

Anlam Belirsizliđi İeren Trke Szcklerin Hesaplamalı Dilbilim Uygulamalarıyla Belirginleřtirmesi

Zeynep Altan
Maltepe niversitesi
Bilgisayar Mhendisliđi Blm
zaltan@maltepe.edu.tr

Zeynep Orhan
Fatih niversitesi
Bilgisayar Mhendisliđi Blm
zorhan@fatih.edu.tr

zet: Birden fazla anlama sahip herhangi bir szcđn tmce ierisindeki yeri nemli olmaksızın, hangi anlamda olduđu ile ilgili genellikle bir belirsizlik mevcuttur. Kiři, anlam belirsizliđi olan bir tmceyi anladığı zaman belirsizliđe neden olan szcđn diđer anlamlarını elemiř, sadece bir anlamını gz nne almıřtır. İnsanođlu anlama iřlemine gerekleřtiren biliřsel bir sisteme sahip olduđu iin, belirsizlik ieren bir tmcenin anlařılması insan dil anlama sisteminde olası anlam kmesi iinden uygun anlamın seilmesi'dir. Bu belirsizlik durumunun szcđn uygulanma řekline gre uygun algoritmaların kullanılması ile zmlenmesi dođal dil iřleme alanındaki alıřmalarla bařlamıř, hesaplamalı dilbilim alıřmalarının yaygınlařması ile daha da nem kazanmıřtır. Szck anlamlarının belirginleřtirilmesi alıřmaları, dilin dođasını inceleyen dilbilim alıřmaları, metin veya konuřmaların evirisinin yapıldığı bilgisayarlı eviri sistemleri, bilgi ynetimi teknolojisi olarak kullanılan bilgi geri getirme ve ıkarımı, İnternet zerinde arama motorlarının tasarımı gibi ok geniř uygulama alanına sahiptir.

Anahtar Szckler: Dođal dil iřleme, hesaplamalı dilbilim, kelime anlamını belirginleřtirme, belirsizlik, btnce.

Summary: Any word including more than one senses which is unimportant its part of speech generally contains ambiguity. When a person comprehends an ambiguous sentence, he or she has eliminated the other meanings of obscure word. Since human beings possess a cognitive system realizing the judgment process, comprehension an ambiguous sentence with human language understanding principle means to choose the appropriate meaning from a probable sense set. This ambiguity situation was partially solved with the initiating of natural language processing studies by applying the appropriate algorithms to the complex words, and became more emphasis as the computational linguistics researches extended. Word sense disambiguation encompasses very broad application areas such as linguistic studies which analyze the features of a natural language, machine translation systems translating the texts or speeches, information retrieval and extraction utilized as the knowledge management technology and the design of research engines on Internet.

Key Words: Natural language processing, computational linguistics, word sense disambiguation, ambiguity, corpus.

1.Giriş

Doğal dil işleme çalışmaları ilk dönemlerinde yapay usun küçük bir uygulama alanı olarak sınıflandırılırken, oldukça kısa bir sürede araştırma konuları genişleyerek tek başına incelenen bir disipline dönüşmüştür. Bu hızlı değişimin nedeni son yıllarda bilginin ön plana çıkarak bilgi toplumu kavramını doğuran bilişim alanındaki önemli gelişmelerdir. İnternet üzerinden bilgiye ulaşmanın ve iletişimin kolaylaşması ile birlikte, erişilecek bilgi miktarı da aynı oranda artmıştır. Öte veri (meta data) olarak adlandırılabilir geniş kapsamlı bilginin düzenlenerek yeniden elde edilmesi, uygun olanlarının çıkarılması, özetlenmesi, hatta belli bir kategoriye göre indekslenmesi doğal diller üzerinde dilbilim çalışmalarının kısmen ya da tamamen otomatikleştirilmesini gerektirmiştir. Ayrıca, farklı toplumların çeşitli ticari ve kültürel ilişkileri çoklu diller arasındaki çeviri sistemleri ile gerçekleşmektedir. Çeviri sistemleri İnternet üzerinden iletişimde de büyük kolaylıklar sağlamaktadır. Birer doğal dil işleme uygulama alanı olan tüm bu dilbilimsel incelemelerde, diğer dillerde olduğu gibi Türkçe için de pek çok farklı anlama sahip sözcüklerin anlamalarının belirginleştirilmesi, çözülmesi gereken önemli dilbilimsel problemlerden biridir. Bu problemin çözümü ise, insan-makine iletişimi arasındaki sorunlar kaldırılabilir ölçüde basitleşecektir. Örneğin, her gün ortalama yirmi milyon sözcük içeren teknik bilginin İnternet ortamına aktarılması ile gittikçe büyüyen bir bilgi okyanusu oluşmaktadır. Kişinin dakikada ortalama bin sözcük okuyabildiği düşünüldüğünde, her gün eklenen bilginin okunabilmesi için günde sekiz saatten bir buçuk ay süre gerekir. Bu süre içerisinde eklenen yeni bilgiler için de beş buçuk yıla gereksinim vardır. Özetle, kişinin elektronik ortamdaki yenilikleri normal koşullarda takip edebilmesi ancak ilave yardımcıları mümkündür (Bird, 1996).

