

1

# VERİ MADENCİLİĞİ

## 0) INTRO (GİRİŞ)

- Predictive  $\rightarrow$  Geçmiş ya da geleceğe yönelik tahminsel (öngörisel)
  - Classification, supervised learning, regression  $\rightarrow$  örnekler
- Descriptive  $\rightarrow$  Veriden formbulayı, yorumlayıcı
  - Clustering (unsupervised), association rules  $\rightarrow$  örnekler
- Data  $\rightarrow$  data preprocessing  $\rightarrow$  data mining  $\rightarrow$  post processing  $\rightarrow$  information } süreç
  - data preprocessing  $\rightarrow$  feature selection, dimensional reduction, normalization, data subsetting
  - post processing  $\rightarrow$  visualization, filtering patterns, pattern interpretation

ANA  
YÖNTEMLER

## 1) DATA PREPROCESSING (ÖN İŞLEM)

Yapılma sebebi  $\rightarrow$  Veri noisy (gürültülü), incomplete (tamamlanmamış), inconsistent (tutaraksız), eksik, tekrarlı, çelişkili olabildiği için

\* Verinin önemli karakteristikleri  $\rightarrow$  Dimensionality, Sparsity, Resolution, Size

\* Veri ön işleme işlemleri:

- Data cleaning  $\rightarrow$  Fill the missing value (eksik, yanlış değerleri doldurma)  $\rightarrow$  naive-bayes, karar ağacı
  - $\rightarrow$  remove outlier  $\rightarrow$  (gürültü yumuşaltma)
  - $\rightarrow$  smoothing noisy data  $\rightarrow$  binning metadları, clustering, regresyon
- Data integration  $\rightarrow$  farklı veri kaynaklarının entegre olması
- Data transformation  $\rightarrow$  normalizasyon (z-score, min-max, decimal scaling)
  - $\rightarrow$  aggregation  $\rightarrow$  birleştirme
- Data reduction (verinin küçültülmesi)  $\rightarrow$  dimensionality reduction (boyutsal)
  - $\rightarrow$  numerosity reduction (sayısal)
  - $\rightarrow$  data compression (verinin sıkıştırılması)
- Data discretization  $\rightarrow$  veri ayrıklaştırılması

\* Data Quality (Veri Kalitesi) Kriterleri

- Kesinlik (Accuracy)  $\rightarrow$  doğru/yanlış
- Tamamlanmışlık (Completeness)  $\rightarrow$  kaydedilmemiş / ulaşılamayan
- Yorumlanabilirlik (Interpretability)  $\rightarrow$  verinin ne kadar kolay anlaşılacağı
- Tutarlılık (Consistency)  $\rightarrow$  verilerin tutarsızlığı / çelişkili olması
- İnanılabilirlik (Believability) - Güncellik (Timeliness)



2

## 2) BINNING METODLARI (KUTULAMA)

⇒ Mean, Median, Boundaries, Mode

SORU: 4, 8, 15, 21, 21, 24, 25, 28, 34 verisindeki görüldüğü şekilde.

1-) Küçükten büyüğe sırala:

Bin1: 4, 8, 15

Bin2: 21, 21, 24

Bin3: 25, 28, 34

} binning } Kaç parça olacağı sordu verilir

2-) Means → Parçaların ortalamaları bulunur, her eleman o değeri alır.

Median → Her eleman, ortadaki elemanın değerini alır. (Veri sayısı çiftse → ortadaki iki değerin ort.)

Boundaries → Her eleman, min ve max olan iki değerden hangisine yakınsa onun değerini alır. Eşit mesafedeysse, min.'u alır.

