

TECHNIQUE & PRATIQUE

Collection dirigée par Emmanuel Cornet et Alexandre Hérault

# Les expressions régulières

*par l'exemple*

VINCENT FOURMOND

Ancien élève de l'École Normale Supérieure

Quinze années d'expérience dans dix langages de programmation



## AVANT-PROPOS

Quiconque ayant déjà écrit un code « à la main » pour extraire des numéros de téléphone, des dates ou même des nombres d'une chaîne de caractères sait que ce n'est ni facile ni intéressant. En conséquence, les programmeurs ont tendance à éviter dans la mesure du possible de tirer des données d'une chaîne de caractères, et quand ils n'ont pas le choix, ils écrivent un code réduit à sa plus simple expression. Cela génère des programmes rigides et des utilisateurs frustrés.

Pourquoi réécrire des routines qui ont déjà été écrites mille fois, alors que, en utilisant des expressions régulières, on peut obtenir toutes les informations dont on a besoin en quelques lignes de code ?

Dans cet ouvrage, nous vous proposons une prise en main rapide et efficace des expressions régulières. Nous vous montrerons notamment, à travers de nombreux exemples, comment :

- construire une expression régulière ;
- effectuer des remplacements ;
- tirer parti de la souplesse des expressions régulières dans un programme.

Les expressions régulières sont des outils indispensables dès que l'on travaille sur des chaînes de caractères contenant des informations. Au moyen d'une notation compacte et efficace, elles permettent de détecter facilement des schémas dans une chaîne et d'en extraire ce dont on a besoin. Munis de ce livre, vous pourrez rapidement vous approprier leur puissance, et oublier les vieux codes longs et difficiles à modifier pour traiter les chaînes.

Ce livre procède en deux mouvements. Dans la première partie (chapitres 1 à 3), vous apprendrez progressivement à construire des expressions régulières, au travers de nombreux exemples et de 44 exercices, tous corrigés en détail (voir page 111). Dans la seconde (chapitres 4 et 5), vous découvrirez comment exploiter les expressions régulières dans un programme. Nous vous invitons à vous référer à l'index (page 124) tout au long de votre lecture.

Nous espérons que vous aurez autant de plaisir à lire cet ouvrage que nous en avons eu à l'écrire ; nous savons qu'il vous sera utile en pratique. Vos critiques comme vos éloges nous aideront à l'améliorer encore : vous pouvez en faire part à l'éditeur, à l'adresse

`contact@H-K.fr`

Si vous rencontrez ce que vous estimez être une erreur ou une imprécision gênante dans l'ouvrage, nous vous serions reconnaissant de nous en faire part également.

Bonne lecture et bon amusement !

Vincent Fourmond

---

# Table des matières

---

<b>Avant-propos</b>	<b>3</b>
<b>1 Qu'est-ce qu'une expression régulière ?</b>	<b>9</b>
Quelques conventions . . . . .	9
Exemple : est-ce une adresse électronique? . . .	10
Définition d'une expression régulière . . . . .	14
Quand utiliser des expressions régulières . . . . .	14
À qui s'adresse ce livre? . . . . .	15
<b>2 Débuter avec les expressions régulières</b>	<b>17</b>
2.1 La notion de correspondance . . . . .	17
2.2 Les caractères usuels et le joker . . . . .	19
a. Le joker : « . » . . . . .	19
b. Pouvez-vous préciser un peu? . . . . .	20
c. Un peu plus? . . . . .	21
d. Un point, pas n'importe quoi! . . . . .	22
2.3 Fromage Dessert? . . . . .	24
2.4 Jean-(Paul Julien) : le groupage . . . . .	25
2.5 Zéro, un, beaucoup . . . . .	27
a. « + » : au moins une fois . . . . .	27
b. Quels éléments au juste? . . . . .	27
c. Combien exactement? . . . . .	29
d. « ? » : au plus une fois . . . . .	30
e. « * » : zéro ou plus . . . . .	31
2.6 Le début et la fin . . . . .	32
Récapitulons! . . . . .	33

