



**ÇUKUROVA UNIVERSITY FACULTY OF ENGINEERING – COMPUTER
ENGINEERING**

COURSE: INTRODUCTION TO DATA MINING

**"Dimensionality Reduction, Feature Selection, and
Classification with Sentiment Labelled Dataset"**

2020555003-ÖMER AKDOĞAN

2020555046-AHMET ÖZALP

2020555060-ALİ TEPELİ

2020555071-EMİRHAN YILMAZOĞLU

2020555072-SEMİH ZENGİNOĞLU

DELIVERED DATE : 02.01.2025

ADVISER = PROF. DR. SELMA AYŞE ÖZEL

INTRODUCTION :

Sentiment analysis is a significant topic in natural language processing and machine learning, aiming to understand and classify emotional expressions in text. In this study, the Sentiment Labelled Sentences Dataset is utilized to examine the impact of different feature selection and dimensionality reduction methods on sentiment classification.

Text data is vectorized using the TF-IDF (Term Frequency-Inverse Document Frequency) method with both unigram and bigram representations. To reduce high dimensionality, methods such as Chi², PCA, Lasso (L1 Regularization), and threshold-based feature selection are applied. Finally, the processed datasets are classified using Naïve Bayes, Decision Tree, and SVM models. The performance of the reduced datasets is compared to that of the original, unreduced datasets.

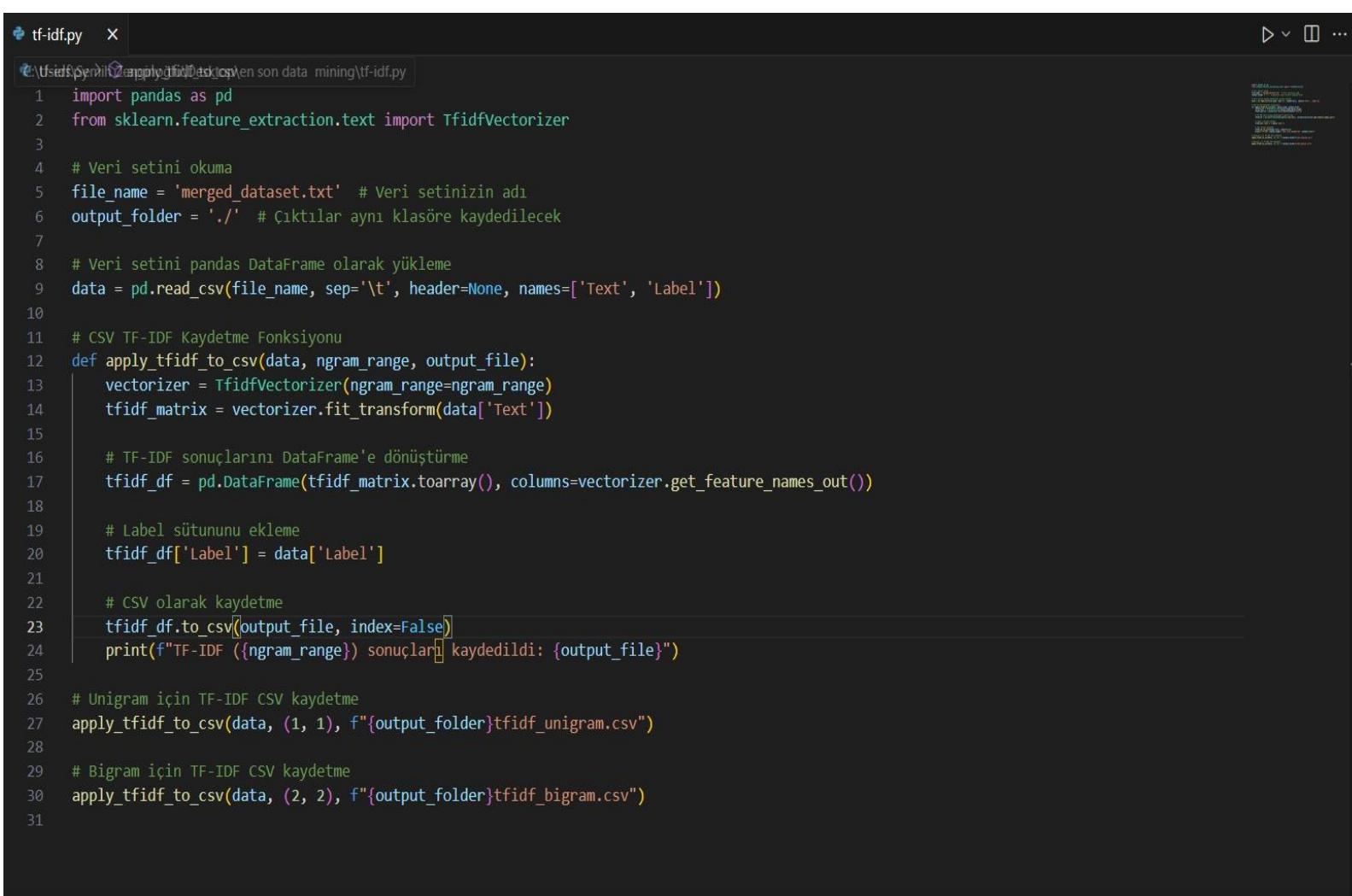
LINK TO OUR DATASET :

<https://archive.ics.uci.edu/dataset/331/sentiment+labelled+sentences>

STEP 1- TF*IDF VECTORIZATION :

First, we combined three datasets into a single comprehensive dataset. Next, we applied unigram and bigram TF-IDF vectorization separately to this dataset. Afterward, we converted the datasets with unigram and bigram applied into Excel tables. The purpose of this step was to enable the use of dimensionality reduction methods, as these methods require the dataset to be numerically represented.

