The Indian Journal of Research

# ANVIKSHIKI

## Bi-monthly International Journal of all Research

## Science

# Anvikshiki
# The Indian Journal of Research

Bi-Monthly International Journal of All Research

# Anvikshiki
# The Indian Journal of Research
Volume 6 Number 4 July 2012

## Science
### Papers

# SOCIAL NETWORK ANALYSIS OF BY MINING ENRON EMAIL DATASET

RAJESH KUMAR*

## *Declaration*

The Declaration of the author for publication of Research Paper in The Indian Journal of Research Anvikshiki  ISSN 0973-9777 Bi-monthly International  Journal of all Research: I, *Rajesh Kumar* the author of the research paper entitled SOCIAL NETWORK ANALYSIS OF BY MINING ENRON EMAIL DATASET declare that , I take the responsibility of the content and material of my paper as I myself have written it and also have read the manuscript of my paper carefully. Also, I hereby give my consent to publish my paper in Anvikshiki journal , This research paper is my original work and no part of it or it's similar version is published or has been sent for publication anywhere else. I authorise the Editorial Board of  the Journal  to modify and edit the manuscript. I also give my consent to the Editor of  Anvikshiki Journal to own the copyright of my research paper.

## *Abstract*

*This paper provides a novel algorithm for automatically extracting social hierarchy data from electronic communication behavior. The algorithm is based on data mining user behaviors to automatically analyze and catalog patterns of communications between entities in a email collection to extract social standing. The advantage to such automatic methods is that they extract relevancy between hierarchy levels and are dynamic over time. We illustrate the algorithms over real world data using the Enron Corporation's email archive. The results show great promise when compared to the corporations work chart and judicial proceeding analyzing the major players.*

## *1 Introduction*

There is a vast quantity of untapped information in any collection of electronic communication records. The recent bankruptcy scandals in publicly held companies, and the subsequent government Act have increased the need to analyze these vast stores of electronic information in order to define risk and identify any conflict of interest among the entities of a corporate household. Corporate household is 'a group of business units united or regarded united within the corporation, such as suppliers and customers whose relationships with the corporation must be captured, managed, and applied for various purposes'. The problem can be broken into three distinct phases; entity identification, entity aggregation, and transparency of inter-entity relationships. Identifying individual entities is straightforward process, but the relationships between entities, or corporate hierarchy is not a straightforward task. Corporate entity

*M.Tech in C.S., Singhania University (Rajasthan) India.

charts sometimes exist on paper, but they do not reflect the day to day reality of a large and dynamic corporation. Corporate insiders are aware of these private relationships, but can be hard to come by, especially after an investigation. This information can be automatically extracted by analyzing the email communication data from within a corporation. Link mining is a set of techniques that uses different types of networks and their indicators to forecast or to model a linked domain. Link mining has been applied to many different areas such as money laundering, telephone fraud detection, crime detection. Here we show that customer modeling is a special case of link mining or relational learning which is based on probabilistic relational models such as those presented by. In general models classify each entity independently according to its attributes. Probabilistic relational models classify entities taking into account the joint probability among them. The application of link mining to corporate communication is of course limited by restrictions to disseminate internal corporate data. Thus testing algorithms against real world data is hard to come by. An exception to this situation is the publicly available Enron email dataset. The Enron Corporation's email collection described in section 2, is a publicly available set of private corporate data released during the judicial proceedings against the Enron corporation. Several researchers have explored it mostly from a Natural Language Processing (NLP) perspective]. Social network analysis (SNA) examining structural features has also been applied to extract properties of the Enron network and attempts to detect the key players around the time of Enron's crisis; studied the patterns of communication of Enron employees differentiated by their hierarchical level; interestingly enough found that word use changed according to the functional position, while conducted a thread analysis to find out employees' responsiveness. The work presented in this paper differs in two major ways. First, the relationship between any two users are calculated based on behavior patterns of each specific user not just links. This allows the algorithm to judge the strength of communication links between users based on their overall communication pattern. Second, we assume a corporate householding perspective and propose a methodology to solve the problem of transparency of inter-entity relationships in an automatic fashion. Our approach determines link mining metrics which can reproduce approximate social hierarchy within an organization or a corporate household, and rank its members. We use our metric to analyze email flows within an organization to extract social hierarchy. We analyze the behavior of the communication patterns without having to take into account the actual contents of the email messages. By performing behavior analysis and determining the communication patterns we are able to automatically:

♦ Rank the major officers of an organization.
♦ Group similarly ranked and connected users in order to accurately reproduce the organizational structure in question.
♦ Understand relationship strengths between specific sets of users.

