĖèŁZăyĭtemăĂĆ æŕRăyĭtem  
 PipelineçzDăzŭăĖŭăôđăŕăſăYŕăyÄäyĭlăôđçŌŕăzĖăyÄäyĭtçôĂă■TăŪzăşTçzDPythonçşzăĂĆăzŪăžnăŌěăRŪă  
 äĭŁçŦĭĭItem PipelineçzDăyŷçŦĭăIJzăZŦĭjZ

- æÿĖǺŒHTMLæȚræ■ō
- éĴNērAècñæŁšăRÚčŽDæȚræ■ō(æčĂæšëitemæỲřăŘęăÑĕăŘñæšŘăžZă■Ůæot)
- éĜ■ad'■æĂgæčĂæšě(çĐúăŘŎäyćaijČ)
- ārEæŁŠăRÚčŽDæȚræ■ōă■ỲăĆlăLræȚræ■ōăžSăy■

## 6.1 çijŨaẸŽèĜłauśçŽĎPipeline

aǒŽazL'äyÄäyIPythončsziijNčDǔāRŌāōđčŎřæŰzæsTprocess\_item(self, item,  
 spider) a■šāRřiiijNēfTāZđāyÄäyIā■ŰāĚyæLŰItemiiijNæLŰēÄĚæLZāGžDropItemāijCāyŷyācāijCēfZāy  
 æLŰēÄĚēfYāRřāzēāōđčŎřāyNéIcāGāāyIæŰzæsTiiijŽ

- `open_spider(self, spider)` èIJYèZZæL'ŞajjĂçŽDæUũæL'gèąN
- `close_spider(self, spider)` èIJYèZZăĖŞéU■æUũæL'gèąN

- `from_crawler(cls, crawler)` `¶`

## 6.2 Item Pipeline `¶`

### äzúæäijélÑérA `¶`

æLSázñéÅŽèfGäyÄäyläzúæäijélÑérAä;Nå■RæIëçIJNçIJNæÄÖæäüä;fçTÍ

```
from scrapy.exceptions import DropItem

class PricePipeline(object):

    vat_factor = 1.15

    def process_item(self, item, spider):
        if item['price']:
            if item['price_excludes_vat']:
                item['price'] = item['price'] * self.vat_factor
            return item
        else:
            raise DropItem("Missing price in %s" % item)
```

### årEitemåEŽåEëjsonæÜGäzú `¶`

äyÑéIcçŽDèfZäyI PipelineårEæL ÄæIJLçŽDitemåEŽåEëåLräyÄäylå■TçNñçŽDjsonæÜGäzúüijNäyÄèaNäyÄ

```
import json

class JsonWriterPipeline(object):

    def __init__(self):
        self.file = open('items.json', 'wb')

    def process_item(self, item, spider):
        line = json.dumps(dict(item)) + "\n"
        self.file.write(line)
        return item
```

### årEitemå■YåCíåLrMongoDBäy■ `¶`

èfZäyIä;Nå■Rä;fçTÍpymongoæIëæijTçd'zæÄÖæäüèöšitemåfIå■YåLrMongoDBäy■äÄC  
MongoDBçŽDåIJräIÄåŠNæTæ■öäzŞåR■åIJlé■ç;öäy■æNĞåöZüijNèfZäyIä;Nå■RäyžèeAæYräRŠä;ääsTç

```
import pymongo

class MongoPipeline(object):
```



```

collection_name = 'scrapy_items'

def __init__(self, mongo_uri, mongo_db):
    self.mongo_uri = mongo_uri
    self.mongo_db = mongo_db

@classmethod
def from_crawler(cls, crawler):
    return cls(
        mongo_uri=crawler.settings.get('MONGO_URI'),
        mongo_db=crawler.settings.get('MONGO_DATABASE', 'items')
    )

def open_spider(self, spider):
    self.client = pymongo.MongoClient(self.mongo_uri)
    self.db = self.client[self.mongo_db]

def close_spider(self, spider):
    self.client.close()

def process_item(self, item, spider):
    self.db[self.collection_name].insert(dict(item))
    return item

```

## éĢāđ'■èĤĢæzd'āZÍ

åAĢèøĳæĹSāznċŽĎiteméĢŇĬcċŽĎidā■ŮāĚÿæŸřāŤrāÿĂçŽĎiĳŇăĴĚæŸřæĹSāznċŽĎèIJŸèŽŽèĤŤāŽđāžĚā

```

from scrapy.exceptions import DropItem

class DuplicatesPipeline(object):

    def __init__(self):
        self.ids_seen = set()

    def process_item(self, item, spider):
        if item['id'] in self.ids_seen:
            raise DropItem("Duplicate item found: %s" % item)
        else:
            self.ids_seen.add(item['id'])
            return item

```

## 6.3 æĤĂæt'zäÿĂäÿItem PipelineçŽĎäZŮ

äĳāāĤĚéazāIJĬé■çĴæŮĢäZŮäÿ■ārĚäĳăĬJĂèçĂæĤĂæt'zçŽĎPipelineçŽĎäZŮæŮzāĹāāĹITEM\_PIPELINESā

ǎŔŎĭcçŽĐæŦŕǎ■Ůeǎłçđ'žǎŏČçŽĐæL'gèǎŇeǎžǎžŦijŇǎžŎǎ;ŎǎŦŕénŸæL'gèǎŇijŇèŇŇčǎžŦŦ0-1000

```
def parse_page1(self, response):
    item = MyItem()
    item['main_url'] = response.url
    request = scrapy.Request("http://www.example.com/some_page.html
↪",
                             callback=self.parse_page2)
    request.meta['item'] = item
    return request

def parse_page2(self, response):
    item = response.meta['item']
    item['other_url'] = response.url
    return item
```

## Requesta Rész

Scrapy a `Request` objektumot használja a weboldalak elérésére. A `Request` objektumot a `FormRequest` osztály segítségével lehet létrehozni, ha a weboldal egy formát tartalmaz, amelyet kitölteni kell.

```
return [FormRequest(url="http://www.example.com/post/action",
                    formdata={'name': 'John Doe', 'age': '27'},
                    callback=self.after_post)]
```

A `FormRequest` objektumot a `from_response` metódussal lehet létrehozni, ha a weboldal egy formát tartalmaz, amelyet kitölteni kell.

```
import scrapy

class LoginSpider(scrapy.Spider):
    name = 'example.com'
    start_urls = ['http://www.example.com/users/login.php']

    def parse(self, response):
        return scrapy.FormRequest.from_response(
            response,
            formdata={'username': 'john', 'password': 'secret'},
            callback=self.after_login
        )

    def after_login(self, response):
        # check login succeed before going on
        if "authentication failed" in response.body:
            self.logger.error("Login failed")
            return

        # continue scraping with authenticated session...
```

## Responsea Rész

A `Response` objektumot a `scrapy.http.Response` osztály segítségével lehet létrehozni. A `Response` objektumot a `TextResponse`, `HtmlResponse` vagy `XmlResponse` osztályokból lehet létrehozni.

- `TextResponse`: A `TextResponse` objektumot a `scrapy.http.TextResponse` osztály segítségével lehet létrehozni. A `TextResponse` objektumot a `body` attribútummal kell megadni.
- `HtmlResponse`: A `HtmlResponse` objektumot a `scrapy.http.HtmlResponse` osztály segítségével lehet létrehozni. A `HtmlResponse` objektumot a `body` attribútummal kell megadni.
- `XmlResponse`: A `XmlResponse` objektumot a `scrapy.http.XmlResponse` osztály segítségével lehet létrehozni. A `XmlResponse` objektumot a `body` attribútummal kell megadni.

## 7 Scrapy a 07- a rész

Scrapy a Python a 07- a rész. A `Scrapy` a 07- a rész. A `Scrapy` a 07- a rész. A `Scrapy` a 07- a rész.

