

Veri Kaynakları

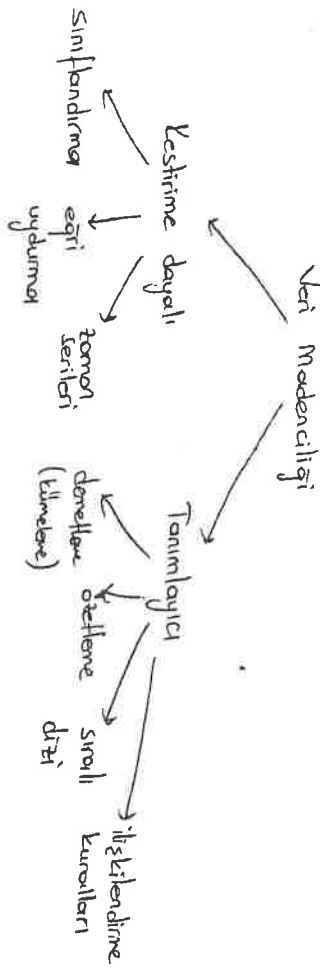
- * Veri dosyaları
- * Veri tabanı kaynaklı veri kümeleri
- * Gelişmiş veri kümeleri (sosyal ağ, www, konumsal veriler....)

Veri Madenciliği Algoritmaları

amaç: Veriyi belli bir modele uydurmak

- * Tanımlayıcı → en iyi müşteri kim?
- * Kestirime dayalı → borsa tahmini

Veri Madenciliği Modelleri



Sınıflandırma: Veriyi önceden belirlenmiş sınıflardan birine dahil eder.
Eğri Uydurma (Regression): Veriyi gerçek değerli bir fonk. dönüştürür.
Zaman serileri inceleme: zaman içinde değişen verinin değişimi belirler.
Denetleme: Benzer verileri aynı grupta toplama.
Özetleme: Veriyi alt gruplara ayırır, her alt grup için özellikler bulunur.
Sıralı dizi: Veri içinde sıralı ilişkileri bulmak için kullanılır.

VERİ MADENCİLİĞİ NOTLARI

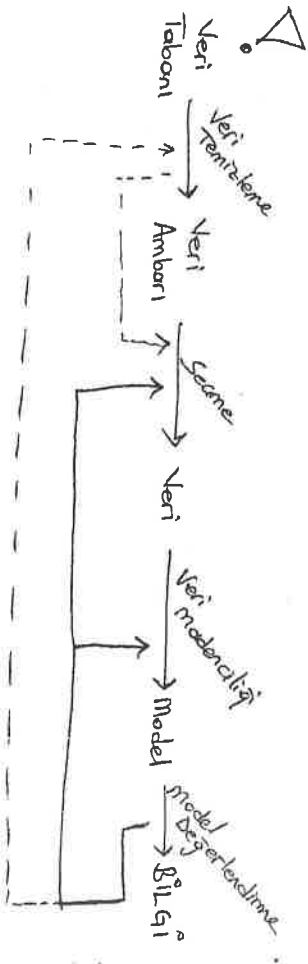
4

- * Teknolojinin gelişimiyle bilgisayar ortamında veri tabanında tutulan veri miktarının artması, kullanıcıların beklentilerinin artması veri madenciliğinin gelişmesine yol açmıştır.

Veri madenciliği: fazla miktardaki veri içinden yararlı bilgiyi bulmak için kullanılır.

- * Bilgi keşfi ise büyük veri kaynaklarından yararlı ve ilginç bilgiyi bulmak için kullanılır. Bulunan bilgi; gribi, yararlı ve önceden bilinmeyen olabilir.

- * Tanımlama ⇒ bilgi keşfi > veri madenciliği
- Pratikte ⇒ bilgi keşfi = veri madenciliği



Bilgi Keşfinin Aşamaları

- Uygulama alanını inceleme
- Amaca uygun veri kümesi oluşturma
- Veri ayıklama ve önizleme
- Veri atılma ve veri dönüşümü
- Veri madenciliği tekniği seçme
- Veri madenciliği algoritmasını seçme
- Model değerlendirme ve bilgi sunumu
- Bulunan bilginin yorumlanması.

Veri madenciliğinde amacı büyük miktardaki veri içerisinde arama yapmak değildir bunu VTYS 'lerinde yapıyor. Amacı; ara-
dığımız veri mevcut ise sonuçlarını anlamaktır.

* Veri tabanı ve veri madenciliği farkı *

Veri Tabanı → Tanımlı bir sorgulama ve SQL dili kullanılır.
→ Canlı veri kullanılır.
→ Belirli bir çıkış elde edilir ve bu verinin bir alt kümesidir.

Veri Madenciliği → Tam tanımlı olmayan sorgulama ve yaygın dil yok.
→ Üzerinde işlem yapılmayan veri kullanılır.
→ Çıkış belirli değildir ve verinin alt-

* Veri madenciliğinin uygulama alanları ise ;

- * Vt analizi ve karar verme desteği
- * Risk analizi 2o Pazar Araştırması
- 3o Sağlıklaırlıkların saptanması gibi uygulamalardır.

1.- Risk Analizi → Finans planlaması ve bilanço değerlendirmesi
→ Kaynak planlaması
→ Rekabet (rakipleri ve pazar eğilimlerini inceleme)

2.- Pazar Araştırması → Uygulamalar için veri kaynağı (anket kartı, kupon)

- Hedef pazarlar bulma
- Müşterilerin davranışlarında zaman içindeki değişiklik.
- Çapraz pazar incelemesi
- Müşteri profili
- Müşterilerin ihtiyaçlarını belirleme

3.- Sağlıklaırlıkların Saptanması → Sigorta, bankacılık ve tele-

komünikasyon alanlarında geçmiş veri kullanılarak sağlıklaırlık raporları için bir model oluşturma ve benzer davranış gösterenleri belirleme.

BÖNEK → Araba Sigortası
Sağlık Sigortası
Kredi Kartı Başvurusu.

*Veri Temizleme : Gerçek uygulamalarda veri temiz, güncel ve tutarlı olabilir ancak bundan kurtulmak..

-Eksik veri : Genellikle eksik nitelik değerleri olan silinir

Ama elle de doldurulabilir ve değer ort. ile de doldurulur.

-Günlüklük Veri : Oluşan bir değerdeki hata gibi. Günlüklük

gök olarak sonucu etkiler. Bu yöntemle yok eder;

①-Bölmeleme : Veri sıralanır buna göre bölümlere ayrılır

②-Kümeleme : Benzer veriler aynı derinlikte gruplanır, ayrılar silinir.

