

Büyük Veri İçin İstatistiksel Öğrenme (Statistical Learning for Big Data)

M. Ali Akcayol
Gazi Üniversitesi
Bilgisayar Mühendisliği Bölümü

Bu dersin sunumları, "The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Trevor Hastie, Robert Tibshirani, Jerome Friedman, Springer, 2017." ve "Mining of Massive Datasets, Jure Leskovec, Anand Rajaraman, Jeffrey David Ullman, Stanford University, 2011." kitapları kullanılarak hazırlanmıştır.

Konular

- Sınıflandırıcıların Değerlendirilmesi
- Skorlar
 - Karışıklık matrisi
 - Accuracy
 - Precision
 - Recall
 - Specificity
 - F-Score
- Eğitim ve Test Kümeleri
 - Hold-out set
 - Multiple random sampling
 - k-fold cross validation
 - Bootstrap

Sınıflandırıcıların Değerlendirilmesi

- Bir sınıflandırıcı geliştirildiğinde doğruluk değerinin belirlenmesi gereklidir.
- Doğruluđu test edilmeden bir sınıflandırıcı gerçek hayattaki problemlerde kullanılamaz.
- Bir sınıflandırıcının değerlendirilmesi için çok sayıda yöntem ve ölçüt vardır.
- **Temel değerlendirme ölçütü sınıflandırmadaki doğruluk oranıdır (accuracy).**

$$Accuracy = \frac{\text{Number of correct classifications}}{\text{Total number of test cases}}$$

3

Sınıflandırıcıların Değerlendirilmesi

- Bazı uygulamalarda **hata oranı (error rate)** değeri kullanılmaktadır.

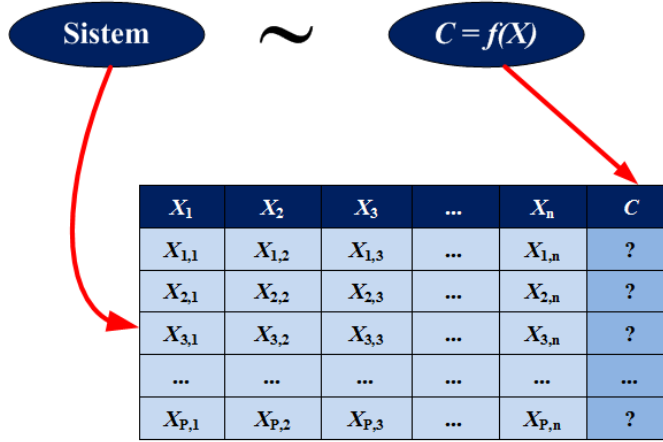
$$error\ rate = 1 - accuracy$$

- Genellikle istatistiksel testler kullanılarak farklı sınıflandırıcıların hangisinin daha uygun olduğuna karar verilebilir.
- Birden fazla sınıflandırıcıyı değerlendirmek için **aynı eğitim verisi ve aynı test verisi kullanılarak doğruluk oranları elde edilir.**

4

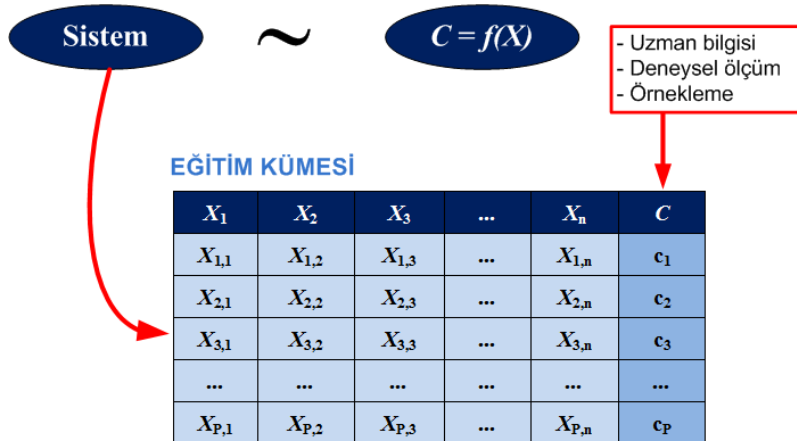
Sınıflandırıcıların Değerlendirilmesi

- Sınıflandırma probleminde, **giriş değerleri ile çıkış sınıfları arasında** ilişkilendirme yapan **bir fonksiyon bulunur**.



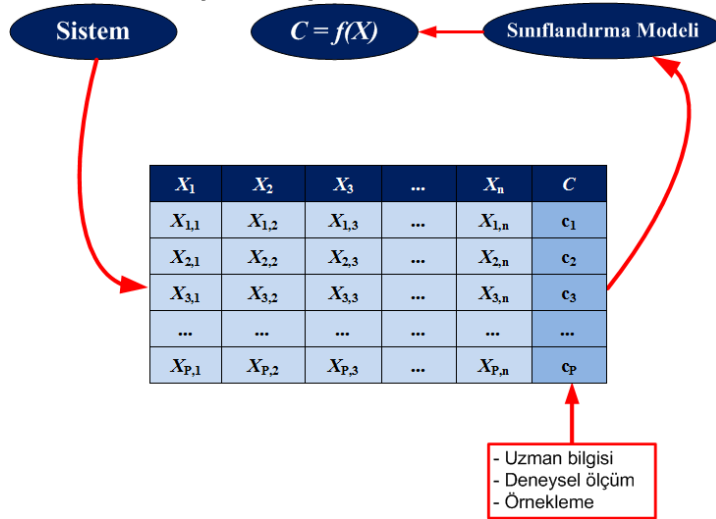
Sınıflandırıcıların Değerlendirilmesi

- Eğitim kümesinde **girişler ve çıkış/lar arasında** ilişkilendirmeler sağlanır.



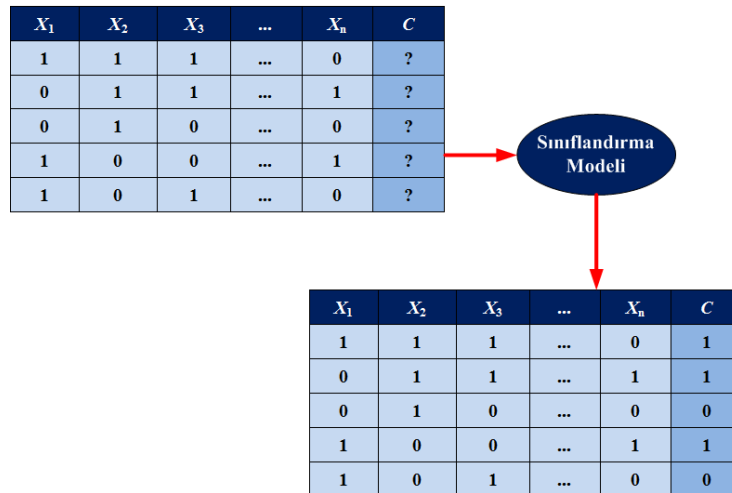
Sınıflandırıcıların Değerlendirilmesi

- Sınıflandırma modeli, eğitim kümesinde giriş/ler ile çıkış/lar arasında bir fonksiyon oluşturur.



