

大数据环境下信息抽取模板自动聚类与发现

计 92 丘骏鹏 2009011282

指导老师：朱小燕 郝宇

研究内容介绍

系统框架与模块实现

- 整体框架

- 预处理模块

- 网页聚类模块

- 模板生成和内容提取模块

实验和分析

总结和展望

研究内容介绍

系统框架与模块实现

- 整体框架

- 预处理模块

- 网页聚类模块

- 模板生成和内容提取模块

实验和分析

总结和展望

背景

- ▶ “大数据”时代来临，我们获得的数据越来越多，研究工作也受到了新的挑战
- ▶ 为了便于计算机处理这些数据，需要从非结构化和半结构化数据中提取出我们关心的内容，存储成结构化的数据
- ▶ 大部分网页是通过查询后台数据库，然后选择合适模板进行渲染方式生成的。即：

背景

- ▶ “大数据”时代来临，我们获得的数据越来越多，研究工作也受到了新的挑战
- ▶ 为了便于计算机处理这些数据，需要从非结构化和半结构化数据中提取出我们关心的内容，存储成结构化的数据
- ▶ 大部分网页是通过查询后台数据库，然后选择合适模板进行渲染方式生成的。即：

```
<html>
<body>
  <h1>{{ news.title }}</h1>
  <p>{{ news.content }}</p>
  <div>
    {% for comment in news.comments %}
      <li>{{ comment }}</li>
    {% endfor %}
  </div>
</body>
</html>
```

背景

- ▶ “大数据”时代来临，我们获得的数据越来越多，研究工作也受到了新的挑战
- ▶ 为了便于计算机处理这些数据，需要从非结构化和半结构化数据中提取出我们关心的内容，存储成结构化的数据
- ▶ 大部分网页是通过查询后台数据库，然后选择合适模板进行渲染方式生成的。即：

```
<html>
<body>
  <h1>{{ news.title }}</h1>
  <p>{{ news.content }}</p>
  <div>
    {% for comment in news.comments %}
      <li>{{ comment }}</li>
    {% endfor %}
  </div>
</body>
</html>
```



生成

- ▶ “大数据”时代来临，我们获得的数据越来越多，研究工作也受到了新的挑战
- ▶ 为了便于计算机处理这些数据，需要从非结构化和半结构化数据中提取出我们关心的内容，存储成结构化的数据
- ▶ 大部分网页是通过查询后台数据库，然后选择合适模板进行渲染方式生成的。即：

```
<html>
<body>
  <h1>{{ news.title }}</h1>
  <p>{{ news.content }}</p>
  <div>
    {% for comment in news.comments %}
      <li>{{ comment }}</li>
    {% endfor %}
  </div>
</body>
</html>
```

生成



工作内容

1. 已经获得大量新闻、博客等网页数据，这些网页可能由不同的模板生成
2. 我们将网页按不同的模板分开，然后从每个集合中提取出所有可能的模板
3. 根据提取出的模板，抽取新的网页中的有用的信息，存储成 XML 格式。

工作内容

1. 已经获得大量新闻、博客等网页数据，这些网页可能由不同的模板生成
2. 我们将网页按不同的模板分开，然后从每个集合中提取出所有可能的模板
3. 根据提取出的模板，抽取新的网页中的有用的信息，存储成XML 格式。



工作内容

1. 已经获得大量新闻、博客等网页数据，这些网页可能由不同的模板生成
2. 我们将网页按不同的模板分开，然后从每个集合中提取出所有可能的模板
3. 根据提取出的模板，抽取新的网页中的有用的信息，存储成XML 格式。



抽取

工作内容

1. 已经获得大量新闻、博客等网页数据，这些网页可能由不同的模板生成
2. 我们将网页按不同的模板分开，然后从每个集合中提取出所有可能的模板
3. 根据提取出的模板，抽取新的网页中的有用的信息，存储成XML 格式。



```
<documents>
<document>
  <title>
    WWDC:2013 苹果仍未解创新和增长之困
  </title>
  <content>
    今日,备受瞩目的苹果全球开发者...
  </content>
</document>
</documents>
```

