



BATCH : **B 84 Data Science**
LESSON : **Machine Learning**
DATE : **12.09.2022**
SUBJECT : **Supervised Learning**



techproeducation



techproeducation



techproeducation



techproeducation



techproedu



MACHINE LEARNING - 3



Makine Öğrenmesi – 3



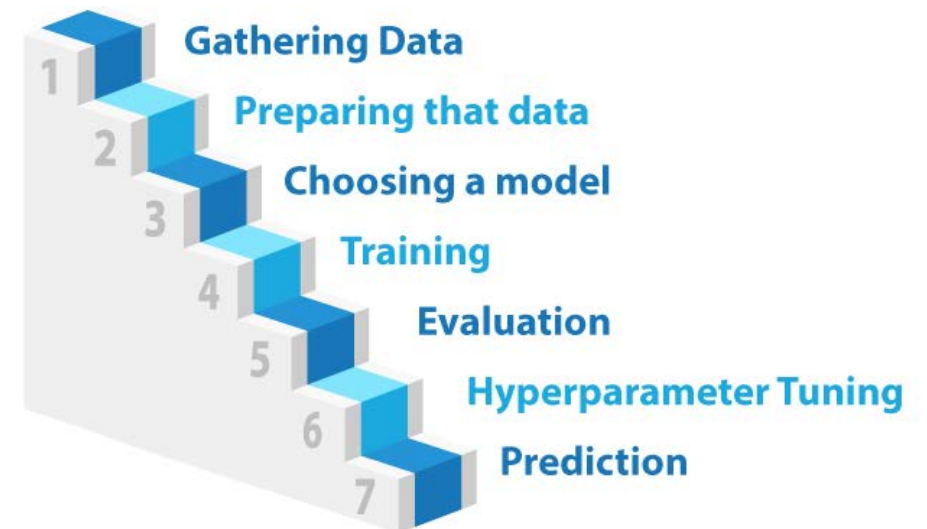
Overall Table of Contents

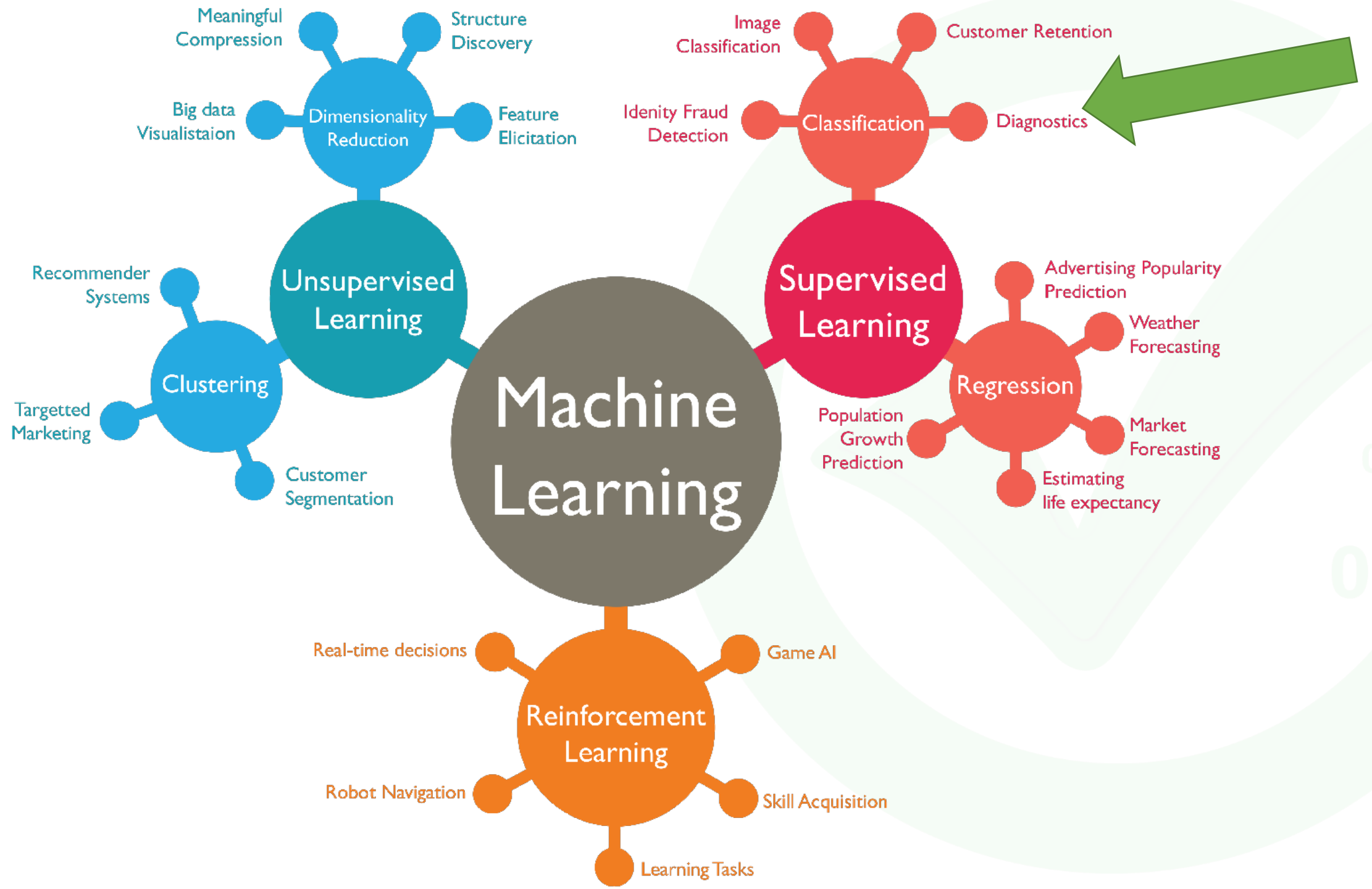


General Content

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

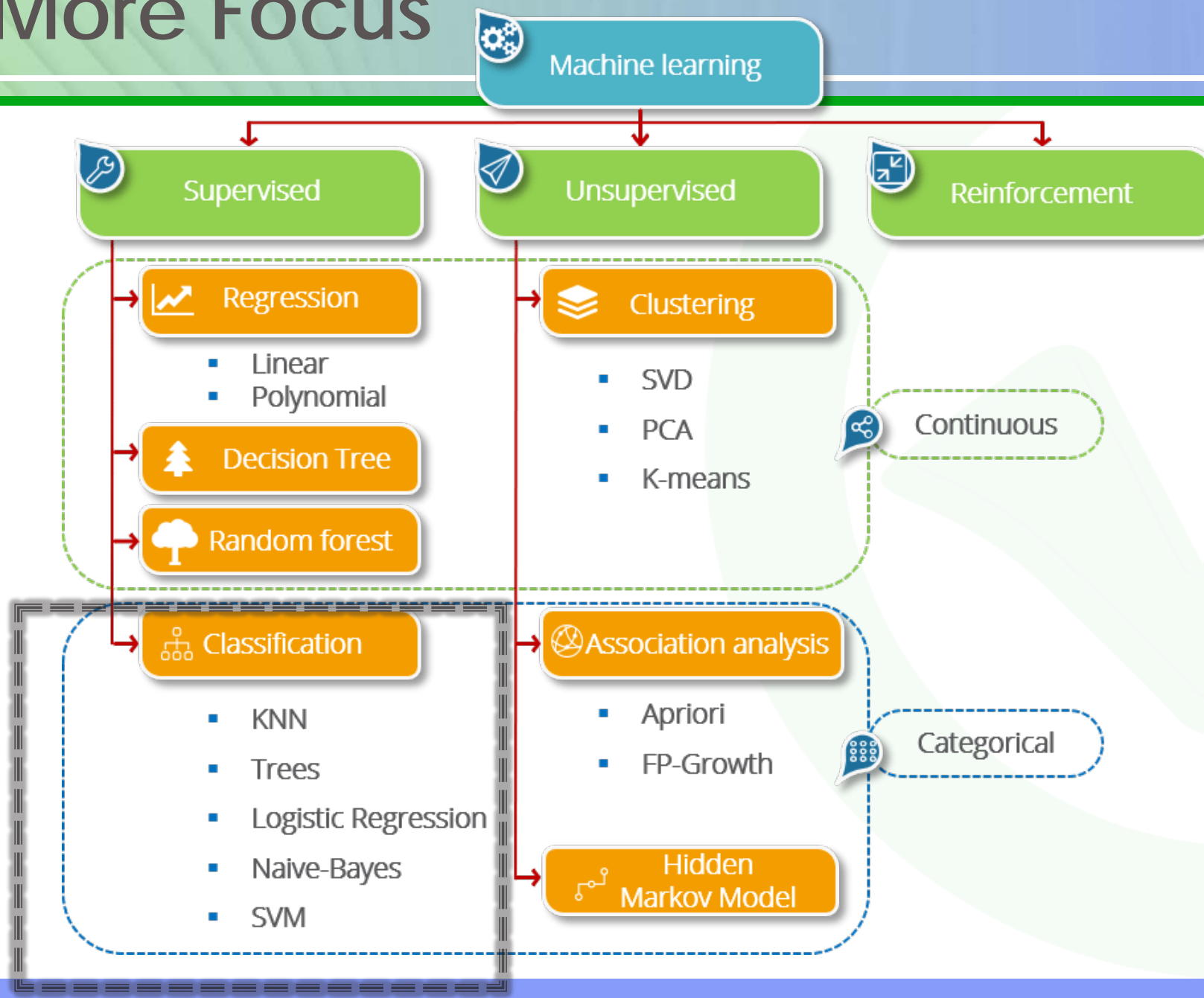
7 steps of Machine Learning







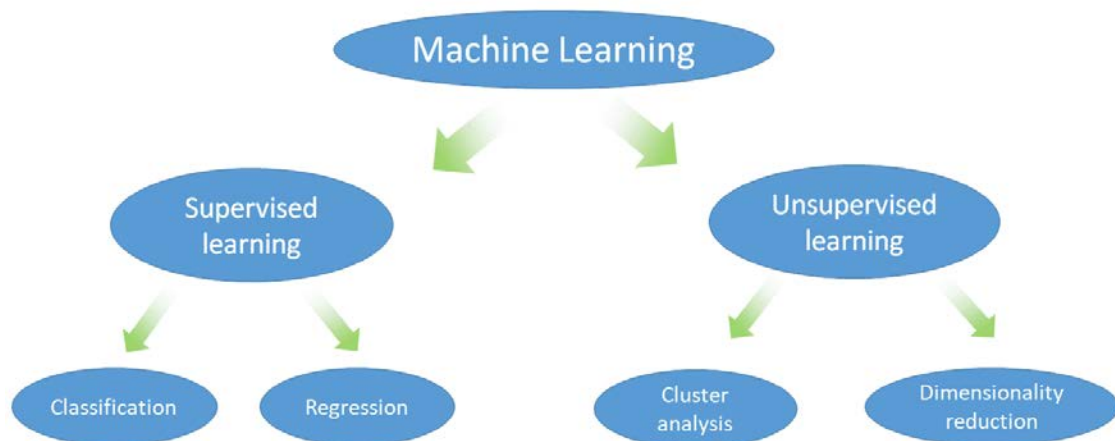
More Focus



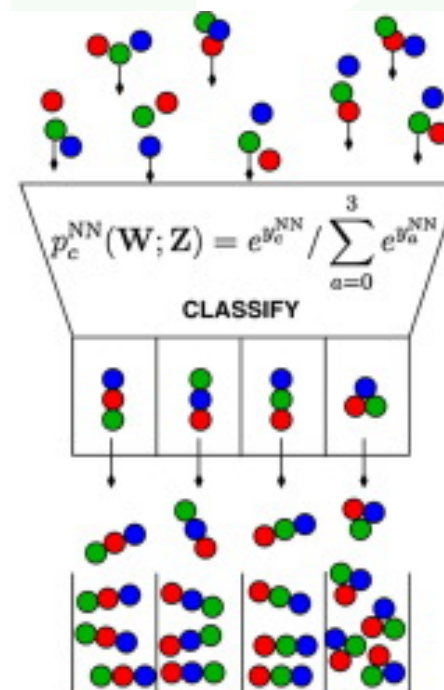
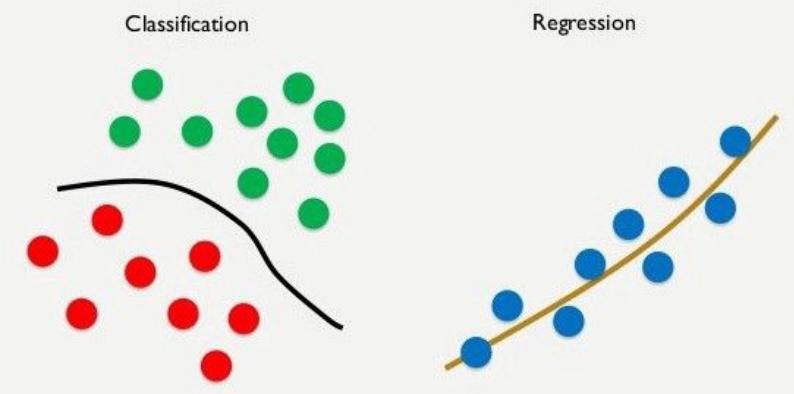


