# Malicious Web Page Detection through Classification Technique: A Survey

[1]Dr. Jitendra Agrawal, [2]Dr. Shikha Agrawal, [3]Anurag Awathe, [4]Dr. Sanjeev Sharma

[1,2,3,4]Rajiv Gandhi Proudyogiki Vishwavidyalaya, Bhopal, Madhya Pradesh, India

## Abstract

A "malicious web page" refers to a web page that contains malicious content that can exploit a client–side computer system. Malicious website may be used as a weapon by cybercriminal to exploit various security threats such as phishing, drive-by-download and spamming. Malicious Web sites are hurdle on the way ofInternet security.  And used as a weapon to mount various security threat like phishing, drive-by-download and spamming. To handle there is need to develop an automatic system to recognized malicious website. This paper gives a bird eye over malicious web site,their vulnerability and recent research to recognize it. In this work various machine learning and graph based technique to detect malicious website are presented. This paper also include various feature extraction technique such as information gain, N-gram, score gram and confidence weighted scheme to study nature of malicious website. The goal of this survey is to provide a comprehensive review of different classification techniques in data mining.

## Keywords

Malicious Web, Blacklisting, Phishing, Machine Learning Technique, Http Response Graph

## I. Introduction

In today's world Internet has become an integral part of one's life. Internet has grown grow from its scholastic origins to become the essential global tool. Internet provides necessary infrastructure that plays a crucial role in various domain for example: communication, finance, commerce, E-governance and education. Some of the prominent application of internet is Hybrid cloud, BYOD, Big Data, E-shopping.

Unfortunately, on the other side, crimes over the Internet (cybercrime) have increased at much faster rate and with high complexity.Cybercrime attacks includes online frauds, cracking into the system, phishing attacks, DNS poisoning, malware attacks, data theft, spamming, scams, blackmailing .A recent example of cybercrime is Sony Cyber Hack [15] in which the computer systems at the corporate network of Sony Pictures were breached and taken offline by a malware-based attack.

The research and development in the field of web security includes network traffic analysis, malware detection and control system, BOT detection, filtering system detection growth, malicious web detection system.

Detection of malicious website is one of the hot research topics in the field of security. The importance of this can be estimated from this statistics which shows that Google finds 9,500 new malicious Web sites a day [10] .This is also important because it prevents the user from being compromised from attacks such as Phishing attack, Drive-by–download attacks, spam attack, Click-jacking, Plug-In and Script-Enabled Attacks, Mal-advertising [5]..

Malicious web detection is defined as the process of identifying those web URLs and web pages which can lead to compromise in user security and affect the users. Similarly malicious web constitute web infrastructure andservices used as intended to affect the users. The different methods used in malicious web are socially engineered theft of passwords or credentials, hide infiltrations, and exploitation of the trust required for economic transactions, government services and social interactions.

These following section elaborate different approaches followed till now. These approaches are based on different types of features of the web pages. The feature selection describes about different features used in malicious web detection process.  Section fourth describe different classification technique section explains about different type of classification methods such as Supervised machine learning, unsupervised machine learning and graph techniques used in detection of malicious web.

## II. Related Work

The detection of malicious web started with introduction of blacklisted URLs. Blacklisted URLs[2] are the recoded list of web address which are earlier analyzed and declared as malicious one. These blacklisted URLs act as repository for cross verification for user request. There is no false positive case in blacklisting URLs. This method of web security has certain limitation. First as the web environment is changing at rapid pace this list need to be updated regularly. Second the attacker can breach this method by switching to new domain. Third this method unable to detect one-day delay malicious web URLs.

To overcome the limitation of blacklisted URLs method recent research focused on the featured based malicious web detection. This malicious URL detection uses various web based features. These include:
1.   URL Lexical features
2.   Domain host based features
3.   Web page content based features
4.   HTTP response graph features

Recent research build statistical models based on above featuresfor classifying URLs into malicious and benign one.

