

Moteurs de recherche multimédias

Introduction & Requêtes

1

Moteur de recherche

A quoi ça sert?

Trouver des informations

De quoi a-t-on besoin?

Données

Moteur

Requêtes

2

Les données

Quel genre de données?

Texte

Numérique

Date

URL

Images

Coordonnées géographiques

Sous quel format?

Idéalement un format générique le plus transparent possible

3

Les requêtes

Requêtes simples

document dont la date de création est le 28 oct. 2009

article dont le numéro d'identification est DK-2093

Requêtes complexes

document créé le 24 octobre 2009 plus ou moins 2 semaines dont le corps principal contient le verbe «créer» et le mot «bonjour», en ne tenant pas compte de la casse.

article valant entre 200 et 300 CHF et stocké dans un rayon de 2km d'une coordonnée géographique donnée.

4

Alternatives

Base de données relationnelle compatible SQL

Avantages: requêtes effectuées directement à la source, syntaxe connue et générique (SQL)

Désavantage: requêtes limitées, sans notions sémantiques et rapidement compliquées

Base de données XML (requêtes XPath)

Moteur de recherche

Avantages: spécialisé dans le traitement des requêtes complexes et utilisant des notions sémantiques

Désavantage: séparé de la source de données initiale

5

Solution proposée

Données (besoin d'un format générique)

XML

Requêtes (besoin d'un système flexible)

Lucene

6

Lucene

C'est une librairie Java de recherche FullText

Fonctionne avec des index composés de documents, stockés sur disque

Un document est un ensemble de champs

Les documents peuvent être cherchés par des requêtes faites sur les champs

Propose de l'analyse de texte flexible (filtres, tokenizers, analyzers)

7

XML + Lucene = Solr

Projet OpenSource

Développé en Java

=> nécessite une machine virtuelle Java.

Serveur de recherche basé sur Lucene

Interfaces XML et HTTP

Définition des types de données et des champs par un schéma flexible

8

Fonctionnalités avancées de Solr

- Recherche par facettes
- Un riche ensemble d'analyseurs sémantiques
- Soulignement des résultats (Highlighter)
- Correcteur orthographique (SpellChecker)
- Proposition de documents similaires (MoreLikeThis)

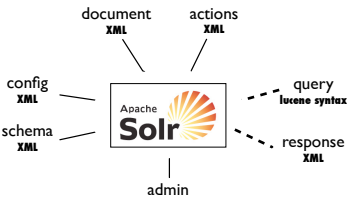
9

Différences SGBD - Solr

	SGBD	Solr
structure	multi-niveaux (database - table - record - field)	un seul niveau (document - field)
commit	instantané	"lent"
update	par champ	par document
recherche	limitée	très rapide
texte	sous-chaînes	termes (mots)

10

Architecture



11

admin

12

admin

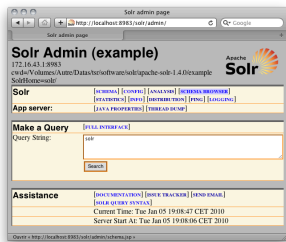
le serveur Solr met à disposition du développeur une interface d'administration, qui s'avère très vite indispensable!

elle permet d'obtenir des informations sur l'instance Solr ainsi que de la tester.

<http://localhost:8983/solr/admin/>

13

admin



14

admin

La zone du haut permet d'accéder, en consultation, aux différents fichiers de configuration et pages d'informations.

La zone centrale permet d'envoyer des requêtes au serveur Solr.

La zone du bas donne accès à des ressources documentaires

15

query
lucene syntax

response
XML

16

queries

Les requêtes sur Solr se font par l'intermédiaire de la syntaxe Lucene

Le but de la syntaxe lucene est de pouvoir exprimer une recherche sur des valeurs spécifiques de champs permettant de choisir et retourner les documents contenant ces champs.