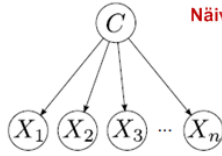
Sınıflandırıcıların Değerlendirilmesi

- Sınıflandırma modeli, yeni girişler için anlamlı sınıflandırma etiketleri belirler.

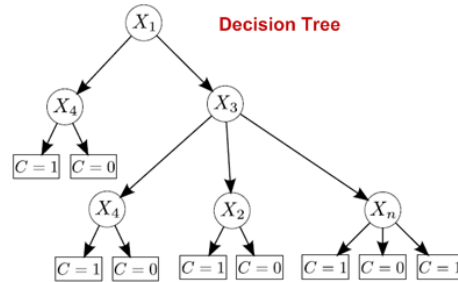
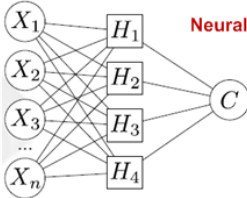


Sınıflandırıcıların Değerlendirilmesi

- Farklı sınıflandırma yaklaşımları kullanılabilir.



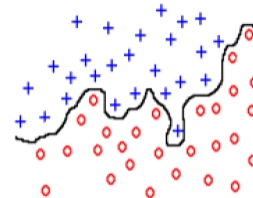
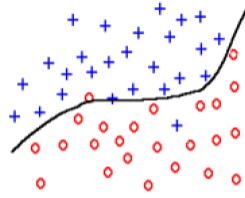
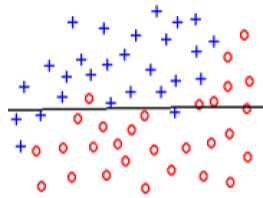
X_1	X_2	X_3	...	X_n	C
1	1	1	...	0	1
0	1	1	...	1	1
0	1	0	...	0	0
1	0	0	...	1	1
1	0	1	...	0	0



Sınıflandırıcıların Değerlendirilmesi

Underfit, fit ve overfit

- Underfit**, eğitim kümesi için iyi düzeyde sınıflandırma yapamaz.
- Fit**, hem eğitim kümesi hem de test kümesi için iyi düzeyde sınıflandırma yapabilir.
- Overfit**, eğitim kümesi için iyi sınıflandırma yapar ancak test kümesi için iyi sınıflandırma yapamaz.



Konular

- Sınıflandırıcıların Değerlendirilmesi
- **Skorlar**
 - Karışıklık matrisi
 - Accuracy
 - Precision
 - Recall
 - Specificity
 - F-Score
- Eğitim ve Test Kümeleri
 - Hold-out set
 - Multiple random sampling
 - k-fold cross validation
 - Bootstrap

11

Skorlar

- Bir sınıflandırıcının ne derece doğru sınıflandırma yaptığının gerçek uygulamada kullanılmadan önce belirlenmesi gereklidir.
- **Farklı eğitim kümeleri için aynı sınıflandırıcı farklı çözümler üretebilir.**
- **Aynı eğitim kümesi için de farklı sınıflandırıcılar farklı çözüm üretebilir** ve bunların en uygun olanının seçilmesi gereklidir.
- **Bir eğitim kümesi için aynı sınıflandırıcı modeli, farklı parametre değerleri için de farklı çözümler üretebilir.** Uygun parametre değerlerinin belirlenmesi gereklidir.
- **Kullanım amaçlarına, problemin büyüklüğüne ve beklenen doğruluk düzeyine göre** farklı skor değerleri kullanılarak sınıflandırıcılar değerlendirilebilir.

12

Konular

- Sınıflandırıcıların Değerlendirilmesi
- Skorlar
 - Karışıklık matrisi
 - Accuracy
 - Precision
 - Recall
 - Specificity
 - F-Score
- Eğitim ve Test Kümeleri
 - Hold-out set
 - Multiple random sampling
 - k-fold cross validation
 - Bootstrap

13

Karışıklık matrisi

- Karışıklık (confusion) matrisi ile, örnek kümedeki gerçek sınıf etiketi ile tahmin edilen sınıf etiketi sayıları gösterilir.

		TAHMİN EDİLEN		TOPLAM
		C ⁺	C ⁻	
GERÇEK	C ⁺	TP True Pozitif (Hits)	FN False Negatif (Miss)	N ⁺ Gerçek Pozitif sayısı
	C ⁻	FP False Pozitif (Miss)	TN True Negatif (Hits)	N ⁻ Gerçek Negatif sayısı
TOPLAM		N̂ ⁺ Tahmin Pozitif sayısı	N̂ ⁻ Tahmin Negatif sayısı	N Toplam Örnek sayısı

14

Karışıklık matrisi

- Çok sayıdaki sınıf için örnek aşağıdadır.

		TAHMİN EDİLEN					TOPLAM
		C ₁	C ₂	C ₃	...	C _n	
GERÇEK	C ₁	T ₁	F ₁₂	F ₁₃	...	F _{1n}	N ₁
	C ₂	F ₂₁	T ₂	F ₂₃	...	F _{2n}	N ₂
	C ₃	F ₃₁	F ₃₂	T ₃	...	F _{3n}	...

	C _n	F _{n1}	F _{n2}	F _{n3}	...	T _n	N _n
TOPLAM		Ñ ₁	Ñ ₂	Ñ ₃	...	Ñ _n	N Toplam Örnek sayısı

15

Konular

- Sınıflandırıcıların Değerlendirilmesi
- Skorlar
 - Karışıklık matrisi
 - Accuracy
 - Precision
 - Recall
 - Specificity
 - F-Score
- Eğitim ve Test Kümeleri
 - Hold-out set
 - Multiple random sampling
 - k-fold cross validation
 - Bootstrap

16

Accuracy

- Bir sınıflandırıcının **doğru sınıflandırdığı örnek sayısının toplam örnek sayısına oranıdır.**

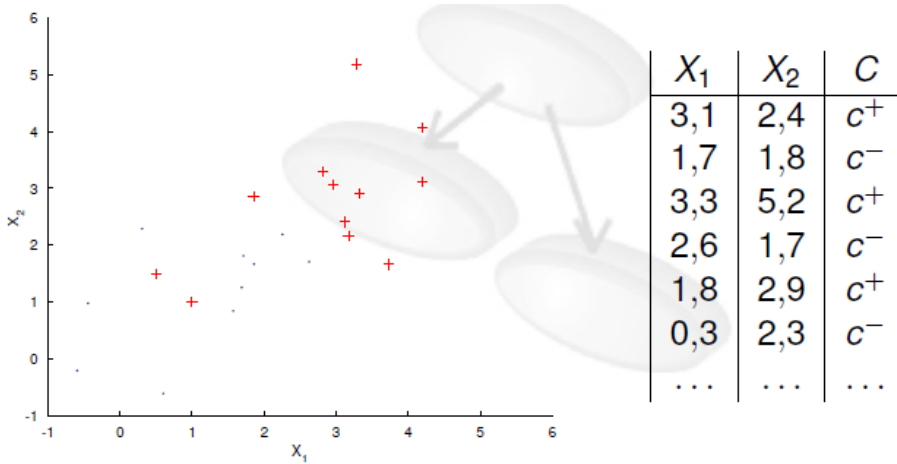
$$Accuracy = \frac{\text{Number of correct classifications}}{\text{Total number of test cases}}$$

- **Doğruluk değerlendirmesi test kümesi kullanılarak hesaplanır.**
- Eğitim sırasında kullanılmayan test verilerinde doğru sınıflandırdığı örnek sayısı alınarak doğruluk düzeyi hesaplanır.

