

ISSN No.-2321-8711

International Journal of Software Engineering and Systems

Volume No. 11

Issue No. 3

September - December 2023



ENRICHED PUBLICATIONS PVT. LTD

**S-9, IInd FLOOR, MLU POCKET,
MANISH ABHINAV PLAZA-II, ABOVE FEDERAL BANK,
PLOT NO-5, SECTOR-5, DWARKA, NEW DELHI, INDIA-110075,
PHONE: - + (91)-(11)-47026006**

International Journal of Software Engineering and Systems

Aims and Scope

Software Engineering has become very important with the ever-increasing demands of the software development to serve the millions of applications across various disciplines. For large software projects, innovative software development approaches are vital importance. In order to gain higher software standards and efficiency, software process adaptation must be derived from social behavior, planning, strategy, intelligent computing, etc., based on various factors. International journals of software engineering address the state of the art of all aspects of software engineering, highlighting the all tools and techniques for the software development process. The journals aims to facilitate and support research related to software engineering technology and the applications. International journals of software engineering welcomes the original research paper, review papers, experimental investigation , surveys and notes in all areas relating to software engineering and its applications. The following list of sample-topics its by no mean to be understood as restricting contributions to the topics mentioned:

Ø Aspect-oriented software development for secure software

Ø Dependable systems

Ø Experience related to secure software system

Ø Global security system

Ø Maintenance and evolution of security properties

Ø Metrics and measurement of security properties

Ø Process of building secure software

Managing Editor
Mr. Amit Prasad

Editorial Board Member

Dr. Pradeep Tomar
School of Information and
Communication Technology,
Gautam Buddha University,
Greater Noida, U.P. INDIA

Dr. O. P. Sangwan
School of Information and
Communication Technology,
Gautam Buddha University,
Greater Noida, U.P. INDIA

Dr. Nasib S. Gill
Department of Computer Science
& Applications, Maharshi Dayanand
University, Rohtak, Haryana, INDIA

Dr. Anurag Singh Baghel
School of Information and
Communication Technology,
Gautam Buddha University,
Greater Noida, U.P. INDIA

Dr. Sanjay Jasola
Graphic Era Hill University,
Dheradhun, Uttrakhand, INDIA

Dr. Bal Kishan
Department of Computer Science
& Applications, Maharshi
Dayanand University,
Rohtak, Haryana, INDIA

Dr. Ela Kumar
School of Information and Communication
Technology, Gautam Buddha University,
Greater Noida, U.P. INDIA

Dr. Sunil Sikka
Department of Computer Science
& Applications, Maharshi
Dayanand University,
Rohtak, Haryana, INDIA

Dr. Rakesh Kumar
Department of Computer Science
Kurukshetra University
Kurukshetra, Haryana, INDIA

Dr. Vijay Kumar
Department of Computer Science
& Engineering and IT, Kautiliya
Institute of Technology and
Engineering, Sitapura, Jaipur,
Rajasthan, INDIA

Dr. Kamal Nayan Aggarwal
Howard University, Howard, USA

Dr. Gurdev Singh
Samsung India Software Center,
Noida, U.P., INDIA

Dr. Dinesh Sharma
University of Maryland, Eastern Shore,
Princess Anne, MD, USA

Dr. Kapil Sharma
Department of Computer Science and
Engineering Delhi Technological
University, New Delhi, INDIA

International Journal of Software Engineering and Systems

(Volume No. 11, Issue No. 3, September - December 2023)

Contents

Sr. No.	Articles / Authors Name	Pg. No.
1	Rean (Robust Efficient Adaptive Network) for Manets: Feasibility and Analysis <i>- Prachi Garg, Shruti Garg</i>	103 - 108
2	Analytical Study on Applications of Web Search using Stochastic Query Covering <i>- Tiruveedula Gopikrishna</i>	109 - 118
3	Detecting Malware by Data Mining <i>- Shafiqul Abidin, Rajeev Kumar, Varun Tiwari</i>	119 - 124
4	Digits in Units and Tens Places of 3-PrimeFactors Numbers till 1 Trillion <i>- Neeraj Anant Pande</i>	125 - 144
5	Skills Required for Web Developer <i>- Dr. Sapna Nagpal</i>	145 - 149

Rean (Robust Efficient Adaptive Network) for Manets: Feasibility and Analysis

Prachi Garg*, Shruti Garg*

*Assistant Professor Dept of Computer Science, Geeta Institute of Technology & Management, Kanipla

ABSTRACT

A mobile ad hoc network (MANET) is a spontaneous network that can be established with no fixed infrastructure. This means that all its nodes behave as routers and take part in its discovery and maintenance of routes to other nodes in the network i.e. nodes within each other's radio range communicate directly via wireless links, while those that are further apart use other nodes as relays. Ad hoc networks have a wide array of military and commercial applications. They are ideal in situations where installing an infrastructure network is not possible or when the purpose of the network is too transient or even for the reason that the previous infrastructure network was destroyed. Security issues are the critical part in such type of networks. DDoS attack is the one of them in MANET. The proposed solution reduces the effect of DDoS attack and the designed network is robust and efficient enough to resolve all the problems related to DDoS. Detection and Prevention of DDoS are two main issues to be tackled.

Keywords – MANET, DDoS.

I. INTRODUCTION

In view of the increasing demand for wireless information and data services, providing faster and reliable mobile access is becoming an important concern. Nowadays, not only mobile phones, but also laptops and PDAs are used by people in their professional and private lives. These devices are used separately for the most part that is their applications do not interact. Sometimes, however, a group of mobile devices form a spontaneous, temporary network as they approach each other. This allows e.g. participants at a meeting to share documents, presentations and other useful information. This kind of spontaneous, temporary network referred to as mobile ad hoc networks (MANETs) sometimes just called ad hoc networks or multi-hop wireless networks, and are expected to play an important role in our daily lives in near future.

A mobile ad hoc network (MANET) is a spontaneous network that can be established with no fixed infrastructure. This means that all its nodes behave as routers and take part in its discovery and maintenance of routes to other nodes in the network i.e. nodes within each other's radio range communicate directly via wireless links, while those that are further apart use other nodes as relays. Its routing protocol has to be able to cope with the new challenges that a MANET creates such as nodes mobility, security maintenance, quality of service, limited bandwidth and limited power supply. These challenges set new demands on MANET routing protocols.

Ad hoc networks have a wide array of military and commercial applications. They are ideal in situations where installing an infrastructure network is not possible or when the purpose of the network is too transient or even for the reason that the previous infrastructure network was destroyed.

Security in mobile ad hoc networks is a hard to achieve due to dynamically changing and fully decentralized topology as well as the vulnerabilities and limitations of wireless data transmissions.

Existing solutions that are applied in wired networks can be used to obtain a certain level of security. Nonetheless, these solutions are not always being suitable to wireless networks. Therefore ad hoc networks have their own vulnerabilities that cannot be always tackled by these wired network security solutions.

One of the very distinct characteristics of MANETs is that all participating nodes have to be involved in the routing process. Traditional routing protocols designed for infrastructure networks cannot be applied in ad hoc networks, thus ad hoc routing protocols were designed to satisfy the needs of infrastructure less networks. Due to the different characteristics of wired and wireless media the task of providing seamless environments for wired and wireless networks is very complicated. One of the major factors is that the wireless medium is inherently less secure than their wired counterpart. Most traditional applications do not provide user level security schemes based on the fact that physical network wiring provides some level of security. The routing protocol sets the upper limit to security in any packet network. If routing can be misdirected, the entire network can be paralyzed. This problem is enlarged in ad hoc networks since routing usually needs to rely on the trustworthiness of all nodes that are participating in the routing process. An additional difficulty is that it is hard to distinguish compromised nodes from nodes that are suffering from broken links.

Recent wireless research indicates that the wireless MANET presents a larger security problem than conventional wired and wireless networks[3]. Distributed Denial of Service (DDoS) attacks has also become a problem for users of computer systems connected to the Internet. A DDoS attack is a distributed, large-scale attempt by malicious users to flood the victim network with an enormous number of packets. This exhausts the victim network of resources such as bandwidth, computing power, etc. The victim is unable to provide services to its legitimate clients and network performance is greatly deteriorated.

II. DISTRIBUTED DENIAL OF SERVICE (DDOS) ATTACK

DoS Attack

A denial of service (DoS) attack is characterized by an explicit attempt by an attacker to prevent legitimate users of a service from using the desired resources [6]. Examples of denial of service attacks include:

- attempts to “flood” a network, thereby preventing legitimate network traffic
- attempts to disrupt connections between two machines, thereby preventing access to a service
- attempts to prevent a particular individual from accessing a service
- attempts to disrupt service to a specific system or person.

DDoS Attack

A DDoS (Distributed Denial-Of-Service) attack is a distributed, large-scale attempt by malicious users to flood the victim network with an enormous number of packets [4]. This exhausts the victim network of resources such as bandwidth, computing power, etc. The victim is unable to provide services to its legitimate clients and network performance is greatly deteriorated. The distributed format adds the “many to one” dimension that makes these attacks more difficult to prevent. A distributed denial of service attack is composed of four elements. First, it involves a victim, i.e., the target host that has been chosen to receive the brunt of the attack. Second, it involves the presence of the attack daemon agents. These are agent programs that actually conduct the attack on the target victim. Attack daemons are usually deployed in host computers. These daemons affect both the target and the host computers.

The task of deploying these attack daemons requires the attacker to gain access and infiltrate the host computers. The third component of a distributed denial of service attack is the control master program. Its task is to coordinate the attack. Finally, there is the real attacker, the mastermind behind the attack. By using a control master program, the real attacker can stay behind the scenes of the attack. The following steps take place during a distributed attack:

- The real attacker sends an “execute” message to the control master program.
- The control master program receives the “execute” message and propagates the command to the attack daemons under its control.
- Upon receiving the attack command, the attack daemons begin the attack on the victim.

III. INTRODUCTION TO REAN (ROBUST EFFICIENT ADAPTIVE NETWORK)

To address different problems in networks we have designed a robust, efficient, adaptive network named as REAN which has following characteristics:-

ROBUST- Robustness is defined as "the ability of a system to resist change without adapting its initial stable configuration".

EFFICIENT- The extent to which time or effort is well used for the intended task or purpose.

ADAPTIVE- Adaptive behavior is a type of behavior that is used to adjust to another type of behavior or situation. This is often characterized by a kind of behavior that allows an individual to change an unconstructive or disruptive behavior to something more constructive.

NETWORK- A network is a telecommunication network that allows computer to exchange data. The physical connection between networked computing devices is established using either cable media or wireless media.

PROTOCOL- A set of rules and regulations that determine how data is transmitted in telecommunications and computer networking.

IV. WORKING OF REAN

Create a network consists of 30 nodes using AODV protocol. Create clusters and make cluster head, gateway nodes using cluster head gateway protocol. For sending data from one node to other we have to select a path that is best over the other paths (i.e. with minimum hopes). After selection of path we have to detect where the DDoS is attacking in the network, then the prevention mechanism is to be applied on that affected area. In the end we apply the maintenance procedure for the nodes in the network. And this maintenance procedure is working till the network is alive.

Parameter used in REAN

In our network use use the following number of parameters which are mentioned below:-

- C_i = Cluster node,
- Ch_i = Head cluster,
- Cg_i = Gateway node,
- NC_i = Centre node,
- B.W = Bandwidth,
- N_i = Node,

P_{Ni} = Participating node, P_i = Path selected,
 W_i = Weight,
 $L.B_i$ = Load balancing Factor, $D.L_i$ = Delay,
 Q = Priority Queue,
 $VAL [N_i]$ = Value of i th Node, $N [RT]$ = Routing table of node, P_f = Path formation,
 Fid = Flow id, Sid = Source id,
 Did = Destination id,
 $P [SR]$ = Packet sending rate,
 $W [f(x)]$ = Weight function of x .

V. PROPOSED ALGORITHM

In this section we discuss about the proposed algorithm.

- Create a network consists of 30 nodes using AODV protocol.
- Create clusters and make cluster head, gateway nodes using cluster head gateway protocol.

i. Condition for node to be cluster head: - Node should be in centre within cluster and bandwidth of that particular node should be the maximum among all nodes within the cluster.

ii. Condition for node to be cluster gateway: - node should lie between two or more clusters.

iii. Rest of the nodes are participating nodes.

• ELECTION OF PATH:-

We would use the path formation and path maintenance procedure of AODV and we will pick best three paths among all paths from source to destination.

Condition of selection of path :- path should contain less hops, weight factor (combination of load balancing + delay rate in respect with network)

Maintain a priority queue to place all 3 paths for communication and pick a particular path only on the basis of priority

• MAINTENANCE OF PATH:-

Periodically all cluster head will flood a status packet to ensure whether all nodes are still within the vicinity of its clusters.

Condition for maintenance of path:- if value of node = '1' then check the IP address of node in the path routing table and remove the path which contains non-participating node (factor value = 1) and refresh the routing table.

• DETECTION OF DDOS:-

Pocket formation: - flow id + source id + destination id + packet sending rate Calculate the weight factor $[W_f(x)]$:-

If $W_f(x)$ of each node lies within the range start the communication and periodically applies the detection of DDOS.

Else: - PREVENTION OF DDOS

Select second path from path routing table on the basis of their priority from priority queue.

- Repeat step 2 to above step.
- Stop communication.

Pseudo Notations REAN (This algorithm is designed for network using these notations :- C_i = Cluster node, CH_i = Head cluster, CG_i = Gateway node, NC_i = Centre node, $B.W$ = Bandwidth, N_i = Node, PN_i = Participating node, P_i = Path selected, W_i = Weight, $L.B_i$ = Load balancing Factor, $D.L_i$ = Delay, Q = Priority Queue, $VAL[N_i]$ = Value of i th Node, $N[RT]$ = Routing table of node, Pf = Path formation, Fid = Flow id, Sid = Source id, Did = Destination id, $P[SR]$ = Packet sending rate, $W[f(x)]$ = Weight function of x)

Step 1. Start.

Step 2. Create a network consists of 50 nodes using AODV.

Step 3. Create $[C_i]$, $[CH_i]$ & $[CG_i]$

For node to be $[CH_i]$

$N_i = NC_i$ & $B.W[N_i] = \text{Max}[N_1, 2, 3, \dots, N]$

For node to be $[CG_i]$

N_i lies b/w 2 or more clusters If $N_i = CH_i$ & $N_i = CG_i$ Then $N_i = PN_i$

Step 4. [Election of path]

Pick best three paths among all paths $[P_i]$ from source to destination. $P_i \rightarrow \text{Min hopes}$

$W_i = [L.B_i + D.L_i]$ Queue = $Q[P_1, P_2, P_3]$

Choose path from the priority queue with highest priority. Step 5. [Maintenance of Path]

Periodically all CH_i will flood a status packet. If $VAL[N_i] = 0$

Then Node is outside the cluster & check $N[IP]$ of Node in Path routing table & remove the path.

Refresh the $N[RT]$ or update. Step 6. [Detection of DDOS] Packet Formation $[Pf]$

$Pf = Fid + Sid + Did + P[SR]$

b) Calculate Wt. Factor ($W[f(x)]$) if $0 < (W[f(x)]) < 7$

Then Start Communication & periodically applies detection of DDOS. Step 7. Else [Prevention of DDOS]

Select P_i from routing table $P_i = Q[P_2, P_3, P_4]$

Step 8. Repeat Step 3 to 6 Step 9. Stop communication.