17

queries

Une requête Lucene est une expression se composant de termes et d'opérateurs.

Un terme peut être de type unique (un seul mot) ou phrase (plusieurs mots entre guillemets).

Il peut débiter par un identifiant de champ, suivi de ':' afin d'indiquer le champ à rechercher. (sans cela, le champ par défaut est utilisé).

exemples: titre:dernier texte:"le journal"

18

queries

Les termes peuvent être organisés entre eux à l'aide de ces 3 opérateurs:

obligatoire	+	AND	&&
interdit	-	NOT	
optionnel		OR	

Attention, les opérateurs AND / && et OR / || n'ont pas la même signification que d'habitude. Leur utilisation équivaut à définir les 2 termes les entourant comme obligatoires / optionnels

19

queries

La définition de termes obligatoires et optionnels permet de faire des requêtes très précises et d'influencer le score des éléments trouvés.

bonjour +comment +va

va retourner tous les documents contenant 'comment' et 'va' mais pas forcément 'bonjour'. Ceux contenant aussi 'bonjour' auront par contre un meilleur score que les autres

bonjour comment AND va

aura le même résultat que l'expression précédente car l'opérateur AND a pour effet d'ajouter un "+" aux deux termes qui l'entourent.

20

queries

Pour faire des expressions plus complexes avec des groupes de termes optionnels par exemple, il faut utiliser des sous-expressions.

(+voici +exemple) (+sous +expression) -ignoré
retourne les documents contenant soit 'voici' et 'exemple', soit 'sous' et 'expression', soit les deux. Mais pas 'ignoré'.

(voici AND exemple) OR (sous AND expression) NOT ignoré
même signification que l'expression précédente.

voici AND exemple OR sous AND expression
retournera les documents contenant les 4 mots
car est équivalent à +voici +exemple +sous +expression.

21

queries

Wildcards

? – 1 caractère
* – 0-n caractères

(attention: un terme ne peut pas commencer par un de ces caractères)

Intervales

[start TO end] – valeur entre start et end - extrémités comprises
{start TO end} – valeur entre start et end - extrémités exclues

Boosting

Permet de mettre du poids sur certains termes, s'utilise en ajoutant le caractère '^' et une valeur après le terme à booster.

Exemple: title:docteur^10 text:docteur – met l'accent sur docteur dans le titre

22

queries

Fuzzy searches

permet de spécifier, pour un terme 'unique', le taux de similarité que doit avoir un champ pour être sélectionné
s'utilise en ajoutant le caractère '~' après le terme.

exemple: music~ – retournerait les champs contenant "music"

Proximity searches

permet de spécifier, pour un terme 'phrase', la distance maximale qui doit séparer les mots de la phrase pour qu'un champ soit sélectionné
s'utilise en ajoutant le caractère '~' après la phrase.

exemple: "iPhone 4S"~2 – retournerait les champs ayant le mot '4S' au maximum à une distance de 2 du mot 'iPhone'.

23

queries

Les valeurs de type 'date' doivent correspondre à un format particulier, selon la norme ISO-8601.

yyyy-mm-ddThh:mm:ssZ
exemple: 2010-01-01T10:45:55Z

Le mot-clef "NOW" peut aussi être utilisé.
Il correspond à l'instant présent.

24

queries

Des dates relatives peuvent aussi être calculées automatiquement grâce à une syntaxe mathématique.

Les unités suivantes sont reconnues:
MILLI - SECOND - MINUTE - HOUR - DAY - MONTH - YEAR

Il est possible de faire des arrondis avec le caractère '/'.

