



## DNA Dizilimindeki Nükleotit Çiftlerinin Frekans Değerlerine Göre Farklı Sınıflandırma Yöntemleri ile Karşılaştırılması

Bihter Daş<sup>1</sup>, İbrahim Türkoğlu<sup>2</sup>

<sup>1</sup>Fırat Üniversitesi, Teknoloji Fakültesi, Yazılım Mühendisliği, 23119, Elazığ/Türkiye.

<sup>2</sup>Bingöl Üniversitesi, Mühendislik ve Mimarlık Fakültesi, Bilgisayar Mühendisliği Bölümü, 12000, Bingöl/Türkiye.

[bihtedas@gmail.com](mailto:bihtedas@gmail.com) [iturkoglu@bingol.edu.tr](mailto:iturkoglu@bingol.edu.tr)

### Özet

DNA dizilim hizalama, biyoinformatiğin en temel problemlerinden biridir. Son yıllarda DNA üzerine yapılan çalışmalarda, DNA'daki nükleotidlerin dizilişlerinin birbiri ardı sıra tek, ikili, üçlü olarak tekrar ettiği belirlenmiştir. Kısa ard arda tekrar eden (STR) olarak adlandırılan ve çeşitli alanlarda kullanılan bu belirteçler genetik hastalıklarda, adli amaçlı kimlik tespitinde, babalık araştırmaları ve tümör biyokimyasal araştırmaları gibi birçok farklı amaçlar için kullanılmaktadır.

Bu makale çalışmasında, bakteri türlerinin farklı uzunluklardaki DNA dizilimleri alınarak, bu dizilimlerde tekrar eden nükleotid çiftlerin frekansı bulunmuş ve elde edilen bu frekans değerlerine YSA, DVM ve KNN yöntemleri uygulanarak bir sınıflandırma yapılmıştır. Sınıflandırma sonucunda KNN yönteminin, YSA ve DVM yöntemlerinden başarımının daha yüksek olduğu görülmüştür.

**Anahtar Kelimeler:** DNA dizilimi, özellik çıkarımı, sınıflandırma, STR

### Abstract

DNA sequence alignment is one of the major problems of bioinformatics. In recent year's studies on DNA, it is determined that sequences of nucleotide in DNA are consecutive as single, double or triple. These markers called STR are used in various fields such as genetic diseases, forensic identification, fatherhood research, biochemical tumor research.

In this paper, it was found the frequency of repetitive nucleotide pairs taking DNA sequences of different lengths of the bacterial species and a classification was made to obtained frequency values by using ANN, SVM and KNN methods. Classification results show performance of KNN method is higher than performance of SVM and ANN methods.

**Keywords:** DNA sequence, feature extraction, classification, STR

### 1. Giriş

Biyoinformatik, biyolojik problemlerin çözümünde bilişim teknolojilerinin kullanılması esasına dayanan bilimsel bir disiplindir. DNA'daki nükleotid diziliminden yola çıkarak hastalıklara sebep olan biyolojik verilerin derlenmesi, analiz edilmesi bu hastalıkların tedavi yöntemlerinin geliştirilmesi biyoinformatiğin en temel amacıdır.

DNA molekülü birbiri etrafında kıvrılmış iki iplikten oluşur. Her bir iplik nükleotid adını verdiğimiz (şeker+fosfat+baz) ünitelerinin tekrarlarından ibaret olup polinükleotid adını alır. DNA'nın yapısında bulunan bazlar adenin(A), timin(T), guanin(G) ve sitozin(C) olmak üzere 4 çeşittir.

Son yıllarda DNA üzerinde yapılan araştırmalarda; DNA'daki bazların ve baz dizilimlerin birbiri ardı sıra tekrar ettiği belirlenmiştir. STR olarak adlandırılan bu tekrarlar genetik kökenli haritalamada, tümör biyokimyasal araştırmalarda, adli bireysel kimlik tespitinde, babalık ve nüfus genetik analizlerde yaygın bir şekilde kullanılmaktadır [4] Çeşitli veri madenciliği tekniklerinin kullanılması DNA dizilimlerinin sınıflandırılması açısından büyük öneme sahiptir.

