# FEATURE SELECTION AND CLASSIFICATION BASED ON URL NAVIGATION

**LALITENDRASINGH S. PAYAL**
PHD Scholar
Shri JJT University
Rajasthan
lalitendrasinghpayal@gmail.com

**Dr.  YOGESH  KUMAR SHARMA**
HOD
Shri JJTU, Rajashthan

**ABSTRACT**

*Site page course of action has been generally analyzed, using various sorts of features that are isolated either from the page content, the page structure or from various pages that interface with that page. Using features from the page itself derives downloading it before its portrayal. We show an examination to check that URL tokens contain information enough to remove features to mastermind site pages. A classifier in perspective of these features can arrange a webpage page without downloading it previously, evading futile downloads.*

*Keywords: Web KB Dataset, Term Frequency-Inverse Document Frequency (TF-IDF), Web page classification.*

## I. INTRODUCTION

The extension in use of Web and its improvement are remarkable. Printed data on the Web is evaluated at one land byte, despite sound and video pictures which powers new challenges to Web lists. Web registries enable customers to look through the Web, by describing Web chronicles into subjects. Pages manual portrayal perseveres as Web reports increase [1].Text gathering hopes to arrange documents into a specific number of predefined classes using record features. Content portrayal has a basic part in recuperation and organization endeavors like information extraction, information recuperation, report filtering, and assembling different leveled records [2]. Exactly when content course of action revolves around site pages, it is called web request or page gathering. Request doles out predefined class imprints to hid or test

data. For this, a game plan of stamped data readies a classifier which by then names hid data. Request is controlled learning. The method isn't various in site page arrange so that there are no less than one predefined class names. Portrayal demonstrate apportions class names to site pages which are hypertext with various features like printed tokens, markup names, URLs and host names in URLs. As site pages have additional properties, this portrayal differs from customary substance arrange. Page gathering has subfields like subject portrayal and viable game plan. In the past, the classifier is stressed over site page content and chooses the site page "subject".

For example, online day by day papers orders like reserve, amusement, and development are instances of subject request. Valuable request oversees limit or kind of a site page. For example, choosing if a site page is an "individual point of arrival" or a "course page" is valuable request. Subject and down to earth arrange are well known portrayal forms. A HTML chronicle's individual portion is a HTML segment made up of a tree of HTML parts and center points like substance centers with all segments having decided qualities. Parts have substance and substance. HTML addresses semantics or centrality. HTML markup has key parts, including character references, character-based data composes, segments (and qualities), and

**Anveshana's International Journal of Research in Engineering and Applied Sciences**
**EMAILID:anveshanaindia@gmail.com,WEBSITE:www.anveshanaindia.com**
104

component references. Another part is file make disclosure, initiating benchmarks mode rendering. Semantic HTML is making HTML emphasizing encoded information's noteworthiness over presentation (looks). HTML joins semantic markup from cause, and presentational markup like , and names. There are in like manner semantically impartial names. HTML and related traditions from root were recognized quickly. In any case, there were no benchmarks in the early years of tongue. Notwithstanding the way that HTML was considered as a semantic tongue without presentation inconspicuous components, helpful use pushed presentational segments and credits to it, driven by changed program dealers. Latest HTML standards are tries to crush turbulent tongue headway and to make a rational foundation to create noteworthy and top notch documents. Feature assurance is a basic course of action step. Pages are in HTML sort out suggesting that site pages are semi composed data, with HTML marks and hyperlinks despite unadulterated substance. Due to this site pages property, incorporate decision in portrayal is exceptional in connection to standard gathering. Feature assurance diminishes data estimation with tens or hundreds or thousands of features which can't be taken care of further. A vital issue of site page portrayal is the high dimensionality of the part space. Best component subsets have least features that most add to arrange exactness.

To improve page game plan execution, various approachs imported from feature decision or substance arrange were associated. Information increment, shared information, record repeat, and term quality are notable part assurance methodology. Information get (IG) measures information in bits about the

class desire, when the fundamental information available is a part and relating class scattering. This examination proposes another component assurance strategy using PSO count for site page gathering.

