# [301] Web 1

Tyler Caraza-Harter

# Learning Objectives Today

Network basics
- IP addresses
- host/domain names
- client/server and request/response

HTTP basics
- URLs
- GET/POST/etc
- headers
- status codes

Requests modules
- downloading data with requests.get
- remote calls with requests.post

# Learning Objectives Today

Motivation

Networking Basics

HTTP (Hypertext Transfer Protocol)

Requests Module

# Data Science and the Internet

There are tons of online sources of data
- Examples: https://tyler.caraza-harter.com/cs301/spring19/datasets.html

Wide range of topics
- healthcare
- roads and city planning
- astronomy
- population
- business
- entertainment
- education
- etc

# Data Science and the Internet

There are tons of online sources of data
- Examples: https://tyler.caraza-harter.com/cs301/spring19/datasets.html

Wide range of topics
- healthcare
- roads and city planning
- astronomy
- population
- business
- entertainment
- education
- etc

Why not just download data by hand?

# Motivation 1: too much data

What if you're analyzing language trends over time?
- Dataset: Project Gutenberg has 57K free books
- Too much work to download one by one

# Motivation 2: data doesn't always come in files

Many datasets are difficult to download complete

Instead, you can make function calls to servers
(we'll learn how) to grab specific data
- Dataset: OpenStreetMap
- You issue calls to get specific data:
  1. specify latitude/longitude rectangle
  2. specify structures of interest (e.g., bike paths)

# Learning Objectives Today

Motivation

Networking Basics

HTTP (Hypertext Transfer Protocol)

Requests Module

# Networking Basics

computer 1

computer 2

Computers communicate over a network (e.g., the Internet)
by sending messages to each other

# Networking Basics

request for data

computer 1

computer 2

Computers communicate over a network (e.g., the Internet)
by sending messages to each other

# Networking Basics

computer 1

request for data

100100101
011001010

sometimes requests
upload data too

computer 2

Computers communicate over a network (e.g., the Internet)
by sending messages to each other

# Networking Basics



response with data

100100101
011001010

computer 1

computer 2

Computers communicate over a network (e.g., the Internet)
by sending messages to each other

# Networking Basics

**server**

**client**

response with data

```
100100101
011001010
```

**computer 1**

**computer 2**

Computers communicate over a network (e.g., the Internet)
by sending messages to each other

# Networking Basics

computer 1

computer 2

**Challenge**: there are millions of computers.
How do we indicate which machine should get our request?

# How do we send a letter?



**1** lookup friend's address in phone book

**2** put address on the envelope

**3** trust postal service to get letter to that address

# Internet Protocol

computer 1

computer 2

**Solution**: every machine* has an IP address (Internet Protocol).
Requests are sent to a specific IP address.

*some machines have more multiple addresses

# Internet Protocol

address: 18.216.110.65

computer 1

computer 2

**Solution**: every machine* has an IP address (Internet Protocol).
Requests are sent to a specific IP address.

*some machines have more multiple addresses

# Internet Protocol

address: `18.216.110.65`

request for data to
`18.216.110.65`

computer 1

computer 2

**Solution**: every machine* has an IP address (Internet Protocol).
Requests are sent to a specific IP address.

*some machines have more multiple addresses

# Internet Protocol

address: `18.216.110.65`

request for data to
`18.216.110.65`

**computer 1**

**computer 2**

**Challenge**: it's hard to remember IP addresses.
Imagine you had to type a number instead of www.google.com!

# Domain Names

domain: tyler.caraza-harter.com
address: 18.216.110.65

request for data to
18.216.110.65

**computer 1**

**computer 2**

**Solution**: use "nicknames" (called domain names)
for IP addresses of machines that serve data

# Domain Names

domain: tyler.caraza-harter.com
address: `18.216.110.65`

if a domain name is used here,
the IP address is looked up here first

request for data to
`18.216.110.65`

computer 1

computer 2

**Solution**: use "nicknames" (called domain names)
for IP addresses of machines that serve data

# Port Numbers

address: `18.216.110.65`

program A

program B `Python`

request for data to
`18.216.110.65`

**computer 1**

`Python` program X

`Python` program Y

NGINX program Z

**computer 2**

**Challenge**: there may be multiple programs running on each computer.
How do we get the messages to the right program?

# Port Numbers

address: `18.216.110.65`

program A

program B | Python

request for data to
`18.216.110.65`

**computer 1**

Python  program X
port 8080

Python  program Y
port 5000

NGINX  program Z
port 80

**computer 2**

**Solution**: give each program a unique ID (called a "port number")

(like apartment numbers)

# Port Numbers

address: `18.216.110.65`

program A

program B  Python

**computer 1**

request for data to
`18.216.110.65`**:80**

Python  program X
port 8080

Python  program Y
port 5000

NGINX

program Z

port 80

**computer 2**

**Solution**: specify port number in request

# Port Numbers

defaults are often used
(e.g., browsers default to
80 or 443)

address: `18.216.110.65`

program A

program B `Python`

**computer 1**

request for data to
`18.216.110.65`**`:80`**

`Python` program X
port 8080

`Python` program Y
port 5000

program Z

port 80

**computer 2**

**Solution**: specify port number in request

what do requests contain?

request for data

computer 1

computer 2

computer 1

response with data

what do responses contain?

computer 2

depends on application!  (video chat, web browsing, etc)

we'll only consider web applications for this semester

response with data

computer 1

computer 2

what do responses contain?

