# Web Scraping and Tools for Scientific Programming
## Computational Data Analytics for Economists

Helge Liebert

University of Basel

March 10, 2016

# What is this lecture about?

- Economics is a computational social science.
- Most research in economics now involves scientific programming.
- Introduce tools and ideas that may make daily research tasks easier.
- Many of these have been developed by other scientists or IT professionals.
- Some are almost as old as computers, others are novel.
- Focus is going to be on working with data.

# Economics and computer science

- Data management is not taught in introductory econometrics.
- Computer science often involves processing data.
- Any common problem you are likely to encounter has been solved.
- CS offers tools and concepts that economists can profit from.
  - ▶ Tools: Remote servers, databases, version control, acessing APIs for data, text processing and analysis, geospatial analysis, OCR, automation of virtually anything (e.g. experiments).
  - ▶ Concepts: Time complexity of algorithms, computational cost.

# Topics

- General points about scientific programming and working with computers.
- Unix tools and the command line.
- Accessing data on the web: APIs.
- Gathering (unstructured) web data and transforming it into structured data ("*web scraping*").
- Use processing of web data to introduce various ideas along the way.
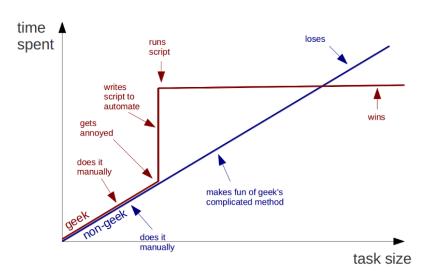
# Goal: Understanding the tools available

- Finding out what to look for is often harder than understanding how a particular tool works.
- Point you towards the resources and explain their general concepts.
- Sometimes these methods may offer the only feasible solution to gather data.
- They may also help you solve tasks efficiently.
- Knowing that a problem can be solved and how is worth a lot.
- Leave you marginally more computer literate.

# Goal: Automation

- Digitalization offers exciting data for research. But: *Data is messy*.
- Gathering or processing data often involves repetitive manual tasks.
- Disadvantages:
  - ▶ Manual tasks are often not well documented or reproducible ex post.
  - ▶ Manual work is frustrating and a huge time-sink.
  - ▶ Manual work may not be feasible with large data.

➥ Automation helps!
  - ▶ Frees you to engage in other work.
  - ▶ You learn new things.
  - ▶ Should you encounter the same class of problem in the future, you already have a solution at hand.
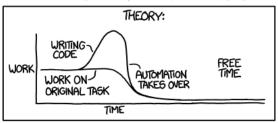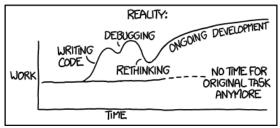
# Automation



**Geeks and repetitive tasks**

time spent ↑

task size →

runs script

loses

writes script to automate

wins

gets annoyed

does it manually

makes fun of geek's complicated method
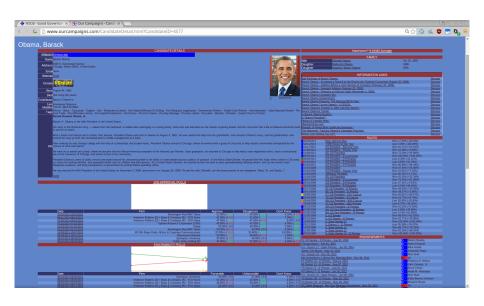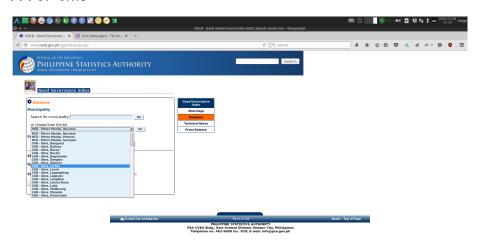
geek

non-geek

does it manually

# Automation

# Example problem

- How to turn unstructured into structured data?

- Consider a situation where
  - You want to get data from the internet.
  - The data is in unstructured/semi-structured form.
  - You want to transform it into a differently structured format for further use.
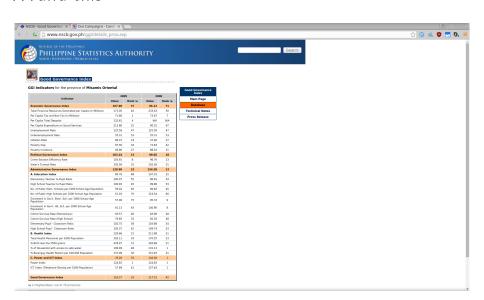  - You need to filter the available information.

