# JWEBMINER : CULTURAL FEATURE EXTRACTOR SOFTWARE PACKAGE OVERVIEW

**Gabriel Vigliensoni**

Music Technology Area, Schulich School of Music, McGill University

gabriel@music.mcgill.ca

## ABSTRACT

The amount of digital media and digital information about it stored in the web is growing exponentially. Data mining researchers studies the extraction of patterns from this data, how to classify and organize it to provide better ways to search, query and access it. In the music field, low-level, high-level and cultural features can be extracted to perform classify and organize musical information. jWebMiner is an open source web package to extract cultural features from the web.

## 1. INTRODUCTION

The amount of music available to users has grown extremely fast during the last twenty years. Digital media and formats, such as compact discs and MPEG1 layer 3, are standard formats to listen, share and carry music from one place to another. In addition, with the increase in the bandwidth available on the Internet, realtime audio transmission with very good consumer quality is available. Furthermore, the number of people and computers connected to the net grows everyday, allowing them to share complete music libraries through peer-to-peer networks.

On the other hand, the amount of information about music grows in the same degree. Online music magazines, weblogs, MP3 blogs, social networks and commercial music sites, stores dynamically huge quantities of musical information written everyday by users and media journalists.

Music information retrieval is a research field that investigates how to deal with big music collections to be categorized, grouped, classified and queried. In order to achieve this goal, a multi-disciplinary approach to music is been done. The study of audio, symbolic and cultural music information is combined by researchers to obtain better results in how we understand, categorize and organize music (McKay 2008).

jMIR, an open source software package developed at McGill University by Cory McKay, gives the possibility to extract low-level features (based on basic signal processing and human physiology), high-level features (based on mu-

sical abstractions) and cultural features (based on cultural information) into one single platform. Each kind of features provides different and complementary musical information, and its combination is useful to gain accuracy and performance in classification and clustering tests (McKay 2007).

## 2. JWEBMINER

jWebMiner is a software package inside the jMIR bundle that allows the extraction of cultural metadata using web services. Basically, it counts hits from different websites using a search text string as a query input (or iTunes XML, ACE XML and Weka ARFF). The number of hits can express the co-occurrence between a number of different text strings (how many times a text string appears in pages that have another text string) or the cross tabulation between two sets of text strings (allowing to measure how often a text string of one class appears in text strings of other class).

### 2.1 WebServices

Web services allows machine interoperability interaction over a network. Through this technology, jWebMiner is capable to communicate with Yahoo! and Google to query strings. Different search engines provides keys to authenticate as users allowing a maximum daily number of request by IP number.

### 2.2 Filtering and Weighting

In order to obtain more clear results, the system has the ability to configure synonyms for different text strings (e.g. allowing the user to control that people can tag in different ways an artist name), filtering hits that does not have any user-definable text strings (e.g. narrowing the huge search to only desired sites that contains some information), and weighting some websites in a bigger or smaller degree depending the ability of the software user to sense how much important is public opinion in one site over other (e.g. giving more weight to sites devoted explicitly to music).

Additional filters for language, region and filetype are offered to narrow the search to some given characteristics and do not count redundant results.

## 2.3 File Output

jWebMiner also offers the possibility to change its statistical functions and the way data is normalized in order to research best ways to extract, weight and see information. Its output data possibilities include ACE XML, Weka ARFF and delimited text files.

## 2.4 Extensibility

jWebMiner is downloadable as standalone user or developer version. The last one allows the possibility to configure or develop special characteristics or new implementations such as the addition of further web services, configure the output visualization or any feature needed for a special project through changes or extension of its application program interface (API).

## 3. USING JWEBMINER

### 3.1 Cross tabulation and co-ocurrence

The GUI and help files of jWebMiner are so clear that its use is straightforward. In its cross tabulation form, the software user must entry a set of strings in the primary field (e.g. musicians names) and in the secondary field another set of strings (e.g. musical styles). If the extract features button is pressed, the system will query automatically the Yahoo site the number of strings of the first field times the number of strings of the second field. Although the process is simple, the number of queries to search engines will be high if we have big data sets. This could be a problem if some limit of queries is given by the provided keys (in fact, Yahoo! has a maximum of 5000 daily queries by unique IP. This will limit our possibilities to a set of five hundred entries classified in ten different classes, but queried only one time)

In co-occurrence mode, the data set must be placed in the primary search string only. The amount of queries will be given by querying each entry with all the other entries. Thus, if $n$ is the number of entries, the system will query search engines $(n*(n+1))/2$. If we have the same of five hundred entries, jWebMiner will query 125250 times an external application. Therefore, if we want to extract cultural metadata using co-occurrence through Yahoo! search engine, out daily data set by unique IP should be limited by 99 different entries.

### 3.2 Search engines

In addition to the limits given by Yahoo!, Google is no longer providing access to its API through web services using SOAP since August 31, 2009. The system migrated to an AJAX API, so to perform a google search into jWebMiner, its API must be changed.

### 3.3 Understanding how we humans refers to music

The results provided by jWebMiner are very sensible to the text strings provided and filters applied. Thus, if we want to perform a music genre classification, artist like The Beatles will be classified as *classical* due to people refers them as *classics* into popular music. On the other hand, depending on the genre categories provided and filters applied, an artist such as Bruce Springsteen can be categorized as jazz, rock, reggae or punk.

## 4. CONCLUSIONS AND FUTURE DEVELOPMENT

jWebMiner is a open source and extensible powerful tool to perform extraction of cultural metadata from the web, allowing MIR researchers to use semantic information about music provided by thousands of anonymous users. This cultural features can complement low-level and high-level features extracted from audio signals and musical abstractions to classify and categorize music more accurate and efficiently. jWebMiner also allows some powerful tools to narrow and filter the search possibilities given by search engines, and the way the output results are given.

However, some future development is needed in order to achieve more accurate results and perform more systematic research. First, new search engines APIs must be programmed in order to perform faster, different and more diary queries. Second, a study of *natural language* must be done in order to understand how human refers to music. This will be important to extract *semantic adjectives* and *noun phrases* from weblogs and electronic publications in general (Whitman 2002). Finally, a improvement in the html results presentation must be implemented allowing a user to better compare results.

## 5. REFERENCES

[1] McKay, C., and I. Fuchinaga. 2007. jWebMiner: a web-based feature extractor. In *Proceedings of the 8th International Conference on Music Information Retrieval*. Vienna. 113-4.

[2] McKay, C., and I. Fuchinaga. 2008. Combining features extracted from audio, symbolic and cultural sources. In *Proceedings of the 9th International Conference on Music Information Retrieval*. Philadelphia. 597-602.

[3] Whitman, B., and P. Smaragdis. 2002. Combining musical and cultural features for intelligent style detection. In *Proceedings of the 3rd International Conference on Music Information Retrieval*. Paris. 47-52.

[4] Whitman, B., and S. Lawrence. 2002. Inferring descriptions and similarity for music from community metadata. In *Proceedings of the 2002 International Computer Music Conference*. Paris. 591-8.

[5] Zadel, M., and I. Fuchinaga. 2004. Web services for music information retrieval. In *Proceedings of 5th International Conference on Music Information Retrieval*. Barcelona. 478-83.