提纲

研究内容介绍

系统框架与模块实现

- 整体框架

- 预处理模块

- 网页聚类模块

- 模板生成和内容提取模块

实验和分析

总结和展望

提纲

研究内容介绍

系统框架与模块实现

- 整体框架

- 预处理模块

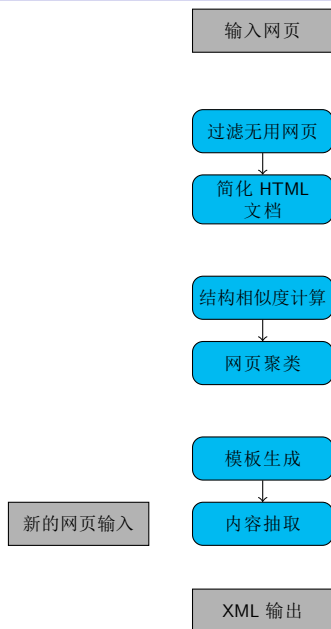
- 网页聚类模块

- 模板生成和内容提取模块

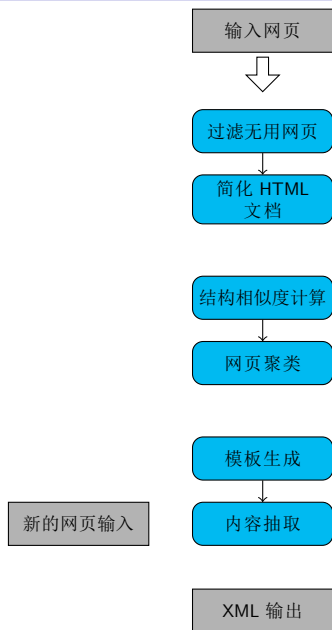
实验和分析

总结和展望

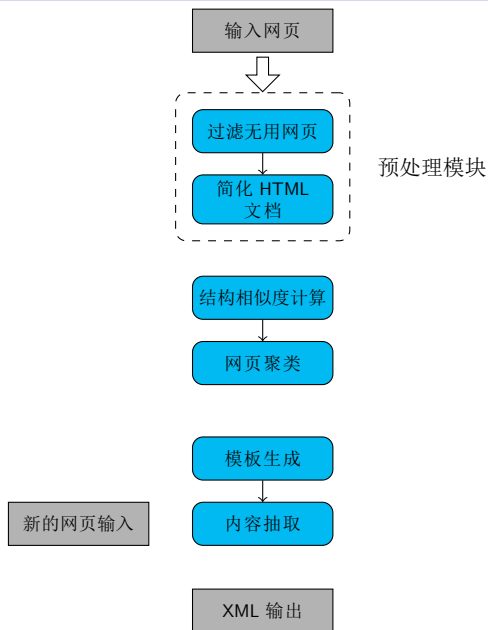
整体框架



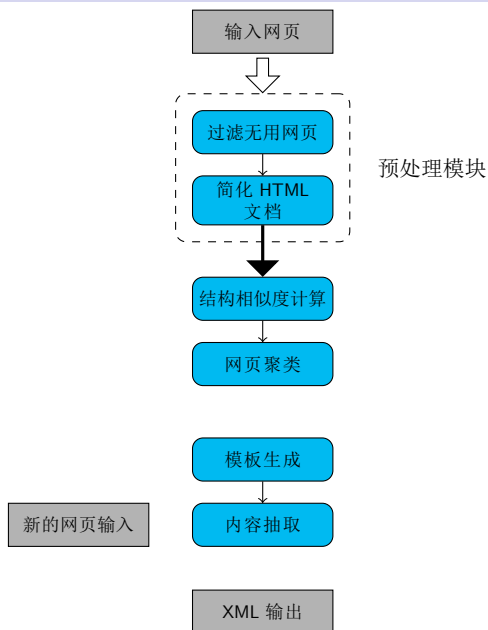
整体框架



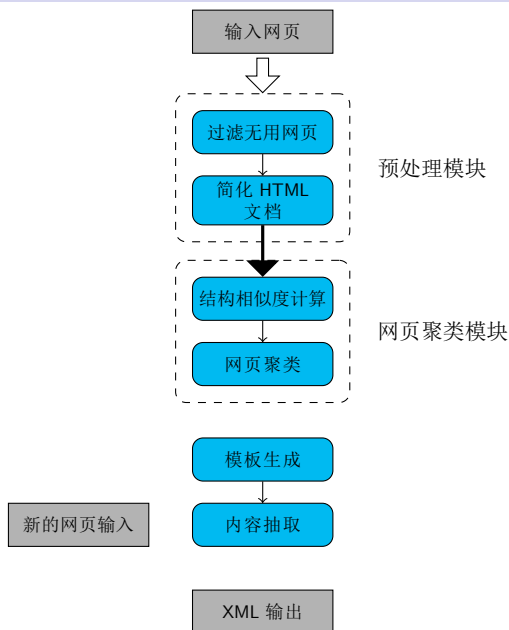
整体框架



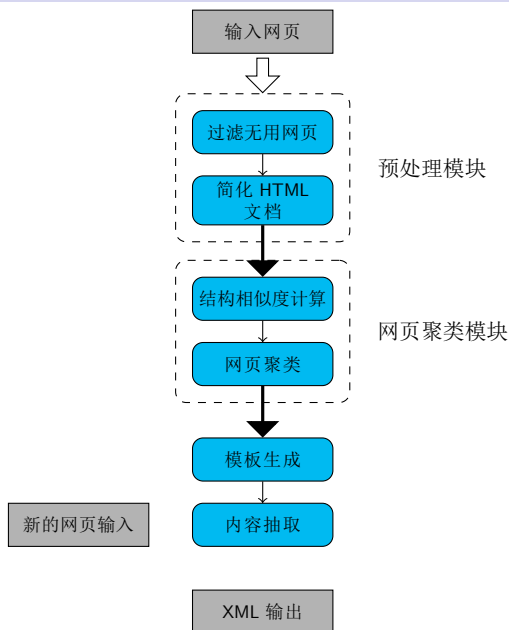
整体框架



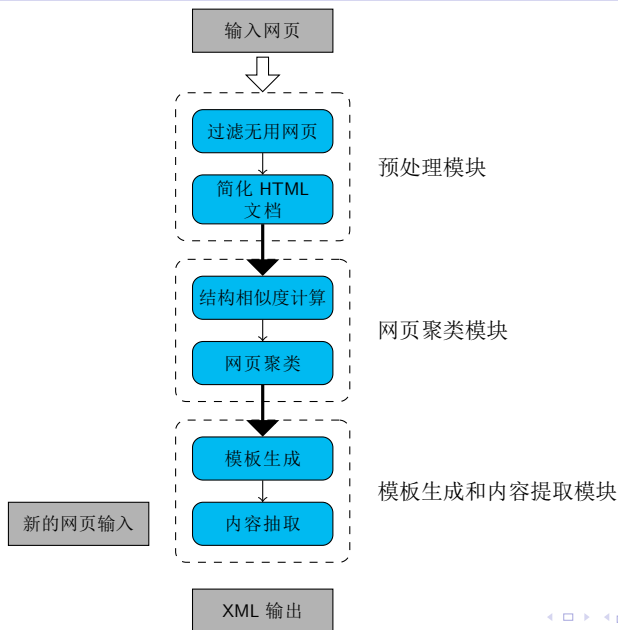
整体框架



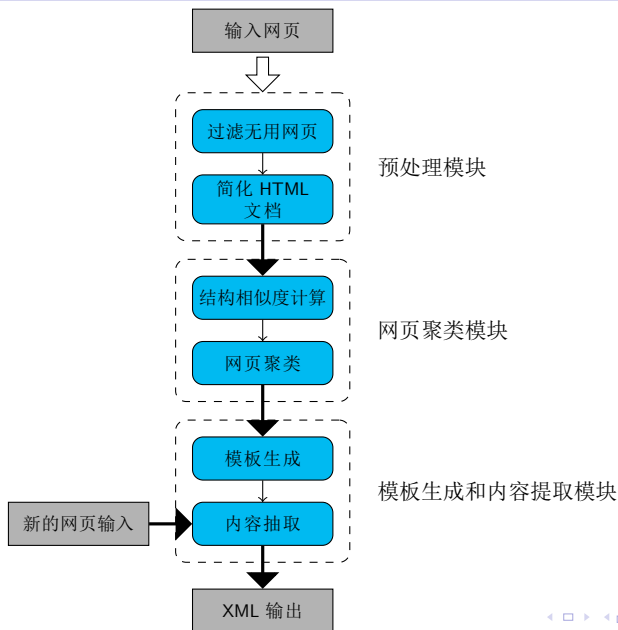
整体框架



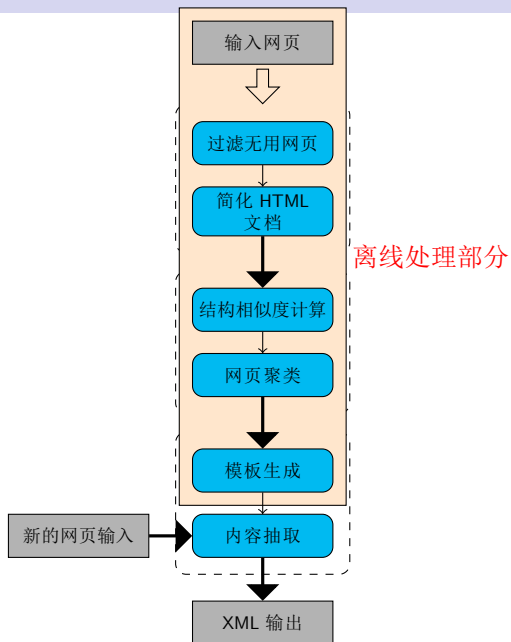
整体框架



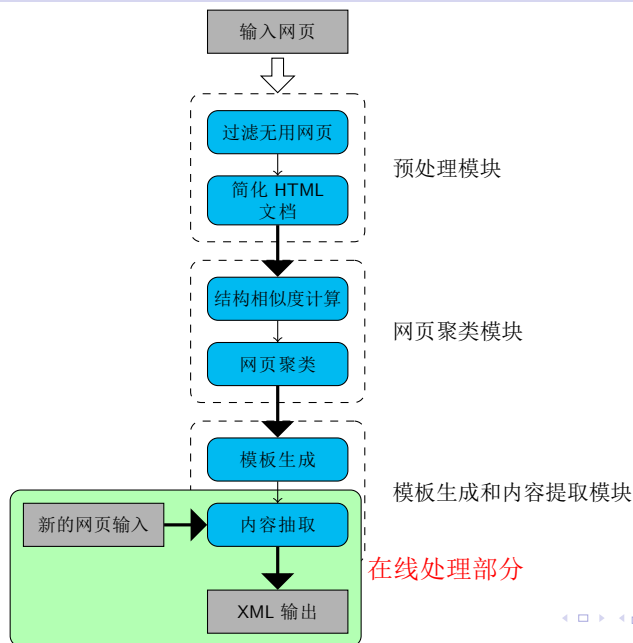
整体框架



整体框架



整体框架



提纲

研究内容介绍

系统框架与模块实现

- 整体框架

- 预处理模块

- 网页聚类模块

- 模板生成和内容提取模块

实验和分析

总结和展望

预处理模块介绍

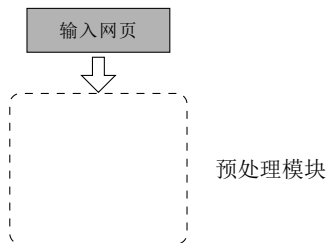
输入网页

预处理模块介绍

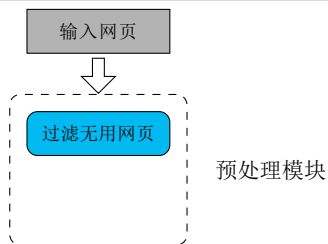
输入网页



预处理模块介绍



预处理模块介绍



预处理模块 -过滤无用网页

错误页过滤

目录页过滤

预处理模块 -过滤无用网页

错误页过滤

- ▶ 404 等错误页大小一般比较小

目录页过滤

预处理模块 - 过滤无用网页

错误页过滤

- ▶ 404 等错误页大小一般比较小
- ▶ 方法：通过设置文件大小阈值过滤

目录页过滤

预处理模块 - 过滤无用网页

错误页过滤

- ▶ 404 等错误页大小一般比较小
- ▶ 方法：通过设置文件大小阈值过滤

目录页过滤

- ▶ 目录页的 URL 有一定的模式：比如对于新浪博客来说，详细页的 URL 都以.html 结尾，而目录页则没有。以下两个 URL：

`http://blog.sina.com.cn/u/1439351555`

`http://blog.sina.com.cn/s/blog_55cac30301016yb1.html`

第一个 URL 对应的是某个博主的文章目录，第二个 URL 则是该博主的某篇文章。

预处理模块 - 过滤无用网页

错误页过滤

- ▶ 404 等错误页大小一般比较小
- ▶ 方法：通过设置文件大小阈值过滤

目录页过滤

- ▶ 目录页的 URL 有一定的模式：比如对于新浪博客来说，详细页的 URL 都以.html 结尾，而目录页则没有。以下两个 URL：

`http://blog.sina.com.cn/u/1439351555`

`http://blog.sina.com.cn/s/blog_55cac30301016yb1.html`

第一个 URL 对应的是某个博主的文章目录，第二个 URL 则是该博主的某篇文章。

- ▶ 方法：通过为每类网页分别建立正则表达式的方法过滤目录页

目录页过滤

- ▶ 为了完备性考虑, 可能存在某些网站, 目录页的 URL 和详细页的 URL 在模式没有太大的区别, 或者人工设置规则的办法太麻烦了

预处理模块 - 过滤无用网页

目录页过滤

- ▶ 为了完备性考虑, 可能存在某些网站, 目录页的 URL 和详细页的 URL 在模式没有太大的区别, 或者人工设置规则的办法太麻烦了
- ▶ 观察百度新闻目录页部分简化的 HTML 代码:

```
<ul>  
  <li><a>红会社监委今日讨论是否重启调查郭美美事件</a></li>  
  <li><a>北京今日14时起堵车 中度以上拥堵将持续7小时</a></li>  
  <li><a>湖北黄冈拒卖地反投8亿建免费公园 市长称不亏</a></li>  
  <li><a>湖南株洲市档案局工会主席杀害局长后跳楼</a></li>  
  <li><a>陕西神木集资大王被查 结婚时曾空运20辆林肯</a></li>  
  <li><a>内地学生扎堆报京沪港高校 在校平均恋爱0.78次</a></li>  
</ul>
```

