

# Web信息处理与应用



## 第二节 爬虫基础

徐童

2023.9.11

- 1965-1989, 两位先驱的成果
- Ted Nelson在1965年提出了超文本的概念。
  - HyperText, 源自于“非连续性著述” (Non sequential writing) 的理念, 即分叉的、允许读者作出选择的、无限延伸扩展的非线性文本。
- 1989年, Tim Berners Lee (万维网之父) 等人首次提出了一个World Wide Web协议, 随后定义了URL、HTML、HTTP等的规范, 使网络能够为大众所使用。



- 1993年，脱胎于爬虫的搜索引擎Wanderer
- Wanderer：最早的爬虫
  - 由MIT的学生Matthew Gray设计
  - 原意用于统计互联网上服务器的数量，而非为搜索引擎所设计
- Wandex：最早的网页索引计划
  - Wanderer后来发展为可以捕获网址，而为这些网址建立索引的计划就是Wandex

Wandex



- 本课程所要解决的问题



- **写在前面：本节课程说明**

- 本节课以爬虫基本概念和基础知识为主

- 计算机课？语文课！（滑稽 😏）

- 有想要了解更多关于爬虫实战内容的同学，可参考以下在线课程：

- Python教程（第21节-第29节部分）

- <https://www.bilibili.com/video/BV1ws411i7Dr/>

- Python爬虫全套课程

- <https://www.bilibili.com/video/BV1Yh411o7Sz>

- **TA倾情推荐**

- 如果觉得听课太累，下面还有更简洁的实战教程

- 零基础入门版

- <https://zhuanlan.zhihu.com/p/21479334>

- 常用工具

- <https://zhuanlan.zhihu.com/p/110448373>

- <https://zhuanlan.zhihu.com/p/57678048>

- **TA倾情推荐**

- 实战教程

- 卷王版

- 一个收集各种爬虫（默认爬虫语言为 python）的集合（如Bilibili 用户、京东商品等），一应俱全。不过有的失效或者不更新了，大佬们自由发挥。

<https://github.com/facert/awesome-spider>

小卷怡情，大卷伤身，强卷灰飞烟灭，内卷需谨慎

- **TA倾情推荐**

- 进阶版：分布式爬虫

- 含豆瓣热门、百度贴吧、百度翻译3个实例。建议循序渐进。

<https://github.com/Kr1s77/Python-crawler-tutorial-starts-from-zero>

- 拓展（移动端数据如何获取？）

- 微信小程序 Charles抓包

<https://blog.csdn.net/HeyShHeyou/article/details/90452656>

- 手机APP fiddler抓包

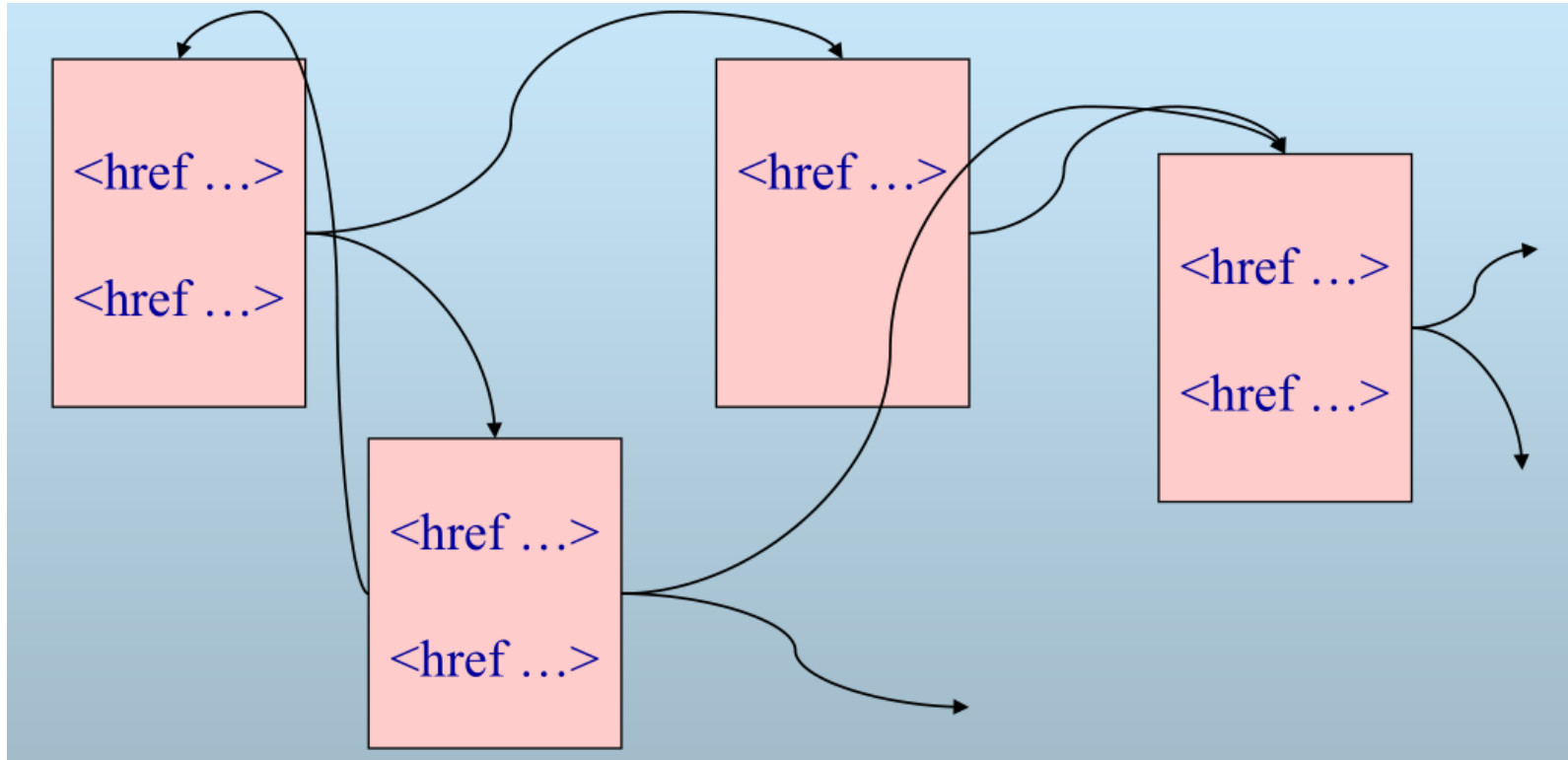
<https://zhuanlan.zhihu.com/p/34430703>



- **网络爬虫的定义与需求**
- 爬虫的基本要素
- 面向API的新爬虫任务
- 常见的爬虫算法
- 常见反爬虫机制与应对策略



- **Web网络的图模型**
- 以网页为节点、超链接（Hyper-Link）为有向边



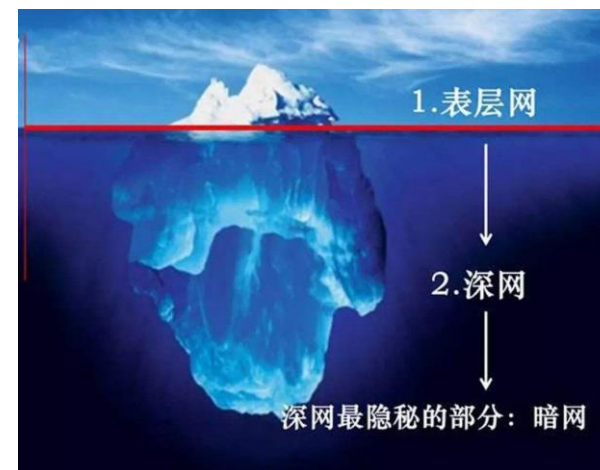
## • 爬虫的任务定义

- 从一个种子站点集合 (Seed sites) 开始, 从Web中寻找并且下载网页, 获取排序需要的相关信息, 并且剔除低质量的网页。
- 常见的爬虫类型
  - 通用网络爬虫: 目标为全网Web信息, 主要为门户网站搜索引擎和大型Web服务提供商采集数据。
  - 聚焦网络爬虫: 选择性爬取与预定主题相关的内容。
  - 增量式网络爬虫: 对已下载内容进行增量式更新并只爬取更新内容。



- **爬虫的任务定义**

- 从一个种子站点集合 (Seed sites) 开始, 从Web中寻找并且下载网页, 获取排序需要的相关信息, 并且剔除低质量的网页。
- 常见的爬虫类型
  - 深度网络爬虫: 专门负责获取搜索引擎无法索引的、超链接不可达的或需提交表单后才可见的网络内容。



