
KMEANS ALGORITMASI ILE TWEET KÜMELEME

DOĞAL DİL İŞLEME DERSİ PROJE ÖDEVİ

Enes Sadi Uysal

Bilgisayar Mühendisliği 4. Sınıf
Yıldız Teknik Üniversitesi
17011041

Ahmet Salih Özmen

Bilgisayar Mühendisliği 4. Sınıf
Yıldız Teknik Üniversitesi
18011081

December 28, 2021

1 Giriş

Twitter günümüzde en çok kullanılan uygulamalardan biri. Her gün atılan binlerce tweet arasından hızla akan gündemi yakalamak, sadece gülmek istemek çok doğal bir ihtiyaç. Fakat atılan tweet'ler çoğu zaman etiket (hashtag) kullanılmadan atılmakta bu ise tweet'leri filtreleme işini zorlaştırmakta. Her ne kadar kelimeler ile filtreleme yapılabilse de, insanların görmek istediği tweet türü kelimelerle ifade edilemeyebilir. İşte bu noktada ayrıştırma işlemi kelimelerle değil tweet'lerin kendileriyle, tamamı kullanılarak yapıp farklı çeşitlerde tweet'lerden oluşan kümeler ortaya konulursa, istediği türü tam olarak bilmeyen kişiler bile beğendiği türdeki tweet'leri, beğendiği mizahı veya sadece siyaset gündemini içeren kümeleri bulabileceklerdir. Projemizin çıkış noktası olan bu fikre dair ilerlemelerimizi, tweet'leri temizlerken izlediğimiz adımları, tweet'leri vektörlere nasıl dönüştürdüğümüzü ve KMeans algoritmasının verimiz üzerindeki başarısını bu raporda anlatacağız.

2 Proje Adımları

Projenin ilk adımı Twitter üzerinden tweet verisinin çekilmesidir. Bu adımda kullanılan kütüphaneden ve tweet'leri alırken yaşanan sorunlardan bir sonraki bölümde bahsedilecektir. Tweet verisi farklı etiketler (hashtag'ler) anahtar kelime olarak verilip ilgili metni içeren tweet'lerin bütün özellikleriyle birlikte çekilmesi ile gerçekleştirilmiştir. Burada her bir anahtar kelime sonucu gelen tweet'ler etiketi o anahtar kelime olmak üzere bir arada etiketli bir şekilde tutulmuştur. Ardından tweet verileri gereksiz özelliklerden temizlenip, tweet metinleri bir sonraki bölümde ayrıntılı açıklanacak olan birçok farklı yöntem ile temizlenmiştir. Temizlenen metinler Türkçe kelime modelleri kullanılarak vektörlere dönüştürülüp sayısal hale getirilmiştir. Sonrasında elde edilen sayısal veri etiket çeşidi sayısı kadar küme belirlenerek kümeleme algoritması olan KMeans'e etiketsiz bir şekilde eğitim verisi olarak verilmiştir. Buradan alınan sonuçların değerlendirmesi ise verilerin etiketleri kullanılarak ve kümelerin entropilerine bakılarak yapılmıştır. Ayrıca Twitter üzerinden toplanan verinin yetersiz olması sebebi ile daha önceden toplanmış tweet verileri veya metin sınıflandırma verileri de model başarısını ölçmede kullanılmıştır.

