# Classification Model of Consumer Question about Motorbike Problems by Using Naïve Bayes and Support Vector Machine

Ekky Wicaksana, Danang Triantoro Murdiansyah, Isman Kurniawan

*School of Computing, Telkom University*
*Bandung, Indonesia*

ekkywicaks@students.telkomuniversity.ac.id
danangtri@telkomuniversity.ac.id
ismankrn@telkomuniversity.ac.id

## Abstract

The motorbike plays an important role in supporting daily activity. The motorbike is known as one of the transportation modes that is frequently used in Indonesia. The number of motorbikes used in Indonesia is continuously increasing time by time. Hence, the occurrence of motorbike problems can affect community activity and disturb the economic condition in society. Since the motorbike problem can occur anytime, a prevention action is required by providing an online consultation platform. However, a classification model is required to handle a wide range of questions about the motorbike problem. By classifying those questions into a specific class of problems, the solution can be delivered to the consumer faster. In this study, we developed prediction models to classify consumer questions. The data set was collected from consumer questions regarding motorbike problems that are commonly occurring. The model was developed using two machine learning algorithms, i.e., Naïve Bayes and Support Vector Machine (SVM). Text vectorization was performed by using the n-gram and term frequency-inverse document frequency (TF-IDF) method. The results show that the SVM model with the uni-trigram model performs better with the value of accuracy and F-measure, which are 0.910 and 0.910, respectively.

**Keywords:** classification, naïve bayes, SVM, n-gram, TF-IDF

## Abstrak

Sepeda motor berperan penting dalam menunjang aktivitas masyarakat sehari-hari. Sepeda motor dikenal sebagai salah satu moda transportasi yang banyak digunakan di Indonesia. Penggunaan sepeda motor di Indonesia terus meningkat dari waktu ke waktu. Sehingga, timbulnya permasalahan sepeda motor dapat mempengaruhi aktivitas dan mengganggu kondisi perekonomian masyarakat. Karena masalah sepeda motor bisa terjadi kapan saja, maka diperlukan tindakan pencegahan dengan menyediakan platform konsultasi online. Namun, diperlukan model klasifikasi untuk menangani berbagai pertanyaan tentang masalah sepeda motor. Dengan mengklasifikasikan pertanyaan-pertanyaan tersebut ke dalam kelas masalah tertentu, solusi dapat disampaikan kepada konsumen lebih cepat. Dalam studi ini, kami mengembangkan model prediksi untuk mengklasifikasikan pertanyaan konsumen. Kumpulan data dikumpulkan dari pertanyaan konsumen mengenai masalah sepeda motor yang sering terjadi. Model dikembangkan menggunakan dua algoritma pembelajaran mesin, yaitu Naïve Bayes dan Support Vector Machine (SVM). Vektorisasi teks dilakukan dengan metode n-gram dan term frequency-inverse document frequency (TF-IDF). Hasil penelitian menunjukkan bahwa model SVM dengan model uni-trigram memiliki kinerja yang lebih baik dengan nilai akurasi dan F-measure masing-masing sebesar 0.910 dan 0.910.

**Kata Kunci:** klasifikasi, naïve bayes, SVM, n-gram, TF-IDF

## I. Introduction

TRANSPORTATION plays an important role in assisting the economic growth of society and the country's economic development. Hence, the success of economic development can be reached with the support of a good transportation system. The transportation system can facilitate the mobility of society and goods, improve the economic value of an area, and maintain price stability. The establishment of a startup company in transportation service points out the need of society on transportation facilities. Among the available transportation mode, the motorbike is one option chosen by most people in society, especially in Indonesia. Many people use motorbikes to assist varied activities, such as working, sending goods, or mobility. In parallel with establishing a transportation startup company, the motorbike service is the most favorite one used by consumers. The main reason is related to the cheap in price compare to the car service. Also, the utilization of motorbikes can improve the efficiency in time when someone should go anyplace.

Regarding the importance of motorbike utilization, the motorbike condition becomes crucial to support the activity of society. The high intensity of motorbike utilization will lead to the decreasing of the motorbike condition. The bad condition of the motorbike can cause an accident and endanger not only the driver but also other people. According to a report published by the Ministry of Communication and Information, three people died every 3 hours caused by a traffic accident, in which 9% of the accident is caused by vehicle factors [1]. The report also tells that 144 motorbike traffic accidents occur from January until February 2020 [2]. The continuously bad condition will make the motorbike become broken. Hence, this condition should be prevented by fix the minor problem as fast as possible.