(means)	(median)	(boundaries)	(mode)
Bin1: 9, 9, 9	Bin1: 8, 8, 8	Bin1: 4, 4, 15	Bin1: 4, 8, 15
Bin2: 22, 22, 22	Bin2: 21, 21, 21	Bin2: 21, 21, 24	Bin2: 21, 21, 21
Bin3: 29, 29, 29	Bin3: 28, 28, 28	Bin3: 25, 25, 34	Bin3: 25, 28, 34

## 3) NORMALİZASYON

a) Z-score Normalizasyonu:

$$v' = \frac{v - \bar{x}}{s_x}$$

v': yeni değer

v: eski "

 $\bar{x}$ : ortalama $s_x$ : standart sapma

\* Standart sapma:

1) Aritmetik ortalamayı bul

2) Her veriden A.O'yu çıkar

3) Sonuçların karelerini al topla

4) Veri sayısının 1 eksikine böl, karekökü al.

b) Min-Max Normalizasyonu:

$$v' = \frac{v - \min(A)}{\max(A) - \min(A)} \cdot [\text{newmax}(A) - \text{newmin}(A)] + \text{newmin}(A)$$

A: ilgili sütun

v: yeni değer

v: eski değer

max(A): sütundaki max

min(A): sütundaki min

newmin(A) &gt; sordu verilen

newmin(A) aralıkların değeri

c) Decimal Scaling Normalizasyonu:

1) Sordu normalizasyon aralığı verilir.

2) Mutlak değerce en büyük olan değeri, verilen normalizasyon aralığında olacağı şekilde, bu sayı j'ın mümkün olan en küçük değerinde olduğu 10<sup>j</sup> sayısı bulunur ve verideki her değer bu sayıya bölünür.

SORU: Veri seti: 10, 20, -52, -700 → [-1, 1] aralığında DS ile normalize et.

- Mutlak değerce en büyük: -700

$$v_1 = 0,01$$

$$v_3 = -0,052$$

$$v_2 = 0,02$$

$$v_4 = -0,7$$

$$10^j : 1000$$



#### 4) CLASSIFICATION TEKNİKLERİ (SINIFLANDIRMA)

##### □ Temel Sınıflandırıcılar:

- Decision Tree: Karar Ağacı
- Nearest-neighbor: KNN
- Neural Networks, Deep Neural Nets (ANN, YSA, Multi Layer Perceptron)
- Naive-Bayes
- Support Vector Machines (SVM)

##### □ Topluluk Sınıflandırıcılar

- Boosting
- Bagging
- Random Forests

#### 5) PRECISION-RECALL-F SCORE, (ALGORİTMA PERFORMANSI)

**SORU:** Varsayalım ki bir doküman koleksiyonunda toplam 50 tane veri madenciliği ile ilişkili doküman yer almaktadır. Koleksiyon toplam 200 dokümandan oluşmaktadır. Otomatik olarak doküman sınıflayan bir veri madenciliği algoritması sonucunda elde edilen durum şu şekildedir. Veri madenciliği ile ilişkili olan 45 adet doküman doğru keşfedilmiştir. Gerçekten veri madenciliği ile ilişkili olmayan 45 doküman ilişkili olarak yakalanmıştır. Bu sınıflama algoritması ile ilgili performans değerleri nedir?

Sınıflama veri madenciliği ile ilişkili kitaplar

$$\text{Precision} = \frac{\text{keşfedilen veri madenciliği kitapları}}{\text{keşfedilen tüm kitaplar}} = \frac{45}{90} = \%50 \quad (0,5)$$

$$\text{Recall} = \frac{\text{keşfedilen veri madenciliği kitapları}}{\text{tüm veri madenciliği kitapları}} = \frac{45}{50} = \%90 \quad (0,9)$$

$$\text{F-score} = \text{Precision ve Recall} = \frac{2 \cdot P \cdot R}{P+R} = \frac{2 \cdot 0,5 \cdot 0,9}{0,5+0,9} = \frac{0,9}{1,4} = 0,64$$

#### 6) ENTROPY (KARAR AĞAÇLARI-1)

⇒ Entropi rastgeleliği, belirsizliği ve beklenmeyen durumun ortaya çıkma olasılığını gösterir.

⇒ Decision tree çizimi ve information gain (bilgi kazancı) hesaplamasında kullanılır.