③-Fgü Uyuma : Veri bir fonksiyona uydurulur. Doğrusal

fgü uydurmada bir değişkenin değeri diğer bir değişken kullanılarak bulunabilir.

*Farklı kaynaklardan olan veriler tutarlı bir şekilde birleştirilmeli ve gereksiz veriler silinmeli.

*Veri, veri madenciliği uygulamaları için uygun olmayabilir. Belirli yöntemlerle uygun hale getirilmeli.

*Veri miktarı fazla olduğunda algoritmalar çok uzun sürebilir. Veri azaltmak gerektirir.

→Nitelik birleştirme işlemi : Sorgulama için gerekli olan bağutlar kullanılıyor, Bu nitelik seçme işlemi ise tüm kümenin bir alt kümesi ile işlem yapılmakta olur. d bağutlu küme k < d olacak şekilde k bağutuna taşınır.

2) Veri Ambarları → Farklı verilerin sorgulandığı ve analizlerinin yapılabilirdiği bir depodur. Ut'unu yormamak için oluşturulmuştur. En

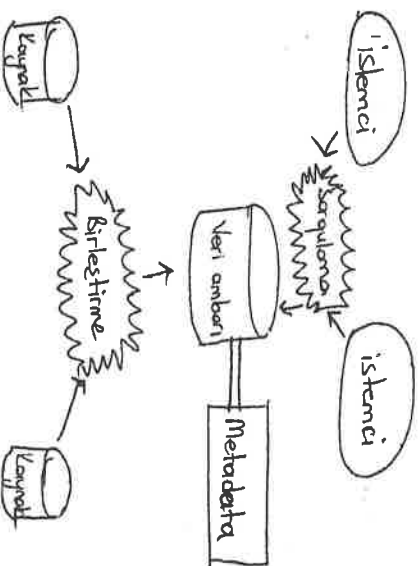
veri ambarı ilgili veriyi kolay, hızlı ve doğru birimde analiz etmek için gerekli işlemleri yerine getirir.

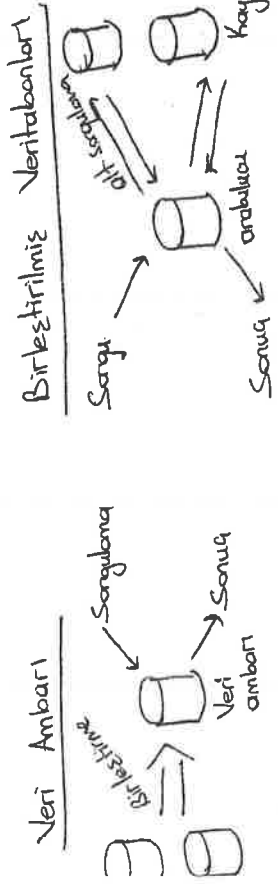
→amaca yönelik : Kayıtlar ilgili karar vermek için gerekli olmayan veriyi kullanılmayarak kayıtlar basit, bilet baltı sağlar.

→birleştirilmiş : Veri kaynaklarının birleştirilmesiyle oluşur. rılır. Dışgılar gibi.

→tamam değışkenli : Tamam değışkeni canlı veritabanlarına göre daha uzundur.

→Değışken değil : Canlı veritabanlarından alınmış verilerin fiziksel olarak başka bir ortamda saklanmasıdır.





Veri Madenciliğinde Sorunlar

* Girdi ile ilgili sorunlar → Veriye ait verilerin toplanarak,

işlerden izinsiz olarak kullanılması.

* Kullanıcı arabelleği → Görüntüleme ve etkileşim

* Veri madenciliği yöntemi → Farklı tipte veriler üzerinde, görüntü-
de, değişimli biçimde çalışma zorunda
kalma.

* Başarım ve ölçeklenebilirlik →

* Veri madenciliği yöntemleri bilimsel kullanımları. Günlük;
bu yöntemler geçmiş olaylara bakarak sonuçları bulur,
gelecekteki verilerin garantisini vermez.

VERİ

Veri → Nesneler ve nesnelerin niteliklerinden oluşan kümedir.
* Kayıt, varlık ve örnek de denilebilir.

Nitelik → Nesnenin bir özelliğidir.

Değer kümesi → Nitelik için saptanmış sayılar veya sembollerdir.

* Nitelik türleri → Belirli aralıkta yer alan değişkenler

→ İktisadi değişkenler

→ Açık ve sıralı değişkenler

* Gerçek uygulamalarda toplanan veri kırılgan; eksik,
günlük ($max = -10 + 1$) ve tutarsız olabilir. Bunların sebepleri;
kurallara uyulmaması, o anda anlık hatalar, kaynak hataları
olabilir. Bu şekilde güvenilir veri ortaya çıkar.

* Veri kalitesi ise veri madenciliği uygulamaları ile yararlı
bilgi bulma sansı daha fazla. Bu yüzden veri ön işleme ile veri

kalitesi hale getirilmelidir.

Veri Ön işleme

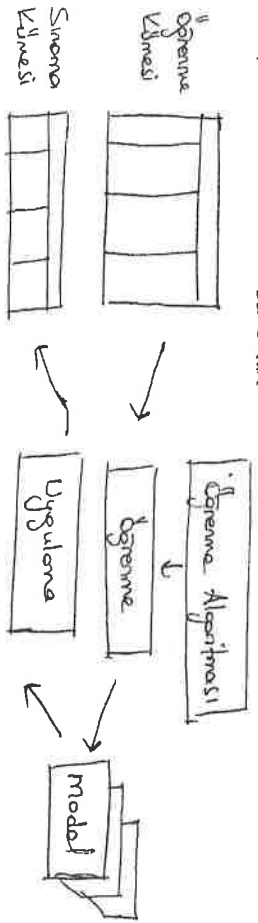
Veri temizleme → Eksik nitelik değerleri tamamlama gibi.

Veri birleştirme → Farklı veri kaynaklarındaki verileri birleştirme

Veri dönüşümü → Normalizasyon ve birleştirme.

Veri azaltma → Aynı veri madenciliği sonuçları elde ederek
şekilde veri miktarını azaltma.

3-Modeli kullanma: Model, daha ömeden görünüş örnekleri sınıflandırmak için kullanılır.



