# DATA THIEVES PROJECT

Pedro Afonso, Rennê Cirqueira, Thais Ternus

Data Analytics
Part-time Course

# GOAL

HELP A PERSON CHOOSING A BOOK FROM
"THE BEST OF THE BEST"

| New York Times Best Sellers | > | Our Project | < | Amazon Books |
|---|---|---|---|---|

IRON
HACK

# PRODUCT

Combined Print & E-Book Fiction

NORA ROBERTS
LEGACY

#1 New York Times

CLICK HERE

AMAZON INFORMATION

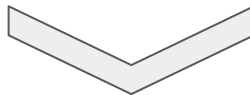⭐⭐⭐⭐½ ▾   7,807 ratings
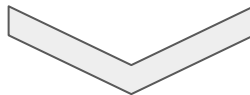
Hardcover
$16.30

The Last Thing

#2 New York Times

Select one Best seller list

The APP will bring the TOP 5 books

Choose the book and click to read more about it

IRON
HACK

# API CONNECTION

## METHODOLOGY

➔ Best Sellers API
comes with several functions

➔ "Overview" retrieves
the top 5 for all categories

➔ Right now, the API returns
more than 50 books
in at least 10 categories

➔ All the categories
were analysed to be included
on the project



IRON
HACK

# WEB SCRAPING

amazon

| Read the CSV file saved on API code | > | Pick the headers to be able to scraping the Amazon website | > | Working on a function to get the informations (Price, review and review count) | > | Merge the TOP 5 API file with scraping file | > | Output to select the books categories |

```
headers= ({'User-Agent':
         'Mozilla/5.0 (Windows NT 6.1) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/41.0.2228.0 Safari/537.36',
         'Accept-Language': 'en-US, en;q=0.5'})
```

IRON
HACK

# WEB SCRAPING

amazon

## FUNCTION METHODOLOGY

Connection to web scrape the amazon website

Then to pick the informations TRY and EXCEPT are used

```python
def book_function(url_data):


    data=[]

    for url_2 in url_data['URL']: #[:4]: slicing only the top 4
        time.sleep(3)
        amazon = requests.get(url_2, headers=headers)
        content = amazon.content
        soup = BeautifulSoup(content)
```
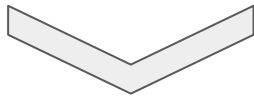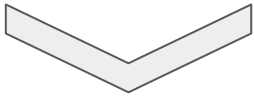
```python
try:
    price =(soup.find('span', attrs = {'class':'a-size-base a-color-price a-color-price'}).get_text()).strip().repl
except:
    #price = 'Nothing'
    try:
        price =(soup.find('span', attrs = {'class':'slot-price'}).get_text()).strip().replace('$', '')
    except:
        try:
            price = (soup.find('span', attrs = {'class':'a-size-base a-color-secondary'}).get_text()).strip().repla
        #price = 'Nothing'
        except:
            price = 'Nothing'
```

# NEXT STEPS

➔  Create a API for easier access to everyone to run the program;

➔  Automatization of the API and the Web Scraping codes to retrieve the updated info from the NYTimes and Amazon Books;

➔  Connect with more data sources;

➔  Create a recommendation system based on machine learning;

➔  Improve Amazon web scraping;

➔  More detailed documentation on the GitHub.

IRON
HACK

# ANY QUESTION?

https://github.com/thaisternus/project-week-3-data-thieves

THANK YOU!

Data Analytics
Part-time Course