# Information Integration on the Web: Homework 1

**Tushar Tiwari**

## Statement of Integrity

I, Tushar Tiwari, declare that the submitted work is original and adheres to all University policies and acknowledge the consequences that may result from a violation of those rules.
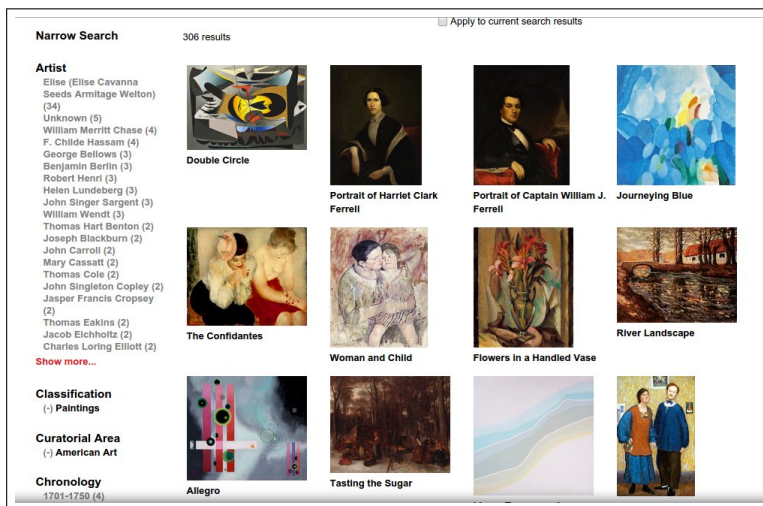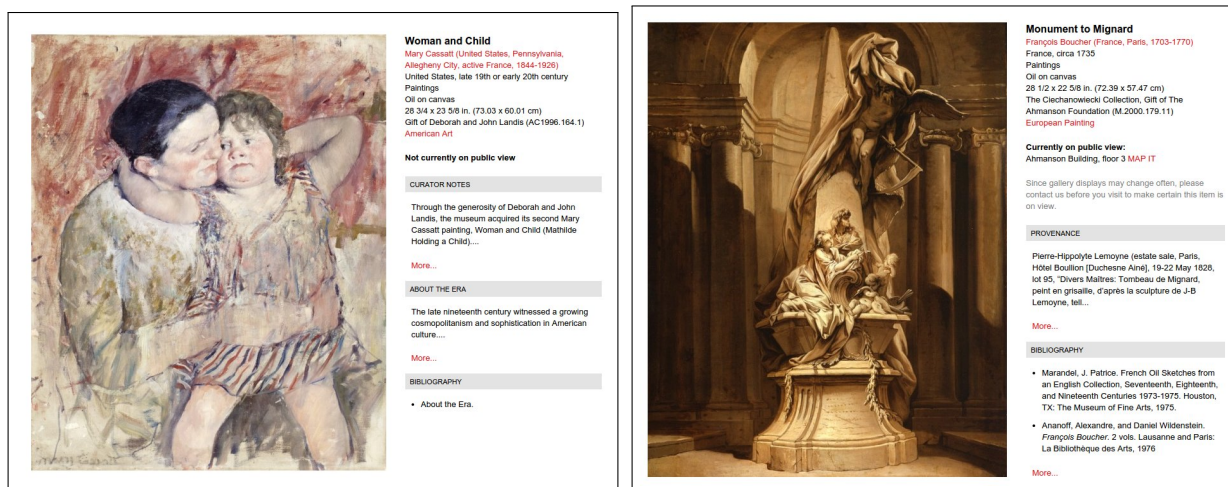
## Screenshots



Figure 1: The search page with the option American Paintings already chosen.
This provides us with a custom url that can be exploited to navigate and then extract all paintings.
`http://collections.lacma.org/search/site/?page=0&f[0]=bm_field_has_image%3Atrue&f[1]=im_field_curatorial_area%3A32&f[2]=im_field_classification%3A22`



(a) An american painting

(b) A european painting

Figure 2: Examples of the paintings on the website

## Sample Record

Let's look at a sample painting and the output that it produces. For purposes of demonstration the link to this sample painting is hard-coded into the program so that the output produced is only for the painting in consideration.



Figure 3: The chosen painting found at `http://collections.lacma.org/node/184481`.



(a) Hardcoding the url into the code



(b) Execution of make command

Figure 4: Sample Run

**Note:**

- The url for the painting is hardcoded only for purposes of demonstration. The usual run will automatically extract image urls from the search page.

- Not all paintings will have all data fields associated with them since the data is unavailable for those paintings. These missing fields have been ignored as they serve no purpose except for increasing the file size which is not desirable.

```
1    [
2        {
3            "_id": 184481,
4            "acquired_from": "Bequest of Fannie and Alan Leslie",
5            "acquired_id": "M.2006.73.2",
6            "artist": {
7                "birthplace": {
8                    "country": "United States"
9                },
10                "dates": {
11                    "born": "1899",
12                    "died": "1978"
13                },
14                "name": {
15                    "first": "Charles",
16                    "last": "Howard",
17                    "middle": "Houghton"
18                }
19            },
20            "category": "american",
21            "copyright": "Estate of Charles Houghton Howard",
22            "created_country": "United States",
23            "created_year": "1950",
24            "image_link": "http://collections.lacma.org/sites/default/files/
25            remote_images/piction/ma-31205657-WEB.jpg",
26            "link": "http://collections.lacma.org/node/184481",
27            "name": "Double Circle",
28            "on_view": "Art of the Americas Building, floor 3",
29            "painting_type": "Oil on canvas"
30        }
31    ]
```

Figure 5: Output JSON of sample run

## Challenges

- Unicode characters are a problem in Python 2 since the default string supports only ascii characters. Due to this certain string operations become very complicated. The problem was oversome by porting the code to Python 3 whose default string supports unicode

- Extracting the artist's details required complex regular expressions since all information pertaining to the artist, like name, date of birth, place of birth, was provided in a single string in the html.

# Tools Used

## Python with BeautifulSoup

Python is one of the most easiest to use yet powerful languages out there which works cross-platform. Installing and distributing dependencies is also very easy. Beautiful Soup is a Python library designed for quick turnaround projects like screen-scraping. It parses anything you give it, and does the tree traversal stuff for you. You can tell it "Find all the links", or "Find all the links of class externalLink", or "Find all the links whose urls match "foo.com", or "Find the table heading that's got bold text, then give me that text.". Some of the features that make it powerful are:

- Provides a few simple methods and Pythonic idioms for navigating, searching, and modifying a parse tree: a toolkit for dissecting a document and extracting what you need. It doesn't take much code to write an application

- Automatically converts incoming documents to Unicode and outgoing documents to UTF-8. You don't have to think about encodings, unless the document doesn't specify an encoding and Beautiful Soup can't detect one. Then you just have to specify the original encoding.

- Sits on top of popular Python parsers like lxml and html5lib, allowing you to try out different parsing strategies or trade speed for flexibility.

# Fields not Extracted

There were no fields that were left untouched. Every field was extracted. But, not all paintings have data for every field. All fields were checked for data. If valid information was present only then it was extracted.