

# Mémo Expressions Régulières

Les **expressions régulières** sont un moyen puissant de recherche de chaîne de caractères proposé dans de nombreux instruments de linguistique de corpus. Elles permettent de rechercher non pas une séquence de caractères figée, telle que la séquence « pomme », mais un motif regroupant plusieurs formes, tel que le **motif** « pommes? » permettant de repérer les formes « pomme » et « pommes » en une seule requête.

La syntaxe des expressions régulières peut varier d'un langage de programmation ou d'un logiciel à l'autre. Lorsqu'elle n'est pas spécifiée dans la documentation du langage ou du logiciel, elle correspond à celle présentée ci-dessous.

Dans cette syntaxe, on distingue les **caractères littéraux** de **caractères spéciaux**, aussi appelés **méta-caractères**. Par caractère littéral, on entend le caractère tel qu'on souhaite le voir apparaître dans le texte. Par exemple, si l'on recherche des mots contenant la lettre « N », on dira que le motif doit comporter le caractère littéral « N ».

Code	Description
\	<b>Marque le caractère suivant comme caractère spécial ou littéral.</b> Par exemple, "n" correspond au caractère "n". "\n" correspond à un caractère de saut de ligne. La séquence "\\" correspond à "\", tandis que "\\(" correspond à "(".
^	Permet de spécifier la position <b>début</b> de la saisie. (bien souvent début de ligne)
\$	Permet de spécifier la position <b>fin</b> de la saisie. (bien souvent fin de ligne)
()	Permet de délimiter un <b>groupe de caractères</b> .
*	Permet de rechercher une chaîne contenant <b>zéro ou plusieurs fois</b> le caractère, ou groupe de caractères, qui précède. Ainsi, "zo*" permet de trouver "z" et "zoo", "(zo)*" permet de trouver "zozo".
+	Permet de rechercher une chaîne contenant <b>une ou plusieurs fois</b> le caractère, ou groupe de caractères, qui précède. Ainsi, "zo+" permet de trouver "zoo", mais pas "z".
?	Permet de rechercher une chaîne contenant <b>zéro ou une fois</b> le caractère, ou groupe de caractères, qui précède. Ainsi, "a?(ve)?" permet de trouver ve dans "lever", mais rien dans "lèvre".
.	Permet de rechercher <b>n'importe quel caractère unique</b> , sauf le caractère de nouvelle ligne.
x y	Permet de rechercher <b>soit x soit y</b> . Par exemple, "z foot" permet de trouver "z" et "foot". "(z f)oot?" permet de trouver "zoo" et "foot".
{n}	n est un nombre entier non négatif. Permet de rechercher un chaîne contenant <b>exactement n fois</b> le caractère, ou groupe de caractères, qui précède. Par exemple, "o{2}" ne permet pas de trouver "Bob", qui ne contient qu'un seul "o" mais correspond aux deux "o" de "foot" et "zoo".

Code	Description
<code>{n,}</code>	n est un entier non négatif. Permet de rechercher un chaîne contenant <b>au moins n fois</b> le caractère, ou groupe de caractères, qui précède. Par exemple, "o{2,}" ne correspond pas à "o" dans "Bob", mais à tous les "o" dans "fooooot". "o{1,}" équivaut à "o+" et "o{0,}" équivaut à "o*".
<code>{n,m}</code>	m et n sont des entiers non négatifs. Permet de rechercher un chaîne contenant <b>au moins n et au plus m fois</b> le caractère, ou groupe de caractères, qui précède. Par exemple, "fo{1,3}t" permet de trouver "fot", "foot" et "fooot". "o{0,1}" équivaut à "o?".
<code>[xyz]</code>	Jeu de caractères. Permet de rechercher un chaîne contenant <b>l'un des caractères indiqués</b> . Par exemple, "[abc]" correspond à "a" dans "plat".
<code>[^xyz]</code>	Jeu de caractères négatif. Permet de rechercher un chaîne contenant <b>tout caractère non indiqué</b> . Par exemple, "[^abc]" correspond à "p", "l" et "t" dans "plat".
<code>[a-z]</code>	Série de caractères. Permet de rechercher un chaîne contenant <b>un caractère de la série spécifiée</b> . Par exemple, "[a-z]" correspond à tout caractère alphabétique minuscule compris entre "a" et "z", "[0-9]" à tout caractère numérique compris entre 0 et 9.
<code>[^m-z]</code>	Série de caractères négative. Permet de rechercher une chaîne contenant <b>un caractère ne se trouvant pas dans la série spécifiée</b> . Par exemple, "[^m-z]" correspond à n'importe quel caractère sauf les caractères alphabétiques qui se trouvent entre "m" et "z" dans l'ordre alphabétique.
<code>\b</code>	Permet de spécifier la position <b>limite de mot</b> , autrement dit, la position entre un mot et une espace, ou un signe de ponctuation. Par exemple, "er\b" correspond à "er" dans "lever", mais pas à "er" dans "verbe".
<code>\B</code>	Permet de spécifier la position <b>autre que limite de mot</b> . Par exemple, "en*\B" correspond à "ent" dans "bien entendu", mais pas à "ent" dans "scient".
<code>\d</code>	Permet de rechercher <b>un caractère représentant un chiffre</b> . Équivaut à [0-9].
<code>\D</code>	Permet de rechercher <b>un caractère ne représentant pas un chiffre</b> . Équivaut à [^0-9].
<code>\n</code>	Correspond à <b>un caractère de saut de ligne</b> .
<code>\s</code>	Correspond à <b>tout espace blanc</b> , y compris l'espace, la tabulation, le saut de page, etc. Équivaut à "[\f\n\r\t\v]"
<code>\S</code>	Correspond à <b>tout caractère qui n'est pas un espace blanc</b> . Équivaut à "[^\f\n\r\t\v]"
<code>\t</code>	Correspond <b>au caractère de tabulation</b> .
<code>\w</code>	Correspond à <b>tout caractère permettant d'écrire un mot, y compris le trait de soulignement</b> . Équivaut à "[A-Za-z0-9_]"
<code>\W</code>	Correspond à <b>tout caractère autre que les caractères permettant d'écrire un mot</b> . Équivaut à "[^A-Za-z0-9_]"