UNIGRAM - BIGRAM TF*IDF PYTHON CODE



```
tf-idf.py x
C:\Users\yusuf\PycharmProjects\deneme\deneme\son data mining\tf-idf.py
1 import pandas as pd
2 from sklearn.feature_extraction.text import TfidfVectorizer
3
4 # Veri setini okuma
5 file_name = 'merged_dataset.txt' # Veri setinin adı
6 output_folder = './' # Çıktılar aynı klasöre kaydedilecek
7
8 # Veri setini pandas DataFrame olarak yükleme
9 data = pd.read_csv(file_name, sep='\t', header=None, names=['Text', 'Label'])
10
11 # CSV TF-IDF Kaydetme Fonksiyonu
12 def apply_tfidf_to_csv(data, ngram_range, output_file):
13     vectorizer = TfidfVectorizer(ngram_range=ngram_range)
14     tfidf_matrix = vectorizer.fit_transform(data['Text'])
15
16     # TF-IDF sonuçlarını DataFrame'e dönüştürme
17     tfidf_df = pd.DataFrame(tfidf_matrix.toarray(), columns=vectorizer.get_feature_names_out())
18
19     # Label sütununu ekleme
20     tfidf_df['Label'] = data['Label']
21
22     # CSV olarak kaydetme
23     tfidf_df.to_csv(output_file, index=False)
24     print(f"TF-IDF ({ngram_range}) sonuçları {output_file} kaydedildi: {output_file}")
25
26 # Unigram için TF-IDF CSV kaydetme
27 apply_tfidf_to_csv(data, (1, 1), f"{output_folder}tfidf_unigram.csv")
28
29 # Bigram için TF-IDF CSV kaydetme
30 apply_tfidf_to_csv(data, (2, 2), f"{output_folder}tfidf_bigram.csv")
31
```

RESULTS FOR TF*IDF VECTORIZATION :

DATASET WITH UNIGRAM TF*IDF APPLIED . The number of unique unigrams (features) obtained: 5155.

	GN	GO	GP	GQ	GR	GS	GT	GU	GV	GW	GX	GY	GZ	HA	HB	HC	HD	HE	HF	HG	HH	HI	HJ	
1	array	arrival	arrived	arrives	arriving	art	article	articulated	artiness	artist	artistic	artless	arts	as	asia	aside	ask	asked	asking	asleep	aspect	aspects	ass	as
2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
4	0	0	0	0	0	0	0	0	0,266548	0	0	0	0	0,130301	0	0	0	0	0	0	0	0	0	
5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
6	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
7	0	0	0	0	0	0,26827	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
9	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0,212782	0	0	0	0	0	0	0	0	0	
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
17	0	0	0	0	0	0	0	0	0	0	0	0	0	0,180929	0	0	0	0	0	0	0	0	0	
18	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0,244548	0	0	0	
20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
21	0	0	0	0	0	0	0	0	0,029304	0	0	0	0	0,060271	0	0	0	0,025187	0	0	0	0	0	
22	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
23	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
24	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
25	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
26	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
27	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
28	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
29	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
30	0	0	0	0	0	0	0	0	0	0	0	0	0	0,358769	0	0	0	0	0	Windows'u etkinleştirmek için Aşağıda girin.	0	0		
31	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	Windows'u etkinleştirmek için Aşağıda girin.	0	0	

DATASET WITH BIGRAM TF*IDF APPLIED . The number of unique bigrams (features) obtained: 20,402.

STEP 2- REDUCTION :

At this step , we applied four (chi_square, pca , l1 _regularization, variance threshold) dimensionality reduction methods to the vectorized dataset to reduce its size.

Chi_square reduction method :

The value of 'k' indicates that we will select the top 100 features with the highest Chi-square scores.

Code :

```
❷ chi-square.py > ...
1  import pandas as pd
2  from sklearn.feature_selection import SelectKBest, chi2
3  import os
4
5  def apply_chi_square(data, target_col, k, output_file):
6      # Hedef klasörün varlığını kontrol et ve yoksa oluştur
7      output_dir = os.path.dirname(output_file)
8      if not os.path.exists(output_dir):
9          os.makedirs(output_dir)
10
11     X = data.drop(columns=[target_col])
12     y = data[target_col]
13
14     chi_selector = SelectKBest(score_func=chi2, k=k)
15     X_new = chi_selector.fit_transform(X, y)
16
17     selected_features = pd.DataFrame(X_new, columns=data.columns[chi_selector.get_support(indices=True)])
18     selected_features[target_col] = y
19
20     selected_features.to_csv(output_file, index=False)
21     print(f"Chi-square ile seçilen özellikler kaydedildi: {output_file}")
22
23     # Veri dosyaları ve klasörler
24     tfidf_files = [
25         ('tfidf_unigram.csv', 'unigram_results/chi_square.csv'),
26         ('tfidf_bigram.csv', 'bigram_results/chi_square.csv')
27     ]
28
29     # Her dosya için Chi-Square uygula ve kaydet
30     for input_file, output_file in tfidf_files:
31         data = pd.read_csv(input_file)
32         apply_chi_square(data, 'Label', k=100, output_file=output_file)
33
```

The table after applying Chi-Square to the dataset with unigram TF-IDF

The table after applying Chi-Square to the dataset with bigram TF-IDF

PCA reduction method :

n is indicates the number that specifies selecting the top 50 principal components with the highest eigenvalues.

