



**BATCH** : **B 150** Data Science  
**LESSON** : **Machine Learning**  
**DATE** : 15.07.2023  
**SUBJECT** : **Supervised Learning  
Classification - KNN**



techproeducation



techproeducation



techproeducation



techproeducation



techproedu



techproeducation.com



info@techproeducation.com



+1 (917) 768-7466



## MACHINE LEARNING - 3



Makine Öğrenmesi – 3



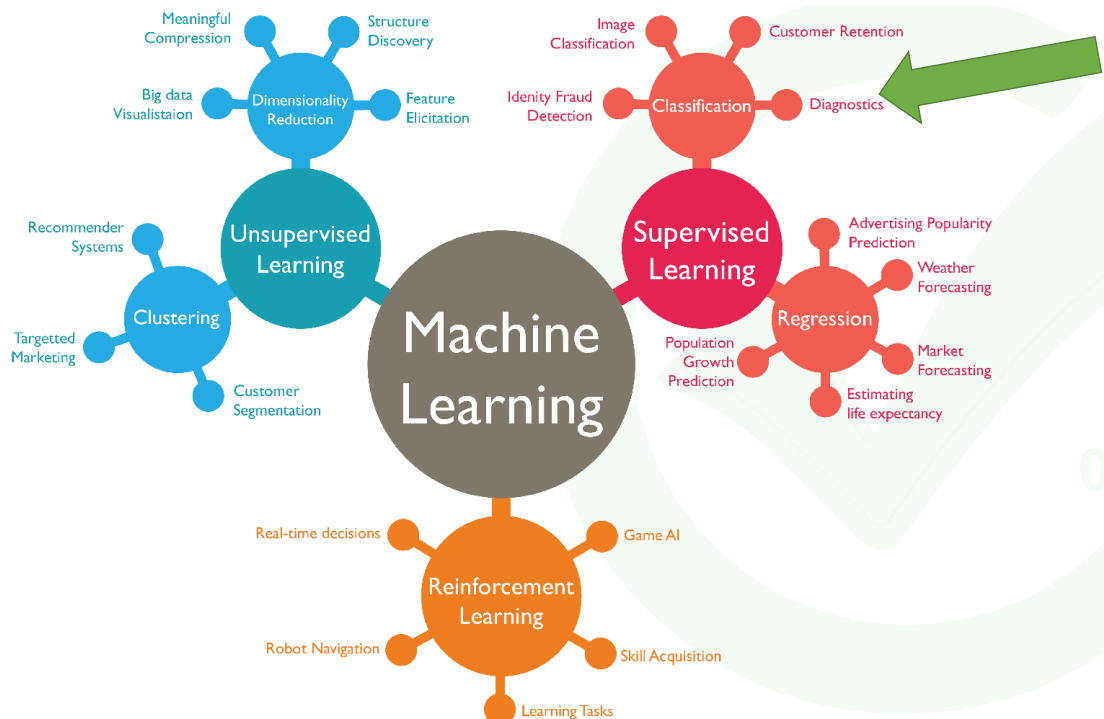
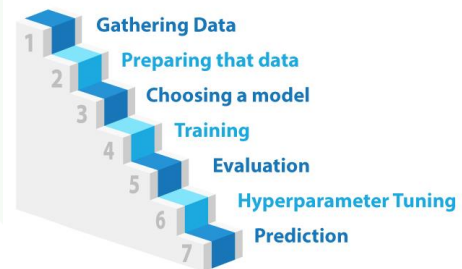
# Overall Table of Contents