One alternative solution to overcome this problem is providing online service to give instant solutions regarding the motorbike problem. Consumers can use the online service to ask about their motorbike problems, and the system will answer their questions. This solution can be developed to complement the motorbike repair shop, which is sometimes limited by distance and time. However, to boost the system's performance, the question from a consumer should be classified to a specific class of motorbike problem. In other words, the system should recognize whether the question is related to a mechanical problem or another problem. This case can be solved by implementing a text classification method to classify the consumer question.

In the term of text classification, many researchers have exploited machine learning algorithm to classify text data. In 2018, Baygin performed a text document classification task using Naïve Bayes (NB) algorithm and obtained a 92% accuracy value [3]. In 2018, Venkatesh and co-workers compared the performance of NB and k-Nearest Neighbor (kNN) algorithms in the case of text classification [4]. They found that the NB algorithm performed better than kNN. In 2018, Bužić and co-workers classified lyric text by using Naïve Bayes algorithms and obtaining satisfying results with the value of precision and F-measure at 93% and 94%, respectively [5]. In 2019, Rahman and co-workers performed topic classification using Decision Tree (DT), kNN, and NB [6]. The results show that the NB algorithm giving the best performance with accuracy is 91.8%. However, despite the success of several text classification studies, there is no research performed concerning the implementation of text classification to predict motorbike problems.

Hence, this research aims to develop a prediction model to classify customer questions regarding the motorbike problem. We defined the class of motorbike problems as engine and non-engine problems. The model was developed by using NB and Support Vector Machine (SVM) methods. Both methods are commonly used in text classification tasks with satisfying results. We utilize the data set consisting of a set of questions in Bahasa Indonesia collected from social media and our questionnaire regarding common problems in the motorbike.

## II. Methodology

This study was conducted by collecting data set regarding consumer questions about the motorbike problem. Then, the data set was pre-processed to obtain clean text and improve the text feature. The process was followed by text vectorization to transform the text representation into a matrix. A prediction model was developed using the NB and SVM method, and their performance was evaluated by calculating several validation parameters.

## A. Dataset

The data set used in this study consists of 505 consumer questions regarding motorbike problems collecting from the QnA forum in social media and own questionnaire. The labeling process was performed by a mechanic that is an expert in the motorbike problem. We classify the consumer question into two motorbike problem classes, i.e., engine and non-engine problems. The number of questions classified as engine and non-engine problems is 265 and 240, respectively. The data set is randomly split into training and test set with a ratio of 4:1.

## B. Text Pre-Processing

Text pre-processing was performed to obtain the clean and appropriate text for the model and reduce the complexity by removing unnecessary features [7]. The text pre-processing stage consists of (i) lowercase, (ii) punctuation removal, (iii) tokenization, (iv) stopword removal, and (v) stemming. Firstly, we converted the text into lowercase to remove letter case bias. Then, we removed symbols and numbers that are appeared in the text to obtain the main message from the text. Afterward, the text is tokenized so that every word is defined as a single entity. Then, we performed stopword removal to remove the word that is commonly found in any text. This stage is conducted to obtain the word that is highly correlated to our case. Finally, we performed a stemming process for each word to obtain the root form of each word. This step is conducted to merge several words that have similar root forms. The comparison of text before and after text pre-processing is presented in Table 1.

TABLE I
THE COMPARISON OF TEXT BEFORE AND AFTER TEXT PRE-PROCESSING

| No. | Before | After |
|---|---|---|
| 1 | *Saya Mau tanya speda saya Suzuki shogun 125r.* | *speda suzuki shogun pengapian dc kal* |
| 2 | *Saya kan punya beat fi 2018. Yang mau saya tanya* | *beat fi gas putar suara kemritik gitu dengerin* |
| 3 | *Apa memang benar bunyi injectornya sprti itu?* | *bunyi injcetornya sprti kedengeran kasar banget* |
| 4 | *Honda spacy fi saya oli gardannya suka netes* | *honda spacy fi oli garden suka netes cek bengkel* |
| 5 | *Sy pynya motor kaisar Ruby. Jika hidup di tempat* | *motor kaisar ruby hidup bagus jalan* |

## C. Feature Extraction

Feature extraction is a process to obtain the most relevant information from the data set to form an appropriate vector for both training and test set. The feature extraction process was performed by using a series step of n-gram and TF-IDF methods.