Code :

```
pca.py > ...
pca.py > ...
1 import pandas as pd
2 from sklearn.decomposition import PCA
3 import os
4
5 def apply_pca(data, target_col, n_components, output_file):
6     # Hedef klasörün varlığını kontrol et ve yoksa oluştur
7     output_dir = os.path.dirname(output_file)
8     if not os.path.exists(output_dir):
9         os.makedirs(output_dir)
10
11     X = data.drop(columns=[target_col])
12     y = data[target_col]
13
14     pca = PCA(n_components=n_components)
15     X_pca = pca.fit_transform(X)
16
17     pca_df = pd.DataFrame(X_pca, columns=[f"PC{i+1}" for i in range(n_components)])
18     pca_df[target_col] = y
19
20     pca_df.to_csv(output_file, index=False)
21     print(f"PCA ile boyutu azaltan veri kaydedildi: {output_file}")
22
23 # Veri dosyaları ve klasörler
24 tfidf_files = [
25     ('tfidf_unigram.csv', 'unigram_results/pca.csv'),
26     ('tfidf_bigram.csv', 'bigram_results/pca.csv')
27 ]
28
29 # Her dosya için PCA uygula ve kaydet
30 for input_file, output_file in tfidf_files:
31     data = pd.read_csv(input_file)
32     apply_pca(data, 'Label', n_components=50, output_file=output_file)
33
```

The table after applying PCA to the dataset with unigram TF-IDF

	PC1, PC2, PC3, PC4, PC5, PC6, PC7, PC8, PC9, PC10, PC11, PC12, PC13, PC14, PC15, PC16, PC17, PC18, PC19, PC20, PC21, PC22, PC23, PC24, PC25, PC26, PC27, PC28, PC29, PC30, PC31, PC32, PC33
1	-0.0044723504809777445, -0.09983232123036741, 0.042160627747927614, 0.08745840539260281, -0.016549152336594136, 0.005056726435301946, 0.11191798284927944, 0.102
2	0.0031741106873912957, 0.09094890485824063, -0.08276061944885453, -0.1911621493927926, 0.15776693858779212, -0.052561373888043954, 0.17309340736601916, -0.03577
3	0.151696520146985, 0.2858815547042354, 0.03752369338172175, 0.08861411104099742, -0.1202389688279117, 0.008147333403210739, 0.10794007091365912, -0.01049832168
4	0.03234794170433043, -0.04121827329797374, -0.03839746528100324, -0.022749073670918966, -0.06529721564281213, -0.02973454506299976, 0.014733830815610981, 0.018
5	0.1848863877325936, 0.3110314208131154, 0.1837503834848383, 0.10453918644108331, -0.008247912039692567, -0.01652441585135748, 0.11463234558348516, 0.0487063
6	-0.0066116315127232865, -0.1292842973972959, 0.03019225489748176, 0.1131173602988854, -0.0612093046573945, -0.0358595332959084, 0.058317091382062806, 0.07603
7	-0.005779752072671274, 0.0606383949411083666, 0.0689170531352892, -0.004648623316527874, -0.09509696221141052, 0.04520072242134111, 0.044728324191756895, 0.00805057655906345, 0.0786
8	0.11524532685593425, -0.11713746004216187, 0.032781102154929645, -0.06539718477304948, -0.04520072242134111, 0.044728324191756895, 0.00805057655906345, 0.0786
9	0.02847840395052757, -0.073487519442646996, -0.03276519021931791, 0.088331652187291929936, -0.020303684014672384, -0.02
10	-0.016530482717700735, -0.055658323718868935, -0.010094481966435833, -0.02683096788710705, -0.022244845881521054, -0.03267843677471044, -0.03267843677471044, -0.02
11	0.14881252790944327, -0.367291315329124, -0.22466476915999013, 0.1044371249944765, 0.03567077232143224, -0.0992555332875087, 0.14278260561092948, 0.12370951770
12	-0.13686862883467418, 0.057156459038546194, 0.0199151054638715, -0.03407186269836006, 0.00216735945331817, -0.043458235304703975, -0.08201012001591212, -0.05
13	-0.0474014713073093, 0.0044829295621025, 0.060363364469744, -0.006943165688326713, -0.03254531658316475, -0.03254531658316475, -0.03254531658316475, -0.0047
14	-0.09161519469492036, 0.13998705499647712, -0.12693117953266786, -0.3069611589369399, 0.3367536023214012, -0.23533488637751368, 0.05376642584915599, -0.0512413
15	-0.0977141786562045, 0.05478823388038281, 0.1667418611797386, -0.10447639066108683, 0.18585110469757272, -0.1851785408834284, -0.012351680051181306, -0.089689
