

Çeviri Yazılımlarında Sözcüklerin Bağlam İçindeki Anlamını Algulamaya Yönelik Bir Öneri

Gökhan Silahtaroğlu¹, Fırat Demircan²

¹Beykent Üniversitesi, Bilgisayar Mühendisliği Bölümü

²Beykent Üniversitesi, Fen Bilimleri Enstitüsü

Özet: Türkçe den başka bir dile ya da başka bir dilden Türkçe'ye çeviri yapan yazılımların en büyük sorunlarından bir tanesi iki dil arasında, birçok sözcüğün iki, üç hatta dört farklı şeyi ifade etmesidir. Yazılımların uygun bir çeviri sağlayabilmesi için cümle içinde geçen sözcüğün o cümlede ya da o bağlamda hangi anlamda kullanıldığının yine yazılım tarafından belirlenmesi gerekir. Bunun içinse çeviri yazılımının uygun bir anlamsal çözümleme işlevine sahip olması gerekir.

Bu çalışmada Türkçe sözcüklerin bir birleriyle olan ilişkilerini, hangi sözcüğün hangi sözcüklerle birlikte geçtiğinin haritası bir bakıma anlamsal sözcük grafini çıkartılması ve buna bağlı olarak İngilizce – Türkçe ve Türkçe – İngilizce Sözlüğün bu anlamsal graf'a göre yeniden düzenlenmesi ve uygun bir algoritmayla sözcüklerin bağlam içindeki gerçek anlamlarının ortaya çıkarılması için bir uygulama ve öneri sunulmaktadır.

1. Giriş

Makine çevirilerinde ortaya çıkan bir çok sorun henüz tam anlamıyla çözülmüş değildir. Türkçe'den İngilizce'ye ya da İngilizce'den Türkçe'ye yapılan makine çevirileri henüz tam olarak güvenilir sonuçlar vermemektedirler. Çevirelerde, dil bilgisi, zaman dönüşümü, deyimler, atasözleri gibi hususların yanı sıra her iki dilde de bulunan eş sesli, çok anlamlı sözcüklerin varlığı da makine çevirilerinde üstesinden gelinmesi gereken bir sorundur. Örneğin, Türkçe'de kullanılan *pas* sözcüğü, futbol dilinde *topu takım arkadaşına kazandırma*, *verme* anlamını taşıırken, başka bir anlamı da *metallerin oksitlenmesi sonucu ortaya çıkan bir maddedir*. Cümle içinde bu sözcüğün hangi anlamda kullanıldığını anlamak için bağlamına ya da gerek cümle gerekse de aynı paragraf içinde birlikte geçtiği diğer sözcüklere bakmak gerekecektir. “*Kaleci Demir geri pası tuttu*” ya da “*Demir pas tuttu*” cümlelerinde geçen pas sözcüğünün gerçek anlamı, yanında, ötesinde berisinde geçen sözcüklere bakılarak bulunabilir.

Bu çalışmada, Türkçe'de anlam kargaşasına yol açan sözcüklerin, birlikte geçtiği sözcükler belirlenerek oluşturulan grafa dayanarak yapılacak bir Türkçe – İngilizce sözlük yardımıyla çok anlamlılıktan kayanaklanan makine çeviri sorunlarına bir çözüm önerilmektedir.

2. Metin Madenciliği

Metin madenciliği veri madenciliğinin bir yan uygulama alanıdır. Bilinen şekliyle veri madenciliği sınıflandırma, kümeleme ve birliktelik analizi modellerini inceler [1]. Birliktelik analizi incelenirken veri setinde kullanıcı tarafından belirtilmiş nesnelerin birbirleriyle ilişkisi ya da birlikte olma sıklıkları ortaya çıkartılır[2]. Veri madenciliğinin bu modeli en çok pazar sepeti analizinde uygulama alanı bulmuştur. Pazar sepeti analizinde müşterinin aldığı ürünlerin bir arada bulunma olasılıkları belirlenir ve bu birliktelik

algoritma tarafından hesaplanan güven ve destek seviyeleriyle birlikte değerlendirilir[1],[2].

