Makine Öğrenmesi Algoritmalarında Veri Setlerine Öznitelik Seçim Tekniklerinin Uygulanması ve Karşılaştırmalı Analizi

Application and Comparative Analysis of Attribute Selection Techniques to Data Sets in Machine Learning Algorithms.

Emrehan EKENTOK Bilgisayar ve Elektrik Elektronik Mühendisliği KTO Karatay Üniversitesi - Konya, Türkiye e.ekentok@gmail.com

Özetçe— Bu çalışmada, Makine öğrenmesi algoritmalarına tabi tutulan veri setlerinde öznitelik seçim tekniğinin uygulanması ile elde edilen sonuçların, tüm özniteliklerin kullanıldığı veri setiyle elde edilen sonuçlar (özgün veri setinin öznitelikleri) ile karşılaştırılması, deneysel gözlemlerin sonuçları ve öznitelik seçim tekniklerinin kullanımının önemi hakkında çıkarımlarda bulunulmuştur.

Anahtar Kelimeler — Makine Öğrenmesi; Öznitelik Seçimi; Naive Bayes;

Abstract— In this study, the results obtained by applying the feature selection technique to the data sets subjected to the machine learning algorithms were compared with the results obtained with the data set of all the attributes (original data sets attributes). the results of the experimental observations and the importance of using the feature selection techniques has been mentioned.

Keywords — Machine Learning; Attribute Selection; Naive Bayes

I. Giris

Günümüzde Makine Öğrenmesi algoritmaları Üretim, Pazarlama-Satış, Sağlık Bilimleri, Finans ve Mali Hizmetler, Enerji, Seyahat veya Sosyal Platformlar gibi birçok alanda sıklıkla kullanılmaktadır. Makine Öğrenmesi Algoritmalarının kullanılmakta olduğu bazı alan ve uygulamalarda kullanılan veri setlerinin sahip olduğu öznitelik sayıları oldukça yüksek sayılardan oluşabilmektedir. Hal böyle olunca araştırmacılar öznitelik seçim yöntemlerine önceki çalışmalara nazaran daha çok ihtiyaç duymaktadırlar. Öznitelik seçimi uygulanmış veri setinde yapılacak işlem sayısı azalmakla beraber gürültülü veya veri setiyle ilgisiz olan öznitelikler veri setinden çıkarılarak,

sınıflandırma başarı oranı yükselmekte, sınıflandırma prosesi (eğitim zamanı, daha az bellek kullanımı) kolaylaşmakta, daha az öznitelik ile veri setinin tüm özniteliklerinin kullanılmasına kıyasla ufak hata oranları, aynı başarı oranları veya daha iyi sonuçlar elde edilmektedir.

Özellikle, Örnekleme sayısının az fakat öznitelik sayısının fazla olduğu veri setlerine sahip uygulamalar ele alındığında, öznitelik seçimi yaparak sınıflandırmaya daha uygun özniteliklerin belirlenebilmesi ve/veya işlem yoğunluğun sadeleştirilmesi konusunda öznitelik seçim tekniğinin önemi ortaya çıkmaktadır.

Bu bildiride iki farklı veri seti üzerinde öznitelik seçimi tekniği uygulanarak elde edilen sonuçların özgün veri setinden elde edilen başarı sonuçları ile karşılaştırılması analiz edilmiş ve öznitelik seçim tekniklerinin öneminin vurgulanması hedeflenmiştir.

II. MATERYAL VE METOT

A. Öznitelik Seçimi

Öznitelik seçim işlemi veri setindeki bağımlı değişken (Supervised Variable) ile hiç ilgisi bulunmayan, gürültülü özniteliklerin elenmesi veya bağımlı değişkeni açıklama unsuru daha yüksek olanların kümelerin belirlenmesi işlemidir.

Genel kanı, veri setinde tüm özelliklerin kullanılması ile daha başarılı sonuçlar elde edileceği yönündedir. Ancak bu yaklaşım çok sayıda öznitelik içeren veri setlerinde her zaman doğru olmayabilir. Veri setindeki her öznitelik algoritmaya uygun açıklayıcı veya doğru bilgiler taşımayabilir. Diğer bir deyişle veri setlerinde bazı öznitelikler algoritmanın işleyiş performasına negatif olarak etki edecek gürültülü bilgi

içerebilir, bu özniteliklerin ayrıştırılması algoritma başarı oranı artırmakla beraber, algoritmada kullanılacak veri boyutunun indirgenmesinin avantajlarınıda beraberinde getirmektedir.