In this paper, we propose an examination to determine features for tokens in URLs to test our past hypothesis, and certify our discernment. We portray features for tokens that perceive these two one of a kind classes of tokens, and we play out a quantifiable examination to test whether these two classes of tokens are discernable in each site. When we legitimize this reality, we can use these features to build the past URL models, and use them to arrange site pages. In examination with substitute sorts of features that can be used to mastermind site pages, token features can be processed without downloading the webpage page, just by researching its URL, which evades trivial downloads. Whatever is left of the article is sorted out as takes after. We propose new features and estimations for use in automated Web arrange errands, for instance, content recommendation and advancement hindering, that help customers adjust to the mass of information on the Web. A prominent method to manage such course of action errands is to use the broadly available library of substance and picture gathering instruments. Be that as it may, in this paper, we fight that two features particular to web reports—their URLs, and the plan of associations with them on a suggesting page—can be used extensively more enough for such request endeavors. We look for after the intuition that substance providers tend to pick URLs and page arranges that coherently structure their

**Anveshana's International Journal of Research in Engineering and Applied Sciences**
EMAILID:anveshanaindia@gmail.com,WEBSITE:www.anveshanaindia.com
105

Figure 1: Screen-shots from an original CNN page (left) and the same page viewed through the Daily You (right). Notice the Daily You's version removes the advertisements, some of the navigation boxes, and also writes the word "pick" (emphasized in picture) near recommended news articles.

content as showed by subject, and that such topical sorting out can be mishandled in portrayal assignments. For example, even with no appreciation of the substance of an every day paper, one may figure associations between articles in perspective of visual groupings alone (for instance, that articles under a comparative heading are about a comparable subject or were of similar essentialness). The goal of this paper is to formalize such senses into a general technique for algorithmically anticipating the properties of the destinations of unvisited joins. Two key walks in game plan are to pick the course of action of features that will be dissected and the decision conclude that will be associated with aggregate in light of those features. In some advancement blocking applications1 , the features join, for example, the estimations of the photo being viewed as and the decision standards ("an ad is a photo 250 by 100 pixels") are troublesomely hand-coded. A noteworthy disserve of such an approach is the prerequisite for human push to make the precepts and to create new ones as advancements progress. To settle this, systems, for instance, AdEater

try to apply machine adjusting, subsequently creating request administers by examining a course of action of named planning cases [13]. In proposition structures, since different decisions rules work for each customer, machine learning is frequently used. Frequently, the Web is managed as a gigantic substance corpus: the different features used are the words in the reports, and standard machine learning estimations, for instance, Naive Bayes or reinforce vector machines are associated [2].

In normal substance portrayal applications, similarity is evaluated by word cover—files are the same to the extent that they join near words (or articulations). In this paper, we receive a substitute procedure to closeness that weights the relative position of things in a tree:

- On various Web regions, page URLs are made in a chain out of significance as demonstrated by subject. For example, the present year's articles about space on the CNN Web site have a URL prefixed by cnn.com/2003/tech/space, which can be deciphered as circumstance in the "space" subtree of the "tech" subtree of the cnn tree. On the normal assumption that a customer is usually propelled by particular subjects anyway not others, the zone of an article in the URL tree is suggestive of the customer's eagerness for it. So additionally, on various Web goals sees every now and again bear joins showing back a single "commercial" subdirectory of the page. Without a doubt, business instruments, for instance, Web

**Anveshana's International Journal of Research in Engineering and Applied Sciences**
EMAILID:anveshanaindia@gmail.com,WEBSITE:www.anveshanaindia.com
106

Washer let customers physically decide certain URL "prefixes" as pointers of advancements that should be blocked.

- We find that Web-goals routinely base the visual plan of their document pages with respect to an issue logical arrangement. This organization is consistently dynamic and reflected in a recursive table plan that can be distinguished in the (different leveled) parse tree of the HTML record. For example, the CNN Web-site first page offers a "section by part manage" distributing stories under different imprints, for instance, "U.S.," "World," "Travel", and "Direction." These circumstances address subject groupings that may well be strong pointers of "captivating quality" for a peruser. Basically, plugs every now and again have a specific circumstance in the page design.

