



REPUBLIC OF TURKEY
ADANA ALPARSLAN TÜRKEŞ SCIENCE AND TECHNOLOGY
UNIVERSITY

FACULTY OF ENGINEERING
DEPARTMENT OF COMPUTER ENGINEERING

Emotion Analysis From Sentiments

ECE NUR BALKAN
BACHELOR DEGREE
SUPERVISOR
ASST. PROF. DR. MÜMINE KAYA KELEŞ

ADANA 2023



REPUBLIC OF TURKEY
ADANA ALPARSLAN TÜRKEŞ SCIENCE AND TECHNOLOGY
UNIVERSITY

FACULTY OF ENGINEERING
DEPARTMENT OF COMPUTER ENGINEERING

Emotion Analysis From Sentiments

ECE NUR BALKAN
BACHELOR DEGREE

SUPERVISOR
ASST. PROF. DR. MÜMINE KAYA KELEŞ

ADANA 2023

ABSTRACT

Sentiment Analysis from Emotions Using Turkish Dataset

Ece Nur BALKAN

Department of Computer Engineering

Supervisor Asst. Prof. Dr. MÜMİNE KAYA KELEŞ

June 2023

The aim of this thesis is to develop a system to automatically classify and analyze emotional content in text data using sentiment analysis methods. Sentiment analysis is a natural language processing method used to identify emotional tone in text data and categorize texts into emotional categories such as positive, negative, or neutral. In this study, we included 6 emotions. These; 'positive', 'negative', 'neutral', 'happiness', 'fear', 'sadness'. For this process, the machine learning model and the Bert model were tested and the best results were tried to be obtained. A total of 6000 data sets were studied on the Türkiye data set. There are 1000 datasets for each emotion. Of these, 800 were used for learning and 200 for testing. Our total test dataset is 1200 and our learning dataset is 4800. The first 6 will be machine learning and the top 3 will be vectorized and embedded in the site. These 6 models are: KNN,NBM,RF,DT,SVM,LR. Results For this, the F-Mesaure and RMSE values will be taken into account. In addition, accuracy, precision, recall values will be examined.

Keywords: sentiment analysis, natural language processing, machine learning, bert modelling.

ACKNOWLEDGEMENTS

I would like to thank my consultant Asst. Prof. Dr. Mümine Kaya Keleş, who always helped and supported me in this thesis journey.

TABLE OF CONTENTS

TABLE OF CONTENTS	iii
LIST OF FIGURES	iv
LIST OF TABLES	v
NOMENCLATURE	vi
FORMULA	vii
PRE- Processing Conditions	viii
1. INTRODUCTION	1
1.1. Description of the problem	1
1.2. AIM	1
2. SENTIMENT ANALYSIS	1
2.1. What is sentiment analysis?	1
2.2. Turkish sentiment analysis	2
3. LITERATURE REVIEW	3
3.1. Turkish sentiment classification studies	3
4. RESULTS USING MACHINE LEARNING AND THE BERT MODEL	7
4.1. Properties of the dataset	7
4.2. Preprocessing Steps	8
4.3. Which models did we use?	14
4.3.1. Methods Used	15
4.3.2 Applications and Libraries Used	16
4.4. What performance metrics did we use?	16
4.5. What is Bert?	18
5. SITE FORMATION	20
6. CONCLUSION AND FUTURE WORK	22
7. REFERENCES	24

LIST OF FIGURES

FIGURE 1 THIS IS THE GRAY WOLF ALGORITHM. THE ALGORITHM GAVE THE BEST RESULTS IN ITS STUDY.....	4
FIGURE 2 STUDIES AND METHODOLOGIES USED IN NATURAL LANGUAGE PROCESSING.....	6
FIGURE 3 WEBSITE	21
FIGURE 4 WEBSITE	21

LIST OF TABLES

TABLE 1 PROJECT_TABLE.....	10
TABLE 2 SVM	11
TABLE 3 LR.....	12
TABLE 4 NBM	13
TABLE 5 SVM,LR SERI1 AND SERI2 = 0000, SERI8 AND SERI9 =0001 , SERI15,16 =0010 , SERI22,23=0011, SERI29,30 =0100, SERI36,37=0101, SERI43,44=0110, SERI50,51 =0111, SERI57,58 =1000, SERI64,65=1001, SERI71,72=1010, SERI78,79 =1011, SERI85,86= 1100, SERI92,93 =1101, SERI99,100 =1110, SERI106,107=1111	15
TABLE 6 TABLE3,4,5 EVALUATE ACCORDINGLY.....	17
TABLE 7 THIS IS 0000.....	17
TABLE 8 THIS IS 0000.....	18
TABLE 9 THIS IS 0000.....	18
TABLE 10 BERTURK.....	19
TABLE 11 BERTURK.....	19
TABLE 12 BERTURK.....	20

NOMENCLATURE

KNN	: K - Nearest Neighbour
NBM	: Naive Bayes Multinomial
RF	: Random Forest
DT	: Decision Table
SVM	: Support Vector Machine
LR	: Logistic Regression
RMSE	: Root Mean Squared Error
BERT	: Bidirectional Encoder Representations from Transformers

FORMULA

Accuracy	$:(TP+TN)/P+N$
Precision	$:TP/TP+FN$
Recall	$:TP/TP+FN$
F1-Score	$:2*precision*recall/precision+recall$
RMSE	$: \sqrt{\text{mean}((y_{\text{true}} - y_{\text{pred}})^2)}$

PRE- Processing Conditions

0000	: Raw
0001	: Stopword
0010	: Stemming
0011	: Stemming+ Stopword
0100	: Punctuation
0101	: Punctuation+ Stopword
0110	: Punctuation+ Stemming
0111	: Punctuation+ Stemming+ Stopword
1000	: Lowercase
1001	: Lowercase+ Stemming
1010	: Lowercase+ Stopword
1011	: Lowercase+ Stemming + Stopword
1100	: Lowercase+ Punctuation
1101	: Lowercase+ Punctuation + Stopword
1110	: Lowercase+ Punctuation + Stemming
1111	: Lowercase+ Punctuation + Stemming+Stopword

1. INTRODUCTION

1.1. Description of the problem

Nowadays, we do most of our work online. Some of us follow apps instead of shopping, studying online, and going to the movies. We also leave a treasure after each transaction. Comments. When we process them, they become very valuable.[1]As the author mentioned in his article. Unfortunately, 80% of the data on the market is unstructured. Building them is a process. However, as we mentioned in our topic, processing Turkish comments makes our job even more difficult. Since Turkish is an agglutinative language, we cannot say it negatively when we see the word 'not' as in English.

1.2. AIM

Machine learning models and BERTurk model will be examined by using Turkish data set. The results will be compared and the model with the best results will be embedded in the background of the site. The sentences entered on the created site will be analyzed accurately and quickly.

2. SENTIMENT ANALYSIS

2.1. What is sentiment analysis?

Sentiment analysis uses various methods and algorithms to determine how a text document or comment is perceived emotionally. These methods include machine learning, deep learning and natural language processing techniques. Text data can be analyzed by these methods and associated with emotion labels such as positive, negative or neutral. In addition, more detailed analyzes such as the intensity of an emotional text, the weights of the emotional components or the emotional tendencies on a particular subject can be made.

Sentiment analysis plays an important role in many areas. Businesses try to understand the emotional reactions of customers with sentiment analysis in areas such as social media analysis, customer feedback, brand management and marketing strategies. Policy makers can assess the impact of policies by monitoring social media and analyzing public reactions. In the social sciences, sentiment analysis is used to monitor emotional changes in society, understand social

trends, and understand human behavior.

It is also widely used in areas such as sentiment analysis, big data analytics and text mining. It helps people process large amounts of text data more quickly and effectively. It also provides valuable information in many application areas such as customer satisfaction analysis, product reviews, social media analytics, and public reaction monitoring. Sentiment analysis is a tool that helps us better understand people's emotional responses and improve informed decision making.

2.2. Turkish sentiment analysis

Turkish sentiment analysis is a technique used to classify Turkish text data into emotional categories. Turkish sentiment analysis aims to understand and analyze emotional tone by applying natural language processing (NLP) methods and algorithms on Turkish text data.

The following steps can be followed to do sentiment analysis in Turkish:

Data Collection: In order to collect Turkish text data, it is necessary to collect data from appropriate sources (social media, websites, surveys, customer feedback, etc.).

Data Pre-Processing: Pre-processing steps should be applied on the collected Turkish text data. These steps include cleaning up text (removing punctuation, removing unnecessary characters, converting to upper-case letters, etc.), removing stop-words, and identifying root words.

Feature Extraction: It is necessary to extract meaningful features from Turkish text data. In this step, text data is converted into numerical representations using methods such as word selection, word frequency, TF-IDF (Term Frequency-Inverse Document Frequency).

Model Training: A sentiment analysis model is trained using Turkish text data and tags from which feature extraction is performed. This model can be a machine learning algorithm or a deep learning model to identify emotional tones in Turkish texts.

Model Validation and Adjustments: The trained model is tested on the validation dataset and its performance is evaluated. If necessary, model adjustments are made and retested.

