

TÜRKÇE METİN ÖZETLEMEDE MELEZ MODELLEME

(A HYBRID MODELLING FOR TURKISH TEXT SUMMARIZATION)

Akif HATİPOĞLU¹ Sevinç İlhan OMURCA²

ÖZ

Orijinal belgelere ait en önemli cümlelerin belirlenmesi için gerçekleştirilen bilgisayar programı ile otomatik metin özetleme işlemi bir doğal dil işleme problemidir. Doğal dil işlemede temel olarak iki çeşit metin özetleme yaklaşımı bulunmaktadır. Bu yaklaşımlar cümle seçerek özetleme ve yorumlayarak özetleme olarak ikiye ayrılmaktadır. Cümle seçerek özetleme iki farklı alt yöntemle ayrılmaktadır. Birincisi özetlenecek metinde cümleleri istatistiksel olarak puanlandırma yöntemidir. İkinci yöntem ise sezgisel olarak gizli anlam çıkarımı yöntemidir. Özetleme çalışmalarında bu iki yöntemi birleştirip melez bir yapı kurularak özetleme gerçekleştirilmektedir. Bu makale kapsamında cümle seçerek özetleme hedeflenmiştir. Türkçenin yapısal özelliklerine göre istatistiksel olarak puanlandırılması ve gizli anlam çıkarım yöntemlerini sezgisel olarak birleştirerek cümle seçimi yapan melez bir model sunulmuştur.

Anahtar Kelimeler: Metin özetleme, Gizli Anlam Analizi, Doğal dil işleme

ABSTRACT

Automatic text summarization with a computer program in order to retain the most important sentences of the original document is a natural language processing problem. There are basically two types of text summarization approaches in natural language processing. These approaches are defined as summarization by selecting sentences and summarization by interpreting sentences. The summarization by selecting sentences method is also separated into two sub-methods. The first one is the method of scoring the sentences statistically. The second one is the method of latent semantic analysis of the sentences. In this study, summarization is realized by selecting sentences. A hybrid model which selects the sentences by combining two methods: statistically scoring sentences due to structural features of Turkish and latent semantic analysis method.

Keywords: Text summarization, Latent semantic analysis, Natural language processing

¹ Kocaeli Üniversitesi, Mühendislik Fakültesi, Bilgisayar Bölümü, Kocaeli, akifhatipoglu.tech@gmail.com (sorumlu yazar)

² Kocaeli Üniversitesi, Mühendislik Fakültesi, Bilgisayar Bölümü, Kocaeli, silhan@kocaeli.edu.tr

1. GİRİŞ

İnternet kullanımının artması ile birlikte bilgi kaynakları da çoğalmıştır. Günümüzde bir konu hakkında araştırma yapıldığında bu konu ile ilgili farklı birçok kaynağa kolayca ulaşılmaktadır. Doğal olarak bu kaynakların hepsini inceleme fırsatı yakalanmamaktadır. Bilgi kaynakları arttıkça incelenmek üzere depolanan metinlerin sayısı ve uzunlukları da artmaktadır. Bu açıdan, metin özetleme ile incelenmesi zaman alan uzun metinleri daha anlamlı kısa metinler haline getirmek büyük önem taşımaktadır. Tüm bunların sonucunda metin özetleme çalışmaları günümüzde büyük önem kazanmıştır.

Literatürde bilgisayar yardımı ile metin özetleme işlemine ait yaygın olarak uygulanan iki farklı yaklaşım bulunmaktadır. Birincisi cümleyi yorumlayarak özetleme, ikincisi ise seçerek özetleme yaklaşımıdır. Yorumlayarak özetleme, metni kelime bazında anlamlandırıp yeni cümleler türetmeye dayalı bir yaklaşımdır. Bu yaklaşımda cümleler anlamlı olarak birleştirilerek yeni bir metin oluşturulmakta veya kelimeler anlamlı olarak birleştirilerek yeni cümleler oluşturulmaktadır. Seçerek özetleme yaklaşımı ise iki farklı alt yöntemden oluşmaktadır. Birincisi özetlenecek metinde cümleleri istatistiksel olarak puanlandırma yoluyla cümle seçimidir. İstatistiksel puanlandırma, cümlelerin doğal dil açısından önemine göre cümleye puan atamaya dayanmaktadır. İstatistiksel puanlandırma çeşitleri şu şekilde açıklanmaktadır. Bir cümlede bulunan kelime sayısı o cümle için bir puan oluşturmaktadır. Aynı şekilde bir cümlede geçen kelime sayısı içerdiği bilgi açısından bir puan oluşturmaktadır. Bir cümlede geçen kelime sayısı o metnin başlığında geçen kelimeler bulunuyor ise o cümle özet cümle olmak için daha güçlü aday kabul edilip belirli bir puan almaktadır. Son olarak cümlede geçen kelimelerin frekansı hesaplanmaktadır. Bu frekans, kelimenin hem cümle içerisinde geçme hem de metin içerisinde kullanım sıklığı ile belirli bir değeri oluşturmaktadır. Bu frekans değerleri ile belge terim matrisi kurulup her cümle için belirli bir puan oluşturulmaktadır [1].

Cümle seçimi ile metin özetlemede ikinci yöntem ise sezgisel olarak gizli anlam çıkarımıdır ya da gizli anlam analizi olarak da adlandırılmıştır. Bu yöntem birliktelik kuralları ile fark edilemeyen bağlantıları çözümleyip gözlemlememizi sağlamaktadır. Bu yöntemde öncelikle kelime cümle matrisi kurulur. Bu matris, kelimelerin cümlelere göre frekanslarından oluşmaktadır. Bu kurulan matris üzerinde tekil değer ayrıştırılması yöntemi uygulanarak daha küçük yönetilebilir bir matris elde edilir. Bu matris bize metindeki kavramları vermektedir. Bu kavramlara göre cümlelere puan verilir.

Otomatik metin özetleme literatürde sık çalışılmış bir alan olmakla birlikte Türkçe metinler üzerinde gerçekleştirilen çok fazla çalışma bulunmamaktadır [2, 3, 4].