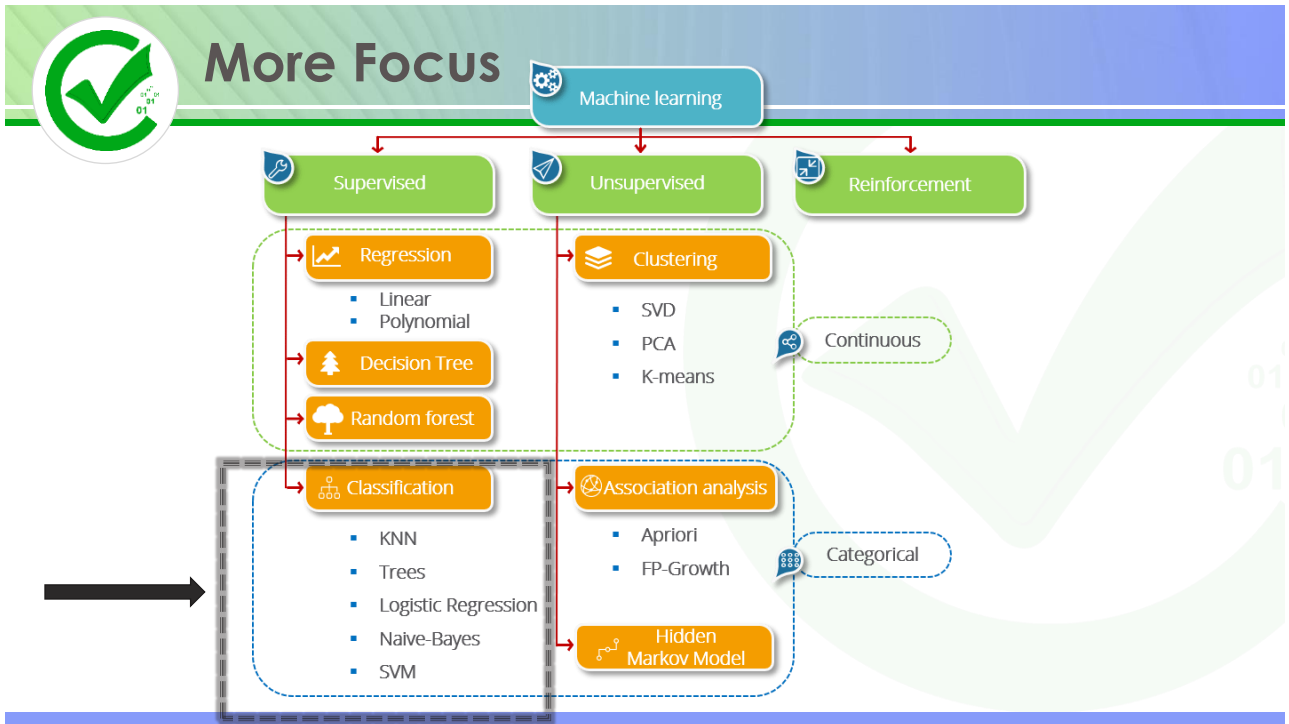


## General Content

- ✓ Supervised Learning Algorithm - **Classification**
- ✓ Supervised Algorithm practices Python application
- ✓ Projects Solutions

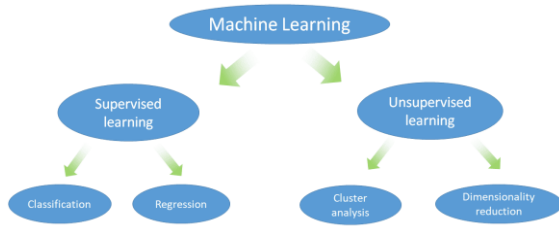
### 7 steps of Machine Learning



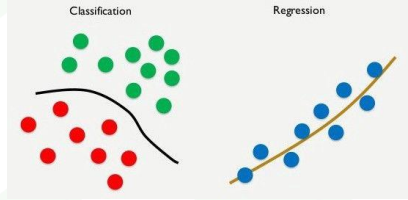


# SUPERVISED LEARNING - Classification

*Sınıflandırma*



## CLASSIFICATION vs REGRESSION



$$P_c^{NN}(\mathbf{W}; \mathbf{Z}) = e^{z_c^{NN}} / \sum_{a=0}^S e^{z_a^{NN}}$$

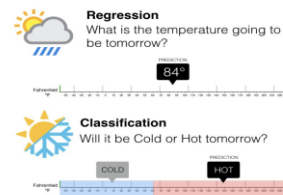
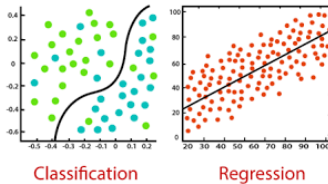
CLASSIFY



## Supervised Learning

### Classification 'a Giriş (Sınıflandırma)

- ✓ Regresyonda neler gördük
- ✓ Supervised Learning in 2. tipi olarak Classification
- ✓ Regression vs Classification
- ✓ Target için Kategorik Sınıflandırma vardır



The Idea of  
**Classification**  
and **Regression**



**!! Regresyonda** target hedef değişkenin **sayısal değerlerini**; **sınıflandırmada** ise target değişkenin ait olduğu **sınıfları (ya da "etiket")** tahmin eden modelleri oluşturmaya çalışırız.