17

Accuracy

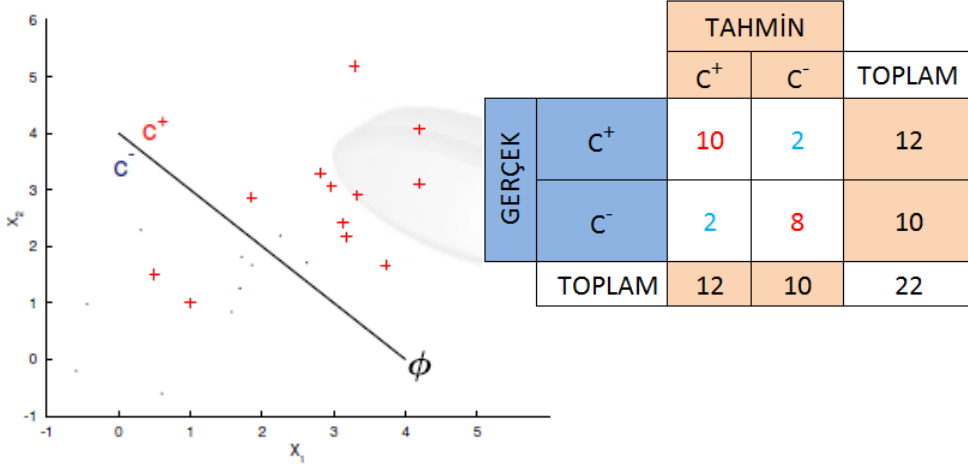
- İyi bir sınıflandırıcının test kümesindeki tüm sınıfları doğru tahmin etmesi beklenir.



18

Accuracy

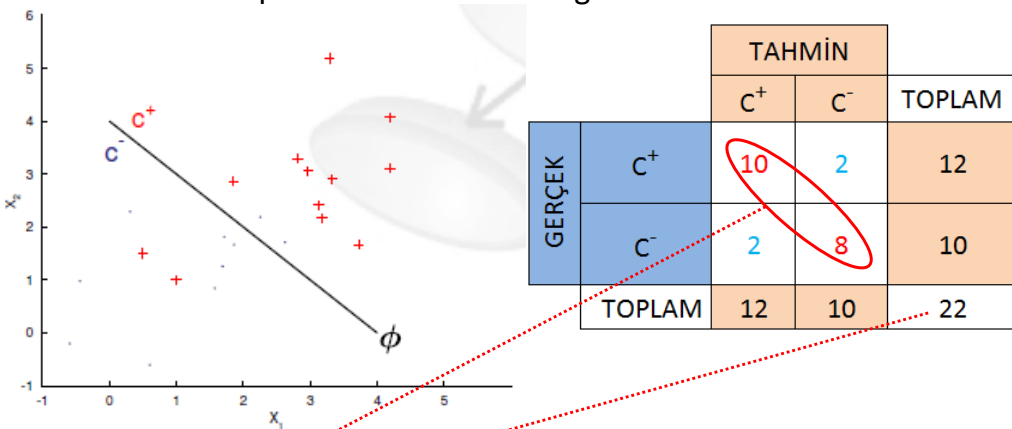
- İki sınıfa sahip bir sınıflandırıcı örneği.



19

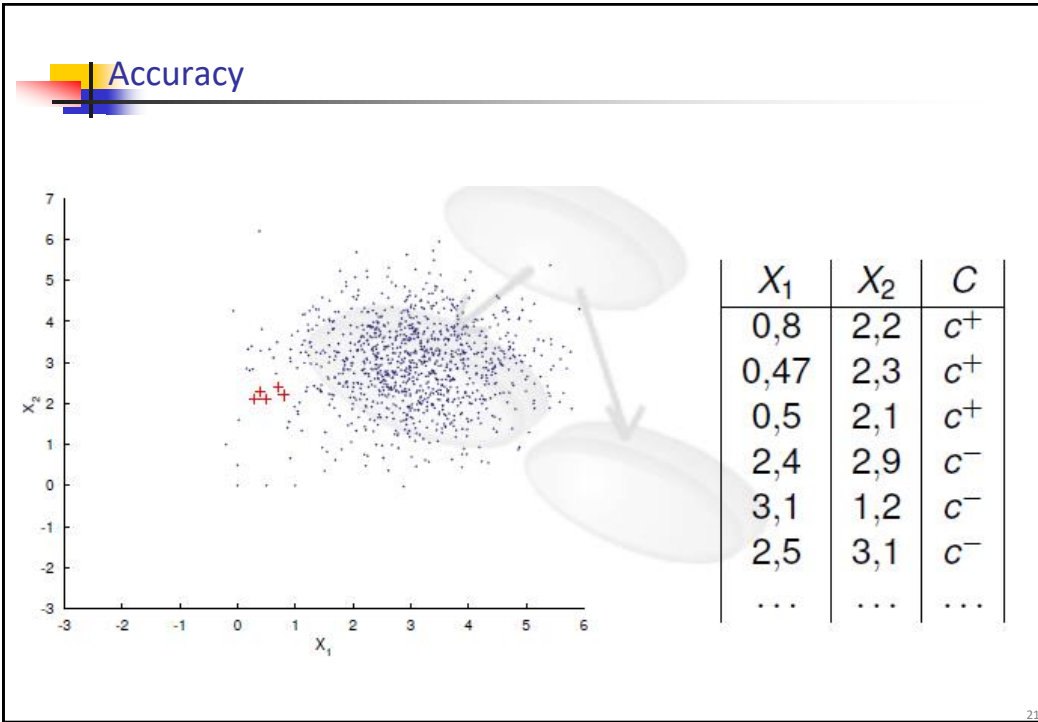
Accuracy

- İki sınıfa sahip bir sınıflandırıcı örneği.