# From this ...

# . . . or this

# . . . and this

# . . . to this.

# Which language to choose?

- *Anything* can be done in *any* language. Convenience varies.
- Concepts and toolkits transfer easily most of the time.
- Trade-off: Prior knowledge vs. task suitability.
- Never re-invent the wheel.

- Choose a high-level, dynamic language. Ideally free and open source.
- Specialized languages: R, MATLAB, Octave, Gauss, Julia, . . .
- General-purpose languages: Python, Perl, Ruby, . . .
- Choice is use case- and taste-specific. Popular is typically better.

# R

- R is the major statistical programming language.
- It is free, used in many sciences and in industry. Good documentation.
- New models are typically first published and implemented in R.
- Having data processing and analysis in the same language is nice.
- Good library support for common tools (e.g. databases, regular expressions).
- Specific tasks for which high-level wrapper functions are not available may be very cumbersome.
- In recent years, R development has been very active and libaries exist for almost anything.

# Python

- General-purpose programming language, supports object-oriented programming.
- Reads like english. Explicit and clear. Whitespace matters, no braces. (*"There should be one obvious way to do it"*.)
- Used extensively in industry and sciences. Good documentation.
- Libraries for almost anything.
- Many science-related libraries exist for other languages, but rarely are they as mature.
- Less support for statistical modeling (but growing).
- Less suited for interactive data work.

# Recommendation

- Research ex ante which libraries are most mature and best for solving your specific problem.

- Focus on getting things done.

- Rule-of-thumb:
  - Simple data processing:
    Stick with R. Augment with other tools where required.
  - More involved projects:
    Go with Python. You can still analyze data in R.

- I am proposing a mix of R, Python and Bash (Unix-Shell),
  . . . and whatever program your co-authors are using.

- R, Python and SQL are highly valued on the job market, knowing your way around a terminal is useful.

- This course uses R, but I will provide some equivalent python code.

# Why not Stata, Matlab, Gauss or similar?

- Advantage: Many domain-specific models supported.
- Less support for almost anything else.
- Much less flexible for anything not to do with data analysis or numerics.
- Difficult to deploy on a server. Often tied to a GUI.
- Less popular, smaller userbase. Proprietary and expensive.
- You can still rely on them for estimation after your data is clean.

# Why not Perl or Ruby?

Perl

- *"There's more than one way to do it."*
- Lots of special cases, reliance on hidden magic, bad readability.
- You may want to work together with somebody else.
- You may want to be able to understand your own code in a few months time.
- Less popular in sciences.

Ruby

- Everything is an object. Really. Intuitive for some people.
- Different focus.
- Even less popular in sciences.
- Less support.

# A few things to get started

- What you need for this course:
  - R.
  - A text editor or an IDE (like RStudio).

- You want to use Python:
  - A Python distribution (use Anaconda) (and possibly a shell).
  - A text editor or an IDE (like Spyder).
  - Which version, 2.7 or 3?
    - Python 3 if you plan to use Python regularly in the future.

- You want to use the command line interface (CLI) and have access to shell tools:
  - Terminal and Bash (Linux, MacOS), package manager (default on Linux, use iterm2 and homebrew on MacOS).
  - Cygwin, Windows Subsystem for Linux, Linux in a virtual machine, dual boot (Windows).

# A note on operating systems

- MacOS or Linux offer built-in access to a Unix shell (Bash).
- Further software is managed via a package management system and distributed via software repositories.
- On Linux, use your package manager to install anything you require.
- On MacOS, familiarize yourself with Homebrew. Install iterm2 if you want a fancier terminal.

- For Windows, many tools are not available or cumbersome to use. Dependency resolution can be a nightmare.
- Windows does not provide proper access to a Unix shell.
- Even reliably installing Python was a chore until recently (now use Anaconda.)

# Command line interpreters and shells

- An interface that lets you interact with your computer.
- A CLI using a programming language that allows you to execute programs and scripts.
- Unix-based operating systems (Linux, MacOS) have Bash pre-installed.
- Windows has cmd (or PowerShell). These are not a viable replacement. Cygwin or WSL may be. Git Bash is incomplete.

