

Makine Öğrenmesi-2

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

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:
 - Hiyerarşik Kümeleme
 - K-means
 - Kendi Kendini Düzenleyen Haritalar (Self Organizing Map -SOM)
 - DBscan
- Regresyon Algoritmaları
- Çok Boyutlu Verilerle Çalışmak
- Veri Sızıntısı
- Pekiştirmeli Öğrenme

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Genel Dataset Yapısı supervised (sınıflandırma / regresyon)

öz1	öz2	...	Öz d	Çıkış
Örnek1 öz1	Örnek1 öz2	...	Örnek 1 öz d	Örnek1 çıkış
Örnek2 öz1	Örnek2 öz2	...	Örnek 2 öz d	Örnek2 çıkış
ÖrnekN öz1	ÖrnekN öz2	...	Örnek N öz d	ÖrnekN çıkış

Özellik tipleri

- Kategorik (K adet ten 1 i)
- Sayı

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.

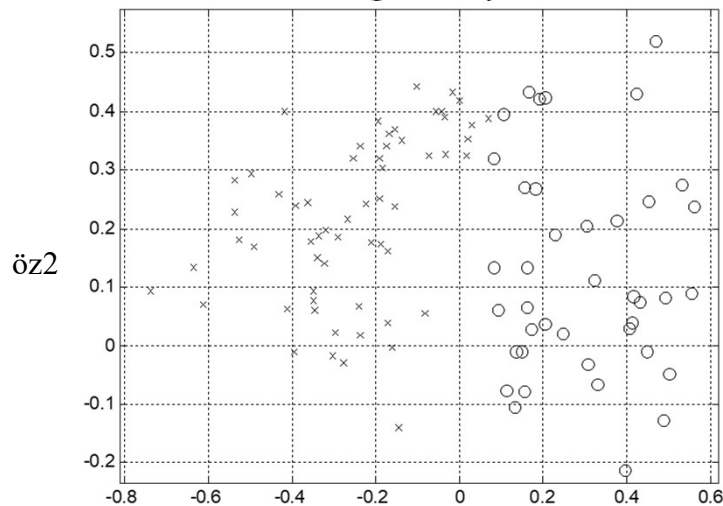
1. Öz	2. Öz	3. Öz	Sınıf
1	3	13	A
2	3	13	B
1	4	13	A
2	3	13	B

If 1.öz==1 then çıkış=A
Else çıkış=B (%100)

If 2.öz==4 then çıkış=A
Else çıkış=B (%) 75

çıkış=B (%50)

Hangi boyut?



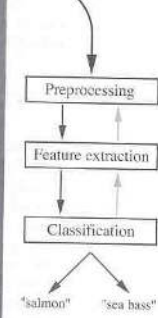
If öz1>0.1 then çıkış=o else çıkış=x

Balık Hali

- Kayan bant üzerindeki balığın türünü belirlemek
- Salmon / Sea Bass
Somon / Levrek



kamera



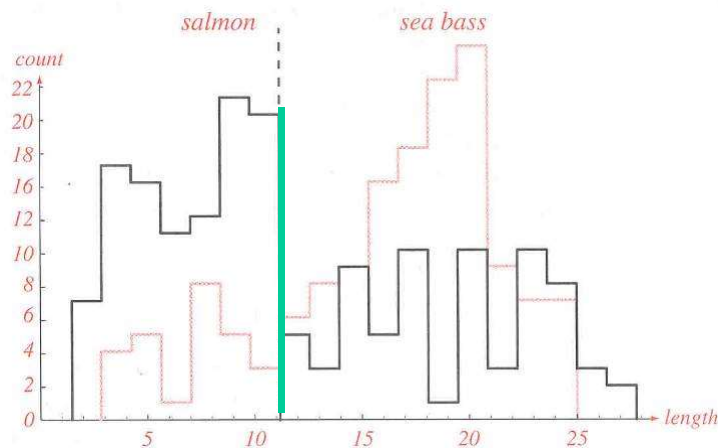
[*] Pattern Classification (2nd Edition) by Richard O. Duda, Peter E. Hart, David G. Stork

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Balık Özellikleri: Uzunluk.

- Salmon lar genelde Sea Bass lardan daha kısalar.



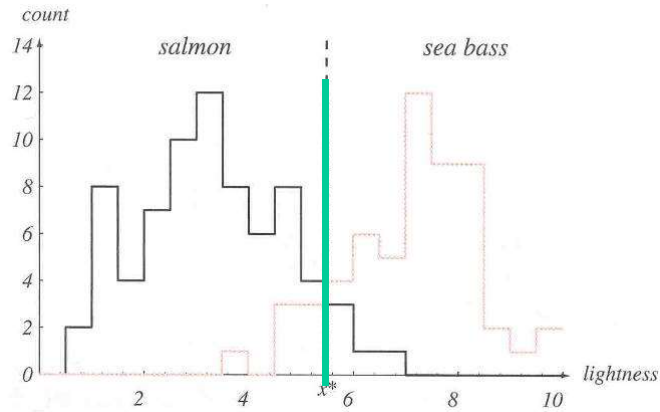
[*] Pattern Classification (2nd Edition) by Richard O. Duda, Peter E. Hart, David G. Stork

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Balık Özellikleri: Parlaklık.

- Sea Bass genelde Salmon lardan daha parlaklar.

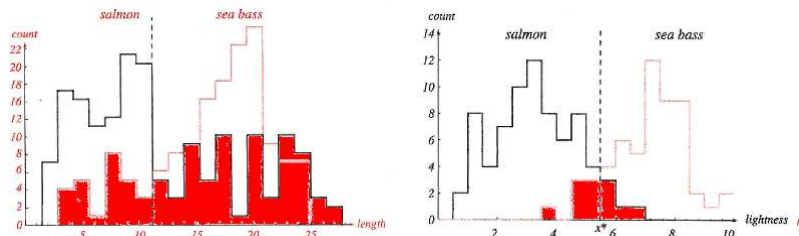


[*] Pattern Classification (2nd Edition) by Richard O. Duda, Peter E. Hart, David G. Stork

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Hangi Özellik?



Kırmızı bölgeler yapılan hataları gösteriyor.

[*] Pattern Classification (2nd Edition) by Richard O. Duda, Peter E. Hart, David G. Stork

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

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.

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

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.