VI. CONCLUSION AND FUTURE PROSPECT**Conclusion**

Detection & Prevention of DDoS attacks is a part of an overall risk management strategy for an organization. Each organization must identify the most important DDoS risks, and implement a cost-effective set of defense mechanisms against those attack types causing the highest risk for business continuity. Studies and news about real-life DDoS attacks indicate that these attacks are not only among the most prevalent network security risks, but that these attacks can also block whole organizations out of the Internet for the duration of an attack. The risk from DDoS attacks should not thus be underestimated, but not overestimated, either. In this dissertation we try to overcome problem of DDoS attack.

Future Scope

In future, we will evaluate our framework for more internet topologies. In particular, we plan to investigate the following issues in more detail.

1. Introduction to load balancing in REAN: This dissertation only uses the load balancing factor but we can't adjust the factor according to our needs. So this issue can be considered for the further research work in future.

2. Quality of service (QOS): In our dissertation the concept of QOS is not introduced, due to bandwidth constraints and dynamic topology of Mobile Ad-hoc Network (MANET), supporting QOS in MANET is a challenging task. So we plan to implement QOS while designing a network.

3. Under water networks: Further we also try to resolve the underwater problem.

REFERENCES

- [1] Ajay Jangra, Sunita Beniwal, Anil Garg, "Co-existence behavior study of Bluetooth & Wi-Fi for 2.4 GHz ISM band" 2006
- [2] Alessandro Mei, Julinda Stefa, "Routing in Outer Space: Fair Traffic Load in Multi-Hop Wireless Networks" *MobiHoc'08, Hong Kong SAR, China in May 26–30, 2008.*
- [3] Binod Vaidya, Sang-Soo Yeo, Dong-You Choi, Seung Jo Han, "Robust and secure routing scheme for wireless multihop network" Published online: 4 April 2009, in Springer-Verlag London Limited 2009.
- [4] Bogdan Carbutar, Ioanis Ioannidis and Cristina Nita-Rotaru, "JANUS: Towards Robust and Malicious Resilient Routing in Hybrid Wireless Networks" *WiSe'04, October 1, 2004.*
- [5] Caroline Gabriel, "WiMax", ARCchart Ltd., London EC2A 1LN
- [6] Carlos A. Flores-Cortés, Gordon S. Blair, Paul Grace, "A Multi-protocol Framework for Ad-hoc Service Discovery" Melbourne, Australia MPAC '06, November 27-December 1, 2006.
- [7] Charikleia Zouridaki, Brian L. Mark, Marek Hejmo and Roshan K. Thomas in October 2008.
- [8] Chien-Chung Shen, Chaiporn Jaikaeo, "Ad Hoc Multicast Routing Algorithm with Swarm Intelligence" *Mobile Networks and Applications 10, Springer Science + Business Media, Inc. Manufactured in The Netherlands, 47–59, 2005.*
- [9] Chin-Yang Tseng, Poornima Balasubramanyam, Calvin Ko, Rattapon Limprasittiporn, Jeff Rowe, Karl Levitt, "A Specification-based Intrusion Detection System for AODV" *Proceedings of the 1st ACM Workshop Security of AdHoc and Sensor Networks Fairfax, Virginia © 2003 ACM-1-58113-783-4/03/0010.*
- [10] Chin-Fu Kuo, Hsueh-Wen Tseng, Ai-Chun Pang, "A Fragment-Based Retransmission Scheme with QoS Considerations for Wireless Networks" *IWCMC'07, August 12{16, 2007, Honolulu, Hawaii, USA}.*
- [11] Claude Castelluccia, Nitesh Saxena, Jeong Hyun Yi, "Robust self-keying mobile ad hoc networks" *Computer Networks 51 (2007) 1169–1182 1389-1286 2006 Elsevier B.V. doi:10.1016/j.comnet.2006.07.009.*
- [12] C.Siva Ram Murthy & B.S Manoj, "Mobile Ad Hoc Networks- Architectures & Protocols", Pearson Education, New Delhi, 2004.
- [13] Danesh, A. and Inkpen K., "Collaborating on AdHoc Wireless network", at www.parc.xerox.com/sl/projects/ubicomp-workshop/positionpapers/danesh.pdf.
- [14] Dr. Sanjeev Sofat, Prof. Divya bansal and Rajinder Kumar, "Security in Mobile Ad Hoc networks", COIT-2008 March 29.

Analytical Study on Applications of Web Search using Stochastic Query Covering

Tiruveedula Gopikrishna*

Doctorate Program in CSE, Research Scholar, Rayalaseema University,
Kurnool, Andhra Pradesh

ABSTRACT

The digital revolution saw amid the most recent decade has resulted in an explosion of accessible online content. An expanding number of individuals utilize the Internet to scan for particular information and to remain informed by perusing news or client generated content. Our research work is identified with the field of web utilization mining. Specifically, we analyze information put away in web crawlers' logs to find utilization patterns, and the point is to upgrade execution of hunt devices and to help clients to discover information on the web. Enhancement of the execution of these frameworks is of paramount significance given that a cutting edge web search tool gets an enormous number of questions each second, and clients expect speedy reactions. As a first commitment, we show a successful approach for choosing a subset of documents to store in a static reserve with the reason for making the query handling speedier.

They are communicated by encasing a multi-word sequence with quotation checks and enforce that the web crawler returns just documents containing the cited expression. In the two commitments we utilize the theoretical aftereffects of the set-cover issue to show the effectiveness of our methodologies. Moreover, we verify the theoretical discoveries with experiments over true datasets.

Keywords: *Web Mining; Stochastic Query Covering; Avaricious Multi-cover; Greedy Multi-cover; Query Log Mining.*

1. INTRODUCTION

The World Wide Web is a prominent medium for disseminating information, staying in contact with companions, and delivering items or administrations. Information accessible on the web is a tremendous source of knowledge, in this manner individuals utilize the overall network each day to satisfy their information needs or to remain educated.

While initially the clients were just consumer of web content, these days they contribute to the quick increment of information and multimedia content accessible on the web. Common cases of client generated content are: (I) thoughts and suppositions Posted on Facebook, Google+, Twitter, and other social websites; (ii) item surveys and tutorials distributed on blog websites; and (iii) photographs or recordings shared on prominent stages, for example, Flickr and YouTube. From one viewpoint, client content accelerates the pace at which information ends up noticeably accessible on the web, yet then again, web clients are suffocating in information and this wonder is otherwise called information overburden.

Over the most recent couple of years, a great deal of consideration has been committed to improving web look and recommender frameworks through information mining. It permits filtering through huge quantities of information for helpful information. Web mining alludes to the way toward gathering knowledge from web information. It is a multidisciplinary exertion that acquires ideas and techniques

from fields, for example, information retrieval, statistics, machine learning, and others. A branch of web mining, called web utilization mining is committed to the extraction of beforehand obscure patterns from information depicting the interaction of clients with the web. The clients give heaps of hints about their interests and purposes through their activities. In this way, the analysis of utilization information is useful for web look, web personalization, and web based business.

Query logs of web crawlers store helpful information about the looking conduct of clients. The query appropriation, clicked comes about, and other information can be exploited to enhance the accuracy of list items to configuration query-result reserving systems and to help propelled look functionalities.

1.1 Web Mining

Information mining is the computational procedure of finding intriguing patterns in vast datasets. It permits the non-unimportant and automatic extraction of implicit and potentially helpful information from huge measure of information. With the blast of content and administrations accessible on the web, the web has as of late turned into a rich zone for information mining. Web mining comprises in the utilization of information mining techniques to information, antiques, and exercises identified with the web. It is a dynamic and wide research region, which draws systems from statistics, database, information retrieval, and some branches of artificial insight, for example, machine learning and characteristic dialect handling (NLP). Web mining is by and large partitioned into three fundamental subareas, comparing to three diverse knowledge-discovery spaces:

1.2 Web content mining

It induces knowledge from content accessible on the web. Web content commonly comprises of content, graphics, and multimedia. Web-content-mining research zone has predominantly centered on unstructured archives (e.g., free content) and semi-organized reports (e.g., HTML records). Content can be spoken to as sack of words, which does not consider the places of the terms in the archives, or utilizing structures which consider likewise the sequences and places of terms (e.g., n-grams or expressions). The primary uses of web content mining are situated to the content arrangement and order and to the occasion detection and following. A different line of research depends on a database perspective of web content mining. It endeavors to demonstrate and coordinate web information as a knowledge base, with the goal that more refined queries can be performed.

1.3 Web structure mining. It separates information from information depicting the association of web content. Web structure mining investigates intra record information and in addition between report information. The previous speaks to any information about the association of content inside a web page (i.e., the courses of action of HTML or XML labels), and the point is to enhance the association of information inside single web reports. The last involves hyperlink associations between web pages that have a place with a similar website or to various websites. Hyperlinks can be utilized for dissecting the structure of the web and its advancement. Besides, connections to a web page can be viewed as an implicit support of pages, so web structure mining has likewise gotten a considerable measure of consideration propelled by applications, for example, finding definitive websites and recognizing noxious movement on the web.

1.4 Web utilization mining

It derives use patterns in information generated from the interactions of the clients with websites, web administrations, and web indexes. Uniquely in contrast to web content and structure mining, which

dissect primary information on the web (e.g., content and connections), web utilization mining investigates secondary information, which catches the use patterns (e.g., queries issued to a web index and navigate information). Web utilization information incorporates server-get to logs, intermediary server logs, program logs, client profiles, registration information, client sessions, transactions, treats, and client queries. Mining utilization information permits finding perusing and seeking patterns.

1.5 Uses of Web Usage Mining

Web use mining has gotten a great deal of consideration in a few regions, including web personalization, web based business, and web based promoting. It is generally utilized by recommendation frameworks for making intriguing proposals for items and web pages. Besides, information put away in query logs of web crawlers can be investigated for improving web look.

Web based business. Information mining has as of late observed a fast increment in commercial intrigue. It permits to plan powerful special battles and to distinguish cross-marketing procedures. The web encourages business transactions, and the web based business is one of the significant powers that enable the web to prosper. The accomplishment of online business relies upon how well the website proprietors comprehend clients' desires. Web use mining is an effective instrument to break down consumers' perusing and purchasing conduct for discovering their inclinations. Commercial websites frequently customize their pages to make them simple to explore. They additionally use item recommendation frameworks to propose things to their clients. Different innovations have been proposed for making recommendations, and a large number of them depend on the things beforehand purchased by the client or by alternate clients who have comparative tastes. Web based Advertising. The analysis of the online movement of customers is additionally critical for internet promoting. Commercial locales and web crawlers have relationship with commercial marketing organizations, for example, Double Click Ad Exchange, for expanding their pay through publicizing. A commercial marketing organization utilizes treats to screen the exercises of guests of web based business destinations. It gathers all the information about a client as a profile in a database, with the goal that when the client visits one of the destinations partnered to the organization, the profile information can be utilized to choose the advertisement to appear on the page. Late investigations have likewise centered on eye development. Given a web page, following eye developments permits to comprehend where clients center their consideration and to determine the best positions for the advertisements.

1.6 Query Log Mining

Query log mining utilizes information put away in logs of web indexes with the reason for dissecting seeking conduct of clients. Measurable highlights extricated from query logs, for example, normal length of queries, query sessions, clicked comes about, have featured that web queries are not quite the same as the queries generally issued by clients of little information retrieval frameworks (e.g., the IR frameworks of computerized libraries). Ordinarily, the web clients sort short queries (i.e., a few terms) and don't utilize Boolean administrators. They take a gander at the main page of results (i.e., top-10 results), and the greater part of the inquiry sessions are short.

1.7 Utilizations of Query Log Mining

Query Expansion, Suggestion, and Spelling Correction. Web queries are issued by clients who need more information about a theme. These queries might be short, seriously defined (i.e., excessively particular, excessively bland, or questionable), and in some cases incorrectly spelled. Information put away in query logs can be utilized for query extension, query proposal, and automatic adjustment of

grammatical errors. Query extension is a strategy generally utilized via web crawlers. It comprises in growing the query by including terms for making the query more expressive and, subsequently, expanding the accuracy of the list items. Queries can be successfully extended by utilizing terms already wrote by clients to enhance the first query or dissecting the content of clicked archives [4,5,6]. Query development is restrictive as far as adaptability. Besides, it is independent of clients' inclinations, in light of the fact that a query is extended similarly for every one of the clients. At last, the clients may feel overpowered by information, since comes about are loaded down with different archives, which might be not intriguing for them.

1.8 Improvement of Performance of Web Search Engines

Query logs have been broadly broke down for improving the productivity of web indexes. Present day web crawlers are expansive scale disseminated frameworks, where the reversed list is apportioned among numerous pursuit modules, running on different bunches of servers. The list can be report parceled, so each segment is with respect to a sub-accumulation of archives, or term-apportioned, to be specific, the record is separated evenly, and segments contain particular subsets of terms happening in the gathering. Once the query is presented, a machine before the bunch communicates the query to the pursuit modules. In record apportioned disseminated designs, to decrease the hunt space, the gathering choice can be utilized. It depends on steering the query to a subset of servers, which are destined to contain significant reports for the query. Methodologies for record dividing and accumulation determination have been exhibited in the writing. Recent examinations have concentrated on utilizing information removed from query logs. Specifically, the creators of proposed archive apportioning and accumulation determination in light of co-grouping of queries and records. Bunches of reports are allotted to various inquiry center modules, while query groups are utilized for accumulation choice.

1.9 Supporting Phrase Queries

A less investigated look into territory concentrates on the analysis of query logs for supporting propelled seek highlights, for example, state queries. Bahle et al. examined logs of commercial web indexes with the motivation behind concentrate the attributes of expression queries. They watched that albeit express expression queries speak to a little level of the queries issued by clients, the vast majority of the non-expression queries which have more than single word can be prepared effectively as expression queries. This is on account of huge numbers of these queries (e.g., titles of melodies, films, or books) are in any case proposed to be phrases. For improving the expression query handling, in the writing a few works have proposed utilizing an upset list enhanced with phrases. Chang and Poon introduced a positional transformed file enhanced with state queries normally issued by the clients. Such sequences of words are listed permitting quick retrieval of those archives which coordinate the expression queries.

1.10 Stochastic Query Covering

With the touchy development of computerized information, individuals find wanted information depending on information retrieval frameworks. Normal cases of these frameworks are expansive scale web search tools, database frameworks, and computerized libraries. Current information retrieval frameworks regularly reserve query results to lessen query preparing and information exchange costs. Database frameworks store queries and their outcomes at the customer side. Storing is likewise generally utilized in web crawlers, which typically reserve consequences of well known queries and posting arrangements of the most continuous query terms.

We will likely show systematically and tentatively that the analysis of the query-archive structure and of factual information separated from query logs can be utilized for choosing a subset of records which, by and large, boosts the quantity of client queries completely served by the store. As a feature of our commitment, we characterize the query-multi-cover issue. We need to make a widespread guide from each query to an arrangement of reports which contribute to the aftereffects of the query. We reserve these records, with the goal that when a query is put together by the client, the framework can utilize stored reports to serve the present query and also potentially future ones. The issue can be viewed as a stochastic speculation of the set-multi-cover issue, in which the components to cover relate to queries and the covering sets are records. The advancement issue is NP-hard. In addition, uniquely in contrast to the conventional issue, we have to characterize a settled mapping from components to covering sets without knowledge of the components to cover, since we have no from the earlier knowledge without bounds queries. We demonstrate that knowing the query circulation gets the job done to give calculations logarithmic approximations of the ideal arrangement. Besides, there exists an arrangement of reports that covers an expansive portion of the queries, and a basic eager approach can discover it [1-2].