Bu çalışmada, Türkçe metinler üzerinde, cümlelerin yapısal ve anlamsal özelliklerini melez şekilde birleştirerek otomatik metin özetlemeyi gerçekleştiren bir model geliştirilmiştir [2]. Cümlelerin uzunluğu, belge içindeki pozisyonu, dokumana ait başlık kelimelerini içerme sıklığı ve içerdiği terimlerin ağırlıkları cümlelerin yapısal özellikleri olarak belirlenmiş; cümlelerin anlamsal özellikleri ise LSA temelli özellikleri ile belirlenmiştir. Özetlenecek metine ait her cümle için yapısal ve anlamsal özellikler skorlanmış ve sezgisel olarak belirlenen özellik katsayıları ile çarpılarak lineer bir model ile birleştirilmiştir. Uygulanan bu melez yöntemin sonucunda en yüksek puan alan cümleden en düşük puan alan cümleye doğru sıralanmış bir liste elde edilmektedir. Sıralanmış cümle listesi içerisinden kullanıcı tarafından belirlenen sayıda cümle seçilip kullanıcıya özet cümleler olarak sunulmaktadır. Deneysel çalışmalar

kısımında geliştirilen sistemi test etmek üzere Türkçe Vikipedi metinleri kullanılmış; seçilen Wikipedia belgesinin özet metni gösterilmiştir.

Liteartürde otomatik metin özetleme sistemlerinin geçerlilik testinin güç bir işlem olduğu yönünde değerlendirmeler yapılmaktadır [5]. Bu çalışmada modelin performansı, özet metin için seçilen cümleler ile 10 ayrı kişinin seçtiği cümleler karşılaştırılarak değerlendirilmiştir.

Makalenin geri kalanı şu şekilde organize edilmiştir: Bölüm 2’de uygulanan metin özetleme yöntemi açıklanmaktadır, bölüm 3’de deneysel sonuçlar raporlanmakta ve bölüm 4’de çalışmadan elde edilen sonuçlara yer verilmektedir.

2. METİN ÖZETLEME

İlk olarak özetlenilecek metin üzerinde veri ön işleme adımları gerçekleştirilir. Bu adımlar gereksiz kelimelerin elenmesi, metin içerisindeki cümlelerin ayrıştırılması ve son olarak kelimelerin köklerinin bulunmasıdır. Veri ön işleme adımları gerçekleştirildikten sonra Yapısal olarak istatistiksel puanlandırma yöntemleri ve sonrasında gizli anlamsal analiz (GAA) aşamaları gerçekleştirilmiştir.

2.1. Metin Ön işleme

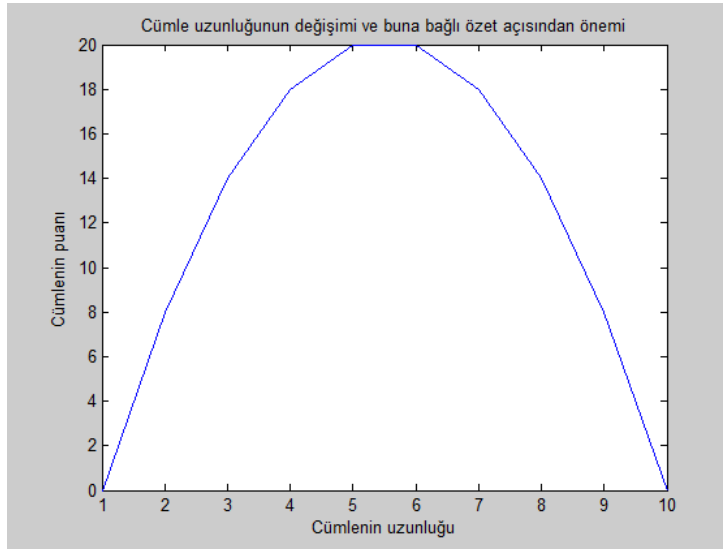
Ön işlemenin ilk adımı olarak metin içerisindeki gereksiz kelimeler elenir. Bunlar genel olarak o dilde sık kullanılan kelimeler, içerikten bağımsız kelimeler, bağlaçlar, imleçler, sayılar, kalıplaşmış kısaltmalardır. Ayrıca noktalama işaretleri ve tek başına duran alfabedeki harfler de gereksiz kelimeler listesine eklenebilir [6, 7].

Sonraki adımda, metin içerisindeki cümleler ayrıştırılır. Bu işlem Zemberek kütüphanesi yardımı ile gerçekleştirilmiştir [8]. Zemberek kütüphanesi açık kaynak kodlu Türkçe doğal dil işleme kütüphanesidir. Platform bağımsız olması ve Türkçe ile ilgili problemlere çözüm getirmesi amacı taşıyan bir projedir. Zemberek dil hatalarına uyarı vermek, ön izleme, düzeltme, heceleme, köklerine ayırma, köklerini inceleme, deyimleri belirleme, ascii dönüşümleri gibi birçok işlevi olan Türkçe ile ilgili en detaylı projedir. Zemberek kütüphanesinin “Sentence boundary detector” fonksiyonu sayesinde bir metin içerisindeki tüm cümleler ayrıştırılabilmektedir. Ancak ayrıştırılan cümlelerde noktalama işaretlerinden kaynaklanan hatalar olmaktadır. Örneğin “.” karakteri cümle sonu gibi algılanarak Zemberek tarafından yanlış cümle ayrıştırmalarına denem olmaktadır. Zemberekten kaynaklanan bu ve benzeri eksiklikler geliştirilen kod ile düzeltilmiştir.

Cümleler ayrıştırıldıktan sonra cümlelerde geçen her bir kelime köklerine ayrıştırılır. Türkçe Ural–Altay dil ailesinden gelmektedir ve sondan eklemeli bir dildir. Türkçede bir kelimenin kökünden sonra farklı ekler gelerek kökten bağımsız farklı anlamlar elde edilebilir. Bu nedenle metin özetleme işlemi gerçekleştirilirken bütün kelimelerin kelime köklerinin bulunması gerekmektedir. Doğal dil işleme işlemlerinin gerçekleştirilmesi için gereken kelime-cümle matrisi kurulurken ve bu matrise bağlı işlemler yapılırken bütün kelimelerin köklerinin bulunmuş olması gerekmektedir. Tüm kelimelerin köklerine ayrıştırılması aşaması Zemberek kütüphanesi yardımı ile yapılmıştır. Her cümleye ait kelime kökleri elde edildikten sonra cümleler üzerinde puanlandırma aşaması gerçekleştirilmiştir.

2.2. Cümle Uzunluğu

Cümlelerin uzunluğunu temsil eden özellik vektörü bir cümle içinde barındırdığı kelime sayısını ifade etmektedir. Veri ön işleme aşamasından sonra cümle içinde cümleye ait kelimelerin kökleri yer almaktadır. Kelime köklerinin sayısı o cümleye ait puanlandırmada kullanılmaktadır [1, 2]. Metin içerisindeki tüm cümleler için ortalama cümle uzunluğu hesaplanır. Bu ortalama cümle uzunluğuna yakın uzunlukta olan cümleler daha önemli cümleler olarak kabul edilir. Ortalamadan daha uzun cümlelerin ayrıntı içeren cümleler olması, daha kısa olan cümlelerin ise yeterli bilgi taşıyor olması olasılıkları bu cümleleri özet metin için daha zayıf adaylar haline getirmektedir.



