



iJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 6

Issue: II

Month of publication: February 2018

DOI:

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Comment Prediction in Facebook Pages using Regression Techniques

V. Pavithra¹, Binil Kuriachan²

¹SITE School, VIT University, Vellore, Tamil Nadu, INDIA

² Verizon, Chennai, Tamil Nadu, INDIA School of Information Technology and Engineering VIT University, Vellore

Abstract: *Data in Social Networks is increasing day by day. It requires highly managing service to handle the large amount of data towards it. This work is about to study the user activity patterns in Social Networks. So, concentrated on active social Networks which is “Facebook” especially in Facebook Page. Here, user comment volume prediction is made based on page category i.e., for a particular category of page’s post will get certain amount of comments. In order to predict the comment volume for each page and to find which page category getting the highest comment. In preliminary work, it has been concluded with decision tree. So, In Further Study, have analyzed with some more Regression Techniques to make the prediction Effective. In this work, modelled the user comment pattern with respect to Page Likes and Popularity, Page Category and Time. Here Decision Tree, LASSO, K-Nearest Neighbor (KNN), Random Forest, and Leaner Regression Techniques are used. The error is found by Root Mean Square Error (RMSE) Metrics. Then, concluded that K-Nearest Neighbor Algorithm performing well and giving the effective prediction.*

Keywords: *Regression, Decision Tree, K-Nearest Neighbor (KNN), Random Forest, Linear Regression, RMSE.*

I. INTRODUCTION

The leading trends towards the Social Networking has been drawn high public attention from past ‘one and half’ decade. The merging and computing with the physical things have been induced the conversion of everyday objects in information applications. These services are acting like a multi-tool (several Category) along with routine applications^[9] for e.g., news, advertisements, communication, commenting, marketing, banking, Entertainment etc. These categories are becoming noticed every day and much more are on the way in this field^[1]. All these services have their daily huge contents generated in common, that is expected to be stored in the Hadoop cluster. As per Facebook feed in daily basis, there are 500+ terabytes of new data are inserted into databases every day, 100+ petabytes of disk space in one of the FB’s largest Hadoop (HDFS)^[9] cluster and there were 2.5 billion of content (feed) items shared per day (news feeds, shares, photos, wall posts, photos, videos, status updates, comments, etc.). For more understanding Twitter went from 5,000 tweets per day in 2007, but in 2013 it become 500,000,000 tweets per day. Flickr app features 5.5 billion images in January 31, 2011 and 3k-5k images are adding every minute.

In this paper, focused on leading Social Networking Application service Facebook, especially ‘Facebook Pages’ (one of the product from Facebook in current Trend), for automatic analysis of trends and patterns of users. So, for this work, Feature Selection has been concentrated more and performed it with different Regression Algorithms to get best predictions. For example, Parametrized and Non Parameterized results are varied based on Target Variables. This research is based towards the comment volume prediction (CVP) that a document is expected to be received in next H hours.

This paper is explained with Section II discussing about the related works, in Section III Feature selection and Regression Techniques. Section IV, Experiment Setting. And finally enclosed with Conclusion and Future work in Section V, followed by References.

II. RELATED WORKS

Predicting the unwanted comments in YouTube by calculating the amount of bad comments based on some key words. For example,

if certain bad words are there in a comment it will be assumed as a bad comment and it will be omitted. This has been done with various regression Techniques E.g. Linear Regression, Decision Tree, Neural Network and used certain metrics like RMSE, R squared and AUC. From the comparison the RMSE worked well to find error in it. Based on that, we have concentrated on certain Regression Technique and chose RMSE to find the error.

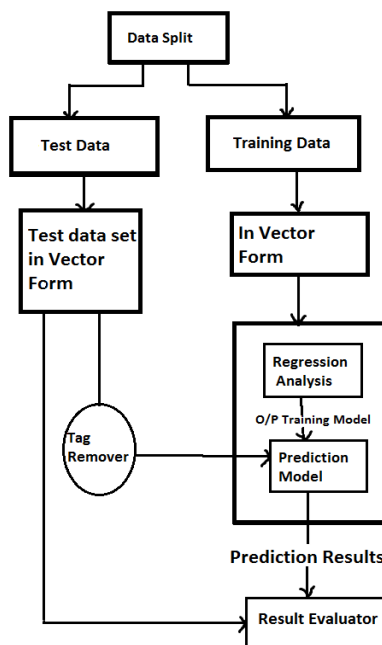


Fig.1. System Architecture

The system Architecture of the model explains that data set split into training and testing before modeling the data and then change into vector form in order to push it for prediction model and the results will be generated with respect to minimal error obtained. It implies with Regression technique. In order to predict the fields expected.