The paper is organized as follows: Section 2 describes the Enron email corpus, section 3 presents the methods used to rank the Enron's officers; section 4 presents the results; section 4 discusses the results, and section 5 presents the conclusions.

## 2 Enron Antecedents and Data

The Enron email data set is a rich source of information showcasing the internal working of a real corporation over a period between 1998-2002. There seems to be multiple versions of the "official" Enron email data set in the literature. In the midst of Enron's legal troubles in 2002, the Federal Energy Regulatory Commission (FERC) made a dataset of 619,449 emails from 158 Enron employees available to the public removing all attachment data. Cohen first put up the raw email files for researchers in 2004, the format was mbox style with each message in its own text file. Following this, a number of research groups around the country obtained and manipulated the dataset in a variety of ways in attempts

SOCIAL NETWORK ANALYSIS OF BY MINING ENRON EMAIL DATASET

to correct inconsistencies and integrity issues within the dataset. The ISI treatment of the Enron corpus consisted of deleting extraneous, unneeded emails and fixing some anomalies in the collection data having to do with empty or illegal user email names and bounced emails messages. In addition duplicates and blank emails were removed. It should be noted that there is indication that a significant number of emails were lost either in converting the Enron data set or through specific deletion of key emails. So although we are working with most of the emails, we will make the assumption that the algorithm is robust although some emails are not part of the analysis. In addition the FERC dataset only covers about 92% of Enron employees at the time.

### 3 SNA Algorithm

The social network analysis algorithm works as follows: For each email user in the dataset analyze and calculate several statistics for each feature of each user. The individual features are normalized and used in a probabilistic framework with which users can be measured against one another for the purposes of ranking and grouping. It should be noted that the list of email users in the dataset represents a wide array of employee positions within the organization or across organizational departments. Two sets of statistics are involved in making the decision about a given user's "importance." First, we collect information pertaining to the flow of information, both volumetric and temporal. Here we count the number of emails a user has sent and received in addition to calculating what we call the average response time for emails. This is, in essence, the time elapsed between a user sending an email and later receiving an email from that same user. An exchange of this nature is only considered a "response" if a received message succeeds a sent message within three business days. This restriction has been implemented to avoid inappropriately long response times caused by a user sending an email, never receiving a response, but then receiving an unrelated email from that same user after a long delay, say a week or two. These elapsed time calculations are then averaged across all "responses" received to make up the average response time. Second, we gather information about the nature of the connections formed in the communication network. Here we rank the users by analyzing cliques (maximal complete subgraphs) and other graph theoretical qualities of an email network graph built from the dataset. Using all emails in the dataset, one can construct an undirected graph, where vertices represent accounts and edges represent communication between two accounts. We build such a graph in order to find all cliques, calculate degree and centrality measures and analyze the social structure of the network. When all the cliques in the graph have been found, we can determine which users are in more cliques, which users are in larger cliques, and which users are in more important cliques. We base it on the assumption that users associated with a larger set and frequency of cliques will then be ranked higher. Finally all of the calculated statistics are normalized and combined, each with an individual contribution to an overall social score with which the users are ultimately ranked.