3 Yöntemlerin Açıklanması

3.1 Twint ile Tweet Verisinin Çekilmesi

Twitter üzerinden veri çekmek için Python dilinde geliştirilen ve kütüphaneleştirilen Twint (Twitter Intelligence Tool) kullanılmıştır. Twint, komut satırı üzerinden çalışabilse de, veri çekerken karşılaşılan kesintiye uğrama sebebiyle ve dosyaların daha düzenli saklanabilmesi sebebi ile bir fonksiyon yazılmıştır.

```

1 def scrape_tw(since, until, keywords, output_dir):
2     """
3     - since (str)      : starting date for scraping. Exp: 2020-01-01 15:30:00
4     - until (str)      : ending date for scraping.   Exp: 2021-01-01 16:00:00
5     - keywords (list)  : which words do you want to scrape on twitter.
6     - output_dir (str) : the path where the csv file will be saved.
7     """
8
9     keys = keywords[0]
10    if len(keywords) > 1:
11        for keyword in keywords[1:]:
12            keys += " OR " + keyword
13    print("Keys: " + keys)
14
15    first_until = until
16    uncut_output = output_dir + '/' + keywords[0] + '_' + since.split()[0] + '_' + since.split()[1] + '_' + first_until.split()[0] + '_' + first_until.split()[1] + '.csv'
17
18    keep = 0
19    while keep == 0:
20        output = output_dir + '/' + keywords[0] + '_' + since.split()[0] + '_' + since.split()[1] + '_' + until.split()[0] + '_' + until.split()[1] + '.csv'
21        call = 'twint -o "' + output + '" --csv --tabs -s "' + keys + '" --since "' + since + '" --until "' + until + "'"
22        a = os.system(call)
23        print(a, call)
24
25        df = pd.read_csv(output, sep = '\t')
26        ctrl = df.iloc[-1]['date'] + ' ' + df.iloc[-1]['time']
27        print(ctrl)
28        if ctrl.split()[0] == since.split()[0] and ctrl.split()[1].split(':')[0] == since.split()[1].split(':')[0] and ctrl.split()[1].split(':')[1] == since.split()[1].split(':')[1]:
29            if output == uncut_output:
30                print("File saved here: " + uncut_output)
31            else:
32                print("File saved here: " + output)
33            keep = 1
34        else:
35            new_name = output_dir + '/' + keywords[0] + '_' + ctrl.split()[0] + '_' + ctrl.split()[1] + '_' + until.split()[0] + '_' + until.split()[1] + '.csv'
36            os.rename(output, new_name)
37            print("File saved here: " + new_name)
38            until = ctrl
39            keep = 1

```

Figure 1: Crawling Function

Figür 1’de gösterilen fonksiyonda *since* ve *until* parametreleri ile hangi tarih ve saatten başlanıp hangi tarih ve saate kadar tweet’lerin çekileceği bilgisi verilir. Ardından çekilecek tweet metinlerinin içermesi gereken anahtar kelimeler *keywords* parametresi ile bir liste olarak verilir. En son ise oluşturulan veri tablolarının nereye kaydedileceği bilgisi *output_dir* parametresi ile verilir. Tweet çekme işlemi sırasında verilen zaman aralığındaki tweet’ler bitmeden çalışma kesintiye uğrarsa, çekilen zamana kadar olan tweet verilerini isimlendirmesini Figür 2’de görüldüğü gibi düzenleyerek kaydeder. Ardından kaldığı yerden tweet çekme işlemine devam eder ve belirlenen zaman aralığındaki bütün tweet’ler bitene kadar çeker.

```

btc_2020-07-28_06:30:11_2020-07-30_08:48:39.csv
btc_2020-07-30_08:48:39_2020-07-30_13:05:26.csv
btc_2020-07-30_13:05:26_2020-07-31_08:45:51.csv

```

Figure 2: Naming Files

3.2 Türkçe GloVe Kelime Modelinin Kullanılması

Tweet metinlerinin vektörlere dönüştürülebilmesi için içindeki her kelimenin vektörünün çıkartılması gerekmektedir. Kelimelerin vektörlere dönüştürülmesi için ise Türkçe veriler ile eğitilmiş bir kelime modeli kullanılması gerekmektedir. Bu projede kullanılan model inzva kapsamında gerçekleştirilen AI Projects etkinliğinde sunulan Türkçe GloVe modelidir.

Projede kullanım kolaylığı olması açısından Figür 3'te görülebildiği gibi GloVe vektörleri Word2Vec modeli haline getirilmiştir.

```
1 glove_file = 'vectors.txt'
2 word2vec_file = 'word2vec.txt'
3
4 glove2word2vec(glove_file, word2vec_file)
5 model = KeyedVectors.load_word2vec_format(word2vec_file)
```

Figure 3: GloVe to Word2Vec

3.3 Metin Ön işleme Adımları

Twitter üzerinden alınan metinler normal metinlerden daha düzensiz ve daha kurallara aykırı ve noktalama işaretleri bakımından yetersiz olduğundan, birçok temizleme işleminden geçirilmiştir.