预处理模块 - 过滤无用网页

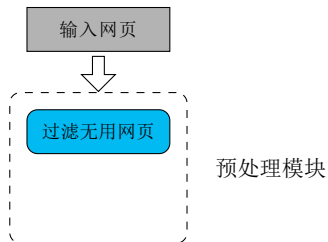
目录页过滤

- ▶ 为了完备性考虑, 可能存在某些网站, 目录页的 URL 和详细页的 URL 在模式没有太大的区别, 或者人工设置规则的办法太麻烦了
- ▶ 观察百度新闻目录页部分简化的 HTML 代码:

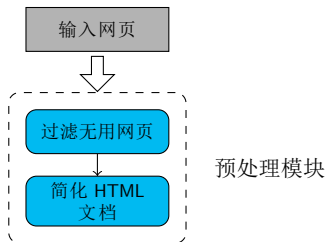
```
<ul>
  <li><a>红会社监委今日讨论是否重启调查郭美美事件</a></li>
  <li><a>北京今日14时起堵车 中度以上拥堵将持续7小时</a></li>
  <li><a>湖北黄冈拒卖地反投8亿建免费公园 市长称不亏</a></li>
  <li><a>湖南株洲市档案局工会主席杀害局长后跳楼</a></li>
  <li><a>陕西神木集资大王被查 结婚时曾空运20辆林肯</a></li>
  <li><a>内地学生扎堆报京沪港高校 在校平均恋爱0.78次</a></li>
</ul>
```

- ▶ 我们提取出网页中的所有 标签中的 <a> 标签, 计算这些 <a> 标签的文本内容占网页所有文本内容的比重, 超过一定的阈值则判定为目录页。

预处理模块介绍



预处理模块介绍



预处理模块 - 简化 HTML 文档

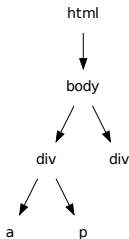
简化 HTML 文档的 3 个步骤

1. 去除无用标签。<script>、<link>、<style>、
 等。

预处理模块 - 简化 HTML 文档

简化 HTML 文档的 3 个步骤

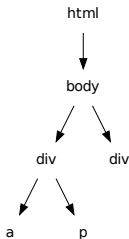
1. 去除无用标签。<script>、<link>、<style>、
 等。
2. 简化解析好的 DOM Tree 树形结构: 采用先序遍历的方式, 将 DOM Tree 转化成一个标签序列。



预处理模块 - 简化 HTML 文档

简化 HTML 文档的 3 个步骤

1. 去除无用标签。<script>、<link>、<style>、
 等。
2. 简化解析好的 DOM Tree 树形结构: 采用先序遍历的方式, 将 DOM Tree 转化成一个标签序列。

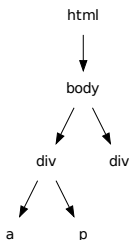


```
<html0><body1>  
<div2><a3><p3>  
<div2>
```

预处理模块 - 简化 HTML 文档

简化 HTML 文档的 3 个步骤

1. 去除无用标签。<script>、<link>、<style>、
 等。
2. 简化解析好的 DOM Tree 树形结构: 采用先序遍历的方式, 将 DOM Tree 转化成一个标签序列。



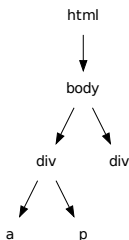
```
<html0><body1>  
<div2><a3><p3>  
<div2>
```

3. 利用后缀树检测并合并重复记录

预处理模块 - 简化 HTML 文档

简化 HTML 文档的 3 个步骤

1. 去除无用标签。<script>、<link>、<style>、
 等。
2. 简化解析好的 DOM Tree 树形结构: 采用先序遍历的方式, 将 DOM Tree 转化成一个标签序列。



```
<html0><body1>  
<div2><a3><p3>  
<div2>
```

3. 利用后缀树检测并合并重复记录

检测重复记录

动机

- ▶ 网页模板中有大量不定个数的重复模式，通常由模板语言的 for 语句生成，我们需要将这些记录合并

```
{% for comment in news.comments %}  
  <li>comment</li>  
{% endfor %}
```

.....

检测重复记录

动机

- ▶ 网页模板中有大量不定个数的重复模式，通常由模板语言的 for 语句生成，我们需要将这些记录合并

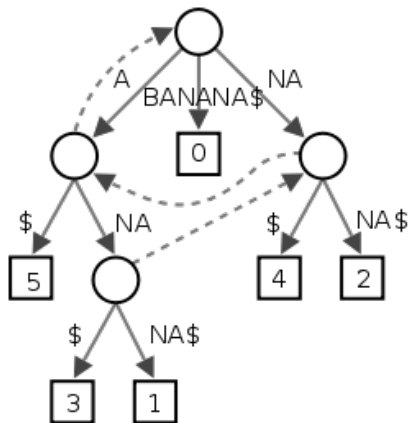
```
{% for comment in news.comments %}  
  <li>comment</li>  
{% endfor %}
```

后缀树简介

- ▶ 一种可以快速完成重复字串的查找的数据结构。
- ▶ 定义：由序列所有的后缀组成的 Trie 树。
- ▶ 后缀树普通的构造算法复杂度很高，系统实现中采用了 Ukkonen 在 1995 年提出了一个 $O(n)$ 时间复杂度的在线构造算法。

后缀树示例

BANANA 对应的的后缀树



后缀树查找重复序列算法

原始算法

- ▶ 任意一条从根节点到内部节点的路径组成的序列都是原序列中重复出现的字串, 且重复的次数是以该内部节点为根的子树的叶子节点个数。
- ▶ 缺点: 没有考虑结构信息, 从串映射回树结构时不一定符合要求。

后缀树查找重复序列算法

原始算法

- ▶ 任意一条从根节点到内部节点的路径组成的序列都是原序列中重复出现的字串, 且重复的次数是以该内部节点为根的子树的叶子节点个数。
- ▶ 缺点: 没有考虑结构信息, 从串映射回树结构时不一定符合要求。

错误情况举例

- ▶ `<div4><a5><div3><div4><a5><div3>` 中的重复串 `<div4><a5><div3>` 横跨了两个不同子树。

后缀树查找重复序列算法

原始算法

- ▶ 任意一条从根节点到内部节点的路径组成的序列都是原序列中重复出现的字串, 且重复的次数是以该内部节点为根的子树的叶子节点个数。
- ▶ 缺点: 没有考虑结构信息, 从串映射回树结构时不一定符合要求。

错误情况举例

- ▶ `<div4><a5><div3><div4><a5><div3>` 中的重复串 `<div4><a5><div3>` 横跨了两个不同子树。
- ▶ `<div3><div4><a5><p5><div3><p4><a5><p5>` 中的两个重复串 `<a5><p5>` 分属不同的子树。

后缀树查找重复序列算法

原始算法

- ▶ 任意一条从根节点到内部节点的路径组成的序列都是原序列中重复出现的字串, 且重复的次数是以该内部节点为根的子树的叶子节点个数。
- ▶ 缺点: 没有考虑结构信息, 从串映射回树结构时不一定符合要求。

错误情况举例

- ▶ `<div4><a5><div3><div4><a5><div3>` 中的重复串 `<div4><a5><div3>` 横跨了两个不同子树。
- ▶ `<div3><div4><a5><p5><div3><p4><a5><p5>` 中的两个重复串 `<a5><p5>` 分属不同的子树。
- ▶ `<div3><div4><div3><div4><a5><p5><div4><a5><p5>` 中的重复串 `<div3><div4>` 和 `<div4><a5><p5>` 有交集 `<div4>`。

新的算法要求

- ▶ 重复序列不能横跨两个子树
- ▶ 重复序列必须有公共的父亲
- ▶ 不同的重复序列不能相交，不能互相包含
- ▶ 重复序列必须尽可能地长

后缀树查找重复序列算法 -1

特点 1: 保证序列尽可能长

后缀树查找重复序列算法 -1

特点 1: 保证序列尽可能长

算法 2 从根节点出发，找出所有的重复子序列

输入: 已经构建好的后缀树，根为 *root*

输出: 该后缀树中所有的重复子序列

- 1: // 从根节点出发，寻找所有的重复子序列
 - 2: **for** *edge* \leftarrow *root.edges* **if** *edge.endNode.isNotLeaf* **do**
 - 3: // 取后缀树根节点的每条边的第一个元素作为每个子树的根节点
 - 4: *subTreeRoot* := *edge.firstElement*
 - 5: // 查找以该节点为根的所有重复子序列
 - 6: findAllRepetitions(*root*, **nil**, *subTreeRoot*)
 - 7: **end for**
-