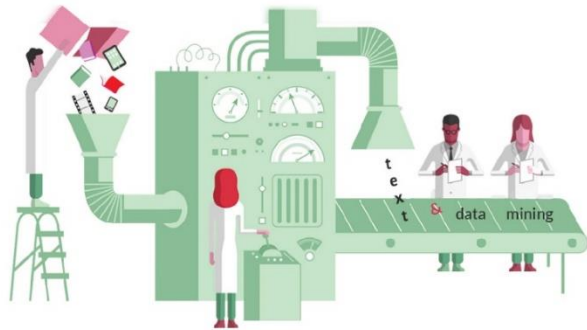
- 爬虫的基本用途



引擎优化



数据展示



数据分析

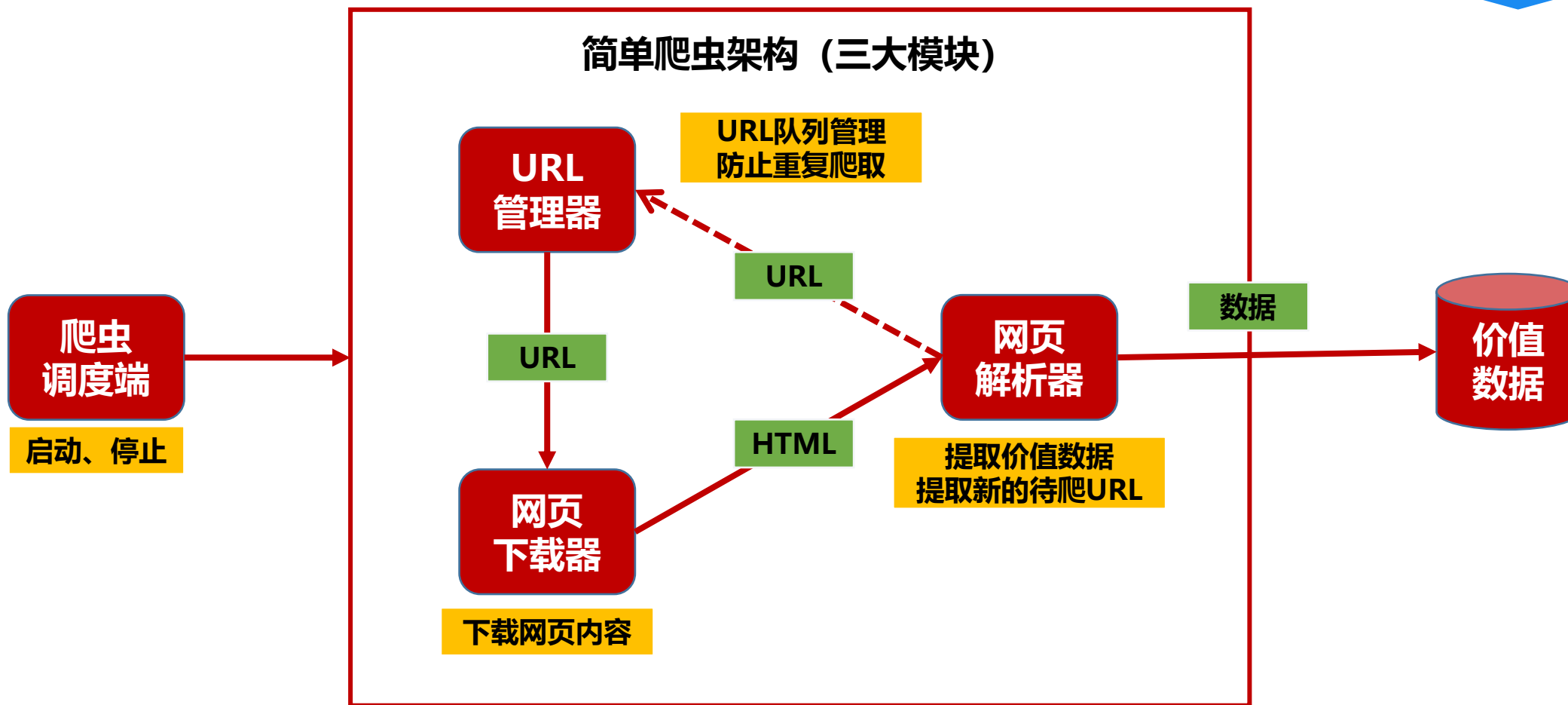
你们尽管抢课



能登进去算我输

特定应用

- 爬虫的基本流程

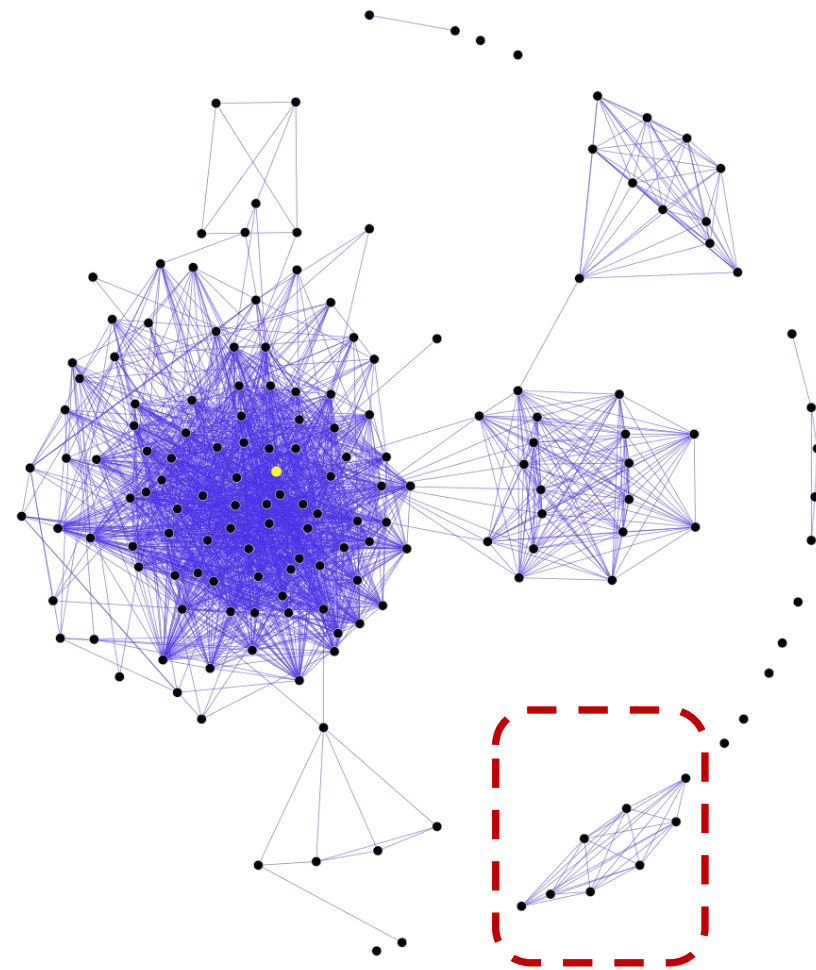


从爬取网页中获取更多URL，充实URL库，进而获取新的网页

- **网络爬虫的性能衡量**
- 数量覆盖率：“全”
  - 搜索引擎索引的网页，占目标区域中所有可能网页数量的百分比
- 质量覆盖率：“好”
  - 搜索引擎索引的网页中，“高质量”网页占目标区域中所有重要网页的百分比

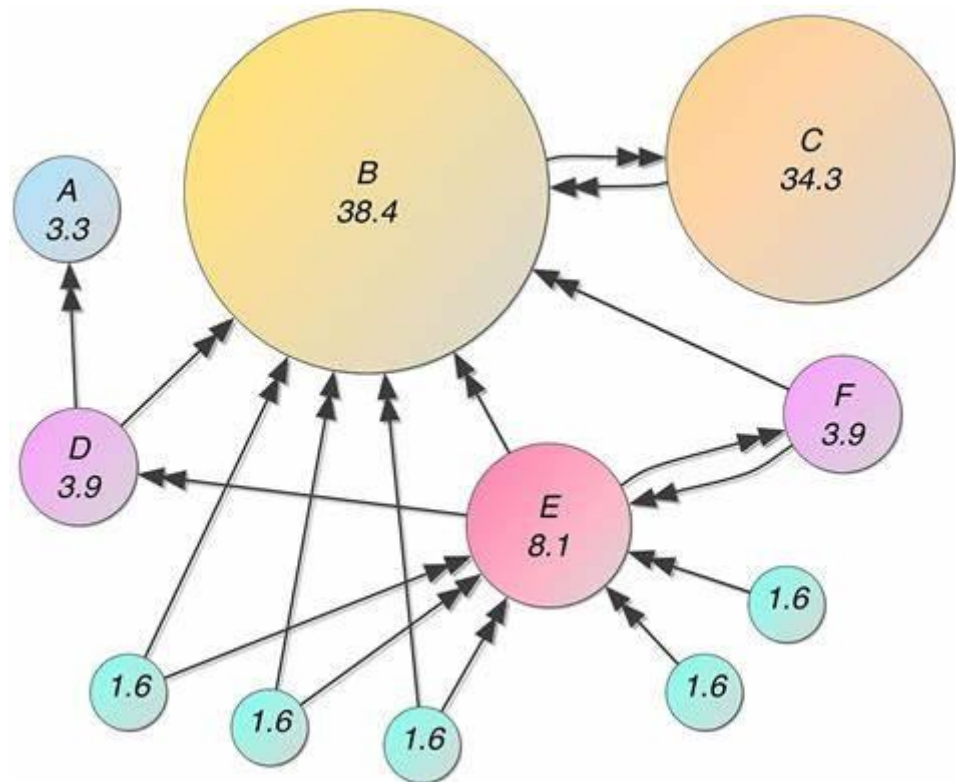


- 爬虫的数量覆盖率问题——“全”
- 遍历完备性无法保证
  - 孤立节点的存在（如右图右半部分）
- 其他影响遍历的因素
  - 部分“偏远”节点难以遍历
  - 网络结构的动态演化
  - IP/Robots等条件的约束





- 爬虫的质量覆盖率问题——“好”
- 如何度量网页的重要性？
  - 启发式方法：出入度...
  - 基于结构：PageRank/HITS...
  - 开放问题：信息密度衡量？
- 不仅重要，还要保证时效
  - “时新性”的要求

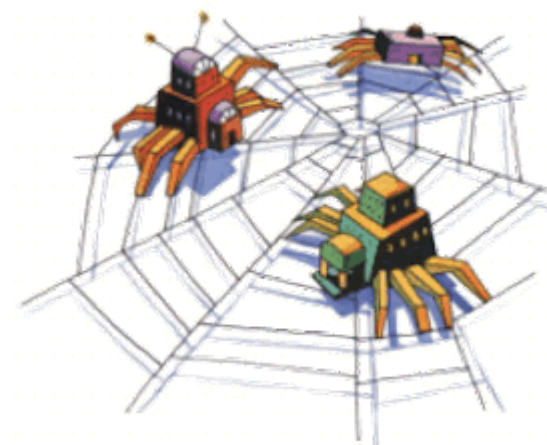


- **爬虫的主要需求**
- 速度 (Speed) : 突破网络瓶颈, 最短时间内获取所需网页
- 可扩展性 (Scalability) : 基于多爬虫机制, 有效打破单爬虫效率天花板
- 健壮性 (Robustness) : 有效应对服务器陷阱等潜在问题
- 时新性 (Freshness) : 保持内容的时效性, 提升用户体验
- 友好性 (Politeness) : 遵循网络秩序, 不影响网络服务正常运转
- 合法性 (Legitimacy) : 遵守相关法律, 不侵犯用户隐私、商业机密等

- **爬虫需求：速度**
- 2012年的统计数据：谷歌每天需要爬取200亿个页面
  - 以每周1400亿个页面计算，约为 $2^{37}$ 个页面
  - 平均每个页面64KB，约为 $2^{16}$ 个bit
    - 实际情况远远不止64KB，大量非文本内容
  - 可知每秒流量 $>100\text{Gb/s}$
- 网络瓶颈的存在对速度的限制
  - 以家用100Mbps网络为例，考虑以太网MTU为1500字节（上限了！），即使达到100%利用率，每秒也仅传输约8333个网页。



- **爬虫需求：可扩展性**
- 单一爬虫工作效率低下，难以突破效率天花板。
- 采用多爬虫机制已成为大型搜索引擎的工作常态。
- 多个爬虫带来新的挑战
  - 如何管理多个并发的连接？
  - 不同爬虫负责不同URL，如何进行划分？
  - 过多的硬件并行好处并不大
    - 抓取的性能瓶颈出现在通讯和硬盘读写



- **爬虫需求：健壮性**
- 如何有效应对爬取过程中所面临的各种风险？
  - 爬取网页时陷入回路怎么处理？
  - URL不规范如何解决？
  - 服务器陷阱如何应对？
  - 系统崩溃如何处理？



• [www.troutbums.com/Flyfactory/hatchline/hatchline/hatchline/flyfactory/flyfactory/flyfactory/flyfactory/flyfactory/flyfactory/flyfactory/hatchline](http://www.troutbums.com/Flyfactory/hatchline/hatchline/hatchline/flyfactory/flyfactory/flyfactory/flyfactory/flyfactory/flyfactory/flyfactory/hatchline)

- **爬虫需求：时新性**

- T时刻的时新性：所抓取的网页内容与T时刻该网页的最新内容一致
- 网页年龄：距离网页最近一次更新的时间
  - 通过对网页更新行为的建模来预测年龄
- 保持爬取内容的时新性能够提升搜索效果
  - 但同时，过度重视时新性将严重增加搜索引擎的负担。
  - 部分网页频繁更新，增加抓取难度



- **爬虫需求：友好性**
- 不能显著影响被爬取的服务器性能
  - 大量DNS查询可能造成类似DOS的效果
- 有些服务器可能不希望被爬取
  - Robots exclusion
  - 反面教材：360搜索、**头条搜索**

头条搜索还没有推出 但派出的ByteSpider爬虫令小网站痛苦不堪

2019年10月24日 08:16 4579 次阅读 稿源：蓝点网 1 条评论

THANK YOU

I'M SORRY

PLEASE

EXCUSE ME

- **爬虫需求：合法性**
- 法律中的规定
  - 我国：《互联网搜索引擎服务自律公约》，无针对爬虫的特定规定，但爬取的数据可能侵犯用户隐私权、违反《反不正当竞争法》等
  - 欧盟：[GDPR](#)
  - 加州：[CCPA](#)
  - .....