Y. Shin et al[1]classified URLs as spam or benign by taking into account the link structure of the hyperlink graph. It relied on the graph metrics of the linked URLs graph. This classification is based on support vector machine learning technique for spam URL classification but quite limited in using the links from only one blog.

J. Ma et al[2]presented an approach automated for URL classification using lightweight statistical methods based on  most predictive tell-tale lexical and host based features. This approcah is quite limited to to handle millions of URLs whose features evolve over time.

M. Cova et al[3]presented a novel approach for detection and analysis ofmalicious JavaScript code. This approach combines anomaly detection techniques and dynamic emulation of the code. The presented method focus only on one type of attack vector i.e. Drive-by-download.

Aaron Blum et al[7]explored confidence weighted classification for classificatio of phishing and benign web using content based

features and produce a dynamic and extensible detection system for present and emerging types ofphishing domains.

## III. Approach

The approach for solving the malicious web detection problem was started by introduction of blacklisting web site.But it has certain limitation in current times. So to eradicate the limitations of blacklisting web there is a need to develop an automated system which can easily identify the malicious web and protect user from being affected. These approaches are based on web featurein real time constraints. These analysis approaches are:

### A. Heuristics Based Approaches

It is based on signature of the known attacks. Attacks payloads are applied in Intrusion Detection System [11, 12, 16] to check the web pages. The detection rate of this technique depend on the rate at which signature are developed for new types of intrusions.
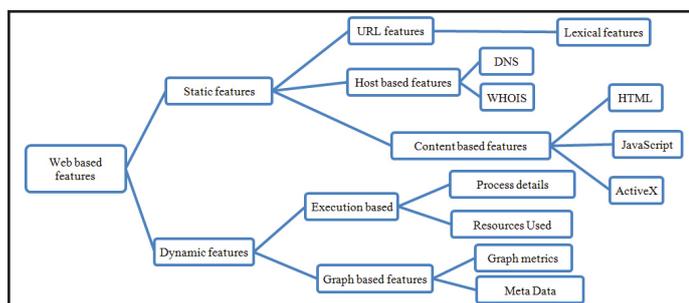


Fig. 1: Different Types of Web Based Features

### B. Static Analysis Approach

Static analysis technique is based on the non-changing characteristics of web pages. Static analysis classifies sites based only on the relationship between URLs and content of the web page. It includes extraction of features from URL, source code features (HTML and JavaScript), host based features (WHOIS and DNS), Linked graph interaction features (HTTP response)

### 1. URL Lexical Feature

Lexical feature is the textual properties of the URL itself. This feature[2, 4] includes the length of the hostname, the length of the entire URL, the number of dots in the URL and token in the hostname (delimited by '.'), in the path URL (strings delimited by '/', '?', '.', '=', '-' and '_'). All these features are integer valued and their values range helps in identifying themalicious web URLs. For example: In Phishing attack the malicious URL has mimic appearance to the benign web page and only differs by these bag-of-words.

1.  http://83.16.123.18/www.paypal.com/update.htm?=
2.  http://signin.paypal.com@10.19.32.3/
3.  www.pay.pal.com

The above Phishing URL appears similar to that of benign URL "https://www.paypal.com/home".
All three URLs above contain "Pay", "Pal", "com" string which makes the malicious URL look alike to benign one and user can easily be trapped by the attacker.

### 2. Host Based Feature:

This feature is based on the WHOIS and DNS record of the URL. Host-based features[2, 4, 7] describe "where" malicious sites are hosted, "who" they are managed by, and "how" they are administered.

