

Railsアプリケーション開発者のためのSolr入門

Solr 1.2対応版

株式会社 ロンウイト

rel 1.0

- このたびは無料セミナー「Railsアプリケーション開発者のためのSolr入門」にご参加くださりまして、まことにありがとうございます。
- 本書はApache Solr 1.2の基本的な使い方について説明しています。
- 本書の権利は株式会社ロンウイトが保有しています。

- Goal
 - Apache Solr 1.2(以下Solr)をインストールし、起動できるようになる
 - Solrの管理画面の基本的な使い方を理解する
 - Solrへの文書登録の基本を理解する
 - Solrを使って基本的な検索ができるようになる
 - Solrを使ったファセットカウントの取得～絞り込み検索の流れを理解する
- Non-Goal
 - Solrのすべての機能を理解する
 - Solrの設定ファイルの項目の意味を理解し、使いこなせるようになる
 - Solrの各種プラグインを開発し、Solrを自在にカスタマイズする
 - パフォーマンスチューニングやメモリチューニングができるようになる

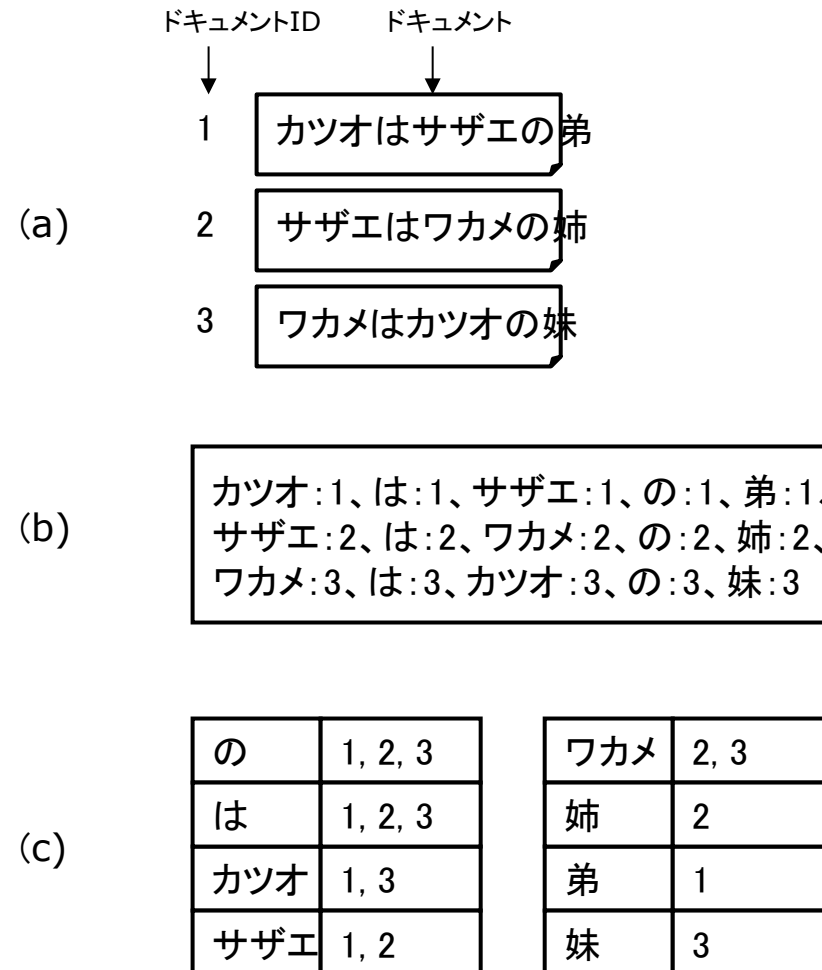
1. 日本語全文検索の基礎
 - 全文検索の方式
 - 転置索引の理解
2. Solrのインストールと起動
 - Solrのインストール
 - Solrの起動
 - Solr管理画面の起動
3. 付属exampleの使い方
 - exampleデータの登録
 - exampleデータの検索
 - 検索リクエストの基本パラメータ
 - solr-rubyの利用
4. 日本語データの登録と検索
 - schema.xmlの変更
 - 管理画面のANALYSIS機能の利用
5. Rails検索アプリケーション
 - ファセットと絞り込みのサンプルRails検索アプリケーション

日本語全文検索の基礎

- 順次検索方式
 - インデックスを作らない
 - ドキュメントの先頭から、検索質問語の文字列と順次比較する
 - 例: UNIXコマンドのgrep、SQLのlike検索

- 転置索引方式
 - あらかじめ検索対象のドキュメントからインデックスを作成
 - 例: Solr、Namazu、Senna、Google、Yahoo!、...