- **爬虫需求：合法性**

- **Linkedin vs hiQ (2019)**

hiQ 作为数据分析公司爬取 Linkedin 的数据，Linkedin 发出要求终止的信函之后限制了 hiQ 访问数据。最终最高法院判决 Linkedin 不得作出限制。

- **杭州魔蝎科技侵犯公民个人信息案 (2021)**

魔蝎科技将开发的前端插件嵌入网贷平台App中，利用各类爬虫技术，爬取用户本人账户内的通话记录、社保、公积金等各类数据，并提供给网贷平台用于判断用户的资信情况

- **爬虫需求：友好 & 合法**

一个合适的爬虫是怎样的

- 只爬取需要的数据
- 满足被爬站点的使用协议
  - 注意：Robots协议被认定为搜索引擎行业内公认的、应当被遵守的商业道德
- 不对被爬网站造成不良影响（如造成DDoS）
- 不收集可定位到个人的信息（数据脱敏/合并）
- 保护好被爬取的数据
- 

[Watch-outs for Legal and Ethical Web Scraping in 2022](#)

- **最后，老师和助教有话说**
- 友情提示：爬虫本中立，数据应保护
  - 遵守协议，不要爬取隐私、涉密数据
  - 爬取方式应友好且合法，不影响网站正常运行
  - 如果不幸吃上大碗牢饭，不把为师说出去就行了



日后你惹出祸来  
不把为师说出来就行了

- 网络爬虫的定义与需求
- **爬虫的基本要素**
- 面向API的新爬虫任务
- 常见的爬虫算法
- 常见反爬虫机制与应对策略

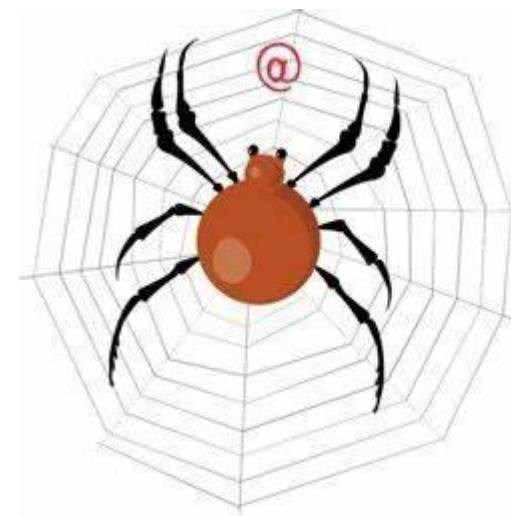
- 爬虫所涉及的协议和要素

- HTML / HTTP

- DNS / URL

- Sitemap

- Robots Exclusion



## • 如何表示与获取网页中的链接结构?

网络课堂 报考科大 科大校友 在校师生

English Version

创蒙学府 育天下英才

学校概况  
院系介绍 师资队伍 本科生教育 研究生教育 理论学习 信息公开  
学校简介 科学研究 发展规划 人才招聘 信息门户 公共服务 电子邮件  
现任领导  
管理机构  
校园地图

科大要闻 更多

专题推荐 更多

包信和校长率代表团访问澳大利亚  
中国科学技术大学召开学习贯彻习近平新时代中国特色社会主义思想...  
中国科大2023级研究生新生入学报到

学习贯彻习近平新时代中国特色社会主义思想  
中国科大2023级研究生新生入学报到

习近平总书记考察中国科大专题网

通知公告 更多

- ▶ 我校举行2023级新生升旗仪式 09-05
- ▶ 学校举行夏培肃先生百年诞辰纪念展开幕式暨夏培肃先生生平档案资料捐赠仪式 09-05
- ▶ 计算机科学与技术学院举办教师节座谈交流会暨校风传承月主题教育活动 09-09
- ▶ 一周会议安排（2023年9月4日—9月10... 09-03
- ▶ 教师节贺词 09-10
- ▶ 关于对2023年职员岗位拟聘人员进行公... 09-07

## • HTML文件示例

```
▶ <div class="container">☰</div>
▼ <nav class="main-menu">
  ▼ <div class="container" style="z-index: 10;">
    <a class="toggle"></a>
    ▼ <ul class="menu">
      ▼ <li>
        <a href="xxgk/xxjj.htm" title="学校概况">学校概况</a>
        ▼ <ul>
          ▼ <li>
            <a href="xxgk/xxjj.htm" title="学校简介">学校简介</a>
          </li>
          ▼ <li>
            <a href="xxgk/xrld.htm" title="现任领导">现任领导</a>
          </li>
          ▼ <li>
            <a href="xxgk/gljg.htm" title="管理机构">管理机构</a>
          </li>
          ▼ <li>
            <a href="xxgk/xydt.htm" title="校园地图">校园地图</a>
          </li>
        </ul>
      </li>
      ▼ <li>
        <a href="yxjs.htm" title="院系介绍">院系介绍</a>
      </li>
      ▼ <li>
        <a href target="_blank" title="师资队伍">师资队伍</a>
      </li>
    </ul>
  </div>
</nav>
```

- **HTML的基本概念**

- HyperText Markup Language, 书写网页的“框架语言”

- 基本组成: “**标记**” (Tags) + “**文本内容**” (Text)

- 例如: `<a href="http://www.ustc.edu.cn/2062/list.htm">学校概况</a>`

- 标记的作用:

- 说明网页元数据 (如“标题”等)

- 说明文本内容的布局和字体、字号等信息

- 嵌入图片、视频、创建超链接等





- **HTML的示例框架**

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 3.2 Final//EN">
```

```
<html>
```

```
<head>
```

```
    <title> This is the title but often omitted </title>
```

```
</head>
```

```
<body>
```

```
    <img src = "url1" alt="text">
```

```
    other text
```

```
    <a href = "url2" title="anchor text"> this is link text </a>
```

```
</body>
```

```
</html>
```

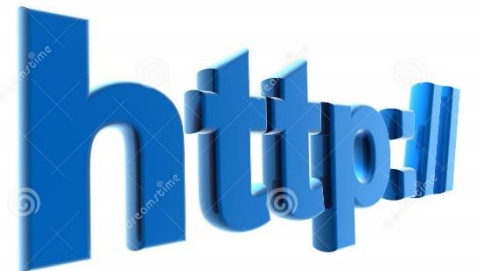


- **HTML中具有特别意义的文字**
- `<head><title> text </text></head>`
- 是搜索服务显示的内容之一（URL，标题，摘要等）
- ``
- 图片描述，可以帮助我们做“从文字到图片”的查询
- `<a href="url" title="text">link text</a>`
- 有助于理解目标网页内容及网页之间在内容上的关系



- **HTTP的基本概念**

- 超文本传输协议 (HyperText Transport Protocol)
  - 工作在TCP 之上 (请求/应答方式)
  - 容许在一个TCP 连接上发多个HTTP请求
- 工作步骤 (从客户端看)
  - 通过域名服务器 (DNS ) 得到服务器主机的IP地址
  - 用TCP和服务器建立联系
    - 发送HTTP 请求 (例如, GET或POST)
    - 接收HTTP应答头
    - 接收HTML网页内容



- HTTP的应用实例

```
telnet www.ics.uci.edu 80
Trying 128.195.1.77...
Connected to lolth.ics.uci.edu.
Escape character is '^]'.
GET http://www.ics.uci.edu/ HTTP/1.1
Host: www.ics.uci.edu

HTTP/1.1 200 OK
Date: Wed, 25 Sep 2002 19:43:12 GMT
Server: Apache/1.3.26 (Unix) PHP/4.1.2 mod_ssl/2.8.10 OpenSSL/0.9.6e
X-Powered-By: PHP/4.1.2
Transfer-Encoding: chunked
Content-Type: text/html

f00
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.0 Transitional//EN">
<html>
<head>
<title>Information and Computer Science at the University of
California, Irvine</title>
...
```

*Refer to [RFC 2616](#)*

User agent's request

Server's response

HTML code of returned webpage

Example of the use of the GET method in an HTTP 1.1 session.

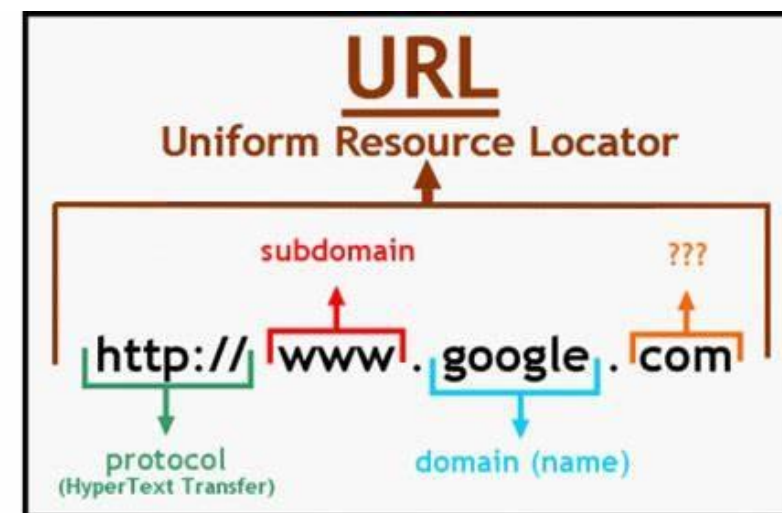
## • HTTP的状态代码分类

分类	分类描述
1**	信息，服务器收到请求，需要请求者继续执行操作
2**	成功，操作被成功接收并处理
3**	重定向，需要进一步的操作以完成请求
4**	客户端错误，请求包含语法错误或无法完成请求
5**	服务器错误，服务器在处理请求的过程中发生了错误

重点记住一些特殊的状态码，如：200 (OK) , 403 (Forbidden) , 404 (Not Found) .....

- **URL的基本概念**

- 统一资源定位符 (Universal Resource Locator)
  - 以某种统一的 (标准化的) 方式标识资源的简单字符串。
- URL一般由四部分组成：
  - 访问资源的模式/协议。
  - 存放资源的主机名。
  - 资源自身的名称, 由路径表示。
  - 被访问的文件名或主页名。



## • URL实例

- URL例子: `http://cs.ustc.edu.cn/3058/list.htm`



- 这是一个可通过HTTP协议获得的文档。
- 存放在名为`cs.ustc.edu.cn`的主机名。
- 通过路径`3058`（对应学院新闻栏目）访问。

- **URL与IP地址的对应关系**

- 一般而言，域名与IP地址处于一一对应的关系
- 实际情况下，域名与IP地址之间存在复杂的对应关系
  - 一对一：基本的对应关系
  - 一对多：可能由于虚拟主机导致，使得多个URL映射到同一IP
    - 例如，[www.pku.edu.cn](http://www.pku.edu.cn) / [www.gh.pku.edu.cn](http://www.gh.pku.edu.cn) 等都映射到162.105.129.12上，但由于各个站点内容不同，所以被认为是不同的Web服务器
  - 多对一：可能由于DNS轮询导致，以应对商业站点的负载问题
    - 例如，多个182.61.200.x均与[www.baidu.com](http://www.baidu.com) 映射。



- **URL链接提取与规范化**

- 目标：得到网页中所含URL的标准型
  - URL的处理与过滤，避免多次获得不同URL指向的相同网页
  - 相对URL：基础URL被省略，需要进行补齐
- URL的不规范现象
  - 不同URL可能指向同一个网页
  - URL的格式存在错误或无用内容
    - 多余的文件名（如index.html），无用的查询变量或空的查询条件
  - 动态网页与动态参数
  - **短链接**

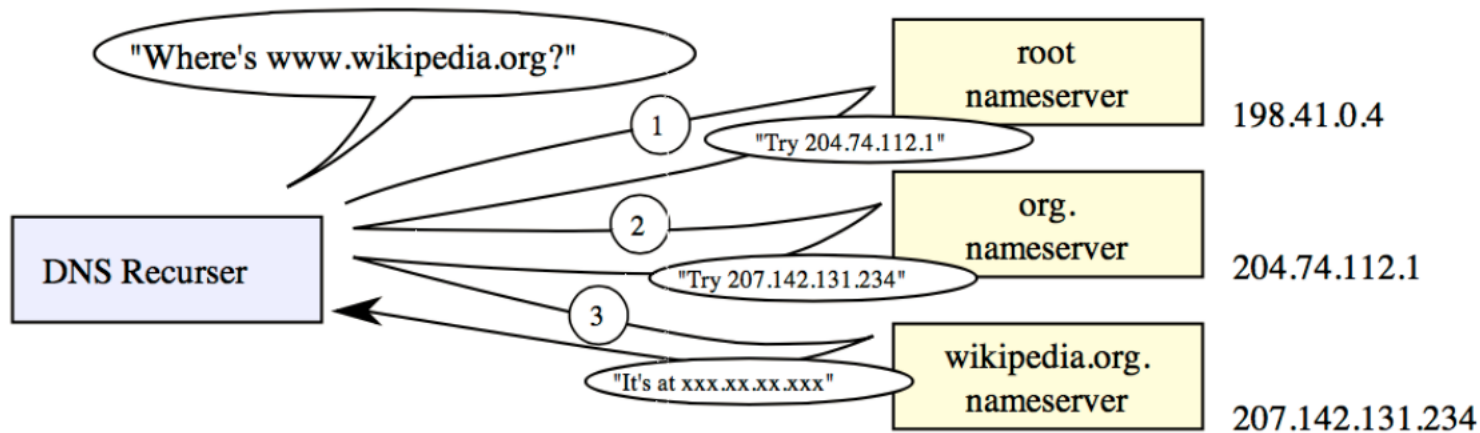
- **一些URL规范化的操作**

- URL的基本组成：协议:// 主机名[: 端口]/ 路径/[: 参数] [? 查询]#Fragment
  - protocol ://hostname[:port]/path/ [;parameters][?query]#fragment
- 一些常见的URL规范化示例
  - URL协议名和主机名的小写化
    - [HTTP://WWW.BAIDU.COM](http://www.baidu.com) -> <http://www.baidu.com>
  - 删除Fragment (#) 或多余的查询串, 如?, &
    - <http://www.example.com#seo> -> <http://www.example.com>
  - 删除默认后缀或多余的www
    - <http://www.example.com/index.html> -> <http://example.com>

- **DNS的基本概念和作用**

- 域名系统 (Domain Name System)

- 将域名和IP地址相互映射的一个分布式数据库。
- DNS地址解析可能成为重要瓶颈，甚至造成类似DOS的攻击效果



- **提升DNS性能的方法**
- 如何提高DNS解析模块的性能?
  - 并行 DNS Client
  - 缓存 cache DNS results
  - 预取 prefetch client

- **并行 DNS Client**
- 并行的地址解析Client
  - 专门对付多个请求的并行处理
  - 容许一次发出多个解析请求
  - 协助在多个DNS server之间做负载分配
    - 例如，根据掌握的URL进行适当调度

- **缓存 cache DNS results**
- 增加DNS缓存的重要性
  - 面向海量URL和主机的搜索任务，如果没有DNS缓存，会造成频繁查询DNS服务器，从而造成类似于拒绝服务攻击（DOS）的副作用。
  - 针对小规模网页搜索（如百万量级），可利用建立在内存中的DNS映射，既能加快网页信息的获取，后能降低对于DNS服务器的压力。
- DNS缓存的内容
  - Internet的DNS系统会定期刷新，交换更新的域名和IP信息
  - 缓存中有部分信息过期影响不大，但注意要适度刷新

- **预取 prefetch client**

- 为了减少等待查找涉及新主机的地址的时间，尽早将主机名投给DNS系统
- 与缓存系统的区别：缓存 – 用了才缓存；预取 – 不用也暂时先缓存
- DNS预取的基本步骤
  - 分析刚刚得到的网页
  - 从<href>属性中提取主机名
  - 向缓存服务器提交DNS解析请求
  - 结果存放在DNS缓存中（也可能用上，也可能用不上）

- **Sitemap, 站点地图**

- 放在服务器根目录中的sitemap.xml, 为爬虫指明抓取的建议
  - 协助爬虫找到隐藏的网页 (如需要查询或表单才能够访问的内容)
  - Changelist属性还提供有关更新频率的说明, 以协助保障时新性

