

International Journal of Advanced Research in Computer and Communication Engineering

Vol. 7, Issue 11, November 2018

# Exemplifying the Significance of Tuning Tf-Idf for Sentiment Mining Online Consumer Review

Nandhini.S<sup>1</sup>, Dr.S.Prema<sup>2</sup>

M. Phil Research Scholar, Department of Computer Science (PG),

K. S. Rangasamy College of Arts and Science (Autonomous), Tiruchengode, Tamilnadu, India<sup>1</sup>

Associate Professor, Department of Computer Science (PG),

K. S. Rangasamy College of Arts and Science (Autonomous), Tiruchengode, Tamilnadu, India<sup>2</sup>

**Abstract**: Text mining have gain huge momentum in recent years, with user-generated content becoming widely available. One keyuse is remark mining, with much attention being given to sentiment analysis and opinion mining. An essential step in the process of comment mining is text pre-processing; a step in which each linguistic term is assigned with a weight that commonly increase with its appearance in the studied text, yet is offset by the occurrence of the term in the domain of interest. A common practice is to use the well-known tf-idf formula to calculate these weights. This paper reveals the bias introduce by between-participants' discourse to the study of comments in social media, and proposes an adjustment. We find that content extract from discourse is often highly correlated, resulting in dependence structures between observations in the study, thus introducing a statistical bias. Ignoring this bias can obvious in a non-robust analysis at best and can lead to an entirely wrong conclusion at worst. We propose a change to tf-idf that accounts for this bias. We show the effects of both the bias and correction with seven Facebook fan pages data, covering different domains, including news, finance, politics, sport, shopping, and entertainment.

Keywords: Sentiment Analysis, Text Mining, Statistical Bias, Discourse, TF-IDF

## I. INTRODUCTION

Data mining is the process of discovering patterns in huge data sets involving methods at the intersection of machine learning, statistics, and database systems. Data mining is an interdisciplinary subfield of computer science with an overall object to mine information (with intelligent methods) from a data set and transform the information into a comprehensible structure for further use. Data mining is the exploration step of the "knowledge discovery in databases" process, or KDD. Aside from the raw analysis step, it also involves database and data management aspect, data preprocessing, model and inference kindness, interestingness metrics, complexity consideration, post-processing of discovered structures, visualization, and online updating. The term "data mining" is in fact a misnomer, because the aim is the extraction of patterns and knowledge from large amounts of data, not the extraction (mining) of data itself. It also is a buzzword and is normally applied to any form of large-scale data or information processing (collection, extrac -tion, warehousing, analysis, and statistics) as well as any claim of computer decision support system, as well as artificial intelligence (e.g., machine learning) and business intelligence. The actual data mining task is the semiautomatic or routine analysis of large quantities of data to extract previously unknown, interesting patterns such as groups of data records (cluster analysis), unusual records (anomaly detection), and dependency (association rule mining, sequential pattern mining). These usually involve using database technique such as spatial indices. These patterns can then be seen as a type of summary of the input data, and may be used in further analysis or, for example, in machine learning and predictive analytics. For example, the data mining step might classify multiple groups in the data, which can then be used to obtain more accurate prediction results by a decision support system. Neither the data collection, data preparation, nor result understanding and reporting is part of the data mining step, but do belong to the overall KDD process as additional steps.

## II. LITERATURE REVIEW

A. TF-IDF: InbalYahavet.a1 [1]reveals the bias introduced by between-participants conversation to the study of comments in social media, and proposes an adjustment. We find that content extracted from conversation is often highly correlated, resulting in dependency structures between observations in the study, thus introducing a statistical bias. Ignoring this bias can manifest in a non-robust analysis at best and can lead to an entirely wrong end at worst. We propose an alter to tf-idf that accounts for this bias. We illustrate the property of both the bias and correction with seven



#### International Journal of Advanced Research in Computer and Communication Engineering

Vol. 7, Issue 11, November 2018

face book fan pages data, covering different domains, including news, finance, politics, sport, shopping, and entertainment. Jitendra kumarroutet et.al [2]with more consumers using online opinion review to notify their service decision making, opinion reviews have an economic impact on the bottom line of businesses. Unsurprisingly, opportunistic those or groups have attempted to abuse or manipulate online opinion review (e.g., spam reviews) to make ports and so on, and that detecting deceiving and fake opinion review is a topic of constant research interest. In this paper, we explain how semi-supervised learning methods can be used to detect spam reviews, prior to representing its utility using a data set of hotel reviews.