*1) N-gram:* The n-gram model is a statistical technique used to examine the next word from the word sequence and predict subsequent words. In sentiment analysis, the n-gram model analyzes the text of document sentiment [8-10]. The output from the n-gram model will be used in the text weighting process. This study utilizes six n-gram models, consisting of unigram, bigram, trigram, uni-bigram, uni-trigram, and bi-trigram. We also evaluate the effect of the N value on the performance of the prediction model.

*2) TF-IDF:* The classification task using a machine learning algorithm requires input data in numerical matrix form. Hence, we performed text vectorization to transform text documents into a numerical feature matrix with lower dimension space [11]. One method that is commonly used in text vectorization is Term Frequency-Inverse Document Frequency (TF-IDF). TF-IDF is one of the word weighting methods that is introduced in information retrieval cases. This method defines a weight that reflects the importance of a word in a corpus. The increase of a word's weight is in parallel with the increase of the word frequency in one class sample and the decrease of

word frequency in another sample. This mechanism improves the input data by decreasing the common word's importance and increasing the rare word [12-13]. The weight calculation in TF-IDF is formulated in (1) – (3).

$$tf_{t,d} = f_{t,d} \tag{1}$$

$$idf_t = log\frac{N}{df_t} \tag{2}$$

$$w = tf_{t,d} \times idf_t \tag{3}$$

The $f_{t,d}$ variable represents the term frequency in the document.. Meanwhile, $N$ and $df_t$ represent the number of documents and document frequency, respectively.

*D. Naïve Bayes*

The Naïve Bayes method is an algorithm that is developed according to the Bayesian theorem founded by Thomas Bayes in the 18th century. In the Bayesian theorem, a conditional probability is expressed as:

$$P(H|X = 1) = \frac{(P(X|H = 1)P(H))}{P(X)} \tag{4}$$

where X and H are events. *P(H/X=1)* and *P(X/H=1)* are the probability of H given X is true, and the probability of X given H is true, respectively. *P(H)* and *P(X)* are the independent probabilities of H and X, respectively. Suppose that *m* class is defined for each input X, Naïve Bayes classifier will classify an input as $i^{th}$ class if and only if $P(Ci|X) > P(Cj|X)$ with $1 \leq j \leq m$, and $j \neq i$. In other word, Naïve Bayes works by looking for the maximum value of $P(C_i/X)$ that is estimated as:

$$P(C_i|X) = \frac{P(X|C_i)P(C_i)}{P(X)} . \tag{5}$$

Since the value of $P(X)$ is similar for each class probability, the calculation of the probability *P(C_i/X)* is only affected by *P(X/Ci))P(Ci)*. If there are many attributes involved in the model, the complexity in calculating *P(C_i/X)* can be reduced by using naïve assumption. By using this approximation, it is assumed that is no dependency amongst those attributes. Hence, the calculation of *P(C_i/X)* can be expressed as:

$$P(Ci|X) = \prod_{k=1}^{n} P(x_k|C_i) = P(x_1|C_i) \times P(x_2|C_i) \times \ldots \times P(x_n|C_i). \tag{6}$$

The $P(x_k|C_i)$ is calculated according to the type of the attribute. As for categorical attribute, the value of $P(x_k|C_i)$ is defined as a ratio of the number of $C_i$ class with $x_k$ attribute to all of $C_i$ class. As for the continued type of attribute, it is assumed that the data distribution follows Gaussian distribution. Hence, the value of $P(x_k|C_i)$ can be calculated as:

$$P(X_k|C_i) = \frac{1}{\sigma_{C_i}\sqrt{2\pi}} e^{\frac{(x-\mu_{C_i})}{2\sigma_{C_i}^2}}, \tag{7}$$

where $\mu_{C_i}$ and $\sigma_{C_i}$ are the value of average and standard deviation, respectively, for data in $C_i$ class with $x_k$ attribute. The implementation of Naïve Bayes in text classification is preferable since this method is considered a simple method. Hence, the utilization of Naïve Bayes in the text classification will give advantages due to its effectiveness and computational simplicity.