Sentiment Analysis Application: The trained and validated model is used to perform sentiment analysis on real-world data. Turkish text data is classified into emotional categories (positive, negative, neutral, etc.) by the trained model.

Turkish sentiment analysis can provide valuable information in areas such as understanding the emotional content of Turkish texts, social media analysis, brand management, customer feedback, marketing strategies and public opinion analysis. By performing sentiment analysis on Turkish texts, it may be possible to understand the emotional expressions specific to the Turkish language and culture and to make effective decisions.

3. LITERATURE REVIEW

3.1. Turkish sentiment classification studies

In the first study in which sentiment analysis was performed for Turkish [16], labeled sentences from Turkish children's tales and the translation of ISEAR [19] dataset were analyzed on 4000 samples in the categories of "joy, sadness, anger and fear". As a result of the study, in which traditional machine learning methods were applied, approximately 80% success was achieved. In another study in Turkish [20], sentiment analysis was conducted on Twitter, one of the social media environments. According to the hashtags used by the users in their posts, as a result of the sentiment analysis using traditional machine learning and classification methods on a total of 6000 tweets in each category, which were collected in Ekman's 6 emotion categories: "fear, anger, disgust, joy, sadness and surprise". 70 achievements have been achieved. In another study, 86% success was achieved as a result of Turkish sentiment analysis [190], which was performed in 7 classes using Naive Bayes. In another study conducted in Turkish [32], a sentiment dictionary was created by using attribute selection and weighting methods over the Turkish sentiment dataset [33], which had 26000 samples previously created by the researchers, and sentiment analysis was carried out with the help of this dictionary. As a result of the feeling analysis made in Ekman feeling categories, approximately 91% success was achieved. In another study, [18] performed sentiment analysis with deep learning algorithms on the TURTED sentiment dataset, which was created from Twitter posts containing keywords belonging to sentiment categories, and achieved 73% success. As can be seen, very few studies have been carried out in the field of Turkish sentiment analysis, and it has been observed that

the studies are generally done by emulating text classification studies and methods with a dictionary have also been tried.

```
Gri kurt popülasyonunu başlat  $X_i$  ( $i=1,2, \dots ,n$ )  
a, A ve C değerlerini başlat  
Her parçacığın komumunu ve uzaklığını hesapla  
 $X_\alpha$ = en iyi komumdaki parçacık  
 $X_\beta$ = ikinci en iyi komumdaki parçacık  
 $X_\delta$ = üçüncü en iyi komumdaki parçacık  
while ( $t < \text{maksimum iterasyon sayısı}$ )  
    for  $p$  in parçacık  
         $p$ 'nin komumunu güncelle (2.20)  
    end for  
     $a, A$  ve  $C$  değerlerini güncelle  
    Her parçacığın komumunu ve uzaklığını hesapla  
     $X_\alpha, X_\beta$  ve  $X_\delta$  değerlerini güncelle  
     $t=t+1$   
end while  
return  $X_\alpha$ 
```

Figure 1 This is the gray wolf algorithm. The algorithm gave the best results in its study.

Historical Development of Sentiment Analysis Studies

Sentiment analysis studies started to be done for English in the early 2000s, as can be seen in Figure 2 [3]. First of all, sentiment analysis studies were carried out with traditional machine learning (DML) methods [3,10,11]. The first sentiment analysis study for Turkish was carried out in 2013 using DME methods [16]. Another study for Turkish using DMR methods was published in 2014 [17]. The success has been increased by adding the attributes obtained from the hiss dictionaries as well as the words as input to the DMI methods [6–9,11]. With the use of deep learning (DL) methods in sentiment analysis in 2017, the prior knowledge obtained from dictionaries has been replaced by word occupants [12,13,14,15]. On the other hand, the creation of a validated and accessible sentiment analysis dataset (TREMO) in Turkish was carried out in 2018 [4]. In 2019, a limited dictionary of sensations was obtained by using the

word association statistics in the TREMO Turkish dataset, and Turkish sentiment analysis was carried out by using it together with the DME methods [5]. The Turkish sentiment analysis approach, on the other hand, could only be done at the end of 2019 [18]. Within the scope of the study with DÖ, the TURTED sentiment dataset consisting of Turkish Twitter posts was created and sentiment analysis was performed for Ekman sentiments. Obtaining more successful results with LL methods than DML methods and dictionary-supported hybrid methods has led to the need to add more prior knowledge to learning. With the pre-trained language model (ODM) approaches aiming to meet this need, the prior knowledge gained through big data has been transferred to deep learning methods.

Yıl	Referans	Dil	Yöntem	Veri Kümesi	Metrik	Sonuç%	Sözlük
2005	Alm ve ekibi [53]	İngilizce	GMÖ	Alm	F	47	WordNet
2007	Aman ve Szpakowicz [3]	İngilizce	GMÖ	Aman	Acc	74	GI + WAL
2008	Aman ve Szpakowicz[42]	İngilizce	GMÖ	Aman	F	59	WAL
2008	Danisman ve Alpkocak [57]	İngilizce	GMÖ	SE2007	F	32	
2010	Binali ve ekibi [23]	İngilizce	AKT	Blog gönderileri	Acc	96	
2010	Ghazi ve ekibi [43]	İngilizce	GMÖ	Alm + Aman	F	50	WAL
2010	Kim ve ekibi [44]	İngilizce	GMÖ	SE2007 + ISEAR + Alm	F	54	WAL + ANEW
2011	Chaffar ve Inkpen [45]	İngilizce	GMÖ	SE2007	Acc	40	WAL
2013	Perikos ve Hatzilygeroudis [24]	İngilizce	AKT	Elle oluşturulmuş	F	89	
2013	Boynukalın ve Karagöz [187]	Türkçe	GMÖ	Masallar + ISEAR çevirisi	Acc	81	
2014	Ghazi ve ekibi [46]	İngilizce	GMÖ	Aman	F	65	WAL
2014	Demirci [189]	Türkçe	GMÖ	Tweetler	Acc	70	
2014	Toçoğlu ve Alpkocak [190]	Türkçe	GMÖ	Açııcı	Acc	86	
2017	Bandhakavi ve ekibi [26]	İngilizce	GMÖ	ISEAR + SE2007+Blog	F	51	Kendi sözlüğü
2018	Toçoğlu ve Alpkocak [32]	Türkçe	GMÖ	TREMO	Acc	86	
2019	Ge ve ekibi [130]	İngilizce	DÖ	SE2019	F	75	
2019	Ma ve ekibi [152]	İngilizce	DÖ	SE2019	F	75	
2019	Basile ve ekibi [155]	İngilizce	ÖDM	SE2019	F	77	
2019	Zhong ve Miao [151]	İngilizce	ÖDM	SE2019	F	74	
2019	Xiao [150]	İngilizce	ÖDM	SE2019	F	77	
2019	Chatterjee ve ekibi [125]	İngilizce	ÖDM	SE2019	F	79	
2019	Huang ve ekibi [154]	İngilizce	ÖDM	Friends + EmotionPush	F	85	
2019	Toçoğlu ve Alpkocak [33]	Türkçe	GMÖ	TREMO	Acc	91	TEL
2019	Toçoğlu ve ekibi [191]	Türkçe	DÖ	TURTED	Acc	74	

*AKT: Anahtar kelime tabanlı, DÖ: Derin Öğrenme, GMÖ: Geleneksel makine öğrenmesi, ÖDM: Ön eğitilmiş dil modeli

Figure 2 Studies and methodologies used in natural language processing.

4. RESULTS USING MACHINE LEARNING AND THE BERT MODEL

Machine learning is a field of artificial intelligence that enables computer systems to gain the ability to solve a specific task or problem by automatically learning and experiencing it using predefined data and algorithms. Machine learning focuses on computer programs learning a particular task or problem through data-driven experiences and algorithms without being manually programmed.

Machine learning can perform complex tasks such as discovering patterns, relationships and trends, making predictions, classifying, clustering, building recommendation systems, and making decisions based on large amounts of data. Machine learning algorithms build models based on data and can use these models to perform tasks such as predicting future data or classifying new data.

It uses methods and techniques from fields such as machine learning, statistics, mathematics, data mining, and artificial intelligence. Data preprocessing includes steps such as feature selection, model training, model validation, and model improvement. Machine learning is widely used to solve real-world problems, make data-driven decisions, and improve automation processes.

It has achieved great success in various application areas such as machine learning, voice and image recognition, natural language processing, automated car driving, financial forecasting, medical diagnostics, product recommendations and target marketing. Thanks to machine learning, computers can process complex data sets, recognize patterns and learn in a similar way to humans.

4.1. Properties of the dataset

The data set consists of Turkish characters. There are 6000 data in total. 4800 is the learning data set. 1200 is the test data set.

4.2. Preprocessing Steps

To process text data for natural language processing applications such as sentiment analysis, the following preprocessing steps are typically applied:

Text Cleaning: It is the process of removing unnecessary elements from text data such as unnecessary characters, punctuation marks, numbers or special characters. In this step, special symbols such as HTML tags, links, emojis in the texts can be cleaned.