Bu makalede DNA dizilimlerini sınıflandırmada Yapay Sinir Ağı(YSA), Destek Vektör Makineleri(DVM) ve K-En Yakın Komşu (KNN) metotları kullanılmıştır.

### 2. Verilerin elde edilmesi

Bu çalışmada deneysel veriler için *National Center for Biotechnology (NCBI)* sitesi gen bankasından [5] gerçek veriler alınmıştır. Alınan *Escherichia coli*, *Bacillus cereus*, *Buchnera aphidicola* ve *Enterobacter cloacae* gibi 4 bakteri türü sınıflandırılma için kullanılmıştır. Tablo 1'de kullanılan bakteri türleri, erişim numaraları, sayıları ve ortalama uzunlukları verilmiştir.

### Biyomedikal Sinyallerde sınıflandırma Uygulamaları

3. Gün 27 Eylül 2014 Cumartesi (09.45-10.45)

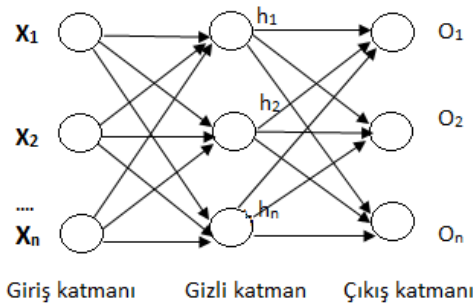
Tablo 1: Kullanılan veri türlerinin açıklaması

Sıra	Bakteri Türü	Erişim Numarası	Bakteri Örnek Sayısı	Dizimdeki Nükleotid Sayısı
1	Escherichia coli	AE14075	31	102
2	Bacillus cereus	NC_016771	40	94
3	Buchnera aphidicola	AF012886	35	96
4	Enterobacter cloacae	EU606203	48	104

### 3. Kullanılan sınıflandırma yöntemleri

#### a. Yapay Sinir Ağları

Yapay Sinir Ağları insan beyninden esinlenerek geliştirilmiş bir bilgi işlem teknolojisidir. YSA ile basit biyolojik sinir sisteminin çalışma şekli benzetim yapılıdır. Yapay sinir ağları, insan beyninin özelliklerinden olan öğrenme yolu ile hafızaya alma, yeni bilgiler oluşturabilme ve keşfedebilme, veriler arasındaki ilişkileri ortaya koyma gibi yetenekleri otomatik olarak gerçekleştirmek amacı ile geliştirilmişlerdir [8]. Yapay sinir ağlarının iki türlü çalışma şekli vardır. Biri eğitim (öğrenme) diğeri test (kullanma) aşamasıdır. Bir yapay sinir ağının kullanılabilmesi için önce eğitilmesi gerekir. Bir YSA'nın birim elemanı nörondur (düğüm). Yapay sinir ağları Şekil 1'de görüldüğü üzere temel olarak girdi, gizli ve çıktı katman olmak üzere üç katmandan oluşmakta ve her katmanda birçok nöron (düğüm) bulunmaktadır.



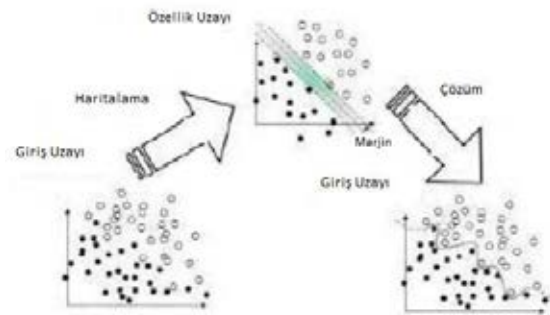
Şekil 1: Üç Katmanlı bir Yapay Sinir Ağı

Bu çalışmada kullanılan YSA modelinde giriş katmanı 16 düğümden oluşmaktadır. Gizli katmandaki düğümlerin sayısının eğitim süresince tespit edilmesi gerekmektedir. Çıkış katmanı ise bakteri türlerini temsil eden 1,2,3,4 şeklinde 4 düğümden oluşmaktadır.