1.11 Proficient Phrase Indexing and Querying

Information retrieval frameworks offer a few hunt functionalities to the clients. One of them is discovering records that contain a correct arrangement of words. The queries, additionally called express queries, are communicated with cited phrases, to be specific; the grouping of words is encased by quotes (e.g., "Bruce Springsteen," "leader of the assembled conditions of America," and "moon stream"). Expression queries are upheld by all the cutting edge web crawlers. They are straightforward and natural to utilize, maintaining a strategic distance from the uncertainty that is regularly taken cover behind a solitary word query or an and-query. In addition, web crawlers implicitly summon express queries, for example, by methods for query division. Expression queries are likewise critical for applications, for example, substance arranged inquiry and copyright infringement detection.

2. RESEARCH METHOD TECHNIQUES

We exhibit the stochastic query covering as an appropriate model of the situation laid out above. What's more, as we examine in Section 1.10, it takes into account shrewd analysis while a most pessimistic scenario approach does not give any instinct [1]. The drawback is that the analysis turns out to be technically more included. In particular, we expect a structure in which clients submit queries to an archive retrieval framework after some time. As the framework utilizes a store of restricted size to hold a subset of the report accumulation. In a perfect world, records in this subset have a tendency to be important for queries that clients are well on the way to submit. At whatever point a client presents a query, the report retrieval framework must restore an arrangement of records pertinent to the query. The framework either develops the outcomes utilizing records put away in reserve, or it flops; in the last case, the framework acquires a store miss, that is, a punishment mirroring the way that building an outcome page will require a tedious operation (e.g., getting to secondary stockpiling). In our model, for any given query, a record has a query-subordinate weight that measures its significance regarding a particular rundown of query comes about (e.g., weights can mirror the level of pertinence of the reports to the queries). At the point when the query q is presented, the record retrieval framework restores the arrangement of reports whose general weight as for the query is in any event some given limit, or as expansive as could reasonably be expected, subject to a cardinality requirement[1].

2.1 Applications

2.1.1 Web Search

Present day web crawlers need to process a huge number of queries every second finished accumulations of billions of records, and clients expect low reaction times. To this end, web indexes utilize an assortment of storing techniques as a way to give comes about auspicious and negligible decrease in quality [1].

2.1.2 Computational Advertising

Another potential application is in the region of computational promoting. Ordinarily, when a client issues a query to a web index (on account of supported hunt) or visits a content page (on account of content match) the online specialist organization (OSP, for example, Google or Yahoo, chooses few advertisements to show to the client from a pool of a few hundred million promotions. Picking the fitting advertisements is a confused method including information retrieval frameworks that recover promotions that are pertinent to the page or query, barter for picking the promotions to show among the applicable ones, and promotion trades.

2.1.3 Different Uses

Query-cover approach is exceptionally broad and can have a considerable measure of utilizations. In situations where the retrieval time exceptionally relies upon the gathering size, a query-mindful reserving methodology can decrease query-handling time and give, in the meantime, quality assurances. A few cases of potential applications are: semantic pursuit (in which measurable NLP techniques may must be connected at query time), picture or video retrieval (which may include tedious picture/video preparing), and querying of huge organic databases [1].

3. ALGORITHMS AND ANALYSIS

The Query Multi-cover(t) issue (Problem 1) is a stochastic speculation of the set-multi-cover issue, in which components compare to queries and archives relate to sets. Specifically, we consider a setting that takes after the one by Grandoni et al., who examined the issue of stochastic widespread set cover. Uniquely in contrast to the customary set-cover issue, we have to characterize a settled mapping from components to covering sets without knowledge of the components to cover, since we don't have from the earlier knowledge of the queries that will be submitted. Likewise, we are thinking about an expansion of the above issue in which sets have related component subordinate weights, and the components have scope necessities [1]. For the double meaning of query-report weights the issue is similar to the set-multi-cover issue, in light of the fact that the point is to cover queries with in any event k records. While for alternate cases, the queries have scope necessities relying upon the general entirety of the weights of archives significant to it. All the more absolutely, the query is secured if the total of record weights is no less than a given edge W .

We say here that in a deterministic setting we can't demonstrate solid, significant outcomes: The most pessimistic scenario limits are free and don't give a considerable measure of understanding. Instead, we demonstrate that knowing the query dissemination gets the job done to furnish algorithms with logarithmic (expected) approximations. We show a basic and effective ravenous algorithm, and we demonstrate that it accomplishes logarithmic estimate proportion under some sensible suppositions. Note that the evidence in for the set-cover issue can't be connected to our setting and that we require more advanced contentions to demonstrate that our algorithm accomplishes a decent estimation.

4. AVARICIOUS MULTI-COVER ALGORITHM

In this segment we introduce a covetous algorithm for the Query Multi-cover(t) issue. Algorithm 1 is called Greedy Multi-cover (GM for curtness), and it is the direct voracious approach: In every cycle it chooses the report that covers the biggest aggregate weight among the revealed queries. We take note of that, despite the fact that the algorithm is straight forward, the analysis is decently nontrivial in our stochastic setting [1].

Continue with the analysis of the execution of the algorithm. For introduction, we expect that for each query q , the i th heaviest record has weight w_i , freely of q , as for the situation that weights depend on the positioning. All things considered, our outcomes hold for a general weight structure with insignificant changes, for the most part to the documentation [1]. By chance, see that a similar report may have distinctive weights for various queries. Characterize ℓ to be the base number with the end goal that $\sum_{i=1}^{\ell} w_i \geq W$, and note that we have $\ell \leq W/w'$, in light of the fact that $w' \leq w_i$, for $i \leq \ell$. Consider a succession of t queries that are inspected from the dissemination $Q[2]$. For a succession $\omega = (q_1, \dots, q_t)$, where the q_i s are freely and indistinguishably appropriated tests from Q , let $C^{opt}(\omega)$ be the ideal cost, and let $C^{-opt} = E * C^{opt}(\omega) +$ be the normal ideal cost. Additionally, characterize $C(\omega)$ and C^{-} as the cost and the normal cost, individually, actuated by the Greedy Multi-cover algorithm. Regarding this setting, we now demonstrate our principle hypothesis:

Hypothesis 2. For any grouping of t queries the mapping made by the Greedy Multi-cover algorithm fulfills.

Evidence. – (Sketch) – Here we give the basic components of the evidence for the case in which Q is the uniform circulation. The total evidence is accounted for in the following passage, where we likewise demonstrate to stretch out the confirmation to the non-uniform case. The confirmation comprises of two sections [1]. The first is given by Lemma 2, where we demonstrate that the algorithm can cover with weight W everything except $8n t C^{-opt} \ln mn$ queries in close to $97 W w' C^{-opt} \ln(nW/w')$ executions of the primary circle. For the second part, let Q_{unc} be the arrangement of components of Q that have been secured with a weight not as much as W , the normal number of queries from Q_{unc} that show up in an irregular succession of length t (potentially with reiterations) is $t n |Q_{unc}| = 8C^{-opt} \ln mn$ [1]. These components are secured by the algorithm without any than ℓ records for each query and, consequently, without any than $8(W/w')C^{-opt} \ln mn$ archives. Along these lines, the cost of the Greedy Multi-cover algorithm is $O(C^{-opt}(W/w') \ln mn)$.

5. ANALYSIS OF ALGORITHMS VARYING PARAMETERS

In this section we dig more into the conduct of the algorithms by concentrate the execution for various estimations of the parameters. Figure 1.1 and 1.2 report review and diminish store for the algorithms parameterized with k , which are Top- k , Bin.GM, and Card.GM. Each bend relates to an alternate estimation of k , and we report the outcomes acquired for $k = 2, 6, \text{ and } 10$.

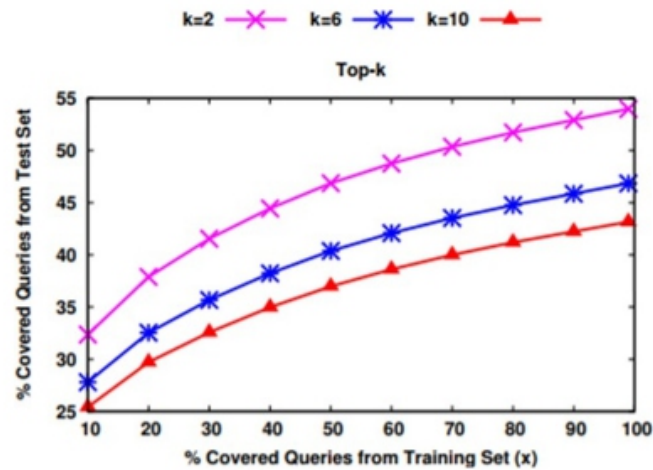


Figure 1.1 Report review and diminish store for the algorithms parameterized with k.

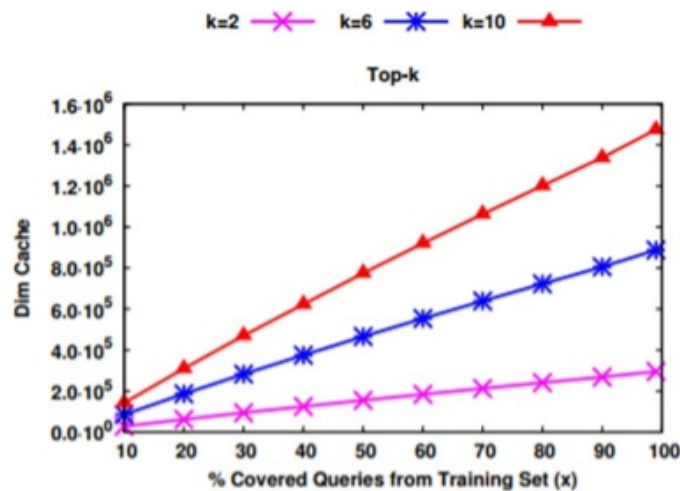


Figure 1.2 Report review

6. ANALYSIS OF ALGORITHMS VARYING TRAINING DATA

The greedy approach is parameterized by x , that is, the level of the queries of the preparation set chose to be secured by the Documents put away in the reserve. In this section we dissect the distinctive practices of the algorithms utilizing diverse preparing information. We differ the determination of the queries of preparing information thinking about two criteria:

6.1 Freq-queries: this algorithm understands the greedy choice over the most incessant queries of the watched period. The algorithm chooses the $x\%$ most incessant queries in the preparation set, and afterward it bit by bit picks the archives to store until the point when each query is secured [2].

6.2 First-queries: this is the same as the past one, with the distinction that it covers avariciously the principal (soonest) $x\%$ of the queries of the watched period [2].

We think about the execution of freq-queries and first-queries, which cut the dispersion of queries and consider the $x\%$ of the most successive queries or the soonest queries in the dissemination, against the execution of the algorithm that watches every one of the queries of the appropriation and chooses records ready to cover $x\%$ of them [2].

7. CONCLUSION

Results accomplished utilizing algorithms GM, Bin.GM, and Card. We watch that freq-queries and first-queries perform correspondingly. This can be ascribed to the way that the appropriation of the queries is sufficiently stationary in a little era, along these lines the statistics gathered as time continues look like those of the whole time frame, prompting the execution of first-queries being like that of freq-queries. As should be obvious, both the methodologies have more regrettable review and the store measure develops straightly when x increments. This leads us to the conclusion that in the event that we truncate the circulation of queries (e.g., considering only a small amount of soonest or most successive queries), we lose critical information.

REFERENCES

- [1] Aris Anagnostopoulos, Luca Becchetti, Ilaria Bordino, Stefano Leonardi, Ida Mele, Piotr Sankowski. "Stochastic Query Covering for Fast Approximate Document Retrieval", *ACM Transactions on Information Systems*, 2015.
- [2] Aris Anagnostopoulos, Luca Becchetti, Stefano Leonardi, Ida Mele, and Piotr Sankowski. Stochastic query covering. In *Proceedings of the 4th International Conference on Web Search and Data Mining (WSDM '11)*, pages 725–734, New York, NY, USA, 2011. ACM.
- [3] www.dblp.dagstuhl.de
- [4] www.doctorat.ubbcluj.ro
- [5] www.user.ceng.metu.edu.tr
- [6] www.junminghuang.com.

Detecting Malware by Data Mining

Shafiqul Abidin¹, Rajeev Kumar², Varun Tiwari³

¹HMR Institute of Technology & Management (GGSIP University) Hamidpur, Delhi, India

²Department of Computer Science, Kalka Institute for Research and Advanced Studies (GGSIP University) Alaknanda, New Delhi, India

³Comm- IT Career Academy (GGSIP University) New Delhi, India

ABSTRACT

The exponentially growth of malware has created number of security threats in IT industry. A large number of viruses are developed and millions of applications are infected and suffered on daily basis. Trojan is one of the fatal and deadly types of malware. But it is often said as legitimate software. They hide themselves within harmless programs. Trojan survived by going unnoticed. They look like just about anything like the computer game as downloaded from different websites. Sometimes even a popup advertisement might try to install something on our computer. Trojan can trick you into using them. In this paper, data mining technique is being proposed to detect Trojan. The technique is based on Naive Bayes – this technique is simple to put into practice and we achieve amazing results in large number of cases. But practically, dependencies exist among variables.

KEYWORDS: Trojan Detection; Data Mining; Decision Tree; Naive Bayesian Network; Naïve Classification Technique.

1. INTRODUCTION

Malware can be said as the collective term for virus, Trojan horse and other malicious that can infect the system. Since, many years these harmful items have evolved and affected smart phones and tablets as well. Malware is sometimes known as computer contaminant, as in the legal codes of several U.S states [1]. Malware comprised of damaging function that is called Payload which has various effects. It bears the quality to be unnoticed. This unnoticed nature is achieved by actively hiding and showing no presence to users. The generic term malware comes from “malicious software”, where malicious describes any code in any part of a software system that is intended to cause damage to a system. The types of malicious codes are Virus, Worms, Trojan and other malicious.

But what is the difference between a virus and a worm? What is the difference between these two and Trojan? Does antivirus apply against Worms, Trojan, Virus and other malicious codes? All these questions come from one source and it's the complex and complicated world of destructive codes [2].

Several types of malicious codes have some kind of behavior which are described below-

Virus

A code which get attached itself to a host program and propagates whenever that infected program executes.

Worms

Unlike virus, a worm does not attach itself to an infected executable program but it spreads itself by transferring via network which includes some connected computers.

Trojan Horse

This includes a hidden program component, which are in a form of pieces of software code which opens a backdoor into the affected computer and thereby allow almost full access to the user noticing.

Trojan horse often referred as Trojans. In 1986, the first Trojan was 'PC-Write'. Trojan is derived from the Ancient Greek story of the wooden horse that used to protect the city of Troy. Trojans are totally different from virus and worms they do not introduce themselves or disseminate themselves into other files, it just represents itself as useful or interesting that tempt a user to install it. Trojans are classified according to the type of actions they can perform on a computer.

Backdoor

This gives malicious users control to do anything they wish on the infected computer, which includes sending, receiving, deleting, launching files, display and rebooting also.

Exploit

This is a piece of data or a sequence of commands that take advantage and attacks within application software running on the system.

Rootkit

This is designed in order to prevent malicious programs being detected. It is difficult to detect because it is activated each time system boots up.

Trojan-Downloadert

This can easily download and install different types of new malicious program into a system. There are also other types of Trojan too like Trojan Banker, Trojan DDOS and many. Trojan can do a lot of harm to a system like - destruction to the system; corrupt data or delete; modify data; spy; use computer resource; infect other connected device etc. In short short it can do a lot of harm to the system.