NOW/DAY -7DAYS
NOW/MONTH -12MONTHS
NOW+18HOURS-5SECONDS
2010-01-01T00:00:00Z + 10YEARS

25

queries

Pour faire une recherche sur tous les documents il faut utiliser la syntaxe:

:
(car '*' seul n'est pas autorisé par Lucene)

Pour rechercher tous les valeurs possibles d'un champ, il faut utiliser la syntaxe:
champ:[* TO *]

Solr étant très flexible, il est possible de définir des types de valeurs spécifiques et donc aussi des manières spécifiques de les rechercher (ex: géolocalisation)

26

queries

Exemples de requêtes:

content-type: video AND title:reportage
publication-date:
[2009-01-01T00:00:00Z TO 2009-12-31T00:00:00Z]
duration: [30 TO 90] AND location:(lausanne OR genève)
type: (maison NOT jumelée)
creation-date: {[* TO NOW-2MONTHS]}
birthday: [NOW/YEAR -10YEARS TO *]
name: [* TO *]
:

27

response

Solr utilise un format XML pour répondre aux requêtes.

Ce format est composé de plusieurs éléments, permettant de déduire les types de données:

<lst/> - liste nommée d'éléments
<arr/> - tableau d'éléments
<str/> - valeur de type String
<int/> <long/> <float/> <double/> - valeurs numériques
<bool/> - valeur booléenne
<date/> - valeur de type Date (format ISO-8601)

28

response

[illegible]

29

paramètres

q – la requête

fq – le filtre de requête

start – le numéro du premier document à retourner

rows – le nombre de documents à retourner

sort – paramètre de tri (ex: city.code asc, city.name desc)

fl – la liste des champs à afficher (ex: score)

wt – le format de sortie (xml, xslt, json, php, etc)

tr – si wt=xslt, indique le fichier xslt à utiliser

debugQuery – si 'true' ajoute des infos de déboggage

paramètres

http://localhost:8983/solr/select?

q=city.code:[2000 TO 2999]&

sort=city.code desc&start=100&rows=5&

```
fl=city.code,city.name,city.region.name
```

va trier les résultats par ordre décroissant des numéros postaux, depuis le résultat numéro 100 par page de 5 éléments. Et n'afficher que les champs numéro postal, nom de localité et acronyme de canton.

paramètres

[illegible]

32

paramètres

Il faut toujours faire bien attention à ce que les URLs soient encodées correctement lors des requêtes Solr.

En effet certains caractères, comme le "+", l'espace, le slash, sont traduits par d'autres dans le navigateur.

Le caractère "+" doit par exemple être inséré dans son format URL-Encodé: %2B

33

paramètres

Certains paramètres peuvent être définis via le formulaire de recherche avancée de Solr (lien "FULL INTERFACE"). Par exemple, q (requête), fl (liste des champs), start (début des résultats), rows (nombre de résultats par page).

MAIS sort, le paramètre pour le tri n'est pas accessible ainsi. Pour l'utiliser, il faut donc l'ajouter "à la main", dans l'URL du navigateur.

34

facettes

La recherche par facettes permet, à partir d'une requête, d'affiner la recherche en se basant sur les valeurs effectives des résultats trouvés.

Par exemple, en faisant une recherche sur toutes les localités contenant 'sur', on va demander à Solr de nous retourner la liste des cantons trouvés.

Il y a 3 types de facettes:
les facettes par champs, par dates et par requêtes.

35

facettes

Pour activer la recherche par facette, il faut ajouter un paramètre 'facet' avec la valeur 'true'.

Pour utiliser les facettes par champ, il faut indiquer à Solr sur quels champs on souhaite calculer les facettes, par le paramètre 'facet.field'. Solr va ensuite faire un total pour chaque valeur possible du champ spécifié.