Tokenization: The process of breaking texts into smaller pieces. Texts are often divided into words or phrases. This means that sentences or paragraphs are segmented as words or symbols.

Upper / Lower Case Conversion: It is the process of converting all letters in the text to uppercase or lowercase letters. This is done to organize the word distribution in the text and to combine different letter cases of the same words.

Removal of Stop Words: It is the removal of commonly used words called grammatical junk words or stop words. Such words are often meaningless and limited in providing information in tasks such as sentiment analysis [2] Author's stopwords.txt is used for this step.

{acaba; altı; altmış; ama; ancak; artık; asla; aslında; az; bana; bazen; bazı; bazıları; bazısı; belki; ben; bende; benden; beni; benim; beş; bile; bin; bir; biri; birçoğu; birçok; birçokları; biri; birisi; birlikte; birkaç; birkaçı; birkez; birşey; birşeyi; biz; bizde; bizden; bize; bizi; bizim; böyle; böylece; bu; bu arada; bu yüzden; buarada; buna; bunda; bundan; bunu; bunun; burada; burda; buyüzden; bütün; çoğu; çoğuna; çoğunu; çok; çünkü; da; daha; dahi; de; defa; değil; demek; diğer; diğeri; diğerleri; diye; doksan; dokuz; dolayı; dört; elbette; elli; en; fakat; falan; felan; filan; gene; gibi; hâlâ; hangi; hangisi; hani; hatta; hem; henüz; hep; hepsi; hepsinde; hepsinden; hepsine; hepsini; hepsinin; her; her biri; herbiri; herkes; herkese; herkesi; hiç; hiç kimse; hiçbir; hiçbirisi; hiçbirine; hiçbirini; hiçkimse; için; içinde; ilk; iki; ile; ise; işte; kaç; kadar; katrilyon; kendi; kendine; kendini; kendisi; kendisine; kez; kırk; ki; kim; kimde; kimden; kime; kimi; kimin; kimisi; kimse; madem; mı; mi; milyar; milyon; mu; mü; nasıl; ne; ne kadar; ne zaman; nekadardır; ne zaman; neden; nedir; nerde; neredede; nereden; nereye; nesi; neyse; niçin; niye; olan; olarak; on; ona; onda; ondan; onlar; onlara; onlarda; onlardan; onları; onların; onu; onun; orada; otuz; oysa; oysaki; öbürü; ön; önce; ötürü; öyle; rağmen; sana; sanki; sekiz; seksen; sen; sende; senden; seni; senin; siz; sizde; sizden; size; sizi; sizin; son; sonra; şayet; şey; şeyde; şeyden;

şeye; şeyi; şeyler; şimdi; şöyle; şu; şuna; şunda; şundan; şunlar; şunu; şunun; tabi; tam; tamam; trilyon; tüm; tümü; üç; üzere; var; ve; veya; veyahut; ya; ya da; yada; yani; yedi; yerine; yetmiş; yine; yirmi; yoksa; yüz; yüzden; zaten; zira}

Root Extraction or Lemmatization: It is the process of obtaining the roots or root forms of words. This includes removing changes such as inflections and plurals while keeping the basic meaning of the word. Root extraction and lemmatization can reduce word diversity in the text and increase the consistency of the analysis.

Text Vectorization: It is the process of converting texts into numerical vector representations. This is necessary so that machine learning algorithms can process text data. Techniques such as TfidfVectorizer can create vectors that reflect the frequency and importance of words in texts.

These preprocessing steps are used to represent text data in a more workable and more consistent way. The preprocessing steps can be customized depending on the particular sentiment analysis application and the dataset to be used.