- Some examples:
  ```
  cd somedir/subdir # navigate to a folder
  cd .. # navigate to parent directory
  cd # return to your home folder
  ls # list directory contents
  R # start the R console
  ```
- *CTRL + c* aborts a process, *CTRL + d* quits.
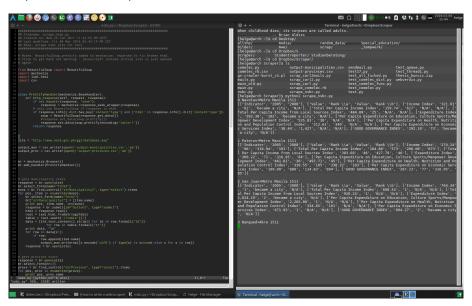
# Examples

- Some examples:
  ```
  vim myscript.r # edit your R script with vim
  R -f myscript.r # execute your R script
  python myscript.py # execute your python script
  git add myscript.py # stage file for version control
  git commit myscript.py # stage file for version control
  man ssh # display manual pages for the ssh program
  ssh myusername@13.438.14.673 # secure shell login to your remote serve
  ```

- Sounds tedious? It is. But it can also be extremely powerful.

- Convert all your pdf files in a folder to text and search them.
  ```
  for file in *.pdf; do pdftotext "$file"; done
  grep -icr "keyword" *.txt
  ```
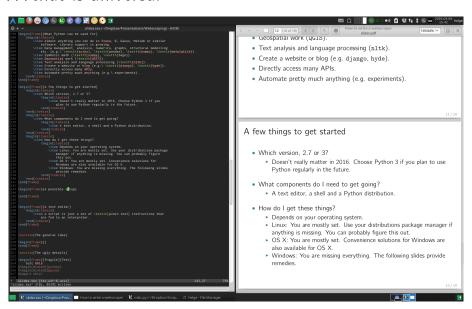
# A note on text editors

- A script is a set of *plain text* instructions, fed to an interpreter.
- R scripts usually have the suffix `.r`, Python `.py`, Shell `.sh`.
- Much of our work involves working with text files.
- *Some* text editor is required. A *good* text editor makes working with text much easier and faster.

- Too many options to list. All are better than Notepad.
- A few suggestions: VS Code, Sublime Text, Atom, Notepad++.
- Learning Vim or Emacs requires you to invest some time.

- Text editors allow you to integrate your work.
- Sometimes IDEs with GUI may be more convenient.
- Features: Efficient text editing, syntax checking, completion, . . .

# A possible setup

# ...that is universal

# Rstudio

# Spyder

# Version control

- Ubiquitous in IT, extremely useful in many purposes.
- Somewhat less useful for statistical data management and analysis due to different workflow.
- Still good to know about and understand the basics.
- Git is the dominant version control software today.
- ProGit is a good and free resource. Skim the first few chapters.
- Sufficient to grasp the concept and know the basic commands.

# Getting started—things to consider before you begin

- Pick up the phone and try to get the data directly.

- Has somebody else done what you intend to do?

- Search if somebody has already faced the same or a similar problem.

- Does the site or service provide an an API that you can access directly?

- Has anybody written a wrapper for it?

- Is the website only online for a limited time? Do you want an original snapshot as a backup? Is it more convenient to filter your data offline?

# Static vs. dynamic websites

Static Website

HTTP request

HTTP response

---

Dynamic Website

HTTP request

HTTP response

Server-side processing

# Web APIs

- Data providers often offer Web APIs (*Application Programming Interface*) to access data.

- Allow programmable access to data via a defined set of HTTP messages. Similar to visiting a website: you specify a URL and information is sent to your machine.

- With a website, you receive code interpreted by your browser (HTML, CSS, JavaScript). With an API, you receive data.

- Usually in JSON (*JavaScript Object Notation*) or XML (*Extensible Markup Language*) format.

# Web APIs

- Often just two steps:
  1. Construct the URL query that serves as the API request.
  2. Process the response message the API sends back.
- Examples:
  - `https://api.kivaws.org/v1/loans/newest.json`
  - `https://api.kivaws.org/v1/loans/search.json?sector=Agriculture& country=VN`
  - `https://www.theyworkforyou.com/api/getMPs?&key=someapikeyhere& output=js`
- Libraries may offer wrappers for APIs: `WDI`, `wbstats`, `twfy`, `pvsR`, Google Maps, OpenStreetMap/OSRM, . . .