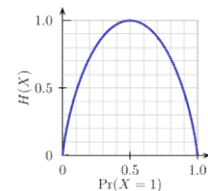
Bilgi Kazancı - Information Gain

S eğitim seti içindeki A özelliğinin

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

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

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



[*] Machine Learning, Tom Mitchell

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

$$\text{Entropi}(3.\text{öz}) = -1 \cdot \log_2(1) = 0$$

$$C=1$$

$$P(13)=1$$

$$\text{Entropi}(\text{sınıf}) = -$$

$$[(1/2) \cdot \log_2(1/2) + (1/2) \cdot \log_2(1/2)]$$

$$= -[(1/2) \cdot -1 + (1/2) \cdot -1] = 1$$

$$C=2$$

$$P(A)=1/2$$

$$P(B)=1/2$$

3. Öz	Sınıf
13	A
13	B
13	A
13	B

K=A A A B B D için

$$C=3$$

$$\text{Entropi}(K) = - (1/2 \cdot \log_2(1/2) + 1/3 \cdot \log_2(1/3) + 1/6 \cdot \log_2(1/6))$$

$$P_a=1/2$$

$$P_b=1/3$$

$$P_d=1/6$$

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

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
- **$E(\text{Piknik}) = -(\frac{3}{4}) \log_2(\frac{3}{4}) - (\frac{1}{4}) \log_2(\frac{1}{4}) = 0.811$ olarak bulunur.**

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

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

- **Gain(Piknik,Nem)**= 0.811- (1/4) (- (1/1) log₂ (1/1) - 0) – (3/4) (- (2/3) log₂(2/3)- (1/3) log₂(1/3))
= 0.811 -0.688= **0.1225**
- **Gain(Piknik,Rüzgar)**= 0.811- (4/4) (- (3/4) log₂(3/4) – (1/4) log₂(1/4))
= 0.811 -0.811= **0**
- **Gain(Piknik,SuSıcaklığı)**= 0.811- (3/4) (- (2/3) log₂(2/3) – (1/3) log₂(1/3)) – (1/4) (- (1/1) log₂ (1/1))
= 0.811 -0.688= **0.1225**
- En büyük bilgi kazancına sahip özellik ‘Hava’dır.
- Gerçek uygulamalarda ise yüzlerce özelliğin bilgi kazançları hesaplanır ve en büyük olanları seçilerek kullanılır.

Signal2Noise Ratio (Sinyalin Gürültüye Oranı)

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

$$S_i = \frac{|m_1 - m_2|}{d_1 + d_2} \begin{matrix} \text{Sinyal} \\ \text{Gürültü} \end{matrix}$$

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

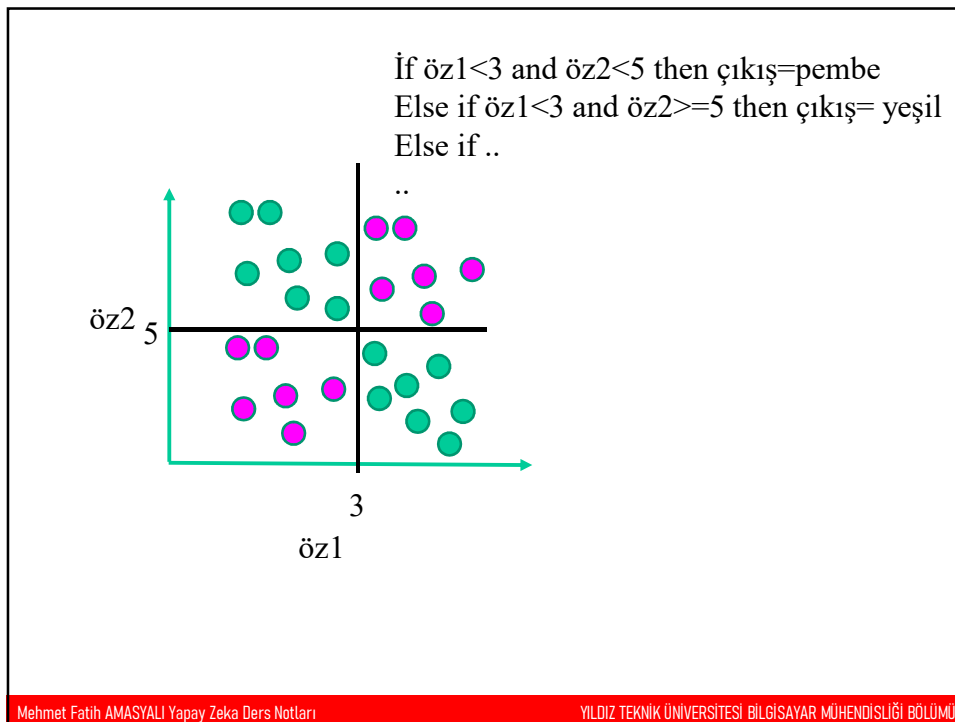
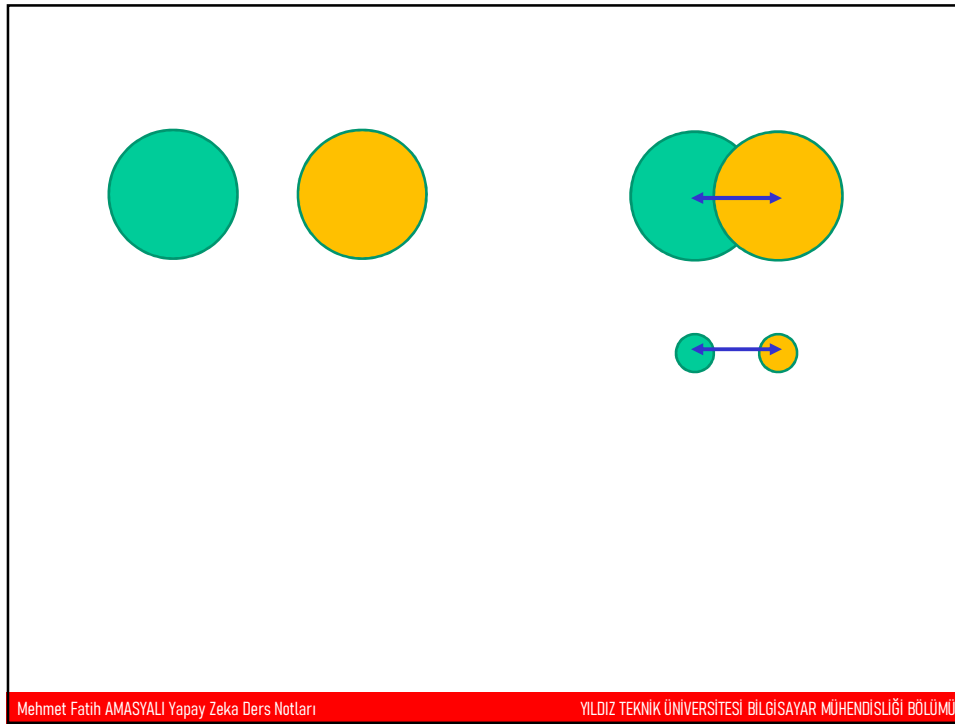
$m_2 \rightarrow$ sınıf2’deki i. özelliklerin ortalaması

