Métadonnée Cours « Document et Web Sémantique »

Nicolas Delestre



Métadonnée - v1.5 1 / 23

Plan...

- Quelques constats
- 2 Vocabulaire : métadonnées, schéma, profil d'applications
- 3 Deux exemples de schéma de métadonnées
 - Dublin Core
 - Schema.org
- 4 Diffusion et utilisation des métadonnées
 - Diffusion au sein des pages HTML
 - Pour les pages
 - Pour une partie d'une page
 - Utilisation par les moteurs de recherche
- Conclusion



Métadonnée - v1.5 2 / 23

Quelques constats 1/3

Google's Algorithm Is Lying to You About Onions http://daringfireball.net/linked/2017/03/07/google-onions

« Not only does Google, the world's preeminent index of information, tell its users that caramelizing onions takes "about 5 minutes" - it pulls that information from an article whose entire point was to tell people exactly the opposite. A block of text from the Times that I had published as a quote, to illustrate how it was a lie, had been extracted by the algorithm as the authoritative truth on the subject. » (Tom Scocca)



Métadonnée - v1.5 3 / 23

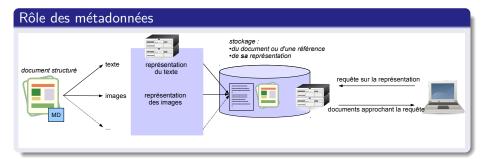
Quelques constats 2/3

Pourquoi?

- Les moteurs de recherche indexent des documents à destination des humains et non des machines
- Non prise en compte du contexte (à l'indexation et au moment de la recherche)
- Les moteurs de recherche utilisent les données (mots) pas l'information (relation entre les mots, les concepts)
- Les moteurs de recherche devraient demander des précisions concernant la recherche (interaction entre l'utilisateur et le moteur de recherche)



Quelques constats 3 / 3





Métadonnée - v1.5 5 / 23

Définitions

- « Une métadonnée (du grec meta "après" et du latin data "informations") est une donnée servant à définir ou décrire une autre donnée quel que soit son support (papier ou électronique). » (Wikipédia)
- « Les métadonnées sont, dans le cadre du Web sémantique, des données signifiantes qui permettent de faciliter l'accès au contenu informationnel d'une ressource informatique, une notice de contenu intégrée en quelque sorte (dans l'en-tête des documents HTML côté code source ou en tant que fichier XML autonome par exemple). » (Wikipédia)



Métadonnée (en pratique)

- C'est une donnée sur une donnée :
 - clairement définie
 - avec une arité (attention au LangString)
 - a un type :
 - type énuméré (vocabulaire fermé)
 - type prédéfini (entier, date, etc.) avec formalisme de représentation

Schéma de métadonnées

- Ensemble de métadonnées défini par un organisme, une entreprise, etc.
- Exemples de schéma de métadonnées
 - Dublin Core
 - Learning Object Metadata

Profil d'application

- Définition
 - « ...métadonnées issues d'un ou plusieurs schémas de métadonnées combinés afin d'améliorer et d'optimiser leur utilisation dans un cadre particulier » [HP00]
 - « ... est d'adapter et de combiner des schémas existants afin d'obtenir un nouveau schéma pour une application particulière tout en gardant l'intéropérabilité avec le ou les schémas de base » [DHSW02]
- Un profil d'application est un schéma de métadonnées
- Exemples de profil d'application :
 - CanCore
 - LOM-FR



Le Dublin Core 1 / 2

Qu'est ce?

- Géré par le (*Dublin Core Metadata Initiative*) : http://dublincore.org/index.shtml
- Schéma de métadonnées permettant de décrire des documents (numériques ou physiques)
 - 15 métadonnées forment le DCMES (Dublin Core Metadata Element Set). C'est aujourd'hui une norme (ISO 15836)
 - 55 (avec les 15 précédents) forment le DCMI Metadata Terms

Les métadonnées du DCMES (Wikipédia)

Elément	Commentaire
Title	Titre principal du document
Creator	Nom de la personne, de l'organisation ou du service à l'origine de la
	rédaction du document
Subject	Mots-clefs, phrases de résumé, ou codes de classement
Description	Résumé, table des matières, ou texte libre
Publisher	Nom de la personne, de l'organisation ou du service à l'origine de la
	publication du document

Le Dublin Core 2 / 2

Les métadonnées (suite)