As signature method is a traditional and usual method to detect malicious program. They are created manually; it matches with at least one byte code pattern of the software. As, researchers have tried to present more reliable methods for malware detection.

2. RELATED WORK

The process of identifying malware is called analyzing, which are roughly divided into static and dynamic analysis.

STATIC ANALYSIS

In this program code is checked. But in actual program code does not execute. It investigates and detects coding flaws, back doors and potentially unwanted codes. In static method, binary codes are checked and detected according to the binary codes given.

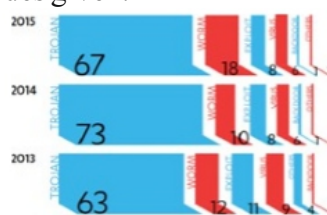


Figure 1. Percentage of Malware Detection Reported

DYNAMIC ANALYSIS

Dynamic Analysis is evaluated in a runtime environment. Its key objective is to find bugs in a program, during run time instead of repeated code examining. Actually they analyze what is being happened behind the scene. Sometimes static and dynamic analysis is considered as glass-box testing. In this article an effort has been made to get ascertain static analysis method by implementing a data mining technique to detect & clean Trojan.

3. LITERATURE SURVEY

Earlier malware was detected using signature methods, and then researchers found a number of classifiers for analysis and detection of malware. Many classifiers used n-gram i.e. a series of bits in some order and extracted from hex dump like mentioned in an International Journal of Intelligent Information Systems paper, presented on “Malware detection using data mining techniques” that has a higher success rate, as it finds whether there is a malware or not using the binary codes and as per rootkit detection success rate is over 97% [3].

In another paper decision tree and Naïve Bayes data mining techniques are used to detect virus. That consists of more than three thousand malicious and more than one thousand benign programs were there, where firstly op-code is used as vector and secondly, op-code as well as first operand were taken where benign and virus programs were mixed which affects the effectiveness of both the classifiers [4]. Proposed surveillance spyware detection system (SSDS), where features were considered as both static and dynamic using information gain method 76 static and 14 dynamic were discovered. According, to their research SSDS is a better performer than other known antivirus like Norton, Kaspersky, etc. A signature based method called as SAVE (Static Analysis of Vicious Executable), represented as API calls and as well as used Euclidean distance to compare signature with API calls [5].



Figure 2. Cyber Security Model

So, besides data mining method other techniques are also used to detect malware as malware detection is a very important part in security. Data preparation is needed for data mining process, where data needs to be collected then its feature needs to be extracted and model should be classified to get the result [6] [11].

4. NAÏVE BAYES

This Naïve Bayes is based on probability and works on independent assumptions. This method is used both for multi class classification and byte sequence data Naive Bayes has been studied. Since, 1950s and is popular for data classification, where document is assigned to one or more categories (can be text, spam, image or music, etc.) [7].

Applications of Naïve Bayes Classification include:

- Text Classification
- Spam Filtering
- Hybrid Recommender System
- Online Applications – Simple Emotion Modeling.

In 2003, Virus detection using data mining method is was published i.e. Multi Naive Bayes. They are quite well in complex problems. Naive Bayes are simple, fast and highly scalable that only requires small amount of data to estimate. It is based on conditional probabilities i.e. calculates a probability by counting the frequency of values and combination of a given data [8].

In this method, we want to compute a certain given text document, as it states-

$$P(x/y) = P(y/x) \cdot P(x)/P(y) \quad (i) \text{ Where,}$$

$P(x/y)$ = posterior probability $P(x)$ = prior probability

Priori: probability of an event before the evidence is observed. **Posterior:** probability of an event after the evidence is observed

Here, x is a vector $x=(x_1, x_2 \dots x_n)$. To use Naive Bayes technique, we assume features which occurs independently. Suppose feature is F , then

$$\begin{aligned} F &= (F_1, F_2 \dots F_n). P(x, F_1, F_2 \dots F_n) = P(x) \cdot P(y_1, \dots, y_n/x) \\ &= P(x) \cdot P(y_1/x) \cdot P(y_2/x, \dots, y_n/x, y_1) \\ &= P(x) \cdot P(y_1/x) \cdot P(y_2/x) \cdot P(y_3, \dots, y_n/x, y_1, y_2) \\ &= P(x) \cdot P(y_1/x) \cdot P(y_2/x, x_1) \dots P(y_n/x, y_1, y_2, y_3, \dots, y_{n-1}) \text{ As, } i \neq j \\ P(x/F) &= \prod_{i=1}^n P(F_i/x) \cdot P(x) / \prod_{j=1}^n P(F_j) \quad (ii) \end{aligned}$$

Since, denominator is same for all the classes. So, we take maximum as computed in (ii) equation, we get $P(x/y_1, \dots, y_n) = \max(P(x) \prod_{i=1}^n P(y_i/x))$

We first, collected data then feature extraction and then we applied the equation for the program.

5. PROPOSED APPROACH

Analysis of program can be carried out in each step of our method like data collection, data preprocessing, feature extraction, and feature selection. At last decision tree and naive Bayesian network algorithms have been suggested and practical are carried to find the effectiveness of proposed technique.

6. DATA COLLECTION AND PROCESSING

We have downloaded collection of Trojan codes at <http://vxheaven.org/vl.php> and benign files were collected from a PC running windows XP includes operating system files and various windows application. Dataset consist of 4722 PE files, where 3000 are Trojan and 1722 benign programs.

The goal was to gather useful features and extract useful features from PEiD that distinguish between malicious and benign files where distribution of different packed, not packed, Trojans and benign programs are there.

FEATURE SELECTION AND EXTRACTION

In this, decision trees are divided into subsets and concurrently a connected decision tree is developed incrementally. Decision tree construction is to find attributes. These attributes return the highest information gain (i.e., the most homogeneous branches). Here, data set one comprises 890 Trojan codes and 150 benign programs. The expected results using the using the equation:

$$= -(|\text{benign}|/|X| \log_2 |\text{benign}|/|X| + |\text{Trojan}|/|X| \log_2 |\text{Trojan}|/|X|)$$

$$= -150/(890+150) \log_2 (150/890+150) + (-890)/890+150 \log_2 (890/890+150).$$

The attribute with the largest information gain is chosen the decision node and negligible information gain can be discarded to reduce number of features to speed up the classification. ID3 algorithm is run recursively, until all data are classified.

CLASSIFICATION AND MODEL TRAINING

Each data set was then fed to Naive Bayes technique; experiments were repeated several times using random sub-sampling holdout method, to obtain the accuracy from the iteration method. So, these results can be obtained [9].

TABLE 1. Results from Iteration Method

Naive	1 byte	76.1	41.2	73.3
Bayesian	2 byte	80.7	41.2	77.1
Decision	1 byte	93.2	29.4	89.5
Tree	2 byte	94.3	23.5	91.4

7. RESULTS

Results have been obtained after testing the data using the obtained data set to evaluate the correctness of the classification model for Trojan detection. The four estimates define the member.

True Positive (TP): Number of programs correctly identified as Trojan codes. **False Positive (FP):** Number of benign programs incorrectly identified as Trojan codes.

True Negative (TN): Number of programs correctly identified as benign programs.

False Negative (FN): Number of Trojan codes incorrectly identified as Trojan codes.

The action of every classifier was evaluated using false alarm, overall accuracy and detection rate: **Detection Rate (DR):** Percentage of correctly identified malicious programs.

Detection Rate = $TP/(TP+FN)$

False Alarm Rate (FAR) or False Positive Rate (FPR): Percentage of wrongly identified benign Programs –

False Alarm Rate = $FP/(TN+FP)$

Overall Accuracy: Percentage of correctly identified Programs Overall Accuracy = $TP+TN/(TP+TN+FP+FN)$

This data set was experimented. The unknown Trojan detection rates 93.2 per cent and 76.1 per cent with accuracies of 89.5 per cent and 73.3 per cent were obtained in first experiments. Each element consists of only the op-code whereas unknown Trojan detection rates are 94.3 per cent and 80.7 per cent and accuracies rise to 91.4 per cent and 77.1 per cent. More information is surfaced in each iteration and therefore Naïve Bayes classifier performs more accurately [10].

8. CONCLUSION

Naïve Bayes technique is simple to put into practice and we achieve amazing results in large number of cases. But practically, dependencies exist among variables. This article examined Trojan Detection using Naive Bayes technique. As Trojan detection is one of the major measures for security. This technique automatically extracts Trojan qualities from Trojan programs. Further, these qualities are used for classification. The obtained results and outcomes indicate that the rate of detection - the Decision Tree and Naive Bayes classifiers computed as 94.3 per cent and 80.7 per cent and the accuracy 91.4 per cent and 77.1 per cent respectively. This shows that Decision tree performs well if compared with Naive Bayes classifier. We need to put into effect suitable policies and checkup the legal aspects and need to undertake privacy from all directions for the security purpose.

REFERENCES

- [1] Matthew G. Schultz, Eleazar Eskin, Erez Ado and Salvatore J. Stolon "Data Mining Methods for Detection of New Malicious Executable".
- [2] Muezzin Ahmed Siddiqui, Morgan Wang, "Detecting Trojans Using Data Mining Techniques", January 2008.
- [3] Tzu-Yen Wang, Shi-Jinn Horn, Ming-Yang Su, Chin-Suing Wu, Pang-Chu Wang and Wei-Zen Su, "A Surveillance Spyware Detection System Based on Data Mining Methods", 2004.
- [4] Yuen Kou, Chang-Tine Lu, Sir rat Sinvongwattana You-Ping Huang " Survey of Fraud Detection Techniques", March, 2004.
- [5] Jay-Hwang WANG, Peter S. DENG, Yi-Sheen FAN, Li-Jing JAW, Yu-Ching LIU, "VIRUS DETECTION USING DATA MINING TECHNIQUES", 2004.
- [6] Sung, A.H., Xu, J., Chavez, P., Mukkamala, S. "Static analyzer of vicious executables", 20th Annual Computer Security Applications Conference, 2004.
- [7] Karta Mathura, Sari Hiranwal, "A Survey on Techniques in Detection and Analyzing Malware Executable", April 2013.
- [8] DauberKauri R. Chakra" Feature selection and clustering for malicious and benign software characterization" August, 2014.
- [9] Sara Najari, Iman Lotfi, "Malware Detection using Data Mining Techniques", International Journal of Intelligent Information Systems, October 20, 2014.
- [10] Mittal A. Saeed, Ali Selma, Ali M. A. Abuagoub, " A Survey on Malware and Malware Detection System" International Journal of Computer Applications, April 2013.
- [11] Shafiqul Abidin, Mohd Izhar, "Attacks on Wireless and its Limitations", International Journal of Computer Science and Engineering, Volume 5, Issue 11, November 2017.

Digits in Units and Tens Places of 3-PrimeFactors Numbers till 1 Trillion

Neeraj Anant Pande

Associate Professor, Department of Mathematics and Statistics, Yeshwant Mahavidyalaya (College),
Nanded-431602, Maharashtra, INDIA

ABSTRACT

Positive integers which have precisely 3 prime divisors are called as '3-PrimeFactors numbers'. As they inherit unidentified distribution pattern from primes from which they are built, their study from all perspectives, like primes, becomes necessary. Using decimal number system, occurrence analysis of all digit combinations in units and tens places of 3-PrimeFactors numbers is presented here. The range chosen covers all numbers having upto as many as 12 significant decimal digits.

Keywords: Prime number; k-PrimeFactors number; 3-PrimeFactors number; Digits in units and tens places.

1. INTRODUCTION

Addition, subtraction and multiplication are the operation which when done with integers, yield again integers. Of these, addition and multiplication exhibit this property with positive integers also. The fourth arithmetic operation, which doesn't show this property with integers, is division. Division of integers need not always give integers. So we are interested in cases when it actually does. What is talked about is divisibility. Positive integers which are primitive in divisibility don't have any non-trivial divisors and are well-known as prime numbers [1]. This is a class of numbers having innocent looking definition but forming most mysterious structure due to hitherto unknown precise properties. Not only they or their types are special, but many other classes of numbers that are based on them become equally so; like the newly defined next one.

Definition (k-PrimeFactors Number) [6] : For any integer $k \geq 0$, a positive integer having k number of prime divisors, not be necessarily distinct, is called as k-PrimeFactors number.

0- PrimeFactors number is unique and it is unity (1). It enjoys many exclusive properties, like being multiplicative identity, being self reciprocal, own factorial, own square, own cube, having unique positive divisor as self and many more so.

1- PrimeFactor numbers are primes themselves, and as mentioned above, are studied in two primary ways : exclusively in specific large ranges [3] or asymptotically in arbitrary ranges. Their types are no exception for first [4] as well as second approach of study.

2- PrimeFactors numbers have been recently considered for their minimum [6] and maximum densities [7], minimum [8] and maximum spacings [9] between their successive members, digits in their units [10] place and units & tens places [11].

Similar analysis of 3-PrimeFactors numbers done till now has thrown light on their minimum [12] and maximum densities [13], minimum [14] and maximum spacings [15] between their successive members as well as digits in their units [16] place.

The approach in getting results about previous two types of numbers was based on first generating usual primes by using efficient algorithms [2] on modern electronic computers running Java programming language [5].

2. DIGITS IN UNITS AND TENS PLACES OF 3-PRIMEFACTORS NUMBERS

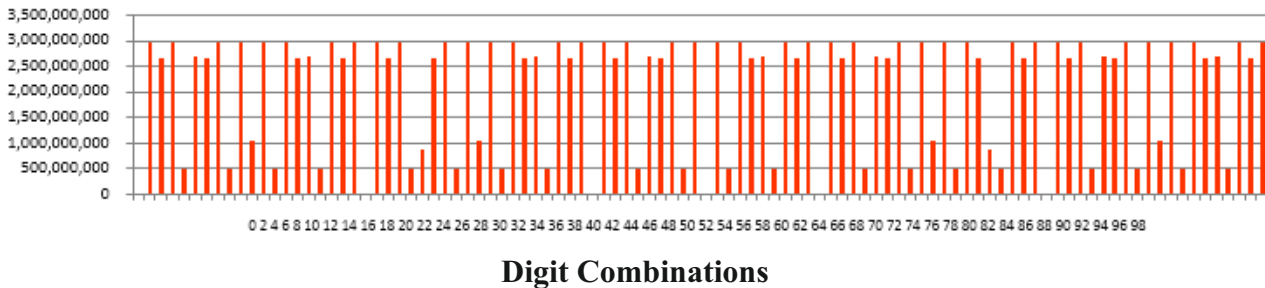
With due respect to decimal number system with base 10 in use almost everywhere, in this work we have exhaustively determined and analysed digits in units and tens places of 3- prime factors numbers till 1 trillion.