Sınıflandırmaya Başlatılmıyın DEĞERLENDİRME

- * Doğru sınıflandırmaya başlatması
- * Hız
- * Kararlı olması
- * Anlaşılabilir olması
- * Ölçülebilirlik
- * Kurallara yapısı

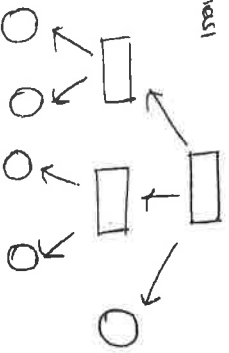
Sınıflandırmaya KONTROLÜ

- * Karar ağaçları
- * Yayıp sınırları
- * Bayes sınıflandırıcılar
- * İstatistik tabanlı sınıflandırıcılar
- * k-en yakın komşu yöntemi
- * Destek vektör makineleri

Karar Ağaçları

Aktif diyagramı şeklinde ağaç yapısı

- * Ara düğümler bir nitelik sınaması
- * Dallar sınama sonucu
- * Yapraklar sınıflar



Benzerlik ve Farklılık

* Nesneler birbirine daha benzer ise daha büyük değer alınır. Bu değer genelde 0-1 aralığında olur.

Öklid Uzaklığı: Nesneler arası farklılığı bulmak için kullanılır.

$$d(i,j) = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2}$$

Manhattan Uzaklığı: Öklid genelleştirilmiş hali.

$$d(i,j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

q → 1 olursa → Manhattan Uzaklığı

q → 2 olursa → Öklid Uzaklığı

*

	Nesne i	Nesne j
Nesne i	0	1
	M ₀₀	M ₀₁
	M ₁₀	M ₁₁

M₀₀: i=0, j=0 olduğu nitelik sayısı.

* Yalın Düzlem Katsayısı → $sim(i,j) = \frac{M_{ii} + M_{00}}{M_{00} + M_{01} + M_{10} + M_{11}}$

* Jaccard Katsayısı → $d(i,j) = \frac{M_{ii}}{M_{01} + M_{10} + M_{11}}$

Ağac oluşturmada birine temel yaklaşımlar

* Bilme kriteri

* Durma kararı

* Bölme kriteri

* Etkileşime kurallı

Karar Ağacı kullanılarak Sınıflandırma

* Doğrudan

* Dolaylı

Ayarlar:

* Karar ağacı oluşturma yöntemi

* Kullanılan kriter

* Sınıflı ve ayrık nit. için kullanılabilir.

Detaylar:

* Sınıflı nit. tahmin etmekte çok başarılıdır.

* Sınıf sayısı fazla ve örneklere birisi ot aldığında başarısız.

Karar Ağaçlarında Aşırı Öğrenme

Örneklere birimindeki örneklerin ağırlığı veya gücüyle ölçülmesi.

Bunu engelleyen 2 yaklaşım:

- İşlemi erken sona erdirmeye.

- Ağac oluşturduktan sonra ağacı kısaltmaya.

Karar ağacı yöntemleri

Genel olarak iki aşamadan oluşur.

1 - ağac oluşturma

2 - ağac budama

Karar ağacı oluşturma

□ * ağac bütün verinin oluşturduğu tek bir ağaçla başlıyor.

① * eğer örnekleri hepsi aynı sınıfta olursa ağacın yaprak olarak sonlanıyor ve sınıf etkisini alıyor.

□ Nit * eğer değil ise örnekleri sınıflara en iyi bilecek olan nitelik seçiliyor.

* İşlem sona eriyor.

∇

Anahtar tipi

a → 10 } örnek var
b → 10 }

o alt üst spes

foratı en yüksek olan kalitesi
yüksek bölünür.

* Ama en iyi bölün nitelik iyi bir fonksiyona ile belirlenir.

1) Bilgi Kazancı: Bütün niteliklerin ayrı değerler aldığı varsayılıyor.

2) Gini İndeksi: Her nitelik için olası bütün ikiye bölünmeler

sınanır.

① Bilgi Kazancı

* Entropi → Rastgelecilik, belirsizlik, bellemeyen durumu ortaya çıkma olasılığını gösterir.

$$H(p_1, p_2, \dots, p_n) = - \sum_{i=1}^n p_i \cdot \log(p_i)$$

→ Sınıfların içindeki olayların olasılığı bir durum olduğuna için entropi = 0 alınmalı.

- Örnekler aynı sınıfa aitse entropi = 0
- Örnekler sınıflar arası eşit dağılımısa entropi = 1
- Örnekler sınıflar arası rastgele dağılımısa $0 < \text{entropi} < 1$ olur.

Örnek: S veri kümesinde 14 örnek : a sınıfında 9 b sınıfında 5 } örnek varsa.

$$\text{Entropi} = -\left(\frac{9}{14}\right) \log_2\left(\frac{9}{14}\right) - \left(\frac{5}{14}\right) \log_2\left(\frac{5}{14}\right) = 0.940$$

BİLGİ KAZANCI:

A niteğinin S veri kümesindeki bilgi kazancı;

$$\text{Bilgi Kazancı}(S, A) = \text{Entropi}(S) - \sum \text{values}(A) \frac{|S_A|}{|S|} \cdot \text{Entropi}(S_A)$$

Örnek: Önce tüm niteğilerin bilgi kazancıları aynı ayı hesaplanır. Bilgi kazancı büyük olan kök değeri seçilir. Aynı işlemler yapıldıkça ulaşılan kadar yapılır.

$$\text{Meseler ; Entropi}(S) = 0.940$$

$$\text{Entropi}(\text{Sweet}) = 0.811 \quad \text{Entropi}(\text{Strong}) = 1.00$$

$$\text{Gain}(\text{wind}) = 0.940 - \left(\frac{5}{14}\right) \cdot 0.811 - \left(\frac{9}{14}\right) \cdot 1.00 = 0.046 \text{ olur.}$$

(weak: 5) (strong: 9)

② GINI INDEX

* Veri kümesi S içinde n sınıf varsa ve p_j c_j sınıfının olasılığı

ise;

$$gini(S) = 1 - \sum_{j=1}^n p_j^2$$

* Eğer veri kümesi S₁ ve S₂ alt kümelerine bölünüyorsa ve her alt kümede srasıyla N₁ ve N₂ örnek varsa;

$$gini_{split}(S) = \frac{N_1}{N} \cdot gini(S_1) + \frac{N_2}{N} \cdot gini(S_2)$$

* Gini index değeri en küçük olan seçilir

Örnek:

C1	0
C2	6

$$P(C1) = 0/6 = 0 \quad P(C2) = 6/6 = 1$$

$$Gini = 1 - P(C1)^2 - P(C2)^2 = 1 - 0 - 1 = 0$$

C1	1
C2	5

$$P(C1) = 1/6 \quad P(C2) = 5/6$$

$$Gini = 1 - \frac{1}{36} - \frac{25}{36} = 1 - \frac{13}{18} = \frac{5}{18} = 0.277$$

BÖLMELEME

Isi	40	40	72	80	80
Tenis Oynamak	No	No	Yes	Yes	No

* Sürekli nitelik A sınıfları. Birbirini izleyen ancak sınıf etkisi farklı olan nitelik değerleri bulunur. En fazla kazancı sağlayan bölme seçilir.

