# SENTIMENT ANALYSIS OF REKSADANA ON BIBIT APPLICATIONS USING THE NAÏVE BAYES METHOD AND K-NEAREST NEIGHBOR (KNN)

**Alisa Fitriyani**[1*] , **Agung Triayudi**[2]

Sistem Informasi, Fakultas Teknologi Komunikasi dan Informatika
Univesitas Nasional
alisafitriyani2018@student.unas.ac.id[1*], agungtriayudi@civitas.unas.ac.id[2]
(*) Corresponding Author

## *Abstrak*

*Kurang Nya minat masyarakat terhadap pasar modal, membuat para petinggi perusahan pasar modal saling berlomba untuk menyediakan layanan agar dapat memberikan kemudahan bagi nasabah di berbagai layanan yang tersedia serta memberikan kemudahan dalam mengakses informasi keuangan. Kemunculan beberapa perusahaan rintisan yang menyediakan produk investasi reksadana bagi investor yaitu PT Bibit Reksadana Tumbuh Bersama, yang menciptakan aplikasi reksadana yaitu Bibit Reksadana dengan pengguna lebih dari satu juta pengguna berdasarkan data unduhan pada play store oleh PT Bibit Tumbuh Bersama yang berlaku sebagai Agen Penjual Efek Reksadana (APERD) dan menjual l34 produk reksadana. Sehingga untuk memberikan informasi kepada masyarakat diperlukan adanya analisis sentimen mengenai bagaimana pendapat dari para pengguna aplikasi bibit reksadana menggunakan metode K-nearest neighbor(KNN) dan Naïve bayes, dengan hasil scraping youtube  sebanyak 33.292 dan scraping ulasan sebanyak 30.708 ulasan, lalu dilakukan tahap text processing, dan labeling menggunakan library textblob, dengan tingkat akurasi pada klasifikasi data youtube dan data ulasan dengan metode K-Nearest Neighbor sebanyak 99%, 99%,  dan Naïve bayes sebanyak 99%,98%,    dapat disimpulkan bahwa sentimen netral lebih banyak dari sentimen positif, dan sentiment positif lebih banyak dari sentiment negative.*

*Kata kunci: Bibit Reksadana, YouTube, Naive Bayes classifier(NBC), K-nearest neighbor(KNN)*

## Abstract

The lack of public interest in the capital market has made the top brass of capital market companies compete with each other to provide services to provide convenience for customers in the various services available and provide convenience in accessing financial information. The emergence of several startup companies that provide reksadana investment products for investors, namely PT Bibit Reksadana Grows Together, which created a reksadana application, namely Bibit Reksadana with more than one million users based on data downloaded on the play store by PT Bibit Grow Bersama which acts as a Reksadana Selling Agent (APERD) and sells 134 reksadana products. So to provide information to the public, it is necessary to have a sentiment analysis on how the opinions of users of the reksadana bibit application use the K-nearest neighbor (KNN) and Naïve Bayes methods, with the results of scraping youtube as much as 33,292 and scraping reviews as much as 30,708 reviews, then the text processing stage is carried out, and labeling using the text blob library, with an accuracy rate of 99%, 99%, 99% accuracy on youtube data classification and review data with the K-Nearest Neighbor method, 99%, and 99%, 98% Naïve Bayes, it can be concluded that neutral sentiment is more than positive sentiment. , and more positive sentiment than negative sentiment.

Keywords: Reksadana, YouTube, Naive Bayes classifier (NBC), K-nearest neighbor(KNN)

## INTRODUCTION

Investment can increase future welfare, which prevents inflation that often occurs. Investment can be carried out in several ways, namely, investing in financial assets and tangible assets. Investing in real assets is a clear investment and form a kind, gold, construction, etc. Financial asset investment is an investment in financial assets in the form of securities in the capital market such as stocks, bonds, and deposits. When viewed from the difference in the number of Indonesian residents and investors. The number of investors in the capital market is 3,022,366 in the capital market as much as 1.125% of the total population of Indonesia, which is 268,835,016 people. This shows that very few Indonesian people are interested in investing in the capital market (Rizal, 2021).

Investment knowledge, lack of investment, and excess investment are some of the factors that can affect people's desire to invest. Information that has been organized in memory to form a unified network of structured information about investment is knowledge of investment (Rizal, 2021).

