

Compte rendu aymane hinane

Jsoup est une librairie Java qui permet d'analyser les pages HTML, en analysant leur code HTML.

Le but de Jsoup est d'extraire les informations d'un code HTML à partir des balises.

On installe Jsoup via le gestionnaire de packages Maven.

```
<dependency> <groupId>org.jsoup</groupId>  
<artifactId>jsoup</artifactId> <version>1.15.3</version>  
</dependency>
```

On intègre Jsoup via la librairie.

```
import org.jsoup.*;
```

On se connecte au site web qu'on veut cibler.

```
Document doc = Jsoup.connect("https://quotes.toscrape.com/").get();
```

```
Document doc =  
Jsoup.connect("https://quotes.toscrape.com/").get();
```

En cas d'échec de connexion, une exception de type [IOException](#) sera déclenchée.

1- première étape

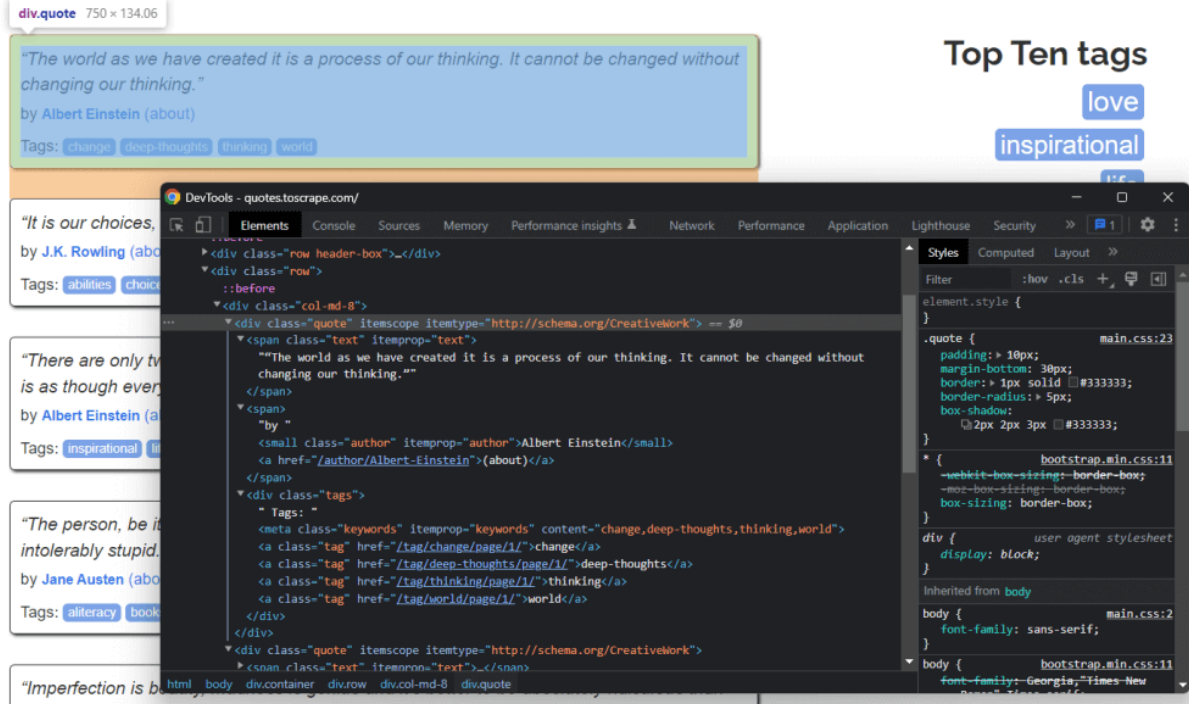
Il faut commencer par analyser le code de la page HTML qu'on veut cibler.

On utilise le développeur ou n'importe quel navigateur.

Et aussi de connaître l'élément HTML qui contient la donnée que l'on veut cibler.

Quotes to Scrape

Login



2-selectionner les elements avec jsoup

Selectionner toute les balises div

```
doc.getElementsByTag("div");
```

selectionner tous les elements html qui ont la class quote

```
Elements quotes = doc.getElementsByClass(".quote");
```

Pour cibler les balise enfant d'un parent

```
for (Element quoteElement: quoteElements) { Element text =
quoteElement.select(".text").first(); Element author =
quoteElement.select(".author").first(); Elements tags =
quoteElement.select(".tag"); }
```

```
//dans cette class nous allons stocker les donner recuperer
List<Quote> quotes = new ArrayList<>();
```

selectionner tous les elements html qui ont la class quote

```

Elements quoteElements = doc.select(".quote");

// iterer sur l'element parent
(Element quoteElement : quoteElements) {

// créer l'objet quote qui servira a stocker les donnees
Quote quote = new Quote();

// on va nettoyer la donner stocker dans la balise text on
// eliminent les quotes

String text = quoteElement.select(".text").first().text()
    .replace("\"", "")
    .replace("'", "");

String author = quoteElement.select(".author").first().text();

List<String> tags = new ArrayList<>();

for (Element tag : quoteElement.select(".tag")) {
tags.add(tag.text());
}

// storing the scraped data in the Quote object
quote.setText(text);
quote.setAuthor(author);
quote.setTags(String.join(", ", tags));
quotes.add(quote);

}

```