### 3.1 Information Flows

First and foremost, we consider the volume of information exchanged, i.e. the number of emails sent and received, to be at least a limited indicator of importance. It is fair to hypothesize that users who communicate more, should, on average, maintain more important placement in the social hierarchy of the organization. This statistic is computed by simply tallying the total number of emails sent and received by each user. Furthermore, in order to rate the importance of user i using the amount of time user j takes to respond to emails from user i, we must first hypothesize that a faster response implies that user i is more important to user j. Additionally, when we iterate and average over all j, we will

user j takes to respond to emails from user i, we must first hypothesize that a faster response implies that user i is more important to user j. Additionally, when we iterate and average over all j, we will assume that the overall importance of user i will be reflected in this overall average of his or her importance to each of the other people in the organization. In other words, if people generally respond (relatively) quickly to a specific user, we can consider that user to be (relatively) important. To compute the average response time for each account x, we collect a list of all emails sent and received to and from accounts y1 through yn, organize and group the emails by account y1 through yn, and compute the amount of time elapsed between every email sent from account x to account yj and the next email received by account x from account $y_j$. As previously mentioned, communication of this kind contributes to this value only if the next incoming email was received within three business days of the original outgoing email.

### 3.2 Communication Networks

The first step is to construct an undirected graph and find all cliques. To build this graph, an email threshold N is first decided on. Next, using all emails in the dataset, we create a vertex for each account. An undirected edge is then drawn between each pair of accounts which have exchanged at least N emails. We then employ a clique finding algorithm. This recursively finds all maximal complete subgraphs (cliques).

a. *Number of cliques* : The number of cliques that the account is contained within.

b. *Raw clique score* : A score computed using the size of a given account's clique set. Bigger cliques are worth more than smaller ones, importance increases exponentially with size. c. Weighted clique score: A score computed using the "importance" of the people in each clique. This preliminary "importance" is computed strictly from the number of emails and the average response time. Each account in a clique is given a weight proportional to its computed preliminary. The weighted clique score is then computed by adding each weighed user contribution within the clique. Here the 'importance' of the accounts in the clique raises the score of the clique. More specifically, the raw clique score R is computed with the following formula:

$R = 2^{n-1}$

where n is the number of users in the clique. The weighted clique score W is computed with the following formula:

$W = t \cdot 2^{n-1}$

where t is the time score for the given user. Finally, the following indicators are calculated for the graph G(V,E) where V = v1, v2, ..., vn is the set of vertices, E is the set of edges, and eij is the edge between vertices vi and vj :

♦ Degree centrality or degree of a vertex vi: $deg(vi) = \sum_{j} a_{ij}$ where $a_{ij}$ is an element of the adjacent matrix A of G

♦ Clustering coefficient:: $C = \frac{1}{n} \sum_{i=1}^{n} CC_i$ where $CC_i = \frac{2|e_{ij}|}{deg(v_i)(deg(v_i)-1)} : v_j \in N_i, e_{ij} \in E$

♦ Each vertex vi has a neighborhood N defined by its immediately connected neighbors: $N_i = \{v_j : e_{ij} \in E\}$

♦ Mean of shortest path length from a specific vertex to all vertices in the graph G: $L = \frac{1}{n} \sum_{i} d_i$, where $d_i \in D$ and D is the geodesic distance matrix (matrix of all shortest path between every pair of vertices) of G, and n is the number of vertices in G.

- Betweenness centrality $Bc(vi) = \sum_k \sum_j \frac{g_{kj}}{g_{kj}}$. This is the proportion of all geodesic distances of all other vertices that include vertex vi where $g_{kj}$ is the number of geodesic paths between vertices k and j that include vertex i, and $g_{kj}$ is the number of geodesic paths between k and j.
- "Hubs-and-authorities" importance: "hub" refers to the vertex $v_i$ that points to many authorities, and "authority" is a vertex $v_j$ that points to many hubs. We used the recursive algorithm that calculates the "hubs-and-authorities" importance of each vertex of a graph G(V,E).