* Statik de ise en fazla bölme yapılır. Bölme esit genişlik, esit derinlik veya denetleme yöntemi ile bulunur.

FARKLI SINIFLANDIRMA YÖNTEMLERİ

(5)

1) Örnek tabanlı yöntemler : Örnekleme kümesi saklanır, sınıflandırılacak yeni bir örnek geldiğinde öğrenme kümesi sınıf etiketini öngörmek için kullanılır.

* k -en yakın komşu yöntemi : Sınıflandırılmak istenen örneğe en yakın örnekleri bul.

Δ Bütün örnekler n -boyutlu uzayda bir noktaya karşı düşürülür. Nesneler arasındaki uzaklık (ölçülür). hesaplanır. Öğrenilen fonksiyon aynı değeri veya genel değeri verir.

* Aynı dağılım fonksiyonlarında;
 k -komşu algoritması K_q örneklere en yakın k öğrenme örneğinde en çok görülen sınıf değeri verir.

* Sınıklı değeri fonksiyonlarında;
en yakın k öğrenme örneğin ortalaması alınır.

2) Genetik Algoritmalar : Optimizasyon amaçlıdır bir başlangıç çözümü öneriyor, tetkiklenen her ana adımda daha iyi çözüm üretilmeye çalışıyor. 5 ana parçadan oluşuyor.

- bireylerden oluşan bir başlangıç kümesi, P
- Çarpılma : Bir ana tabandan yeni bireyler üretmek için yapılan işlemler
- Mutasyon : Bir bireyi rastgele değiştirme.
- Uygunluk : En iyi bireyleri belirleme.
- Çarpılma ve mutasyon tekniklerini uygulayan ve uygunluk fonk. değeri tablosu içindeki en iyi bireyi seçen algoritma.

Ayarlar → Paralel çalışabilir.

NP karmaşık problem çözümlerine uygun

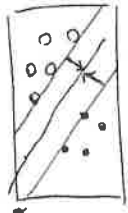
Dezavantajlar → Son kullanıcının modeli anlaması zor.

Problemi GA ile çözmeye uygun hale getirmek zor.

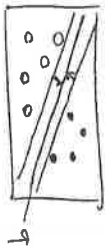
Uygunluk fonksiyonunu belirlemek zor.

3) Destek Vektör Makineleri: Hem doğrusal olarak ayırt edilebilen hem de edilemeyen veri kümesini sınıflandırabilir.

Verileri; a → Bir çözüm önerisi,



→ Başka bir çözüm önerisi.



Hangisi daha iyi?

a, b'den daha iyi...

4) Bulanık Küme Sınıflandırıcılar:

Bulanık mantık 0.0 ve 1.0 arasında genel değerler kullanılarak üyelik dereceleri hesaplar. Nitelik değerleri bulanık değerlere dönüştürülür. Kurallar kümesi oluşturulur.

Yeni bir örneği sınıflandırmak için birden fazla kural kullanılır. Her kuraldan gelen sonuç toplanır.

5) Öğreni:

Sınıflandırmaya problemleriyle aynı yaklaşım... Model oluştur, bililmeyen leğeri hesaplamak için modeli kullan. Ama;

Sınıflandırmaya → Ayrıntı değeri

Öğreni → Silineli değeri.

Eğri Uydurma

Doğrusal eğri uydurma → En basit eğri uydurma yönteminin Veri doğrusal bir eğri ile modellenir.

$$\hat{y} = w_0 + w_1 \cdot a_1 + w_2 \cdot a_2 + \dots + w_k \cdot a_k$$

İLİŞKİLENDİRME KURALLARI

* Veri kimesi içindeki yaygın düğümlerin, nesneleri oluşturan ni-tekliler arasındaki ilişkilerin bulunması için vardır.

* Kural Zetli : "Görde \rightarrow Baş[destek, güven]"

Satır alma ($x_i, etme$) \rightarrow Satır alma ($x_i, "sif"$) [0,0,0,0, %65]

Yaygın Nitelikler

Nitelikler kimesi \rightarrow 1 veya daha fazla nitelikten oluşan kime.

Destek Değeri \rightarrow Bir nitelikler kimesinin veri kimesinde görülme

sıklığı \rightarrow 5 kisten 3'ü kula alıyorsa $d.d=3$

Güven Değeri \rightarrow Destek Değeri \rightarrow $g.d = \frac{3}{5}$

Kendisi

* Amaç; D hareket kimesinden kurallar çıkarmak;

Kurallardan d.d, belirleyen en küçük d.d'nden büyük olmalı. (minup)

Genellik g.d, " " " g.d, " " (minup)

esitlikle kabulünüz \rightarrow bu eşit değerleri dışardan verilen

Route-force Yaklaşım

Okası bütün kuralları listele

Her kural için destek ve güven değeri hesapla.

minup ve minof eşit değerlerinden büyükleri sil.

Hesaplamaya maliyeti yüksek.

$$D = 3^d - 2^{d+1} + 1$$

d = Nitelik sayısı

İstiklamlendirme Kural Oluşturma: 2 aşamadan oluşur

1- Yayıgın nitelikli betimlere, (minsup \neq a olmalı.) $\rightarrow abc$

2- Kural oluşturma. $\rightarrow ab \rightarrow c, b \rightarrow ac, \dots$ hepsi yapılır.

d nitelik için 2^{de} yaygın nitelikli oluşturulabilir.



TID	Öğeler
1	
2	
3	
4	
5	

TID	Öğeler
1	
2	
3	
4	
5	

Yayıgın nitelikli adayları.

* Yayıgın nitelikli oluşturma yöntemleri *

\rightarrow Yayıgın nitelikli aday sayısını (m) azaltma (Apriori Algoritması).

\rightarrow Hareket sayısını (n) azaltma.

\rightarrow Veri kümesi taranma sayısını azaltma (mv)

APRIORI ALGORİTMA

Temel Yıkılma: $\forall x, y: (x \subseteq y) \Rightarrow s(x) \geq s(y)$

* Bir nitelikler kümesinin destek değeri alt kümesinin destek değeri den büyük demektir.

