# IT497 Week 3 Lab

## Web Scraping using R and `rvest`
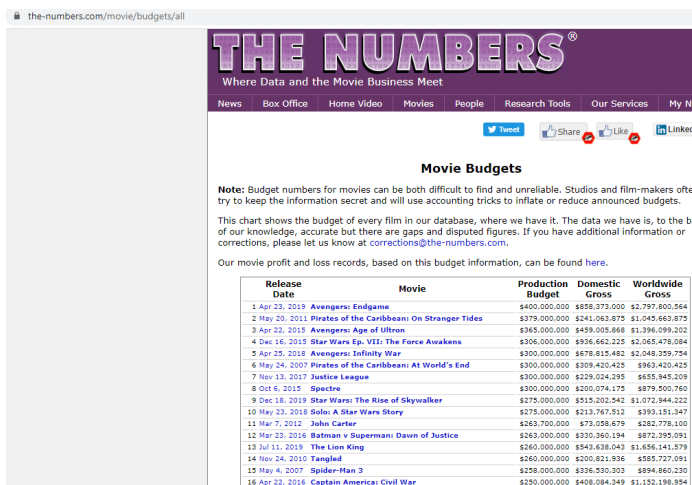
### About the Data

Nash Information Services' (https://www.the-numbers.com/) publishes revenue estimates for past movies, analyses for future releases, and many other reports that can be critical to planning an investment or creating a compelling business plan for a movie.

For our lab, we will be using Nash's Movie Budgets (https://www.the-numbers.com/movie/budgets/all)

This chart shows the budget of every film in Nash's database.



For our lab, we are only going to scrape the top 100.

**Note: for extra credit, you may scrape additional movies. Instructions for the extra credit follow the lab instructions.**

# Following are the steps you need to follow:

Step 1. Open R Studio Server and go to `File -> New File -> R Notebook`



**Step 2. R Studio will create a template file like this. Read through the text and pay special attention to everything in blue. That is how to "Run" code chunks and how to "Insert Chunk".**

**Step 3. Delete everything below the 3 dashes:**



**Step 4. Change the title to `IT497 Lab 3` and save the file as `yourlastnamelab3.Rmd`**



**Step 5. Insert a new R Code Chunk**

**If you have not already installed the `tidyverse`, do it in the 1st code chunk like this:**

```
# Install from CRAN

install.packages("tidyverse")
```

**The tidyverse packages are the packages at heart of data science/research using R. We will learn much more about many of these packages throughout the semester. However, if you want to do a little reading on your own, a great place to start is here:**
https://www.tidyverse.org/

**If you have not already, you will also need to install `rvest`.**

```
install.packages("rvest")
```

rvest helps you scrape information from web pages. It is designed to make it easy to express common web scraping tasks, inspired by libraries like beautiful soup. For more information about rvest, please visit: https://github.com/tidyverse/rvest

A nice rvest tutorial can be found here: https://www.datacamp.com/community/tutorials/r-web-scraping-rvest

**Please note that you only need to do this ONE TIME THIS SEMESTER.**
After this time, you will just load the tidyverse and revest packages using the library command. Insert the following R code into your code chunk. In code chunk #1, you should have the following R code:

```
library(tidyverse)
library(rvest)
```

Step 6. Insert a new R Code Chunk. In code chunk #2, you should have the following R code:

```
base_url <- "https://www.the-numbers.com/movie/budgets/all"
base_webpage <- read_html(base_url)
```

Step 7. Insert a new R Code Chunk. In code chunk #3, you should have the following R code:

```
new_urls<- "https://www.the-numbers.com/movie/budgets/all/%s"

table_base <- rvest::html_table(base_webpage)[[1]] %>%
  tibble::as_tibble(.name_repair = "unique") # repair the repeated columns
```

Step 8. Insert a new R Code Chunk. In code chunk #4, you should list the first 10 movies (you know how to do this).

Step 9. Insert a new R Code Chunk. In code chunk #5, write your data to a *.csv file. You will submit this file via ReggieNet.

In code chunk #5, you should have the following R code:

```
write.csv(table_base,"moviesData.csv")
```

**YOU SHOULD BE ABLE TO KNIT YOUR RMD FILE TO EITHER HTML OR WORD.**

# Extra Credit:

There are <u>2 extra credit opportunities </u>this week

 You may do both to earn up to 5 points extra credit,

**Up to 3 additional extra credit points**
In this week's lab, we scraped 100 movie records. Our lab was based off of this R Blogger's Tutorial: https://www.r-bloggers.com/tutorial-web-scraping-of-multiple-pages-using-r/

**The bottom of the tutorial explains how to read 100's of additional movie records.**

For up to 3 additional extra credit points, scrape an additional 600 movie records. Submit the 600 records following the lab's instructions (there should just be more records).

**Up to 2 additional extra credit points**
This video https://vimeo.com/110804387 explains how to use a Word Template to format your knitted output.

For up to two points, create a word template that makes the R code in your knitted word document **9-point Consolas.** Using the template when creating your output and submit your template via ReggieNet along with the other lab materials.

Deliverables: Upload the following to ReggieNet before the due date/time.

1. Your `moviesData.csv file`
2. `Your *.Rmd file`
3. `Either a Word or HTML file knitted from your *.Rmd file`
4. `Any extra credit files you want to include,`