`http://localhost:8983/solr/select?
q=sur&facet=true&facet.field=city.region.name`

36

facettes

```
- <response>
- <clt name="responseHeader">
- <clt name="status"><clt>
- <clt name="QTime"><clt>
- <clt name="params">
- <clt name="facet.query"><clt>
- <clt name="facet.field">city.region.name<clt>
- <clt name="q">date<clt>
- <clt name="facet.limit">5<clt>
- <clt>
- <clt>
+ <result name="response" numFound="46" start="0"><result>
- <clt name="facet_counts">
- <clt name="facet_queries">
- <clt name="facet_fields">
- <clt name="city.region.name">
- <clt name="Vaud">35<clt>
- <clt name="Fribourg">7<clt>
- <clt name="Gruyère">1<clt>
- <clt name="Jura">1<clt>
- <clt name="Neuchâtel">1<clt>
- <clt>
- <clt>
- <clt name="facet_dates">
- <clt>
</response>
```

37

facettes

Pour utiliser les facettes par requête, il faut indiquer à Solr des requêtes qu'il doit exécuter pour calculer les différents totaux de facettes.

http://localhost:8983/solr/select?
q=*&facet=true&facet.query=city.code:[* TO
2999]&facet.query=city.code:[3000 TO
5999]&facet.query=city.code:[6000 TO *]

38

facettes

```
- <response>
- <clt name="responseHeader">
- <clt name="status"><clt>
- <clt name="QTime"><clt>
- <clt name="params">
+ <clt name="facet.query"><clt>
- <clt name="facet">true<clt>
- <clt>
+ <result name="response" numFound="8014" start="0"><result>
- <clt name="facet_counts">
- <clt name="facet_queries">
- <clt name="city.code[* TO 2999]">1420<clt>
- <clt name="city.code[3000 TO 5999]">141<clt>
- <clt name="city.code[6000 TO *]">2183<clt>
- <clt>
- <clt name="facet_fields">
- <clt name="facet_dates">
- <clt>
</response>
```

39

facettes

Pour utiliser les facettes par date, il faut indiquer à Solr le champ de type 'date' à utiliser pour les facettes, via le paramètre 'facet.date', ainsi que les intervalles désirés.

Les intervalles se définissent à l'aide des paramètres 'facet.date.start', 'facet.date.end' et 'facet.date.gap'.

http://localhost:8983/solr/select?
q=*&facet=true&facet.date=timestamp&facet.date.start=NOW/
YEAR&facet.date.end=NOW&facet.date.gap=%2B1MONTH

40

facettes

D'autres paramètres des facettes permettent d'influencer la manière dont elles s'affichent.

facet.limit - permet de définir le nombre maximal de facettes affichées (par défaut 100).

facet.offset - permet de définir le numéro de la première facette à afficher (par défaut 0), équivalent à 'start' pour les requêtes.

facet.mincount - permet de définir le nombre de résultats minimum requis pour qu'une facette soit affichée (par défaut 0).

41

spellchecker

Le correcteur orthographique fonctionne en établissant un dictionnaire lors de l'indexation des documents, avec les valeurs des champs.

Il retourne ensuite une liste de termes semblables.

Pour l'activer, il faut ajouter un paramètre 'spellcheck' avec la valeur 'true'.

<http://localhost:8983/solr/select?q=rosé&spellcheck=true>

42

spellchecker

```
- <response>
- <lst name="responseHeader">
  <lst name="status">0</lst>
  <lst name="QTime">10</lst>
  <lst name="params">
    <str name="spellcheck">true</str>
    <str name="q">rosé</str>
  </lst>
</lst>
+ <result name="response" numFound="2" start="0"><result>
- <lst name="spellcheck">
- <lst name="suggestion">
  <lst name="sug">
    <lst name="numFound">3</lst>
    <lst name="startOffset">0</lst>
    <lst name="endOffset">4</lst>
    <arr name="suggestion">
      <str>Rosé</str>
      <str>Arna</str>
      <str>Gros</str>
    </arr>
  </lst>
</lst>
</result>
</response>
```

43

highlighting

La fonction de soulignement automatique permet de demander à Solr de nous indiquer où le mot-clé a été trouvé dans les champs retournés.

