

# Simple way of minimum redundancy maximum relevance feature selection approach for Korean review data

*Ju yeon Heo, Kun chang Lee*

[heojuyeon12@gmail.com](mailto:heojuyeon12@gmail.com)

[kunchanglee@gmail.com](mailto:kunchanglee@gmail.com)

*Sungkyunkwan University*

---

## Abstract

Feature selection plays an important role in sentimental analysis to perform well. In this paper, a supervised feature selection approach is presented, which is considering not only feature-feature and feature-target variable relationship but also trying to include only mutually exclusive and completely exhaustive features for efficiency. Applying a clustering, the two methods searches for a best set of features based on hybrid approach which is combination between filter and wrapper methods. Eventually, several experiments using 4 Korean datasets are presented to show the effectiveness of the features selected with two methods from the point of view of the number of features needed to reach certain classification accuracy.

---

## 1. Introduction

### 1.1 Importance of feature selection

As social network services become more and more popular, we can obtain the information of social tendencies and preferences [1,2]. Consumer's evaluations of products or service are getting important because they are not only used for products or brand development, but also they impact other potential customers. The aim of sentiment analysis is to classify opinions and feelings from enormous reviews into two-class or multi-class polarity by creating accurate machine learning models.

In text datasets, both useful and useless high dimensional features which are contained. Because fruitless features are un-evenly distributed, irrelevant and noisy, enormous features and dimensions decrease the efficiency and performance [3]. Thus, the performance of sentiment analysis in text mining highly depends on feature selection. It is important in both supervised and unsupervised sentiment analysis by reducing dimensions and making models focus on most informative features [4]. Also, it can beat the problem of overfitting and can lower the cost of time and computational resources.

### 1.2. Feature selection study about sentiment analysis in Korean

Because of grammatical uniqueness by language, it is often difficult to stipulate them with a consistent rule [23]. However, the studies on Korean sentiment analysis have the difficulties and there are still fewer studies on Korean sentiment analysis than English sentiment analysis, especially on the models

that considering both official reviews and microblogging. Thus, this paper will be meaningful in studying feature selection methods in Korean data including both reviews and microblogging.

## 1.2. Previous feature selection methods

For the feature selection, there are 2 frameworks that are search-based framework and correlation-based framework.

### 1.2.1. search-based feature selection framework

For the search-based feature selection framework, 3 approaches are suggested; filtering, wrapper and embedded approaches. The hybrid method that are the mixture of filter and wrapper methods is proposed for complimenting the cons of both filter and wrapper methods. It can decrease cost compared to wrapper method and get rid of the limitation of the dependence of the feature subset selection in filter method [5]. This paper takes hybrid approaches.

### 1.2.2. correlation-based framework

The second framework, correlation-based framework, considers both feature–feature correlation and feature–class correlation. Generally, the feature–feature correlation is called as feature redundancy, while the correlation between features and class is viewed as feature relevance. The chi-square method that are used most for weighting features focus on how to select high feature relevance. If the target variable, for example class label, is independent of the feature variable, we can discard that feature variable. If they are dependent, the feature variable is important. However, they are not considering feature redundancy and they include all features that are relevant to target variables even if they are overlapped. If the features are sharing similar characteristics, we don't have to include all these features for feature vectors. For example, 'Love story' and 'Romance' in movie domain have really similar meanings and they have high possibilities to occur together in a same document. Even if both features have high feature relevance and highly dependent to the target variables, selecting both features can be inefficient. For considering both correlation, mRMR(minimum redundancy maximum relevance) is studied actively [13,14]. However, there's only a few works on considering both redundancy and relevance in feature selection with hybrid approaches in sentiment analysis. Also, for the best of our knowledge, few articles considered the meaning of features that have high redundancy. Especially these kind of researches about Korean are scarce in sentiment analysis study fields. Studying in the field of not only official reviews but also microblogging in Korean is valuable when there are not many researches. Thus, this paper is going to compare the performance of three hybrid approaches. First, our two methods, ToCchi and ToCknn, second, chi-square methods hybrid approaches and last, other methods presented in 2017(OIFV) [5].

## 1.3. Proposals

This paper takes a different look at mutually exclusive and completely exhaustive feature selection based on document prospects.

### 1.3.1. clustering the features based on document appearance.

We assume that features that appear similarly related to the documents share the similar characteristics. The things that features can share can be the similar topics, similar meaning, and similar emotions that can be used almost together <Figure 1>. For example, ‘ㅌ’ and ‘ㅌㅌㅌ’ are the emoticons used in Korean that mean sad and they occurred almost in a same sentence or document. Thus, selecting the features equally balanced from the whole datasets not only can deal with the diverse and smaller features of the whole datasets harmoniously. But also it can be efficient and reduce the feature dimension without sacrifice representing the data.

**Figure 1. Characteristics that features in the same cluster share**

Similar topics		Similar meaning		Appearance almost together	
'캐릭터'	'우체국'	'로맨스'	(sad emoticon)	'명예훼손'	(exclamation for something)
'character'	'post office'	'Romance'	'ㅌㅌㅌ'	'defamation'	
'배우'	'택배'	'러브스토리'	'ㅌㅌ'	'침해'	'완전'
'actor'	'delivery'	'Love story'	'ㅌ'	'invasion'	'대박'
'연기'	'배송'	'사랑'	'ㅌ'	'타인'	'미치다'
'act'	'deliver'	'Love'	(happy emoticon)	'others(formal)'	'진짜'
'액션'	'도착'	'애정선'	'ㅎㅎ'	'욕설'	'정말'
'action'	'arrive'	'Affection'	'ㅋㅋ'	'swear word'	
'감정'	'상태'				
'emotion'	'condition'				

Also, our methods can play role as the goal of TF-IDF method. This is important because the number of reviews about the certain topics are the dependent to the training data, so choosing the features based on term frequency can lead to overfitting problem. The data is only collected for certain periods and it is only the small part of the whole reviews. Thus, when we depend on the frequency of the features, it can lead to wrong feature selection. However, by using our proposed methods, in the movie reviews, we can separate reviews of actors and the one of directions and then consider both topics equally balanced even if the number of reviews of actors are much larger than the one of directions. In other word, our methods can work as the role of TF-IDF unexpectedly.

