

Time Series Analysis and Forecasting Methods for Temporal Mining of Interlinked Documents

Prasanna Desikan and Jaideep Srivastava
Department of Computer Science
University of Minnesota.
<[desikan,srivastava](mailto:desikan,srivastava@cs.umn.edu)>@cs.umn.edu

Abstract

The need to study and understand the evolving content, structure and usage of interlinked documents has gained importance recently, especially with the advent of the World Wide Web and online Citation Indices. In this report we have surveyed the various time series models and forecasting methods that could be used as an effective tool to capture the evolving data that is constituted by interlinked documents. Our dataset consists of research papers from a high-energy physics archive of journals that can be represented in the form of a graph with the papers as the nodes and the citations as the edges. The evolving nature of interlinked documents has been posed as a challenging problem in the KDD Cup 2003. We have applied some of the existing forecasting techniques to this graph to be able to predict the number of citations a paper would receive in the future. This report reflects our finding on how the various time series models can be applied and the effective ways of forecasting the link structure in such interlinked documents.

1. Introduction

Time series models have been the basis for any study of a behavior of process or metrics over a period of time. The applications of time series models are manifold, including sales forecasting, weather forecasting, inventory studies etc. In decisions that involve factor of uncertainty of the future, time series models have been found one of the most effective methods of forecasting. Most often, future course of actions and decisions for such processes will depend on what would be an anticipated result. The need for these anticipated results has encouraged organizations to develop forecasting techniques to be better prepared to face the seemingly uncertain future. Also, these models can be combined with other data mining techniques to help understand the behavior of the data and to be able to predict future trends and patterns in the data behavior.

The evolving structure of interlinked documents, such as the World Wide Web or online citation indices, and the usage of these documents over a period of time has been of interest to both the researchers and the industry. These set of documents form a graph, with the nodes representing the documents and the edges representing the hyperlinks or the citations. Research has been carried out in extracting information from the pure structure of such graphs and also on the usage of these documents, especially with respect to the World Wide Web. The stability of the Web structure has led to the more research related to Hyperlink Analysis and the field gained more recognition with the advent of Google [1]. A survey on Hyperlink Analysis is provided in [2]. Usage aspects of such documents have also received wide attention and Srivastava et al [3] provide a good overview of Web usage mining research ideas and its applications.

Most research has thus focused more recently on mining information from structure and usage of such graphs. In this study we focus on another important dimension of mining such graphs as identified [4] - the *Temporal Evolution*. The Web is changing fast over time and so is the users interaction in the Web suggesting the need to study and develop models for the evolving *Web Content*, *Web Structure* and *Web Usage*.

Also, of interest have been how the citation structure of research papers changes over time and how the access patterns of these papers vary over time.

The need to study the *Temporal Evolution* of the interlinked document has motivated us to analyze the various time series models that can be applied to them. We study the various the time series models that exist and the kind of data they are suitable to apply to. We also discuss some of the forecasting methods that are currently used. **Figure 1** depicts the idea of change of such interlinked documents over time.