Yapay us, dilbilim, mantık, psikoloji ve bilişsel bilimler gibi farklı bilim dalları arasında yer alan hesaplamalı dilbilim ise, doğal dilleri mantıksal olarak modeller ve hesaplamalarını gerçekleştirir. Hesaplamalı dilbilim, dilin işlenmesini matematiksel olarak ifade etmek üzere anlam ve anlatım arasındaki hesaplama yeteneğini araştırır. Bu özelliği ile dilbilim ve bilişsel bilimin bir parçasıdır. Diğer taraftan hesaplamalı dilbilim, anlam ve biçim arasındaki dönüşümü gerçekleştiren bilgisayar programlarının kullanılması ile mühendislik biliminin bir alt dalı olarak sınıflandırılır. Ayrıca hesaplamalı dilbilimin matematiksel dilbilim ve teorik bilgisayar bilimi ile yakın ilişkisi, anlam ve biçim arasındaki dönüşüm hesaplamalarının çeşitli özelliklerinin incelenmesi gerekliliğini ve bunun genel hesaplama teorisi ile bağlantısını açıklar (Uszkoreit, 2000).

Aklın çalışmasının araştırılması bilişsel bilim ve dilbilim arasındaki ilişki olarak tanımlanabilir. Burada, dinamik ve disiplinler arası bir yaklaşım olarak dil ve düşüncenin kökenlerinin araştırılmasına önem verilir. Bilişsel bilim insanları bir taraftan bilgisayarların düşünüp düşünmediği sorusunu cevaplamaya çalışırken, diğer taraftan öğrenme ve hatırlamanın nasıl gerçekleştiği, çevreyi algılama yetisi, beyin ve us ilişkisi, usun gelişimi gibi zihinsel prosesleri de araştırır. Dilbilimciler ise dilin yapısı, tarihi, felsefe ve psikolojisini inceler. Bu alandaki araştırmalara örnek olarak dillerin özellik ve edinimleri, dillerin gelişimi, zaman içerisindeki değişimleri ve dilin beyindeki örgütlenişi verilebilir.

Doğal dil işleme uygulamalarında da bu yaklaşımlar göz önüne alınarak, her bir uygulamanın ait olduğu disiplinin gözetiminde, farklı teknik ve yaklaşımlarla çözümleme gerçekleştirilir. Araştırma yöntemlerindeki farklılıklara rağmen, bilişsel bilim insanları usun beynin bir fonksiyonu olduğu, düşüncenin bir hesaplama türü olduğu, dil ve bilişin uzmanlaşmış bir dizi işlem ve betimlemelerle anlaşılabilceği şeklinde ortak bir fikre sahiptir. Kelime anlamlarının açıklatırılmasının gerekli olduğu uygulamalarda ise, örneğin özelliklerine göre farklı yaklaşımlar içeren algoritmalarından yararlanılır.

