

Doğal Dil İşlemede Eğilimler

Önceden: Yapay Zeka Tabanlı, tam olarak anlama

Şimdiki: Külliyyat(Corpus)-tabanlı, İstatistiki, makine öğrenmesi içeren

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Akış

- Makine Öğrenmesi nedir?
- Günlük Hayatımızdaki Uygulamaları
- Verilerin Sayısallaştırılması
- Özellik Belirleme
 - Özellik Seçim Metotları
 - Bilgi Kazancı (Informaiton Gain-IG)
 - Sinyalin Gürültüye Oranı: (S2N ratio)
 - Alt küme seçiciler (Wrappers)
 - Yeni Özelliklerin Çıkarımı
 - Temel Bileşen Analizi (Principal Component Analysis)
 - Doğrusal Ayırteden Analizi (Linear Discriminant Analysis)
- Sınıflandırma Metotları
 - Doğrusal Regresyon
 - Karar Ağaçları (Decision Trees)
 - Yapay Sinir Ağları
 - En Yakın K Komşu Algoritması (k - Nearest Neighbor)
 - Öğrenmeli Vektör Kuantalama (Learning Vector Quantization)
- Kümeleme Algoritmaları:
 - K-Ortalama (K-Means)
 - Kendi Kendini Düzenleyen Haritalar (Self Organizing Map -SOM)

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Makine Öğrenmesi

Çok büyük miktardaki verilerin elle işlenmesi, analizinin yapılması mümkün değildir.

Bu tür problemlere çözüm bulmak amacıyla makine öğrenmesi metotları geliştirilmiştir.

Bu metotlar

geçmişteki verileri kullanarak

veriye en uygun **modeli** bulmaya çalışırlar.

Yeni gelen verileri de bu modele göre analiz edip sonuç üretirler.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



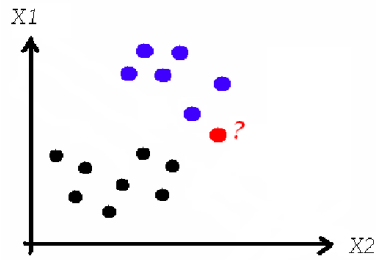
Metot türleri

- Farklı uygulamaların analizlerden farklı beklentileri olmaktadır.
- Makine öğrenmesi metotlarını bu beklentilere göre sınıflandırmak mümkündür.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Sınıflandırma



Kırmızı hangi sınıftan?

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Kümeleme

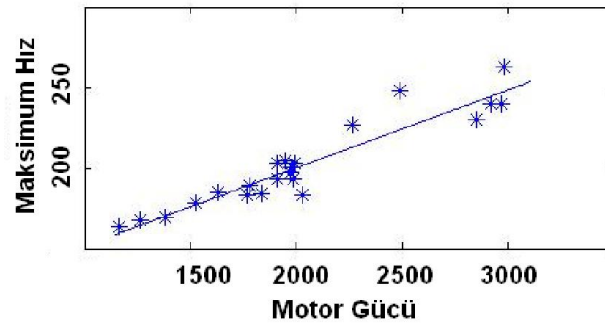
- 256 rengi 16 rene nasıl indiririz?



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Regresyon Eğri Uydurma



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Birliktelik Kuralları Keşfi

- Sepet analizi
 - hangi ürünler birlikte satılıyor?
- Raf düzenlemesi
 - hangi ürünler yan yana konmalı?
- Promosyonlar
 - neyin yanında ne verilmeli?

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Akış

- Makine Öğrenmesi nedir?
- **Günlük Hayatımızdaki Uygulamaları**
- Verilerin Sayısallaştırılması
- Özellik Belirleme
 - Özellik Seçim Metotları
 - Bilgi Kazancı (Informaiton Gain-IG)
 - Sinyalin Gürültüye Oranı: (S2N ratio)
 - Alt küme seçiciler (Wrappers)
 - Yeni Özelliklerin Çıkarımı
 - Temel Bileşen Analizi (Principal Component Analysis)
 - Doğrusal Ayırteden Analizi (Linear Discriminant Analysis)
- Sınıflandırma Metotları
 - Doğrusal Regresyon
 - Karar Ağaçları (Decision Trees)
 - Yapay Sinir Ağları
 - En Yakın K Komşu Algoritması (k - Nearest Neighbor)
 - Öğrenmeli Vektör Kuantalama (Learning Vector Quantization)
- Kümeleme Algoritmaları:
 - K-Ortalama (K-Means)
 - Kendi Kendini Düzenleyen Haritalar (Self Organizing Map -SOM)

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Makine Öğrenmesinin

Günlük Hayatımızdaki Uygulamaları

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



El yazısı / Kitap Yazısı Tanıma
HCR /OCR



İşlem: Şekillerin hangi harf olduğunu tahmin etme

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Kredi Taleplerini Değerlendirme

- Birisi bankadan borç ister.
- Banka borcu versin/vermesin.
- Nasıl?

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



e-ticaret

- Birisi Amazon.com dan bir kitap yada ürün alıyor.

Görev ne olabilir?

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



e-ticaret

- Birisi Amazon.com dan bir kitap yada ürün alıyor.

Görev ne olabilir?

Müşteriye alması muhtemel kitapları önerelim.

Ama nasıl?

Kitapları

konularına

yazarlarına

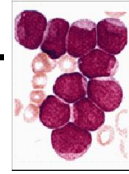
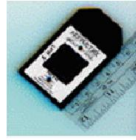
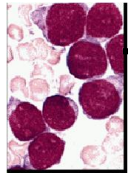
birlikte satışlarına

göre kümelemek?

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



ALL



AML

Gen Mikrodizilimleri

100 kişinin (hasta/sağlam) elimizde gen dizilimleri var.

Bu dizilimleri analiz ederek hasta olup olmadığı bilinmeyen birisinin hasta olup olmadığını yada hastalığının türünü öğrenebilir miyiz?

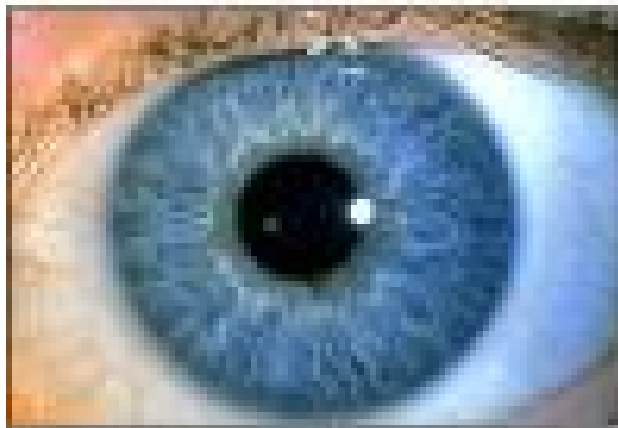
En iyi tedaviyi önerebilir miyiz?

Nasıl? Elimizde hangi bilgiler olmalı?

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



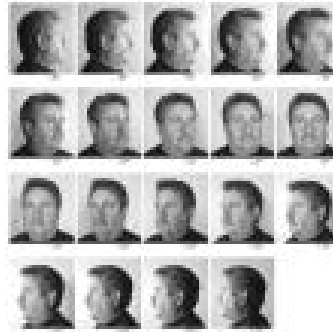
Bu adam kim? İçeri girsin mi?



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



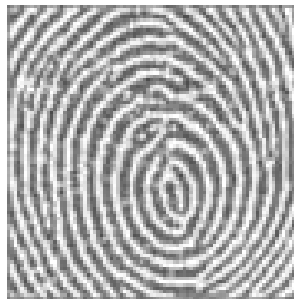
Bu adam kim?
Bu adam havaalanında mı?



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Bu parmak izi kimin?
Bu adamı tutuklayalım mı?



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Bu ses kimin?
Bu ses ne diyor?



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Bu imza kimin? Yoksa taklit mi?

Taklit olup olmadığını nasıl anlarız?
Zaman bilgisi ?

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Bu metnin konusu nedir? Bu mail spam mi?