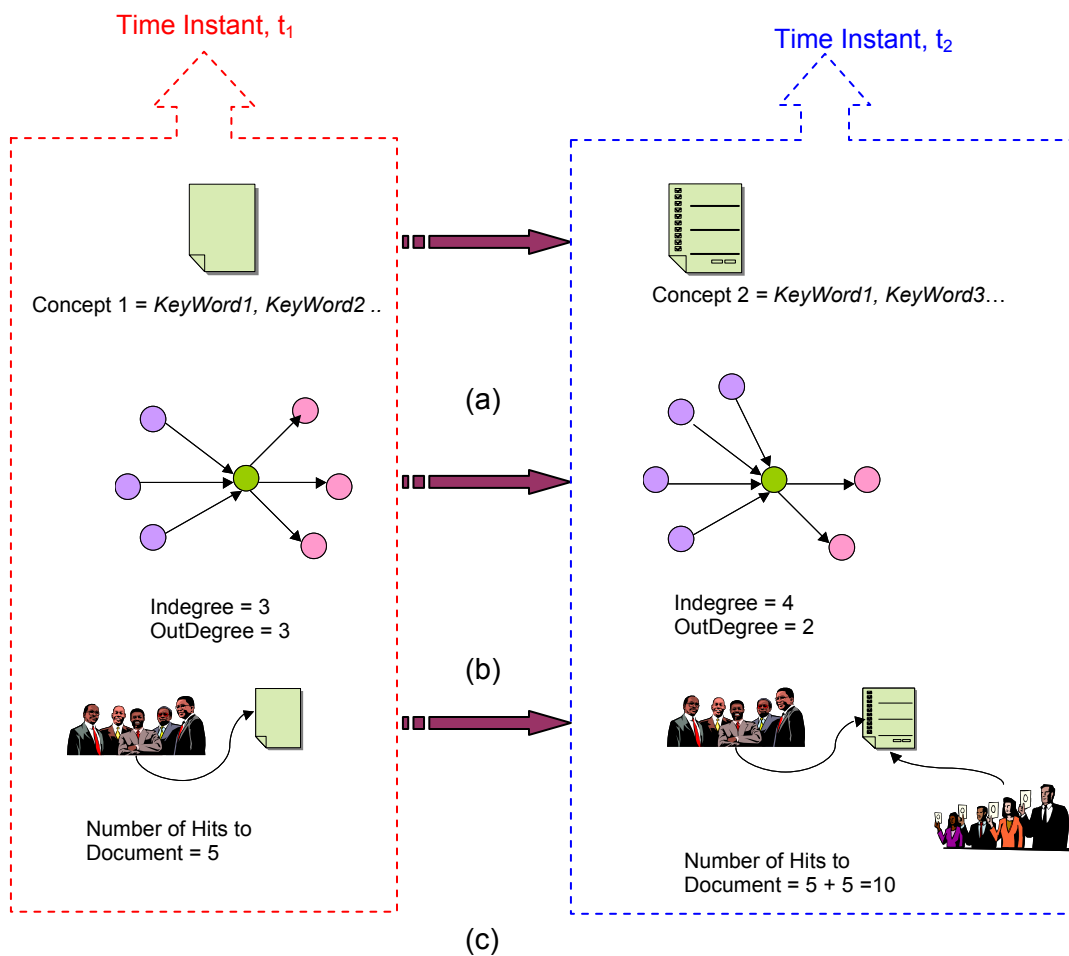


Figure 1: Temporal Evolution of a single document.

(a) Change in the *Content* of a document over time.

(b) Change in the *Structure* i.e. number of inlinks and outlinks; of a document over time.

(c) Change in the *Structure* i.e. number of inlinks and outlinks; of a document over time.

As part of our experiment, we try some of fitting in some of these models to the publicly available KDD Cup data that consists of research papers from high-energy physics. The rest of the document is organized as follows. In the section 2 we briefly describe the various time series methods that exist. Time Series Forecasting models and forecasting methods are discussed in this section 3. In section 4 we describe the data and the experimental set up. The results of the experiments are presented in section 5. Finally we provide the some conclusions and future directions.

2. Time Series Analysis Techniques

Time Series can be defined as *an ordered sequence of values of a variable at equally spaced time intervals* [5]. The motivation to study time series models is twofold:

- Obtain an understanding of the underlying forces and structure that produced the observed data
- Fit a model and proceed to forecasting, monitoring or even feedback and feedforward control.

Time Series Analysis can be divided into two main categories depending on the type of the model that can be fitted. The two categories are:

- *Kinetic Model*: The data here is fitted as $x_t = f(t)$. The measurements or observations are seen as a function of time.
- *Dynamic Model*: The data here is fitted as $x_t = f(x_{t-1}, x_{t-2}, x_{t-3} \dots)$.

The classical time series analysis procedures decomposes the time series function $x_t = f(t)$ into up to four components [6]:

1. **Trend**: a long-term monotonic change of the average level of the time series.
2. **The Trade Cycle**: a long wave in the time series.
3. **The Seasonal Component**: fluctuations in time series that recur during specific time periods.
4. **The Residual component** that represents all the influences on the time series that are not explained by the other three components.