# Save an offline copy

- Use the shell utilities `wget` or `curl` to download the complete site.
- Also useful if you just want a set of files (e.g. pdf documents) from the same site directory.
- Convenient for static sites of limited size.
- Infeasible for large sites or sites that create content dynamically.

# Examples

- **Simple http GET request.**

  `wget http://www.google.com`

- **Recursively download a website.**

  `wget -r http://www.some-site.com/some-subdir/`

- **Download all pdfs from a site.**

  `wget -r -A.pdf http://url-to-webpage-with-pdfs/`

- **Mirror a site offline and convert links for local browsing.**

  `wget --mirror -p --convert-links -p ./local-dir`
  `    http://target-website.com`

# HTML and the Document Object Model

- Extracting information from the web requires a basic understanding of HTML and the associated Document Object Model (DOM).

- HTML elements provide the structure and content of web pages.

- Typically consist of *start* and *end* tags, with content in between.

    `<tagname>Content here</tagname>`

- A page consists of nested elements.

- The `html` element is the outer-most element, nesting the `head` and `body` elements, which in turn have nested elements.

- Nesting structure of elements can be represented by a tree (DOM).

# Document Object Model

- The DOM is a programming interface for HTML and XML documents.

- Provides a structured representation of the document.

- A document as a group of nodes, each node representing a part of the document.

- Allows programmatic access to the tree to change the structure, style and content of the document.

- Connects web pages to scripts or programming languages.

# A simple HTML page

A simple HTML page:

```html
<html>
  <head>
    <title>My Web Page</title>
  </head>
  <body>
    <h1>Welcome To My Web Page</h1>
  </body>
</html>
```

How a browser renders this page:

# HTML and the DOM

A simple HTML page:

```
<html>
  <head>
    <title>My Web Page
    </title>
  </head>
  <body>
    <h1>Welcome To My Web Page
    </h1>
  </body>
</html>
```
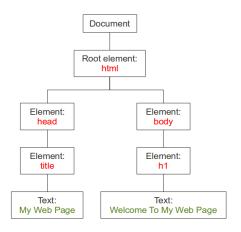
Corresponding node tree:

# Example: An HTML table element



- Tables are represented by a top-level `table` element.
- The `table` element nests `tr` (*table row*) elements.
- These nest `th` (*table header*) and `td` (*table data*) element cells.

# HTML and the DOM

- HTML tags can have attributes and text content.

  ```
  <tag attribute="value" attribute2="value">Text content.</tag>
  ```

- Example page:

  ```
  <html>
    <head>
      <title>My Web Page</title>
    </head>
    <body>
      <h1>Welcome To My Web Page</h1>
      <img src="picture.png" alt="example photo" width="100">
        <a href="pagelink.html" id="pageid">Check this other
        ↪  page.</a>
    </body>
  </html>
  ```
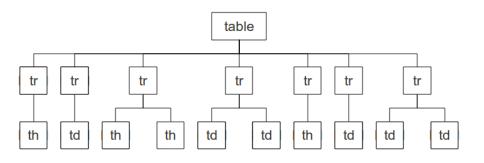
# DOM node trees

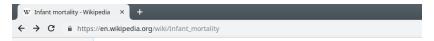- HTML DOM views a document as a tree structure called node tree.
- Everything in an HTML document is a node.
  - The entire document is a document node
  - Every HTML element is an element node
  - Every HTML attribute is an attribute node
  - Text content in the HTML elements is a text node
- Nodes can be accessed through the tree.
- Nodes may be assigned unique id attributes through the tree.