2. Sözcük Anlamının Belirginleştirilmesi

Doğal dil işleme çalışmaları ana ve ara uygulamalar olarak iki gruba ayrılabilir. Ana uygulamalar bilgisayarla çeviri, otomatik özetleme, bilgi çıkarımı, bilginin yeniden eldesi gibi kendi başına bir uygulama oluşturan örneklerdir. Ara uygulamalar ise, tümceyi öğelerine ayırma, çözümleme, biçimbilimsel analiz (sözcük ek ve köklerini bulma), sözcük anlamını belirginleştirme gibi ana uygulamalar için gerekli işlemleri gerçekleştirirler. Bir sözcüğün tümce içinde hangi öğeye karşılık kullanıldığı bilmesi ayırt edici bir özelliktir (Agirre ve diğerleri, 2001). Örneğin, “yüz” sözcüğü “yüz lira” ve “denizde yüz” kullanımlarında ad ve eylem olmasına göre anlamlandırılacaktır. Kök sözcüklerle türemiş sözcükler arasındaki ilişkiler bir başka ara uygulama olarak, sözcüğün anlamının belirginleştirilmesine katkı sağlar. “Git” sözcüğünün “gittik” şeklinde kullanıldığında eylem olduğu, “gideri” sözcüğünün ise ad olarak kullanıldığı biçimbilimsel analiz sonucunda çıkarılacaktır. Farklı anlamları olan sözcüklerden “kara” ise, “karaya çıkmamıza çok az kaldı” ya da “kara kara düşünmek” şeklinde kullanıldığında, sözcüğün hangi anlama geldiğini açık olarak belirler. Bu tür sözcükler yardımcı sözcükler olarak sınıflandırılır. Anlamsal sözcük birliktelikleri ise, “yüz-sayı” birlikteliği ile bir taksonomi, “yüz-deniz” birlikteliği ile durum, “yüz-spor” birlikteliği ile konu tanımlar.

Sözcük anlamının belirginleştirilmesi farklı bir uygulama alanı olarak mesaj anlama, insan-makine iletişimi gibi amacın anlama olduğu uygulamalarda mutlaka gerçekleşmelidir. Ayrıca, amacın anlama olmadığı çalışmalar da belirginleştirmeyi gerektirir (Ide ve Veronis, 1998). Örneğin, bilgisayarlı çeviri sistemlerinde sözcük kaynak dilde belirsiz olabilir veya hedef dile birden fazla şekilde çeviri yapılabilir. “Yüz” sözcüğü kullanıldığı yere göre, İngilizce’ye “swim”, “float”, “skin”, “face”, “surface”, “cheek”, “hundred” gibi farklı şekillerde çevrilebilir. Doğru olan sözcük belirsiz sözcüğün anlamının belirginleştirilmesi ile seçilecektir. Bilgi çıkarımının yapıldığı bir uygulamada önceden belirlenmiş bir anahtar sözcük taranırken, sözcüğün farklı anlamlarını elemek arama sonuçlarının kalitesini arttıracaktır. Örneğin “fare” sözcüğü bilgisayar terimi olarak arandığında, hayvan olarak kullanıldığı anlamın elenmesi çözüme ulaşmayı kolaylaştırır. Sözcük anlamlarını belirginleştirmeye bir başka yeni yaklaşım biçimi ise, metnin içinde incelenen sözcükten önceki ve sonraki sözcüklerin kavramsal olarak sınıflandırılmasıdır. Sözcük kategorilerini oluşturan bu ontolojik sıradüzen aynı zamanda bir anlamsal ağ yapısı oluşturur. Bu tür sözcük anlamını açıklatırma yaklaşımları günümüzde pek çok doğal dil uygulamasında kullanılmaktadır (Mihalcea R. ve diğerleri 2004, Altıntaş ve diğerleri 2005, Ide N. ve diğerleri 1998).