<b>3</b>	<b>Enrichir ses expressions régulières</b>	<b>35</b>
3.1	Remplacements . . . . .	35
3.2	B[ée]r[æ]ng[èe]re . . . . .	37
	a. Caractères spéciaux et classes . . . . .	38
	b. Et a-z ? . . . . .	40
3.3	La Disparition : [^e] . . . . .	41
3.4	Des raccourcis bien pratiques . . . . .	43
3.5	« ♪ Donna ♪ » {7,12} . . . . .	46
3.6	IgNOrER lA cAsSE . . . . .	48
3.7	« Je descends à la prochaine » . . . . .	49
3.8	Mettez-le-moi de côté . . . . .	51
	a. Numérotation des groupes . . . . .	52
	b. Remplacements . . . . .	52
	c. Extraction . . . . .	54
	d. Les captures et les quantificateurs . . . . .	55
3.9	Une impression de déjà-vu . . . . .	55
3.10	Les mots : le début et la fin . . . . .	58
	Récapitulons ! . . . . .	60
<b>4</b>	<b>Bien utiliser les expressions régulières</b>	<b>61</b>
4.1	Utilisations courantes . . . . .	61
	a. Les tests simples . . . . .	61
	b. Utilisation des groupes capturés . . . . .	62
	c. Toutes les correspondances . . . . .	64
	d. Et les remplacements ? . . . . .	64
4.2	N'en demandez pas trop ! . . . . .	65
	a. Fragmentez ! . . . . .	65
	b. Commentez si possible ! . . . . .	67
	c. Utilisez des structures appropriées . . . . .	68
	d. Une piste pour les imbrications . . . . .	69
4.3	Les lignes . . . . .	70
4.4	L'international . . . . .	71
	a. Les raccourcis . . . . .	72
	b. Les classes de type POSIX . . . . .	72
	c. Unicode . . . . .	73
	Récapitulons : les modificateurs . . . . .	73

<b>5</b>	<b>Syntaxe dans dix langages</b>	<b>75</b>
5.1	Perl .....	76
	a. Les opérations élémentaires.....	76
	b. Le code .....	78
5.2	JavaScript .....	79
	a. Les opérations élémentaires.....	79
	b. Le code .....	81
5.3	PHP .....	81
	a. Les opérations élémentaires.....	82
	b. Le code .....	84
5.4	Java .....	85
	a. Les opérations élémentaires.....	85
	b. Le code .....	87
5.5	C# .....	88
	a. Les opérations élémentaires.....	88
	b. Le code .....	90
5.6	Python .....	91
	a. Les opérations élémentaires.....	91
	b. Le code .....	93
5.7	C++/Qt .....	94
	a. Les opérations élémentaires.....	95
	b. Le code .....	96
5.8	Emacs Lisp .....	97
	a. Les opérations élémentaires.....	97
	b. Le code .....	98
5.9	C/POSIX .....	99
	a. Les opérations élémentaires.....	100
	b. Le code .....	102
5.10	Grep et Sed .....	104
	a. Les opérations élémentaires.....	105
	b. Le code .....	106
	Tableau récapitulatif .....	108

## Annexes

<b>A Solutions des exercices</b>	<b>111</b>
A.1 Exercices du chapitre 2.....	111
A.2 Exercices du chapitre 3.....	116
<b>B Que faire quand ça ne marche pas ?</b>	<b>122</b>
B.1 Règles générales .....	122
B.2 Règles dépendant d'un langage .....	123
B.3 Et sinon ? .....	123
<b>Index</b>	<b>124</b>

# Qu'est-ce qu'une expression régulière ?

---

Les ordinateurs n'ont pas du tout la même conception des textes que nous : pour nous, un texte est un ensemble d'idées couchées sur papier. Nous nous en servons pour transmettre des sentiments, des émotions ou, plus prosaïquement, des informations. Un ordinateur a une vision beaucoup plus pragmatique : un texte, c'est une suite de lettres ou plus exactement une suite de nombres qui représentent chacun une lettre. Il ne connaît ni langage ni signification.

Les ordinateurs ne savent pas naturellement repérer des informations dans un texte : il est souvent laborieux et rarement intéressant d'extraire des données d'une chaîne de caractères « à la main ». De plus, dès que l'on change un tant soit peu le type de données auxquelles on s'intéresse (des dates, des numéros de téléphone, des comptes bancaires, etc.), il faut récrire un nouveau code, à grand renforts de « copier/coller », de fautes de frappe et d'oublis.

À quoi bon réinventer la roue ? Les *expressions régulières* fournissent un moyen simple et particulièrement élégant de repérer des schémas dans une chaîne de caractères, en extraire des données et y apporter des modifications.