## Supervised Learning

### Classification 'a Giriş

- ✓ Classification un hayattaki kullanım alanları
- ✓ Binary Classification

Fraud Protection



Spam Mail Detection



#### What is Classification



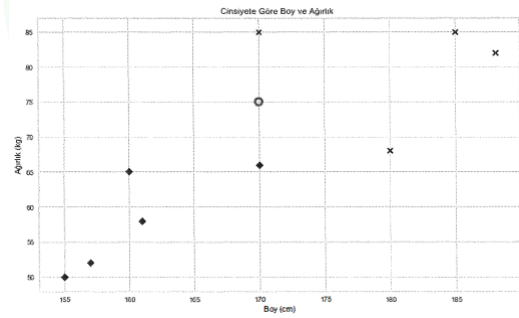
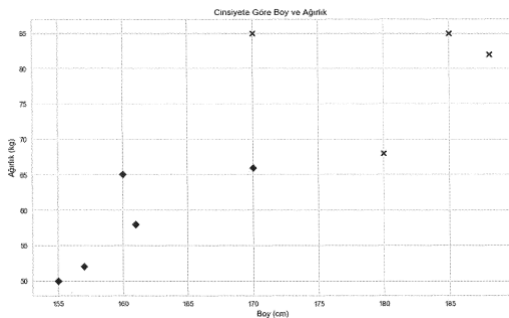
Boy	Kilo	Cinsiyet
160	65	K
170	85	E
185	85	E
188	82	E
155	50	K
161	58	K
180	68	E
157	52	K
170	66	K



## Supervised Learning

### Classification 'a Giriş

- ✓ Binary Classification
- ✓ Kategori tahmini
- ✓ «Yakınlık» kavramı - 'ara mesafe ölçümü'

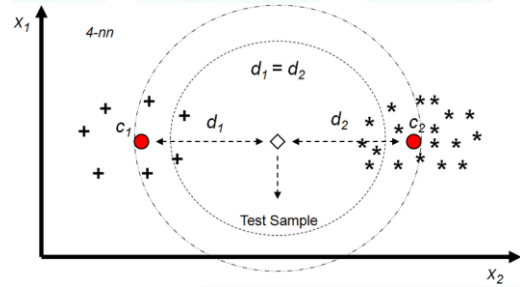
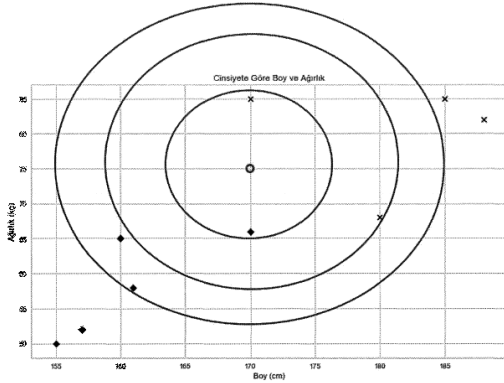