# Learning Objectives Today

Motivation

Networking Basics

HTTP (Hypertext Transfer Protocol)

Requests Module

# HTTP

Protocol for communicating web data
- downloading a specific webpage, image, etc

domain: <u>example.com</u>
address: `12.34.56.78`

program B `Python`

**computer 1**

`Python` program Y
port 80

**computer 2**

**Note**: we won't talk about HTTPS today, which is HTTP with encryption

# HTTP

Protocol for communicating web data
- downloading a specific webpage, image, etc

domain: <u>example.com</u>
address: `12.34.56.78`

program B `Python` —— please send **home page** ——→ `Python` program Y
port 80

**computer 1**

**computer 2**

# HTTP

Protocol for communicating web data

- downloading a specific webpage, image, etc

domain: <u>example.com</u>
address: `12.34.56.78`

program B `Python`

**computer 1**

please send **/index.html**

`Python` program Y

port 80

**computer 2**

# HTTP

Protocol for communicating web data
- downloading a specific webpage, image, etc

domain: <u>example.com</u>
address: `12.34.56.78`

program B **Python**

please send **/about.html**

**computer 1**

**Python** program Y
port 80

**computer 2**

# HTTP

Protocol for communicating web data
- downloading a specific webpage, image, etc

domain: <u>example.com</u>
address: `12.34.56.78`

program B `Python`

please send **/logo.gif**

`Python` program Y

port 80

**computer 1**

**computer 2**

# HTTP

Protocol for communicating web data
- downloading a specific webpage, image, etc

domain: <u>example.com</u>
address: `12.34.56.78`

program B **Python** ———→ **Python** program Y
                                       port 80

please send **/logo.gif**
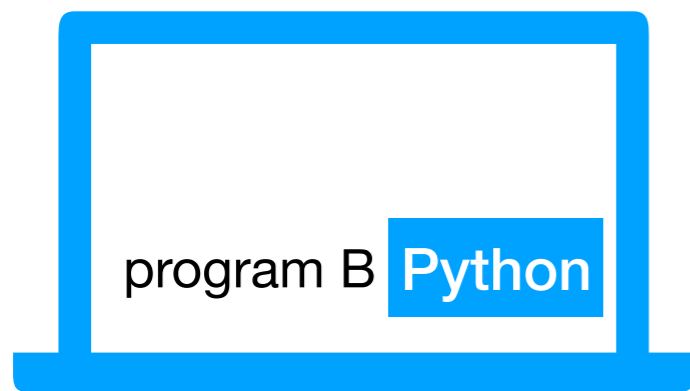
**computer 1**

**computer 2**

Note we need three things:
1. domain name
2. port number
3. resource (file name)

which computer?

which program
on that computer?

which resource
from that program?

getting specific
about what we want

Note we need three things:
1. domain name
2. port number
3. resource (file name)

which computer?

which program
on that computer?

which resource
from that program?

**getting specific
about what we want**

**URL**

Note we need three things:
1. domain name
2. port number
3. resource (file name)

# URLs

`https://en.wikipedia.org:443/wiki/URL`



URL

Note we need three things:
1. domain name
2. port number
3. resource (file name)

# URLs

`https://en.wikipedia.org:443/wiki/URL`



**URL** ◁ Note we need three things:
1. domain name
2. port number
3. resource (file name)

# URLs

domain name

`https://en.wikipedia.org:443/wiki/URL`

port



URL ◁ Note we need three things:
1. domain name
2. port number
3. resource (file name)

# URLs

domain name resource

`https://en.wikipedia.org:443/wiki/URL`

port



**URL**

Note we need three things:
1. domain name
2. port number
3. resource (file name)

# URLs

domain name          resource

`https://en.wikipedia.org/wiki/URL`

port would have defaulted to 443 if not specified

URL ⟨ Note we need three things:
1. domain name
2. port number
3. resource (file name)

# HTTP

Protocol for communicating web data
- downloading a specific webpage, image, etc

domain: <u>example.com</u>
address: `12.34.56.78`

program B **Python**    →    **Python** program Y
please send /about.html    port 80

**computer 1**

**computer 2**

**URL** ⟨ Note we need three things:
1. domain name
2. port number
3. resource (file name)

# HTTP

Protocol for communicating web data
- downloading a specific webpage, image, etc

domain: <u>example.com</u>
address: `12.34.56.78`

program B `Python`

please send /about.html

**computer 1**

`Python` program Y
port 80

**computer 2**

**URL** Note we need three things:
1. domain name
2. port number
3. resource (file name)

how do we specify this?

# HTTP

Protocol for communicating web data
- downloading a specific webpage, image, etc

domain: <u>example.com</u>
address: `12.34.56.78`

program B **Python**

please send /about.html

**Python** program Y
port 80

**computer 1**

**computer 2**

**HTTP Request:**

```
GET /about.html HTTP/1.1
Host: example.com
User-Agent: ...
Accept: */*
```

# HTTP

Protocol for communicating web data
- downloading a specific webpage, image, etc

domain: example.com
address: `12.34.56.78`

program B **Python** ← **response** ← **Python** program Y

port 80

**computer 1**

**computer 2**

**HTTP Response:**

```
HTTP/1.0 200 OK
Content-Type: text/html; charset=utf-8
Content-Length: 74
Server: Werkzeug/0.14.1 Python/3.6.6
Date: Sun, 11 Nov 2018 17:00:29 GMT

all the contents
```

# Request and Response Headers

**HTTP Request:**

```
GET /about.html HTTP/1.1
Host: example.com
User-Agent: ...
Accept: */*
```

**HTTP Response:**

```
HTTP/1.0 200 OK
Content-Type: text/html; charset=utf-8
Content-Length: 74
Server: Werkzeug/0.14.1 Python/3.6.6
Date: Sun, 11 Nov 2018 17:00:29 GMT

all the contents
```