The emergence of several startups that offer Reksadana investment products to investors, one of which is PT Bibit Reksadana Growing Together, which makes a Reksadana application, namely Reksadana Bibit, with more than one million users, according to data provided by PT Bibit Grow Bersama on the Play Store which acts as a Reksadana selling agent. (APERD) by selling l34 Reksadana products.

The government made several efforts to foster public willingness to invest in the Indonesian capital market, namely the "Yuk Nabung Saham" campaign, which was launched in 2015. Today's high-speed technology has even become a necessity for some people(Indah Ramadhani, 2021)

With the accelerated development of information technology, especially in the financial industry, various financial institutions are trying to provide clients with access to various services and easy access to financial information. According to data from Digital 2021, the growth of social media users in Indonesia in January 2021 was 170.0 million social media users. This number increased by 6.3% or 10 million between 2020 and 2021(Gunawan, Pratiwi, & Pratama, 2018)(Diniyati, Triayudi, & Solehati, 2021).

Therefore, to find out public opinion, the author will conduct a sentiment analysis(Gunawan et al., 2018)(Park & Seo, 2018). There have been many studies conducted using sentiment analysis to find out the opinions of users of a product, political opinions, and opinions about movies that have been launched(Sandoval-Almazan & Valle-Cruz, 2018)(Ikoro, Sharmina, Malik, & Batista-navarro, 2020).

One way to get a public opinion in research is through social media. Because social media is usually used as a platform for buying and selling, providing information, opinions and as a medium to devote oneself. Based on these conditions, to conduct sentiment analysis, the researchers used two data sources, namely data from youtube, and reviewed data from the play store using scrapping methods.

Youtube users have now reached 30,000,000, people use Youtube to provide opinions/information about something. This opinion can be used to obtain information(Muhammad, 2019). Youtube is a social media that can be used as a medium of information in the form of videos, music, or writing, YouTube can also be used to share videos and films. Likewise, comments or opinions regarding an object such as reksadana. Writing youtube comments can be done by writing 200 characters into a message. An API (Application Programming Interface) facility that can make it easier for users to obtain data from social media. Data from youtube is used by several researchers, namely to solve problems such as the limitations of the sentiment dictionary(Balya, 2019). Sentiment analysis of the teacher's room on YouTube to find out the strengths and weaknesses of the Ruangguru application(Firdaus, Rizki, Gaus, & Susanto, 2020). Sentiment analysis of film reviews obtained from IMDB divides viewers' reactions to the films they watch into two groups, namely negative and positive reactions. Text mining is used in the analysis process to extract the information obtained and uses Naive Bayes classification. The reaction emotion will be tested by Chi Square(Amrullah, Sofyan Anas, & Hidayat, 2020).

There are many techniques and methods in data mining. To get maximum results, the researchers will use two methods for comparison, namely Naïve Bayes and K-nearest neighbor (KNN). The Naïve Bayes method is a method that can be used to group opinions properly (Wisnu, Afif, & Ruldevyani, 2020).

Naive Bayes can classify people's opinions into positive or negative. K-nearest neighbor (KNN) is a method that can form a classification of objects according to the assessment data, which is the closest step to the object. Learning data is predicted into multiple-dimensional sections, each dimension showing the initial features of the data. This section is divided into sections based on the description of the learning data. The best value of k for this algorithm depends on the data. In general, a high value of k will reduce the effect of noise on the classification but make the boundaries between each translation more blurred (Wisnu et al., 2020).

The research uses sentiment analysis that has been carried out by several researchers, sentiment analysis research on application reviews using the Naive Bayes method, the naive Bayes method can estimate the sentiment class on online application reviews by the system created, but the system created is not yet competent in carrying out the relevance and predicting sentiment class, this research can still be developed by using other algorithms to get more relevant analysis results(Olabenjo, 2016). To see the level of consumer satisfaction in digital payments, sentiment analysis was carried out using the Naïve Bayes classifier (NBC) and K-Nearest neighbor (KNN) classification methods(S et al., 2021). Two nave Bayes classification groupings were made with

review data from the google play store. This grouping was evaluated with several evaluation methods for comparison. The results prove that the Naïve Bayes classifier (NBC) method is good for structured distribution problems regarding software on the Google Play Store and for automating the categorization of android software on the Google Play Store (Fairuz, 2017)

**RESEARCH METHODS**

The stages of the research method for sentiment analysis, start from data collection until the results of the analysis are found.