Kim Schoutenet.a1[4]using online consumer review as electronic word ofmouth to assist purchase-decision making has become more and more popular. The Web provide an extensive source of consumer review, but one can hardly read all reviews to obtain a fair estimate of a product or service. A text processing framework that can summarize review would therefore be desirable. A sub-task to be performed by such a framework would be to find the general aspect category addressed in review sentences, for which this paper presents two methods. In contrast to most existing approaches, the first method presented is an unsupervised method that applies organization rule mining on co-occurrence frequency data obtained from a corpus to find these aspect category.

Yuanlin Chenet.a1 [5]online transaction platforms provide a extremely convenient channel for consumers to generate and retrieve product reviews. In addition, consumers can also vote review perceived to be helpful in making their conclusion. However, due to various characteristics, consumers can have different preferences on products and review. Their voting behavior can be influenced by reviews and accessible review votes. To explore the authority mechanism of the reviewer, the review, and the existing votes on review helpfulness, we propose three hypotheses based on the consumer perception and perform statistical tests to confirm these hypotheses with actual review data from Amazon. Our empirical study indicates that review helpfulness has important correlation and trend with reviewers, review valance, and review votes.

Jo Mackiewicz et.a1[3]increasingly, professional and technical communicators analyse, synthesize, andreply to usergenerated content, including online consumer review of products, as the influence of user-generated content onconsumers' purchasing decision grows. But product review differ in the degree to which people perceive them to be credible. The analysed summary of existing work is given in table 1.

S. No	Title	Author, Publisher and Year	Working Platform	Objects	Future Scope
1	Comments Mining With TF-IDF: The Inherent Bias and Its	InbalYahav et.a1[1] IEEE [2015]	TF-IDF	The Inherent Bias and Its Removal.	Propose an adjustment to tf- idf that accounts for this bias.
	Removal				
2	Revisiting Semi- Supervised Learning for Online Deceptive Review Detection	JitendraKumar routet.a1[2] IEEE[2017]	Data Mining Online Review	Learning for Online Deceptive Review Detection.	Direction includes implementing and evaluating the proposed approach in the real-world.
3	Supervised and Unsupervised Aspect Category Detection for Sentiment Analysis with Co-occurrence Data	Kim Schouten et.a1[4] IEEE [2017]	Data Mining Online Review	Sentiment Analysis with Co- occurrence Data.	Proposed unsupervised method performs better than several simple baselines, a similar but supervised method, and a super-vised baseline, with an F <sub>1</sub> -score of 67%.
4	The Impact of Review Environment on Review Credibility	Jo Mackiewicz et.a1[3] IEEE [2018]	Data Mining Online Review	Review Environment on Review Credibility.	Managing user-generated product reviews, they should try to make credible content more noticeable to review users.
5	Analysis of Review Helpfulness Based on Consumer Perspective	Yuanlin Chen et.a1[5] IEEE [2018]	Data Mining Online Review	Review Helpfulness Based on Consumer Perspective.	More hypotheses and parameters to construct a quantitative helpfulness model.

#### Table I: Summary of the Literature Review



International Journal of Advanced Research in Computer and Communication Engineering

Vol. 7, Issue 11, November 2018

#### III. PROPOSED WORK

A. Classification Algorithm: Two classification algorithms are deployed for the selected pre-processed dataset.

- Naïve Bayesian algorithm
- SVM

Evaluation Metrics: The aforementioned dataset listed is applied for the two Classifiers. The classifier considered for this work is Naïve Bayes and SVM. Comparison between class detection accuracy was carried out. Evaluation metrics used are listed below,

- 1. Precision =  $\frac{\text{tp}}{\text{tp+fp}}$
- 2. Recall =  $\frac{\text{tp}}{\text{tp+fn}}$