Host based features include:
- IP address properties: It is used to identify whether the IPs of the A, MX, or NS records belong to same Autonomous Systems or not. If the hosting malicious URL belongs to a specific IP prefix or AS belonging to an Internet Service Provider (ISP), then it will be classified as disreputable ISP.
- WHOIS properties: It contains details about what is the date of registration, update, expiration, who is the registrar, who is the registrant. If WHOIS features of malicious domains belong to some other URLs then it would be more chances that they are also malicious one. Moreover, if web sites are taken down frequently, then there is possibility it may be malicious.
- Geographic properties: This determines country, city a particular IP address belongs. As with IP address properties, probability for malicious activity in specific geographic regions.
- Connection speed:The connection speed also important factor for determining malicious. Some malicious sites tend to reside on compromised machines (connected via cable or DSL).
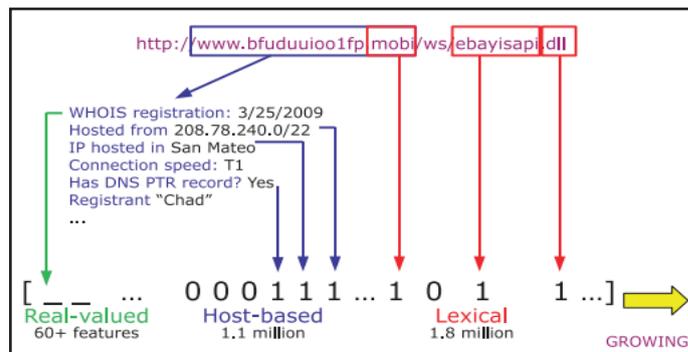


Fig. 2: WHOIS Records

### 3. Content Based Feature

Source code of the web page is important source of feature for identifying malicious web pages. The contents based feature of web pages includesnumber of functions, objects, I-Frames, data streams, hyperlinks, In-bound, Out-bound, #f hidden elements. These features are based on the vulnerabilities and loopholes in the language used for scripting purpose. The scripting languages include the DHTML, CSS, JavaScript, and ActiveX.Content based features [3] are as:
- **DHTML or HTML features:** DHTML or HTML feature includes word count, Average word length, distinct word count, and size of I-frame. These features are exploited by the attacker for obfuscating the malicious code or script into the web page.
- **JavaScript features:** JavaScript is one of the popular scripting languages which is used as the validation code for the web page. Some functions in the JavaScript are quite vulnerable and are exploited by the attacker for injecting the malicious script. These malicious scripts get executed just by clicking on the page. Some of JavaScript functions which are most exploited are eval(), escape(), unescape(), exec(), ubound() etc
- **ActiveX:** ActiveX provide high user end facility. ActiveX used list of objects for various facilities for example: Scripting->FileSystem Object which has the ability of I/O to file system, WScript.Shell can execute shell commands on the client's computer, Adodb.Stream is used to download files from the Internet.

## 4. Graph Based Features

Linked graph based feature includes the structural properties of the hyperlinks. Linked graph helps in determining how redirection path are followed to reach the target URL. The link graph metrics [1] include:

### (i). Degree of the node

Degree is the measure of number of links or hosts are connected to that particular node. It includes the number of links directed to the URL i.e. In-degree and the number of links directed out of the URL i.e. Out-degree.

$D(x_i) = IN(x_i) + OUT(x_i)$

Where $IN(x\_i)$ is the In-degree and $OUT(x\_i)$ is the Out-degree of the node $x_i$

### (ii). Centrality

Centrality of a node within a graphdenotes how many shortest paths in the graph traverse through a given node. A node with high centrality denotes a hub or central node in the graph.

### (iii). Clustering Coefficient

It measures the degree to which all neighbors of a given node in a graph tend to interconnect. This metric show the degree of local interconnectivity among URLs.The clustering coefficient [18] for the whole graph is the average of the local values $C_i$

$$C = \frac{1}{n}\sum_{i=1}^{n} C_i$$

Where n is the number of nodes in the network and $0<C\_i\leq1$ and $0<C\leq1$

### (iv). Metadata of the Link Graphs

Metadata feature includes the number of nodes, edges, domains, and hosts; the Hypertext Transfer Protocol (HTTP) status codes for each URL; and type of edges.