転置索引の作成方法





全文検索の各方式の特徴

方式	長所	短所	適用例
順次検索方式	<ul style="list-style-type: none">•インデックスを使わないので余計なメンテナンス作業が不要•インデックスを使わないので「今あるドキュメント」の内容をリアルタイムに検索できる	<ul style="list-style-type: none">•大量のドキュメントの検索には不向き•多数のユーザから繰り返し検索される状況下ではかなり非効率	<ul style="list-style-type: none">•UNIXのgrepコマンド•SQLのlike検索
転置索引方式	<ul style="list-style-type: none">•大量ドキュメントを保有し、多数のユーザから繰り返し検索される状況下でも効率的に処理できる•大規模な検索に向く	<ul style="list-style-type: none">•インデックスをメンテナンスしなければならないため、「今あるドキュメント」とインデックスの内容に差異が生じる場合がある•インデックスのサイズが巨大になる	<ul style="list-style-type: none">•多数のユーザから利用されるアプリケーション•インターネットやイントラネットなどのコンテンツ検索機能

- インデックスのサイズ
- インデックスの作成にかかる時間
- ヒットしすぎる
 - 検索結果一覧の表示順序(ランキング)のスコアを計算
 - Google ⇒ PageRank
 - Lucene ⇒ tf*idf (*Similarity*抽象クラスのJavadoc参照)
- 日本語テキスト処理
 - 英語などと違って、単語を識別するのが困難
 - 形態素解析
 - 辞書を用いる方式が主流 ⇒ 流行語に弱い
 - 辞書にない単語が検索されにくい(「東京」で「東京大学」がヒットしない)
 - JapaneseAnalyzer
 - N-gram
 - Nの大きさによりノイズや検索漏れ
 - CJKAnalyzer (「東京都」⇒「東京」「京都」)
 - 表記の揺れ
 - 「インタフェース」「インターフェイス」、「引っ越し」「引越し」
 - 半角・全角、新旧漢字(「亜」と「亞」)