Şekil 1. Cümle uzunluğunun özet açısından önemi

Cümlelerin uzunluğunun cümle puanına etkisi Şekil 1'de gösterilmiştir. Bu kabulü şu şekilde formüle edebiliriz. Cümle uzunlukları bulunduğundan sonra cümle uzunluklarının ortalaması bulunur. Bu ortalamaya sahip ya da bu ortalamaya yakın cümleler metin özetleme için önemlidir. Çok kısa cümleler veya çok uzun cümleler metin özetleme açısından daha az öneme sahiptir.

$D = \{S^{(1)}, S^{(2)}, \dots, S^{(m)}\}$ belgesinde geçen bir $S^{(i)}$ cümlesi için "cümle uzunluğu" özelliği Eşitlik 1'deki gibi formüle edilmiştir.

$$S^{(i)}_{\text{puan}_{f_{i1}}} = (\text{uzunluk}_{S^{(i)}} - 1/2m (\sum_{j=1}^m \text{luzunluk}_{S^{(j)}}) - 1/2)^2 \quad (1)$$

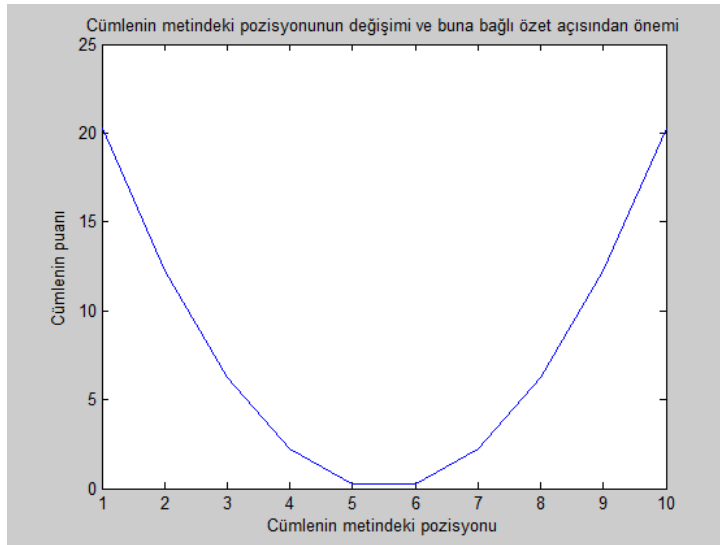
$$S^{(i)}_{\text{puan}_{f_{i1}}} = -1 \times S^{(i)}_{\text{puan}_{f_{i1}}} + \max_{S^{(i)}_{\text{max}_{f_{i1}}}} \{S^{(1)}_{\text{puan}_{f_{i1}}}, S^{(2)}_{\text{puan}_{f_{i1}}}, \dots, S^{(m)}_{\text{puan}_{f_{i1}}}\} \quad (2)$$

Burada $S^{(i)}_{\text{puan}_{f_{i1}}}$ $S^{(i)}$ cümlesine ait uzunluk özelliğinin hesaplanan değeridir. m belgede geçen cümle sayısını ifade etmektedir.

2.3. Cümlelerin Metindeki Pozisyonu

Cümlelerin pozisyonunu temsil eden özellik vektörü bir cümlelerin metin içerisindeki pozisyonunu ifade etmektedir. Türkçe metinler giriş, gelişme ve sonuç bölümlerinden oluşmaktadır. Verilen bir konuyu açıklayabilmek için önce açıklanması gereken düşünce bulunur ve bu düşünce giriş bölümünde belirtilir. Giriş bölümünde sadece açıklanması gereken düşünce kısa ve öz bilgiler ile belirtilir; olay örneklenmez veya açıklama yapılmaz. Gelişme bölümünde, giriş bölümünde belirtilen düşünce daha genişletilmiş ve gerekirse örneklerle zenginleştirilmiş şekilde açıklanır. Gelişme bölümü giriş bölümünden daha uzun yazılır, konu ayrıntıları ile anlatılır. Sonuç bölümünde, giriş ve gelişmenin ortak düşüncesi yani ana düşünce yazılır. Bu bölüm yazılı anlatımın diğer bir kısa bölümüdür, sadece çıkarılan ana düşünceyi anlatır.

Bu çalışmada cümleye ait pozisyon özelliği, her cümlelerin özetlenmesi gereken metnin başlangıcına olan uzaklığını temsil etmektedir. Cümlelerin başlangıcına olan uzaklığı değiştikçe cümlelerin ağırlığı da değişmektedir. Türkçede metinlerin giriş-gelişme-sonuç şeklinde yazıldığı düşünüldüğünde; çalışma kapsamında cümlelerin pozisyonu metnin başlangıcına yaklaştıkça ya da metnin sonuna yaklaştıkça özet açısından önemi artmaktadır şeklinde bir varsayım ile hareket edilmiştir.