### 1.3.1. chi-square and K-nearest neighbors within the cluster

We used both chi-square and KNN for ordering the features within the cluster. Because the chi-square method is focusing on feature relation, chi-square scoring after clustering the features can consider feature relevance and feature redundancy at once. On the other hand, the clustering - KNN method scores the features that are in core of the clusters high. Even if it can reflect high feature relevance and MECE, they can't represent the whole cluster because they put center features before other features. Also, it doesn't

take account of feature redundancy. Thus, we predict that chi-square after clustering method will show better performance than KNN after clustering.

As a result, we expect considering selecting features balanced related to the features can increase the performance by reducing the number of unnecessary features. In other words, we can reach high accuracy with only small number of features by mutually exclusive and completely exhaustive feature selection. This paper proposed two method with hybrid feature selection approaches; Topic clustering with CHI(ToCchi) and Topic clustering with kNN(ToCknn) model. Comparative experiments were done on 4 datasets with different themes written in Korean.

#### 1.4. Research Questions

RQ1 : Do ToCchi and ToCknn outperform in feature selection compare to base or previous hybrid method?

RQ2 : Which model would show better accuracy or stable performance in feature selection?

The rest of paper is arranged accordingly. Section 2 indicates several related works in sentiment analysis and Section 3 illustrates the aspects of 4 datasets and the way of representation. In Section 4, proposed methodologies for feature selection are described with algorithms. The comparative experiments and evaluation is showed in Section 5. Lastly, conclusion and future works are given in Section 6.

## 2. Related works

### 2.1. Sentiment analysis

Sentiment analysis is studied well from early 2000 in NLP [15]. Numerous approaches has been reported to classify sentiment from languages, both in supervised and unsupervised method [16].

The main process of sentiment analysis can be divided into three steps. The first step includes representation of the data and second one involves extraction and selection of the features from large sets. The last is classifying each sample into binary or multi classes with single-classifiers such as NB, LR and ensemble-classifiers like RF, AB [6].

In this process, feature selection is the key steps because high-dimensional feature space is a significant challenge in text mining. It can reduce lots of cost occurred in computational time and power and increase efficiency of text analysis [7]. Also, it can improve the performance by eliminating unnecessary text features that increases noisy [9].

#### 2.2.1. Feature selection study about sentiment analysis in Korean

Unlike general structured data, texts in different languages have lexical and grammatical uniqueness by language. And as the forms of expression are varied and complex, it is often difficult to stipulate them with a consistent rule [23]. The studies on Korean sentiment analysis have the difficulties. The use of artificially modified Korean makes lexical analysis difficult. The combination of consonants and vowels

made up the Korean language. However, the texts in SNS often distorted, such as using consonants or vowels only and not the combination for emoticons. In that case, the features will be exploding and it will be hard to tell the emoticons and frequently used typos. There are still fewer studies on Korean sentiment analysis than English sentiment analysis, especially on the models that considering both official reviews and microblogging. The research of emotion analysis in Korean microblog texts [22] applied the machine learning model based on Korean documents and classified human sentiments into seven emotions.

## 2.2. Previous feature selection methods

### 2.2.1 search-based feature selection framework

Feature selection techniques are categorized into filter, wrapper, and embedded methods [8]. Filter methods are for statistical scoring of text feature excluding the consideration of performance of learning mechanisms. Wrapper methods are trying to find the best subset of features by considering the interaction with learning algorithms. In the embedded approach, the feature selection is done in the process of training algorithms and the best subset of features are found by a classifier used. The hybrid method that are the mixture of filter and wrapper methods are also included in the selection categories [10]. It is the complement of the filter and wrapper that reduces computational cost compared to wrapper method and overcomes the limitation of the dependence of the feature subset selection in filter method. Yousefpour, Ibrahim, 2017 investigated the hybrid methods of feature selection for text classification. They scored features with 5 filter methods and tested the performance of each subsets from the vector generated by ordinal-based and frequency-based integration (OIFV) [5]. Wang, Suge, et al also researched the hybrid feature selecting method based on category distinguishing ability of words and information gain [11]. They concluded that hybrid methods is superior to the one with directly using information gain. Le Nguyen Hoai Nam et al proposed a hybrid filter feature selection, called FCFS and related filter feature selection methods as CMFS, OCFS, CIIC, IG, CHI with two datasets about news and medicine [21].

These papers focus on the hybrid method by integrating the process used in filter and wrapper methods. In the first step, features within the topics are ordered to create feature vectors based on the score of relation between the feature and the each topic. The next is making subset of the vector and evaluating the performance of each subset and finding the best one which is the method of wrapper approach.

### 2.2.2 correlation-based framework

The correlation-based framework considers both feature–feature correlation and feature–class correlation. First in the feature-feature correlation study, Vinh, Nguyen X., and James Bailey, a method for supervised feature selection based on clustering the features into groups is proposed, using a conditional mutual information based distance measure. They find that there is a reasonable condition, namely when all features are independent given the class variable (as assumed by the popular naive Bayes classifier) [12]. Second, in the feature-class correlation, the method by Sotoca, José Martínez, and Filiberto Pla builds a dissimilarity space using information theoretic measures, in particular

conditional mutual information between features with respect to a relevant variable that represents the class labels [20].

For considering both correlation, mRMR which means minimum redundancy maximum relevance is studying actively. In the work of Assi, E. Bou, et al, based on a Support Vector Machine and an Adaptive Neuro Fuzzy inference system, data reduction was performed by mRMR features selection approach for electrodes selection and a genetic algorithm. The selected subset of features performed equally and sometimes even better than the whole features set [13]. In the article of Niphath Claypo and Saichon Jaiyen, the mRMR feature selection is used to select the features of data in order to reduce the number of features in the data set. Consequently, the computational times of learning algorithms are reduced for neural networks with mRMR approaches based on Thai restaurant reviews [14].