Metin madenciliğinde ise müşterinin aldığı ürünlerin bir arada bulunma sıklığı yerine sözcüklerin belli koşullar altında bir birini izleme sıklığı ortaya çıkartılır. Bu uygulamada veri seti bir düz dosyadaki metindir. İnternetin yaygınlaşmasıyla birlikte elektronik ortamdaki metinlerin birçoğu artık web sayfalarında yer almaktadır [3]. Metin madenciliğini besleyen metinlerde bu yüzden çoğunlukla web sayfalarından sağlanmaktadır. BU yüzden metin madenciliği zaman zaman İngilizce'de Web Mining (Web Madenciliği) ismiyle de anılmaktadır. Ancak işin özü ve kullanılan algoritmalar değişmemektedir. Metin madenciliğinde de birliktelik analizleri dışında sınıflandırma, kümeleme gibi model uygulamaları yapılabilir.

Herhangi bir kavramın metinde bir kez ya da belirli sayıda geçme olasılığının belirlenmesi bununla ilgili kurallar türetme çalışmaları metin madenciliğinde sınıflandırma çalışmalarına tipik bir örnektir. Örneği biraz daha açarsak, *bulut* sözcüğünün bir metinde üç ve üzerinde geçmesi için gerekli kurallar çeşitli analizlerin yapılması ve metinler üzerinde algoritmaların koşturulmasından sonra şu veya benzer bir şekilde oluşacak ya da algoritma tarafından üretilip kullanıcıya sunulacaktır.

```
IF title INCLUDE bulut AND body INCLUDE bilişim( > 3) AND body INCLUDE veri (>0) THEN TRUE (Support %3, Confidence %78, Lift 236).
```

Yukarıdaki sonuç cümlesi şu şekilde açıklanabilir: *Bulut* sözcüğünün bir metinde üç defa ve üzerinde geçmesi için ilk şartın başlıkta *bulut* sözcüğünün geçmesidir, bunu yanı sıra metin içinde *bilişim* sözcüğünün en az dört defa ve *veri* sözcüğünün en az bir defa geçmesi gerekmektedir. Bu kural tüm metinler göz önüne alındığında veri setinde %3 sıklıkla rastlanılmış ve her bir rastlantıda *bulut* sözcüğünün üç ve üzerinde gerçekleşme olasılığı % 78'dir. Bu kural

ağacın tümü göz önüne alındığında verilere ulaşma açısından 2.36 kat daha etkindir. Metin madenciliğinde kullanılabilir bir başka veri madenciliği modeli kümeleme değildir.

Metin madenciliğinde kümeleme metinlerin içinde barındırdıkları tüm sözcüklere göre kaç ayrı kümeye ayrılması gerektiği konusunu ilgilendirir[1],[4]. Örneğin İnternetteki tüm metinler metin madenciliği ile doğal kümeler ayrılıp daha sonra ortaya çıkartılan her bir kümenin, yapısal ve anlamsal özellikleri belirlenebilir. Daha sonra da bununla ilgili kurallar türetilerek, öğrenme evreni dışından karşımıza çıkacak herhangi bir metnin hangi kümeye, yapısal ya da anlamsal olarak ait olduğu söylenebilir.

Bu çalışmanın da ana konusunu oluşturan üçüncü modelde birliktelik kuralları analizdir[3]. Metin madenciliğinde bu model, hangi şartlarda ve konularda örneğin, hukuk, spor, magazin vs hangi sözcüklerin hangi sözcüklerle birlikte geçtiği ve bu birlikteliğin mutlak ve göreceli gücünü ortaya koyan çalışmadır [5]. Bu yöntemle ortaya çıkartılan kurallar ve örüntüler web semantiği alanında temel kaynak veya kullanılacak olan algoritmanın temel bileşeni olarak kullanılabilir. Bağlantı analiz için kullanılan algoritmalar en yaygın olanları APRIORI, AIS ve SETM algoritmalarıdır [1],[3],[5]. Üç algoritma da en kısa süre içinde nesne kümeleri yaratarak birliktelik kurallarını ortaya çıkarır.