* Bir yaygın nitelikler kümesinin alt kümesinde yaygın nitelikler kümesidir.

ÖRNEK

TID	Öğeler
1	Ekmek, Süt
2	Ekmek, Bez, Bira, Yumurta
3	Süt, Bez, Bira, Kola
4	Ekmek, Süt, Bez, Bira
5	Ekmek, Süt, Bez, Kola
6	Şakaa

Görüm

minsup = 3 \Rightarrow bu eşik değeri dışından verilir. altındaki öğeler silinir.

Şimdi tüm 2'li kombinasyonları yazarız.

C1 \rightarrow	Öğeler	Sayı
	Ekmek	4
	Süt	4
	Bez	4
	Bira	3
	Yumurta	1
	Kola	2

C2 \rightarrow	Öğeler	Sayı
	(ekmek, süt)	3
	(ekmek, bez)	3
	(ekmek, bira)	2
	(süt, bez)	3
	(süt, bira)	2
	(bez, bira)	3
	"	"
	"	"
	"	"
	"	"

C3 \rightarrow	Öğeler	Sayı
	(ekmek, süt, bez)	2
	(ekmek, süt, bira)	"
	(süt, bez, bira)	"
	"	"

Destek

$$C_1 + C_2 + C_3 = 6 + 15 + 20 = 41$$

Yayıgın sayı azalınca

$$C_1 + C_2 + C_3 = 6 + 6 + 1 = 13$$

DEMETLEME = KÜMELEME

* Nesneleri gruplara ayırma işlemidir. Aynı gruptaki nesneler birbirine daha çok benzer, farklı gruptaki nesneler birbirine daha az benzerdir.

* Eğitim verisi kullanılmaz, sınıflandırmadan farklı budur.

* Sınıflandırmada başta bilgi varken kümelere de yoktur.

* Gözetimsiz öğrenmeye girer, yani hangi nesnenin hangi

sınıfa ait olduğu ve sınıf sayısı belli değil.

* Bu metod; görselliği işleme, ekonomi ve ayrıntınlıkları

belirleme gibi uygulamalarda kullanılır.

İyi bir demetleme yöntemi

Aynı demet içindeki nesneler arası fazla benzerlik olmalı.

Farklı demet içindeki nesneler arası az benzerlik olmalı.

Veri içindeki gürültüleri düşürmek mümkün.

Gruplama için uygun demetleme kriteri bulunmalı.

Temel Demetlerne Yaklaşımları

'Bölünmeli' yöntemler → Veriyi bölerek, her grubu belirlenmis

ir kriteri göre değerlendirir. (K-means)

- Hiyerarşik yöntemler → Veri kümelerini "önceden belirlenmiş"

15 kritere göre hiyerarşik olarak ayırır.

2 Yogunluk tabanlı yöntemler → Nesnelerin yoğunluğuna göre

emelleri olustur.

Model tabanlı yöntemler → Her modelin bir modeli

İyidegi varsayılır. Ama bu modellere uygun verileri gruplamak.

K-Means Demetlene

Bilinen bir k değeri için k -means denetleme algoritması -

in 4 asamas varir:

1-Veri kümesi K alt kümeye ayrılır.

2- Her demetin ortalaması hesaplanır : merkezi nokta:

3- Her nesne en yakın market noktasının olduğu denetle dahil
dilir.

4- Nesnelerin demetlenmesinde değişiklik olmayana kadar

dim 2'ye geri döndür.

1421551

* $\text{Karmasiklig} \rightarrow \text{yer karmasiklig} (o(n+k).d)$

Zaman " " (o.k.t.n.d) (p.v.t.n.d)

K-means

- Gerçekleşmesi kolay

- Karmaşıklık diğerlerine göre daha az

- K-means bazı durumlarda iyi sonuç veremeyebilir.

→ veri grupları farklı boyutlarda ise

→ Veri gruplarının yoğunlukları farklı ise

→ Veri gruplarının şekli küresel değilse

→ Ven tände
öytriliktat varsa en iji
sarsa vermayabilir.

- baslangici sartlari cok baglidir. sartlar farklı durur

senavlar da farklı olur.

- Nesnelerin merteye uzaklıkları ne kadar yakın ise cözüm
- kadar iyidir.

- K-means'in kötü yanı tüm verilerin illaki bir kümeye dahil edilme zorudur.

• Eğer veri sayısı kadar merkezi olsun her bir oranı en az olur ama bu istenen bir şey değildir.

- Eğer tek küme ederse hatasız olur en fazla olur.

Hata Oranı \rightarrow Noktaların merkezlerine uzaklıklarını $\frac{||d_k||}{||d_k||}$

Sonra hepsini toplar bulduğın sonucu hatla orandır.

*Teoreme göre teter teter tüm sınıflara bakılır hangi sınıf değeri yüksek ise o sınıfa yerleştirir.

$$P(v_1, v_2, \dots, v_n | S_j) \rightarrow \text{Sınıf nitelikler}$$

ÖRNEK :

v_1	v_2	v_3	v_4	Sınıf
Yes	No	No	Yes	B
Yes	No	Yes	No	B
No	Yes	Yes	No	M
No	No	Yes	Yes	M
Yes	No	No	Yes	B
Yes	No	No	No	M
Yes	Yes	Yes	No	M
Yes	Yes	No	Yes	M
No	No	No	Yes	B
No	No	Yes	No	M

Yandaki veri kümesi için sırasıyla gelen $\langle \text{Yes, No, Yes, Yes} \rangle$ örneğin sınıfını bulalım.

$$\text{Formül: } S_{\text{değer}} = \arg \max \left[\prod_{i=1}^n P(v_i | S_j) \right]$$

$$P(S | \text{Yes, No, Yes, Yes}) = ?$$

→ tüm sınıfları teter teter dene

$$P(S=B | v_1, v_2, v_3, v_4) = \frac{1}{10} \cdot \left[\frac{3}{4} \cdot \frac{1}{4} \cdot \frac{3}{4} \right] = \frac{9}{160}$$

$$P(S=M | v_1, v_2, v_3, v_4) = \frac{6}{10} \cdot \left[\frac{3}{6} \cdot \frac{3}{6} \cdot \frac{4}{6} \cdot \frac{2}{6} \right] = \frac{1}{20}$$

$$\frac{9}{160} > \frac{1}{20} \text{ olduğu için bu veri B sınıfında dendir.}$$

* Eğer nitelikler sürekli değerler dursa yani evet-hayır değil de 1-2-3... gibi değerler olursa gauss dağılımını kullanırız.