### 2.3. Related works about proposals

When generating the feature vectors in the first step in hybrid feature selection approaches, several methods are used and studied for its performance. According to Yousefpour et al, feature selection methods such as the IG and CHI methods were found to achieve better accuracy than other methods in filter approaches. Compare to the CHI method, kNN is less studied in feature selection process [5]. kNN score the features according to how they share the similar patterns or characteristics. If we set one feature which is strongly related to the topic as a center and use the kNN method within the specific topic group of features, we can pick features that have strong patterns within themselves and they can represent the topic group well. In other word, kNN method also can be an indicator that can show the relationship between features and the topic as CHI can be. In this paper, we compare the CHI and kNN method in terms of indicator that can score the features with strong association to the topic before using wrapper method in hybrid approach. This article proposed two method in scoring the relevance which is the **ToCchi** and **ToCknn**.

## 3. Datasets and feature preprocessing

The detection of features from the datasets should be done before selecting feature for classifier models. Feature representation are done with several methods such as bag-of-words [16], lexicon etc. The unigram BOW is adopted for feature extraction.

**Table1. Information of datasets**

Datasets	subject	# of positive	# of negative	total	# of features After CountVectoizer	Language	Gathering period
Google App Store	Various app review	1688	1845	3533	3374	Korean	2017.7~8

11Street Shopping	Various products	9551	10449	20000	6378	Korean	2017.5~6
Twitter	Idol (Wanna One)	6792	7751	14543	6382	Korean	2017.7~8
Watcha-movie	Movie	4215	3481	7696	6352	Korean	2011

## (1) Datasets

### Google App Store

We collected the data from Google App store by crawling the reviews of various apps. The Google App is the biggest market of the apps on Android OS. For each download page, there are reviews and the ratings of the apps. The problems and the benefits of the apps are pointed by numerous active users.

### Twitter

We collected the data from Twitter about ‘Wanna One’ which is popular idol singer group in Korea. It is picked out as a search word because it is a hot topic in July, 2017 so that there are various sentiment reactions on Twitter. We gather the results of the topic except for sentences repeated meaninglessly, advertisement and sentences which are too short. For the rating of the text data, we employed the students and ask them for making polarity consisted of positive and negative.

### 11 street shopping

In ‘11 street shopping’<sup>1</sup>, consumers can compare the price of products from various categories and read the reviews and ratings from 1 to 5. We save the reviews of several categories with active reviews from May to June, 2017 and divided into positive and negative based on ratings. Due to the characteristic of the culture in this shopping site, customers almost rate the products from 3 to 5. Thus, we give positive polarity to under 4 and negative polarity to the upper 4 after the relationship between reviews and ratings are checked by humans.

### Watcha movie

‘Whatcha’<sup>2</sup> is the movie recommend site in Korea. After accumulating the data of preference movie of users, they advise some movies to watch. It has lots of movie fan users and active review system. Diverse users write down their feelings and evaluations and they write down opinion, impression and deep thoughts of movies and its topics. We choose movies among

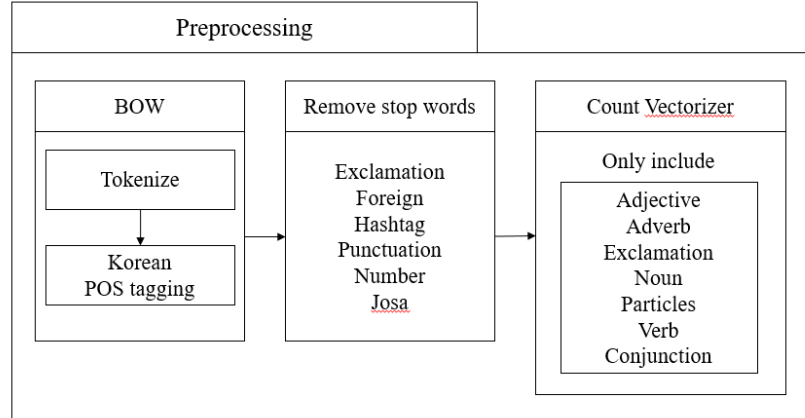
---

<sup>1</sup> [www.11st.co.kr/](http://www.11st.co.kr/)

<sup>2</sup> <https://watcha.net/>

movies passed over million viewers mark and place sentences with score 5 to positive and 0.5, 1 to negative.

## (2) Feature preprocessing



To annotate sentence with POS, this article used the ‘KoNLPy’. It splits the word into morphemes. After being excluded stop words which is above, features in datasets are selected including POS mentioned above by ‘Count Vectorizer’ in Sklearn open source<sup>3</sup>. The particles are contained because they are used as emoticons widely in Korean. For example, ‘ㅋㅋ’, ‘ㅎㅎ’ show laughing and ‘ㅠㅠ’ shows crying.

## (3) Methods used in feature selection

### Chi-square

The Chi-square calculates the degree of the relationship between the feature and the category.

$$CHI(f, c_i) = \frac{N \times (AB - CD)^2}{(A + C) \times (A + D) \times (B + C) \times (B + D)}$$

$$i = 1, 2, 3, \dots, M$$

Where f: feature, c: category, N: number of all documents, M: the number of the categories (data in this paper use M=2 for binary polarity), A: the number of times f and c occur together, B: the number of time neither c nor f occurs, C: the number of times f occurs without c and D is the number of times c occurs without f.

<sup>3</sup> [http://scikit-learn.org/stable/modules/generated/sklearn.feature\\_extraction.text.CountVectorizer.html](http://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html)



## **kNN**

K-Nearest Neighbor is a widely applied text classifier due to its simplicity and efficiency. Its training-stage consists of storing all training examples as classifier, thus it has often been called as lazy learner since ‘it defers the decision on how to generalize beyond the training data until each new query instance is encountered’ [18].

## **Kmeans**

K-means clustering [19] is a method commonly used to automatically partition a data set into  $k$  groups. After selecting  $k$  initial cluster centers, this model iteratively refine them. 1. Each instance  $di$  is assigned to its closest cluster center. 2. Each cluster center  $Cj$  is updated to be the mean of its constituent instances. The algorithm converges when there is no further change in assignment of instances to clusters.