## C. Dynamic Analysis Approach

Dynamic analysis technique inspects features based on the execution the web page in the controlled environment for example: Honey-pot.Honey-pot are the real machine which mimic the characteristic of multiple operating system at the same time. These are designed to lure the attacker without knowing that they are actually monitored and decoy the purpose of identifying new attack tactics. Honey-pots are divided into two types based on their interaction level i.e. Low interaction honey-clients[12] for example HoneyC and High level honey-clients for example Capture-HPC, MITRE Honey-client.

The dynamic approaches attempt to capture unusual behaviors of the malicious web such as launched processes, registry changes, memory heap allocation, list of open ports in the victim system. This analysis technique is useful for identifying those attacks that require wait for time, logic bombs, and user interactions. Dynamic analysis has high detection rate and require resource intensive environment.

## D. Hybrid Analysis Approach

Hybrid Analysis technique makes use of the above two analysis approach to remove the limitations of correctly identifying and time constraints. Hybrid detection method [14] has a two hierarchy steps. First well-known malicious web pages can be detected in the misuse detection step. Second unknown malicious web pages, which are obtained from the suspicious web pages that

were identified in the first detection component, can be detected in the anomaly detection step.

Table 1: Type of Features in Relation to Type of Attack

| No | Types of attacks | Features | Ref. |
|---|---|---|---|
| 1 | Phishing | URL lexical features | 5,11,12,13 |
| 2 | Drive by download | Content based features | 3,11,12 |
| 3 | Click jacking | Http response | 1 |
| 4 | Spamming | Content based and execution based | 1 |
| 5 | Mal-advertising | Graph and content based | 12 |
| 6 | Session High jacking | Execution based features | 1 |

Above table give the relationship between types of attacks to various types of features. Identification of the most appropriate set of features help in efficiently distinguishing web pages and URLsinto malicious benign. Identifying most appropriate feature is done by applying different feature selection technique.

## IV. Feature Selection Technique

Features selection technique is the mathematical models which are used to identify features that are needed to differentiate between malicious and benign web pages. The quality of classification technique depends on the feature selection technique.

### A. Entropy

Entropycalculates the variation or measure of impurity. High value of entropy shows the uniform distribution of the featured value and show the normal behavior, on the contrary low entropy shows the varied distribution of the featured value and helps in determining the malicious behavior based on those features. Entropy [6, 9]can be represented as:

$$H(x) = -\sum_{i=1}^{n} p(x_i)log_b p(x_i)$$

where $x_i$ is the $i^{th}$ attribute.

### B. Information Gain

Information gain (IG) measures the amount of information about the class prediction. This measurement helps in determining whether a new instance should be classified malicious or benign one.

The greater information gain of an attribute, the higher value it contributes to the process to identify malicious web pages.

Information gain [9, 17] is a measurement method to choose high valuable features only. Information gain for an attribute 'a' is defined as follows:

$$IG(S, a) = Entropy(S) - \sum_{v\in a} \frac{|S_v|}{|S|} * Entropy(S_v)$$

Where S is collection of instances, $S_v$ is a subset of S with relevant value v of attribute 'a'.

### C. Confidence Weighted Algorithm

Confidence-weighted method for URL classification is based on linear classification by utilizing individual confidence weighted values to improve the flexibility and accuracy of classification. The relationship of each parameter's confidence weight factor helps in

automatically classifying into malicious and benign class. Further, this allows the model to automatically adjust when parameters change in significance, as often occurs throughout the course of a phishing campaign [7, 17].

It uses mean μ and standard deviation σrepresentative of the class and confidence for each feature.The class of a new data member represented by a feature vector x is determined by computing the sign of w • x,

Where w ~N (μ, Σ) and Σ represents the covariance matrix with a diagonal of σ, and zero for all off-diagonal elements.