These two cases prescribe the probability of requesting a report in perspective of its circumstance in some past logical order (the class name itself isn't a bit of the logical characterization).

Most portrayal estimations deal basically with features that have been reduced to numbers, for instance, the amount of occasions of "apple" in a record or the length and width of a photo. To rather execute the likelihood of collection using such tree-based features, we need to deal with a couple of issues. To begin with, we need to develop a gathering model that empowers us to get ready and make conjectures using the tree-based segment. Second, we need to show how arrange

using that model can be executed gainfully.

## II. RELATED WORK

We have perceived a couple of suggestions in the region of page portrayal using features isolated from URLs. Some of them are controlled, presented one of the fundamental approaches to manage page subject portrayal using just URLs. Their recommendation included on tokenizing URLs into tokens, and a while later isolating features from those, like the words they contain, their sort or length, among others. These features are used to create a Maximum Entropy arrange show. They work with a subset of the URLs to fabricate a portrayal show, so they play out a diminished crawling to procure their arrangement set. Regardless, this is an overseen suggestion that requires a named set of planning URLs. Moreover, they go for requesting pages having a place with more than one site, and relies upon words being human-sensible, which isn't for the most part substantial.

Vidal et. al propose a methodology to dismember a singular site, and therefore find pages that resemble a case page that is given. To achieve in this manner, the site is mapped, and URL plans are made for those pages that lead (particularly or over the long haul) to those pages. To recognize likeness between pages, the Tree Edit Distance is used. To develop the arrangement set, they have to as of now crawl the entire site, download each page and after that strategy them, which takes a great deal of time. In like manner, it doesn't arrange a page according to its substance, yet to its fundamental likeness to the given page. That infers that pages containing information about different subjects may be named a comparable class. Zhu et. al. [16] propose an

**Anveshana's International Journal of Research in Engineering and Applied Sciences**
EMAILID:anveshanaindia@gmail.com,WEBSITE:www.anveshanaindia.com
107

association classifier, as opposed to a site page classifier, despite the way that we join it in this framework on account of their examination of association features. They go for gathering joins as showed by their ability inside the site. They propose a logical arrangement of predefined interface classes, dependent upon the limit that every association performs in the page, to be particular: investigating, requesting, refering to, endorsing and publicizing. They separate associations with independent visual, substance and essential features, among others, and they create two sorts of regulated classifiers: SVM and decision trees. It is controlled suggestion in which the classes are predefined in a rigid logical arrangement. Moreover, their goal isn't particularly related to information extraction, like our own, so they are not point arranged; rather, their classifier may be used for customer interface proposition. Baykan et. al. [3] proposed a classifier that resembles [12]. They tokenise URLs too to evacuate features, yet they apply unmistakable regulated portrayal estimations, as SVM, Naïve-Bayes or Maximum Entropy, and consider the results. These computations ought to be empowered an once-over of words trademark for each topic that will be requested. On a past work [4], the makers used a comparable idea, anyway this time remembering the true objective to arrange pages as showed by their tongue, as opposed to their subject. Much the same as [12], it is a coordinated framework that organizes pages from different districts. Finally, Blanco et. al. [6] consider that each site is made by populating HTML groups with data from a database. They will probably cluster site pages with the objective that each gathering contains pages following a particular arrangement. They watched that URLs made from a

comparable design have a relative case, much the same as pages delivered from a comparable configuration contain near terms, so they proposed a computation for unsupervised gathering that solidifies page substance and its URL as features, by strategies for the base depiction length procedure (MDL). They require a considerable getting ready set, so they sneak the entire site in their investigations. Additionally, to improve the portrayal capability, features from the page itself are fused into development to the association based features, which suggests that it must be downloaded as of now. As opposed to the past, our suggestion isn't coordinated, and it doesn't require to crawl the whole site to build the course of action appear, as in [15] or [6]. We use a little subset of pages and URLs from the site, and we apply a quantifiable framework to expel gathering features from those URLs.

