

# İstanbul Teknik Üniversitesi Elektronik ve Haberleşme Mühendisliği Bölümü

İsaret İşleme için Makine Öğrenmesi Donem Odevi

2022 Sonbahar

Oğrenci: Giray Arat

Oğretmen: Bilge Günsel

**Proje Adı: Sosyal Medya'dan Alınan Verileri Kullanarak Olay Tanımlama**

## Proje Amacı

Sosyal medya'dan alınan verilerin yakınsaklık vektörlerinin oluşturulması, öbeklenmesi ve sonrasında kelime öğrenme yapılması ile "olay tanımlama" yapmak.

## Proje Referansları

Proje makalesi DOI:

<https://doi.org/10.1007/s10994-021-05988-7>

Proje gerçekleştirildiği içeren GitHub depositesi:

<https://github.com/HHansi/Embed2Detect>

Projede kullanılan sosyal medya verisinin işlenmiş halini içeren GitHub depositesi:

<https://github.com/HHansi/Twitter-Event-Data-2019>

## Projede Kullanılacak Giriş Verisi

[https://github.com/HHansi/Twitter-Event-Data-2019/blob/master/BrexitVote/ids\\_7.30-17.30.txt](https://github.com/HHansi/Twitter-Event-Data-2019/blob/master/BrexitVote/ids_7.30-17.30.txt)

## Projede Yapılacaklar

### Genel bilgi

Kelimeler arasında syntax ile alakalı ve semantic benzerlikler vardır. Syntax benzerlikleri sadece aynı kelime yapısına sahip olmakla alakalıdır. Semantic yani anlamsal benzerlikler ise aynı

anlama gelen farklı söyleyislerin kullanılmasıdır. Makale bunun için verdiği örnekte kendi ürettikleri veriden iki tweet arasında aynı anlama gelen fakat farklı syntax'e sahip iki deyişi seçmiş:

hashtags. In addition, different word phrases such as *impact assessments* and *economic analysis* were used to mention the same subject discussed in them. In such cases, semantics are needed to understand the relationships between terms to extract valuable information.

Bu sebeple verilerin semantic yakınsaklıklar kurularak işlenmesi gerekiyor. Bunun için kullanılabilecek birçok "Word2vec" algoritması var. Word2vec algoritmaları, "text" yani metni vektörlere çevirmek üzere geliştirilmiş bazı algoritmaların adı. Syntax üzerine semantic benzerlikleri de yakalayabileceği için seçilen Skip-gram modeli kullanılıyor makalede. Skip-gram modeli bir denklem.

#### 3.1.1 Skip-gram model

Skip-gram model is a log-linear classifier which is composed by a 3-layer neural network with the objective to predict context/surrounding words of a centre word given a sequence of training words  $w_1, w_2, \dots, w_n$  (Mikolov et al. 2013b). More formally, it focuses on maximizing the average log probability of context words  $w_{k+j} | -m \leq j \leq m, j \neq 0$  of the centre word  $w_k$  by following the objective function in Eq. 1. The length of the training context is represented by  $m$ .

$$J = \frac{1}{n} \sum_{k=1}^n \sum_{-m < j < m, j \neq 0} \log p(w_{k+j} | w_k) \quad (1)$$

## Kelime gömülmesini (Word Embedding) öğrenme

Sonraki adım Skip-gram modeli kullanarak alınan ham veriden sayısal vektörel veri oluşturmak. Benim anladığım kadarıyla bu adımda sadece ham veriden sayısal vektörler oluşturma işlemi yapılıyor. Fakat bunu yapan fonksiyonun "learn word embeddings" diye bir başlığı var. Neden sadece "word embedding" değil de "learn word embeddings" yazıyor bilemiyorum.

## Olay Pencerelerini (Event Windows) öğrenme

Twitter'dan alınan veriler kullanarak olay pencereleri belirleme. Sadece belirli token üzerinden filtrelenmiş ve başka ön-işleme geçirmemiş veriyi alarak ilk önce sayısal değerlere sahip vektörleri oluşturmuşlardı. Şimdi kodun kalan kısmında bu vektörleri kullanarak olay pencerelerini belirleyecekler.

