https://schoolofdata.org/extracting-data-from-pdfs/

Lesson 04:
Data Scraping

# Extracting Data from PDFs



Professor:
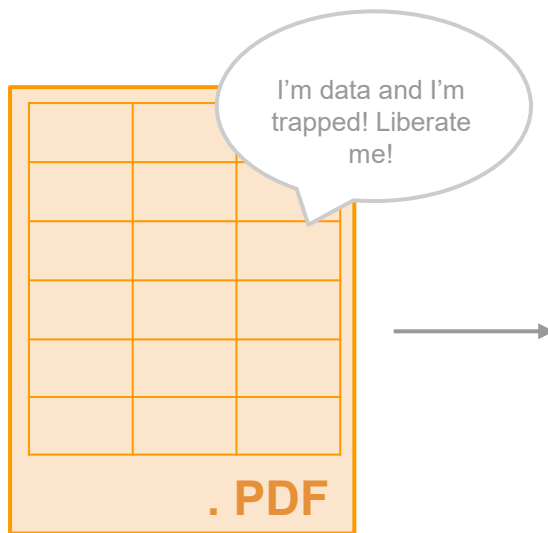Amr Eleraqi

aeleraqi@aucegypt.edu

# Unlock it!

**Your problem**
PDFs are great for humans to read. But not for machines. When researchers want to analyze the data trapped in a PDF, what tools can we recommend?

**Your solution**
Tabula



PDF (It's a trap!)



https://tabula.technology/

# Unlock it!

Tabula magic:



1. upload file

**Tabula**
**Processing File**
9%: generating page thumbnails....



2. select stuff



3. download data!

## About Tabula

- The program is available in different versions that are compatible with various operating systems.

- It requires an internet connection.

- No one can view the files and documents that are loaded inside the program, nor does the program keep a copy of them.

- The program works well with both structured and unstructured data.

- The program handles data effectively in Arabic language.

**Pros**

- Free, open source.

- Easy to install and use.

- Pretty accurate, especially for PDFs with simple table layout.

- Also has options for Python or R.

**Cons**

- No OCR, so that must be done separately.

- Had some issues with complex table layout.

- Minimal documentation.

- Not in active development.

## Installing Tabula

Here are the instructions to download and install Tabula:

*Ensure Java is installed on your computer. You can download Java here: https://www.java.com/en/download/

- Open the Tabula website: http://tabula.technology/
- Download the version of Tabula for your operating system

# Tabula



Tabula is a tool for liberating data tables locked inside PDF files.

View the Project on GitHub
tabulapdf/tabula

| Download for **Windows** | Download for **Mac** | View source on **GitHub** |
|---|---|---|

**Current Version:** 1.2.1

**Other Versions:** pre-releases & archives

**Need help?** Open an issue on Github.

**Donate:** Help support this project by backing us on OpenCollective.

We'd love to hear from you! Say hi on Twitter at @TabulaPDF

## Latest Version: Tabula 1.2.1

June 4, 2018

Tabula 1.2.1 fixes several bugs in the user interface and processing backend. (You can read about all the changes in the release notes.)

Download Tabula below, or on the release notes page.

Special thanks to our OpenCollective backers for supporting our work on Tabula; if you find Tabula useful in your work, please consider a one-time or monthly donation.

## How Can Tabula Help Me?

If you've ever tried to do anything with data provided to you in PDFs, you know how painful it is — there's no easy way to copy-and-paste rows of data out of PDF files. Tabula allows you to extract that data into a CSV or Microsoft Excel spreadsheet using a simple, easy-to-use interface. Tabula works on Mac, Windows and Linux.

## Who Uses Tabula?

Tabula is used to power investigative reporting at news organizations of all sizes, including ProPublica, The Times of London, Foreign Policy, La Nación (Argentina), The New York Times and the St. Paul (MN) Pioneer Press.

Grassroots organizations like SchoolCuts.org rely on Tabula to turn clunky documents into human-friendly public resources.

