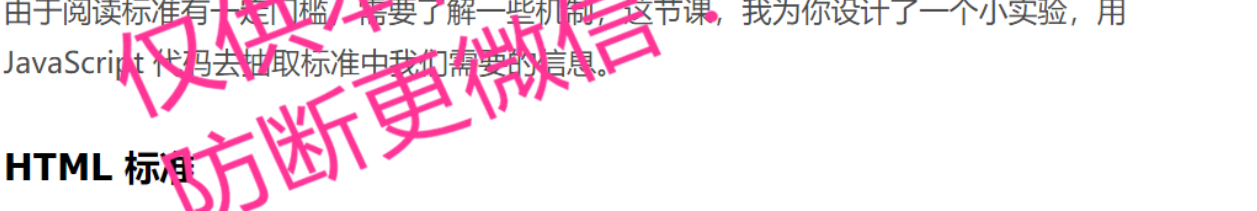


34 | HTML小实验：用代码分析HTML标准

winter 2019-04-11



你好，我是 winter。

前面的课程中，我们已经讲解了大部分的 HTML 标签。

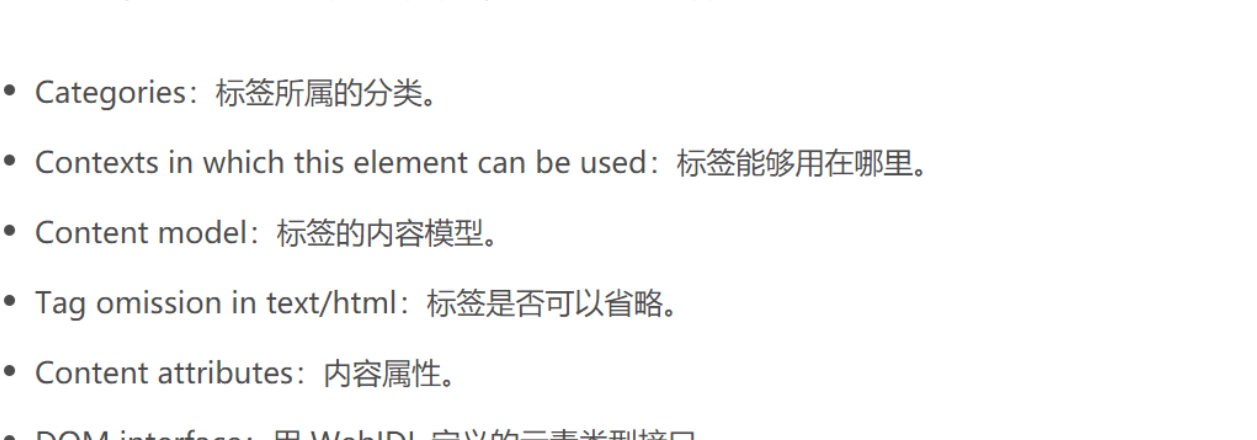
然而，为了突出重点，我们还是忽略一些标签类型。比如表单类标签和表格类标签，我认为只有少数前端工程师用过，比如我在整个手机淘宝的工作生涯中，一次表格类标签都没有用到，表单类则只用过 input，也只有几次。

