

# Corporate Communication Network and Stock Price Movements Insights from Data Mining

SRI. CH. RATNA BABU<sup>1</sup>, M. NAVEEN KUMAR<sup>2</sup>, V. SRAVANI<sup>3</sup>

<sup>1</sup> Associate Professor, Department of Computer Science and Engineering, RVR&JC College of Engineering, India

<sup>2,3</sup> UG Student, Department of Computer Science and Engineering, RVR&JC College of Engineering, India

**Abstract-** To detect communication patterns with in a company and find out how these patterns are related to performance of a company. In particular, the study concentrates on whether or not there exists any association relationships between the frequency of e-mail exchange of key employees in company and performance of company that reflects in stock prices. If there are such relationships, then they find out whether or not the company's stock price could be accurately predicted based on detected relationships.

**Data Mining algorithm is proposed to mine the e-mail communication record and historical stock prices to detect the association relationships and based on the detected relationships, rules that can predict changes in stock prices can be constructed.**

## I. INTRODUCTION

Recent research the existence of interesting communication patterns among different participants of different social network platforms. These patterns have been shown to be useful in predicting product sales and stock prices. Compared to social network, which can be considered as representing connections among people in the public, a corporate network can be express opinions on any issues of interest, members of a corporate communication network are expected to mainly talk about company-specific business. In a company, employee communications can mean the success or failure of any major change program resulting from a merger, acquisition, new venture, new process improvement approach, or other management issues. In other words, employee communication can serve a critical “business function that drives performance and contributes to a company’s financial success. Therefore, prediction of movements of stock

price is the main focus of this paper.

Unlike social networks, in a corporate communication network, e-mails have long been used as a tool for interorganizational and information exchange. In the same way, a social network platform is able to capture participants’ behaviour and their opinions about various issues and events. Thus, we argue that a corporate communication network in the form of an ecosystem also contains insightful information, such as organizational stability and robustness, about a company’s development. We believe our argument is in line with corporate communication theory, which suggests that “employee communications can mean the success or failure of any major change program” resulting from a merger, acquisition, new venture, new process improvement approach, or other management issues.

In this paper, We propose that a company’s performance, in terms of its stock price movements, can be predicted by internal communication patterns. To obtain early warnings signals, we believe that it is important for patterns in corporate communication networks to be detected earlier for the prediction of significant stock price movement to avoid possible adversities that a company may face in the stock market so that stakeholders’ interests can be protected as much as possible. Despite the potential importance of such knowledge about corporate communication, little work has been done in this important direction.

## II. PREVIOUS METHODS

- Support vector Machines: “Support Vector Machine” (SVM) is a supervised machine learning algorithm which can be used for both classification or regression challenges. However,

it is mostly used in classification problems. In this algorithm, we plot each data item as a point in  $n$ -dimensional space (where  $n$  is number of features you have) with the value of each feature being the value of a particular coordinate. Then, we perform classification by finding the hyper-plane that differentiate the two classes.

• Decision Trees:

The general motive of using Decision Tree is to create a training model which can use to predict class or value of target variables by learning decision rules inferred from prior data (training data). As one of the traditional classifiers, a decision tree is a form of multiple variable analysis for classification and prediction. It has been applied in various fields including stock market prediction and stock trading rule generation

III. PROPOSED METHOD

• Modules to be implemented


The modules to be implemented for prediction of movement in stock prices are:

- discretizing the weight matrices and the value of movement.
- discovering patterns for describing the relationship between communication frequency and stock price
- constructing prediction rules based on patterns;
- predicting the movement of stock price using prediction rules.

• Module description

Notations and Definitions

Definition 1 (communication networks): Let  $G = \{G_1, G_2, \dots, G_t, \dots, G_p\}$  represents the communication networks ( $p$  is the total time points). The one of networks is represented as a graph with a 2-element tuple  $G_t = (V_t, R_t)$ , with the following characteristics.



$V_t = \{v_1, v_2, \dots, v_n\}$  represents a finite set of nodes

(employees in the communication network),  $n$  is the total number of nodes in the network, and  $v_c$  represents the central node in the network. In the Enron case, we choose the CEO as the central node and examined the communication network between the CEO with other employees.

$R_t = \{v_i^j, t\}$  denotes the set of edges representing all the links in the network. If there is an edge  $R_t^{ij} \in R_t$ , it means that the two vertices  $v_i$  and  $v_j$  are connected through the link of  $R_t$  at  $t^{\text{th}}$  time point.

$1, t \in \mathbb{Z}^+$  and  $t = 1, 2, \dots, p$ .