There are **LOTS** of details here we don't care about right now

# Request and Response Headers

we want the about.html page

**HTTP Request:**

```
GET /about.html HTTP/1.1
Host: example.com
User-Agent: ...
Accept: */*
```

**HTTP Response:**

```
HTTP/1.0 200 OK
Content-Type: text/html; charset=utf-8
Content-Length: 74
Server: Werkzeug/0.14.1 Python/3.6.6
Date: Sun, 11 Nov 2018 17:00:29 GMT

all the contents
```

data in about.html

There are **LOTS** of details here we don't care about right now

# Request and Response Headers

we want the about.html page

**HTTP Request:**

```
GET /about.html HTTP/1.1
Host: example.com
User-Agent: ...
Accept: */*
```

**status code.** 200 is good. 404, 500, others are various errors or other more complicated states

**HTTP Response:**

```
HTTP/1.0 200 OK
Content-Type: text/html; charset=utf-8
Content-Length: 74
Server: Werkzeug/0.14.1 Python/3.6.6
Date: Sun, 11 Nov 2018 17:00:29 GMT
```

data in about.html

```
all the contents
```

There are **LOTS** of details here we don't care about right now

**method.** *GET* is simple download. *POST* means we are uploading data as part of our request. We won't talk about others today.

we want the about.html page

**HTTP Request:**
```
GET /about.html HTTP/1.1
Host: example.com
User-Agent: ...
Accept: */*
```

**status code.** 200 is good. 404, 500, others are various errors or other more complicated states

**HTTP Response:**
```
HTTP/1.0 200 OK
Content-Type: text/html; charset=utf-8
Content-Length: 74
Server: Werkzeug/0.14.1 Python/3.6.6
Date: Sun, 11 Nov 2018 17:00:29 GMT

all the contents
```

data in about.html

There are **LOTS** of details here we don't care about right now

# Learning Objectives Today

Motivation

Networking Basics

HTTP (Hypertext Transfer Protocol)

Requests Module

# Requests module

Purpose
- easily send requests to a server and parse the response
- *"HTTP for Humans™"*
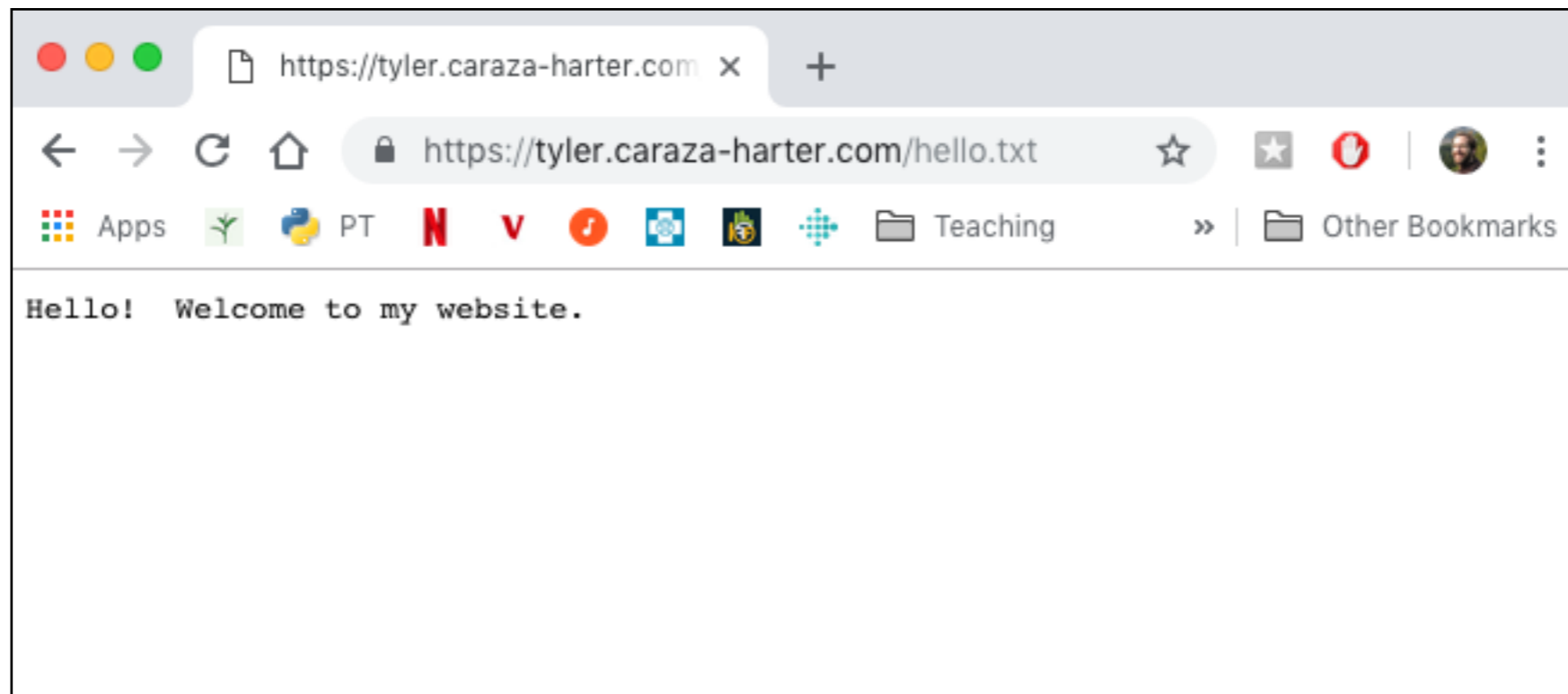
Installation
- install:
  `pip install requests`