Preprocessing Code	Algorithms	Evaluation Criteria				
		F-Measure	Accuracy	RMSE	Precision	Recall
0000	KNN [k=3]	0.2025458436745291	0.2544444444444444		0.5568838711972174	0.26517059381692065
	NaiveBayesMultinomial	0.8983407399520978	0.8961111111111111		0.8771343493072553	0.8688414480429264
	Random Forest	0.7188018256594041	0.7944444444444444		0.7112061607401579	0.72027081114497
	Decision Table	0.7134590188440265	0.7133333333333334		0.7915223808580111	0.7977712533964189
	Support Vector Machine	0.8955555555555555	0.8955555555555555	0.672887641	0.8972435193145606	0.8981820657546349
	Logistic Regression	0.889241311648938	0.8877777777777778	0.672887641	0.889460926060963	0.8900713108518552
0001	KNN [k=3]	0.20120041655751078	0.2538888888888889		0.5552332880577155	0.2646223482028856
	NaiveBayesMultinomial	0.8700656870225313	0.8677777777777778		0.8771343493072553	0.8688414480429264
	Random Forest	0.7831851663151608	0.7872222222222223		0.7844236726251138	0.7924378536308215
	Decision Table	0.7167227013328731	0.7166666666666667		0.7141644024765125	0.7234553595308692
	Support Vector Machine	0.8972842542035222	0.8955555555555555	0.677823314	0.8972435193145606	0.8981820657546349
	Logistic Regression	0.8892413116489383	0.8877777777777778	0.672887641	0.889460926060963	0.8900713108518552
0010	KNN [k=3]	0.230192997689212	0.230192997689212		0.5821077144652045	0.2855873818432311
	NaiveBayesMultinomial	0.8771636469900764	0.8755555555555555		0.8837982088847814	0.8764771145831646
	Random Forest	0.7833977619806577	0.7877777777777778		0.7850654930832701	0.79293230314672
	Decision Table	0.7196109822824047	0.7205555555555555		0.7167342745105332	0.7272200020443763
	Support Vector Machine	0.8938213669682734	0.8922222222222222	0.679869268	0.8947200692784956	0.8947559836070967
	Logistic Regression	0.8919958298962708	0.8905555555555555	0.7	0.8925351339546571	0.8928060090674865
0011	KNN [k=3]	0.230192997689212	0.2738888888888889		0.5821077144652045	0.2855873818432311
	NaiveBayesMultinomial	0.8771636469900764	0.8755555555555555		0.8837982088847814	0.8764771145831646
	Random Forest	0.7803061056981822	0.7844444444444445		0.7815159060050961	0.7895275194321096
	Decision Table	0.722566944122834	0.7238888888888889		0.7182861604522226	0.7307752081809132
	Support Vector Machine	0.8938213669682734	0.8922222222222222	0.679869268	0.8947200692784956	0.8947559836070967
	Logistic Regression	0.8919958298962708	0.8905555555555555	0.7	0.8925351339546571	0.8928060090674865
0100	KNN [k=3]	0.2403257189770034	0.2827777777777778		0.5839366364822548	0.29490355093790144
	NaiveBayesMultinomial	0.8749343216325381	0.8733333333333333		0.8820392911070877	0.8742084102287239
	Random Forest	0.7764100468905507	0.7822222222222223		0.7780460928660674	0.787667972758332
	Decision Table	0.70219359045996	0.7038888888888889		0.7000007873414006	0.7109961394524631
	Support Vector Machine	0.8955555555555555	0.89692757302259	0.681909085	0.8970426907861918	0.8978841623377112
	Logistic Regression	0.8903426294330576	0.8888888888888888	0.717247826	0.8909028945431663	0.8911261046098043
0101	KNN [k=3]	0.2403257189770034	0.2827777777777778		0.5839366364822548	0.29490355093790144
	NaiveBayesMultinomial	0.8749343216325381	0.8733333333333333		0.8820392911070877	0.8742084102287239
	Random Forest	0.7764100468905507	0.7822222222222223		0.7780460928660674	0.787667972758332
	Decision Table	0.70219359045996	0.7038888888888889		0.7000007873414006	0.7109961394524631
	Support Vector Machine	0.8955555555555555	0.89692757302259	0.681909085	0.8970426907861918	0.8978841623377112
	Logistic Regression	0.8903426294330576	0.8888888888888888	0.717247826	0.8909028945431663	0.8911261046098043
0110	KNN [k=3]	0.21960717259323825	0.2738888888888889		0.6476022390130414	0.2776802987654326
	NaiveBayesMultinomial	0.8664345717229379	0.8655555555555555		0.8733306026278201	0.8673955794888789
	Random Forest	0.7789029145239578	0.7827777777777778		0.7832557418731119	0.7856803690939663
	Decision Table	0.6869963279475427	0.6877777777777778		0.6858757975845272	0.6918499963922383
	Support Vector Machine	0.875	0.8757952349270427	0.747217059	0.8753666002918087	0.8769251322343073
	Logistic Regression	0.8826955475629034	0.8816666666666667	0.757554545	0.8841355744892475	0.8827268156766364
0111	KNN [k=3]	0.21960717259323825	0.2738888888888889		0.6476022390130414	0.2776802987654326
	NaiveBayesMultinomial	0.8664345717229379	0.8655555555555555		0.8733306026278201	0.8673955794888789
	Random Forest	0.7756532198470895	0.7794444444444445		0.7811654187291458	0.7822449807385142
	Decision Table	0.6946349891815661	0.695		0.69392933758085819	0.6988679100598828
	Support Vector Machine	0.8757952349270427	0.875	0.747217059	0.8753666002918087	0.8769251322343073
	Logistic Regression	0.8826955475629034	0.8816666666666667	0.757554545	0.8841355744892475	0.8827268156766364
1000	KNN [k=3]	0.20234818052755418	0.2594444444444444		0.5986440088629343	0.2631993888160669
	NaiveBayesMultinomial	0.8742323701085927	0.8722222222222222		0.8801192594275608	0.87383929179525
	Random Forest	0.7803140338689646	0.7816666666666666		0.7838100240412164	0.7838891115257711
	Decision Table	0.7028464106085647	0.7016666666666667		0.7046974154886304	0.7049312566192635
	Support Vector Machine	0.8833712516645592	0.8827777777777778	0.707106781	0.8829188435000631	0.8845525165484549
	Logistic Regression	0.8838018689011878	0.8827777777777778	0.708676388	0.884583108573432	0.8841136758409034
1001	KNN [k=3]	0.2403257189770034	0.2827777777777778		0.5839366364822548	0.29490355093790144
	NaiveBayesMultinomial	0.8749343216325381	0.8733333333333333		0.8820392911070877	0.8742084102287239
	Random Forest	0.7872576273122558	0.7905555555555556		0.7880154200859124	0.7957670961514084
	Decision Table	0.708962450369226	0.7094444444444444		0.706358935561669	0.7160348526531776
	Support Vector Machine	0.89692757302259	0.8955555555555555	0.681909085	0.8970426907861918	0.8978841623377112
	Logistic Regression	0.8903426294330576	0.8888888888888888	0.717247826	0.8909028945431663	0.8911261046098043
1010	KNN [k=3]	0.2403257189770034	0.2827777777777778		0.5839366364822548	0.29490355093790144
	NaiveBayesMultinomial	0.8749343216325381	0.8733333333333333		0.8820392911070877	0.8742084102287239
	Random Forest	0.7765612248582757	0.7811111111111111		0.7779456400862776	0.7868524883339312
	Decision Table	0.6967241047505869	0.6977777777777778		0.6938216626014063	0.704935517818129
	Support Vector Machine	0.89692757302259	0.8955555555555555	0.681909085	0.8970426907861918	0.8978841623377112
	Logistic Regression	0.8903426294330576	0.8888888888888888	0.717247826	0.8909028945431663	0.8911261046098043
1011	KNN [k=3]	0.2827777777777778	0.2403257189770034		0.5839366364822548	0.29490355093790144
	NaiveBayesMultinomial	0.8749343216325381	0.8733333333333333		0.8820392911070877	0.8742084102287239
	Random Forest	0.7774748918057734	0.7822222222222223		0.7786699690046371	0.7876781214075605
	Decision Table	0.708283252279609	0.7088888888888889		0.7051005390928622	0.715487066814708
	Support Vector Machine	0.89692757302259	0.8955555555555555	0.681909085	0.8970426907861918	0.8978841623377112
	Logistic Regression	0.8903426294330576	0.8888888888888888	0.717247826	0.8909028945431663	0.8911261046098043
1100	KNN [k=3]	0.2244822081546037	0.2244822081546037		0.6196585653512974	0.2821697785838598
	NaiveBayesMultinomial	0.8643041245292089	0.8633333333333333		0.8709874573007607	0.8651444155227029
	Random Forest	0.7796568054509887	0.7827777777777778		0.7842661919334745	0.7858730667248266
	Decision Table	0.6958504317467199	0.6944444444444444		0.7001056291644675	0.6979871624614705
	Support Vector Machine	0.8795575858216521	0.8788888888888889	0.722264956	0.8791810321440149	0.8807521855158532
	Logistic Regression	0.885335575261615	0.8844444444444445	0.748331477	0.8867861913134956	0.8856679345229744
1101	KNN [k=3]	0.2244822081546037	0.2783333333333333		0.6196585653512974	0.2821697785838598
	NaiveBayesMultinomial	0.8643041245292089	0.8633333333333333		0.8709874573007607	0.8651444155227029
	Random Forest	0.7721224901056951	0.7761111111111111		0.7757577606280016	0.77937096197375
	Decision Table	0.6981865241644233	0.6983333333333334		0.6976162660094768	0.7022321646095754
	Support Vector Machine	0.8795575858216521	0.8788888888888889	0.722264956	0.8791810321440149	0.8807521855158532
	Logistic Regression	0.885335575261615	0.8844444444444445	0.748331477	0.8867861913134956	0.8856679345229744
1110	KNN [k=3]	0.2244822081546037	0.2783333333333333		0.6196585653512974	0.2821697785838598
	NaiveBayesMultinomial	0.8643041245292089	0.8633333333333333		0.8709874573007607	0.8651444155227029
	Random Forest	0.7776245032505241	0.7805555555555556		0.7819059059892411	0.7834079888037643
	Decision Table	0.6963029524662297	0.6955555555555556		0.6972732748954147	0.6994053381417041
	Support Vector Machine	0.8795575858216521	0.8788888888888889	0.722264956	0.8791810321440149	0.8807521855158532
	Logistic Regression	0.885335575261615	0.8844444444444445	0.748331477	0.8867861913134956	0.8856679345229744
1111	KNN [k=3]	0.2244822081546037	0.2783333333333333		0.6196585653512974	0.2821697785838598
	NaiveBayesMultinomial	0.8643041245292089	0.8633333333333333		0.8709874573007607	0.8651444155227029
	Random Forest	0.7743793533025242	0.7783333333333333		0.7772315141406771	0.7810197495032346
	Decision Table	0.6907899491867472	0.6911111111111111		0.6906015425483717	0.6947368486780902
	Support Vector Machine	0.8795575858216521	0.8788888888888889	0.722264956	0.8791810321440149	0.8807521855158532
	Logistic Regression	0.885335575261615	0.8844444444444445	0.748331477	0.8867861913134956	0.8856679345229744

Table 1 Project_table

Preprocessing Code	SVM Algorithms	Evaluation Criteria			
		F-Measure	Precision	Recall	Support
0000	hüzün	0,9	0,98	0,83	204
	korku	0,95	0,99	0,92	189
	mutluluk	0,91	0,85	0,99	194
	negatif	0,84	0,84	0,83	215
	nötr	0,93	0,91	0,95	176
	pozitif	0,82	0,81	0,82	222
	macro avg	0,89	0,89	0,9	1200
0001	hüzün	0,9	0,98	0,83	204
	korku	0,95	0,99	0,92	189
	mutluluk	0,91	0,85	0,99	194
	negatif	0,84	0,84	0,83	215
	nötr	0,93	0,91	0,95	176
	pozitif	0,81	0,81	0,82	222
	macro avg	0,89	0,89	0,9	1200
0010	hüzün	0,9	0,98	0,84	204
	korku	0,95	0,98	0,92	189
	mutluluk	0,9	0,82	0,99	194
	negatif	0,84	0,83	0,84	215
	nötr	0,94	0,93	0,95	176
	pozitif	0,82	0,83	0,81	222
	macro avg	0,89	0,89	0,9	1200
0011	hüzün	0,9	0,98	0,84	204
	korku	0,95	0,98	0,92	189
	mutluluk	0,9	0,92	0,99	194
	negatif	0,84	0,83	0,84	215
	nötr	0,94	0,93	0,95	176
	pozitif	0,82	0,83	0,81	222
	macro avg	0,89	0,89	0,9	1200
0100	hüzün	0,9	0,97	0,84	204
	korku	0,95	0,98	0,92	189
	mutluluk	0,9	0,82	0,99	194
	negatif	0,84	0,83	0,84	215
	nötr	0,94	0,93	0,95	176
	pozitif	0,81	0,82	0,8	222
	macro avg	0,89	0,89	0,89	1200
0101	hüzün	0,9	0,97	0,84	204
	korku	0,95	0,98	0,92	189
	mutluluk	0,9	0,82	0,99	194
	negatif	0,84	0,83	0,84	215
	nötr	0,94	0,93	0,95	176
	pozitif	0,81	0,82	0,8	222
	macro avg	0,89	0,89	0,89	1200
0110	hüzün	0,88	0,96	0,81	197
	korku	0,93	0,97	0,88	189
	mutluluk	0,89	0,8	1	205
	negatif	0,86	0,9	0,82	215
	nötr	0,93	0,92	0,94	186
	pozitif	0,81	0,79	0,84	208
	macro avg	0,88	0,88	0,89	1200
0111	hüzün	0,88	0,96	0,81	197
	korku	0,93	0,97	0,88	189
	mutluluk	0,89	0,8	1	205
	negatif	0,86	0,9	0,82	215
	nötr	0,93	0,92	0,94	186
	pozitif	0,81	0,79	0,84	208
	macro avg	0,88	0,88	0,89	1200
1000	hüzün	0,89	0,96	0,83	197
	korku	0,94	0,98	0,9	189
	mutluluk	0,9	0,83	0,99	205
	negatif	0,85	0,9	0,8	215
	nötr	0,94	0,9	0,97	186
	pozitif	0,8	0,78	0,83	208
	macro avg	0,88	0,88	0,89	1200
1001	hüzün	0,9	0,97	0,84	204
	korku	0,95	0,98	0,92	189
	mutluluk	0,9	0,82	0,99	194
	negatif	0,84	0,83	0,84	215
	nötr	0,94	0,93	0,95	176
	pozitif	0,81	0,82	0,8	222
	macro avg	0,89	0,89	0,89	1200
1010	hüzün	0,9	0,97	0,84	204
	korku	0,95	0,98	0,92	189
	mutluluk	0,9	0,82	0,99	194
	negatif	0,84	0,83	0,84	215
	nötr	0,94	0,93	0,95	176
	pozitif	0,81	0,82	0,8	222
	macro avg	0,89	0,89	0,9	1200
1011	hüzün	0,9	0,97	0,84	204
	korku	0,95	0,98	0,92	189
	mutluluk	0,9	0,82	0,99	194
	negatif	0,84	0,83	0,84	215
	nötr	0,94	0,93	0,95	176
	pozitif	0,81	0,82	0,8	222
	macro avg	0,89	0,89	0,89	1200
1100	hüzün	0,88	0,96	0,81	197
	korku	0,93	0,97	0,88	189
	mutluluk	0,89	0,8	1	205
	negatif	0,86	0,9	0,81	215
	nötr	0,93	0,91	0,94	186
	pozitif	0,81	0,78	0,83	208
	macro avg	0,88	0,88	0,88	1200
1101	hüzün	0,88	0,96	0,81	197
	korku	0,93	0,97	0,88	189
	mutluluk	0,89	0,8	1	205
	negatif	0,86	0,9	0,81	215
	nötr	0,93	0,91	0,94	186
	pozitif	0,81	0,78	0,83	208
	macro avg	0,88	0,88	0,89	1200
1110	hüzün	0,88	0,96	0,81	197
	korku	0,93	0,97	0,88	189
	mutluluk	0,89	0,8	1	205
	negatif	0,86	0,9	0,81	215
	nötr	0,93	0,91	0,94	186
	pozitif	0,81	0,78	0,83	208
	macro avg	0,88	0,89	0,88	1200
1111	hüzün	0,88	0,96	0,81	197
	korku	0,93	0,97	0,88	189
	mutluluk	0,89	0,8	1	205
	negatif	0,86	0,9	0,81	215
	nötr	0,93	0,91	0,94	186
	pozitif	0,81	0,78	0,83	208
	macro avg	0,88	0,89	0,88	1200