Let  $S = \{s_1, s_2, \dots, s_{t-\tau}, \dots, s_t, \dots, s_p\}$  represents a time series of data values collected over a period of time and  $p$  is the total time points.  $S_t$  represents the value of movement of stock price in the  $t^{\text{th}}$  time point, where  $s_t \in [L_t, U_t]$ ;  $L_t$  represents the lower bound and  $U_t$  represents the upper bound of the values,

$1 \leq t \leq t-1, t \in \mathbb{Z}^+$  and  $t = 1, 2, \dots, p$ .

In order to better visualize the data we used, a synthetic data set shown in Fig. 3 is used to describe the format of input graph stream. Six nodes in the example represent six persons in the communication network, respectively. The value of the edge between two nodes represents communication frequency between two persons. The whole input data describes the changes of communication network during five time points. In this paper, the proposed algorithm can discover the patterns for describing the relationship between communication frequency network and the change of stock price. For testing the performance of the proposed algorithm, we predict the change of stock price of Enron using its e-mail corpus data set.




Fig. Prediction process

To better analyze the graphical data and link

information of  $G_t$ , we use a  $n \times n$  matrix  $\mathbf{D} = [d_{ij}] (1 \leq i, j \leq n)$  to represent the graph and the value of  $d_{ij}$  can capture the communication frequency that occurred between two nodes, as shown in the following:

$$d_{ij} = \begin{cases} e_{ij}, & e_{ij} \in E \\ 0, & \text{other} \end{cases}$$

By using the same synthetic data set, the graph set shown in Fig. 3 is transformed into a set of matrices as Fig. 5 demonstrated. In this case, we capture five matrices as shown in Fig. 5 to represent the communication network for all five time points.




Fig. Weighted matrices for every time point

• Proposed Algorithm

Discretizing Communication Matrix and Stock Price: The value of edge ( $d_{ij}$ ) between two different nodes for representing communication frequency between two nodes is numerous. In order to reduce and simplify the original data, numerous values are always replaced by a small number of interval label, which leads to a concise easy-to-use, knowledge-level representation of mining results. The existing and mature methods on unsupervised discretization are primarily equal frequency discretization and equal width discretization. The equal width method is typically used in every statistic program to produce regular histograms. However, equal width discretization can hardly handle the situation if outliers exist in the data set. Equal frequency can overcome the limitations of the equal width discretization by dividing the domain in intervals with the same distribution of data points. Hence, considering different people may have different habits to exchange e-mails with each other, we use equal frequency, the most represented algorithm, to discretize other nonzero value for the value of weight in the communication network.

Besides the differences between people, we also consider the differences of networks in different time points. Hence, we use  $GS_t$  represent the total number of communications occurrences, denoted as  $\sum_{j=1}^n \sum_{i=1}^n$  which can measure the overall graph sparsity in the  $t$ th

time point.

Let a set of event-based values, denoted as  $E_{ij} = \{E_{ij,1}, E_{ij,2}, \dots, E_{ij,t}, \dots, E_{ij,p}\}$ ,  $E_{ij,t} \in E_{ij}$  correspond to  $v_{ij,t} \in R_t$ , where the domain of the value of  $E_{ij}$  is denoted as  $\text{dom}(E_{ij}) = [L_{ij}, U_{ij}]$ , in which  $L_{ij}$  represents the lower bound and  $U_{ij}$  represents the upper bound of the values. Hence, the set of discrete states for  $E_{ij,t}$  is denoted as  $D(E_{ij,t}) = \{dr_1, dr_2, \dots, dr_k, \dots, dr_n\}$ , where  $dr_k$  represents the discrete states for the value of ( $R_t / GS_t$ ),  $n$  is the total number of states.

In terms of the Enron case, a secretary needs to send a great number of e-mails for coordination in course of the daily work. Therefore, the mean level, for instance 10 e-mails per day, is considered as the normal number for a secretary according to his/her communication pattern. In a relative sense, 10 e-mails per day might be already well above the mean level of the communication frequency for a front line employee. We standardize the communication level according to each person/node's mean of his/her communication frequency, and all continuous data are transformed into discrete data, such as  $dr_1, dr_2, dr_3$ . Three levels are applied for Enron case which represent "no communication," "weak relationship," and "strong relationship".