Anti spam yazılımları nasıl çalışır?

Spamciler nasıl çalışıyor?

Yeni nesil spam mailler: Mesaj resimde,
metinde ise anti spamlardan kaçmak için gereken kelimeler var.

Makine öğrenmesi metodlarını hem spamciler hem anti spamciler kullanıyor.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Olağan dışı bir durum var mı? Güvenlik kamerası kayıtları



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Kamera kaydındaki kişi ne anlatıyor?



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Sonuç: İletişimin artması

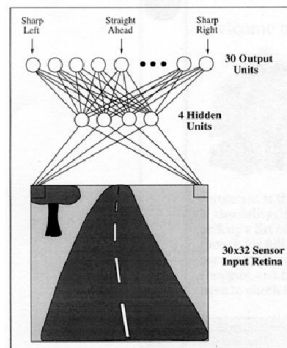


ALVIN

Otobanda saatte 70 mil
hızla **sürücüsüz**
gidebilen bir otomobil

Bütün denemeler
trafiğe kapalı alanlarda
gerçekleştirilmiştir 😊

Neden şehiriçi değil?
Neden otoban?
Neden diğer arabalar yok?
Araba birine çarparsa suçlu
kim?



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Adalet

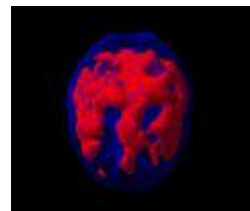
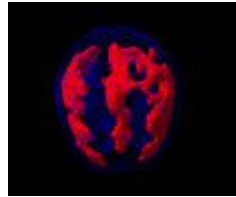
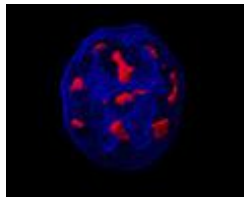
- Çin’de pilot uygulama:
 - bir şehrin mahkeme hakimleri bir bilgisayar programı
 - Amaç: Daha adil bir dünya
 - Aynı özelliklere sahip davalarda aynı kararların alınması
 - Sistemin eğitimi için neler gerekli?
 - Milyonlarca/Milyarlarca (orası Çin) davaya ait verilerin kategorilenmesi

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Beyin Aktiviteleri

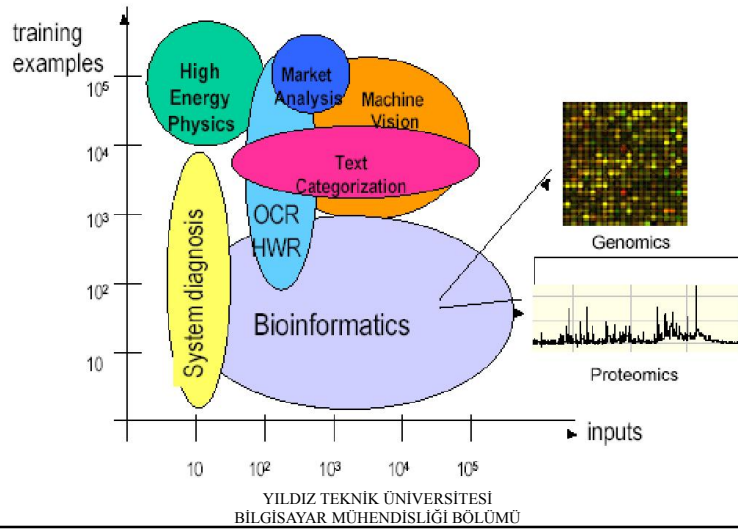
- İnsanların
 - değişik şeyler düşünürkenki,
 - değişik duygulara sahipkenki,
 - problem çözerken ki
 beyin aktiviteleri kaydedilir.
- Görev?



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Uygulamalardaki boyut örnek sayıları



Learning is the future

- Learning techniques will be a basis for every application that involves a connection to a real world
- Basic learning algorithms are ready for use in limited applications today
- Prospects for broader future application make for exciting fundamental research and development opportunities
- Many unresolved issues – Theory and Systems



Akış

- Makine Öğrenmesi nedir?
- Günlük Hayatımızdaki Uygulamaları
- **Verilerin Sayısallaştırılması**
- Özellik Belirleme
 - Özellik Seçim Metotları
 - Bilgi Kazancı (Informaiton Gain-IG)
 - Sinyalin Gürültüye Oranı: (S2N ratio)
 - Alt küme seçiciler (Wrappers)
 - Yeni Özelliklerin Çıkarımı
 - Temel Bileşen Analizi (Principal Component Analysis)
 - Doğrusal Ayırteden Analizi (Linear Discriminant Analysis)
- Sınıflandırma Metotları
 - Doğrusal Regresyon
 - Karar Ağaçları (Decision Trees)
 - Yapay Sinir Ağları
 - En Yakın K Komşu Algoritması (k - Nearest Neighbor)
 - Öğrenmeli Vektör Kuantalama (Learning Vector Quantization)
- Kümeleme Algoritmaları:
 - K-Ortalama (K-Means)
 - Kendi Kendini Düzenleyen Haritalar (Self Organizing Map -SOM)

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Verilerin Sayısallaştırılması

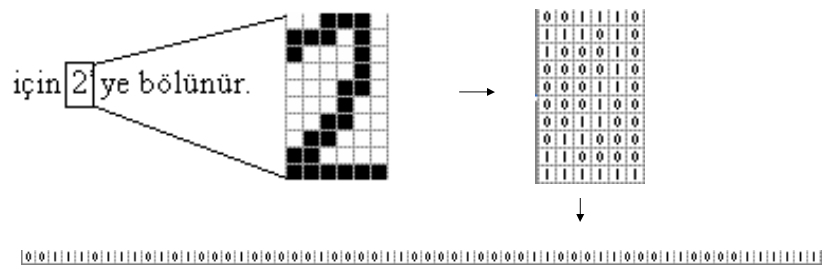
Resim	Resmin her bir pikselinin renkli resimlerde R,G,B değerleri, siyah-beyaz resimlerde 1–255 arası gri seviyesi kullanılarak sayılara çevrilir. Renkli resimler 3 adet, siyah beyazlar 1 adet en*boy büyüklüğünde matrisle ifade edilir.
Metin	Metindeki harfler, heceler ve kelimeler genelde frekanslarına göre kodlanarak sayılara çevrilir.
Hareketli görüntü	Resim bilgisine ek olarak resmin hangi resimden sonra geldiğini gösteren zaman bilgisini de içerir. Bu ek bilgi haricinde yapılan işlem resim ile aynıdır.
Ses	Ses, genlik ve frekansın zaman içinde değişimiyle kodlanır.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Örnek sınıflandırma uygulaması

- Sistem: bir kitap fotokopisinin içindeki yazılarının metne dönüştürülmesi
- Öncelikle metindeki satırlar bulunur.
- Her bir satırdaki harfler bulunur. Her harfe ait onlarca örnek resimden etiketlenmiş bir veritabanı oluşturulur. Her bir resim için



- Bu şekilde tanınmak istenen harf için çeşitli fontlarla yazılmış birçok örneği temsil eden 60 boyutlu vektörler elde edilir.
- Bu uygulamamız için özellik sayımız 60'tır. Diğer bir deyişle örneklerimiz 60 boyutlu bir uzayda temsil edilmektedirler.
- Elimizde 10 rakama ait farklı fontlarla yazılmış 10'ar resim olursa veri kümemiz 100 örnek* 60 boyutluk bir matris olacaktır.
- Elimizde her örneğin hangi harf olduğunu gösteren sınıf bilgiside bulunmaktadır.
- Bu matris eğitim ve test kümesi oluşturmak için 2'ye bölünür.
- Eğitim kümesi bir sınıflandırıcıya verilir.
- Modellenir.
- Modelin başarısını ölçmek için sınıflandırıcının daha önce görmediği, modelini oluşturmakta kullanmadığı test kümesi için tahminde bulunması istenir.
- Bu tahminlerle gerçek sınıfların ayınlığının ölçüsü sınıflandırıcının başarı ölçüsüdür.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Akış