⇒  $p_1, p_2, \dots, p_n$  toplamları 1 olan olasılıklardır.

$$\text{Entropy} = E(p_1, p_2, \dots, p_n) = - \sum_{i=1}^n p_i \cdot \log_2(p_i)$$

⇒ S veri kümesinde 14 örnek: 9'u A sınıfına, 5'i B sınıfına ait.

$$E(p_1, p_2, \dots, p_n) = - \sum_{i=1}^n p_i \cdot \log_2(p_i) \Rightarrow E(p_1, p_2) = - \frac{9}{14} \cdot \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \cdot \log_2\left(\frac{5}{14}\right) = 0,940$$



4

⇒ Sınıflandırmada; olayın olması beklenen bir durumda entropi = 0

↓

⇒ Örnekler aynı sınıfa aitse, entropi = 0

⇒ Örnekler sınıflar arasında eşit dağıtılmışsa, entropi = 1

⇒ Örnekler sınıflar arasında rastgele dağılmışsa,  $0 < \text{entropi} < 1$

### (7) BİLGİ KAZANCI (KARAR AĞAÇLARI-2)

⇒ Ağaç oluşturmada çoklu dallanma yaklaşımıdır.

⇒ Bilgi kazancı en büyük olan node, root (kök) düğümdür. (Bilgi kazancına göre nitelik seçilir)

⇒ Sınıflandırma sonucu için en az sayıda karşılaştırma yapmayı hedefler.

⇒ Bilgi Kazancı = Genel Entropi - Sütun Entropi

⇒ Genel Entropi, en sağdaki sütun ile hesaplanır.

⇒ Sütun Entropi =  $E_1 \cdot p_1 + E_2 \cdot p_2 + \dots + E_n \cdot p_n$

Örnek:

$$\text{Genel Entropi} = -\frac{4}{10} \cdot \log_2\left(\frac{4}{10}\right) - \frac{6}{10} \cdot \log_2\left(\frac{6}{10}\right) = 0,9709$$

A	B	Label
T	F	+
T	T	+
T	T	+
T	F	-
T	T	+
F	F	-
F	F	-
F	F	-
T	T	-
T	F	-

	(+)	(-)		(+)	(-)
A ⇒ T	4	3	B ⇒ T	3	1
F	0	3	F	1	5

$$E(A, \text{label}) = E(T) \cdot p(T) + E(F) \cdot p(F)$$

$$= \left[ \left( -\frac{4}{7} \cdot \log_2\left(\frac{4}{7}\right) - \frac{3}{7} \cdot \log_2\left(\frac{3}{7}\right) \right) \cdot \frac{7}{10} \right] + \left[ \left( -\frac{0}{3} \cdot \log_2\left(\frac{0}{3}\right) - \frac{3}{3} \cdot \log_2\left(\frac{3}{3}\right) \right) \cdot \frac{3}{10} \right]$$

$$= 0,6897$$

$$A'nın bilgi kazancı = 0,9709 - 0,6897 = 0,2812$$

# B için de aynı yapılır ve B.K en büyük olan root seçilir.

### (8) GINI İNDEKSİ (KARAR AĞAÇLARI-3)

⇒ Ağaç oluşturmada, tekli dallanma yaklaşımıdır.

⇒ Gini index değeri en küçük olan nitelik seçilir.

$$\Rightarrow \text{GINI}(S) = 1 - \sum [p_i]^2$$

Örnek:

C <sub>1</sub>	2
C <sub>2</sub>	4

$$p(C_1) = \frac{2}{6}, p(C_2) = \frac{4}{6} \Rightarrow \text{Gini} = 1 - \left(\frac{2}{6}\right)^2 - \left(\frac{4}{6}\right)^2 = 0,444$$



## 9) KARAR AĞAÇLARI (4)

⇒ Akış diyagramı şeklinde bir ağaç yapısıdır.

⇒ Veri kümesinden, belli durumlarda nasıl davranacağı belirtilen bir ağaç çıkarılır.

⇒ Decision tree (karar ağacı) oluşturma algoritmaları; ağacı daha küçük, kompakt, sıkıştırılmış ve verimli hale getirmeyi hedefler.