Table 2 SVM

Algorithms	Evaluation Criteria				
	Precision	Recall	F1-Score	Support	
hüzün		0,97	0,85	0,91	204
korku		0,97	0,94	0,95	189
mutluluk		0,88	0,98	0,92	194
negatif		0,82	0,85	0,83	215
nötr		0,91	0,95	0,93	176
pozitif		0,82	0,8	0,81	222
macro avg	0,89		0,89	0,89	1200
hüzün		0,97	0,85	0,91	204
korku		0,97	0,94	0,95	189
mutluluk		0,88	0,98	0,92	194
negatif		0,82	0,85	0,83	215
nötr		0,91	0,95	0,93	176
pozitif		0,82	0,8	0,81	222
macro avg	0,89		0,89	0,89	1200
hüzün		0,97	0,87	0,92	204
korku		0,98	0,93	0,95	189
mutluluk		0,85	0,99	0,91	194
negatif		0,83	0,84	0,84	215
nötr		0,92	0,94	0,93	176
pozitif		0,83	0,79	0,81	222
macro avg	0,9		0,9	0,89	1200
hüzün		0,97	0,87	0,92	204
korku		0,98	0,93	0,95	189
mutluluk		0,85	0,99	0,91	194
negatif		0,83	0,84	0,84	215
nötr		0,92	0,94	0,93	176
pozitif		0,83	0,79	0,81	222
macro avg	0,9		0,9	0,89	1200
hüzün		0,97	0,86	0,91	204
korku		0,98	0,93	0,95	189
mutluluk		0,85	0,99	0,91	194
negatif		0,82	0,85	0,83	215
nötr		0,92	0,94	0,93	176
pozitif		0,83	0,78	0,8	222
macro avg	0,89		0,89	0,89	1200
hüzün		0,97	0,86	0,91	204
korku		0,98	0,93	0,95	189
mutluluk		0,85	0,99	0,91	194
negatif		0,82	0,85	0,83	215
nötr		0,92	0,94	0,93	176
pozitif		0,83	0,78	0,8	222
macro avg	0,89		0,89	0,89	1200
hüzün		0,95	0,83	0,89	197
korku		0,97	0,9	0,93	189
mutluluk		0,82	1	0,9	205
negatif		0,89	0,84	0,86	215
nötr		0,93	0,95	0,94	186
pozitif		0,81	0,81	0,81	208
macro avg	0,89		0,89	0,89	1200
hüzün		0,95	0,83	0,89	197
korku		0,96	0,93	0,94	189
mutluluk		0,86	0,98	0,91	205
negatif		0,88	0,81	0,85	215
nötr		0,92	0,97	0,95	186
pozitif		0,79	0,81	0,8	208
macro avg	0,89		0,89	0,89	1200
hüzün		0,97	0,86	0,91	204
korku		0,98	0,93	0,95	189
mutluluk		0,85	0,99	0,91	194
negatif		0,82	0,85	0,83	215
nötr		0,92	0,94	0,93	176
pozitif		0,83	0,78	0,8	222
macro avg	0,89		0,89	0,89	1200
hüzün		0,97	0,86	0,91	204
korku		0,98	0,93	0,95	189
mutluluk		0,85	0,99	0,91	194
negatif		0,82	0,85	0,83	215
nötr		0,92	0,94	0,93	176
pozitif		0,83	0,78	0,8	222
macro avg	0,89		0,89	0,89	1200
hüzün		0,97	0,86	0,91	204
korku		0,98	0,93	0,95	189
mutluluk		0,85	0,99	0,91	194
negatif		0,82	0,85	0,83	215
nötr		0,92	0,94	0,93	176
pozitif		0,83	0,78	0,8	222
macro avg	0,89		0,89	0,89	1200
hüzün		0,95	0,84	0,89	197
korku		0,97	0,9	0,93	189
mutluluk		0,82	1	0,9	205
negatif		0,89	0,82	0,85	215
nötr		0,92	0,95	0,93	186
pozitif		0,8	0,8	0,8	208
macro avg	0,89		0,89	0,89	1200
hüzün		0,95	0,84	0,89	197
korku		0,97	0,9	0,93	189
mutluluk		0,82	1	0,9	205
negatif		0,89	0,82	0,85	215
nötr		0,92	0,95	0,93	186
pozitif		0,8	0,8	0,8	208
macro avg	0,89		0,89	0,89	1200