Önce bütün metin küçük harflere dönüştürülmüştür. Ardından kelimeler arasında kullanılabilecek '_' gibi karakterler metinlerden silinip boşluk karakteri ile doldurulmuştur. Sonrasında tweet içerisinde bulunabilecek olan linkler (http...) veya etiketler (@username) RegEx sorguları kullanılarak temizlenmiştir (Figür 4). Ardından bütün noktalama işaretleri boş string ile değiştirilip metin bütün noktalama işaretlerinden temizlenmiştir.

```
1 for k in range(len(df)):
2     df['text'][k] = re.sub(r'http\S+', '', str(df['text'][k])) # remove links
3     df['text'][k] = re.sub(r'@\S+', '', str(df['text'][k])) # remove usertags
```

Figure 4: RegEx for Removing Links and Tags

Türkçe "stopword" olarak sayılan kelimeler metinlerden temizlenmek için **nlTK** modülünden alınan stopwords listesine bakılmıştır. Bütün metinler taranarak Figür 5'teki kod ile temizlenmiştir. Ardından bütün kelimeleri stopwords olan, yani bu temizleme işlemi bittiğinde tamamı silinen metinler veri setinden çıkarılmıştır.

```

1 docs = df['text']
2 docs_list = []
3
4 WPT = nltk.WordPunctTokenizer()
5 stop_word_list = nltk.corpus.stopwords.words('turkish')
6
7 for doc in docs:
8     doc = re.sub("\d+", " ", doc) # remove numbers
9     doc = WPT.tokenize(doc)
10    filtered_tokens = [item for item in doc if item not in stop_word_list]
11    lemma = nltk.WordNetLemmatizer()
12    lemma_word = [lemma.lemmatize(word) for word in filtered_tokens]
13
14    doc = " ".join(lemma_word)
15
16    docs_list.append(doc)
17
18 df['text'] = pd.DataFrame(docs_list, columns=['text'])

```

Figure 5: Removing Stopwords

Kullanıcı isteğine bağlı olarak metinlere "stemming" işlemi uygulanabilmektedir. Stemming kelimelerin gövde veya kök halleri ile tutulması işlemidir. Türkçe kelimeler için stemming yapan bir Python kütüphanesi olan TurkishStemmer kullanılmıştır. Figür 6'daki şekilde kullanılan kütüphane ile kelimeler kök hallerine döndürülebilmektedir.

```

1 if stemming:
2     stemmer = TurkishStemmer()
3     for i, text in enumerate(df['text']):
4         word_list = text.split()
5         new_sentence = ""
6         for word in word_list:
7             new_word = stemmer.stem(word)
8             new_sentence += (" " + new_word)
9
10    df['text'][i] = new_sentence

```

Figure 6: Turkish Stemming

Temizleme işlemleri bittiğinde kelimelerin vektörleri bulunmuştur. Kelimelerin vektör ortalamalarını alıp metinlerin vektörleri ortaya çıkartılmıştır. Bu aşamada kelimelerin tamamının vektörlerini almak yerine **Bag of Words** yöntemi ile (Figure 7) metinden eşsiz kelimeleri elde edip, onların vektör ortalamaları alınmıştır. Metindeki hiçbir kelime sözlükte yani kelime modelinde bulunamazsa o metin veri setinden çıkartılmıştır.

```

1 def find_bow(sentence):
2     tokenizer = Tokenizer()
3     tokenizer.fit_on_texts(sentence)
4     sequences = tokenizer.texts_to_sequences(sentence)
5     word_index = tokenizer.word_index
6     bow = {}
7     for key in word_index:
8         bow[key] = sequences[0].count(word_index[key])
9
10    return bow, len(word_index)

```

Figure 7: Bag of Word

3.4 KMeans Model Eğitimi ve Değerlendirmesi

Kümeleme işlemi yapılacağından, bilinen ve iyi çalışan kümeleme algoritmalarından biri olan KMeans modeli kullanılmıştır. Burada oluşturulacak küme sayısı toplanan tweet'lerin sınıf sayısı olarak belirlenmiştir. (Figür 8)

```

1 def train_kmeans(df):
2     X = np.array([np.array([*arr]) for arr in [vector for _, vector in enumerate(df["vectors"], 0)]]])
3     y = np.array(df['labels'])
4
5     kmeans = KMeans(n_clusters=len(df['labels'].unique()))
6     kmeans.fit(X)
7
8     df['clusters'] = kmeans.labels_
9
10    avg_ent = average_entropy(df)
11
12    return kmeans, avg_ent

