



T.C.

**SİVAS CUMHURİYET ÜNİVERSİTESİ
TEKNİK BİLİMLER MESLEK YÜKSEKOKULU
BÜYÜK VERİ ANALİSTLİĞİ BÖLÜMÜ**

DOĞAL DİL İŞLEME UYGULAMA RAPORU

GÖKAY İLYAS KOÇ

Danışman

Öğretim Görevlisi İstatistik Bölüm Başkanı Fatih GÖKÇE

Sivas

25. 2025

Twitter Metin Sınıflandırması

Hazırlayan: Gökay İlyas Koç

Öğrenci Numarası: 2024295008

Tarih: 25 Aralık 2025

Özet

Bu projede Twitter tweet verileri kullanarak tweetlerin **Nefret**, **Saldırgan** ve **Nötr** olmak üzere üç farklı sınıfa ayırmayı amaçlıyoruz. Projede Python ve makine öğrenmesi yöntemlerini kullanıyoruz. Tweetler üzerinde çeşitli ön işleme adımlarını deneyip ardından TF-IDF yaparak metinleri sayısal değerlere çeviriyoruz. Farklı sınıflandırma algoritmalarını kullanarak modelimizin performanslarını karşılaştırıp elde ettiğimiz sonuçlara göre en yüksek doğruluk oranı veren modeli belirleyip, performans analizini gerçekleştiriyoruz.

1. Amaç

Projemizin amacı, Tweetler'den elde edilen metinleri kullanarak nefret söylemi mi yoksa saldırgan içerik mi otomatik olarak tespit etmeyi amaçlıyoruz. Bu süreç sonunda makine öğrenmesi yöntemlerini deneyerek, oranları karşılaştırıyoruz ve en başarılı modeli belirliyoruz.

2. Kullanılan Veri ve Ön İşlemler

Projede kaggle üzerinden aldığımız **veri.csv** adlı içinde tweetler bulunan veri setini indirip kullanacağız. İşlemler öncesinde metinleri bir takım ön işleme adımları uygulayıp deneyeceğiz.

Uygulanan Ön İşleme Adımları

- tweet metinlerini küçük harfe dönüştürdüm.
- Linkleri, kullanıcı hashtag , sayıları ve noktalama işaretlerini temizleyip tweetleri analize uygun hale getiriyorum.
- Temizlenen tweetleri **clean_text** adında oluşturduğumuz yeni bir sütuna aktarıyoruz.

Temizlik Kontrolü

Orijinal Tweetler

!!! RT @mayasolovely: As a woman you shouldn't...
!!!! RT @mleew17: boy dats cold...
!!!!!! RT @UrKindOfBrand Dawg!!!!
!!!!!!! RT @C_G_Anderson:
!!!!!!!!!!!! RT @ShenikaRoberts:

Temizlenmiş Tweetler

rt as a woman you shouldnt complain about...
rt boy dats coldtyga dwn bad...
rt dawg rt you ever fuck...
rt she look like a tranny
rt the shit you hear about me...

Bu tweetler, temizleme işlemini doğru uyguladığımızı gösteriyor.

3. Özellik Çıkarımı

Tweetleri sayısal hale dönüştürmek için **TF-IDF (Term Frequency – Inverse Document Frequency)** kullanıyoruz. Bu işlemde **5000 kelimelik** bir özellik oluşturuyoruz.

- **TF-IDF matris boyutu:** (19.826, 5000)

Bu boyut, verinin yeterli kelime bilgisiyle eğitildiğini kanıtlar niteliktedir.

4. Kullanılan Modeller ve Performansları

Projemizde beş tane makine öğrenmesi yöntemini kullanarak eğitiyoruz ve test verisi ile doğruluk (accuracy) değerini hesaplıyoruz.

Model	Doğruluk (%)
Lojistik Regresyon	89.75
Naive Bayes	83.50
Destek Vektör Makineleri (SVM)	89.39
Karar Ağacı	87.37
Rastgele Orman	89.51

Sonuçlara göre en yüksek oranı **Lojistik Regresyon** modelinde yakalıyoruz.