The *Trend* and *Trade Cycle* correspond to the smoothing factor and the *Seasonal* and *Residual* component contribute to the cyclic factor. Often before time series models are applied, the data needs to be examined and if necessary, it has to be transformed to be able to interpret the series better. This is done to stabilize the variance. For example, if there is a trend in the series and the standard deviation is directly proportional to the mean, then a logarithmic transformation is suggested. And in order to make the seasonal affect additive, if there is a trend in the series and the size of the seasonal effect tends to increase with the mean then it may be advisable it transform the data so as to make the seasonal effect constant from year to year. Transformation is also applied sometimes to make the data normally distributed. The fitting of time series models can be an ambitious undertaking. There are many methods of model. These models have been well discussed in [7, 8]. The user's application and preference will decide the selection of the appropriate technique. We will now discuss some of the existing methods in time series analysis.

2.1 Smoothing Methods

Inherent in the collection of data taken over time is some form of random variation. There exist methods for reducing or canceling the effect due to random variation. An often-used technique in industry is "smoothing". This technique, when properly applied, reveals more clearly the underlying trend, seasonal and cyclic components. There are two distinct groups of smoothing methods: *Averaging Methods* and *Smoothing Methods*.

2.1.1 Averaging Methods

The "simple" average or mean of all past observations is only a useful estimate for forecasting when there are no trends. The average "weighs" all past observations equally. In general:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \left(\frac{1}{n}\right) x_1 + \left(\frac{1}{n}\right) x_2 + \dots + \left(\frac{1}{n}\right) x_n$$

The $(1/n)$ are the weights associated with each value of x . As we can see, these weights are normalized and sum upto 1.

An alternative way to summarize the past data is to compute the mean of successive smaller sets of numbers of past data. This smoothing process is continued by advancing one period and calculating the next average of t numbers, dropping the first number. Such type of averaging is called *Single Moving Average* and the general expression for the moving average is

$$M_t = [X_t + X_{t-1} + \dots + X_{t-N+1}] / N$$

There exists a variation on the MA procedure that often does a better job of handling trend. It is called Double Moving Averages for a Linear Trend Process. It calculates a second moving average from the original moving average, using the same value for M . As soon as both single and double moving averages are available, a computer routine uses these averages to compute a slope and intercept, and then forecasts one or more periods ahead.

2.1.2 Exponential Smoothing Methods

This is a very popular scheme to produce a smoothed Time Series. Whereas in Single Moving Averages the past observations are weighted equally, Exponential Smoothing assigns exponentially decreasing weights as the observation get older. In other words, recent observations are given relatively more weight in forecasting than the older observations. In the case of moving averages, the weights assigned to the observations are the same and are equal to $1/N$. In exponential smoothing, however, there are one or more smoothing parameters to be determined (or estimated) and these choices determine the weights assigned to the observations.

This smoothing scheme begins by setting S_2 to y_1 , where S_t stands for smoothed observation or EWMA, and y stands for the original observation. The subscripts refer

to the time periods, 1, 2, ..., n . For the third period, $S_3 = \alpha y_2 + (1-\alpha) S_2$; and so on. There is no S_1 ; the smoothed series starts with the smoothed version of the second observation.

For any time period t , the smoothed value S_t is found by computing

$$S_t = \alpha y_{t-1} + (1-\alpha)S_{t-1}, \quad 0 < \alpha \leq 1, \quad t \geq 3$$

This is the *basic equation of exponential smoothing* and the constant or parameter α is called the *smoothing constant*. The speed at which the older responses are dampened (smoothed) is a function of the value of α . When α is close to 1, dampening is quick and when α is close to 0, dampening is slow. We choose the best value for α so the value which results in the smallest Mean Squared Error.

3. Time Series Models and Forecasting

Time series Models and forecasting methods have been studied by various people and detailed analysis can be found in [9, 10,12]. Time Series Models can be divided into two kinds. Univariate Models where the observations are those of single variable recorded sequentially over equal spaced time intervals. The other kind is the Multivariate, where the observations are of multiple variables. A common assumption in many time series techniques is that the data are stationary. A stationary process has the property that the mean, variance and autocorrelation structure do not change over time. Stationarity can be defined in precise mathematical terms, but for our purpose we mean a flat looking series, without trend, constant variance over time, a constant autocorrelation structure over time and no periodic fluctuations. There are a number of approaches to modeling time series. We outline a few of the most common approaches below.