Solrのインストールと起動

(演習) Solrのインストール～起動

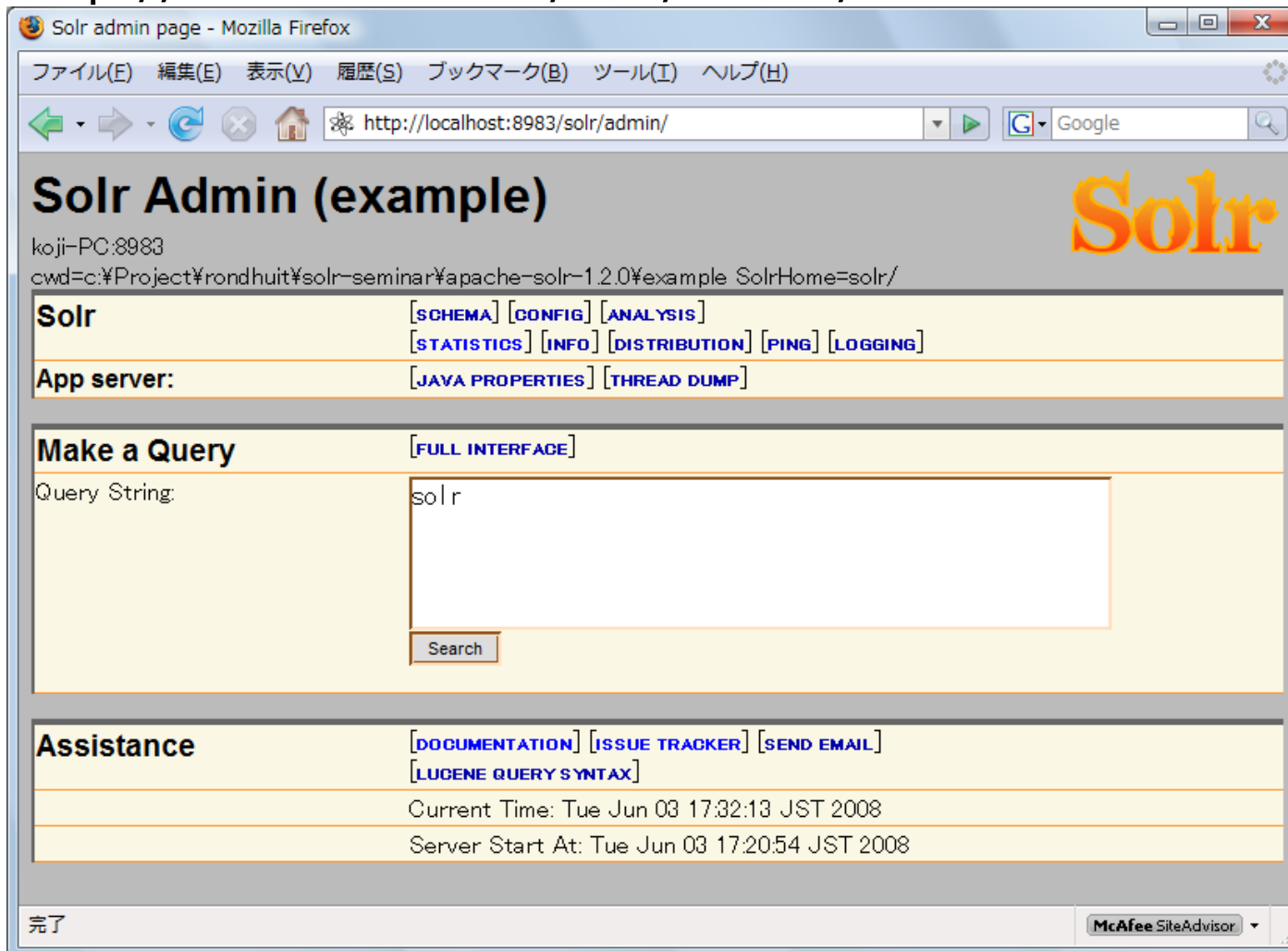
- インストール
 - apache-solr-1.2.0.zipを展開します (Subversionから作成したパッケージ)
 - これ以降、展開したディレクトリを\$SOLRとして参照します
- 起動
 - \$SOLR/exampleに移動します
 - start.jarを使って起動します

```
$ cd $SOLR/example
$ java -jar start.jar
2008-06-03 17:20:52.050::INFO: Logging to STDERR via org.mortbay.log.StdErrLog
2008-06-03 17:20:52.376::INFO: jetty-6.1.3
      :
      :
2008/06/03 17:20:55 org.apache.solr.servlet.SolrDispatchFilter init
情報: SolrDispatchFilter.init() done
2008/06/03 17:20:55 org.apache.solr.servlet.SolrServlet init
情報: SolrServlet.init()
2008/06/03 17:20:55 org.apache.solr.servlet.SolrServlet init
情報: SolrServlet.init() done
2008/06/03 17:20:55 org.apache.solr.servlet.SolrUpdateServlet init
情報: SolrUpdateServlet.init() done
2008-06-03 17:20:55.675::INFO: Started SocketConnector @ 0.0.0.0:8983
```

- 停止は[Ctrl]+[c]

(演習) Solr管理画面の起動

- ブラウザを使って管理画面を起動します
http://localhost:8983/solr/admin/



付属exampleの使い方

(演習) 付属exampleデータの登録と検索 RONDHUIT

- 付属exampleデータの登録

- \$SOLR/example/exampledocsに移動します
- post.jarを使って登録します

```
$ cd $SOLR/example/exampledocs
$ java -jar post.jar *.xml
SimplePostTool: version 1.2
SimplePostTool: WARNING: Make sure your XML documents are encoded in UTF-8, other
encodings are not currently supported
SimplePostTool: POSTing files to http://localhost:8983/solr/update..
SimplePostTool: POSTing file hd.xml
      :
      :
SimplePostTool: POSTing file vidcard.xml
SimplePostTool: COMMITting Solr index changes..
```