Figure 1 Research Methode framework

**Data collection**
This study uses two data sources, namely Youtube data and application user review data obtained from the Playstore. Retrieval of data from youtube and play store review data using the scraping method with the python programming language.

Table 1 Data

| YouTube data | Review Data |
|---|---|
| 33,292 | 30,708 |

Table 1 results from scraping youtube data as much as 33,292 data, while review data is taken as much as 30,708 data.

**Text Processing**

Text processing is a very important step to carry out the classification process. At this stage, the attributes on the data that do not affect the classification process will be removed so that the result of this process is in the form of data that is clean and ready to be used.


Figure 2 Stages of processing text

Figure 2 flow in doing text processing with the Python programming language.


Figure 3 Example Of Data Result

Figure 3 is an example of the results of data retrieval in python on Youtube data.

**Case Folding**
At this stage, the letters in the text data will be changed to lowercase letters.


Figure 4 Case Folding Results

Figure 4 is an example of the case folding process, changing the sentence to lowercase, before "Penjualan…." After "penjualan….".

**Filtering**
In this process, the character, URL, hashtag, and punctuation stage is carried out because these components do not affect the sentiment value.

129

| | text |
|---|---|
| 0 | mantap jiwa |
| 1 | penjualan kok ada potongan rugi dong klu ada p... |
| 2 | mohon kembalikan dana saya atau negara anda yg... |
| 3 | muaaaaaannnntaaaappppppp |
| 4 | min kasih fitur dark mode dong klo buka app sa... |

Figure 5 Filter Result

Figure 5 is an example of character deletion filtering results before there were emojis in some data after the emojis in the data were lost.

**Tokenizing**

The process of cutting sentences into pieces of text is accompanied by the removal of numbers, characters, and punctuation marks.

| | text |
|---|---|
| 0 | [mantap, jiwa] |
| 1 | [penjualan, kok, ada, potongan, rugi, dong, kl... |
| 2 | [mohon, kembalikan, dana, saya, atau, negara, ... |
| 3 | [muaaaaaannnntaaaappppppp] |
| 4 | [min, kasih, fitur, dark, mode, dong, klo, buk... |

Figure 6 Tokenizing Result

Figure 6 is an example of tokenizing word cuts in each sentence.

**Stopword Removal**

*A stopword* is a list of general terms that have no significance and are not used. In this process, the common words will be removed to reduce the number of terms stored by the system.

| | text |
|---|---|
| 0 | mantap jiwa |
| 1 | penjualan potongan rugi klu potongannya |
| 2 | mohon kembalikan dana negara terus menerus mem... |
| 3 | muaaaaaannnntaaaappppppp |
| 4 | kasih fitur dark mode klo buka app mlam hari b... |

Figure 7 Stopword Removal Results

Figure 7 is an example of the removal process before "penjualan kok ada...." after "penjualan potongan...."

**Stemming**
So that the words in the sentiment sentence turn into basic words, the *stemming* process is used.

| | text |
|---|---|
| 0 | mantap jiwa |
| 1 | jual potong rugi klu potong |
| 2 | mohon kembali dana negara terus terus panen be... |
| 3 | muaaaaaannnntaaaappppppp |
| 4 | kasih fitur dark mode klo buka app mlam hari b... |

Figure 8 The results of an example of the stemming process

Figure 8 is an example of the stemming process of changing common words into basic words, before "penjualan potongan..." after "jual potong..."

**Labels**
*Labeling* is the process of giving positive and negative labels to sentiment data, and this *labeling* process uses the text blob method. Textblob is *a* processing *library* in NLP *(Natural Language Processing).*

| | label | text |
|---|---|---|
| 30498 | Positif | so far so good |
| 30505 | Positif | ok |
| 30509 | Positif | easy to use |
| 30525 | Positif | really helpfull untuk mulai investasi waktu verifikasi rasional lumrah registrasi mudah cepat |
| 30529 | Positif | nice |
| 30532 | Positif | good |
| 30582 | Positif | kode sayakaya dapet cashback k pakai nya user friendly mula takut robo bantu |
| 30583 | Positif | ok |
| 30592 | Positif | apps modern mudah aplikasi proses cepat edukatif |
| 30593 | Positif | mantap ajang ajar reksadana bikin tarik untuk coba baru top dech nyesel kalo gak nyoba rasa iuta |

Figure 9 Labeling Results

Figure 9 examples of positive and negative labeling in each sentence.