The structure of each process is carried out in each phase that is shown in the Fig.1. architecture diagram.

A. Data Mining Techniques

Supervised learning is defined as where the input variables(x) and output(y) is given. With that any of algorithm to be used in order to map the data in required format. It has two types,

- 1) *Classification*: A Classification is defined as any variable is categorized or having certain limits like data up to 100 or “True” or “False”.
- 2) *Regression*: Regression is defined as the variable is continuous, mostly real values termed to be target values. Unsupervised learning is where there is no output variable (y) and having only input variable to predict the outcome. It mostly based on
- 3) *Clustering*: Clustering is basically termed as grouping the similar data into one group and so on. Such as cluster of customer bought same product X
- 4) *Association*: It's from the Rule mining Technique where the cluster of frequent item set has been formed a rule to identify the best match for a product A in product B and C.

B. Regression

Regression Analysis is the predictive modelling technique which help in dealing with continuous variables and determines the relationship between the target value and predictors in it. Facebook data is continuous so regression analysis will work. This model is used to Forecast any results and finding the causal effects among variables. For example, possibilities of a person getting cancer due to smoking compared with junk food can be determined using regression technique ^[5].

- 1) *Linear Regression*: It's a common regression technique that helps in the forecasting the results. It has been widely used. In this model the dependent value will be continuous and the independent variable may be of discrete or continuous depends on the values given. So based on the line equation it has been calculated along with mean square error (difference of Residuals and observation).
- 2) *K-Nearest Neighbours*: KNN is the other effective algorithm that takes for analysis without specifying the parameters and calculates based on the data similarity.
- 3) *Decision Tree*: Decision Tree is the tree based structured model. It selects the node by itself, from the input given and forms tree. From classification it differs in regression i.e., it takes average of every parameter and forms the root node with the highest influencing node. Then the remaining nodes follows ^[8].
- 4) *Random Forest*: It clearly for the large set of data that picks the variables randomly which fits for it. If the response is a factor, random Forest performs Classification; if the response is continuous, random Forest performs Regression. The unsupervised data is generally called as unlabelled data. It randomly picks up the predictors (i.e., group of decision trees) to form the model ^[10]. Categorical predictor variables must be specified as factors (or else they will be wrongly treated as continuous). In a Random Forest, each node is formed based upon the best predictors accordingly and form tree based structure like in decision tree with multiple combinations. This makes the Algorithm to perform well compared to other classifiers like support vector machines, Neural Networks, Discriminant analysis and its robust performance against overfitting.

III. FEATURE SELECTION

Feature Selection is the process before mapping the data into model. Feature Selection is done based on two ways with parameters and without parameters. The following processes explain about the basic Feature Selection in this work ^[2]. Features are classified as follows:

A. *Relevant*

These kind of features are highly influence towards the output. This feature cannot be replaced by any of the other.

B. *Irrelevant*

It does not relate at all. These kind of features will be eliminated by feature selection Technique. If this has been used it will give random results each time.

C. *Redundant*

The meaning or purpose of the parameter is repeating is termed to be redundant. It will make the model weak.

D. *Different Kinds Feature Selection Techniques (FST)*

There are several Feature Selection Techniques available, based on dataset the selection of each technique varies. In general, the most common Feature Selection Techniques is Forward and backward feature Selection process ^[2].

- 1) In corresponding to Forward FS, deleting each parameter one by one from beginning which leads to find the best features influencing the Target Variable.
- 2) Where the Backward FS represents the reverse form of Forward FS which eliminates from backwards.

In both the case addition of removed feature does not make any changes to get best feature. It will lead to random feature selection Process.

E. Feature's used for Prediction (with Parameter)

Identified almost 53 features, with one as target value^[3] for each post and categorized the features based on relation between Target Variable. Here 4 features have been identified as more influencing than other features.