ÖRNEK :

v_1	v_2	v_3	v_4	Yas	Cinsiyet
Evet	Hayır	Evet	Hayır	38	K
Evet	Evnet	Evnet	Hayır	40	K
Evnet	Evnet	Evnet	Hayır	41	K
Evnet	Evnet	Evnet	Hayır	55	K
Hayır	Hayır	Hayır	Hayır	27	E
Hayır	Evnet	Hayır	Hayır	30	E
Evnet	Hayır	Evnet	Evnet	35	E
Hayır	Hayır	Hayır	Hayır	42	E
Evnet	Hayır	Hayır	Hayır	43	E
Evnet	Hayır	Hayır	Hayır	45	E

Yandaki veri kümesi için $\langle \text{Evnet, Evnet, Hayır, Hayır, 45} \rangle$ Cinsiyet = ?

$$\text{Normal Dağılımı göre } \Rightarrow P(v) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{1}{2} \left(\frac{v-\mu}{\sigma} \right)^2}$$

σ: standart sapma.

$$\mu = \frac{x_1 + x_2 + \dots + x_{10}}{10} = 39,6$$

$$\sigma = \sqrt{\frac{1}{10} \cdot [(28-39,6)^2 + (40-39,6)^2 + \dots + (45-39,6)^2]}$$

$$P(\text{Yas} | \text{Cinsiyet} = K) = \left\{ \begin{array}{l} \text{aynı ayın hesapların Bulduğumuz değerleri;} \\ P(\text{Yas} | \text{Cinsiyet} = E) = \end{array} \right\}$$

buraya yerleştir.

$$P(C=E | v_1, v_2, v_3, v_4, \text{Yas}) = \frac{2}{3} \left(\frac{2}{3} \cdot \frac{1}{3} \cdot \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{3}{100} \right) = 0,0016$$

* 0 problemi

Eğer sonuçlardan bir 0 dursa çarpım durumunda olduğu için hiç bir veri alamayız. Tüm olasılıkları 1 ile toplarsak;

$$V_1 = \frac{2}{5} \quad V_2 = \frac{0}{5} \quad V_3 = \frac{2}{5} \quad \text{olursa;}$$

$$V_1 = \frac{4}{6} \quad V_2 = \frac{1}{6} \quad V_3 = \frac{1}{6} \quad \text{der işlem yaparız}$$

BAYES TEOREMİ

Birbirinden bağımsız ve rastgele iki olayın birbiri ardı sıra gerçekleştiği durumlarda bu iki olayın birinin gerçekleşmesi durumunda ikinci olayın gerçekleşmesi olasılığı $P(A|B)$, $P(B|A)$ yada $P(A \cap B)$ ifadesiyle gösterilebilir.

$$P(A \cap B) = \frac{P(A|B) \cdot P(B)}{P(A)} = \frac{P(B|A) \cdot P(A)}{P(B)}$$

* Bayes teoremi rastgele bir A olayı ile diğer bir rastgele B olayı için koşullu olasılıklar ve marginal olasılıklar arasındaki ilişkiyi tanımlar.

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

* Naive Bayes Sınıflayıcı ile bayes teoremi hesaplarında dikkat edilmesi gereken en önemli fark sınıflayıcıların olasılık değerinden ziyade hedef sınıfı bulmaya odaklanmasıdır. Bu yüzden paydada bulunan değer, tüm hedef sınıflara ait olasılık hesaplarında ortak olduğunda ihmal edilebilir.

K-NN ÖRNEK

X1	X2	Y
2	4	bad
3	6	good
3	4	good
4	10	bad
5	8	bad
6	3	good
7	9	good
9	7	bad
11	7	bad
10	2	bad

X1=3 ve X2=4 değerleri için k=4 ile

hangı sınıfı ait olduğunu bulalım... k=4

$$d(i,j) = \sqrt{(2-3)^2 + (4-4)^2} = 1$$

$$d(i,j) = \sqrt{(3-3)^2 + (6-4)^2} = 2$$

$$d(i,j) = \sqrt{(3-3)^2 + (4-4)^2} = 0$$

Hesaplamaları yaptıktan sonra en yakın 4 noktayı alıp

tablodan karşılıklarına bakmaktır. Good, Bad, Bad, Bad ver diye

BAD sınıfına dahil ederiz...

Algoritma Adımları

* K değeri dışardan isteriz.

* Sınıf bilinmeyen noktaya sınıfı veri noktalarının uzaklık-

ları hesaplanır. Genelde 3'lük uzaklığı hesaplanır.

* Hesaplanan uzaklıklar sıralanır. En küçük uzaklığa sahip

K adet sınıfı veri noktası tespit edilir.

* Seçilen K adet sınıfı veri noktası içerisinde çoğunluğu

güç sahipli sınıf belirlenir. Bu sınıf bilinmeyen noktaya için sonuçtır.

↓ ↑

(K-NN)



a) derinlik
genişlik



b) derinlik genişlik

$$\min = 2$$

$$\max = 20$$

$$20 - 2 = 18$$

$$\frac{18}{2} = 9$$

$$\begin{pmatrix} 2 & -11 \\ 11 & -20 \end{pmatrix} \rightarrow$$

VERI MADENCILIGI

BÖNÜT

Kayıt No	Yaş	Medeni Durum	Araba Sayısı
1	23	Bekar	1
2	25	Evlü	1
3	29	Bekar	0
4	34	Evlü	2
5	38	Evlü	2
...			

Tablodan İstatistiklere
Kuralları Derleme -

① * Bekar bakmak; * yaş = 30-39, medeni durumu: evli \Rightarrow a. sayısı: 2
dört

② * Yaş: 29, Medeni durum: Bekar \Rightarrow Araba sayısı: 0

③ * Yaş: 30-39 \Rightarrow Medeni Durum: Evli

④ * arabası: 0-1 \Rightarrow medeni durumu: bekar

$$\textcircled{1} \frac{(Yaş: 30-39, m.d: evli, a.s=2)}{(Yaş: 30-39, m.d: evli)} = \frac{2}{2} = \%100 \rightarrow \text{güven değeri}$$

$$\frac{(Yaş: 30-39, m.d: evli, a.s=2)}{5} = \frac{2}{5} = \%40 \rightarrow \text{destek değeri}$$

$$\textcircled{2} \frac{(Yaş: 29, m.d: bekar, a.s: 0)}{(Yaş: 29, m.d: bekar)} = \frac{1}{1} = \%100 \rightarrow \text{güven değeri}$$

$$\frac{(Yaş: 29, m.d: bekar, a.s: 0)}{5} = \frac{1}{5} = \%20 \rightarrow \text{destek değeri}$$