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

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

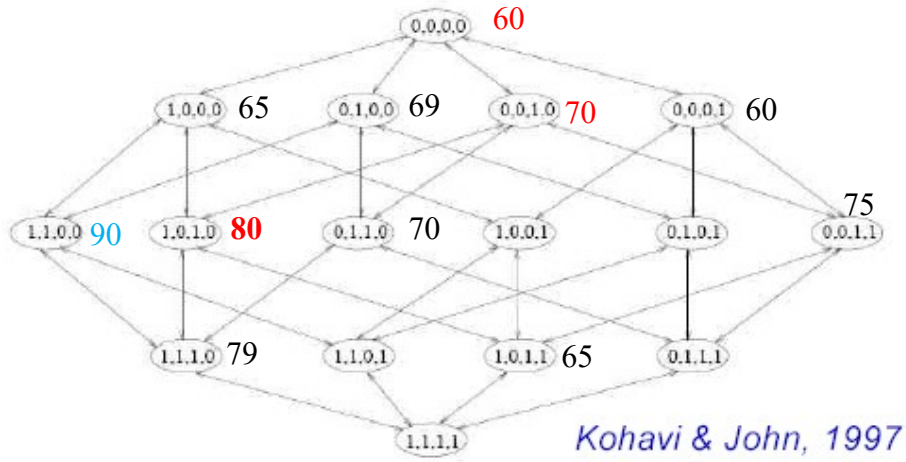
S değeri en yüksek olan özellikler

seçilerek sınıflandırmada kullanılırlar.



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

Wrappers (Özellik altkümesi seçiciler)



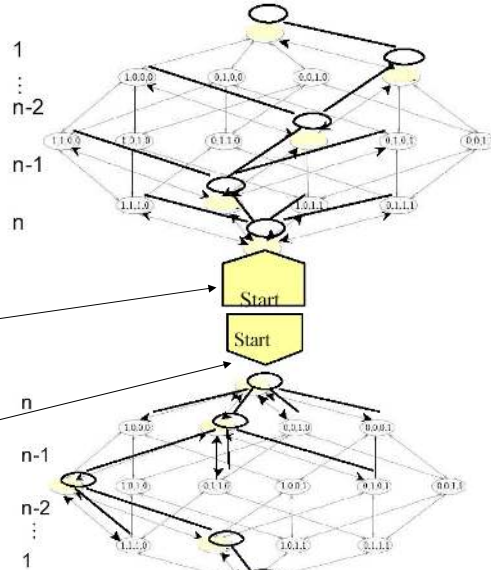
Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Özellik altkümesi seçiciler

- Hızlandırmak için tüm olasılıkları denemek yerine

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



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

[*] <http://clopinet.com/isabelle/Projects/ETH/lecture9.pdf>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

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.

Orj özellikler: x_1, x_2, x_3, x_4

Yeni özellikler:

$$Y_1 = w_{11} * x_1 + w_{12} * x_2 + w_{13} * x_3 + w_{14} * x_4 + w_{15}$$

$$Y_2 = w_{21} * x_1 + w_{22} * x_2 + w_{23} * x_3 + w_{24} * x_4 + w_{25}$$

$$Y_3 = w_{31} * x_1 + w_{32} * x_2 + w_{33} * x_3 + w_{34} * x_4 + w_{35}$$

$$Y_4 = w_{41} * x_1 + w_{42} * x_2 + w_{43} * x_3 + w_{44} * x_4 + w_{45}$$

- PCA
- LDA

öz1	öz2
3	1
2	5
3	6

Y1	Y2
18	7
25	22
33	27

$$w_{11}=5 \quad w_{12}=3 \quad w_{13}=0$$

$$w_{21}=1 \quad w_{22}=4 \quad w_{23}=0$$

$$Y_1 = w_{11} * \text{öz1} + w_{12} * \text{öz2} + w_{13}$$

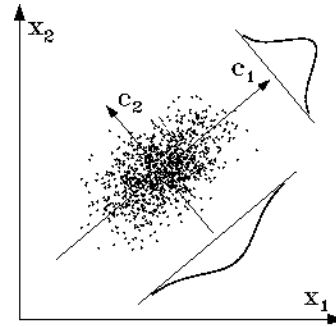
$$Y_2 = w_{21} * \text{öz1} + w_{22} * \text{öz2} + w_{23}$$

$$Y_1 = 5 * \text{öz1} + 3 * \text{öz2} + 0$$

$$Y_2 = 1 * \text{öz1} + 4 * \text{öz2} + 0$$

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

- Bu metotta örneklerin en fazla değişim gösterdikleri boyutlar bulunur. Yanda 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 iz düşürülerek indirgenmiş olur.

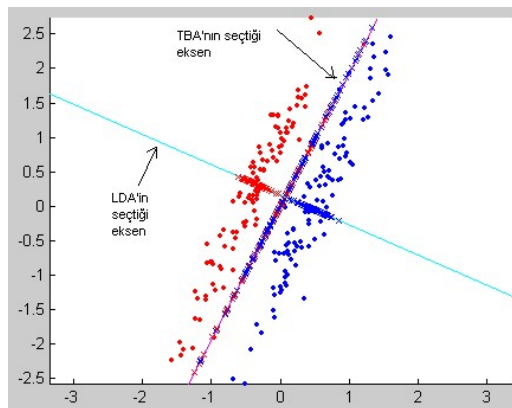


Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Doğrusal Ayırte Den 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.



Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

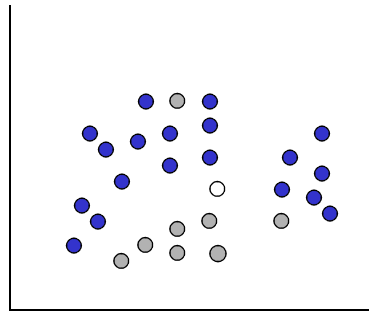
YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Hangisi

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

Sınıflandırma Metotları

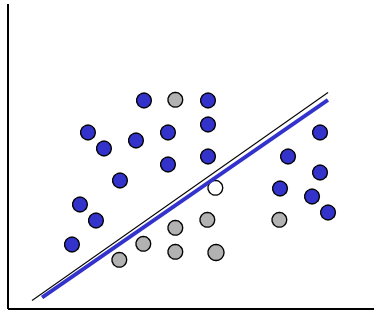
Görev: Önceden etiketlenmiş örnekleri kullanarak yeni örneklerin sınıflarını bulmak



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

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

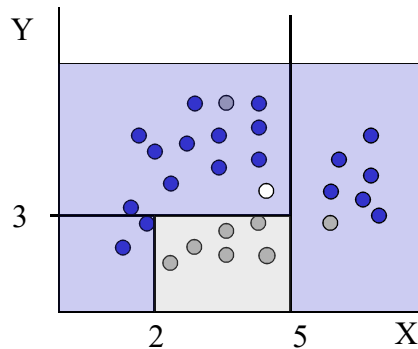
Doğrusal Regresyon



- $w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d \geq 0$
- En az hata yapan w_i leri bulmaya çalışır.
- Basit bir model
- Yeterince esnek değil

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

Ürettikleri kurallar anlaşılır.