Using it
- just import:
  `import requests`

# GET Request

```python
import requests

url = "https://tyler.caraza-harter.com/hello.txt"

requests.get(url)
```
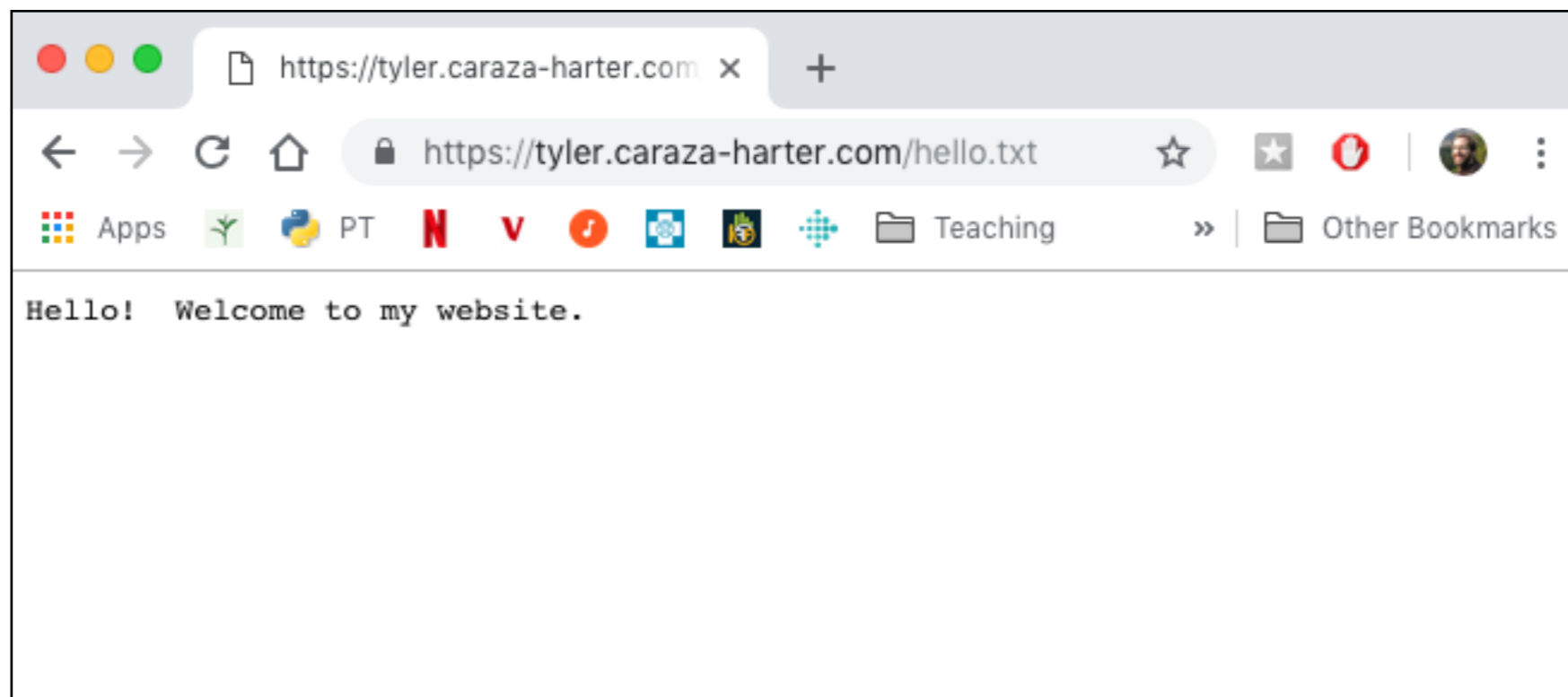
# GET Request

```python
import requests

url = "https://tyler.caraza-harter.com/hello.txt"

requests.get(url)
```

sends a **GET** request to tyler.caraza-harter.com,
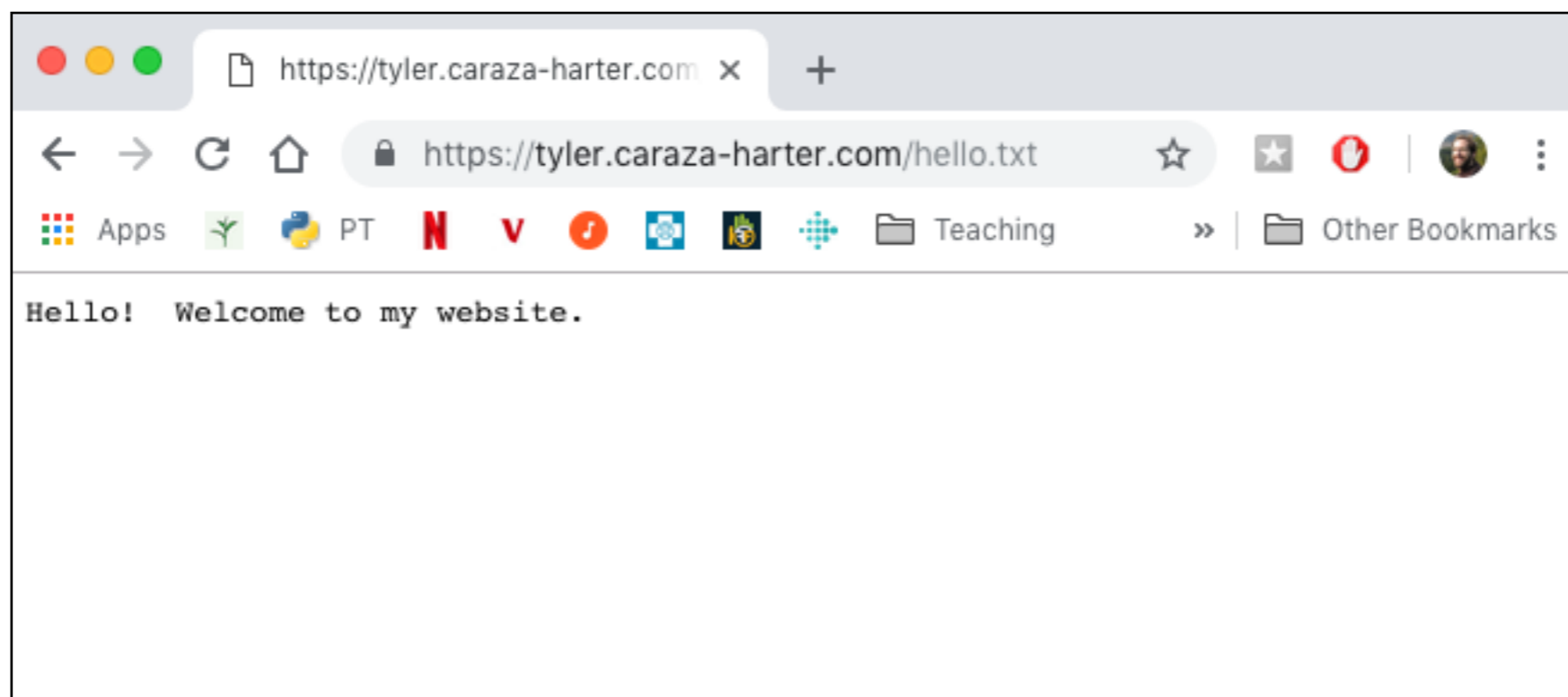asking for the contents of the **/hello.txt** page