And researchers of all kinds use Tabula to turn PDF reports into Excel spreadsheets, CSVs, and JSON files for use in analysis and database applications.
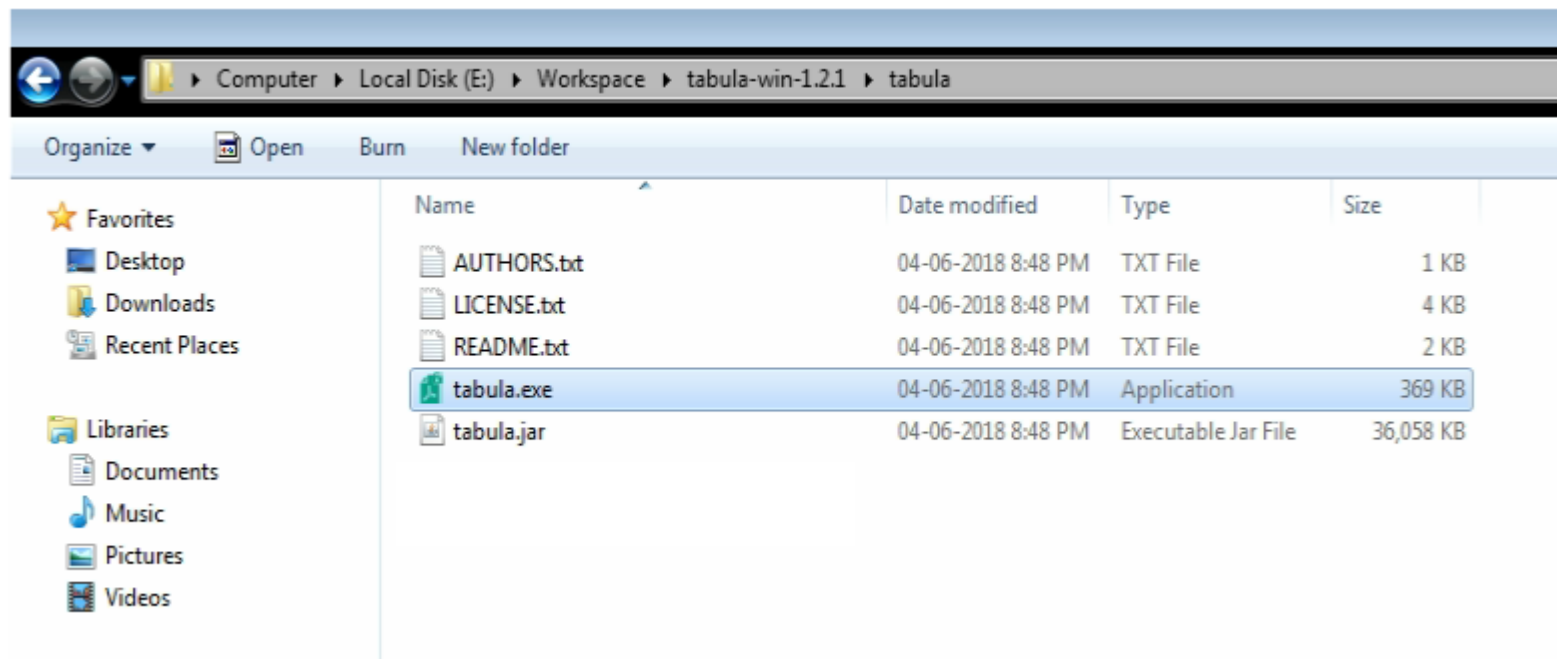
## Download & Install Tabula

Windows & Linux users will need a copy of Java installed. You can download Java here. (Java is included in the Mac version.)

1. Download the version of Tabula for your operating system:
    - **Windows:** tabula-win.zip
    - **Mac OS X:** tabula-mac.zip
    - **Linux/Other:** tabula-jar.zip, view README.txt inside for instructions

2. Extract the zip file. (Instructions: Windows, Mac)
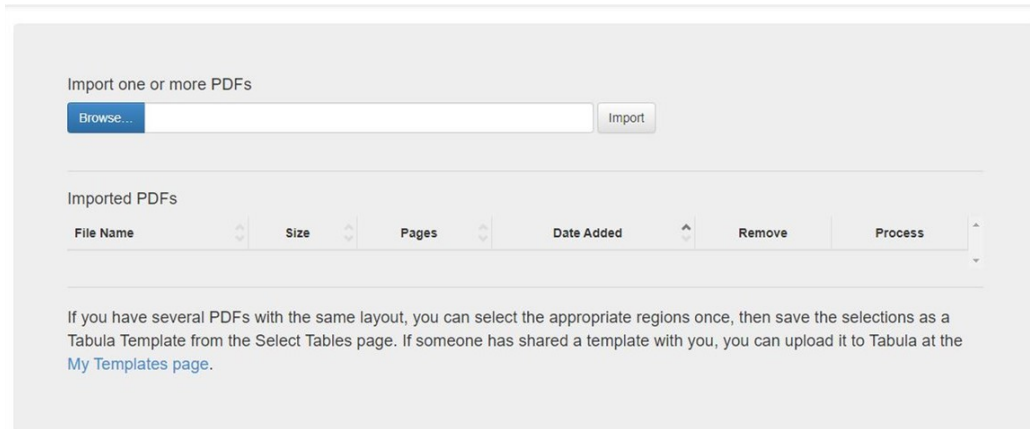3. Go into the folder you just extracted. Run the "Tabula" program inside.

## Installing Tabula

- Tabula downloads as a zip file on your computer. Extract the downloaded zip file – this creates a folder called "tabula" on your computer.