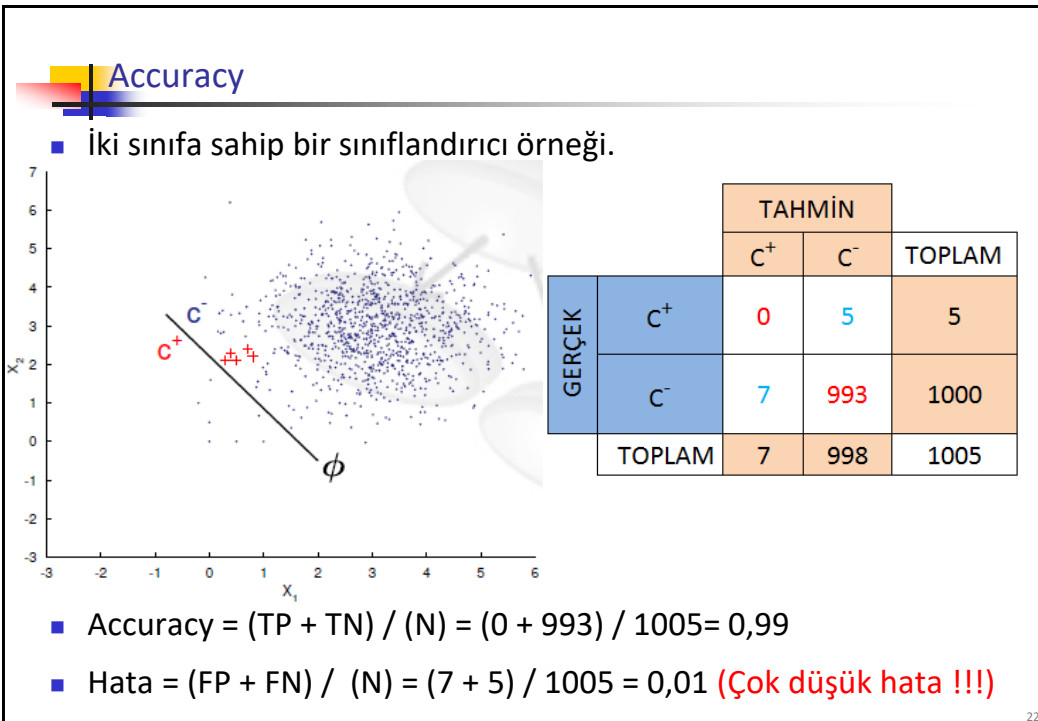


- Accuracy = $(TP + TN) / (N) = (10 + 8) / 22 = 0,82$
- Hata = $(FP + FN) / (N) = (2 + 2) / 22 = 0,18$

20



21



22

Accuracy

- İki sınıfa sahip bir sınıflandırıcı örneği.

		TAHMİN		TOPLAM
		C ⁺	C ⁻	
GERÇEK	C ⁺	0	5	5
	C ⁻	0	1000	1000
TOPLAM		0	1005	1005

- Accuracy = $(TP + TN) / (N) = (0 + 1000) / 1005 = 0,995$
- Hata = $(FP + FN) / (N) = (0 + 5) / 1005 = 0,005$ (Çok daha düşük hata !!!)

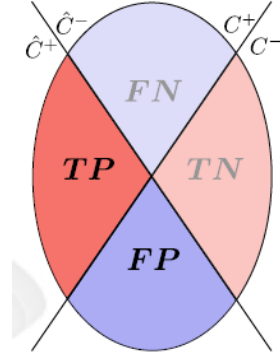
Konular

- Sınıflandırıcıların Değerlendirilmesi
 - Skorlar
 - Karışıklık matrisi
 - Accuracy
 - Precision
 - Recall
 - Specificity
 - F-Score
 - Eğitim ve Test Kümeleri
 - Hold-out set
 - Multiple random sampling
 - k-fold cross validation
 - Bootstrap

Precision

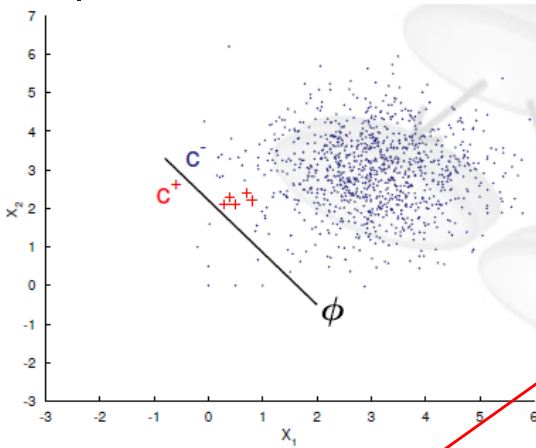
- Precision, gerçek değeri pozitif olup pozitif değere sınıflandırılan sayısının, pozitif değere sınıflandırılanların toplamına oranıdır.

$$\text{Precision} = \frac{TP}{TP + FP}$$



25

Precision



		TAHMİN		TOPLAM
		C ⁺	C ⁻	
GERÇEK	C ⁺	0	5	5
	C ⁻	7	993	1000
TOPLAM		7	998	1005

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{0}{0 + 7} = 0 \quad (\text{Çok kötü precision değeri !!!})$$

26

Konular

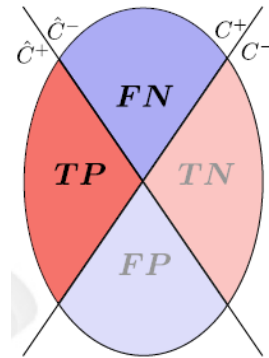
- Sınıflandırıcıların Değerlendirilmesi
- Skorlar
 - Karışıklık matrisi
 - Accuracy
 - Precision
 - **Recall**
 - Specificity
 - F-Score
- Eğitim ve Test Kümeleri
 - Hold-out set
 - Multiple random sampling
 - k-fold cross validation
 - Bootstrap

27

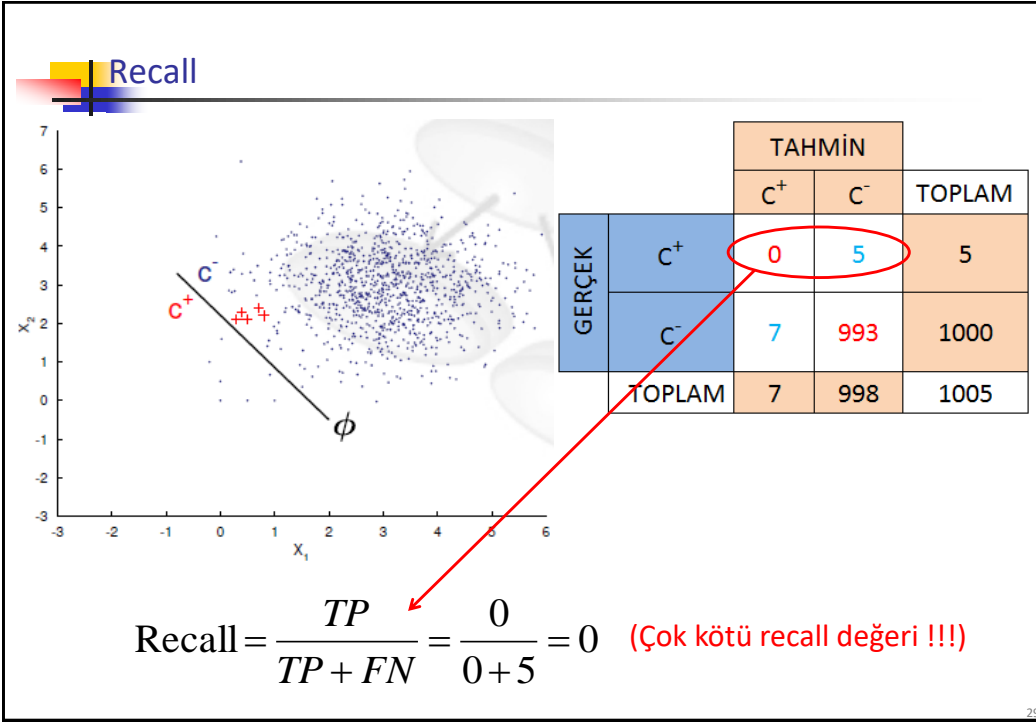
Recall

- Recall, gerçek değeri pozitif olup pozitif değere sınıflandırılan sayısının, gerçek değeri pozitif olanların tümüne oranıdır.