- 1) *Page Features*: It defines about popularity/Likes of a page, check-in's, category of a page. *Page Likes*: This feature describes about the user specific interest related to page category such as Status, wall posts, Photos, Profile pic, shares or pages. *Pag*
- 2) *Category*: It specifies about the particular page that varies from other page like, Entertainment, Politics, brand or product, artist, music bands, place, tourism, medicine, company or institute etc.,
- 3) *Page Check-in's*: It tells about the person presence and the act on liking the post, pages and how many shares on particular pages etc., *Page Talking About*: It tells about the user who all are interested in each Facebook page category and 'engaged' with it. Specially about the user who coming back to that page after liking it. This includes some of activities like page share and likes to a post and comments to that post etc.,
- 4) *Essential Features*: The pattern of comment from different users on the post at various time interval with respect to randomly selected base time/date.
- 5) *Weekday Features*: It is for the complete week that indicates in binary values (0,1). This is used to pick the post that got published on selected base time/date.
- 6) *Other basic Features*: The remaining features that help to predict the volume of comment for each page category and that includes to document about the source of the page and date/time for about next H hours, document status volume (0,1) and the count of the post share. The remaining five feature are identified. From the above mentioned parameters of 53 are not completely required. Here after the feature selection found that 8 parameters which are highly prioritized. But performing with parameter in the algorithm didn't give expected results so worked without using parameter.
- 7) *Without using Parameter*: The prediction comes with expected way when performing without specifying any parameters in it. The regression gives the result which expected and termed as best prediction Results among the results that with specified parameters. `glm (formula = train_sacle$Target.Variable ~ ., data = train_sacle)` The above code represents the model with dot representing that Non-Parametrized form.

IV. EXPERIMENT SETTINGS

For this experiment, Data of Facebook page with user pattern in each page is taken for training and testing. In total there were 2,770 pages are modelled for 57,000 posts and 4,120,532 of comments using JQuery and Facebook Query Language(FQL). The sorted data (cleaned data) adds up to certain Giga bytes and this process takes up to certain weeks to model. After, the data is cleaned (After cleansing 5,892 posts and taken 51,108 posts remaining)^{[1][3]}. Then dividing the cleaned corpus into two different subsets using temporal splits, i.e., (1) Training data (80%, 40988) and (2) Testing data (20%, 10120) and then these data are sent to pre-processing modules where it has two divisions of datasets.



A. Training Dataset

Training Dataset is from the variant selection and calculation of it then vectorising it termed to be pre- processing.

B. Testing Dataset

In Testing also data are vectorised i.e., from 10,120 it had formed as 100 in each vector of total 10 as modelled.

In previous work [1], concluded that decision tree performs better compared to neural networks. Based on those results as part, decision tree results with other Algorithms like Linear Regression, Random Forest, Lasso and KNN has been made. Here, Decision Tree results are compared in both cases with parameter and without parameter towards other Algorithms. The prediction results in each algorithm of non-parametric form gives better results compared to parametric form. So, results here are for non-parameterized form. Here Random Forest Algorithm RMSE value is pretty high, but its prediction results came closer to the expected results but, it's not giving the constant results as well, due to its random FS process.

Here results of these algorithms are made for the Facebook dataset. First,

```
> summary(dtree)
```

Call

```
n= 40949
```

	CP	nsplit	rel error	xerror	xstd
1	0.11856017	0	1.0000000	1.0000339	0.08565948
2	0.04939177	2	0.7628797	0.7929352	0.07076855
3	0.04115737	3	0.7134879	0.7602787	0.06974433
4	0.04013157	5	0.6311731	0.7001560	0.06674591
5	0.02915479	6	0.5910416	0.6562073	0.06152217
6	0.02292093	7	0.5618868	0.6165994	0.05988566
7	0.01847200	8	0.5389658	0.5859780	0.05704875
8	0.01741446	9	0.5204938	0.5783039	0.05650934
9	0.01593535	10	0.5030794	0.5766578	0.05650324
10	0.01190385	11	0.4871440	0.5456269	0.05377161
11	0.01139054	12	0.4752402	0.5418436	0.05381691
12	0.01055486	13	0.4638496	0.5328157	0.05381680
13	0.01000000	15	0.4427399	0.5205916	0.05186179

Variable importance

Base.time	CC1
52	36
Page.talking.about	Page.Popularity.likes
6	3
Page.Category	Page.Checkins
2	1

In decision tree, the Root node error is the one to compute two measures of predictive variables, when according to the values shown in the relror and xerror column, and its depending on the CP (complexity of parameter in first column):

$0.71348 \times 0.76027 = 0.5424$ (54.2%) is the *re-substitution error rate in a 3rd row* (i.e., error rate computed on the training sample).