后缀树查找重复序列算法 -2

特点 2: 保证序列不会横跨多个子树, 且不互相包含

后缀树查找重复序列算法 -2

特点 2: 保证序列不会横跨多个子树，且互相包含

算法 4 简化的 findAllRepetitions 实现

输入: 一个内部节点 *node*，当前已经找到的重复序列 *prefix*，要找的子树的根节点 *subTreeRoot*

输出: 所有经过该内部节点的符合要求的重复序列

```
function FINDALLREPETITIONS(node, prefix, subTreeRoot)
    // 定义一个空集合
    results := Collection.empty
    // 对于该内部节点的每一条不连接叶子节点的边
    for edge ← node.edges if edge.endNode.isNotLeaf do
        // 依次取出该条边上属于该根节点子树上的点
        seq := edge.takeWhile(element inSubTreeOf subTreeRoot)
        if seq.length == edge.length then
            // 遍历完了该条边上所有元素，
            // 则取该条边连接的下一个内部节点进行递归查找
            findAllRepetitions(edge.endNode, prefix + seq, subTreeRoot)
        else
            // 否则，将当前得到的序列加入到结果集合中
            addToResults(prefix + seq, results)
        end if
    end for
    return results
end function
```

后缀树查找重复序列算法 -3

特点 3: 保证重复序列相交均在同一子树中, 且无相交的情况

后缀树查找重复序列算法 -3

特点 3: 保证重复序列相交均在同一子树中，且无相交的情况

- ▶ 所有序列按实际的父节点分组，然后对每个分组进行合并，保证重复串有公共的父亲。
- ▶ 对于重复串存在交集的情况，如之前的例子：

`<div3><div4><div3><div4><a5><p5><div4><a5><p5>`

对于 `<div3><div4>` 和 `<div4><a5><p5>`，只取

`<div4><a5><p5>`

提纲

研究内容介绍

系统框架与模块实现

整体框架

预处理模块

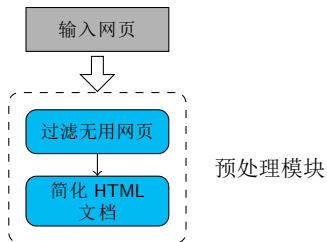
网页聚类模块

模板生成和内容提取模块

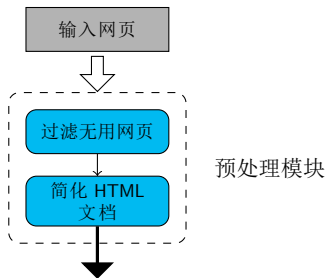
实验和分析

总结和展望

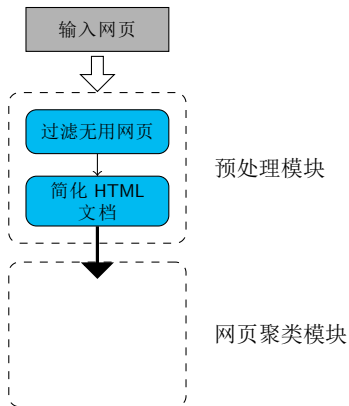
网页聚类模块介绍



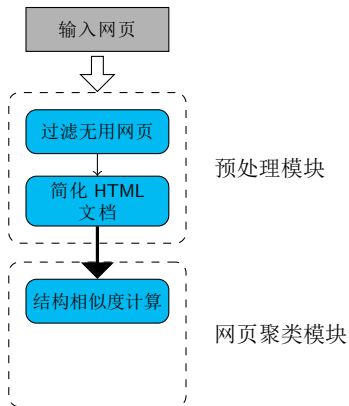
网页聚类模块介绍



网页聚类模块介绍



网页聚类模块介绍



网页聚类模块 - 计算结构相似度

- ▶ 序列 S_1, S_2 , x_i 和 y_j 分别表示 S_1 的第 i 个元素和 S_2 的第 j 个元素, $f(\text{depth})$ 是根据深度加权的函数, 我们认为深度越深的节点, 成为模板的概率就越小。

$$t(i)(j) = \begin{cases} 0 & i = 0, j = 0 \\ t(i-1)(j-1) + f(x_i.\text{depth}) & i, j > 0, x_i = y_j \\ \max(t(i)(j-1), t(i-1)(j)) & i, j > 0, x_i \neq y_j \end{cases} \quad (1)$$

网页聚类模块 - 计算结构相似度

- ▶ 序列 S_1, S_2 , x_i 和 y_j 分别表示 S_1 的第 i 个元素和 S_2 的第 j 个元素, $f(\text{depth})$ 是根据深度加权的函数, 我们认为深度越深的节点, 成为模板的概率就越小。

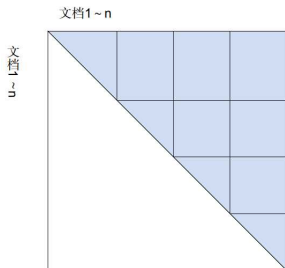
$$t(i)(j) = \begin{cases} 0 & i = 0, j = 0 \\ t(i-1)(j-1) + f(x_i.\text{depth}) & i, j > 0, x_i = y_j \\ \max(t(i)(j-1), t(i-1)(j)) & i, j > 0, x_i \neq y_j \end{cases} \quad (1)$$

- ▶ 结构相似度计算公式

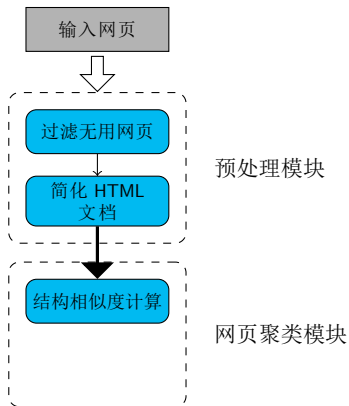
$$\text{Sim}(D_1, D_2) = \frac{|elcs(S_1, S_2)|}{\max(\sum_{n \in S_1} f(n.\text{depth}), \sum_{n \in S_2} f(n.\text{depth}))}$$

计算时间优化

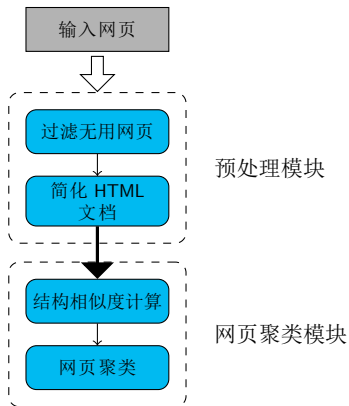
- ▶ 文档数量多，计算量很大，需要一定的优化。
- ▶ 采用 Actor 库实现多线程计算，加快计算速度。将任务分割后交给每个 Actor 进行计算，Actor 的调度的算法采用 RoundRobin。
- ▶ 任务分割示意图



网页聚类模块介绍

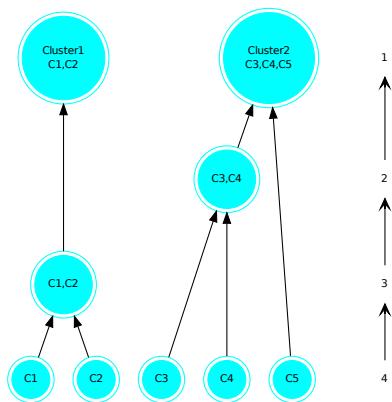


网页聚类模块介绍



聚类算法

1. 初始时，让每个文档实例都单独为一类
2. 迭代时，每次选择距离最近的两个类合并
3. 直到任意两个类的距离都大于阈值时程序退出
4. 算法的结束条件由阈值决定，无需实现设定类的个数。