Karar düğümleri ve yapraklardan oluşan hiyerarşik bir yapı.

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.

Karar ağacı oluşturalım

5	o	x	x	x	$X1 < 1.5$ e1	$X2 < 1.5$ e8
4	o	x	x	x	$X1 < 2.5$ e2	$X2 < 2.5$ e9
3	o	o	o	x	$X1 < 3.5$ e3	$X2 < 3.5$ e10
2	o	o	o	x	$X2 < 1.5$ e4	$X2 < 4.5$ e11
1	x	x	o	o	$X2 < 2.5$ e5	$X2 < 3.5$ e6
	1	2	3	4	$X2 < 4.5$ e7	$X1 < 2.5$ e26
						$X1 < 3.5$ e36
						$X2 < 1.5$ e46
						$X2 < 2.5$ e56
					$X1 < 2.5$ e27	$X2 < 2.5$ e57
					$X1 < 3.5$ e37	$X2 < 3.5$ e66
					$X2 < 1.5$ e47	$X2 < 4.5$ e76
					$X1 < 2.5$ e29	$X1 < 2.5$ e28
					$X1 < 3.5$ e39	$X1 < 3.5$ e38
						$X2 < 2.5$ e58

Karar Düğümleri Nasıl Bulunur?

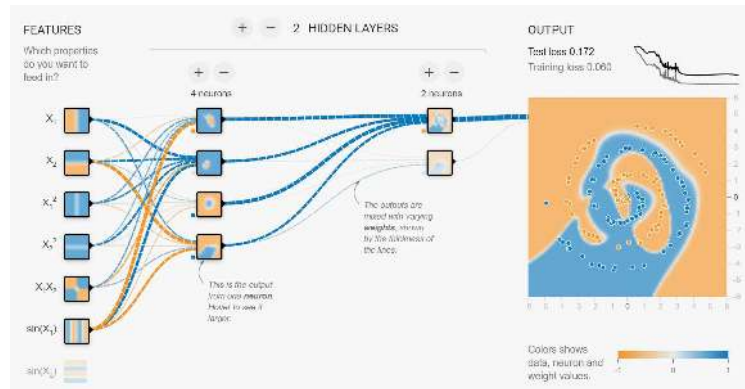
- Karar düğümlerinde yer alan özelliğin ve eşik değerinin belirlenmesinde genel olarak **entropi** kavramı kullanılır.
- Eğitim verisi her bir özelliğin her bir değeri için ikiye bölünür. Oluşan iki alt kümenin entropileri toplanır. En **düşük** entropi toplamına sahip olan özellik , değer ikilisi karar düğüme yerleştirilir.

Karar Ağaçlarıyla Sınıflandırma

- En tepedeki kök karar düğüminden başla.
- Bir yaprağa gelinceye kadar karar düğümlerindeki yönlendirmelere göre dallarda ilerle. (Karar düğümlerinde tek bir özelliğin adı ve bir eşik değeri yer alır. O düğüme gelen verinin hangi dala gideceğine verinin o düğümdaki özelliğinin eşik değerinden büyük ya da küçük olmasına göre karar verilir.)
- Verinin sınıfı, yaprağın temsil ettiği sınıf olarak belirle.

Yapay Sinir Ağları

- Daha kompleks karar sınırlar üretebilirler*
- Elimizdeki Eğitim seti Girişler ve Çıkışlar ı içerir. Bu girişler verildiğinde bu çıkışları verecek ağırlık değerlerini (W) bulmaya çalışır.



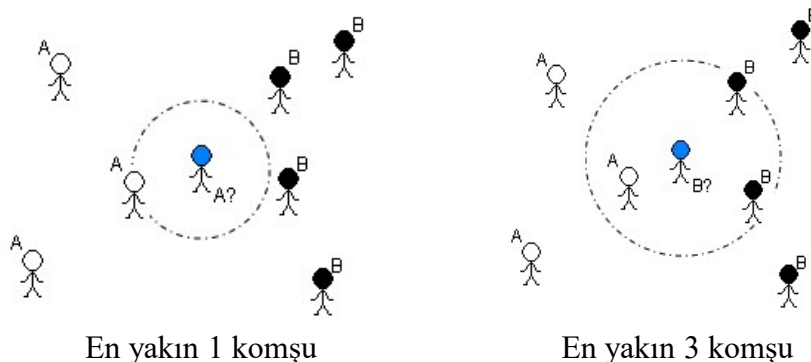
[*] <https://playground.tensorflow.org/>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

En Yakın K Komşu

- Bana Arkadaşını söyle, sana kim olduğunu söyleyeyim.



Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

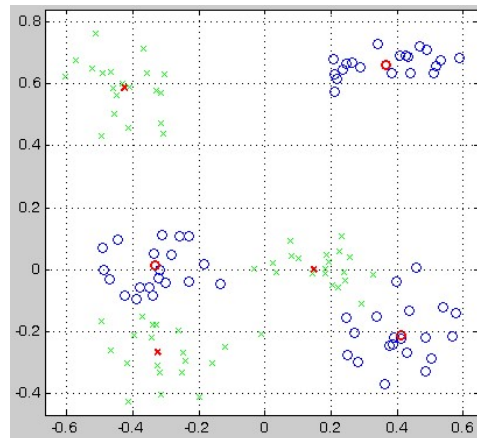
YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

En Yakın K Komşu

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

1NN'deki Tüm Noktaları kullanmasak?

- Tüm noktalar yerine yeni/ seçilmiş noktalar
- Bu noktalar nasıl belirlenebilir?



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 datasından bir örnek

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

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

2. Aşağıdaki işlemleri n defa tekrar et

2.1 X eğitim datasını 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 nın güncellenmesi:

Eğer x doğru sınıfsa ($S(x)=S(\mu_k)$ sınıfları aynı ise)

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

değilse

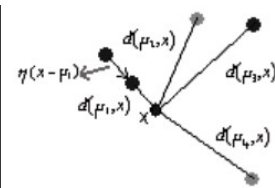
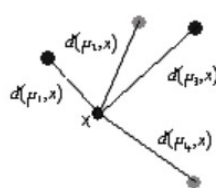
$\mu_k \leftarrow \mu_k - \eta(x - \mu_k)$ cezalandır x 'den uzaklaştır

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

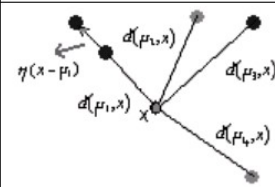
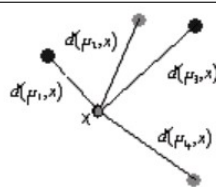
YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

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)