$0.76027 \times 0.06974 = 0.5302$ (53.02%) is the *cross-validated error rate* (using 10-fold CV, see xval in rpart.control(); but also see the results of xpred.rpart() and plotcp() which relies on the kind of measure). This measures are a more objective indicator of predictive accuracy.

Then the result of linear regression is explained below,

```
>summary(lm)
```

```
Call:glm(formula = train_sacle$Target.Variable ~ ., data = train_sacle)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-173.20	-8.82	-2.88	4.33	1280.70

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	7.32289	0.15983	45.818	< 2e-16 ***
Page.Popularity.likes	-1.07511	0.20714	-5.190	2.11e-07 ***
Page.Checkins.Ŷ	-0.74871	0.16276	-4.600	4.24e-06 ***
Page.talking.about	3.12843	0.22157	14.120	< 2e-16 ***
Page.Category	-0.58454	0.16290	-3.588	0.000333 ***
CC1	11.14411	0.17213		
64.741				< 2e-16 ***
Base.time	-8.42235	0.16029	-52.545	< 2e-16 ***
Post.length	0.08752	0.15994	0.547	0.584269

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Null deviance: 51588873 on 40948 degrees of freedom

Residual deviance: 42825449 on 40941 degrees of freedom

AIC: 400927

Number of Fisher Scoring iterations: 2

Std. Error is standard deviation error; that is, a point which is deviating from the regression line(residuals), it's called as Std. Error. Otherwise, the value which is deviating from the point at where the regression line is termed to be zero error state. Estimating the coefficient under standard regression model of the corresponding quantity (the coefficient estimated).

t value is the value of the t-statistic for which defined as the value obtained is differs from Zero.

Pr. is the p-value from hypothesis testing of statistic model along with t value. It refers alpha that termed to have 0.5 as constant. If the value is greater than 0.5 it termed to be likely and if it is less than 0.5 it termed to be failure and the results obtaining will be of unusual, otherwise if the null hypothesis were true.

Least Absolute Shrinkage and selection operator (LASSO), it creates a regression model that will penalize with the L1-norm which is the sum of the absolute coefficients. It makes the effect of shrinking the coefficients^{[11][12]}. This graph Fig.2., resembles a sine curve but not exactly because of the noise present in this model towards the data. So, it shows that LASSO won't be a good choice. There were three more graphs has been formed with the representation of Residuals but this sine wave graph shown below is bit more effective in representing the exact situation of data in the model.

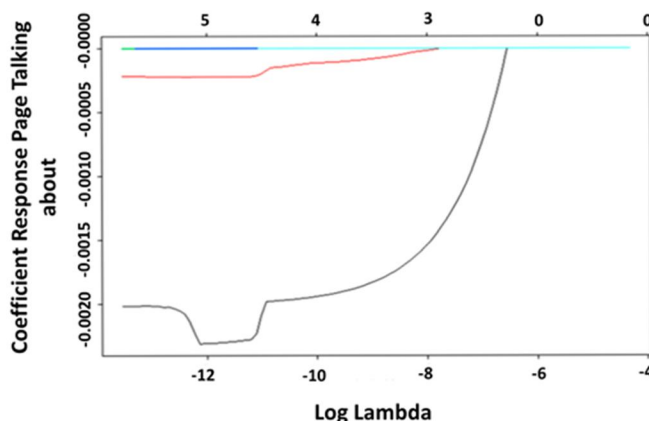


Fig.2. Graph of Lasso model

The above graph of LASSO model is formed with the Facebook data that originally deals with regression technique. Here the curve which trying to form sin curve but due to noise in the model towards the data it started increasing as well. So, the model has been failed and it shows the similar results of Ridge regression model. Result of RMSE values for different regression techniques which clearly shows that KNN gives the least error compared to other algorithms.