提纲

研究内容介绍

系统框架与模块实现

- 整体框架

- 预处理模块

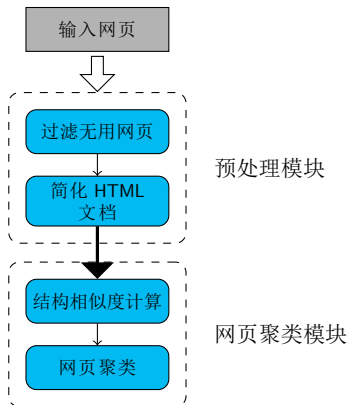
- 网页聚类模块

- 模板生成和内容提取模块

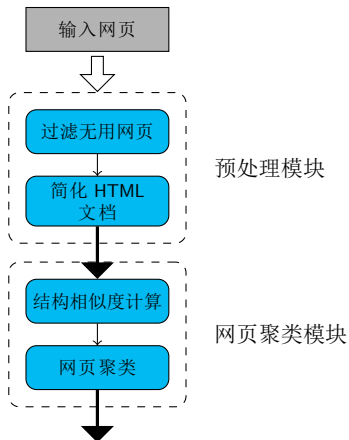
实验和分析

总结和展望

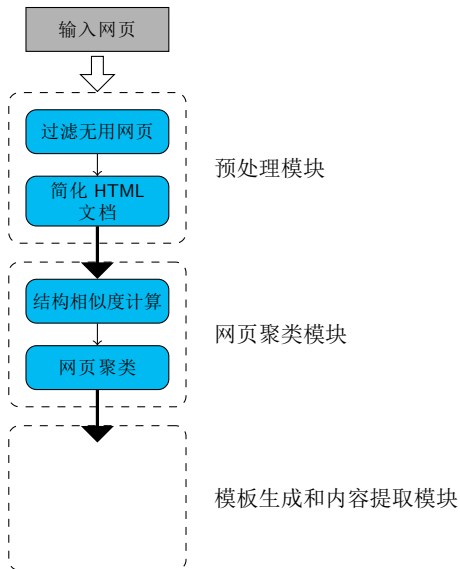
模板生成和内容提取模块介绍



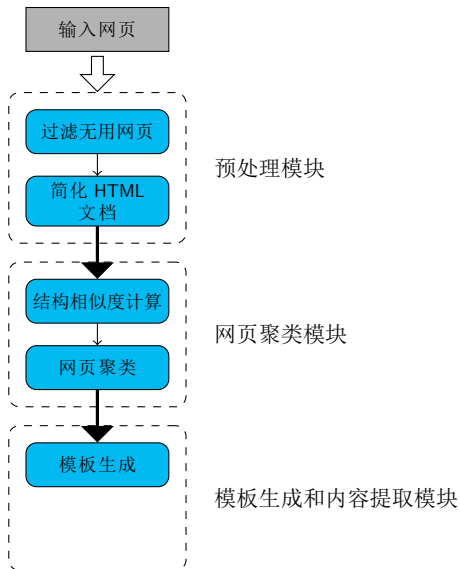
模板生成和内容提取模块介绍



模板生成和内容提取模块介绍



模板生成和内容提取模块介绍



模板形式化定义

基本节点

- ▶ 两种组成方式

模板形式化定义

基本节点

- ▶ 两种组成方式

1. 单个不重复的 HTML 标签，即 `< tag >`

模板形式化定义

基本节点

► 两种组成方式

1. 单个不重复的 HTML 标签，即 $\langle tag \rangle$
2. 由一个或多个 HTML 标签组成的序列，这些序列可以出现一次或多次

模板形式化定义

基本节点

- ▶ 两种组成方式
 1. 单个不重复的 HTML 标签，即 `< tag >`
 2. 由一个或多个 HTML 标签组成的序列，这些序列可以出现一次或多次
- ▶ 第二种形式通过合并重复记录得到，对应的模板语言为：

```
{% for comment in news.comments %}  
  <li>comment</li>  
{% endfor %}
```

模板形式化定义

基本节点

► 两种组成方式

1. 单个不重复的 HTML 标签，即 `< tag >`
2. 由一个或多个 HTML 标签组成的序列，这些序列可以出现一次或多次

► 第二种形式通过合并重复记录得到，对应的模板语言为：

```
{% for comment in news.comments %}  
  <li>comment</li>  
{% endfor %}
```

► 还有一种模板生成形式：

```
{% if news.has_subtitle %}  
  <h1>news.title</h1>  
  <h2>news.sub_title</h2>  
{% else %}  
  <h1>news.title</h1>  
{% endif %}
```

必选和可选节点

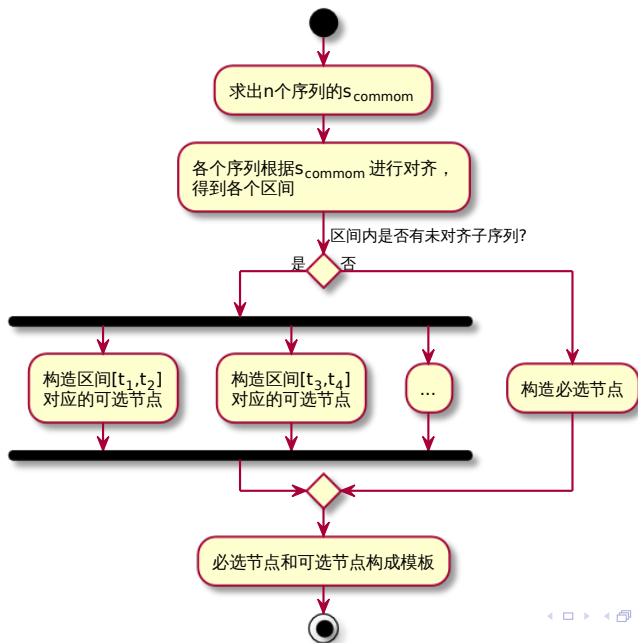
- ▶ 必选节点对应着由基本节点组成的一个序列
- ▶ 可选节点同时对应多个序列，每个序列由不同的基本节点组成，同时每个序列还对应着一个出现概率
- ▶ 我们规定：模板 Tp 必须由必选节点 TN_i 和 ON_i 交替组成，第一个和最后一个可选节点是可选的：

$$Tp = [ON_0]EN_1ON_1EN_2ON_2.....EN_n[ON_n]$$

模板生成算法

- ▶ 找出所有组成必选节点的基本节点：对于一个聚类的全部序列，从聚类中心点开始，依次计算一次最长公共子序列，得到 n 个序列的公共子序列 S_{common}
- ▶ 每个序列和 S_{common} 对齐，得到一些未对齐的区间，每个未对齐的区间计算一个可选节点
- ▶ 已对齐的基本节点组成必选节点
- ▶ 必选节点和可选节点交替出现，组成最终的模板

模板生成流程图



简单的例子

- ▶ 设有 3 个序列，分别为：

$s_1 = aorzbcdlx e$

$s_2 = athubeatcdlx e$

$s_3 = athubeatcdpkue$

- ▶ 根据最长公共子串算法，得到 s_{common} 为 *abcde*
- ▶ 每个序列同 s_{common} 进行对齐，得到

s_1	:	a	orz	b		cd	lx	e
s_2	:	a	thu	b	eat	cd	lx	e
s_3	:	a	thu	b	eat	cd	pkue	e
s_{common}	:	a		b		cd		e

简单的例子

- ▶ 根据上述对齐结果，生成以下可选节点：

$$ON_{a,b} = (thu, 2/3) \mid (orz, 1/3)$$

$$ON_{b,c} = eat, 2/3$$

$$ON_{d,e} = (lx, 2/3) \mid (pku, 1/3)$$

- ▶ 对齐的部分则对应生成必选节点，注意将连续的基本节点合并，分别为： a, b, cd, e

可选节点生成的实际实现

- ▶ 实际情况要比上面的例子复杂。区间 $[t_1, t_2]$ 中未对齐的那些标签子序列, 有些可能差异很大, 有些则可能非常相近。

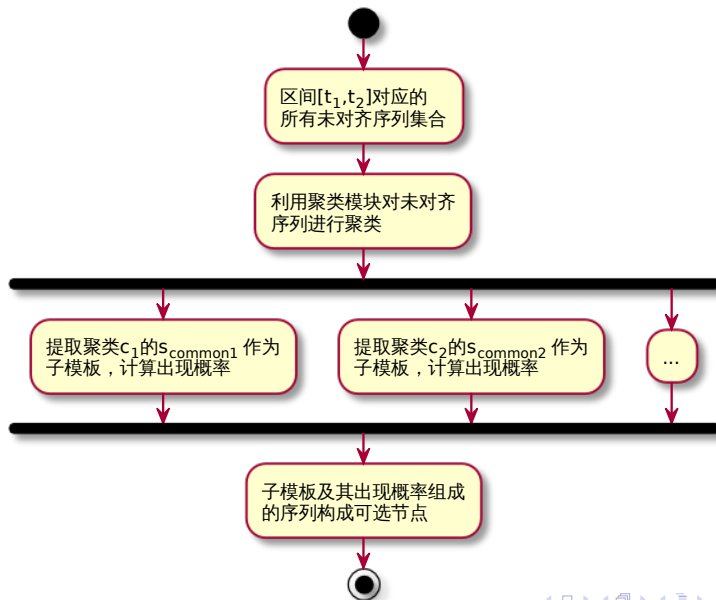
可选节点生成的实际实现

- ▶ 实际情况要比上面的例子复杂。区间 $[t_1, t_2]$ 中未对齐的那些标签子序列, 有些可能差异很大, 有些则可能非常相近。
- ▶ 做法: 对于每个区间 $[t_1, t_2]$, 先将几种不同的子序列分开, 然后针对每个子序列集合提取“模板”。

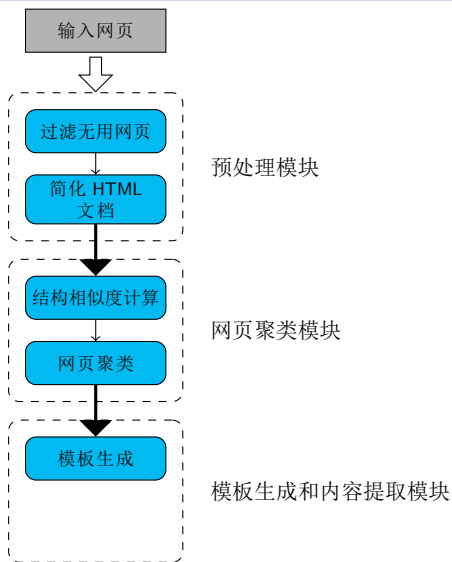
可选节点生成的实际实现

- ▶ 实际情况要比上面的例子复杂。区间 $[t_1, t_2]$ 中未对齐的那些标签子序列, 有些可能差异很大, 有些则可能非常相近。
- ▶ 做法: 对于每个区间 $[t_1, t_2]$, 先将几种不同的子序列分开, 然后针对每个子序列集合提取“模板”。
- ▶ 我们发现, 这和整体的框架很相似。因此, 可选节点生成问题可以使用一个“自相似”的框架解决。