## **4. Methodologies**

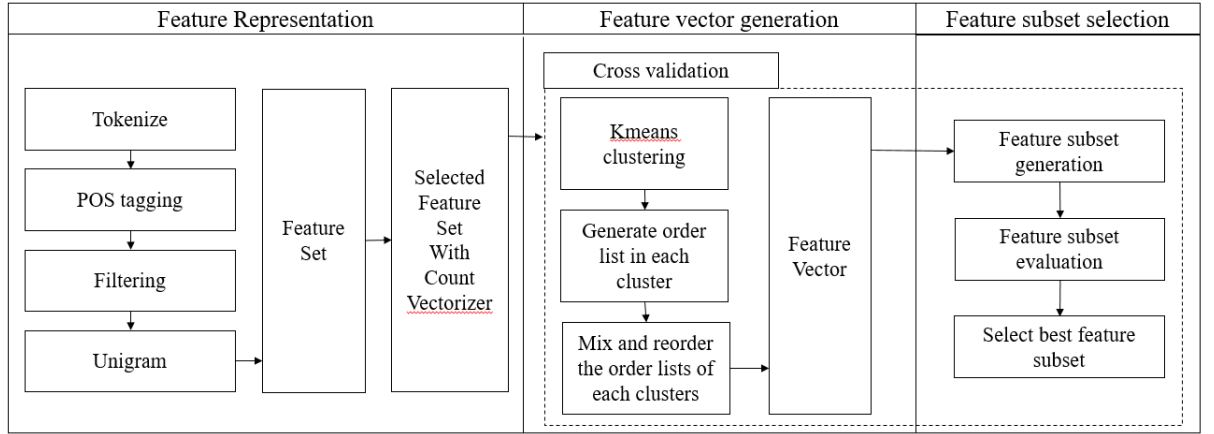
This paper creates feature vectors based on scoring and ordinal feature subsets as a hybrid approach. This is better to overcome the limitation occurred when choosing the size of feature subsets in filter methods and reduce the computational cost compared to wrapper methods because it decreases the number of possible feature subsets.

The research proposed two methods to generate feature vectors which are the core stage in the whole feature selection process. The first method is Topic clustering with CHI(ToCchi) and the second one is Topic clustering with kNN(ToCknn)for feature selection.

These methods consider different clusters of features to select the features equally balanced. They assume that selecting the features equally balanced can increase the performance by reducing the number of unnecessary features with feature redundancy, feature relevance, and MECE. Selecting features without considering the characteristics of the features is not efficient way because it can only choose features about one or two specific topics or contain similar features sharing same topic, same meaning or same emotions that can be used almost together. The data is only collected for certain periods and it is only the small part of the whole reviews. If the model depends on the number of appearance of the features in the training datasets, it can be overfitting because the frequency of the features is the characteristic of the data.

For the ToCknn, the center of the feature is selected based on the nearest centroid scoring. Then k-nearest neighbor method is used to make an ordinal list of features. For the ToCchi, the feature vector is made based on the chi-square score of the features. The unigram BOW with counter vectorized process for feature selection which are mentioned above in the section 3 are used as the way of representing the raw text sets in both methods.

***Figure2. Process of ToCchi and ToCknn models***



### Stages in method

#### Stage 1: Clustering based on the sentence or documents.

The dimensions are each documents and data points represent original features which are unigram BOW. By using Kmeans, features are clustered into several groups as the parameter of the number of clustering changes. Because it is hard to decide hyper parameter for the number of clusters, cross validation is used to find best set.

#### Stage 2: Creating feature vectors within each clusters.

$Features = [f_1, f_2, f_3, f_4, \dots, f_N]$  are the set of features, where n is the number of features. For the ToCchi method, features are ordered based on its chi-square scores. On the other hand, for the ToCknn method, first center feature is selected based on the Chi-square among the features. Next, for each clusters, the list of features are generated as the order of relevance with the center by using K-nearest neighbors.

##### (1) ToCchi

Features order by chi-square in cluster A:

$$F_{CHI-a} = [f_{100}, f_2, f_{17}, f_{20}, \dots, f_N]$$

$$= [f_{CHI-1}, f_{CHI-2}, f_{CHI-3}, \dots, f_{CHI-N}]$$

##### (2) ToCknn

First center feature based on  $NearestCentroid = f_{KNN1}$

Features order by kNN of the center ( $= f_{KNN1}$ ):  $\mu_X f_{KNN1}$

$$F_a = [f_{Cluster\_a1}, f_{Cluster\_a2}, f_{Cluster\_a3}, \dots, f_{Cluster\_aN}]$$

$$= [f_{KNN-a1}, f_{KNN-a2}, f_{KNN-a3}, f_{KNN-a4}, \dots, f_{KNN-aN}]$$

fCHI1 is the center and fCHI100 is the nearest features from the center and fCHI3 is the furthestmost.

### Stage 3: Generating feature subsets and evaluation.

Each vectors are combined as the sequence of indexes.

$$F_a = [f_{Cluster\_a1}, f_{Cluster\_a2}, f_{Cluster\_a3}, \dots, f_{Cluster\_aN}]$$

$$F_b = [f_{Cluster\_b1}, f_{Cluster\_b2}, f_{Cluster\_b3}, \dots, f_{Cluster\_bN}]$$

$$F_c = [f_{Cluster\_c1}, f_{Cluster\_c2}, f_{Cluster\_c3}, \dots, f_{Cluster\_cN}]$$

↓

$$FV_{final} = [f_{Cluster\_a1}, f_{Cluster\_b1}, f_{Cluster\_c1}, f_{Cluster\_a2}, f_{Cluster\_b2}, f_{Cluster\_c2} \dots f_{Cluster\_aN}, f_{Cluster\_bN}, f_{Cluster\_cN}]$$

### Stage 4: Generating feature subsets and evaluation

After making the feature vectors, feature subsets are made as follows:

$$ToChi, ToCknn = [x_1, x_2, \dots, x_N], \quad \forall_{i,j} i < j \rightarrow rank(x_i) \geq rank(x_j)$$

$$Feature\ subsets = \{ \{x_1\}, \{x_1, x_2\}, \{x_1, x_2, x_3\}, \dots, \{x_1, x_2, \dots, x_N\} \}$$

Where  $x_i$  is the feature and  $N$  is the total number of features.  $x_1$  has the highest rank and  $x_2$  has second highest rank. After creating the feature vectors, feature subsets are generated and the subsets with highest accuracy is chosen as a result. The algorithm of two methods is shown in Algorithm 1.