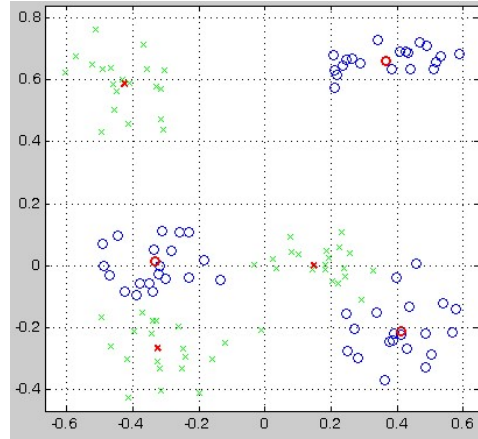


Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

LVQ- Test İşlemi

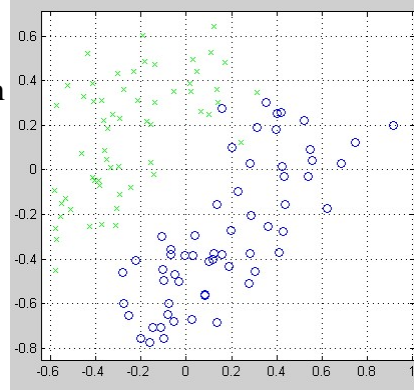
- Eğitim sonucu bulunan 2 sınıfa ait 3'er betimleyici vektör.
- Test işlemi, test örneğinin bu 6 vektörden en yakın olanının sınıfına atanmasıdır.



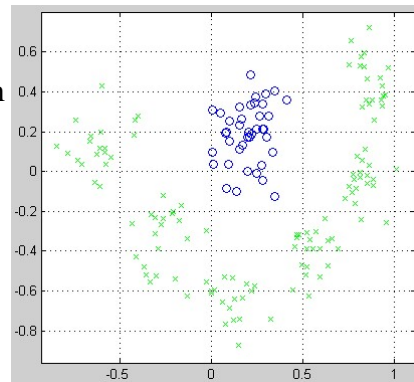
- LVQ ne zaman lineer/doğrusal karar sınırı üretir?

Hangisi?

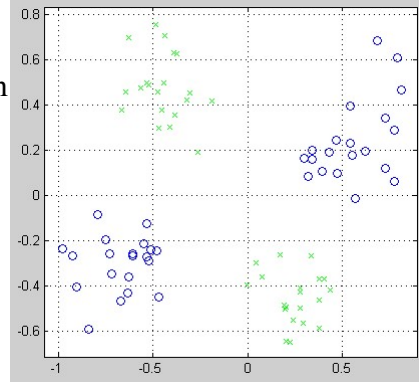
Lineer Regresyon
Karar Ağaçları
KNN
YSA
LVQ



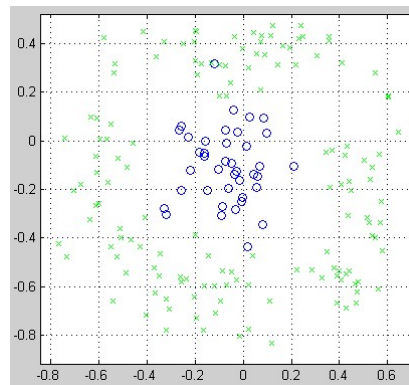
Lineer Regresyon
Karar Ağaçları
KNN
YSA
LVQ



Lineer Regresyon
Karar Ağaçları
KNN
YSA
LVQ



Başka bir model?



Eğitim kümesi üzerindeki performans neyi ifade eder?

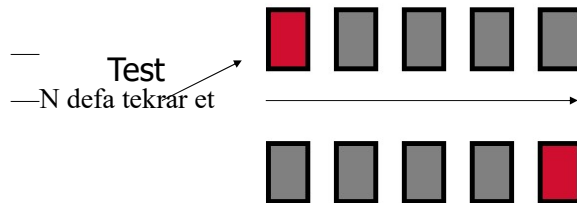
- Borsa oynayan maymunlar*
- 1000 kişiye hisse senedi tahminleri göndermek*
- * Hatasız Düşünme Sanatı'ndan

Hata nasıl ölçülür? Çapraz Geçerleme

—Tüm dataseti eşit boyutlu N gruba böl



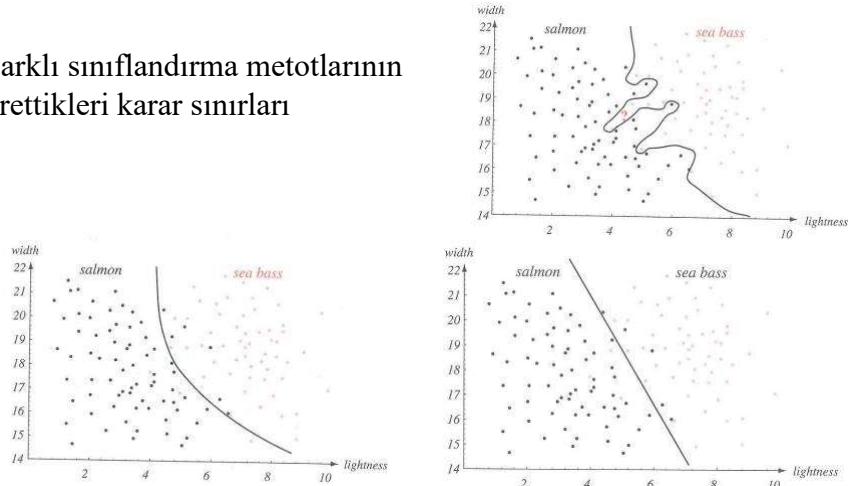
—Bir grubu test için geriye kalanların hepsini eğitim için kullan



—N defa tekrar et

Sınıflandırma Metotları- Sonuç

Farklı sınıflandırma metotlarının
ürettikleri karar sınırları



[*] Pattern Classification (2nd Edition) by Richard O. Duda, Peter E. Hart, David G. Stork

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Sınıflandırma Metotları- Sonuç

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

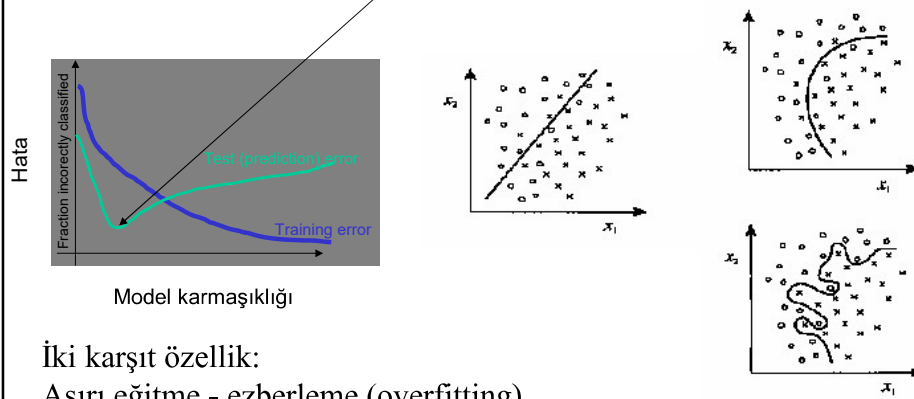
dataset	amlall	ann	bi75ds3	derma	gkanser	Hava
Özellik sayısı	7129	21	470	34	30	34
Sınıf sayısı	2	3	9	6	2	2
Örnek sayısı	72	3772	315	286	456	281
NB	97,14	95,55	68,49	77,97	94,29	89,31
SVM	92,86	93,74	62,11	79,37	96,26	86,48
1NN	94,29	93,4	63,19	76,26	96,26	89,72
C45	83,39	99,58	65,01	75,2	93,62	91,82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

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Modelin karmaşıklığı arttığında eğitim kümesindeki hata sürekli düşerken, test kümesindeki hata bir noktadan sonra yükselir.