**RESULTS AND DISCUSSION**

This study has two data sources, namely data retrieval from YouTube and from user review data on the Playstore which was taken using the scraping method with the keywords "bibit" and "reksadana bibit" with a time range from January 2020 - to February 2022. The results of the python program for data retrieval and labeling are in the form of a CSV file containing several sentiments with tweet, text, polarity, subjectivity, and negative and positive labels.

Table 2 Label

|  | Amount | Positive | Negative | Neutral |
|---|---|---|---|---|
| YouTube | 33,292 | 1415 | 400 | 31,477 |
| Review | 30,708 | 3673 | 259 | 26,766 |

Table 2 results from labeling data on 33,292 YouTube data with 1415 positive sentiment results, 31,477 neutral, 400 negative, and from user reviews data on the Playstore as much as 3 0.708 with 3673 positive sentiment results, 26,766 neutral, and 259 negatives. The sentiment analysis process for reksadana has been carried out using the python program. For the labeling process, each data uses the text blob library.



Figure 10 Visualization of Labeling Tweet Bibit Data

Figure 10 Graph of the number of negative, neutral, positive sentiment data on YouTube comments.



Figure 11 Spreading the polarity of Bibit tweet data.

Figure 11 distribution of sentiment data based on polarity and subjectivity values on YouTube comment data.



Figure 12 Results of labeling data

Figure 12 Graph of the number of negative, neutral, positive sentiment data for the community in the Playstore review with the Bibit application.
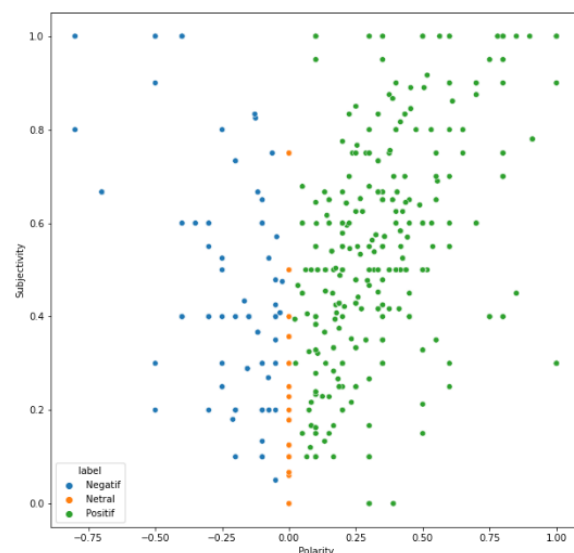


Figure 13 Polarity Spread of Bibit Review Data

Figure 11 distribution of sentiment data based on polarity and subjectivity values in the Playstore review in the Bibit application.

**Sentiment Analysis**
To determine the accuracy level of sentiment analysis of Reksadana users, two classification methods are used for comparison, namely the Naïve Bayes Classifier (NBC) and K-nearest neighbor (KNN), Support Vector Machine (SVM) methods.

**K-Nearest Neighbor(KNN)**
The accuracy obtained from testing YouTube data and review data with the bibit keyword is 99% and 99%.

131

```
Tingkat Akurasi :99 persen
              precision    recall  f1-score   support

        -1        1.00      1.00      1.00        87
         0        1.00      1.00      1.00      7858
         1        1.00      0.99      1.00       378

  accuracy                            1.00      8323
 macro avg        1.00      1.00      1.00      8323
weighted avg      1.00      1.00      1.00      8323
```

Figure 14 Accuracy Results For Bibit of Data
YouTube

Figure 14 is the result of the accuracy test of the k - nearest neighbor method on YouTube bibit data, with 24,969 test data and 8323 test data, the results of the test data accuracy test on YouTube data with the number of negative sentiments 87, neutral 7858 and positive 378 the results are 99%.

```
Tingkat Akurasi :99 persen
              precision    recall  f1-score   support

        -1        1.00      0.98      0.99        65
         0        1.00      1.00      1.00      6741
         1        1.00      1.00      1.00       871

  accuracy                            1.00      7677
 macro avg        1.00      0.99      1.00      7677
weighted avg      1.00      1.00      1.00      7677
```