- Makine Öğrenmesi nedir?
- Günlük Hayatımızdaki Uygulamaları
- Verilerin Sayısallaştırılması
- **Özellik Belirleme**
 - Özellik Seçim Metotları
 - Bilgi Kazancı (Informaiton Gain-IG)
 - Sinyalin Gürültüye Oranı: (S2N ratio)
 - Alt küme seçiciler (Wrappers)
 - Yeni Özelliklerin Çıkarımı
 - Temel Bileşen Analizi (Principal Component Analysis)
 - Doğrusal Ayırteden Analizi (Linear Discriminant Analysis)
- Sınıflandırma Metotları
 - Doğrusal Regresyon
 - Karar Ağaçları (Decision Trees)
 - Yapay Sinir Ağları
 - En Yakın K Komşu Algoritması (k - Nearest Neighbor)
 - Öğrenmeli Vektör Kuantalama (Learning Vector Quantization)
- Kümeleme Algoritmaları:
 - K-Ortalama (K-Means)
 - Kendi Kendini Düzenleyen Haritalar (Self Organizing Map -SOM)

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Özellik Belirleme

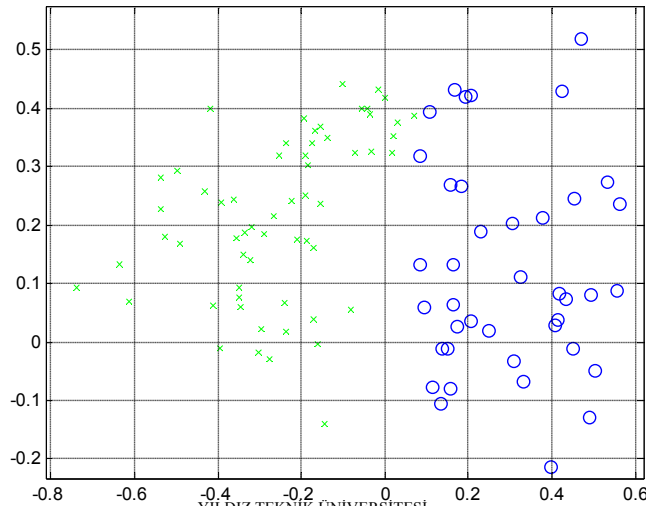
- Bir doktor
- Veri: Kişi bilgilerini içeren dosyalar
- Görev: Kimler hasta bul.
- Hangi bilgilere bakılır?
 - Ad soyad
 - Doğum yeri
 - Cinsiyet
 - Kan tahlili sonuçları
 - Röntgen sonuçları
 - vs.

1. Özellik	2. Özellik	Sınıf
1	3	A
2	3	B
1	4	A
2	3	B

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Hangi boyut?

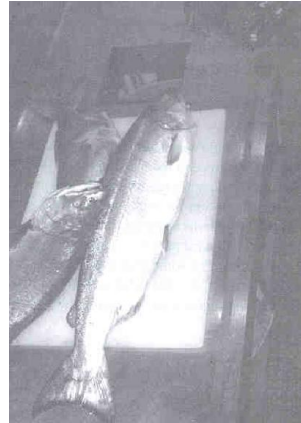


YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

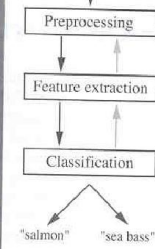


Balık Hali

- Kayan bant üzerindeki balığın türünü belirlemek (Salmon? Sea Bass?)



kamera

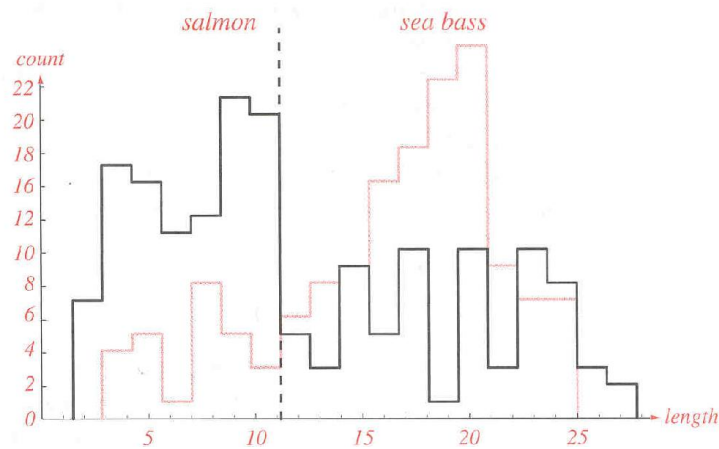


YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Balık Özellikleri. Uzunluk.

- Salmon lar genelde Sez Bass lardan daha kısalar.

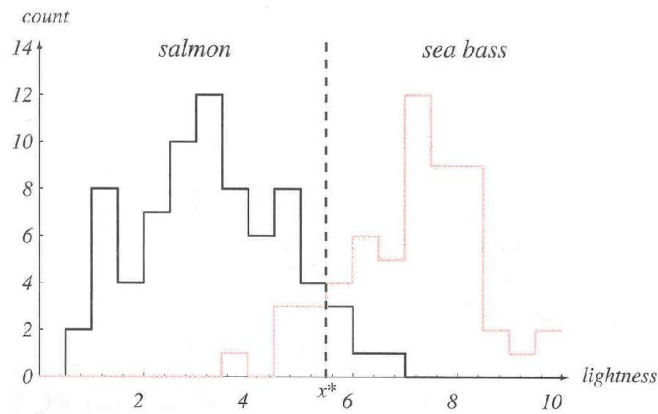


YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Balık Özellikleri. Parlaklık.

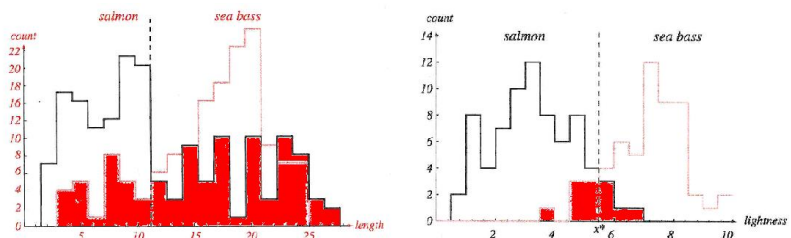
- Sea Bass genelde Salmon lardan daha parlaklar.



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Hangi Özellik?



Kırmızı bölgeler yapılan hataları gösteriyor.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Doktoru yoralım ☺

- Hastalık dosyasında 5000 adet özellik olsaydı?
Örneğin kişinin DNA dizisine bakarak hasta olup olmadığına karar verecek olsaydık ne yapardık?
Nerelere bakacağımıza nasıl karar verirdik.
- Burada devreye makineleri sokmamız gerekiyor gibi gözükmemekte.
- Bu olay bir insanın hesap yapma kabiliyetiyle, bir hesap makinesininkini karşılaştırmaya benziyor.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Özellik seçimi

- Bu problem makinelerle iki farklı metotla çözülebilir.
 - Var olan özelliklerden bazılarını seçmek
 - Özellikleri tek tek değerlendirmek (Filter)
 - Özellik alt kümeleri oluşturup, sınıflandırıcılar kullanıp performanslarını ölçüp, bu alt kümeleri en iyilemek için değiştirerek (Wrapper)
 - Var olan özelliklerin lineer birleşimlerinden yeni özelliklerin çıkarımı

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Özellikleri birer birer inceleme (Filters)

- Eğitim bilgilerindeki her bir özellik teker teker ele alınır.
- Örnek ile ilgili sadece o özellik elimizde olsaydı ne olurdu sorusunun cevabı bulunmaya çalışılır.
- Seçilen özellikle sınıf ya da sonucun birlikte değişimleri incelenir.
- Özellik değiştiğinde sınıf ya da sonuç ne kadar değişiyorsa o özelliğin sonuca o kadar etkisi vardır denilir.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Bilgi Kazancı - Informaiton Gain