那么，剩下的标签我们怎么样去了解它们呢？当然是查阅 HTML 标准。

由于阅读标准有门槛，需要了解一些机制，这节课，我为你设计了一个小实验，用 JavaScript 代码去抽取标准中我们需要的信息。

HTML 标准

我们采用 WHATWG 的 living standard 标准，我们先来看看标准是如何描述一个标签的，这里我们看到，有下面这些内容。



我们看到，这里的描述分为 6 个部分，有下面这些内容。

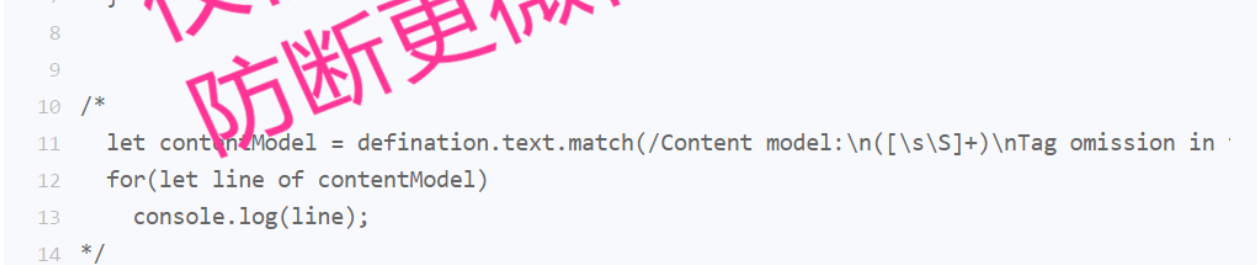
- Categories: 标签所属的分类。
- Contexts in which this element can be used: 标签能够用在哪里。
- Content model: 标签的内容模型。
- Tag omission in text/html: 标签是否可以省略。
- Content attributes: 内容属性。
- DOM interface: 用 WebIDL 定义的元素类型接口。

这一节课，我们关注一下 Categories、Contexts in which this element can be used、Content model 这几个部分。我会带你从标准中抓取数据，做一个小工具，用来检查 X 标签是否能放入 Y 标签内。

代码角度分析 HTML 标准

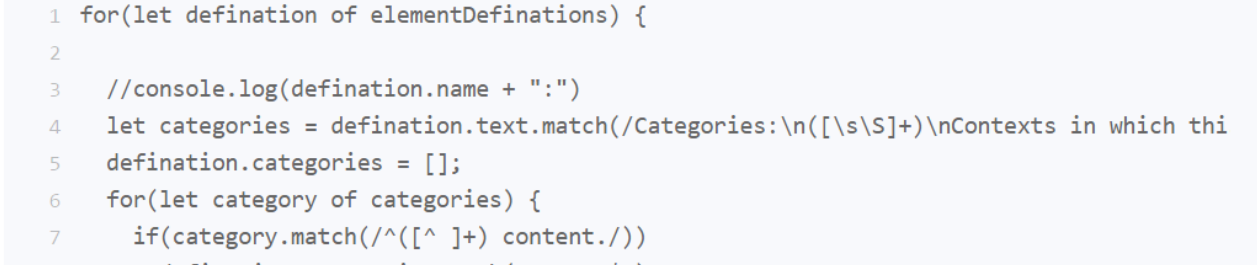
HTML 标准描述用词非常的严谨，这给我们抓取数据带来了巨大的方便，首先，我们打开单页面版 HTML 标准 <https://html.spec.whatwg.org/>