## Quelques conventions

Avant de passer à un premier exemple, détaillons quelques conventions que nous avons adoptées dans ce livre.

Un texte entre "guillemets américains" représente une chaîne de caractères : c'est une entité qu'en tant que programmeur vous ne taperez quasiment jamais. De telles chaînes

sont le substrat des expressions régulières et seront, pour la plupart, connues uniquement au moment de l'exécution du programme.

En revanche, les textes entre « guillemets français » représentent des éléments que vous serez amenés à utiliser directement dans vos programmes.

Beaucoup de langages incorporent des fonctionnalités plus ou moins évoluées d'expressions régulières, parfois grâce à des bibliothèques externes ; chaque langage a sa manière de les écrire et, surtout, sa manière de les utiliser. Nous avons choisi de nous appuyer sur la syntaxe de Perl, car c'est le langage par excellence des expressions régulières : c'est là qu'elles sont les plus puissantes et aussi les plus faciles à utiliser. Nous verrons de plus que la syntaxe varie très peu d'un langage à un autre ; la majorité des exemples que nous proposons est directement utilisable dans les autres langages. Dans de nombreux langages on définit d'ailleurs les expressions régulières par rapport à celles de Perl. Ce choix n'aura donc presque pas d'influence sur la conception des expressions régulières.

Nous prendrons l'habitude de noter les expressions régulières entre barres obliques, par exemple /*expression*/, pour bien les démarquer du reste du texte.

Enfin, pour lever toute ambiguïté, nous matérialiserons les espaces par le symbole `␣`, aussi bien dans les expressions régulières que dans les chaînes de caractères.

### **Exemple : est-ce une adresse électronique ?**

Entrons dès maintenant dans le vif du sujet au moyen d'un exemple complet de construction d'une expression régulière. Nous introduirons très brièvement les symboles les plus courants, mais, ne vous inquiétez pas, nous reviendrons en détail sur chacun d'eux, de façon progressive, à partir du chapitre 2.

Certains sites Internet proposent des formulaires en ligne, à remplir directement dans une page web ; il est fréquent que

des utilisateurs y saisissent des données incomplètes ou inexactes suite à une erreur de manipulation ou à une faute de frappe. Pour alléger la charge des serveurs, il est donc important de savoir filtrer les informations qui y parviennent.

Supposons que vous travaillez sur un formulaire en ligne dans lequel l'utilisateur doit entrer une adresse électronique. Vous souhaitez faire rapidement le tri entre les adresses invalides, comme "a", "a@b", etc., et les adresses correctes, comme "contact@h-k.fr". Nous allons construire pas à pas une expression régulière qui permet de vérifier si une adresse électronique est valide. Notons dès maintenant que les expressions régulières ne peuvent pas faire la différence entre une adresse réelle, comme "contact@h-k.fr" et une adresse inexistante mais « syntaxiquement correcte », par exemple "bidule@biniou.edu". Pour cela, il faut recourir à d'autres mécanismes, comme l'envoi d'un courrier de confirmation.

Pour commencer, posons-nous la question : « À quoi ressemble une adresse électronique ? » La réponse la plus élémentaire que l'on puisse donner est : « C'est une chaîne de caractères qui contient une arobase (@). » On peut représenter de telles chaînes au moyen de l'expression régulière

/@/

Elle se lit : « Détecter toutes les chaînes de caractères qui contiennent un arobase. » Si l'adresse fournie par l'utilisateur ne satisfait pas à ce critère, nous pouvons le lui signaler par un message d'erreur.

Le test précédent permet de repérer les erreurs les plus grossières ; imposons maintenant qu'il y ait quelque chose avant et après l'arobase. « . » est un *joker* : il peut représenter n'importe quel caractère. « + » s'applique au caractère qui le précède et signifie « au moins une fois » (rappelons que nous reviendrons en détail sur tous ces symboles à partir du chapitre 2). On peut les assembler en « .+ », soit « une suite quelconque d'au moins un caractère » ; nous nous en servi-

rons pour représenter « quelque chose ». Nous pouvons alors combiner cet élément avec l'arobase pour former

$$/.+@./$$

qui représente « une chaîne d'au moins un caractère (.+) suivie d'un arobase, puis d'une autre chaîne d'au moins un caractère (.+). » Nous avons progressé : cette dernière expression règle le sort de "a@" et "@bcd" car il leur manque l'une ou l'autre des chaînes autour de l'arobase. En revanche, elle laisse encore passer "a@a" qui n'est manifestement pas plus correcte que les deux précédentes.