S eğitim seti içindeki A özelliğinin

$$Gain(S, A) \equiv Entrophy(S) - \sum_{v \in Value(A)} \frac{|S_v|}{|S|} Entrophy(S_v)$$

N kavramının c farklı değeri varsa N'in entropisi,
N'in aldığı her değer olasılıkları kullanılarak

$$Entropy(N) = \sum_{i=1}^c -p_i \log_2 p_i$$

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



daha önceki hava, nem, rüzgar, su sıcaklığı gibi değerlere göre pikniğe gidip gitmeme kararı verilmiş 4 olay

Olay No	Hava	Nem	Rüzgar	Su sıcaklığı	Pikniğe gidildi mi?
1	güneşli	normal	güçlü	ılık	Evet
2	güneşli	yüksek	güçlü	ılık	Evet
3	yağmurlu	yüksek	güçlü	ılık	Hayır
4	güneşli	yüksek	güçlü	soğuk	Evet

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Her bir özelliğin piknik kavramı için bilgi kazancını bulalım

- Pikniğe gidildi mi? sorusunun iki cevabı vardır.
- Evet cevabının olasılığı $\frac{3}{4}$
- Hayır cevabının olasılığı $\frac{1}{4}$
- Dolayısıyla Pikniğin Entropi'si
- $E(\text{Piknik}) = -(3/4) \log_2(3/4) - (1/4) \log_2(1/4)$
= **0.811** olarak bulunur.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



- **Gain(Piknik,Hava)**= $0.811 - \left(\frac{3}{4}\right) \left(-\left(\frac{3}{3}\right) \log_2 \left(\frac{3}{3}\right) - 0\right) - \left(\frac{1}{4}\right) \left(0 - \left(\frac{1}{1}\right) \log_2 \left(\frac{1}{1}\right)\right) = \mathbf{0.811}$
- Hava özelliğinin IG'si hesaplanırken bulunan rakamların açıklamaları:
 $0.811 \rightarrow$ Pikniğe gitme olayının Entropisi
 $\left(\frac{3}{4}\right) \rightarrow$ havanın güneşli olma oranı
 $\left(\frac{3}{3}\right) \rightarrow$ hava güneşli iken pikniğe gidilme oranı
 $0 \rightarrow$ hava güneşli iken pikniğe gidilmeme oranı
 $\left(\frac{1}{4}\right) \rightarrow$ havanın yağmurlu olma oranı
 $0 \rightarrow$ hava yağmurlu iken pikniğe gidilme oranı
 $\left(\frac{1}{1}\right) \rightarrow$ hava yağmurlu iken pikniğe gidilmeme oranı

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



- **Gain(Piknik,Nem)**= $0.811 - \left(\frac{1}{4}\right) \left(-\left(\frac{1}{1}\right) \log_2 \left(\frac{1}{1}\right) - 0\right) - \left(\frac{3}{4}\right) \left(-\left(\frac{2}{3}\right) \log_2 \left(\frac{2}{3}\right) - \left(\frac{1}{3}\right) \log_2 \left(\frac{1}{3}\right)\right)$
 $= 0.811 - 0.688 = \mathbf{0.1225}$
- **Gain(Piknik,Rüzgar)**= $0.811 - \left(\frac{4}{4}\right) \left(-\left(\frac{3}{4}\right) \log_2 \left(\frac{3}{4}\right) - \left(\frac{1}{4}\right) \log_2 \left(\frac{1}{4}\right)\right)$
 $= 0.811 - 0.811 = \mathbf{0}$
- **Gain(Piknik,SuSıcaklığı)**= $0.811 - \left(\frac{3}{4}\right) \left(-\left(\frac{2}{3}\right) \log_2 \left(\frac{2}{3}\right) - \left(\frac{1}{3}\right) \log_2 \left(\frac{1}{3}\right)\right) - \left(\frac{1}{4}\right) \left(-\left(\frac{1}{1}\right) \log_2 \left(\frac{1}{1}\right)\right)$
 $= 0.811 - 0.688 = \mathbf{0.1225}$
- En büyük bilgi kazancına sahip özellik 'Hava'dır.
- Gerçek uygulamalarda ise yüzlerce özelliğin bilgi kazançları hesaplanır ve en büyük olanları seçilerek kullanılır.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



S2N

- sınıflar arası ayrılıkların fazla sınıf içi ayrılıkların az olan özellikler seçilir.

$$S_i = \frac{m_1 - m_2}{d_1 - d_2}$$

$m_1 \rightarrow$ sınıf1'deki i. özelliklerin ortalaması

$m_2 \rightarrow$ sınıf2'deki i. özelliklerin ortalaması

$d_1 \rightarrow$ sınıf1'deki i. özelliklerin standart sapması

$d_2 \rightarrow$ sınıf2'deki i. özelliklerin standart sapması

S değeri en yüksek olan özellikler

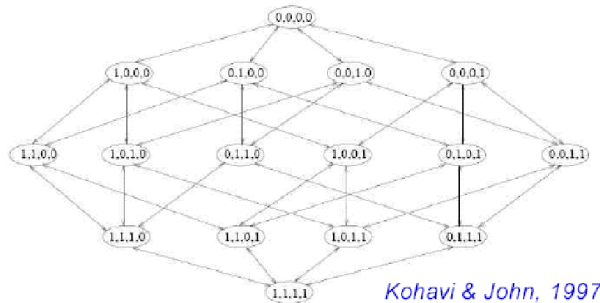
seçilerek sınıflandırmada kullanılırlar.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Özellik altkümesi seçiciler (Wrappers)

N özellik için olası 2^N özellik alt kümesi = 2^N eğitim



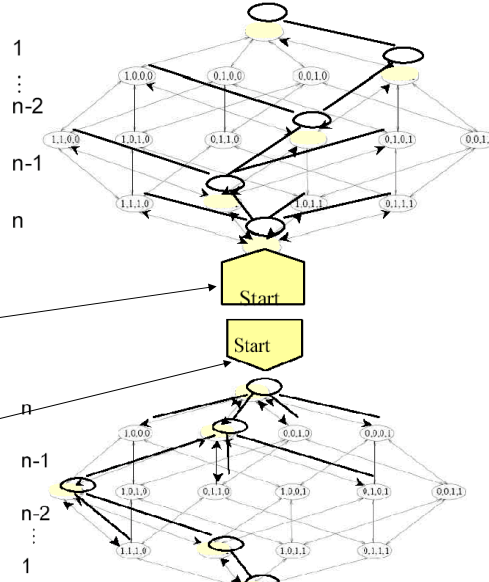
YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Özellik altkümesi seçiciler

- Hızlandırmak için tüm olasılıkları denemek yerine

- Hepsiyle başlayıp her seferinde bir tane elemek
- Tek özellekle başlayıp her seferinde bir tane eklemek



Hangi yoldan gidileceğine o özellik kümesinin sınıflandırmadaki performansına bakılarak karar verilir.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Yeni Özelliklerin Çıkarımı

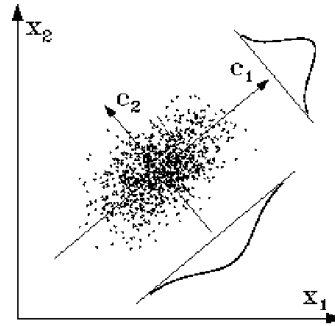
- Var olan özelliklerin lineer birleşimlerinden yeni bir özellik uzayı oluşturulur ve veriler bu uzayda ifade edilirler. Yaygın olarak kullanılan 2 metot vardır.
- PCA
- LDA

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Temel Bileşen Analizi-TBA (Principle Component Analysis - PCA)

- Bu metotta örneklerin en fazla değişim gösterdikleri boyutlar bulunur. Yansa veriler c_1 ve c_2 eksenlerine izdüşümü yapıldığındaki dağılımları gösterilmiştir.
- C_1 eksenindeki değişim daha büyüktür. Dolayısıyla veriler 2 boyuttan bir boyuta C_1 eksenine üzerine iz düşürülerek indirgenmiş olur.