äı£çŤĺăżşăıŁçőĂă■Ŧ

æCædIJåIJSpideréGÑeİcă;ǣǾTliijÑeĆčārsæZt'çõÄ■ǾZæijÑāZāäyzloggerårşæYřaoČčZǾDäyÄäyłaōđäçNā

## 7.1 aRŠéÅemail

åŖŠéĂÄäÿ■åŇĚåŘnéŽĎäzú

éĚ■ç;ó

```
MAIL_FROM = 'scrapy@localhost'
MAIL_HOST = 'localhost'
MAIL_PORT = 25
MAIL_USER = ""
MAIL_PASS = ""
MAIL_TLS = False
MAIL_SSL = False
```

## 7.2 aRÑäyÄäyIè£ŽçÍNè£RèaÑad'ŽäyISpider

```
import scrapy
from scrapy.crawler import CrawlerProcess

class MySpider1(scrapy.Spider):
    # Your first spider definition
    ...

class MySpider2(scrapy.Spider):
    # Your second spider definition
    ...

process = CrawlerProcess()
process.crawl(MySpider1)
process.crawl(MySpider2)
process.start() # the script will block here until all crawling_
                ↳ jobs are finished
```

## 7.3 aLÈäyČaijRçŁñèŽn

ScrapyāzūæšqæIJL'æRŘä;ŽaEĚç;øçŽDāLÈäyČaijRæLŞaRŪāLşèČ;iiijNäy■è£GæIJL'ä;Łād'ŽæÚzæşTāRřazē  
æČæđIJä;äæIJL'ä;Łād'ŽäyIspideriiijNæIJÄçõÅa■TçŽDæŪzaijRārsæYřaRřaŁād'ŽäyIscrapydāōđä;NiiijNç  
æČæđIJä;äæČşad'ŽäyIæIJžāZÍè£RèaÑaRÑäyÄäyIspideriiijNāRřazēařEurlāLÈçL'ĞaRŌäzd'çZæfRäyIæIJžāZ

```
http://somedomain.com/urls-to-crawl/spider1/part1.list
http://somedomain.com/urls-to-crawl/spider1/part2.list
http://somedomain.com/urls-to-crawl/spider1/part3.list
```

çDūāRŌè£RèaÑ3äyIscrapydāōđä;NiiijNāLÈāLñāRřaŁāōČazñiiijNāzūaijæÄŞpartāRCæTř

```
curl http://scrapy1.mycompany.com:6800/schedule.json -d_
↳project=myproject -d spider=spider1 -d part=1
curl http://scrapy2.mycompany.com:6800/schedule.json -d_
↳project=myproject -d spider=spider1 -d part=2
curl http://scrapy3.mycompany.com:6800/schedule.json -d_
↳project=myproject -d spider=spider1 -d part=3
```

## 7.4 éYšæ■céçnářAçŽDç■ŪçTě

äyÄäzŽç;ŚçnŽāōđçŌřazEäyÄäzŽç■ŪçTěæIèççAæ■çŁñèŽnæIèçŁñāRŪāōČazñçŽDç;ŚéatāÄČæIJL'çŽDæřT  
äyNéIcæYřařazēŌè£ZäzŽç;ŚçnŽçŽDäyÄäzŽæIJL'çTÍçŽDāzžèōōiiijŽ

- ä;£çTluser agentæšāāÄČzşāřsæYřæfRæñqāRŚéÄAçŽDæŪūāÄŽéŽRæIJžāZŌæšāy■éÄL'æNl'äy■äyÄ

- `çAæ■cCookieiijNæšRäzZç;ŠçñZäijŽéÄŽè£GCookieèrEāLñçTlæLûèžñäz;iiijNççAçTlāRŌä;fā;UæIJ`
  - `èøç;ôdownload_delayäyNè;;āzûè£šiiijNæTřā■Uèøç;ôäyž5çgšiiijNèuLād'gèuLāōL'āĚÍ`
  - `āēĆædIJæIJL'ārřēČ;çŽĐērīār;éGRā;£çTlGoogle cacheēŌuāRŪç;ŠéaŋiiijNèĀNäy■æYřçŽt'æŌēøē£éŪō`
  - `ä;£çTlāyĀäyĽè;øē;ñIPæšāiiijNā;NāēĆāĚ■et'ççŽĐTor projectæLŪèĀĚæYřāzYèt'ççŽĐProxyMesh`
  - `ä;£çTlād'gādNāLEāyČāijRāyNè;;āZlriijNè£ZæāuārseČ;āōNāĚléA£āĚ■ēcnārAžEriijNāRlēIJĀēAāĚšæ`
- `āēĆædIJē£ZāžZè£YæYřæUāæšTēA£āĚ■ēcnççAiiijNāRřāzēēĀČēŽSāTēäyŽæTřæNā`

## 8 ScrapyæTŽčlN08- æŮGäzúäyŌāŽçL'Ĝ

ScrapyäyžæLŠāznæRŘä;ZāžEāRřéĜ■çTlçŽĐitem pipelinesäyžæšRäyĽçL'zāōŽçŽĐItemāŌzäyNè;;æŮGäzūā  
éĀŽāyŷæIēērt'ä;āāijŽéĀL'æNl'ä;£çTlFiles PipelineæLŪImages PipelineāĀČ

è£Žāyđ'äyĽçŌæAšéČ;āōđçŌřāžEriijŽ

- `éA£āĚ■ēĜ■ād'■äyNè;;`
- `ārřāzēæNĜāōŽāyNè;;ārŌā£Iā■YçŽĐāIJræŮž(æŮGäzúçšçžçŽōā;Täy■,Amazon S3äy■)`

Images Pipelineäyžād'DçRĚāŽçL'ĜæRŘä;ZāžEēciād'ŮçŽĐāLšèČ;iiijŽ

- `ārĚæL'ĀæIJL'äyNè;;çŽĐāŽçL'ĜæāijāijRè;ñæ■cæLRæŽōéĀŽçŽĐJPGāžūā;£çTlRGBēciJèL'sæIāāijR`
- `çTšæLRçijl'çTēāŽ;`
- `æČĀæšēāŽçL'ĜçŽĐāō;āžēāšNénYāžēçāŌā£IāōČāžñæzāēŭšæIJĀārRçŽĐāržāryéŽRāLŪ`

çŌæAšāRŊNæŮūāijŽāIJlāEĚēČlā£Iā■YäyĀäyĽēcnērČāžēäyNè;;çŽĐURLāLŪēāĽiiijNçĐūāRŌārĚāNĚāRñçŽy

### 8.1 ä;£çTlFiles Pipeline

äyĀèLñæLŠāznäijŽæNl'çĚgäyNéIççŽĐæ■ēēld'æIēä;£çTlæŮGäzúçŌæAšiiijŽ

1. `āIJlæšRäyĽSpideräy■iiijNā;āçLñāRŪäyĀäyĽitemārŌōiiijNārĚçŽyāžTçŽĐæŮGäzūURLæT;āĚēfile_u`
2. `itemēcnē£TāŽđāzNāRŌāršāijŽè;ñāžđ'ççŽitem pipeline`
3. `ā;Šè£ŽāyĽitemāLřè;çFilesPipelineæŮūiiijNāIJlfile_urlsā■Ůæŏtäy■çŽĐURLāLŪēāĽiiijŽéĀŽ`
4. `ā;ŠæŮGäzūēcnäyNè;;āōNāžNāRŌōiiijNçžSædIJāijŽēcnētNāĀijçžŽāRēäyĀäyĽfilesā■ŮæŏtāĀČē£Žāy`

### 8.2 ä;£çTlImages Pipeline

ImagesPipelineèŭ\$FilesPipelineçŽĐä;£çTlāŭŏäy■ād'ŽiiijNäy■è£Gä;£çTlçŽĐā■ŮæŏtāR■äy■äyĀ  
ä;£çTlImagesPipelineçŽĐāē;ād'DæYřā;āārřāzēēĀŽè£ĜēĚ■ç;ŏæIēæRŘä;Zēciād'ŮçŽĐāLšèČ;iiijNærT  
ImagesPipelineä;£çTlPillowæIēçTšæLRçijl'çTēāŽ;āžēāRĽè;ñæ■cæLRæāĜāĜĚçŽĐJPEG/RGBæāijāij