Similarly, the value of stock price is discretized. First, the original value of "ups and downs" is represented using the price fluctuation ratio, which can be calculated by  $(C_{Pt-1} - C_{Pt}) / C_{Pt-1}$  where  $C_{Pt-1}$  represents "closing price" in the day of "t-1" and  $C_{Pt}$  represents "closing price" in the day of "t." Then the movement of stock price are classified into different states, such as  $ds_1, ds_2$ , which, respectively, represents "the ratio of increases of stock price is higher than 0%" and "the ratio of decreases of stock price is lower than 0%."

In summary, the original matrix is discretized as  $R_t = \{E_{ij,t} | E_{ij,t} \in D(L_{ij,t}), D(L_{ij,t}) = \{dr_i\}, 1 \leq i, j \leq n\}$  and the value of movement of stock price is discretized as  $st = E_t$ , where  $E_t \in D(L_t)$  and  $D(L_t) = \{ds_i\}$ .

2-D TABLE WITH J ROWS AND I COLUMNS

$o_{ik}$ ( $e_{ik}$ )		Stock Price (in the certain position later)					Total
		$ds_1$	...	$ds_l$	...	$ds_l$	
$E_{j,t-1}^c$	$dr_1$ ( $e_{1j}$ )	$o_{11}$ ( $e_{11}$ )	...	$o_{l1}$ ( $e_{l1}$ )	...	$o_{l1}$ ( $e_{l1}$ )	$o_{+1}$
	.	.		.		.	.
	.	.		.		.	.
	$dr_j$ ( $e_{1j}$ )	$o_{1j}$ ( $e_{1j}$ )	...	$o_{lj}$ ( $e_{lj}$ )	...	$o_{lj}$ ( $e_{lj}$ )	$o_{+j}$
	.	.		.		.	.
$dr_j$ ( $e_{1j}$ )	$o_{1j}$ ( $e_{1j}$ )	...	$o_{lj}$ ( $e_{lj}$ )	...	$o_{lj}$ ( $e_{lj}$ )	$o_{+j}$	
<b>Total</b>	$o_{1+}$	...	$o_{l+}$	...	$o_{l+}$	$o_{+}$	

FREQUENCY TABLE FOR  $E_{j,t-1}^c$  AND STOCK PRICE

$o_{ij}$ ( $e_{ij}$ )		Stock			Total
		Up ( $ds_1$ )	Middle ( $ds_2$ )	Down ( $ds_3$ )	
$E_{j,t-1}^c$	None ( $dr_1$ )	0 (0)	0 (0)	0 (0)	0
	Strong ( $dr_2$ )	0 (0)	1 (0.5)	0 (0.5)	1
	Weak ( $dr_3$ )	0 (0)	1 (1.5)	2 (1.5)	3
	<b>Total</b>	0	2	2	4

WEIGHT VALUES OF EVIDENCE FOR  $E_{j,t-1}^c$  AND  $E_t$

$W_{ij}$		Stock		
		Up	Middle	Down
$E_{j,t-1}^c$	None	N/A	N/A	N/A
	Strong	N/A	N/A	N/A
	Weak	N/A	-1	1

In a simple example,  $z_{ij}$  ( $E_{j,t-1}^c$ =strong→stock=middle) can be calculated as  $z_{ij}=1 - 0.5/(0.5)1/2=0.71$  and statistical significances then tested based on the use of the following residual measure to determine the graph-stock pattern according to (3).  $d_{ij} = 0.71/((1 - 1/4)(1 - 2/4))1/2 = 1.154$ . And Table III shows all values of  $d_{lk}$  to measure the correlations between  $E_{j,t-1}^c$  and  $E_t$ . N/A in Table III represents nonexistence when  $o_{lk}$  is zero for counting the number of relationship takes  $o_{lk}$  on corresponded value.

The weight of evidence for or against a certain prediction of the attribute values of future objects by using the same information measure can be assessed quantitatively as follows.

Suppose that  $vil$  of stock at one time point later is dependent on  $vjk$  in the records stream. The weight of evidence measures the amount of positive or negative

evidence that is provided by the value of relationship between central node and destination node supporting or refusing the stock price being observed together. Hence,  $vjk$  that provides positive evidence supporting stock at one time point later in the graph stream.

The value  $vil$  and  $vjk$  is the characteristics of weight for  $R(Nc, NJ)$  in the preceding time point. Otherwise,  $vjk$  provides negative evidence.