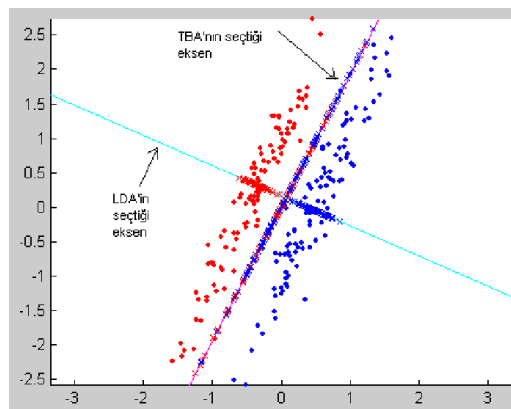


YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Doğrusal Ayırteden Analizi (Linear Discriminant Analysis - LDA)

Yandaki gibi durumlar için LDA önerilmiştir. LDA varyanslara ek olarak sınıf bilgisini de kullanarak boyut indirgene yapar. Sadece varyansa değil sınıflandırabilmeye de bakar.



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Hangisi

- Niye bu kadar çok metot var?
- Ne zaman hangisini kullanacağız?

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Akış

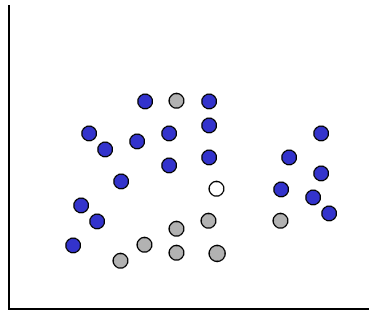
- Makine Öğrenmesi nedir?
- Günlük Hayatımızdaki Uygulamaları
- Verilerin Sayısallaştırılması
- Özellik Belirleme
 - Özellik Seçim Metotları
 - Bilgi Kazancı (Information Gain-IG)
 - Sinyalin Gürültüye Oranı: (S2N ratio)
 - Alt küme seçiciler (Wrappers)
 - Yeni Özelliklerin Çıkarımı
 - Temel Bileşen Analizi (Principal Component Analysis)
 - Doğrusal Ayırtmadan Analizi (Linear Discriminant Analysis)
- **Sınıflandırma Metotları**
 - Doğrusal Regresyon
 - Karar Ağaçları (Decision Trees)
 - Yapay Sinir Ağları
 - En Yakın K Komşu Algoritması (k - Nearest Neighbor)
 - Öğrenmeli Vektör Kuantalama (Learning Vector Quantization)
- Kümeleme Algoritmaları:
 - K-Ortalama (K-Means)
 - Kendi Kendini Düzenleyen Haritalar (Self Organizing Map -SOM)

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Sınıflandırma Metotları

Görev: Önceden etiketlenmiş örnekleri kullanarak yeni örneklerin sınıflarını bulmak



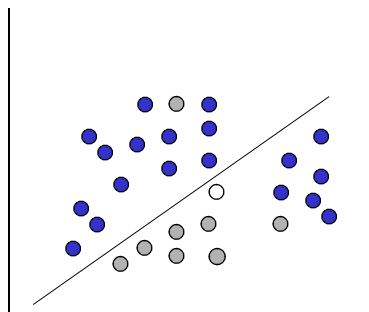
Metotlar:
Regresyon,
Karar Ağaçları,
LVQ,
Yapay Sinir Ağları,
...

Mavi ve gri sınıftan örnekler ● ●
Beyaz, mavi mi gri mi? ○

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Doğrusal Regresyon



- $w_0 + w_1 x + w_2 y \geq 0$
- Regresyon en az hata yapan w_i leri bulmaya çalışır.
- Basit bir model
- Yeterince esnek değil

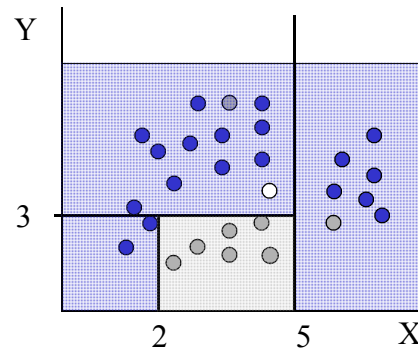
YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Karar Ağaçları

Böl ve yönet stratejisi

Nasıl böleceğiz?



if $X > 5$ then blue
 else if $Y > 3$ then blue
 else if $X > 2$ then green
 else blue

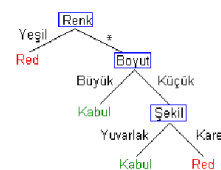
YILDIZ TEKNİK ÜNİVERSİTESİ
 BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Karar Ağaçları

- Ürettikleri kurallar anlaşılır.
- Karar düğümleri ve yapraklardan oluşan hiyerarşik bir yapı.

Şekil	Renk	Boyut	Sınıf
Yuvarlak	Yeşil	Küçük	Red
Kare	Siyah	Büyük	Kabul
Kare	Sarı	Büyük	Kabul
Yuvarlak	Sarı	Küçük	Red
Kare	Yeşil	Büyük	Red
Kare	Sarı	Küçük	Kabul



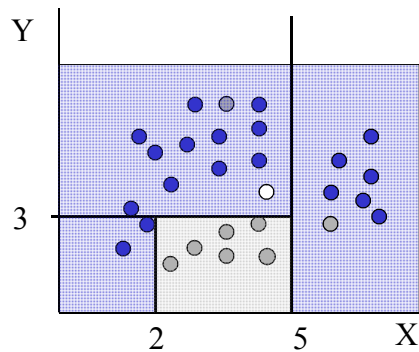
YILDIZ TEKNİK ÜNİVERSİTESİ
 BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Karar Ağaçları Oluşturma

- Tüm veri kümesiyle başla.
- Bir özelliğin bir değerlerine göre veri kümesi iki alt kümeye böl. Bölmede kullanılan özellikler ve değerleri karar düğüme yerleştir.
- Her alt küme için aynı prosedür her alt kümede sadece tek bir sınıfa ait örnekler kalıncaya kadar uygula.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Karar D ğ mleri Nasıl Bulunur?

- Karar d ğ mlerinde yer alan  zelliğ n ve eřik değ rinin belirlenmesinde genel olarak **entropi** kavramı kullanılır.
- Eđitim verisi her bir  zelliğ n her bir değ ri i in ikiye b l n r. Oluřan iki alt k menin entropileri toplanır. En d ř k entropi toplamına sahip olan  zellik ve değ ri karar d ğ m ne yerleřtirilir.

YILDIZ TEKNİK  NİVERSİTESİ
BİLGİSAYAR M HENDİSLİĐİ B L M 



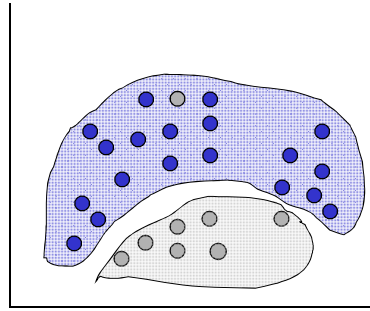
Karar Ađa larıyla Sınıflandırma

- En tepedeki k k karar d ğ m nden bařla.
- Bir yaprađa gelinceye kadar karar d ğ mlerindeki y nlendirmelere g re dallarda ilerle. (Karar d ğ mlerinde tek bir  zelliğ n adı ve bir eřik değ ri yer alır. O d ğ me gelen verinin hangi dala gideceğ ne verinin o d ğ mdeki  zelliğ nin eřik değ rinden b y k ya da k   k olmasına g re karar verilir.)
- Verinin sınıfı, yaprağ n temsil ettiğ  sınıf olarak belirle.