SUPERVISED LEARNING - Classification

Sınıflandırma



CLASSIFICATION vs REGRESSION

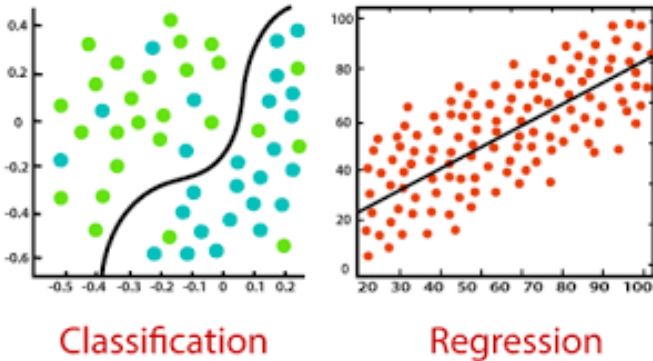




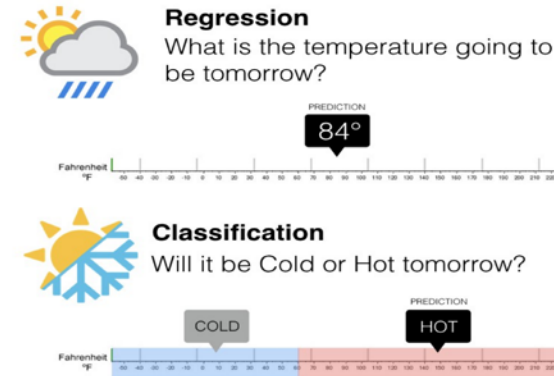
Supervised Learning

Classification 'a Giriş (Sınıflandırma)

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



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



The Idea of
Classification
and **Regression**



Supervised Learning

Classification 'a Giriş

- ✓ Classification un hayattaki kullanım alanları
- ✓ Binary Classification

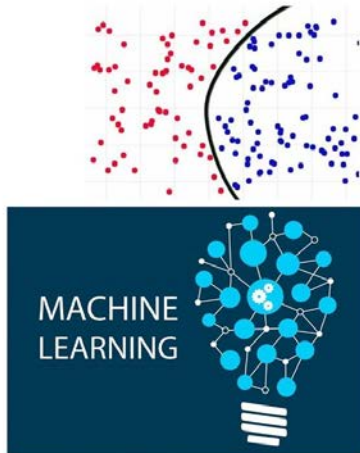
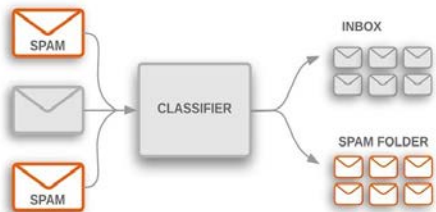
Fraud Protection