在这个页面上，我们执行一下以下代码：



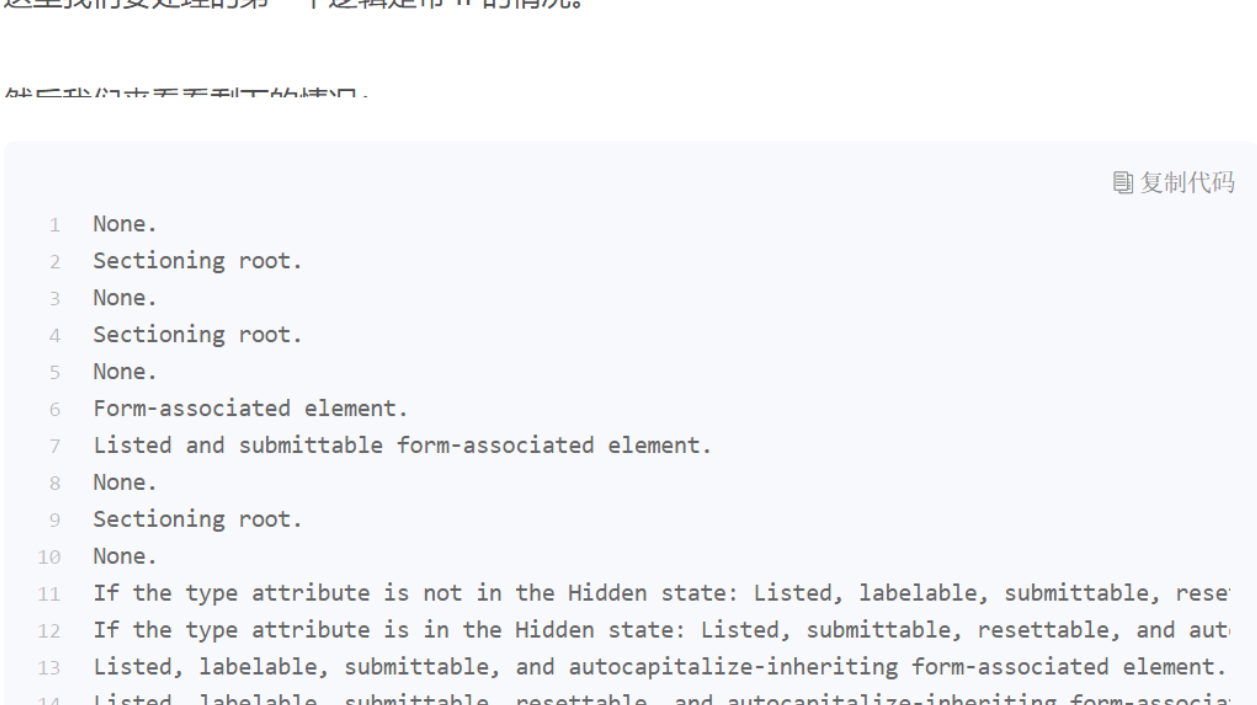
这样我们就得到了所有元素的定义了，现在有 107 个元素。

不过，比较尴尬的是，这些文本中并不包含元素名，我们只好从 id 属性中获取，最后代码类似这样：



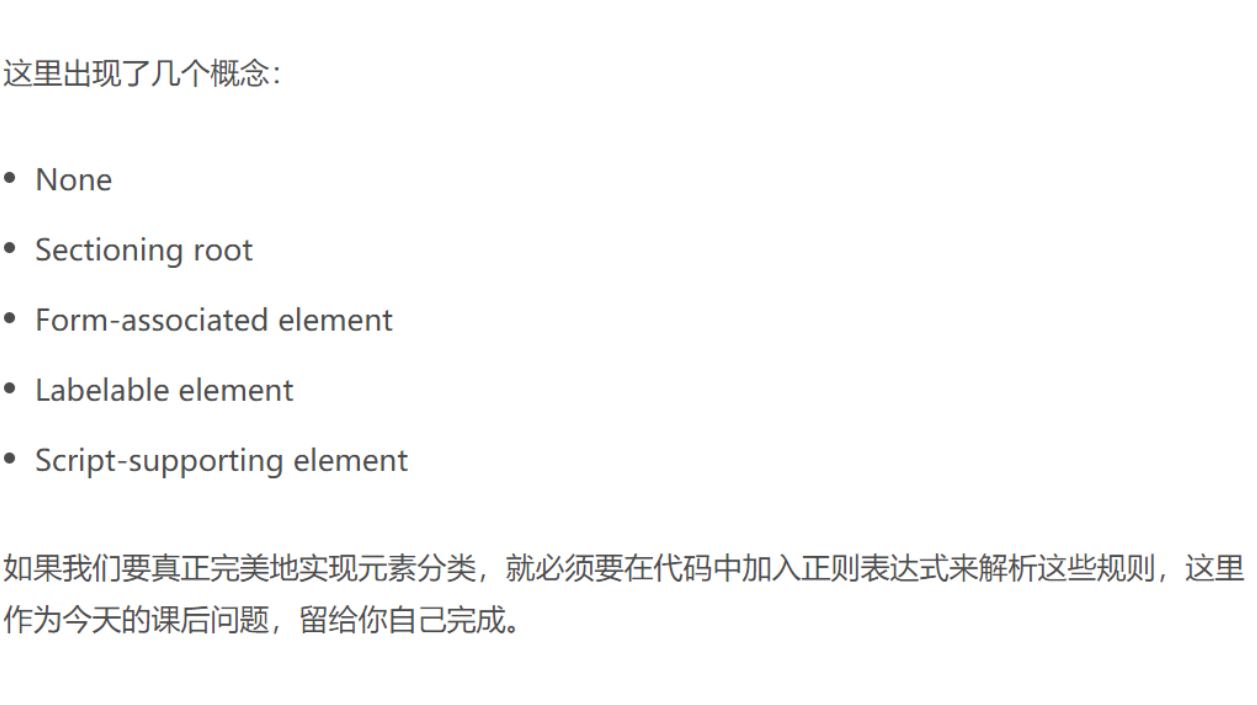
接下来我们用代码理解一下这些文本。首先我们来分析一下这些文本，它分成了 6 个部分，而且顺序非常固定，这样，我们可以用 JavaScript 的正则表达式匹配来拆分六个字段。