YILDIZ TEKNİK  NİVERSİTESİ
BİLGİSAYAR M HENDİSLİĐİ B L M 



Yapay Sinir Ağları

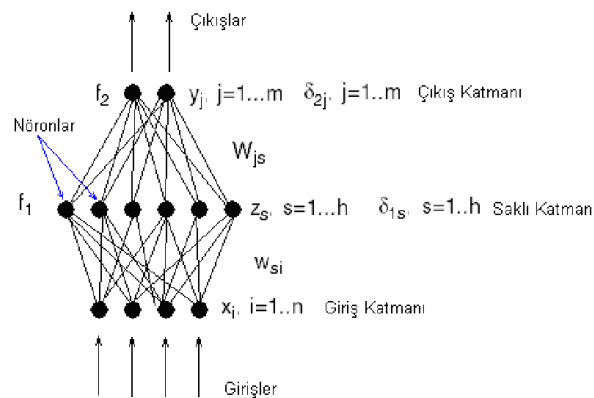


- Daha kompleks karar sınırlar üretebilirler.
- Can be more accurate
- Also can overfit the data – find patterns in random noise

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Yapay Sinir Ağları



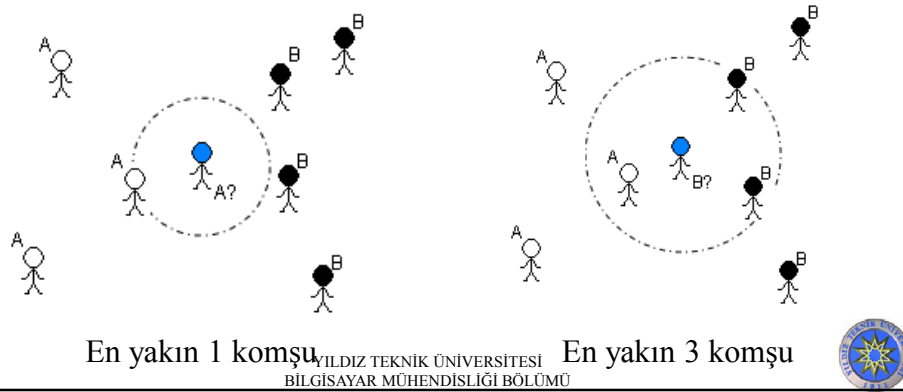
Elimizdeki Eğitim seti Girişler ve Çıkışlar ı içerir.
Bu girişler verildiğinde bu çıkışları verecek
Ağırlık değerlerini (W) bulmaya çalışır.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



En Yakın K Komşu

- Bana Arkadaşını söyle, sana kim olduğunu söyleyeyim.



En Yakın K Komşu

- Eğitim yok.
- Test verileri en yakınındaki K adet komşularının sınıf değerlerine bakılarak sınıflandırılırlar.



Öğrenmeli Vektör Kuantalama (Learning Vector Quantization)

[η] öğrenme oranı

[n] maximum eğitim sayısı

[c] betimleyici vektör sayısı

[μ_1, \dots, μ_c] betimleyici vektörler (centroids)

[x] eğitim verisinden bir örnek

[S(x)] x vektörünün ait olduğu yada betimlediği sınıf olmak üzere

1. $\eta, \mu_1, \mu_2, \dots, \mu_c$ için ilk değer atamalarını gerçekleştir

2. Aşağıdaki işlemleri n defa tekrar et

2.1 x eğitim verisini al

2.2 x'e en yakın betimleyici vektörü bul

$(\mu_k) : k \leftarrow \arg \min_j \|x - \mu_j\| \quad j=1..c$

2.3 μ_k nin güncellenmesi:

Eğer x doğru sınıfta ($s(x)=s(\mu_k)$) sınıfları aynı ise)

$\mu_k \leftarrow \mu_k + \eta(x - \mu_k)$ ödüllendir x'e yaklaştır

değilse

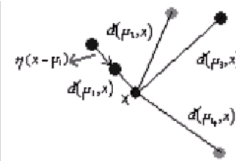
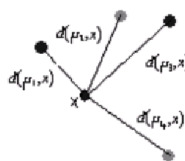
$\mu_k \leftarrow \mu_k - \eta(x - \mu_k)$ cezalandır x'den uzaklaştır

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

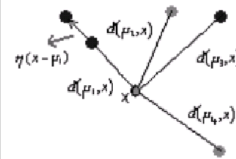
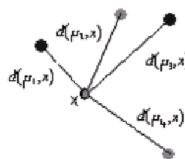


LVQ'da eğitim

LVQ'da ödüllendirme
Kazanan vektörle, örnek aynı sınıftan (ikisi de siyah sınıftan)



LVQ'da cezalandırma
Kazanan vektörle, örnek farklı sınıflardan (kazanan siyah, örnek gri sınıftan)

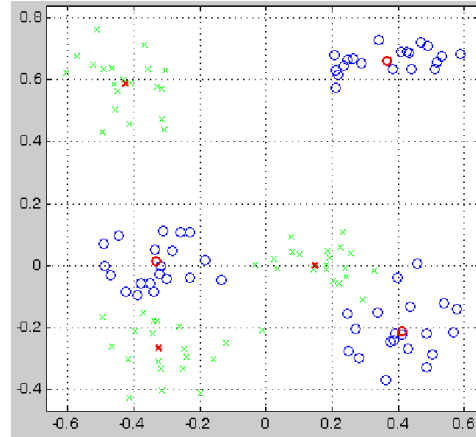


YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



LVQ- Test İşlemi

- Eğitim sonucu bulunan 2 sınıfa ait 3'er betimleyici vektör.
- Test işlemi, test örneğinin bu 6 vektörden en yakın olanının sınıfına atanmasıdır.



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

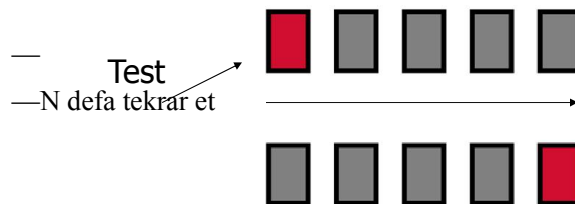


Çapraz Geçerleme

—Tüm dataseti eşit boyutlu N gruba böl



—Bir grubu test için geriye kalanların hepsini eğitim için kullan

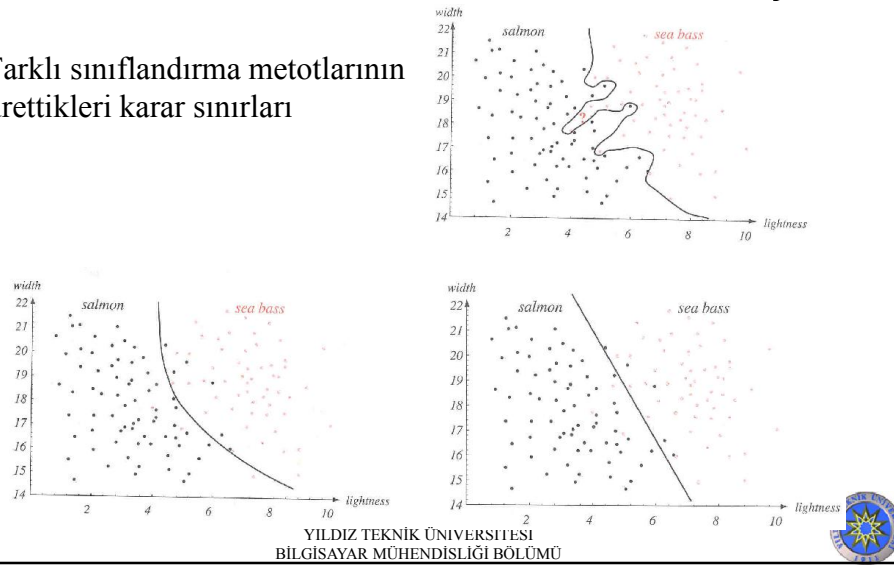


YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Sınıflandırma Metotları- Sonuç

Farklı sınıflandırma metotlarının ürettikleri karar sınırları



Sınıflandırma Metotları- Sonuç

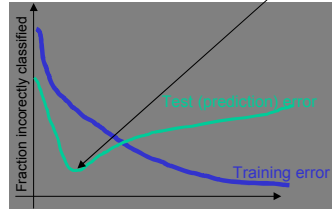
- Neden bu kadar çok algoritma var?
- Ne zaman hangisini seçeceğiz?

