LESSON 15 : Part 01 of 04 in the Visualizing Data module

# Google Sheets: Scraping data from the internet

Build your own data sets using Google Sheets.

# Lesson overview

## Learn to build your own data sets using Google Sheets.

There is a massive amount of data available on the internet that you can use to research and visualize stories. Finding the data, and getting it into a format you can work with is the first step.



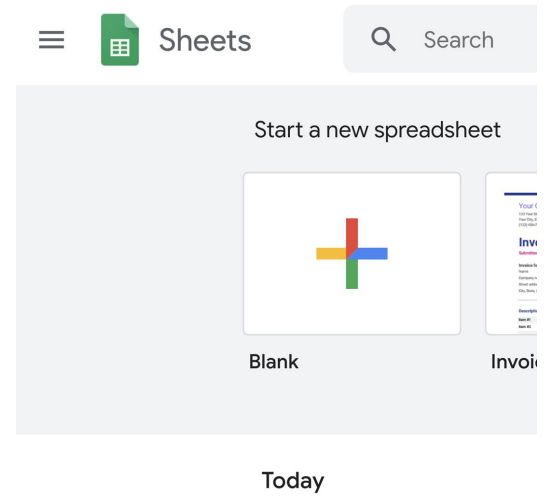| 1 | Starting a new spreadsheet. |
| 2 | Finding reliable data. |
| 3 | Importing data to Google Sheets. |
| 4 | Troubleshooting and error messages. |
| 5 | Displaying your data. |

For more Data Journalism lessons, visit:

newsinitiative.withgoogle.com/training/course/data-journalism

# Starting a new spreadsheet.
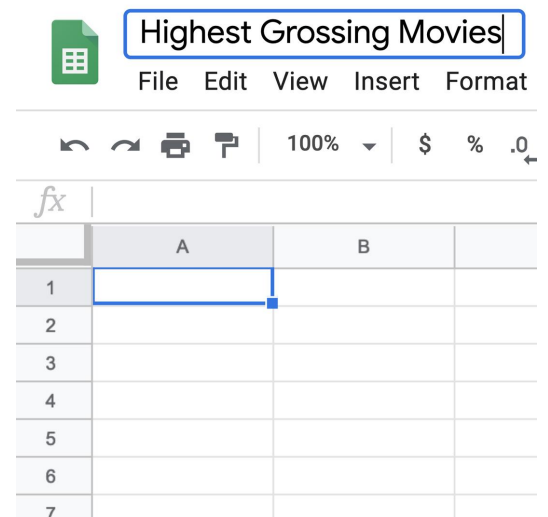
First, you need to create a blank spreadsheet. Go to sheets.google.com. Under **Start a new spreadsheet**, click the **+** icon.

To name your spreadsheet, click the text in the top left corner. Let's name this one "Highest Grossing Movies."

# Finding reliable data.

By sourcing data from government sites, scientific publications, Wikipedia, Google Public Data Explorer and more, you can tell data stories on almost any topic. In this lesson, we'll practice with data about movies.

---

Go to **google.com** and **search highest grossing films**. One of the first links should be a Wikipedia entry with multiple tables. One list, called "the top 50 highest-grossing films of all time" cites multiple references, so we will use that one. Always check to make sure you're scraping data from reliable sources.



---

To import this table to Google Sheets, copy the address of the Wikipedia page by highlighting the URL, right clicking on it, and selecting **copy**.
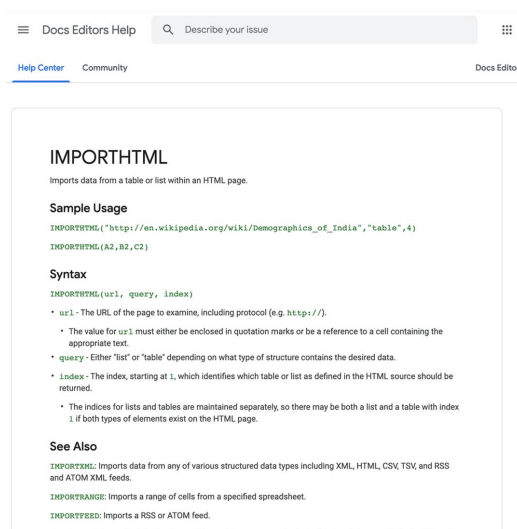
# Importing data to Google Sheets.

We'll use  importHTML to import the table from Wikipedia to our spreadsheet. This powerful formula is built into Google Sheets to help you import tables or lists from web pages. To learn more about how importHTML works and see examples, read the Google Sheets documentation pages.

The `importHTML` tool needs three parameters to work:
1) a URL
2) the type of data we're collecting, either a table or list
3) the number representing the position of the table or list in the HTML code.
In this example, the first instance of a table would be numbered as one, as the table we want is the first one that shows up in the HTML. You can use trial and error to find what the position of the table is (1, 2, 3, etc.) or right click the webpage, select **Inspect** > **Find** to locate the table in the code.

Go  to the blank sheet you created and navigate to cell A1. Type:

```
=importHTML("https://en.wikipedia.org/wiki/List_of_highest-grossing_films", "table", 1)
```

Notice that the URL and the element type (in our case, table) go between quotes — this will make the parameters green. The last parameter is a number not within quotes and it will be colored blue.
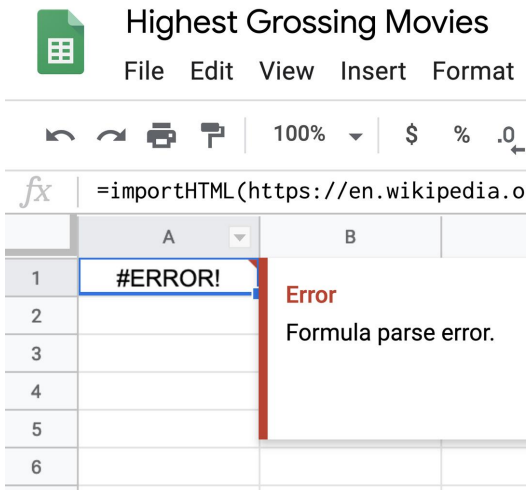
```
g_films", "table", 1)
```

```
_films", "table", 1)
```

# Troubleshooting and error messages.

If you get an ERROR! Message, check to make sure the quotes are double quotes as shown in the example.

If you get a VALUE! error, check to make sure you don't have extra parentheses or quotation marks in the cell.

# Displaying your data.

Once your ImportHTML formula is correct, press **enter** and give Google Sheets a couple of seconds. The table should load with all the rows and columns formatted.

Notice that there are some elements we need to remove so that we can visualize this data. We will learn this in the next lesson, "Google Sheets: Cleaning data."

# Congratulations!

You completed "Google Sheets: Scraping data from the internet."

To continue building your digital journalism skills and work toward Google News Initiative certification, go to our Training Center website and take another lesson:



## Google Sheets: Cleaning data

Prepare your data for analysis and visualization.

For more Data Journalism lessons, visit:

newsinitiative.withgoogle.com/training/course/data-journalism