```

Figure 8: KMeans Training

Eğitim gerçekleştikten sonra bulunan kümelerin değerlendirilmesi için kümelerin entropileri hesaplanmıştır ve her sınıfa ait ne kadar veri olduğu çizdirilerek kontrol edilmiştir. Kümelerin entropilerinin ortalaması alınarak ise ortalama entropi bulunmuş ve model başarısı bu şekilde ölçülmüştür. (Figür 9)

```

1 def ent(data):
2     """
3     Calculates entropy of the passed 'pd.Series'
4     """
5     p_data = data.value_counts() # counts occurrence of each value
6     entropy = scipy.stats.entropy(p_data) # get entropy from counts
7
8     return entropy
9
10 def average_entropy(df):
11     sum_ent = 0
12     n_clusters = len(df['clusters'].value_counts())
13
14     for i in range(n_clusters):
15         sum_ent += ent(df[df['clusters'] == i]['labels'])
16
17     return sum_ent/n_clusters

```

Figure 9: Entropy Average

3.5 Tweet Tahminlemesi

Girdi olarak bir metin alan fonksiyon alınan metne ön işleme adımlarını uyguladıktan sonra kelimelerine ayırıp metnin vektörünü hesaplar. Bu vektör modele girdi olarak verilir ve alınan küme numarası çıktısı daha önceden atanan sınıf etiketi olarak fonksiyondan geri döndürülür.(Figür 10)

```

1 def predict_text(text, assign_dict, kmeans, model):
2     '''
3     - text (str)           : string for prediction
4     - dict_for_clusters (dict) : dictionary for cluster labels of hashtags
5     - kmeans              : model for predict the cluster of given vector
6     - model               : model for taking vector values of words
7     '''
8     dict_for_clusters = assign_dict_clusters(assign_dict, kmeans, df)
9     text = text.lower()
10    text = remove_punctuation(text)
11
12    WPT = nltk.WordPunctTokenizer()
13    stop_word_list = nltk.corpus.stopwords.words('turkish')
14
15    text = re.sub("\d+", " ", text) # remove numbers
16    text = WPT.tokenize(text)
17    filtered_tokens = [item for item in text if item not in stop_word_list]
18    lemma = nltk.WordNetLemmatizer()
19    lemma_word = [lemma.lemmatize(word) for word in filtered_tokens]
20
21    text = " ".join(lemma_word)
22
23    not_found = 0
24    bow, count = find_bow([text])
25    word_list = list(bow.keys()) # finding bag of words for a sample
26    avg_vector = np.zeros(300)
27    for word in word_list:
28        try:
29            avg_vector += model.get_vector(word) # finding vectors of these words
30        except:
31            not_found += 1
32    if count != 0:
33        avg_vector = avg_vector/count # taking average of all words in bag of words for this sample
34
35    avg_vector = avg_vector.reshape(1,-1)
36    prediction = dict_for_clusters[kmeans.predict(avg_vector)[0]]
37
38    return prediction

```

Figure 10: Prediction Function

4 Sonuçlar

Bu projede 3 farklı veri seti üzerinde çalışıldı. Veri setlerinin ikisi Kaggle üzerinden hazır elde edilmiş olup, bir tanesi Twitter üzerinden çekilerek elde edilmiştir.

Projede kullanılan veri setleri:

- A Benchmark Data for Turkish Text Categorization @YTU Kemik Group
- Turkish Tweets Dataset
- Twitter üzerinden <https://github.com/twintproject/twint> ile toplanan veri seti

4.1 "A Benchmark Data for Turkish Text Categorization" Verisi ile Alınan Sonuçlar

Bu veri seti ekonomi, siyaset, dünya, teknoloji, spor, sağlık, kültür olarak 7 sınıftan oluşmaktadır. Bu sınıflara ait veriler düzenli dağılmış olup veri setinde model performansını etkileyecek bir dengesizlik yoktur.