- Go into the "tabula" folder. Run the tabula.exe program inside. A control window may open; allow this window to run.

- Next, a web browser will open – this is Tabula. If your web browser does not open, use your web browser to go to: http://localhost:34555

**Installing Tabula**

- After running the program, a web browser will open – this is Tabula. If your web browser does not open, use your web browser to go to: http://localhost:34555
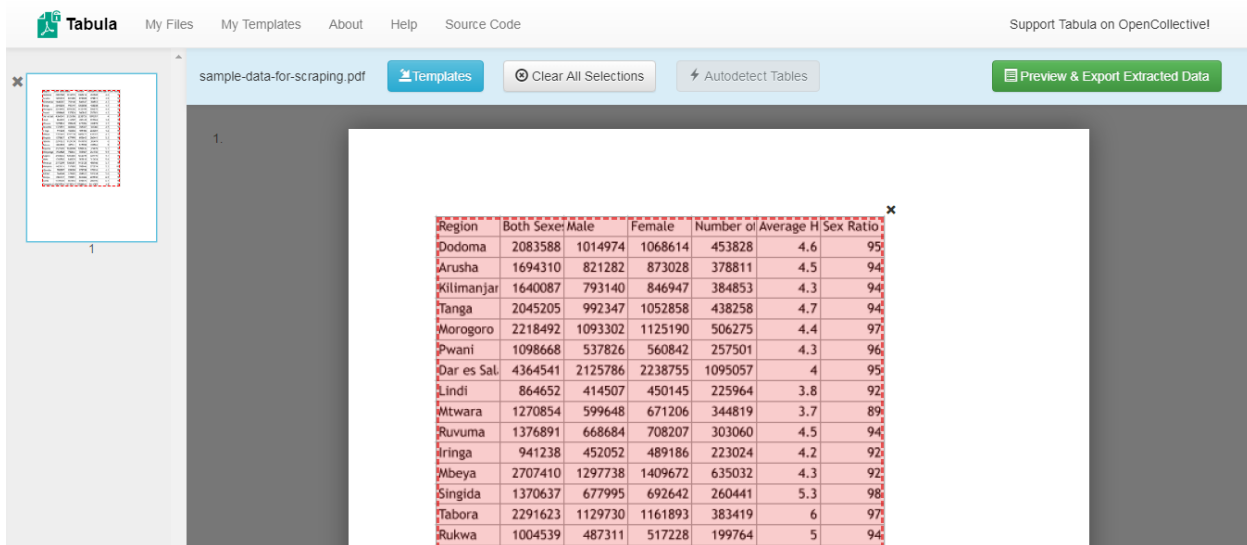
- There's Tabula!

## Tabula in Action

**1. Upload PDF File:** In home screen you will find file selection option where you need to browse and upload PDF file from which you want to extract data. After selecting the file, click on the Import button. After submission, you will be shown uploaded PDF file as shown in the screenshots below:





Depending on the size of the PDF, it may take some time for the file to upload.
Be patient!

## Tabula in Action

**2. Make table selection:** Now you need to move to the page from where you want to extract table data, then select the table by clicking and dragging to draw a box around the table. *You can also click on "Autodetect Tables" option which will select the tabular data automatically.*

# Tabula in Action

**3. Preview Data:** After table selection "Preview & Export Extracted Data" button will be enabled. Click on that button to preview the data. Preview is shown in the below screenshot:

**Double-check your data by cross-referencing your table**

Double-check your Tabula preview of your table with the original PDF. We use another program, like Preview or Adobe Acrobat, to compare. This way, you'll make sure no data has been lost or misread.