*E. Support Vector Machine*

Support vector machine (SVM) is a supervised learning algorithm commonly used for classification and regression tasks [14]. SVM works by finding the best hyperplane to maximize the distance between two classes, as illustrated in Figure 1. The object data closest to the hyperplane is defined as a support vector, which is the only object considered in finding the optimal hyperplane. In the case of 2 dimensions and 3 dimension attributes, the separating function is defined as line and plane. As for the higher dimension, the function is defined as hyperplane [15].
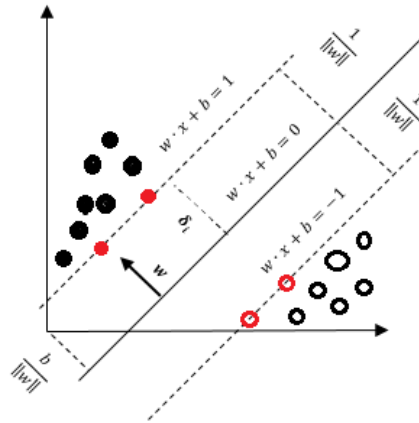


Fig. 1. Hyperplane of Support Vector Machine

As for an example, suppose the data set consists of *x* attribute and *y* label with $y \in \{-1,1\}$. For two-dimension case, the hyperplane formula can be expressed as:

$$w.x + b, \tag{8}$$

where *w* and *b* represent vector weight and bias, respectively. After the optimal hyperplane is obtained, a sample will be classified as -1 class if $w.x + b \leq -1$, and will be classified as +1 class if $w.x + b \geq +1$. Meanwhile, the distance between two margins is formulated as $\frac{2}{\|w\|}$. In this study, we use a linear kernel with the value of the regularization parameter as 1.0. SVM is one of the machine learning techniques that is commonly used for text classification due to its ability to handle multidimensional features. The total unique words found in all text are considered as features, and as consequence the prediction will involve high number of features.

*F. Validation Parameter*

Several validation parameters were calculated to evaluate the performance of the classification model. Those parameters consist of accuracy, recall, precision, and F-1 score, in which the calculation is derived from the confusion matrix [16]. The formulation of those parameters is expressed in (9) – (12).

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}, \tag{9}$$

$$Precision = \frac{TP}{TP+FP}, \tag{10}$$

$$Recall = \frac{TP}{TP+FN} \times 100\%, \tag{11}$$

$$F - 1\ score = 2 \times \frac{precision \times recall}{precision + recall}, \tag{12}$$

where TP, TN, FP and FN mean true positive, true negative, false positive and false negative, respectively. Here, we consider the F-1 score as an overall measurement to evaluate the model performance.
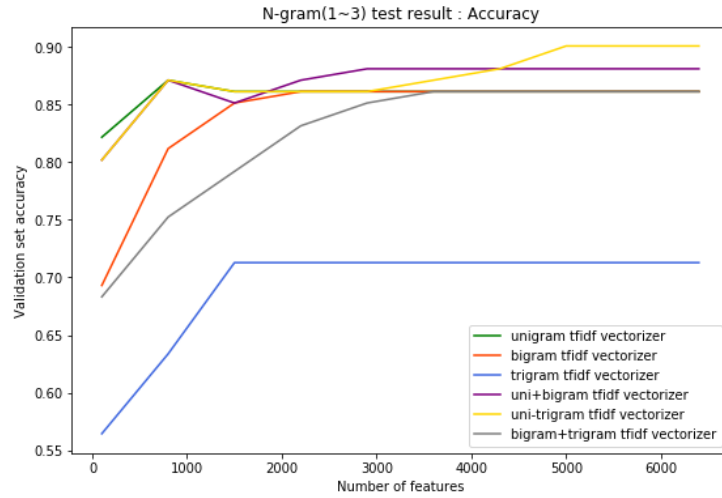
## III. RESULTS AND DISCUSSION