### III. URL AND TABLE FEATURES

In this section we discuss in greater detail two tree-structured features that are particularly relevant to certain Web classification tasks.

### *3.1 URL trees*

The World Wide Web Consortium battles that report URLs should be dull (http://www.w3.org/Axioms.html#opaque) .

On this Web page, Tim Berners-Lee makes his proverb out of cloud URIs: "... you should not look at the substance of the URI string to increment other information...". Rather than those style rules, most URLs nowadays have human-masterminded suggestions that are important for proposition issues. No ifs ands or buts, the run's URL contain semantics including creation (w3.org), that the page is made in HTML, and that the subject relates to a "Saying in regards to

**Anveshana's International Journal of Research in Engineering and Applied Sciences**
EMAILID:anveshanaindia@gmail.com,WEBSITE:www.anveshanaindia.com
108

Opaqueness." As the file's URL delineates (to some degree startlingly), URLs are more than basically pointers: makers and editors dole out basic ramifications to URLs. They do this to make internal affiliation more direct (creation rights, record approvals, self-arrange), and now and again to make that affiliation plot clear to perusers. Perusers frequently make surmisings from URLs, which is the reason projects and web records generally demonstrate URLs close by the substance delineation of an association. We can accumulate from a URL that a report serves a particular limit (a specific Web inventory may constantly serve ads); or relates to a subject ('business' stories might be under one registry); or has a particular starting point. Or then again, we may delete an expansion of a URL attempting to move to a more "general" page still related to our starting stage. Essentially, near chronicles (as described by the site's designers) every now and again live under similar URLs. A fair URL structure gives obliging coherent signs to the peruser. URLs are extraordinarily awesome features for learning. To begin with, they are definitely not hard to remove and for the most part consistent. Each URL maps strangely to a file, and any fetchable record must have a URL. Then again, other Web features like hook content, alt marks, and picture sizes, are optional and not stand-out to a record. Clearly, URLs can be scrambled, concealed or changed in robotized shape; yet such changes in the meantime make it troublesome for customers and web files to find and return to information. Second, URLs can be examined without downloading the goal record, which allows us to perform gathering more quickly. This is a fundamental condition for realtime portrayal errands like advancement

blocking. Third, as we fight underneath, URLs have an instinctual and clear mapping to certain gathering issues. For example, we give trial confirmation in Section 4 that the URL is significantly related with whether an association is a notice or not. Most notice clicks are finished couple of undertakings; these activities are regularly contained in subtrees of the URL tree, like http://doubleclick.net or http://nytimes.com/adx/... . To change over a URL into a tree-shape, we tokenized the URL by the characters/, ? additionally, and. The/is a standard delimiter for lists that was continued into Web records; ? in addition, and are standard delimiters for passing elements into a substance. The farthest left thing (http:) transforms into the root center point of the tree. Dynamic tokens in the URL (i.e. nytimes.com) transform into the posterity of the past token. Note that our advancement guarantees we end up with a tree, paying little respect to whether the site itself isn't tree shaped (two pages may point to a comparative URL, anyway it is essentially the URL that portrays the tree zone).

### 3.2 HTML

Table Trees Similarly, the visual plan of a page is generally dealt with to empower a customer to perceive how to use a site. This outline tends to be organize—most pages will hold a 'look and feel' regardless of the way that the major substance might be dynamic. For example, unprecedented articles on one particular point may appear in a comparative place on the page for quite a while. The page configuration is ordinarily controlled by HTML table names, identifying with rectangular groupings of substance, pictures and associations. Routinely, one table at the edge or best of a page will contain a noteworthy piece of the site's course. The

