

(Geliş Tarihi / ReceivedDate: 23.01.2020, Kabul Tarihi/ AcceptedDate: 19.04.2020)

Twitter'daki Verilere Metin Madenciliği Yöntemlerinin Uygulanması

Hatice Kübra ÜÇÜKKARTAL*¹

¹Eskişehir Osmangazi Üniversitesi, Mühendislik-Mimarlık Fakültesi, Bilgisayar Mühendisliği Böl., 26480, Eskişehir

Anahtar Kelimeler:

Sosyal Medya,
Metin Madenciliği,
Veri Madenciliği,
Veri Görselleştirme,
Duygu Analizi

Özet: Sosyal medya insanların birbirleriyle iletişimini daha etkileşimli hale getirmiştir. Bu etkileşimler sayesinde çok büyük miktarlarda veriler üretilmektedir. Bu sayede gün geçtikçe artan verilerin işlenmesi ve analiz edilmesine yönelik ihtiyaçlar ortaya çıkmaktadır. Bu çalışmanın amacı, sosyal medya platformu olan twitter üzerinden insanların tartıştığı konular hakkında fikir sahibi olmaktır. İnsanların görüşlerinin sonuçlarını göstermek metin madenciliği konusunda alt yapı oluşturmak ve veri görselleştirme yöntemleriyle de verileri daha anlamlı hale getirmek istenmiştir. Bu çalışmada Twitter üzerinden toplanan verilere metin madenciliği yöntemleri uygulanmıştır. Metin madenciliğini, veri bilimi ve veri görselleştirme araçlarıyla birleştirerek veriler kolay anlaşılır hale getirilmiştir. Duygu analizi ile insanların yaptıkları paylaşımların pozitif, negatif veya nötr olma durumu analiz edilmiştir.

Applying Text Mining Methods to Twitter Data

Keywords:

Social Media,
Text Mining,
Data Mining,
Data Virtualization,
Sentiment Analysis

Abstract: Social media has made people's communication more interactive. Thanks to these interactions, huge amounts of data are produced. In this way, the needs for processing and analyzing increasing data are emerging. The aim of this study is to have an idea about the issues that people discuss via the social media platform twitter. It was aimed to show the results of people's opinions, to create a substructure on text mining and to make the data more meaningful with data visualization methods. In this study, text mining methods were applied to the data collected on Twitter. By combining text mining with data science and data visualization tools, data is made easy to understand. Sentiment analysis was used to analyze whether people's posts were positive, negative or neutral.

1. GİRİŞ

Sosyal medya insanların birbirleriyle iletişimini daha etkileşimli hale getirmiştir. Bu etkileşimler sayesinde çok büyük miktarlarda veriler üretilmektedir. Bu sayede gün geçtikçe artan verilerin işlenmesi ve analiz edilmesine yönelik ihtiyaçlar ortaya çıkmaktadır.

Twitter'da paylaşılan içerikler yapısal olmayan veri türlerindedir. Metin tabanlı olan bu yapısal olmayan veriler analiz etmek için, veri madenciliği algoritmalarında kullanılabilecek yapısal formata dönüştürülmesi gerekir. Bu verileri yapısal veri formatına dönüştürmek ve analiz etmek için metin madenciliği yöntemleri uygulanır.

Metin madenciliği, doğal dil işleme ve veri madenciliğinin birlikte kullanılmasıdır[1]. Doğal dil işleme, Natural Language Processing (NLP) olarak bilinen yapay zeka ve dilbilimin alt kategorisidir. İnsanların konuştuğu dillerin işlenmesi ve kullanılması amacı ile araştırma yapan bilim dalıdır[2]. Veri madenciliği, büyük veri yığınlarını anlamlı ve faydalı bilgiye dönüştürme işlemidir[3].

Bu çalışmanın amacı, sosyal medya platformu olan twitter üzerinden insanların tartıştığı konular hakkında fikir sahibi olmaktır. İnsanların görüşlerinin sonuçlarını göstermek metin madenciliği konusunda alt yapı oluşturmak ve veri görselleştirme yöntemleriyle de verileri daha anlamlı hale getirmek istenmiştir.

*İlgili yazar: Hatice Kübra ÜÇÜKKARTAL, haticekubra26@gmail.com)

2. MATERYAL VE METOT

2.1. Projede Kullanılan Araçlar

Bu projede kullanılacak olan araçlar şu şekildedir;

- Twitter Api
- Python
- Jupyter Notebook

Twitter Api, twitter'dan verileri çekmek için kullanılacak olan arayüzdür.

Python, veri analizi, veri işleme makine öğrenmesi ve veri görselleştirme için kullanılacak olan programlama dilidir.

Jupyter Notebook, twitter'dan çekeceğimiz metinleri analiz edeceğimiz, python kodlarını çalıştıracığımız geliştirme ortamıdır.