$$\text{Recall} = \frac{TP}{TP + FN}$$



28



Precision ve Recall Karşılaştırma

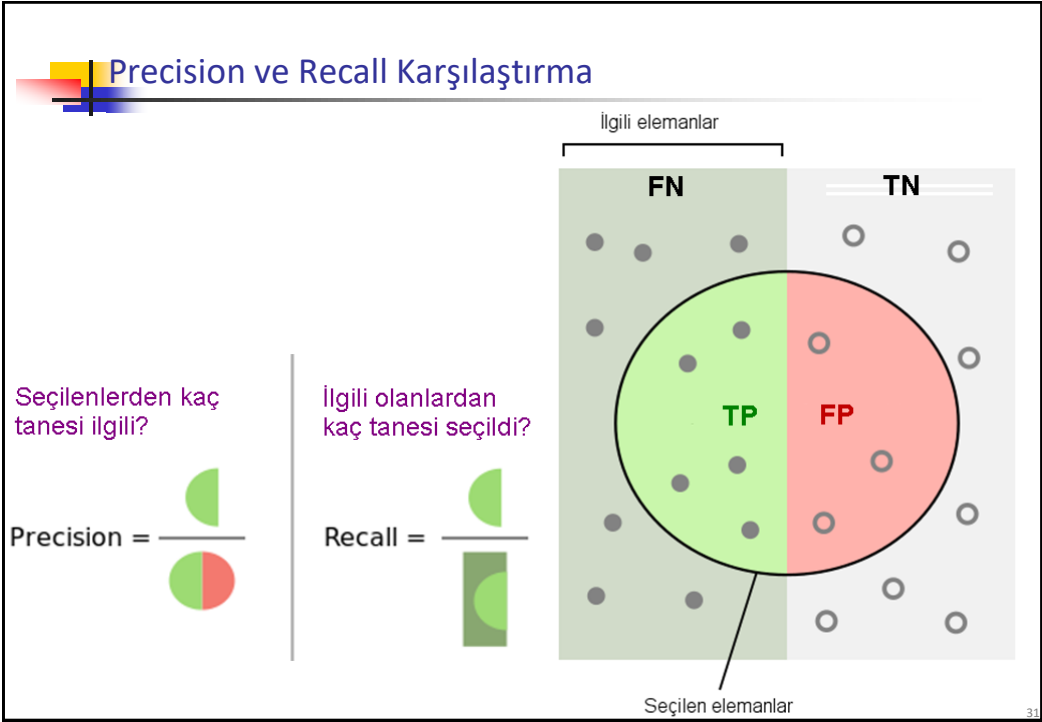
Spam filtreleme

- Precision: gerçekte spam olup da, spam kutusuna alınanların **spam kutusundaki tüm mesajlara oranıdır.**
- Recall: gerçekte spam olup ta spam kutusuna alınan mesajların, **gerçekte spam olan tüm mesajlara oranıdır.**

Duygu analizi

- Precision: gerçekte pozitif olup da pozitif sınıflandırılanların sayısının **pozitif olarak sınıflandırılanların tümüne oranıdır.**
- Recall: gerçekte pozitif olup da pozitif sınıflandırılanların **gerçekte pozitif olanların tümüne oranıdır.**

$$\text{Precision} = \frac{TP}{TP + FP} \qquad \text{Recall} = \frac{TP}{TP + FN}$$

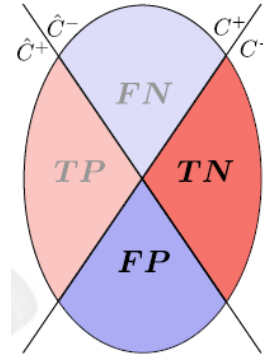


- ## Konular
- Sınıflandırıcıların Değerlendirilmesi
 - Skorlar
 - Karışıklık matrisi
 - Accuracy
 - Precision
 - Recall
 - **Specificity**
 - F-Score
 - Eğitim ve Test Kümeleri
 - Hold-out set
 - Multiple random sampling
 - k-fold cross validation
 - Bootstrap
- 32

Specificity

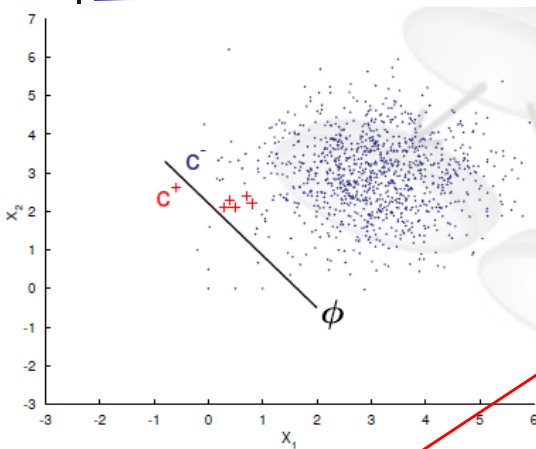
- Specificity, gerçek değeri negatif olup negatif sınıflandırılan sayısının, gerçek değeri negatif olanların tümüne oranıdır.

