

VERİ ANALİZİ OKULU

2025-2026

Makine Öğrenmesi ve Yapay Zekaya Giriş

Dr. Öğr. Üyesi M. Fuat KINA

Marmara Üniversitesi
Nüfus ve Sosyal Araştırmalar Enstitüsü





Dersin Çerçevesi

- Makine öğrenmesi nedir?
- Makine öğrenmesi türleri (gözetimli, gözetimsiz)
- Gözetimli öğrenme: *sınıflandırma* (k-NN), *regresyon* (düzenlileştirme-regularizasyon)
- Hiperparametre ayarı
- Ön-işleme ve pipeline
- Gözetimsiz öğrenme: **k-means** kümeleme



Makine Öğrenmesi Nedir?

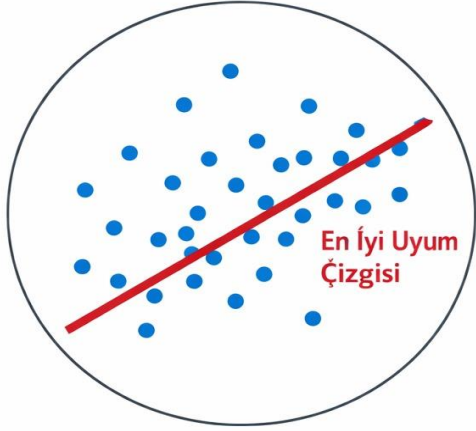
- Bilgisayarlara, açıkça programlamadan, veriden öğrenerek iş yapmayı öğretme teknolojisi.
 - Örnek: Bir e-postanın içerik ve gönderen bilgisine bakarak spam / değil diye karar verme.
 - Örnek: Wikipedia girdilerini, içerdiği kelimelere göre kümelere ayırma.
- **İstatistikten makine öğrenmesine** geçiş:
 - Veriden anlamlı neden-sonuç ilişkileri geliştirmek 
 - İçinde yaşadığımız dünyayı tahmin etmek (*nowcasting / forecasting*) 
 - **Temel fark:** Makine öğrenmesi modeli (bir istatistik modelinin aksine) geliştirildiği yerdeki değil, bunun dışında yer alan uzaydaki ilişkileri öğrenebildiği ölçüde başarılı kabul edilir.
- Etiketler varsa: gözetimli öğrenme (supervised).
- Etiketler yoksa: gözetimsiz öğrenme (unsupervised).



Makine Öğrenmesi Nedir?

İstatistiksel Modelin Mantığı

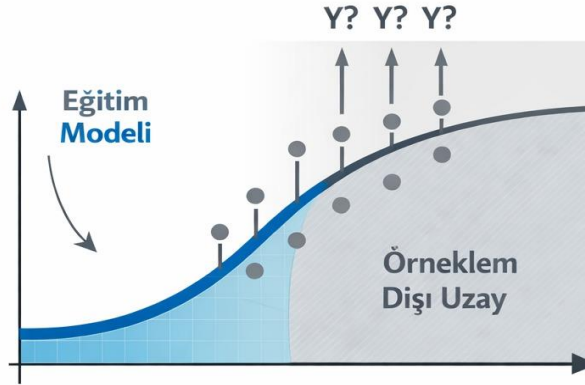
Veriyi Anlamlandırmak



$X \rightarrow Y$ İlişkisini Anlamak

Makine Öğrenmesi (Gözetimli Öğrenme)