- `get_media_requests(self, item, info)` efTāZdāyÄyhlRequestārZēsa

• `item_completed(self, results, item, info)`,  
ImagesPipeline

```
import scrapy
from scrapy.pipelines.images import ImagesPipeline
from scrapy.exceptions import DropItem

class MyImagesPipeline(ImagesPipeline):

    def get_media_requests(self, item, info):
        for image_url in item['image_urls']:
            yield scrapy.Request(image_url)

    def item_completed(self, results, item, info):
        image_paths = [x['path'] for ok, x in results if ok]
        if not image_paths:
            raise DropItem("Item contains no images")
        item['image_paths'] = image_paths
        return item
```

## 9 Scrapy

Scrapy

- Scrapy
- Scrapy Cloud

### 9.1 Scrapy

Scrapy

Scrapy

Scrapy

#### Scrapy

Scrapy

```
pip install scrapy
```

Scrapy



```
apt-get install scrapyd
```

## éĚ■ç;ő

éĚ■ç;őæŰĜäzúaIJřaĪÄrijŇäijŸăĚĹçžgäzŎä;ŎäĹřénŸ

- /etc/scrapyd/scrapyd.conf (Unix)
- /etc/scrapyd/conf.d/\* (in alphabetical order, Unix)
- scrapyd.conf
- ~/.scrapyd.conf (users home directory)

ăĚüä;ŞăŔĆæŢřăŔĆèĂĈscrapydéĚ■ç;ő

çőĂă■ŢçŽĎă;Ňă■Ř

```
[scrapyd]
eggs_dir      = eggs
logs_dir      = logs
items_dir     =
jobs_to_keep  = 5
dbs_dir       = dbs
max_proc      = 0
max_proc_per_cpu = 4
finished_to_keep = 100
poll_interval = 5
bind_address  = 0.0.0.0
http_port     = 6800
debug         = off
runner        = scrapyd.runner
application   = scrapyd.app.application
launcher      = scrapyd.launcher.Launcher
webroot       = scrapyd.website.Root

[services]
schedule.json      = scrapyd.webservice.Schedule
cancel.json        = scrapyd.webservice.Cancel
addversion.json    = scrapyd.webservice.AddVersion
listprojects.json  = scrapyd.webservice.ListProjects
listversions.json  = scrapyd.webservice.ListVersions
listspiders.json   = scrapyd.webservice.ListSpiders
delproject.json    = scrapyd.webservice.DeleteProject
delversion.json    = scrapyd.webservice.DeleteVersion
listjobs.json      = scrapyd.webservice.ListJobs
daemonstatus.json  = scrapyd.webservice.DaemonStatus
```

## éČlčjš

ä;ŁçŤİscrapyd-clientæIJÆŮžä;ŁřijŇ Scrapyd-clientæŸřscrapydçŽDäŸÄäŸŁăóċæŁŭčŇřijŇăóČæŖŖä;ŽăžEs  
éĂŽăŸŸăŖEä;ăçŽDăŭċĹŇéČlčjšăĹŖScrapydéIJĂċĚAăŸd'ăŸŁæ■ċĹd'ĲijŽ

1. âŖEăŭċĹŇæŁŸăŇĚæŁŖpythonċŽŇĲijŇă;ăĉIJĂċĚAăŖŁċċĚsetuptools

2. éĂŽċĚĜaddversion.jsonçžŁčŇŖăŖEċŸŠċŽĜċŽŇăŸŁăĲijăċĚĜŸScrapdæIJ■ăŁăăŽĹ

ă;ăăŖŖăžċăIJă;ăçŽDăŭċĹŇéĚ■ç;ŖăŮĜăžŮscrapy.cfgăŖŖăžăžŁScrapydçŽŖăăĜ

```
[deploy:example]
url = http://scrapyd.example.com/api/scrapyd
username = scrapy
password = secret
```

ăĹŮăĜžæŁŸæIJĹăŖŖçŤlčŽŖăăĜă;ŁçŤĹăŸ;ăžd'

```
scrapyd-deploy -l
```

ăĹŮăĜžæŸŖăŸlčŽŖăăĜăŸĹĹċæŁŸæIJĹăŖŖċĚŖĚăŇçŽDăŭċĹŇĲijŇæŁŸġċăŇăŸ;ăžd'

```
scrapyd-deploy -L example
```

ăĚĹċdăĹŖăŭċĹŇæăžçŽŖă;ŤĲijŇçDŭăŖŖŖă;ŁçŤĹăċăŸŇăŸ;ăžd'æĹċċĹčjšĲijŽ

```
scrapyd-deploy <target> -p <project>
```

ă;ăċĚŸăŖŖăžċăŖăžăžŁŸċŽŸċŖăŖçŽDtargetăŇŇprojectĲijŇçIJĂçŽDă;ăăŖŖăŇăċĹăŖŖăŤŸăžċçăĂ

```
[deploy]
url = http://scrapyd.example.com/api/scrapyd
username = scrapy
password = secret
project = yourproject
```

ċĚŖăăŭă;ăăŖŖçŽŤæŖŖăŖŮŮæŁŸġċăŇ

```
scrapyd-deploy
```

ăċĹċăđIJă;ăæIJĹăđŸăŸŁtargetĲijŇċĹăžăžăŖŖăžăă;ŁçŤĹăŸŇĹĹăŸ;ăžd'ăŖEprojectċĹčjšăĹŖăđŸăŸŁtargetæIJ■ă

```
scrapyd-deploy -a -p <project>
```

## 9.2 éČlčjšăĹŖScrapy Cloud

Scrapy CloudæŸŖăŸĂăŸŁæŁŸçŖăçŽDăžŖæIJ■ăŁăăŽĹĲijŇçŤŖScrapyeĹŇăŖŖŖăžDăĚŇăŖŸScrapinghubçŤŖæŁă  
ăŖČăĚ■ċŽd'ăžEăŖŁċċĚăŇŇçŽŖæŖĜăæIJ■ăŁăăŽĹčŽDĹIJĂċĚAĲijŇăžŮăŖŖă;ŽăžEĹĹdăŸŸç;ŖċġČçŽDUIăĹċçŖă

Scrapy CloudŠŃScrapydaŸřřĚijăőžčŽĎijŇä;ăăŔřăžěăăžă■őéIJĂëęAăIJläyď'èĂĚăžŇăL'■ăĹGă■ćijŇéĚ■  
cfġijŇëüşscrapyd-deployëržăŔŮčŽĎăŸřăŸĂăăŮčŽĎăĂĆ

æIJL'ăĴŁăd'ŽæŮŭăĂŽæĹŚăžněIJĂèĕAăžŌăd'Žăylĕ;ŚĉnŽĴĹnăRŮæĹ'ĂéIJĂèĕAĉŽDăŦræ■ōiijŊărfŦăĕĈăĹŚ  
ăĚŭăōădăy■éIJĂèĕAŦiijŊăĹŚăžnăRfăžéěĂŽèĴĠĉzt'æĹd'ăyĂăylĕğĎăĹŽéĒĉ;ōeăĴăĹŮĕĂĖăyĂăylĕğĎăĹŽéĒ  
èĕAĕĴŽăăŭăĂŽiijŊăĹŚăžnăŕsăy■èĈ;ăĒ■ă;ĴĉŦĴăĹ'■éĴĉŽDscrapy crawl  
testèĴŽĉğ■ăŚ;ăzd'ăžĒiijŊăĹŚăžněIJĂèĕAă;ĴĉŦĴiijŮĴĴĴŊĉŽDăŮăiŕĴĕĴĴăŕăŊScrapy  
spideriijŊăŔĈĕĂĈăōŸăŮăăŮĴĴăăĴ

aRraZealL'cTlscrapyaeRRaZcZcZDaeayafCAPIeAZefGcijUclNaeUzaijRaRraLlscrapyijNazcaeZfaijaczscZDs  
 crawlRaRraLlaUzaijRaAC  
 ScrapyadDazZazOTwistedaijCaec;SczIJaeaEaduaSzcaAazNayLijNaZaa'd'ajaeIJAeeAaIJlTwisted  
 reactorGNeIcefReaNaAC  
 eeUaELa;aaRraZeaj;cTlscrapy.crawler.CrawlerProcessefZaylcszaIeefReaNa;aczZDspiderijNef  
 reactorijNazueC;eE;c;oa;aczZDaUeaUaSNshutdownad'DcReaZlaACaL'AaeIJl'cZDscrapyaS;azd'eC;aj;c;

```
from twisted.internet import reactor
import scrapy
from scrapy.crawler import CrawlerRunner
from scrapy.utils.log import configure_logging
```

```

class MySpider(scrapy.Spider):
    # Your spider definition
    ...

configure_logging({'LOG_FORMAT': '%(levelname)s: %(message)s'})
runner = CrawlerRunner()

d = runner.crawl(MySpider)
d.addBoth(lambda _: reactor.stop())
reactor.run() # the script will block here until the crawling is_
↳ finished