<i>Digit in Units & Tens Places</i>	<i>Number of 3-PrimeFactors Numbers < 10¹²</i>	<i>Digit in Units & Tens Places</i>	<i>Number of 3-PrimeFactorsNumbers < 10¹²</i>
0	0	10	1,02,95,17,130
1	2,97,92,77,532	11	2,97,92,20,747
2	2,63,69,78,935	12	49,60,02,799
3	2,97,92,57,921	13	2,97,92,65,188
4	49,60,06,640	14	2,63,69,92,085
5	2,69,43,69,122	15	2,69,43,51,460
6	2,63,69,88,732	16	49,60,04,725
7	2,97,92,56,767	17	2,97,92,33,988
8	49,60,06,148	18	2,63,69,86,216
9	2,97,92,51,044	19	2,97,93,44,677
20	1	30	1,02,95,09,448
21	2,97,92,69,874	31	2,97,92,45,251
22	2,63,69,91,893	32	49,60,12,087
23	2,97,92,85,991	33	2,97,92,57,037
24	49,60,02,832	34	2,63,69,85,706
25	85,59,72,440	35	2,69,43,56,492
26	2,63,69,99,010	36	49,60,07,407
27	2,97,92,79,986	37	2,97,92,61,670
28	49,60,04,412	38	2,63,69,78,877
29	2,97,92,51,241	39	2,97,92,54,721
40	0	50	1
41	2,97,92,47,090	51	2,97,92,85,045
42	2,63,69,99,481	52	49,60,00,355
43	2,97,92,71,742	53	2,97,92,34,797
44	49,60,01,886	54	2,63,69,65,210
45	2,69,43,62,312	55	2,69,43,55,069
46	2,63,69,73,256	56	49,60,05,188
47	2,97,92,68,772	57	2,97,92,18,076
48	49,60,05,686	58	2,63,69,92,557
49	2,97,92,41,332	59	2,97,92,70,725
60	0	70	1,02,95,18,337
61	2,97,92,38,768	71	2,97,92,70,931
62	2,63,70,04,656	72	49,59,98,098
63	2,97,92,42,871	73	2,97,92,49,073
64	49,59,98,951	74	2,63,69,74,877
65	2,69,43,51,957	75	85,59,82,992
66	2,63,69,67,593	76	49,60,04,785
67	2,97,92,68,785	77	2,97,92,71,043
68	49,59,97,535	78	2,63,69,76,708
69	2,97,91,78,247	79	2,97,92,28,844
80	0	90	1,02,95,09,896
81	2,97,92,06,905	91	2,97,92,34,449
82	2,63,69,91,954	92	49,60,06,703
83	2,97,92,44,360	93	2,97,92,66,680
84	49,60,05,945	94	2,63,69,95,736

85	2,69,43,64,378	95	2,69,43,62,295
86	2,63,69,75,992	96	49,59,97,643
87	2,97,92,73,217	97	2,97,93,00,256
88	49,60,09,779	98	2,63,70,41,411
89	2,97,92,85,478	99	2,97,92,47,971

Their graphical comparison follows.

Number of 3-PrimeFactors Numbers less than 1 Trillion with Different Digits in Tens and Units Places



3. RANGE-WISE DIGITS IN UNITS & TENS PLACES OF 3-PRIMEFACTORS NUMBERS

In the earlier section the values given were for whole range of 1 trillion. Here, instead for more details, the same values are given in gradually increasing ranges.

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		0	1	2	3	4
1	$<10^1$	0	0	0	0	0
2	$<10^2$	0	0	0	0	0
3	$<10^3$	0	1	4	2	2
4	$<10^4$	0	17	38	19	18
5	$<10^5$	0	240	399	247	140
6	$<10^6$	0	2,689	3,831	2,679	1,083
7	$<10^7$	0	28,335	35,899	28,434	9,143
8	$<10^8$	0	2,93,157	3,35,739	2,94,075	78,304
9	$<10^9$	0	29,85,017	31,43,990	29,86,202	6,83,786
10	$<10^{10}$	0	3,00,13,100	2,95,38,802	3,00,10,034	60,71,825
11	$<10^{11}$	0	29,97,05,140	27,86,31,343	29,97,09,710	5,45,99,594
12	$<10^{12}$	0	2,97,92,77,532	2,63,69,78,935	2,97,92,57,921	49,60,06,640

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		5	6	7	8	9
1	$<10^1$	0	0	0	1	0
2	$<10^2$	0	0	0	1	0
3	$<10^3$	4	5	2	3	2
4	$<10^4$	45	42	21	19	20
5	$<10^5$	412	416	239	142	231
6	$<10^6$	3,956	3,823	2,657	1,106	2,673
7	$<10^7$	36,943	35,778	28,286	9,177	28,252
8	$<10^8$	3,44,865	3,35,393	2,93,978	78,378	2,93,127

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		5	6	7	8	9
9	$<10^9$	32,25,275	31,42,120	29,86,417	6,84,137	29,84,813
10	$<10^{10}$	3,02,70,258	2,95,43,136	3,00,14,436	60,71,966	3,00,13,125
11	$<10^{11}$	28,50,81,376	27,86,34,770	29,97,10,850	5,45,99,424	29,97,25,113
12	$<10^{12}$	2,69,43,69,122	2,63,69,88,732	2,97,92,56,767	49,60,06,148	2,97,92,51,044

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		10	11	12	13	14
1	$<10^2$	0	0	1	0	0
2	$<10^3$	5	1	3	0	2
3	$<10^4$	40	18	18	21	36
4	$<10^5$	306	234	130	246	409
5	$<10^6$	2,387	2,687	1,096	2,621	3,842
6	$<10^7$	19,617	28,435	9,202	28,340	35,950
7	$<10^8$	1,66,104	2,93,533	78,429	2,93,982	3,35,541
8	$<10^9$	14,40,298	29,85,357	6,84,117	29,85,638	31,42,202
9	$<10^{10}$	1,27,11,386	3,00,13,064	60,72,391	3,00,13,048	2,95,39,039
10	$<10^{11}$	11,37,61,519	29,97,08,797	5,45,98,574	29,97,14,164	27,86,22,227
11	$<10^{12}$	1,02,95,17,130	2,97,92,20,747	49,60,02,799	2,97,92,65,188	2,63,69,92,085

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		15	16	17	18	19
1	$<10^2$	0	0	0	1	0
2	$<10^3$	3	4	1	5	0
3	$<10^4$	40	18	15	38	22
4	$<10^5$	408	136	238	403	241
5	$<10^6$	3,889	1,088	2,687	3,844	2,714
6	$<10^7$	36,798	9,144	28,382	36,011	28,435
7	$<10^8$	3,44,437	78,228	2,93,572	3,35,928	2,93,704
8	$<10^9$	32,24,399	6,84,132	29,86,880	31,42,906	29,84,113
9	$<10^{10}$	3,02,65,960	60,73,001	3,00,14,499	2,95,42,065	3,00,07,825
10	$<10^{11}$	28,50,75,931	5,45,98,030	29,96,94,477	27,86,23,162	29,97,11,299
11	$<10^{12}$	2,69,43,51,460	49,60,04,725	2,97,92,33,988	2,63,69,86,216	2,97,93,44,677

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		20	21	22	23	24
1	$<10^2$	1	0	0	0	0
2	$<10^3$	1	0	4	1	3
3	$<10^4$	1	22	42	25	20
4	$<10^5$	1	246	425	244	139
5	$<10^6$	1	2,633	3,854	2,692	1,098
6	$<10^7$	1	28,381	35,981	28,371	9,109
7	$<10^8$	1	2,93,726	3,35,392	2,93,641	78,235
8	$<10^9$	1	29,84,032	31,41,916	29,84,496	6,83,920
9	$<10^{10}$	1	3,00,04,755	2,95,44,076	3,00,16,057	60,72,530
10	$<10^{11}$	1	29,97,16,340	27,86,39,418	29,97,13,131	5,45,98,675
11	$<10^{12}$	1	2,97,92,69,874	2,63,69,91,893	2,97,92,85,991	49,60,02,832

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		25	26	27	28	29
1	$<10^2$	0	0	1	1	0
2	$<10^3$	5	2	3	3	1
3	$<10^4$	37	42	13	17	20

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		25	26	27	28	29
4	$<10^5$	269	405	252	140	243
5	$<10^6$	2,085	3,886	2,705	1,107	2,614
6	$<10^7$	16,900	35,844	28,469	9,166	28,279
7	$<10^8$	1,41,502	3,34,948	2,93,682	78,289	2,93,569
8	$<10^9$	12,16,687	31,42,448	29,85,056	6,84,091	29,85,466
9	$<10^{10}$	1,06,67,607	2,95,40,298	3,00,12,609	60,72,928	3,00,14,422
10	$<10^{11}$	9,49,79,097	27,86,32,110	29,97,05,115	5,45,97,945	29,96,98,723
11	$<10^{12}$	85,59,72,440	2,63,69,99,010	2,97,92,79,986	49,60,04,412	2,97,92,51,241

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		30	31	32	33	34
1	$<10^2$	1	0	0	0	0
2	$<10^3$	7	3	2	2	3
3	$<10^4$	42	20	20	18	40
4	$<10^5$	310	245	140	248	417
5	$<10^6$	2,402	2,670	1,111	2,692	3,878
6	$<10^7$	19,665	28,386	9,158	28,281	35,788
7	$<10^8$	1,66,230	2,93,353	78,294	2,93,875	3,35,362
8	$<10^9$	14,40,474	29,85,203	6,84,066	29,84,204	31,42,949
9	$<10^{10}$	1,27,12,499	3,00,14,591	60,71,678	3,00,14,728	2,95,38,893
10	$<10^{11}$	11,37,65,625	29,97,20,078	5,45,99,070	29,97,18,838	27,86,34,822
11	$<10^{12}$	1,02,95,09,448	2,97,92,45,251	49,60,12,087	2,97,92,57,037	2,63,69,85,706

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		35	36	37	38	39
1	$<10^2$	0	0	0	0	0
2	$<10^3$	2	2	1	6	2
3	$<10^4$	35	20	20	47	20
4	$<10^5$	398	139	248	407	245
5	$<10^6$	3,920	1,107	2,710	3,822	2,696
6	$<10^7$	36,823	9,145	28,431	35,757	28,534
7	$<10^8$	3,44,547	78,245	2,93,468	3,35,778	2,93,133
8	$<10^9$	32,24,127	6,83,951	29,83,402	31,43,638	29,83,803
9	$<10^{10}$	3,02,66,540	60,71,882	3,00,12,294	2,95,42,612	3,00,16,452
10	$<10^{11}$	28,50,66,725	5,46,00,654	29,97,21,067	27,86,33,694	29,97,07,484
11	$<10^{12}$	2,69,43,56,492	49,60,07,407	2,97,92,61,670	2,63,69,78,877	2,97,92,54,721

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		40	41	42	43	44
1	$<10^2$	0	0	1	0	1
2	$<10^3$	0	1	6	1	3
3	$<10^4$	0	16	40	23	19
4	$<10^5$	0	227	417	239	134
5	$<10^6$	0	2,647	3,807	2,666	1,102
6	$<10^7$	0	28,444	35,891	28,518	9,158
7	$<10^8$	0	2,93,366	3,35,557	2,93,500	78,297
8	$<10^9$	0	29,84,511	31,43,676	29,84,133	6,84,173
9	$<10^{10}$	0	3,00,11,088	2,95,43,339	3,00,16,060	60,71,869
10	$<10^{11}$	0	29,97,09,648	27,86,29,054	29,97,15,997	5,45,99,873
11	$<10^{12}$	0	2,97,92,47,090	2,63,69,99,481	2,97,92,71,742	49,60,01,886

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		45	46	47	48	49
1	$<10^2$	1	0	0	0	0
2	$<10^3$	5	3	3	2	1
3	$<10^4$	45	36	23	20	15
4	$<10^5$	427	402	232	139	229
5	$<10^6$	3,963	3,870	2,659	1,092	2,658
6	$<10^7$	36,954	35,890	28,489	9,151	28,417
7	$<10^8$	3,44,829	3,35,711	2,93,643	78,376	2,93,857
8	$<10^9$	32,26,240	31,42,259	29,84,197	6,83,798	29,83,802
9	$<10^{10}$	3,02,66,847	2,95,40,268	3,00,20,501	60,72,543	3,00,11,626
10	$<10^{11}$	28,50,78,692	27,86,22,118	29,97,27,187	5,46,00,232	29,97,03,401
11	$<10^{12}$	2,69,43,62,312	2,63,69,73,256	2,97,92,68,772	49,60,05,686	2,97,92,41,332

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		50	51	52	53	54
1	$<10^2$	1	0	1	0	0
2	$<10^3$	1	1	3	1	5
3	$<10^4$	1	25	19	20	47
4	$<10^5$	1	250	135	233	406
5	$<10^6$	1	2,647	1,096	2,699	3,776
6	$<10^7$	1	28,360	9,157	28,470	35,821
7	$<10^8$	1	2,93,496	78,282	2,93,741	3,35,262
8	$<10^9$	1	29,84,194	6,84,072	29,84,757	31,42,919
9	$<10^{10}$	1	3,00,12,768	60,72,952	3,00,13,376	2,95,43,227
10	$<10^{11}$	1	29,97,26,685	5,46,01,634	29,97,04,480	27,86,29,425
11	$<10^{12}$	1	2,97,92,85,045	49,60,00,355	2,97,92,34,797	2,63,69,65,210

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		55	56	57	58	59
1	$<10^2$	0	0	0	0	0
2	$<10^3$	3	3	3	2	1
3	$<10^4$	41	18	22	42	20
4	$<10^5$	419	140	232	413	244
5	$<10^6$	3,943	1,110	2,649	3,882	2,678
6	$<10^7$	36,889	9,167	28,395	35,927	28,299
7	$<10^8$	3,44,460	78,360	2,93,722	3,35,570	2,93,741
8	$<10^9$	32,24,156	6,83,632	29,84,900	31,43,305	29,85,358
9	$<10^{10}$	3,02,65,273	60,71,725	3,00,12,991	2,95,41,593	3,00,12,667
10	$<10^{11}$	28,50,67,855	5,46,00,491	29,97,23,888	27,86,36,125	29,97,16,786
11	$<10^{12}$	2,69,43,55,069	49,60,05,188	2,97,92,18,076	2,63,69,92,557	2,97,92,70,725

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		60	61	62	63	64
1	$<10^2$	0	0	0	1	0
2	$<10^3$	0	3	2	4	3
3	$<10^4$	0	18	43	16	15
4	$<10^5$	0	230	412	236	138
5	$<10^6$	0	2,669	3,846	2,705	1,097
6	$<10^7$	0	28,361	35,964	28,375	9,164
7	$<10^8$	0	2,93,543	3,35,262	2,94,202	78,318
8	$<10^9$	0	29,85,859	31,42,592	29,85,533	6,83,766
9	$<10^{10}$	0	3,00,13,173	2,95,39,772	3,00,16,963	60,71,966
10	$<10^{11}$	0	29,97,10,957	27,86,31,964	29,97,06,326	5,45,96,962

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		60	61	62	63	64
11	$<10^{12}$	0	2,97,92,38,768	2,63,70,04,656	2,97,92,42,871	49,59,98,951

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		65	66	67	68	69
1	$<10^2$	0	1	0	1	0
2	$<10^3$	3	3	1	3	2
3	$<10^4$	39	36	22	18	20
4	$<10^5$	412	394	237	136	229
5	$<10^6$	3,919	3,793	2,627	1,102	2,648
6	$<10^7$	36,855	35,935	28,532	9,155	28,560
7	$<10^8$	3,44,465	3,35,631	2,93,944	78,408	2,93,704
8	$<10^9$	32,24,647	31,43,246	29,84,678	6,84,215	29,83,806
9	$<10^{10}$	3,02,66,983	2,95,39,849	3,00,14,444	60,72,128	3,00,10,161
10	$<10^{11}$	28,50,67,237	27,86,29,545	29,97,19,667	5,45,98,568	29,97,09,784
11	$<10^{12}$	2,69,43,51,957	2,63,69,67,593	2,97,92,68,785	49,59,97,535	2,97,91,78,247