# GET Request

```python
import requests

url = "https://tyler.caraza-harter.com/hello.txt"

resp = requests.get(url)
```

put response from tyler.caraza-harter.com in the resp variable

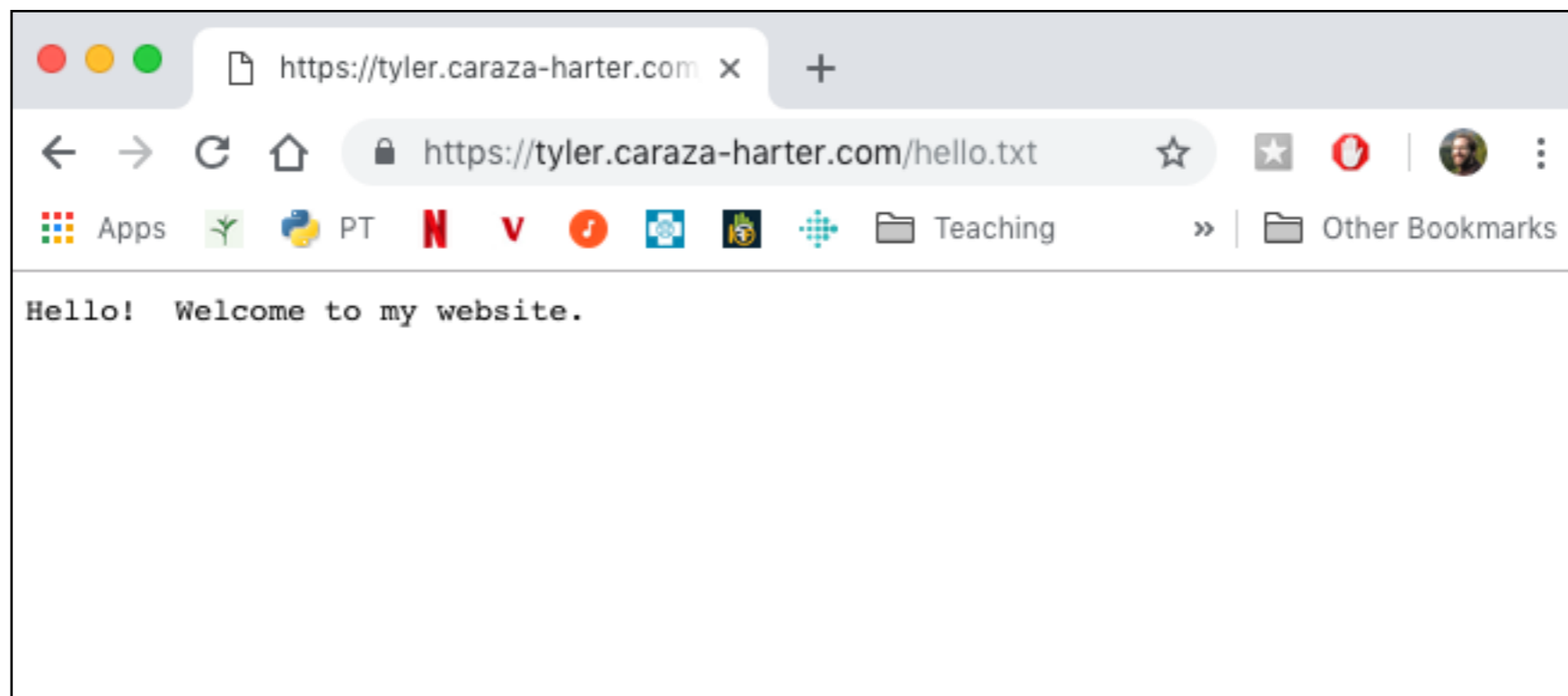# GET Request

```python
import requests

url = "https://tyler.caraza-harter.com/hello.txt"

resp = requests.get(url)

# make sure we got 200 (success) back
assert(resp.status_code == 200)
```
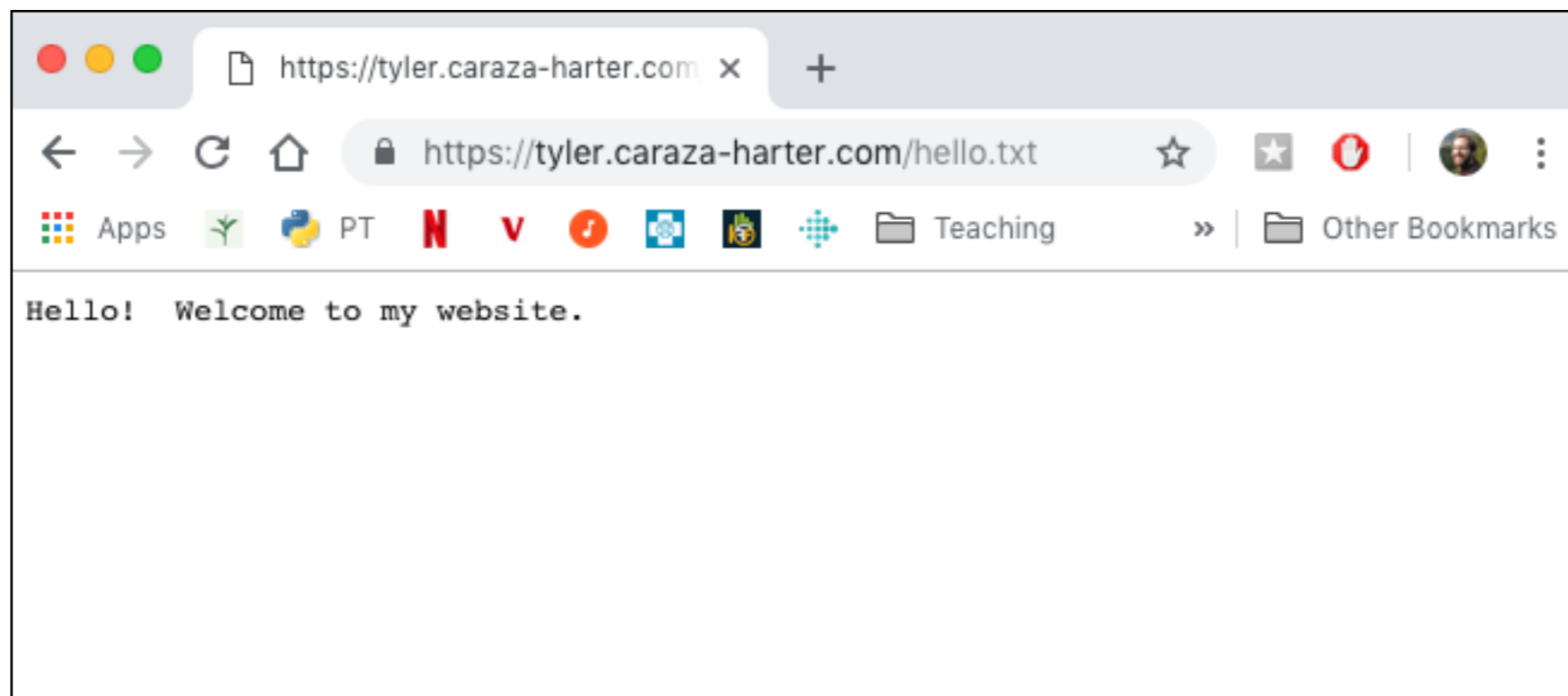
# GET Request

```python
import requests

url = "https://tyler.caraza-harter.com/hello.txt"