#### b. Destek vektör makinesi

Son yıllarda, sınıflandırma problemlerinin çözümü için geliştirilmiş en başarılı makine öğrenimi algoritmalarından biri Destek Vektör Makineleri'dir. Destek Vektör Makineleri, değişkenler arasındaki örüntülerin bilinmediği veri setlerindeki

birçok sınıflandırma probleminin çözümünde başarıyla uygulanmış, performansı yüksek ve etkin makine öğrenimi algoritmalarından biri olarak veri madenciliği uygulamalarındaki yerini almıştır [9]. Bu yöntem, sınıflandırmayı bir doğrusal ya da doğrusal olmayan bir fonksiyon yardımıyla yerine getirir. Doğrusal olmayan dönüşümlerde kernel fonksiyonu kullanılmakta ve verilerin daha yüksek boyutta doğrusal olarak ayırılmasına imkân sağlanmaktadır ve Şekil 2' de destek vektör algoritmasının genel yapısı görülmektedir. [7]

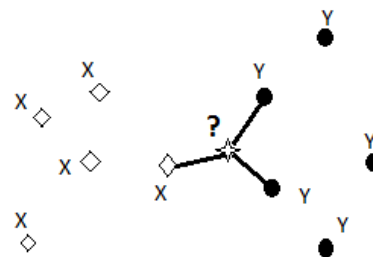


Şekil 2: Destek Vektör Makine Algoritması [7]

Bu çalışmada destek vektör makineleri (DVM) ile gerçekleştirilecek dört sınıflı bir sınıflandırma için radyal tabanlı fonksiyon (RBF) kerneli kullanılmıştır. Uygulamada kernel parametresi (RBF kerneli için band genişliği değeri) 2, düzenleme parametresi (C) ise 100000 olarak belirlenmiştir.

#### c. K-En yakın komşu algoritması

K-En Yakın Komşu Algoritması (K-Nearest Neighbor Algorithm), yeni bir veri geldiğinde varolan öğrenme verisi üzerinde sınıflandırma yapan eğitilmiş öğrenme algoritmasıdır. Algoritma, yeni bir veri geldiğinde, onun en yakın K komşusuna bakarak bu verinin sınıfına karar verir. Her sınıfın özelliklerinin önceden belirlenmiş olması çok önemlidir. Yeni gelen verinin daha önceki verilerden k tanesine yakınlığına bakılır. Bu iki veri arasındaki mesafe çeşitli uzaklık fonksiyonları kullanılarak hesaplanır. Manhattan Uzaklık Fonksiyonu, Minkowski Uzaklık Fonksiyonu, Öklid Uzaklık Fonksiyonu içerisinde en çok tercih edilen fonksiyon oklid uzaklık fonksiyonudur. En yakın mesafe neresi ise yeni veri o sınıfa atanır. Bu çalışmada kullanılan KNN algoritmasında k değeri 5 olarak alınmıştır.



Şekil 3: K-En Yakın Komşu Algoritması





***Biyomedikal Sinyallerde sınıflandırma Uygulamaları***

3. Gün 27 Eylül 2014 Cumartesi (09.45-10.45)

- [4] Zhou Q., Jiang Q. Ve Wei D. ,”A new method for classification in dna sequence”, the 6th international conference on computer science&education (iccse 2011) august 3-5, 2011.
- [5] <http://www.ncbi.nlm.nih.gov/Genbank/genomes/bacteria>.
- [6] Ayhan S., Erdoğan Ş., “Destek vektör makineleriyle sınıflandırma problemlerinin çözümü için çekirdek fonksiyonu seçimi “,iibf dergisi, 9(1), 175- 198, 2014.
- [7] Karagülle F., “Destek vektör makinelerin kullanarak yüz bulma” ,yüksek lisans tezi, 2008
- [8] A.şengür, i.türkoğlu ve m.c.ince, wavelet packet neural networks for texture classification, expert systems with applications, 32(2), mart 2007.
- [9] Sengur, a., “multiclass least-squares support vector machines for analog modulation classification”, expert systems with applications, 36(3), 6681-6685 (2009).