Veri setleri içerisindeki etkin öznitelikleri seçme işlemi WEKA programı ile gerçekleştirilmiştir. Öznitelik seçim yöntemi olarak Correlation-based Feature Subset Selection (CfsSubsetEval), arama metodu olarak bestFirst tercih edilmiştir.

CFS, 1999 yılında Hall tarafından geliştirilmiştir [1]. CFS, alt öznitelik kümelerini korelasyon-bazlı değerlendirerek en iyi alt öznitelik kümeyi bulmayı hedefleyen filtre modelli öznitelik seçme algoritmalarından biridir. Temel prensip olarak kendi aralarında korelasyonu az, sınıf etiketleri ile korelasyonu fazla öznitelik alt kümesini seçmeye çalışan bir algoritmadır.

B. Veri Setleri

Bu çalışmada UCI veritabanından, *Waveform Database Generator (version 1) v*eri seti ve Waikato üniversitesinde açık kaynak kodlu olarak JAVA dili üzerinde geliştirilmiş Weka yazılımının (versiyon 3.8.1) kütüphanesinde bulunan *Vote* veri seti kullanılmıştır.

Waveform Database Generator veri setinde 3 adet sınıf, 40 öznitelik ve 5000 adet örnek bulunmaktadır. Vote veri setinde ise 16 öznitelik ve 465 örnek bulunmaktadır.

C. Sınıflandırma

Bildiriye konu olan bu çalışmada her iki veri setindede makine öğrenmesi algoritmalarından Naive Bayes algoritması kullanılmıştır.

Naive Bayes Sınıflandırma algoritması Bayes teoremine dayanan temel bir olasılıksal sınıflandırma yöntemidir. Hali hazırda sınıflandırılmış durumdaki örnek verileri kullanarak yeni bir verinin mevcut sınıflardan herhangi birine ait olma olasılığını hesaplayan bir yaklaşımdır. Bu sınıflandırma algoritmasında öznitelikler birbirinden bağımsız olarak değerlendirilir. Bir Bayes yaklaşımı olarak, n boyutlu uzayda tanımlı olan X vektörü $(x_1,...,x_2)$, m adet sınıf bulunan C_k $(C_1,..,C_n)$ veri kümesinde son olasılığı maksimize eden bir sınıf etiketi C arar.

III. DENEYSEL SONUÇLAR

Deniz dalga formlarının sınıflandırılması için elde edilen *Waveform Database Generator* veri seti ve *Vote* veri seti Naive Bayes algoritması, çapraz doğrulama (k10) tekniği ile çalıştırılmıştır. Çapraz doğrulama test yönteminde veri seti 10 eşit kümeye bölünmekte, 9 küme eğitim için 1 küme test için kullanılmaktadır. Test sonunda 10 adet performans metriği elde edilmekte ve elde edilen her bir metriğin aritmetik ortalaması alınmaktadır.

Waveform Database Generator veri seti özgün olarak 40 adet özniteliğe ve 1 adet sınıflandırma niteliğine sahiptir (class) sahiptir. 40 öznitelik ile Naive Bayes sınıflandırıcı ile 10k çapraz doğrulama test tekniği kullanılarak gerçekleştirilen çalışmanın sonucunda, % 80 'lik bir başarı oranı elde edilmiştir.

Correctly Classified Instances Kappa statistic Mean absolute error Root mean squared error Relative absolute error Root relative squared error Total Number of Instances === Confusion Mar	4000 0.7005 0.1357 0.3369 30.5282 % 71.56 5000	80	%
a b c 890 394 408 16 1547 90 2 90 1563	< classified as a = 0 b = 1 c = 2		

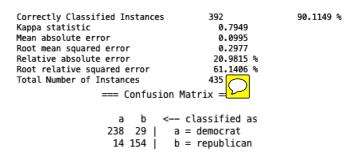
Resmit; Özgün Öznitelikler ile Sınıflandırma Sonucu

Özgün veri setine uygulanan Öznitelik Seçim tekniği Correlation-based Feature Subset Selection (CfsSubsetEval) olarak seçilmiş, arama metodu ise bestFirst tercih edilmiştir. Bu işlem sonucun da veri setinde 15 adet öznitelik kalmış 25 adet öznitelik çıkartılarak veri setinde kullanılmamıştır. 15 öznitelikle aynı şekilde Naive Bayes ile 10k çapraz doğrulama test tekniği kullanılarak sınıflandırma algoritması çalıştırılmıştır. Başarı oranı % 80.12 'ye yükselerek sadece 15 adet öznitelik ile bu sonuç elde edilmiştir. Sonuçlara ilişkin diğer Detaylar ise aşağıdaki Resim 2 görselinde verilmiştir.