5. Detaylı Performans Analizi

En yüksek oran veren modelimiz **Lojistik Regresyon** için sınıf bazlı precision, recall ve F1-score değerlerini hesaplıyoruz ekte:

	precision	recall	f1-score	support
Nefret	0.62	0.20	0.31	295
Saldırgan	0.92	0.96	0.94	3814
Nötr	0.84	0.84	0.84	848
accuracy			0.90	4957

Değerlendirme

- Modelimizin yüksek bir doğruluk oranına sahip olduğunu görüyoruz. (%89.75).
- Özellikle saldırgan sınıfında başarılı sonuçlar elde ettiğimizi gördük.
- Nefret sınıfında recall değerinin düşük olması ise bu sınıftaki istisna verilerin doğru tespit edilmediğini görüyoruz.
- Nötr sınıfında ise dengeli bir performans dağılımı gözlemliyoruz.

6. Karışıklık Matrisi

Modelimizin hata analizinde ise karışıklık matrisini inceliyoruz:

Gerçek \ Tahmin	Nefret	Saldırgan	Nötr
Nefret	60	204	31
Saldırgan	37	3673	104
Nötr	0	132	716

Tabloda, nefret içerikli tweetlerin genellikle saldırgan sınıfı ile karıştırıldığını görüyoruz.

7. Görselleştirme

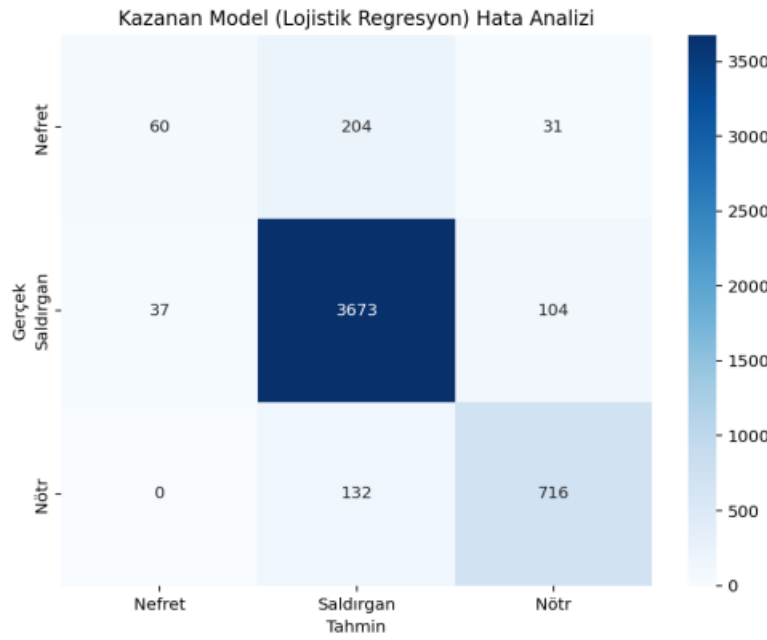
Karışıklık matrisi, **seaborn** kütüphanesiyle ısı haritası (heatmap) şeklinde görselleştiriyoruz. Bu görsel grafik, modelin hangi sınıflarda hata yaptığını daha net bir şekilde gösteriyor.

8. Sonuç

Bu analiz projemizde Twitter tweetleri üzerinde başarılı bir sınıflandırma süreci gerçekleştirdik ve Lojistik Regresyon modelinin en iyi performansını yakaladık.

EKLER

Ek-1: Karışıklık Matrisi Isı Haritası



Ek-2: Variable Explorer Çıktıları

Name	Type	Size	Value
cm	Array of int64	[3, 3]	[[60 204 31] [37 3673 104] [0 0 0]]
df	DataFrame	[24783, 8]	Column names: Unnamed: 0, count, hate_speech, offensive_language, neit ...
en_yuksek_puan	float	1	0.8975186604801291
isim	str	14	Rastgele Orman
kazanan_model_ismi	str	18	Lojistik Regresyon
kazanan_tahminler	Array of int64	[4957]	[1 1 1 ... 0 1 2]
model	ensemble._forest.RandomForestClassifier	50	RandomForestClassifier object of sklearn.ensemble._forest module
modeller	dict	5	{'Lojistik Regresyon':LogisticRegression, 'Naive Bayes':MultinomialNB, ...}
puan	float	1	0.8957030461972968
tahmin	Array of int64	[4957]	[1 1 1 ... 1 1 2]

tfidf	feature_extraction.text.TfidfVectorizer	1	TfidfVectorizer object of sklearn.feature_extraction.text module
X_test	Series	[4957]	Series object of pandas.core.series module
X_test_vec	sparse._csr.csr_matrix	[4957, 5000]	csr_matrix object of scipy.sparse._csr module
X_train	Series	[19826]	Series object of pandas.core.series module
X_train_vec	sparse._csr.csr_matrix	[19826, 5000]	csr_matrix object of scipy.sparse._csr module
y_test	Series	[4957]	Series object of pandas.core.series module
y_train	Series	[19826]	Series object of pandas.core.series module

Confusion Matrix (cm)

cm değişkeni, 3x3 boyutunda Confusion Matrix'i temsil etmektedir.