Şekil 2. Cümle pozisyonunun özet açısından önemi

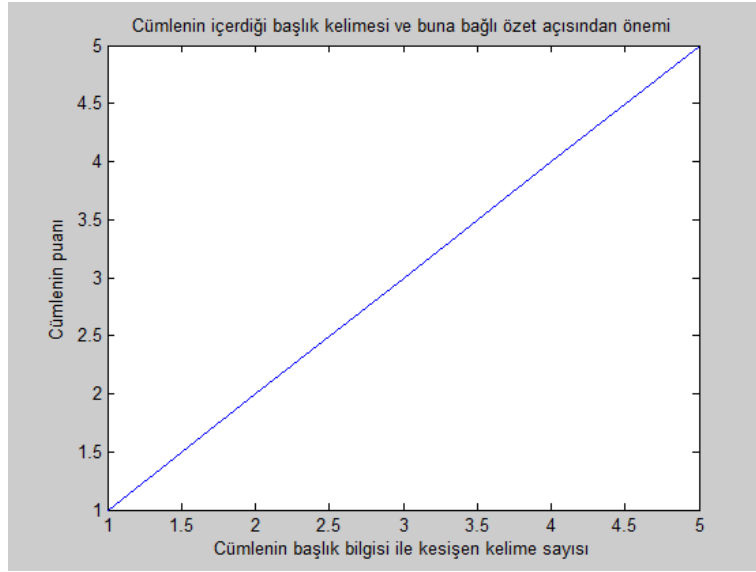
Şekil 2'de gösterildiği gibi cümlelerin pozisyonu değiştikçe özet açısından önemi de değişmektedir. Bütün cümlelerin metnin başlangıcına olan uzaklığı o cümlelerin pozisyonunu vermektedir. Bu pozisyon değişimi cümlelerin giriş, gelişme, sonuç, bölümlerinde olabileceğini göstermektedir. Cümlelerin pozisyonu metnin ortasına göre arttıkça veya azaldıkça özet açısından önemi değişmektedir. Eğer cümlelerin metin içerisindeki pozisyonu metnin ortalarına karşılık geliyor ise cümlelerin gelişme bölümünde bulunma ihtimali artmaktadır. Cümlelerin pozisyonu metnin başlangıcına yaklaştıkça o cümlelerin giriş bölümünde bulunma ihtimali artmaktadır. Cümlelerin pozisyonu metnin sonuna doğru yaklaştıkça cümlelerin sonuç bölümünde bulunma ihtimali artmaktadır [1, 2]. Bu durum Eşitlik 3'te formüle edilmiştir.

$$S^{(l)}_{\text{puan}_{f_{i2}}} = (\text{pozisyon}_{S^{(l)}} - 1/2 (\text{pozisyon}_{S^{(m)}}) - 1/2)^2 \quad (3)$$

Burada $S^{(i)}_{\text{puan}_{f_{iz}}} S^{(i)}$ cümlesine ait pozisyon özelliğinin hesaplanan değeridir. $\text{pozisyon}_{S^{(m)}}$ belgenin son cümlesinin pozisyon değeridir.

2.4. Başlık Kelimeleri

Başlık kelimelerini temsil eden özellik vektörü her cümle için sayfanın başlık bilgisi ile kesişen kelimelerin sayısını ifade etmektedir. Eğer cümle başlıkta geçen kelimeyi içeriyorsa bu cümle metin özetleme açısından önemlidir ve cümle bu varsayıma göre puanlandırılır [1, 2]. Bir cümle içerisinde birden fazla başlık ile kesişen kelime bulunabilir. Her cümle için, başlık bilgisi ile kesişen kelime sayısı ile cümle için puanı doğru orantılıdır.



Şekil 3. Cümlelerin içerdiği başlık kelimesi ve özet açısından önemi

Şekil 3'de cümlelerin başlık bilgisi ile kesişen kelime sayısının, cümlelerin özeti açısından önemini gösterilmektedir.

2.5. Cümlelerin Ağırlıklandırılması

Cümlelerin ağırlığını temsil eden özellik vektörü her terimin her cümle için ağırlık bilgisini temsil etmektedir. Terimlerin cümledeki ağırlık bilgisi Tf-tbf (Terim frekansı – ters belge frekansı) yöntemi ile hesaplanmakta ve belge – terim matrisi kurulmaktadır [1, 2, 9, 10]. Bütün cümlelerin içerdikleri terimlerin belirli bir frekansı bulunmaktadır. Tf-tbf ağırlıklandırma yöntemi metinde geçen terimlerin çıkarılması ve bu terimlerin frekanslarına göre çeşitli hesaplamaların yapılması üzerine kurulmuştur. Tek bir belgenin özetlenmesi yapıldığından dolayı belge yerine cümleler kullanılmıştır. Bütün metinde geçen kelimeler üzerinde Tf-tbf ağırlıklandırılması yapılmakta ve terim – cümle matrisi oluşturulmaktadır. Metinde bir kelime kökü birden çok cümlede geçebilmektedir. Bu yüzden ilk önce metin içerisinde birden çok yerde geçse dahi tekil kabul edilen kelimeler bulunmakta ve bir listede tutulmaktadır. Bu kelimelerin listesi terim–belge matrisindeki terimleri oluşturmaktadır. Terim–belge matrisinde geçen belge bilgisi ise cümleler olarak belirlenmektedir. Terim–belge matrisinin hesaplanmasında her terimin cümle ile ilgili olan ağırlık bilgisi hesaplanmaktadır. Ağırlık bilgisi iki değişkene bağlıdır. Ağırlık bilgisi hesaplanırken ilk değişken bir kelimenin

cümlede geçme sıklığının bulunmasıdır; ikinci değişken ise bir kelimenin metinde geçme sıklığının bulunmasıdır. Bu işlem tüm kelimeler için gerçekleştirilir. Bir cümlede sıkça geçen ancak diğer cümlelerde pek bulunmayan terimin ağırlığı daha yüksek olur. Yüksek ağırlıklı kelimeler içeren cümleler özet için daha önemlidir. Bu varsayım Eşitlik 4'de formüle edilmiştir.

$$\begin{aligned} tf &= \text{terimin cümlede geçme} \div \text{cümledeki toplam terim sayısı} \\ tbf &= \log_{10}(\text{toplam cümle sayısı} \div \text{terimi bulunduran cümle sayısı}) \\ Tf - tbf &= tf \times tbf \end{aligned} \quad (4)$$

Her tekil kelime ve bağlı olduğu cümle için $tf - tbf$ değeri hesaplanmaktadır [2, 9]. Terim-cümle matrisinin sütunundaki değerlerin toplamı her cümle için hesaplanır. Yani cümleye ait kelimelerin ağırlıklarının toplamı o cümlenin toplam ağırlığını vermektedir. Bu toplam ağırlık puanlandırma ile doğru orantılıdır. Ağırlık yöntemi metninin içerdiği cümle sayısı ya da cümlelerin içerdiği kelime sayısı arttıkça daha iyi sonuçlar vermektedir.

2.6. Gizli Anlamsal Analiz

GAA birliktelik kuralları ile fark edilemeyen bağlantıları çözümleyip gözlemlememizi sağlamaktadır [2, 9, 11]. Bu yöntem belgelerde yer alan her cümle içerisindeki terimlerin benzerliklerinden kavramlar çıkartmaya dayanmaktadır. Kavramlar, belgenin konusu hakkında bilgi veren doğal dildeki kelime (terim) ya da kelime öbekleridir. GAA ise terim benzerliklerine dayalı matris ayrıştırma tekniğidir. Cümle-terim matrisi GAA aşamasında da kullanılmaktadır. Cümle-terim matrisi tekil değer ayrışımı (TDA) yöntemi ile daha küçük yönetilebilir bir matrise dönüştürülmektedir. İndirgenmiş matris, metin ile ilişkili terim gruplarını kapsamaktadır. Bu matriste bulunan terimler gruplanarak birbirleri ile ilişkili terim grupları saptanabilmektedir. GAA girdi matrisinin kurulması, TDA ve cümle seçimi olmak üzere üç aşamadan oluşmaktadır.