```
<?xml version="1.0" encoding="UTF-8"?>
<urlset xmlns="http://www.sitemaps.org/schemas/sitemap/0.9">
  <url>
    <loc>http://www.company.com/</loc>
    <lastmod>2008-01-15</lastmod>
    <changelist>monthly</changelist>
    <priority>0.7</priority>
  </url>
  <url>
    <loc>http://www.company.com/items?item=truck</loc>
    <changelist>weekly</changelist>
  </url>
  <url>
    <loc>http://www.company.com/items?item=bicycle</loc>
    <changelist>daily</changelist>
  </url>
</urlset>
```



- Sitemap, 站点地图

- 案例：中央政府门户网站的站点地图-百度版本 (<http://www.gov.cn/baidu.xml>)

```
▼<urlset>
  ▼<url>
    <loc>http://www.gov.cn/xinwen/2021-09/13/content_5636996.htm</loc>
    <lastmod>2021-09-13T11:34:06-08:00</lastmod>
  </url>
  ▼<url>
    <loc>http://www.gov.cn/xinwen/2021-09/13/content_5636985.htm</loc>
    <lastmod>2021-09-13T10:34:16-08:00</lastmod>
  </url>
  ▼<url>
    <loc>http://www.gov.cn/xinwen/2021-09/13/content_5636984.htm</loc>
    <lastmod>2021-09-13T10:31:16-08:00</lastmod>
  </url>
  ▼<url>
    <loc>http://www.gov.cn/xinwen/2021-09/13/content_5636981.htm</loc>
    <lastmod>2021-09-13T10:28:06-08:00</lastmod>
  </url>
  ▼<url>
    <loc>http://www.gov.cn/xinwen/2021-09/13/content_5636980.htm</loc>
    <lastmod>2021-09-13T10:26:34-08:00</lastmod>
  </url>
  ▼<url>
    <loc>http://www.gov.cn/xinwen/2021-09/13/content_5636979.htm</loc>
    <lastmod>2021-09-13T10:26:09-08:00</lastmod>
  </url>
  ▼<url>
    <loc>http://www.gov.cn/xinwen/2021-09/13/content_5636977.htm</loc>
    <lastmod>2021-09-13T10:20:15-08:00</lastmod>
  </url>
```

- **Robots Exclusion, 排斥协议**

- Sitemap是允许协议, 而Robots是排斥协议
- 在服务器文档根目录中的文件robots.txt, 包含一个路径前缀表, 描述了服务器给出的抓取限制
- 例如, 腾讯新闻的Robots.txt ([news.qq.com/robots.txt](http://news.qq.com/robots.txt))

---

```
User-agent: *  
Disallow:  
Sitemap: //www.qq.com/sitemap_index.xml  
Sitemap: http://news.qq.com/topic_sitemap.xml
```

- **Robots Exclusion, 排斥协议**

- 案例：中央政府门户网站的Robots.txt ([www.gov.cn/robots.txt](http://www.gov.cn/robots.txt))

```
User-agent: *
Allow: /1
Sitemap:http://www.gov.cn/baidu.xml
Sitemap:http://www.gov.cn/google.xml
Sitemap:http://www.gov.cn/bing.xml
Sitemap:http://www.gov.cn/sogou.xml
Disallow:/2016gov/
Disallow:/2016shuju/
Disallow:/2016guoqing/
Disallow:/2016zhengce/
Disallow:/2016hudong/
Disallow:/2016fuwu/
Disallow:/2016xinwen/
Disallow:/premier/
Disallow:/2016guowuyuan/
Disallow:/2016public/
Disallow:/2016ducha/
Disallow:/guowuyuan/yangjing/
Disallow:/guowuyuan/yj.htm
Disallow:/c16629/gwyld_yj.htm
Disallow:/guowuyuan/yj_hy.htm
Disallow:/guoqing/2013-03/16/content_2583121.htm
Disallow:/guowuyuan/2015-12/24/content_5027563.htm
```

## • Robots Exclusion, 排斥协议

- 该限制只针对爬虫，一般浏览不受影响
- 该限制属于“君子协定”，没有强制执行力
  - 但目前，绝大多数搜索引擎都遵守该限制



tmall.com/robots.txt

```
User-agent: *  
Disallow: /
```



[淘宝商城 天猫](https://www.tmall.com/)

<https://www.tmall.com/> - 2909条评价

由于该网站的robots.txt文件存在限制指令（限制搜索引擎抓取），系统无法提供该页面的内容描述 - [了解详情](#)

[爱淘宝-淘宝网购物分享平台](https://ai.taobao.com/)

<https://ai.taobao.com/> - 305条评价

由于该网站的robots.txt文件存在限制指令（限制搜索引擎抓取），系统无法提供该页面的内容描述 - [了解详情](#)

[淘宝网 - 淘!我喜欢](https://mai.taobao.com/)

[mai.taobao.com/](https://mai.taobao.com/) - 8483条评价

由于该网站的robots.txt文件存在限制指令（限制搜索引擎抓取），系统无法提供该页面的内容描述 - [了解详情](#)

- 网络爬虫的定义与需求
- 爬虫的基本要素
- **面向API的新爬虫任务**
- 常见的爬虫算法
- 常见反爬虫机制与应对策略

- **API的基本概念**

- 应用程序接口（Application Programming Interface）
- 2000年，Salesforce和eBay推出了自己的API，程序员可以访问并下载一些公开数据。从那时起，许多网站都提供API用于访问公共数据库。
- 通过开放的API来吸引更多的用户和更多的创意，与具备分享、标准、去中心化、开放、模块化的Web 2.0时代相得益彰，为创造者和平台都带来价值。
- 同时，网页API也为开发人员提供了一种更友好的网络爬虫方式：
  - 直接获取资源列表，清晰、明确
  - 接收JSON或XML等格式化数据的反馈



## • 基于API的数据爬取示例

```
In [28]: import urllib.request as request
import json
url = 'https://api.douban.com/v2/movie/top250'
crawl_content = request.urlopen(url).read()
top20 = json.loads(crawl_content.decode('utf8'))['subjects']
for movie in top20:
    url = 'https://api.douban.com/v2/movie/' + movie['id']
    movieContent = request.urlopen(url).read()
    print(json.loads(movieContent.decode('utf8')))
```