- Satırlar: Gerçek sınıflar
- Sütunlar: Model tarafından tahmin edilen sınıflar

Veri Seti (df)

- Satır sayısı: 24.783
- Sütun sayısı: 8

En Yüksek Başarı Skoru (en_yuksek_puan)

- Değer: %89.75

Kazanan Model Bilgileri

- Kazanan Model: Lojistik Regresyon

Ek-3: DataFrame Sütun Açıklamaları

Index	Unnamed: 0	count	hate_speech	offensive_language	neither	class	tweet	clean_text
0	0	3	0	0	3	2	!!! RT @mayasolely: As a woman you shouldn't compl-	rt as a woman you shoudnt complain about cleaning-
1	1	3	0	3	0	1	!!!!!! RT @blowm17: boy datz col6f...tyga dem bad for	rt boy datz coldtyga dem bad for cuffin dat hoe in
2	2	3	0	3	0	1	!!!!!! RT @xindofbrand Dmg!!!! RT @soobabyalife-	rt dmeg rt you ever fuck a bitch and she start to
3	3	3	0	2	1	1	!!!!!! RT @C_G_Anderson: @viva_based she look like	rt she look like a tranny
4	4	6	0	6	0	1	!!!!!! RT @Shenkaaberts: The shit you hear-	rt the shit you hear about me might be true or it
5	5	3	1	2	0	1	!!!!!! RT @Madison_x: The shit just blows-	the shit just blows meclain you so faithful and down
6	6	3	0	3	0	1	!!!!!! @_BrighterDays: I can not just sit up and HA-	i can not just sit up and hate on another bitch i
7	7	3	0	3	0	1	!!!!!! @0220/RealQueenBri: cause I'm tired of you	got too much shit going on
8	8	3	0	3	0	1	" ∓ you might not get ya bitch back ∓ thats	cause im tired of you big bitches coming for us
9	9	3	1	2	0	1	" @rhythmicx...hobbies include fighting Marlan"	skiny girls
10	10	3	0	3	0	1	hitch	asp you might not get ya bitch back asp thats that
11	11	3	0	3	0	1	" Keeks is a bitch she curves everyone " lol I walked	hobbies include fighting marlan
12	12	3	0	2	1	1	into a conversation like this. Sm	hitch
13	13	3	0	3	0	1	" Murda Gang bitch its Gang Land "	keeks is a bitch she curves everyone lol i walked
14	14	3	1	2	0	1	" So hoes that smoke are losers ? " yea ... go on IG	into a conversation like this sm
15	15	3	0	3	0	1	" bad bitches is the only thing that i like "	murda gang bitch its gang land
16	16	3	0	3	0	1	" bitch get up off me "	so hoes that smoke are losers yea go on ig
17	17	3	1	2	0	1	" bitch nigga miss me with it "	bad bitches is the only thing that i like
18	18	3	0	3	0	1	" bitch plz whatever "	bitch get up off me
19	19	3	0	3	0	1	" bitch who do you love "	bitch nigga miss me with it
20	20	3	0	3	0	1	" bitches get cut off everyday S "	bitch plz whatever
21	21	3	0	3	0	1	" black bottle ∓ a bad bitch "	bitch who do you love
		3	0	3	0	1	" broke bitch cant tell me nothing "	bitches get cut off everyday b
		3	0	3	0	1	" cancel that bitch like Nino "	black bottle asp a bad bitch
		3	0	3	0	1		broke bitch cant tell me nothing
		3	0	3	0	1		cancel that bitch like nino

1. **Unnamed: 0** – Otomatik indeks sütunu
2. **count** – Metin içerisindeki kelime sayısı
3. **hate_speech** – Nefret söylemi etiketi
4. **offensive_language** – Saldırgan dil etiketi
5. **neither** – Nötr içerik etiketi
6. **class** – Nihai sınıf etiketi
7. **Metin sütunu** – Ham tweet metni
8. **Ön işlenmiş metin sütunu** – Temizlenmiş metin

Ek-4: Kullanılan Kütüphaneler

- Pandas
- Scikit-learn
- Matplotlib
- Seaborn