# Supervised Learning

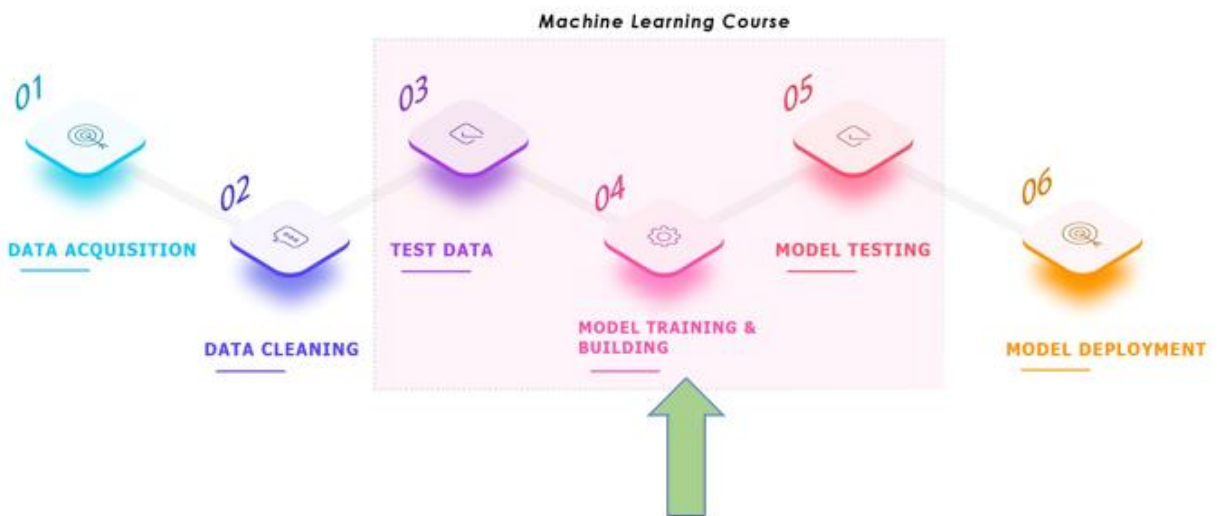
## Classification 'a Giriş

- «Yakınlık» kavramı - «ara mesafe ölçümü»
- En yakın komşu kavramı



## Where are we?

DATA SCIENCE





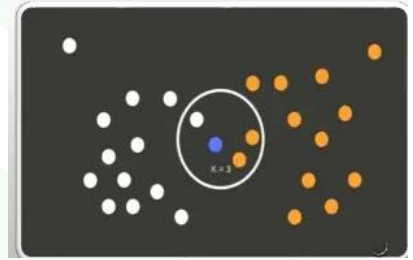
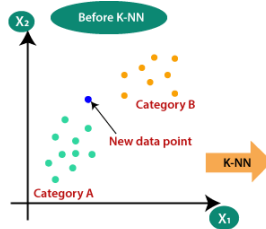
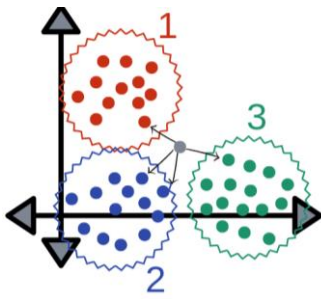


## Supervised Learning

### K Nearest Neighbour-KNN Algoritması

#### K-EN YAKIN KOMŞU Algoritması

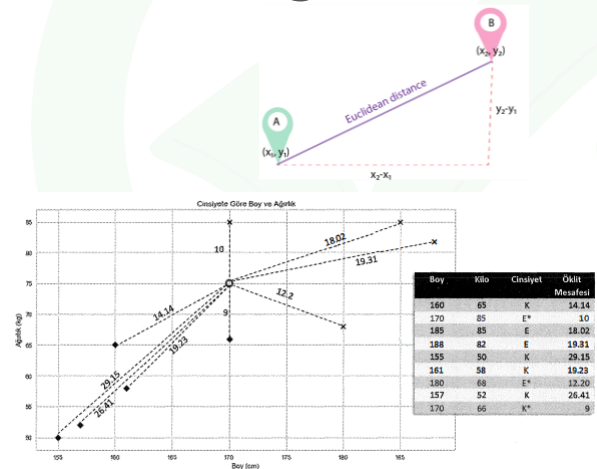
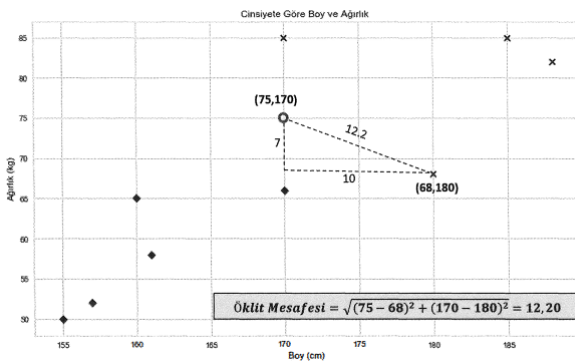
- ✓ Logistic Regression ' dan bir adım öncesinde
- ✓ Classification için en basit yol olarak KNN
- ✓ En yakın komşu sayısı



## Supervised Learning

### K Nearest Neighbour-KNN Algoritması

- ✓ Öklid mesafesi





# Supervised Learning

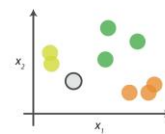
## K Nearest Neighbour-KNN Algoritması

- ✓ KNN avantajları
- ✓ KNN dezavantajları



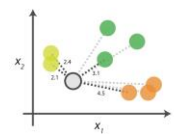
### kNN Algorithm

#### 0. Look at the data



Say you want to classify the grey point into a class. Here, there are three potential classes - lime green, green and orange.

#### 1. Calculate distances



Start by calculating the distances between the grey point and all other points.

#### 2. Find neighbours

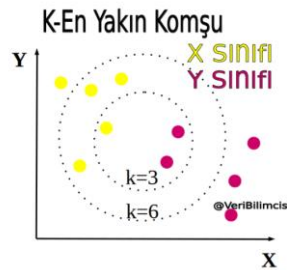
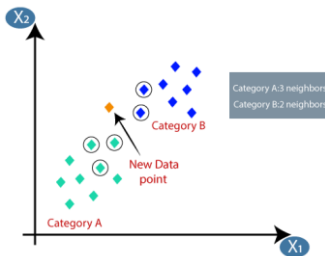
Point	Distance	Rank
●	2.1	1st NN
●	2.4	2nd NN
●	3.1	3rd NN
●	4.5	4th NN

Next, find the nearest neighbours by ranking points by increasing distance. The nearest neighbours (NNs) of the grey point are the ones closest in dataset.

