# 网络爬虫
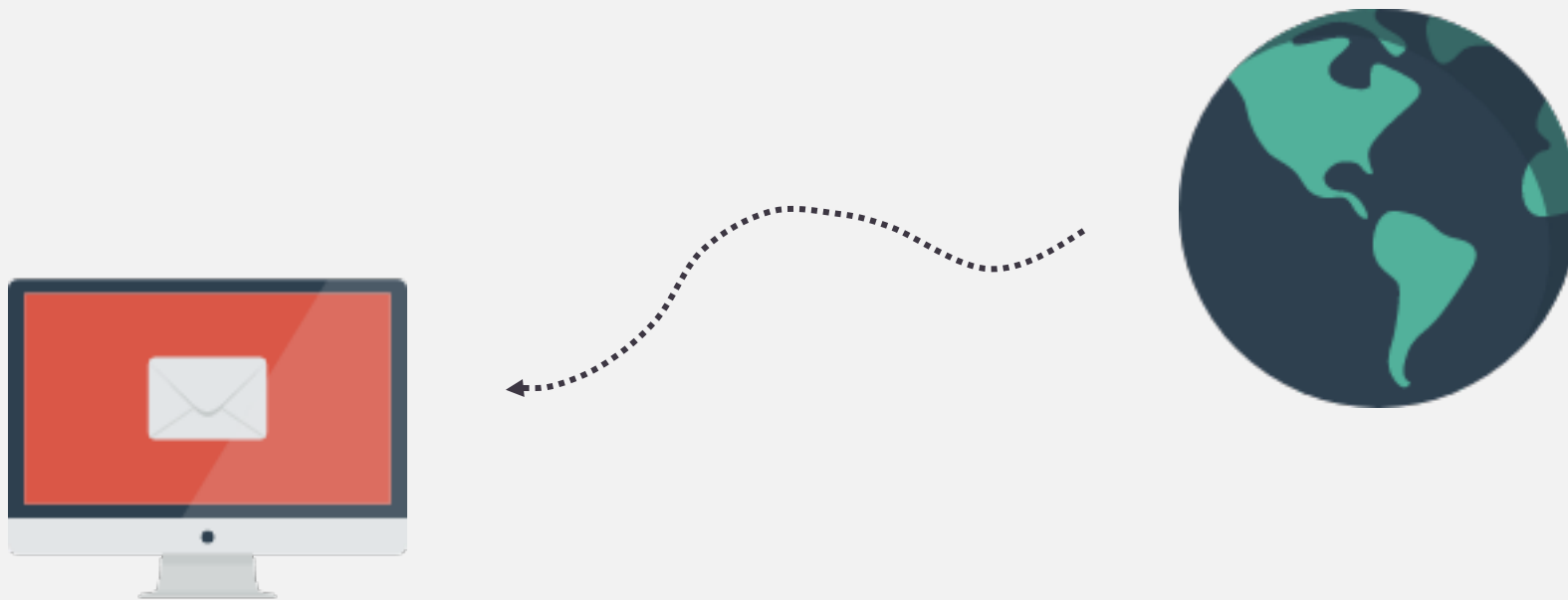
数据科学引论（Python之道）

# 爬虫是什么

爬虫crawler，即网络爬虫Spider。是去**自动化**获取网络上的内容，是一个能够自动化地访问互联网并将网站内容下载下来的的程序或脚本。

# 为什么需要爬虫？

高效自动化地从网络获取收集数据，后续可做数据处理。

# 基本流程



## URL
### 地址

统一资源定位符(网址)

根据url指定需要爬取的网页

## HTML
### 网页原内容

网页内容编码

我们要从中提取我们需要的关键信息

## TOOL
### 工具

爬虫工具

访问网页并从html解析信息的具体工具，如Scrapy

## FORMAT
### 格式

统一化格式

得到信息将数据整理成统一格式导出

# Scrapy – URL指定

```python
def __init__(self):

    self.file = open('demo1_quotes.json', 'w');

    #设置待爬取网站列表
    self.urls = []
    for i in range(1,3):
        self.urls.append('http://quotes.toscrape.com/page/' + str(i) )

#   初始化效果  效果等同
#     self.urls = [
#           'http://quotes.toscrape.com/page/1/',
#           'http://quotes.toscrape.com/page/2/',
#       ]

    print(self.urls)
```

quotes.toscrape.com/page/1/

quotes.toscrape.com/page/2/

# Quotes to Scrape

"This life is what you make it. No matter what, you're going to
sometimes, it's a universal truth. But the good part is you get to
how you're going to mess it up. Girls will be your friends - they
anyway. But just remember, some come, some go. The ones t