The graph showing KNN RMSE for different K-Values. In this Algorithm we can iterate the model with n number of K values. So that, it will help to pick any K value which giving minimum error rate. Here, at K=5; which is shown in the Fig.3. Since KNN will not allow to specify parameters, takes as whole data and performs. The variable selection is done with in the model like Lasso and Ridge does. So that it will assign the high influencing parameters from the data and also it has the K-value specified to give the prediction which all matches the K-values.

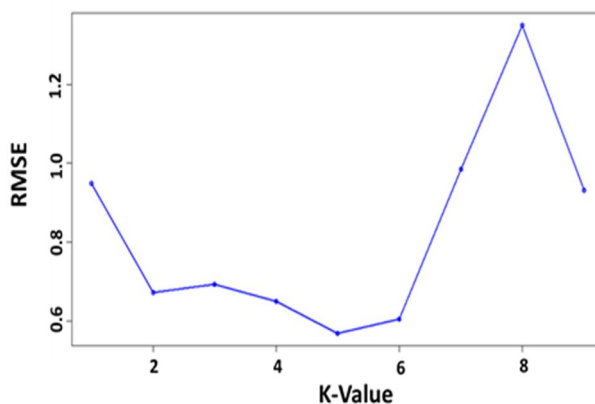


Fig.3. RMSE with respect to K-Values

The 3D plot which represents values with respect to main features (Base Time, Page category). So that the plot which represent the comments based on page and time(Hours) as well. Shown in. Fig.4. The colour from red to black represents the data. First value starts with red and finally ends at black. With Respect to Parameters Base time of each post based on page category, comments are predicted accordingly. So that the dependency of both parameters plays the major role in it.

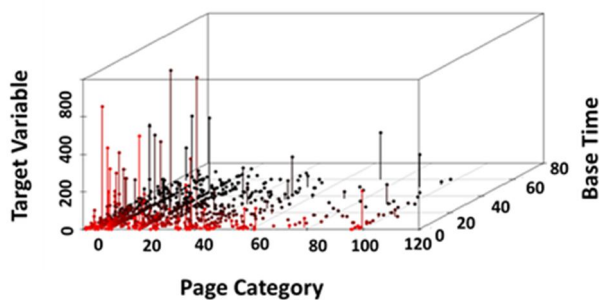


Fig.4. 3D plot of comment volume prediction.

Know, by selecting the different base date/tim randomly for each post with different variants, can be chosen to get good results. Then the clear plot of possible number of comments for each post category is displayed. Fig.5. The Comments for page category 18(Artists) received maximum of comments in next H hours among all. It’s a period of work done describing the efficiency of the model in it. This processes includes the time taken to train the data and taken for regression process and conclude the validation with test cases. This work can be made as a software application with some of image processing technique included to give the comment volume prediction.

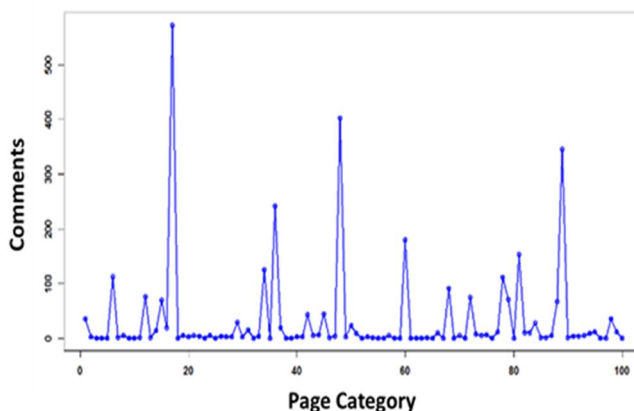


Fig.5. Plot of comment for each page category

V. RESULTS AND COMPARISONS

Thus the Regression Model for the Facebook dataset is concluding with KNN Regression Techniques, that is a Non Parametrized Model gives the accurate prediction results compared to other Algorithms. Here proving that with RMSE Results; Compared with the Algorithms like Linear, KNN, Random Forest (RF), Decision Tree (DT), Lasso and Ridge. In this bar chart below Fig.6, KNN RMSE is less compared to others. So, it clearly shows KNN works better.

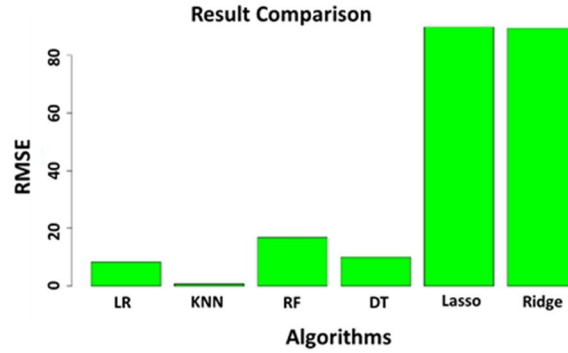


Fig.6. RMSE Result comparison

A. Comparing D-Tree and KNN with R-Squared

The graph in Fig.7. shows the comparison of decision tree with KNN results with respect to R-Squared. R-Squared metric is used to prove D-Tree in the previous work [1] So, tried to perform KNN with R-Squared, but the result is still better than D-Tree.

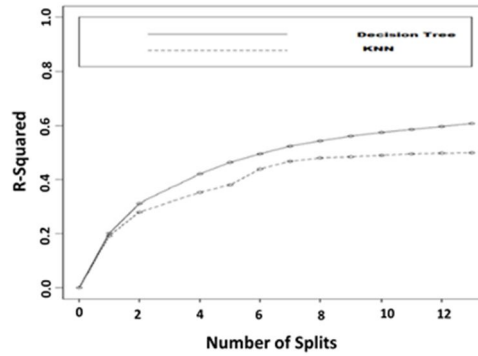


Fig.7. Decision Results Vs KNN

B. Statistical Testing

Statistical testing is for choosing best predictors among others. It is also called as hypothesis testing. Comparing with observed and expected values with Probability Function Constant called alpha value. If the fit is less than the alpha value observed value will be taken otherwise expected values will be taken [11]. From this process best predictors will obtain.

