

Informatique et Sciences du Numérique

Représentations structurées des données

Alexandre Termier / Cyril Labbé

2017-2018

- Différent type d'information
 - données / méta-données
 - Information de présentation
- Web
 - HTML
 - CSS
- Modèle de données : semi-structuré (XML)
- Modèle de données : relationnel
 - Algèbre relationnelle
 - Language SQL

De l'information et des caractères "cachés"

- Pour nous : texte en langage écrit (ici le français)
- Pour la machine : une **chaîne de caractères**
- Jeux de caractères contiennent des **caractères de contrôle** qui permettent le formatage du texte : `\n` : retour à la ligne, `\t` : tabulation, ...

Exemple (Ex: texte formaté)

Vocabulaire et énoncés

Un triangle rectangle est un triangle admettant un angle droit (c'est-à-dire de mesure 90°).

Les deux côtés adjacents sont appelés cathètes et le côté opposé est l'hypoténuse.

Théorème

La forme la plus connue du théorème de Pythagore est la suivante :

Théorème de Pythagore - Dans un triangle rectangle, le carré de la longueur de l'hypoténuse est égal à la somme des carrés des longueurs des côtés de l'angle droit.

En particulier, la longueur de l'hypoténuse est donc toujours supérieure à celle de chaque autre côté.

- Format de fichier pour stocker des lignes

Format d'un fichier texte

```
ligne 1EOL
ligne 2EOL
...
dernière ligneEOL
EOF
```

- Caractéristiques
 - EOL : **fin de ligne** (**E**nd **O**f **L**ine), `\r\n` DOS, `\n` Unix, `\r` vieux Macs
 - EOF : **fin de fichier** (**E**nd **O**f **F**ile)
 - Lisible par les humains
 - Standard tant que l'on utilise des caractères ASCII / UTF-8

ASCII CONTROL CODE CHART

b7 b6 b5 BITS	0 0 0 0		0 1 0 1		1 0 0 0		1 1 0 1	
	CONTROL		SYMBOLS NUMBERS		UPPER CASE		LOWER CASE	
b4 b3 b2 b1	0	16	32	48	64	80	96	112
0 0 0 0	0 NUL	16 DLE	32 SP	48 0	64 @	80 P	96 ' .	112 p
0 0 0 1	1 SOH	17 DC1	33 !	49 1	65 A	81 Q	97 a	113 q
0 0 1 0	2 STX	18 DC2	34 " .	50 2	66 B	82 R	98 b	114 r
0 0 1 1	3 ETX	19 DC3	35 #	51 3	67 C	83 S	99 c	115 s
0 1 0 0	4 EOT	20 DC4	36 \$	52 4	68 D	84 T	100 d	116 t
0 1 0 1	5 ENQ	21 NAK	37 %	53 5	69 E	85 U	101 e	117 u
0 1 1 0	6 ACK	22 SYN	38 &	54 6	70 F	86 V	102 f	118 v
0 1 1 1	7 BEL	23 ETB	39 ' .	55 7	71 G	87 W	103 g	119 w
1 0 0 0	8 BS	24 CAN	40 (56 8	72 H	88 X	104 h	120 x
1 0 0 1	9 HT	25 EM	41)	57 9	73 I	89 Y	105 i	121 y
1 0 1 0	10 LF	26 SUB	42 *	58 :	74 J	90 Z	106 j	122 z
1 0 1 1	11 VT	27 ESC	43 +	59 ;	75 K	91 [107 k	123 {
1 1 0 0	12 FF	28 FS	44 ,	60 <	76 L	92 \	108 l	124
1 1 0 1	13 CR	29 GS	45 -	61 =	77 M	93]	109 m	125 }
1 1 1 0	14 SO	30 RS	46 .	62 >	78 N	94 ^	110 n	126 ~
1 1 1 1	15 SI	31 US	47 /	63 ?	79 O	95 _	111 o	127 DEL

LEGEND:



Victor Eijkhout
Dept. of Comp. Sci.
University of Tennessee
Knoxville TN 37996, USA

Bits of code point	Bytes in sequence	Byte 1	Byte 2	Byte 3	Byte 4
7	1	0xxxxxxx			
11	2	110xxxxx	10xxxxxx		
16	3	1110xxxx	10xxxxxx	10xxxxxx	
21	4	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx

Métadonnées d'un document

Données fournissant des informations sur un ou plusieurs aspects du document, sans faire explicitement partie du contenu du document.

- Quelques métadonnées
 - Titre et auteur d'un document
 - Date de création du document
 - Données EXIF d'une photo
 - ...

- Utilisations type
 - Catalogage/indexation de grandes quantités de documents
 - Filtrage et recherche dans ces documents

- Utilisations type
 - Catalogage/indexation de grandes quantités de documents
 - Filtrage et recherche dans ces documents
- Notion ancienne !
 - Un livre a des métadonnées sur sa couverture
 - Catalogue d'une bibliothèque
 - Système Dewey

- Utilisations type
 - Catalogage/indexation de grandes quantités de documents
 - Filtrage et recherche dans ces documents
- Notion ancienne !
 - Un livre a des métadonnées sur sa couverture
 - Catalogue d'une bibliothèque
 - Système Dewey
- Difficiles à repérer dans un fichier texte

- Fichiers texte
- Points positifs
 - Faciles à stocker
 - (Assez) standards
 - Quelques informations de formatage
 - Faciles à lire par un humain
- Points négatifs
 - Formatage insuffisant pour des documents complexes (**gras**, *italique*, . . .)
 - Pas de métadonnées
 - Pas d'images
 - Informations redondantes : perte de place en stockage

- Origine d'HTML
 - CERN : beaucoup de chercheurs, beaucoup de systèmes différents
 - Communications par échange de documents et mails (via Internet)
- HTML : réponse à deux besoins cruciaux
 - Pouvoir faire **facilement** un formatage avancé, "comme sur papier"
 - cela inclut l'intégration d'images
 - ...tout en devant rester portable entre systèmes différents
 - → Tim Berners-Lee : faire un sous-dialecte de **SGML**
 - Pouvoir "relier" des documents les uns aux autres
 - **hyperliens**
 - → reprise d'idées anciennes en recherche documentaire !

- **Hypertext Markup Langage**

- **hypertext** : documents peuvent être reliés entre eux
- **markup** : texte est annoté, en HTML les annotations décrivent la présentation du document

- Forme des annotations : **balises**

Théorème de `<bold>Pythagore</bold>`

- Intégration des métadonnées dans le document

- balises `<title>`, `<author>`, `<meta>`,...

- Informations de structuration

- balises `<title>`, `<body>`, `<h1>`, `<h2>`, ``, ``...

- Informations de présentation

- balises `<bold>`, `<italic>`,...

- Liens

- balise `<a>`