

# EdmondFrank's 时光足迹

この先は暗い夜道だけかもしれない　それでも信じて進むんだ。星がその道を少しでも照らしてく  
れるのを。  
或许前路永夜，即便如此我也要前进，因为星光即使微弱也会我为照亮前途。  
——《四月は君の嘘》

Blog   Archives



## 比较网页结构相似度

- -

### Table of Contents

- [1. 总体介绍](#)
- [2. 最长公共子序列](#)
- [3. 递归式展示](#)
- [4. 算法实现（python实现）](#)
- [5. 网页相似度计算](#)

## 总体介绍

网页网页结构相似度计算通常是网页自动分类的基础，在一般的网页信息提取中，判断网页片断是“噪声”还是“有效信息”通常是个两类分类问题。 简单地，我们可以把一般网页分为三个类，即：

- 目录导航式页面（List\Index Page）
- 详细页面（Detail Page）
- 未知页面（Unknown Page）

由于网页本身就可以抽象成串行的节点或者是DOM树，那么对于串行序列，就可以常用最长公共子序列来衡量相似度

# 最长公共子序列

最长公共子序列是动态规划的基本问题：

序列a共有m个元素，序列b共有n个元素，如果a[m-1]==b[n-1]，

那么a[:m]和b[:n]的最长公共子序列长度就是a[:m-1]和b[:n-1]的最长公共子序列长度+1；

如果a[m-1]!=b[n-1]，那么a[:m]和b[:n]的最长公共子序列长度就是

MAX (a[:m-1]和b[:n]的最长公共子序列长度， a[:m]和b[:n-1]的最长公共子序列长度)

## 递归式展示

$$c[i, j] = \begin{cases} 0 & \text{if } i = 0 \text{ or } j = 0, \\ c[i - 1, j - 1] + 1 & \text{if } i, j > 0 \text{ and } x_i = y_j, \\ \max(c[i, j - 1], c[i - 1, j]) & \text{if } i, j > 0 \text{ and } x_i \neq y_j. \end{cases}$$

## 算法实现（python实现）

```
1  #params:
2  # - a : str
3  # - b : str
4  #return
5  # - c : 过程处理矩阵
6  # - c[x][y] : the lcs-length(最长公共子序列长度)
7  def lcs(a, b):
8      lena=len(a)
9      lenb=len(b)
10     c=[[0 for i in range(lenb+1)] for j in range(lena+1)]
11     for i in range(lena):
12         for j in range(lenb):
13             if a[i]==b[j]:
14                 c[i+1][j+1]=c[i][j]+1
15             elif c[i+1][j]>c[i][j+1]:
16                 c[i+1][j+1]=c[i+1][j]
17             else:
18                 c[i+1][j+1]=c[i][j+1]
19     return c, c[lena][lenb]
```

# 网页相似度计算

```
# -*- coding: utf-8 -*-
1 import lxml.html.soupparser as soupparser
2 import requests
3 headers = {
4     "User-Agent": "Mozilla/5.0 (X11; Linux x86_64) AppleWebKit/537.36 (KHTML, like
5     Gecko) Chrome/55.0.2883.87 Safari/537.36"
6 }
7 def get_domtree(html):
8     dom = soupparser.fromstring(html)
9     for child in dom.iter():
10         yield child.tag
11
12 def similar_web(a_url, b_url):
13     html1 = requests.get(a_url, headers=headers).text
14     html2 = requests.get(b_url, headers=headers).text
15     dom_tree1 = ">".join(list(filter(lambda e:
16     isinstance(e, str), list(get_domtree(html1)))))
17     dom_tree2 = ">".join(list(filter(lambda e:
18     isinstance(e, str), list(get_domtree(html2)))))
19     c, flag, length = lcs(dom_tree1, dom_tree2)
20     return 2.0*length/(len(dom_tree1)+len(dom_tree2))
21
22 percent = similar_web(
23 'http://edmondfrank.github.io/blog/2017/04/05/qi-an-tan-mongodb/',
24 'http://edmondfrank.github.io/blog/2017/03/27/emacsshi-yong-zhi-nan/')
25 print(percent) #相似度 (百分比)
```

Posted by EdmondFrank • • [python](#)

[Tweet](#)

[« 浅谈MongoDB](#)

[Web正文提取\(偏純文本类\) »](#)

## About Me

## Recent Posts

## Popular Posts



GitHub: [@EdmondFrank](#)

Twitter: [@EdmondFrank4](#)

Blog:

<https://edmondfrank.github.io>

この町、冗談と気まぐれと偶然でて  
きっているらしい。