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

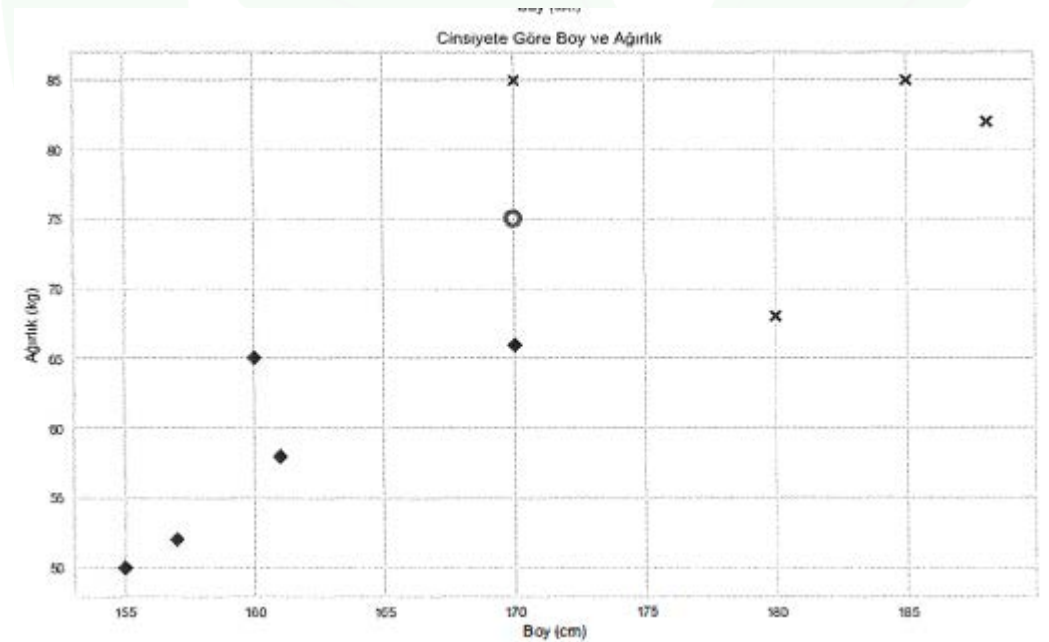
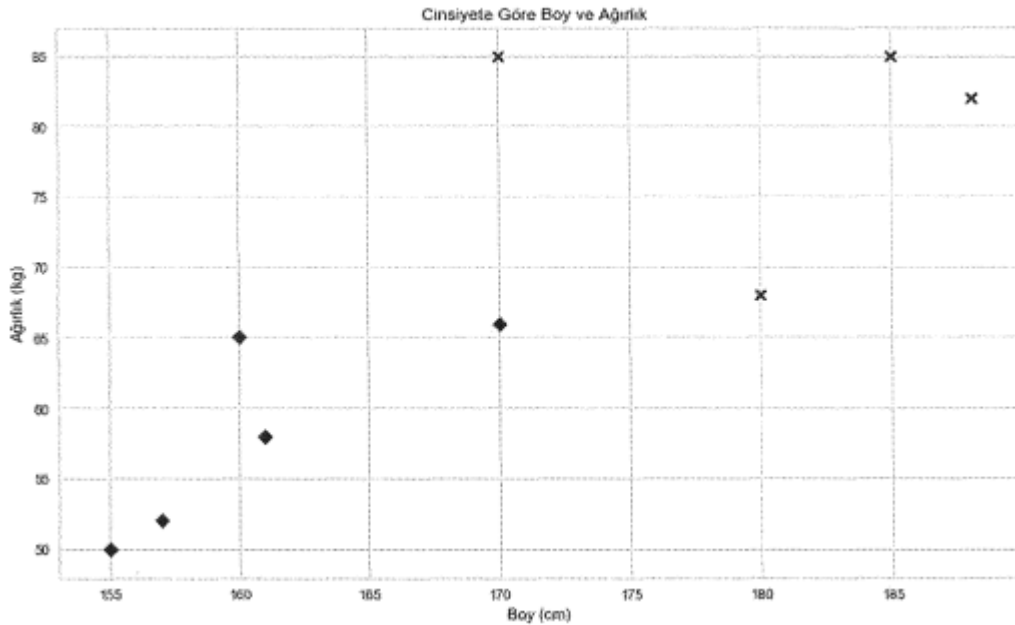




Supervised Learning

Classification 'a Giriş

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

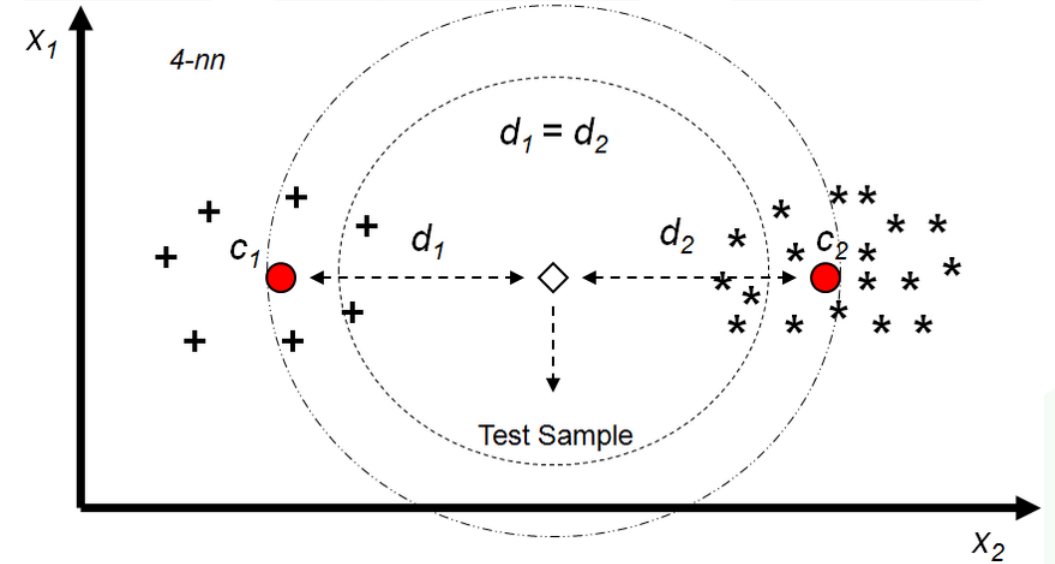
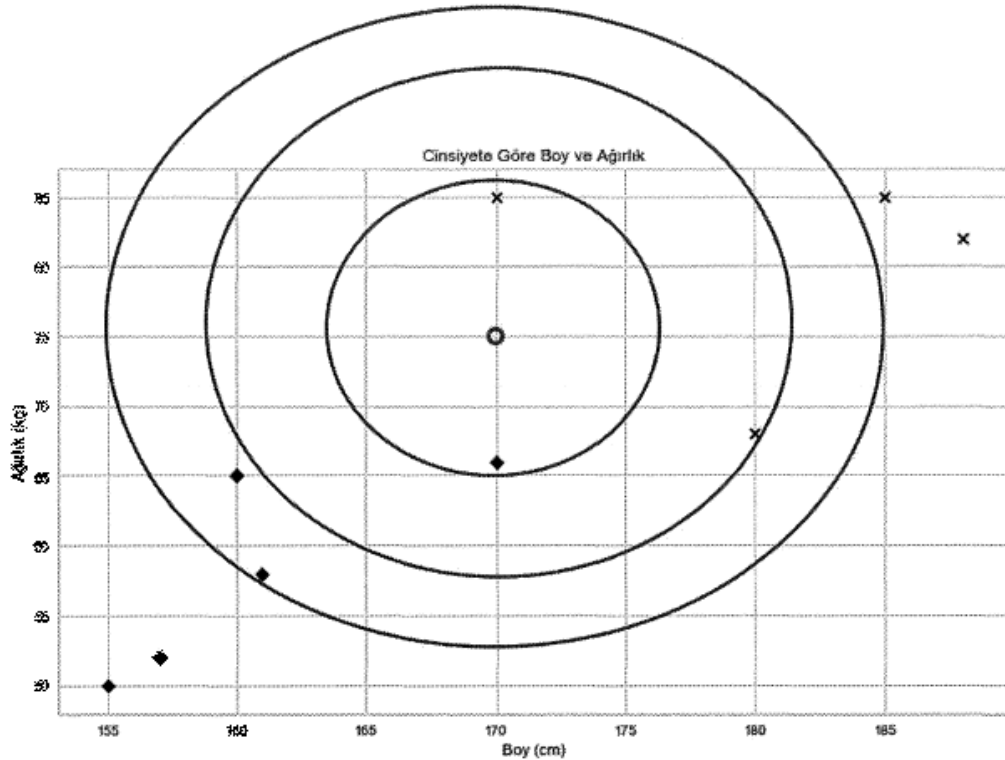