### 3.3 The Social Score

We introduce the social score S, a normalized, scaled number between 0 and 100 which is computed for each user as a weighted combination of the number of emails, response score, average response time, clique scores, and the degree and centrality measures introduced above. The breakdown of social scores is then used to:

i. Rank users from most important to least important
ii. Group users which have similar social scores and clique connectivity
iii. Determine n different levels (or echelons) of social hierarchy within which to place all the users. This is a clustering step, and n can be bounded. The rankings, groups and echelons are used to reconstruct an organization chart as accurately as possible. To compute S , we must first scale and normalize each of the previous statistics which we have gathered. The contribution, C , of each metric is individually mapped to a [0, 100] scale and weighted with the following formula

$$W_i \cdot C_i = W_i \cdot 100 \left[ \frac{x_i - \inf x}{\sup x - \inf x} \right]$$

where x is the metric in question, is the respective weight for that metric, the sup x and inf x are computed across all I users and xi is the value for the ith user. This normalization is applied to each of the following metrics:

1. number of emails
2. average response time
3. response score
4. number of cliques
5. raw clique score
6. weighted clique score
7. degree centrality
8. clustering coefficient
9. mean of shortest path length from a specific vertex to all vertices in the graph
10. betweenness centrality
11. "Hubs-and-Authorities" importance

Finally, these weighted contributions are then normalized over the chosen weights $w_x$ to compute the social score as follows:

$$S = \sum_{all} \frac{W_i \cdot C_i}{W_i}$$

This gives us a score between 0 and 100 with which to rank every user into an overall ranked list. Our assumption is that although the number of emails, average response time, number and quality of cliques, and the degree and centrality measures are all perfectly reasonable variables in an equation for

### *3.4 Visualization*

As part of this research, we developed a graphical interface for EMT, using the JUNG library, to visualize the results of social hierarchy detection by means of email flow. After the results have been computed, the statistics calculated and the users ranked, the option to view the network is available. When this option is invoked, a hierarchical, organized version of the undirected clique graph is displayed. Nodes represent users, while edges are drawn if those two users have exchanged at least m emails. Information is provided to the user in two distinct ways, the qualities of a user are reflected in the look of each node, where the relative importance of a user is reflected in the placement of each node within the simulated organization chart. Although every node is colored red, its relative size represents its social score. The largest node representing the highest ranked individual, the smallest representing the lowest. The transparency of a given node is a reflection of the user's time score. A user boasting a time score near to 1 will render itself almost completely opaque where a user with a very low time score will render almost entirely transparent. The users are divided into one of n echelons using a grouping algorithm, we use n = 5 in this paper. Currently, the only grouping algorithm which has been implemented is a straight scale level division. Users with social scores from 80-100 are placed on the top level, users with social scores from 60-80 are placed on the next level down, etc. If the weights are chosen with this scale division in mind, only a small percentage of the users will maintain high enough social scores to inhabit the upper levels, so a tree-like organizational structure will be manifested. Different, more sophisticated, ranking and grouping algorithms have been considered and will be implemented, and will be discussed in the following section on future work. When a node is selected with the mouse, all users connected to the selected user through cliques are highlighted and the user, time score and social score populate a small table at the bottom of the interface for inspection. Nodes can be individually picked or picked as groups and rearranged at the user's discretion. If the organization is not accurate or has misrepresented the structure of the actual social hierarchy in question, the user can return to the analysis window and adjust the weights in order to emphasize importance in the correct individuals and then can recreate the visualization. If the user would prefer to analyze the network graphically with a non-hierarchical structure, more traditional graph/network visualization is available by means of the Fruchterman-Reingold node placement algorithm. This node placement algorithm will emphasize the clique structure and the connectedness of nodes in the graph rather than the hierarchical ranking scheme in the first visual layout.

### *4. Results and Discussion*

We have performed the data processing and analysis using EMT. EMT is a Java based email analysis engine built on a database back-end. The Java Universal Network/ Graph Framework (JUNG) library] is used extensively in EMT for the degree and centrality measures, and for visualization purposes. In order to showcase the accuracy of our algorithm we present the analysis of the North American West Power Traders division of Enron Corporation. As one can see in Table 1 and Figure 1, when running the code on the 54 users contained with the North American West Power Traders division we can reproduce the very top of the hierarchy with great accuracy. The transparency of the vertices in the graph visualization (Figure 1) denotes the response score of the user, a combination of the number of responses and the average response time. By our assumptions made in section three, we have determined that lower average response times infer higher importance, and appropriately, Tim Belden and Debra Davidson have fast average response times, causing more opaque colored node representations. Once we turn to