dataset	amlall	ann	bi75ds3	derma	gkanser	Hava
Özellik sayısı	7129	21	470	34	30	34
Sınıf sayısı	2	3	9	6	2	2
Örnek sayısı	72	3772	315	286	456	281
NB	97,14	95,55	68,49	77,97	94,29	89,31
SVM	92,86	93,74	62,11	79,37	96,26	86,48
1NN	94,29	93,4	63,19	76,26	96,26	89,72
C45	83,39	99,58	65,01	75,2	93,62	91,82
RF	95,71	99,5	72	76,96	95,38	95,02

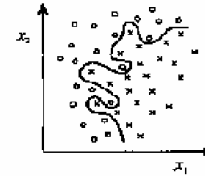
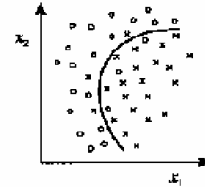
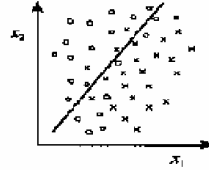
YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Modelim karmaşıklığı arttığında eğitim kümesindeki hata düşerken test kümesindeki hata yükselir.

Her veri kümesi için optimum nokta (optimum karmaşıklık) farklıdır.



Model karmaşıklığı



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Akış

- Makine Öğrenmesi nedir?
- Günlük Hayatımızdaki Uygulamaları
- Verilerin Sayısallaştırılması
- Özellik Belirleme
 - Özellik Seçim Metotları
 - Bilgi Kazancı (Information Gain-IG)
 - Sinyalin Gürültüye Oranı: (S2N ratio)
 - Alt küme seçiciler (Wrappers)
 - Yeni Özelliklerin Çıkarımı
 - Temel Bileşen Analizi (Principal Component Analysis)
 - Doğrusal Ayırtıcı Analizi (Linear Discriminant Analysis)
- Sınıflandırma Metotları
 - Doğrusal Regresyon
 - Karar Ağaçları (Decision Trees)
 - Yapay Sinir Ağları
 - En Yakın K Komşu Algoritması (k - Nearest Neighbor)
 - Öğrenmeli Vektör Kuantalama (Learning Vector Quantization)
- **Kümeleme Algoritmaları:**
 - K-Ortalama (K-Means)
 - Kendi Kendini Düzenleyen Haritalar (Self Organizing Map -SOM)

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Kümeleme Algoritmaları

- Kümeleme algoritmaları eğitimcişiz öğrenme metotlarıdır.
- Örneklere ait sınıf bilgisini kullanmazlar.
- Temelde verileri en iyi temsil edecek vektörleri bulmaya çalışırlar.
- Verileri temsil eden vektörler bulunduğtan sonar artık tüm veriler bu yeni vektörlerle kodlanabilirler ve farklı bilgi sayısı azalır.
- Bu nedenle birçok sıkıştırma algoritmasının temelinde kümeleme algoritmaları yer almaktadır.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Kümeleme Algoritmaları

- Bir boyutlu (özellikli) 10 örnek içeren bir veri
12-15-13-87-4-5-9-67-1-2
- Bu 10 farklı veriyi 3 farklı veriyle temsil etmek istersek:
12-12-12-77-3-3-3-77-3-3
- şeklinde ifade edebiliriz.
- Kümeleme algoritmaları bu 3 farklı verinin değerlerini bulmakta kullanılırlar.
- Gerçek değerlerle temsil edilen değerler arasındaki farkları minimum yapmaya çalışırlar.

Yukarıdaki örnek için 3 küme oluşmuştur.

- 12-15-13 örnekleri 1. kümede
- 87-67 örnekleri 2. kümede
- 4-5-1-2-9 örnekleri 3. kümede yer almaktadır.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Renk Kümeleme

Quantization process



Resimdeki farklı renk sayısı 106846'dan 55'e indirilmiş.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Resim Kümeleme



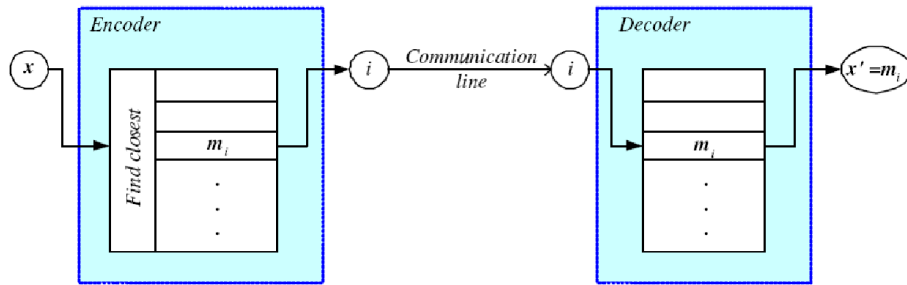
10*10 luk blokları ifade eden
vektörler kümelanmış

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Nasıl Kullanılır?

Bulunan (renkleri yada blokları temsil eden) küme merkezlerinden bir kod kitabı oluşturulur. Bu kitap her iki merkeze verilir. Vektörlerin kendileri yerine sadece indisler kullanılır. İndisin maximum büyüklüğü kodlanması için gereken bit sayısını artırır. Bu yüzden farklı vektör sayısının az olması istenir.



ETHEM ALPAYDIN © The MIT Press, 2004

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



K-means

Works with numeric data only

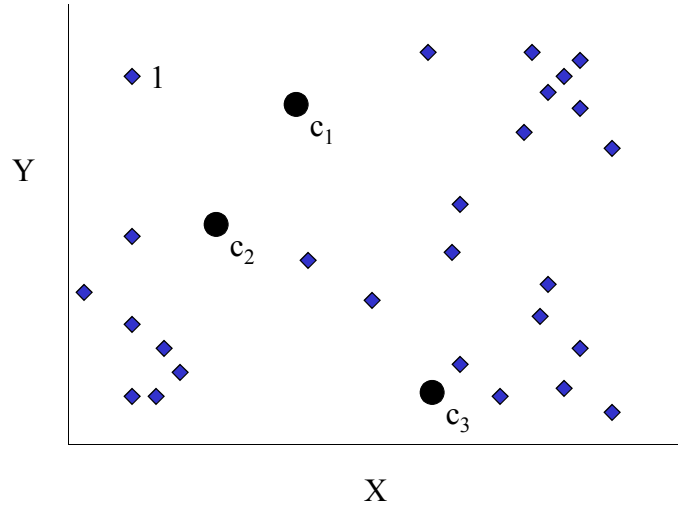
- 1) Rasgele K adet küme merkezi ata
- 2) Her örneği en yakınındaki merkezin kümesine ata
- 3) Merkezleri kendi kümelerinin merkezine ata
- 4) 2. ve 3. adımları küme değiştiren örnek kalmayıncaya kadar tekrar et.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



K-means örnek adım 1

Rasgele
3 küme
merkezi
ata.

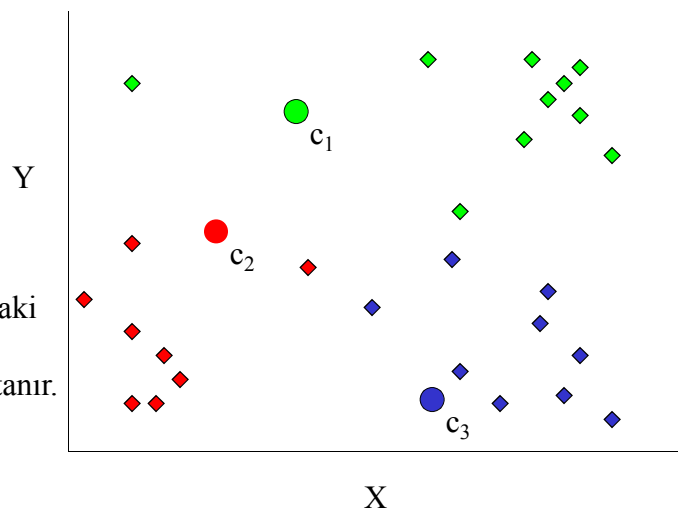


YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



K-means örnek adım 2

Her örnek
en yakınındaki
merkezin
kümesine atanır.