# Data from the web

# Wikipedia on infant mortality



Infant mortality - Wikipedia

https://en.wikipedia.org/wiki/Infant_mortality

## Epidemiology [ edit ]

*See also: List of countries by infant mortality rate*

For the world, and for both less developed countries (LDCs) and more developed countries (MDCs), IMR decli

However, IMR was, and remains, higher in LDCs. In 2001, the IMR for LDCs (91) was about 10 times as large
countries are, on average, much less than those among the more developed countries.[clarification needed]

A factor of about 67 separate countries with the highest and lowest reported infant mortality rates. The top and

| Rank | Country | Infant mortality rate (deaths/1,000 live births) |
|------|---------|--------------------------------------------------|
| 1 | Afghanistan | 121.63 |
| 2 | Niger | 109.98 |
| 3 | Mali | 109.08 |
| 4 | Somalia | 103.72 |
| 5 | Central African Republic | 97.17 |
| 218 | Sweden | 2.74 |
| 219 | Singapore | 2.65 |
| 220 | Bermuda | 2.47 |
| 221 | Japan | 2.21 |
| 222 | Monaco | 1.80 |

According to Guillot, Gerland, Pelletier and Saabneh "birth histories, however, are subject to a number of error

# Fetching a table from Wikipedia

```r
library(rvest)

# 1) fetch and parse the website
page <- read_html("https://en.wikipedia.org/wiki/Infant_mortality")
# 2) extract the html node containing the table
table <- html_node(page, xpath = "//*[@id='mw-content-text']/div/table[2]")
# 3) extract the table as a data frame
mrates <- html_table(table)
```

# Inspecting the HTML source

- Convenient with modern browsers: Use the developer tools.
- Right-click *Inspect* (for Chrome there is also SelectorGadget).
- Look at the HTML source to grasp the structure.
- Find out how to navigate the site.
- Find the element(s) you want to extract.
- Get the Xpath expression or CSS selector to extract elements.

# HTML elements visualized

# Infant mortality rates from Wikipedia

```html
<table class="wikitable" style="text-align:left">
  <tbody>
    <tr>
      <th>Rank</th>
      <th>Country</th>
      <th>Infant mortality rate <br> (deaths/1,000 live births)</th>
    </tr>
    <tr>
      <td>1</td>
      <td style="text-align:left;"><a href="/wiki/Afghanistan" title="Afghanistan">Afghanistan</a></td>
      <td>121.63</td>
    </tr>
    <tr>
      <td>2</td>
      <td style="text-align:left;"><a href="/wiki/Niger" title="Niger">Niger</a></td>
      <td>109.98</td>
    </tr>
    <tr>
      <td>3</td>
      <td style="text-align:left;"><a href="/wiki/Mali" title="Mali">Mali</a></td>
      <td>109.08</td>
    </tr>
    <tr>
      <td>4</td>
      <td style="text-align:left;"><a href="/wiki/Somalia" title="Somalia">Somalia</a></td>
      <td>103.72</td>
    </tr>
    ...
  </tbody>
</table>
```

# Fetching a table from Wikipedia

```r
# fetch and parse the website
page <- read_html("https://en.wikipedia.org/wiki/Infant_mortality")
# list the table nodes
html_nodes(page, "table")
# using xpath expressions or css selectors is equivalent
table <- html_node(page,
xpath = "//*[@id='mw-content-text']/div/table[2]")
table <- html_node(page,
css = "#mw-content-text > div > table:nth-child(121)")
```
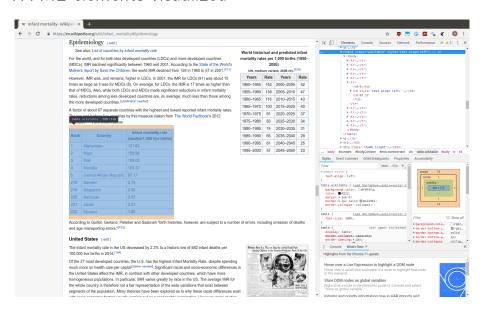
# The general structure

- There is no universal recipe. But most programs follow a certain structure.
  1. Open a website mimicking a browser and navigate it (optional).
  2. Get the page source HTML and feed it to a parser.
  3. Extract the elements you need.
  4. Filter and arrange them as needed and save them.
  5. Repeat 1.–4. until you have everything you want.
  6. Output your data.