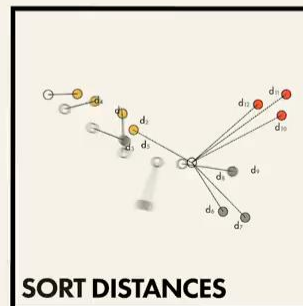
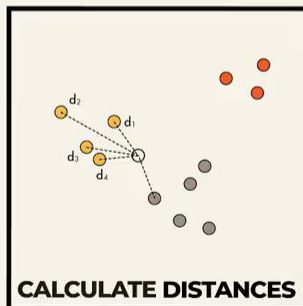
#### 3. Vote on labels

Class	# of votes	Result
●	2	Class ● wins the vote! Point ● is therefore predicted to be of class ●.
●	1	
●	1	

Vote on the predicted class labels based on the classes of the k nearest neighbours. Here, the labels were predicted based on the k=3 nearest neighbours.



# K NEAREST NEIGHBORS

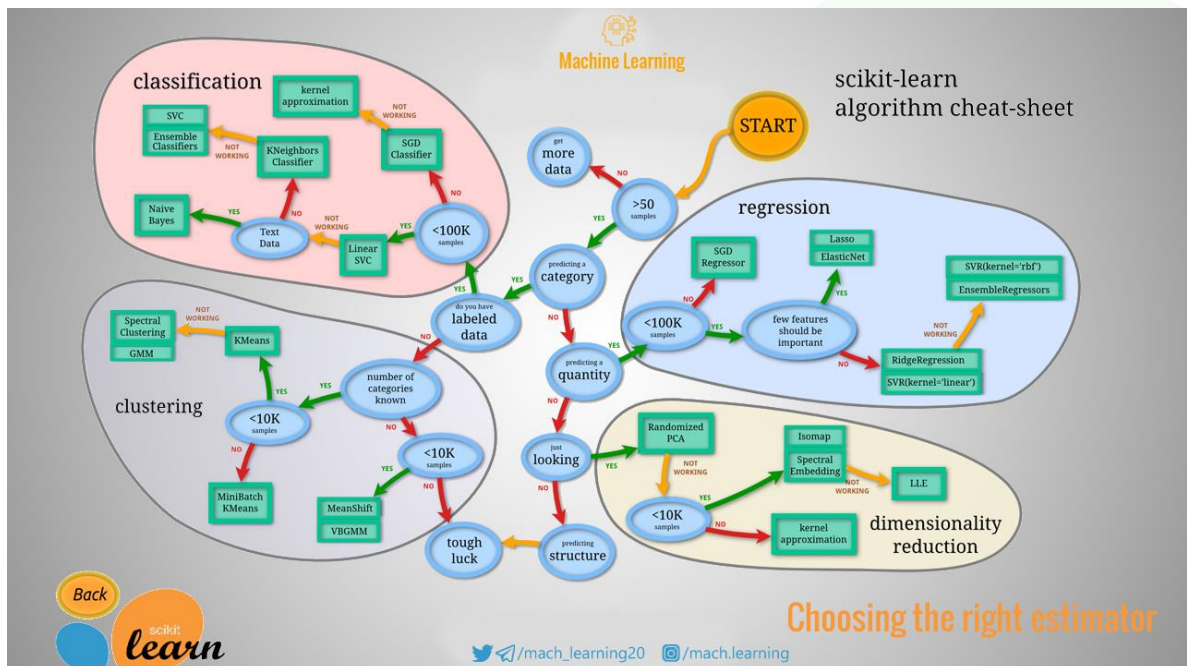
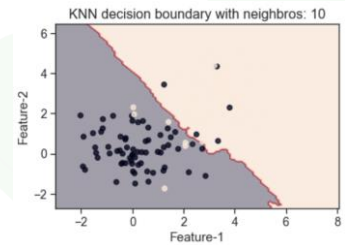
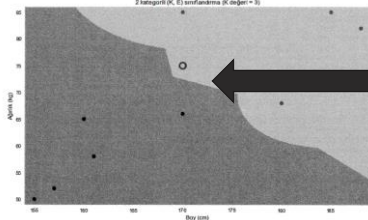
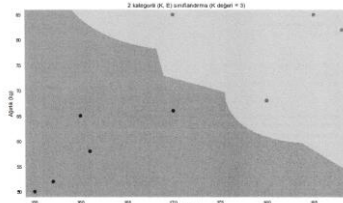
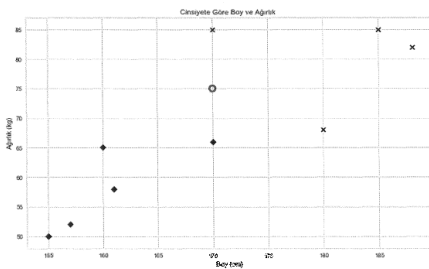






