# HW#4 - Web Archiving, Part 2
Ethan Landers
Due: Sunday, October 27, 2024, by 11:59 PM

# Q1

*Q: What can you say about the relationship between the age of a URI-R and the number of its mementos?*

This question was challenging to answer because the URI-Rs I collected TimeMaps for had either 0 or 3 mementos each. I excluded the URI-Rs with 0 mementos from my analysis, resulting in plots for those with only 3 mementos. A scatter plot was not effective in this case due to the lack of variability in the number of mementos. Instead, I opted to create a box plot, which effectively illustrates the age of each memento across all the URI-Rs for which I grabbed TimeMaps.
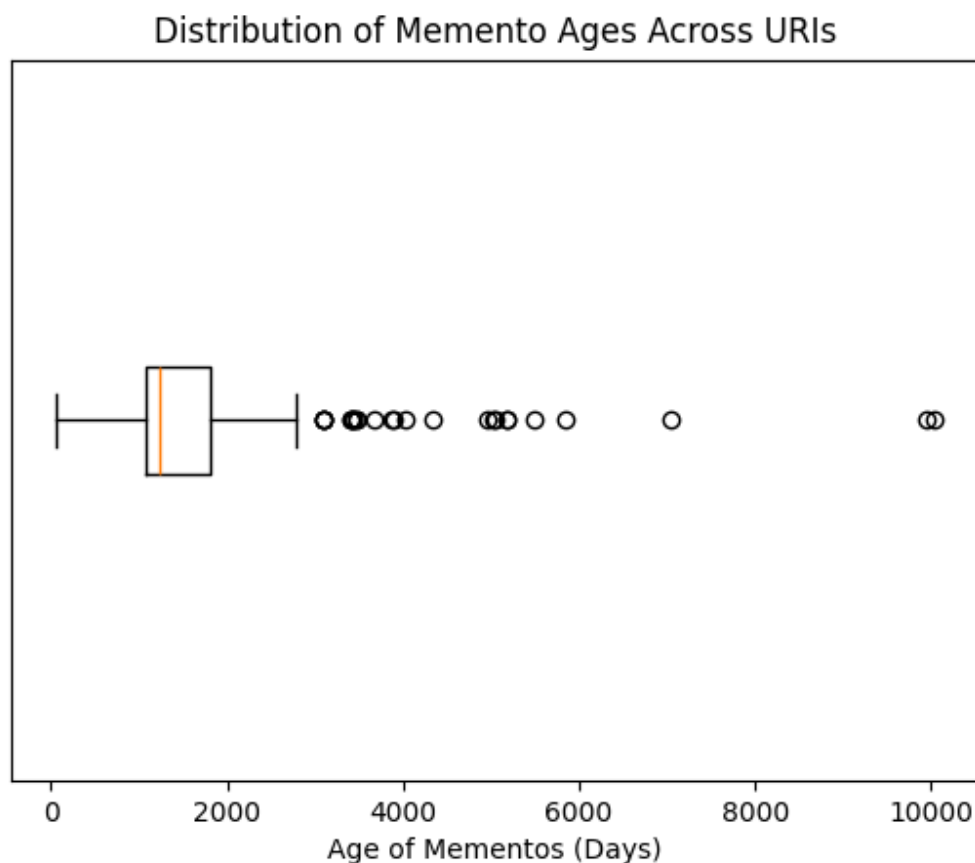


**Figure 1:** Distribution of Mementos Ages Across URIs

The results indicate that the majority of mementos for archived URI-Rs are between 1,000 and 2,000 days old, with a few outliers exceeding 2,000 days.

*Q: What URI-R had the oldest memento? Did that surprise you?*

The URI-R that has the oldest memento from the URI-Rs that were analyzed was https://www.unc.edu, with the earliest memento date being 1997-04-27 05:36:51. I'm not that surprised by the result as the University of North Carolina Chapel Hill is a prestigious state university that is known for conducting research.
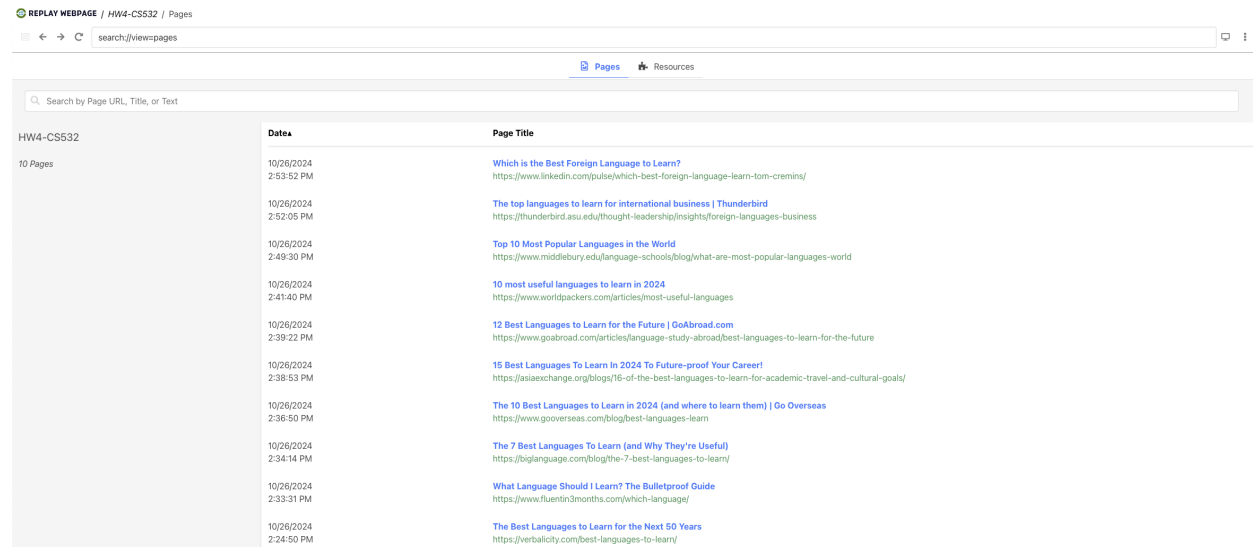
*Q: How many URI-Rs had an age of < 1 week, meaning that their first memento was captured the same week you collected the data?*