Girdi matrisi üzerinde yapılacak olan işlemler her kelimenin cümleye göre frekansına bağlı olduğundan, çalışma kapsamında bu aşamada $tf-tbf$ matrisi kullanılmıştır. Her tekil kelimenin, geçtiği cümlelere göre $tf-tbf$ değerleri hesaplanarak terim-cümle matrisi kurulmuştur.

Bir önceki adımda kurulmuş olan girdi matrisi $m \times n$ boyutlu bir cümle-terim matrisidir. Matris boyutları $m > n$ şeklindedir yani matrisin kelime sayısı cümle sayısından fazladır. TDA bu girdi matrisinin öz değer matrisi ve öz vektörleri matrislerini elde etmemizi sağlamaktadır. Bu aşamadan sonra girdi matrisi, A matrisi olarak adlandırılmıştır. A matrisinin rankı $r \leq \min(m, n)$ olarak gösterilmektedir. A matrisi,

$$A = USV^T \quad (5)$$

olarak TDA algoritması tarafından çözümlenir. Burada U matrisi $m \times r$ sütun ortogonal matrisi olup sol tekil vektörü bize verir. V^T matrisi $r \times n$ ortonormal matrisi verir ve bu matris sağ tekil vektördür. Bu sağ ve sol vektörel matrisler A matrisimizin öz vektör matrislerini oluşturmaktadır. S matrisi ise $r \times r$ tekil değerlerin diagonal matrisidir [12].

Diğer bir açıdan bakıldığında TDA A matrisinin indirgenmiş matrislerini vermektedir. A matrisi r tane kavram içermektedir. S matrisi seçilmiş konuların önemini belirtmekte, bir diğer

deyişle kavramları vermektedir. V^T matrisinin her sütunu, metindeki cümleleri; satırları ise kavramları vermektedir. U matrisi ise r tane konu için m adet terim arasındaki eşleşmeyi vurgulamaktadır [2, 9].

TDA algoritmasının uygulanmasında “Apache Commons Math” açık kaynak kodlu kütüphanesi kullanılmıştır. Bu kütüphanenin seçilme nedeni kıyaslama testlerinde diğer kütüphanelere göre daha performanslı çıkmasıdır. “Apache Commons Math” kütüphanesi büyük matrisler üzerinde TDA uygulamasının daha performanslı yürütüldüğü gözlemlendiğinden seçilmiştir.

GAA'nın ilk aşamasında girdi matrisi oluşturulmuştur. Girdi matrisi kullanılarak TDA aşaması gerçekleştirilmiştir. Bu aşama sonucunda üç farklı matris elde edilerek girdi matrisi bu üç matrise indirgenmiştir. Bu aşamalar sonrasında indirgenmiş matrisleri kullanarak cümle seçimi gerçekleştirilmiştir. Cümle seçimi aşaması cümlelere puan vermeye dayalı bir aşamadır. En yüksek puana sahip cümleler metin özetleme açısından daha önemlidir. Cümle seçimi aşamasında çapraz yöntem ve benzerlik yöntemi melez olarak kullanılmıştır. Bu iki yöntem cümleleri puanlandırmada özellik vektörünün devamı niteliğindedir.

2.6.1. Çapraz Yöntem

Çapraz yöntem TDA aşamasından sonra elde edilen V^T matrisini kullanarak hesaplama yapmaktadır. V^T satırları kavramları ve sütunları cümleleri temsil eden kavram-cümle matrisidir. Bu yöntemde öncelikle her satırın ortalaması hesaplanmaktadır. Hesaplanan ortalama değer kavram için eşik değeri kabul edilmektedir. Bu ortalamanın altında kalan kavramların değerleri sıfır yapılmaktadır. Sonrasında yeni oluşturulan V^T matrisinin her sütunu için yani her cümle için o cümlenin uzunluğu hesaplanmaktadır. Bu işlemler Eşitlik 6'da formüle edilmiştir.

$$S^{(i)}_{\text{puan}_{f_{i5}}} = \sqrt{V^T_{(i,j)} \times S_{(j,j)}} \quad (6)$$

Her cümle için çapraz yöntemden elde edilen ağırlık bilgisi, V^T matrisin o sütunundaki değerleri ile S matrisinin ilgili sütununa ait öz değeri kullanılarak hesaplanmaktadır [2, 8, 9, 10].

2.6.2. Benzerlik Yöntemi

Benzerlik yöntemi TDA aşamasından sonra elde edilen U, S, V^T matrislerini kullanarak hesaplama yapmaktadır. TDA aşamasından sonra elde edilen indirgenmiş matrisler, k maksimum tekil değerler kadar yeniden indirgenir. Eski U, S, V^T matrisleri indirgenerek yeni bir $A_k = U_k S_k V_k^T$ matrisi elde edilir.

İndirgeme işlemi tamamlandıktan sonra her cümle için cümle ilişki bağlantısı kurularak her cümle arası bağlantı oluşturulur. A_k matrisinin sütunları cümle vektörleridir; cümlelerin bağlantıları eşit benzerlikte olurlarsa geçersiz sayılırlar. A_k matrisine ait bütün cümle vektörleri arasındaki benzerlikler Eşitlik 7'de gösterildiği gibi kosinüs benzerliğine göre hesaplanmaktadır.

$$S_{puan_{f_{i6}}^{(i)}} = Benzerlik (S_i, S_j) = \frac{\bar{s}_i \times \bar{s}_j}{|\bar{s}_i| \times |\bar{s}_j|} \quad (7)$$

Bütün i ve j cümleleri için benzerlik puanları hesaplandıktan sonra bu benzerlik değerleri cümleler arasındaki bağlantıya ait olan değerlerdir [2, 9, 13, 14].

3. Deneysel Sonuçlar

Uygulamamızda özetlenecek metni oluşturan cümlelere yapısal analiz ve GAA açısından puan ataması gerçekleştirilmiştir. Cümlenin genel özellik vektörü kurulurken melez bir yapı esas alınmıştır. Kurulan özellik vektörü içerisinde altı farklı özellik bulunmaktadır. Özellik f_{ij} i. cümleye ait j. özelliği temsil etmektedir. Çalışmada kullanılan özellikler ve açıklamaları Çizelge 1’de verilmektedir.