```

## 10.2 ăŖŇăŷĂèŁŻćÍŇèŁŖèąŇăđ'Žăŷłspider

ézŸèóđ'æČĚăĖġă;ŞăġăăŖŖăŋăăĽ'ğèąŇscrapy crawlăŚġăzd'æŮŭăġŽăĽŽăzzăŷĂăŷłæŮŖčŽĎĕŹćÍŇăĂĆ

```

import scrapy
from twisted.internet import reactor
from scrapy.crawler import CrawlerRunner
from scrapy.utils.log import configure_logging

class MySpider1(scrapy.Spider):
    # Your first spider definition
    ...

class MySpider2(scrapy.Spider):
    # Your second spider definition
    ...

configure_logging()
runner = CrawlerRunner()
runner.crawl(MySpider1)
runner.crawl(MySpider2)
d = runner.join()
d.addBoth(lambda _: reactor.stop())

reactor.run() # the script will block here until all crawling jobs_
↳ are finished

```

## 10.3 ăŏŽăžĽ'èğĎăĽŽèąĹ

ăĕġăžĖĹĹăġ;Şă■ăġġăġġŇăĲĽ'ăžĖăĽ'■ĹćčŽĎĎŽăĲŇăŖŖăĽăŝžčăĂġġŇăŖŖăŖăžĕăġġĂăğŇăĽŖăžŋčŽĎăĽăĕ  
ăĽŖăžŋčŽĎĖĲĂăŖćăŸŖĕŹăăŭčŽĎġġŇăžŮăŷđ'ăŷłăŷ■ăŖŇčŽĎč;ŚćŋŽćĽŇăŖŮăĽŖăžŋăĽ'ĂĖĲĂĕĕĂčŽĎă  
ĕĕŮăĖĽăĽŖăžŋĖĲĂăĕĕĂăŏŽăžĽ'èğĎăĽŽèąĹăŖŇăŮĠćŋăĕăĹġġŇăĂŽĕĕĠăĽăĂĂčŽĎăĽŽăžžĖĲŸĖŽŽćŝġġŇă

```
#!/usr/bin/env python
# -*- encoding: utf-8 -*-
"""
Topic: ǎŏžǎžL'æŧřæ■ŏǎžšæĺǎǎđŇǎŏđǎ;š
Desc :
"""
import datetime

from sqlalchemy.engine.url import URL
from sqlalchemy.ext.declarative import declarative_base
from sqlalchemy import create_engine, Column, Integer, String, Text,
↳ DateTime
from coolscrapy.settings import DATABASE

Base = declarative_base()

class ArticleRule(Base):
    """èġĎǎĹžǎř■çġř
    __tablename__ = 'article_rule'

    id = Column(Integer, primary_key=True)
    # èġĎǎĹžǎř■çġř
    name = Column(String(30))
    # èĹŘèǎŇçŽĎǎššǎř■ǎĹŮèǎĹīījŇéǎŮǎŘŮéŽŤǎījǎ
    allow_domains = Column(String(100))
    # ǎījǎǎġŇŮŮŮǎĹŮèǎĹīījŇéǎŮǎŘŮéŽŤǎījǎ
    start_urls = Column(String(100))
    # ǎŷŇǎŷǎĕǎŧçŽĎxpath
    next_page = Column(String(100))
    # æŮĜçŇǎĕšçæŎĕæ■čǎĹžèǎĹèççǎījŘ(ǎ■Řǎŷš)
    allow_url = Column(String(200))
    # æŮĜçŇǎĕšçæŎĕæŘǎřŮǎŇžǎššxpath
    extract_from = Column(String(200))
    # æŮĜçŇǎæǎĜéçŷxpath
    title_xpath = Column(String(100))
    # æŮĜçŇǎǎĕĚǎŏžxpath
    body_xpath = Column(Text)
    # ǎŘšǎŷčæŮŮéŮŧ'xpath
    publish_time_xpath = Column(String(30))
    # æŮĜçŇǎæĹèæžŘ
    source_site = Column(String(30))
    # èġĎǎĹžæŷřǎřçŧšæŧĹ
    enable = Column(Integer)

class Article(Base):
    """æŮĜçŇǎçšž"""
    __tablename__ = 'articles'

    id = Column(Integer, primary_key=True)
```

```
url = Column(String(100))
title = Column(String(100))
body = Column(Text)
publish_time = Column(String(30))
source_site = Column(String(30))
```

## 10.4 ǎŏŽǎžL'æŰĞçńǎltem

èƒZäyĭā;ŁçőĀă■ȚăžEĭijŃæšăüzĂăžLéIJĂèęAęért' æYŎçŽĐ

```
import scrapy

class Article(scrapy.Item):
    title = scrapy.Field()
    url = scrapy.Field()
    body = scrapy.Field()
    publish_time = scrapy.Field()
    source_site = scrapy.Field()
```

## 10.5 ĀŽāŽL'ArticleSpider

æŌëäyNæIëæLŚäznârEåöZäZL'çLñåRŪæŪĞçnäçŽĐëIJYëZZijŇNefZäyIspideräijŽä;ŁçTlāyĀäyIRuleåöđäçN

```
from scrapy.utils import parse_text
from scrapy.spiders import CrawlSpider, Rule
from scrapy.linkextractors import LinkExtractor
from coolscrapy.items import Article

class ArticleSpider(CrawlSpider):
    name = "article"

    def __init__(self, rule):
        self.rule = rule
        self.name = rule.name
        self.allowed_domains = rule.allow_domains.split(",")
        self.start_urls = rule.start_urls.split(",")
        rule_list = []
        # æůžăŁă`äÿŇäÿĂëąť`çŽĎëğĎăĹŽ
        if rule.next_page:
            rule_list.append(Rule(LinkExtractor(restrict_
→xpaths=rule.next_page)))
            # æůžăŁăŁă;ăŘŮæŮĞńăéŞ;æŎěçŽĎëğĎăĹŽ
            rule_list.append(Rule(LinkExtractor(
                allow=[rule.allow_url],
                restrict_xpaths=[rule.extract_from]),
```

```

        callback='parse_item'))
    self.rules = tuple(rule_list)
    super(ArticleSpider, self).__init__()

    def parse_item(self, response):
        self.log('Hi, this is an article page! %s' % response.url)

        article = Article()
        article["url"] = response.url

        title = response.xpath(self.rule.title_xpath).extract()
        article["title"] = parse_text(title, self.rule.name, 'title')

        body = response.xpath(self.rule.body_xpath).extract()
        article["body"] = parse_text(body, self.rule.name, 'body')

        publish_time = response.xpath(self.rule.publish_time_xpath).
        extract()
        article["publish_time"] = parse_text(publish_time, self.
        rule.name, 'publish_time')

        article["source_site"] = self.rule.source_site

        return article