我们这个小实验的目标是计算元素之间的包含关系，因此，我们先关心一下 categories 和 contentModel 两个字段。

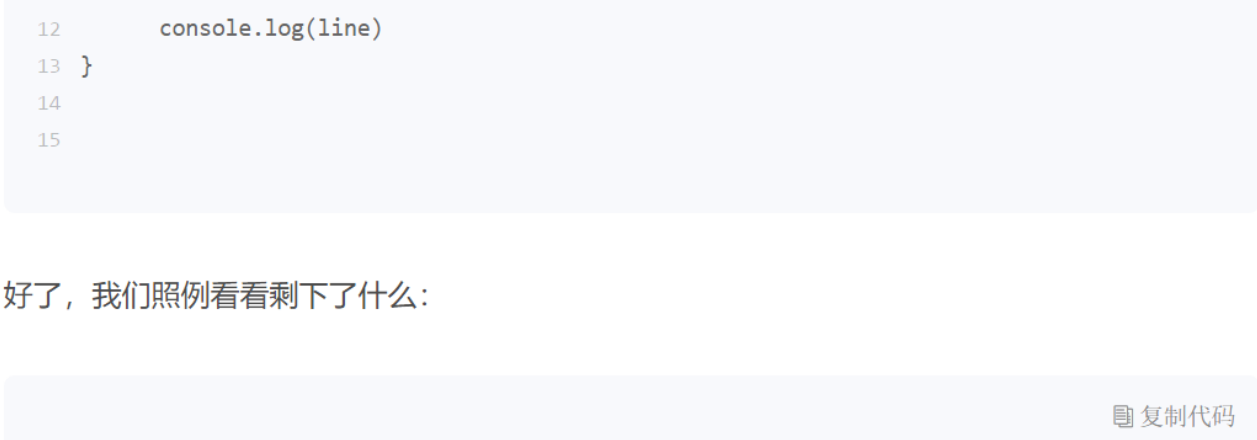
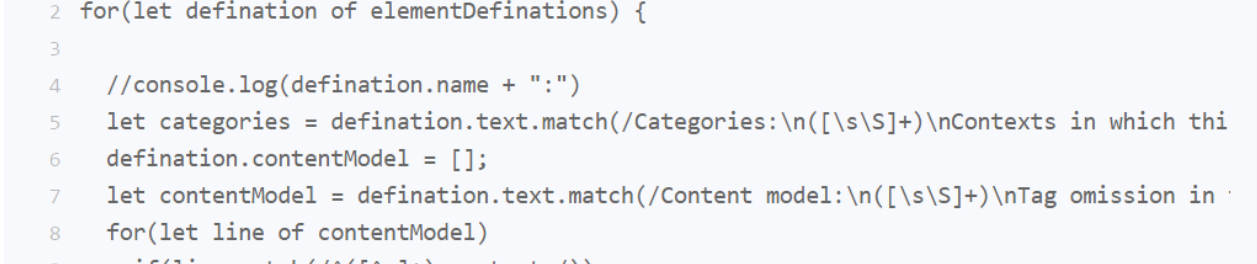


接下来我们来处理 category。

首先 category 的写法中，最基本的就是直接描述了 category 的句子，我们把这些不带任何条件的 category 先保存起来，然后打印出来其它的描述看看：