# K Nearest Neighbour-KNN Algorithması

Decision boundary kavramı



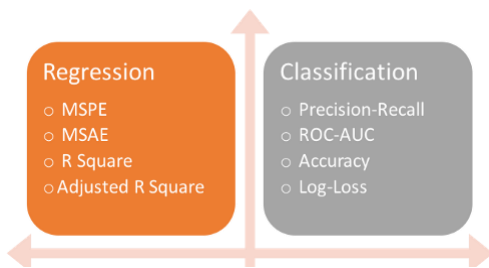


## Supervised Learning

### Evaluation Metrics for Classification Problems (Performans Ölçütleri)

*very Important*

- ✓ Regresyon kriterleri nasıl olurdu burada?
- ✓ Confusion metrics kavramı



		Predicted Class		
		Positive	Negative	
Actual Class	Positive	True Positive (TP)	False Negative (FN) Type II Error	<b>Sensitivity</b> $\frac{TP}{(TP + FN)}$
	Negative	False Positive (FP) Type I Error	True Negative (TN)	<b>Specificity</b> $\frac{TN}{(TN + FP)}$
		<b>Precision Value</b> $\frac{TP}{(TP + FP)}$	<b>Negative Predictive Value</b> $\frac{TN}{(TN + FN)}$	<b>Accuracy</b> $\frac{TP + TN}{(TP + TN + FP + FN)}$

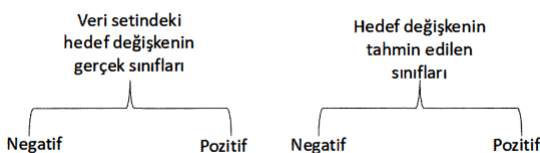


## Supervised Learning

### Evaluation Metrics for Classification Problems

- ✓ Confusion metrics kavramı: TN-FN-FP-TP

#### 2-SINIFLI (CLASS) SINIFLANDIRMA MODELİ İÇİN HATA MATRİSİ



#### 2-SINIFLI (CLASS) SINIFLANDIRMA MODELİ İÇİN HATA MATRİSİ

Gerçek Sınıflar	Tahmin Edilen Sınıflar	
	Negatif (0)	Pozitif (1)
Negatif (0)	DN	YP
Pozitif (1)	YN	DP





# Supervised Learning

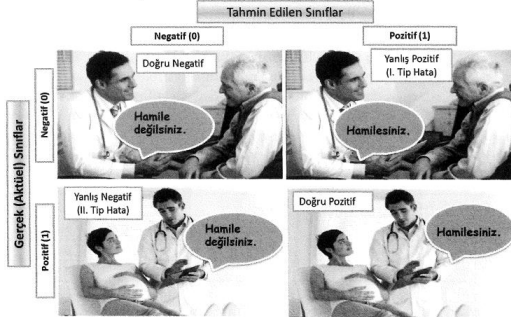
## Evaluation Metrics for Classification Problems



### Confusion metrics kavramı: TN-FN-FP-TP

1 291-5020 • +1 (814) 325-1026 • +1 (817) 771-0931 • +1 (848) 213-2921 • +1 (880) 997-9075 • +1 (800) 610-

### SINIFLANDIRMA MODELLERİ İÇİN PERFORMANS DEĞERLENDİRME ÖLÇÜTLERİ



Actual \ Predicted	Predicted	
	0	1
0	30	12
1	8	56

Gerçek Sınıflar

	Tahmin Edilen Sınıflar	
	Negatif (0)	Pozitif (1)
Negatif (0)	DN	YP
Pozitif (1)	YN	DP

**Doğruluk (Accuracy):** Doğru tahmin edilen hedef değişkenlerin tüm hedef değişkenlerine oranıdır. Model hedef değişkenleri ne kadar doğrulukla tahmin ediyor?

$$\text{Doğruluk (Accuracy)} = \frac{DP + DN}{DP + DN + YP + YN}$$

**Keskinlik (Precision):** Doğru pozitif olarak tahmin edilen gözlemlerin tüm pozitif gözlemlere oranıdır. Doğru bir şekilde pozitif tahmin edilen gözlemlerin gerçekte ne kadarı doğrudur?

$$\text{Keskinlik (Precision)} = \frac{DP}{DP + YP}$$

**Duyarlılık (Recall):** Doğru bir şekilde pozitif olarak tahmin edilen gözlemlerin ne kadar başarılı tahmin edildiğini gösterir.

$$\text{Duyarlılık (Recall)} = \frac{DP}{DP + YN}$$

**F1 Skoru:** Keskinlik (precision) ve Duyarlılık (Recall) harmonik ortalamasıdır.

$$F1 = 2 * \frac{(\text{Keskinlik} * \text{Duyarlılık})}{(\text{Keskinlik} + \text{Duyarlılık})}$$

### Performance metrics associated with Class 1

Predicted Labels	Actual Labels	
	1	0
1	True Positive	False Positive
0	False Negative	True Negative