$$\text{Specificity} = \frac{TN}{TN + FP}$$



33

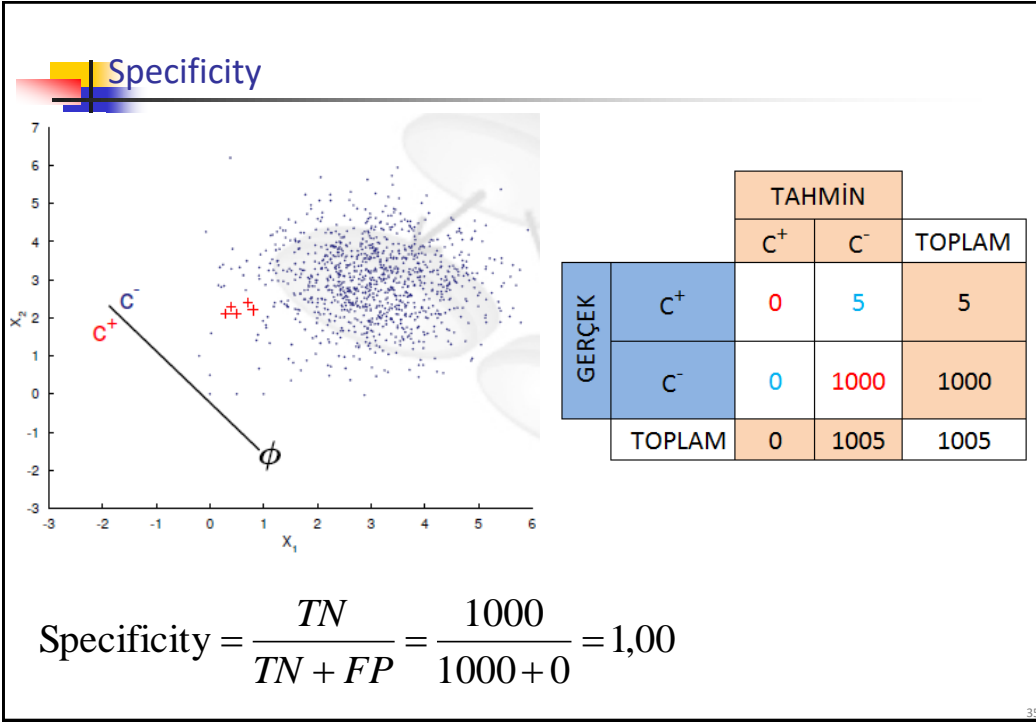
Specificity



		TAHMİN		TOPLAM
		C ⁺	C ⁻	
GERÇEK	C ⁺	0	5	5
	C ⁻	7	993	1000
TOPLAM		7	998	1005

$$\text{Specificity} = \frac{TN}{TN + FP} = \frac{993}{993 + 7} = 0,99$$

34



- ## Konular
- Sınıflandırıcıların Değerlendirilmesi
 - Skorlar
 - Karışıklık matrisi
 - Accuracy
 - Precision
 - Recall
 - Specificity
 - **F-Score**
 - Eğitim ve Test Kümeleri
 - Hold-out set
 - Multiple random sampling
 - k-fold cross validation
 - Bootstrap

F-Score

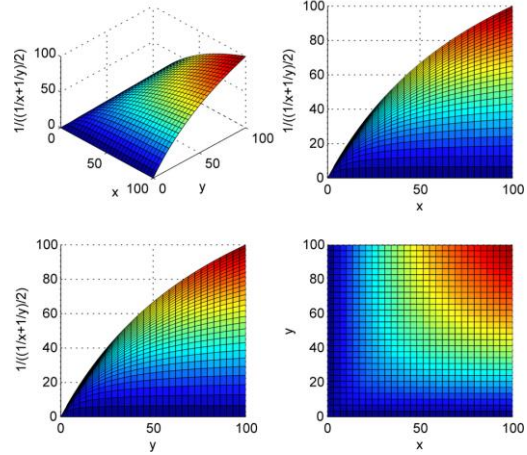
- F₁-score (**Harmonic mean**), iki sınıflandırıcının tek ölçüt ile değerlendirilmesi için kullanılır.

$$F\text{-score} = \frac{(\beta^2 + 1)\text{Precision} \cdot \text{Recall}}{\beta^2 (\text{Precision} + \text{Recall})}$$

$$F_1\text{-score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$x_{\text{harmonic mean}} = \frac{n}{\sum_{k=1}^n \frac{1}{x_k}}$$

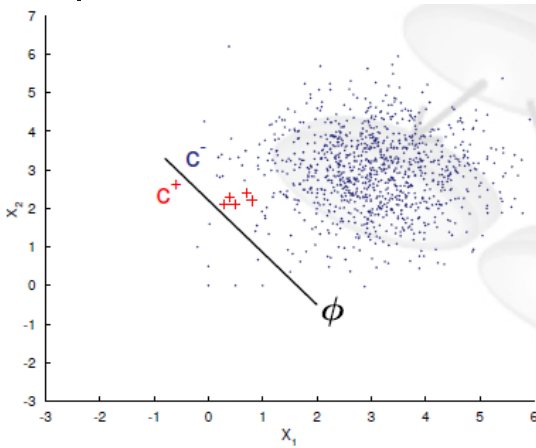
$$\begin{aligned} \text{Harmonic mean} &= \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}} \\ &= \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \end{aligned}$$



x ve y için f₁-score değişimi

37

F-Score



		TAHMİN		TOPLAM
		C ⁺	C ⁻	
GERÇEK	C ⁺	0	5	5
	C ⁻	7	993	1000
TOPLAM		7	998	1005

$$F_1\text{-score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2 \cdot 0 \cdot 0}{0 + 0} = 0$$

38

F-Score

		TAHMİN		TOPLAM
		C ⁺	C ⁻	
GERÇEK	C ⁺	10	0	10
	C ⁻	0	5	5
TOPLAM		10	5	15

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{10}{10 + 0} = 1$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{10}{10 + 0} = 1$$

$$F_1 \text{-score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2 \cdot 1 \cdot 1}{1 + 1} = 1$$

39

F-Score

		TAHMİN		TOPLAM
		C ⁺	C ⁻	
GERÇEK	C ⁺	0	10	10
	C ⁻	5	0	5
TOPLAM		5	10	15

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{0}{0 + 5} = 0$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{0}{0 + 10} = 0$$

$$F_1 \text{-score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2 \cdot 0 \cdot 0}{0 + 0} = 0$$

40

F-Score

		TAHMİN		TOPLAM
		C ⁺	C ⁻	
GERÇEK	C ⁺	5	5	10
	C ⁻	5	5	10
TOPLAM		10	10	20

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{5}{5+5} = 0,5$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{5}{5+5} = 0,5$$

$$F_1 \text{ - score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2 \cdot (0,5) \cdot (0,5)}{0,5 + 0,5} = 0,5$$

41

F-Score

		TAHMİN		TOPLAM
		C ⁺	C ⁻	
GERÇEK	C ⁺	5	0	5
	C ⁻	15	0	15
TOPLAM		20	0	20

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{5}{5+15} = 0,25$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{5}{5+0} = 1$$

$$F_1 \text{ - score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2 \cdot (0,25) \cdot 1}{0,25 + 1} = 0,4$$

42

Konular

- Sınıflandırıcıların Değerlendirilmesi
- Skorlar
 - Karışıklık matrisi
 - Accuracy
 - Precision
 - Recall
 - Specificity
 - F-Score
- Eğitim ve Test Kümeleri
 - Hold-out set
 - Multiple random sampling
 - k-fold cross validation
 - Bootstrap

43

Eğitim ve Test Kümeleri

- Bir sınıflandırıcının eğitimi ve testi için sınırlı sayıda veri elde edilebilir.
- **Eğitim verisi arttıkça daha iyi öğrenme ve genelleme sağlanabilmektedir.**
- **Test verisi arttıkça sınıflandırıcının hata olasılığı daha iyi tahmin edilebilmektedir.**
- Eğitim ve test verileri birbirinden farklı oluşturulmalıdır.
- Bir eğitim döngüsü içerisinde aynı değerlere sahip eğitim ve test veri kümesi kullanılmamalıdır.