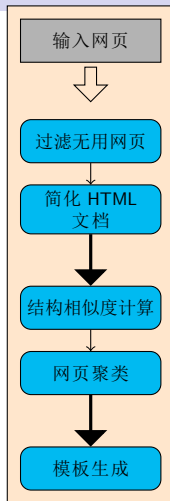
构造可选节点流程



模板生成和内容提取模块介绍

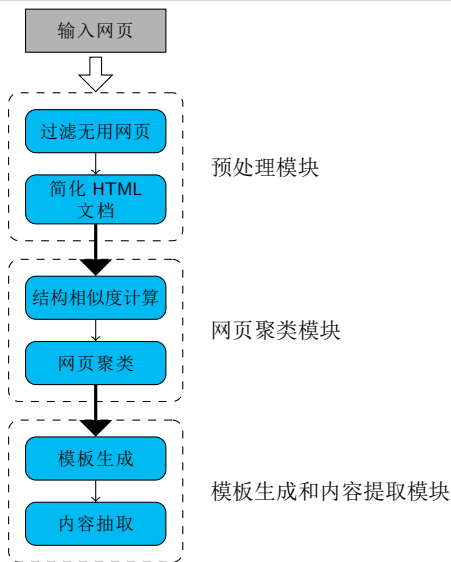


模板生成和内容提取模块介绍

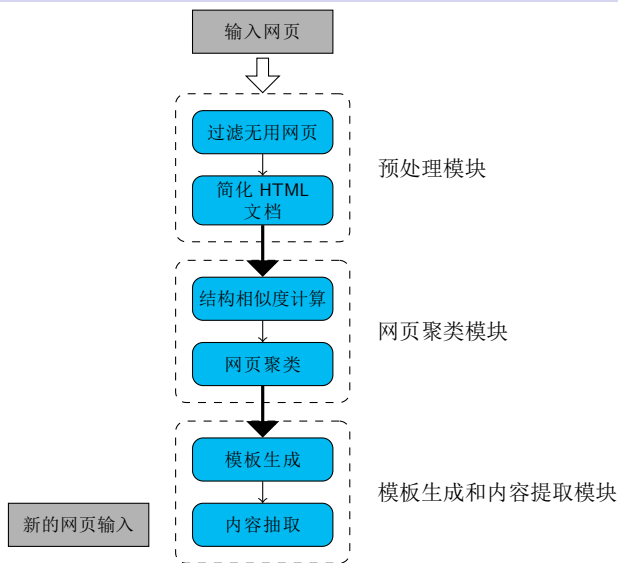


离线处理部分

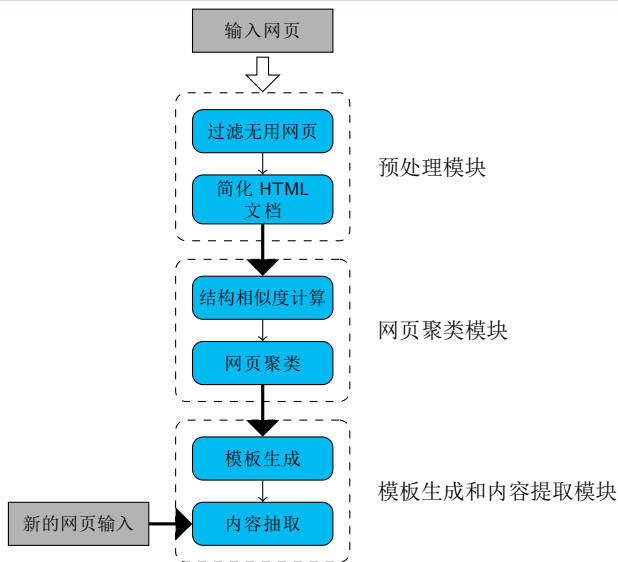
模板生成和内容提取模块介绍



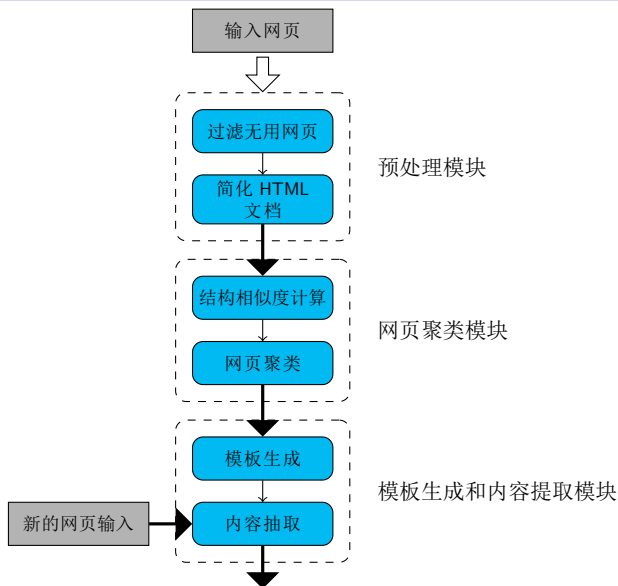
模板生成和内容提取模块介绍



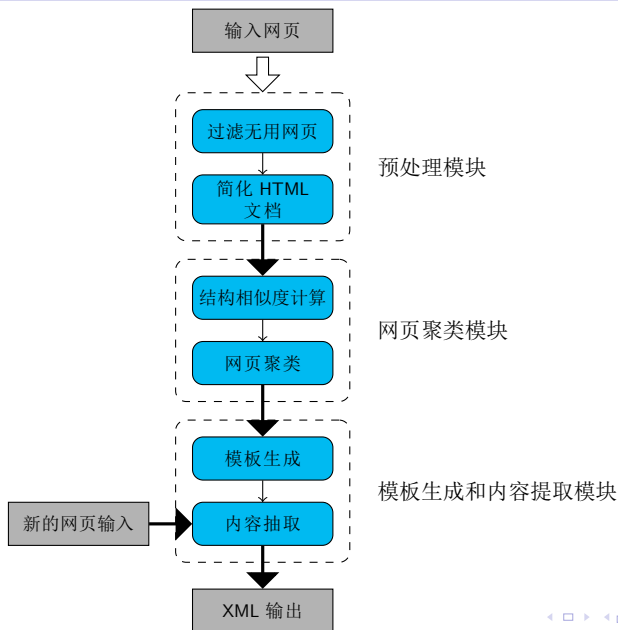
模板生成和内容提取模块介绍



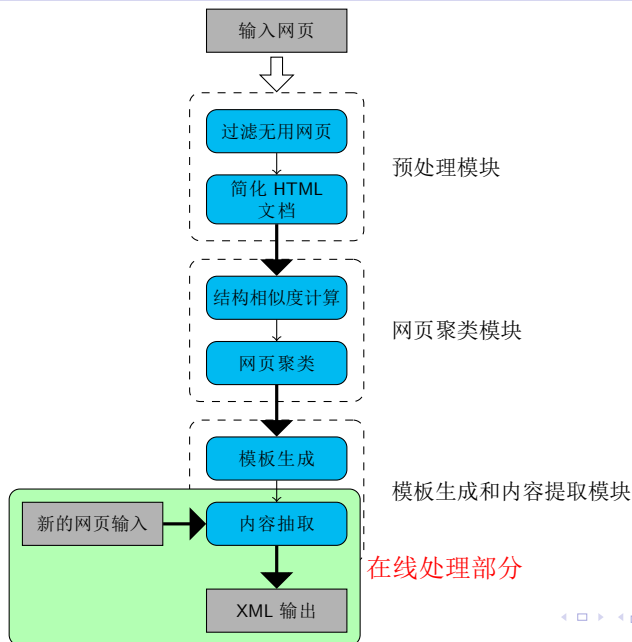
模板生成和内容提取模块介绍



模板生成和内容提取模块介绍

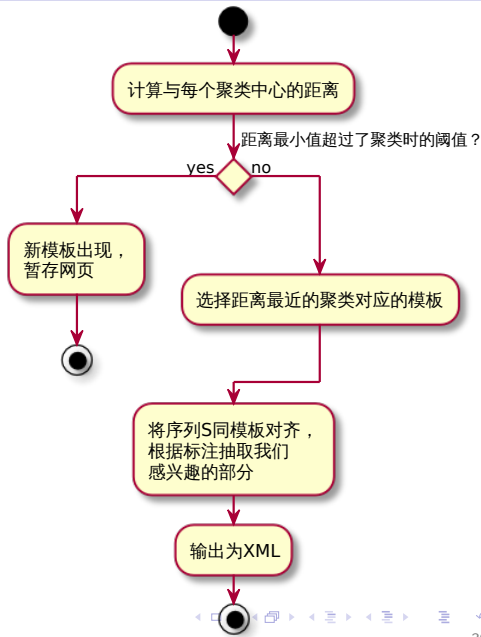


模板生成和内容提取模块介绍



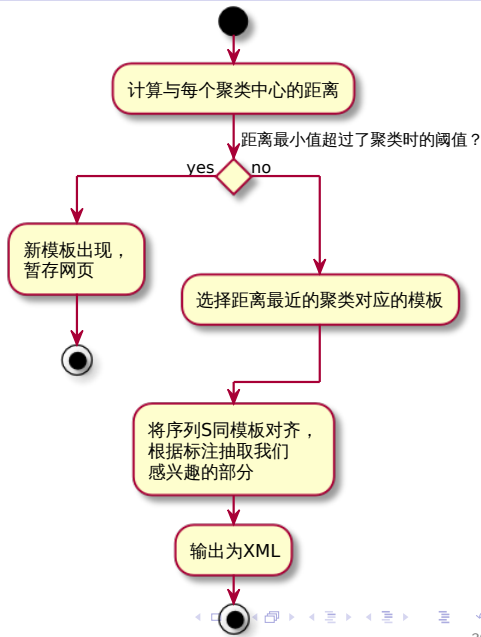
内容提取模块

1. 人工标注感兴趣的内容，如新闻标题，正文等。



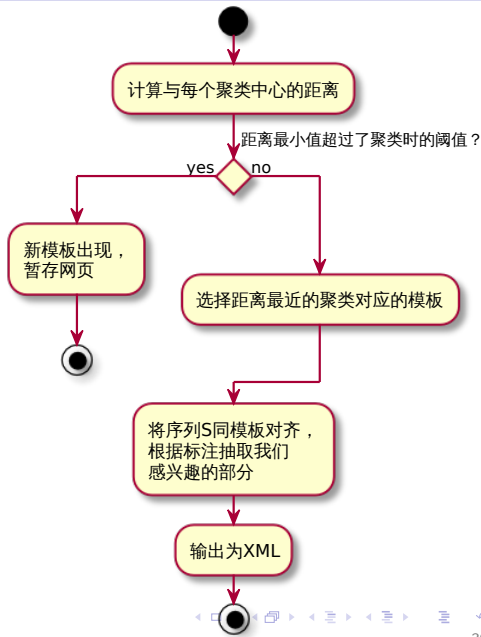
内容提取模块

1. 人工标注感兴趣的内容，如新闻标题，正文等。
2. 计算距离，判断使用哪个模板



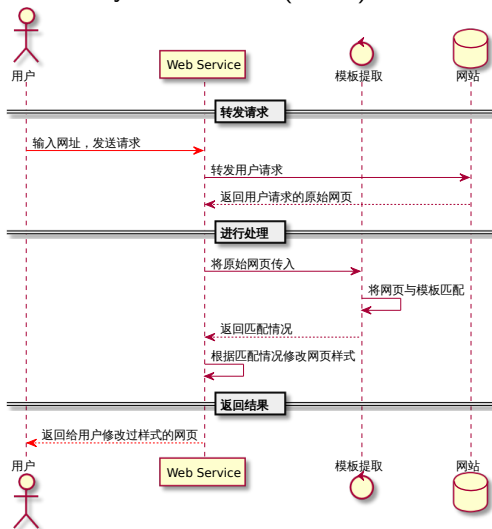
内容提取模块

1. 人工标注感兴趣的内容，如新闻标题，正文等。
2. 计算距离，判断使用哪个模板
3. 与模板对齐，抽取内容



模板匹配演示系统

- 基于 Play! Framework(Scala) 实现了一个 Web Service



提纲

研究内容介绍

系统框架与模块实现

- 整体框架

- 预处理模块

- 网页聚类模块

- 模板生成和内容提取模块

实验和分析

总结和展望

- ▶ 机器的配置情况是：16 个逻辑核的 CPU，24G 内存，操作系统为 64 位的 CentOS。
- ▶ 数据统计

	blog	news
文件个数	59998	81561
总大小	5.4G	7.9G
来源	blog.sina.com.cn	ent.sina.com.cn

► 过滤目录页的规则

	blog	news
目录页 URL 规则	.*(?<!\.html?)\$.*(?<!\.s?html)\$
错误页最大长度	6000	6000
原来总文件个数	59998	81561
过滤后详细页个数	35466	65655

► 过滤目录页的规则

	blog	news
目录页 URL 规则	<code>.*(?<!\.html?)\$</code>	<code>.*(?<!\.s?html)\$</code>
错误页最大长度	6000	6000
原来总文件个数	59998	81561
过滤后详细页个数	35466	65655

► 均选取 10000 个样本作为训练集。

无用标签去除规则

正则表达式模式	对应的标签名
(?is)<tag.*?>.*?</tag>	style,script
(?is)<[/]?tag.*?>	link,input,br,img,meta,wbr
(?is)<[/]?tag.*?>	strong,em,font,b,p,table

重复记录长度分布统计

- 在两个个数为 10000 的训练集 blog 和 news 上分别做统计

重复记录长度	blog	news