Cette fonction s'active en ajoutant un paramètre 'hl' avec la valeur 'true', puis en spécifiant quels champs doivent être soulignés via le paramètre 'hl.fl' (spécifier * pour souligner tous les champs).

http://localhost:8983/solr/select?q=le&hl=true&hl.fl=*

44

highlighting

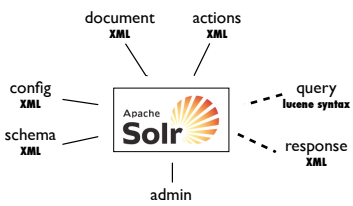
```
-<response>
  -<del name="responseHeader">
    <del name="status"><del>
    <del name="QTime"><del>
  -<del name="params">
    <del name="id"><del>
    <del name="id"><del>
    <del name="q"><del>
    <del name="rows"><del>
  </del>
+<result name="response" numFound="132" start="0"><result>
  -<del name="highlighting">
    -<del name="city.1040-13">
      -<del name="city.name">
        <del>Villars</del><del>Terre</del><del>
      </del>
    </del>
  -<del name="city.1041-13">
    -<del name="city.name">
      <del>Polze</del><del>Grand</del><del>
    </del>
  </del>
</del>
</response>
```

Moteurs de recherche multimédias

Schéma et Indexation

1

Architecture



2

configuration

3

configuration

La configuration de Solr passe par deux fichiers principaux:

`solrconfig.xml`

`schema.xml`

Ils se trouvent dans le répertoire "solr/conf"

Ils peuvent être consultés directement via l'interface d'administration de Solr.

4

configuration



solrconfig.xml

solrconfig

Ce fichier contient tous les paramètres de configuration de l'instance Solr.

Solr étant un projet hautement modulaire, car basé sur des interfaces Java très génériques, une grande quantité de modules peuvent être activés pour modifier son comportement.

C'est précisément dans ce fichier de configuration que les modules peuvent être activés et paramétrés.

solrconfig

Le principe général de configuration est que Solr est construit autour de composants (SearchComponent).

Ces composants sont des classes Java qui sont configurées dans le fichier solrconfig.xml

Chaque composant peut ensuite être utilisé dans un gestionnaire de requête (RequestHandler) qui est aussi configuré dans ce fichier.

Les pages 'admin' et 'select' sont 2 RequestHandler.

schema.xml

9

schema

Ce fichier contient la définition des champs et types de données utilisés dans l'index Solr.

Il est composé de trois parties:

définition des types
définition des champs
paramètres généraux

10

types

Solr met à disposition plusieurs types de données, par l'intermédiaire de classes java.

La plupart des types sont des types simples, ne nécessitant pas une logique de traitement particulière:

int, float, long, double, boolean, string,
binary, date, random

Ces types simples sont indexés tels-quels.

11

type text

En plus des types simple, Solr propose des champs de type 'text'. Ce type est assez complexe à gérer car il peut être configuré de manière très précise.

C'est grâce à ce type qu'il est possible de faire de la recherche fulltext avec Solr.

Sa configuration se base sur un principe de filtres.

12

type text

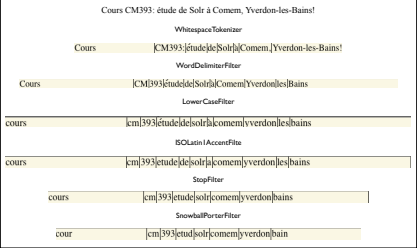
Lors de l'indexation, Solr aura pour objectif de générer une liste de termes distincts à partir d'un texte.

Pour générer ces termes distincts, il va faire passer le texte par une chaîne de filtres.

Il stockera ensuite dans l'index les termes résultant de l'application de cette chaîne de filtres.

13

type text



14

type text

Ensuite, lors de la recherche, Solr fera passer la requête de l'utilisateur par les mêmes filtres que ceux utilisés lors de l'indexation.