Supervised Learning

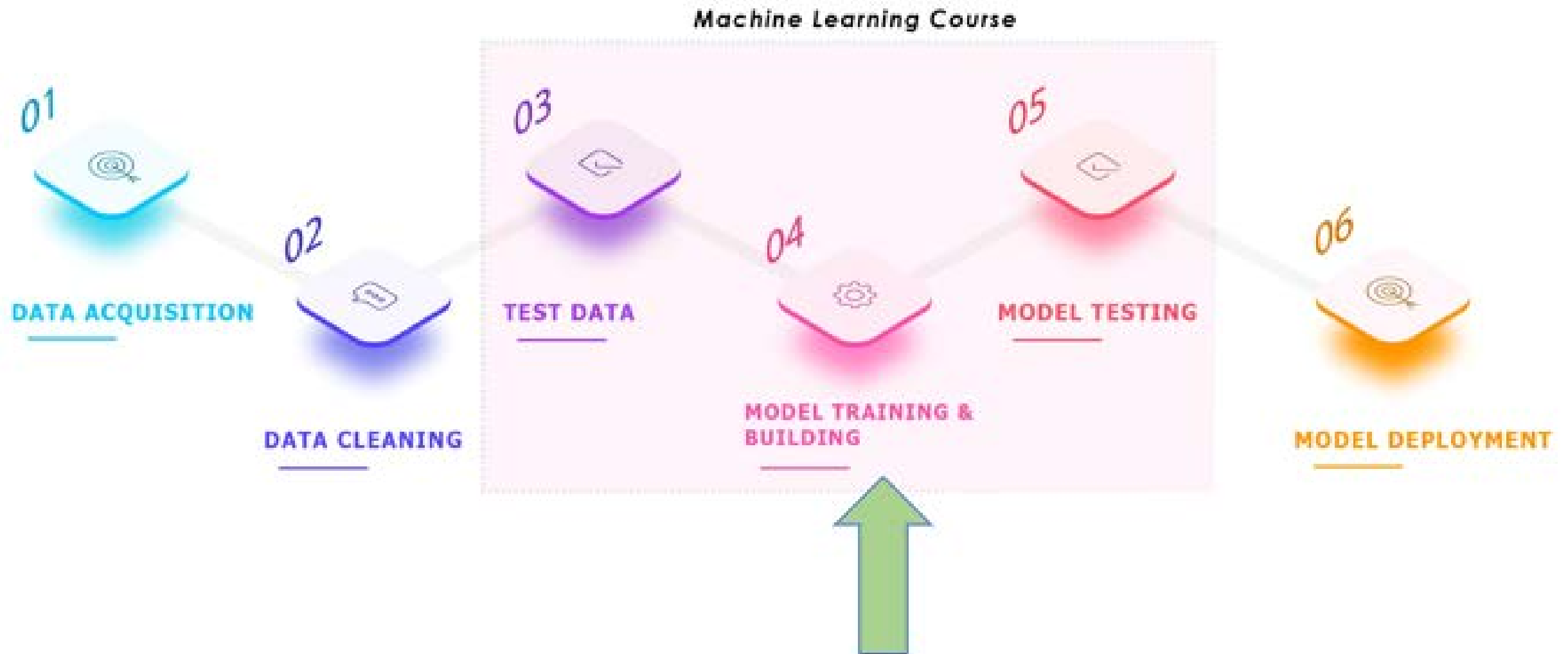
Classification 'a Giriş

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



Where are we?

DATA SCIENCE



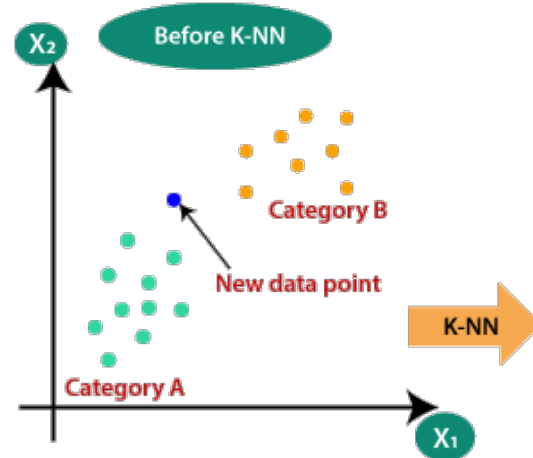
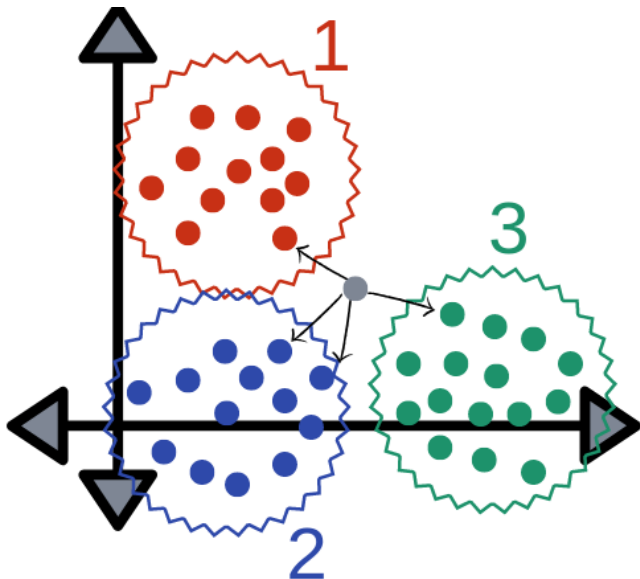