Çizelge 1. i. Cümle özellikleri

Özellik	Açıklaması
f_{i1}	Uzunluk
f_{i2}	Pozisyon
f_{i3}	Başlık
f_{i4}	Ağırlık
f_{i5}	Çapraz
f_{i6}	Benzerlik

Sonraki aşamalarda özellikler analiz edilerek metindeki her cümlenin özet cümle olmaya ne kadar uygun olduğu hesaplanmaktadır. Cümlelerin özet cümle için puanları, Türkçenin dil özellikleri göz önüne alınarak aşağıdaki Eşitlik 8’de olduğu gibi hesaplanmaktadır.

$$F(x) = 0.150 \times f_{i1} + 0.135 \times f_{i2} + 0.140 \times f_{i3} + 0.160 \times f_{i4} + 0.090 \times f_{i5} + 0.070 \times f_{i6} \quad (8)$$

Burada, i cümlesinin toplam puanı, cümlenin tüm özelliklerinin belirli bir ağırlık değeri ile çarpılıp toplama eklenmesi sonucunda elde edilmiştir [2, 9]. Eşitlikteki katsayılar literatürde daha önceden benzer şekilde gerçekleştirilen bir çalışmada ‘2’ belirlenmiş katsayılar olarak seçilmiştir. Katsayıların optimizasyonu için farklı yöntemler denenebilir yada uzman görüşüne bağlı olarak farklı değerler seçilebilir.

Çalışmamızda, tüm cümleler Eşitlik 7 ile puanlandırıldıktan sonra, cümle seçimi aşaması gerçekleştirilmektedir. Cümleler yüksek puan alandan az puan alana doğru sıralanmaktadır. Bu sıralı listenin % 40’ı oranında cümle seçilip özet metine eklenmiştir. Son olarak özet metin çıktısı sorguyu oluşturan kullanıcıya gösterilmektedir. Deneylerde kullanılan metinlere ait özet bilgiler Çizelge 2’de verilmektedir.

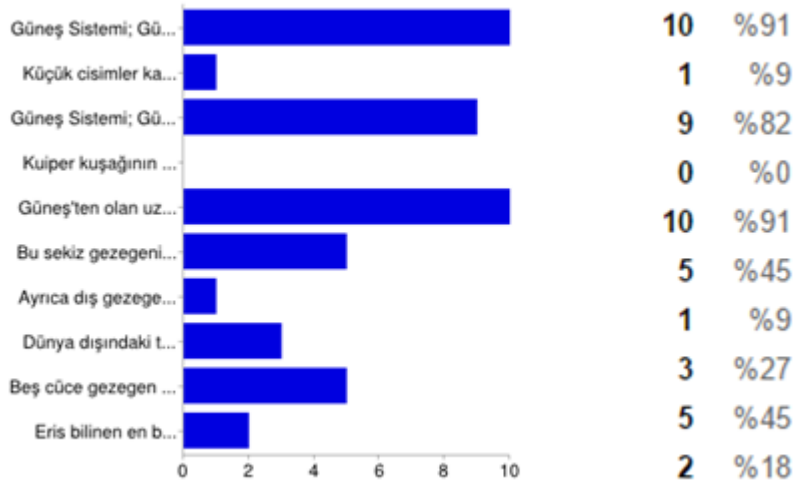
Çizelge 3’de “Güneş Sistemi” belgesine ait özetleme örneği gösterilmektedir.

Cümlelerin puanlarına göre sıralanmış özet cümleleri, cümlelerin orijinal metindeki sıralamasına göre otomatik olarak düzenlendiğinde elde edilen sonuç yukarıdaki çizelgede “Metindeki cümle sırasına göre sıralanmış özet cümleleri” başlığında belirtilmiştir. Metnin anlam bütünlüğü açısından daha doğru bir sıralama bu şekilde elde edilmiştir.

Çizelge 3’de verilen Güneş Sistemi metni proje hakkında bilgisi olmayan 10 farklı kullanıcıya verilmiş ve bu metinden 4 farklı cümle seçerek özetlemeleri istenmiştir. Bu çalışmanın önermiş olduğu cümleler ile kullanıcıların seçtikleri cümleler % 77.5 eşleşmiştir. Şekil 4’de kullanıcıların özet için seçtikleri cümleler ve bu cümlelerin kaç kişi tarafından seçildiğine dair yüzdeler dilim gösterilmiştir.

Çizelge 2. Deneilerde kullanılan metinler

Belge No	Başlık	Cümle Sayısı	Kaynakça
1	Güneş Sistemi	10	http://tr.wikipedia.org/wiki/Güneş_Sistemi
2	Charles Bukowski	6	http://tr.wikipedia.org/wiki/Charles_Bukowski
3	Güneş Ocağı	6	http://tr.wikipedia.org/wiki/Güneş_ocağı
4	Shine On You Crazy Diamond	5	http://tr.wikipedia.org/wiki/Shine_On_You_Crazy_Diamond
5	Fil	32	http://tr.wikipedia.org/wiki/Fil
6	Gandalf	22	http://tr.wikipedia.org/wiki/Gandalf
7	Yıldız Savaşları	13	http://tr.wikipedia.org/wiki/Yıldız_Savaşları
8	Java Programlama Dili	13	http://tr.wikipedia.org/wiki/Java_(programlama_dili)
9	Bilgisayar	13	http://tr.wikipedia.org/wiki/Bilgisayar
10	Afrika	10	http://tr.wikipedia.org/wiki/Afrika



Şekil 4. Metindeki her cümle için kullanıcılar tarafından seçilme yüzdesi

Çizelge 4’de verilen Charles Bukowski metnini aynı şekilde proje hakkında bilgisi olmayan 10 farklı kullanıcıya verilmiş ve bu metinden 2 farklı cümle seçerek özetlemeleri istenmiştir. Bu çalışmanın önermiş olduğu sonuçlar ile kullanıcıların seçtikleri cümleler % 82

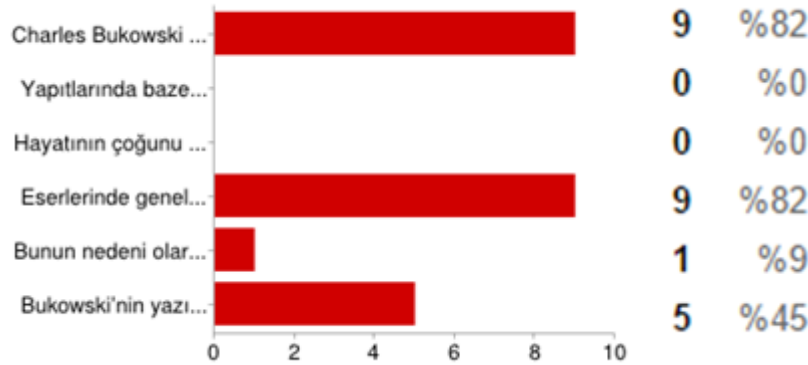
eşleşmiştir. Şekil 5’de kullanıcıların özet için seçtikleri cümleler ve bu cümlelerin kaç kişi tarafından seçildiğine dair yüzdelik dilim gösterilmiştir.