Pour affiner notre tri, remarquons que le membre de droite d'une adresse électronique se compose d'au moins deux parties séparées par un point, comme dans "h-k.fr". Nous savons déjà exprimer « n'importe quelle chaîne d'au moins un caractère », grâce à « .+ ». La première idée pour vérifier que la partie de droite est conforme serait d'employer « .+.+ ». Mais on sait déjà que « . » représente non pas un point, mais n'importe quel caractère. Il faut donc « protéger » le point pour conserver sa signification d'origine en l'associant au caractère antislash : « \. » correspond spécifiquement au caractère « . ». Ce que l'on est en train de construire peut donc se représenter au moyen de « .+\.+. ». Nous pouvons alors modifier l'expression régulière précédente en

$$/.+@.\.+. /$$

pour tenir compte de la forme attendue pour le membre de droite. Cette expression se lit : « N'importe quelle chaîne non vide de caractères (.+) suivie d'une arobase (@), puis de n'importe quelle chaîne (.+), d'un point (\.) et enfin d'une troisième chaîne (.+). »

Continuons notre progression : en réalité, on ne peut pas utiliser n'importe quel caractère autour de l'arobase. En effet, les adresses électroniques ne peuvent comporter qu'une seule arobase ; cette dernière doit donc être proscrite dans les parties de gauche et de droite. Dans ce but, nous pouvons utiliser

en lieu et place des jokers (les points) la construction «  $[\hat{\text{C}}]$  » qui signifie « n'importe quel caractère *sauf* une arobase ». En effectuant les remplacements dans l'expression régulière précédente, nous obtenons :

$$/[\hat{\text{C}}]+\text{C}[\hat{\text{C}}]+\backslash.[\hat{\text{C}}]+/$$

Elle se lit : « Une chaîne d'au moins un caractère ne comprenant pas d'arobase («  $[\hat{\text{C}}]+$  ») suivie d'une arobase («  $\text{C}$  »), puis d'une autre chaîne d'au moins un caractère sans arobase («  $[\hat{\text{C}}]+$  »), d'un point («  $\backslash.$  ») et d'une troisième chaîne sans arobase («  $[\hat{\text{C}}]+$  »). »

De même que l'expression  $/\text{C}/$  détecte les chaînes qui contiennent une arobase, toutes les expressions régulières que nous avons introduites jusqu'à présent permettent de vérifier que la chaîne fournie par l'utilisateur *contient* une adresse syntaxiquement valide ; s'il envoie "a@contact@h-k.fr@b", l'expression se contente de détecter la partie du milieu (ici, "contact@h-k.fr"). Notre test laisse passer cette chaîne, ce qui n'est pas satisfaisant. Il faut s'assurer que la chaîne *entière* passe le test. On peut dans ce but utiliser «  $\hat{\text{C}}$  » pour repérer le début d'une chaîne de caractères et «  $\text{\$}$  » pour repérer sa fin. En les plaçant de part et d'autre de l'expression précédente, nous obtenons :

$$/\hat{\text{C}}[\hat{\text{C}}]+\text{C}[\hat{\text{C}}]+\backslash.[\hat{\text{C}}]+\text{\$}/$$

Cette nouvelle expression permet de s'assurer que la chaîne entière correspond : le début et la fin de la chaîne à vérifier coïncident exactement avec ceux de notre expression. Elle se lit « Détecter une chaîne qui

- commence ( $\hat{\text{C}}$ )
- par une chaîne sans arobase ( $[\hat{\text{C}}]+$ ),
- suivie d'un arobase ( $\text{C}$ ),
- puis d'une autre chaîne sans arobase ( $[\hat{\text{C}}]+$ ),
- d'un point ( $\backslash.$ ),

- et d'une troisième chaîne sans arobase (`[^@]+`)
- allant jusqu'à la fin (`$`). »

Nous pourrions encore améliorer cette expression en prenant en compte d'autres facteurs, comme le nombre de caractères de la dernière chaîne, entre deux et quatre (`fr`, `com`, `info`, etc.), la présence éventuelle d'espaces de part et d'autre de l'adresse, de la liste exacte des caractères autorisés dans une adresse électronique, etc. Vous pourrez compléter cette expression après la lecture du chapitre 3.