It was concluded that zero URI-Rs had an age of less than 1 week after analysis. All the URI-Rs that I looked at had mementos that were captured more than a week ago (as of 10/26/2024).

# Q2

Click here to access my Conifer public collection for Q2 of this assignment.

Figure 2 shows the list of archived pages as well as the browser address bar after uploading the WARC file (created by archiving 10 webpages using Conifer) to ReplayWeb.page (`https://replayweb.page/`).



**Figure 2:** ReplayWeb.page WARC file "Pages" Tab

*Q: Why did you choose this particular topic?*

I am a big language fan, and I speak Spanish self-taught. I want to learn a new language, but I have the hardest time choosing the next language to learn. Therefore, I chose to archive websites discussing the best foreign languages to learn.

*Q: Did you have any issues in archiving the webpages?*

For some websites, when I started the archiving process, a 404 error message would appear, and I couldn't archive the website that I wanted to. Otherwise, no issues. While capturing, I would scroll to the bottom of every webpage so the whole page could be captured.

*Q: Do the archived webpages look like the original webpages?*

For the most part, the archived webpages look quite similar to their original webpages. Sometimes there is a minor formatting difference, but nonetheless very similar.

*Q: How many URLs were archived in the WARC file? How does this compare to the number of Pages?*

A total of 1,140 URLs were archived in the WARC file, corresponding to 10 webpages. This indicates that each archived webpage contained a diverse array of resources, including images, scripts, and other types of content.
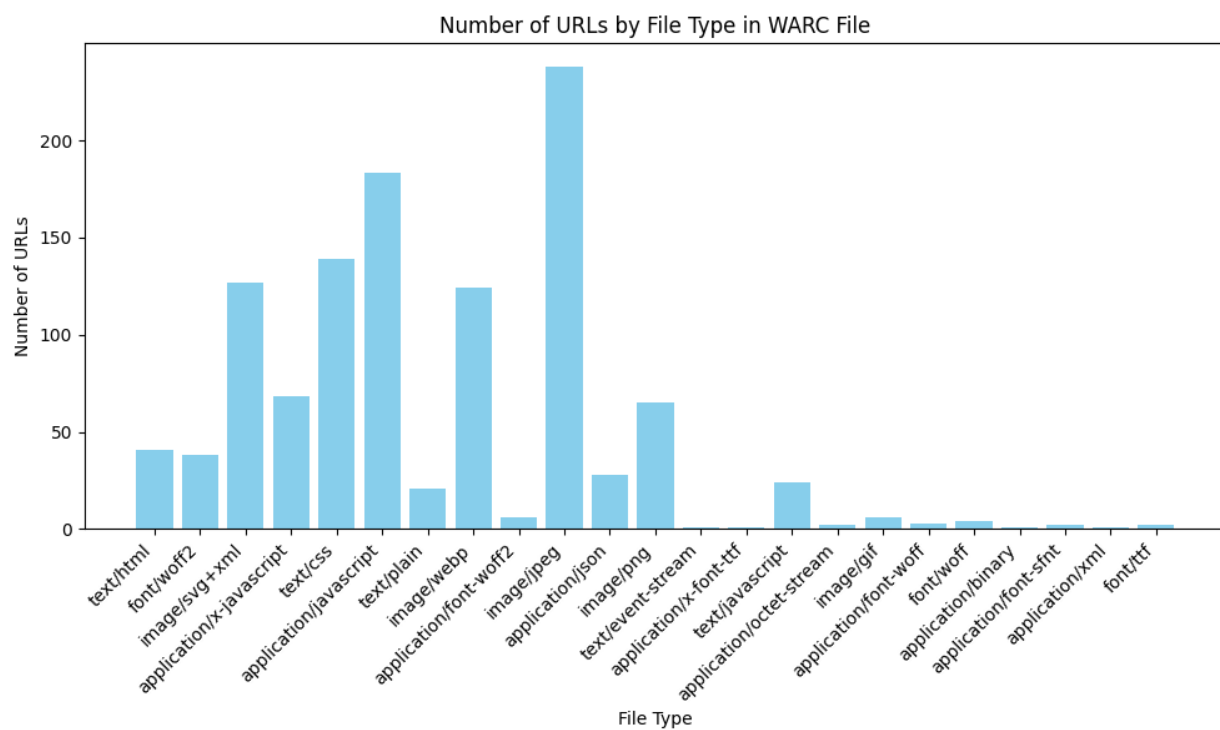


**Figure 3:** Number of URLs by File Type in WARC File

*Q: Which file type had the most URLs? Were you surprised by this?*

JPEG was the file type that had the most URLs. This doesn't surprise me as many web pages contain several JPEG image files.

# References

- Counters in Python, `https://www.geeksforgeeks.org/counters-in-python-set-1/`

- Python — datetime.timedelta() function, `https://www.geeksforgeeks.org/python-datetime-timedelta-function/`

- Warcio, `https://github.com/webrecorder/warcio`