⇒ Karar ağacı oluşturma algoritmaları:

- Hunt's Algorithm
- CART
- ID3, C4.5
- SLIQ, SPRINT

⇒ Karar ağacı modelleri, eğitim verilerindeki detayları öğrenir ve uygular.

⇒ Hem kategorik hem de sayısal verileri işleyen tek yöntemdir.

⇒ En iyi bölme belirlenmesinde "Greedy yaklaşımı" kullanılır (çoğunlukla aynı sınıfa ait örneklerin bulunduğu dallar tercih edilir. (homojen sınıf dağılımına sahip olan dallar))

⇒ Düğüm sarılığı ölçümü 3 yolla yapılabilir:

- Entropy ⇒ information gain (bilgi kazancı) ile → çoklu dallanma
- Gini index ⇒ ikize dallanma (binary)
- Misclassification error

⇒ Ağaç Oluşturma Yaklaşımları:

- Kökten yapraklara ulaşılır (parçala ve çöz = divide and conquer)
- Dügümü dallara ayırmada Greedy yaklaşımı kullanılır. Her adımda en iyi bölme bulunur. Her düğümden dallanmak için en iyi özellik seçilir. Her dal için aynı yaklaşım uygulanır.
- Yaprak düğümler en çok örneği olan sınıfla etiketlenir.

⇒ Karar ağaçlarının uygulanması kolaydır.

⇒ Karar ağaçları, bilinmeyen kayıtları sınıflandırmada son derece hızlıdır.

⇒ Küçük boyutlu ağaçların yorumlanması kolaydır.

## 10) KNN (K-Nearest Neighbor)

⇒  $k=3, 5, \dots$  gibi tek sayılar kullanılarak belirlenir. Verilen verinin, eğitim verisindeki değerlerle öklid uzaklığı ölçülür ve en yakın  $k$  komşuya göre benzerlik tayin edilir.

⇒ Supervised bir tekniktir ⇒ label (etiketler) olduğu için

⇒ Feature similarity measure yöntemi kullanır ⇒ parameter tuning

⇒ Avantajları:

- Uygulanması basittir.
- Gürültülü verilere karşı etkilidir.
- Eğitim dokümanlarının sayısı fazla ise daha etkilidir.



⇒ Dezavantajları:

- Algoritma başlangıçta  $k$  parametresine ihtiyacı duyar.
- Bu ryi sonucu alınabilmesi için hangi özelliklerin uygulanacağı ve hangi özelliklerin alınacağı bilgisi açık değildir.
- Hesaplama maliyeti yüksektir.

⇒ Lazy learner (tembel) bir tekniktir. Decision tree ve, eager learner tekniğidir.

⇒ Tüm eğitim dataseti kullanır ⇒ eğitim adımı yok

⇒ Generalized model oluşturmaz.

## 11) NAIVE BAYES

⇒ Bayes teoremini dayanak alan bir yöntemdir.  $P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$

⇒ Class labeler ile sınıflandırılarak olasılıklar birbirleri ile çarpılır ve son olarak class label olasılığı ile çarpılarak son durum bulunur.

⇒ Örneklerin hangi sınıfa, hangi olasılıkla ait olduklarını belirtir.

⇒ Nitelikler birbirinden bağımsızdır ve hepsi aynı derecede önemlidir.

⇒ Koşullu olasılıklardan birinin 0 olduğu büyük veri setlerinde, ilgili koşullu olasılığın payına 1 eklenerek sonuç 0'dan kurtarılır ⇒ Laplace düzeltmesi

⇒ Avantajları:

- Gerçekleşmesi kolaydır.
- Çoğu durumda ryi sonuçlar verir.
- Alakasız niteliklere ve fazla girili noktalarına dirençlidir.

⇒ Dezavantajları:

- Niteliklerin bağımsız olduğunu kabul eder ancak gerçek hayatta değişkenler birbirine bağımlıdır. Bu sebepten dolayı değişkenler arası ilişki modellenemez.

⇒ Çözüm ⇒ bayes ağları

## 12) REGRESYON

⇒ Classification ⇒ kesikli veriler  
Regression ⇒ sürekli veriler

⇒ Regression, verilerden yola çıkarak ortaya bir takım ilişkiler koymakta kullanılır. Bağımsız bir değişkenin değerleri için, bağımlı değişken tahmin edilmek istenir.