16	0.12961929873265893, -0.0228141465067822, -0.011251094937046407, 0.001040580248668456, -0.02584506998373993, 0.012659749256022026, -0.0447452911785382, -0.08
17	0.041608676874179316, -0.09392433263773816, -0.10723001153904803, -0.10231829072600293, -0.04616652346421397, 0.04497125656228304, 0.09123015930148552, -0.0777
18	-0.06064146148133899, -0.023450093095948977, -0.116037987768316, -0.09366500342077258, 0.157244450384492, 0.0044904144977940687, 0.1957534236570428, 0.0015964241
19	0.41059834660563427, 0.48959457745825407, -0.14306448129415575, 0.18149594571226668, -0.18686243861207139, -0.05079155630367532, 0.0785466719681708, 0.0370778
20	0.05202321394159326, -0.1204319313016385, -0.02013798721910007, 0.05687547036744453, 0.023682736793544993, 0.044245180979274275, 0.0026597605632076708, 0.0462
21	0.004283309719807825, -0.105234529141182, -0.0964198820502053, -0.036386294866706, -0.0206529403739601826, -0.000936795739601826, 0.0658464808241536, -0.04580167
22	0.05479599854781587, -0.09437841364838971, -0.050612850376564614, 0.10772816047321898, 0.043735728497671746, -0.02206986818590959, -0.027989300557546, -0.0085
23	-0.0953988747253539, 0.0366912403012499, 0.13298808715291485, -0.01096594179309866, -0.0928624395380795, -0.032278853958702155, 0.054004758905158115, -0.0538
24	0.0568445012741671, -0.01585710246650394, 0.02107964930179281, 0.0009430108412159309, -0.023613962458311974, 0.05259914084832296, -0.04381475216084843, 0.01
25	-0.04660939386731194, -0.01593094712299063, -0.016435400260058745, -0.1276214451216442, -0.016927120941706452, -0.044803052824765, -0.0022069291935080484, -0.0
26	0.2215022356137919, 0.20945691458683632, -0.05696234986230923, -0.0499598355250163, -0.02629509067383067, -0.018923691773045655, -0.021336334679053916, -0.0991
27	0.0972061923840264, 0.0020530278435959503, 0.05142963178984104, 0.1419111288191913562, -0.018822759748823397, 0.023370561766652975, -0.056108270385345264, -0.120
28	-0.049346312039030005, -0.03386813632898175, 0.04228204811796534, 0.01120326168530481, -0.09415095929547654, -0.01763699771251844, 0.028067139227597935, 0.006
29	-0.03147389513690063, 0.025837172657820897, 0.0630396070748358, 0.0714618670082038, -0.10254175421650691, -0.03644420639663023, 0.012055405290931936, -0.05477
30	0.027349380646141324, -0.021418619884028926, -0.04574201564351659, -0.03700274717522683, -0.046742330490172214, -0.051436365991649295, -0.01200463261431771, -0.0
31	0.12794319055427403, -0.01595537274559991, 0.2027868072231544, -0.044715011377763625, 0.09931696468243256, -0.019170113433925277, -0.008923117523673686, 0.015249
32	0.08593678746652272, -0.00933085125298589, 0.04861974066886629, -0.10252462792513815, 0.11577062447925898, -0.07460725057606411, -0.06362006533238991, 0.0252141
33	0.026012331834885088, -0.04530250243156589, 0.011527299131486373, 0.06593153967621249, 0.05702802731862344, -0.05499099659596979, 0.02450285312817014, 0.039473
34	0.16100123204219394, -0.03255503608675444, 0.07533468654746298, -0.0005482477446741627, 0.038758468513489105, 0.09895458156320739, -0.13844605524951945, -0.094
35	0.011664308799538576, -0.0674398443458514, -0.08869188851145768, -0.019867957163625034, -0.07942255957017389, -0.05519453079526104, 0.04458270811410894, -0.00
36	0.06307186651406733, -0.024492730105212548, 0.11914987911835581, 0.047531562138265274, -0.08480443599922094, -0.03688221466606005, 0.048657754815737427, -0.043
37	-0.0288495034028842, -0.027768755687018663, -0.13153846537934635, 0.05052808420366883, 0.16972964244921757, -0.26787200762110375, -0.14143465552922987, 0.005209
38	-0.021135604107012317, -0.0223204125512328604, -0.05010464827854050, -0.08172678257852377, -0.053708166604004006, -0.06828732241652443, -0.0247811661004401