resp = requests.get(url)

resp.raise_for_status() # shortcut
```
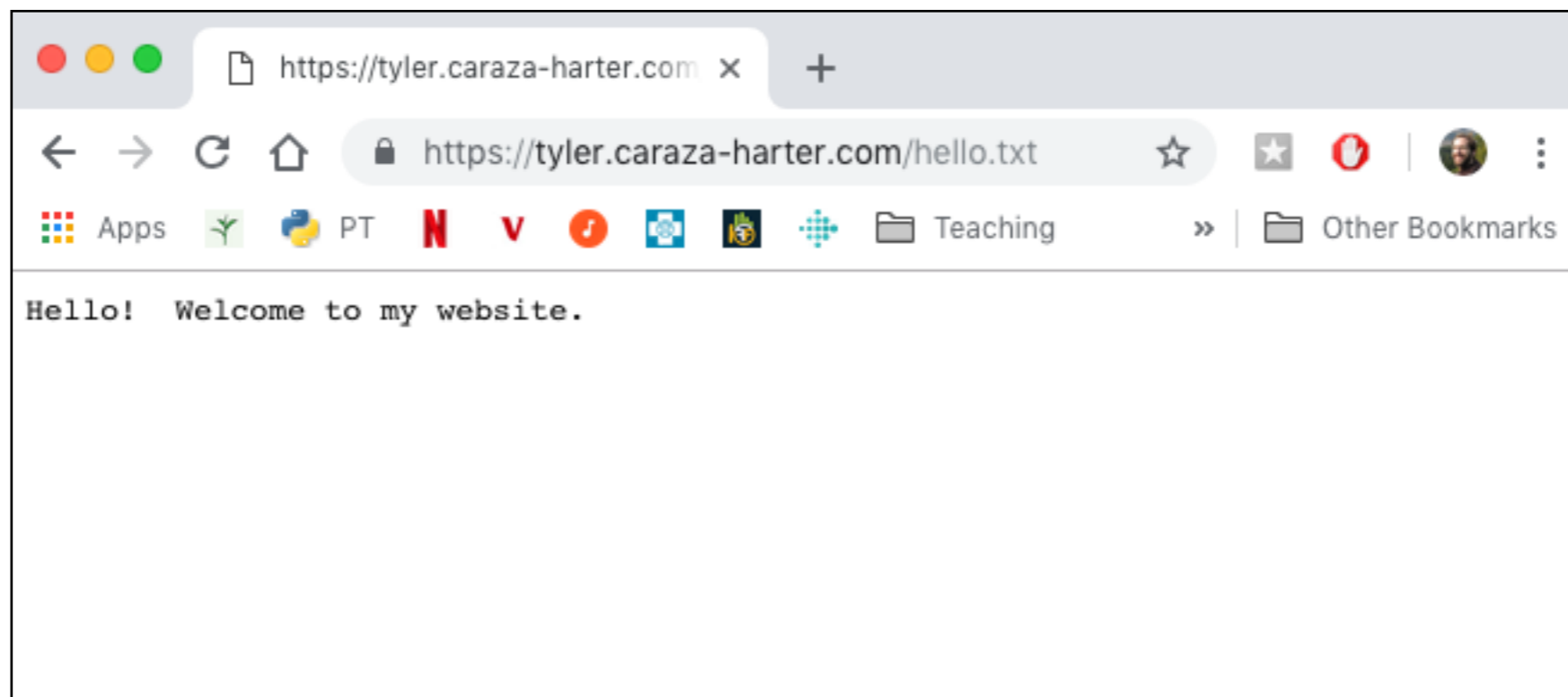
# GET Request

```python
import requests

url = "https://tyler.caraza-harter.com/hello.txt"

resp = requests.get(url)

resp.raise_for_status() # shortcut
print(resp.text) # "Hello! Welcome to my website."
```
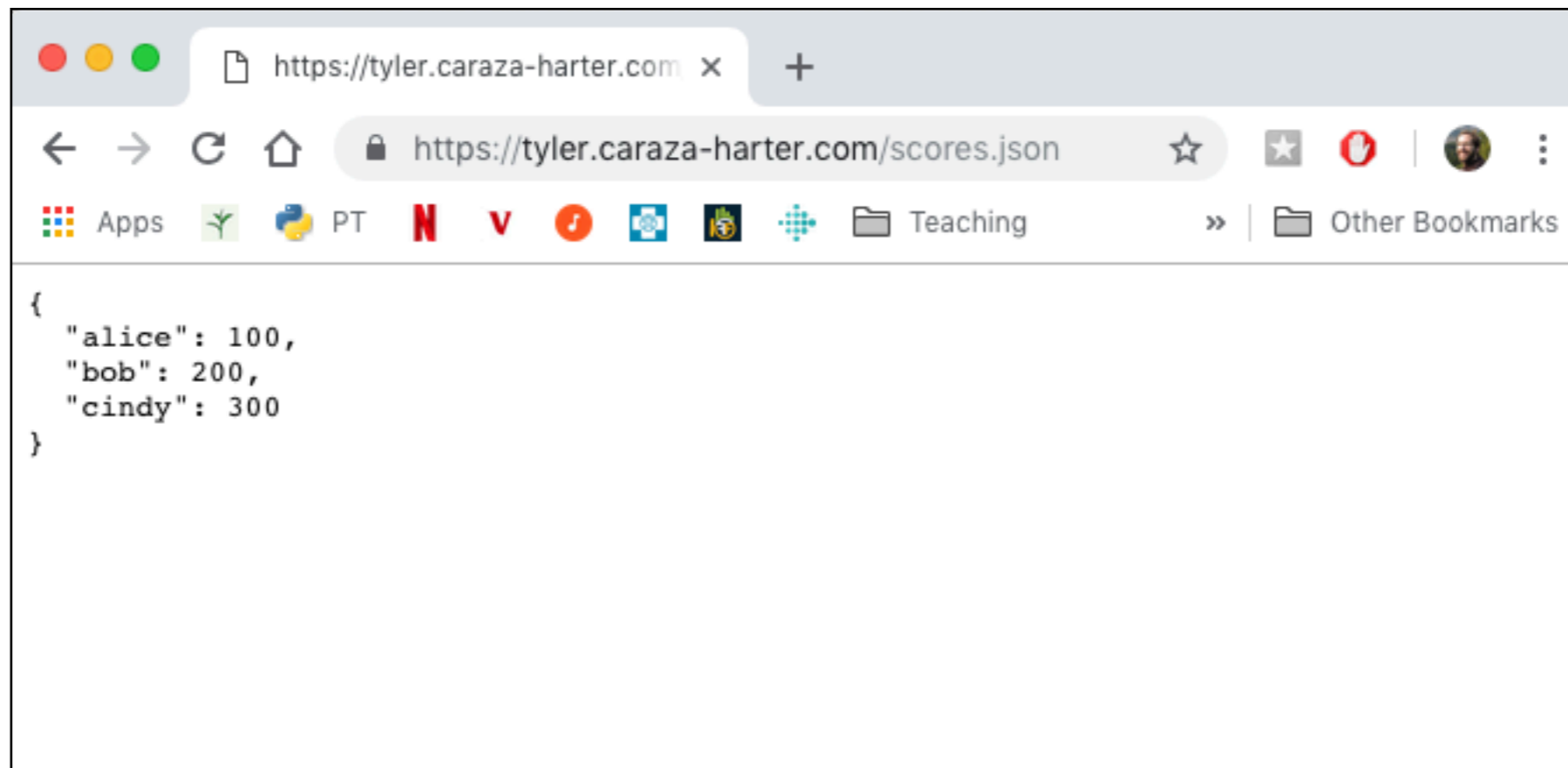
# JSON Responses

```python
import requests, json

url = "https://tyler.caraza-harter.com/scores.json"
resp = requests.get(url)

scores = json.loads(resp.text)
```