**Anveshana's International Journal of Research in Engineering and Applied Sciences**
EMAILID:anveshanaindia@gmail.com,WEBSITE:www.anveshanaindia.com
109

substance of a site may use tables to gather together articles by noteworthiness (the component news fragment of a news-magazine), by subject, or successively (freshest things generally at the best). Like the URL, this page configuration can be used to wipe out certain substance, (for instance, the principles at the most astounding purpose of the page); or to revolve around other substance (the highlights, or the recreations fragment). Like the URL feature, tables make incredible features for machine learning. For a page to indicate truly in programs, the names need to conform to a standardized HTML sentence structure; this furthermore makes the table component easy to expel. In the accompanying portion, we give the instance of a Chinese Web site that might be seen despite with the exception of understanding the specifics of the substance on the page. To change over the HTML table structure into a tree-shape, we used a translated Perl program that removed the HTML table names
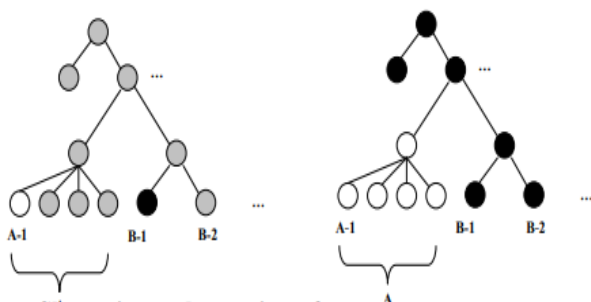


Figure 2: Shown is an abstraction of Web problem to the domain of "tree learning."
The top-left shows the original Website. The top-right shows that visual portions of the page are collected in chunks of HTML, which are indented to show the HTML's tree structure.

The base left shows the fascinated tree, getting a fragmentary naming of the page: white ("notice"), dull ("substance") and

decrease ("cloud") focus focuses. The base right shows one potential speculation of the tree which proposes everything in box A will be an advertisement.(and ). The base of the tree is the whole page's HTML. The successors of a middle point are the going with bring down level of table parts. Note that while an equivalent story may be incorporated in excess of one place on the page (e.g., an article may show up in both "world" and "getting ready"), with a definitive goal of this paper, we see those two appearances as two separate leaves on the legitimate request.

## IV. EXPERIMENTS AND RESULTS

We used 10 times 10-wrinkle cross-endorsement (unless for the most part communicated) to evaluate the test exactness. The examinations were continue running on a machine with 2 twofold focus 2 GHz Intel processors with 4 GB memory. To lead each one of the examinations, we used WEKA (Waikato Environment for Knowledge Analysis) data mining structure [14] with default parameter regards where legitimate. We take a gander at incorporate assurance and request frameworks using 3 generally used machine learning figurings in issues like our own: Naïve Bayes (NB) [17], Logistic Regression (LR) [28] and Random Forests (RF) [16]. We furthermore tried evaluating C4.5 [24] and Multilayer Perceptron, anyway the Wrapper incorporate decision technique was prohibitively slower setting aside quite a while for these slower classifiers.

### 4.1 Data Sets

For phishing site pages, we used attested phishing URLs from PhishTank [11]. PhishTank, worked by OpenDNS, is an aggregate clearing house for data and information about phishing on the Internet. A phish once submitted is checked by

**Anveshana's International Journal of Research in Engineering and Applied Sciences**
EMAILID:anveshanaindia@gmail.com,WEBSITE:www.anveshanaindia.com
110