### A. Optimal N-gram and Feature Number Analysis

We examined the contribution of n-gram and feature number variation on model performance. Six n-gram models were utilized in this study, in which the number of features for each n-gram model is presented in Table 2. We found that the unigram model contains the least number of features, while uni-trigram contains the most feature. Here, the number of features reflects the complexity of the model. This means that more features lead to a risk of overfitting conditions. Hence, we analyzed the contribution of the feature number of each n-gram model on the model performance to obtain the optimal number of features, as shown in Figure 2.
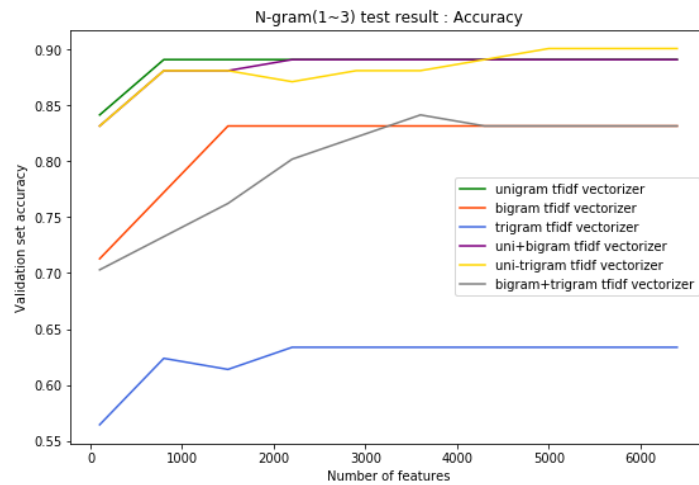
TABLE II
FEATURE NUMBER FOR EACH N-GRAM MODEL

| N-gram Model | Feature Number |
|--------------|----------------|
| Unigram | 910 |
| Bigram | 2147 |
| Trigram | 2188 |
| Uni-bigram | 3057 |
| Uni-trigram | 5245 |
| Bi-trigram | 4335 |

As for the NB method, we found that the uni-trigram model gives the best accuracy after using all available features. This indicates that utilization of all features in the uni-trigram model improves the ability of the model to classify the proper class. Meanwhile, the worst result is obtained from the trigram model with the lowest value of accuracy. This is caused not only by the feature number but also by the feature quality derived from the model. Also, we found that the increase of feature number improves the model performance for all n-gram models. This points out the dependency of model performance on feature numbers. The flat horizontal line in the figure indicates that all features of the n-gram model have been used.

(a)



(b)

Fig. 2. The N-gram Model and Feature Number Analysis for (a) NB and (b) SVM Method

As for the SVM model, we found that the uni-trigram model also performed better than other n-gram models. This indicates that the feature quality of the uni-trigram model is suitable to be used by both NB and SVM models. Meanwhile, the worst performance was also obtained from the trigram model. This confirms that the trigram model is not suitable for obtaining the feature for this classification task. We also found that the increase of feature number improves the model performance, except for the bi-trigram model.

## B. *Model Validation*

The prediction models for the classification task were developed by using the optimized feature number for each n-gram model obtained from the previous step. We calculated several validation parameters to evaluate those models. The result of the NB method evaluation is presented in Table 3. As for the training set, we found that the uni-trigram model shows the best performance with the value of accuracy, precision, recall, and F-1 score are 0.985, 0.995, 0.976, and 0.985, respectively. The trigram model gives the worst result with the lowest values of accuracy and F-1 score.

As for the test set, we also found that the uni-trigram model performs better than other models with the value of accuracy, precision, recall, and F-1 score are 0.900, 0.890, 0.924, and 0.900, respectively. Like the training set, the trigram model also gives the worst performance in the test set with the lowest accuracy and F-1 score. This result indicates that the quality of features obtained from the uni-trigram model is well enough to identify the characteristic of each class by using the NB method. Meanwhile, the features obtained from the trigram model are not good in quality to represent the character of each class.

The evaluation result of the SVM method is presented in Table 4. As for the training set, we found that the uni-bigram and uni-trigram models perform better than other models with a similar value of accuracy, precision, recall, and F-1 score as 0.990, 0.990, 0.990, and 0.990, respectively. This indicates that the addition of bigram and trigram features to the unigram model gives a similar contribution in improving the model. Also, we found that the bigram model gives the worst performance with the lowest value of F-1 score. As for the test set, we found that the uni-trigram model performs better than the other model with the value of accuracy, precision, recall, and F-1 score, which are 0.910, 0.912, and 0.910 and 0.910, respectively. This indicates that the feature obtained from the uni-trigram model is also suitable to be used in the SVM method. The trigram model is found to gives the worst performance with the lowest value of F-1 score. Also, this confirms that the quality of features obtained from the trigram model is not enough to recognize the characteristic of each class.