;

nicel deęer olan yaşı b l mleyelim; 20-29, 30-39 olsun;

Kayıt No	Yaş	M. durum	A. sayısı
1	23	Bekar	1
5	38	Evl�	2

K. No	Yaş 20-29	Yaş 30-39	m.d. evli	m.d. bekar	Araaba	A ₁	A ₂
1	1	0	0	1	0	1	0
5	0	1	1	0	0	0	1

Yaş. 20-29, m.d. bekar $\Rightarrow A_1$

Yaş 30-39, m.d. evli $\Rightarrow A_2$

  NEK : Faturaların tutulduęu veri tabanı tablosundaki kayıt-

lar ařaędaki sekildedir;

TID	���	Frekans
1	{1,2,3,4}	3
2	{1,2}	6
3	{2,3,4}	4
4	{2,3}	5
5	{1,2,4}	
6	{3,4}	
7	{2,4}	

\Rightarrow

\Rightarrow

���	Frekans
1,2	3
1,3	1
1,4	1
2,3	3
2,4	4
3,4	3

\Leftarrow

mins p deęeri 3 olduęu

i in 3 l  deęerleri

bir deęer alınamaktadır.

$1 \Rightarrow 2$; d.d $\rightarrow 3$

$g.d \Rightarrow 3/3 = 1$

$2 \Rightarrow 1$; d.d $\Rightarrow 3$

$g.d \Rightarrow 3/6 = 0.5$

D	F
2,3,4	2

Kümelere (Denetleme)

* eğitim verisi çok
tüm verileri dışardan
veren k sayısına göre
kümelere, k -means

algoritması kullanılarak.

* Uzaklığa dayalı kümelere
* k -means avantajları

* Rasgele noktalar ile başlan-

dıgından elde edilen her sonuç

kararı değildir,

* Küme sayısı broad kümelere

algoritmasında olduğu gibi dışardan
verilir.

* Konkar olmayan küme setlerinde
başarılı değildir,

* Keskin kümelere yapar.

* Aykırı değerlere karşı hassastır,

* Başlangıç noktasına bağlıdır.

Sınıflandırma

* eğitim verisi kullanılarak
test verisi uygun yere yer-
leştirilir k -nn algoritması
kullanılarak.

* k -NN avantajları

* Bellek tabanlı bir sınıflayıcı-

dır. Sınıflı veri noktaları sürekli
hafızada tutulmalıdır. Veri

kümesi çok büyük olduğunda

hesaplama süresi de artacaktır.

* Hesaplama için özellikler
kaldırılabilir için gereksiz veya

ilgisiz özellikler sınıflandırıcı
mayı kötü yönde etkileye-
bilir.

* Sınıflandırma başarısı

ocısından genellikle yapay

sınırlı ağlar gibi gelişmiş

sınıflandırma tekniklerinin

gerisinde kalır.

Uzayın belli bir kısmı umman sınıftır.))

GINI İLE KARAR VERME GİRİ

k -nn

k -NN

k -means
eğitim verisi kullanılmaz.

* Sınıflandırmada başka bilgi vardır kümelere ise bilgi yoktur.

* Kümelere olduğu için eğitim verisi kullanılmaz, sınıflandırmada farkı vardır.

disordan verilir.

* Kümelere yöntemidir.

k -MEANS



- $1.2 \Rightarrow 3 \Rightarrow \frac{1.2}{1.2, 3} \Rightarrow \frac{1.2}{2} \Rightarrow \%50 < \%70$
- $1.3 \Rightarrow 2 \Rightarrow \frac{1.3}{1.2, 3} \Rightarrow \frac{1.3}{4} \Rightarrow \%50 < \%70$
- $2.3 \Rightarrow 1 \Rightarrow \frac{2.3}{1.2, 3} \Rightarrow \frac{2.3}{4} \Rightarrow \%50 < \%70$

- $1.2 \Rightarrow 5 \Rightarrow \frac{1.2}{1.2, 5} \Rightarrow \frac{1.2}{4} \Rightarrow \%50 < \%70$
- + $1.5 \Rightarrow 2 \Rightarrow \frac{1.5}{1.5, 2} \Rightarrow \frac{1.5}{2} \Rightarrow \%100 > \%70$
- + $2.5 \Rightarrow 1 \Rightarrow \frac{2.5}{1.5, 2} \Rightarrow \frac{2.5}{2} \Rightarrow \%100 > \%70$

C_3

Itemset	Sup
$(1.2, 3)$	2
$(1.2, 5)$	2

L_2

Itemset	Sup
(1.2)	4
(1.3)	4
(1.5)	4
(2.3)	4
(2.4)	2
(2.5)	2

C_2

Itemset	Sup
(1.2)	4
(1.3)	4
(1.4)	1
(1.5)	2
(2.4)	2
(2.5)	2
(3.4)	0
(3.5)	1
(4.5)	0

Items

$1.2, 5$
2.4
2.3
$1.2, 4$
1.3
2.3
$1.2, 3, 5$
$1.2, 3$

$min_{sup} = 2$
 $min_{conf} = \%70$

C_1

Itemset	Sup
1	6
2	7
3	3
4	4
5	5
Sold	2

"ÖRNEK 1"

Gül, nur, ilk.	ilk, con	nur, ilk	gül, ilk	gül, ilk, nur, con
----------------	----------	----------	----------	--------------------

⇒ Tamamı nitelik kümesi
 gülün destek değeri ⇒ 3
 gülün güven değeri ⇒ 3/5

ilk	gül	3
2	nur	3
3	ilk	5
4	con	2

ilk	gül	2	(gül, ilk)	3
1	(gül/nur)	2	(nur, ilk)	3

minsup = 3 dediyimz anda "con" silinir
 3'le kombinasyonlarını yazarız.

minsup = 3'ü di
 (gül/nur) yitirir

"ÖRNEK 2"

10	A, C, D
20	B, C, E
30	A, B, C, E
40	B, E

Supmin = 2

ilk	A	2
ilk	B	3
ilk	C	3
ilk	D	1
ilk	E	3

ilk	A, B	1
ilk	A, C	2
ilk	A, E	1
ilk	B, C	2
ilk	B, E	3
ilk	C, E	2

(2)

(1,2,5) için;

$$\cancel{(1,2)} \Rightarrow 5 \rightarrow \frac{(1,2,5)}{(1,2)} \Rightarrow \frac{2}{4} = \%50 < \%70 \text{ (elmez)}$$

$$\checkmark (1,5) \Rightarrow 2 \rightarrow \frac{(1,2,5)}{(1,5)} \Rightarrow \frac{2}{2} = \%100 \neq \%70 \text{ (elur.)}$$

$$\checkmark (2,5) \Rightarrow 1 \rightarrow \frac{(1,2,5)}{(2,5)} \Rightarrow \frac{2}{2} = \%100 \neq \%70 \text{ (elur.)}$$

* Yani sonda elimizde kalan 2 kural var bunlar; $(1,5) \Rightarrow 2$ ve $(2,5) \Rightarrow 1$ 'dir.