different enrolled customers to attest it as phishing. We accumulated first course of action of phishing URLs from June 1 to October 31, 2010. Phishing procedures used by scalawags progress after some time. With a particular ultimate objective to look at these propelling techniques and to immovably copy our examinations as a general rule circumstance, we accumulated second gathering of avowed phishing URLs that were submitted for check from January 1 to May 3, 2011. We used substance [13] to subsequently recognize and expand the truncated URLs gave by online organization longurl.org. We accumulated our true blue site pages from two open data sources. One is the Yahoo! directory1 , the web interfaces in which are aimlessly given by Yahoo's server redirection advantage [10]. We used this organization to subjectively pick a URL and download its page substance close by server header information. To cover more broad URL structures and varieties in page substance, we moreover made an once-over of URLs of most typically phished targets. We by then downloaded those URLs, parsed the recuperated HTML pages, and gathered and crawled the hyperlinks in that to in like manner use as liberal site pages. We made the supposition, which we accept is sensible, to see those site pages as generous, since their URLs were removed from a true blue sources. These site pages were crawled between September 15 and October 31 of 2010. The other wellspring of true blue site pages is the DMOZ Open Directory Project 2 . DMOZ is a registry whose areas are audited physically by editors. In light of the date on which phishing URLs were submitted to PhishTank for check, we created two educational files. The principle enlightening file, we suggest it as DS1, contains 11,240 phishing webpage pages submitted before October 31, 2010 and 21,946 good 'ol fashioned site pages from Yahoo! additionally, seed URLs. The second instructive list, we imply it as DS2, contains 5,454 phishing pages submitted for affirmation between January 1 and May 3 of 2011 and 9,635 heedlessly picked true blue site pages from DMOZ. We discarded the URLs that were no more considerable as the page couldn't be gotten the opportunity to remove features from their substance.

### 4.2 Features

We start with a course of action of 177 features of which 38 are content-based and the rest are URL based. Content-based features are generally gotten from the particular (HTML) substance of site pages e.g., counting outside and internal associations, counting IFRAME marks, and checking whether IFRAME name's source URLs are accessible in blacklists and web seek instruments, checking for watchword field and testing how the casing data is transmitted to the servers (paying little heed to whether Transport Layer Security is used and whether ―GET‖ or ―POST‖ procedure is used to transmit shape data with mystery key field), et cetera. URL-based features fuse lexical properties of URLs, for instance, counting number of ―.‖, ―-―, ―_‖, et cetera in various parts of URLs, checking whether IP address is used and what kind of documentation is used to address the IP address set up of a territory name. URLs and zone some portion of the URLs are checked against top 3 web lists (Google, Yahoo, and Bing) records to check whether the URLs are recorded. Features in like manner fuse checking IPs and space name of the URLs against the best once-over of IPs and zones certainly surely understood for encouraging phishing and distinctive pernicious destinations.

**Anveshana's International Journal of Research in Engineering and Applied Sciences**
EMAILID:anveshanaindia@gmail.com,WEBSITE:www.anveshanaindia.com
111

Features also fuse a summary of eye-getting watchwords (e.g., log, click, pay, free, remunerate, bank, customer, et cetera.) that are more routinely used as a piece of phishing URLs to trick the end customers.

### 4.3 Feature Selection

Table 1 shows the demand correctnesses of Naïve Bayes, Logistic Regression and Random Forests classifiers with and without highlight choice utilizing CFS on DS1 edifying collection. Acquired pursue methodology understood a subset of 42 consolidates out of 177 highlights; anyway insatiable forward demand (Greedy FS) picked every single one of the highlights (happens as intended are not appeared as they are same as without include affirmation, grayed push). Natural ask for strategy overhauled Naïve Bayes classifier's outcomes the most with its blunder enhancing from 2.2% to 1.7% with the immense reducing in both FPR and FNR

.

## V. CONCLUSIONS AND FUTURE WORK

In this paper, we familiarize a trial with insistence that data contained in URL tokens is attractive to sort out those URLs, and to make classes of URLs. We have made highlights that experience the data contained in the tokens. These highlights are probabilities estimators, in light of token frequencies acquired from the examination. The part respects histograms demonstrate that we can see two sort of tokens, subordinate upon their segment respects, being some of them dependant on the demand made to get the URL, while others are independent from it, and they overall have a place with URLs of a near kind. As a determination, URL tokens in actuality contain enough data to gather

URL models, and in this manner, to mastermind pages, with the benefit of not downloading the page starting at now. We have perceived unmistakable recommendations in the composed work that go for utilizing the data contained in URLs to organize pages.