Table 3 LR

Preporocessing Code	Algorithms	Evaulation Criteria			
		Precision	Recall	F1-Score	Support
0000	hüzün	0,93	0,87	0,9	204
	korku	0,87	0,99	0,93	189
	mutluluk	0,99	0,79	0,88	194
	negatif	0,79	0,9	0,84	215
	nötr	0,99	0,92	0,96	176
	pozitif	0,79	0,92	0,81	222
	macro avg	0,89	0,88	0,88	1200
0001	hüzün	0,93	0,87	0,9	204
	korku	0,87	0,99	0,93	189
	mutluluk	0,99	0,79	0,88	194
	negatif	0,79	0,9	0,84	215
	nötr	0,99	0,82	0,96	176
	pozitif	0,79	0,82	0,81	222
	macro avg	0,89	0,88	0,88	1200
0010	hüzün	0,91	0,9	0,9	204
	korku	0,87	0,99	0,93	189
	mutluluk	0,99	0,79	0,88	194
	negatif	0,81	0,89	0,85	215
	nötr	0,99	0,92	0,95	176
	pozitif	0,81	0,82	0,82	222
	macro avg	0,9	0,82	0,89	1200
0011	hüzün	0,91	0,9	0,9	204
	korku	0,87	0,99	0,93	189
	mutluluk	0,99	0,79	0,88	194
	negatif	0,81	0,89	0,85	215
	nötr	0,99	0,92	0,95	176
	pozitif	0,81	0,82	0,82	222
	macro avg	0,9	0,82	0,89	1200
0100	hüzün	0,9	0,89	0,89	204
	korku	0,86	0,99	0,92	189
	mutluluk	0,99	0,79	0,88	194
	negatif	0,8	0,88	0,84	215
	nötr	0,99	0,92	0,95	176
	pozitif	0,8	0,8	0,8	222
	macro avg	0,89	0,88	0,88	1200
0101	hüzün	0,9	0,89	0,89	204
	korku	0,86	0,99	0,92	189
	mutluluk	0,99	0,79	0,88	194
	negatif	0,8	0,88	0,84	215
	nötr	0,99	0,92	0,95	176
	pozitif	0,8	0,8	0,8	222
	macro avg	0,89	0,88	0,88	1200
0110	hüzün	0,93	0,86	0,89	197
	korku	0,85	1	0,92	189
	mutluluk	0,99	0,78	0,87	205
	negatif	0,83	0,87	0,85	215
	nötr	0,99	0,9	0,94	186
	pozitif	0,75	0,88	0,81	208
	macro avg	0,89	0,88	0,88	1200
0111	hüzün	0,93	0,86	0,89	197
	korku	0,85	1	0,92	189
	mutluluk	0,99	0,78	0,87	205
	negatif	0,83	0,87	0,85	215
	nötr	0,99	0,9	0,94	186
	pozitif	0,75	0,88	0,81	208
	macro avg	0,89	0,88	0,88	1200
1000	hüzün	0,93	0,85	0,89	197
	korku	0,86	1	0,92	189
	mutluluk	0,99	0,79	0,88	205
	negatif	0,86	0,86	0,86	215
	nötr	0,99	0,88	0,93	186
	pozitif	0,72	0,88	0,79	208
	macro avg	0,89	0,88	0,88	1200
1001	hüzün	0,9	0,89	0,89	204
	korku	0,86	0,99	0,92	189
	mutluluk	0,99	0,79	0,88	194
	negatif	0,8	0,88	0,84	215
	nötr	0,99	0,92	0,95	176
	pozitif	0,8	0,8	0,8	222
	macro avg	0,89	0,88	0,88	1200
1010	hüzün	0,9	0,89	0,89	204
	korku	0,86	0,99	0,92	189
	mutluluk	0,99	0,79	0,88	194
	negatif	0,88	0,84	0,84	215
	nötr	0,99	0,92	0,95	176
	pozitif	0,8	0,8	0,8	222
	macro avg	0,89	0,88	0,88	1200
1011	hüzün	0,9	0,89	0,89	204
	korku	0,86	0,99	0,92	189
	mutluluk	0,99	0,79	0,88	194
	negatif	0,8	0,88	0,84	215
	nötr	0,99	0,92	0,95	176
	pozitif	0,8	0,8	0,8	222
	macro avg	0,89	0,88	0,88	1200
1100	hüzün	0,93	0,86	0,89	197
	korku	0,84	1	0,91	189
	mutluluk	0,99	0,77	0,87	205
	negatif	0,84	0,87	0,85	215
	nötr	0,99	0,9	0,94	186
	pozitif	0,75	0,88	0,81	208
	macro avg	0,89	0,88	0,88	1200
1101	hüzün	0,93	0,86	0,89	197
	korku	0,84	1	0,91	189
	mutluluk	0,99	0,77	0,87	205
	negatif	0,84	0,87	0,85	215
	nötr	0,99	0,9	0,94	186
	pozitif	0,75	0,88	0,81	208
	macro avg	0,89	0,88	0,88	1200
1110	hüzün	0,93	0,86	0,89	197
	korku	0,84	1	0,91	189
	mutluluk	0,99	0,77	0,87	205
	negatif	0,84	0,87	0,85	215
	nötr	0,99	0,9	0,94	186
	pozitif	0,75	0,88	0,81	208
	macro avg	0,89	0,88	0,88	1200
1111	hüzün	0,93	0,86	0,89	197
	korku	0,84	1	0,91	189
	mutluluk	0,99	0,77	0,87	205
	negatif	0,84	0,87	0,85	215
	nötr	0,99	0,9	0,94	186
	pozitif	0,75	0,88	0,81	208
	macro avg	0,89	0,88	0,88	1200

Table 4 NBM

4.3. Which models did we use?

Naive Bayes: Naive Bayes classifier is a probabilistic model used to predict emotional labels of texts. It assumes that the words in the texts are independent and calculates the probabilities of these words between classes.

Support Vector Machines (SVM): Support Vector Machines are a widely used classification method for classifying texts. By converting texts to vector representations, it creates a decision boundary to predict belonging to a particular emotional class.

Decision Trees: Decision trees use a tree structure to predict the emotional classes of texts. Each node represents a division by a particular property and its values. Using the features of the texts, a series of decision structures are created and emotional classes are determined at the end.

K-Nearest Neighbor (KNN):

KNN is a simple and popular classification and regression algorithm.

In the case of classification, it uses the k nearest neighbors whose labels are known to classify a data point.

In the case of regression, it uses the mean of the k nearest neighbors to estimate the output of a data point.

KNN is mainly based on distance measurements of data points (usually the Euclidean distance is used).

The user-specified parameter k determines how many neighbors are to be considered.

KNN is a widely used algorithm with its simple and understandable structure.

Logistic Regression:

Logistic regression is a linear regression algorithm used in classification problems.

It is used to separate data points into two or more classes.

Logistic regression attempts to classify data points according to a linear decision boundary.

Logistic regression uses the logistic function (sigmoid function) to estimate the output.

Logistic regression generates the probability values of the predictions and classifies by setting a cutoff threshold.

Logistic regression is a widely used algorithm because of its simplicity, interpretability and speed.

Random Forest:

A random forest is an ensemble (combination) algorithm created by combining many decision trees.

Each tree is trained by random sampling and random feature selection.

The random forest takes an estimate of each tree and determines the result by majority voting or average value.

Random forest is resistant to overfitting and generally provides high performance.

Random forest can be used in both classification and regression problems.

Random forest is a computationally efficient algorithm that can run in parallel on large datasets.

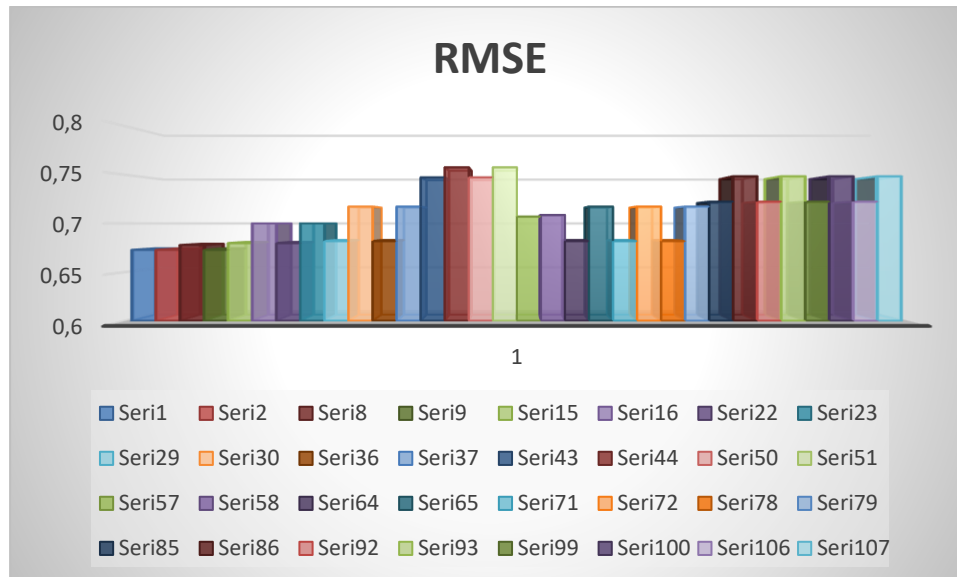


Table 5 svm,lr seri1 and seri2 = 0000, seri8 and seri9 =0001 , seri15,16 =0010 , seri22,23=0011, seri29,30 =0100, seri36,37=0101, seri43,44=0110, seri50,51 =0111, seri57,58 =1000, seri64,65=1001, seri71,72=1010, seri78,79 =1011, seri85,86= 1100, seri92,93 =1101, seri99,100 =1110, seri106,107=1111

4.3.1. Methods Used

Machine Learning Based Classification Models: These approaches use machine learning algorithms to identify emotional labels of texts. A classification model is trained using a pre-labeled dataset and then predictions are made on new texts. For example, algorithms such as Support Vector Machines (SVM), Decision trees, Logistic Regression, Naive Bayes can be used.

Different from Figure 1, it was evaluated on 6 titles. These are ‘pozitif’, ‘negatif’, ‘nötr’, ‘hüzün’ ‘korku’ , ‘mutluluk’.

Models used: KNN, Naive Bayes Multinomial, Random Forest, Decision Table, SVM, LR.
Models used on the website, BERTurk, SVM, LR, NBM.
TfidfVectorizer was used for vectorization.

4.3.2 Applications and Libraries Used

Data extraction: beautifulsoup,selenium /jupyter notebook
Preprocessing: nltk.download(),WPT = nltk.WordPunctTokenizer() stop_word_list
nltk.corpus.stopwords.words('turkish'),nltk.re,snowballstemmer import
TurkishStemmer/Google Colab
Vectorization:pickle,pandas,sklearn.feature_extraction.text,TfidfVectorizer, sklearn.svm
import SVC,sklearn.model_selection import train_test_split/Google Colab
Site: flask import Flask, render_template, request,HTML/ Spyerdar

4.4. What performance metrics did we use?

F1-Score (F1-Score): F1-Score is a criterion used to evaluate the performance of classification models. It is the harmonic mean of the precision and recall metrics. The F1-Score tends to minimize both false positives and false negatives of a model. That is, it evaluates the overall performance of the model, taking into account both precision and recall.

