

Tests & Résultats

) Scrapy ☐) parse »
) Dans l'exemple ci-dessous, la page est constituée d'une séquence de citations, dans des blocs « <div class=)quote> ». Dans chacun de ces blocs, la citation elle-même est placée dans un bloc « », et le nom de l'auteur dans un bloc « ».

Cette fonction « parse » retournera principalement (via un « **yield** ») deux)

- Soit des « données collectées »,)

Python (ici: {'text' : ..., 'author' : ...}). Chacun de ces objets est rajouté au résultat)

recupère typiquement sous forme) json ».

- Soit une ou plusieurs URLs, qui seront rajoutées au scheduler pour un traitement ultérieur, sauf si la page en question a déjà été traitée préalablement (ici : **yield** response.follow(...))

)) au choix via deux standards : CSS ou XPath.

) exemple ci-dessus, on parcourt (**for** quote **in** ...) div » de classe « quote » (response.css('div.quote')). Dans chacun de ces blocs, on extraira le texte contenu dans un « span » de classe « text » (via le standard CSS : quote.css('span.text::text')), et le contenu d'un bloc « small » contenu dans un bloc « span » (via le standard XPath : quote.xpath('span/small/text()')).

) précisera à chaque « yield » renvoyant une URL (plus précisément un objet de type scrapy.Request ») la fonction à utiliser lors du parsing (ici response.follow(..., self.parse)).

)) possible, ou au : utilisation de proxies, nombre de requêtes en parallèle, délai entre les requêtes, respect ou non des fichiers « robots.txt »

```
def parse(self, response):
    for quote in response.css('div.quote'):
        yield {
            'text': quote.css('span.text::text').get(),
            'author': quote.xpath('span/small/text()').get(),
        }

    next_page = response.css('li.next a::attr("href")').get()
    if next_page is not None:
        yield response.follow(next_page, self.parse)
```

Exemple extrait de : <https://docs.scrapy.org/>

Conditions d'utilisation & Budget

Scrapy est gratuit et Open Source.