The table after applying PCA to the dataset with bigram TF-IDF

```
bigram_results > pca.csv
1 [PC1,PC2,PC3,PC4,PC5,PC6,PC7,PC8,PC9,PC10,PC11,PC12,PC13,PC14,PC15,PC16,PC17,PC18,PC19,PC20,PC21,PC22,PC23,PC24,PC25,PC26,PC27,PC28,PC29,PC30,PC31,PC32,PC33]
2 0.005254820801315349,-0.0055101887189825375,-0.018154884419891378,0.0016679282727952635,-0.01024168309104478,-0.006870014576673043,-0.014130914607036367
3 0.0004196959972610621,0.005479707930836288,-0.009290969376485863,-0.003257390554233253,-0.011242914982681448,-0.004500605563344882,-0.00709104820435705
4 -0.009168918321842452,-0.017474807117252283,-0.004784172614712596,0.0027342074491334824,-0.004549924994029902,-0.01071692518470274,-0.022058669643172345,
5 0.00023864190956300374,-0.006061530887588352,-0.007718697252588374,-0.002346850499024322,-0.0117491166694015,-0.004591278917633504,-0.00752078581255636
6 -0.007731572685383603,0.08393002557429045,-0.009370628532221764,-0.077808888643326,0.115890677773326,0.0444901385769344,0.05246801999086373,0.033955
7 -0.004788113565301059,-0.006788984979333621,0.009950440368602001,-0.0023552613551863204,-0.004910912026889002,-0.00651066486265194
8 0.005037738353464614,-0.003943002520136171,-0.011234718117558476,-0.005574311454997458,-0.011075084825709496,-0.003949463156296442,-0.007828421537287307
9 0.014611561361151234,0.001416465782402019,-0.011741805173012286,-0.00691789985085724,-0.012090875080257737,-0.0024208135676434006,-0.007984172883888432
10 -0.00233207851360595,-0.008462780040586397,-0.003031943876284621,0.0063178506405536,-0.0113373273260607,-0.00938320426592867,-0.009789950988001,-0.01
11 -0.002542280958763872,-0.011207682299031383,-0.027721063164974715,0.006915795437034554,-0.03786170118075765,-0.01524266123054545,0.07495749624391937,0
12 -0.022061537732915082,0.061332652626052565,0.0465934113020268,-0.1258726895463193,0.13963239356078674,0.0532629915575084,0.05373925397749205,0.093184
13 -0.012467515214548011,-0.02628216442458532,0.03837974468637076,-0.008743767525912678,-0.00934231546774674,-0.002527336097429426,-0.0076292131691990984,-0
14 -0.001448772312715538,-0.00757391275403316,-0.0109953007530307347,-0.003637881044476387,-0.00990716729377673,-0.003627281790293118,-0.00838822440435479
15 -0.017728313769731874,-0.03265266558020426,0.06734361259041535,-0.010680599794328007,-0.007849527566176986,-0.002444002594253462,-0.005586709384756049,-0
16 3.959759320816179e-05,-0.009171388060031501,-0.00619524287349102,-0.011453723205814205,0.0016491282769978027,-0.0062530053489705615,-0.009700294169793
17 0.011948191069450703,0.000759792203629682,-0.008693605548995339,-0.0007264775372683335,0.01427807901375755,-0.008923248109754076,-0.014947937025796014
18 0.006038242450698447,-0.006329379436133407,-0.0099617799255885939,-0.027759393746278557,-0.018934787192140274,-0.001975863350311587
19 0.0023193520323970965,-0.007640663701684185,-0.014250791924034323,-0.002835605170173885,-0.01686188184975649,-0.006922847010476939,-0.0141911723236139
20 0.006187770155614013,-0.020779470495407876,-0.08413354350976049,-0.06627476055456405,-0.341574531732721,0.822511211479161,0.26440812953809223,-0.05324994
21 -1.5294414350842764e-06,-0.005239931725152428,-0.00908691038487415,-0.0020924565363534126,-0.009682807064852545,-0.004280550529133969,-0.006689517008420
22 0.0024533264203815282,-0.005091527168808789,-0.01044238083943701,-0.0029233962605668973,-0.012571073118379754,-0.00665182674515501,-0.00932744917243215
23 -0.00105738203756834,-0.008486055726714514,-0.017876624952145825,-0.004213470231439749,-0.010425794728236628,-0.0064604243567038057,-0.015714806481792947
24 -0.001529909780908686,-0.005735552319123644,-0.016237637787965192,-0.011655276935863048,-0.004638867986872372,-0.0014015732121463267,-0.0055076851757532
25 -0.0001980908574981255,-0.005499809793023333,-0.009704832533807288,-0.011257365905205806,-0.013016261303799183,0.0048402761103149994,-0.003543735494040
26 -0.002544771815724009,-0.01374025837290704,-0.008977403398869032,-0.005437552096448491,-0.020044489183432164,0.01345706735709247,-0.010223745888361635,-0
27 0.0003846945750736072,-0.0049747479795237155,-0.008364823780941433,-0.0029233962605668973,-0.009966536877080184,-0.0039662444348532765,-0.00624399898825
28 0.012607376619555709,0.00429562865716297,-0.021456850613228957,-0.0036841288456013455,-0.0006857295301845849,0.026261117380375964,0.006227153857732988,0
29 4.6491611863259514e-05,-0.006845738561136572,-0.009982787041169222,-0.00348004285100106,-0.01004210501883723,-0.002423977023417228,-0.0074387739741782
30 -0.0036937090687410305,-0.011144232408566377,-0.01985897677287744,-0.0013355078189526504,-0.01811059852720502,-0.004944565990709763,-0.0155911526992158
31 0.00038479644123941704,-0.00494747921978312406,-0.00836505839570302,-0.002922166643251307,-0.00996615859050353,-0.003965958353467768,-0.00624375197344309
32 -0.023236608277259947,0.03128079192207384,-0.026405213530031137,-0.037637914393503319,0.10913860264674588,0.04002164175540324,0.039286017303965255,0.0237
33 0.015066980224928465,0.0005384003273337963,-0.008446637476578658,0.0009598303580977339,-0.025480563428970024,0.00023373958401440143,-0.01703710674749956
34 0.0003379070192271707,-0.0068418962787927615,-0.02055082254874565,-0.0023418087856640095,-0.011539563845562043,-0.0078450281101944,-0.016679014826753578,
35 0.019056132695416633,-0.0046199764711224115,-0.009625796978291109,0.00886012334003873,-0.028405797001773776,0.0240595996018147,-0.0035762657789531085,0.
36 -0.0012678759359064336,-0.008292045427369046,-0.013218544250251752,-0.00672916075199048,-0.010040699599285274,-0.0021278400663061296,-0.007188442028039
37 -0.005126365785335731,-0.011567179328048542,-0.021623747297505803,-0.003107181142973053,-0.00670081410593358,-0.002558227073275758,-0.0133130549318362
38 -0.000801507950362285,-0.007991765192846226,-0.01420221700787814,-0.000691356891130895,-0.013924833802594084,-0.004904913020781531,-0.01088768960258949
```