Supervised Learning

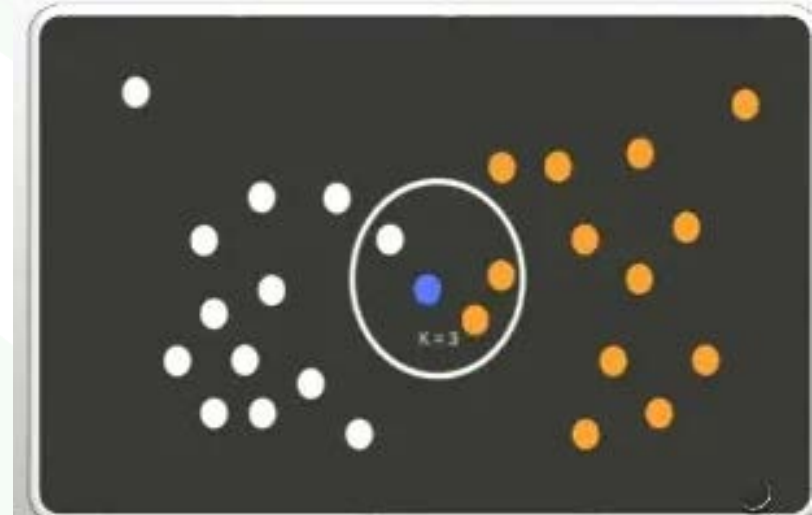
K Nearest Neighbour-KNN Algoritması

K-EN YAKIN KOMŞU Algoritması

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



K-NN

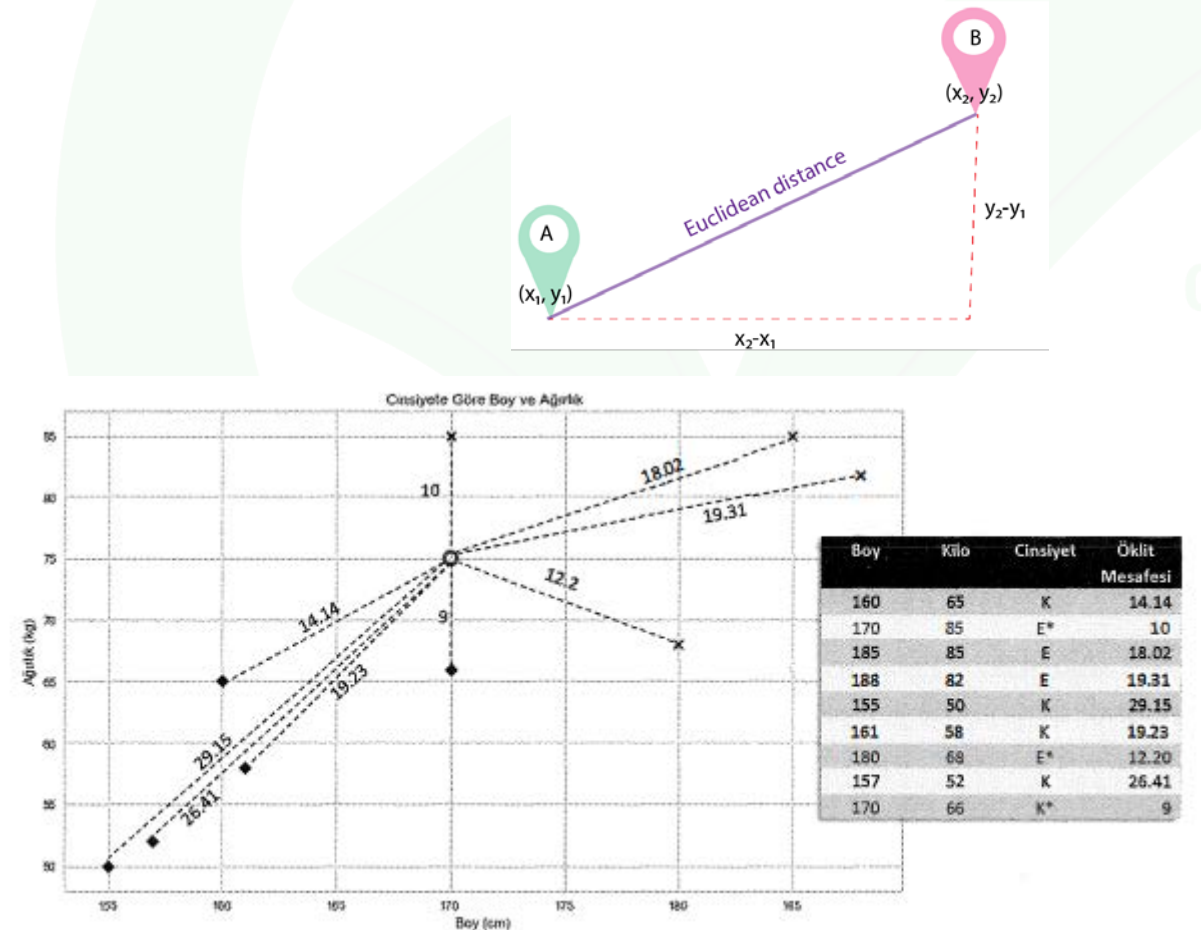
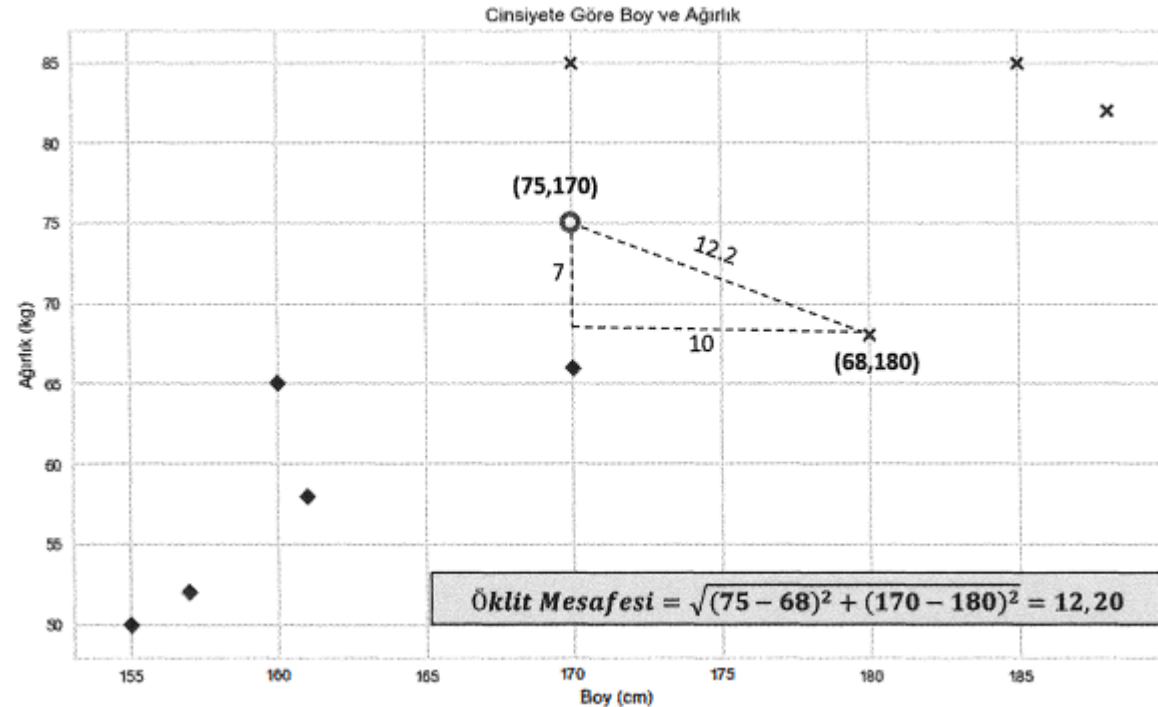




Supervised Learning

K Nearest Neighbour-KNN Algoritması

✓ Öklid mesafesi

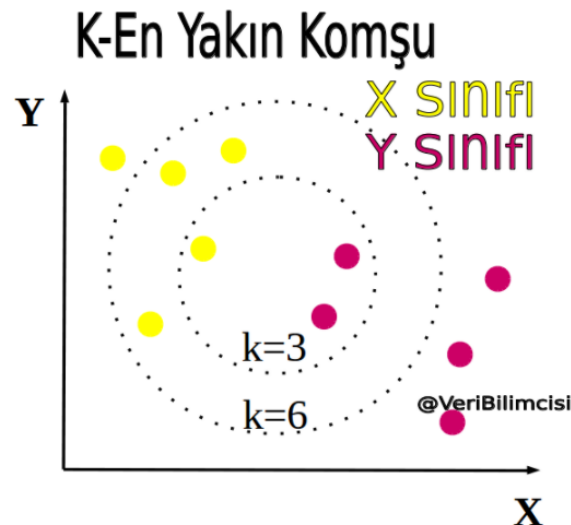
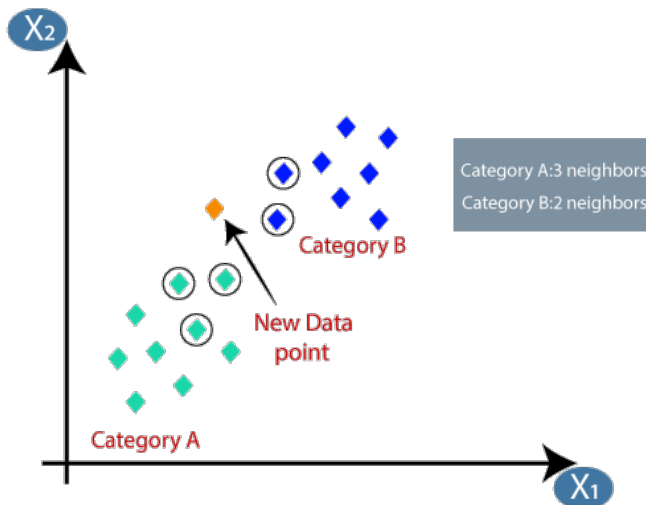




Supervised Learning

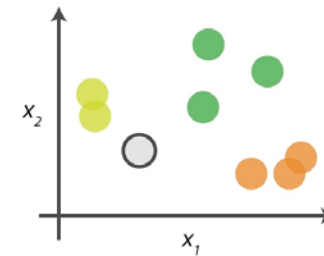
K Nearest Neighbour-KNN Algoritması

- ✓ KNN avantajları
- ✓ KNN dezavantajları



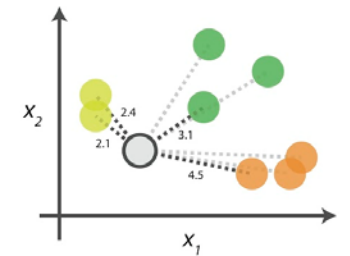
kNN Algorithm

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			
		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 dataspace.