```
{'rating': {'max': 10, 'average': '9.6', 'numRaters': 947593, 'min': 0}, 'author': [{'name': '弗兰克·德拉邦特 Frank Darabont'}], 'alt_title': '肖申克的救赎 / 月黑高飞(港)', 'image': 'https://img3.doubanio.com/view/photo/s_ratio_poster/public/p480747492.jpg', 'title': 'The Shawshank Redemption', 'summary': '20世纪40年代末,小有成就的青年银行家安迪(蒂姆·罗宾斯 Tim Robbins 饰)因涉嫌杀害妻子及她的情人而锒铛入狱。在这座名为肖申克的监狱内,希望似乎虚无缥缈,终身监禁的惩罚无疑注定了安迪接下来灰暗绝望的人生。未过多久,安迪尝试接近囚犯中颇有声望的瑞德(摩根·弗里曼 Morgan Freeman 饰),请求对方帮自己搞来小锤子。以此为契机,二人逐渐熟稔,安迪也仿佛在鱼龙混杂、罪恶横生、黑白混淆的牢狱中找到属于自己的求生之道。他利用自身的专业知识,帮助监狱管理层逃税、洗黑钱,同时凭借与瑞德的交往在犯人中间也渐渐受到礼遇。表面看来,他已如瑞德那样对那堵高墙从憎恨转变为处之泰然,但是对自由的渴望仍促使他朝着心中的希望和目標前进。而关于其罪行的真相,似乎更使这一切朝前推进了一步.....\n本片根据著名作家斯蒂芬·金(Stephen Edwin King)的原著改编。', 'attrs': {'pubdate': ['1994-09-10(多伦多电影节)', '1994-10-14(美国)'], 'language': ['英语'], 'title': ['The Shawshank Redemption'], 'country': ['美国'], 'writer': ['弗兰克·德拉邦特 Frank Darabont', '斯蒂芬·金 Stephen King'], 'director': ['弗兰克·德拉邦特 Frank Darabont'], 'cast': ['蒂姆·罗宾斯 Tim Robbins', '摩根·弗里曼 Morgan Freeman', '鲍勃·冈顿 Bob Gunton', '威廉姆·赛德勒 William Sadler', '克兰西·布朗 Clancy Brown', '吉尔·贝罗斯 Gil Bellows', '马克·罗斯顿 Mark Rolston', '詹姆斯·惠特摩 James Whitmore', '杰弗里·德曼 Jeffrey DeMunn', '拉里·布兰登伯格 Larry Brandenburg', '尼尔·吉恩托利 Neil Giuntoli', '布赖恩·利比 Brian Libby', '大卫·普罗瓦尔 David Proval', '约瑟夫·劳格诺 Joseph Ragno', '祖德·塞克利拉 Jude Ciccolella'], 'movie_duration': ['142 分钟'], 'year': ['1994'], 'movie_type': ['犯罪'], 'id': 'https://api.douban.com/movie/1292052', 'mobile_link': 'https://m.douban.com/movie/subject/1292052/', 'alt': 'https://movie.douban.com/movie/1292052', 'tags': [{'count': 194779, 'name': '经典'}, {'count': 165432, 'name': '励志'}, {'count': 147828, 'name': '信念'}, {'count': 133384, 'name': '自由'}, {'count': 98525, 'name': '美国'}, {'count': 94751, 'name': '人性'}, {'count': 73431, 'name': '人生'}, {'count': 60696, 'name': '剧情'}]}
```

基于豆瓣提供的API, 获得评分最高的20部电影信息

- 基于API的数据爬取一般流程



其中，第3-4步可能存在循环，即通过获得某个资源ID对应的详细信息，获得更多其他资源ID，并继续抓取更多信息。



- **Token的概念、作用与取得方式**

- API的调用需要用户身份认证（用户授权），而Token相当于授权后的通行证
- 使用Token认证而非用户名/密码认证方式的优点
  - Token 的生成完全独立于帐号密码，往来通信不会影响账号安全
  - 即使 Token 丢失或泄露，只需登录后台删除该Token即可



目前微博的API认证授权机制

- **结构化数据反馈 (1) : XML**

- 可扩展标记语言 (Extensible Markup Language)
  - 是一种允许用户对自己的标记语言进行定义的源语言。
  - 提供统一的方法来描述和交换独立于应用程序的结构化数据。
- XML相比于HTML的优势所在：
  - 可扩展性方面：HTML不允许用户自行定义标识或属性，而在XML中，用户能够根据需要自行定义新的标识及属性名
  - 结构性方面：HTML不支持深层的结构描述，XML的文件结构嵌套可以复杂到任意程度，能表示面向对象的等级层次。
  - 可校验性方面：HTML没有提供规范文件以支持结构校验，而XML文件可以包括一个语法描述，使应用程序可以对此进行结构校验。



## XML文件的实例与基本要求

```
<?xml version="1.0" encoding="UTF-8"?>
- <item>
- <venue>
  <address_1>162 Winn St</address_1>
  <state>MA</state>
  <zip>01803</zip>
  <lat>42.504240</lat>
  <repinned>False</repinned>
  <name>American Legion Hall</name>
  <city>Burlington</city>
  <id>486621</id>
  <country>us</country>
  <lon>-71.185790</lon>
</venue>
- <fee>
  <label>Price</label>
  <accepts>amazon</accepts>
  <currency>USD</currency>
  <description>per person</description>
  <amount>10.0</amount>
  <required>0</required>
</fee>
<status>past</status>
<description><b>Looks like the storm predicte
confirm the band. More details will follow.
holidays with old friends and new at a spe
American Legion Hall in Burlington (right c
be a cash bar. You are welcome to invite fa
done so on your reply, it will help me keep
on your reply, &quot;paying by check&quo
soon as you know you can attend. If you h
we had a gift exchange which was alot of f
receiving yourself. <br />Hope it will be a
<how_to_find_us>We have rented the hall...so
```

必须有XML声明语句作为文档开头

良好的XML文档有且只有一个根节点

例如，在本实例中为item

可根据需要嵌套文件结构

例如，<venue>中包含<city>信息

所有的标记必须有相应的结束标记

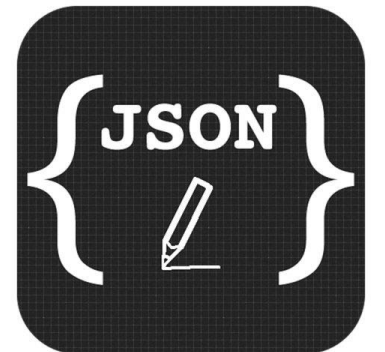
例如，<zip>与</zip>的成对出现

注意，所有的空标记也必须被关闭

## • 结构化数据反馈 (2) : JSON

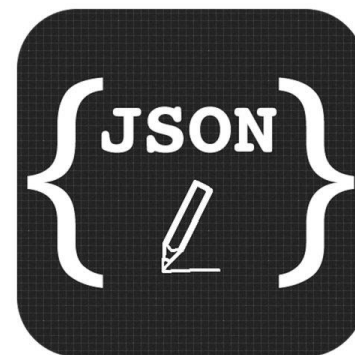
- JS 对象表示法 (JavaScript Object Notation)
  - 轻量级的数据交换格式，采用完全独立于编程语言的文本格式。
  - 简洁和清晰的层次结构使得 JSON 成为理想的数据存储和交换语言。
  - 易于人阅读和编写，同时也易于机器解析，并有效地提升网络传输效率。

```
"name":"Michael",  
"address":  
{  
  "city":"Beijing",  
  "street":" Chaoyang Road ",  
  "postcode":100025  
}
```



## • JSON与XML的区别

- XML的优点：格式统一、标准，容易与其他系统进行远程交互，方便共享
- XML的缺点：文件庞大，文件格式复杂，解析复杂且方式繁多，工作量大
- JSON的优点
  - 数据格式比较简单，易于读写，格式经过压缩，占用带宽小；
  - 支持多种语言，易于解析，可以简单的进行数据的读取；
  - 能直接为服务器端代码使用，大大简化了开发和维护工作。



## • JSON的基本元素

```
1 <?xml version="1.0" encoding="utf-8" ?>
2 <root>
3   <resultcode>200</resultcode>
4   <reason>Return Successd!</reason>
5   <result>
6     <area>江苏省苏州市</area>
7     <location>电信</location>
8   </result>
9 </root>
```

## XML与JSON 格式对比