Tahmin Yapmak



X Değerlerinden Y Tahmini

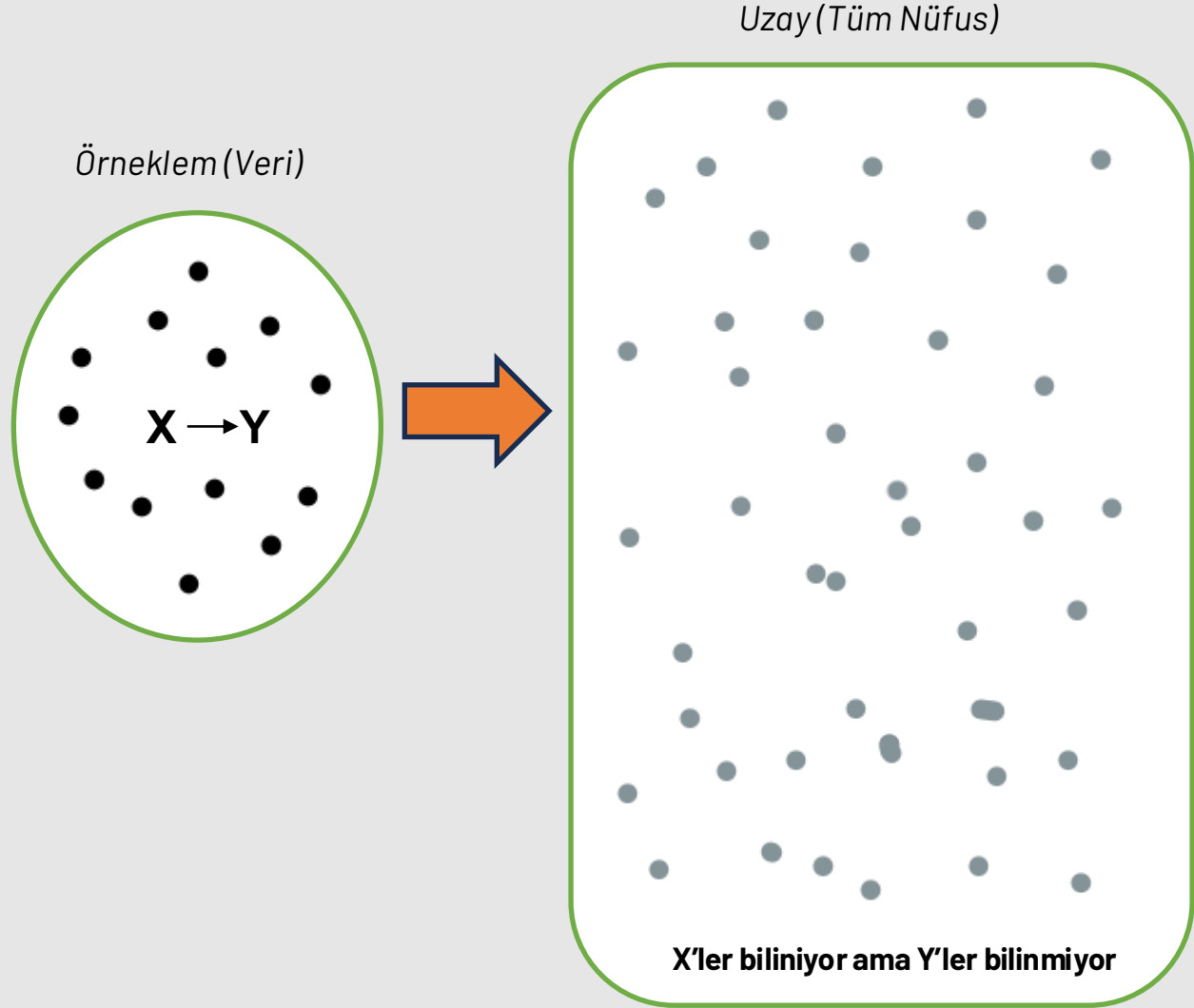
İstatistiksel çıkarım: 1000 kişilik bir anket datası içerisinde eğitimin gelire etkisi

- Amaç: Neden-sonuç ilişkisi kurup hipotezleri test etmek
- Örneklem dışı uzak (out-of-sample) önemsiz

Makine öğrenmesi: Haber metnine bakıp protesto-eylem var mı yok mu karar veren bir model.

- Amaç: Henüz tedavülde olmayan bir veriyi yaratmak
- Örneklem dışı uzak (out-of-sample) modelin temel ilgisi

Makine Öğrenmesi Nedir?



Örneklem (Veri): Hem X'leri hem Y'leri bildiğim yer. Makine Öğrenmesi için eğitim alanı, tahminlemenin başladığı yer. İstatistik içinse kurulacak olan iddianın hem başladığı hem bittiği yer.