### D. Scoring Mechanism

Scoring mechanism is a filtering technique which classifies suspicious web pages into classes: benign web pages and potential malicious web pages. Scoring algorithm[11]is based on the concept ofstandard score which measure how many standarddeviations a value of observed attribute is far from the mean (Carroll and Carroll 2002). It provide an balance between accuracy and performance to classification techniques.Scoring mechanism has three types of scores based on contents of web pages:

- Exploit content score.
- Script content score
- Foreign content score

A group score of instance x is calculated as follows :

$$GS_{g \in G}(x) = \sum_{a \in g} \frac{|x_a - \mu_a|}{\delta_a}$$

Where g is an attribute group which can be foreign content group, script content group or exploit content group; a is an attribute of g; $x_a$ is value of attribute a of instance x; $\delta_a$ is a standard deviation of attribute a which is estimated during training a set of benign instances; $\mu_a$ is mean of attribute a which is estimated during training a set of benign instances.

The greater score an instance x has in eachtype; themore likely it is classified as potential malicious class. Threshold value $T_g$ for each content type"g"is calculated toidentify potential malicious instances. The rule of classification is as follows:

$$x = \begin{cases} Potentialymaliciousif \exists_g \in G : GS_g(x) > T_g \\ Otherwise, \quad xisbenign \end{cases}$$

A web page will be classified as potential maliciousone thathas score greater than the threshold value for thattype.

### E. N-Gram Technique

N-gram[5] is a contiguous series of n items from a given sequence of text. It is a probabilistic method of predicting the next item in a sequence. N-gram technique can be used to extract the features without knowledge about a DHTML webpage. The n-gram model is used in statistical web content based processing.

## V. Classification Techniques

### A. Machine Learning

#### 1. Naïve Bayes

Naïve Bayes classification is a supervised machine learning technique. This classification technique is based on the condition probabilistic approach. It is assumed that all the features are independent from each other.The parameters in the Naive Bayes[2, 7] classifier are estimated to maximize the joint log-likelihood of URL features and labels, as opposed to the accuracy of classification.

Let X = {$x_1, x_2, ....x_n$} g be the features of web pages. x and y be the class of classification. The posterior probability of class y given x is calculated as:

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} = \frac{P(x|y)P(y)}{\sum_j P(x|y_j)P(y_j)}$$

Where

$$P(x|y) = \prod_{i=1}^{n} P(x_i|y)$$

The probability of each class $P(y_j)$ and the probability of each feature given a class $P(x_i|y)$ are obtained from the frequencies of the training dataset.

### 2. Decision Tree

Decision tree classification technique (Quinlan, 1992) is a classifier that models the data into a rooted tree. Each internal node of the tree corresponds to a feature of the web pages and edges from the node separate the data based on the value of the feature. By tracking down the nodes from the root of the tree based on the feature values of an example, we can predict class for that.

The feature selection for each node is based on the information gain of the features on the training data. The posterior probability of a class for an item can be calculated based on the number of training data of each class on the leaf nodes.

### 3. Boosted Decision Tree

Boosted decision tree classification is the improved version of decision tree. This classification technique is combination of multiple classifiers. The final prediction of an example is calculated as a weighted sum of the prediction of all the classifiers.During training process, small decision trees are created, and weights are assigned to the decision trees by minimizing the exponential loss of the boosted classifier on the training data. The posterior probability of a class given a web page is calculated as weighted sum of the posterior probability of the decision trees with normalization.

### 4. Support Vector Machine

A support vector machine or SVM [8] is a supervised method of machine learning for classification. Ittries to determine a decision boundary that separates the web pages into malicious and benign class.SVM maximizes the distance of the hyper plane and the resulting decision boundaries are robust to slight change of the feature vectors.

When input data is not separable by the linear function, SVMs use a kernel function to map the data into a higher dimensional space, and separate the data on the mapped dimension.

### 5. Logistic Regression [2]

This is a simple parametric model for binary classification where examples are classified based on their distance from a hyper plane decision boundary.