```

æŁAæşŁæĐRçŽĐæŸrstart\_urlsiiŸNrulesç■ŁéČĵāĹiāġŃāŃŮæŁŘäžEārŷzèšaçŽĐāśđæĂġiiŸŃéČĵčŤśaijāāĖēçŽĐ

## 10.6 çijŮåĖŽpipelineā■ŸāĆĹāĹræŤræ■ōāžŞäŸ■

æŁŚäžñèŁŸæŸřäĵčŤŤSQLAlchemyæĹæřĖæŮĠçñăăĹræŤræ■ōā■ŸāĆĹāĹræŤræ■ōāžŞäŸ■

```

@contextmanager
def session_scope(Session):
    """Provide a transactional scope around a series of operations."""
    session = Session()
    try:
        yield session
        session.commit()
    except:
        session.rollback()
        raise
    finally:
        session.close()

class ArticleDataBasePipeline(object):
    """äĹĹā■ŸæŮĠçñăăĹræŤræ■ōāžŞ"""

```

```
def __init__(self):
    engine = db_connect()
    create_news_table(engine)
    self.Session = sessionmaker(bind=engine)

def open_spider(self, spider):
    """This method is called when the spider is opened."""
    pass

def process_item(self, item, spider):
    a = Article(url=item["url"],
                title=item["title"].encode("utf-8"),
                publish_time=item["publish_time"].encode("utf-8"),
                body=item["body"].encode("utf-8"),
                source_site=item["source_site"].encode("utf-8"))
    with session_scope(self.Session) as session:
        session.add(a)

def close_spider(self, spider):
    pass
```

## 10.7 æŁóæŤzrun.pyåŘráŁìĎŽæIĴň

æŁśäznârĖäyŁēlċcŽĎrun.pyčl■ä;IĲăŃŃăŤză■șăŔrăŃŃăŤăLŭæŁśäznčŽĎæŮĠčnăcĹnĕŽnăŔrăĹlĕĎŽăIĲñ

```
import logging
from spiders.article_spider import ArticleSpider
from twisted.internet import reactor
from scrapy.crawler import CrawlerRunner
from scrapy.utils.project import get_project_settings
from scrapy.utils.log import configure_logging
from coolscrapy.models import db_connect
from coolscrapy.models import ArticleRule
from sqlalchemy.orm import sessionmaker

if __name__ == '__main__':
    settings = get_project_settings()
    configure_logging(settings)
    db = db_connect()
    Session = sessionmaker(bind=db)
    session = Session()
    rules = session.query(ArticleRule).filter(ArticleRule.enable == 1).all()
    session.close()
    runner = CrawlerRunner(settings)

    for rule in rules:
```



```

        # stop reactor when spider closes
        # runner.signals.connect(spider_closing, signal=signals.
→ spider_closed)
        runner.crawl(ArticleSpider, rule=rule)

    d = runner.join()
    d.addBoth(lambda _: reactor.stop())

    # blocks process so always keep as the last statement
    reactor.run()
    logging.info('all finished.')

```

OK iijNäy ÄäL GäRđăôŽăĂĆçŎřăIJăĽSăžnăRřăžěă; ĂArticleRuleăĹăy■ăĽăăĔăĔăĽRčŽ; äyĽă■ČăyĽç; ŠçñŽç; Ĺ; ŠçDŭă; äăŏNăĔĹăRřăžěăĂŽăyĂăyĽWebăĽ■çnřăĹăăŏNăĽRčžt' æĽd' ArticleRuleăĹçŽDăžžăĽăăĂČă; ŠçDŭăă; äăRřăžěăIJĹGitHubăyĽçIJNăĽřăIJăĽŮĞçŽDăŏNăTř' éąžçŽŏăžRčăĂăĂĆ

## 11 ScrapyæŤŽçĹN11- æĹăæNšçŽzăĽŤ

æIJĽæŮăăĂŽçĽnăRŮç; ŠçñŽçŽDăŮăăĂŽéIJĂăĔAçŽză; ŤiijNăIJĹScrapyăy■ăRřăžěăĂŽèĽGăĹăæNšçŽză; ŤăæAçšăŏđçŎřçŽză; ŤăřséIJĂăĔAăĹă■ŤăRŘăžd' iijNăĔĹăĂŽèĽGăŤRèĝĽăŽĹèŏĽéŮŏgithubçŽDçŽză; ŤéąŤéĹă//github.com/loginiijN çDŭăRŎă; ĽçŤĹăŤRèĝĽăŽĹèŤçŤăŭăăĔăĔăĹă; ŮăĽřçŽză; ŤăŮăĔĹĂăĔAăRŘăžd' äžĂă

æĽSèĽŽéGŇă; ĽçŤĹchromeăŤRèĝĽăŽĹçŽDèŤçŤăŭăăĔăĔăiijNŤ12æĽŠăijĂăRŎéĂĽæNĹNetworkiijNăžŭăŤEPlogăN; äyĽăĂĆ æĽSæŤĔăDŘè; ŠăĔăĔŽèŤŤçŽDçŤĹăĽăăR■ăŠNăŤEçăĂiijNă; ŮăĽŤăŏČăRŘăžd' çŽDformăĹă

▼ Form Data	view source	view URL encoded
commit: Sign in		
utf8: ✓		
authenticity_token: bvS09fDUV/Q0oXyfXAtgrbidJ+GnkuNQGEW		
login: yidao620c		
password: dddd		

åŒæšëçIJNhtmlæžŘčăĀiijŽăŔŚçŎřëăĹ■ŤëĜŇéíçæIJL'äyĭéŽŖëŮŔçŽĐauthenticity\_tokenăĀiijjŇ

```

<!-- </textarea> -->
<!-- ' ' -->
<form accept-charset="UTF-8" action="/session" data-form-nonce="2da80bbd
<div style="margin:0;padding:0;display:inline">
  <input name="utf8" type="hidden" value="✓">
  <input name="authenticity_token" type="hidden" value="bvS09fDUV/Q0oXyfXAtgrbidJ+GnkuNQGEW" />
  <input type="text" value="yidao620c" />
  <input type="password" value="ddd" />
</div>
<div class="auth-form-header">...</div>
<div id="js-flash-container">
</div>
<div class="auth-form-body">

```

## 11.1 éĜ■ăĒŽstart\_requestsæŮžæšŤ

ëçĀă;ŁçŦícookieiijŇçñňăyĂæ■ěă;ŮæL'ŞăijĂăŏČăŚĀiijŇézŸëŏd'scrapyă;ŁçŦíCookiesMiddlewareăy■

```
COOKIES_ENABLED = True
```

æĹŚăžňăĒĹëçĀăL'ŞăijĂçŽžă;ŤëąŧéíçiijŇëŎŭăŔŮauthenticity\_tokenăĀiijjŇëçŽéĜŇæĹŚéĜ■ăĒŽăž

```

# éĜ■ăĒŽăžĒçĹŇëžŇçşşçŽĐæŮžæšŤ, åŏđçŎřăžĒçĒĒăŏžăžL'èŕŭæśç,
→ēĹŖëăŇăĹŔăĹŚăŔŎăiijžëŕčçŦícallbackăžđëŕčăĜ;æŤŕ
def start_requests(self):
    return [Request("https://github.com/login",
                    meta={'cookiejar': 1}, callback=self.post_login)]

# FormRequestet
def post_login(self, response):
    # åĒĹăŎžăŇĒéžŖëŮŔçŽĐëăĹă■ŤăŔçăŤŕauthenticity_token
    authenticity_token = response.xpath(
        '//input[@name="authenticity_token"]/@value').extract_
    →first()
    logging.info('authenticity_token=' + authenticity_token)
    pass

```

start\_requestsæŮžæšŤæŇĜăŏžăžĒçăžđëŕčăĜ;æŤŕiijŇçŦíæĹëëŎŭăŔŮéžŖëŮŔëăĹă■ŤăĀiijauthenti