Elément	Commentaire
Contributor	Nom d'une personne, d'une organisation ou d'un service qui contribue ou
	a contribué à l'élaboration du document. Chaque contributeur fait l'objet
	d'un élément Contributor séparé
Date	Date d'un évènement dans le cycle de vie du document
Type	Genre du contenu
Format	Type MIME, ou format physique du document
Identifier	Identificateur non ambigu : il est recommandé d'utiliser un système de
	référencement précis, afin que l'identifiant soit unique au sein du site, par
	exemple les URI ou les numéros ISBN
Source	Ressource dont dérive le document : le document peut découler en totalité
	ou en partie de la ressource en question. Il est recommandé d'utiliser une
	dénomination formelle des ressources, par exemple leur URI
Language	La langue du document
Relation	Lien avec d'autres ressources. De nombreux raffinements permettent
	d'établir des liens précis, par exemple de version, de chapitres, de stan-
	dard, etc.
Coverage	Couverture spatiale (point géographique, pays, régions, noms de lieux)
	ou temporelle
Rights	Droits de propriété intellectuelle, Copyright, droits de propriété divers

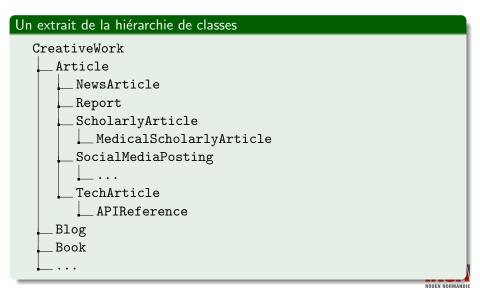
Schema.org 1/3

Qu'est-ce?

- http://schema.org
- Proposé par Google, Yahoo, Bing et Yandex
- Schéma de métadonnées généraliste (version 28.1 de novembre 2024)
 - environ 800 classes (organisées en hiérarchie héritage -)
 - environ 1450 propriétés
 - plusieurs centaines valeurs énumérées

Premier niveau de la hiérarchie Thing Action CreativeWork Event Thing Person Person Place Product

Schema.org 2 / 3



Métadonnée - v1.5

Schema.org 3 / 3

Propriétés

• Les propriétés sont attachées aux classes :

Thing additionalType, alternateName, description, image, name, potentialAction, sameAs, url

Person additionalName, address, affiliation, alumniOf, award,

. . .

- Les valeurs des propriétés peuvent de type simples (Datatype) ou de classes :
 - Datatype: Boolean, Date, DateTime, Number (deux sous types: Float et Integer), Text (un sous type URL), Time
 - ex : addtionalName : Text
 - ex : address : PostalAddress
- Une sous classes peut utiliser les propriétés des super classes



Métadonnées pour les pages (HTML5)

- Métadonnée par défaut : balise title
- Métadonnées additionnelles, utilisation de la balise meta avec les attributs:

name pour le nom de la métadonnée : *author*, *description*, *generator*, *keywords*

content pour la valeur

Exemple: http://www.w3.org/(2015)

Métadonnée - v1.5 14 / 23

JSON-LD

- Notation JSON avec des clés réservées (commençant par un @ pour le Linked Data)
- Format recommandé par Google

Exemple (http://schema.org/Movie)

```
<script type="application/ld+json">
  "@context": "http://schema.org",
  "@type": "Movie",
  "name": "Pirates of the Carribean: On Stranger Tides (2011)",
  "actor": [{"@tvpe": "Person", "name": "Johnny Depp"}, {"@tvpe": "Person", "name": "
       Penelope Cruz"}],
  "aggregateRating": {
    "@type": "AggregateRating", "bestRating": "10", "ratingCount":
"200", "ratingValue": "8", "reviewCount": "50"
  "author": [{"@type": "Person", "name": "Ted Elliott"}, {"@type": "Person", "name": "
       Terry Rossio"}].
  "description": "Jack Sparrow and Barbossa embark on a quest to find the elusive
       fountain of youth, only to discover that Blackbeard and his daughter are after
       it too.".
  "director": {"@type": "Person", "name": "Rob Marshall"}
  </script>
```

Métadonnée - v1.5

Microformat

Microformat (http://microformats.org)

- Spécifie une syntaxe et des schémas à utiliser
- Utilisation de l'attribut class pour spécifier la classe et le rôle de d'une valeur
- Plusieurs classes sont proposées (hCard, hCalendar, rel-licence, hRecipe, etc.). Par exemple marmiton.org utilise hRecipe.
- La version 2 propose une nomenclature de nommage (h-* pour les noms de classes, p-* pour les propriétés simples, etc.)