(Is your prediction correct?) (What did you predict)  
True Negative (You predicted 0)  
True Positive (Your prediction is correct)

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

$$\text{False +ve rate} = \frac{FP}{TN + FP}$$

$$\text{F1 score} = 2x \frac{(\text{Prec} \times \text{Rec})}{(\text{Prec} + \text{Rec})}$$

$$\text{Accuracy} = \frac{TP + TN}{TP + FN + FP + TN}$$

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

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

Actual class	Predicted class	
	+	-
+	TP True Positives	FN False Negatives Type II error
-	FP False Positives Type I error	TN True Negatives

In Python  

```

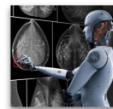
15] confusion_matrix(y, y_pred)
16] array([[448, 52],
        [121, 147]], dtype=int64)

```

Actual Cancer: 1 1 0 0 1 0 0 0 0 1 ...  
 Predicted Cancer: 0 1 1 0 1 0 0 0 0 0 ...  
 X ✓ X ✓ ✓ ✓ ✓ ✓ ✓ X  
 FN TP FP TN TP TN TN TN FN ...

### Why is Accuracy not a good metric?

#### Cancer Detection Example:



Actual Cancer: 1 1 0 0 1 0 0 0 0 1 ...  
 Predicted Cancer: 0 1 1 0 1 0 0 0 0 0 ...  
 X ✓ X ✓ ✓ ✓ ✓ ✓ ✓ X

All correctly predicted values (7) / All predicted values (10) X 100 → **Accuracy = 70 %**

All correctly predicted values (60) / All predicted values (63) X 100 → **Accuracy = 95 % (PERFECT ????)**

Accuracy is very high, but missed 2 actual patient.



# Supervised Learning

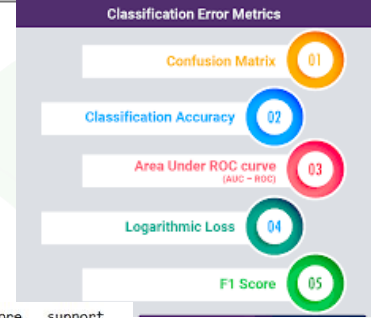
## Evaluation Metrics for Classification Problems

✔ Confusion metrics kavramı: TN-FN-FP-TP

Gerçek/Actual Veri				Tahmin		Hata Türü
Boy	Kilo	Cinsiyet	Etiket	Tahmin Edilen Cinsiyet	Etiket	
170	75	E	0	E	0	DN
180	95	E	0	E	0	DN
160	50	K	1	K	1	DP
165	62	K	1	K	1	DP
167	88	K	1	E	0	YN

Hata Matrisi		Tahmin Edilen Etiketler	
Negatif (0)		Pozitif (1)	
Gerçek/Aktüel	Negatif (0)	DN=2	YP=0
Veri	Pozitif (1)	YN=1	DP=2

	precision	recall	f1-score	support
E	0.67	1.00	0.80	2
K	1.00	0.67	0.80	3
micro avg	0.80	0.80	0.80	5
macro avg	0.83	0.83	0.80	5
weighted avg	0.87	0.80	0.80	5

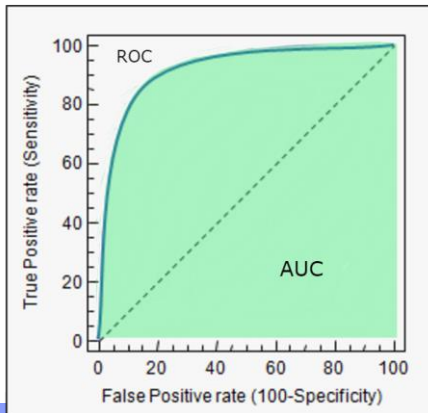




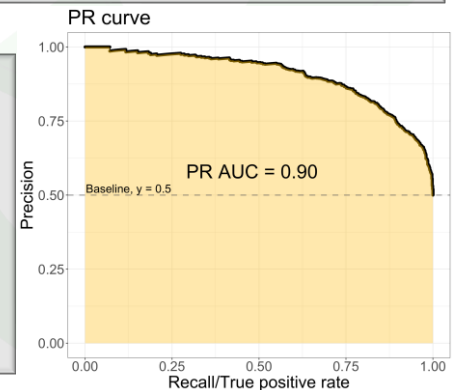
# Supervised Learning

## Evaluation Metrics for Classification Problems

- ✓ ROC Curve (Receiver Operating Characteristics)
- ✓ AUC Area (Area Under the Curve)
- ✓ (ROC Eğrisi ve AUC Alanı)



AUC'un olabildiğinde yüksek olması (1'e yakın olmasını) istiyoruz  
AUC ne kadar yüksekse model o kadar negatif (0) durumları negatif; pozitif (1) durumları da pozitif öngörüyor demektir.