```
1 {
2   "resultcode": "200",
3   "reason": "Return Successd!",
4   "result": {
5     "area": "江苏省苏州市",
6     "location": "电信"
7   }
8 }
```

- JSON中的六大构造字符
- “[ ” ]” 表示数组开始/结束
- “{ ” }” 表示对象开始/结束
- “: ” 表示名称的分隔符, “ , ” 表示值的分隔符
- 从内容表示上看, JSON更为轻量, 但在无缩进情况下层次化解析相对困难

- 网络爬虫的定义与需求
- 爬虫的基本要素
- 面向API的新爬虫任务
- **常见的爬虫算法**
- 常见反爬虫机制与应对策略

- 最最基础的算法

PROCEDURE SPIDER 1 (G)

Let ROOT := any URL from G

Initialize STACK <stack data structure>

Let STACK := push(ROOT, STACK)

Initialize COLLECTION <big file of URL-page pairs>

While STACK is not empty,

*URLcurr* := pop(STACK)

PAGE := look-up(*URLcurr*)

STORE(<*URLcurr*, PAGE>, COLLECTION)

For every *URLi* in PAGE,

    push(*URLi*, STACK)

Return COLLECTION

潜在风险:

重复收集的问题?

回路或不连通的解决方法?

如何控制搜集特定的一部分?



- 最最基础的改进算法

PROCEDURE SPIDER(G, {SEEDS})

Initialize COLLECTION <big file of URL-page pairs>

Initialize VISITED <big hash-table>

For every ROOT in SEEDS

    Initialize STACK <stack data structure>

    Let STACK := push(ROOT, STACK)

    While STACK is not empty,

        Do *URLcurr* := pop(STACK)

            Until *URLcurr* is not in COLLECTION

                insert-hash(*URLcurr* , VISITED)

                PAGE := look-up(*URLcurr* )

                STORE(<*URLcurr* , PAGE>, COLLECTION)

                For every *URLi* in PAGE,

                    push(*URLi* , STACK)

Return COLLECTION

解决方案:

利用BigTable排除重复部分

基于种子集合控制爬取内容

通过更换种子重新启动解决回路等

- **重复性的监测问题**

- 不仅URL本身可能重复，即使不同URL，也可能导致相似的文档内容
- 完全重复文档的检测方法
  - 检验和（Checksumming），所有字节进行加和，相同文档有相同检验和
  - 缺陷：相同检验和不代表相同（如打乱顺序）
- 近似重复文档的检测方法：指纹表示法
  - 1) 对文档进行分词处理，并进行n-gram组合
  - 2) 挑选部分n-gram用于表示这一文档
  - 3) 对被选中的n-gram进行散列，以提升检索效率和减少指纹大小
  - 4) 存储散列值作为文档指纹，通常存储在倒排索引中

- 重复性的监测问题

Tropical fish include fish found in tropical environments around the world, including both freshwater and salt water species.

a) 原始文本

tropical fish include, fish include fish, include fish found, fish found in, found in tropical, in tropical environments, tropical environments around, environments around the, around the world, the world including, world including both, including both freshwater, both freshwater and, freshwater and salt, and salt water, salt water species

b) 3-gram

938 664 463 822 492 798 78 969 143 236 913 908 694 553 870 779

c) 散列值

664 492 236 908

d) 使用  $0 \bmod 4$  选择的散列值

图3-14 指纹生成过程实例

# 更进一步的遍历策略，从某博士读论文开始说起

tion of benchmark datasets, it is also encouraged that statistics information including geography, gender, ethnicity and other demographic information should be provided, for those datasets containing information about people.

### 5.2 Intersectional Fairness

The investigation of intersectional fairness, i.e., combination of multiple sensitive attributes, is relatively lacking in current research. Take bias mitigation for example, current work generally focus on one kind of bias. Although this may increase model fairness in terms of a specific bias, it is highly possible that the model is still biased from the intersectional perspective. For instance, a DNN classifier is fair to women, while exhibiting discrimination towards a subminority group, e.g., African American women or women over the age of 60. Similarly for DNN based job recruiting tool, even if the biased model is free of gender bias, it is hard to guarantee that the model is not biased towards other protected attributes, e.g., race, age. More work is needed to figure out methods which are effective for identification and mitigation of intersectional bias.

### 5.3 Fairness and Utility Trade-off

The removal of bias could possibly hurt the model's ability for main prediction task. For instance, adversarial training could increase fairness with respect to demographic parity measurement. A possible deficiency of this mitigation solution is that it could compromise overall prediction accuracy, especially the accuracy for non-protected groups. Thus this might undermine the principle of beneficence. It remains a challenge to simultaneously reduce unintended bias and maintain satisfactory model prediction performance.

### 5.4 Formalization of Fairness

As the field of fairness machine learning is evolving quickly, there is still no consensus about the measurements of fairness. In certain cases, some measurements could be conflicting with others. A model may be fair in terms of one metric, but may lead to other sorts of unfairness. For instance, a loan approval tool may satisfy demographic parity measurement, while violating equality of opportunity measurement. There is no silver bullet, and each application domain calls for a fairness measurement or combination of measurements which meet its specific requirements [2].

### 5.5 Fairness in Large-scale Training

Large-scale training is employed in some domains to boost model performance. Take NLP domain for instance, current paradigm is to pre-train language models (e.g. BERT [3] and XLNet [24]) on large-scale text corpus, which will be further fine-tuned on downstream tasks such as machine translation. These powerful language models could capture biases and propagate them to other tasks. Since these models need to be trained on corpus with billion-scale words and are typically trained for days, bias mitigation either through data preprocessing or training regularization remains a challenge and more research is needed in this direction.

### 6. CONCLUSIONS

With increasing adoption of DNNs in high-stake real world applications, e.g., job hunting, criminal justice and loan approval, the undesirable algorithmic unfairness problem has

attracted much attention recently. We present a clear categorization of unfairness and introduce the most widely used measurements of fairness. By introducing interpretability as an essential ingredient, we also give a comprehensive overview of existing bias detection and mitigation techniques from the computational perspective. The model bias to some extent exposes biases present in our society. To really benefit our society, DNN models are supposed to reduce these biases instead of amplifying biases. In future, endeavors from different disciplines, including computer science, statistics, cognitive science, should be joined together to eliminate disparity and promote fairness. In this way, DNN systems could be readily applied for fairness sensitive applications and really improve benefits of all groups.

### 7. REFERENCES

- [1] A. Beutel, P. Cui, S. G. Jiang, F. He, A. He, W. Wen, C. Lu, F. Koehnemann, J. Biegel, and E. H. Chi. Putting fairness in perspective: Challenges, metrics, and techniques. *arXiv preprint arXiv:1808.05671*, 2018.
- [2] A. Beutel, J. Chen, Z. Zhao, and E. H. Chi. Data decision and theoretical implications when adversarially learning fair representations. *Fairness, Accountability, and Transparency in Machine Learning (FAITML)*, 2017.
- [3] S. L. Biderman, L. Green, and B. O'Connor. Demographic dialectal variation in social media: A case study of afro-american english. *2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2018.
- [4] T. Bolukbasi, K.-W. Chang, J. Y. Zou, V. Saligrama, and A. T. Kabaj. Man is to computer programmer as woman is to housewife? obtaining word embeddings in *Flaircraft Conference on Neural Information Processing Systems (NIPS)*, 2016.
- [5] J. Boudouani and T. Cohen. Gender shades: Intersectional neutrality diagnosis in convolutional gender classification. In *Conference on Fairness, Accountability and Transparency (FAIT)*, pages 77–91, 2018.
- [6] T. Calder and S. Verwer. Three naive bayes approaches for discriminative classification. *Data Mining and Knowledge Discovery*, pages 277–292, 2010.
- [7] Y. Chen, P. Sclafone, and M. Ghahramani. Can it help reduce disparities in general medical and mental health care? *AMA journal of ethics*, 2019.
- [8] A. Das, A. Dashtnoveh, and F. Brnoomeh. Mitigating bias in gender, age and ethnicity classification: a multi-task, convolution neural network approach. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [9] J. Devlin, M. W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *2019 Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2019.
- [10] M. Du, N. Liu, and X. Hu. Techniques for interpretable machine learning on downstream tasks. *arXiv preprint arXiv:1808.00273*, 2018.
- [11] M. Du, N. Liu, Q. Sang, and X. Hu. Towards explanation of model-based prediction with guided feature importance. *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, 2018.
- [12] M. Du, N. Liu, F. Yang, and X. Hu. Learning credible deep neural networks with relational regularization. In *IEEE International Conference on Data Mining (ICDM)*, 2019.
- [13] M. Du, N. Liu, F. Yang, S. Ji, and X. Hu. On attribution of recurrent neural network predictions: An additive decomposition. In *The World Wide Web Conference (WWW)*, 2019.
- [14] C. Dhanraj and A. Wong. Auditing inquest: Towards a model-driven framework for auditing demographic

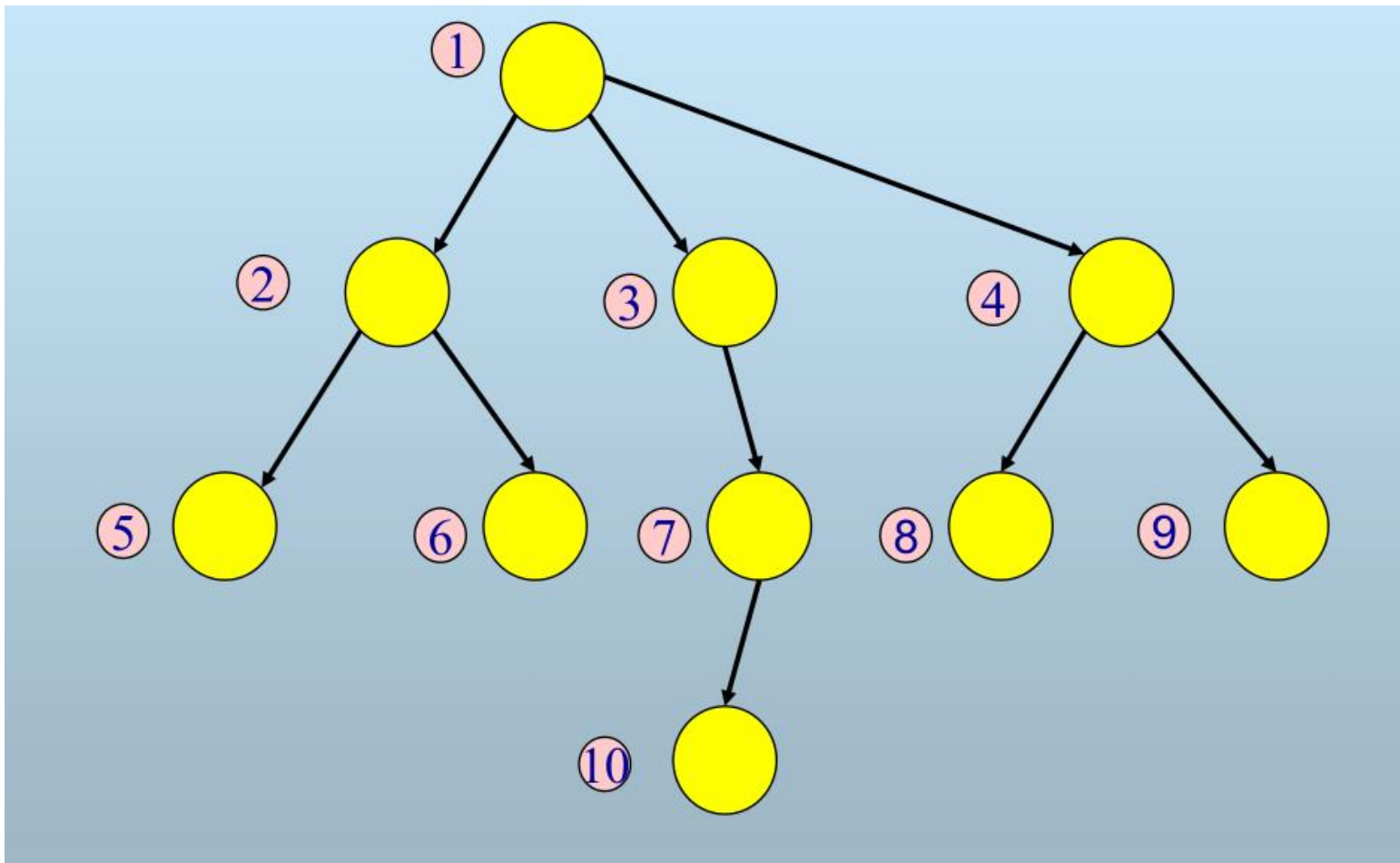


### SIGIR '17, August 07-11 2017, Shinjuku, Tokyo, Japan/Ku Chen\*, Yongfeng Zhang\*, Qingyao Ai, Hongfeng Xu, Junchi Yan, and Zheng Qin†

- [1] B. H. Fan, and Z. Zhang. Siftnet: Small-scale exploring ratings and reviews for recommendation. In *ACM paper*, page 1–6, 2014.
- [2] J. Cao, J. Li, X. Guo, S. Bao, H. Ji, and J. Tang. Stone hit only? target user identification in short texts. In *EMNLP paper*, page 1–6, 2013.
- [3] Y. Chen, Z. Qin, Z. Zhang, and Y. Tang. Learning to rank features for recommendation. In *Proceedings of the 20th International ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 119–127, 2014.
- [4] Y. Chen, F. Wang, Z. Qin, and Y. Zhang. Hyper: A hybrid local relevance personal ranking method. In *Proceedings of the 20th International Conference Companion on World Wide Web*, pages 31–32, International World Wide Web Conference Steering Committee, 2014.
- [5] K. Cho, E. Van Den Broeck, C. Galambos, D. Baheti, F. Bagnato, H. Schleich, and Y. Bergin. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1607.05934*, 2016.
- [6] P. Compton, J. Adams, and E. Sarguz. Deep neural networks provide a statistical machine translation. *arXiv preprint arXiv:1611.00148*, 2016.
- [7] Q. Dai, M. Cao, C. Wu, A. Loshchilov, J. Jia, and C. Wang. Jointly modeling aspects, ratings and reviews for movie recommendation (joint). In *Proceedings of the 2018 AAAI Conference on Artificial Intelligence*, pages 4916–4922, 2018.
- [8] X. Fan, T. Tang, M.-K. Wu, and S. Wang. Adaptive feature extraction for video recommendation using an aggregation mechanism. *Journal of Visual Communication and Image Representation*, 157(10):1046–1052, 2018.
- [9] G. Gama, N. Elhadad, and A. Mannes. Beyond the state improving rating prediction using review text content. In *WWW*, volume 3, pages 1–6, 2009.
- [10] C. Glaser, A. Borden, and Y. Bergin. Deep sparse recurrent neural networks. In *Empirical Video and Action Recognition, 2009 Proceedings. IEEE Conference on Computer Vision and Pattern Recognition*, pages 174–181, 2009.
- [11] Y. Gong and X. Liu. Video recommendation using singular value decomposition. In *Empirical Video and Action Recognition, 2009 Proceedings. IEEE Conference on Computer Vision and Pattern Recognition*, pages 178–189, 2009.
- [12] H. Guo, Y. Zhang, S. Fan, and Y. Tang. Deep neural learning for image recognition. *arXiv preprint arXiv:1212.0410*, 2012.
- [13] H. Guo and Y. Zhang. Why visual features pervaded ranking from implicit feedbacks. *arXiv preprint arXiv:1304.0744*, 2013.
- [14] H. Guo, T. Chen, H. Li, Fan, and C. Chen. Textual feature-aware relevance recommendation by modeling aspects. In *Proceedings of the 20th ACM International Conference on Information and Knowledge Management*, pages 1041–1048, 2011.
- [15] H. Guo, M. Guo, M. Cui, Fan, and Y. Wang. Break: Towards ranking on bipartite graphs. In *Proceedings of the 2014 ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 97–105, 2014.
- [16] X. He, L. Lian, H. Zhang, L. Nie, X. Hu, and T. S. Chua. Neural collaborative filtering. In *Proceedings of the 20th International Conference Companion on World Wide Web (WWW)*, pages 173–182, International World Wide Web Conference Steering Committee, 2012.
- [17] H. He, H. Zhang, M. Kuo, and T. S. Chua. Fast matrix factorization for online recommendation with implicit feedback. In *Proceedings of the 20th International ACM SIGKDD Conference on Knowledge and Information Discovery*, pages 145–154, 2012.
- [18] J. Heisterkamp and J. Scholte. Long short-term memory. *Neural computation*, pages 1416–1439, 1997.
- [19] J. H. Holland. *Cognitive Psychology*. Long R. Gordin, G. Gendreau, and S. T. Darnell. CAFE: Computational architecture for fast feature combining. In *Proceedings of the 2nd ACM International Conference on Knowledge Discovery and Data Mining*, pages 679–691, ACM, 2004.
- [20] H. M. King, R. F. W. McKenzie, and J. Crocker. Advances in video recommendation and sharing. In *Recent Advances in Multimedia Information Processing and Communication*, pages 16–26, Springer, 2009.
- [21] W. Jiang, C. Cortes, and A. C. Liou. Automatic content video summarization on mobile and visual analysis. In *ICCV*, pages 2099–2106, 2011.
- [22] A. Khosh, R. Hariri, C. Li, and N. Soudarshan. Large-scale video summarization using multi-stage analysis. In *CVPR*, pages 2699–2706, 2011.
- [23] G. Kim, I. Song, and E. Yang. Joint summarization of large-scale collections of web images and their text descriptions. In *CVPR*, 2014.
- [24] M. Li, J. Dai, and Y. Zhang. Efficiently learning neural network models for ad-hoc relevance video summarization. In *CVPR*, pages 1346–1353, 2012.
- [25] L. Li, J. Dai, and Y. Zhang. Efficiently learning neural network models for ad-hoc relevance video summarization. In *CVPR*, pages 1347–1354, 2012.
- [26] L. Li, J. Dai, Zhang, and Y. Zhang. An overview of video summarization techniques. Technical Report HPL-2001-191, HP Laboratory, 2001.
- [27] L. Li, J. Dai, Zhang, W. Song, and X. Du. Combining user preferences and user opinions for accurate recommendation. *Electronic Commerce Research and Applications*, 13(1):1–13, 2013.
- [28] F. Li, X. X. He, L. Wang, and H. J. Zhang. A generic framework of user-centric model and its application in video summarization. *IEEE Transactions on Multimedia*, 13(5):947–953, 2011.
- [29] J. McInerney and L. Jordan. Hidden factors and hidden topics: understanding rating dimensions with review text. In *Proceedings of the 20th International Conference on Knowledge Discovery and Data Mining*, pages 174–183, ACM, 2014.
- [30] M. Wang, C. Tang, P. Shi, and A. van den Broeck. Using neural networks to detect user feedback. In *Proceedings of the 20th International ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 191–200, 2014.
- [31] A. Mahajan and K. Bharadwaj. Probabilistic matrix factorization. In *Advances in neural information processing systems*, pages 1275–1284, 2009.
- [32] T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- [33] A. Nagandhi and V. Tanaka. Adaptive video ranking and full video search for three applications. *Journal of Information Processing*, 20(2):109–122, 2012.
- [34] J. Peng, Y. Zhou, and Q. Tang. Learning neural network users and modeling recommendation. In *2017 20th International Conference on Modeling, Identification and Control (ICMIC)*, pages 1–6, IEEE, 2017.
- [35] Z. Ren, S. Liang, H. J. Wang, and M. de Rube. Social collaborative viewpoint prediction with explicit recommendations. *China Communications*, page 30, 2016.
- [36] S. Rendle, C. Freudenthaler, C. Gantner, and L. Schmidt-Thieme. Ray: Bayesian personalized ranking from implicit feedback. In *Proceedings of the 24th Conference on Artificial Intelligence*, pages 411–414, ACM Press, 2009.
- [37] J. Rendle, C. Freudenthaler, and Y. Yang. Sequence recommender learning with neural networks. In *Advances in neural information processing systems*, pages 1034–1042, 2011.
- [38] Y. Tan, M. Zhang, L. Liu, and Y. Ma. Rating-based latent topic understanding users and items with rating and reviews. In *ECCV*, 2016.
- [39] L. Tang, B. Qin, and L. Lian. Learning matrix representations of users and products for document-level recommendation. In *ACM*, pages 1014–1023, 2012.
- [40] D. Tang, B. Qin, Liu, and Y. Yang. User modeling with neural network for review topic prediction. In *ICML*, pages 1040–1048, 2013.
- [41] D. Vetro, A. Foahey, S. Bergin, and D. Ebrahimi. Show and tell: A neural image caption generator. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3136–3144, 2015.
- [42] H. Wang, S. Wang, and D. S. Wang. Collaborative deep learning for recommendation systems. In *AAAI*, pages 1191–1194, 2013.
- [43] H. Wang, S. Wang, and D. S. Wang. Collaborative recurrent autoencoder for recommendation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 413–419, 2014.
- [44] M. Wang, H. Guo, C. Li, Z. He, Y. Yao, and T. S. Chua. Cross-domain user recommendation. In *Proceedings of the 2013 IEEE International Conference on Data Mining*, pages 1049–1058, 2013.
- [45] F. Wang, M. Kuan, J. Luo, and B. L. Tseng. Deep neural-based graph representation for automatic video summarization. In *Proceedings of the 2017 AAAI Conference on Artificial Intelligence and Empowering people*, pages 7–9, 2017.
- [46] B. Wu, F. Zhang, H. Tan, A. He, and Q. Yang. Cross-modal time-wise video tagging using temporal and perceptual topic modeling. In *SIGIR*, pages 770–784, 2014.
- [47] C. A. Borden, A. Borden, A. Ahmad, and J. J. Smalls. Explaining reviews and ratings with topic-driven auto-encoding. In *Proceedings of the 20th International Conference Companion on World Wide Web*, pages 127–138, International World Wide Web Conference Steering Committee, 2014.
- [48] W. Wu and M. Ester. Hara: A probabilistic model combining aspect-based opinion mining and collaborative filtering. In *WWW*, pages 1089–1100, ACM, 2014.
- [49] Y. Xia, J. Li, C. Zhang, and L. Lian. Video highlight shot extraction with time-view coattention. In *Workshop on the Paper on Pattern and Multi-media Computing and Video Search*, pages 31–36, 2013.
- [50] H. Xu, Y. Zhou, and H. Zhu. Trade-off between a great general-based video summarization model. In *ICCV*, pages 2198–2204, 2015.
- [51] Y. Yan, G. Li, L. Fan, and L. Li. A multiple visual models based perspective analysis framework for multi-view video summarization. In *IEEE International Conference on Multimedia and Expo*, 2007.
- [52] J. You, M. H. Heisterkamp, S. Vijayarajasekaran, O. Vetro, R. Hariri, and G. Gendreau. Combining user and implicit feedback. Deep networks for video classification. *arXiv preprint arXiv:1406.1817*, 2014.
- [53] W. Zhang, C. Tang, P. Shi, and A. van den Broeck. Collaborative multi-level context-based learning for video rating prediction. In *ICCV*, 2013.
- [54] Z. Zhang, C. Li, L. Lian, and Y. Zhang. A deep neural network model for the applicable recommendation based on phrase-level sentiment analysis. In *Proceedings of the 19th International Conference on Knowledge Discovery and Data Mining in Information Systems*, pages 61–72, ACM, 2013.
- [55] Y. Zhang, Y. Ren, T. S. Chua, and H. J. Zhang. Adaptive feature extraction using supervised clustering. In *Image Processing, 2008. ICCP 2008. Proceedings. 2008 International Conference on*, volume 4, pages 17–20, IEEE, 2008.

# 遍历所有参考文献，然后从第一篇文章开始，遍历参考文献的参考文献

- 第一种遍历策略：广度优先算法



- **第一种遍历策略：广度优先算法**