这里我们要处理的第一个逻辑是带 if 的情况。

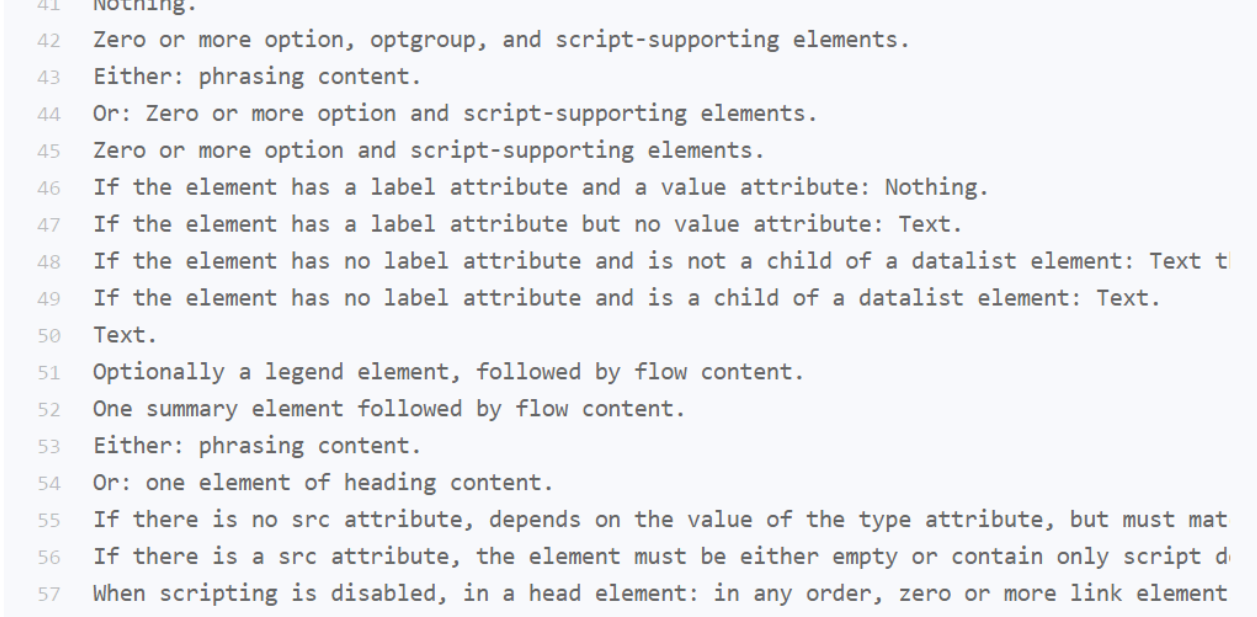


这里出现了几个概念：

- None
- Sectioning root
- Form-associated element
- Labelable element
- Script-supporting element