3. Sözcük Anlamını Belirginleştirmede Yararlanılan Kaynaklar

Sözcük anlamlarını belirginleştirmede incelenen dilde düzenlenmiş elektronik sözlüklerden, sözcüklerin kavramsal ilişkilerine göre düzenlendiği ontolojik sözlüklerden, herhangi bir konuda analizi kısmen ve tamamen oluşturulmuş derleme metinlerden yararlanılabilir ya da elektronik sözlük ve ontoloji, elektronik sözlük ve derlem, ontoloji ve derlem gibi farklı kaynakların birlikte kullanımı gerçekleştirilebilir. Derleme metinler terimi derlem ya da bütüncü olarak ta adlandırılmaktadır. Bu konudaki ilk araştırmalar İngilizce için yapılmış olduğu için, bu dilde pek çok kaynağa ulaşmak mümkündür. Özellikle Princeton Üniversitesi Bilişsel Bilimler Laboratuvarı'nda 1985 yılında Prof. A.G. Miller tarafından başlatılan WordNet projesi anlamsal bir sözlük olarak İngilizce sözcükleri eşanlamlılar kümelerinde (synsets) sınıflandırır ve sözcüklerin kısa, genel tanımlamalarını yaparak bu eşanlamlılar kümeleri arasındaki çeşitli anlamsal ilişkileri oluşturur (Fellbaum C., 1998). Burada amaç iki farklı işlevi yerine getirmektir: İlki, sözcüklerin tanımlarının verildiği bir listeyi (dictionary), sözcüklerden kavramları, kavramların özelliklerini ve kavramlar arasındaki ilişkileri (thesaurus-ontology) oluşturmak iken, ikincisi yapay us uygulamalarını ve özellikle otomatik metin analizini desteklemektir. WordNet başlangıcından itibaren ücretsiz olarak kullanılabilir; yeni sürümü olan WordNet 2.1 ise, ad, eylem, sıfat, belirteç olarak sınıflandırılmış toplam 155327 farklı ögenin üst kavram (hyperonym), alt kavram (hyponym), eşanlamlılık (synonym), zıtlıklılık (antonym), parçanın (bölümün-üyenin) bütünü (holonomy) ve bütünün parçası (bölümü -üyesi) gibi çeşitli anlamsal sınıflandırma sonuçlarını vermektedir.

Senseval Projeleri pek çok dilde sözcük anlamlarının açıklanması çalışmalarının yaygınlaşmasına neden olmuştur. İlk Senseval Projesinde (1998) İngilizce, İtalyanca ve Fransızca için çalışma grupları oluşturmuştur. Senseval 2 ise 2001 yılında incelendiği dil sayısını arttırarak Baskça, Çince, Çekçe, Danimarkaca, Hollandaca, İngilizce, Estçe, Japonca, Korece, İspanyolca ve İsveççe dillerinde farklı kategorilerde düzenlenmiştir. Senseval 3 çalışmayı 2004 yılında Barselona'da yapılmış ve kapsamına diğer çalıştaylara ek olarak anlamsal rollerin tanınması, çok dilli açıklamalar, mantıksal biçimler, alt sınıflandırma edinimi gibi konular daha fazla dili kapsayacak şekilde eklenmiştir. Senseval 4 çalışmayı için duyurular yapılmaya başlanmış olup, 2007 yılında gerçekleştirilecektir. Bu çalışmaya Türk dilinin de takım olarak katılımı için öneri verilmiştir.

4. Türkçe Sözcük Anlamını Belirginleştirme Araştırmaları

Sabancı Üniversitesi'nde BalkaNet Projesi'nin bir parçası olarak Türkçe bir kavramsal sözlük hazırlanmıştır. (Bilgin O. ve diğerleri, 2004). Bulgarca, Çekce, Yunanca, Romence, Türkçe ve Sırpça olarak 6 farklı Balkan dilinde uygulanan BalkaNet projesi temel olarak Princeton WorldNet modelini kullanmıştır. BalkaNet projesi için kurulan

konsorsiyum projenin ilk aşamasında EuroWordNet¹ projesinin 1310 temel kavramını her bir çalışma takımının diline çevirmiştir. Bu kavramlar sıradüzendeki düzey sayısının yüksekliği ve pek çok alt kavram içermesi nedeni ile tüm dillerde oldukça önemli bir yapı taşı olmuştur. Birinci aşama Türkçe için eş anlamların, zıt anlamların ve alt kavramların elektronik Türkçe dilbilgisi sözlüğünden otomatik çıkarımı şeklinde gerçekleşmiştir. Daha sonra konsorsiyum incelenecek kavramların sayısının beşbine çıkarılmasını kararlaştırmış; böylece Türkçe dahil tüm takımlar bütüncü sıklıkları (corpus frequencies), sözcük dağılımının tanımlanması, tek dilli sözlükler, çoklu anlamlar (polysemy) gibi farklı kriterleri de ekleyerek alt kümelerini genişletmişlerdir.