L1 Regularization (Lasso) : The penalty coefficient in L1 is set to its default value of C=1.0.

Code:

```
l1reg.py > ...
1 import pandas as pd
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.feature_selection import SelectFromModel
4 import os
5
6 def apply_l1_regularization(data, target_col, output_file):
7     # Hedef klasörün varlığını kontrol et ve yoksa oluştur
8     output_dir = os.path.dirname(output_file)
9     if not os.path.exists(output_dir):
10         os.makedirs(output_dir)
11
12     X = data.drop(columns=[target_col])
13     y = data[target_col]
14
15     # L1 Regularization uygulama
16     l1_model = LogisticRegression(penalty='l1', solver='liblinear', max_iter=1000)
17     l1_model.fit(X, y)
18
19     selector = SelectFromModel(l1_model, prefit=True)
20     X_new = selector.transform(X)
21
22     selected_features = pd.DataFrame(X_new, columns=data.columns[selector.get_support(indices=True)])
23     selected_features[target_col] = y
24
25     selected_features.to_csv(output_file, index=False)
26     print(f'L1 Regularization ile seçilen özellikler kaydedildi: {output_file}')
27
28     # Veri dosyaları ve klasörler
29     tfidf_files = [
30         ('tfidf_unigram.csv', 'unigram_results/l1_regularization.csv'),
31         ('tfidf_bigram.csv', 'bigram_results/l1_regularization.csv')
32     ]
33
34     # Her dosya için L1 Regularization uygula ve kaydet
35     for input_file, output_file in tfidf_files:
36         data = pd.read_csv(input_file)
37         apply_l1_regularization(data, 'Label', output_file)
38
```

The table after applying L1 Regularization (Lasso) to the dataset with unigram TF-IDF

The table after applying L1 Regularization (Lasso) to the dataset with bigram TF-IDF

VARIANCE THRESHOLD : The threshold value in this reduction method has been set to k=0.001.

Code :

```
threshold.py ×
C:\Users\Semih Zenginoğlu\Desktop\en son data mining\threshold.py (preview ⓘ)
1 import pandas as pd
2 from sklearn.feature_selection import VarianceThreshold
3 import os
4
5 def apply_variance_threshold(data, target_col, threshold, output_file):
6     # Hedef klasörün varlığını kontrol et ve yoksa oluştur
7     output_dir = os.path.dirname(output_file)
8     if not os.path.exists(output_dir):
9         os.makedirs(output_dir)
10
11     X = data.drop(columns=[target_col])
12     y = data[target_col]
13
14     # Variance Threshold uygulama
15     selector = VarianceThreshold(threshold=threshold)
16     X_new = selector.fit_transform(X)
17
18     selected_features = pd.DataFrame(X_new, columns=data.columns[selector.get_support(indices=True)])
19     selected_features[target_col] = y
20
21     selected_features.to_csv(output_file, index=False)
22     print(f"Variance Threshold ile seçilen özellikler kaydedildi: {output_file}")
23
24 # Veri dosyaları ve klasörler
25 tfidf_files = [
26     ('tfidf_unigram.csv', 'unigram_results/variance_threshold.csv'),
27     ('tfidf_bigram.csv', 'bigram_results/variance_threshold.csv')
28 ]
29
30 # Her dosya için Variance Threshold uygula ve kaydet
31 for input_file, output_file in tfidf_files:
32     data = pd.read_csv(input_file)
33     apply_variance_threshold(data, 'Label', threshold=0.001, output_file=output_file)
34
```

The table after applying threshold to the dataset with unigram TF-IDF

The table after applying threshold to the dataset with bigram TF-IDF

STEP 3 – CLASSIFICATION :

In the third step, we classified the datasets reduced using four different methods by applying classification models and compared the results using four different metrics. The classification models we used are Naive Bayes, SVM, and Decision Tree.

1 - NAIVE BAYES MODEL :

We are sharing the classification of the operations we performed using the Naive Bayes algorithm and the corresponding code. First, in the Naive Bayes model training, we used 20% of the data as test data and 80% as training data. The Python code for training this model is as follows:

Code:

```
naive.py > evaluate_model
1 import os
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.naive_bayes import MultinomialNB
5 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
6
7 # Performans değerlendirme fonksiyonu
8 def evaluate_model(data_file, method_name, results):
9     if not os.path.exists(data_file):
10         print(f"Dosya bulunamadı: {data_file}")
11         return
12
13     data = pd.read_csv(data_file)
14     X = data.drop(columns=['Label'])
15     y = data['Label']
16
17     # Negatif değerleri sıfıra çevir
18     X[X < 0] = 0
19
20     # Eğitim ve test setlerine ayırma
21     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
22
23     # Naive Bayes model eğitimi
24     model = MultinomialNB()
25     model.fit(X_train, y_train)
26
27     # Test setinde tahmin
28     y_pred = model.predict(X_test)
29
30     # Performans metriklerini hesapla
31     accuracy = accuracy_score(y_test, y_pred)
32     precision = precision_score(y_test, y_pred, average='weighted')
33     recall = recall_score(y_test, y_pred, average='weighted')
34     f1 = f1_score(y_test, y_pred, average='weighted')
35
36     # Sonuçları kaydet
37     results.append({
38         'Method': method_name,
39         'Accuracy': accuracy,
40         'Precision': precision,
41         'Recall': recall,
42         'F1 Score': f1
43     })
44
45 # Sonuçları karşılaştırmak için fonksiyon
46 def compare_methods(result_folders, output_file):
47     results = []
48     for folder in result_folders:
49         for method_file in ['chi_square.csv', 'pca.csv', 'variance_threshold.csv', 'l1_regularization.csv', 'tfidf.csv']:
50             method_path = os.path.join(folder, method_file)
51             method_name = f'{folder.split('_')[0]}_{method_file.split('.')[0]}'
52             evaluate_model(method_path, method_name, results)
53
54     # Sonuçları bir DataFrame'e dönüştür
55     results_df = pd.DataFrame(results)
56
57     # Excel dosyasına kaydet
58     results_df.to_excel(output_file, index=False)
59     print(f"Sonuçlar kaydedildi: {output_file}")
60
61 # Unigram ve Bigram klasörleri
62 result_folders = ['unigram_results', 'bigram_results']
63
64 # Çıkıtı dosyası
65 output_excel = 'results_comparison.xlsx'
66
67 # Yöntemleri karşılaştır ve Excel'e yaz
68 compare_methods(result_folders, output_excel)
69
```