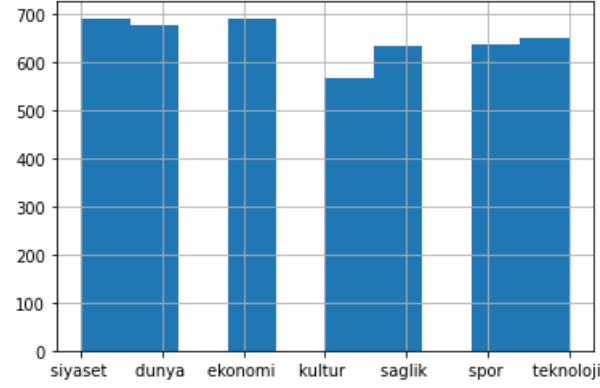


Figure 11: Histogram of Dataset

Eğitilen model **0.7378** ortalama entropiye sahip olup, diğer veri setleri ile oluşturulan modellerden daha iyi çalışmaktadır. Kümelerin ortalama sınıf dağılımları Figür 12'dekine benzer olduğundan, verileri doğru bir biçimde ayırdığı söylenebilir.

	index	labels
0	teknoloji	433
1	ekonomi	25
2	kultur	14
3	dunya	4
4	siyaset	4
5	saglik	1

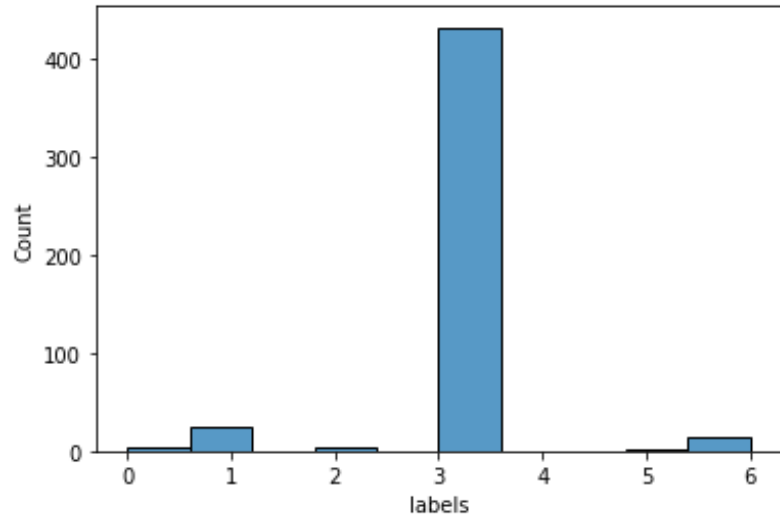


Figure 12: Distribution of Cluster

Figure 13: Histogram of Cluster

4.2 "Turkish Tweets Dataset" Verisi ile Alınan Sonuçlar

Bu veri seti mutlu, kızgın, üzgün, korku ve sürpriz olmak üzere 5 sınıftan oluşmaktadır. Bir önceki verideki gibi bu veride de sınıflara ait örnekler eşit dağılmıştır.

Eğitilen model **1.3968** ortalama entropiye sahiptir ve bir önceki veride çalıştığı kadar iyi çalışmamaktadır.

	index	labels
0	korku	279
1	surpriz	190
2	kızgın	153
3	üzgün	105
4	mutlu	93

Figure 14: Histogram of Cluster

Figür 14'te verilen kümenin "korku" sınıfını temsil ettiği söylenebilse de, entropisinin yüksek olduğu görülmektedir. Ayrıca bu veri seti ile gerçekleştirilen eğitim sonucunda her sınıf tek bir küme tarafından temsil edilememekte, bazı sınıflar iki küme tarafından temsil edilmektedir. Çok başarılı bir model olmasa da tahmin sonuçları, verilen girdinin basitliğine göre doğru sonuçlar vermektedir.

4.3 Twitter Üzerinden Twint ile Çekilen Tweet Veri Seti ile Alınan Sonuçlar

Twitter üzerinden toplanan tweet'ler 5 farklı anahtar kelimeyle çekilmiştir. Elde edilen veriler "#asgariucet", "#hekimlerhaklariniistiyor", "#sigarazammi", "#karyagiyor" ve "#KilicdaroglundanOnlineEgitim" etiketlerini içeren tweet'ler içermektedir.