Figure 15 Accuracy Results For Review Bibit

Figure 15 is the result of the accuracy of the classification test of the k - nearest neighbor method on bibit review data, with 23,031 train data and 7677 test data, the results of the test data accuracy test on review data with the number of sentiments 65 negatives, 871 positive, and 6741 neutral, the results are 99%.
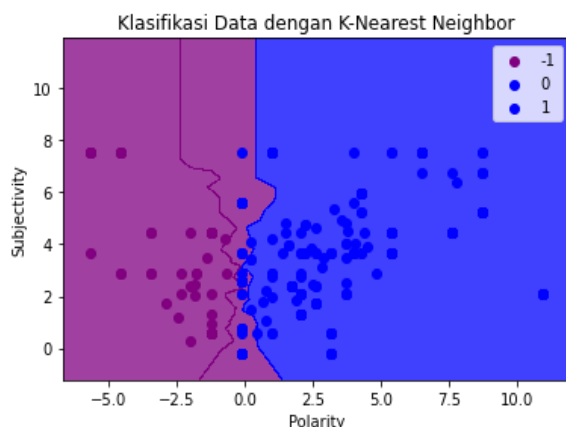


Figure 16 Youtube Bibit Data Visualization

Figure 16 visualization of YouTube comment data classification results with the k-nearest neighbor method based on sentiment negative, neutral and positive.
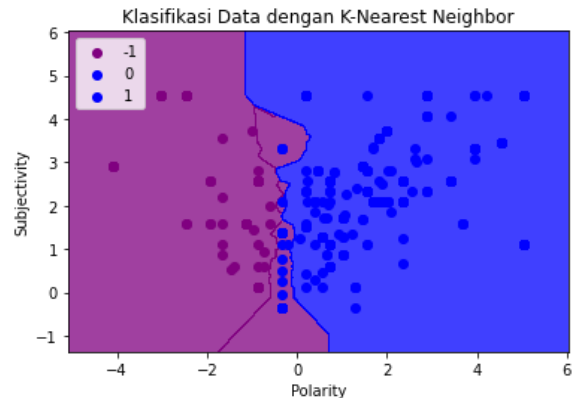


Figure 17 Visualization KNN Bibit Review Data

Figure 17 Visualization of bibit review data distribution with the k- nearest neighbor algorithm based on sentiment negative, neutral and positive.

**Naïve Bayes Classifier (NBC)**

```
Akurasi Naive Bayes :99 persen
              precision    recall  f1-score   support

        -1        0.92      1.00      0.96        80
         0        1.00      1.00      1.00      6237
         1        0.97      1.00      0.99       342

  accuracy                            1.00      6659
 macro avg        0.96      1.00      0.98      6659
weighted avg      1.00      1.00      1.00      6659
```

Figure 18 Accuracy Result NBC Bibit YouTube Data

Figure 18 results of the accuracy test of the naive Bayes method on bibit YouTube data, with 26,633 train data and 6659 test data, the results of the test data accuracy test on YouTube data with a negative sentiment count of 80, neutral 6237 and positive 342 the results are 99%.

```
Akurasi Naive Bayes :98 persen
              precision    recall  f1-score   support

        -1        0.28      1.00      0.43        23
         0        1.00      0.99      0.99      5282
         1        1.00      1.00      1.00       837

  accuracy                            0.99      6142
 macro avg        0.76      1.00      0.81      6142
weighted avg      1.00      0.99      0.99      6142
```

Figure 19 Accuracy Results of NBC Bibit Data
Reviews

Figure 19 is the result of the accuracy test of the naive Bayes method on bibit review data, with 24,566 train data and 6142 test data, the results of the test data accuracy test on YouTube data with the number of negative sentiments 23, neutral 5282, and positive 837 the results are 98%.

**Support Vector Machine (SVM)**

```
model accuracy is: 99.9249136506983 %
              precision    recall  f1-score   support

         -1       1.00      0.99      0.99        82
          0       1.00      1.00      1.00      6306
          1       1.00      0.99      0.99       271

   accuracy                           1.00      6659
  macro avg       1.00      0.99      1.00      6659
weighted avg      1.00      1.00      1.00      6659
```