RESULT OF NAIVE BAYES MODEL :

A	B	C	D	E	F
Method	Accuracy	Precision	Recall	F1 Score	
unigram_chi_square	0,785454545	0,80555711	0,785454545	0,783856903	
unigram_pca	0,598181818	0,632475936	0,598181818	0,583492448	
unigram_variance_threshold	0,763636364	0,769391207	0,763636364	0,76358323	
unigram_l1_regularization	0,836363636	0,842709562	0,836363636	0,836326852	
unigram_tfidf (without reduction)	0,821818182	0,822883567	0,821818182	0,821959736	
bigram_chi_square	0,583636364	0,730715584	0,583636364	0,522739868	
bigram_pca	0,523636364	0,613042298	0,523636364	0,446510344	
bigram_variance_threshold	0,48	0,513195803	0,48	0,379865759	
bigram_l1_regularization	0,578181818	0,662935357	0,578181818	0,534841598	
bigram_tfidf (without reduction)	0,709090909	0,729702556	0,709090909	0,706084467	

2 - SVM MODEL : In this model, we used 20% of the data for testing and 80% for training. A linear kernel SVM model was created and trained using the dataset. We compared the applied TF-IDF and reduction methods based on four metrics and presented the results in an Excel table. The Python code for the SVM model is as follows:

Code :

```
svm.py > evaluate_model
1 import os
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.svm import SVC
5 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
6
7 # Performans değerlendirme fonksiyonu
8 def evaluate_model(data_file, method_name, results):
9     if not os.path.exists(data_file):
10         print(f"Dosya bulunamadı: {data_file}")
11         return
12
13     data = pd.read_csv(data_file)
14     X = data.drop(columns=['Label'])
15     y = data['Label']
16
17     # Negatif değerleri sıfıra çevir
18     X[X < 0] = 0
19
20     # Eğitim ve test setlerine ayırma
21     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
22
23     # SVM model eğitimi
24     model = SVC(kernel='linear') # SVM with linear kernel
25     model.fit(X_train, y_train)
26
27     # Test setinde tahmin
28     y_pred = model.predict(X_test)
29
30     # Performans metriklerini hesapla
31     accuracy = accuracy_score(y_test, y_pred)
32     precision = precision_score(y_test, y_pred, average='weighted')
33     recall = recall_score(y_test, y_pred, average='weighted')
34     f1 = f1_score(y_test, y_pred, average='weighted')
```

```

35
36     # Sonuçları kaydet
37     results.append({
38         'Method': method_name,
39         'Accuracy': accuracy,
40         'Precision': precision,
41         'Recall': recall,
42         'F1 Score': f1
43     })
44
45 # Sonuçları karşılaştırmak için fonksiyon
46 def compare_methods_svm(result_folders, output_file):
47     results = []
48     for folder in result_folders:
49         for method_file in ['chi_square.csv', 'pca.csv', 'variance_threshold.csv', 'l1_regularization.csv', 'tfidf.csv']:
50             method_path = os.path.join(folder, method_file)
51             method_name = f"{folder.split('_')[0]}_{method_file.split('.')[0]}"
52             evaluate_model(method_path, method_name, results)
53
54     # Sonuçları bir DataFrame'e dönüştür
55     results_df = pd.DataFrame(results)
56
57     # Excel dosyasına kaydet
58     results_df.to_excel(output_file, index=False)
59     print(f"SVM sonuçları kaydedildi: {output_file}")
60
61 # Unigram ve Bigram klasörleri
62 result_folders = ['unigram_results', 'bigram_results']
63
64 # Çıkıtı dosyası
65 output_excel_svm = 'svm_results_comparison.xlsx'
66
67 # Yöntemleri karşılaştır ve Excel'e yaz
68 compare_methods_svm(result_folders, output_excel_svm)
69

```

RESULT OF SVM MODEL :

1	Method	Accuracy	Precision	Recall	F1 Score
2	unigram_chi_square	0,783636364	0,78673522	0,783636364	0,781933833
3	unigram_pca	0,681818182	0,690640005	0,681818182	0,673450865
4	unigram_variance_threshold	0,765454545	0,765627885	0,765454545	0,764803506
5	unigram_l1_regularization	0,834545455	0,83447994	0,834545455	0,834487791
6	unigram_tfidf(without reduction)	0,827272727	0,827741088	0,827272727	0,827372859
7	bigram_chi_square	0,585454545	0,749975702	0,585454545	0,521276313
8	bigram_pca	0,516363636	0,666355015	0,516363636	0,410842975
9	bigram_variance_threshold	0,483636364	0,65742621	0,483636364	0,334852604
10	bigram_l1_regularization	0,58	0,676552285	0,58	0,532813988
11	bigram_tfidf(without reduction)	0,709090909	0,716029406	0,709090909	0,708783105

3 – DECISION TREE :

In this study, 80% of the dataset was used for training and 20% for testing. A classic decision tree model, DecisionTreeClassifier, was used. The model utilized the default Gini Impurity criterion for making branching decisions. The Python code for the Decision Tree model is as follows:

Code :

```
decissiontree.py > ...
1 import os
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
6
7 # Performans değerlendirme fonksiyonu
8 def evaluate_model(data_file, method_name, results):
9     if not os.path.exists(data_file):
10         print(f"Dosya bulunamadı: {data_file}")
11         return
12
13     data = pd.read_csv(data_file)
14     X = data.drop(columns=['Label'])
15     y = data['Label']
16
17     # Negatif değerleri sıfır'a çevir
18     X[X < 0] = 0
19
20     # Eğitim ve test setlerine ayırma
21     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
22
23     # Decision Tree model eğitimi
24     model = DecisionTreeClassifier(random_state=42) # Karar ağacı modeli
25     model.fit(X_train, y_train)
26
27     # Test setinde tahmin
28     y_pred = model.predict(X_test)
29
30     # Performans metriklerini hesapla
31     accuracy = accuracy_score(y_test, y_pred)
32     precision = precision_score(y_test, y_pred, average='weighted')
33     recall = recall_score(y_test, y_pred, average='weighted')
34     f1 = f1_score(y_test, y_pred, average='weighted')
```

```

35
36     # Sonuçları kaydet
37     results.append({
38         'Method': method_name,
39         'Accuracy': accuracy,
40         'Precision': precision,
41         'Recall': recall,
42         'F1 Score': f1
43     })
44
45 # Sonuçları karşılaştırmak için fonksiyon
46 def compare_methods_decision_tree(result_folders, output_file):
47     results = []
48     for folder in result_folders:
49         for method_file in ['chi_square.csv', 'pca.csv', 'variance_threshold.csv', 'l1_regularization.csv', 'tfidf.csv']:
50             method_path = os.path.join(folder, method_file)
51             method_name = f"{folder.split('_')[0]}_{method_file.split('.')[0]}"
52             evaluate_model(method_path, method_name, results)
53
54 # Sonuçları bir DataFrame'e dönüştür
55 results_df = pd.DataFrame(results)
56
57 # Excel dosyasına kaydet
58 results_df.to_excel(output_file, index=False)
59 print(f"Decision Tree Sonuçları kaydedildi: {output_file}")
60
61 # Unigram ve Bigram klasörleri
62 result_folders = ['unigram_results', 'bigram_results']
63
64 # Çıktı dosyası
65 output_excel_tree = 'decision_tree_results_comparison.xlsx'
66
67 # Yöntemleri karşılaştır ve Excel'e yaz
68 compare_methods_decision_tree(result_folders, output_excel_tree)
69

```

RESULT OF DECISION TREE MODEL:

	A	B	C	D	E
1	Method	Accuracy	Precision	Recall	F1 Score
2	unigram_chi_square	0,718182	0,722424242	0,718181818	0,7144022
3	unigram_pca	0,610909	0,610340057	0,610909091	0,610495302
4	unigram_variance_threshold	0,638182	0,638260429	0,638181818	0,638219015
5	unigram_l1_regularization	0,716364	0,716649994	0,716363636	0,716468925
6	unigram_tfidf(without reduction)	0,692727	0,695572476	0,692727273	0,692937595
7	bigram_chi_square	0,576364	0,725913212	0,576363636	0,511514404
8	bigram_pca	0,592727	0,593712024	0,592727273	0,593023946
9	bigram_variance_threshold	0,525455	0,508011698	0,525454545	0,43591626
10	bigram_l1_regularization	0,567273	0,644643806	0,567272727	0,522811639
11	bigram_tfidf(without reduction)	0,578182	0,600193146	0,578181818	0,568061584

RESULT AND COMPARISONS :

In this section, we will compare the methods, models, and types of TF*IDF used. Based on this study, I am sharing the most successful methods, models, and types of TF*IDF.

1) Comparison of TF-IDF Methods

Unigram TF-IDF provided the most effective representation in all models and delivered high accuracy rates.

2) Comparison of Dimensionality Reduction Methods

L1 Regularization has demonstrated the highest performance across all models, making it the most successful dimensionality reduction method. It minimized data loss, enhancing the model's overall performance.

Chi-Square produced effective results, particularly with the SVM model, but it was not as strong as L1 Regularization.

PCA showed the lowest performance, making it an ineffective dimensionality reduction method.

Non-reduced TF-IDF achieved results close to the dimensionality reduction techniques in some cases. Not applying dimensionality reduction can be effective, especially for smaller datasets.

3) Model Comparison

Naive Bayes exhibited the best performance in terms of Accuracy and F1 Score, ranking first overall.

SVM produced results close to Naive Bayes, ranking second. It is competitive in terms of F1 Score.

Decision Tree showed lower performance compared to the other models and was less effective for high-dimensional datasets.

4. Best Combinations

Naive Bayes + L1 Regularization + Unigram TF-IDF

Accuracy: 83.64%, F1 Score: 83.63%

SVM + L1 Regularization + Unigram TF-IDF

Accuracy: 83.45%, F1 Score: 83.44%

SVM + Non-reduced TF-IDF + Unigram

Accuracy: 82.73%, F1 Score: 82.73%

Conclusion

L1 Regularization is the most effective dimensionality reduction method.

Naive Bayes ranks first in overall performance.

Unigram TF-IDF, when combined with dimensionality reduction methods, produces the strongest results.

Non-reduced TF-IDF is an effective alternative for smaller datasets.

AS A RESULT :

We worked on making the given datasets more understandable and simplified by performing specific operations. As a result of these operations, we used certain algorithms to classify the outcomes, trained them, and compared their performance. We aimed to demonstrate the operations, methods, and algorithms we used step by step and shared their results in this report. Everything we did from the start to the end of this project has been explained in this report.

REFERENCES:

-) <https://tr.wikipedia.org/wiki/Anasayfa>
-) <https://archive.ics.uci.edu/dataset/331/sentiment+labelled+sentences>
-) <https://www.geeksforgeeks.org/>
-) <https://www.geeksforgeeks.org/understanding-tf-idf-term-frequency-inverse-document-frequency/>
-) <https://www.ibm.com/think/topics/classification-models>
-) ChatGPT