Trend, Seasonal, Residual Decompositions: One approach is to decompose the time series into a trend, seasonal, and residual component. Triple exponential smoothing is an example of this approach. Another example, called seasonal loess, is based on locally weighted least squares.

Frequency Based Methods: Another approach, commonly used in scientific and engineering applications, is to analyze the series in the frequency domain. An example of this approach in modeling a sinusoidal type data set is shown in the beam deflection case study. The spectral plot is the primary tool for the frequency analysis of time series.

Autoregressive (AR) Models: A common approach for modeling univariate time series is the autoregressive (AR) model:

$$X_t = \delta + \phi_1 X_{t-1} + \phi_2 X_{t-2} + \dots + \phi_p X_{t-p} + A_t$$

where X_t is the time series, A_t is white noise, and

$$\delta = \left(1 - \sum_{i=1}^p \phi_i\right) \mu$$

with μ denoting the process mean.

An autoregressive model is simply a linear regression of the current value of the series against one or more prior values of the series. The value of p is called the order of the AR model. AR models can be analyzed with one of various methods; including standard linear least squares techniques. They also have a straightforward interpretation.

Moving Average (MA): Models another common approach for modeling univariate time series models is the moving average (MA) model:

$$X_t = \mu + A_t - \theta_1 A_{t-1} - \theta_2 A_{t-2} - \dots - \theta_q A_{t-q}$$

where X_t is the time series, μ is the mean of the series, A_{t-i} are white noise, and $1, \dots, q$ are the parameters of the model. The value of q is called the order of the MA model. That is, a moving average model is conceptually a linear regression of the current value of the series against the white noise or random shocks of one or more prior values of the series. The random shocks at each point are assumed to come from the same distribution, typically a normal distribution, with location at zero and constant

scale. The distinction in this model is that these random shocks are propagated to future values of the time series. Fitting the MA estimates is more complicated than with AR models because the error terms are not observable. This means that iterative non-linear fitting procedures need to be used in place of linear least squares. MA models also have a less obvious interpretation than AR models. Note, however, that the error terms after the model is fit should be independent and follow the standard assumptions for a univariate process.

Box-Jenkins Approach: The Box-Jenkins ARMA model is a combination of the AR and MA models:

$$X_t = \delta + \phi_1 X_{t-1} + \phi_2 X_{t-2} + \dots + \phi_p X_{t-p} + A_t - \theta_1 A_{t-1} - \theta_2 A_{t-2} - \dots - \theta_q A_{t-q}$$

where the terms in the equation have the same meaning as given for the AR and MA model [9].

The Box-Jenkins model assumes that the time series is stationary. Box and Jenkins recommend differencing non-stationary series one or more times to achieve stationarity. Doing so produces an ARIMA model, with the "I" standing for "Integrated". Some formulations transform the series by subtracting the mean of the series from each data point. This yields a series with a mean of zero. Whether you need to do this or not is dependent on the software you use to estimate the model. Box-Jenkins models can be extended to include seasonal autoregressive and seasonal moving average terms. Although this complicates the notation and mathematics of the model, the underlying concepts for seasonal autoregressive and seasonal moving average terms are similar to the non-seasonal autoregressive and moving average terms. The most general Box-Jenkins model includes difference operators, autoregressive terms, moving average terms, seasonal difference operators, seasonal autoregressive terms, and seasonal moving average terms. As with modeling in general, however, only necessary terms should be included in the model.

3.1 Forecasting Methods

The main objective of forecasting for a given series $x_1, x_2, x_3, \dots, x_N$; to estimate future values such as x_{N+k} , where the integer k is called the lead time [7]. The forecast of x_{N+k} made at a time N for k steps ahead is denoted by $\hat{x}(N, k)$.

Figure 2 depicts a classification of Forecasting Methods based on the kind of approach.