Figure 20 YouTube data distribution accuracy results

Figure 20 results of the accuracy of the SVM classification test on YouTube bibit data, with 26,633 train data and 6659 test data, the results of the test data accuracy test on YouTube data with the number of negative sentiments 82, neutral 6306 and positive 271 the results are 99%.

```
model accuracy is: 99.9348746336698 %
              precision    recall  f1-score   support

         -1       1.00      0.98      0.99        55
          0       1.00      1.00      1.00      5340
          1       1.00      1.00      1.00       747

   accuracy                           1.00      6142
  macro avg       1.00      0.99      1.00      6142
weighted avg      1.00      1.00      1.00      6142
```

Figure 21 The results of the accuracy of the distribution of review data

Figure 19 is the result of the accuracy of the Naive Bayes method of classification test on bibit review data, with *24,566* train data and 6142 test data, the results of the test data accuracy test on YouTube data with the number of negative sentiments 55, neutral 5340, and positive 747 the results are 99%.
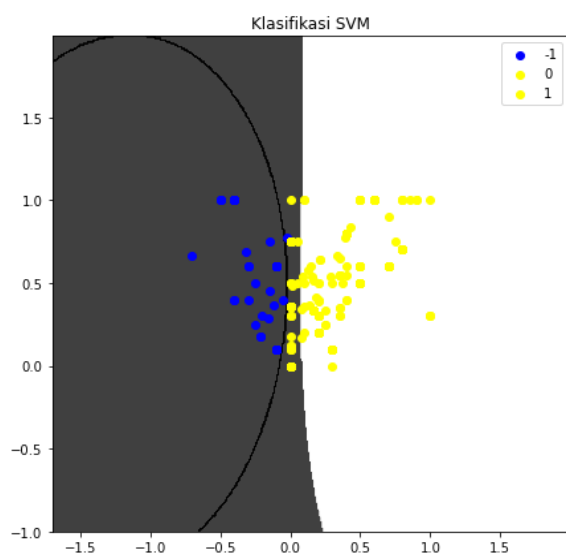

Figure 22 visualization SVM YouTube data

Figure 22 visualization of YouTube comment data classification results with Svm on sentiment negative, neutral, and positive.
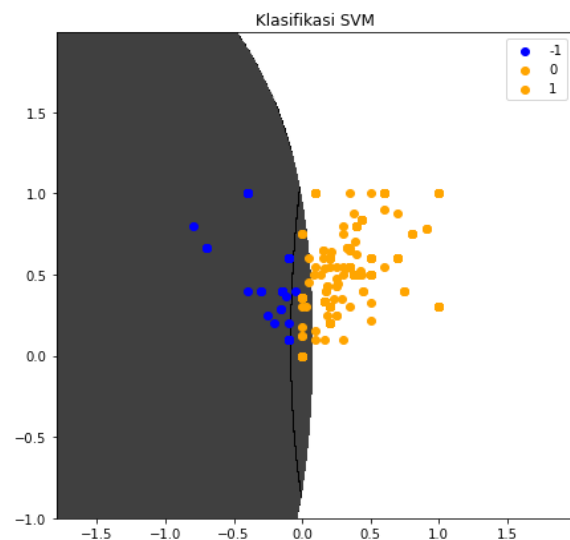

Figure 23 visualization SVM review data

Figure 23 visualization of review data classification results with Svm on sentiment negative, neutral, and positive.

Table 3 Accuracy

|  | NBC | | | KNN | | | SVM | | |
|---|---|---|---|---|---|---|---|---|---|
|  | Acc | + | - | Acc | + | - | Acc | + | - |
| Yt | 99% | 342 | 80 | 99% | 378 | 87 | 99% | 271 | 82 |
| Review | 98% | 837 | 23 | 99% | 871 | 65 | 99% | 747 | 55 |

Table 3 shows the accuracy results of the Naive *Bayes* and *k-nearest neighbor* and Svm methods.

**CONCLUSION**

From the results of sentiment analysis carried out with YouTube data sources and review data, using the *k-nearest neighbor (KNN) classification method*, *nave Bayes classifier,* and SVM on sentiment towards reksadana applications, it can be concluded that the 3 types of classification have the maximum level of accuracy for performing sentiment analysis, Based on the results of the YouTube data visualization and review data, it can be concluded that there are more neutral sentiments than positive sentiments, and more positive sentiments than negative opinions. The negative opinion obtained from the results of this analysis is on the buying, selling, and registration processes in conducting verification which takes a lot of time. The neutral opinion

obtained is because many people know about bibit but still lack understanding about bibit reksadana, while the positive opinion obtained is that the features are easy to understand so that people who want to start investing can use the reksadana application.