the lower ranked individuals, differences in our computed hierarchy and the official hierarchy are quite noticeable in Figure 3. As we move down the corporate ladder, the conversational flows of dissimilar employees can in fact be quite similar. Despite the discrepancies of our selections with the lower ranked officers, we find that consistently we are able to pick out the most important 2 or 3 individuals in any given subset, affording us the power to build a hierarchy from small groups up. Not only does the head of Enrons Western trading operation, Tim Belden, appear on the top of our list, both his administrative assistants appear with him. Additionally, in the first fourteen positions we are also able to identify the majority of directors, and an important number of managers and specialists. Assume the positions and their key role in the organizational structure. The placement of accounts other than the top two or three is in fact giving us insight into the true social hierarchy of this particular Enron business unit over the course of time from which the emails were gathered. This differs noticeably from the official corporate hierarchy, which can be expected as the data reflects the reality of the corporate communication structure. With this sort of technique, it may be possible to view a snapshot of a corporate community (or any number of sub-communities) and effectively determine the real relationships and connections between individuals, a set of insights an official corporate organization chart simply could not offer.

## 5 Conclusions and Future work

Although real world data is hard to come by, the Enron dataset provides an excellent starting point for these tools. When we analyzed the algorithm on our own email data the social hierarchy of our lab was very apparent. Figure 2 clearly shows professor, PhD, lab students, and outsiders The next immediate concern is to apply these tools to the Enron dataset in a comprehensive and formal manner over time based data sets. The dataset contains enough email volume and generality to provide us with very useful results if we are interested in knowing how social structure changes over time. By varying the feature weights it is possible to use the mentioned parameters to:

a. Pick out the most important individual(s) in an organization
b. Group individuals with similar social/email qualities, and
c. Graphically draw an organization chart which approximately simulates the real social hierarchy in question

