

北 京 林 业 大 学

2022 学年—2023 学年第 2 学期 Python 应用 实验任务书

专业名称：计算机、大数据、物联网、信息 实验学时： 2

课程名称：Python 应用 任课教师： 王春玲

实验题目：实验 3 爬取中国工程院院士信息

实验环境：Python、PyCharm 等

实验目的：

1. 熟练使用标准库 urllib 读取网页内容。
2. 熟练使用正则表达式提取文本中感兴趣的信息。
3. 熟练使用内置函数 open() 创建文本文件和二进制文件。
4. 熟悉 HTML 语法以及常见的 HTML 标签。

实验内容：

爬取中国工程院网页，把每位院士的简介保存为本地文本文件，并把每位院士的照片保存为本地图片，文本文件和图片文件都以院士的姓名为主文件名。

实验步骤如下：

(1) 使用 Google Chrome 或其他浏览器打开下面的网址，然后在页面上右击，在弹出的菜单中选择“查看网页源代码”。

http://www.cae.cn/cae/html/main/col48/column_48_1.html

(2) 分析网页源代码，确定每位院士的姓名和链接所在的 HTML 标签，为后面编写正则表达式做准备，如图 1 所示。

```
<li class="name_list"><a href="/cae/html/main/colys/63775817.html" target="_blank">曹喜滨</a></li>
<li class="name_list"><a href="/cae/html/main/colys/35791989.html" target="_blank">陈学东</a></li>
<li class="name_list"><a href="/cae/html/main/colys/01567139.html" target="_blank">邓宗全</a></li>
<li class="name_list"><a href="/cae/html/main/colys/25235806.html" target="_blank">丁荣军</a></li>
<li class="name_list"><a href="/cae/html/main/colys/46604755.html" target="_blank">董春鹏</a></li>
<li class="name_list"><a href="/cae/html/main/colys/42793697.html" target="_blank">樊会涛</a></li>
<li class="name_list"><a href="/cae/html/main/colys/15623030.html" target="_blank">冯培德</a></li>
<li class="name_list"><a href="/cae/html/main/colys/48072917.html" target="_blank">冯煜芳</a></li>
<li class="name_list"><a href="/cae/html/main/colys/71311121.html" target="_blank">甘晓华</a></li>
<li class="name_list"><a href="/cae/html/main/colys/15801057.html" target="_blank">高金吉</a></li>
```

图 1 每位院士的链接