3. Vote on labels

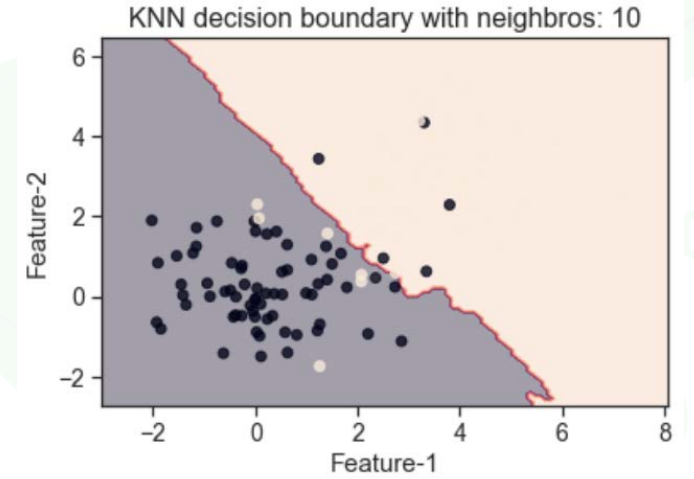
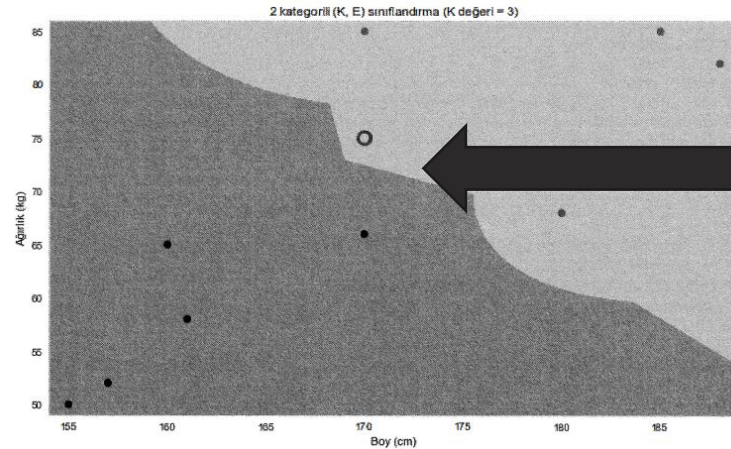
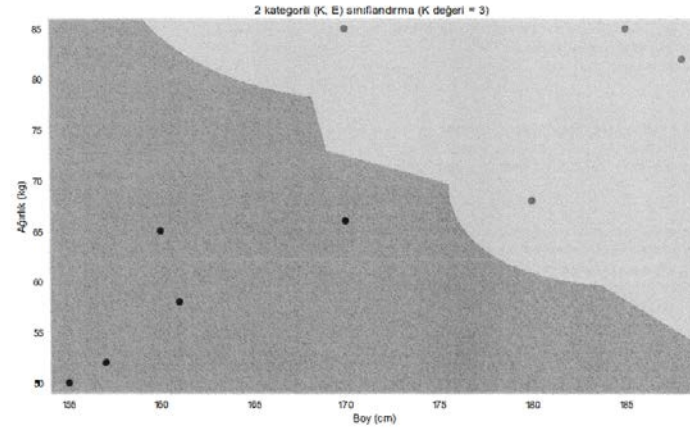
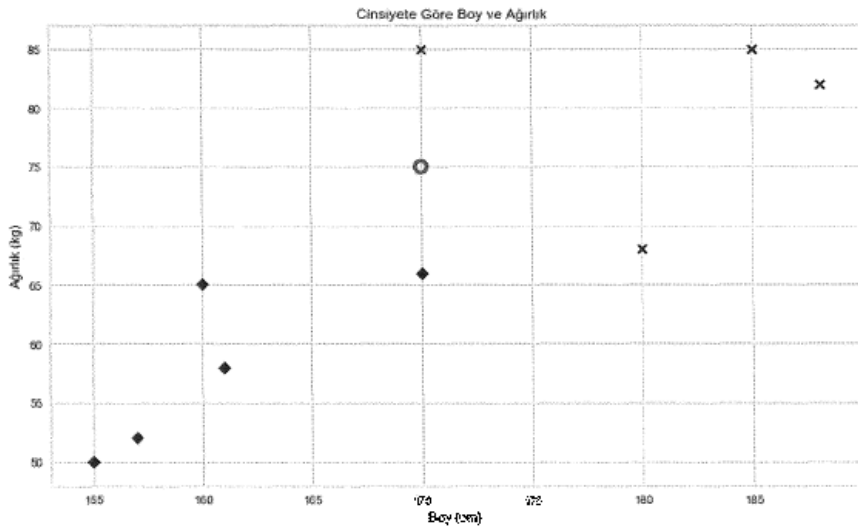
Class	# of votes	
	2	→ Class wins the vote! Point is therefore predicted to be of class .
	1	
	1	

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 Neighbour-KNN Algoritması

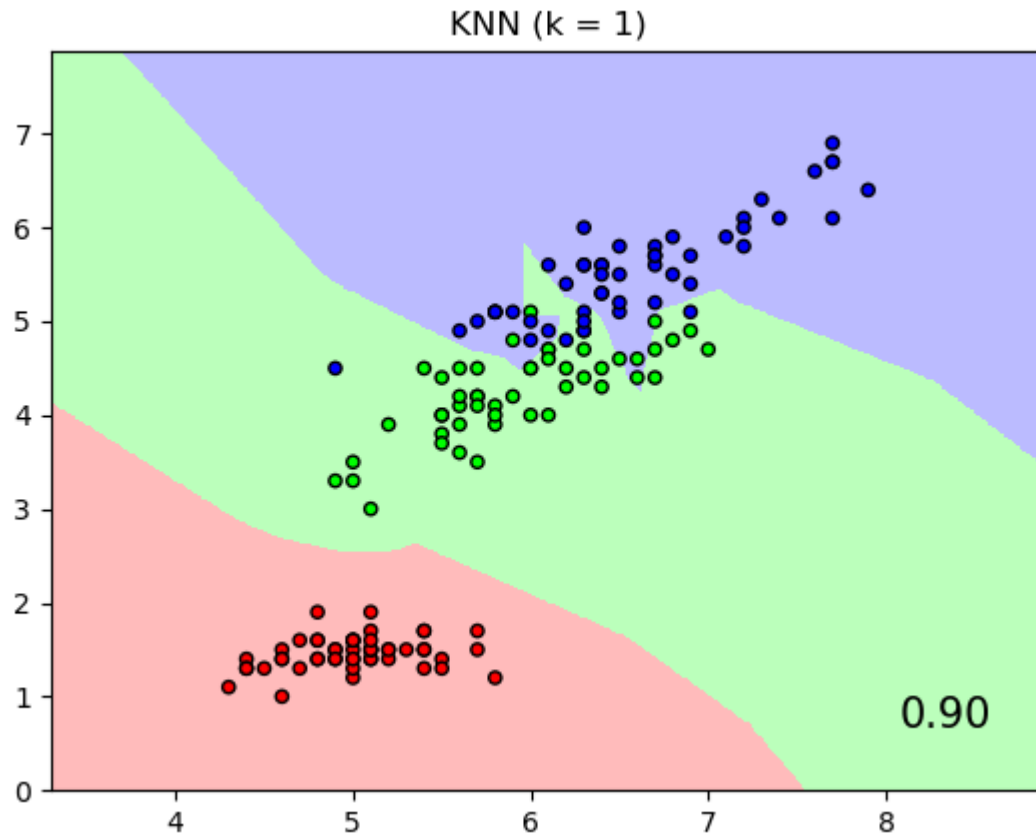
Decision boundary kavramı

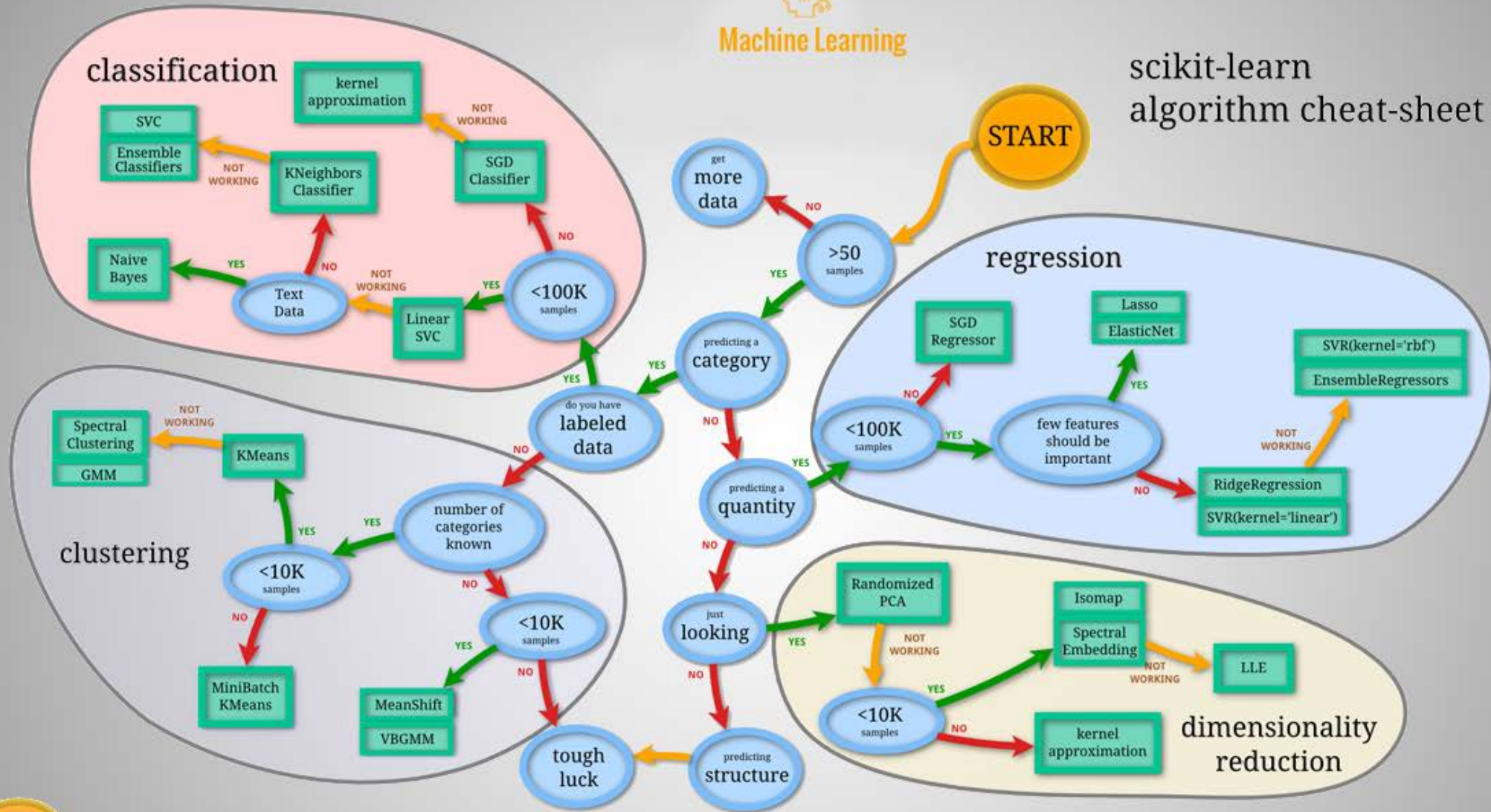




Supervised Learning

Multi Label KNN (Çok sınıflı KNN)