3. Anlam ve Bağlam

Anlam sözcüklerin ilk okunduğu ya da duyulduğu anda zihinde meydana getirdikleri olgudur[6]. Konuya başlamadan önce bir sözcüğün ilk duyulduğu anda bizde ne hissettirdiğini anlamamız gerekmektedir. İnsan beynini bir bilgisayar olarak düşündüğümüz zaman sözcüğü gördüğü ya da duyduğu anda o sözcüğü daha önce bağdaştırdığı bir sürü anlam gelir aklına ama cümleye veya konuya göre o anlamlardan bir tanesini seçer ve o anlama göre karşısındaki kişiye cevap verir [7]. Şimdi düşünecek olursak biz İngilizce bir sözcüğü duyduğumuzda ne yaparız? Arkadaşlarımız ya da yakınlarımızdan birine şu sözcüğün anlamı ne diye sorarız bize cevap verirken birbirinden farklı bir kaç anlam söylediklerini ya da bunları söylemeden önce bize o sözcüğü nereden duyduğumuzu veya nasıl bişer cümle içinde duyduğumuzu sorduklarını görürüz. Peki karşımızdaki kişi bunu niye yapar? Anadili İngilizce olmadığı için mi? Hayır, sözcüğüm anlamını bilmediği ve bunu cümleden yararlanarak bulmak istediği için mi? Hayır, peki niye? Çünkü sözcüğü duyduğu anda beynine birçok anlam gelir ve bunların hangisinin doğru olduğunu bilemez işini şansa bırakmamak için ve karşısındaki kişiye en iyi şekilde yardımcı olabilmek için sorar. Sözcüğün o cümle içinde kullanılan başka bir deyişle o bağlam içindeki anlamını karşısındaki kişiye verir. Türkçe’de ve diğer bütün dillerde de bu böyledir. Eş sesli kelimeler, eş anlamlılar, anlam

genişlemeleri ve anlam daralmaları kelimenin birden çok anlam ifade etmesini veya birden fazla kelimenin sadece bir anlam ifade etmesini sağlarlar[8]. Bu da çevirilerde işimizi bir hayli zorlaştırır. Bunu aynı dili konuşan kişiler fark etmeyebilir ya da önemsemeyebilirler ama çeviri gibi hassa konularda çok önemlidirler.

Konuşma dili yazım dilinden farklı olduğu için yazarak karşımızdaki hissettirdiğimiz şeyi konuşurken farklı şekilde açıklarız[9]. Çünkü yazı dilinde vurgu yoktur ve bunun için çoğu insan yazarken duygularını tam olarak ifade etmediğini söyler. Yazım dilinde kurallar, kullanılacak sözcükler ve vurgular konuşma diline göre çok daha önemlidir ve daha büyük bir yer kaplar. Eş anlamlıların çokluğu zaman kavramları bütünlük derinlik ilişkileri gerçek anlamda önemlidir. Bunun için insanlar okudukları zaman dili daha iyi kullandıklarını görürler. Çünkü birçok kavramı takip etmesi gerekir. Önce sözcüğün gerçek anlamını, daha sonra deyimleri, daha sonra cümleyi ve belki paragrafın tümünü hesaba katarak o sözcüğün gerçekte ne anlam ifade ettiğini görmeye çalışır. Bizim üzerinde çalışacağımız konu işte kelimelerin gerçekte ya da bulunduğu bağlam içinde ne anlam ifade ettiğini bu özelliklerden birkaç tanesini kullanarak bilgisayara öğretmek ve çevirilerin daha rahat, daha sağlıklı olmasını sağlamaktır.

4. Uygulama Adımları ve Öneri

4.1 Amaç

Makine çevirilerinde karşılaşılan temel güçlüklerden bir tanesi iki dil arasında oluşan sözcükler arası çok anlamlılıktır. Örneğin, İngilizce’de **strong** sözcüğü güçlü kuvvetli anlamına geldiği gibi **‘strong tea’** olarak kullanıldığında **‘demli’** anlamını taşımaktadır. Bu nedenle iki dil arasında da makine çevirisi yaparken bu anlam farklılıkları, sözcüğün kullanıldığı bağlam da göz önünde bulundurularak değerlendirilmeli ve çeviri işlemi makine tarafından bu şekilde yürütülmelidir.

Bu çalışmadaki temel çözüm önerisi Türkçe – İngilizce çevirilerde hedef dil İngilizce olduğu durumlarda , iki dil arasındaki sözcükler arasındaki çok anlamlılık ya da anlam kargaşasını gidermektir. Bunun başarılabilmesi için aşağıdaki hususların yerine getirilmesi gerekmektedir:

1. Tercihen her iki dilin de ama en azından hedef dilin anlam haritasının çıkartılması.
2. Bu anlam haritasına uygun Türkçe – İngilizce yeni bir sözcük hazırlanması
3. Makine çeviri programının yeni sözlüğü okuyacak şekilde düzenlenmesi.

bir sözlük prototipi ürettik. Bu sözlükte kelimelerimiz, İngilizce karşılıkları ve onlarla birlikte en çok kullanılmış sözcükler bulunmaktadır. Bu sözcükleri ise kullanılma sıklıklarına göre derecelendirdik. Programın yapacağı şey öncelikle İngilizceye çevrilecek sözcüğün içinde bulunduğu paragrafı bulup, algoritma yardımıyla o paragrafta en çok geçen sözcükleri sıralayıp, daha sonra önceden oluşturduğumuz sözlüğümüzde o kelimenin farklı anlamlarındaki onlarla birlikte geçmiş (daha önce derecelendirdiğimiz) kelimelere bakarak hangi anlamlarındaki sözcükler ile incelediğimiz paragraftaki sözcükler birbirine daha çok uyduğuna bakarak sözcüğüm gerçek anlamını bulunacaktır.