Finally, we compared the performance of the NB and SVM method for the classification task. As for the NB method, the best model obtained using uni-trigram features with the value of F-1 score is 0.985 and 0.900 for training and test set, respectively. As for the SVM method, the best model was also obtained using uni-trigram features with the value of F-1 score is 0.990 and 0.910 for training and test set, respectively. From the comparison, we found that the uni-trigram model developed by using the SVM method performed better than that developed by the NB method. This indicates the better ability of SVM in recognizing the character of each class. This ability is assisting by the complexity level of SVM that is higher than NB.

TABLE III
THE VALIDATION RESULTS OF THE NB METHOD

| Gram Methods | TP | FP | FN | TN | A | P | R | F |
|---|---|---|---|---|---|---|---|---|
| **Training Set** | | | | | | | | |
| Unigram | 187 | 5 | **3** | **209** | 0.980 | 0.976 | **0.985** | 0.981 |
| Bigram | 187 | 5 | 4 | 208 | 0.977 | 0.976 | 0.981 | 0.978 |
| Trigram | 163 | 29 | 4 | 208 | 0.918 | 0.877 | 0.981 | 0.926 |
| Uni-bigram | 190 | 2 | 5 | 207 | 0.982 | 0.990 | 0.976 | 0.983 |
| Uni-trigram | **191** | **1** | 5 | 207 | **0.985** | **0.995** | 0.976 | **0.985** |
| Bi-trigram | 187 | 5 | 4 | 208 | 0.977 | 0.976 | 0.981 | 0.978 |
| **Test Set** | | | | | | | | |
| Unigram | 39 | 9 | 5 | 48 | 0.861 | 0.842 | 0.905 | 0.872 |
| Bigram | 36 | 12 | 2 | 51 | 0.861 | 0.809 | 0.962 | 0.879 |
| Trigram | 20 | 28 | **1** | **52** | 0.712 | 0.650 | 0.981 | 0.781 |
| Uni-bigram | 40 | 8 | 4 | 49 | 0.881 | **0.859** | 0.924 | 0.890 |
| Uni-trigram | **42** | **6** | 4 | 49 | **0.900** | 0.890 | 0.924 | **0.900** |
| Bi-trigram | 36 | 12 | 2 | 51 | 0.861 | 0.809 | **0.962** | 0.879 |

A: Accuracy, P: Precision, R: Recall, F: F-1 Score.

TABLE IV
THE VALIDATION RESULTS OF THE SVM METHOD

| Gram Methods | TP | FP | FN | TN | A | P | R | F |
|---|---|---|---|---|---|---|---|---|
| **Training Set** | | | | | | | | |
| Unigram | 188 | 4 | 3 | 209 | 0.982 | 0.982 | 0.982 | 0.982 |
| Bigram | 185 | 7 | 1 | 211 | 0.980 | 0.980 | 0.980 | 0.980 |
| Trigram | 160 | 32 | **0** | **212** | 0.920 | 0.931 | 0.920 | 0.919 |
| Uni-bigram | 191 | 1 | 3 | 209 | **0.990** | **0.990** | **0.990** | **0.990** |
| Uni-trigram | **191** | **1** | 3 | 209 | **0.990** | **0.990** | **0.990** | **0.990** |
| Bi-trigram | 188 | 4 | 3 | 209 | 0.982 | 0.982 | 0.982 | 0.982 |
| **Test Set** | | | | | | | | |

| Unigram | 40 | 8 | 4 | 49 | 0.881 | 0.883 | 0.881 | 0.880 |
|---|---|---|---|---|---|---|---|---|
| Bigram | 33 | 15 | **1** | **52** | 0.841 | 0.868 | 0.841 | 0.837 |
| Trigram | 19 | 29 | **1** | **52** | 0.702 | 0.788 | 0.702 | 0.672 |
| Uni-bigram | 41 | 7 | 4 | 49 | 0.891 | 0.892 | 0.891 | 0.890 |
| Uni-trigram | **42** | **6** | 3 | 50 | **0.910** | **0.912** | **0.910** | **0.910** |
| Bi-trigram | 40 | 8 | 4 | 49 | 0.881 | 0.883 | 0.881 | 0.880 |

A: Accuracy, P: Precision, R: Recall, F: F-1 Score.