#### **Algorithm 1.**

Cross validation is used for finding hyper parameter about the number of clusters made in Kmeans model.

---

**Input :** unigram or bigram of text datasets

**Output :** feature subset which gain the highest accuracy

Shuffle the datasets

Weight data basis on TF-IDF

**For** fold in numRepetition:

**For** CV in CrossValidation:

Clustering based on documents with Kmeans method when number of clusters=CV

**If** ToCchi model:

**For** cluster in each clusters:

score features basis on Chi-square

---

---

```

        make feature vectors by ordering them according to the score
    End cluster
If ToCknn model:
    For cluster in each clusters:
        select the feature with Nearest Centroid as a center
        (e.g. If 3 clusters, 3 centers are elected.)
        Get the feature vectors in order of degree of nearness to the center by using kNN
    End cluster
    Generate the set of feature subsets incrementally as Equation 7
    For sub in the set of feature subsets:
        Create data with only features in sub from original datasets
        Separate data into test and train set
        Classification: LR, RF
    End sub
    Save feature subset with highest accuracy among the set of feature subsets
End CV
    Save feature subset with highest accuracy among the CV in CrossValidation
End fold
Average the accuracy, precision, recall and f-value on all saved feature subsets.

```

---

## 5. Evaluation

### 5.1. Benchmarks for sentiment analysis

We compare the performance of the two method with the results of base hybrid model and chi model. For base model, feature vectors are generated by the order of scoring based on Chi-square. After creating the feature vectors, same process of producing subsets of feature vector are employed. All the subsets created are evaluated its performance. Also OIFV method in Yousefpour, Ibrahim, 2017 are used to compare the result of our two methods from other hybrid feature selection methods, which is mentioned in section 2 [5].

### 5.2 Performance measure

Usually, four indexes are adopted to evaluate the performance of sentiment classification, so called Accuracy, Precision, Recall and F1-score. Accuracy is the ratio of all true predictions to all predicted samples. Precision is the ratio of true predicted samples against all predicted samples. The ratio of true predicted samples against all actual samples is called Recall. F1 is the average of recall and precision.

### 5.3 Repetition

We repeated 2 times for performing the whole process to generate trustworthy results. For the characteristic of the method, repetition is adopted rather than cross-validation. After the repetition, all evaluating measures are averaged for the results.

### 5.4 Classification used

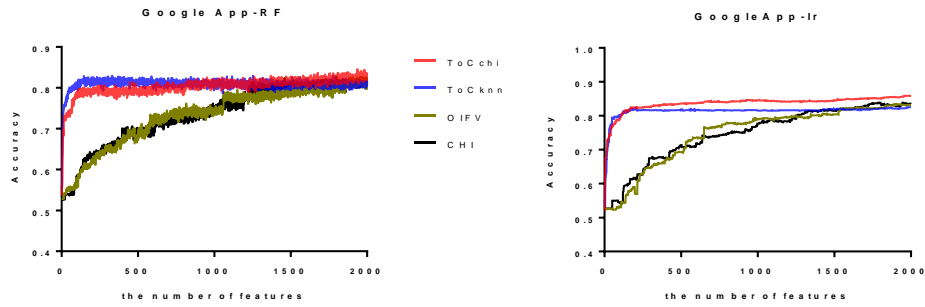
Among the various classifiers for text sentiment analysis, this paper employees LR and RF. Both single classifiers and ensemble classifiers are applied. Because LR is the basic component used in neural network, we believe that it will be helpful to the future study. Random Forest which is performed well is a

general term for ensemble methods using tree-type classifiers  $\{h(\mathbf{x}, \Theta_k), k = 1, \dots, \}$  where the  $\{\Theta_k\}$  are independent identically distributed random vectors and  $\mathbf{x}$  is an input pattern [17]. These tools are from ‘Sci-kit learn’ open source<sup>4</sup>.

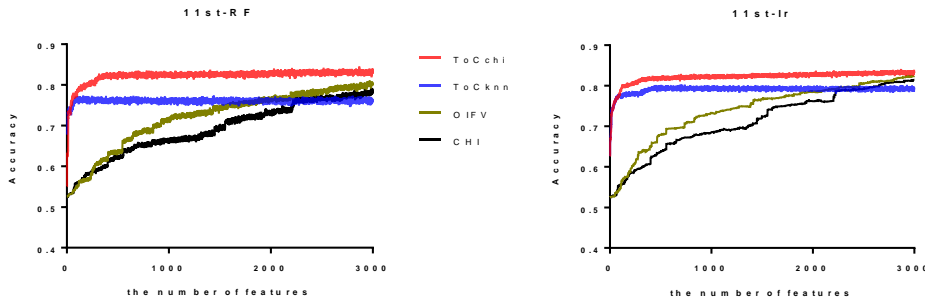
## 5.5 Results

We obtained the performance results of the classifiers by the proposed methods in *Algorithm 1*. We tested the performance of our proposed methods on four Korean review datasets. For ToCchi model, the results of the accuracy and the number of features to reach some level of accuracy are presented on *Table 2*. It also compare the result of the base hybrid feature selection model and OIFV model which is mentioned above in RQ and section 2. The changes of the accuracy as the length of the sub feature vector is increased can be checked on *Fig 2,3,4,5*. The same results of the ToCknn model is presented on *Table 3 and Fig 2,3,4,5*.