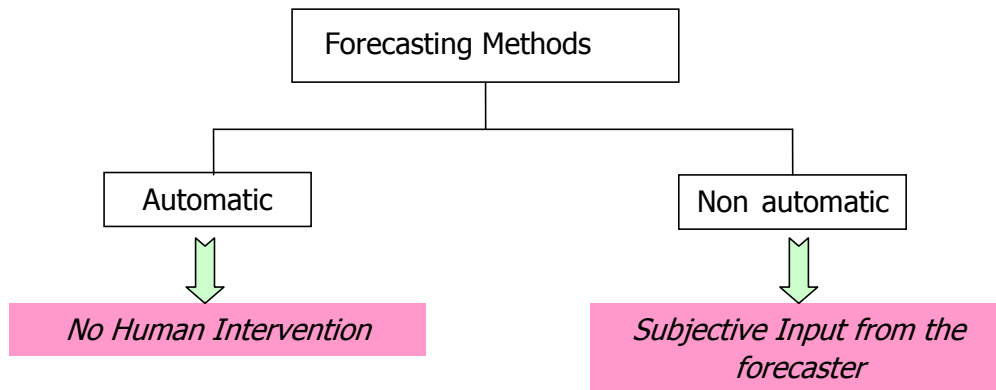


Figure 2: Forecasting Methodology

A straight line model is used to relate the time series, Y_t , to time, t , and the least squares line is used to forecast the future values of Y_t .

$$E(Y_t) = \beta_0 + \beta_1 t.$$

It is risky to use a least squares regression model outside the experimental region, especially for prediction purposes. Cyclical or Trade effects like the effects of an inflation or recession are not included.

Among the other forecasting techniques that are based on time series models and techniques discussed earlier we present only the one based on exponentially smoothing as it fits in most cases and we have also implemented the model, in our experiment.

The exponential smoothed forecast for Y_{t+j} is the smoothed value at time t .

$$F_{t+1} = E_t$$

where F_{t+1} is the forecast of Y_{t+1} .

$$\begin{aligned} F_{t+1} = E_t &= wY_t + (1-w)E_{t-1} \\ &= wY_t + (1-w)F_t \\ &= F_t + w(Y_t - F_t). \end{aligned}$$

Exponential smoothed forecast are appropriate only when trend and seasonal components are relatively insignificant. Smoothed values will tend to lag behind when a long-term trend exists. Averaging tends to smooth any seasonal component.

The most popular forecasting technique is the Holt-Winters Forecasting Technique [12]. It consists of both exponential component (E_t) and a trend component (T_t).

The calculation time begins at $t=2$, because the first two observations are needed to obtain the first estimate of trend T_2 .

$$\begin{aligned} E_2 &= Y_2 \\ T_2 &= Y_2 - Y_1 \\ E_t &= wY_t + (1-w)(E_{t-1} + T_{t-1}), 0 < w < 1. \\ T_t &= v(E_t - E_{t-1}) + (1-v)T_{t-1}, 0 < v < 1. \end{aligned}$$

Note: 'v' closer to zero suggests more weight to past estimates of trend, and 'v' value closer to one suggests more weight to current change in level.

Firstly, the exponentially smoothed and trend components, E_t and T_t , for each observed value of $Y_t (t \geq 2)$ are calculated. The one-step-ahead forecasting is determined using.

$$F_{t+1} = E_t + T_t$$

And the k-step-ahead forecast using:

$$F_{t+k} = E_t + kT_t$$

4. Experimental Task and Analysis

The goal of this task is to predict changes in the number of citations to individual papers over time [11]. The data available consisted of all papers and their unique ids, their publishing dates and the paper cited from and the paper cited to. We now built a graph using this information with the nodes as the papers and the links as the citations. We also noted that time in terms of the month in which a paper received a citation. There were about 30,000 papers and 300,000 citations in all. Out of these papers about 8000 papers were published but were not referenced at all. We noticed that the indegree distribution followed a power-law distribution with a constant of about 1.6. The plot of the indegree distribution is depicted in Figure 3.