3. Accuracy = 
$$\frac{\text{tp+tn}}{\text{tp+tn+fp+fn}}$$

B.Naïve Bayesian Considerations: Naive Bayes classifier is a simple probabilistic classifier based on relatingBayes' theorem with strong (naive) independence assumptions. This classifier assumes that the presence (or absence) of a particular feature of a class is unconnected to the presence (or absence) of any other feature, given the class variable. Even if these features depend on each other or upon the existence of the other structures, a naive Bayes classifier considers all of these properties to independently contribute to the probability of class. The possibility model for a classifier is a conditional modelp( $C|F_1, \ldots, F_n$ ) over a dependent class variable C with a small number of results or classes, conditional on several feature variables  $F_1$  through  $F_n$ . The problem is that if the number of features n is large or when a feature can take on a large number of values, then basing such a model on probability tables is infeasible.

C. SVM: Support Vector Machines are based on the model of decision planes that define decision boundaries. A decision plane is one that splits between a set of objects having different class memberships. SVM is a classifier method that performs classification tasks by constructing hyper planes in a multidimensional space that separates cases of different class labels. In Online Consumer Review OCR context, SVM classify {positive, negative} acts as a categorical value for classification. Linear kernel is used with epsilon 0.01, gamma 0.0, loss 0.1, and nu 0.5. Probability estimate and normalize set to false and shrinking based on function set to true.



Figure 1: Working Method of the Proposed System

## IV. ALGORITHM IMPLEMENTATION

- A. IDF Calculation
- 1. definverseDocumentFrequency(term, allDocuments):
- 2. numDocumentsWithThisTerm =0
- 3. fordoc inallDocuments:
- 4. ifterm.lower() inallDocuments[doc].lower().split():
- 5. numDocumentsWithThisTerm =numDocumentsWithThisTerm +1
- 6.ifnumDocumentsWithThisTerm> 0:
- 7. return1.0+log(float(len(allDocuments)) /numDocumentsWithThisTerm)
- 8. else
- 9. return1.0



International Journal of Advanced Research in Computer and Communication Engineering

Vol. 7, Issue 11, November 2018

B. TF Calculation deftermFrequency(term, document): normalizeDocument =document.lower().split()returnnormalizeDocument.count(term.lower()) /float(len(normalizeDocument))

## V. RESULT AND DISCUSSION

A. Term Frequency (TF): Term Frequency also known as TF measures the number of times a term (word) arises in a document.

B. Inverse Document Frequency (IDF): The main purpose of show a search is to find out relevant documents matching the query. In the first step all terms are measured equally important. In fact definite terms that occur too frequently have little power in determining the relevance. We need a way to weigh down the properties of too frequently occurring terms. Also the terms that occur less in the document can be more applicable. We need a way to weigh up the properties of less frequently occurring terms.

C. TF \* IDF: Remember we are annoying to find out relevant documents for the query: **life learning** For each term in the query increase its normalized term frequency with its IDF on each document. In Document1 for the term **life** the normalized term regularity is 0.1 and its IDF is 1.405507153. Multiplying them composed we get **0.140550715**(0.1 \* 1.405507153).

D. Vector Space Model – Cosine Similarity: From each document we originate a vector. If you need some review on vector refer here. The set of documents in a collection then is noticed as a set of vectors in a vector space. Each term will have its own alignment. Using the formula given below we can find out the similarity among any two documents. Cosine Similarity (d1, d2) = Dot product(d1, d2) / ||d1|| \* ||d2||

Dot product (d1, d2) = Dot product <math>(d1, d2) / ||d1|| + ||d2||  $||d1|| = square root (d1[0]^2 + d1[1]^2 + ... + d1[n]^2)$  $||d2|| = square root (d2[0]^2 + d2[1]^2 + ... + d2[n]^2)$ 