44

Konular

- Sınıflandırıcıların Değerlendirilmesi
- Skorlar
 - Karışıklık matrisi
 - Accuracy
 - Precision
 - Recall
 - Specificity
 - F-Score
- Eğitim ve Test Kümeleri
 - **Hold-out set**
 - Multiple random sampling
 - k-fold cross validation
 - Bootstrap

45

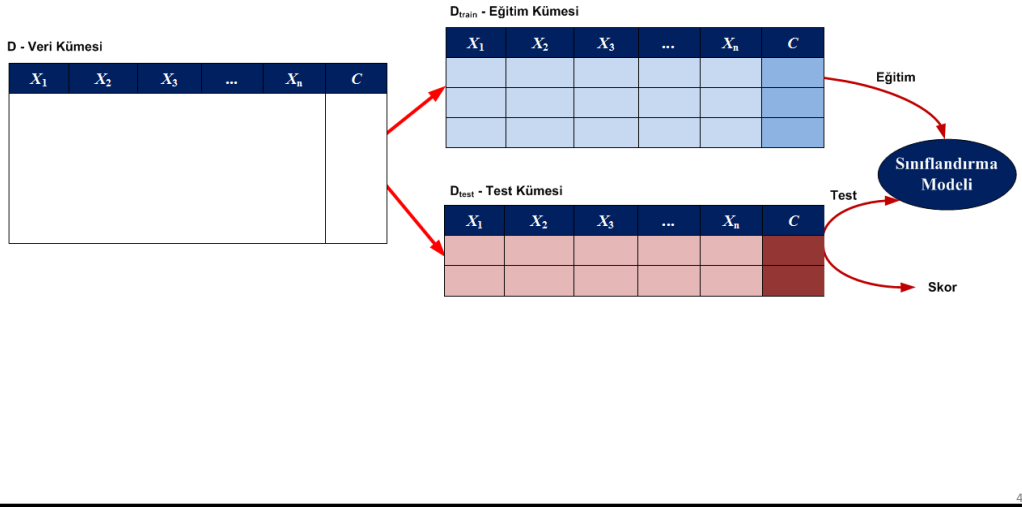
Hold-out set

- Kullanılabilir veri kümesi D, **eğitim (D_{train})** ve **test (D_{test})** olmak üzere iki ayrı kümeye ayrılır.
$$D = D_{\text{train}} \cup D_{\text{test}}, D_{\text{train}} \cap D_{\text{test}} = \emptyset$$
- Oluşturulan test kümesi **holdout set** olarak adlandırılır.
- Bu yöntem veri kümesi D büyükse kullanılabilir.
- **D veri kümesindeki tüm veriler bir sınıf etiketine atanmıştır.**
- **Training set** sınıflandırıcının **eğitimi için**, **test set** ise sınıflandırıcının **değerlendirilmesi için kullanılır.**
- **Training set için yüksek doğruluk** oranına sahip olan sınıflandırıcı, **test set için düşük doğruluk** düzeyine sahipse **overfit** yapılmıştır.
- İki küme için %50'şer veya 2/3 train ve 1/3 test için alınabilir.

46

Hold-out set

- Kullanılabilir veri kümesi D , eğitim (D_{train}) ve test (D_{test}) olmak üzere iki ayrı kümeye ayrılır.



Konular

- Sınıflandırıcıların Değerlendirilmesi
- Skorlar
 - Karışıklık matrisi
 - Accuracy
 - Precision
 - Recall
 - Specificity
 - F-Score
- Eğitim ve Test Kümeleri
 - Hold-out set
 - Multiple random sampling
 - k-fold cross validation
 - Bootstrap

Multiple random sampling

- Kullanılabilir veri kümesi D çok küçük boyutta ise, test kümesi çok daha küçük boyutta olacağından güvenilir sonuç vermez.
- Veri kümesinin küçük boyutta olduğu bu gibi durumlarda **n kez rastgele örnekleme** ile **eğitim** ve **test kümesi** oluşturulur.
- Bu durumda n tane doğruluk değeri elde edilir.
- **Sonuç doğruluk değeri**, elde edilen **doğruluk değerlerinin ortalaması alınarak hesaplanır**.

49

Konular

- Sınıflandırıcıların Değerlendirilmesi
 - Skorlar
 - Karışıklık matrisi
 - Accuracy
 - Precision
 - Recall
 - Specificity
 - F-Score
 - Eğitim ve Test Kümeleri
 - Hold-out set
 - Multiple random sampling
 - **k-fold cross validation**
 - Bootstrap

50

k-fold cross validation

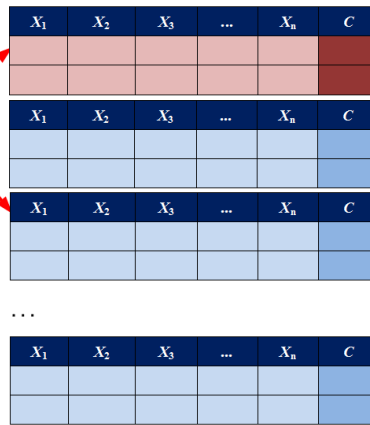
- **Veri kümesi D küçük boyutta** olduğunda sık kullanılan yöntemdir.
- Bu yöntemde veri kümesi D, k adet eşit boyutta disjoint alt kümeye bölünür.
- Her (k-1) küme eğitim için kullanılırken kalan bir küme test için kullanılır.
- Bu işlem k kez tekrarlanır ve k adet doğruluk değeri elde edilir.
- Sonuç doğruluk değeri tüm doğruluk değerlerinin ortalaması alınarak hesaplanır.
- 5-fold ve 10-fold cross-validation literatürde farklı uygulamalarda yaygın kullanılmaktadır.
- **Hold-out set yöntemine göre daha iyi sonuç vermektedir.**