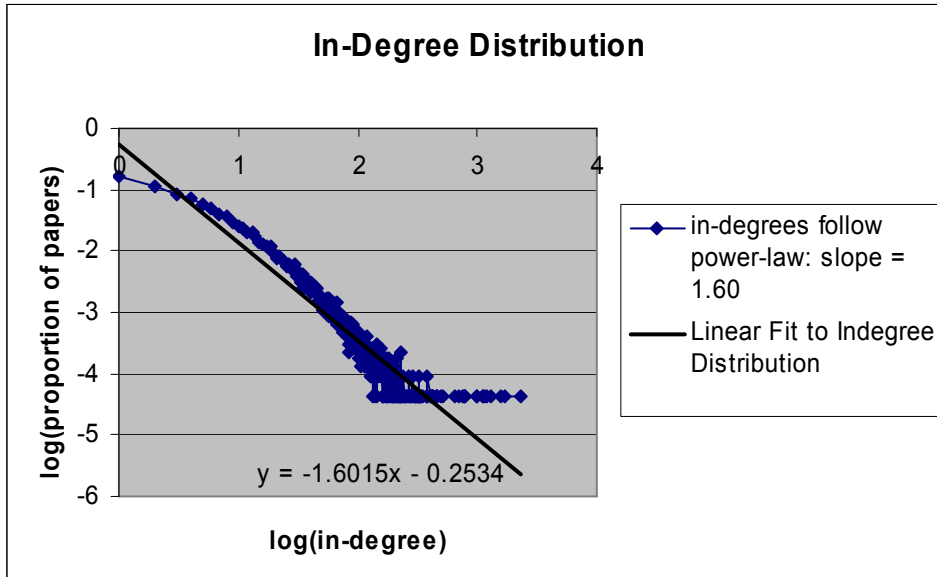


Figure 3: Indegree Distribution of Papers

We divided the data into training and test sets where only the last month was considered for testing. We tried using linear regression models for forecasting, but due to the nature of the data, the predicted values turned out to be negative in certain cases as the fitted line had a negative slope. The best-fit method seemed to be the exponential smoothing method. For this method we tried to fit in two models, one for the time series as such with the citations per month as the function value and the other with the cumulative citations as the function value. The model for forecasting was:

$$Predicted_{t+1} = Predicted_t + 0.3*(Actual_t - Predicted_t)$$

We noticed that the error as interpreted by the L 1 norm of the difference between the predicted value and the actual value as greater for the cumulative citations case. When we tried to predict for a period of three months, using single and double exponentially smoothing methods. There was not much difference between the L 1 norm in both the cases and in our case the single exponential method seem to perform better. We then decided to count the citations on a quarterly basis, which means we kept a count for the number of citations a paper received in every three months. We then used single exponential smoothing method to predict the number of citations the paper would receive for the next three months. This proved to be slightly better which indicated that as the granularity increases, there is chance for better prediction. Figure 4 shows the predicted citations and the actual citations for each paper on the data.