Supervised Learning

Evaluation Metrics for Classification Problems (Performans Ölçütleri)

very
Important



Regresyon kriterleri nasıl olurdu burada?



Confusion metrics kavramı

Regression

- MSPE
- MSAE
- R Square
- Adjusted R Square

Classification

- Precision-Recall
- ROC-AUC
- Accuracy
- Log-Loss

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

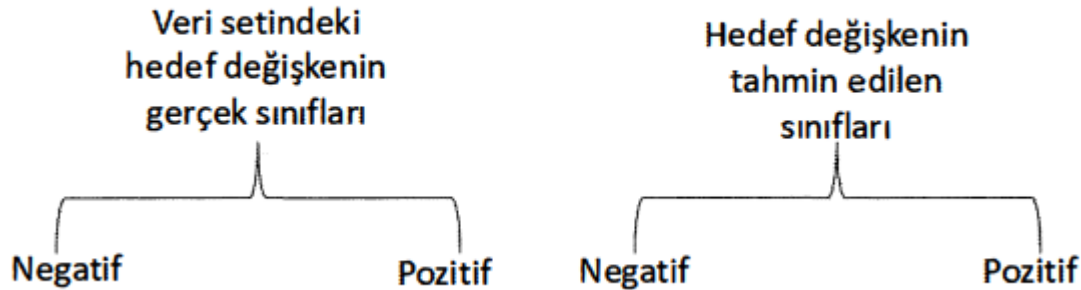


Supervised Learning

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





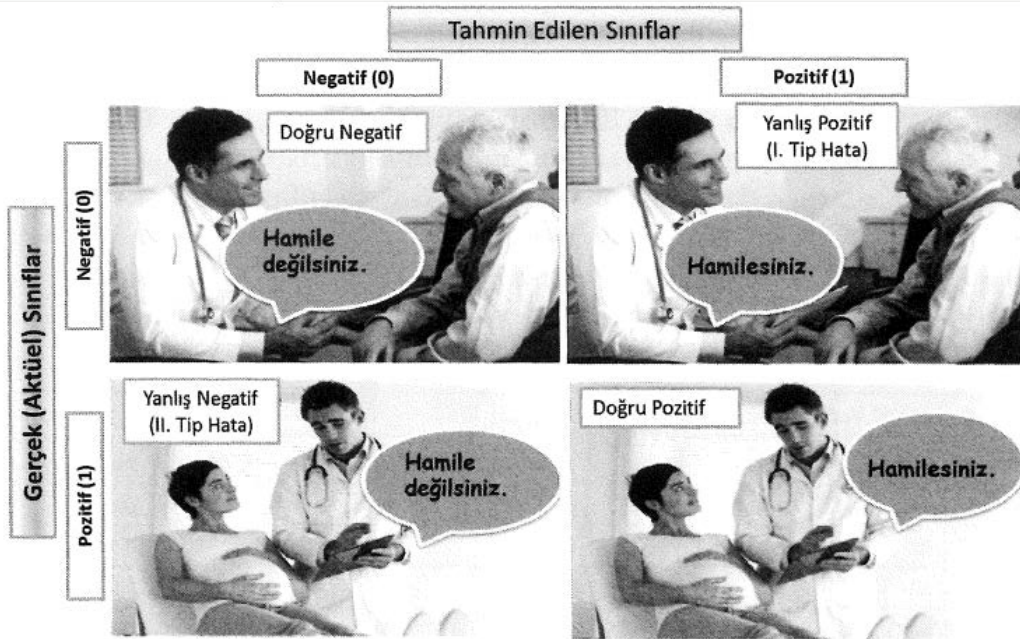
Supervised Learning

Evaluation Metrics for Classification Problems



Confusion metrics kavramı: TN-FN-FP-TP

291-5020, +1 (814) 325-1026, +1 (817) 771-9951, +1 (848) 213-2921, +1 (860) 997-9075, +1 (980) 616-



Tahmin Edilen Sınıflar

Negatif (0)

Pozitif (1)

Doğru Negatif

Yanlış Pozitif
(I. Tip Hata)

Yanlış Negatif
(II. Tip Hata)

Doğru Pozitif

Predicted

Actual		Predicted	
		0	1
Actual	0	30	12
	1	8	56

SINIFLANDIRMA MODELLERİ İÇİN PERFORMANS DEĞERLENDİRME ÖLÇÜTLERİ

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?

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

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

F1 Skoru: Kesinlik (precision) ve Duyarlılığın (Recall) harmonik ortalamasıdır.

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

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

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

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

Performance metrics associated with Class 1

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

(Is your prediction correct?) (What did you predict)

True Negative (You predicted 0)

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

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

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

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

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

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

True +ve rate

		Predicted class	
		+	-
Actual 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)
15]: 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

✗ ✓ ✗ ✓ ✓ ✓ ✓ ✓ ✓ ✗

FN TP FP TN TP TN TN TN TN FN...

Actual Labels: 1 1 0 1 0 0 0 1 1 0

Predicted Labels: 0 1 1 1 0 0 0 1 1 1

✗ ✓ ✗ ✓ ✓ ✓ ✓ ✓ ✓ ✗

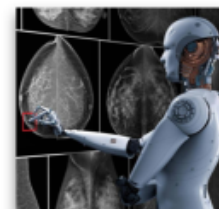
All correctly predicted values (7)

All predicted values (10)

X 100 → **Accuracy = 70 %**

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

✗ ✓ ✗ ✓ ✓ ✓ ✓ ✓ ✓ ✗

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 Veri	Negatif (0)	DN=2	YP=0
	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

Classification Error Metrics

Confusion Matrix

01

Classification Accuracy

02

Area Under ROC curve
(AUC - ROC)