如果我们要真正完美地实现元素分类，就必须要在代码中加入正则表达式来解析这些规则，这里作为今天的课后问题，留给你自己完成。

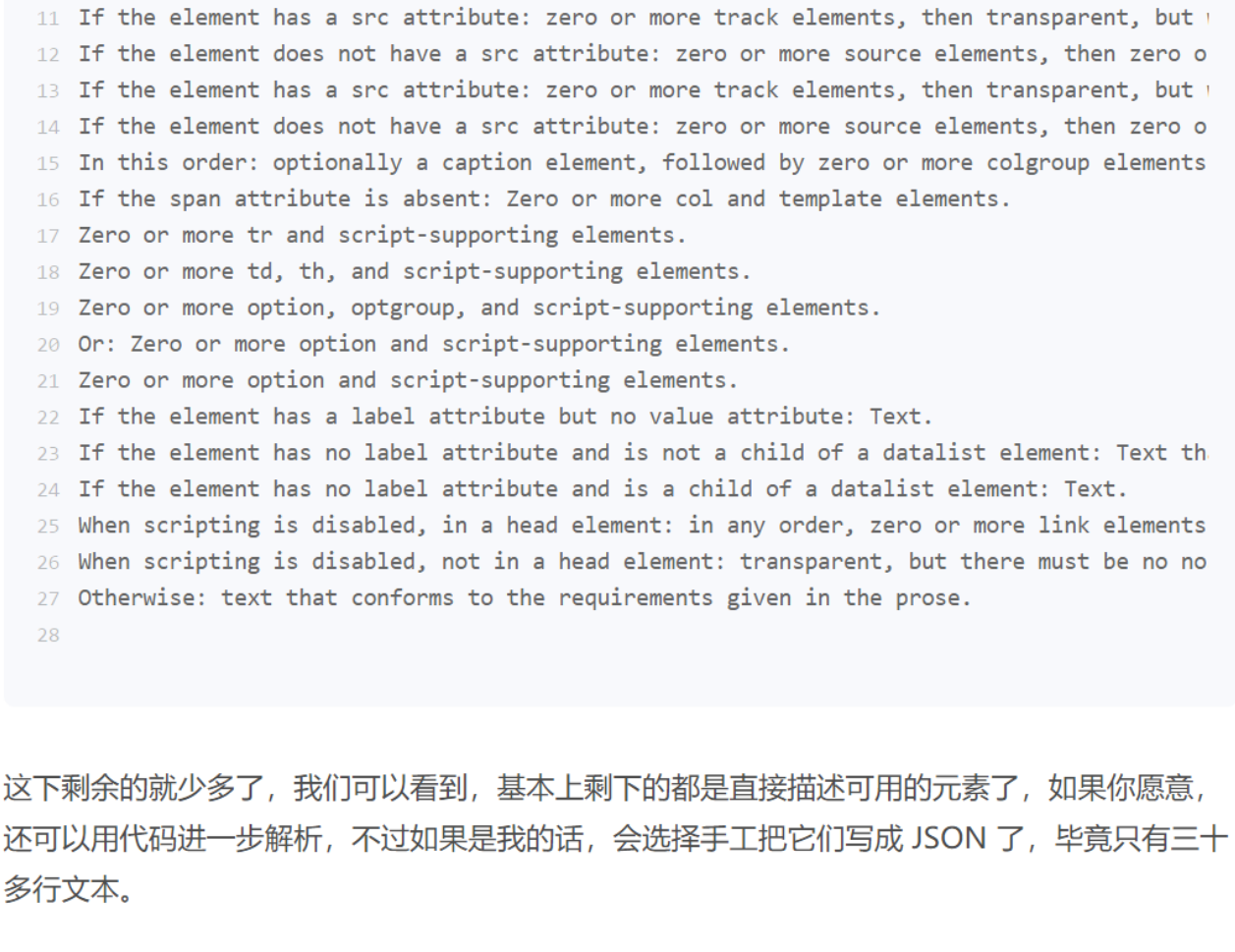
接下来我们看看 Content Model，我们照例先处理掉最简单点的部分，就是带分类的内容模型：



好了，我们照例看看剩下有什么：



这有点复杂，我们还是把它做一些分类，首先我们过滤掉带 If 的情况、Text 和 Transparent。

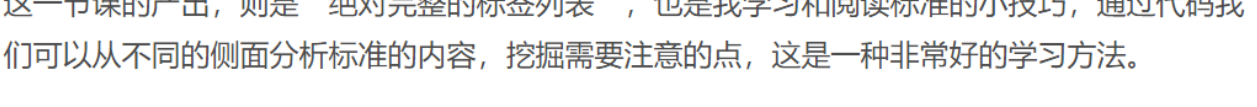


这时候我们再来执行看看：



这下剩余的也就少了，我们可以看到，基本上剩下的都是直接描述可用的元素了，如果你愿意，还可以用代码进一步解析，不过如果是我的话，会选择手工把它们写成 JSON 了，毕竟只有三十多行文本。

好了，有了 contentModel 和 category，我们要检查某一元素是否可以作为另一元素的子元素，就可以判断一下两边是否匹配啦，首先，我们要做个索引：



然后我们编写一下我们的 check 函数：

总结

这一节课，我们完成了一个小实验：利用工具分析 Web 标准文本，来获得元素的信息。

通过这个实验，我希望能够传递一种思路，代码能够帮助我们从 Web 标准中挖掘出来很多想要的信息，编写代码的过程，也是更深入理解标准的契机。

我们前面的课程中把元素分成了几类来讲解，但是这些分类只能大概地覆盖所有的标签，我设置课程的目标也是讲解标签背后的知识，而非每一种标签的细节。具体每一种标签的属性和细节，可以留给大家自己去整理。

这一节课的产出，则是“绝对完整的标签列表”，也是我学习和阅读标准的小技巧，通过代码我们可以从不同的侧面分析标准的内容，挖掘需要注意的点，这是一种非常好的学习方法。