# JSON Responses

```python
import requests, json

url = "https://tyler.caraza-harter.com/scores.json"
resp = requests.get(url)

scores = json.loads(resp.text)
scores = resp.json() # shortcut
```
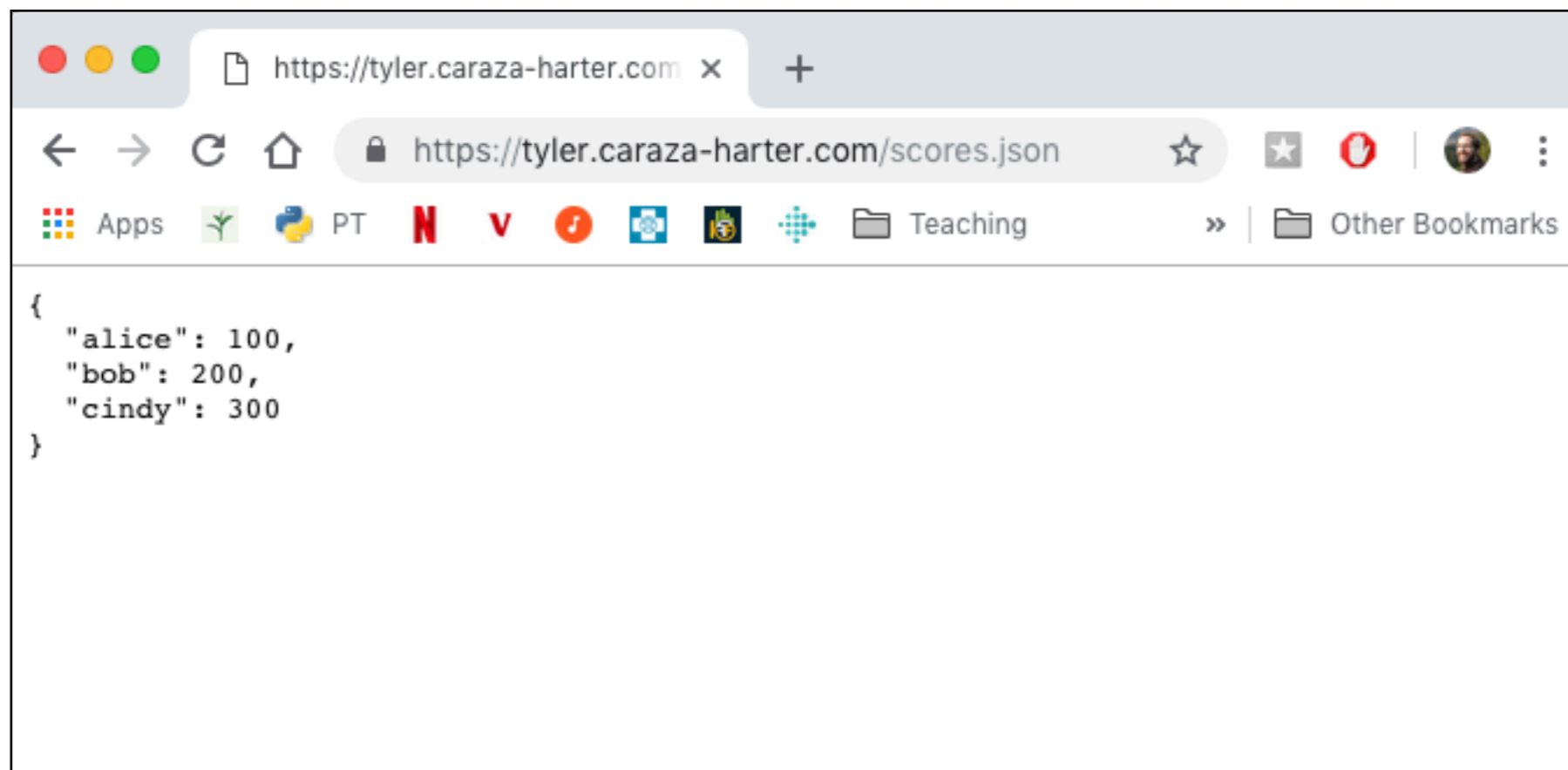