⇒ Linear Regression: Bağımlı bir değişken ( $Y$ ) ile bağımsız bir değişken ( $X$ ) serisinin arasındaki ilişkinin doğrusal bir fonksiyonel biçimidir ve veriyi lineer bir eğriye uydurma çalışmasıdır. (Linear modelleme)

⇒ Regression Amaçları:

- $Y$ 'nin  $X$ 'e bağlı ilişkisinin ortaya çıkarılması ⇒ istatistiksel olarak anlamlı olması
- Yeni  $X$  değerleri için  $Y$ 'nin tahmin edilmesi ⇒ daha önce verilmemiş veriler verildiğinde, ortaya veya geleceğe yönelik tahmin ortaya koyma



⇒ Lojistik Regresyan: Bağımlı=hedef=çıktı değişkeninin binary (0,1) olması söz konusuysa kullanılır.

Classification problemlerinde kullanılmak üzere geliştirilmiş ve tüm kesikli değerler ile çalışabilmektedir.

⇒ Linear regresyon, en fazla kullanılan ve en basit regresyon çeşididir. Bağımlı ve bağımsız değişken arasında doğu bir ilişki çizerek aradaki bağıntıyı ortaya çıkarır. Sadece sürekli verilerde çalışır. Kategorik verilerde çalışmaz.

⇒ Lojistik regresyon ⇒ Kesikli (belirli sayıda değer alabilme) ve sürekli (bir aralıktaki tüm değerleri alabilme) verilerle çalışabilir.

### [13] SVM (Destek Vektör Makineleri)

⇒ Sınırlandırmayı doğrusal veya doğrusal olmayan bir fonksiyon yardımıyla yerine getirir.

⇒ Veriyi birbirinden ayırmak için en uygun fonksiyonun tahmin edilmesi esasına dayanır.

⇒ SVM; böl, parçala, yönet metoduyla çalışır ⇒ Her problemi ikili (binary) şekilde böler.

⇒ Kümelendirme ve regresyon analizinde de kullanılır.

⇒ Sınırlar lineer olmadığında, kernel trick yoluyla problem bir boyut yukarı taşınarak çözülür.

⇒ Doğrusal olmayan bir işlem ile  $n$  boyutlu veri kümesi  $m > n$  olacak şekilde dönüştürülür.

⇒ Yüksek boyutta doğrusal sınıflandırma işlemi yapılabilir.

⇒ Uygun bir dönüşüm ile her zaman veri bir hiper düzlem ile iki sınıfa ayrılabilir.

⇒ Hiper düzleme en yakın öğrenme verileri destek vektörleri olarak adlandırılır.

⇒ Overfitting durumunda  $C$  parametresi küçültülür.

⇒ Öğrenme problemi dışbükey bir optimizasyon problemi olarak formüle edilmiştir.

### [14] YSA (Yapay Sinir Ağları)

⇒ Lineer olmayan veriler arasındaki ilişkileri bulabilen güçlü bir yapıdır.

⇒ Perceptron: Yapay sinir hücresi ⇒ en temel (1 girdi, 1 çıktı katmanı)

⇒ Girdi tüm değerlerin istenilen formata dönüştürülmesi için bir aktivasyon fonksiyonu gerekir.

⇒ Aktivasyon fonksiyonları kullanılarak, çıktı değerinin istenen sınıflara ait olması sağlanır.

⇒ Sık kullanılan aktivasyon fonksiyonları: Sigmoid, Tanh, ReLU, Linear, Sign

⇒ Sigmoid, Tanh, ReLU aktivasyon fonksiyonları; doğrusal olmayan (nonlinear) karar düzlemlerinin öğrenilmesine yardım eder.

⇒ Öğrenme katsayısı ( $\lambda$ =learning rate) gözlem altında tutulması gereken bir değerdir.

⇒ YSA, genellikle birden fazla gizli (orta) katmandan oluşur.