5. Sonuç

Çeviri; farklı dilleri konuşan insanlığın birbirleriyle daha rahat iletişim kurabilmesi için gerekli olup gün geçtikçe gelişen, üzerinde çok fazla araştırma yapılan konuların başında gelir. İnternetin tüm dünyada yaygınlaşmasından sonra, farklı ülkelerdeki farklı dilleri konuşan insanların birbirleriyle daha fazla iletişim kurma istekleri, farklı dillerdeki kaynaklardan yararlanma istekleri çeviri yapan programların yaygınlaşmasını ve yeni çözüm önerileri çıkmasını sağlamıştır. Biz de en önemli konulardan biri olan eş sesli ve anlam daralması yaşamış kelimeler sorununu çözmeye çalıştık. Daha önce yapılan araştırmaları inceleyerek ve onların veri ve sonuçlarından yararlanarak bu sorun için bir algoritma geliştirdik. Öncelikle kelimelerin paragraflar içindeki birbirlerine yakınlıklarından yararlanıp yeni bir tür Türkçe İngilizce sözlük üretilip bu sözlük ile kelimenin kullanıldığı cümle, paragraf ya da başka bir deyişle kullanıldığı bağlam içindeki anlamını bulmaya çalıştık. Bildiğiniz gibi çeviri gibi hassas konular, çok farklı dillerde çok farklı yapılar olduğundan dolayı tüm dünyada aynı hızla ilerlememektedir. Ancak, farklı ülkelerde kendi dillerinin yapısını çözümlenmek için birçok araştırma yürütülmektedir. Kendi dilimizin yapısını çözdüğümüz ve dilimizin yapısı hakkında yeteri kadar fikir sahibi olduğumuz durumlarda çeviri gibi hassas konular daha hızlı ilerleme gösterebilir. Gelişen dünyada yeni teknikler sayesinde gelecekte insanların birbirlerinin dillerini öğrenmek zorunda

kalmadan birbirleriyle çok rahat bir şekilde iletişim kurabilecekleri görebiliyoruz. Bu süreci kısaltmak insanların kendi anadillerine verdikleri özenden ve kendi anadillerinin yapısını ne derece iyi bilmelerinden kaynaklandığı bir gerçektir..

6. Kaynaklar.

[1] Silahtaroglu, Gökhan, Kavram Ve Algoritmalarıyla Temel Veri Madenciliği, , İstanbul, Papatya Yayıncılık, 2008.

[2] Kantardzic, Mehmed, Data Mining Concepts, Models, Methods, And Algorithms, A John Wiley & Sons, Inc., Publication, United States Of America, 2003.

[3] Çelikyay, E. Kübra, Metin Madenciliği Yöntemiyle Türkçe'de En Sık Kullanılan Ve Birbirini Takip Eden Harflerin Analizi Ve Birliktelik Kuralları, Beykent Üniversitesi Fen Bilimleri Enstitüsü , 2010, (Yayınlanmamış Yüksek Lisans Tezi).

[4] Cabena, Peter- Hadjuna, Pablo -Rolf- Verhees, Jaap – Zanası, Alessandro – Hall, Prentice, Discovering Data Mining: From Concept To Implementation, Usa, International Business Machines Corporation, 1998.

[5] Bayer, Harun, Verimadenciliğinde Bir Metin Madenciliği Uygulaması, Beykent Üniversitesi Fen Bilimleri Enstitüsü , 2011, (Yayınlanmamış Yüksek Lisans Tezi).

[6] Türk Dil Kurumu, Anlam Sözcüğünün Türkçe Karşılığı

[7] Condon, John (1998), Sözcüklerin Büyülü Dünyası (Anlambilim Ve İletişim)-, İstanbul, İnsan Yayınları,

[8] Muzaffer Barın, Dinleme Ve Konuşma Becerilerinin Önemi, Elazığ, 1997

[9] Feridun Karaarslan, Konuşma Ve Yazma Eğitiminde Beyin Fırtınası Tekniğinin Etkinliği, Sakarya 2010