### B. Graph Analysis

Graph analysis is a classification technique which uses graph properties to determine malicious web site. Linkgraphs uses

differing depths and aggregating sub-graphs of the link graphs for determine the Spam type of malicious nature of web page. Support vector machine (SVMlight) technique[1] is used in this approach. SVM-light classifier useshigh-dimensional graph and metadata metrics features for train classifier and uses linear kernels for decision boundary. In order to decide what values of the SVM output correspond to spam (class label +1) or legitimate URLs (class label −1), one needs to apply a decision threshold, which is then used to control the fraction of false positive versus false negative predictions.

## VI. Conclusion

Detection of malicious web has become a necessary and hot topic of research as numbers of internet users are increasing at a high pace. There are lots of challenges regarding this detection process. First the number of online URL is very large. Second web environment uses diverse platform and difficult to find security solution for them. Third now threats are become more and more complex and used various obfuscation techniques to bypass detection techniques. The existing detection techniques are focused only on single type of attacks only. New generated malicious web pages exploit multiple types of attacks for targeting the client. Cloaking type of attacks is difficult to detect because these web respond differently to browser and crawler. Size of web is a big challenge in the process.

## References

[1] Y. Shin, S. Myers, M. Gupta, P. Radivojac, "A link graph-based approach to identify forum spam," Security and Communication Networks, Vol. 2, pp. 71–81, 2009.

[2] J. Ma, L. K. Saul, S. Savage, G. M. Voelker, "Beyond Blacklists : Learning to Detect Malicious Web Sites from Suspicious URLs," In proceedingsof 15th ACM SIGKDD international conference on Knowledge Discovery and Data mining, 2009. pp. 1245-1254.

[3] M. Cova, C. Kruegel, G. Vigna, "Detection and analysis of drive-by-download attacks and malicious JavaScript code," In Proceedings of the 19th International Conference on World wide web - WWW '10, pp. 281-290, 2010.

[4] C. Curtsinger, L. Benjamin, G. Z. Benjamin, C. Seifert, "ZOZZLE: Fast and Precise In-Browser JavaScript Malware Detection," In proceedings of 20th USENIX conference on Security Symonposium, pp. 33-48, 2011.

[5] D. A. T. Holz, C. F. Felix, "ADSandbox: Sandboxing JavaScript to fight malicious websites," In Proceedings of the ACM Symposium on Applied Computing, pp. 1859-1864, 2010.

[6] Aaron Blum, B. Wardman, T. Solorio, G. Warner, "Lexical feature based phishing URL detection using online learning," In Proceedings of 3rd ACM Workshop on Artificial Intelligence and Security - AISec '10, pp. 54-60, 2010

[7] Y. T. Hou, Y. Chang, T. Chen, C. S. Laih, C. M. Chen, "Malicious web content detection by machine learning", published in journal Expert Systems with Applications, Vol. 37, Issue 1, pp. 55–60, Jan. 2010.

[8] T. M. Qassrawi, Z. Hongli, "Detecting malicious web servers with honeyclients.," Journal of Networks, Vol. 6, No. 1, pp. 145-152, 2011.

[9] J. Ma, K. S. Lawrence, S. Stefan, M. V. Geoffrey, "Identifying suspicious URLs: an application of large-scale online learning," In Proceedings of the 26th Annual International Conference on Machine Learning, ACM, pp. 681-688, 2009.

[10] K. Alexandros, M. Cova, C. Kruegel, G. Vigna, "Escape from monkey island: Evading high-interaction honeyclients.," In Detection of Intrusions and Malware, and Vulnerability Assessment, Springer Berlin Heidelberg, pp. 124-143, 2011

[11] V. L. Le, I. Welch, X. Gao, P. Komisarczuk, "Identification of potential malicious web pages," In proceedings of ninth australasian information security Conference-2011, Vol. 116, pp. 33-40.

[12] B. Eshete, "Effective analysis, characterization, and detection of malicious web pages," in Proceedings of the 22nd International Conference on World Wide Web companion, 2013, pp. 355–360.

[13] M. K. Wanawe, M. S. Awasare, N. V Puri, "An Efficient Approach to Detecting Phishing A Web Using K-Means and Naïve-Bayes Algorithms," Vol. 2, No. 3, pp. 106–111, 2014.