⇒ Epoch: Öğrenmenin gerçekleştiği dönüye denir.

⇒ Her bir epochta öğrenmenin kalitesi veya hata(loss) miktarı ölçülür.

⇒ Hata ve kayıpları azaltmak için iterasyon yapılır. (Öğrenme, süreçlerin sürekli tekrarlendiği bir döngüdür)

⇒ Optimizasyonda gradient descent algoritması kullanılır.

⇒ YSA ile öğrenme aşaması:

1) Yapay sinir ağı oluşturma

- giriş verisini modelleme; gizli katman sayısını ve bu katmanlardaki nöron sayısını belirleme

2) Yapay sinir ağı eğitme

3) Yapay sinir ağıni küçültme

4) Sonucu yorumlama

⇒ YSA oluşturma

• giriş nöron sayısı = öğrenme kümesindeki verilerin nitelik sayısı

• gizli nöron sayısı ⇒ öğrenme sırasında ayarlanır

• çıkış nöron sayısı = sınıf sayısı

⇒ YSA eğitme: Amaç, veri kümesindeki örneklerin hepsini doğru sınıflandıracak ağırlıkları belirlemektir.

1) Ağırlıklara rastgele değerler ata

2) Öğrenme kümesindeki giriş değerlerini teker teker sinir ağına uygula

3) Çıkışı hesapla

4) Hata değerini hesapla

5) Ağırlıkları hata fonksiyonunu en küçültmek şeklinde yor.

⇒ YSA küçültme ⇒ Ağırlıklardan bazıları sınıflandırma sonucunu etkilemeyecek şekilde silinir

⇒ Avantajları:

- Doğru sınıflandırma oranı genelde yüksek

- Kararlı - öğrenme kümesinde hata olduğu durumda da çalışır

- Çıkış ayık, sürekli ya da ayık veya sürekli değişkenlerden oluşan bir vektör olabilir.

⇒ Dezavantajları:

- Öğrenme süresi uzun

- Öğrenilen fonksiyonun anlaşılması zor

⇒ Perceptron ⇒ ara katman içermez ⇒ doğrusal olmayan ayırma yapamaz.

⇒ Multi-layer Neural Network (Çok katmanlı Sinir Ağı) ⇒ En az bir ara katman içerir ve doğrusal olmayan karar yüzeylerini içeren her türlü sınıflandırma görevini çözebilir.

⇒ Gizli katmandaki her düğüm, önceki katmandan gelen aktivasyonlar üzerinde çalışır ve aktivasyonları bir sonraki katmanın düğümlerine iletilir. Heri baslemeli sinir ağıları olarak da adlandırılır.

⇒ Her gizli katman bir soyutlama düzeyini temsil eden karmaşık özellikler, daha basit özelliklerin bileşimleridir. Katman sayısı, YSA derinliği olarak bilinir. Daha derin ağılar, karmaşık özellikler hiyerarşisini ifade eder.

⇒ YSA, kendisine örnekler gösterdikçe bu ağırlık değerlerini değiştirir. Hedef, ağıya gösterilen örnekler için doğru çıkışı verecek ağırlıkları bulmaktır.



⇒ Biyolojik sinir hücresi → Yapay sinir hücreleri

Akson → Çıktı  
Dendrit → Toplama fonksiyonu  
Çekirdek → Aktivasyon fonksiyonu  
Sinaps → Ağırlıklar

⇒ YSA'nın öğrenme süresi uzundur fakat öğrenilen durumun yeni örneklerde test edilmesi yarıda-  
gerlendirilmesi hızlıdır.

⇒ YSA'da bilgi:

- Bilgi ağı bağlantı ağırlıkları ile temsil edilir.
- YSA zeka ağırlıklıdır.
- Ağırlık sahip olduğu ağırlık değerlerinin dağılımı olduğu ölçüde ağırlık performansı fazladır.
- Bilgi dağıtılmıştır, ağırlık değerlerinin bazıları kayıbsa dahi ağırlık çalışmasını sürdürebilir.