Exemple (http://microformats.org/wiki/microformats2-fr)

Métadonnée - v1.5 16 / 23

RDFa

- Mécanisme proposé par le w3c
- Définit la syntaxe, ne spécifie pas les schémas à utiliser
- Deux versions: RDFa v1.0 pour XHTML 1.0 et HTML 4.01 et RDFa v1.1 pour HTML5 (deux versions HTML+RDFaLite et HTML+RDFa)
- Propriétés : vocab, typeof, property, etc.

Exemple HTML+RDFaLite (http://www.w3.org/TR/html-rdfa/)

```
<!DOCTYPE html>
<html lang="en">
<head>
    <title>Example Document</title>
</head>
<body vocab="http://schema.org/">

        Welcome to my <a property="url" href="http://example.org/">blog</a>.

        </body>
</html>
```

Métadonnée - v1.5 17 / 23

Micro données (microdata)

- Syntaxe principalement associée à Schema.org
- Compatible avec HTML5, trois attributs principaux :
 - itemscope : Crée un élément et indique que les descendants de cette balise HTML contiennent des informations à son sujet.
 - itemtype : Un URL pointant vers un vocabulaire qui décrit l'élément et ses propriétés (http://schema.org/XX)
 - itemprop : Indique que la balise contient la valeur de la propriété indiquée.

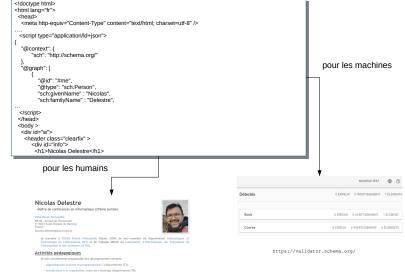
Exemple (http://schema.org/docs/gs.html)

```
<div itemscope itemtype="http://schema.org/Movie">
  <\1 itemprop="name"&g;Avatar</h1>
  <div itemprop="director" itemscope itemtype="http://schema.org/Person">
    Director:
        <span itemprop="name">James Cameron</span>
        (born <span itemprop="birthDate">August 16, 1954)</span>
  </div>
  <span itemprop="genre">Science fiction</span>
        <a href="../movies/avatar-theatrical-trailer.html" itemprop="trailer">Trailer</a>
</div>
</div>
```

Métadonnée - v1.5 18 / 23

Un exemple :

https://delestre.pages.insa-rouen.fr/siteweb/

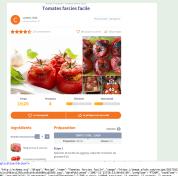


Rich Snippet 1 / 2

Qu'est ce?

- Informations ajoutées dans la page des résultats
- Peut être des images ou du texte

Utilisation des métadonnées de www.marmiton.org par Google



The state of the s



Recette Tomates farcies simples - Cuisine A7

Rich Snippet 2 / 2

Classes prises en compte par Google

- Product
- Recipe
- Review
- Event
- SoftwareApplication
- Video
- NewsArticle



Conclusion

 Les métadonnées permettent à des programmes de récupérer et d'interpréter des informations structurées



- C'est grâce à ces métadonnées que les moteurs de recherche peuvent ajouter des *Rich Snippets*
- Mais:
 - ces métadonnées ne sont pas vraiment liées
 - mélanger métadonnées et document peut ralentir les processus de chargement, d'extraction et d'interprétation des documents (le navigateur charge les métadonnées contenues dans les pages alors qu'il ne va pas en tirer parti...)
 - les représentations syntaxiques sont diverses
 - il n'y a pas de langage permettant d'interroger ces métadonnées (à l'image du SQL pour le monde relationnel)

ROUEN NORMANDIE

Métadonnée - v1.5 22 / 23

Références

[DHSW02] E. Duval, W. Hodgins, S. Sutton, and S.L. Weibel.

Metadata principles and practicalities.

 $\hbox{D-Lib Magazine, http://www.dlib.org/dlib/april02/weibel/~04weibel.html,~2002.}$

[HP00] R. Heery and M. Patel.

Application profiles: mixing and matching metadata schemas.

Ariadne issue 25, 2000.



Métadonnée - v1.5