Correctly Classified I Kappa statistic Mean absolute error Root mean squared erro Relative absolute erro Root relative squared Total Number of Instar	r or error		4006 0.7023 0.1349 0.333 30.3476 % 70.6464 %	80.12	%
		ion Mat	trix ===		
a	b	C 404	< classi lica as		
894	394	404	a = 0		
15	1546	92	b = 1		
4	85	1566	c = 2		

Azaltılmış Öznitelik ile Sınıflandırma Sonucu

Weka 3.8.1 versiyonlu yazılımın kütüphanesinde bulunan *Vote (oy)* veri seti özgün olarak 16 adet öznitelik ve 1 adet sınıflandırma (class) niteliğine ve 465 adet örneğe sahiptir. 16 öznitelik özgün olarak Naive Bayes sınıflandırıcı algoritması ile 10k çapraz doğrulama test tekniği kullanılarak çalıştırıldığında % 90.1149'luk başarı oranı elde edilmiştir. Sonuçlar ile ilgili detaylar Resim 3' de gösterilmiştir.





VOTE veri setine *CfsSubsetEval* — *bestFirst* öznitelik seçim tekniği uygulandığında ise veri setinde 4 adet öznitelik seçilmiş 12 adet öznitelik veri setinden çıkartılarak kullanılmamıştır. 4 öznitelikle aynı şekilde Naive Bayes ile 10k çapraz doğrulama test tekniği kullanılarak sınıflandırma algoritması çalıştırılmıştır. Başarı oranı % 96.092 'ye yükselmiştir.

Correctly Classified Instances 418 96.092 % Kappa statistic 0.9177 Mean absolute error 0.0575 Root mean squared error 0.1768 Relative absolute error 12.1285 % Root relative squared error 36.3023 % Total Number of Instances 435 === Confusion classified as a = democrat 258 9 8 160 b = republican

Re ; Azaltılmış Öznitelik ile Sınıflandırma Sonucu

Öznitelik Seçim Tekniğinin Veri Setlerine Uygulanması ile Elde Edilen Sonuçların Karşılaştırılması							
Veri Seti	Özgün Öznitelik	Azaltılmış Öznitelik	Algoritma	Başarı Oranı	Az.Öz.Baş arı Oranı		
WaveForm	40	15	Naive Bayes	%80	% 80.12		
Vote	16	4	Naive Bayes	% 90.1149	%96.092		

Tabl Başarı Sonuçlarının Karşılaştırılması

IV. SONUÇLAR VE DEĞERLENDİRME

Çalışmaya Konu olan Makine Öğrenme Algoritmalarında kullanılan veri setlerinde filtre olarak da değerlendirilebilinen Öznitelik seçim tekniği iki adet veri setine uygulanmıştır. Bu veri setlerinin her ikisinde de öznitelik sayıları CFS Tekniği ile azaltılmış ve Naive Bayes algoritması 10 katlı Çapraz Doğrulama test tekniği ile çalıştırılmıştır. Elde edilen sonuçlar karşılaştırılmıştır.

Her iki veri setindede bildiriye konu olan öznitelik azaltma işleminin önemi ve başarısı açıkça görülmüştür. Daha az öznitelik ile elde edilen başarı oranı analiz edilerek tespit edilmiştir. Ayrıca daha çok özniteliğe sahip olan veri setlerinde ise Öznitelik seçim tekniklerinin başarı oranlarında azaltılmamış özgün öznitelik sayılı veri seti ile yapılan çalışmalardan daha başarılı olabileceğine dair sonuçlar gözlemlenmiştir.

Sonuç olarak günümüzde birçok alan ve uygulamada kendine yer bulan makine öğrenmesi algoritmalarında ayrıştırılmış ve/veya azaltılmış öznitelikler sayesinde başarı oranlarının artırılması, işlem sürelerinin azaltılması, en uygun öznitelikler ile algoritmaların çalıştırılması, daha az özniteliğe sahip olunması gibi yararları gözlemlenmiştir.

Bu çalışmada yeni bir yöntem veya metot önerimi yapılmamış makine öğrenmesi algoritmalarında öznitelik seçim tekniklerinin kullanılmasının önemi ve sonuçları tecrübe edilerek aktarılmıştır. ileride bu konu ile ilgili yapılacak çalışmalara bir temel olması amaçlanmıştır.

KAYNAKLAR

