

题目 102

要求:

1、用 python 将所给数据集中的样本数据, 针对不同的网页特征进行特征提取, 分别统计出特征数据集;

2、特征数据预处理: 对所有样本进行特征提取后, 对每个特征的特征值进行归一化处理(这里可采用线性归一化处理)和向量化(文本特征可采用 TF/TF-IDF 等向量化)处理, 将特征值缩放至区间[0,1]。

3、经过归一化处理后, 将不同类的数据集进行特征融合, 所有网页样本均表示为一个 N 维的特征向量 $f=[u_1,...,u_7,h_1,...,h_{10},j_1,...,j_{18},...]^T$ 。

一、URL 特征提取规则:

U1: URL 的长度 (含 http://)

U2: URL 中点的个数

U3: URL 中数字的个数

U4: URL 中特殊字符的个数 (特殊字符: “#、@、_、-、&、/、=”)

U5: URL 中子域名的个数

U6: URL 路径深度 (层数)

U7: URL 中是否包含 IP 地址 (为二值特征, 包含 IP 地址时为 1, 否则为 0。)

Flag 为: n, d, p。 (可将其表示: n 为 0, d 为 1, p 为 2)。

表 1: URL 特征:

Id	U1	U2	U3	U4	U5	U6	U7	Flag
1								
2								
3								
...								
n								

二、HTML 特征提取规则:

H1: HTML 长度

H2: 网页文本内容占整个 HTML 文件长度的比重

H3: URL 出现次数

H4: 隐藏标签的数量 (主要考察标签的 size 属性(包括 width 属性和 height 属性)、hidden 属性、display 属性以及 visible 属性)

H5: <iframe>标签数量

H6: <meta>标签个数

H7: <image>标签数量

H8: <script>标签数量

H9: <object>标签数量

H10: <embed>标签数量

Flag 为: n, d, p。 (可将其表示: n 为 0, d 为 1, p 为 2)。

表 2:: HTML 特征:

Id	H1	H2	H3	H4	H5	H6	H7	H8	H9	H10	Flag
1											
2											
3											
...											
n											

三、JavaScript 特征提取规则:

J1: eval 函数调用次数

J2: 长字符串数量(长度大于 30)

J3: JavaScript 代码长度

J4: 可疑文件出现次数。(可疑文件主要指.exe、.ini、.dll、.tmp 等后缀文件)

J5: setTimeout 函数数量

J6: setInterval 函数数量

J7: window.location 函数数量

J8: Window.open 函数数量

J9: scriptObject.src 数量

J10: scriptObject.setAttribute 数量

J11: scriptObject.innerHTML 数量

J12: document.location 函数数量

J13: document.cookie 函数数量

J14: document. Write()函数个数

J15: unescape()和 escape()函数个数

J16: split()和 replaoe()函数个数

J17: JavaScript 代码占 HTML 文档的比率

J18: Navigator 字符串出现次数

Flag 为: n, d, p。(可将其表示: n 为 0, d 为 1, p 为 2)。

表 3: JavaScript 特征:

Id	J1	J2	J3	...	J16	J17	J18	Flag
1								
2								
3								
...								
n								

四、文本特征:

提取页面文本特征 W1, W2, ..., Wn。

(1) 文本抽取:首先对网页进行 HTML 解析, 并利用 Xpath 分别提取 Web 网页中<title> 标签里面的网页标题、<meta>标签里面的 keyword 和 description 以及<script>标签中 alert 方法中的指定消息, 并将抓取出来的文本信息按行存入 txt 文本中。

如: 抽取的部分文本如下所示。

[illegible]