**Fig 2. Google App Store (half number of features) - RF(left), LR(right)**

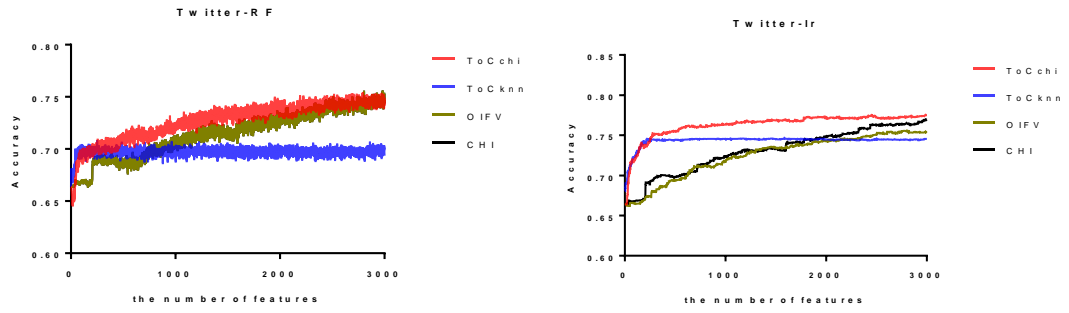


**Fig 3. 11 Street Shopping (half number of features) - RF(left), LR(right)**

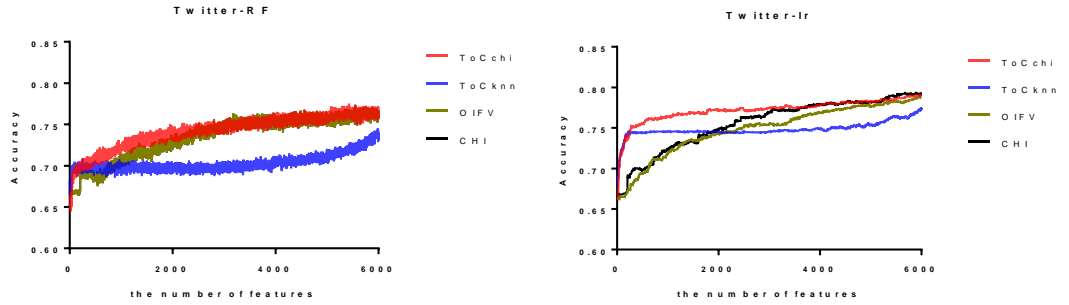


**Fig 4. Twitter (half number of features) - RF(left), LR(right)**

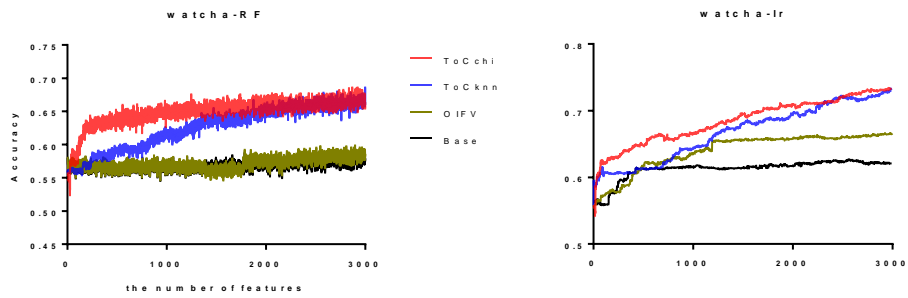
<sup>4</sup> <http://scikit-learn.org/stable/index.html>



**Fig 5. Twitter (Full feature version) - RF(left), LR(right)**



**Fig 6. Watcha Movie (half number of features) - RF(left), LR(right)**



**Table2. Results with RF classifier**

Dataset	Model	# of features for certain level of Accuracy				Performance of subset with best Accuracy				
		start	0.6	0.7	0.8	Accuracy	Features	Precision	Recall	F1
Google App Store	Base	0.52546	140.50	620.00	1418.50	0.765211	6380.00	0.757317	0.689927	0.703633
	OIFV	0.528289	145.50	580.00	1853.00	0.847949	3370.00	0.83239	0.83239	0.83239
	ToCchi	0.531471	7.00	17.00	573.00	0.853253	3332.05	0.843774	0.822167	0.823653

	ToCknn	0.539604	6.50	8.50	337.00	0.842999	3373.05	0.831774	0.822167	0.823653
11street Shopping	Base	0.52546	139.50	621.00	1425.00	0.845827	7534.00	0.856691	0.855392	0.855796
	OIFV	0.525375	283.00	901.50	3017.50	0.856	6373.00	0.843375	0.843375	0.843375
	ToCchi	0.701862	0.00	4.50	119.50	0.851997	7539.05	0.85951	0.85806	0.858497
	ToCknn	0.626563	5.00	15.00	272.00	0.843688	6378.00	0.835942	0.831205	0.83227
Twitter	Base	0.663802	1	3916.50	None	0.774235	6378.05	0.766011	0.686295	0.700435
	OIFV	0.663802	1	3471.50	None	0.771915	6377.05	0.763149	0.763149	0.763149
	ToCchi	0.66363	1	208.00	None	0.77329	6379.00	0.753228	0.679599	0.692597
	ToCknn	0.663802	1	1147.00	None	0.76968	6379.00	0.75612	0.685086	0.698555
Watcha movie	Base	0.580383	4825.50	6267.50	None	0.727184	6348.05	0.68961	0.68953	0.68961
	OIFV	0.557973	4764.00	6238.50	None	0.716142	6352.05	0.68951	0.68951	0.68951
	ToCchi	0.557973	165.50	5887.00	None	0.72215	6348.00	0.697623	0.695346	0.696171
	ToCknn	0.557973	750.40	5108.50	None	0.721825	6350.00	0.706399	0.701482	0.702869