# Scrapy – HTML分析

```
def parse(self, response):
    #提取名言列表
    quotes = response.css("div.quote");
    for quote in quotes:
        #提取每条名言中的作者名
        author = quote.css("small.author::text").extract_first();
        #提取名言的文字内容
        text = quote.css(".text::text").extract_first();
        #提取名言标签
        tags = quote.css(".tags .tag::text").extract();
```

```html
▼<div class="quote" itemscope itemtype="http://schema.org/CreativeWork">
    ▼<span class="text" itemprop="text">
        ""The world as we have created it is a process of our thinking. It cannot be changed without
        changing our thinking.""
    </span>
    ▼<span>
        "by "
        <small class="author" itemprop="author">Albert Einstein</small>
        <a href="/author/Albert-Einstein">(about)</a>
    </span>
    ▼<div class="tags">
        "
                    Tags:
        "
        <meta class="keywords" itemprop="keywords" content="change,deep-thoughts,thinking,world">
        <a class="tag" href="/tag/change/page/1/">change</a>
        <a class="tag" href="/tag/deep-thoughts/page/1/">deep-thoughts</a>
        <a class="tag" href="/tag/thinking/page/1/">thinking</a>
        <a class="tag" href="/tag/world/page/1/">world</a>
    </div>
</div>
▶<div class="quote" itemscope itemtype="http://schema.org/CreativeWork">…</div>
▶<div class="quote" itemscope itemtype="http://schema.org/CreativeWork">…</div>
```

# Scrapy – 格式化导出

```python
#parse方法会在每个request收到response之后调用
def parse(self, response):
    #提取名言列表
    quotes = response.css("div.quote");
    for quote in quotes:
        #提取每条名言中的作者名
        author = quote.css("small.author::text").extract_first()
        #提取名言的文字内容
        text = quote.css(".text::text").extract_first();
        #提取名言标签
        tags = quote.css(".tags .tag::text").extract();
        #构建字典对象
        item = {"author":author, "text": text, "tags":tags };
        #将字典转换成json字符串
        line = json.dumps(dict(item))
        #将每个条目写入文件
        self.file.write(line + "\n")

    #及时将内容写入文件，否则可能会出现少许延迟
    self.file.flush()
    os.fsync(self.file)
    #输出当前解析完成的网页网址，可以当做爬取进度来看待,与程序逻辑无关
    print("over: " + response.url)
```

```
1  {"author": "Albert Einstein", "tags": ["change", "d
2  {"author": "J.K. Rowling", "tags": ["abilities", "c
3  {"author": "Albert Einstein", "tags": ["inspiration
4  {"author": "Jane Austen", "tags": ["aliteracy", "bo
5  {"author": "Marilyn Monroe", "tags": ["be-yourself"
6  {"author": "Albert Einstein", "tags": ["adulthood"
7  {"author": "Andr\u00e9 Gide", "tags": ["life", "lov
8  {"author": "Thomas A. Edison", "tags": ["edison",
9  {"author": "Eleanor Roosevelt", "tags": ["misattri
10 {"author": "Steve Martin", "tags": ["humor", "obvio
11 {"author": "Marilyn Monroe", "tags": ["friends", "h
12 {"author": "J.K. Rowling", "tags": ["courage", "fri
13 {"author": "Albert Einstein", "tags": ["simplicity"
14 {"author": "Bob Marley", "tags": ["love"], "text"
15 {"author": "Dr. Seuss", "tags": ["fantasy"], "text"
16 {"author": "Douglas Adams", "tags": ["life", "navig
17 {"author": "Elie Wiesel", "tags": ["activism", "apa
18 {"author": "Friedrich Nietzsche", "tags": ["friends
19 {"author": "Mark Twain", "tags": ["books", "content
20 {"author": "Allen Saunders", "tags": ["fate", "life
21
```

# 常见问题

IP被封杀

重构

网页更新

# Thanks for Watching