#### IV. IMPLEMENTATION

**Input:**  $D=\{R_1, R_2, \dots, R_p\}$  (A set of matrices to represents communication network)

$S = \{S_1, S_2, \dots, S_p\}$  (A set of stock prices),

$v_c$  (central node)

$sl$  (the number of states of stock prices)

$rl$  (the number of states of communication frequency)

**Output:** Sub-Res (List of Adjusted Residuals Value)

Res (Integrated all Sub-Res)

Discretize  $S$  into  $E_t \subseteq (ds_1, ds_2, \dots, ds_l, \dots, ds_{sl})$

**for**  $r=1: n, r \neq c$  ( $n$  is the total number of nodes)

Discretize  $E_{j,t-1}^c$  into  $D_r \subseteq (dr_1, dr_2, \dots, dr_j, \dots, dr_{rl})$

**for**  $i=1:sl$

**for**  $j=1:rl$

$$\text{calculate } d_{ij} = \frac{o_{ij} - e_{ij}}{\left(1 - \frac{\sum_{u=1}^j o_{iu}}{M^j}\right) \left(1 - \frac{\sum_{u=1}^l o_{uj}}{M^l}\right) \cdot e_{ij}}$$

Add  $d_{ij}$  into Sub-Res

**end**

**end**

Add Result into Sub-Res

**end**

#### V. RESULTS

- Actual Results of the Work.

Firstly, the dataset containing tweets is pre-processed. Firstly, we test whether the dataset is being loaded or not and if is taken in the write format or not. If wrong format of dataset is given is given as input that the system is expected to return an error.

And here it is checked whether the system returns the error or not. If error is returned that the test case is passed else failed. If right file format is given, the system proceeds that the next module without any error. The probability of stock price moving upward or downward is obtained if done correctly, test case is passed. If there are no correct number of frequencies for calculating probabilities then the system should display no output to the user.

Next the input is entered in between the specified dates. Now the values entered are tested. The values of dates are entered as in the specified range. If the values entered tested. The values of dates are entered as in the specified range. If the values entered are out of the range then the specified range. If the values entered out of the range then the system should return an error else proceed to next step. The probability of stock price moving upward or downward is obtained if done correctly, test case is passed. If there are no correct number of frequencies for calculating probabilities then the system should display no output to the user. Input is entered in between the specified dates. Now the values entered are tested. The values of dates are entered as in the specified range.

```
Microsoft Windows [Version 10.0.17134.706]
(c) 2018 Microsoft Corporation. All rights reserved.

C:\Users\Waveen Chowdary>cd Desktop
C:\Users\Waveen Chowdary\Desktop>cd regression
C:\Users\Waveen Chowdary\Desktop\regression>javac Preprocess.java
C:\Users\Waveen Chowdary\Desktop\regression>java Preprocess
Preprocessing done successfully
C:\Users\Waveen Chowdary\Desktop\regression>javac Input.java
C:\Users\Waveen Chowdary\Desktop\regression>java Input
Example: 2013-12-25
Enter startdate: 2016-04-17
Enter enddate: 2016-05-06
132
C:\Users\Waveen Chowdary\Desktop\regression>
```

Screenshot of entering the values

```
Microsoft Windows [Version 10.0.17134.706]
(c) 2018 Microsoft Corporation. All rights reserved.

C:\Users\Waveen Chowdary>cd Desktop
C:\Users\Waveen Chowdary\Desktop>cd regression
C:\Users\Waveen Chowdary\Desktop\regression>javac Preprocess.java
Preprocessing done successfully
C:\Users\Waveen Chowdary\Desktop\regression>
```

Screenshot of Importing the dataset

```
C:\Users\Waveen Chowdary\Desktop\regression>javac SentimentAnalysis.java
C:\Users\Waveen Chowdary\Desktop\regression>java SentimentAnalysis
Minimum stock value:186.21
Maximum stock value:186.45
Min value:10.119999
C:\Users\Waveen Chowdary\Desktop\regression>
```

Screenshot of analysis

Next the input is entered in between the specified dates. Now values entered are tested. The values of dates are entered as in the specified range. If the values entered are out of range

```
--upgrade-module-path (path)
    Override location of upgradeable modules
--verbose                    Output messages about what the compiler is doing
--version, -version          Version information
-Merror                       Terminate compilation if warnings occur

C:\Users\Waveen Chowdary\Desktop>cd regression
C:\Users\Waveen Chowdary\Desktop\regression>javac Preprocess.java
C:\Users\Waveen Chowdary\Desktop\regression>java Preprocess
Example: 2013-01-25
Enter dates between 2016-04-14 to 2016-07-02
Enter startdate: 2016-04-17
Enter enddate: 2016-05-15
No of tweets between entered dates:285
Minimum stock value:86.21
Maximum stock value:186.45
Min value:10.119999
Frequency table for tweets and stock values
32 89
8 76
Probability that stock value can move upwards when frequency is positive will be:0.46
Probability that stock value can move upwards when frequency is negative will be:0.2
C:\Users\Waveen Chowdary\Desktop\regression>
```

Screenshot of prediction

• Analysis of Results Obtained:

We analyze the complexity of proposed algorithm. Suppose that we have N nodes in a community network with the p total time points s represents states of stock price and r states of communication

frequency. After preprocessing and discretizing the stock price and communication between central nodes, the proposed algorithm needs to iterate in order to calculate a the adjusted residual to measure the correlations between communication frequency and stock prices.