**Table 3. Results with LR classifier**

Dataset	Model	# of features for certain level of Accuracy				Performance of subset with best Accuracy				
		start	0.6	0.7	0.8	Accuracy	Features	Precision	Recall	F1
Google App Store	Base	0.525459689	170.00	478.00	1277.50	0.850426	3369.00	0.850127	0.857250	0.850832
	OIFV	0.526874116	220.00	527.00	1354.00	0.860679	3369.00	0.860679	0.860679	0.860679
	ToCchi	0.523338048	7.50	20.50	117.00	0.873409	3370.50	0.877108	0.871443	0.872506
	ToCknn	0.552333805	10.00	25.00	110.00	0.861386	3370.00	0.866584	0.858087	0.859766
11street Shopping	Base	0.52475	329.50	1373.00	2660.00	0.849000	6374.50	0.839399	0.845029	0.840029
	OIFV	0.525375	226.00	741.50	2469.50	0.823875	6072.50	0.871632	0.871632	0.869873
	ToCchi	0.626056	1.00	17.50	164.50	0.848375	6374.50	0.839279	0.844679	0.839679
	ToCknn	0.626563	1.00	15.00	317.00	0.847625	6372.00	0.835739	0.837939	0.836139
Twitter	Base	0.663630113	1.00	556.50	None	0.794001	6378.00	0.794111	0.794111	0.794111
	OIFV	0.663801994	1.00	224.00	None	0.794431	6377.50	0.794431	0.794431	0.794431
	ToCchi	0.663801994	1.00	44.50	None	0.794775	6378.50	0.792913	0.729731	0.746165
	ToCknn	0.68047439	1.00	32.00	None	0.766243	6378.50	0.756	0.693731	0.707346
Watcha movie	Base	0.557973368	356.00	5353.50	None	0.730107	6128.00	0.730111	0.730097	0.730121
	OIFV	0.562861	395.50	5210.50	None	0.74797	6114.50	0.74811	0.74811	0.74873
	ToCchi	0.555699903	38.00	1715.00	None	0.774602	6034.50	0.775162	0.775186	0.775186
	ToCknn	0.563819422	65.50	2150.50	None	0.788892	6122.00	0.785326	0.783522	0.787723

### 5.5.1. Results for RQ1

First, do ToCchi and ToCknn outperform in feature selection compare to base hybrid method with CHI for filtering or previous hybrid feature selection method(OIFV) mentioned in section 2? As shown in Fig 2,3,6,

both models get the certain level of accuracy with only few number of features compare to the base and OIFV hybrid method in three datasets with both classifiers. For example, in Google Store review datasets, ToCchi reaches the accuracy 0.7 with only 17 features and ToCknn does with 9 features when base model requires 618 features and OIFV needs 578 features with RF. Because both methods make models to be trained with small cost of time and space, they can be useful when the few number of features needed for models need to be trained faster and have limited space. However, for the Twitter datasets, only ToCchi outperform the base and OIFV method and ToCknn shows even lower performance than comparison targets. Considering the characteristic of the datasets, Twitter is unofficial space to express users' feelings easily and quickly. Comparing to the other datasets, it is not for rating and evaluating, but for just instant emotions the moment users write it down. Thus, there are lots of typing error and no word spacing, which make POS tagging difficult and harsh. Also, there are many neologism used only for certain small scope of groups and not widely used. We believe that's the reason the result of the Twitter dataset is different from the others. ToCknn model come out to be not good at unofficial data with large noise.

### 5.5.2. Results for RQ2

Second, which model would show better accuracy or stable performance in feature selection? ToCchi seems to be better than ToCknn due to its stable performance in several kind of datasets. This result is driven from the fact that ToCchi outperform both base and previous hybrid method in all four datasets when ToCknn just do in only three datasets. When datasets seem to have lots of noisy such as type errors and neologisms used only in small scale, ToCchi is better in performance. Because chi-square scoring after clustering the features can consider feature relevance, MECE, and feature redundancy at once, it can outperform even if the datasets have lots of noise like Twitter. On the other hand, the clustering - KNN method can reflect high feature relevance and MECE, but they can't represent the whole cluster because they put center features before other features. Also, it doesn't take account of feature redundancy. When the features are not in a good shape such as containing typing error and unofficial words, it seems like strong relevance between features are rather shows worse performance. As we expected, chi-square after clustering method gives better results than KNN after clustering.

## 6. Discussion

### 6.1. Academic Implications

1) The loan words, they are used mixed with the words in the other language that have the same meaning. For example, '러브' which is the loan word meaning 'love' is used mixed with '사랑' that is the Korean with the same meaning. For the researchers, it will be useful to filter these words before applying to the classifiers. In Figure1, it is shown that the two methods in this paper, ToCchi and ToCknn have functions to distinct and not include all the words with similar meanings. Even if, two methods are not selected as the main feature selection methods, they can adopt them to roughly solve these kinds of loan word problems.

2) For the unofficial resources such as Twitter, feature selection methods that emphasize feature relevance seem to be better rather than feature redundancy. When we compare the performance between two feature selection



methods, ToCchi show the better performance with the datasets having lots of noisy such as type errors and neologisms used only in small scale. Thus, with the unformal datasets or sources, it would be better to use methods considering the relationship between features and target variables like ToCchi.

3) For the clustering, this paper used K-means and it require 2 main hyper-parameters which are the number of clusters and the number of initialization features within each clusters. Because there are more than 1 parameters, grid search is useful for finding proper hyper-parameters. Even if it can pass the optimal point because of grid, this method is good for the researchers with limited time and computing powers.

## 6.2. Practical Implications

1) Even if the domain of the datasets is the same, the characteristic of the data and the words used can be completely different. For example, the reviews of the movie in Twitter include lots of interjection and swear words used for expressing strong emotions. However, the movie reviews in Watcha or Naver that are the Korean movie review webpages, contain more formal and professional words. Thus, when hands on workers want to balance considering the characteristics of sources, it will be useful to apply the ToCchi and ToCknn in this paper. Because they select small number of features with high performance, users can gather features with them from each sources and apply for the mixed datasets with various sources.

2) Because hyper feature selection method which is the complement between filter and wrapper still need lots of time and computing power, this paper used counter-vectorizer to pick out the features roughly before applying the ToCchi and ToCknn. However, if the practical users have enough time and computing powers, it would be better to use those two methods without any processing of picking out.

3) If the agents have enough computing power and time, it would be better to try both grid search and random search for finding the hyper-parameter. Grid search will help the workers to narrow range of the hyper-parameters and random search can complement a weakness of the grid search that it can overlook the area with optimal solutions.