03

Logarithmic Loss

04

F1 Score

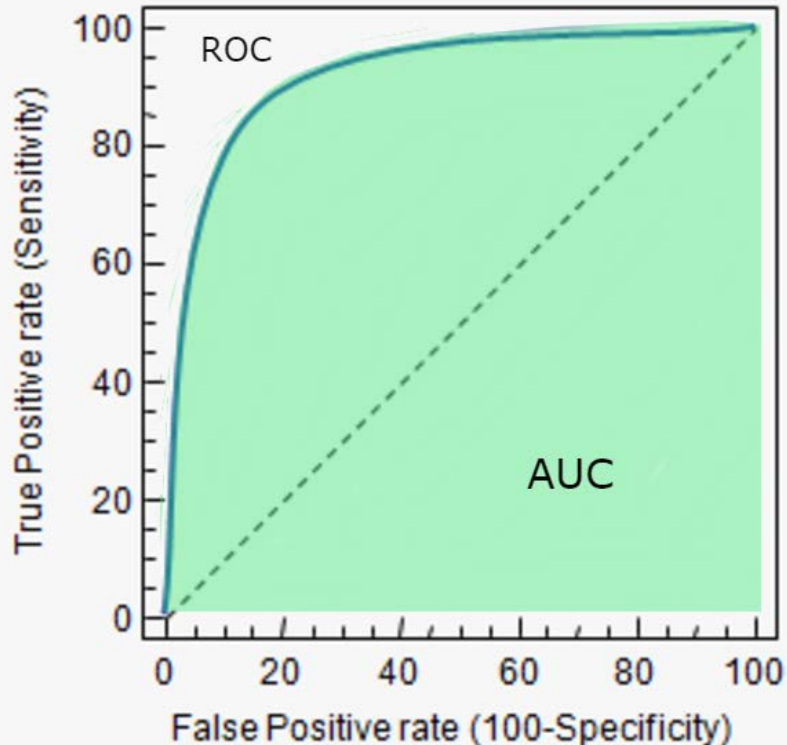
05



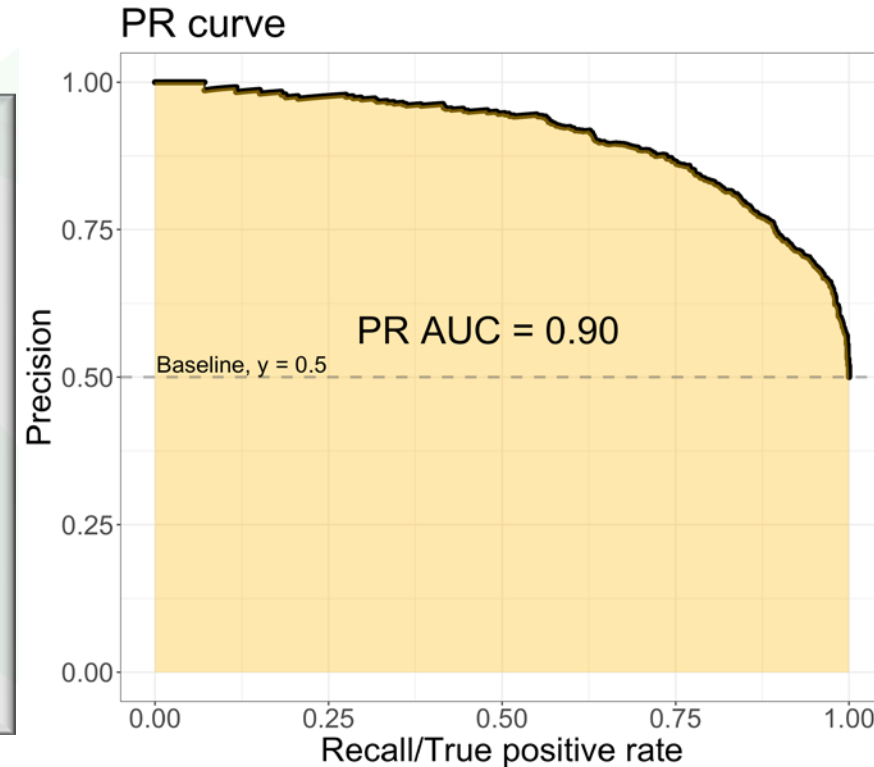
Supervised Learning

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.

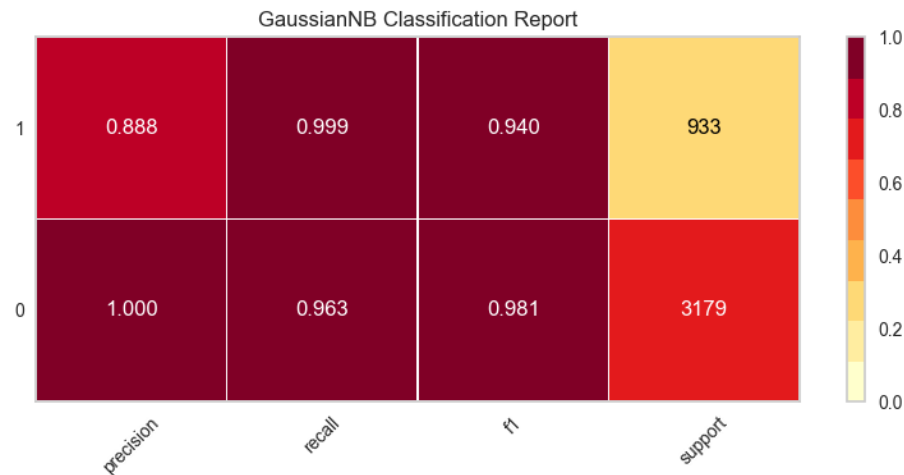
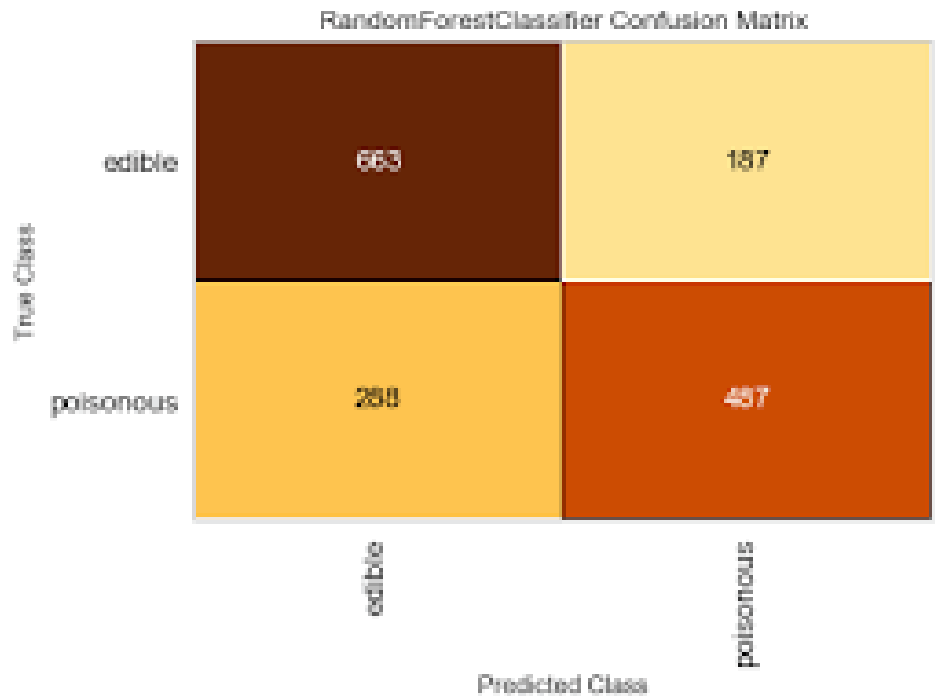




Supervised Learning

Evaluation Metrics for Classification Problems

✓ Yellowbrick ile confusion matrices



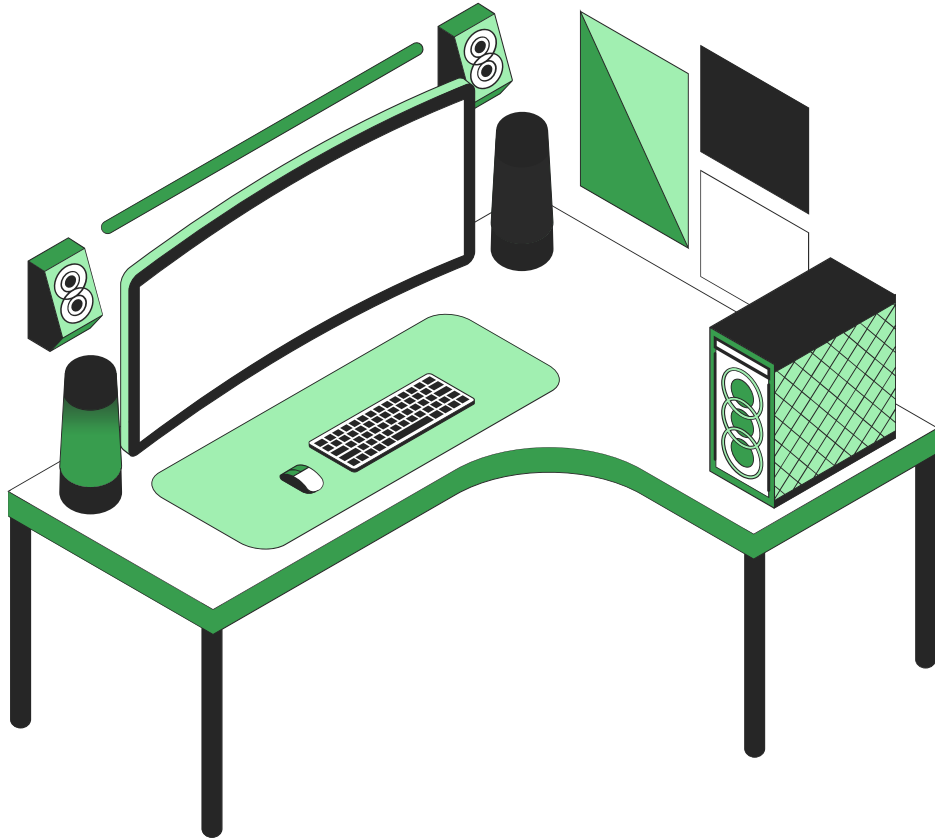


Bu dersi anladım..



Everything is
clear ?





Do you
have any
questions?

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