### Définition d'une expression régulière

La démarche précédente nous a montré ce qu'est une expression régulière : il s'agit d'un mini-programme, écrit dans un langage rudimentaire très compact, qui permet de décrire efficacement toute une série de chaînes de caractères qui se ressemblent. Elles servent à formaliser suffisamment ces ressemblances pour qu'un ordinateur puisse les reconnaître.

On parle souvent de *regex*<sup>1</sup> ou *regex*, qui sont les abréviations de l'anglais *regular expression*, « expression régulière ». « Régulière » fait référence à la régularité<sup>2</sup> des chaînes que les expressions régulières permettent de reconnaître.

### Quand utiliser des expressions régulières

Les expressions régulières apparaissent naturellement dès que l'on souhaite appliquer des traitements automatiques à du texte. On s'en sert pour :

- vérifier si celui-ci est syntaxiquement correct, comme dans le cas des adresses électroniques ;

---

1. Prononcer « règle-expe ».

2. La régularité est une notion mathématique de la théorie des langages. C'est de cette notion qu'il est question dans « expression régulière ».

- en extraire des informations, par exemple pour filtrer le résultat d'une commande particulièrement bavarde dont seule une partie nous intéresse ;
- effectuer des remplacements, par exemple pour « nettoyer » une chaîne issue de l'interaction avec l'utilisateur avant de la passer à une fonction qui exige qu'une syntaxe rigide soit respectée.

Nous verrons aussi qu'une bonne expression régulière ne peut pas remplacer un programme bien conçu !

### **À qui s'adresse ce livre ?**

Ce livre s'adresse aux programmeurs qui ne connaissent pas ou peu les expressions régulières, et qui manquent d'outils pour traiter les entrées textuelles de leurs programmes.

Nous vous montrerons progressivement au cours des chapitres 2 et 3 comment exprimer des idées de plus en plus complexes en termes d'expressions régulières ; nous verrons ensuite au chapitre 4 comment les incorporer dans vos programmes et en tirer parti.

Le chapitre 4 s'appuie sur des codes en Perl. Que les lecteurs qui ne sont pas familiers avec ce langage ne s'en inquiètent pas : nous n'utiliserons pas les spécificités syntaxiques de Perl dans nos programmes. En effet, on peut écrire des codes en Perl dans une syntaxe « universelle », proche des autres langages de haut niveau. Ils seront compréhensibles par tous les programmeurs. Nous expliquerons en outre tout point de syntaxe inhabituel que nous ne pouvons pas contourner.

Les expressions régulières s'acquièrent par l'expérience ; ce n'est qu'en les manipulant que l'on comprend vraiment tout les bénéfices que l'on peut en tirer. C'est pourquoi nous nous appuyerons sur maints exemples que vous pourrez réutiliser tels quels ou adapter à votre guise. Nous vous proposerons également de nombreux exercices dont la correction se trouve

dans l'annexe A (page 111). Ils peuvent se résoudre de tête, et nous vous invitons à y réfléchir avant de regarder leurs solutions. Nous vous invitons aussi à consulter systématiquement ces dernières même si vous pensez avoir trouvé : un détail peut vous avoir échappé. De plus, nous y présenterons souvent des prolongements qui peuvent vous intéresser.

Enfin, le chapitre 5 expose dans le détail les manières d'utiliser les expressions régulières dans dix langages différents : Perl, Python, PHP, C#, Java, JavaScript, C, C++, Emacs Lisp et les scripts shell. Un exemple complet y est traité dans les dix langages ; vous pourrez adapter ces codes pour tester les nouveaux éléments au fur et à mesure de leur découverte. Ces programmes sont en téléchargement libre sur la page

<http://www.h-k.fr/liens/tp/regexprs.html>

À la fin de ce livre, vous vous demanderez comment vous avez pu vous passer des expressions régulières jusqu'à présent !

# Débuter avec les expressions régulières

---

Nous allons découvrir dans ce chapitre le principe de base des expressions régulières : la correspondance. Nous verrons ensuite plusieurs constituants élémentaires des expressions régulières qui servent à exprimer des idées relativement simples. Certains manques se feront ressentir ; ils seront presque tous comblés dans le chapitre suivant.