In order to more completely answer our question, as previously mentioned, a number of additions and alterations to the current algorithms exist and can be tested. First, the concept of average response time can be reworked or augmented by considering the order of responses, rather than the time between responses, For example, if user a receives an email from user b before receiving an email from user c, but then promptly responds to user c before responding to user b, it should be clear that user c carries more importance (at least in the eyes of user a). Either replacing the average response time statistic with this, or introducing it as its own metric may prove quite useful. Another approach is to consider common email usage times for each user and to adjust the received time of email to the beginning of the next common email usage time. For example, if user a typically only accesses her email from 9-11am and from 2-5pm, then an email received by user a at 7pm can be assumed to have been received at 9am the next morning. We hypothesize that this might correct errors currently introduced in the average response time calculations due to different people maintaining different work schedules. In addition to the continued work on the average response time algorithms, new grouping and division algorithms are being considered. Rather than implementing the straight scale division algorithm, a more statistically sophisticated formula can be used to group users by percentile or standard deviations of common distributions. Furthermore, rather than ignoring the clique connections between users at this step, the graph edges could very well prove important in how to arrange users into five different levels of social ranking, by grouping users with respect to their connections to others.
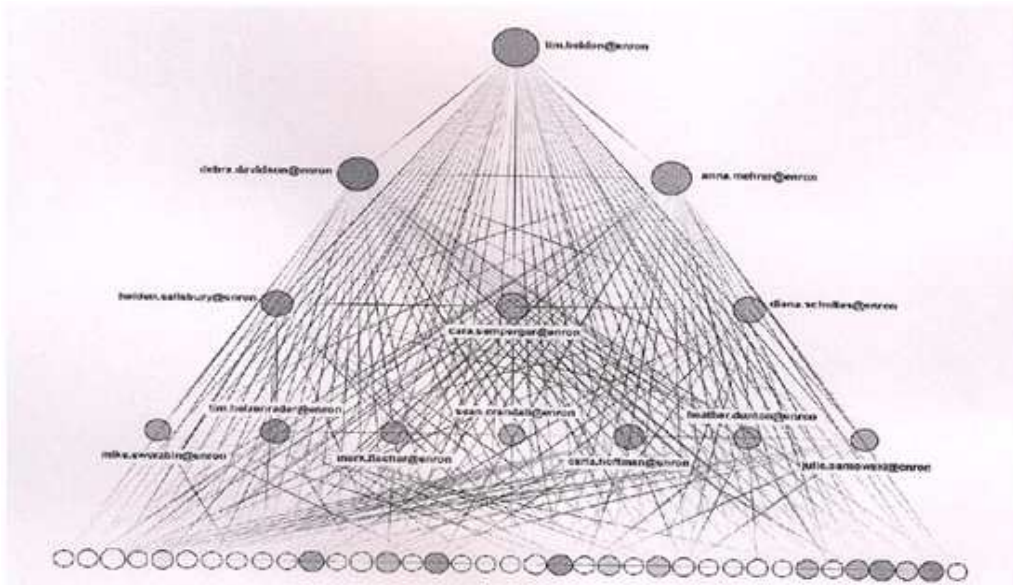
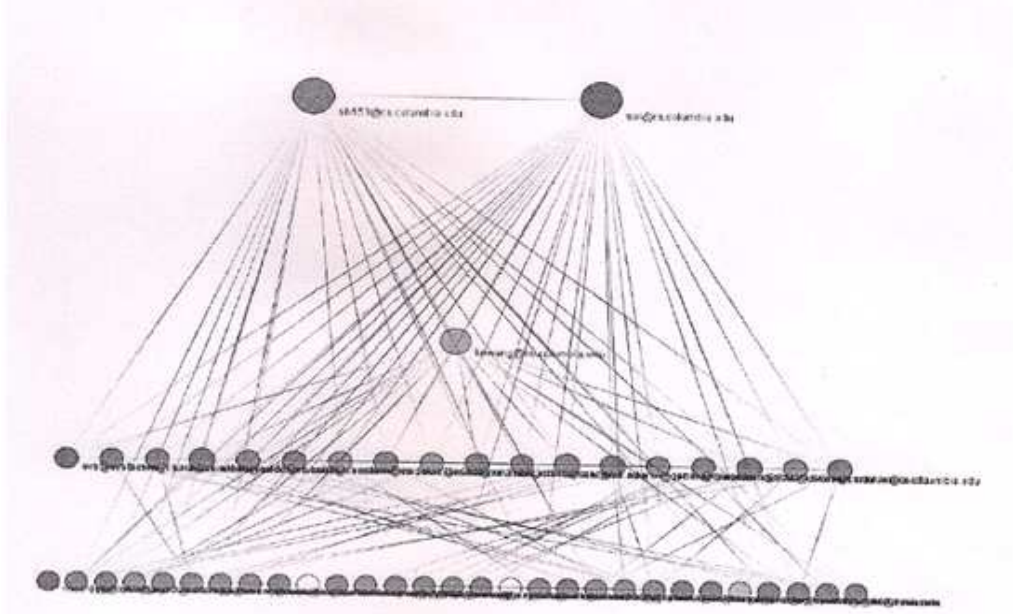Figure 1: Enron North American West Power Traders Extracted Social Network



Figure 2: Analysis of our own emails

REFERENCES

C. BRON & J. KERBOSCH (1973), Algorithm 457: finding all cliques of an undirected graph. Commun. ACM, 16(9):575–577.

D. G. DEEPAK P & V. VARSHNEY (2007),Analysis of enron email threads and quantification of employee responsiveness. In Proceedings of the Text Mining and Link Analysis Workshop on International Joint Conference on Artificial Intelligence, Hyderabad, India.

G. CARENINI, R. T. NG & X. ZHOU. Scalable discovery of hidden emails from large folders. In KDD '05: Proceeding of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining, pages 544–549, New York, NY, USA, 2005. ACM Press.

J. Diesner & K. Carley (2005),Exploration of communication networks from the enron email corpus. In Proceedings of Workshop on Link Analysis, Counterterrorism and Security, Newport Beach CA.

L. Freeman (1979), Centrality in networks: I. conceptual clarification. Social networks, 1:215–239.

T. Elsayed & D. W. Oard (July 2006),Modeling identity in archival collections of email: a preliminary study. In Third Conference on Email and Anti-spam (CEAS), Mountain View, CA.