Tweet verisi bütün özellikleriyle birlikte geldiğinden, bu özelliklerin çıkarılması için ve veri setinin ön işleme fonksiyonu için hazır hale getirilmesi için birkaç adım uygulanmaktadır. (Figür 15)

```

1 def pre_processing(df, keywords):
2     remove_features = ['id', 'conversation_id', 'created_at', 'date', 'time', 'timezone',
3                         'user_id', 'username', 'name', 'place', 'language', 'mentions',
4                         'urls', 'photos', 'replies_count', 'retweets_count', 'likes_count',
5                         'cashtags', 'link', 'retweet', 'quote_url', 'video',
6                         'thumbnail', 'near', 'geo', 'source', 'user_rt_id', 'user_rt',
7                         'retweet_id', 'reply_to', 'retweet_date', 'translate', 'trans_src',
8                         'trans_dest']
9     df = df.drop(remove_features, axis=1)
10    df['tweet'] = df['tweet'].str.lower()
11
12    filter = []
13    for i in df['hashtags']:
14        if (i.count(',') + 1) <= len(keywords):
15            filter.append(True)
16        else:
17            filter.append(False)
18    df = df[filter]
19    df = df.reset_index(drop=True)
20
21    df['tweet'] = [i.replace(keywords[0] and keywords[1], '') for i in df['tweet']]
22    df['tweet'] = [i.replace(keywords[0] or keywords[1], '') for i in df['tweet']]
23
24    filter = []
25    for i in df['tweet']:
26        if len(i) == 0:
27            filter.append(False)
28        else:
29            filter.append(True)
30    df = df[filter]
31    df = df.reset_index(drop=True)
32
33    return df

```

Figure 15: Preprocessing for Twitter Data

Tweet metinleri vektör haline getirilmeden önce, toplanırken kullanılan etiket veya anahtar kelimeler tweet metninden çıkartılır. Eğer sadece etiket veya anahtar kelimedenden oluşan tweet'ler varsa bunlar veri setinden çıkartılır.

Elde edilen verilerden sigara zammı, hekimler haklarını istiyor ve asgari ücret ile ilgili olanlar kullanılarak bir model oluşturulmuştur. Dengeli bir dağılıma sahip olan 3 sınıf ile eğitim gerçekleştirilmiş ve ortalama entropi **0.9716** olarak bulunmuştur.

Kümelerin sınıf dağılımları ilk veri setindeki kadar başarılı olmamıştır. (Örnek küme dağılımı Figür 16'da görülebilir) Bunun sebebi toplanan tweet'lerin içerdiği konuların birbirine yakın olması olarak yorumlanabilir. Sigara zammı ekonominin kötü olduğundan şikayet eden tweet'ler içerirken, asgari ücret ile ilgili tweet'ler de bu konudan bahsetmektedir. Hekimler ile ilgili tweet'ler ise nispeten daha iyi ayrıştırılabilmektedir.

	index	labels
0	#sigarazammi	3046
1	#asgariucret	1243
2	#hekimlerhaklariniistiyor	945

Figure 16: Distribution of Cluster

5 Demo

Demo için *Next.js* ile yapılmış bir web uygulaması kullanılmıştır. Uygulamanın front-end kısmında kullanıcının tweeti girebileceği bir alan ve tweetin kategorisini öğrenmek için bastığı bir buton bulunmaktadır. Back-end kısmında ise kullanıcının gönderdiği tweet parametre olarak kullanılarak *demo.py* adlı dosya çalıştırılmaktadır. Bu dosya çalıştıktan sonra tweetin kategorisini döndürülür. Demoda "A Benchmark Data for Turkish Text Categorization" verisi ile eğitilen model kullanılmıştır.



Figure 17: Demo Example

6 Referanslar

- TurkishGloVe, Link: <https://github.com/inzva/Turkish-GloVe>
- Turkish Stemmer for Python, Link: <https://github.com/otuncelli/turkish-stemmer-python>
- A Benchmark Data for Turkish Text Categorization, Link: <https://www.kaggle.com/savasy/ttc4900>
- Turkish Tweets Dataset, Link: <https://www.kaggle.com/anil1055/turkish-tweet-dataset/version/1>
- Kemik Doğal Dil İşleme Grubu Yıldız Teknik Üniversitesi, Link: <http://www.kemik.yildiz.edu.tr>