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



İki karşıt özellik:
Aşırı eğitme - ezberleme (overfitting)
Genelleştirebilme

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Kümeleme Algoritmaları

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

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

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.

Renk Kümeleme



256 Farklı renkten ... 2 renge dönüşüm

*] http://fourier.eng.hmc.edu/e161/lectures/digital_image/node2.html

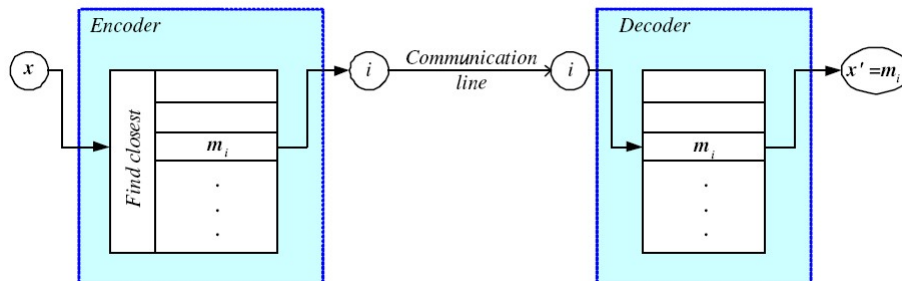
Resim Kümeleme



10*10 luk blokları ifade eden vektörler kümelenmiş

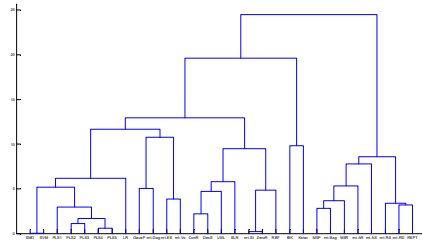
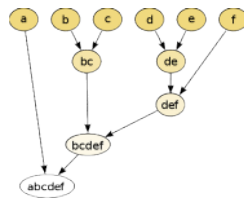
Nasıl Kullanılır?

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.



Hiyerarşik Kümeleme

- Tek kümeden çok kümeye (bölmeli)
- Çok kümeden tek kümeye (eklemeli)



Algoritmaların başarılarına göre kümelenmesi

[*] https://en.wikipedia.org/wiki/Hierarchical_clustering

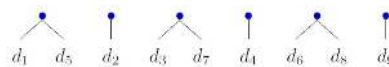
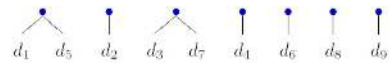
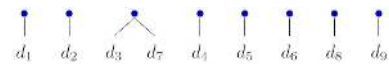
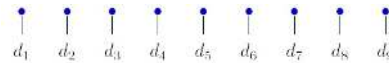
Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Hiyerarşik Eklemeli Kümeleme

- Birbirine en benzeyen iki kümeyi birleştir
- Tekrar et

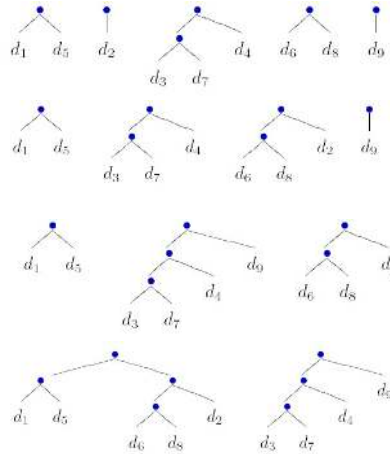
Başlangıçta küme sayısı =
örnek sayısı



Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Hiyerarşik Eklemeli kümeleme



Sonuçta küme sayısı=2

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Kümelerin birbirine benzerliği

- İki kümenin benzerliği
 - En benzer elemanları (Single link)
 - En benzemeyen elemanları (Compete link)
 - Ortalamaları (Group average)

kullanılarak bulunabilir.

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Bölmeli Hiyerarşik Kümeleme

- Tek bir kümeyle başla
- Küme içinde birbirine en az benzeyen iki elemanı bul.
- Kümeyi bu iki elemana yakınlığa göre böl.
- Oluşan her alt küme için bu işlemi tekrar et.

Ne zaman duracağız

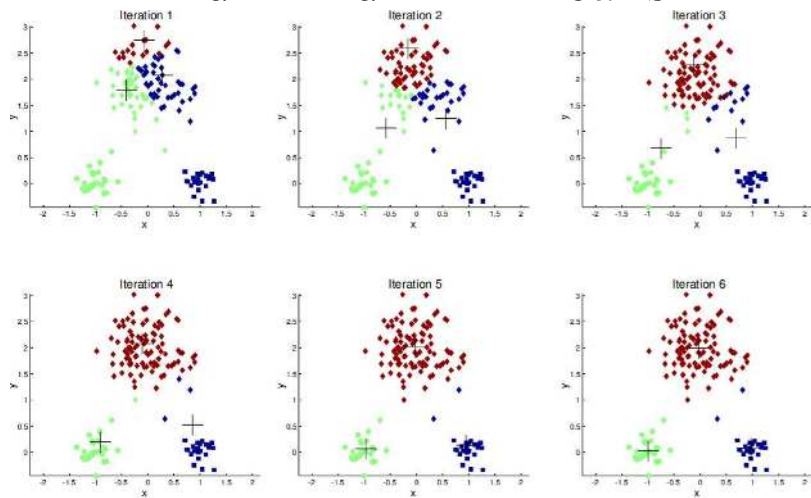
- İstenen küme sayısına ulaşıncaya
- Önceden belirlenmiş bir toplam benzerlik eşik değerine göre