There are certain test cases for hypothesis testing for example, Chi Square, Anova, T-test, F-test, Pearson Correlation etc., For large set of data with different group of parameters and values of continues as well as categorical will be going for Anova Test.

C. Anova Testing

Anova tests significant among two or mor groups so that the predictors whic influencing more to the outcome is determined. Anova testing which can be implemented when the values are categorical as well as continuous. The Model helps to find the best predictors by generating four different graphs called Residuals Vs Fitted, Normal Vs Q-Q, Scale-Location and Residuals Vs Leverage.

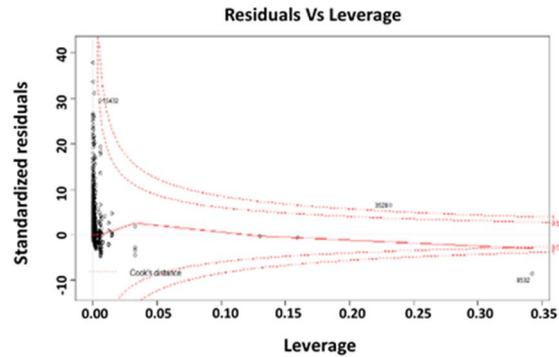


Fig.8. Plot of Residuals Vs Leverage in statistical testing

Above graph shows the result of residuals and Leverage i.e. which parameter is influencing more among other parameters. Residuals Vs Leverage graph is taken to get the clear image of observed and residuals data. Leverage explains that data which are having large difference. Shown in Fig.8.

Generally, in Hypothesis Testing F-statistics are the ratio of 2 different measures of data variance. If it is null hypothesis, then it leads to estimate the same value of ratio of around 1.

- 1) Here, the numerator is computed based on the variance mean and along with the true mean of each identical variance of data measured^[11].
- 2) But in other case, if null hypothesis is false and mean values are not at all equal, then this model measures larger.
- 3) Then coming to the denominator is an average of each sample variance group, which will estimate the overall population that assuming all groups having equal variance.

VI. CONCLUSION

This paper examines the Decision Tree, KNN, linear Regression and Random Forest results and concluding that KNN gives expected results compared to other Algorithms in comment volume prediction model. Moreover, with this examination, also shown that this model can be used for forecasting the comment volume perhaps choosing up of right variant is must. There is further a room for improvement using more features and with other regression techniques. The outcome of this process is a software application specially for a prediction of comment volume which can be enhanced further more using category based on predictors and by including some image processing features etc.

REFERENCES

- [1] Kamal jot Singh*, Ramjet Kaur, Dinesh Kumar "Comment Volume Prediction using Neural Networks and Decision Trees",2015 17th UKSIM-AMSS International Conference on Modelling and Simulation
- [2] M. Kara Giannopoulos, D. A Yfantis, S. B. Kotsiantis, P. E. Pintelas "Feature Selection for Regression Problems" Educational Software Development Laboratory,2000
- [3] UCI Machine Learning Dataset Forum. <https://archive.ics.uci.edu/ml/datasets/Facebook+Comment+Volume+Dataset>
- [4] K. Shvachko, H. Kuang, S. Radia, R. Chansler, "Prediction of bad comment for a video in Youtube" ,2010 IEEE 26th Symposium on, 2010, pp. 1–10.
- [5] Alan O. Sykes, The Inaugural Coase Lecture "An Introduction to Regression Analysis"1999
- [6] Eibe Frank University of Waikato New Zealand, "Machine Learning Techniques for Data Mining",2000



- [7] K. Buza, M. Spiliopoulou, L. Schmidt-Thieme, and R. Janning "Feedback prediction for blogs," in Data Analysis, Machine Learning and Knowledge Discovery, ser. Studies in Classification, Eds. Springer International Publishing, 2014, pp. 145–152.