[14] S. Yoo, S. Kim, "Two-Phase Malicious Web Page Detection Scheme Using Misuse and Anomaly Detection," International Journal of Reliable Information and assurance, Vol. 2, No. 1, pp. 1–9, 2014.

[15] R. B. Basnet, H. S. Andrew, L. Quingzhong, "Feature selection for improved phishing detection," In Advanced Research in Applied Artificial Intelligence, Springer Berlin Heidelberg, pp. 252-261, 2012

[16] C. Seifert, I. Welch, P. Komisarczu, "Identification of malicious web pages with static heuristics", In Proceedings of the Australasian Telecommunication Networks and Applications Conference, 2008.

[17] J. Ma, L. K. Saul, S. Savage, G. M., "Voelker.Identifying suspicious urls: An application of large-scale online learning", In Proceedings of the 26th International Conference on Machine Learning (ICML-2009), pp. 681–688, Montreal, Quebec, Canada, 2009.

[18] P. Zhao, H. C. Steven, "Cost-sensitive online active learning with application to malicious URL detection," In Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining, ACM, pp. 919-927, 2013.

Dr. Jitendra Agrawal was born in 1974, is an Assistant Professor in the Department of Computer Science & Engineering at the Rajiv Gandhi Proudyogiki Vishwavidyalaya, MP, India. He earned his Master Degree from Samrat Ashok Technology Institute, Vidisha (M.P.) in 1997 and awarded Doctor of Philosophy in Computer & Information Technology in 2012. He has published more than 50 publications in International Journals and Conferences. He has published two books named Data Structures and Advanced Database Management System. He is the recipient of the Best Professor in Information Technology award by the World Education Congress in 2013. He is a senior member of the IEEE (USA), Life member of Computer Society of India (CSI), Life member of Indian Society of Technical Education (ISTE).

Dr Shikha Agrawal is an Assistant Professor in Department of Computer Science & Engineering at University Institute of Technology, Rajiv Gandhi Proudyogiki Vishwavidyalaya, Bhopal (MP) India. She obtained B.E., M.Tech. and Ph.D in Computer Science & Engineering from Rajiv Gandhi Proudyogiki Vishwavidalaya Bhopal. She has more than twelve years of teaching experience. She has been awarded as "Young Sciencetist" by Madhya Pradesh Council of Science and Technology, Bhopal in 2012. Her other extraordinary achievements include "ICT Rising Star of the Year Award 2015" and "Young ICON Award 2015". Her area of interest is Artificial Intelligence, Soft Computing and Particle Swarm Optimization and Database. She has published more than 30 research papers in different reputed international journals and 9 chapters. She is a Madhya Pradesh State Students' Coordinator, Computer Society of India (CSI), 2016-2017.

ANURAG AWATHE He received Bachelor of Engineering Degree in INFORMATION TECHNOLOGY From INDORE INSTITUTE OF SCIENCE AND TECHNOLOGY INDORE. He is currently pursuing the M.Tech at School of Information technology, RGPV Bhopal (Madhya Pradesh). His research interest MALICIOUS WEB PAGE DETECTION.

Dr. Sanjeev Sharma was born in 1970, working as Associate Professor and Head of School of Information Technology at Rajiv Gandhi Proudyogiki Vishwavidyalaya Bhopal India. Dr Sharma is receipt of Best Teacher Award in Information Technology awarded by World Education Congress. He has graduated in Electrical & Electronics from Samrat Ashok Technical Institute, India and post graduated in Microwave and Millimeter from Maulana Azad College of Technology, India. He completed his Doctorate in Information Technology from Rajiv Gandhi Proudyogiki Vishwavidyalaya. He possesses teaching and research experience of more than 25 years. His areas of interest are Mobile Computing, Adhoc Network, Data Mining, Image processing and Information Security. He has edited proceedings of several national and international conferences and published more than 120 research papers in reputed journals. He is member of various professional bodies like IEEE, CSI, ISTE.