- 登録されるXMLをエディタで開いて見てみよう

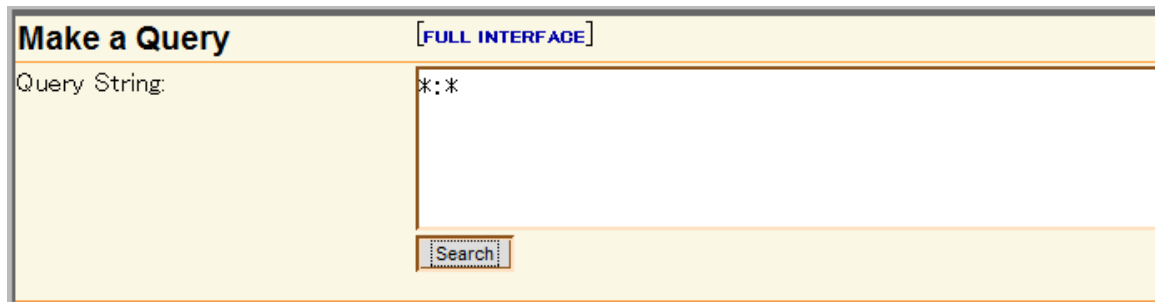
文書登録

```
<add>
  <doc>
    <field name="id">12345</field>
    <field name="name">マッサージチェア</field>
    <field name="cat">家電>健康器具>マッサージ</field>
  </doc>
  <doc/>
</add>
```

登録する文書(検索対象文書)

(演習) 付属exampleデータの登録と検索 RONDHUIT

- 付属exampleデータの検索
 - 管理画面から適当な検索語で検索します
 - レスポンスXMLを観察しよう
 - 検索語に"*:*"と指定すると「すべて検索」の意味になります



The screenshot shows a web interface titled "Make a Query" with a sub-link "[FULL INTERFACE]". Below the title, there is a "Query String:" label followed by a text input field containing the text "*:*". At the bottom right of the input area, there is a "Search" button.

- 基本パラメータ
 - q=検索語(検索式)
 - start=開始文書番号(デフォルト=0)
 - rows=ページあたりの文書件数(デフォルト=10)
 - fl=フィールド指定
 - wt=ライター名(xml(デフォルト),json,ruby,python,...)
 - indent=on(出力の字下げを実行)
 - debugQuery=on(デバッグ出力)
- ファセットカウント取得 & 絞り込み検索関連パラメータ
 - facet=on
 - facet.field=ファセットフィールド指定
 - facet.mincount=ファセットフィールド指定時最小カウント指定
 - facet.query=任意の検索式(範囲検索など)
 - fq=絞り込み検索条件

http://localhost:8983/solr/select?q=%3A*&indent=on&facet=on&facet.field=cat

(演習) solr-rubyの利用

- solr-rubyクライアントライブラリ
 - \$SOLR/client/ruby/solr-ruby/lib/solr.rb
- Solr接続の取得
 - solr=Solr::Connection.new("http://localhost:8983/solr")
- Solr検索リクエスト
 - request=Solr::Request::Standard.new(:query=>'*:*')
 - response=solr.send(request)
- irbから使ってみよう

```
$ irb -I $SOLR/client/ruby/solr-ruby/lib -r solr
```

```
irb(main):001:0> solr=Solr::Connection.new("http://localhost:8983/solr")  
=> #<Solr::Connection:0x3b6eb00 @url=#<URI::HTTP:0x1db74d6 URL:http://localhost:  
8983/solr>, @connection=#<Net::HTTP localhost:8983 open=false>, @autocommit=false>
```

```
irb(main):002:0> request=Solr::Request::Standard.new(:query=>'*:*')  
=> #<Solr::Request::Standard:0x3b6b270 @query_type="standard", @params={:query=>  
"*:*", :field_list=>["*", "score"]}>
```

```
irb(main):003:0> pp solr.send(request)
```