## IV. Conclusion

We developed prediction models to classify customer questions regarding motorbike problems using Naïve Bayes (NB) and Support Vector Machine (SVM) method. The feature extraction was performed by using the n-gram model and TF-IDF method. We utilized six n-gram models and evaluated the contribution of the feature number of each n-gram model to the model performance. We found that uni-trigram produce the best feature for both NB and SVM method. This is indicated by the highest F-1 score that is achieved by using this n-gram model. From the comparison, the SVM model is found to perform better than the NB model with the accuracy and F-1 score are 0.990 and 0.910, respectively.

## References

[1]　KOMINFO, "Setiap Jam Rata-rata 3 Orang Meninggal Akibat Kecelakaan Jalan di Indonesia," 2017. [Online]. Available: https://kominfo.go.id/index.php/content/detail/10368/rata-rata-tiga-orang-meninggal-setiap-jam-akibat-kecelakaanjalan/0/artikel_gpr.

[2]　Lokadata, "Kecelakaan Lalu Lintas Menurut Jenis Kendaraan," 2020. [Online]. Available: https://lokadata.id/data/kecelakaan-lalu-lintas-menurut-jenis-kendaraan-2020-1582708742.

[3]　M. Baygin, "Classification of text documents based on naive bayes using N-gram features," International Conference on Artificial Intelligence and Data Processing (IDAP), Malatya, Turkey, 2018, pp. 1-5.

[4]　Venkatesh and K. V. Ranjitha, "Classification and optimization scheme for text data using machine learning naïve bayes classifier," IEEE World Symposium on Communication Engineering (WSCE), Singapore, 2018, pp. 33-36.

[5]　D. Bužić and J. Dobša, "Lyrics classification using naive bayes," 2018 41st International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO), Opatija, 2018, pp. 1011-1015.

[6]　M. A. Rahman and Y. A. Akter, "Topic classification from text using decision tree, K-NN and multinomial naïve bayes," 1st International Conference on Advances in Science, Engineering and Robotics Technology (ICASERT), Dhaka, Bangladesh, 2019, pp. 1-4.

[7]　G. Singh, B. Kumar, L.Gaur, and A.Tyagi, "Comparison between multinomial and bernoulli naïve bayes for text classification," International Conference on Automation, Computational and Technology Management (ICACTM), India, 2019.

[8]　A. Nugroho, "Analisis sentimen pada media sosial twitter menggunakan naive bayes classifier dengan ekstrasi fitur N-gram," J-SAKTI, vol. 2, no. 2, 2018, p. 200.

[9]　F. Peng and D. Schuurmans, "Combining naive bayes and N-gram language models for text classification," Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics), vol. 2633, 2003, pp. 335–350.

[10]　L. Kobyliński and A. Przepiórkowski, "Definition extraction with balanced random forests," Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics), vol. 5221, 2008, pp. 237–24.

[11]　M. Hakiem and M. A. Fauzi, "Klasifikasi ujaran kebencian pada twitter menggunakan metode naïve bayes berbasis N-gram dengan seleksi fitur information gain," vol. 3, no. 3, 2019, pp. 2443–2451.

[12]　A. Tripathy, A. Agrawal, and S. K. Rath, "Classification of sentiment reviews using n-gram machine learning approach," Expert Syst. Appl., vol. 57, 2016, pp. 117–126.

[13]　M. Shirakawa, T. Hara, and S. Nishio, "N-gram IDF: A global term weighting scheme based on information distance," WWW 2015 - Proc. 24th Int. Conf. World Wide Web, 2015, pp. 960–970.

[14]　Suyanto, "Data Mining: Untuk klasifikasi dan klasterisasi data," Informatika, 2017, pp. 196-210.

[15]　P. A. Octaviani, Y. Wilandari, and D. Ispriyanti, "Penerapan metode klasifikasi support vector machine pada data akreditasi sekolah dasar di kabupaten magelang," Jurnal Gaussian, vol. 3, no. 4, 2014, pp. 811-820.

[16]　X. Zhou and A. Del Valle, "Range based confusion matrix for imbalanced time series classification," 6th Conference on Data Science and Machine Learning Applications (CDMA), Riyadh, Saudi Arabia, 2020, pp. 1-6.