## 11.2 ä;£çŤÍFormRequest

ScrapyäyžæÍáæŇSætŔèġŁăŹÍiijŇæÍŚăžňăŏŽăžŁ' httpheader

```
# äyžăžEæÍáæŇSætŔèġŁăŹÍiijŇæÍŚăžňăŏŽăžŁ' httpheader
post_headers = {
    "Accept": "text/html,application/xhtml+xml,application/xml;q=0.
→9,image/webp,*/*;q=0.8",
    "Accept-Encoding": "gzip, deflate",
    "Accept-Language": "zh-CN,zh;q=0.8,en;q=0.6",
    "Cache-Control": "no-cache",
    "Connection": "keep-alive",
    "Content-Type": "application/x-www-form-urlencoded",
    "User-Agent": "Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/
→537.36 (KHTML, like Gecko) Chrome/49.0.2623.75 Safari/537.36",
    "Referer": "https://github.com/",
}
# ä;£çŤÍFormRequesetæÍáæŇSætŔèġŁăŹÍiijŇæÍŚăžňăŏŽăžŁ'
def post_login(self, response):
    # äĖÍŁăŐžæŇÉéŽŔèŮŔçŽĎèáÍá■ŤăŔĆăŤřauthenticity_token
    authenticity_token = response.xpath(
        '//input[@name="authenticity_token"]/@value').extract_
→first()
    logging.info('authenticity_token=' + authenticity_token)
    # FormRequeset.from_responseæŸŔScrapyæŔŔă;ŽçŽĎăŸĂăŸłăĢ;æŤř,
→çŤÍăžŐpostèáÍá■Ť
    # çŽžéŽEæÍŔăŁŚăŔŐ, äijžèŕČçŤÍafter_
→loginăžĎèŕČăĢ;æŤřiijŇăęĆăĎIJurleŭSRequestéáŧéÍçŽĎăŸĂăŸłăŕŝçIJAçŤěăŐŁ'
    return [FormRequest.from_response(response,
        url='https://github.com/
→session',
        meta={'cookiejar': response.
→meta['cookiejar']},
        headers=self.post_headers, #
→æŖÍăĎŔă■ăĎ'ĎçŽĎheaders
        formdata={
            'utf8': '✓',
            'login': 'yidao620c',
            'password': '*****',
            'authenticity_token':
→authenticity_token
        },
        callback=self.after_login,
        dont_filter=True
    )]

def after_login(self, response):
    pass
```

FormRequest.from\_response() æŰžæŖŤèŐŦ'ä;ăæŇĢăŏŽăŔŔăžĎ'çŽĎurliijŇèŕŭæŝĆăĎ't'èĖŸæIJL'for

→ăžăăÿžæŁŚăznăÿŁéİcǎőŽăzL'ăžĖRuleiijŃæL'ĂăžěăŔléIJĂèęAçóĂă■ŦçŽĐçŦŠæŁŔăĹİăğŇçĹăă

æŁŻēĠNæŁSéĀZēŁĠstart\_urlsāōZāzL'āžEāijĀāğNēāŁēŁcīijNçDūāRŌçTšæŁRRequestīijNāEūā;ŠçŁñāK

### 1.3 éGåEZ\_requests\_to\_follow

eIJL'äyléÜóécYǎŁŻaiiÄägŃăZřæL'řæŁŚâ;ŁäzĚåršæYřèŁZéGŃæŁŚâőŽäzL'čŽǾspiderčzǧæL'èĚłCrawlSpi

""éĞāEŽāŁāāĚěcookiejarcŽĎæŽt'æŮř""

## 1.4 éaṭélcād'DčŘEæÚzæşŢ

IJlègDāLZRuleéGÑélcæLŠăoŽzāL'āžEærRäyléS;æOěčŽDāZdërČăĜ;æTřparse\_pageiijNārśæYræIJĀcz

" "

→ "ēĚŽăŷłæŸřä;ĚÇŤĹLinkExtractorēĜłāŁłād'DĲŘĚéŞĲæŎěăžěăŘĹ`ăŷNăŸĂéąŧ`"  
→ ""

```
logging.info(u'-----æúŁæAråŁEåL' šćžŁ-----')
```

```
logging.info(response.url)
```

```
issue_title = response.xpath(
```

```
'//span[@class="js-issue-title"]/text()').extract_first()
```

```
logging.info(u'issue_titleïijŽ' + issue_title.encode('utf-8'))
```

## 11.5 11.5

```
#!/usr/bin/env python
# -*- encoding: utf-8 -*-
"""
Topic: çŻżâ;ŦçĹñèŽń
Desc : æĴæŦSçŻżâ;Ŧhttps://github.
→comăŘŌăřĒèĜĴăŭşçŻĎIssueăĒĴéĴĴĴĴăĜžæĴě
tipsiijžă;ŧçŦĴchromeërĴĴerŦpostèăĴă■ŦçŻĎŪŭăĂžăŦ;éĂĴ'Preserve_
→logăŠŦDisable cache
"""
import logging
import re
import sys
import scrapy
from scrapy.spiders import CrawlSpider, Rule
from scrapy.linkextractors import LinkExtractor
from scrapy.http import Request, FormRequest, HtmlResponse

logging.basicConfig(level=logging.INFO,
                    format='%(asctime)s %(filename)s[line:
→%(lineno)d] %(levelname)s %(message)s',
                    datefmt='%Y-%m-%d %H:%M:%S',
                    handlers=[logging.StreamHandler(sys.stdout)])

class GithubSpider(CrawlSpider):
    name = "github"
    allowed_domains = ["github.com"]
    start_urls = [
        'https://github.com/issues',
    ]
    rules = (
        # æŪĴæĴăĴŪèăĴ
        Rule(LinkExtractor(allow=('/issues/\d+',),
→ restrict_xpaths='//ul[starts-with(@class,
→ "table-list")]//li/div[2]/a[2]'),
            callback='parse_page'),
        # äŸŦăŸĂéăŦ, If callback is None follow defaults to True,
→otherwise it defaults to False
        Rule(LinkExtractor(restrict_xpaths='//a[@class="next_page"]
→')),
    )
    post_headers = {
        "Accept": "text/html,application/xhtml+xml,application/xml;
→q=0.9,image/webp,*/*;q=0.8",
        "Accept-Encoding": "gzip, deflate",
        "Accept-Language": "zh-CN,zh;q=0.8,en;q=0.6",
        "Cache-Control": "no-cache",
        "Connection": "keep-alive",
```

```

        "Content-Type": "application/x-www-form-urlencoded",
        "User-Agent": "Mozilla/5.0 (Windows NT 6.1; WOW64)
→AppleWebKit/537.36 (KHTML, like Gecko) Chrome/49.0.2623.75 Safari/
→537.36",
        "Referer": "https://github.com/",
    }

    # éĜ■āEĴžāžEçĹñèŽńçśżçŽĎæŮzæşŦ, āōđçŎřžEèĜĹāōŽžāzL'èrŭæśĆ,
→èĹŘèāNăĹĹăĹSăŔŎăijžèřČĤĹcallbackăžđèřČăĜ;æŦř
    def start_requests(self):
        return [Request("https://github.com/login",
                        meta={'cookiejar': 1}, callback=self.post_
→login)]

    # FormRequestSet
    def post_login(self, response):
        # āĖĹāŎžæNĹéŽĹŔçŽĎeāĹā■ŦăŔĆæŦřauthenticity_token
        authenticity_token = response.xpath(
            '//input[@name="authenticity_token"]/@value').extract_
→first()
        logging.info('authenticity_token=' + authenticity_token)
        # FormRequestSet.from_responseæŦřScrapyæŔŘăžçŽĎăÿĂăÿĹăĜ;æŦř,
→çĤĹăžŎpostèāĹā■Ŧ
        # çŽzéŽEæĹĹăĹSăŔŎ, äijžèřČĤĹafter_
→loginăžđèřČăĜ;æŦřiiijNăeĆăđIJurlèŭSRequestéaťéĪćçŽĎăÿĂăŭăŕśçIJAçŦěăŎĹ'
        return [FormRequest.from_response(response,
                        url='https://github.com/
→session',
                        meta={'cookiejar':
→response.meta['cookiejar']},
                        headers=self.post_headers,
→ # æşĹăĎŔă■d'ăđ'ĎçŽĎheaders
                        formdata={
                            'utf8': '✓',
                            'login': 'yidao620c',
                            'password': '*****',
                            'authenticity_token':
→authenticity_token
                        },
                        callback=self.after_login,
                        dont_filter=True
                    )]

    def after_login(self, response):
        for url in self.start_urls:
            #
→ăžăăÿžæĹSăžnăÿĹéĪćăōŽžāzL'ăžEŔRuleiijNăL'ĂăžěăŔĹéIJĂèeAçōĂă■ŦçŽĎçŦSæĹĹăĹĹăĜNçĹĹă
            yield Request(url, meta={'cookiejar': response.meta[
→'cookiejar']})