```
PROCEDURE SPIDER(G, {SEEDS})
```

```
  Initialize COLLECTION <big file of URL-page pairs> // 结果存储
```

```
  Initialize VISITED <big hash-table> // 已访问URL 列表
```

```
For every ROOT in SEEDS
```

```
  Initialize QUEUE <queue data structure> // 待爬取URL 队列
```

```
  Let QUEUE := EnQueue(ROOT, QUEUE)
```

```
  While QUEUE is not empty,
```

```
    Do URLcurr := DeQueue(QUEUE)
```

```
      Until URLcurr is not in VISITED
```

```
    insert-hash(URLcurr , VISITED)
```

```
    PAGE := look-up(URLcurr ) // 爬取页面
```

```
    STORE(<URLcurr , PAGE>, COLLECTION)
```

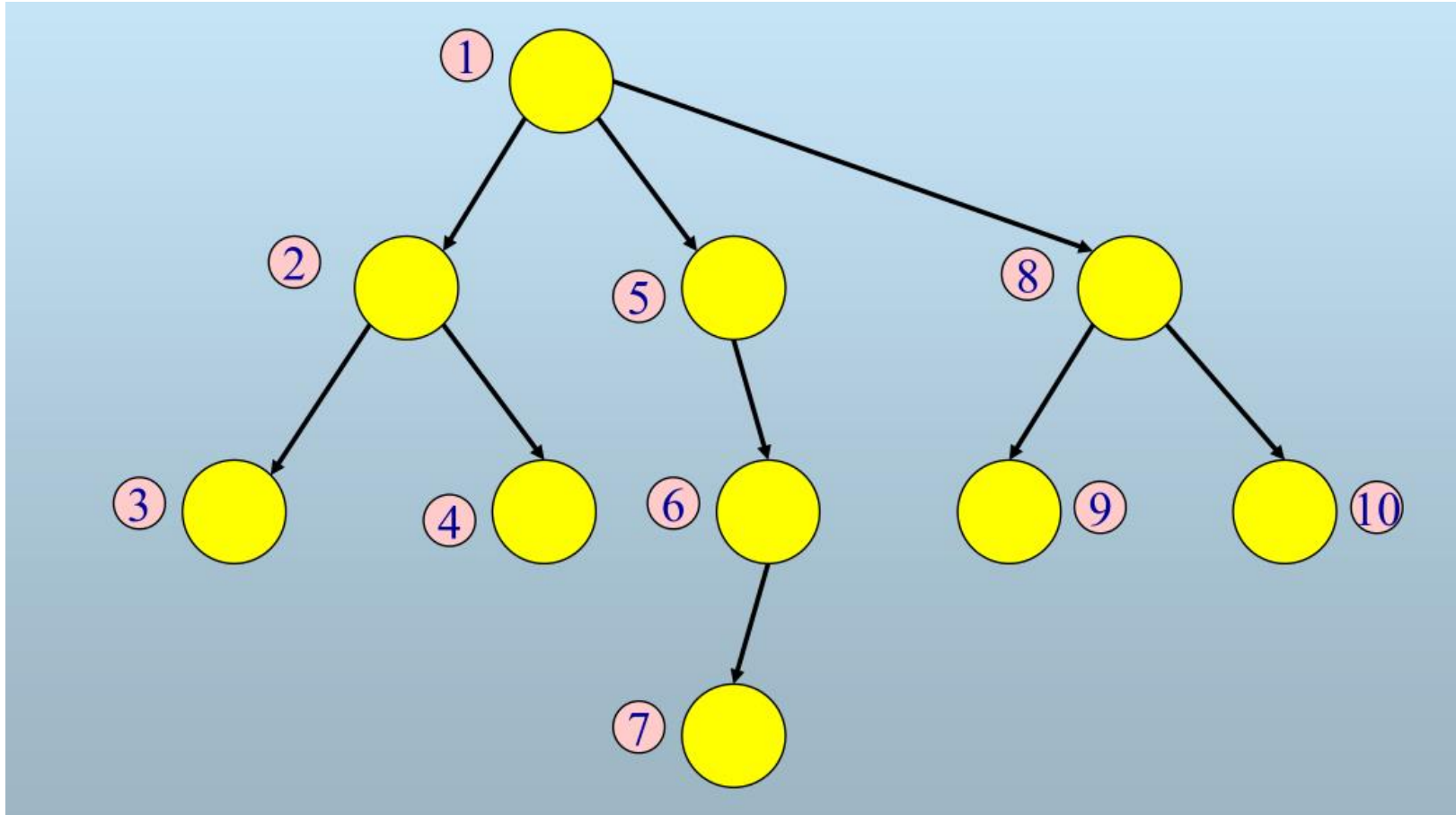
```
    For every URL i in PAGE, // 链接提取
```