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		70	71	72	73	74
1	$<10^2$	1	0	0	0	0
2	$<10^3$	6	1	2	2	5
3	$<10^4$	46	19	17	22	44
4	$<10^5$	308	242	141	225	391
5	$<10^6$	2,411	2,728	1,108	2,662	3,811
6	$<10^7$	19,621	28,380	9,161	28,463	35,813
7	$<10^8$	1,66,211	2,93,790	78,298	2,93,115	3,35,546
8	$<10^9$	14,40,495	29,84,416	6,83,625	29,86,184	31,42,985
9	$<10^{10}$	1,27,12,314	3,00,14,889	60,71,760	3,00,17,483	2,95,42,055
10	$<10^{11}$	11,37,64,039	29,97,10,908	5,45,99,362	29,97,09,136	27,86,34,153
11	$<10^{12}$	1,02,95,18,337	2,97,92,70,931	49,59,98,098	2,97,92,49,073	2,63,69,74,877

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		75	76	77	78	79
1	$<10^2$	1	1	0	1	0
2	$<10^3$	6	1	2	4	1
3	$<10^4$	40	15	22	42	21
4	$<10^5$	280	139	221	391	238
5	$<10^6$	2,117	1,106	2,676	3,810	2,631
6	$<10^7$	16,959	9,172	28,490	36,044	28,415
7	$<10^8$	1,41,643	78,233	2,93,993	3,35,684	2,93,741
8	$<10^9$	12,16,966	6,83,816	29,84,014	31,43,961	29,85,146
9	$<10^{10}$	1,06,68,718	60,72,519	3,00,12,640	2,95,42,264	3,00,16,320
10	$<10^{11}$	9,49,82,714	5,45,97,432	29,97,15,203	27,86,33,071	29,97,19,177
11	$<10^{12}$	85,59,82,992	49,60,04,785	2,97,92,71,043	2,63,69,76,708	2,97,92,28,844

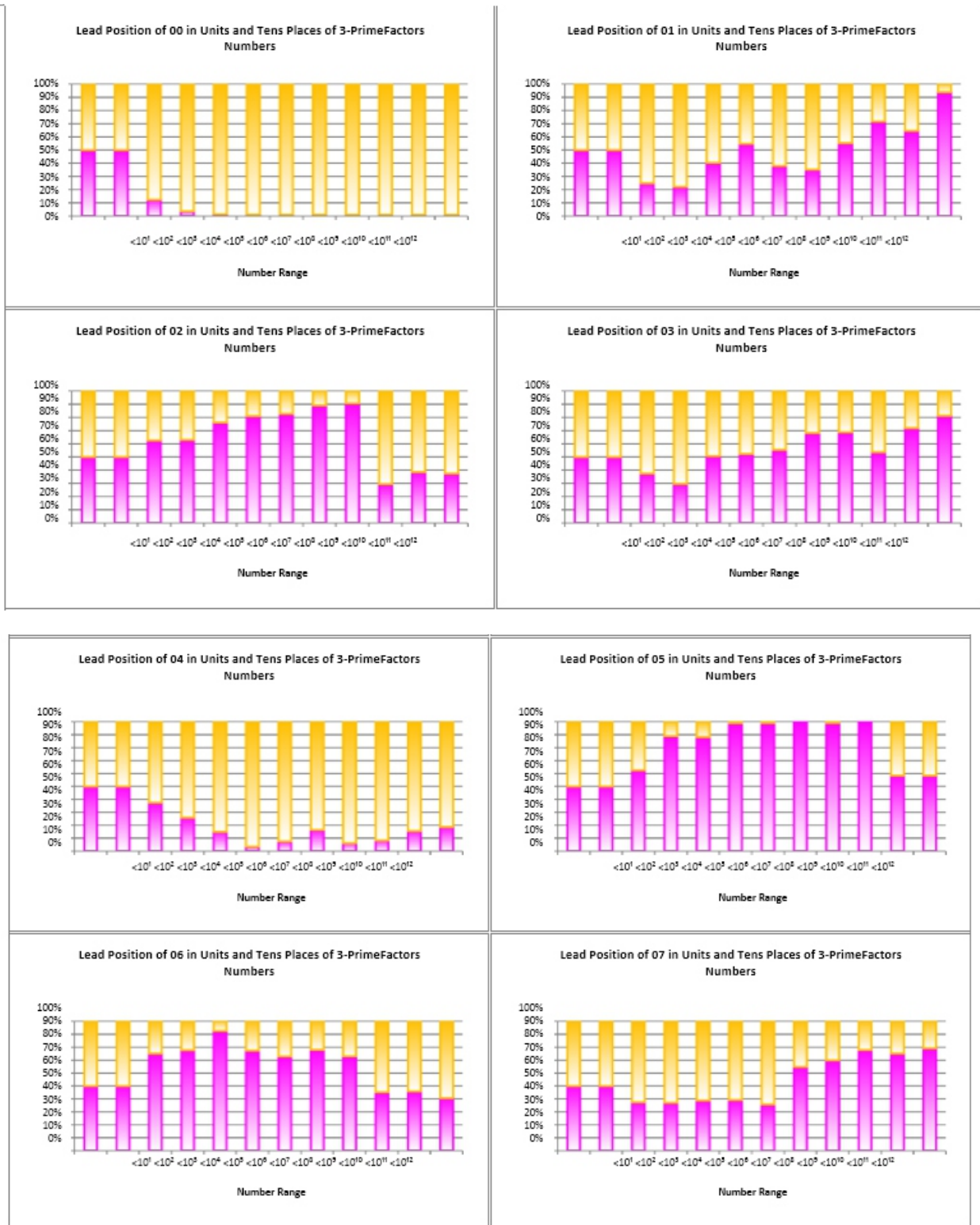
Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		80	81	82	83	84
1	$<10^2$	0	0	0	0	0
2	$<10^3$	0	1	5	1	1
3	$<10^4$	0	23	49	23	17
4	$<10^5$	0	233	397	256	138
5	$<10^6$	0	2,617	3,826	2,703	1,115
6	$<10^7$	0	28,256	35,881	28,412	9,187

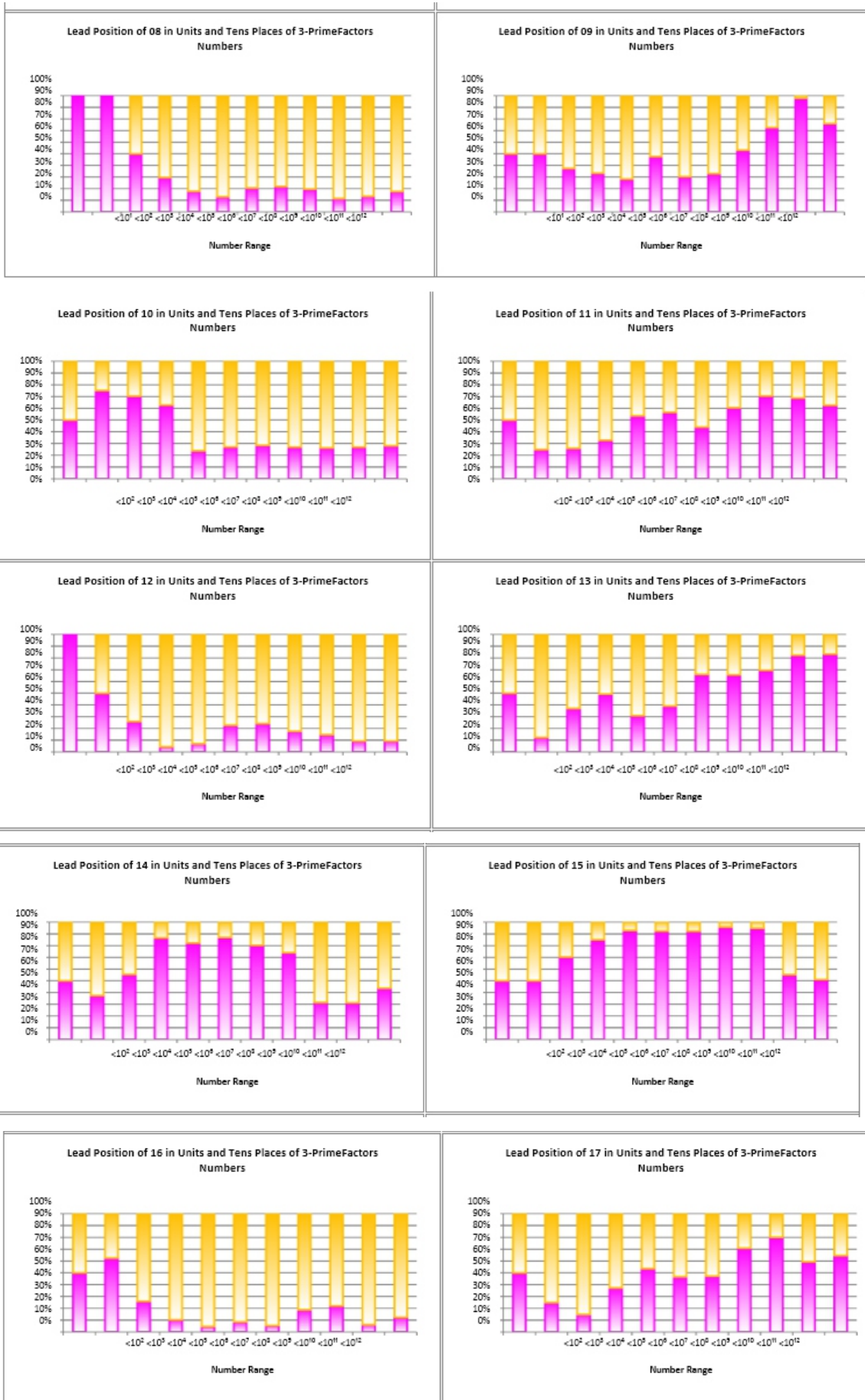
Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		80	81	82	83	84
7	$<10^8$	0	2,93,861	3,35,337	2,93,498	78,175
8	$<10^9$	0	29,84,111	31,44,680	29,83,980	6,83,891
9	$<10^{10}$	0	3,00,17,106	2,95,43,547	3,00,13,680	60,71,432
10	$<10^{11}$	0	29,97,15,517	27,86,40,230	29,97,21,210	5,45,99,539
11	$<10^{12}$	0	2,97,92,06,905	2,63,69,91,954	2,97,92,44,360	49,60,05,945

Sr. No.	Range	Number of 3-PrimeFactors Numbers with Following Digits in Units & Tens Places				
		85	86	87	88	89
1	$<10^2$	0	0	0	0	0
2	$<10^3$	3	4	2	3	0
3	$<10^4$	46	44	21	21	19
4	$<10^5$	415	395	250	138	244
5	$<10^6$	3,919	3,812	2,712	1,111	2,669

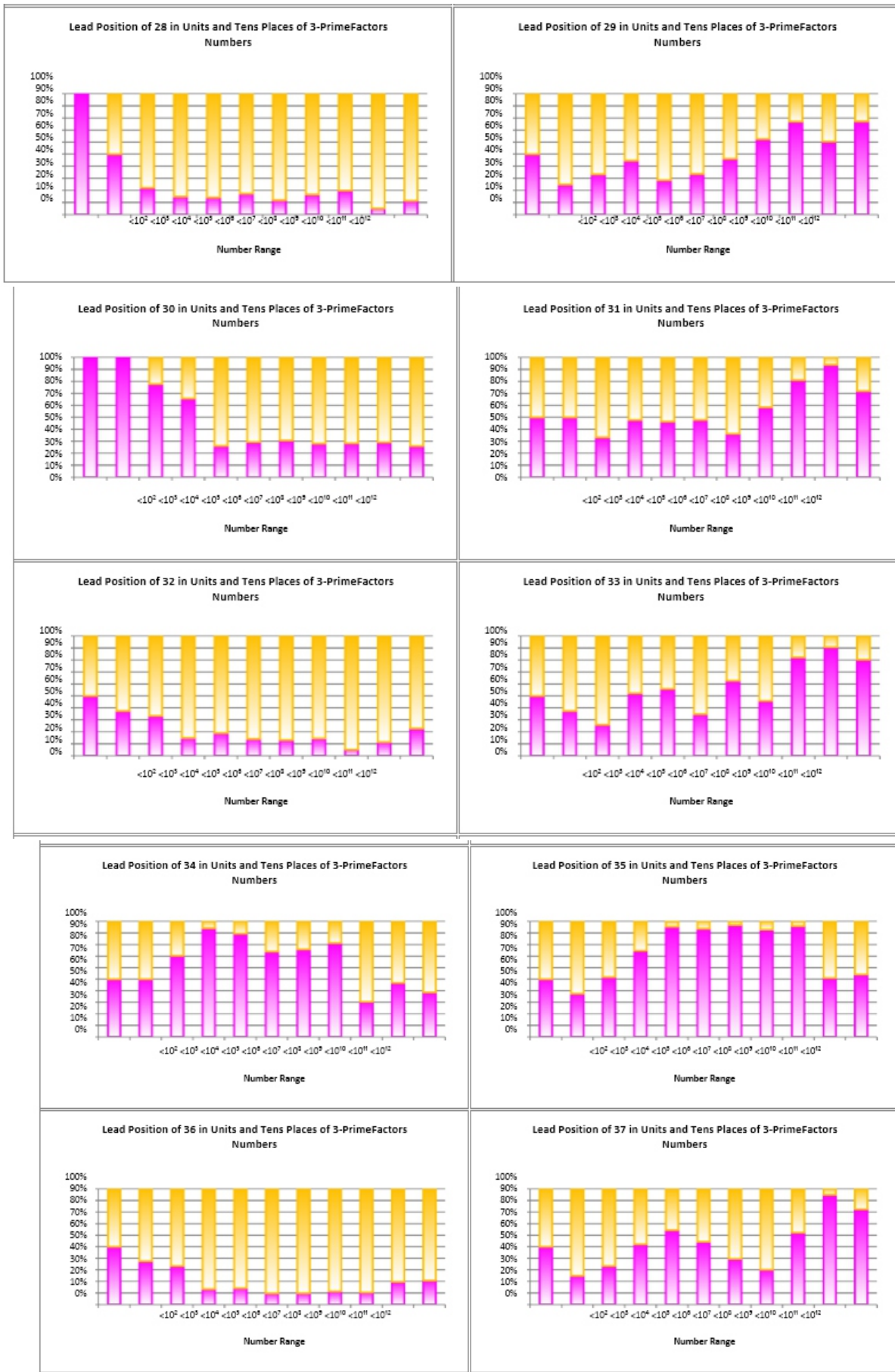
4. RANGE-WISE LEAD POSITIONS OF DIGITS IN UNITS & TENS PLACES OF 3-PRIMEFACTORS NUMBERS

The 1110 values determined rigorously above give the range-wise lead positions in percentages of each combination of digits in units and tens places of 3-Prime Factors numbers till 1 trillion.