Çizelge 3. Örnek metin özetlemesi

Metin başlığı: Güneş Sistemi
Ana Metin
<p>Güneş Sistemi; Güneş ve onun çekim etkisi altında kalan sekiz gezegen ile onların bilinen 166 uydusu, altı cüce gezegen (Ceres, Plüton, Eris, Haumea, Makemake, Sedna) ile onların bilinen altı uydusu ve milyarlarca küçük gök cisiminden oluşur.</p> <p>Küçük cisimler kategorisine asteroitler, Kuiper kuşağı nesnelere, kuyruklu yıldızlar, gök taşları ve gezegenlerarası toz girer.</p> <p>Güneş Sistemi; Güneş, dört Yer benzeri iç gezegen, küçük, kaya ve metal içerikli asteroitlerden oluşan bir asteroit kuşağı, dört gaz devi dış gezegen, ve Kuiper kuşağı denen buzlu cisimlerden oluşan ikinci bir kuşaktan ibarettir.</p> <p>Kuiper kuşağının ötesinde ise seyrek disk, gündurgun (heliopause) ve en son olarak da varsayımsal Oort bulutu bulunur.</p> <p>Güneş'ten olan uzaklıklarına göre gezegenler sırasıyla Merkür, Venüs, Dünya, Mars, Jüpiter, Satürn, Uranüs ve Neptün'dür.</p> <p>Bu sekiz gezegenin altısının çevresinde doğal uydular döner.</p> <p>Ayrıca dış gezegenlerin her birinin toz ve diğer parçacıklardan oluşan halkaları vardır.</p> <p>Dünya dışındaki tüm gezegenler adlarını Yunan ve Roma mitolojisi'nin tanrılarından alır.</p> <p>Beş cüce gezegen ise; Kuiper kuşağında yer alan Plüton, Haumea ve Makemake; asteroit kuşağındaki en büyük cisim olan Ceres ve seyrek diskte yer alan Eris'tir.</p> <p>Eris bilinen en büyük cüce gezegendir.</p>
Cümlelerin puanlarına göre sıralanmış özet cümleleri
<p>Güneş'ten olan uzaklıklarına göre gezegenler sırasıyla Merkür, Venüs, Dünya, Mars, Jüpiter, Satürn, Uranüs ve Neptün'dür.</p> <p>Güneş Sistemi; Güneş ve onun çekim etkisi altında kalan sekiz gezegen ile onların bilinen 166 uydusu, beş cüce gezegen (Ceres, Plüton, Eris, Haumea, Makemake) ile onların bilinen altı uydusu ve milyarlarca küçük gök cisiminden oluşur.</p> <p>Beş cüce gezegen ise; Kuiper kuşağında yer alan Plüton, Haumea ve Makemake; asteroit kuşağındaki en büyük cisim olan Ceres ve seyrek diskte yer alan Eris'tir.</p> <p>Güneş Sistemi; Güneş, dört Yer benzeri iç gezegen, küçük, kaya ve metal içerikli asteroitlerden oluşan bir asteroit kuşağı, dört gaz devi dış gezegen, ve Kuiper kuşağı denen buzlu cisimlerden oluşan ikinci bir kuşaktan ibarettir.</p>
Metindeki cümle sırasına göre sıralanmış özet cümleleri
<p>Güneş Sistemi; Güneş ve onun çekim etkisi altında kalan sekiz gezegen ile onların bilinen 166 uydusu, beş cüce gezegen (Ceres, Plüton, Eris, Haumea, Makemake) ile onların bilinen altı uydusu ve milyarlarca küçük gök cisiminden oluşur.</p> <p>Güneş Sistemi; Güneş, dört Yer benzeri iç gezegen, küçük, kaya ve metal içerikli asteroitlerden oluşan bir asteroit kuşağı, dört gaz devi dış gezegen, ve Kuiper kuşağı denen buzlu cisimlerden oluşan ikinci bir kuşaktan ibarettir.</p> <p>Güneş'ten olan uzaklıklarına göre gezegenler sırasıyla Merkür, Venüs, Dünya, Mars, Jüpiter, Satürn, Uranüs ve Neptün'dür.</p> <p>Beş cüce gezegen ise; Kuiper kuşağında yer alan Plüton, Haumea ve Makemake; asteroit kuşağındaki en büyük cisim olan Ceres ve seyrek diskte yer alan Eris'tir.</p>

Çizelge 4. Örnek metin özetlemesi

Metin başlığı: Charles Bukowski
Ana Metin
Charles Bukowski (16 Ağustos 1920 – 9 Mart 1994), asıl adı Heinrich Karl Bukowski olan Amerikalı yazar ve şair. Yapıtlarında bazen Henry Chinaski ismini de kullanmıştır. Hayatının çoğunu ABD'nin Los Angeles şehrinde geçirmiştir. Eserlerinde genellikle toplum dışı insanları ve depresyonu konu alması ve alkolizme yakın bir hayat tarzını anlatmasıyla ünlüdür. Bunun nedeni olarak kendisinin bu hayatı yaşaması gösterilebilir. Bukowski'nin yazılarında kendi hayatını yazıp yazmadığı tartışma konusu olmuştur; hayranlarının bir kısmı bunları kurguladığını, çoğunluğu ise yaşamadan bu tip kurguları yapmasının mümkün olmayacağını ve o karakterde bir insanın bu hayatı sürmesinin zaten doğal olduğu görüşünü savunmaktadır.
Cümlelerin puanlarına göre sıralanmış özet cümleleri
Eserlerinde genellikle toplum dışı insanları ve depresyonu konu alması ve alkolizme yakın bir hayat tarzını anlatmasıyla ünlüdür. Charles Bukowski (16 Ağustos 1920 – 9 Mart 1994), asıl adı Heinrich Karl Bukowski olan Amerikalı yazar ve şair.
Metindeki cümle sırasına göre sıralanmış özet cümleleri
Charles Bukowski (16 Ağustos 1920 – 9 Mart 1994), asıl adı Heinrich Karl Bukowski olan Amerikalı yazar ve şair. Eserlerinde genellikle toplum dışı insanları ve depresyonu konu alması ve alkolizme yakın bir hayat tarzını anlatmasıyla ünlüdür.