⇒ YSA'daki her hata, çıkış değerlerinin ve bias değerlerinin çıkışa yansımalarıdır.

⇒ YSA eğitimi batch processing (cevrimici) moda yapılır.

⇒ Büyük öğrenme oranında; sistem veriyi çok hızlı öğrenir, toplam hata artar. Düşük öğrenme oranı  
sistem çok yavaş öğrenir, eğitim zamanı artar ancak hata azalır.

⇒ YSA için lokal minimum hata tuzuguna düşmek doğaldır ve bu durum global minimuma ulaşma-  
sına engel olabilir. Momentum, öğrenme algoritmasını lokal minimumdan kurtarıp sekilde önceki yönde  
tutmaya çalışır.

⇒ Bu tür sonuçlar alınmadan birçok deneme yapılmalıdır. 0,6-0,8 aralığındaki değerler önerilir.

⇒ YSA durdurma kriterleri:

- Zaman eşik değere ulaştığında bitebilir.
- Önceden tanımlı epoch değere ulaştığında bitebilir.
- Önceden tanımlı hata değere ulaştığında bitebilir.
- İki epoch arasındaki hata değeri azaldığında eğitim bitebilir.

⇒ Örnekler eğitim kümesinden rastgele seçilir.

⇒ Bias:

- Sistem performansını etkiler.
- Öğrenmede bir ağırlık gibi alınır.

- Giriş sinyallerinin toplamı 0 olduğunda öğrenme gerçekleşmez. Çıkış değerleri hep 1 olan bi-  
nöröler, nöronların giriş sinyallerinin sürekli sıfırdan farklı olmasını sağlar.

- Öğrenmeyi hızlandırırken, yerel optimum değerlere takılmayı gösterir.

⇒ YSA yapılarına göre ileri beslemeli ve geri beslemeli olarak ikiye ayrılır.

⇒ Backpropagation: İleri beslemeli YSA'nı eğitmek için yaygın olarak kullanılan bir algoritma.  
Geri yayılım algoritması, zincir kurallı tarafından her ağırlığa göre kayıp fonksiyonunun  
hesaplayıcı, gradyanı tek bir eğerde bir katman olarak hesaplanır. Zincir kurallıdaki ara testler  
geriye hesaplanmaları önlemek için son katmandan geriye doğru yineliyerek çalışan bir dinamik  
programlama örneğidir.



⇒ Backpropagation'daki ağırlık güncelleme dengesi binlerce kez iteratif olarak devam edebilir. Belirli bir iterasyon sonucunda dengesi sonlandırılabilir ya da hata değeri göre sonlandırılabilir. Bu seçim kritiktir; çünkü az sayıda iterasyon hatayı yeterince düşürmez, çok sayıda iterasyon ise overfitting'e sebep olabilir.

### [15] MODEL SEÇİMİ, DEĞERLENDİRMESİ VE ENSEMBLE

⇒ Model başarısı değerlendirilirken, validation test dikkate alınır (training = eğitim seti değil)

⇒ Modelin başarısını ölçmede kullanılan metodlar:

- Confusion matrix (karışıklık matrisi)
- F-score
- Cross validation (k-fold)
- Bootstrap
- Confidence intervals } iki modelin başarısını karşılaştırma/ölçme
- ROC eğrisi
- Holdout
- Random subsampling

⇒ Bias → isabet, variance → farklılık

↓

high bias ve low variance → underfitting

low bias ve high variance → overfitting

⇒ Yüksek variance → eğitim verilerinde düşük hata, dağıtım verilerinde yüksek hata durumu

- model eğitim verisi eklemek } çözüm
- regularization (düzenleştirme)

⇒ Yüksek bias → Hem eğitim hem de dağıtım verilerinde yüksek hata durumu

- model eğitiminde iterasyonlara devam etmek } çözüm
- öğrenme algoritmasını farklılaştırmak

⇒ Regularization yöntemleri:

- Lasso regularization } ağırlık azaltımı
- Ridge regularization
- Early stopping regularization
- Dropout regularization
- Data augmentation
- Batch Normalization
- Ensemble Learnings

⇒ Curse of dimensionality (boyut laneti): Çok yüksek boyutlu veri kümelerinde karışıklık problemi

- veri seyrekliği
- mesaj konsantrasyonu

⇒ Ensemble: Tek bir hipotez yerine hipotez uzayından bir hipotez topluluğu seçmek ve tahminleri birleştirmektir. Tek bir güçlü sınıf yerine çok sayıda zayıf sınıflayıcı kullanılır. Hipotezler birbiriyle farklılaştıkça hataları arasındaki korelasyon azalır ve topluluk öğrenmesi daha kullanışlı olur.