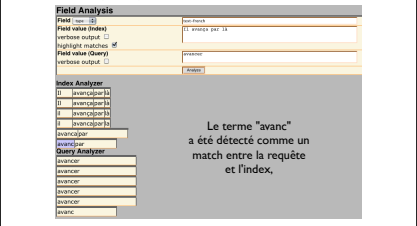
Tous les termes finaux similaires entre la requête et les termes indexés permettront à Solr de savoir quels documents correspondent à la requête.

L'enchaînement des filtres et leur configuration sont intégralement définis dans le schema.

15

type text

L'outil "Analysis" permet de tester les chaînes de filtres



16

type text

La plupart du temps, plusieurs types de champs textes vont être définis dans le schéma. Chacun d'eux permettant des recherches tenant compte de critères différents.

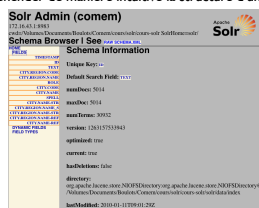
Par exemple, un type texte pour les contenus en français, un autre pour les contenus en anglais, etc.

Cela signifie qu'un même contenu texte se retrouvera indexé de plusieurs manières dans des champs différents. Cela permet une grande flexibilité.

17

schema

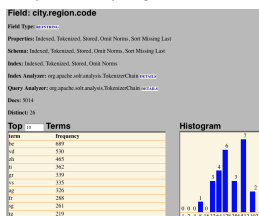
Solr offre un outil très pratique, le "Schema browser", pour appréhender de manière intuitive la structure d'un index.



18

schema

Chaque champ du schéma peut, grâce à cet outil, être analysé.



19

schema

Pour définir effectivement le schema d'une instance Solr, il faut éditer le fichier "schema.xml" qui se trouve dans le répertoire "config".

En plus de la définition des types et des champs, il contient aussi des paramètres généraux:

uniqueKey - champ utilisé comme identifiant unique
defaultSearchField - champ de recherche par défaut
defaultOperator - l'opérateur par défaut (OR ou AND)
copyField - permet de copier un champ dans un autre

20

schema

La définition d'un champ dans le schema de Solr peut avoir plusieurs paramètres:

name (obligatoire) - le nom du champ
type (obligatoire) - le nom du type du champ
indexed - true si le champ doit être indexé (recherchable et triable)
stored - true si la valeur du champ doit être stockée
compressed - true si le champ est compressé (ralenti les traitements)
multiValued - true si le champ peut être multivalué
default - une valeur par défaut

21

indexation

22

indexation

L'indexation se fait par l'intermédiaire de fichiers XML qui sont envoyés à Solr, via une requête HTTP POST sur l'url `http://localhost:8983/solr/update`

Le fichier XML contient des actions qui sont ensuite interprétées par Solr:

```
<add>
<commit> <rollback> <optimize>
<delete>
```

23

indexation

Pour effectuer la requête HTTP Post, n'importe quel logiciel peut être utilisé (par exemple, curl, wget, ...)

Un petit programme Java pratique est fourni avec Solr pour cela. Il s'agit du programme "post.jar".

```
java -jar post.jar fichier-a-indexer.xml
```

Aura pour effet d'envoyer le fichier à Solr et d'effectuer un commit automatique.

24

add

Permet d'ajouter un document dans

```
<?xml version="1.0" encoding="UTF-8" standalone="yes"?>
<add>
  <doc>
    <field name="role">city</field>
    <field name="ld">city.1000-5</field>
    <field name="city.code">-1000</field>
    <field name="city.name">Lausanne</field>
    <field name="city.region.code">VD</field>
    <field name="city.region.name">Vaud</field>
  </doc>
</add>
```

25

actions

L'action 'commit' permet de dire à Solr que les documents ajoutés peuvent être effectivement indexés.

```
<?xml version="1.0" encoding="UTF-8" standalone="yes"?>
<commit/>
```

L'action 'rollback' permet d'annuler les dernières additions.

```
<?xml version="1.0" encoding="UTF-8" standalone="yes"?>
<rollback/>
```