Bir doğal dil işleme alanı olarak bilgi-tabanlı tekniklerle olası modellerin bütünlüğü, veri tabanı sorgulamalarıyla sınırlı dil uygulamalarını zenginleştirmiştir. Böylece metinlere uygulanan istatistiksel yöntemlerle en olası yorumun tahmini mümkün olmaktadır. Bunun için de ayrıntılı olarak işlenmiş derleme metinlere (bütüncü) gereksinim vardır. Metin örnekleri kullanarak birden fazla anlama sahip kelimelerin, özellikle eylem türündeki kelimelerin anlamlarını çıkarabilmek için, bu metinler üzerinde sözcüksel ve anlamsal bilginin doğru olarak işaretlenmiş olması önemlidir. Örneğin Türkçe için her biri yaklaşık 25000 sözcükten oluşan 7 farklı metin koleksiyonunun² bi-gram model üzerinde test edildiği bir çalışma yapılmıştır (Altan Z., Yanık E., 2001). Burada tümcelerın sözdizimsel ve anlamsal sınıflandırmasında sadece incelenmek istenen sözcükten önceki sözcük işaretlenmiştir. Bütüncü üzerinde uygulanan olası dil modeli, elle tanımlanan kurallara ek bir öğrenme bileşeni olarak en olası çözümü tahmin edebilmekte ve dili işlemedeki belirsizlikleri de büyük ölçüde azaltmaktadır. Bu çalışmada eylemlere ait kavramsal sınıflandırma yol alma, yönelme ve terk etmeden biri şeklinde devinim (motion), kavrama (perception), duygu (emotion), fonksiyon (bodily care and functions), bağlantı (contact) gibi WordNet'in eylemler için grupladığı kavramsal özelliklerinden yararlanarak gerçekleştirilmiş ve eylemlerin sözcük anlamları bunlara göre numaralandırılmıştır. Tümcelerın işaretlenmeleri "git" eylemi için Tablo 1'de görüldüğü gibidir. Artık herhangi bir eylemin anlamı olasıya bağlı olarak tahmin edilebilir. Tahmin için kullanılacak yöntem en olası maksimumun kestirimi (Maximum Likelihood Estimation- MLE) olabilir. MLE, işlenmiş bütüncü içinde aranan sözcüğün eğitilme sayısını hesaplar. Eğitim sadece bir önceki sözcüğe göre yapıldığı için araştırılan kelimedenden önceki kelimenin öğelerine ayrılmış olması önemlidir. Bu sınıflandırmadan elde edilen değerler, işaretlenmiş bu bütüncü üzerinde farklı anlamların belirlenmesi için bir Bayes sınıflandırması oluştururlar. Bu bütüncü tümce öğeleri elle işaretlenerek elde edilmiştir.

Benzeri problemler Türkçe dil işleme çalışmalarının pek çoğunda mevcuttur. Fakat ODTÜ derleme metninin kullanıma açılması ile birlikte gerek sözdizimsel, gerekse biçimbirimsel olarak çözümlenmiş; bütüncü bulma problemi kısmen de olsa çözümlenmiştir.

¹ EuroNet Projesi 1996 yılında başlayıp üç yıl süren WordNet'in Avrupa dillerine uyarlanması şeklinde bir konsorsiyumun gerçekleştirdiği ortak bir çalışmadır

² Dünya klasiklerinden örnek hikayeler: Guliver Devler Ülkesinde , Candide , Ivan Nikiforoviç, Tours Papazı , Mozart Prag Yolunda , Mektuplar, Kır Atlı

Tablo 1: “git” eylemi için örnek işaretleme

No	Tümce
1	ülkeyi keşfetmek için yazar da [birlikte] (DuZ) \$gidiyor\$ {1} ve karada kalıyor
2	gemi, [Suratya] (YeZ) \$gidiyordu\$ {1}
3	Umut Burnuna kadar [rüzgâr] (ÖZ) [çok iyi] (DuZ) \$gitti\$ {3}
4	Biraz daha kuzeye dönerek Tataristanın kuzeybatısına ve [Buzdenizine] (YeZ) \$gitmek\$ {1} [olasılığı karşısında] (ZaZ), bulunduğumuz rotayı izlemenin daha iyi olacağını düşündük
5	[merakımı] (KEyl) \$giderecek\$ {5} bir şey de göremediğimden
6	[olanca hızımla] (DuZ) [o önce] (ZaZ) \$gittiğim\$ {1} [yana] (YeZ) koşmuştum
7	[sesim ve işaretlerim] (ÖZ) [hoşuna] (KEyl) \$gitmiş\$ {4} [gibiydi] (Eyl)
8	Fakat nasıl davrandığımı, kocasının işaretlerine göre ne kadar iyi davrandığımı görünce bana alıştı ve \$gitgide\$ {12} [artan bir sevgi] (DoT) beslemeye başladı
9	Bu yaptığım [pek] (MiZ) [hoşlarına] (KEyl) \$gitmişti\$ {4}
10	Ben de \$gittim\$ {1}, elini öptüm

