# 爬虫案例6(随机代理池+随机usragent+时间间隔)

https://www.cnblogs.com/cnkai/p/7401343.html

https://www.cnblogs.com/cnkai/p/7401526.html

https://blog.csdn.net/mouday/article/details/81512748 随机时间间隔

先固定间隔 在setting中 设置
```
DOWNLOAD_DELAY=3
```

————————

1.创建项目

```
MichaelYun:PycharmProjects Yun$ scrapy startproject randommm
New Scrapy project 'randommm', using template directory '/usr/local/lib/python3.7/site-package
s/scrapy/templates/project', created in:
    /Users/Yun/PycharmProjects/randommm

You can start your first spider with:
    cd randommm
    scrapy genspider example example.com
MichaelYun:PycharmProjects Yun$ cd randommm
MichaelYun:randommm Yun$ scrapy genspider  random http://www.ip138.com/
Created spider 'random' using template 'basic' in module:
  randommm.spiders.random
MichaelYun:randommm Yun$
```

2.在setting.py 中设置user_agent 和 代理ip池

注意:以后的思路可以自己爬，读取json或者自己爬的json文件

```
1 MY_USER_AGENT = [
2     "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; AcooBrowser; .NET CLR 1.1.4322;
  .NET CLR 2.0.50727)",
3     "Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 6.0; Acoo Browser; SLCC1; .NET CLR
  2.0.50727; Media Center PC 5.0; .NET CLR 3.0.04506)",
4     "Mozilla/4.0 (compatible; MSIE 7.0; AOL 9.5; AOLBuild 4337.35; Windows NT 5.1; .NET CLR
  1.1.4322; .NET CLR 2.0.50727)",
5     "Mozilla/5.0 (Windows; U; MSIE 9.0; Windows NT 9.0; en-US)",
6     "Mozilla/5.0 (compatible; MSIE 9.0; Windows NT 6.1; Win64; x64; Trident/5.0; .NET CLR
  3.5.30729; .NET CLR 3.0.30729; .NET CLR 2.0.50727; Media Center PC 6.0)",
7     "Mozilla/5.0 (compatible; MSIE 8.0; Windows NT 6.0; Trident/4.0; WOW64; Trident/4.0;
  SLCC2; .NET CLR 2.0.50727; .NET CLR 3.5.30729; .NET CLR 3.0.30729; .NET CLR 1.0.3705; .NET
  CLR 1.1.4322)",
8     "Mozilla/4.0 (compatible; MSIE 7.0b; Windows NT 5.2; .NET CLR 1.1.4322; .NET CLR
  2.0.50727; InfoPath.2; .NET CLR 3.0.04506.30)",
9     "Mozilla/5.0 (Windows; U; Windows NT 5.1; zh-CN) AppleWebKit/523.15 (KHTML, like Gecko,
```

```
     Safari/419.3) Arora/0.3 (Change: 287 c9dfb30)",
10      "Mozilla/5.0 (X11; U; Linux; en-US) AppleWebKit/527+ (KHTML, like Gecko, Safari/419.3)
   Arora/0.6",
11      "Mozilla/5.0 (Windows; U; Windows NT 5.1; en-US; rv:1.8.1.2pre) Gecko/20070215 K-
   Ninja/2.1.1",
12      "Mozilla/5.0 (Windows; U; Windows NT 5.1; zh-CN; rv:1.9) Gecko/20080705 Firefox/3.0
   Kapiko/3.0",
13      "Mozilla/5.0 (X11; Linux i686; U;) Gecko/20070322 Kazehakase/0.4.5",
14      "Mozilla/5.0 (X11; U; Linux i686; en-US; rv:1.9.0.8) Gecko Fedora/1.9.0.8-1.fc10
   Kazehakase/0.5.6",
15      "Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/535.11 (KHTML, like Gecko)
   Chrome/17.0.963.56 Safari/535.11",
16      "Mozilla/5.0 (Macintosh; Intel Mac OS X 10_7_3) AppleWebKit/535.20 (KHTML, like Gecko)
   Chrome/19.0.1036.7 Safari/535.20",
17      "Opera/9.80 (Macintosh; Intel Mac OS X 10.6.8; U; fr) Presto/2.9.168 Version/11.52",
18      "Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/536.11 (KHTML, like Gecko)
   Chrome/20.0.1132.11 TaoBrowser/2.0 Safari/536.11",
19      "Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.1 (KHTML, like Gecko)
   Chrome/21.0.1180.71 Safari/537.1 LBBROWSER",
20      "Mozilla/5.0 (compatible; MSIE 9.0; Windows NT 6.1; WOW64; Trident/5.0; SLCC2; .NET CLR
   2.0.50727; .NET CLR 3.5.30729; .NET CLR 3.0.30729; Media Center PC 6.0; .NET4.0C; .NET4.0E;
   LBBROWSER)",
21      "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; QQDownload 732; .NET4.0C;
   .NET4.0E; LBBROWSER)",
22      "Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/535.11 (KHTML, like Gecko)
   Chrome/17.0.963.84 Safari/535.11 LBBROWSER",
23      "Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 6.1; WOW64; Trident/5.0; SLCC2; .NET CLR
   2.0.50727; .NET CLR 3.5.30729; .NET CLR 3.0.30729; Media Center PC 6.0; .NET4.0C;
   .NET4.0E)",
24      "Mozilla/5.0 (compatible; MSIE 9.0; Windows NT 6.1; WOW64; Trident/5.0; SLCC2; .NET CLR
   2.0.50727; .NET CLR 3.5.30729; .NET CLR 3.0.30729; Media Center PC 6.0; .NET4.0C; .NET4.0E;
   QQBrowser/7.0.3698.400)",
25      "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; QQDownload 732; .NET4.0C;
   .NET4.0E)",
26      "Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 5.1; Trident/4.0; SV1; QQDownload 732;
   .NET4.0C; .NET4.0E; 360SE)",
27      "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; QQDownload 732; .NET4.0C;
   .NET4.0E)",
28      "Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 6.1; WOW64; Trident/5.0; SLCC2; .NET CLR
   2.0.50727; .NET CLR 3.5.30729; .NET CLR 3.0.30729; Media Center PC 6.0; .NET4.0C;
   .NET4.0E)",
29      "Mozilla/5.0 (Windows NT 5.1) AppleWebKit/537.1 (KHTML, like Gecko) Chrome/21.0.1180.89
   Safari/537.1",
30      "Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.1 (KHTML, like Gecko)
   Chrome/21.0.1180.89 Safari/537.1",
31      "Mozilla/5.0 (iPad; U; CPU OS 4_2_1 like Mac OS X; zh-cn) AppleWebKit/533.17.9 (KHTML,
   like Gecko) Version/5.0.2 Mobile/8C148 Safari/6533.18.5",
32      "Mozilla/5.0 (Windows NT 6.1; Win64; x64; rv:2.0b13pre) Gecko/20110307
   Firefox/4.0b13pre",
33      "Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:16.0) Gecko/20100101 Firefox/16.0",
34      "Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.11 (KHTML, like Gecko)
   Chrome/23.0.1271.64 Safari/537.11",
35      "Mozilla/5.0 (X11; U; Linux x86_64; zh-CN; rv:1.9.2.10) Gecko/20100922 Ubuntu/10.10
   (maverick) Firefox/3.6.10",
36      "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko)
   Chrome/58.0.3029.110 Safari/537.36",
```

```
37      ]
38
39
40  PROXIES = [
41      'http://218.2.80.59:8080',
42      'http://60.208.44.228:80',
43      'http://218.202.219.82:80',
44      'http://139.180.220.37:2333'
45  ]
```