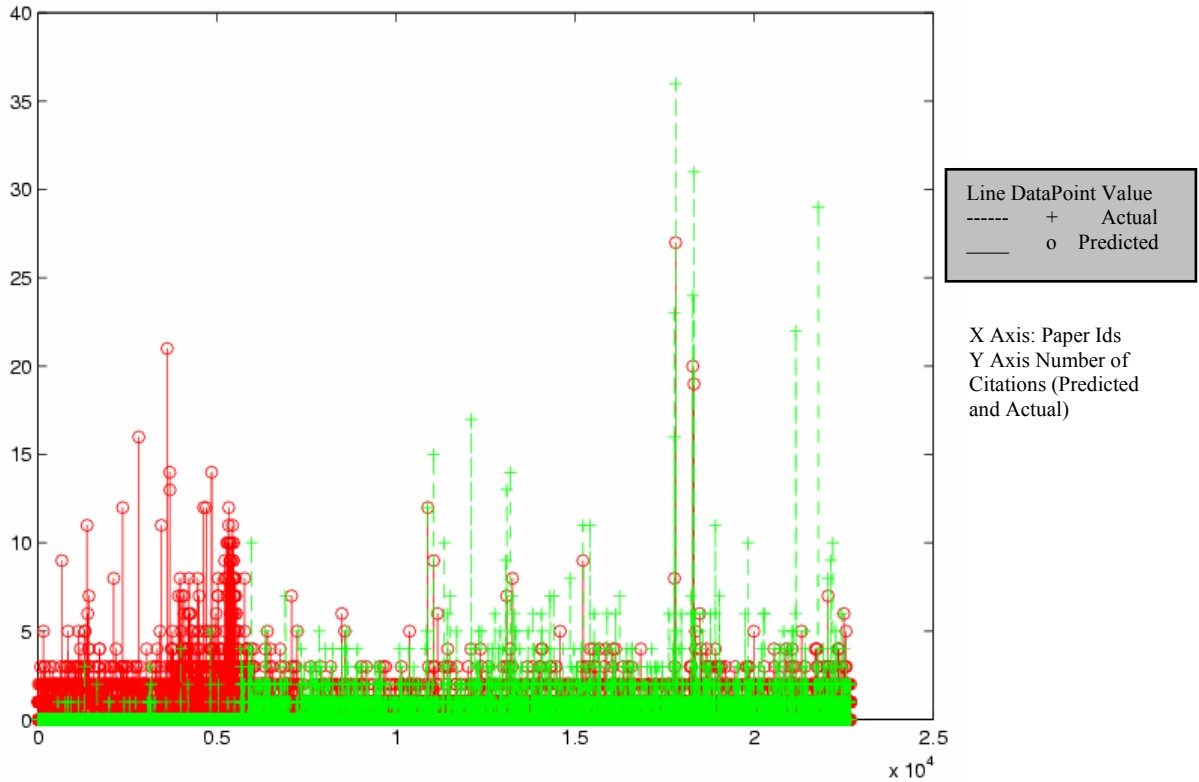


Figure 4: Predicted Number of Citations versus Actual Number of Citations for the papers in the data set.

5. Conclusions and Future Work

It is an interesting problem to predict the number of citations a paper would receive. Though forecasting and prediction is not very accurate, it would be good if we could achieve higher percentage of accuracy. Another interesting idea would be to cluster the papers based on the time series and see if that could reveal the kind of paper, such as a survey paper, a new idea in the paper etc. Other issues could be to segment the time series depending on the age of the paper and trying to fit different models. Though we tried doing it, given the short frame of time we could not perceive it actively and so further exploration is necessary. Also the next immediate task is to estimate the number of downloads a paper would receive in the first three months of its publications. We believe that the “popularity” of an author could play a role in predicting the number of downloads. The evolvement of content, structure and usage could thus reveal very interesting information.

6. Acknowledgements

We would like to thank Sandeep Mopuru and Praveen Vutukuru for providing valuable feedback and their efforts in data pre-processing. This work was partially supported by Army High Performance Computing Research Center contract number DAAD19-01-2-0014. The content of the work does not necessarily reflect the position or policy of the government and no official endorsement should be inferred. Access to computing facilities was provided by the AHPCRC and the Minnesota Supercomputing Institute.

References

1. Google Inc, <http://www.google.com/>
2. P. Desikan, J. Srivastava, V. Kumar, P.-N. Tan, “Hyperlink Analysis – Techniques & Applications”, Army High Performance Computing Center Technical Report, 2002.
3. J. Srivastava, R. Cooley, M. Deshpande and P-N. Tan. “Web Usage Mining: Discovery and Applications of usage patterns from Web Data”, SIGKDD Explorations, Vol 1, Issue 2, 2000.
4. J. Srivastava, P. Desikan and V. Kumar, “Web Mining – Accomplishments and Future Directions”, Invited paper in National Science Foundation Workshop on Next Generation Data Mining, Baltimore, MD, Nov. 1-3, 2002.
5. *NIST/SEMATECH e-Handbook of Statistical Methods*,
<http://www.itl.nist.gov/div898/handbook/>
6. J.T. McClave, P.G. Benson and T. Sincich, “Statistics for Business and Economics”, Prentice Hall, 2001
7. C. Chatfield, “The Analysis of Time Series – An Introduction”, Chapman and Hall , 1996.
8. P.J. Brockwell and R.A. Davis. (1987). *Time Series: Theory and Methods*, Springer-Verlang
9. G. E. P Box, G. M. Jenkins, and G. C. Reinsel. (1994). *Time Series Analysis, Forecasting and Control*, 3rd ed. Prentice Hall, Englewood Clifs, NJ.
10. P.J. Brockwell and R.A. Davis, (2002). *Introduction to Time Series and Forecasting*, 2nd. ed., Springer-Verlang.
11. KDD Cup 2003, <http://www.cs.cornell.edu/projects/kddcup/index.html>.
12. R.Yaffee “A. Introduction to time series analysis and forecasting with applications of SAS and SPSS /.” San Diego: Academic Press, c2000.