

医療情報 (MML3.0) Hadoopプログラミング

2011年10月21日

キヤノンソフト情報システム 上村準也

発表の経緯



- パッケージ製品の開発者です
 - JavaでWebアプリケーションを開発
 - PaaS/IaaSや仮想マシンを使うことに興味あり
- Hadoopについては
 - 技術的なニュースは度々聞くが使ったことはない
 - ・ 仕事でも案件として話を聞くことがあるように
- 2011年1月にHadoopのハンズオン
 - クラウド・テクノロジー分野第四回勉強会で初めて
 - 本業とは関係なく勉強することに

発表の内容



- 医療情報(MML3.0)を題材にHadoopでプログラミング
 - Hadoop Core 0.21.0 を使用
 - MapReduceアルゴリズムの話ではなく…
- 1. 開発環境の準備
- 2. テストコードの記述
- 3. 圧縮データの入力
- 4. XMLデータの入出力
 - XMLデータをどう扱うのが良い?
 - <u>象本とHadoop徹底入門</u>を読みながら試したこと

開発環境の準備(1)



- 開発環境 Windows7 および Windows XP SP3
- ダウンロード&インストール
 - JDK Oracle 1.6.0_27
 - Eclipse Indigo Service Release 1
 - Gnu On Windows (Gow) 0.4.0
 - <u>Hadoop</u> 0.21.0

開発環境の準備(2)



• Hadoopの設定ファイルを編集

conf/core-site.xml

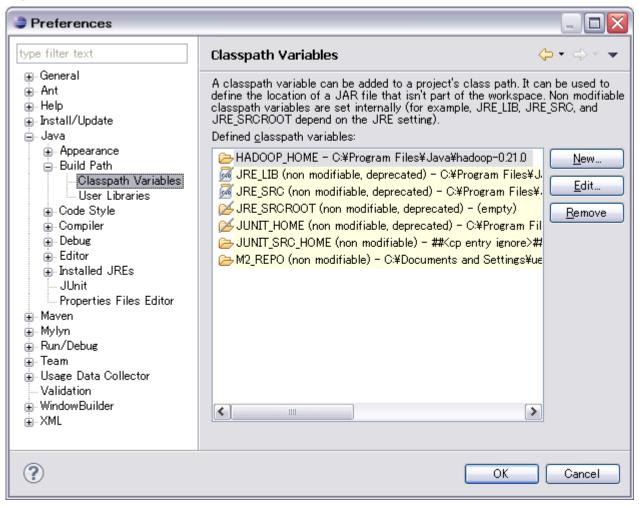
conf/mapred-site.xml

```
<name>mapreduce.jobtracker.address
```

開発環境の準備(3)



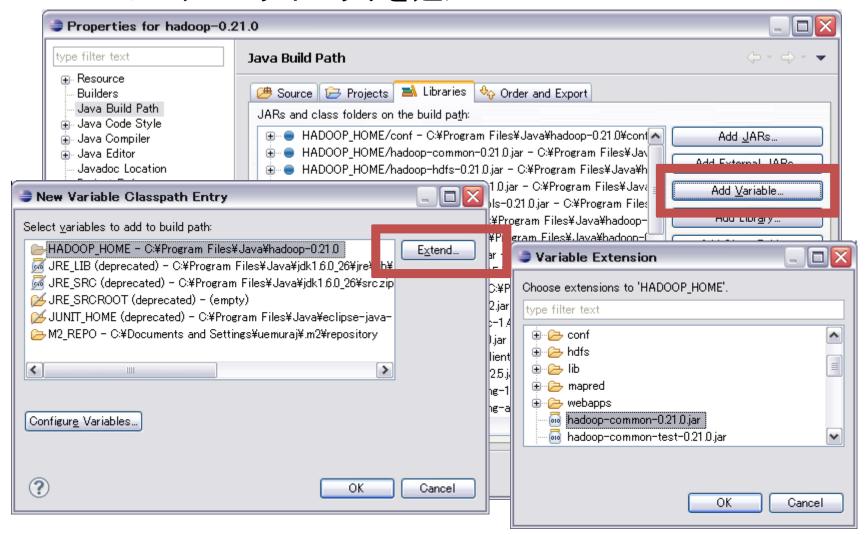
• Eclipseのクラスパス変数を追加



開発環境の準備(4)



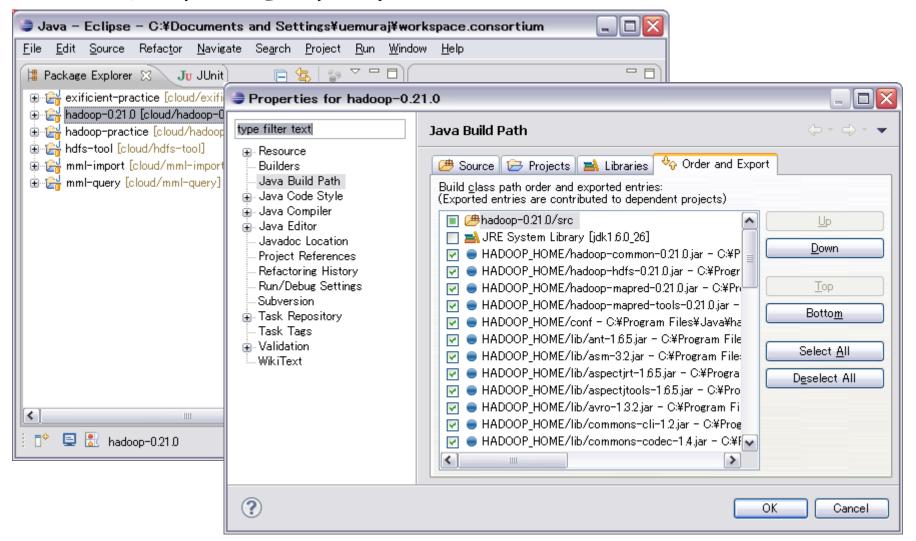
プロジェクトにライブラリを追加



開発環境の準備(5)