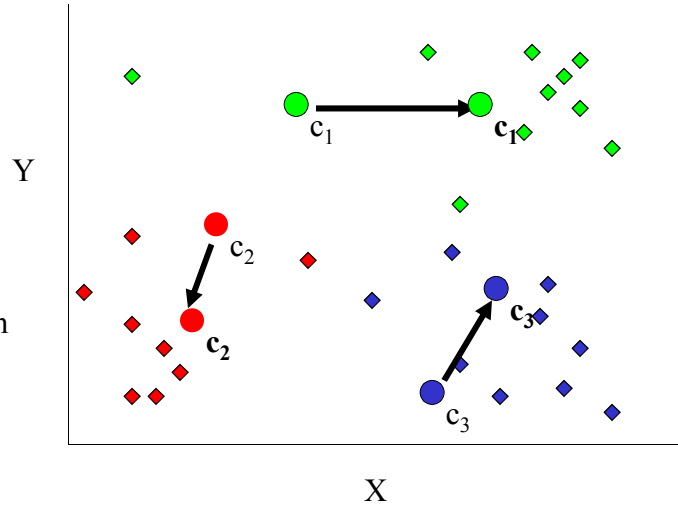


YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



K-means örnek adım 3

Merkezleri
kendi
kümelerinin
merkezine
götür.



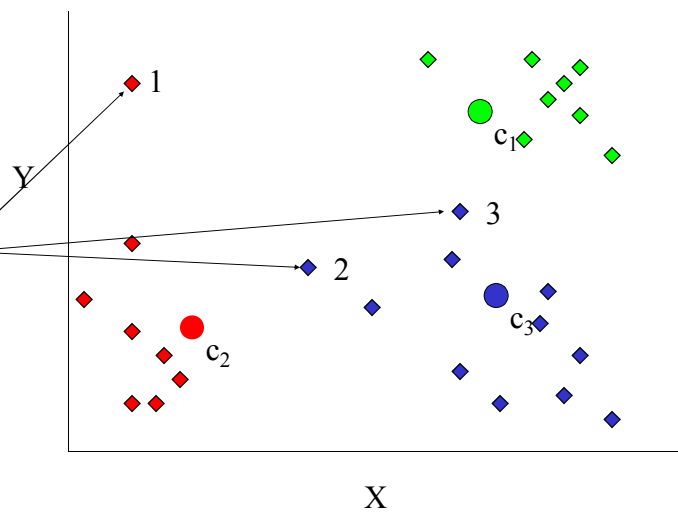
YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



K-means örnek adım 4

Her örneği
yeniden en
yakınındaki
merkezin
kümesine
ata.

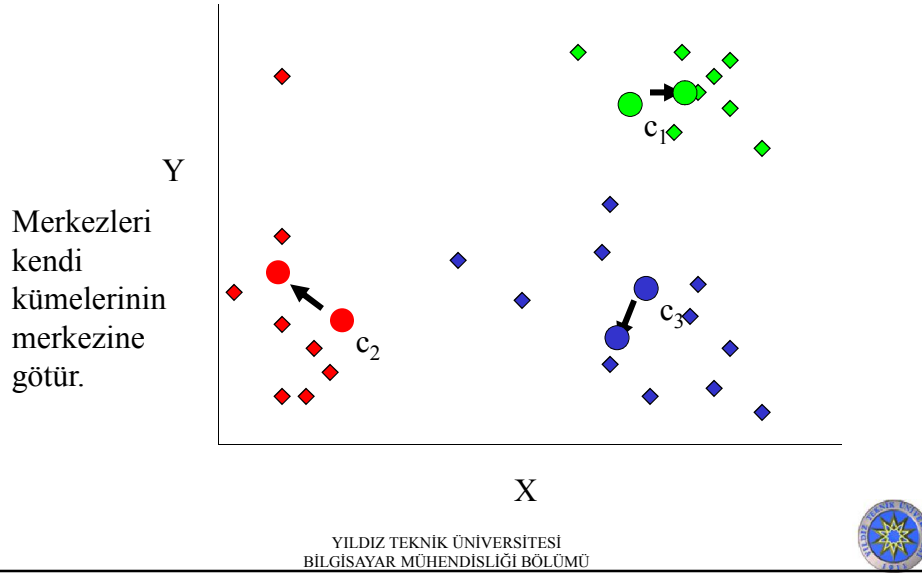
*Q: Hangi
örneklerin
kümesi
değişti?*



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



K-means örnek adım 5



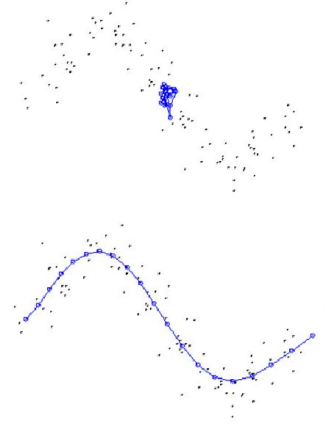
Kendi Kendini Düzenleyen Haritalar Self Organizing Maps

- Kmeans algoritmasında merkez noktalar arasında herhangi bir ilişki yoktur. SOM'da ise merkez noktalar 1 ya da 2 boyutlu bir dizi içinde yer alırlar. Buna göre birbirlerine 1 ya da 2 boyutlu uzayda komşudurlar.
- Kmeans algoritmasında sadece kazanan (en yakın) merkez güncellenirken SOM'da bütün merkezler kazanan nörona komşuluklarına göre güncellenir. Yakın komşular uzak komşulara göre daha fazla hareket ederler (güncellenirler).
- Merkezlerin birbirlerine bağlı oluşu verinin 1 ya da 2 boyutlu uzaydaki yansımasının da elde edilmesini sağlar.



SOM

- SOM merkezleri 1 boyutlu bir dizide birbirlerine komşudurlar. Başlangıçtaki durumları rasgele atandığı için bir yumak şeklindedirler. Eğitim tamamlandığında ise SOM merkezleri verinin şeklini almıştır.



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Sonuç olarak

- makineler insanlığın işgücüne sağladıkları katkıyı, makine öğrenmesi metotları sayesinde insanlığın beyin gücüne de sağlamaya başlamışlardır.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Bir gün bilgisayarlar bizi anarlarsa?

Ve bütün bunları mükemmel bir şekilde yaparlarsa Nasıl bir dünya

- Bir sürü işsiz bilgisayar mühendisi ☺
- Bir sürü işsiz insan
- ???

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Kaynaklar

- Alpaydın E. (2004) "Introduction to Machine Learning", The MIT Press, 3-6
- <http://www.autonlab.org/tutorials/infogain11.pdf>
- http://www.kdnuggets.com/dmcourse/data_mining_course/assignments/assignment-4.html
- http://pespmc1.vub.ac.be/asc/SENSIT_ANALY.html
- http://csnet.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf
- http://www.cavs.msstate.edu/hse/ies/publications/reports/isip_internal/1998/linear_discrim_analysis/lda_theory.pdf
- <http://www.kernel-machines.org>
- T.Kohonen, " Self-Organization and associative Memory", 3d ed, 1989, Berlin :Springer-Verlag.
- <http://www.willamette.edu/~gorr/classes/cs449/Classification/perceptron.html>
- O. T. Yıldız, E. Alpaydın, Univariate and Multivariate Decision Trees, Tainn 2000
- <http://www.ph.tn.tudelft.nl/PHDTheses/AHoekstra/html/node45.html>
- <http://mathworld.wolfram.com/K-MeansClusteringAlgorithm.html>
-

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Weka



Copyright: Martin Kramer (mkramer@wxs.nl)

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