### Olay pencereleri

Bu makine öğrenmesi uygulamasının amacı olan olay tanımlamayı yapmak için vektörel verinin içerisinde bazı olay pencereleri belirlemek gerekiyor. Olay penceresinin ne olduğunu anlamak için makaledeki örneği tekrar edebiliriz. Makalede kullanılan veri dizilerinden bir tanesi bir futbol maçıyla ilgili. Diğer veri seti ise Brexit ile ilgili.

Bir veri seti politika ile ilgili iken diğeri sporla ilgili. Makaleyi yazan doktora öğrencisinin beraber çalıştığı profesörün bununla ilgili makalesinde gösteriyormuş ki, politika ile ilgili verilerin olay penceresi daha uzun ve sporla ilgili verilerin olay penceresi daha kısa olur. Bu demektir ki sporla ilgili olaylar olduğunda kısa bir zaman dilimi içinde çok sık ve çok fazla Tweet atılır. Buna kıyasla politikayla ilgili Tweetler daha uzun bir zaman dilimine yayılır. O zaman Brexit’le ilgili olabilecek olayları tanımlayabilmek için daha uzun olay pencereleri tanımlamak gerekir. Futbol maçıyla ilgili olabilecek olaylar bir kişinin gol atması, birisinin ceza kartı görmesi olabilir. Bu olayları tanımlayabilmek için daha kısa olay pencereleri tanımlamak gerekir.

### Kod’un yaptığı: (Cluster) Öbeklerin Arasındaki Değişimleri Hesaplamak

Kod’un bu bölümünde olay pencerelerini tanımlayacak işlemler yapılır. Önceki adımlarda “word embeddings” üretilmişti. Bu word embeddingsdeki semantic ve syntax ile ilgili değişimlere bakarak, öbekler arasındaki değişimler hesaplanır.

Clustering için hierarchical agglomerative clustering (HAC) seçilmiş. İki cluster arasındaki mesafe için bir sonraki denklemi kullanmışlar:

$$D(C_i, C_j) = \frac{1}{|C_i||C_j|} \sum_{w_p \in C_i} \sum_{w_q \in C_j} d(w_p, w_q),$$

Kelime öbeklerinde tokenler vardı. Tokenler Tweetlerde geçen bir kaç kelime için geçerli, o kelimenin benzer kelimelerinin toplandığı kelimeler. Makale mesela futbol maçından örnek veriyor. Futbol maçıyla ilgili Tweetlerde “goalıı” ve “goalııı” kelimeleri geçiyorsa, bu kelimeler tek bir token altında toplanıyor. Bu iki kelime de “goalıı” sayılıyor.

Her zaman penceresi için bu tokenler arasındaki değişimleri içeren bir matrix kuruluyor. Bu matrislerin tokenleri arasında karar vermek için Dendrogram seviyesinde benzerlik kullanılıyor.

$$DL\ Similarity_{(w_i, w_j)} = \frac{dl_{(w_i, w_j)}}{\max(dl_{r \rightarrow x} : x \in L) + 1}$$

### Kelime Damıtma (Word Extraction) yapma

Önceki adımda olay pencereleri tanımlanmıştı. Eğer bir olay penceresi, içinde olay gerçekleşmiş bir olay penceresi olarak belirlenirse, “event word extractor”, yani olay kelime damıtıcısı, olayla ilgili kelimeleri o pencereden damıtıyor.

Olaylar, metin dağılıcısında değişimler yaratıyor. Yani sosyal medyada akan, yaşayan metin havuzunda, olayların gerçekleşmesiyle beraber ani değişimler oluyor. Bu sebeple önceki olay penceresine göre öbek farkları taşıyan kelimeler, olay kelimeleri olarak seçiliyor. Öbek farkları yani “cluster changes” gösteren kelimeler.

Yukarıda gösterilen  $DL\ Similarity_{(w_i, w_j)}$  değeri sıfırdan farklı olan çiftlere ait olan kelimelerin

hepsi, geçici “temporal” öbekler değişimine sahip kelimeler sayılıyor.

Benim anladığım kadarıyla olay tanımlama bu şekilde gerçekleşiyor. Makine öğrenmesi algoritmaları olay penceresi ve sonrasında olay kelimeleri belirlemede kullanılmış oluyor.