⇒ Ensemble sınıflayıcıları oluşturma:

- Eğitim setini manipüle ederek
  - bagging → ortalama, majority voting
  - boosting → ağırlık (özellik) ile, weighted voting
- Girdi özelliklerini (input features) manipüle ederek
  - random forests
- Sınıf etiketlerini (class label) manipüle ederek
  - hata düzeltici çıktı kodlaması
- Öğrenme algoritmasını manipüle ederek:
  - YSA veya karar ağacında rastgelelik eklemek

⇒ Random forest → Kapsula göre hangi algoritmanın çalışacağını, hangi şartta hangi algoritmanın daha başarılı olacağını belirleyen algoritmadır.

⇒ Adaboost → Algoritmalar basarılarına göre bir ağırlığa sahip olur ve ağırlık zaman içinde değişir. Toplam hatanın kontrol altına alınması amaçlanır.

## 16) CLUSTERING

⇒ Clustering algoritmaları verileri gruplayarak kümeleme işlemi yapar.

⇒ Unsupervised bir tekniktir. Sınıf etiketi (class label) kullanmaz, feature similarity yapar.

⇒ Similarity matrix hesaplayarak benzer özelliklere sahip verileri kümeler. Bu hesaplamada çeşitli metriklerden yararlanılır.

↳ Öklid, Minkowski, Manhattan

⇒ Şekil, yoğunluk, konum bağıli kümeleme yapılabilir.

⇒ Kümeleme (clustering) algoritmaları:

- K-means ve varyantları
- Hierarchical clustering
- Density-based clustering

⇒ Hierarchical clusteringte küme arası uzaklığı belirleme:

- MIN (tek bağ)
- MAX (tam bağ)
- Ortalama
- Merkezler arası uzaklık

## 17) K-MEANS

⇒ Algoritma aşamaları:

- 1) Veri kümesi  $k$  altkümeğe ayrılır.
- 2) Her kümenin ortalaması hesaplanır merkez nokta
- 3) Her nesne en yakın merkez noktanın olduğu kümeğe dahil edilir.
- 4) Nesnelerin kümelenebilirliği doğruluk alıncaya kadar 2. adıma geri dönlür.

⇒ K-means değerlendirme: Hataların karesi toplamını (SSB) azaltmak için;

- $k$  küme sayısı artırılabilir
  - Başlangıç için farklı merkez noktaları seçilebilir
- } en az SSB değerine sahip olan kümeleme seçilir



→ Karmasıklığa etkili olan faktörler:

- Nokta sayısı
- Kume sayısı
- İterasyon sayısı
- Nitelik sayısı

→ Soru çözümü:

- k kadar rastgele nokta seç.
- Seçilen noktalar ile diğer noktaların mesafelerini ayrı ayrı hesapla
- En küçük mesafeye göre I'ları kümele.
- k kadar noktadan, new c değerleri, kümeledeki I'ların aritmetik ortalamasıyla bulunur.
- Grafik çizilir, I'lar kümeden, new c değerleri belirtilir.
- İterasyon sayısı kadar devam edilir. (c → new c)

### 18) ASSOCIATION RULE VE APRIORI ALGORİTMASI

$$A \rightarrow B \quad \text{Support} = \frac{A, B \text{ sayısı}}{\text{satır sayısı}} \quad ; \quad \text{Confidence} = \frac{A, B \text{ sayısı}}{A \text{ sayısı}}$$

$$\Rightarrow \text{support count} = \text{support} \times \text{satır sayısı}$$