# Navigating to another page

```r
# Open infant mortality page
session <- html_session("https://en.wikipedia.org/wiki/Infant_mortality")
# Goto page on Somalia
session <- follow_link(session, "Somalia")
# Read the source
page <- read_html(session)
# Extract html
table <- html_node(page, xpath = "//*[@id='mw-content-text']/div/table[4]")
regions <- html_table(table)
```

# Filtering links

```
# read wiki page
page <- read_html("https://en.wikipedia.org/wiki/Infant_mortality")
# get the links
wikilinks <- html_attr(html_nodes(page, "a"), "href")
# use regex to filter internal links:
#   select only articles, no files or category pages,
#   matching with mortality or somalia
links <- grep("^(?!.*:)(/wiki/.*Mortality)|(/wiki/.*Somalia)", wikilinks,
ignore.case = TRUE, value = TRUE, perl = TRUE)
links <- unique(links)
# go to first selected article page and process it
session <- jump_to(session, links[1])
page <- read_html(session)
html_nodes(page, "title")
```

# A note on regular expressions

- Regular expressions are sequences of characters that define a search pattern.
- Usually used for find/replace operations on strings, or for validation.
- Regexes are an *extremely* helpful tool.
- Easy to grasp, complex to master.
- Pin a cheatsheet to your office wall.
- But: Regular expressions are not parsers. Always use a dedicated HTML parser to extract elements.

# Trial project

- Phillippine Statistics Authority *Good Governance Index*.
- Available at `http://nap.psa.gov.ph/ggi/default.asp`.
- Get all GGI data tables for all municipalities.
- Save them in a local data file for further processing and analysis.
- Try for yourself. How would you go about this?

# General remarks

- Start simple and expand your program incrementally.
- Keep it simple. Do not overengineer the problem.
- Do not repeat yourself. Code duplication implies bug reuse.
- Limit the number of iterations for test runs. Use print statements to inspect objects.
- Write tests to verify things work as intended.
- If the web page cannot be navigated easily or has hidden javascript, look into Selenium (`library(rselenium)`).
- If you scraper requires complex monitoring/validation procedures or threading for performance, look into Python.

# Final remarks

- Sometimes a little program can help alot.
- Never lose sight of your goal. Your time is valuable.
- Do not engage in perfectionism, focus on GTD.
- Maybe identify everyday tasks that you can optimize.
- It might even be fun.

Appendix

# Why Python?

```
>>> import this
The Zen of Python, by Tim Peters

Beautiful is better than ugly.
Explicit is better than implicit.
Simple is better than complex.
Complex is better than complicated.
Flat is better than nested.
Sparse is better than dense.
Readability counts.
Special cases aren't special enough to break the rules.
Although practicality beats purity.
Errors should never pass silently.
Unless explicitly silenced.
In the face of ambiguity, refuse the temptation to guess.
There should be one-- and preferably only one --obvious way to do it.
Although that way may not be obvious at first unless you're Dutch.
Now is better than never.
Although never is often better than *right* now.
If the implementation is hard to explain, it's a bad idea.
If the implementation is easy to explain, it may be a good idea.
Namespaces are one honking great idea -- let's do more of those!
>>>
```

# What (else) can Python be used for?

- Almost anything you can do in Stata, R, Gauss, Matlab or similar software. Library support is growing.
- Data management, analysis, numerics, graphs, structural modelling etc. (e.g. `scipy`, `pandas`, `numpy`, `matplotlib`, `seaborn`).
- Symbolic math (e.g. `sympy`, `Sage`).
- Geospatial work (e.g. `QGIS`).
- Text analysis and language processing (e.g. `nltk`, `spacy`, `gensim`).
- Create a website or blog (e.g. `django`, `hyde`, `sphinx`).
- Directly access many APIs (e.g. Twitter).
- Automate pretty much anything (e.g. experiments, data collection).

- In recent years, R has been extended to many of these domains.

# Python resources

- Relevant modules.
  - ▶ `requests`, `bs4/BeautifulSoup`, `mechanize/mechanicalsoup`, `selenium`
  - ▶ `Scrapy` provides a complete framework for more complex projects.
  - ▶ `csv`, `re`, `pickle`, `pprint`, `pandas`, `random`, `itertools`, `pickle`, …

- Learning the language.
  - ▶ *A byte of Python* is free. *The Quick Python Book* or *Dive into Python* offer a denser treatment.
  - ▶ O'Reilly: *Learning Python/Programming Python/Pocket reference*, *Web scraping with Python*.
  - ▶ *Automate the boring stuff with Python* for inspiration.
  - ▶ Plenty of video lectures and courses online. Stackoverflow helps.
  - ▶ (For Git: *ProGit* is free and really all you need.)

- Read about basic types, syntax, look at a few examples, then just have a go.