51

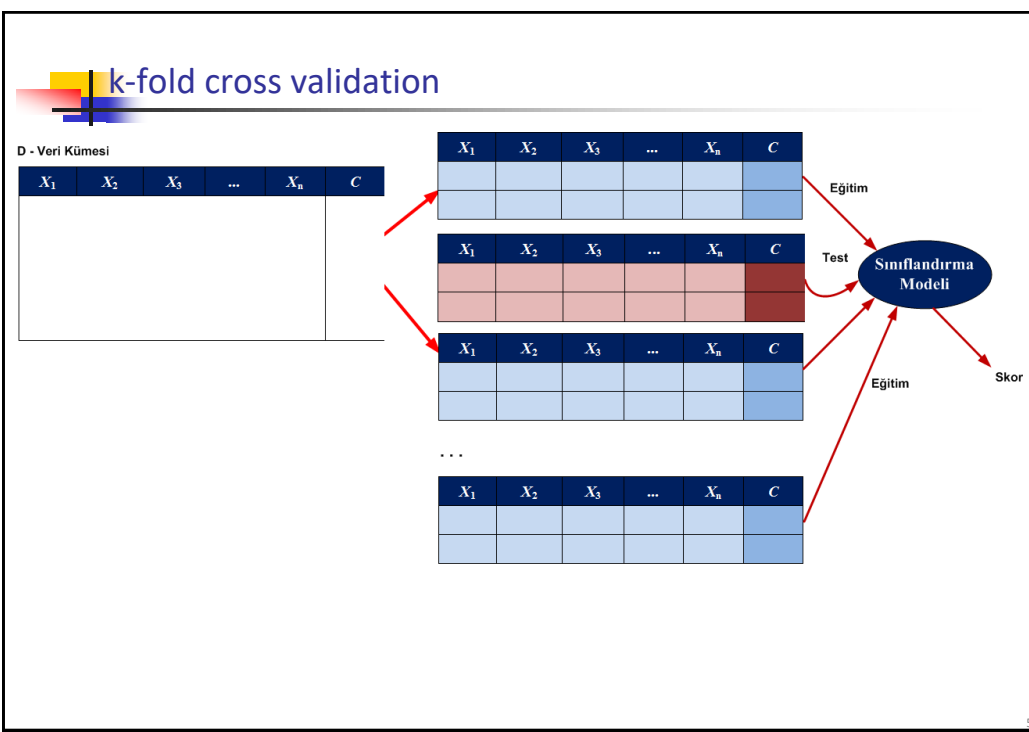
k-fold cross validation

D - Veri Kümesi

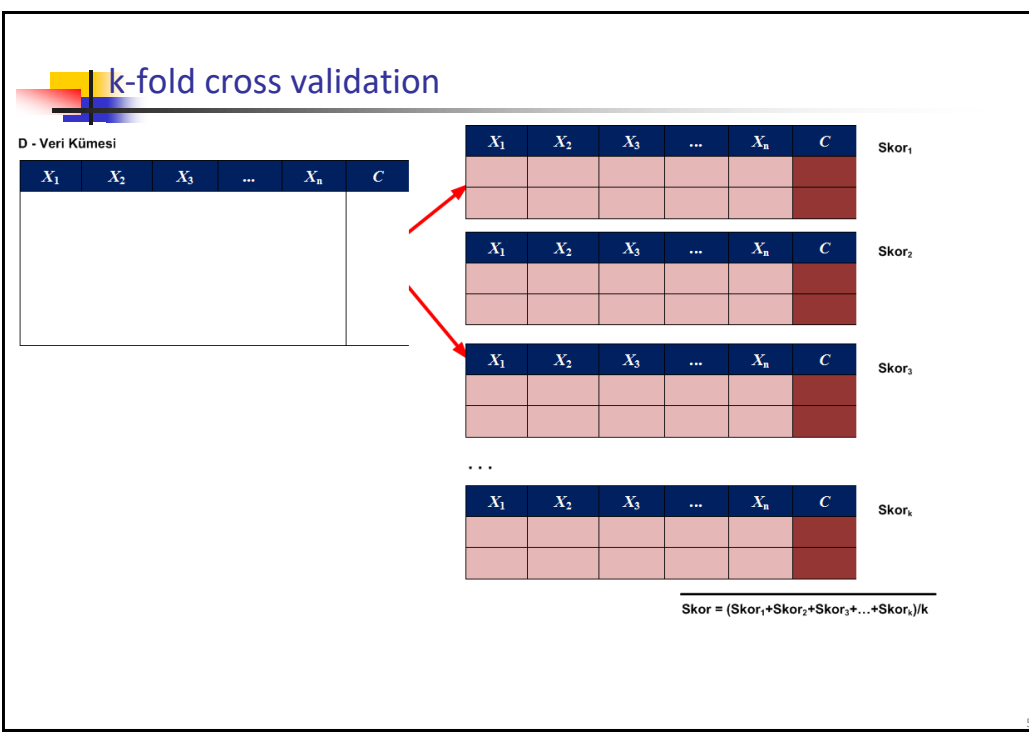
X_1	X_2	X_3	...	X_n	C



52



53



54

k-fold cross validation

Leaving-one-out cross validation

- k-fold cross validation yönteminin özel durumudur.
- **Eğitim kümesi, örnek sayısı kadar alt kümeye ($k=N$) bölünür.**
- Her iterasyonda **1 eleman dışarıda tutularak diğerlerinin tümü eğitim için kullanılır.**
- **Kalan 1 eleman ise test için kullanılır.**
- Sonuç doğruluk değeri tüm doğruluk değerlerinin toplamı alınarak hesaplanır.

55

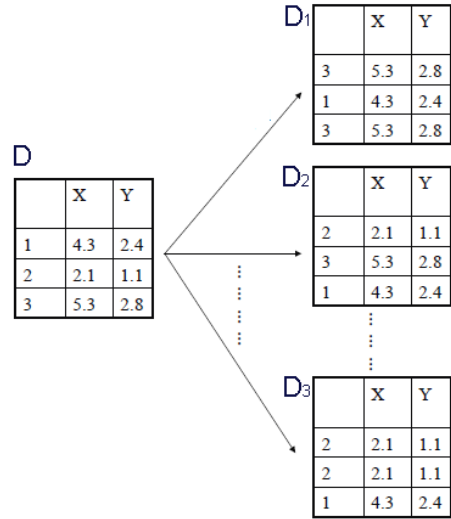
Konular

- Sınıflandırıcıların Değerlendirilmesi
- Skorlar
 - Karışıklık matrisi
 - Accuracy
 - Precision
 - Recall
 - Specificity
 - F-Score
- Eğitim ve Test Kümeleri
 - Hold-out set
 - Multiple random sampling
 - k-fold cross validation
 - **Bootstrap**

56

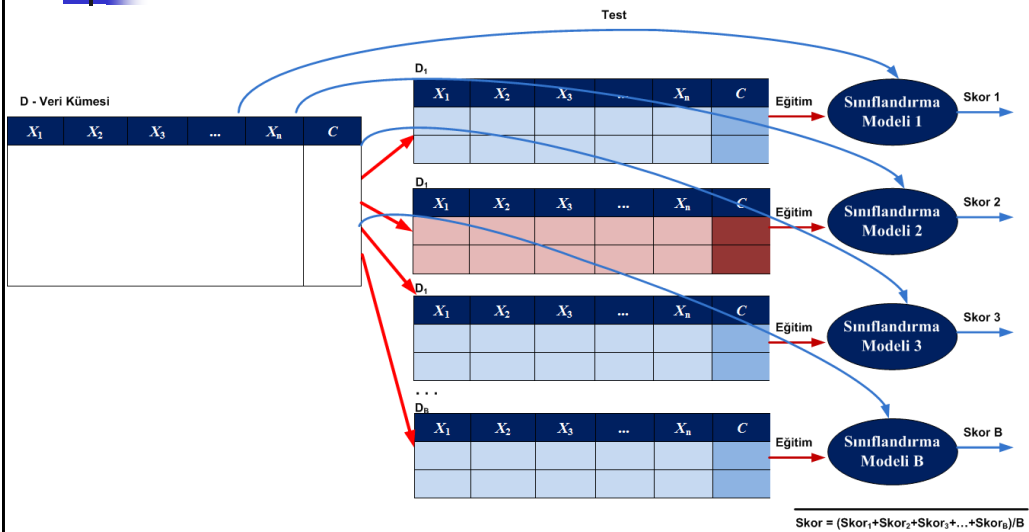
Bootstrap

- D veri kümesinden belirli sayıda değer rastgele seçilerek bir altküme elde edilir.
- Elde edilen her altküme için ayrı model oluşturulur ve her birisinin doğruluk değeri ayrı ele alınır.
- **Sonuç doğruluk değeri ise tüm modellerden elde edilen doğruluk değerlerinin ortalaması alınarak hesaplanır.**
- Test ve eğitim verileri tekrarlı olabilir.
- Overfit olma olasılığı vardır.



57

Bootstrap



58

Ödev

- Veri artırma (data augmentation) hakkında bir araştırma ödevi hazırlayınız.