日本語データの登録と検索

1. インデックスの削除

- Solrサーバを停止して、\$SOLR/example/solr/data/indexディレクトリを削除します

2. schema.xmlの変更

- \$SOLR/example/solr/conf/schema.xmlを入れ替えてSolrサーバを起動します
- schema.xmlをエディタで開いて見てみよう

3. 日本語サンプルデータの登録

- post.jar(P.15)を使ってblog.xmlを登録します

4. 管理画面からの検索

- 管理画面を起動し、適当な単語で検索してみます
- (注意) fl=* (デフォルト)を指定した場合、巨大な文書がヒットするとレスポンスが遅くなる場合があります

(演習) CJKAnalyzerの動作を理解する

The screenshot shows the Solr Admin interface in a Mozilla Firefox browser window. The address bar shows the URL: http://localhost:8983/solr/admin/analysis.jsp?nam. The page title is "Solr Admin (example)".

On the left side, there is a navigation menu with the following items: [SCHEMA], [CONFIG], [ANALYSIS], [STATISTICS], [INFO], [DISTRIBUTION], [PING], [LOGGING], [JAVA PROPERTIES], and [THREAD DUMP]. The [ANALYSIS] item is highlighted with a red dashed box.

Below the menu, there is a text box: "管理画面メニューから [ANALYSIS]を選択"

The main content area is titled "Field Analysis" and contains the following form:

Field name	text
Field value (Index)	今日はいい天気です。
Field value (Query)	天気

There are checkboxes for "verbose output" and "highlight matches" (both checked) under the Field value (Index) section, and "verbose output" (checked) under the Field value (Query) section. An "Analyze" button is located at the bottom of the form.

Below the form, there is a section titled "Index Analyzer" with the class name "org.apache.lucene.analysis.cjk.CJKAnalyzer {}". It contains a table with the following data:

term position	1	2	3	4	5	6	7	8
term text	今日	日は	はい	いい	い天	天気	気で	です
term type	double	double	double	double	double	double	double	double
source start,end	0,2	1,3	2,4	3,5	4,6	5,7	6,8	7,9

Below the table, there is a section titled "Query Analyzer" with the class name "org.apache.lucene.analysis.cjk.CJKAnalyzer {}". It contains a table with the following data:

term position	1
term text	天気
term type	double
source start,end	0,2

At the bottom of the page, there is a status bar with the text "完了" and a "McAfee SiteAdvisor" logo.

Rails検索アプリケーション

(演習) Rails検索アプリケーションの準備 RONDHUIT

1. Railsアプリケーションの雛形生成

```
$ rails sample
```

2. environment.rbの修正

- ロードパスに"\$SOLR/client/ruby/solr-ruby/lib"を加えます
 - 代わりに"\$ gem install solr-ruby"でも可

```
$.unshift '$SOLR/client/ruby/solr-ruby/lib'
```

- ActiveRecord不使用の設定を行います

```
Rails::Initializer.run do |config|  
  config.frameworks -= [ :active_record ]  
end
```

3. 既成ファイルのコピー

- app/controllers/search_controller.rb
- app/helpers/search_helper.rb
- app/views/search/*.rhtml
- lib/solr_context.rb

起動時`load_missing_constant': uninitialized constant ActiveRecord (NameError)"のエラーになる場合は、config/initializers/new_rails_defaults.rbのActiveRecord定数を参照している2つの行をコメントアウトする。

4. Railsアプリケーションサーバの起動

```
$ ruby script/server
```

1. トップページが表示

<http://localhost:3000/search/index>

- しくみを観察しよう
 - qパラメータはどのようになっていますか？
 - ファセットカウント取得をどのように行っていますか？
 - リンクはどのようになっていますか？

2. 検索結果ページの表示

- しくみを観察しよう
 - qパラメータはどのようになっていますか？

3. 絞り込み検索の実行

- しくみを観察しよう
 - fqパラメータはどのようになっていますか？

4. 検索結果[次ページ][前ページ]の表示

- しくみを観察しよう
 - startパラメータはどのようになっていますか？