Şekil 5. Metindeki her cümlelerin kullanıcılar tarafından seçilme yüzdesi

4. SONUÇLAR

Bu çalışmada bir doğal dil işleme uygulaması olan metin özetleme problemi üzerinde çalışılmıştır. Özet cümlelerin seçimi, özetlenecek metinlerin Türkçenin dil özelliklerine dayalı istatistiksel puanlandırılması ve anlamsal puanlandırılması yöntemlerinin melez şekilde değerlendirilmesi ile gerçekleştirilmiştir. Cümlelerin değerlendirilmesinde ele alınacak yapısal ve anlamsal özellikler, literatürdeki çalışmalara benzer şekilde gerçekleştirilmiş ve önerilmiştir. Özetlenecek metinlerde yer alan cümlelerin özet cümle adaylığı için aldıkları puanlar, yapısal ve anlamsal özelliklerin sezgisel bir ağırlıklandırma yöntemi ile birleştirilmesi ile belirlenmiştir. Çalışma kapsamında veri ön işleme, yapısal olarak istatistiksel puanlandırma, GAA analiz ve melez cümle seçimi aşamaları Türkçe yazılmış metinler üzerinde başarıyla gerçekleştirilmiştir.

Elde edilen sonuçların değerlendirilmesi için geliştirilen özetleme sistemine ve farklı kullanıcılara aynı metinler verilmiştir. Özetleme sisteminin önerdiği cümleler ile kullanıcıların önermiş olduğu cümleler karşılaştırılmıştır. Bu karşılaştırma sonucunda “Güneş Sistemi” metninin özeti ile kullanıcıların bu metinden seçtiği cümleler % 77.5, “Charles Bukowski” metninin özeti ile kullanıcıların bu metinden seçtiği cümleler % 82 oranında eşleşmiştir. Burada gösterilemeyen benzer sonuçlar değerlendirildiğinde kısa metinlerde sistemin başarımı daha yüksek iken metinlerin uzunluğu arttıkça sistemin başarımı düşmektedir. Uzun metinlerde başarımı artırmak için metin özellikleri için önerilmiş ağırlıkların optimizasyonu gerçekleştirilebilir.

Çalışmanın bundan sonraki aşamalarında değerlendirilen belge sayısı artırılarak, özetleme için hesaba katılan metin özellikleri zenginleştirilmeye çalışılacaktır.

KAYNAKLAR

- [1] Uzundere E, Dedja E, Diri B, Amasyalı MF. Türkçe Haber Metinleri İçin Otomatik Özetleme, *Akıllı Sistemlerde Yenilikler ve Uygulamaları Sempozyumu*, 2008, s.1-3.
- [2] Güran A, Güler BN, Gürbüz ZM. Efficient Feature Integration with Wikipedia-Based Semantic Feature Extraction for Turkish Text Summarization, *Turkish Journal of Electrical Engineering & Computer Sciences*, 2013, s.3-11.
- [3] Çığır C, Kutlu M, Çicekli I. Generic Text Summarization for Turkish, *The Computer Journal*, 2010, s.1315-1323.
- [4] Güran A, Bayazıt NG, Bekar B. Automatic Summarization of Turkish Documents Using Non-negative Matrix Factorization, *Innovations in Intelligent Systems and Applications (INISTA)*, 2011, s.480-484.
- [5] Das D, Martins AFT, A Survey on Automatic Text Summarization, *Literature Survey for the Language and Statistics II course at CMU*, 2007.
- [6] [http://en.wikipedia.org/wiki/Stop_words], Erişim tarihi: 20.01.2015.
- [7] [<http://www.konumuzseo.com/stop-words-ve-kullanim-amaclari/>], Erişim tarihi: 20.01.2015.
- [8] [[http://tr.wikipedia.org/wiki/Zemberek_\(yazılım\)](http://tr.wikipedia.org/wiki/Zemberek_(yazılım))], Erişim tarihi: 20.01.2015.
- [9] Özsoy GM. Text Summarization Using Latent Semantic Analysis, *Proceeding of the 23rd International Conference on Computational Linguistics*, 2011, s.22-41.
- [10] Yohei S. Sentence Extraction by tfidf and Position Weighting from Newspaper Articles, *Third NTCIR Workshop*, 2003, s.2-6.
- [11] Ünalı İ, Kırkgöz Y. Latent Semantic Analysis: An Analytical Tool for Second Language Writing Assessment, *Mustafa Kemal University Journal of Social Sciences Institute*, 2011, s.2-9.
- [12] Golub HG, Reinsch C. Singular Value Decomposition and Least Squares Solutions, *Numerische Matematik*, 1970, s.1-14.
- [13] Steinberger J, Jezek K. Using Latent Semantic Analysis in Text Summarization and Summary Evaluation, *ISIM2004*, s.2-21.
- [14] Özsoy GM, Çicekli İ, Alpaslan F. Text Summarization of Turkish Texts Using Latent Semantic Analysis, *Proceeding of the 23rd International Conference on Computational Linguistics*, 2010, s.3-6.

ÖZGEÇMİŞ / CV

Akif Hatipoğlu

Akif Hatipoğlu, 27 Mayıs 1993 yılında Balıkesir'in Bandırma ilçesinde doğdu. 18 yaşına kadar yaşamını Bandırma'da sürdürdü. Şu an ki yaşamını Kocaeli'nde sürdürmektedir. İlköğretim ve ortaöğretimi Bandırma Evyapan İlk Öğretim Okulu'nda okudu. Liseyi Bandırma Anadolu Lisesi'nde okudu. 2011 LYS sınavı ile şu an okumanda olduğu Kocaeli Üniversitesi Bilgisayar Mühendisliği bölümünü kazandı. Bu bölümde 4.sınıf öğrencisi olarak halen öğrenim hayatına devam etmektedir.

Akif Hatipoglu, was born on May 27,1993 in Bandırma/Balıkesir. He continued living up to 18 years in Bandırma. He continues his current life in Kocaeli. He studied at the elementary and secondary education Bandırma Evyapan Elementary School. He studied high school in Bandırma Anatolian High School. He won the Kocaeli University Computer Engineering Department in 2011. He still continues his education career as a 4th year student of Computer Engineering at the Kocaeli University.