L'action 'optimize' permet d'optimiser l'index.

```
|
<?xml version="1.0" encoding="UTF-8" standalone="yes"?>
<optimize/>
```

26

delete

Permet de supprimer un document de l'index.
Soit par query, soit par id.

```
<?xml version="1.0" encoding="UTF-8" standalone="yes"?>
<delete>
  <id>city-1000-11</id>
</delete>
```

```
<?xml version="1.0" encoding="UTF-8" standalone="yes"?>
<delete>
  <q>city_postcode:1000</q>
</delete>
```

27

références

Solr 1.4 Enterprise Search Server by David Smiley and
Eric Pugh - Packt Publishing Ltd. ISBN
978-1-847195-88-3

(disponible à la bibliothèque HEIG-VD)

<http://lucene.apache.org/solr/tutorial.html>

<http://people.apache.org/~yonik/presentations/Solr.pdf>

28

Moteurs de recherche multimédias

Interfaçage

1

Formats de sortie



2

Formats de sortie

Plusieurs formats différents peuvent
être définis à l'aide du paramètre
'wt'

Exemples: xml, xslt, csv, json, ruby,
python php, phps, ...

3

Utilisation en XSLT

Le format d'export XSLT permet de traiter
le résultat XML avec une feuille de style.

Cette feuille de style peut être
programmée pour afficher une interface de
recherche.

Cette méthode implique un traitement
côté serveur et ne permet pas beaucoup
d'interactivité.

4

Utilisation en javascript

En utilisant javascript, on peut offrir une plus grande interaction.

Le format le plus pratique pour une utilisation avec Javascript est le format 'json'.

JSON est un format d'échange de données structurées entre différentes langages.

Il permet de transférer d'un langage à l'autre un tableau associatif.

5

Syntaxe JSON

JSON = JavaScript Object Notation

Un "objet" JSON peut contenir différents éléments:

nombre: 5
texte: "bonjour"
tableau: [item1, item2, item3]
objet: {"key1": value1, "key2": value2}

Les éléments structurés peuvent être imbriqués.

6

Exemple JSON

http://localhost:8983/solr/select/?q=role:audio&rows=1&indent=on&wt=json

```
{
  "responseHeader": {
    "status": 0,
    "QTime": 4,
    "params": {
      "q": "role:audio",
      "start": 0,
      "rows": 1,
      "wt": "json",
      "indent": "on"
    },
    "response": {
      "numFound": 29, "start": 0, "docs": [
        {
          "id": "audio-nicolas-allumer le feu",
          "role": "audio",
          "date_creation": "2010-12-31T00:00:00Z",
          "date_modification": "2010-12-31T00:00:00Z",
          "url": "http://www.nicolas.audio/allumer le feu",
          "album": "Ce que je sais",
          "auteur": "Mallory Johnny",
          "taille": 8.0,
          "titre": "Allumer le feu",
          "type": "mp3",
          "timestamp": "2011-12-05T13:29:55.045Z",
          "genre": ["Rock"]
        }
      ]
    }
  }
}
```

7

Utilisation de JSON

Avec javascript, la construction de la page html va se faire côté client, à partir des données obtenues de Solr.

Pour simplifier la programmation, on va utiliser la librairie JQuery.

Cette librairie nous offre la méthode 'getJSON' qui permet d'obtenir un objet JS à partir d'une requête retournant un JSON.

8

Exemple avec JQuery

9

```
var request = {};  
request['q'] = "**:*";  
request['rows'] = 10;  
request['start'] = 0;  
request['wt'] = "json";  
var url="http://localhost:8983/solr/select";  
$.getJSON(url, request, function(result, status, data) {  
    var documents = result.response.docs;  
    for ( var i = 0; i < documents.length; i++) {  
        displayDocument(documents[i]);  
    }  
});  
});
```