```
      EnQueue(URL, QUEUE)
```

```
Return COLLECTION
```

本质是维持一个URL队列，  
先进先出

- 第二种遍历策略：深度优先算法



- 第二种遍历策略：深度优先算法

```
PROCEDURE SPIDER(G, {SEEDS})
```

```
    Initialize COLLECTION <big file of URL-page pairs> // 结果存储
```

```
    Initialize VISITED <big hash-table> // 已访问URL 列表
```

```
For every ROOT in SEEDS
```

```
    Initialize STACK <stack data structure> // 待爬取URL 栈
```

```
    Let STACK := push(ROOT, STACK)
```

```
    While STACK is not empty,
```

```
        Do URLcurr := pop(STACK)
```

```
            Until URLcurr is not in VISITED
```

```
                insert-hash(URLcurr , VISITED)
```

```
                PAGE := look-up(URLcurr ) // 爬取页面
```

```
                STORE(<URLcurr , PAGE>, COLLECTION)
```

```
                For every URLi in PAGE, // 链接提取
```

```
                    push(URLi , STACK)
```

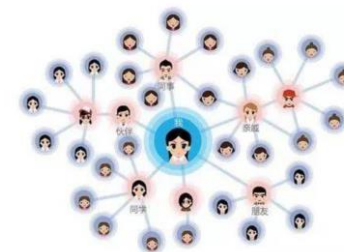
```
Return COLLECTION
```

本质是维持一个URL栈，  
先进后出



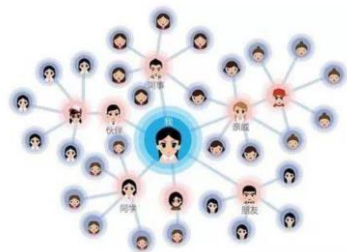
- 广度 or 深度优先算法?

- 对于一般的网页爬虫而言，在时间无限的情况下，两种算法可以视作等价
- 假如在**有限步**条件下，两者的优劣如何？
  - 如果希望获得更为多样化的内容，应采用广度优先
  - 如果希望获得更为深层次的信息，应结合路径选择采用深度优先
    - 部分站点需要深层次浏览以判断站点信息质量



- **广度 or 深度优先算法?**

- 对于一般的网页爬虫而言，在时间无限的情况下，两种算法可以视作等价
- 此外，**不同的应用场景**，对两种算法的适应度可能不尽相同
  - 以社交网络场景为例，广度优先有着以下特点：
    - ✓ 通常而言，选择度数较高的起点，可以快速到达大量节点
    - 但全局性不佳，难以反映网络全貌
    - 另外，此类采样所得节点的度数往往虚高
    - ◆ Facebook网络，随机抽样平均度94，广度优先抽样竟达324



- **再进一步，从遍历要求到网页本身重要性的要求**
- 需要对网页的重要性进行评估，进而优先搜集重要的网页
- 根据经验，体现网页重要度的常见特征类别：
  - 启发式特征（简单统计指标）
    - 如，入度、父网页入度、镜像度、较小的目录深度（易于浏览）
  - 结构性特征
    - 如，PageRank/HIT值，各类Betweenness值等
  - 主题类特征（反应网页与特定需求的契合程度）

- 网络爬虫的定义与需求
- 爬虫的基本要素
- 面向API的新爬虫任务
- 常见的爬虫算法
- **常见反爬虫机制与应对策略**

- **有爬虫，就有反爬虫**
- 使用任何技术手段，阻止别人批量获取自己网站信息
  - 顺便，大家也一起来当一下免费的数据标注员



Please show you're not a robot

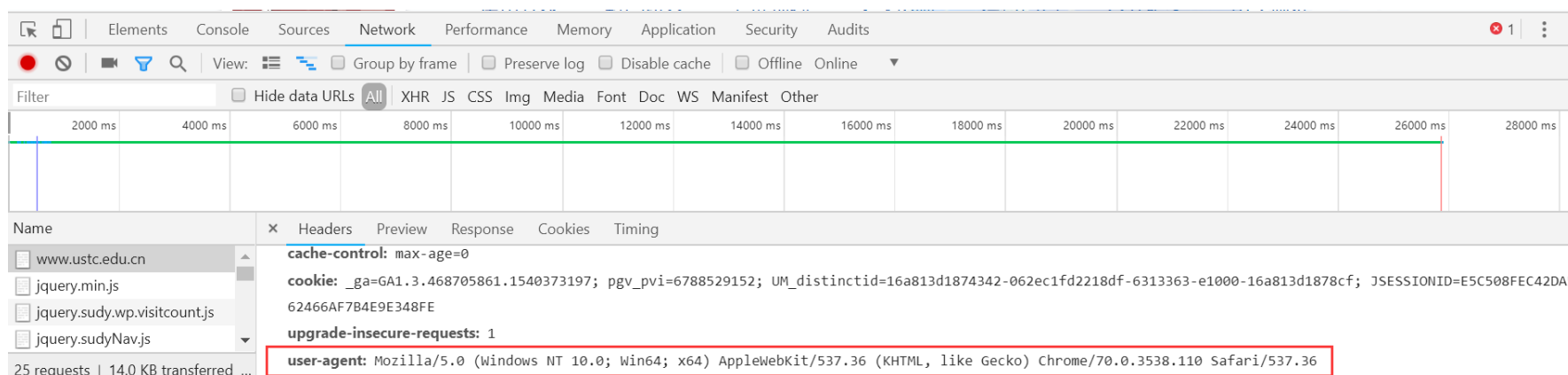
I'm not a robot

 reCAPTCHA  
Policy - Terms

We don't think. We know.

## • 常见反爬虫策略 (1) : User Agent

- 最常见的反爬虫策略，利用访问网站时浏览器发布的Request Headers信息中的User-Agent信息，判断用户使用何种方式浏览



```
UA.py x
1 import requests
2
3 html = requests.get('https://zhihu.com').content
4 print(html.decode())
5
```

爬虫往往U-A  
部分为空

- **常见反爬虫策略 (1) : User Agent**

- 应对策略：利用Python的Request库允许用户自定义请求头信息的手段
  - 在请求头信息中将 User-Agent 的值改为浏览器的请求头标识，从而  
绕开反爬虫机制

```
import requests
# 伪造请求头信息 欺骗服务器
headers = {"User-Agent": "Mozilla/5.0 (Macintosh; Intel Mac OS X 10.13; rv:9527.0) Gecko/20100101 Firefox/9527.0"}
resp = requests.get("http://127.0.0.1", headers=headers)
print(resp.status_code)
```

- **常见反爬虫策略 (2) : IP/账号访问次数/频率**
- 通过限制特定IP地址/账号访问频率和次数进行反爬
  - 其本质在于判断浏览行为是否是**人类行为**
- 应对手段：
  - 构造 IP 代理池，然后每次访问时随机选择代理
    - Github中有相关服务，通过各种提供免费IP的网站来提供代理池
  - 每次爬取行为后间隔一段时间
  - 注册多个账号以保障数据收集
    - 挑战：账号本身的成本问题、账号被查封的危险



- 常见反爬虫策略 (3) : 验证码

- 通过各类验证码, 判断浏览者属于人类还是机器



从简单的字符识别到  
复杂的逻辑推理

验证码:

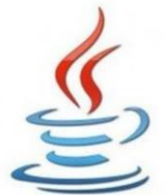


- 应对手段:

- 简单的字符识别: 基于机器学习与模式识别相关技术
- 复杂的逻辑推理: 人工辅助破解

- **常见反爬虫策略 (4) : 动态网页**

- 从网页的 url 加载网页的源代码之后，会在浏览器里执行JavaScript程序
  - 网页内容由脚本加载，而直接抓取则只能得到空白页面
  - 此类情况在抓取在线播放的音视频文件时尤为常见
- 应对手段：
  - 核心思路：模拟调用请求
  - 使用审查元素分析ajax请求，如此循环直到获得包含数据信息的json文件



## • 常见反爬虫策略 (5) : 蜜罐技术

- 网页上会故意留下一些人类看不到或者绝对不会点击的链接。由于爬虫会从源代码中获取内容，所以爬虫可能会访问这样的链接
- 只要网站发现有IP访问这个链接，立刻永久封禁访问者，从而难以继续爬取
- 应对手段：
  - 核心思路：干涉爬虫路径
  - 通过工具库判断页面上的隐含元素，使爬虫避开这些元素，可以部分回避蜜罐的诱导。



## • 常见反爬虫策略 (6) : 加密-解密技术

- 对网页中的关键信息进行加密或混淆, 来增加
- 一类常见案例: 页面上显示为正常文本, 而源代码中为编号

- 网页加载时需要借助字体库, 将编码转为文字

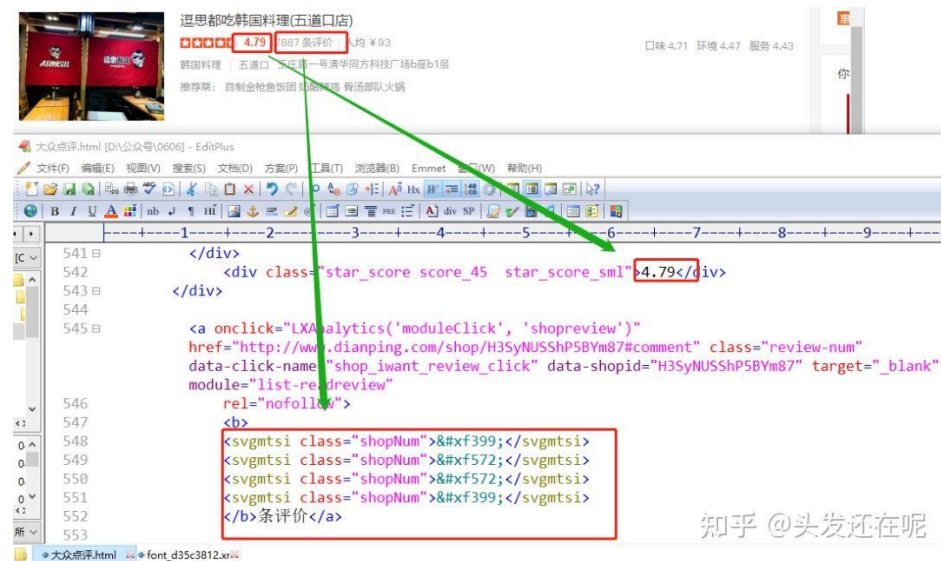
- 字体库/映射可以动态更新!

- 解决方法:

- 通过抓包获取字体库/映射

- 基于映射进行解码

参考链接: <https://zhuanlan.zhihu.com/p/406615464>



- **常见反爬虫策略 (7) : 用户权限限制**
- 不同类型/级别的用户给予不同的内容权限
  - VIP、SVIP、蓝钻、红钻、绿钻、各种钻.....
- 应对手段:
  - 氪金, 就可以变强



- **其他的反爬虫策略**

- 不同的网页结构：每个相同类型的页面的源代码格式均不相同
- 同样也是双刃剑：增加爬取难度的同时降低用户浏览体验（[逼死强迫症](#)）
- 多模态的呈现方式：文字转为图像或视频
- 应对策略：OCR、语音识别、图像/视频标签技术



## 年度锦鲤！太原一彩民中奖1.34亿,扮"猪猪侠"领奖当场捐300万



# 本章小结

## 网络爬虫

- 网络爬虫的定义与需求
- 爬虫的基本要素
- 面向API的爬虫任务
- 常见的爬虫算法
- 常见反爬虫机制与应对策略