1	1365336	1278686
2	177588	63475
3	5921	11733
4	1224	391
5	243	2872
6	95	22
7	18958	23
8	50	40
9	60	356
≥ 10	745	3575

聚类和模板提取

- ▶ 选择 news 数据集做实验，调整聚类时的距离的阈值，观察聚类结果及模板抽取结果的变化

距离阈值	聚类个数	必选节点 平均长度	可选节点 平均长度	模板平均长度
0.3	16	5.18	2.14	21.8
0.4	9	4.66	2.59	22.6
0.5	7	4.68	2.71	24.9
0.6	5	3.76	2.90	24.8

提取效果

- 我们在 5000 个测试集上运行我们的程序，抽取出我们需要的信息，并将其保存为 XML 格式输出。对于博客，我们抽取标题，作者和正文 3 个内容。

```
</document>
<document name="/home/qiu junpeng/Data/blog_detail/http%3A%2Fblog.sina.com.cn%2Fs%2Fblog_10bc233401
019ao6.html">
  <TITLE>缘分_夏诺的宝盒_新浪博客</TITLE>
  <AUTHOR>夏诺的宝盒的博客</AUTHOR>
  <CONTENT>缘分 (2012-10-13 14:50:18) 转载 ▼原文地址: 缘分 (来自 @轻博客 ) 作者: 芸芸的云_0baea3
  _ _ 缘分,一个美丽而又飘渺的词语,多少人因它而走到了一起,又有多少人因它而一次次 陷入迷途。 _ 缘分就是
  缘到自然分的无赖,是别离的痛苦,是无情得分手 _ 但缘分也是美丽的 _ 当缘分未来时。我们应该学会等候 _ 当它
  来时,我们要更加好好的去珍惜 _ 不要无视缘分,因为前世五百年的回眸才 _ 换来今生的一次cha肩而过 _ 缘分虽
  然不存在于现实 _ 但它是感情的一种表达 _ 相信缘分 _ 也许它就在你身边 _ ---芸子 _ 分享: 分享到新浪Qing 喜
  欢阅读 ! 评论 ! 收藏 ! 转载 ! 喜欢 ▼ ! 打印 ! 举报 已投稿到: 排行榜</CONTENT>
</document>
<document name="/home/qiu junpeng/Data/blog_detail/http%3A%2Fblog.sina.com.cn%2Fs%2Fblog_0435c17801
02drbt.html">
  <TITLE>地方政府或减持股票还债_36家公司最可能被减持_百闻_新浪博客</TITLE>
  <AUTHOR>百闻博客</AUTHOR>
  <CONTENT>地方政府或减持股票还债_36家公司最可能被减持 (2011-07-26 12:41:39) 转载 ▼地方政府或减持股
  票还债 36家公司最可能被减持 2011年07月26日07:25 瑞银证券 对于部分市场人士关心的“地方政府是否会减持股票
```

- 模板匹配演示系统的运行效果如下：



提纲

研究内容介绍

系统框架与模块实现

- 整体框架

- 预处理模块

- 网页聚类模块

- 模板生成和内容提取模块

实验和分析

总结和展望

总结

1. 设计并实现了一套利用后缀树在树的先序遍历序列中查找重复子树的算法。

总结

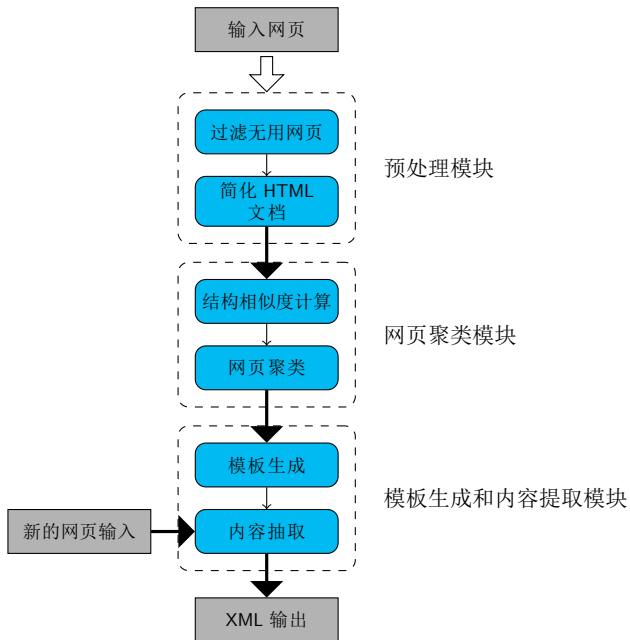
1. 设计并实现了一套利用后缀树在树的先序遍历序列中查找重复子树的算法。
2. 实现并改进了最长公共子序列算法，将其用于计算文档的结构相似度；实现了简单的凝聚层次聚类算法，并将其用于文档的聚类