## 2.1 La notion de correspondance

Un *outil* est un « objet fabriqué, utilisé manuellement ou sur une machine pour réaliser une opération déterminée<sup>1</sup>. » En s'appuyant sur cette définition, il est aisé de vérifier qu'une *pioche* est un *outil* car c'est un objet fabriqué que l'on utilise manuellement pour réaliser l'opération « creuser un trou ». On peut dire que l'objet « pioche » *correspond* à la définition d'un « outil ».

Les expressions régulières sont aux chaînes de caractères ce que les mots génériques, comme « outil », sont à la langue française : elles désignent un ensemble plus ou moins grand de chaînes en donnant leur définition. Une chaîne de caractères fait partie de cet ensemble si elle vérifie point par point chacun des critères de la définition. Dans ce cas, on dit que la chaîne *correspond*<sup>2</sup> à l'expression régulière. Nous avons vu dans le premier chapitre que "contact@h-k.fr" correspond à /.+@.+/ parce que "contact" correspond à « .+ », "@" à « @ » et "h-k.fr" à « .+ ».

---

1. Larousse, dictionnaire encyclopédique en 4 volumes, édition 1998.

2. On utilise *match* en anglais.

# Enrichir ses expressions régulières

---

Nous verrons dans ce chapitre des éléments pour mieux cibler les correspondances de nos expressions régulières. Ainsi, à l'issue de ce chapitre, vous saurez tout sur la conception des expressions régulières, à l'exception de quelques constructions spécifiques que nous verrons dans le chapitre suivant.

## 3.1 Remplacements

Nous nous sommes limités dans le chapitre précédent à observer des chaînes de caractères et à y chercher des correspondances. Nous allons aussi dans ce chapitre apprendre à les modifier, à l'aide de la fonction de remplacement :

---

`s/a/b/` remplace la première correspondance de `/a/` dans une chaîne par `"b"`<sup>1</sup>.

---

`s/chat/chien/` permet donc de remplacer par `"chien"` la première occurrence de `"chat"` dans une chaîne. Si nous l'appliquons à `"Le chat vient, le chat s'en va"`, nous obtenons `"Le chien vient, le chat s'en va"`.

EXERCICE 20. Quelle transformation a-t-on effectué pour changer `"65535"` en `"65536"` (plusieurs possibilités) ?

Dans la plupart des cas, il est beaucoup plus utile de remplacer toutes les correspondances non coïncidentes<sup>2</sup> d'une expression régulière dans une chaîne. On peut le faire en une seule opération à l'aide du modificateur `g`.

---

1. Le `s` initial est un raccourci pour « substitution ».

2. Ce sont toutes les correspondances disjointes d'une chaîne (p. 18).

# Bien utiliser les expressions régulières dans un programme

---

Après avoir appris à construire des expressions régulières dans les chapitres précédents, le temps est venu de voir comment s'en servir. Nous allons vous montrer comment utiliser les tests dans vos programmes, comment extraire des données de chaînes de caractères ainsi que des utilisations courantes des remplacements. Nous vous donnerons ensuite quelques règles de « savoir programmer » avec les expressions régulières, en particulier ce qu'il faut éviter. Enfin, ce chapitre se terminera par deux sujets importants : la gestion des lignes et quelques pistes pour travailler dans l'international.

Tous les codes de ce chapitre sont écrits en Perl. Nous avons cependant porté une attention particulière à la syntaxe, de manière à la rendre claire et compréhensible pour tous, quitte à s'éloigner de la manière dont un « perliste » les aurait écrits<sup>1</sup>.

## 4.1 Utilisations courantes

### a. Les tests simples

On utilise naturellement les expressions régulières pour vérifier si une entrée de l'utilisateur ou une partie d'un fichier sont dans un format bien déterminé. Ainsi, dans le programme qui gère le formulaire en ligne de l'introduction, on pourrait trouver le code suivant :

```
if($adresse =~ /\s*[\^@]+@[^\^]+\.[^\^]+\s*$/)
  { utiliseAdresse($adresse); }
else {afficheErreur("Adresse incorrecte");}
```

---

1. Qu'ils nous en excusent par avance !

## Syntaxe dans dix langages

---

Au contraire du chapitre précédent, où nous avons utilisé Perl comme langage de référence, nous allons maintenant détailler quelques exemples sur une sélection de dix langages. Vous y apprendrez quelles bibliothèques il faut utiliser et quels objets ou fonctions implémentent les expressions régulières. Nous vous montrerons dans chacun de ces langages comment vérifier si une chaîne correspond à une expression régulière, comment réutiliser les groupes capturés, comment progresser dans une chaîne pour trouver toutes les correspondances non coïncidentes d'une expression et enfin comment effectuer des remplacements.

