

# jWebMiner: A Web-Based Feature Extractor

## Cory McKay and Ichiro Fujinaga



Social Sciences and Humanities  
Research Council of Canada



McGill



Schulich School of Music  
École de musique Schulich



Centre for Interdisciplinary Research  
in Music Media and Technology



### Overview

jWebMiner is a software package for extracting cultural features from the web. At its most basic level, it operates by using web services to extract hit counts for sets of query strings from search engines. It then processes these hit counts to calculate feature values denoting classifications or similarity measurements. jWebMiner emphasizes extensibility, generality and an easy-to-use interface.

### Cross-Tabulation Extraction

One of the two main types of feature extraction that can be performed by jWebMiner is "cross tabulation extraction." This involves finding relatively how often each query string in one set of query strings appears on the same web page as each string in another set of query strings. This is useful for classifying instances from one list into classes denoted in the other list, as shown in Figures 1 and 2:

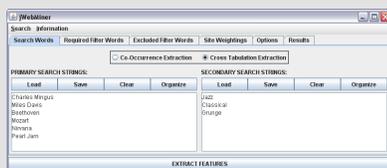


Figure 1: Query terms used for a sample cross tabulation extraction that will classify the six musical artists amongst the three musical genres.

	Classical	Genre	Jazz
Beethoven	<b>0.57685729661511</b>	0.083023385112237	0.3489208840274252
Charles Mingus	0.22287765011872584	0.16881377497486955	<b>0.626386274984692</b>
Miles Davis	0.22418619120572363	0.0948076872578883	<b>0.677468584283897</b>
Mozart	<b>0.56678526249966</b>	0.08883525320532933	0.3645692107573874
Nirvana	0.11754652389738058	<b>0.7657867962584234</b>	0.17674676169942264
Pearl Jam	0.0894814752561329	<b>0.792247789649911</b>	0.17177007559940205

Figure 2: The results for the queries shown in Figure 1. The figures represent the normalized weighted relative frequencies with which each artist appears on the same web page with each genre. The highest value for each artist is in bold.

### Synonyms

Multiple terms are often used to refer to the same entity. For example, "Charles Mingus" and "Charlie Mingus" are often used interchangeably. jWebMiner therefore makes it possible to define synonyms so that hits for corresponding synonyms are combined.

### Co-Occurrence Extraction

The other main type of feature extraction that can be performed by jWebMiner is "co-occurrence extraction." This involves finding relatively how often each query string in a set of query strings appears on the same web page as every other query string in the same set. This is useful for similarity measurement, as demonstrated in Figures 3 and 4:

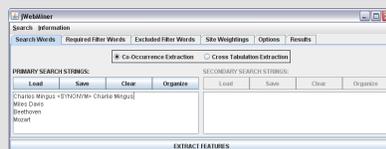


Figure 3: Query terms used for a sample co-occurrence feature extraction that will measure the similarity between each of the four musical artists. A synonym for "Charles Mingus" is included.

	Beethoven	Charles Mingus Charlie Mingus	Miles Davis	Mozart
Beethoven	-	0.0646173807491334	0.20758391144428176	<b>0.4857967812958482</b>
Charles Mingus Charlie Mingus	0.030543728177661813	-	<b>0.814289489876489</b>	0.02526582388465726
Miles Davis	0.06688074739193561	<b>0.873962496322882</b>	-	0.0642563526164477
Mozart	<b>0.4916275964976749</b>	0.0797323501164333	0.22612194938589172	-

Figure 4: The results for the queries shown in Figure 3. The figures represent the normalized weighted relative frequencies with which each artist appears on the same web page with each other artist. The highest value for each artist's row is in bold.

### Required String Filters

There is usually some feature noisiness due to hits not related to the topic being studied. There are many Hinduism-related sites that contain the string "Nirvana," for example, but that have nothing to do with the band. jWebMiner allows users to define required string filters (e.g., "music") that must be found on a site in addition to the query strings for it to be counted as a hit.

### Excluded String Filters

Excluded string filters offer another way for jWebMiner users to reduce feature noisiness. There are user-definable strings that may not appear on a web site if it is to be included in feature counts. For example, one might wish to exclude sites containing the words "paradise" or "zen" in order to avoid false hits relating to the band "Nirvana."

### Site Weightings

Another way to reduce feature noisiness is to limit searches to sites that are known to be relevant to music, such as the All Music Guide, Pitchfork, etc. jWebMiner therefore allows users to assign relative weightings to both the Internet as a whole and to individual sites. These weightings control the effect that hits from certain sites have relative to hits from other sites on the final feature calculations.

### Additional Options

jWebMiner offers a variety of additional ways that users can customize their feature extractions. Users may choose among a variety of formulas and normalization schemes for calculating feature values. Users may also choose to use either Google, Yahoo! or both to perform searches, and may configure these engines using a variety of options, as shown in Figure 5:

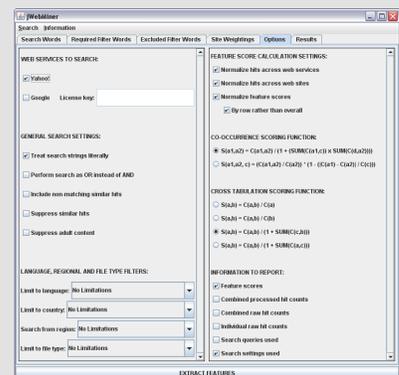


Figure 5: Additional ways of customizing jWebMiner feature extractions.

### Usability & Extensibility

jWebMiner includes a simple and intuitive GUI as well as a full user manual. The software is open source and the Java code is well-documented and implemented using a modular plug-in design to facilitate the addition of new functionality such as additional web services and feature calculation schemes. jWebMiner can parse query terms from text, iTunes XML, ACE XML, or Weka ARFF files. Feature values can be saved as text, ACE XML, Weka ARFF or HTML files. jWebMiner is one component of the jMIR software suite, which includes audio and symbolic feature extractors and machine learning software.