K-means

Sadece sayısal veriyle çalışır

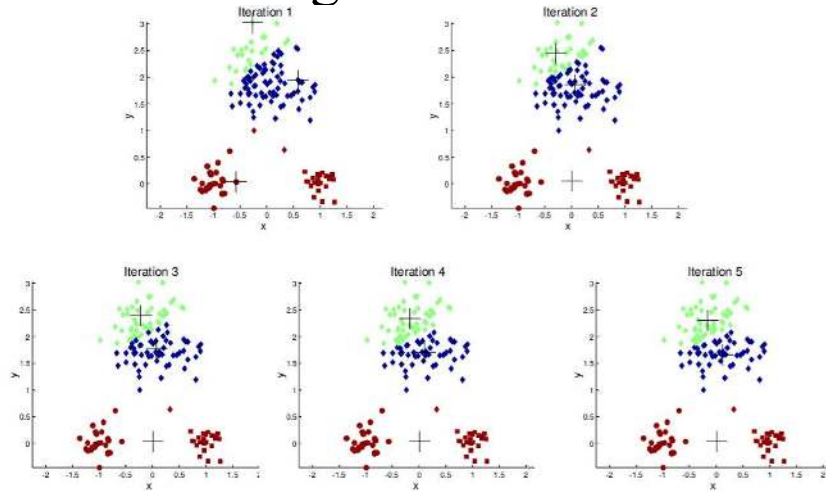
- 1) Rasgele K adet küme merkezi ata
- 2) Her örneği en yakınındaki merkezin kümesine ata
- 3) Merkezleri kendi kümelerinin merkezine ata
- 4) 2. ve 3. adımları küme değiştiren örnek kalmayıncaya kadar tekrar et.

Adım Adım K-means



[*] <https://cs.wmich.edu/alfuqaha/summer14/cs6530/lectures/ClusteringAnalysis.pdf>

İlk değerlerin önemi



[*] <https://cs.wmich.edu/alfuqaha/summer14/cs6530/lectures/ClusteringAnalysis.pdf>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

K-means Hata Fonksiyonu

- k adet küme: C_1, C_2, \dots, C_k
- c_i : i. küme merkezi
- x : örnekler
- Farklı ilk değerlerin sonuçları bu hataya göre karşılaştırılır.

$$Cost(C) = \sum_{i=1}^k \sum_{x \in C_i} (x - c_i)^2$$

[*] <https://cs.wmich.edu/alfuqaha/summer14/cs6530/lectures/ClusteringAnalysis.pdf>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Tekli (Online) vs. Toplu (Batch) K-means

- Önceki örnekte yapılan işlem toplu (batch) k-means
- Eğer her bir örnek için mean'ler güncellenirse → tekli k-means

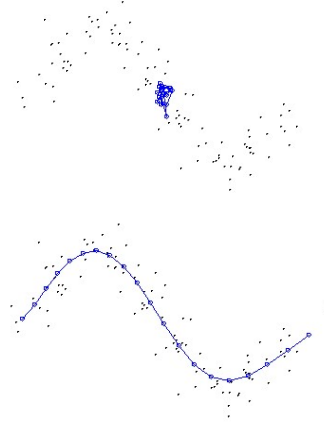
Kendi Kendini Düzenleyen Haritalar Self Organizing Maps*

- 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.

[*] Kohonen, Teuvo. "The self-organizing map." *Proceedings of the IEEE* 78.9 (1990): 1464-1480.

SOM

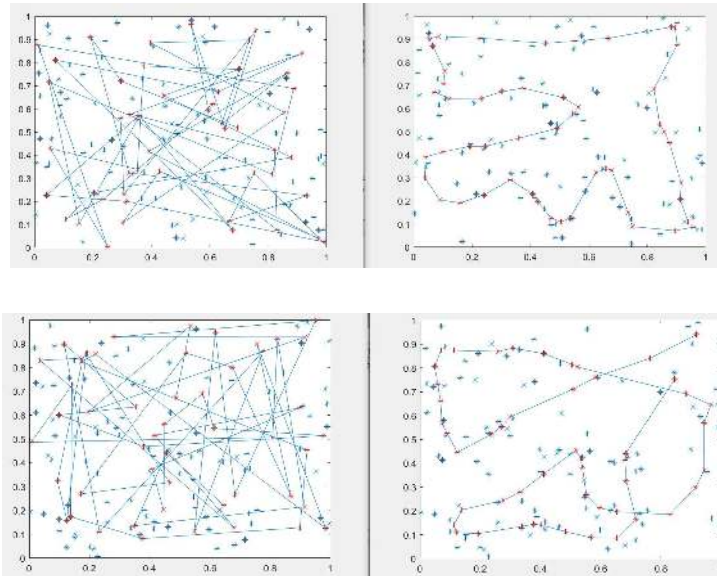
- Örnekte SOM merkezleri 1 boyutlu bir dizide birbirlerine komşudurlar. Başlangıçtaki durumları rasgele atandığı için bir yumak şeklindedirler. Eğitim tamamlandığında ise SOM merkezleri verinin şeklini almıştır.



Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Kod: my_som.m

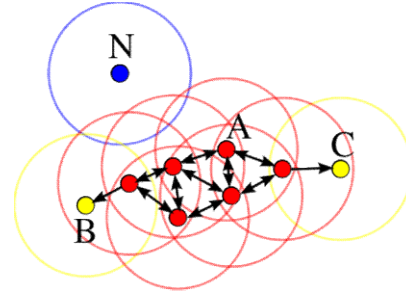


Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

DBscan

- Yoğunluk tabanlı kümeleme
- Hiper-parametreleri: ϵ , \min_p
- Önce tüm noktaları çekirdek, sınır, gürültü olarak belirle
- Çekirdek (kırmızı): ϵ komşuluğunda en az \min_p adet nokta bulunan noktalar
- Gürültü (mavi): ϵ komşuluğunda hiç nokta bulunmayan noktalar
- Sınır (sarı): ϵ komşuluğunda \min_p 'den az nokta bulunan noktalar



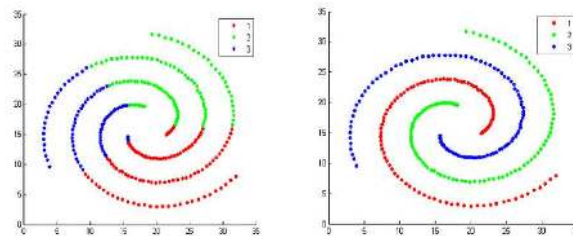
[*] <https://en.wikipedia.org/wiki/DBSCAN>

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

DBscan

- Noktalar etiketlendikten sonra gürültüler elenir
- Çekirdeklerden birbirine ϵ mesafede olanlar birleştirilir.
- X (hiper parametrelere bağlı) adet küme oluşur.
- Sınır noktalar en yakın çekirdeğin kümesine atanır



[*] https://www.researchgate.net/figure/Results-of-a-k-means-b-DBSCAN-on-Jain-dataset_fig2_291831757

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Kümeleme algoritmalarında

- Ne zaman dururuz?
- Ölü nöron kavramı (SOM, K-means)
- Başarı ölçütü nedir?
- **Küme sayısı nasıl belirlenir?**
 - Mecburiyet (verilmiştir)
 - ?

Regresyon Algoritmaları

- Basit lineer regresyon
- kNN
- Regresyon ağaçları
- YSA

Basit lineer regresyon

X: girişler ($n \times d$)

Y: çıkışlar ($n \times 1$)

Model: $Y = XW$

W: katsayılar

Amaç W yi bulmak

Her 2 tarafı X^{-1} ile çarpalım.