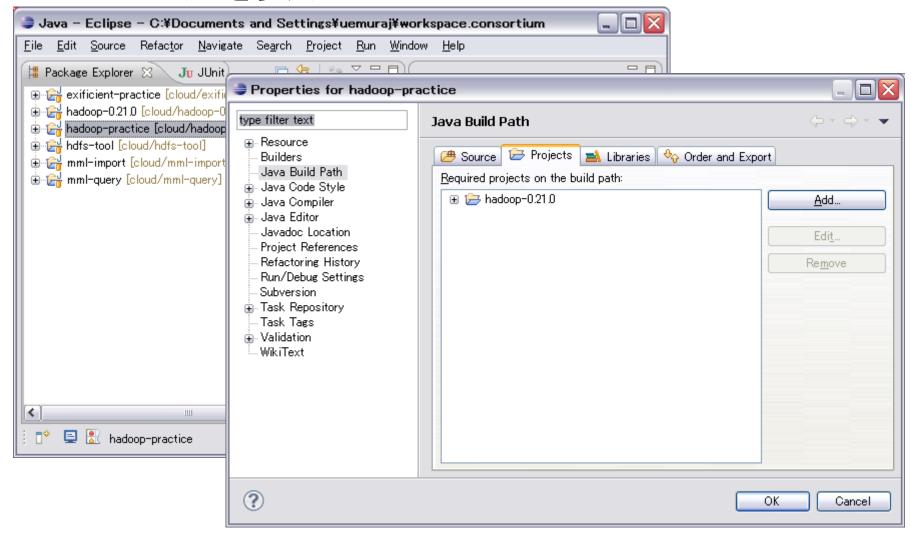
・ プロジェクトからエクスポート



開発環境の準備(6)



プロジェクトを参照



テストコードの記述



- Hadoopで記述した初めてのコード
 - Javaソースコードのimport文を数える
- テストされるコード PackageCount.java
- テストするコード PackageCountTest.java
- テスト結果 PackageCountTest

圧縮データの入力(1)



- ZIP/TAR形式のデータを処理するInputFormatを実装
 - 入手したテスト用データをそのまま処理したい
 - MMLはXMLでありテキストデータ、まとめて圧縮すると 圧縮率は高く、取り回しには便利
- 標準でもサポートされているが少し扱いが難しい
 - ZIPエントリで切り分けられたバイトデータがそのまま渡 されるように見える
- それらしい名前の <u>ZipFileRecordReader.java</u> で検索してみる
 - 独自に実装されたものも見つかる(0.20系)
 - トレーニングとしてZIP/TARそれぞれ実装を試みる

圧縮データの入力(2)



- ZIP/TAR形式のデータを処理するInputFormatを実装
 - ZIPとTARどちらも非常に似た実装になる
 - GZIP圧縮されたTAR形式の場合の処理を紹介
- 実装したコード TgzFileInputFormat.java
- テストするコード <u>TgzFileInputFormatTest.java</u>
- テストするコード XmlValidateMapper.java
- テスト結果 <u>TgzFileInputFormatTest</u>

XMLデータの入出力(1)



- StreamXmlRecordReader を試してみる
 - mapred/contrib/streaming に用意されているのだが…
- テストするコード <u>StreamXmlRecordReaderTest.java</u>
- 実行結果 <u>StreamXmlRecordReaderTest</u>

XMLデータの入出力(2)



- org.w3c.dom.Documentを処理するSerializationを実装
 - Mapperの入力としてDOMオブジェクトが欲しい
 - Mapperを重ねれば可能
 - だがSerializationを実装するとノード間の通信にも使用 される
 - XMLを良く圧縮できる実装があれば全体の効率UPが見込める
 - Writableを実装する方法よりシンプルに
- バイナリXMLを使用したSerializationを実装
 - XMLコンソーシアムに多数の公開資料
 - <u>EXIficient</u> W3C Efficient XML Interchange (EXI) フォーマットのオープンソース実装を試す

XMLデータの入出力(3)



- org.w3c.dom.Documentを処理するSerializationを実装
 - WritableSerializationのソースコードを参考に
- 実装したコード <u>DocumentHolder.java</u>
- 実装したコード DocumentSerialization.java
- テストするコード <u>DocumentSerializationTest.java</u>
- テストするコード XmlDocumentMapper.java
- 実行結果 <u>DocumentSerializationTest</u>

XMLデータの入出力(4)



- <u>DocumentSerialization.java</u>を使用して勉強会の課題を解く
 - 年齢と性別をキーに、その患者さんのBMI値を出力
- 実装したコード BmiMapper.java
- テストするコード BmiTest.java
- 実行結果 BmiTest

まとめ



- まだまだ勉強すること、試すことが多くあります
 - 測定して実装方法による特性を調べたり
 - もっと他の手段がありそう、Hadoopを利用する他のソフトウェアも見てみたい
 - 実際にはそれらのソフトウェアを使うことになると思うが、 このレベルの勉強は後々自分を助けると考える
- Hadoopをどのように稼動して利用するか
 - 「Amazon EC2」「GMOクラウド」「さくらVPS」での環境構築も試行錯誤中
 - 大規模な環境やセキュリティの話を聞くが、個人的には 中小規模でも面白いことがありそうだと感じる

さいごに



クラウド部会で一緒に勉強しましょう