Browser showing `https://tyler.caraza-harter.com/scores.json`:

```json
{
    "alice": 100,
    "bob": 200,
    "cindy": 300
}
```

# Demo 1: State Populations

Goal: fetch population data for all states and provide summary stats

**Input:**

- List of state files: https://tyler.caraza-harter.com/cs301/spring19/materials/code/lec-30/data/state_files.txt
- The 50 JSON files

**Output:**

- Stats about population: mean, max, min, etc

**Bonus!** "cache" results to make reruns of notebook faster

In [19]: `df.describe().astype(int)`

Out[19]:

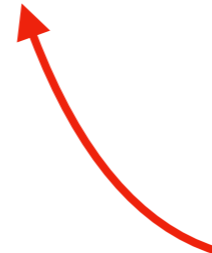|       | 2000 | 2010 | 2015 |
|-------|------|------|------|
| count | 50 | 50 | 50 |
| mean | 5616996 | 6162876 | 6364951 |
| std | 6185579 | 6848235 | 7152085 |
| min | 493782 | 563626 | 584304 |
| 25% | 1735533 | 1833004 | 1857308 |
| 50% | 4026890 | 4436369 | 4530803 |
| 75% | 6281944 | 6680312 | 6986155 |
| max | 33871648 | 37253956 | 38792291 |

# POST Request

```python
import requests

url = "…"

requests.post(url, data)
```

# POST Request

```python
import requests

url = "…"

requests.post(url, data)
```

data to upload

# Demo 2: Score Keeper

Goal: use POSTs and GETs to keep track of scores

**Server Setup**:

- `pip install flask`
- download https://raw.githubusercontent.com/tylerharter/caraza-harter-com/master/tyler/cs301/spring19/materials/code/lec-30/scores.py
- run this: `python scores.py`
- open `http://127.0.0.1:8080/` in web browser
- see code examples

**to view scores:**

`GET IP:PORT/scores`

**to record score:**

`POST` a player name to `IP:PORT/scores`