3.编写自定义中间件，随机

```
1
2  from scrapy.downloadermiddlewares.useragent import UserAgentMiddleware
3  import random
4
5
6  #设置随机代理
7  class MyUserAgentMiddleware(UserAgentMiddleware):
8
9
10     def __init__(self, user_agent):
11         self.user_agent = user_agent
12
13     @classmethod
14     def from_crawler(cls, crawler):
15         return cls(
16             user_agent=crawler.settings.get('MY_USER_AGENT')
17         )
18
19     def process_request(self, request, spider):
20         agent = random.choice(self.user_agent)
21         request.headers['User-Agent'] = agent
22
23
24  #设置随机ip
25  class ProxyMiddleware(object):
26
27     def __init__(self, ip):
28         self.ip = ip
29
30     @classmethod
31     def from_crawler(cls, crawler):
32         return cls(ip=crawler.settings.get('PROXIES'))
33
34     def process_request(self, request, spider):
35         ip = random.choice(self.ip)
36         request.meta['proxy'] = ip
```

4.添加配置文件中

```
1  DOWNLOAD_DELAY=3
2
3  DOWNLOADER_MIDDLEWARES = {
4      'scrapy.downloadermiddleware.useragent.UserAgentMiddleware': None,
5      'randommm.middlewares.MyUserAgentMiddleware': 400,
6      'randommm.middlewares.ProxyMiddleware': 401,
```

```
 7 }
```

5.逻辑代码参考

```python
 1 # -*- coding: utf-8 -*-
 2 import scrapy
 3
 4
 5
 6 class RandomSpider(scrapy.Spider):
 7     name = 'random'
 8     allowed_domains = []
 9
10     def start_requests(self):
11
12         url = 'https://www.ip.cn/'
13
14         for i in range(4):
15             yield scrapy.Request(url=url, callback=self.parse, dont_filter=True)
16
17     def parse(self,response):
18         print('IP地址是:',response.css('#result > div > p:nth-
   child(3)::text').extract_first())
19
```

```
r: None)
IP地址是: GeoIP: Nanning, Guangxi, China
2019-02-22 12:27:57 [scrapy.core.engine] DEBUG: Crawled (200
r: None)
IP地址是: GeoIP: Nanning, Guangxi, China
2019-02-22 12:27:58 [scrapy.core.engine] DEBUG: Crawled (200
r: None)
IP地址是: GeoIP: Nanning, Guangxi, China
2019-02-22 12:28:00 [scrapy.core.engine] DEBUG: Crawled (200
r: None)
IP地址是: GeoIP: Wuhan, Hubei, China
2019-02-22 12:28:00 [scrapy.core.engine] INFO: Closing spid
```

可以看到有随机切换IP


注意参考:

```python
def judge_ip(self, ip, port):
    #判断ip是否可用
    http_url = "http://www.baidu.com"
    proxy_url = "http://{0}:{1}".format(ip, port)
    try:
        proxy_dict = {
            "http":proxy_url,
        }
        response = requests.get(http_url, proxies=proxy_dict)
    except Exception as e:
        print ("invalid ip and port")
        self.delete_ip(ip)
        return False
    else:
        code = response.status_code
        if code >= 200 and code < 300:
            print ("effective ip")
            return True
        else:
            print  ("invalid ip and port")
            self.delete_ip(ip)
            return False


def judge_ip(self, ip, port):
    #判断ip是否可用
    http_url = "http://www.baidu.com"
    proxy_url = "http://{0}:{1}".format(ip, port)
    try:
        proxy_dict = {
            "http":proxy_url,
        }
        response = requests.get(http_url, proxies=proxy_dict)
```