```

```

def parse_page(self, response):
    """
    """
    logging.info(u'-----æúŁæAřăŁEăL' ščžf-----
    --')
    logging.info(response.url)
    issue_title = response.xpath(
        '//span[@class="js-issue-title"]/text()').extract_
    first()
    logging.info(u'issue_titleïijž' + issue_title.encode('utf-8
    '))

def _requests_to_follow(self, response):
    """éĜ■ăEŁŻăŁăăĚăcookiejarčŽĐăŽt'æŮř"""
    if not isinstance(response, HtmlResponse):
        return
    seen = set()
    for n, rule in enumerate(self._rules):
        links = [l for l in rule.link_extractor.extract_
        links(response) if l not in seen]
        if links and rule.process_links:
            links = rule.process_links(links)
        for link in links:
            seen.add(link)
            r = Request(url=link.url, callback=self._response_
            downloaded)
            # äÿŇéİcèŁŻăŘăæŸřăŁŚéĜ■ăEŁŻčŽĐ
            r.meta.update(rule=n, link_text=link.text,
            cookiejar=response.meta['cookiejar'])
            yield rule.process_request(r)

```

ä;ăăŘřăžăăĬĬGitHubäÿŁçĬJŇăĹřăĬJňăŮĜčŽĐăŏŇăæŤt' éąžčŽŏăžŘčăAĭijŇěŁŸăĬJĹăŘăăđ' ŮäÿÄäÿŁěĜăĹăĹč

## 12 ScrapyæŤŽčĬŇ12- æŁŞăRŮăĹăăĀăçĭŞčŇŽ

ăĹ■éİăĹŚăžŇăžŇčž■čŽĐéČĭæŸřăŌžăæŁŞăRŮéİŽăăĀăçŽĐçĭŞčŇŽéąŁéİčĭijŇăžşăřşăæŸřěŕt' æĹŚăžŇăĹ'ŞăĭjĂăĭEăŸřăéČăžŁčŽĐăžŞăAŤçĭŞăđ'ġéČĹăĹčŽĐwebéąŁéİcéČĭæŸřăĹăăĀăçŽĐĭijŇčžŘăÿÿéĂŽčŽĐçĭŞčŇŽăĭŇăÿŇăĭĭæşŘăÿŁéŞĭæŌăăŮăĹăŁřčŽĐéąŁéİcéĜŇéİăăŘăăĬJĹăăĭjČă■ăăĹăăĭĭčŽĐăEĚăăŏžĭijŇěŁŹăăăăăE■ăĭŁçŤĹăžĭăĭŁçŤĹJavascriptăÿşăşŞăŇăđ'ĐčŘEçĭŞéąŁăŸřčġ■éİđăÿÿăÿÿăġăĀăçŽĐăĂăşŤĭijŇăéČăĭŤăđ'ĐčŘEăÿĂăÿĹăđăĚŹčřĜăŮĜčŇăăřEĭŕt'æŸŌăéČăĭŤăĬĬScrapyçĹŇăžŇăÿ■ăĭŁçŤĹscrapy-splashăĹăăđ'ĐčŘEéąŁéİčăÿ■ăĭŮJavascriptăĀĀ

## 12.1 scrapy-splash

scrapy-splash is a **Splash** web browser extension for **Scrapy**. It allows you to use **Scrapy** to scrape websites that require JavaScript. It is a **Scrapy** extension that uses **Splash** to render JavaScript and return the HTML of the page. It is a **Scrapy** extension that uses **Splash** to render JavaScript and return the HTML of the page. It is a **Scrapy** extension that uses **Splash** to render JavaScript and return the HTML of the page.

## 12.2 Docker

Docker is a platform for containerizing applications. It allows you to run applications in containers, which are isolated from the host operating system. It is a platform for containerizing applications. It allows you to run applications in containers, which are isolated from the host operating system. It is a platform for containerizing applications. It allows you to run applications in containers, which are isolated from the host operating system.

Install Docker on Ubuntu 12.04 LTS

Update the package index and install Docker

```
$ sudo apt-get update
$ sudo apt-get install linux-image-generic-lts-trusty
$ sudo reboot
```

Configure Docker

```
$ sudo apt-get install apt-transport-https ca-certificates
```

Generate GPG key

```
$ sudo apt-key adv --keyserver hkp://p80.pool.sks-keyservers.net:80
--recv-keys 58118E89F3A912897C070ADB76221572C52609D
```

Add Docker repository to sources.list

```
deb https://apt.dockerproject.org/repo ubuntu-precise main
```

Update APT

```
$ sudo apt-get update
$ sudo apt-get purge lxc-docker
$ apt-cache policy docker-engine
```

Install Docker

```
$ sudo apt-get install docker-engine
```

Start Docker

```
$ sudo service docker start
```

Test Docker

```
$ sudo docker run hello-world
```



äyŁÉİcèŁZæİaŞjāzd'āijZāyNèj;āyĀāyŁæŁNèrŁéŁIJāČŘāzŭāIJlāōzāZlāy■èŁRèaŃāōČiijŃāōČāijZæL'Şā■řāyĀ

## 12.3 āóŁ'èčĚSplash

æŃL'āŘŮéŁIJāČŘāyŃæİè

```
$ sudo docker pull scrapinghub/splash
```

āŘřāŁlāōzāZl

```
$ sudo docker run -p 5023:5023 -p 8050:8050 -p 8051:8051_
↳scrapinghub/splash
```

çŎřāIJlāŘřāzèčĚZèŁĜ0.0.0.0:8050(http),8051(https),5023 (telnet)æİèèōŁéŮōSplashžĚāĀČ

## 12.4 āóŁ'èčĚscrapy-splash

äjŁçŁlŁpīpāōŁ'èčĚ

```
$ pip install scrapy-splash
```

## 12.5 éĚ■çjőscrapy-splash

āIJlājāçZĎscrapyāŭèçlŃçZĎéĚ■çjőæŮĜāzŭsettings.pyāy■æŭzāŁă

```
SPLASH_URL = 'http://192.168.203.92:8050'
```

æŭzāŁăSplashäy■éŮŁ'āzŭiijŃèŁŸæŸřāIJlāsettings.pyāy■éĀZèŁĜDOWNLOADER\_MIDDLEWARESæŃĜā

```
DOWNLOADER_MIDDLEWARES = {
    'scrapy_splash.SplashCookiesMiddleware': 723,
    'scrapy_splash.SplashMiddleware': 725,
    'scrapy.downloadermiddlewares.httpcompression.
↳HttpCompressionMiddleware': 810,
}
```

ézŸèōđ'æČĚāĚjāyŃiijŃHttpProxyMiddlewareçZĎāijŸāĚŁçžgæŸř750iijŃèçAæŁŁāōČæŤ;āIJlāSplashäy■éŮ

èōŁçjőSplashèĜlāŭsçZĎāŎzéĜ■èŁĜæzd'āZl

```
DUPEFILTER_CLASS = 'scrapy_splash.SplashAwareDupeFilter'
```

āçČæđIJā;āäjŁçŁlāSplashçZĎHttpçijŞā■ŸiijŃéČčāzŁèŁŸèçAæŃĜāōZāyĀāyŁèĜlāōZāzL'çZĎçijŞā■ŸāŘŎāŘ
splashæŘŘājZāžĚāyĀāyŁscrapy.contrib.httpcache.
FilesystemCacheStorageçZĎā■Řçşz



```

        'dont_send_headers': True,      # optional, default is False
        'magic_response': False,       # optional, default is True
    }
})

