```
全称XML Path Language,即XML路径语言
                基础
                                 最初是用来搜寻XML文档的,但是它同样适用于HTML文档的搜索
                                     nodename 。 选取此节点的所有子节点
                                   ▶/ ◎ 从当前节点选取直接子节点
                                   // ◎ 从当前节点选取子孙节点
                         重点规则
                规则
                                    · ② 选取当前节点
                                    .. ◎ 选取当前节点的父节点
                                   @ 。 选取属性
                         导包 。 import lxml.etree
                                                    ★html = lxml.etree.HTML(text) #text是html_str
                         读取数据,构造element对象 html的类型 <type 'lxml.etree_Element'>
                                               ►lxml.etree.tostring(html)
                         element对象转化成str 但是结果是bytes类型
                                           lxml.etree.parse(html_file_path,解析器)
                         从文件中读取html 解析器: etree.HTMLParser()
                                      html.xpath(规则) #html是上面的element对象
                                      //* ◎ 提取所有节点
                                      [<Element li at 0x1359f248>, <Element li at 0x1359f208>, <Element li at 0x1359c208>, <Element li at 0x1359c208>, <Element li at 0x1359ac8>]里面的都是element元素,还可//li 。 所有li节点 。 以继续调用xpath方法进行提取元素
                                                   取出li的所有结点的class的值,class是属于li的
                                      //li/@class o/
                                                  若其中的属性不存在,则跳过不显示
                                      ★//li/a/@href ○ 最常用的提取链接
                                      //li/a/@href=\"link3.html\" 。 提取li标签下的a标签中href属性为link3.html的元素
                                      取出li的所有结点的class的值,class是不一定属于li的,
//li//@class © 也可能是li下其他标签的
                                               //li[1] ◎ 所有li标签中的第一个
                         ★规则提取
                                               //li[last()] ◎ 所有li标签中的最后一个
                                              //li[last()-1] 。 所有li标签中的倒数第二个
                                      //li[last()-1]/a/@href 。 所有li中倒数第二个下的a标签的href属性
                                      //*[@href=\"link4.html\"] ◎ 所有标签中href属性为link4.html
                                                                            所有标签中href属性为link4.html的标签之间的文本内
                                              //*[@href=\"link4.html\"]/text() 。 容
                                            ★//div[@class=\"haha\"]//p//text() ◎ 所有div标签中class属性为haha的所有p标签之间的文本内容
                                                                        contains的使用
                                      ★//li(contains(@class,"li")]/a/text() 
第一个參数传入属性名称,第二个参数传入属性值,只
要此属性包含所传入的属性值,就可以完成匹配了
                                      ★//li[contains(@class, "li") and @name="item"]/a/text() ◎ 一个标签多个属性值
Xapth
                                      position函数的使用
//li[position()<3]/a/text()
位置小于3的li元素
               代码
                                                 //li[1]/ancestor::*
                                      ancestor //ii[1]/ancestor::div
                                      attribute o //li[1]/attribute::*
                                      child o //li[1]/child::a[@href="link1.html"]
                                      descendant o //li[1]/descendant::span
                         节点轴选择
                                      following o //li[1]/following::*[2]
                                      following-sibling o //li[1]/following-sibling::*
                                                 http://www.w3school.com.cn/xpath/
                                      更多用法 o xpath_axes.asp
                                           or ⊚ 或
                                           and ∘ 与
                                           mod 🌼 计算余数
                                           ○ 计算两个节点集
                                           + • 加法
                                           - 🌼 减法
                                           * 🌼 乘法
                         XPath中的运算符
                                           div 🌼 除法
                                            = 🌼 等于
                                           != ○ 不等于
                                           < ⊙ 小于
                                           > ◎ 大于
                                           <= ◎ 小于或等于
                                           >= ◎ 大于或等于
                                                          def download(url):
                                                          der download(ur);
headers = {"User-Agent": "Mozilla/5.0 (compatible; MSIE 9.0; Windows NT 6.1; Trident/5.0);"}
request = urllib2.Request(url,headers = headers)
response = urllib2.urlopen(request)
                                                           data = response.read()
                                                           html = lxml.etree.HTML(data)
e_html = html.xpath("//*[@class=\"emphasis\"]/text()")
                                                           print(e html)
                                                          download("https://www.autohome.com.cn/166/#pvareaid=3311284")
                         结合urllib2解析数据 ⊙ 标准代码 ⊙
                                                        response.read()可以直接放入到lxml.etree.HTML()中
                         etree模块可以自动修正HTML文本
                                  当我们使用这种方式进行数据提取并没有提取到数据的时候,可能是a并不是ul的直接子节点,而是子孙节点,使用//ul//a 进行提取
                         //ul/a o
                                  典型犯错的: //li[@class="item-0"]/text() 。 text在a中不在li中, /提取不到
                补充
                                        //li[@class="item-0"]/a/text()
                                       //li[@class="item-0"]//text()
                                               第一个只会将a标签下面的文本数据挖掘出来
                         ▶★经典问题
                                       区別 

第二个会将a下的所有文本数据挖掘出来,并且li下不属于a标签的其他文本也都会被挖掘出来
```

联系 ○ 第一个挖掘的文本信息是第二个的子集

安装