## 7. Concluding Remarks

In this paper, two methods is presented which are considering not only feature-feature and feature-target variable relationship but also trying to include only mutually exclusive and completely exhaustive features for efficiency. From the view of the previous researches, it is meaningful because there are only a few researches about feature selection considering MECE. Eventually, the results with 4 Korean datasets show the effectiveness of the features selected with two methods from the point of view of the number of features needed to reach certain classification accuracy. From RQ1, Both models get the certain level of accuracy with only few number of features compare to the base and OIFV hybrid method in three datasets with both classifiers. From RQ2, ToCchi seems to be better than ToCknn due to its stable performance even in unofficial and noisy datasets.

However, this article also has the limitations. Because of the cost of time and computing power, counter vectorizer is used to pick features out roughly first before using our two methods. Thus the full features within each datasets are about 6000 and because of that, the effectiveness of the two methods are not shown effectively. Also the classifiers used in this article are only LR and RF, which are the basic of the linear classifier and ensemble classifier. That's also because of the time and memory cost, and it would be better to show the results using other classifiers such as SVM and AB(AdaBoost). For active research in this field, progressing the experiment with known English data will be useful because it is easy to compare with other methods. Therefore, for the further research, experiments with well-known English data without any counter-vectorizer process using various classifiers should be proceeded. Also, it would be better to use various BOW such as POS and n-gram because this paper used only unigram for reducing cost of time and memory.

## References

- [1] Hu, Nan, Noi Sian Koh, and Srinivas K. Reddy. "Ratings lead you to the product, reviews help you clinch it? The mediating role of online review sentiments on product sales." *Decision support systems* 57 (2014): 42-53.
- [2] Cheung, Christy MK, Bo Sophia Xiao, and Ivy LB Liu. "Do actions speak louder than voices? The signaling role of social information cues in influencing consumer purchase decisions." *Decision Support Systems* 65 (2014): 50-58.
- [3] Zheng, Ling, Ren Diao, and Qiang Shen. "Self-adjusting harmony search-based feature selection." *Soft Computing* 19.6 (2015): 1567-1579.
- [4] Liu, Luying, et al. "A comparative study on unsupervised feature selection methods for text clustering." *Natural Language Processing and Knowledge Engineering, 2005. IEEE NLP-KE'05. Proceedings of 2005 IEEE International Conference on.* IEEE, 2005.
- [5] Yousefpour, Alireza, Roliana Ibrahim, and Haza Nuzly Abdel Hamed. "Ordinal-based and frequency-based integration of feature selection methods for sentiment analysis." *Expert Systems with Applications* 75 (2017): 80-93.
- [6] Aggarwal, U., and G. Aggarwal. "Sentiment Analysis: A Survey." *International Journal of Computer Sciences and Engineering* 5.5 (2017): 222-225.
- [7] Nassirtoussi, Arman Khadjeh, et al. "Text mining of news-headlines for FOREX market prediction: A Multi-layer Dimension Reduction Algorithm with semantics and sentiment." *Expert Systems with Applications* 42.1 (2015): 306-324.
- [8] Saeys, Yvan, Iñaki Inza, and Pedro Larrañaga. "A review of feature selection techniques in bioinformatics." *bioinformatics* 23.19 (2007): 2507-2517.
- [9] Dash, Manoranjan, and Huan Liu. "Feature selection for classification." *Intelligent data analysis* 1.1-4 (1997): 131-156.

- [10] Lin, Kuan-Cheng, et al. "Feature selection based on an improved cat swarm optimization algorithm for big data classification." *The Journal of Supercomputing* 72.8 (2016): 3210-3221.
- [11] Wang, Suge, et al. "A hybrid method of feature selection for Chinese text sentiment classification." *Fuzzy Systems and Knowledge Discovery, 2007. FSKD 2007. Fourth International Conference on*. Vol. 3. IEEE, 2007.
- [12] Vinh, Nguyen X., and James Bailey. "Comments on supervised feature selection by clustering using conditional mutual information-based distances." *Pattern Recognition* 46.4 (2013): 1220-1225.
- [13] Assi, E. Bou, et al. "A hybrid mRMR-genetic based selection method for the prediction of epileptic seizures." *Biomedical Circuits and Systems Conference (BioCAS), 2015 IEEE*. IEEE, 2015.
- [14] Claypo, Niphat, and Saichon Jaiyen. "Opinion mining for Thai restaurant reviews using neural networks and mRMR feature selection." *Computer Science and Engineering Conference (ICSEC), 2014 International*. IEEE, 2014.
- [15] Pang, Bo, and Lillian Lee. "Opinion mining and sentiment analysis." *Foundations and Trends® in Information Retrieval* 2.1–2 (2008): 1-135.
- [16] Pang, Bo, Lillian Lee, and Shivakumar Vaithyanathan. "Thumbs up?: sentiment classification using machine learning techniques." *Proceedings of the ACL-02 conference on Empirical methods in natural language processing-Volume 10*. Association for Computational Linguistics, 2002.
- [17] Breiman, Leo. "Random forests." *Machine learning* 45.1 (2001): 5-32.
- [18] Sebastiani, Fabrizio. "Machine learning in automated text categorization." *ACM computing surveys (CSUR)* 34.1 (2002): 1-47.
- [19] MacQueen, James. "Some methods for classification and analysis of multivariate observations." *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*. Vol. 1. No. 14. 1967.
- [20] Sotoca, José Martínez, and Filiberto Pla. "Supervised feature selection by clustering using conditional mutual information-based distances." *Pattern Recognition* 43.6 (2010): 2068-2081.
- [21] Nam, Le Nguyen Hoai, and Ho Bao Quoc. "A Combined Approach for Filter Feature Selection in Document Classification." *Proceedings of the 2015 IEEE 27th International Conference on Tools with Artificial Intelligence (ICTAI)*. IEEE Computer Society, 2015.
- [22] Lee, C., Choi, D., Kim, S., Kang, J.: Classification and analysis of emotion in korean microblog texts. *KIISE* 40(3), 159–167 (2013)
- [23] Jung, Younghee, et al. "A corpus-based approach to classifying emotions using Korean linguistic features." *Cluster Computing* 20.1 (2017): 583-595.
-