# Tabula in Action

**4. Export Data:** Now you can copy data to the clipboard and paste it to anywhere you want or you can export the data to a variety of file formats like CSV, TSV, JSON etc. The following screenshot shows exported CSV file:
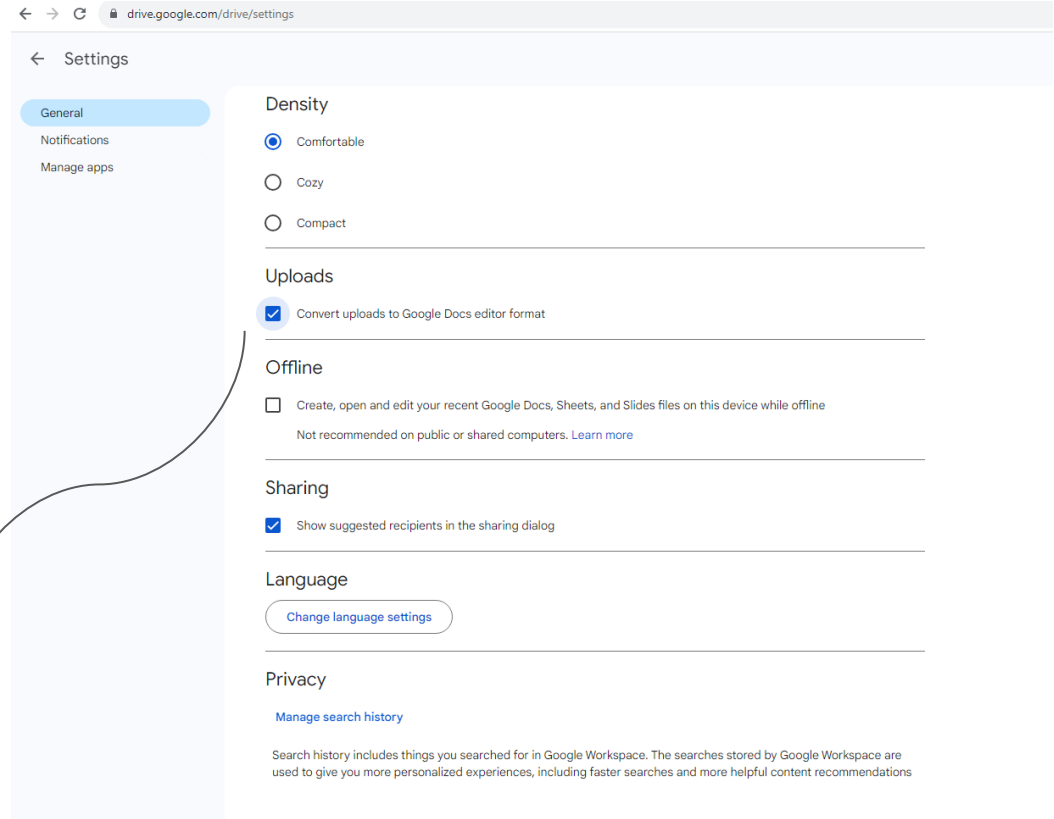
## Extract text from image in Google Drive

Google Docs has a powerful Optical Character Recognition feature built right in..

## Enabling the feature

Open up Google Drive and then click on the gear icon. From the drop-down, click Settings. In the resulting window, make sure Convert uploads is checked.
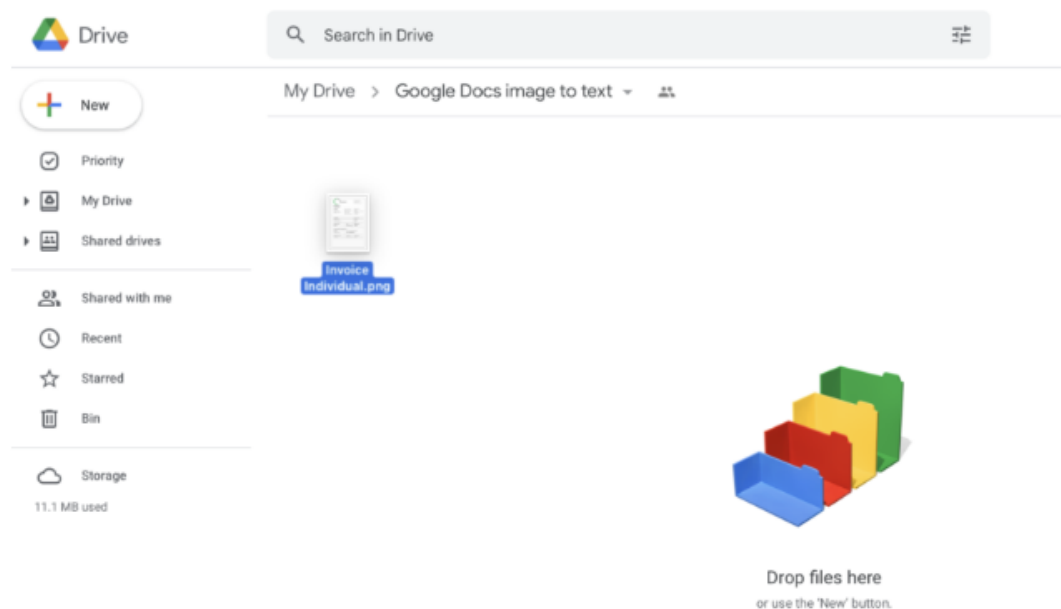
## Extract text from image in Google Drive

### Using the feature

Upload either an image to Google Drive. The uploaded image doesn't automatically convert. In fact, it will remain exactly as uploaded.

*It is important to remember that Google Drive supports documents with .jpg, .png, .gif extensions and PDF files with a maximum size of 2 MB.*

## Extract text from image in Google Drive

**Using the feature**

Once the file is in your Google Drive account, right-click it and select Open with | Google Docs.

After selecting to open your file with Google Docs, wait for a bit and a new document will open, containing the newly converted text.