* Min sup değerini büyük seçersek veri kümesinden bazı örüntüler elde edilmeyebilir, önemli bilgi taşıyan veriler silinebilir.

* Min sup değeri küçük seçersek yöntem karmaşıktır, çok fazla sayıda yoğun nitelik kümesi elde edilir.

(1)

BİRLİKTELİK HAREKETİ

Soru :

Tid	Items
1	1,2,5
2	2,4
3	2,3
4	1,2,4
5	1,3
6	2,3
7	1,3
8	1,2,3,5
9	1,2,3

Tablodan $\text{minsup} = 2$ ve $\text{minconf} = \%70$ değerleri için kuralları elde ediniz.

Not \Rightarrow Burada 2'nin ;

destek değeri = ~~7~~ \rightarrow kaç satırda geçtiği
güven değeri = ~~7~~/9 \rightarrow destek değeri tamamı

Çözüm :

Itemset	Sup
1	6 \rightarrow 1'in tablodaki kaç defa geçtiği.
2	7
3	6
4	2
5	2

\Rightarrow

Itemseti 2'li ler için oluşturalım ;

Itemset	Sup
(1,2)	4
(1,3)	4
(1,4)	1
(1,5)	2
(2,3)	4
(2,4)	2
(2,5)	2
(3,4)	0
(3,5)	1
(4,5)	0

support değeri 2'nin altında olanları eleyelim..
ve buradan 3'lü oluşturalım.

Itemset	Sup
* (1,2,3)	2 ✓
* (1,2,5)	2 ✓

\Leftarrow

3'lü ye sadece bu kadar kaldı.

Çünkü (1,2,3) şartı sağlıyor

yani (1), (2), (3), (1,2), (1,3), (2,3)'ü silmedik

o yüzden (1,2,3) oluşturabiliriz. Şimdi kuralları

incelleyelim..

~~(1,2,3)~~ \Rightarrow ~~(1,2)~~ \Rightarrow 3 \rightarrow $\frac{(1,2,3)}{(1,2)} = \frac{2}{4} = \%50$ \rightarrow ilk tablodan 1,2,3'ün birlikte olduğu durumlara bak.

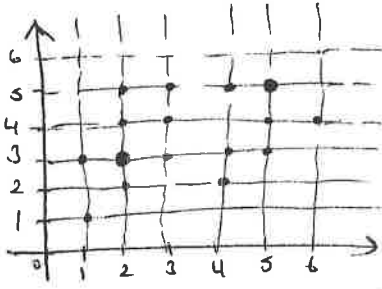
\rightarrow ilk tablodan 1,2'nin birlikte olduğu durum.

~~(1,3)~~ \Rightarrow 2 \rightarrow $\frac{(1,2,3)}{(1,3)} = \frac{2}{4} = \%50$

* Hepsi minconf değerinin altında kalıyor, bu kuralları oluşturamayız.

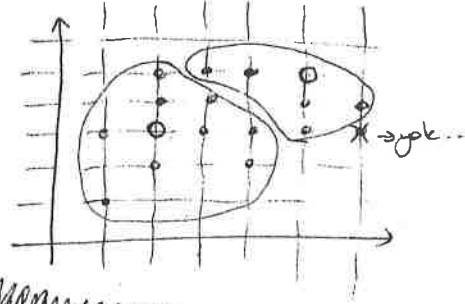
~~(2,3)~~ \Rightarrow 1 \rightarrow $\frac{(1,2,3)}{(2,3)} = \frac{2}{4} = \%50$

K-means Örnek



$k=2$ olsun.

Rastgele 2 nokta
koyduk her nesneyi
en yakın gruba
dahil edelim



~~...~~

* Of veri moderciligi güzel olansın ama kötü derssin "

~~...~~

Soru

	X_1	X_2	D
1 →	2,3	3,4	a
2 →	5,3	4,2	a
3 →	1,1	5,7	a
4 →	-3,1	-0,3	b
5 →	-4,1	-1,5	b
6 →	-0,5	2,3	b

K-means yöntemiyle 3 küme merkezi seç;

$$c_1 = (5,4)$$

$$c_2 = (0,1)$$

$$c_3 = (3,3)$$

1. verinin kümesini bulalım;

$$c_1: \sqrt{(5-2,3)^2 + (4-3,4)^2} = \sqrt{7,29 + 0,36}$$

$$c_1 = 2,765$$

$$c_2: \sqrt{(0-5,3)^2 + (1-4,2)^2} = \sqrt{24,01 + 10,24}$$

$$c_2 = 6,711 \rightarrow 2 \text{ veri}$$

$$c_3: \sqrt{(0-2,3)^2 + (1-3,4)^2} = \sqrt{5,29 + 5,76}$$

$$c_3 = 3,324$$

$$c_3 = \sqrt{0,49 + 1,44} = 1,359$$

✓ Derlen

Toplam maliyeti hesapla = ?

Noktalardan merkezleri çıkar

Soru

Tablo 1'e göre K-nn yöntemiyle sınıflandırma yapılacaktır. Buna göre sınıfı bilinmeyen (1,2) veri noktasının $k=1$, $k=3$, $k=5$ için sınıflarını bularak sınıflandırma hakkında yorumunuzu yazınız.

Soru

x_1	x_2	x_3	D
K	P	G	H
K	N	G	H
B	E	G	H
B	P	G	S
B	N	T	S
K	E	T	S

Naive Bayes sınıflandırıcıya göre (K,P,T) ve (B,P,T) veri noktalarının sınıflarını bulunuz

$$P(S=H | K, P, T) = \frac{3}{6} \left[\frac{3}{4} \cdot \frac{2}{4} \cdot \frac{1}{2} \right] = \frac{3}{16}$$

$$P(S=S | K, P, T) = \frac{3}{6} \left[\frac{1}{3} \cdot \frac{1}{3} \cdot \frac{2}{3} \right] = \frac{1}{27}$$

$$\frac{3}{16} > \frac{1}{27}$$