Kullanılan kütüphaneler;

- Tweepy (Twitter Api)
- Pandas (Veri analiz aracı)
- Matplotlib (Veri görselleştirme)
- Seaborn (Veri görselleştirme)
- WordCloud (Veri görselleştirme)
- TextBlob (Metin verilerini işleme)
- NLTK (Doğal dil işleme paketi)
- Scikit-learn (Yapay öğrenme)

2.2. Veri Seti

Twitter'dan anlık olarak veri çekme işlemi Twitter API (Twitter Apps, 2019) ile gerçekleştirilmiştir. Twitter API kullanmak için Twitter Application oluşturulmuştur. Tweetlere erişim sağlayabilmek için gerekli izinler (key, secret) alınmıştır. Anahtar kelime olarak kullanıcının girdiği arama terimini içeren tweetler toplanmıştır. Bu işlemler için Python dilinin tweepy kütüphanesi kullanılmıştır. Twitter'dan çıkarılan yaklaşık 10000 İngilizce tweet üzerinde işlem yapılmıştır.

Ham tweetler Python dilinin Pandas kütüphanesi kullanılarak csv dosyasına yazılır. Daha sonra bu csv dosyasındaki verilere temizleme işlemi uygulamak için bazı metin madenciliği yöntemleri uygulanır.

2.3. Metin Madenciliği ile Ön İşleme

Metin madenciliği çalışmaları metni veri kaynağı olarak kabul eden veri madenciliği çalışmasıdır. Diğer bir tanımla metin üzerinden yapılandırılmış veri elde etmeyi amaçlar[12]. Metin Madenciliği, doküman koleksiyonlarının önışlemden geçirilmesi, ara sonuçların saklanması, ara sonuçların analiz edilmesi için çeşitli tekniklerin kullanılması ve nihai sonuçların görselleştirilmesi gibi aşamalardan oluşmaktadır[13].

Twitter'dan toplanan veriler metin tabanlı olduğu için metin madenciliği yöntemleri uygulanarak kullanılabilir formata dönüştürülür.

Tweetler, sözdizimsel olarak iyi oluşturulmamıştır. Bu yüzden bazı ön işlemlerden geçirilmiştir. Bunlar;

- Tokenization
- Stopwords
- Lemmatization
- Twitter metinlerindeki büyük harfler küçük harflere dönüştürülmüştür.
- Metinler, noktalama işaretlerinden arındırılmıştır.
- Metin içeriğindeki URL, hashtag ve kullanıcı isimleri kaldırılmıştır.
- Metinlerdeki sayısal ifadeler kaldırılmıştır.
- TF-IDF (Term Frequency - Inverse Document Frequency) ağırlıklandırma işlemi yapılmıştır.

Tokenization, metinleri istenilen özelliklere göre parçalama işlemidir. Stopwords, bir dilde sık kullanılan etkisiz kelimeleri filtreleme işlemidir. Lemmatization, metindeki kelimeleri köküne indirgeme işlemidir. Twitter metinleri içerisinde '#' ile başlayan hashtag, '@' ile başlayan kullanıcı isimleri, URL formatında web sitesi adresleri, sayısal ifadeler ve noktalama işaretleri barınabilmektedir. Metinler bu karakterlerden temizlenmiştir. Bu işlemler için Python dilinin NLTK kütüphanesi kullanılmıştır.

TF-IDF, bir kelimenin doküman içerisinde ne kadar önemli olduğunu değerlendirmede kullanılan istatistiksel bir ölçüdür. Terim sıklığı (Term Frequency), bir doküman içerisinde geçen terim ağırlıklarını hesaplamak için kullanılan yöntemdir[14]. Ters doküman sıklığı (Inverse Document Frequency), birden fazla dokümanda kelimenin geçme sayısını bularak bu kelimenin terim olup olmadığını bağlaç vb (Stop Words) olduğu anlamaya çalışır. Bunun için Terimin Geçtiği Doküman Sayısı / Doküman Sayısı'nın logaritmasının mutlak değeri alınmaktadır[15]. Bu projede her bir tweet doküman olarak baz alınmıştır. Bu işlemler için Python dilinin Scikit-Learn kütüphanesi kullanılmıştır.

2.4. Veri Görselleştirme

Veri görselleştirme, soyut bilgilerin grafik biçiminde sunulmasını tanımlar. Normalde klasik formatta sunulan karmaşık ve dağınık haldeki verileri, kolay algılanabilir görseller ile rahatça anlaşılır ve yorumlanabilir hale getirmektedir.

Görselleştirme metotlarından biri olan WordCloud, türkçe karşılığı kelime bulutu olarak geçmektedir. Kelime bulutları genellikle kelime sıklıklarına göre analiz sonrasında, farklı renk ve desenlerle oluşturulan metin görselleridir[16].

Twitter'dan çekilen tweetler temizleme ve ön işlemden geçirildikten sonra tweetlerde en çok geçen kelimeler, Python içerisinde bulunan WordCloud kütüphanesi kullanılarak veriler görselleştirilmiştir.

- [14] İnternet: Onur, D. 2015. Term Frequency.
<https://medium.com/algorithms-data-structures/tf-idf-term-frequency-inverse-document-frequency-53feb22a17c6>, (Erişim Tarihi: 22.01.2020)
- [15] İnternet: Onur, D. 2015. Inverse Term Frequency.
<https://medium.com/algorithms-data-structures/tf-idf-term-frequency-inverse-document-frequency-53feb22a17c6>, (Erişim Tarihi: 22.01.2020)
- [16] İnternet: Hüseyin, D. 2014. Kelime Bulutu.
<https://huseyinemirtas.net/kelime-bulut/>, (Erişim Tarihi: 22.01.2020)