1: daes	אמיניין איז	inter deliniters ###(d));((0): max s m
Nominal	2. CEAL Nominal	
1	Now all Apple has to do is get swype on the johone and it will be crack. Iphone that is	
	Apple will be adding more carrier support to the iPhone 45 just announced	
5	Hilarious youtube video quy does a duet with apple s Siri Pretty much sums up the love affair http t co 8ExbnOjY	
)S	RIM you made it too easy for me to switch to Apple iPhone See ya	
os	I just realized that the reason I got into twitter was ios5 thanks apple	
s	I m a current Blackberry user little bit disappointed with it Should I move to Android or Apple iphone	
20	The 16 strangest things Siri has said so far I am SOOO glad that Apple gave Siri a sense of humor http t co TWAEUDBp via HappyPlace	
os	Great up dose personal event Apple tonight in Regent St store	
DS	From which companies do you experience the best customer service aside from zappos and apple	
os	Just apply for a job at Apple hope they call me lol	
os	RT JamaicanIdler Lmao I think apple is onto something magical I am DYING haha Siri suggested where to find whores and where to h	
os	Lmao I think apple is onto something magical I am DYING haha Siri suggested where to find whores and where to hide a body lolol	
os	RT PhillipRowntree Just registered as an apple developer Here s hoping I can actually do it Any help greatly appreciated	
os	Wow Great deals on refurbed iPad first gen models RT Apple offers great deals on refurbished 1st gen iPads http t co ukWOKBGd Apple	
os	Just registered as an apple developer Here s hoping I can actually do it Any help greatly appreciated	
os	Currently learning Mandarin for my upcoming trip to Hong Kong. I gotta hand it to Apple iPhones their uber useful flashcard apps	
os	Come to the dark side gretchenedark Hey apple if you send me a free iPhone I will publidy and ceremoniously burn my BlackBerry	
os	Hey apple if you send me a free iPhone any version will do I will publicly and ceremoniously burn my BlackBerry	
os	Thank you apple for Find My Mac just located and wiped my stolen Air smallvictory thievingbastards	
os	Thanks to Apple Covent Garden GeniusBar for replacing my MacBook keyboard cracked wristpad during my lunch break today out of warranty	
os	DailyDealChat apple Thanks	
Pos	Pads Replace Bound Playbooks on Some N F L Teams http: t co 2UXAWKwf apple nytimes	
os	apple good ipad	
os	apple siri is efffing amazing	
os	Amazing new Apple iOs 5 feature http://to.jatFVfpM	
os	RT TripLingo We re one of a few Featured Education Apps on the Apple Website today sweet http t co 0yWvbe1Z	
los	We re one of a few Featured Education Apps on the Apple Website today sweet http t co 0yWvbe1Z	
os	When you want something done right you do it yourself or go to Apple ATT you re useless these days yourdaysarenumbered	
los	We did an unexpected workshop for the iPhone4S at apple yesterday and we got an awesome amount of info notjustaboutthephone gamerchik 16	
os	It 3 ios5 apple	
los	RT Apple No question bro RT AintEeenTrippin Should I get dis iPhone or a EVO 3D?	
Pos	RT inightbewrong I m OVER people bitching about the iPhone45 I think its the smartest phone I ve ever had and I m very happy	
os	I m OVER people bitching about the iPhone45 I think its the smartest phone I ve ever had and I m very happy Way to go Apple	
los	Twitter CEO points to Apple as corporate mentor as iOS signups triple http://tio.co.gcY8iphN	
os	At the bus with my iPhone thxx apple	
los	azee1v1 apple umber AppStore is well done so is iTunes on the mobile devices I was talking about desktop app	
os	NYTimes Coach Wants to See You And Bring Your iPad http: t co J2FTIEnG iPad apple set red 42 red 42 hut hut NFL wish I had an iPad	
		· · · · · · · · · · · · · · · · ·

Figure 2: Sample Data Set



#### International Journal of Advanced Research in Computer and Communication Engineering

Vol. 7, Issue 11, November 2018



#### Figure 3:Generic Object Editor

🥥 Weka Explorer		Street, Street			1.00				- 0 X		
Preprocess Classify Cluster Associ	iate Select attributes Visualize										
Open file	Open URL	Open DB	Gener	ate		Undo	Edit	Save.			
Filter											
Choose StringToWordVector -R first-last -W 1000 -prune-rate -1.0 -N 0 -S -stemmer weka.core.stemmers.lterated.ovinsStemmer -M 1 -tokenizer "weka.core.tokenizers.NGramTokenizer -delimiters \" \ /\ n\\\.,:\\\\\\\\?\?\" -max 3 -min 1"											
Current relation Relation: twitter-sanders-applenew	v-weka.filters.unsupervised.attribute.St	iringToWordVector-R-W1000-p	Selected attribute Name: dass Type: Nominal								
Instances: 479		Su	ım of weights: 479	Missing	: 0 (0%)	Distinct	: 2	Jnique: 0 (0%)			
Attributes				No.	Label	Co	punt	Weight			
All	None	Invert P	attern		2 Neg	31	6	316.0			
No Name						I					
1 dass											
2 text											
				Class: tex	t (Nom)			•	Visualize All		
				163							
	Remove										
Status OK								Log	×0		
🚱 🖉 🚞	D 🔮 🧿 🛛		YA.	1		10.10	n	- 隆 🛱 🔶	10:35 AM 9/27/2018		