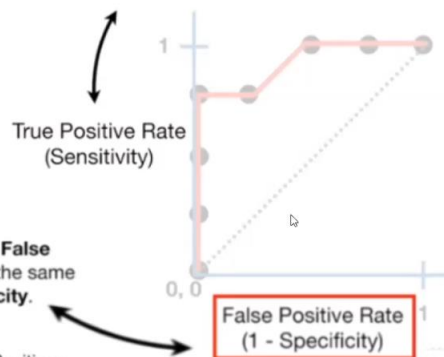


$$\text{True Positive Rate} = \text{Sensitivity} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

Receiver  
Operator  
Characteristics

Area  
Under  
Curve

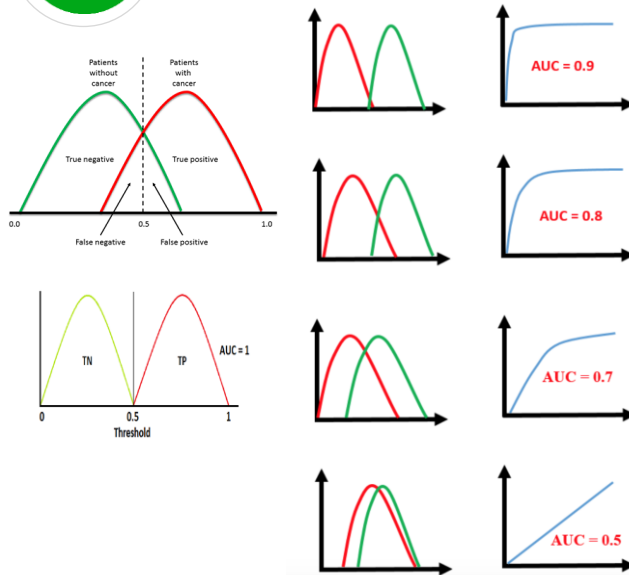
The x-axis shows the **False Positive Rate**, which is the same thing as **1 - Specificity**.



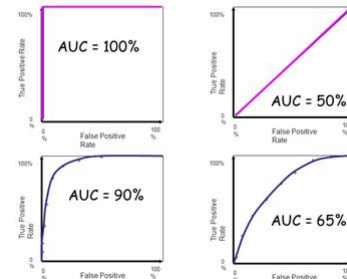
$$\text{False Positive Rate} = (1 - \text{Specificity}) = \frac{\text{False Positives}}{\text{False Positives} + \text{True Negatives}}$$

Picture credit : StatQuest





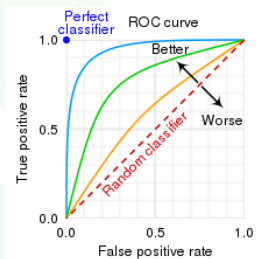
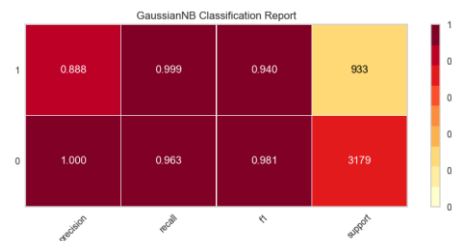
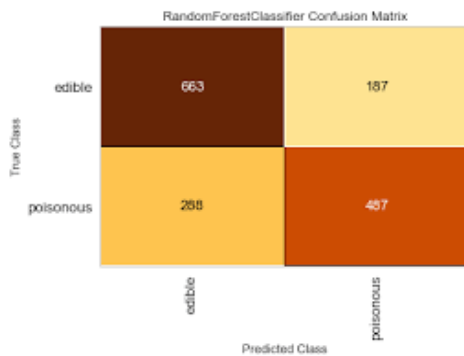
AUC for ROC curves



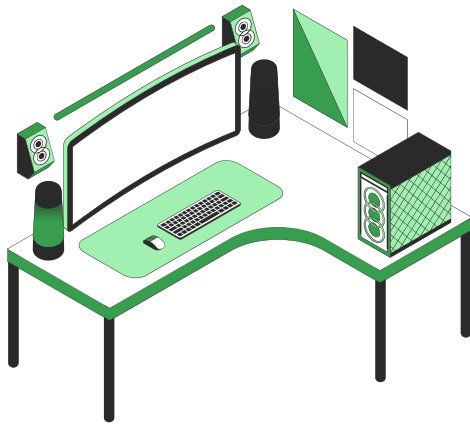
## Supervised Learning

### Evaluation Metrics for Classification Problems

✓ Yellowbrick ile confusion matrices







Do you  
have any  
questions?

Send it to us! We hope you learned  
something new.