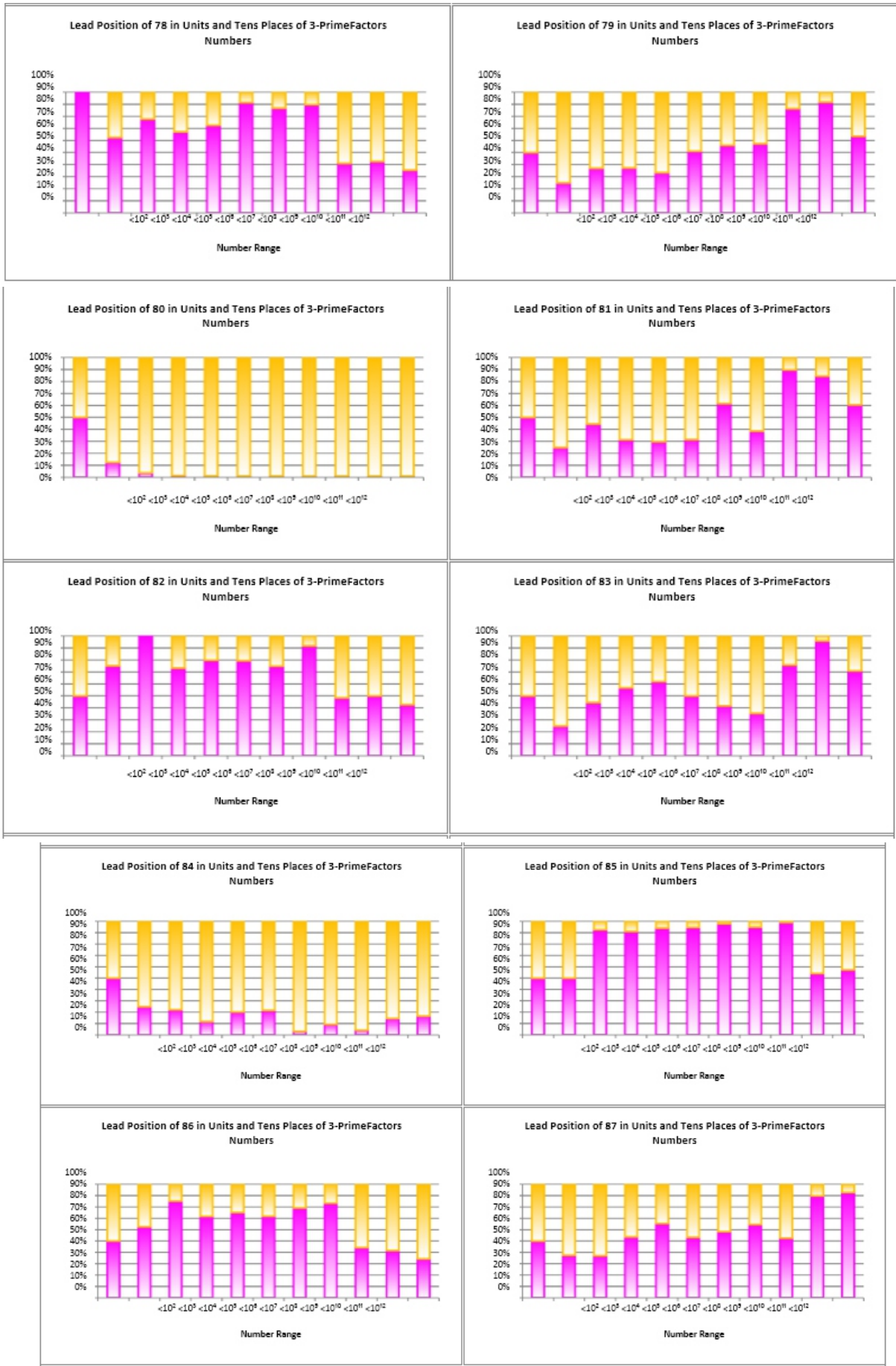


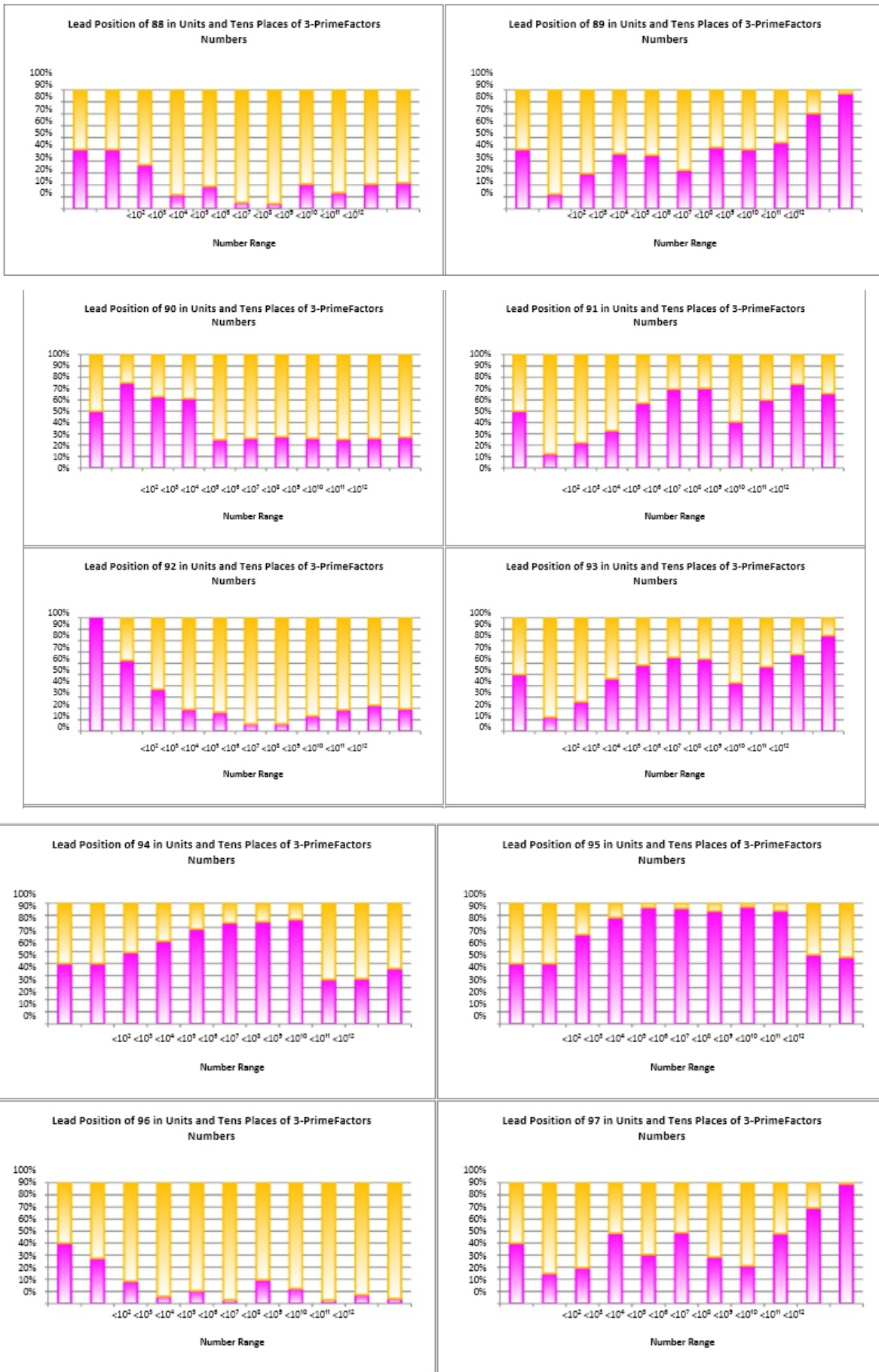


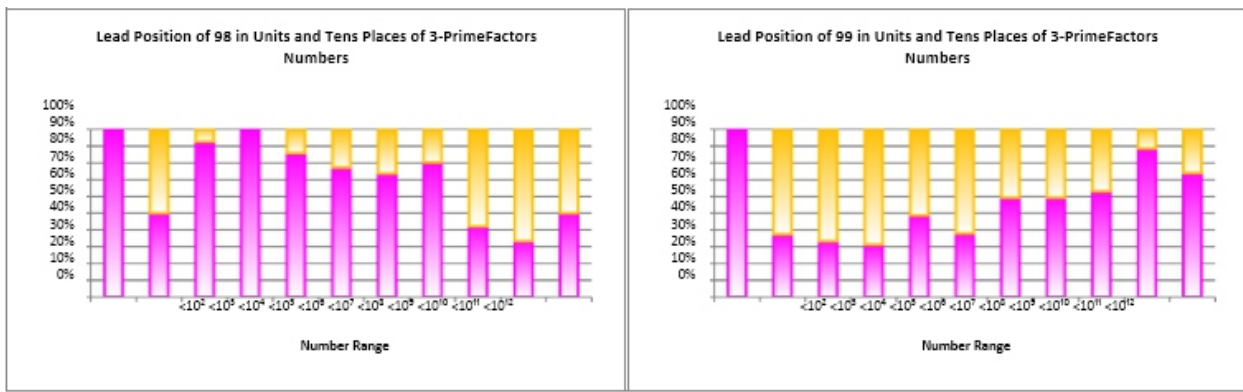












5. PATTERNS OF DIGITS IN UNITS & TENS PLACES OF 3-PRIMEFACTORS NUMBERS

All the above values indicate following properties so long as the decimal digits in units and tens places of 3-PrimeFactors numbers are concerned.

The 4 digit combinations 00, 40, 60 and 80 never occur in units and tens places of 3- PrimeFactors numbers.

The digit combinations 20 & 50 occur only once, that too right at these numbers and not later. The other digit combination which are perfectly divisible by 4 viz., 04, 08, 12, 16, 24, 28, 32, 36, 44, 48, 52, 56, 64, 68, 72, 76, 84, 88, 92 and 96 occur at next higher class about 0.049% of total numbers.

The other 2 digit combinations which are multiples of 25, viz., 25 and 75 stand at next higher level with percentage of around 0.0855, near double of their above value.

The remaining digit combinations of multiples of 10, viz., 10, 30, 70 and 90 are having next higher percentage of occurrence of about 0.10.

All other remaining digit combinations dominate occurrence percentage with values from 0.26 to 0.29.

ACKNOWLEDGEMENT

The author gratefully mentions continuous uninterrupted use of all computers in Laboratory of the Department of Mathematics & Statistics of his host institute for several months to get these results. The software tools of Java programming language, NetBeans IDE and Microsoft Excel were instrumental in making computers work and their development teams are deeply acknowledged.

Finally the author will like to thank anonymous referee(s) of this paper.

REFERENCES

- [1] Benjamin Fine, Gerhard Rosenberger, "Number Theory: An Introduction via the Distribution of Primes", Birkhauser, 2007.
- [2] Neeraj Anant Pande, "Improved Prime Generating Algorithms by Skipping Composite Divisors and Even Numbers (Other Than 2)", *Journal of Science and Arts*, Year 15, No.2 (31), 135-142, 2015.
- [3] Neeraj Anant Pande, "Analysis of Primes Less Than a Trillion", *International Journal of Computer Science & Engineering Technology*, Vol. 6, No. 06, 332-341, 2015.
- [4] Neeraj Anant Pande, "Analysis of Twin Primes Less Than a Trillion", *Journal of Science and Arts*, Year 16, No.4 (37), 279-288, 2016.
- [5] Herbert Schildt, "Java : The Complete Reference, 7th Edition", Tata Mc-Graw Hill, 2007.

- [6] Neeraj Anant Pande, "Low Density Distribution of 2-PrimeFactors Numbers till 1 Trillion", *Journal of Research in Applied Mathematics*, 2017, Vol. 3, Issue 8, 35-47, 2017.
- [7] Neeraj Anant Pande, "High Density Distribution of 2-PrimeFactors Numbers till 1 Trillion", *American International Journal of Research in Formal, Applied & Natural Sciences, Communicated*, 2017.
- [8] Neeraj Anant Pande, "Minimum Spacings between 2-PrimeFactors Numbers till 1 Trillion", *Journal of Computer and Mathematical Sciences*, Vol. 8 (12), 769-780, 2017.
- [9] Neeraj Anant Pande, "Maximum Spacings between 2-PrimeFactors Numbers till 1 Trillion", *International Journal of Mathematics Trends and Technology*, Volume 52, Issue 5, 311-321, 2017.
- [10] Neeraj Anant Pande, "Digits in Units Place of 2-PrimeFactors Numbers till 1 Trillion", *International Journal of Mathematics And its Applications*, Accepted, 2017.
- [11] Neeraj Anant Pande, "Digits in Units and Tens Place of 2-PrimeFactors Numbers till 1 Trillion", *International Journal of Engineering, Science and Mathematics*, Volume 6, Issue 8, 254-273, 2017.
- [12] Neeraj Anant Pande, "Low Density Distribution of 3-PrimeFactors Numbers till 1 Trillion", *International Journal of Latest Engineering Research and Applications*, Accepted, 2017.
- [13] Neeraj Anant Pande, "High Density Distribution of 3-PrimeFactors Numbers till 1 Trillion", *International Journal of Mathematics and Statistics Invention*, Communicated, 2017.
- [14] Neeraj Anant Pande, "Minimum Spacings between 3-PrimeFactors Numbers till 1 Trillion", *Journal of Research in Applied Mathematics, Communicated*, 2017.
- [15] Neeraj Anant Pande, "Maximum Spacings between 3-PrimeFactors Numbers till 1 Trillion", *Journal of Computer and Mathematical Sciences, Communicated*, 2017.
- [16] Neeraj Anant Pande, "Digits in Units Place of 3-PrimeFactors Numbers till 1 Trillion", *International Journal of Mathematics Trends and Technology, Communicated*, 2017.

Skills Required for Web Developer

Dr. Sapna Nagpal

Assistant Professor in Computer Science PG Govt.
College, Tigaon (Faridabad)

ABSTRACT

There is huge unemployment in India. People are getting degrees but not jobs. This paper emphasizes the area of web development which has huge job requirements due to increasing web drive in the present time. Web presence today has become a vital element for an organization striving to fulfil its objectives in an efficient manner. And equally important, if not more, in the mantra of success is the role of the web service provider. Setting up a powerful web presence is a voyage of discovery. This research work focuses on the web technologies used for coding the websites at the front-end as well as the back-end. While implementing a web application, designing part has various options, this work has endeavoured to find the mostly used language platform for web application development.

1. INTRODUCTION

In today's vast range of technology, language and platform choices, it can be very difficult to figure out where to best invest time in training your skills as a web developer. To secure a tech-hefty graduate job such as web developer, web designer or graphics designer, you'll need to ensure you have the web designing languages and other technical/Interface skills, the employers want. You may have a degree in computer science, but you may still not able to fit the expectations of your prospective employer. Degree alone is not sufficient to skill a person but one should invest time to develop skills that are required as per the latest technologies.

Web development can range from developing the simplest static single page of plain text to the most complex web-based internet applications, electronic businesses, and social network services. The different areas of web design include web graphicdesign, interface design, authoring, including standardised code and proprietary software; user experience design; and search engine optimization and so on. A Web Designer earns an average salary of Rs 231,555 per year in India as per a survey done on 1572 individuals in January 2017 [3]. A skill in Web Development is associated with high pay for this job. People in this job generally don't have more than 10 years' experience. Experience strongly influences income for this job. Web design is as much a science as it is an art form. While half of the job is based on sound coding and design know-how, the other half is based on just having an intuitive sense of what looks good and what doesn't.

2. WEB DEVELOPMENT PROCESS AND EXPERIMENT RESULTS

This research paper is focused as to which web application platforms are used for designing of websites, website hosting server and other development services. Data was collected from various web development companies working in Delhi and NCR. The people contacted were the Director or HR person who were kind enough to share their data for this research work. The main focus was to find which technologies are sought most while recruiting new web designers/programmers. The results found from that data are represented here.