Figure 4:Weka Explorer



International Journal of Advanced Research in Computer and Communication Engineering

Vol. 7, Issue 11, November 2018



Figure 5:ROC Curve



Figure 6:Cost Benefit Analysis

## CONCLUSION

This research deploys the TF-IDF tuning along with specified preprocessing techniques which is exemplified as the proposed method working scenario in the next chapter. This chapter illustrates the significance of the TF-IDF tuning in document. The proposed method of preprocessing works well to attain the accuracy. The Naïve Bayesian and SVM both work well with the pre-processed dataset. Hence in terms of the evaluation metrics SVM performs well which has been implemented in MATLAB. It is concluded that the projected method accomplishes well when equated with the prevailing methods.

#### **FUTURE WORK**

Text mining is the potential area, which seek lots of importance due to the nature of the data it handles. Since the social Medias seek large place in our daily lives, the importance of sentiment mining also hikes. Incorporating N-Gram method and Bag of Words method along with fuzzy text categorization would elevate the performance of the research to the next level. That would be the further enhancement for the research.



International Journal of Advanced Research in Computer and Communication Engineering

Vol. 7, Issue 11, November 2018

#### REFERENCES

- [1]. Amrita Kaur, NeelamDuhan. "A Survey on Sentiment Analysis and Opinion Mining." IEEE International Journal of Innovations & Advancement in Computer Science IJIACS IEEE, 2015.
- [2]. DivyaBohra, Sanjay Deshmukh. "A Survey on Sentiment Analysis in NLP." IEEE International Journal of Advanced Research in Computer and Communication Engineering. IEEE,2015.
- [3]. Horakova, Marketa. "Sentiment Analysis Tool using Machine Learning." IEEE Global Journal on Technology. IEEE, 2015.
- [4]. InbalYahav, OnnShehory, and David Schwartz." Comments Mining With TF-IDF: TheInherent Bias and Its Removal." IEEE Transactions on Knowledge and Data Engineering, VOL.14, NO.8, pp.1-14.IEEE, August 2015.
- [5]. Jitendra Kumar Rout, AnmolDalmia, Kim-Kwang Raymond Choo, SambitBakshi, and Sanjay Kumar Jena. "Revisiting Semi-Supervised Learning for Online Deceptive Review Detection." Vol.5,pp.1319-1327.IEEE, January 2017.
  [6]. Jo Mackiewicz andDave Yeats. "Product Review Users' Perceptions of Review Quality: The Role of Credibility, Informativeness, and
- Readability." IEEE Transactions on Professional Communication, Vol.57, No.4, pp.309-324. IEEE, December 2014.
- Jo Mackiewicz, Dave Yeats, and Thomas Thornton."The Impact of Review Environment on Review Credibility." IEEE Transactions on [7]. Professional Communication, VOL.59, NO. 2, pp.71-88. IEEE June 2016.
- [8]. Kim Schouten, Onne van der Weijde, Flavius Frasincar, and Rommert Dekker."Supervised and Unsupervised Aspect Category Detection for Sentiment Analysis with Co-occurrence Data."IEEE Transactions on Cybernetics, VOL.48, NO.4, pp.1263-1275. IEEE, April 2018
- Nikhil R, Nikhil Tikoo, SukritKurle; HariSravanPisupati, Dr. Prasad G R. "A Survey on Text Mining and Sentiment Analysis for Unstructured [9]. Web Data."IEEE Journal of Emerging Technologies and Innovative Research (JETIR).IEEE, 2015.
   [10]. Saurin Dave, Prof.HiteishiDiwanji. "Trend Analysis in Social Networking using Opinion Mining a Survey." IJSRSET IEEE, 2015.