Bu bütünce, ODTÜ-BAP ve TÜBİTAK tarafından desteklenmiş ve ODTÜ-Sabancı Üniversiteleri işbirliği ile gerçekleştirilmiştir. Çalışmada bir ana derleme metin oluşturulmuş; ayrıca farklı kullanımlar için bu ana derleme metinden bazı farklı özellikleri olan bir de ağaç bankası derleme metni geliştirilmiştir (Ofлаzer ve diğerleri, 2003). Derlemede kullanılan metinler 1990 yılı sonrası basılan eserlerden seçilmiştir. Derlemede yaklaşık olarak 2.000.000 sözcük bulunmaktadır. 201 kitap, 87 makale ve 3 tane günlük gazeteden seçilmiş haberlerden oluşan 999 farklı yazılı metin kullanılmıştır. Derlemede bulunan metinlerin çoğunluğu biçimbirimsel olarak çözümlenmiştir. Fakat yapısal belirsizlikler tamamen çözümlenmemiş olduğu için kullanımda bazı problemlerle karşılaşmaktadır.

5 . Sözcük Anlamını Belirginleştirmede Kullanılan Yöntemler

Sözcük anlamının belirginleştirilmesinde en etkili çalışmalardan biri, olasılı başka deyişle bilgisayarla öğrenme (machine learning) algoritmalarının kullanılmasıdır. Bu yaklaşımlardan öngörmeli (supervised) yaklaşımların, öngörmesiz (unsupervised) yaklaşımlardan daha iyi sonuçlar verdiği gözlenmiştir³.

Standart bilgisayarla öğrenme algoritmalarından pek çoğu öngörmeli öğrenme yaklaşımında kullanılabilir. Sonuçlar genellikle başarılı olmasına rağmen, öngörmeli

³ Öngörmeli öğrenmede çalışılan verinin her parçasının gerçek durumu bilinmesine rağmen, öngörmesiz öğrenmede eğitim örneği içindeki verinin sınıflandırılması bilinmez. Öngörmesiz öğrenme çoğunlukla kümelendirme olarak bilinirken, öngörmeli öğrenme sınıflandırma olarak adlandırılır

yöntemler anlamsal olarak işaretlenmiş derlemlerin az olmasından veya hiç olmamasından dolayı dezavantaj oluşturlar.

Tablo 1’de küçük bir örneği verilmiş olan bütüncü, sözcük anlamını belirginleştirmek üzere kullanılmış ve öngörmeli öğrenme gerçekleştirilmiştir. Sonuçlar Tablo 2’de istatistiksel bir yaklaşım olan Naïve Bayes (NB) ve örnek tabanlı Exemplar Based algoritmalarının uygulamaları olarak görülebilir (Altan ve Orhan, 2003). Tablo 3 yöntemlerde kullanılan özellikler açıklamaktadır. Küçük bir bütüncü üzerinde uygulanan bu algoritmalarından örnek tabanlı yaklaşım, Naïve Bayes’e göre biraz daha iyi sonuç vermiştir .

Tablo 2: Bütünceden çıkarılan özelliklere bir örnek

No	L1P	L2P	R1P	L1M	L2M	R1M	Git Kok	Git Ek	Anlam	CL	CR	NB	FB
1	duz	-	-	-	-	-	gidiyor	gidiyor	1	birlikte	-	1	1
2	yez	-	-	e	-	-	gidiyordu	gidiyordu	1	surat	-	1	1
3	duz	öz	-	-	-	-	gitti	gitti	3	iyi	-	3	3
4	yez	-	zaz	e	-	de	gitmek	gitmek	1	buzdenizi	olasılık	1	1
5	keyl	-	-	i	-	-	giderecek	giderecek	5	merak	-	5	5
6	zaz	duz	yez	-	-	e	gittiğim	gittiğim	1	önce	yan	1	1
7	keyl	öz	eyl	e	-	-	gitmiş	gitmiş	4	hoş	gibi	4	4
8	-	-	dt	-	-	-	gitgide	gitgide	12	-	artan	1	12
9	keyl	miz	-	e	-	-	gitmişti	gitmişti	4	hoş	-	4	4
10	-	-	-	-	-	-	gittim	gittim	1	-	-	1	1