Web development is all about communication and exchange of data between client and server over HTTP protocol. Client-side programming is writing code that will run on the client, it deals with the user interface with which the user interacts in the web. The coding is mainly done in any scripting language to design interactive web pages, interact with temporary storage, works as an interface between user and server, sends requests to the server and retrieval of data from server. On the other side, the sever side programming is the kind of program that runs directly on the server and deals with dynamic content as most web pages deal with searching databases. The programming on this end deals with processing the user input, displaying the requested pages, structuring web applications, interaction with servers/storages, interaction with databases, querying the database, encoding of data into HTML and operations over databases like delete, update.

There are various web coding languages available for web programming for client as well as server. The data was collected for both client and server-side application development technologies. The responses received are shown in the figure1 and figure 2. The main platforms used for client-side web development are HTML, Java Script and PHP. CSS is another choice for client-side programming. In very few cases Java, ASP.NET and Web 2.0 are used for website development at client-side. On the other hand, PHP is the main choice for sever-side web application development. Java and ASP.NET are the next choices for the server-side and Java Script follows these. HTML, CSS and Web 2.0 are used in very rare case for this end.

1. Which technology is used to develop Client side Web application?

8 responses

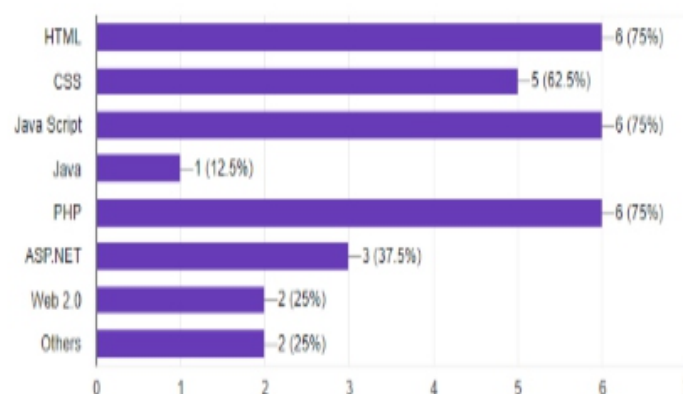


Figure 1: Languages used to develop Client-side Web applications.

2. Which technology is used to develop Server side Web application?

8 responses

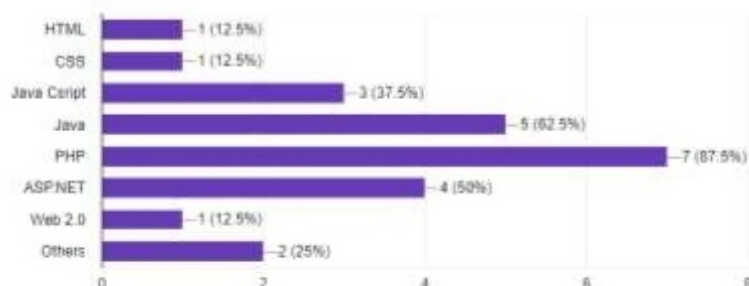


Figure 2: Languages used to develop Server-side Web applications.

The back-end database is another platform used along with dynamic websites which deals with security, structure and content management. MySQL is the first choice for database management in back end [Figure 3]. Oracle and other database management language are not among the mostly used platforms.

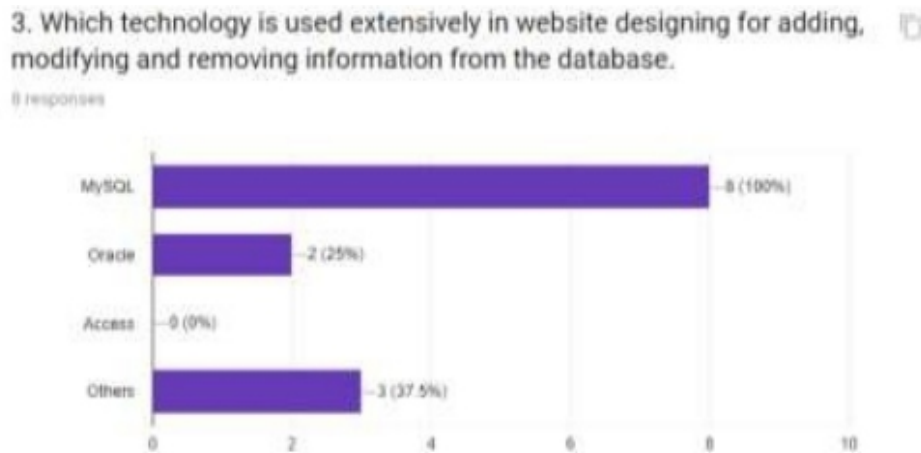


Figure 3: Technology used for Database Management in Websites.

Having choices is always a good thing. But sometimes it's supportive to have a chaperon, so the next part is to help cut the noise i.e. which guides the beginners to learn the skills they need to get their first job in web development. The person interested in making their career in website development should invest in mastering which skill which is sought most to improve their employability. The survey results yield the output that PHP is at the peak among all other technologies for hiring web programmers. HTML, Java Script, Java come after that. CSS, Web 2.0 and others stand at the last in securing the web designer job (Figure 4).

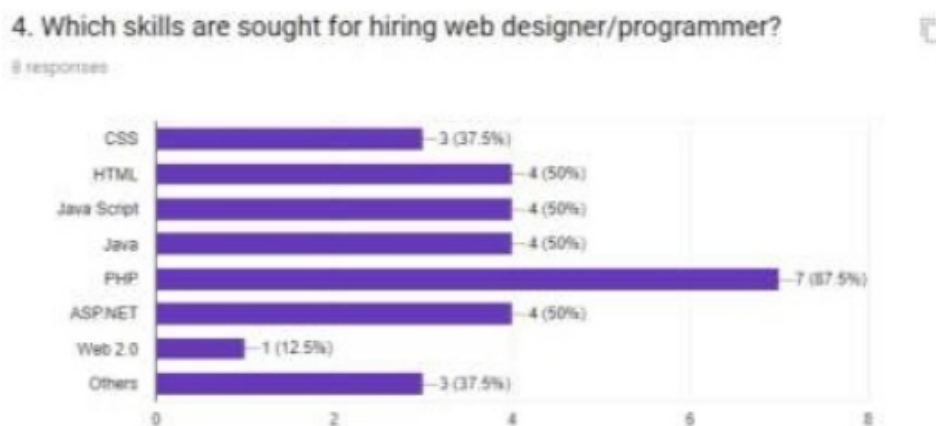


Figure 4: Major Technologies sought while hiring Web Designers.

Major criteria for hiring the web designers is Interview including both technical and HR interview, which is conducted in 75% of the cases whereas written test is another choice to sieve the candidates (Figure 5). Generally, persons with less than 5 years of experience are most sought in this profession. Sometimes freshers (25% of cases) are also considered for new recruitments. Since, in this profession there are very few people with more than 10 years of experience, so the demand is for less experienced persons which can be trained as per the demands of the company.

5. What is the criteria for selection of web designer.

8 responses

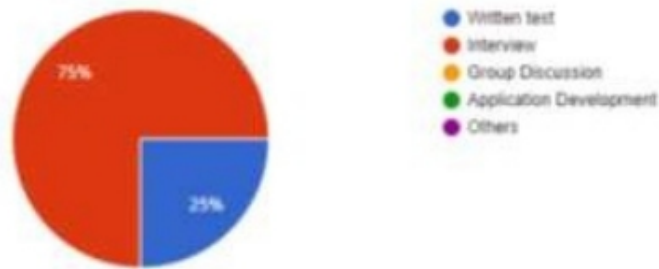


Figure 5: The Criteria for hiring Web Designers.

6. Which type of person is preferred?

8 responses

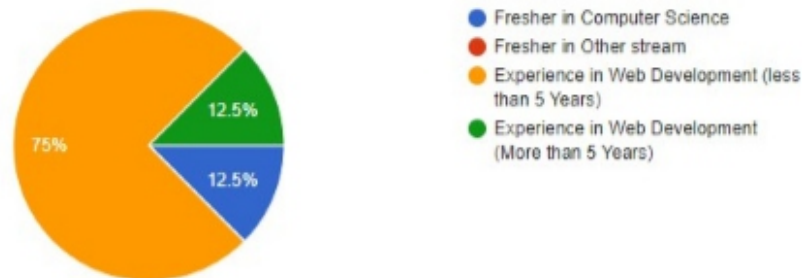


Figure 6: The Experience in Recruitment of Web Designer/ Programmer.

3. DISCUSSION

Web designers use a combination of design and IT skills to create, develop and maintain websites. They must find a balance between the visual appeal of a website and its functionality as per the specific requirements of the clients to design the underlying architecture. This research work was initiated to know as to which technologies are mostly used for website development and related solutions. A total of 8 companies in Delhi and NCR were contacted including NIC Lodhi Road, which is extending comprehensive WWW services to the Central and State Governments, ministries and departments in the areas of consultancy, web design and development, web hosting, value added web services for promotion of websites, enhancement of websites & training. The other companies which were contacted were the companies working and recruiting in Delhi and NCR region to provide web services. The results found in this research work are highlighting the major technologies/skills considered for recruiting web designers. There are definitely other skills such as empathy in communication, time and stress management, perspective and design sense which are also required in web designers. The purpose of this work was to find out the core languages/technologies used for hiring web professionals.

REFERENCES

- [1] <http://webservices.nic.in>
- [2] H. Wang, J. Yang, "Research and application of web development based on ASP.NET 2.0+Ajax", 3rd IEEE Conference on Industrial Electronics and Applications, 2008, ICIEA 2008.
- [3] https://www.payscale.com/research/IN/Job=Web_Designer_%26_Developer/Salary
- [4] Sabah Al-Fedaghi, "Developing Web Applications", International Journal of Software Engineering and Its Applications, Vol. 5 No. 2, April, 2011.

- [5] E. Mendes, "Web Development Versus Software Development", *Practitioner's Knowledge Representation*, DOI 10.1007/978-3-642-54157-5_2, #Springer-Verlag, Berlin, Heidelberg 2014.
- [6] Ch Rajesh and K S V Krishna Srikanth, "Research on HTML5 in Web Development", *International Journal of Computer Science and Information Technologies*, Vol. 5 (2), 2014, 2408-2412, ISSN:0975-9646.
- [7] Manya Sharma, "Web Development Technology-PHP. How it is related to Web Development Technology ASP.NET", *International Journal of Scientific & Technology Research Volume 4, Issue 01, January 2015*, ISSN 2277-8616.

Instructions for Authors

Essentials for Publishing in this Journal

- 1 Submitted articles should not have been previously published or be currently under consideration for publication elsewhere.
- 2 Conference papers may only be submitted if the paper has been completely re-written (taken to mean more than 50%) and the author has cleared any necessary permission with the copyright owner if it has been previously copyrighted.
- 3 All our articles are refereed through a double-blind process.
- 4 All authors must declare they have read and agreed to the content of the submitted article and must sign a declaration correspond to the originality of the article.

Submission Process

All articles for this journal must be submitted using our online submissions system. <http://enrichedpub.com/> . Please use the Submit Your Article link in the Author Service area.

Manuscript Guidelines

The instructions to authors about the article preparation for publication in the Manuscripts are submitted online, through the e-Ur (Electronic editing) system, developed by **Enriched Publications Pvt. Ltd.** The article should contain the abstract with keywords, introduction, body, conclusion, references and the summary in English language (without heading and subheading enumeration). The article length should not exceed 16 pages of A4 paper format.

Title

The title should be informative. It is in both Journal's and author's best interest to use terms suitable. For indexing and word search. If there are no such terms in the title, the author is strongly advised to add a subtitle. The title should be given in English as well. The titles precede the abstract and the summary in an appropriate language.

Letterhead Title

The letterhead title is given at a top of each page for easier identification of article copies in an Electronic form in particular. It contains the author's surname and first name initial, article title, journal title and collation (year, volume, and issue, first and last page). The journal and article titles can be given in a shortened form.

Author's Name

Full name(s) of author(s) should be used. It is advisable to give the middle initial. Names are given in their original form.

Contact Details

The postal address or the e-mail address of the author (usually of the first one if there are more Authors) is given in the footnote at the bottom of the first page.

Type of Articles

Classification of articles is a duty of the editorial staff and is of special importance. Referees and the members of the editorial staff, or section editors, can propose a category, but the editor-in-chief has the sole responsibility for their classification. Journal articles are classified as follows:

Scientific articles:

1. Original scientific paper (giving the previously unpublished results of the author's own research based on management methods).
2. Survey paper (giving an original, detailed and critical view of a research problem or an area to which the author has made a contribution visible through his self-citation);
3. Short or preliminary communication (original management paper of full format but of a smaller extent or of a preliminary character);
4. Scientific critique or forum (discussion on a particular scientific topic, based exclusively on management argumentation) and commentaries. Exceptionally, in particular areas, a scientific paper in the Journal can be in a form of a monograph or a critical edition of scientific data (historical, archival, lexicographic, bibliographic, data survey, etc.) which were unknown or hardly accessible for scientific research.

Professional articles:

1. Professional paper (contribution offering experience useful for improvement of professional practice but not necessarily based on scientific methods);
2. Informative contribution (editorial, commentary, etc.);
3. Review (of a book, software, case study, scientific event, etc.)

Language

The article should be in English. The grammar and style of the article should be of good quality. The systematized text should be without abbreviations (except standard ones). All measurements must be in SI units. The sequence of formulae is denoted in Arabic numerals in parentheses on the right-hand side.

Abstract and Summary

An abstract is a concise informative presentation of the article content for fast and accurate Evaluation of its relevance. It is both in the Editorial Office's and the author's best interest for an abstract to contain terms often used for indexing and article search. The abstract describes the purpose of the study and the methods, outlines the findings and state the conclusions. A 100- to 250-Word abstract should be placed between the title and the keywords with the body text to follow. Besides an abstract are advised to have a summary in English, at the end of the article, after the Reference list. The summary should be structured and long up to 1/10 of the article length (it is more extensive than the abstract).

Keywords

Keywords are terms or phrases showing adequately the article content for indexing and search purposes. They should be allocated heaving in mind widely accepted international sources (index, dictionary or thesaurus), such as the Web of Science keyword list for science in general. The higher their usage frequency is the better. Up to 10 keywords immediately follow the abstract and the summary, in respective languages.

Acknowledgements

The name and the number of the project or programmed within which the article was realized is given in a separate note at the bottom of the first page together with the name of the institution which financially supported the project or programmed.

Tables and Illustrations

All the captions should be in the original language as well as in English, together with the texts in illustrations if possible. Tables are typed in the same style as the text and are denoted by numerals at the top. Photographs and drawings, placed appropriately in the text, should be clear, precise and suitable for reproduction. Drawings should be created in Word or Corel.

Citation in the Text

Citation in the text must be uniform. When citing references in the text, use the reference number set in square brackets from the Reference list at the end of the article.

Footnotes

Footnotes are given at the bottom of the page with the text they refer to. They can contain less relevant details, additional explanations or used sources (e.g. scientific material, manuals). They cannot replace the cited literature.

The article should be accompanied with a cover letter with the information about the author(s): surname, middle initial, first name, and citizen personal number, rank, title, e-mail address, and affiliation address, home address including municipality, phone number in the office and at home (or a mobile phone number). The cover letter should state the type of the article and tell which illustrations are original and which are not.

Address of the Editorial Office:

Enriched Publications Pvt. Ltd.
S-9, IInd FLOOR, MLU POCKET,
MANISH ABHINAV PLAZA-II, ABOVE FEDERAL BANK,
PLOT NO-5, SECTOR -5, DWARKA, NEW DELHI, INDIA-110075,
PHONE: - + (91)-(11)-45525005