- [8] Y. Zhao and Y. Zhang, "Comparison of decision tree methods for finding active objects", Advances in Space Research, vol. 41, no. 12, pp. 2008.
- [9] K. Shvachko, H. Kuang, S. Radia, and R. Chansler, "The Hadoop distributed file system," in Mass Storage Systems and Technologies(MSST), 2010 IEEE 26th sep 2010.
- [10] Tao Shi and Steve Horvath, "Unsupervised Learning with Random Forest Predictors" Journal of Computational and Graphical Statistics, Volume 15, Number 1, Pages 118–138
- [11] Hastie, T., Tibshirani, R., and Friedman, J. (2011). The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Springer Series in Statistics.
- [12] Hoerl, A.E. and Kennard, R. (2015). Ridge regression: Biased estimation for non-orthogonal problems. Technometrics, 12:55-67
- [13] T. Reuter, P. Cimiano, L. Drumond, K. Buza, L. Schmidt-Thieme, Scalable event-based clustering of social media via record linkage, techniques., in: ICWSM, 2016.
- [14] I. Polato, R. R' e, A. Goldman, F. Kon, A comprehensive view of hadoop researcha systematic literature review, Journal of Network and Computer Applications 46 (2016) 1–25.
- [15] Guy Lewis Steele, Jr. "Debunking the 'Expensive Procedure Call' Myth, or, Procedure Call Implementations Considered Harmful, or, Lambda: The Ultimate GOTO". MIT AI Lab. AI Lab Memo AIM-443. October 2015.
- [16] Knuth, Donald (1974), "Structured Programming with go-to Statements" (PDF), Computing Surveys, ACM, 6 (4), archived from the original (PDF) on 24 August 2009, retrieved 19 May 2013
- [17] Floating Point Benchmark: Comparing Languages (Fourmilog: None Dare Call It Reason)". Fourmilab.ch. 4 August 2005. Retrieved 14 December 2011
- [18] Box, G. E. P. (2010). "Some Theorems on Quadratic Forms Applied in the Study of Analysis of Variance Problems, II. Effects of Inequality of Variance and of Correlation Between Errors in the Two-Way Classification". The Annals of Mathematical Statistics. 25 (3): 484. doi:10.1214/aoms/1177728717
- [19] Caliński, Tadeusz; Kageyama, Sanpei (2000). Block designs: A Randomization approach, Volume I: Analysis. Lecture Notes in Statistics. 150. New York: Springer-Verlag. ISBN 0-387-98578-6.
- [20] Hettmansperger, T. P.; McKean, J. W. (2014). Edward Arnold, ed. Robust nonparametric statistical methods. Kendall's Library of Statistics. Volume 5 (First ed.). New York: John Wiley & Sons, Inc. pp. xiv+467 pp. ISBN 0-340-54937-8. MR 1604954
- [21] Wichita, Michael J. (2016). The coordinate-free approach to linear models. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge: Cambridge University Press. pp. xiv+199. ISBN 978-0-521-86842-6. MR 2283455.
- [22] Phadke, Madhav S. (2009). Quality Engineering using Robust Design. New Jersey: Prentice Hall PTR. ISBN 0-13-745167-9.
- [23] Montgomery (2011, Section 3-5.8: Experiments with a single factor: The analysis of variance; Practical interpretation of results; Comparing means with a control)
- [24] Anderson, David R.; Sweeney, Dennis J.; Williams, Thomas A. (2016). Statistics for business and economics (6th ed.). Minneapolis/St. Paul: West Pub. Co. pp. 452–453. ISBN 0-314-06378-1.



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)