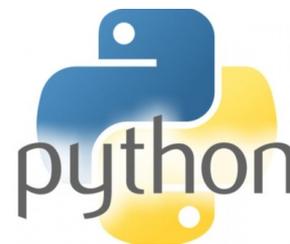


# Le module `browser_snt.py`

(sources : [python.org](http://python.org), [crummy.com](http://crummy.com) )



Ce module comporte des fonctions qui permettent d'aborder quelques notions de programmation au sujet des robots web (spiders, crawlers) et la recherche de mots clés dans une page (Term Frequency). Il est conçu à partir des modules **urllib** et **BeautifulSoup 4** ainsi que des fonctions de traitement des chaînes de caractères de Python. Il faut donc installer le module **BeautifulSoup 4** pour pouvoir l'utiliser.

**python -m pip install beautifulsoup4** dans une invite de commande Windows  
**python pip install beautifulsoup4** dans une invite de commande Linux

Enfin, il suffit de copier le fichier **module\_browser\_snt.py** dans le même dossier que les programmes Python qui l'utilisent.

## [module\\_browser\\_snt.find\\_tags\(page, tag\) :](#)

**page** est une chaîne de caractères. **tag** est une chaîne de caractères du type 'a' pour une balise <a>, 'img' pour une balise <img>, ... La fonction retourne une liste. Elle permet de obtenir toutes les balises du type considéré contenues dans la page.

## [module\\_browser\\_snt.get\\_attributes\(tags\\_list, attribute\) :](#)

**tags\_list** est une liste de balises du même type. **attribute** est une chaîne de caractères du type 'href', 'alt', ..., spécifiant l'attribut à considérer. La fonction retourne une liste. Elle permet de obtenir la valeur des attributs du type considéré, pour chaque balises.

## [module\\_browser\\_snt.get\\_page\(url\\_name\) :](#)

**url\_name** est une chaîne de caractères du type 'https://home.cern/fr', ... La fonction retourne une chaîne de caractères. Elle permet de obtenir le code html d'une page web d'adresse **url\_name**.

## [module\\_browser\\_snt.get\\_text\(page\) :](#)

**page** est une chaîne de caractères. La fonction retourne une chaîne de caractères. Elle permet de obtenir tout le texte contenu dans une page html.

## [module\\_browser\\_snt.keep\\_attributes\\_with\(attributes\\_list, term\) :](#)

**attributes\_list** est une liste de valeurs d'attributs. **term** est une chaîne de caractères. La fonction retourne une liste. Elle permet de ne conserver que les valeurs d'attributs contenant le terme considéré.

## [module\\_browser\\_snt.keep\\_attributes\\_without\(attributes\\_list, term\) :](#)

**attributes\_list** est une liste de valeurs d'attributs. **term** est une chaîne de caractères. La fonction retourne une liste. Elle permet de ne conserver que les valeurs d'attributs ne contenant pas le terme considéré.

#### module\_browser\_snt.lowercase(text) :

**text** est une chaîne de caractères. La fonction retourne une chaîne de caractères. Elle permet de passer tout un texte en minuscules.

#### module\_browser\_snt.remove\_duplicates\_words(words\_list) :

**words\_list** est une liste de chaînes de caractères, correspondant à une liste de mots. La fonction retourne une liste de mots (chaînes de caractères). Elle permet de supprimer les doublons dans une liste de mots.

#### module\_browser\_snt.remove\_small\_words(words\_list, taille) :

**words\_list** est une liste de chaînes de caractères, correspondant à une liste de mots. **taille** est un entier correspondant à la taille. La fonction retourne une liste de mots (chaînes de caractères). Elle permet de supprimer des mots dans une liste de mots.

#### module\_browser\_snt.remove\_symbols(text, old\_symbols) :

**text** est une chaîne de caractères. **old\_symbols** est une chaîne de caractères contenant les caractères à enlever. La fonction retourne une chaîne de caractères. Elle permet de supprimer des caractères d'un texte. Les caractères sont en fait remplacés par des espaces.

#### module\_browser\_snt.remove\_words(words\_list, old\_words) :

**words\_list** est une liste de chaînes de caractères, correspondant à une liste de mots. **old\_words** est une liste de chaînes de caractères contenant les mots à enlever. La fonction retourne une liste de mots (chaînes de caractères). Elle permet de supprimer des mots dans une liste de mots.

#### module\_browser\_snt.replace\_symbols(text, old\_symbols, new\_symbols) :

**text** est une chaîne de caractères. **old\_symbols** est une chaîne de caractères contenant les caractères à remplacer. **new\_symbols** est une chaîne de caractères contenant les nouveaux caractères, dans le même ordre. La fonction retourne une chaîne de caractères. Elle permet de remplacer des caractères d'un texte.

#### module\_browser\_snt.split\_into\_words(text) :

**text** est une chaîne de caractères. La fonction retourne une liste de mots (chaînes de caractères). Elle permet de découper un texte en liste de mots.

#### module\_browser\_snt.uppercase(text) :

**text** est une chaîne de caractères. La fonction retourne une chaîne de caractères. Elle permet de passer tout un texte en majuscules.

#### module\_browser\_snt.words\_frequencies(text, words\_list) :

**text** est une chaîne de caractères. **words\_list** est une liste de chaînes de caractères, correspondant à une liste de mots. La fonction retourne une liste de doublets (chaînes de caractères, entier). Elle permet d'obtenir le nombre de mots identiques contenus dans un texte. La liste est triée par ordre croissant d'occurrence des mots.