Pour illustrer ces notions, nous avons implémenté dans les dix langages un exemple complet dont voici la motivation : une association a décidé de saisir toutes les informations sur ses adhérents dans une base de données. Elle a confié cette tâche à un stagiaire américain qui l'a très bien accomplie, au détail près qu'il a tapé toutes les dates (comme les dates de naissance et d'adhésion) au format américain, soit « mois-jour-année ». Depuis, les secrétaires de l'association se trompent systématiquement : il est urgent de les convertir dans un format plus habituel pour des Français, comme « jour.mois.année ».

En tant qu'adhérent de cette association, vous avez décidé de prendre vous-même en charge la conversion. Vous commencez par exporter la base de données dans un format texte, comme :

```
"Fourmond: Sophie:07-28-2004:02-14-2005:etc."
```

À l'aide d'un petit programme utilisant des expressions régulières, il est facile de modifier ce fichier texte pour ensuite le réimporter dans la base de données.

## Solutions des exercices

---

### A.1 Exercices du chapitre 2

**Solution 1** (*énoncé p. 18*). Quand on cherche des correspondances de /aba/ dans "abababa", on trouve la première correspondance, soit "abababa". Ensuite, il faut chercher les suivantes à partir de la fin de la première; on ne trouve que "abababa". En résumé, les correspondances sont :

"abababa"

**Solution 2** (*énoncé p. 20*). Le mot que l'on recherche est constitué d'un "b" suivi d'un caractère inconnu, d'un "l", de deux autres caractères inconnus et finalement de "re". On peut donc le représenter au moyen de l'expression

/b.l..re/

**Solution 3** (*énoncé p. 21*). On cherche des mots qui commencent par deux minuscules, soit « [a-z] [a-z] », puis "re", puis une minuscule, « [a-z] ». En combinant cela, on obtient

/[a-z] [a-z]re[a-z]/

Notons que cette expression ne nous assure pas que les correspondances sont des mots entiers. Nous verrons plus loin plusieurs manières d'y remédier.

**Solution 4** (*énoncé p. 21*). Nous savons maintenant exprimer spécifiquement un chiffre, au moyen de « [0-9] ». En remplaçant chaque point dans /...-...-.../ par cette construction, on obtient :

/[0-9] [0-9] - [0-9] [0-9] - [0-9] [0-9] [0-9] [0-9]/

Cette expression élimine bien les correspondances contenant des lettres ou des espaces.

# Que faire quand ça ne marche pas ?

---

Construire une expression régulière n'est pas toujours une tâche facile, d'autant qu'il n'est pas aisé d'obtenir un diagnostic. Pour vous aider, nous avons rassemblé quelques pistes à essayer quand quelque chose ne fonctionne pas.

Dans un premier temps, écrivez un petit programme de test à la manière de ceux du chapitre 5 avec l'expression régulière fautive et un exemple de chaîne sur laquelle vous souhaitez l'utiliser.

Si cela ne permet pas de cerner le problème, reconstruisez l'expression régulière morceau par morceau, en vérifiant à chaque étape les correspondances de chaque élément à l'aide de groupes capturés.

## B.1 Règles générales

- Changez les « `□` » en « `\s` » si votre implémentation le supporte ; sinon, essayez « `[□\t\n]` » tel quel (sans doubler les antislashes).
- Essayez des quantificateurs non avides, surtout si vous utilisez « `.*` » ou « `.+` ».
- Si vous utilisez « `[a-z]` » dans des correspondances avec du texte accentué, préférez « `\w` » (voir page 71).
- Attention, « `\w` » inclut aussi les chiffres : "78" correspond à « `\w+` ». S'il vous faut vraiment des lettres, utilisez `/[^W\d_]/` ou bien des constructions de type POSIX ou Unicode si c'est possible (voir page 71).
- Vérifiez la numérotation de vos groupes capturés.