T. Fawcett & F. Provost (1999),Activity monitoring: noticing interesting changes in behavior. In Proceedings of the Fifth ACM SIGKDD International conference on knowledge discovery and data mining (KDD-99), pages 53–62.

W. Cohen (March 2004), Enron data set.

Z. Bar-Yossef, I. Guy, R. Lempel, Y. S. Maarek & V. Soroka (2006), Cluster ranking with an application to mining mailbox networks. In ICDM '06: Proceedings of the Sixth International Conference on Data Mining, pages 63–74, Washington, DC, USA,  IEEE Computer Society.

**SUBMISSION OF PAPERS**

Contributions should be sent by email to Dr. Maneesha Shukla Editor-in-Chief, Anvikshiki, The Indian Journal of Research (maneeshashukla76@rediffmail.com)*.* www.onlineijra.com

　　Papers are reviewed on the understanding that they are submitted solely to this Journal. If accepted, they may not be published elsewhere in full or in part without the Editor-in-Chief's permission. Please save your manuscript into the following separate files-***Title; Abstract; Manuscript; Appendix.*** To ensure anonymity in the review process, do not include the names of authors or institution in the abstract or body of the manuscript.

*Title***:** This title should include the manuscript, full names of the authors, the name and address of the institution from which the work originates the telephone number, fax number and e-mail address of the corresponding author. It must also include an exact word count of the paper.

*Abstract***:** This file should contain a short abstract of no more than 120 words.

*MANUSCRIPT:* This file should contain the main body of the manuscript. Paper should be between 5 to 10 pages in lenth,and should include only such reviews of the literature as are relevant to the argument. An exact word count must be given on the title page. Papers longer than 10 pages (including *abstracts, appendices and references*) will not be considered for publication. Undue length will lead to delay in publication. Authors are reminded that Journal readership is abroad and international and papers should be drafted with this in mind.

　*References should be listed alphabetically* at the end of the paper, giving the name of journals in full. Authors must check that references that appear in the text also appear in the References and *vice versa.* Title of book and journals should be italicised.

*Examples***:**

BLUMSTEIN,A.and COHEN,J.(1973),'A Theory of Punishment' *Journal of Criminal Law and Criminology,*64:198-207

GUPTA,RAJKUMAR(2009),*A Study of The Ethnic Minority in Trinidad in The Perspective of Trinidad Indian's Attempt to Preserve Indian Culture,* India: Maneesha Publication,

RICHARDSON,G(1985),Judicial Intervention in Prison Life', in M. Maguire ,J. Vagg and R. Morgan, eds., *Accountability and Prisons,*113-54.London:Tavistocs.

SINGH,ANITA.(2007),*My Ten Short Stories,*113-154.India:Maneesha Publication.

In the text,the name of the author and date of publication should be cited as in the Harvard system(e.g.Garland 1981: 41-2;Robertson and Taylor 1973;ii.357-9)If there are more than two authors, the first name followed by *et al.* is manadatory in the text,but the name should be spelt out in full in the References. Where authors cite them as XXXX+date of publication.

*Diagrams and tables* are expensive of space and should be used sparingly. All diagrams, figures and tables should be in black and white, numbered and should be referred to in the text.They should be placed at the end of the manuscript with there preferred location indication in the manuscript(e.g.Figure 1 here).

*Appendix:* Authors that employ mathematical modelling or complex statistics should place the mathematics in a technical appendix.

*NOTE :* Please submit your paper either by post or e-mail along with your photo, bio-data, e-mail Id and a self-addressed envelop with a revenue stamp worth Rs.51 affixed on it. One hard copy along with the CD should also be sent.A self-addressed envelop with revenue stamp affixed on it should also be sent for getting the acceptance letter. Contributors submitting their papers through e-mail, will be sent the acceptance letter through the same. Editorial Board's decision will be communicated within a week of the receipt of the paper. For more information, please contact on my mobile before submitting the paper. All decisions regarding members on Editorial board or Advisory board Membership will rest with the Editor. Every member must make 20 members for Anvikshiki in one year. For getting the copies of 'Reprints', kindly inform before the publication of the Journal. In this regard, the fees will be charged from the author.