The benefits of our highlights in association with different recommendations are 1) client mediation is kept to a base, which spares a fundamental resource as is client time; 2) pages are accumulated for highlights that are outside them, which avoids downloading a page so as to sort out it; 3) it is tongue independent, since it depends upon the URL setup paying little respect to the specific words or groupings of characters that make every token; 4) it doesn't envision that affiliations will be consolidated by words gainful for depiction; and 5), we don't have to slither broadly a website to hoard a demand demonstrate that works legitimately, rather, we play out a lightweight creeping that recovers a little subset of pages. On account of the legitimate idea of the recommendation, we can ensure that the classifier is as right as it would be on the off chance that it had been constructed utilizing the entire course of action of pages. Later on, we hope to make a site page classifier utilizing these highlights. We give some cognizance about how to deliver such a classifier in [10]. In addition, we recognize such a classifier can be utilized to enhance web crawlers sufficiency, which we uncover in [9]. We ought to watch that we should investigate how to apply our highlights to the charged inviting URLs, which don't fit our speculation.

## REFERENCES

**Anveshana's International Journal of Research in Engineering and Applied Sciences**
EMAILID:anveshanaindia@gmail.com,WEBSITE:www.anveshanaindia.com
112

1. Arasu, A., Garcia-Molina, H.: Extracting structured data from web pages. In: SIGMOD, pp. 337–348 (2003)

2. Bar-Yossef, Z., Rajagopalan, S.: Template detection via data mining and its applications. In: WWW, pp. 580–591 (2002)

3. Baykan, E., Henzinger, M.R., Marian, L., Weber, I.: Purely URL-based topic classification. In: WWW, pp. 1109–1110 (2009)

4. Baykan, E., Henzinger, M.R., Weber, I.: Web page language identification based on URLs. PVLDB 1(1), 176–187 (2008)

5. Blanco, L., Crescenzi, V., Merialdo, P.: Structure and semantics of Data-IntensiveWeb pages: An experimental study on their relationships. J. UCS 14(11), 1877–1892 (2008)

6. Blanco, L., Dalvi, N., Machanavajjhala, A.: Highly efficient algorithms for structural clustering of large websites. In: WWW, pp. 437–446. ACM, New York (2011)

7. Cohen, W.W.: Improving a page classifier with anchor extraction and link analysis. In: NIPS, pp. 1481–1488 (2002)

8. F¨urnkranz, J.: Hyperlink ensembles: a case study in hypertext classification. Information Fusion 3(4), 299–312 (2002)

9. Hern´andez, I., Sleiman, H.A., Ruiz, D., Corchuelo, R.: A Conceptual Framework for Efficient Web Crawling in Virtual Integration Contexts. In: Gong, Z., Luo, X., Chen, J., Lei, J., Wang, F.L. (eds.) WISM 2011, Part II. LNCS, vol. 6988, pp. 282–291. Springer, Heidelberg (2011)

10. Hern´andez, I., Rivero, C.R., Ruiz, D., Corchuelo, R.: A Tool for Link-Based Web Page Classification. In: Lozano, J.A., G´amez, J.A., Moreno, J.A. (eds.) CAEPIA 2011. LNCS, vol. 7023, pp. 443–452. Springer, Heidelberg (2011)

11. Hotho, A., Maedche, A., Staab, S.: Ontology-based text document clustering. KI 16(4), 48–54 (2002)

12. Kan, M.-Y., Thi, H.O.N.: Fast webpage classification using URL features. In: CIKM, pp. 325–326 (2005)

13. Pierre, J.M.: On the automated classification of web sites. CoRR, cs.IR/0102002 (2001)

14. Selamat, A., Omatu, S.: Web page feature selection and classification using neural networks. Inf. Sci. 158, 69–88 (2004)

15. Vidal, M.L.A., da Silva, A.S., de Moura, E.S., Cavalcanti, J.M.B.: Structure-based crawling in the hidden web. J. UCS 14(11), 1857–1876 (2008).

**Anveshana's International Journal of Research in Engineering and Applied Sciences**
EMAILID:anveshanaindia@gmail.com,WEBSITE:www.anveshanaindia.com
113