```

Splash API `æYÖrijNä;ŁçTÍ SplashRequest æYřäYÄäyléÍđäyyä;ŁáL'çŽDăũăĚũălēăąăăĚrequest`  
`meta['splash']éĜNçŽDăTřă■ō`

- `meta['splash']['args']` āNĚăRňăžĚăRŚă;Ă SplashçŽDăRĆăTřăĂĆ
- `meta['splash']['endpoint']` æNĜăōŽăžĚ SplashăL'Ăă;ŁçTÍçŽDendpointiijNézYèōđ' æYřrender.html
- `meta['splash']['splash_url']` ěĚçŽŮăžĚĚsettings.pyăŮĜăžũăy■ēĚ■ç;ōçŽD Splash URL
- `meta['splash']['splash_headers']` ěĚRĚăNă;ăăćđăŁăăL'ŮăŁăŮăTřăRŚă;Ă SplashăIJ■ăŁăăŽÍçŽDHTTPă
- `meta['splash']['dont_send_headers']` âĚĆăđIJă;ăăy■ăĚšăijăăĂšheadersçŽŽ SplashiijNăřĚăōĚĚō;ç;ōăL'
- `meta['splash']['slot_policy']` ěōL'ă;ăĚĜăōŽăžL' SplashĥĚũăśĚçŽDăRňă■ēĚō;ç;ō
- `meta['splash']['dont_process_response']` â;Šă;ăĚō;ç;ōăL'RTrueăRÖrijN SplashMiddlewareăy■ăijŽ ResponseĥĚũăśĚăĚĆăžYèōđ' æYřäijŽĚĚTăŽD SplashResponseă■RçšăăŠăăžTăřTăĚĆ SplashTe
- `meta['splash']['magic_response']` éžYèōđ'ăyžTrueiijNS-plashăijŽĚĜăL'Ěō;ç;ōResponseçŽDăYăăžŽăśđăĂĝiijNăřTăĚĆresponse.headers,response.bodyç■L'

âĚĆăđIJă;ăăĚšăĂžĚĚĜ SplashălēăRĚăžđ' FormĥĚũăśĚiijNăRřăžăă;ŁçTÍ scrapy\_splash.  
 SplashFormRequestiijNăōĚĚũ\$ SplashRequestă;ŁçTÍăYřäYĂăũçŽDăĂĆ

## Responses

ăržăžŮăy■ăRňçŽD SplashĥĚũăśĚiijN scrapy-splashĥĚTăŽđăy■ăRňçŽD Responseă■Rçšă

- `SplashResponse` äžNĚĚŽăLŮăŠ■ăžTrijNăřTăĚĆărž/render.pngçŽDăŠ■ăžT
- `SplashTextResponse` æŮĜăĚIJňăŠ■ăžTrijNăřTăĚĆărž/render.htmlçŽDăŠ■ăžT
- `SplashJsonResponse` JSONăŠ■ăžTrijNăřTăĚĆărž/render.jsonăL'Ůă;ŁçTÍLuaĚĎŽăIJňçŽD/executeçŽDă

âĚĆăđIJă;ăăRĬăĚšă;ŁçTÍăĜăĜĚçŽD ResponseăržĚăqijNăřśĚō;ç;ōmeta['splash']['dont\_proces  
 ăL'ĂăĚIJĬĚĚŽăžŽ ResponseăijŽăĚLresponse.urlĚō;ç;ōăL'RăŮšăĜNĥĚũăśĚCURL(ăžšăřśăYřă;ăĚĚĂăyš  
 endpointçŽDURLăIJřăĬĂăĚĆăōđĚŽĚăIJřăĬĂăĂžĚĚĜresponse.real\_urlă;ŮăLř

## SessionçŽDăđ'DçRĚ

SplashăIJňĚžăYřăŮăçŁăăĂăçŽDijNĚĆăžLăyžăžĚăTřăNă scrapy-  
 splashçŽD sessionăĚĚăžçijŮăĚŽLuaĚĎŽăIJňiijNă;ŁçTÍ/execute

èĀÑæǻĞăĜĖçŽĐscrapy sessionȳŔĆæȚřāŔřäzëǻ|ŁçŦlSplashRequestǻřĖcookieæũżǻŁǻǻŦřǻŞǻL'■Splash  
cookiejarǻÿ■

æŌëäyNæIëæĹŚéĀŽèĤGäyÄäyġaôðéZĖçŽDä;Nā■RæIëæijTçd'žæĀŌæäüä;ĤçTīiijNæĹŚéĀĹ'æNĲçĹnāRŪāā  
 äžnāyIJç;SæĹŠāijĀéçŪéatçŽDæUūāĀŽāRġaijŽāEārijeĹĹēRĲJā■TāĹæj;āGžæIëriijNāĒŪāzŪāĒŪā;ŠéçŪéatāE  
 æĹŚçŌrāIJġārséĀŽèĤGçĹnāRŪēĤŽäyP'çNĲJā;āāŪIJæñç"èĤŽāZŽäyġa■ŪæIëèft'æŸŌäyNæŽôéĀŽçŽDScrapy  
 éçŪāĒĹæĹŚäzñāEŽäyġçōĀā■TçŽDæTNeftSpideriijNäy■ä;ĤçTīsplashiijŽ

çĐúãŔŎè£ŔèąŇçz\$æđIjijŽ

æŁŚæŁ; äy■āŁřćČčäyŁ”čŃIJä; āāŪIJæñć”èŁŻāZZäyŁā■Ū

æŌěäÿŊæİěæŁŚä;£çŦİsplashæİěçŁňåŦŮ

```

import scrapy
from scrapy_splash import SplashRequest

class JsSpider(scrapy.Spider):
    name = "jd"
    allowed_domains = ["jd.com"]
    start_urls = [
        "http://www.jd.com/"
    ]

    def start_requests(self):
        splash_args = {
            'wait': 0.5,
        }
        for url in self.start_urls:
            yield SplashRequest(url, self.parse_result, endpoint=
→ 'render.html',
                                args=splash_args)

    def parse_result(self, response):
        logging.info(u'-----
→ ä;ŁçŤlsplashçŁňăŔŮăžňăÿIJç;ŚéęŮéatâijĆă■ěăŁăè;ĭăĖĚăăőź-----')
        guessyou = response.xpath('//div[@id="guessyou"]/div[1]/h2/
→ text()).extract_first()
        logging.info(u"findĭijŽ%ŝ" % guessyou)
        logging.info(u'-----success-----')

```

èŁŔăĤçzŞăđIJĭjŽ

```

2016-04-18 14:42:51 js_spider.py[line:36] INFO -----
→ ä;ŁçŤlsplashçŁňăŔŮăžňăÿIJç;ŚéęŮéatâijĆă■ěăŁăè;ĭăĖĚăăőź-----
2016-04-18 14:42:51 js_spider.py[line:38] INFO _
→ findĭijŽçŇIJă;ăăŮIJăňć
2016-04-18 14:42:51 js_spider.py[line:39] INFO -----
→ success-----

```

ăŔŕăžčIJŇăĠçzŞăđIJéĠŇéĬăŭşçzŔăĤ;ăĤŔăžĖèŁăÿŦ'çŇIJă;ăăŮIJăňć"ĭijŇêŦ'ăŸŎăijĆă■ěăŁăè;ĭăĖĚăă

## 13 èĤŦçşzæĹŚ

- EmailĭijŽ yidao620@gmail.com
- ă■ŽăăćĭijŽ <https://www.xncoding.com/>
- GitHubĭijŽ <https://github.com/yidao620c>

ăĤŎăŁăæĤŦŇăÿĂăĤŦĭijŇêŦăă■ŽăÿzăŮĭăĤŦăŦăŦă



掃描上面的QR Code，加我WeChat。