### CONCLUSION

The findings and theoretical implications from this paper are twofold. On one hand, we captured the communications among nodes in Enron's major corporate communication network and identified employees' communication patterns. This paper demonstrates that a corporate e-mail ecosystem contains meaningful information about employees' communication patterns. Even if we only focus on the communication frequency, a company (Enron in our case) has identifiable patterns of e-mail exchange. Such identifiable patterns can reveal important information about major corporate activities and organizational stability that may subsequently influence the focal company's performance in the stock market. Therefore, cooperate communication patterns can serve as a good proxy to predict a company's stock performance. Our experimental results demonstrated the existence of dependence between e-mail communication network and stock price for Enron.

This paper extended the existing communication theories to capture the patterns of corporate communication and the focal company's stock performance.

On the other hand, social networks have become a hot topic in the field of data mining recently. In this paper, we not only provided an innovative idea on using data-mining algorithms but also constructed the relationship between social network and finance. Hence, this paper demonstrated great potential to predict the amounts of increases and decreases of stock price based on the weighted rules.

### ACKNOWLEDGMENT

We have taken a lot of effort into this project. However, completing this project would not have been possible without the support and guidance of our

faculty of Computer Science and Engineering department at RVR&JC College Of Engineering and our fellow students. We would like to extend our sincere thanks to all of them.

We are highly indebted to Sri. Ch. Ratna Babu(Ass.Professor) for her guidance and supervision. We would like to thank her for providing the necessary information and resources for this project.

We would like to express our gratitude towards our parents & our friends for their kind co-operation and encouragement which help us a lot in completing this project. Our thanks and appreciations also go to our colleague in developing the project. Thank you to all the people who have willingly helped us out with their abilities.

### REFERENCES

- [1] B. Collingsworth, R. Menezes, and P. Martins, "Assessing organizational stability via network analysis," in Proc. IEEE Symp. Comput. Intell. Financial Eng., Mar. 2009, pp. 43–50.
- [2] D. J. Barrett, "Change communication: Using strategic employee communication to facilitate major change," *Corporate Commun., Int. J.*, vol. 7, no. 4, pp. 219–231, 2002.
- [3] M.-C. Wu, S.-Y. Lin, and C.-H. Lin, "An effective application of decision tree to stock trading," *Expert Syst. Appl.*, vol. 131, no. 2, pp. 270–274, 2006.
- [4] C. W. Down, G. C. Phillip, and A. L. Pfeiffer, "Communication and organizational outcomes," in *Handbook of Organizational Communication*, G. Goldhaber and G. Barnett, Eds. Norwood, NJ, USA: Ablex, 1988.
- [5] P. G. Clampitt and C. W. Downs, "Employee perceptions of the relationship between communication and productivity: A field study," *J. Bus. Commun.*, vol. 30, no. 1, pp. 5– 28, 1993.
- [6] G. S. Hansen and B. Wernerfelt, "Determinants of firm performance: The relative importance of economic and organizational factors," *Strategic Manage. J.*, vol. 10, no. 5, pp. 399–411, 1989.
- [7] W. W. Burke and G. H. Litwin, "A causal model of organizational performance and change," *J. Manage.*, vol. 18, no. 3, pp. 523–545, 1992.

- [8] R. Katz, "The effects of group longevity on project communication and performance," *Admin. Sci. Quarterly*, vol. 27, no. 1, pp. 81–104, 1982.
- [9] V. A. Zeithaml, L. L. Berry, and A. Parasuraman, "Communication and control processes in the delivery of service quality," *J. Marketing*, vol. 52, no. 2, pp. 35–48, 1988.
- [10] C. Hargreaves and Y. Hao, "Prediction of stock performance using analytical techniques," *J. Emerg. Technol. Web Intell.*, vol. 5, no. 2, pp. 136–142, 2013.
- [11] R. B. Higgins and B. D. Bannister, "How corporate communication of strategy affects share price," *Long Range Planning*, vol. 25, no. 3, pp. 27–35, 1992.