Precision: Precision is the rate at which samples that a classification model predicts positively are actually positive. That is, it shows how accurately a model predicts true positives. Precision is important in applications with the goal of minimizing false positives. For example, it is important to have as few false positives as possible in spam filtering applications.

RMSE (Root Mean Square Error): RMSE is an error metric used to measure how far the predictions of a regression model are from the true values. The RMSE is the square root of the mean of the squares of the prediction errors. Smaller RMSE values indicate that the model's predictions are closer to the true values.

Recall: Recall is a metric that shows how accurately a classification model predicts true positives. Recall is important in applications with the goal of minimizing false negatives. For example, it is important to have as few false negatives as possible in disease diagnosis practices.

These metrics are used to measure and evaluate the performance of a model. When assessing the success of the model, it is important to consider the class balance, the characteristics of the dataset, and the intended application. The combination of metrics such as F1-Score, precision, RMSE and recall helps to comprehensively evaluate the performance of the model.

	<i>Emation</i>
1	Hüzün
2	Korku
3	Mutluluk
4	Negatif
5	Nötr
6	Pozitif

Table 6 Table3,4,5 evaluate accordingly.

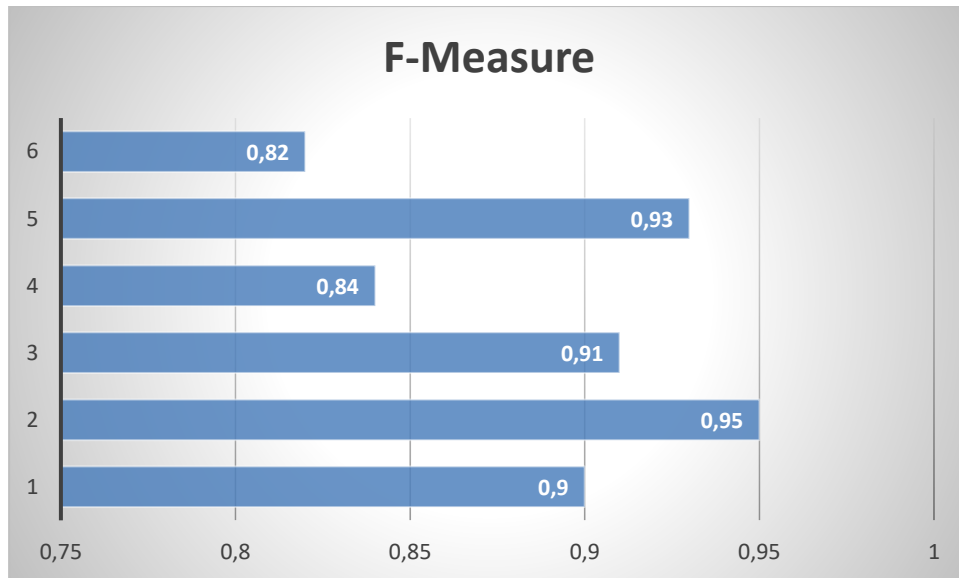


Table 7 This is 0000

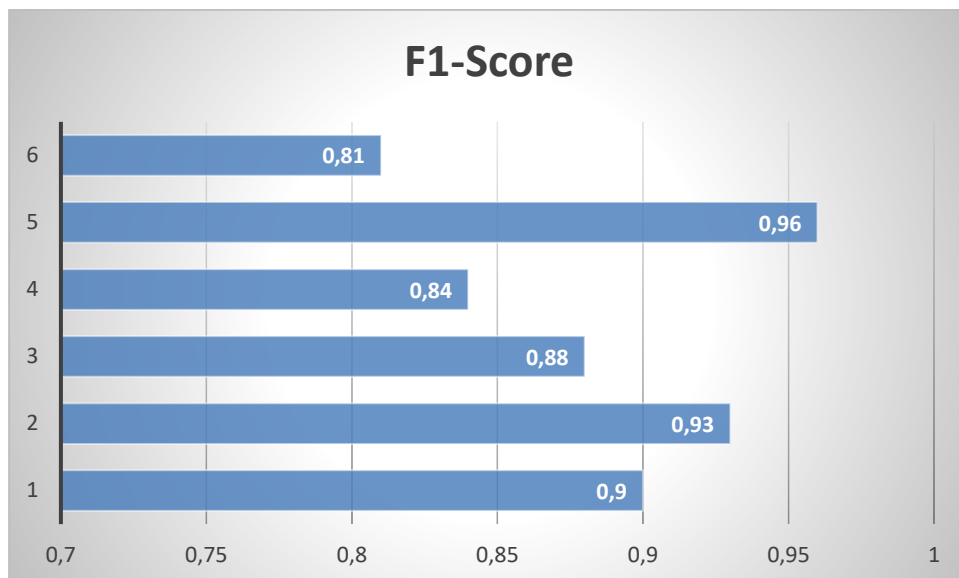


Table 8 This is 0000

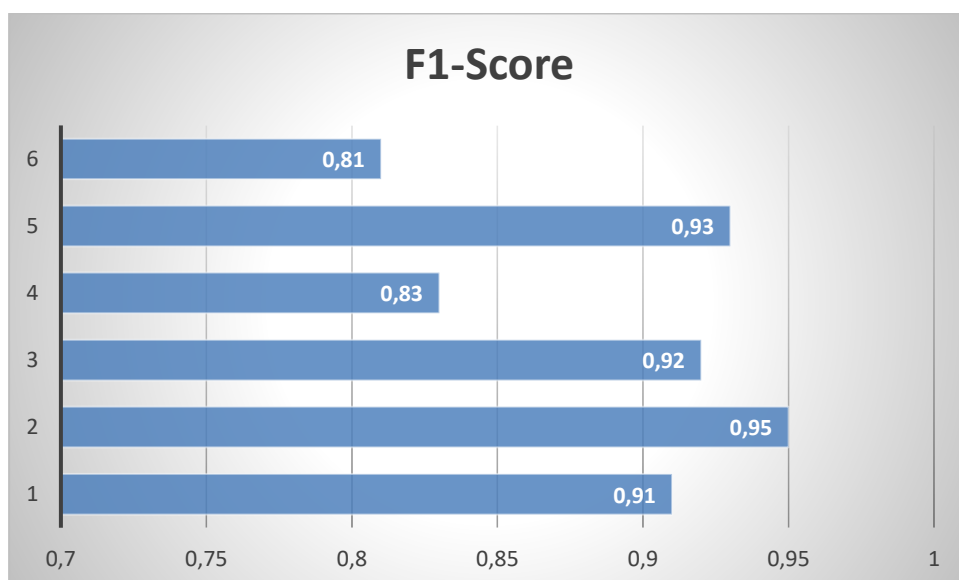


Table 9 This is 0000

4.5. What is Bert?

BERT (Bidirectional Encoder Representations from Transformers) is a deep learning based language model developed by Google. BERT has had great success in natural language processing (NLP) and has delivered best-in-class results on many NLP tasks.

BERT uses the Transformer architecture and is pre-trained as a general language model for many language processing tasks. Using large amounts of text data, the model learns word-

level relationships and grammatical patterns. The most important feature of BERT is that it can process text in both directions (forward and backward). This allows him to better understand the meaning of the word and the context in which it is used.

BERTurk is a pre-trained BERT model for use in the Turkish language. It is a BERT model trained with Turkish dataset to understand grammatical patterns and relationships in Turkish texts. BERTurk provides high performance in Turkish language by using sentiment analysis, text classification, word similarity and other NLP tasks.

BERT and BERTurk are very effective and successful models in the field of natural language processing. They stand out for their general language comprehension and performance in text processing tasks.

The use of the BERTurk model is a very long study. On the other hand, this model could not be included because we could not get good results from our own data set.

F1-Score	0.1019623060944757
Recall	0.13240059540447713
Precision	0.1692532266627289

Table 10 BERTurk

	negatif	nötr	pozitif	Hüzün	Korku	mutluluk
precision	0.272727	0.009804	0.375000	0.250000	0.0	0.107988
recall	0.338462	0.010309	0.015306	0.065327	0.0	0.365000
f1-score	0.302059	0.010050	0.029412	0.103586	0.0	0.166667
support	195.000000	194.000000	196.000000	199.000000	0.0	200.000000

Table 11 BERTurk

	accuracy	macro avg	weighted avg
precision	0.132713	0.169253	0.169004
recall	0.132713	0.132401	0.132713
f1-score	0.132713	0.101962	0.101913
support	0.132713	1183.000000	1183.000000

Table 12 BERTurk

5. SITE FORMATION

The site is connected with the HTML method, which is a basic structure. The spider ide is used for the site creation. The analysis file has been created and links with HTML templates have been added. The site creation phase has been completed using flash in Python. The studies were carried out by embedding the files vectorized with the TDFIDF method into the site.

Duygu Analizi

Metin girin

Analiz Et

Figure 3 WEBSITE

Duygu Analizi

Bugün hava çok güzeldi.Ama içim de bir hüznün var.