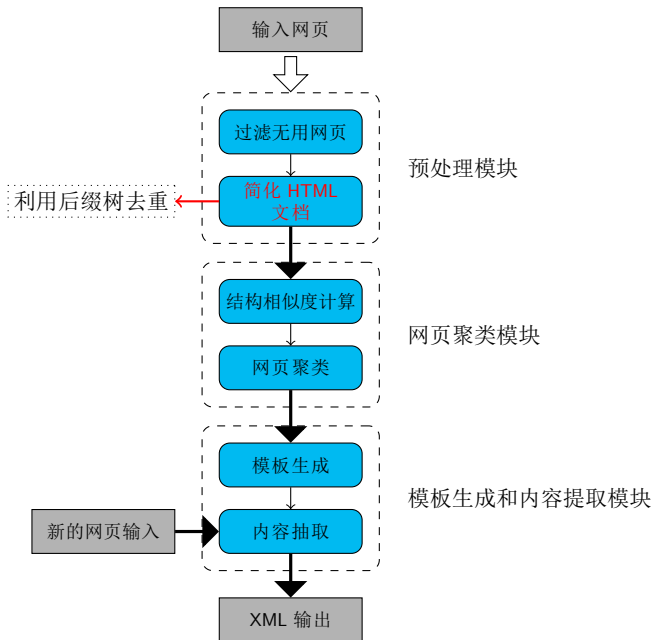
总结

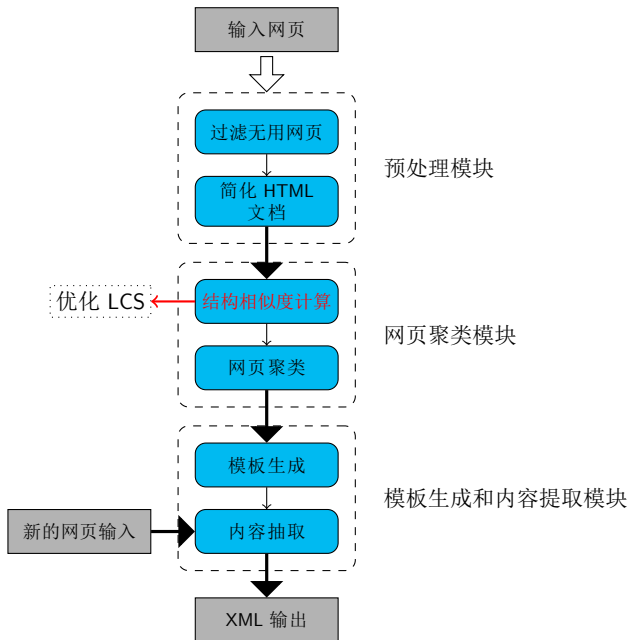
1. 设计并实现了一套利用后缀树在树的先序遍历序列中查找重复子树的算法。
2. 实现并改进了最长公共子序列算法，将其用于计算文档的结构相似度；实现了简单的凝聚层次聚类算法，并将其用于文档的聚类
3. 实现了一个无监督的模板生成算法，通过人工指定模板中某些部分对应的语义，可以使用模板提取新的由该模板生成的网页中对应的信息。

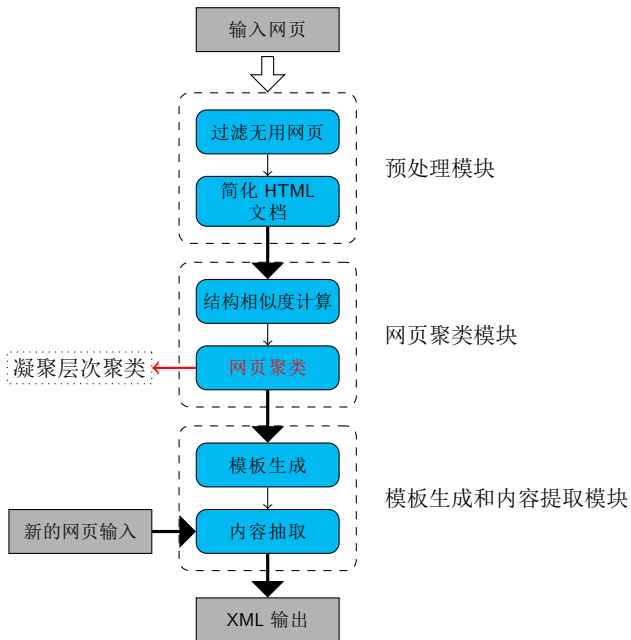
总结

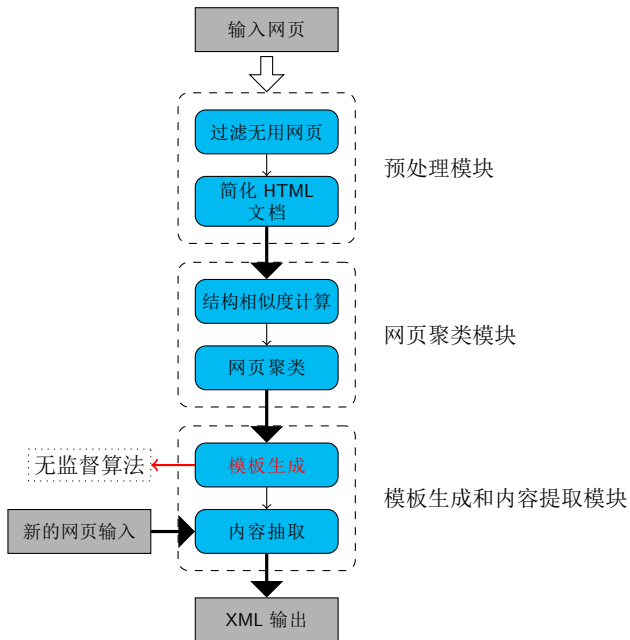
1. 设计并实现了一套利用后缀树在树的先序遍历序列中查找重复子树的算法。
2. 实现并改进了最长公共子序列算法，将其用于计算文档的结构相似度；实现了简单的凝聚层次聚类算法，并将其用于文档的聚类
3. 实现了一个无监督的模板生成算法，通过人工指定模板中某些部分对应的语义，可以使用模板提取新的由该模板生成的网页中对应的信息。
4. 实现了一个 Web Service，可以直观地看到网页和模板的匹配情况。

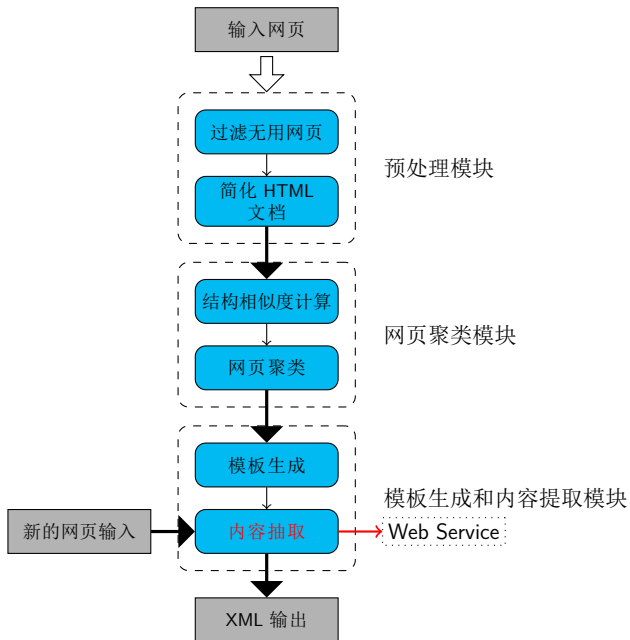












未来的工作

1. 预处理模块：改进 Ukkonen 原有的在线构造算法，在构造树的时候就考虑一些原本的结构信息。查找重复时利用后缀树中其他的一些信息，比如后缀链。

未来的工作

1. 预处理模块：改进 Ukkonen 原有的在线构造算法，在构造树的时候就考虑一些原本的结构信息。查找重复时利用后缀树中其他的一些信息，比如后缀链。
2. 网页结构相似度计算和网页聚类模块：可以进一步改进计算相似度的算法，加入一些其他的除了深度之外的更多结构信息；可以考虑使用更复杂但是鲁棒性更好的一些聚类算法。

未来的工作

1. 预处理模块：改进 Ukkonen 原有的在线构造算法，在构造树的时候就考虑一些原本的结构信息。查找重复时利用后缀树中其他的一些信息，比如后缀链。
2. 网页结构相似度计算和网页聚类模块：可以进一步改进计算相似度的算法，加入一些其他的除了深度之外的更多结构信息；可以考虑使用更复杂但是鲁棒性更好的一些聚类算法。
3. 模板生成和内容提取：对模板生成算法做一些修改，设计一些更高级的机器学习算法。在内容提取的时候，可以采用树的匹配算法而不是序列的匹配算法。

致谢

衷心感谢朱小燕老师和郝宇老师在整个毕设过程的悉心指导。他们不仅在整个系统的实现和优化上提出了很多宝贵的意见和建议，而且在科研的态度和方法上也深深影响了我。

感谢黄民烈老师在方法上提供的一些建议，感谢实验室师兄师姐的帮助。实验室良好的科研氛围给我的毕设提供了巨大的动力。