---

# Index

---

## Symboles

"chaîne" .....	10
« élément » .....	10
\$ .....	13, <b>32</b>
\$& .....	64
(ab)+ .....	27
(b c d) .....	25
* .....	31
*? .....	49
+ .....	11, <b>27</b>
\+ .....	28
+? .....	49
- .....	40
. .....	11, <b>19</b>
\. .....	12, <b>23</b>
.+ .....	11, <b>27</b>
\.+ .....	28
/expression/ .....	10
/ .....	23
? .....	30
[0-9] .....	20
[1b-d7P-S] .....	40
[[:alpha:]] .....	72
[[:class:]] .....	72
[^f-h2] .....	13, <b>41</b>
[A-Z] .....	20
[a-z] .....	20
[a-z]+ .....	28
[A-Za-z] .....	21
[cgz] .....	37

[f-h] .....	21
\ .....	12, <b>22</b>
□ .....	10
^ .....	13, <b>32</b> , <b>41</b>
{5,12} .....	46
{5,12}? .....	49
{5,} .....	47
{5} .....	47
.....	24

## Nombres

\$1 .....	52
\1 .....	55

## A

Alternance .....	24
Avidité .....	29

## B

\B .....	59
\b .....	59

## C

C# .....	88
C+ + /Qt .....	94
C/POSIX .....	99
Capture .....	51
Caractère d'échappement .....	23

Caractères spéciaux 19, 22  
 Casse des caractères ... 48  
 Chiffre ..... 20  
 Classes de caractères .. 37  
   Caractères spéciaux . 38  
   Négatives ..... 41  
   POSIX ..... 72  
   Raccourcis ..... 43  
 Conventions ..... 9  
 Correspondance .... 17–18  
   Non coïncidente . 18, 35  
   Notion de ..... 17

## D

\D ..... 44  
 \d ..... 43  
 Début  
   d'un mot ..... 59  
   d'une chaîne ..... 32  
   d'une ligne ..... 71

## E

Échapper ..... 23  
 Emacs ..... 97  
*Escape* .... voir Échapper  
 Étendues de caractères 40  
 Expression régulière 9, 14  
 Extraction ..... 54

## F

Factoriser ..... 25  
 Fin  
   d'un mot ..... 59  
   d'une chaîne ..... 32  
   d'une ligne ..... 71

Fromage | Dessert .. voir |

## G

g ..... 36, 64  
 Gourmand ... voir Avidité  
*Greedy* ..... voir Avidité  
 Grep ..... 104  
 Groupage ..... 25  
 Groupes capturés .... 51  
   Extraction ..... 54  
   Numérotation ..... 52  
   Références internes .. 55  
   Remplacements ..... 52  
   Utilisation ..... 62

## I

i ..... 48

## J

Java ..... 85  
 JavaScript ..... 79  
 Joker ..... 11, 19

## L

Lignes ..... 70  
 Limite de mot ..... 59

## M

m ..... 71  
 Majuscule ..... 20  
*Match* ..... voir  
   Correspondance  
 MediaBox ..... 44  
 Minuscule ..... 20  
 Modificateurs ..... 48

g ..... 36, **64**  
 i ..... 48  
 m ..... 71  
 s ..... 70  
 x ..... 65, **67**

## N

Négation ..... 13  
 Nettoyer ..... 64  
 Nombres premiers ..... 57  
 Non avidité ..... 49

## P

Parenthèses ..... 25  
   capturantes ..... 51  
 Perl ..... 10, **76**  
 PHP ..... 81  
 POSIX ..... 72, 99  
 Protéger ... voir Échapper  
 Python ..... 91

## Q

Qt ..... 94  
 Quantificateurs ..... 27–31  
   \*? ..... 49  
   \* ..... 31  
   +? ..... 49  
   + ..... 11, **27**  
   ? ..... 30  
   {5,12}? ..... 49  
   {5,12} ..... 46  
   {5,} ..... 47  
   {5} ..... 47  
 Non avides ..... 49

## R

Référence interne ..... 55  
 Regex ..... 14  
 Regexp ..... 14  
*Regular expression* .... 14  
 Remplacement 35, 52, 62,  
   64  
 Retour à la ligne ..... 70

## S

\S ..... 44  
 \s ..... 43  
 s ..... 70  
 s/a/b/ ..... 35  
 s/a/b/g ..... 36  
 Sed ..... 104  
 Substitution ..... voir  
   Remplacement

## T

Tableau associatif ..... 69  
 Tests simples ..... 61

## U

Unicode  

{classe} ..... 73

## W

\W ..... 44  
 \w ..... 43

## X

x ..... 65, **67**