Analiz Et

Duygu: {'negatif': 6.36, 'pozitif': 12.89, 'nötr': 27.21, 'korku': 0.0, 'mutluluk': 19.47, 'hüzün': 34.06}

Figure 4 WEBSITE

6. CONCLUSION AND FUTURE WORK

Our study is based on sentiment analysis and its transfer to a site. Our study, which consists of 6000 data in total, is divided into 4800 train and 1200 test data. There are 1000 datasets for each emotion and they are 800 train, 200 test datasets. These are added to the machine learning models by going through certain preprocessing steps. Work started with 6 machine learning models. These are: K - Nearest Neighbor, Decision Tree, Logistic Regression, Random Forest, Support Vector Machine, Naive Bayes Multinomial. Algorithms with the highest results will be selected. The algorithms with the highest results will be compared with the Bert model. It should not be overlooked that the pre-processing stage is also of great importance. 16 different combinations were applied. These are sub, punctuation, single word, root word combinations. The best results were obtained with the 1110. So alt+punctuation+keyword.

In this study, Support Vector Machine (SVM), Naive Bayes (Multinomial) and Logistic Regression models were selected as machine learning models and compared with the BERT model. The aim of the study is to evaluate the performance of these models in the sentiment analysis task and to obtain the best result.

According to the results obtained, although the SVM model has a high accuracy rate, it worked slower than other models in terms of computation time. The Naive Bayes (Multinomial) model ran fast, but the accuracy rate was slightly lower than the SVM model. The Logistic Regression model, on the other hand, provided a good balance in terms of both computation time and accuracy. The preference of the SVM model for Model Selection is that the RMSE values are higher than the others.

On the other hand, although the BERT model is more complex than other machine learning models, it was expected to have the highest accuracy, but this could not be achieved. This is surprising. Because BERT is a tool that can better understand grammatical patterns and better represent texts. Thus, while it should outperform text-based tasks such as sentiment analysis, it yielded the opposite result. The BERT model has a longer computation time than other models, and this is a factor to consider when working with large datasets. In my own work, the model took 13 hours to train and used most of the computer's CPU.

As a result, in our study, when SVM, Naive Bayes (Multinomial) and Logistic Regression models were compared with the BERT model, it was seen that the BERT model did not

provide the highest accuracy rate and was not successful in sentiment analysis. In addition, considering the computation time and complexity of the model, it seems risky to choose this model in the context of the application. The high F1-Score and RMSE values of the SVM model also enabled this model to be selected. The SVM model and vectors are embedded in the 'Sense Analysis' website.

It cannot be said that our business is completely trouble-free. But it should not be overlooked that there are good lessons in their mistakes. Some notes should be taken from the study. If the size of our dataset was sufficient, artificial neural networks and deep learning would be preferred, which are the most successful in language learning. Even if machine learning is reintroduced, a few changes need to be made. Since the biggest problem in this study is experienced in the data preprocessing phase, ZEMBEREK library should be included for root analysis. For the frequency range, Word2Vect or Doc2Vec can be used. In this way, more successful results can be obtained. On the other hand, easier tagging studies can be done by considering emojis.

It should be known that working with Turkish data is more difficult than languages such as English. It should be known that it complicates our work because it is an agglutinative language. For this reason, it is of great importance that every work done is successful and unsuccessful. Every unsuccessful work will actually lead us to the truth. Further work in the future will take us further in Turkish language processing.

7. REFERENCES

- [1] TÜRKÇE HİS ANALİZİNDE OPTİMİZASYON VE ÖNEĞİTİMLİ MODELLERİN KULLANIMI. Alaettin UÇAN Aralık 2020
- [2] ÇUKUROVA ÜNİVERSİTESİ, TÜRKİYE / FEN BİLİMLERİ ENSTİTÜSÜ ELEKTRİK-ELEKTRONİK MÜHENDİSLİĞİ (DR) (İNGİLİZCE) / (2011-2015) Mümine Kaya Keleş
- [3] S. Aman, S. Szpakowicz, Identifying expressions of emotion in text, in: Int. Conf. Text, Speech Dialogue, 2007: pp. 196–205.
- [4] M.A. Tocoglu, A. Alpkocak, TREMO: A dataset for emotion analysis in Turkish, J. Inf. Sci. 44 (2018) 848–860.
- [5] M.A. Tocoglu, A. Alpkocak, Lexicon-based emotion analysis in Turkish, Turkish J. Electr. Eng. Comput. Sci. 27 (2019) 1213–1227.
- [6] S. Aman, S. Szpakowicz, Using roget’s thesaurus for fine-grained emotion recognition, in: Proc. Third Int. Jt. Conf. Nat. Lang. Process. Vol., 2008.
- [7] D. Ghazi, D. Inkpen, S. Szpakowicz, Hierarchical versus flat classification of emotions in text, in: Proc. NAACL HLT 2010 Work. Comput. Approaches to Anal. Gener. Emot. Text, 2010: pp. 140–146.
- [8] S. Mac Kim, A. Valitutti, R. a Calvo, Evaluation of Unsupervised Emotion Models to Textual Affect Recognition, Proc. NAACL HLT 2010 Work. Comput. Approaches to Anal. Gener. Emot. Text. (2010) 62–70.
- [9] S. Chaffar, D. Inkpen, Using a heterogeneous dataset for emotion analysis in text, in: Can. Conf. Artif. Intell., 2011: pp. 62–67.
- [10] C.O. Alm, D. Roth, R. Sproat, Emotions from text: machine learning for text-based emotion prediction, Proc. Hum. Lang. Technol. Conf. Conf. Empir. Methods Nat. Lang. Process. (2005) 579–586. doi:10.3115/1220575.1220648.
- [11] T. Danisman, A. Alpkocak, Feeler: Emotion Classification of Text Using Vector Space Model, Aisb. (2008) 53–59.
- [12] . Naderalvojoud, A. Ucan, E. Akcapinar Sezer, HUMIR at IEST-2018: Lexiconsensitive and left-right context-sensitive bi-lstm for implicit emotion recognition, in: Proc. 9th Work. Comput. Approaches to Subj. Sentim. Soc. Media Anal., Association for Computational Linguistics, 2018: pp. 182–188.
- [13] M. Abdul-Mageed, L. Ungar, EmoNet: Fine-Grained Emotion Detection with Gated Recurrent Neural Networks, Proc. 55th Annu. Meet. Assoc. for Comput. Linguist. (Volume 1 Long Pap. (2017) 718–728. doi:10.18653/v1/P17-1067.

- [14] S. Ge, T. Qi, C. Wu, Y. Huang, THU_NGN at SemEval-2019 Task 3: Dialog Emotion Classification using Attentional LSTM-CNN, in: Proc. 13th Int. Work. Semant. Eval., 2019: pp. 340–344.
- [15] L. Ma, L. Zhang, W. Ye, W. Hu, PKUSE at SemEval-2019 task 3: emotion detection with emotion-oriented neural attention network, in: Proc. 13th Int. Work. Semant. Eval., 2019: pp. 287–291.
- [16] . Boynukalin, P. Karagoz, Emotion analysis on Turkish texts, in: Inf. Sci. Syst. 2013, Springer, 2013: pp. 159–168.
- [17] M.A. Tocoglu, A. Alpkocak, Emotion extraction from turkish text, Proc. - 2014 Eur. Netw. Intell. Conf. ENIC 2014. (2014) 130–133. doi:10.1109/ENIC.2014.17.
- [18] M.A. Tocoglu, O. Ozturkmenoglu, A. Alpkocak, Emotion analysis from Turkish tweets using deep neural networks, IEEE Access. 7 (2019) 183061–183069. doi:10.1109/ACCESS.2019.2960113.
- [19] H.G. Wallbott, K.R. Scherer, How universal and specific is emotional experience? Evidence from 27 countries on five continents, Inf. (International Soc. Sci. Counc. 25 (1986) 763–795.
- [20] S. Demirci, Emotion analysis on Turkish tweets, Middle East Technical University, 2014.
- [21] Revision Date: 20.06.2022 An alternative word embedding approach for knowledge representation in online consumers' reviews Ekin EKİNCİ^{1*}, Sevinç İLHAN OMURCA²
- [22] Natural language processing methods for knowledge management—Applying document clustering for fast search and grouping of engineering documents Ivar O’rn Arnarsson¹ , Otto Frost², Emil Gustavsson², Mats Jirstrand² and Johan Malmqvist¹
- [23]Instructional Technology and Lifelong Learning Vol. 3, Issue 2, 129-143 (2022)Modeling Education Studies Indexed in Web of Science Using Natural Language Processing Tuncer AKBAY^{*1}
- [24] Revealing the Reflections of the Pandemic by Investigating COVID-19 Related News Articles Using Machine Learning and Network Analysis Ulya BAYRAM
(Geliş/Received:11.06.2021; Kabul/Accepted:12.04.2022)
- [25](Arrived/Received: 01.03.2022; Accepted/Accepted in Revised Form: 16.11.2022)
ACADEMIC TEXT CLUSTERING USING NATURAL LANGUAGE PROCESSING
1Salimkan Fatma TAŞKIRAN, 2Ersin KAYA