$X^{-1}Y = X^{-1}XW$ ($X^{-1}X = I$)

$W = X^{-1}Y$ (tamam, ama X kare matris değilse?)

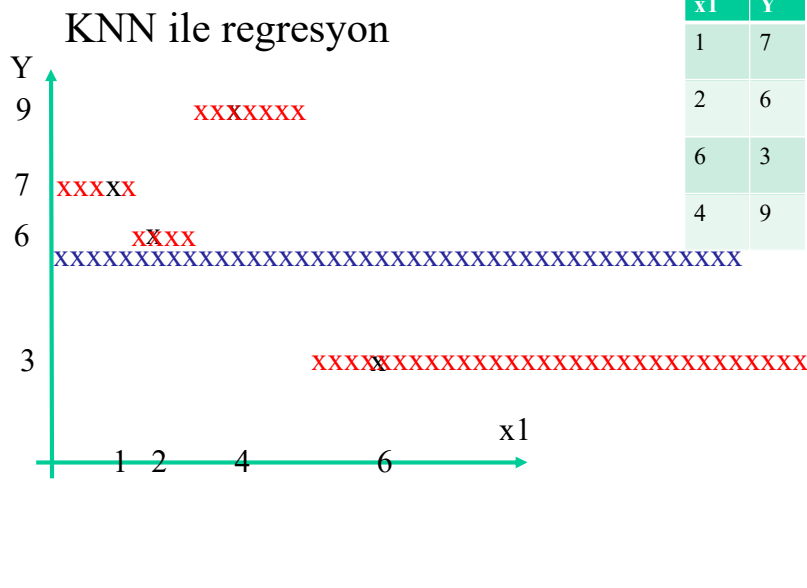
$X^T Y = X^T X W$ ($X^T X$ her zaman kare matristir)

$(X^T X)^{-1} (X^T Y) = (X^T X)^{-1} (X^T X) W$ [$(X^T X)^{-1} (X^T X) = I$]

$W = (X^T X)^{-1} (X^T Y)$

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

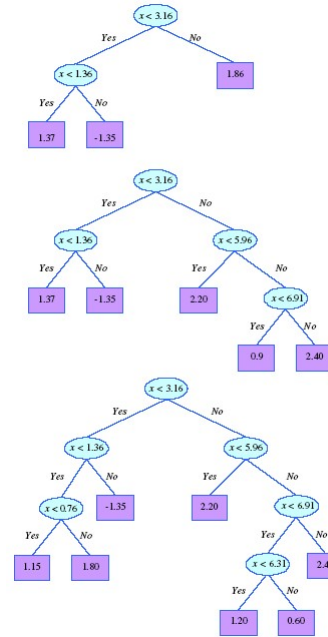
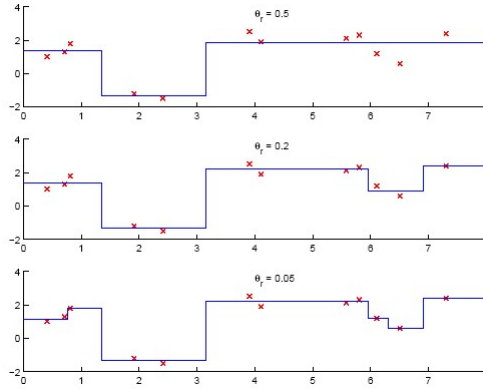
YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Regresyon ağaçlarında kabul edilen hata değerinin etkisi

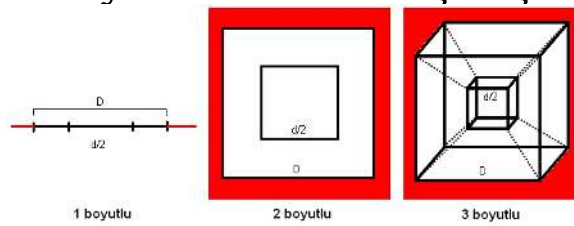


ETHEM ALPAYDIN © The MIT Press, 2004

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Çok Boyutlu Verilerle Çalışmak-1



Boyut Sayısı	Merkeze daha yakın noktaların oranı (%)
1	50
2	25
3	12,50
...	...
P	$(\frac{1}{2})^P$

Boyut sayısı arttığında verilerin çok büyük bir kısmı sınıfları ayıran sınırlara çok yakın yerlerde bulunacağından sınıflandırma yapmak zorlaşmaktadır.

Mehmet Fatih AMASYALI Yapay Zeka Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

Çok Boyutlu Verilerle Çalışmak-2

- Tek boyutlu uzayda $[0,1]$ aralığı temsil eden 10 nokta
- Rastgele bir noktanın, uzayı temsil eden noktalardan en yakın olanına ortalama uzaklığı = 0.5
- İki boyutlu uzayda rasgele bir noktanın en yakın noktaya olan ortalama uzaklığının düşey ya da dikey (manhattan) 0.5 olması için gerekli temsilci nokta sayısı = 100

Boyut Sayısı	Gerekli temsil eden nokta sayısı
1	10
2	100
3	1000
...	...
p	10^p

Doğru sınıflandırma yapmak için gereken örnek sayısı artıyor.

Veri Sızıntısı (Data Leakage)

- Tahmin sonuçlarının çok iyi görünmesine sebep olabilir, ama gerçekler uygulamada ortaya çıkar
- Test kümesindeki bilgilerin eğitim sürecine karışması
 - Verileri normalize ederken, özellik seçimi yaparken test kümesini de kullanmak
- Test zamanı elde olmayacak özelliklerin kullanılması
 - $x(t)=f(x(t-1),x(t+1))$
- Zaman serisi sınıflandırmada eğitim ve test kümelerini oluşturmada hata
 - Rasgele seçim yapılmamalı, bir t anından öncesi eğitim, sonrası test olmalı ki ardışık süreçleri değil sınıfı tanısin.
- Verilerde çıkışla korelasyonu çok yüksek olan özelliklerin olması (kişi tanırken id, telefon no vb.)

Sonuç olarak

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

Bir gün bilgisayarlar
bütün bunları mükemmel bir şekilde
yaparlarsa Nasıl bir dünya

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

Kaynaklar

- Alpaydın E. (2004) "Introduction to Machine Learning", The MIT Press, 3-6
- <https://www.autonlab.org/resources/tutorials/information-gain>
- http://www.cs.otago.ac.nz/cosc453/student_tutorials/principal_components.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://mathworld.wolfram.com/K-MeansClusteringAlgorithm.html>
- Pattern Classification (2nd Edition) by Richard O. Duda,Peter E. Hart,David G. Stork
- Machine Learning, Tom Mitchell, McGraw Hill, 1997

Weka



Copyright: Martin Kramer (mkramer@wxs.nl)