## REFERENCE

Amrullah, A. Z., Sofyan Anas, A., & Hidayat, M. A. J. (2020). Analisis Sentimen Movie Review Menggunakan Naive Bayes Classifier Dengan Seleksi Fitur Chi Square. *Jurnal*, *2*(1), 40–44. https://doi.org/10.30812/bite.v2i1.804

Balya. (2019). Analisis Sentimen Pengguna Youtube Di Indonesia Pada Review Smartphone Menggunakan Naïve Bayes.

Diniyati, D., Triayudi, A., & Solehati, I. D. (2021). *Analisa Interaksi Pengguna Media Sosial Perusahaan Sekuritas di Indonesia Saat Covid-19 menggunakan Social Network Analysis ( Studi Kasus : Indopremier dan Bursa Efek Indonesia )*. (January). https://doi.org/10.35870/jtik.v5i1.166

Fairuz, F. (2017). *Klasifikasi Review Software Pada Google Play Menggunakan Pendekatan Analisis Sentimen* (Universitas Gajah Mada). Universitas Gajah Mada. Retrieved from http://etd.repository.ugm.ac.id/penelitian/detail/112855

Firdaus, M. R., Rizki, F. M., Gaus, F. M., & Susanto, I. K. (2020). Analisis Sentimen Dan Topic Modelling Dalam Aplikasi Ruangguru. *J-SAKTI (Jurnal Sains Komputer Dan Informatika)*, *4*(1), 66. https://doi.org/10.30645/j-sakti.v4i1.188

Gunawan, B., Pratiwi, H. S., & Pratama, E. E. (2018). Sistem Analisis Sentimen pada Ulasan Produk Menggunakan Metode Naive Bayes. *Jurnal Edukasi Dan Penelitian Informatika (JEPIN)*, *4*(2), 118. https://doi.org/10.26418/jp.v4i2.27526

Ikoro, V., Sharmina, M., Malik, K., & Batista-navarro, R. (2020). 2020 7th International Conference on Social Network Analysis, Management and Security, SNAMS 2020. *2020 7th International Conference on Social Network Analysis, Management and Security, SNAMS 2020*, 95–98.

Indah Ramadhani, N. (2021). *Pengaruh Pengetahuan, Manfaat Dan Risiko Investasi Terhadapminat Investasi Padamahasiswafakultas Ekonomidan Bisnisuniversitassumatera Utara*.

Muhammad, A. N. (2019). Analisis Sentimen Positif Dan Negatif Komentar Video Youtube Menggunakan Metode Naïve Bayes - Support Vector Machine (Nbsvm) Classifier. *Skripsi*, 17.

Olabenjo, B. (2016). Applying Naive Bayes Classification to Google Play Apps Categorization. Retrieved from arXiv preprint website: https://arxiv.org/abs/1608.08574

Park, C. W., & Seo, D. R. (2018). Sentiment analysis of Twitter corpus related to artificial intelligence assistants. *2018 5th International Conference on Industrial Engineering and Applications, ICIEA 2018*, 495–498. https://doi.org/10.1109/IEA.2018.8387151

Rizal, S. (2021). Fenomena Penggunaan Platform Digital Reksa Dana Online dalam Peningkatan Jumlah Investor Pasar Modal Indonesia. *Humanis: Humanities, Management and Science Proceedings*, *1*(2), 851–861.

S, F. F., Si, M. N. S., Twitter, A., Giovani, A. P., Haryanti, T., Kurniawati, L., … Indah Ramadhani, N. (2021). Sistem Analisis Sentimen pada Ulasan Produk Menggunakan Metode Naive Bayes. *Proceeding SNTEI*, *1*(2), 113. https://doi.org/10.1088/1742-6596/1444/1/012034

Sandoval-Almazan, R., & Valle-Cruz, D. (2018). Facebook impact and sentiment analysis on political campaigns. *ACM International Conference Proceeding Series*. https://doi.org/10.1145/3209281.3209328

Wisnu, H., Afif, M., & Ruldevyani, Y. (2020). Sentiment analysis on customer satisfaction of digital payment in Indonesia: A comparative study using KNN and Naïve Bayes. *Journal of Physics: Conference Series*, *1444*(1). https://doi.org/10.1088/1742-6596/1444/1/012034