Tablo 3: Öngörmeli eğitimin yapıldığı algoritmaların özellikleri

Özellik	Açıklama
L1P	Hedef sözcüğün solundaki birinci tümce ögesi
L2P	Hedef sözcüğün solundaki ikinci tümce ögesi
R1P	Hedef sözcüğün sağındaki birinci tümce ögesi
L1M	Hedef sözcüğün solundaki birinci tümce ögesinin hal eki
L2M	Hedef sözcüğün solundaki ikinci tümce ögesinin hal eki
R1M	Hedef sözcüğün sağındaki birinci tümce ögesinin hal eki
GK	Git sözcüğünün kökü
GM	Git sözcüğünün eki
CL	Soldan birinci kalıp sözcük
CR	Sağdan birinci kalıp sözcük

İkinci bir Türkçe sözcük anlamını belirginleştirme çalışması, Bölüm 4’de açıklandığı gibi, ODTÜ-Sabancı Türkçe ağaç bankası kullanılarak gerçekleştirilmiştir. Bütüncü olarak ODTÜ derleme metinlerinin kullanıldığı çalışmada, anlamları belirginleştirilecek sözcük sayısı ve tipi arttırılmış; bu sözcükler için yeni anlam sınıfları eklenmiştir (Tablo 4).

Tablo 4: ODTÜ ağaç bankası metinlerinden seçilen sözcüklerin anlam sayıları

Sözcük	Metinlerdeki tümce sayısı	Anlam sayısı
Yan	104	9
Git	189	10
Gör	133	9
Çık	231	15
Al	250	10
Gel	281	12
Yap	328	6
Ol	941	4

6 Sonuç

Türkçe sözcüklerin ortalama anlam sayısı bu alanda üzerinde çok fazla çalışma İngilizce gibi dillere göre çok daha fazladır. Ayrıca tüm gereksinmelere cevap verebilecek Türkçe bir bütüncenin de olmaması çalışmaları daha da güçleştirmektedir.

Kaynaklar

Agirre, E., Ansa, O., Martinez, D., Hovy, E., 2001, Enriching Wordnet Concepts with Topic Signatures, Proc. of the NAACL Workshop on Wordnet and other Lexical Resources: Applications, Extensions And Customizations, Pittsburg, USA, 123-132.

Altan, Z., Orhan Z., 2003, Disambiguation of Turkish Word Senses by Supervised Statistical Methods, International Journal of Computational Intelligence, Vol:1, 16-21

Altan Z. ve Yanık E., 2001 Kelime Anlamlarının İstatistiksel Çıkarımı için Metin örneklerinin İşlenmesi, İstanbul üniversitesi Elektrik& Elektronik Dergisi 1-2, 287-295

Altıntaş E., Karşlıgil E., Coskun, V., New Semantic Similarity Measure Evaluated In Word Sense Disambiguation , 15th Nordic Conf. of Computational Linguistics, 2005

Bilgin O., Çetinoğlu Ö., Oflazer K., 2004, Building a Wordnet for Turkish, Romanian Journal of Information Science and Technology Vol: 7, Num: 1-2, 163-172.

Bird, M., 1996, System Overload. Excess Information Is Clogging The Pipes Of Commerce - And Making People Ill, In Time Magazine, December 9th, 1996, 46-47.

Fellbaum, C., 1998, WordNet: An Electronic Lexical Database, The MIT Press

Ides N., Véronis J., Word Sense Disambiguation: The State of the Art , Computational Linguistics, 1998, 24(1)

Mihalcea R., Tarau P., Figa E., 2004, Pagerank on Semantic Networks with Application Toward Sense Disambiguation, The 20th International Conference on Computational Linguistics, 1126-1132.

Oflazer, K., Say, B., Tur, D. Z. H., Tur, G., 2003, Building A Turkish Treebank, Invited Chapter In Building and Exploiting Syntactically-Annotated Corpora, Anne Abeille Editor, Kluwer Academic Publishers, 2003

Uszkoreit H., 2000, Language Technology for Knowledge Management, Proceedings of Japanese-German Workshop Comp. Linguistics, Yokohama, 26 May 2000, 1-10.