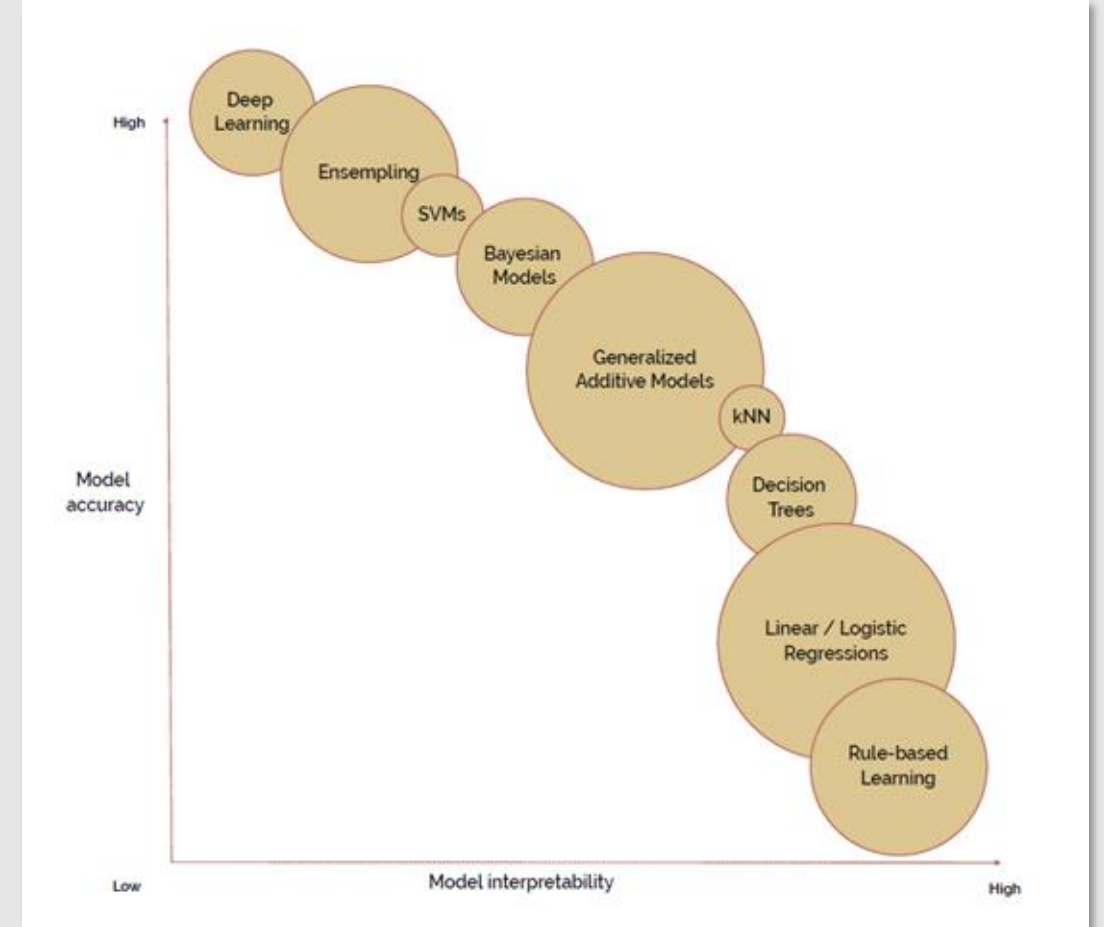
Uzay (Tüm Nüfus): Tahminleme yaptığım alan. Tüm insanlık, tüm Türkiye, tüm haber metinleri, tüm tweetler

- İstatistiksel çıkarımda örneklemden uzaya geçiş temsiliyet problemiyle ilgilidir.
- Makine öğrenmesi ise bu ilişkiyi temsiliyete dayalı bilimsel varsayımla değil test sonuçlarıyla edinir. Eğitim datası dışında (out-of-sample) X'ler Y'leri doğru bir şekilde tahmin edebiliyorsa model başarılıdır.

Makine Öğrenmesi Nedir?

- Modeller büyük ölçüde 1980'lerde geliştirildi, ancak 2000'lerde uygulamaya konulabildi.
- Modellerin karmaşıklığı zaman içerisinde arttıkça yorumlanabilirlik azalırken, tahminleme kapasitesi yükseldi.
 - Makine öğrenmesi (2000-2010)
 - Derin öğrenme ve sinir ağları (2010-2020)
 - Üretken yapay zeka (2020'ler)

(Şekil için: Toosi et al. 2022)



Gözetimli ve Gözetimsiz Öğrenme

Gözetimsiz, etiketsiz

- Etiketlenmemiş datanın içerisindeki saklı desen ve yapıları ortaya çıkartır
- Kümeleme, boyut eksiltme, yapılandırılmamış konu modellemesi

Gözetimli, etiketli

- Amaç, doğru çıktısını bildiğimiz veriden öğrenerek, çıktısını bilmediğimiz yeni veriler için tahmin yapabilmektir
Zaman alıcı veya maliyetli bir manuel işi otomatikleştirmek; gelecek hakkında tahminlerde bulunmak vb.
Etiketli veriye ihtiyaç duyar (tarihsel veriler, kitle kaynaklı etiketleme vb.)

1. Gözetimli Öğrenme

Gözetimli Öğrenme

- Tahmin değişkenlerini (özellikleri) ve bir hedef değişkeni kullanarak, hedef değişkeni tahmin edebilen bir model kurmak
 - Sürekli hedef değişken
 - **Regresyon**
 - İkili (binary) hedef değişken
 - **Sınıflandırma (classification)**
- **Feature** ya da bağımsız değişken = *predictor* (tahmin değişkeni)
- **Target** ya da bağımlı değişken = *response / outcome* (hedef / çıktı)
- Python: scikit-learn kütüphanesi (**sklearn**)



Bunlardan hangisi bir
sınıflandırma problemidir?

Yandaki makine
öğrenmesine ilişkin 4 örnek
uygulama verilmiştir.

Bunlardan hangisi
gözetimli bir sınıflandırma
problemidir?



Scan Me

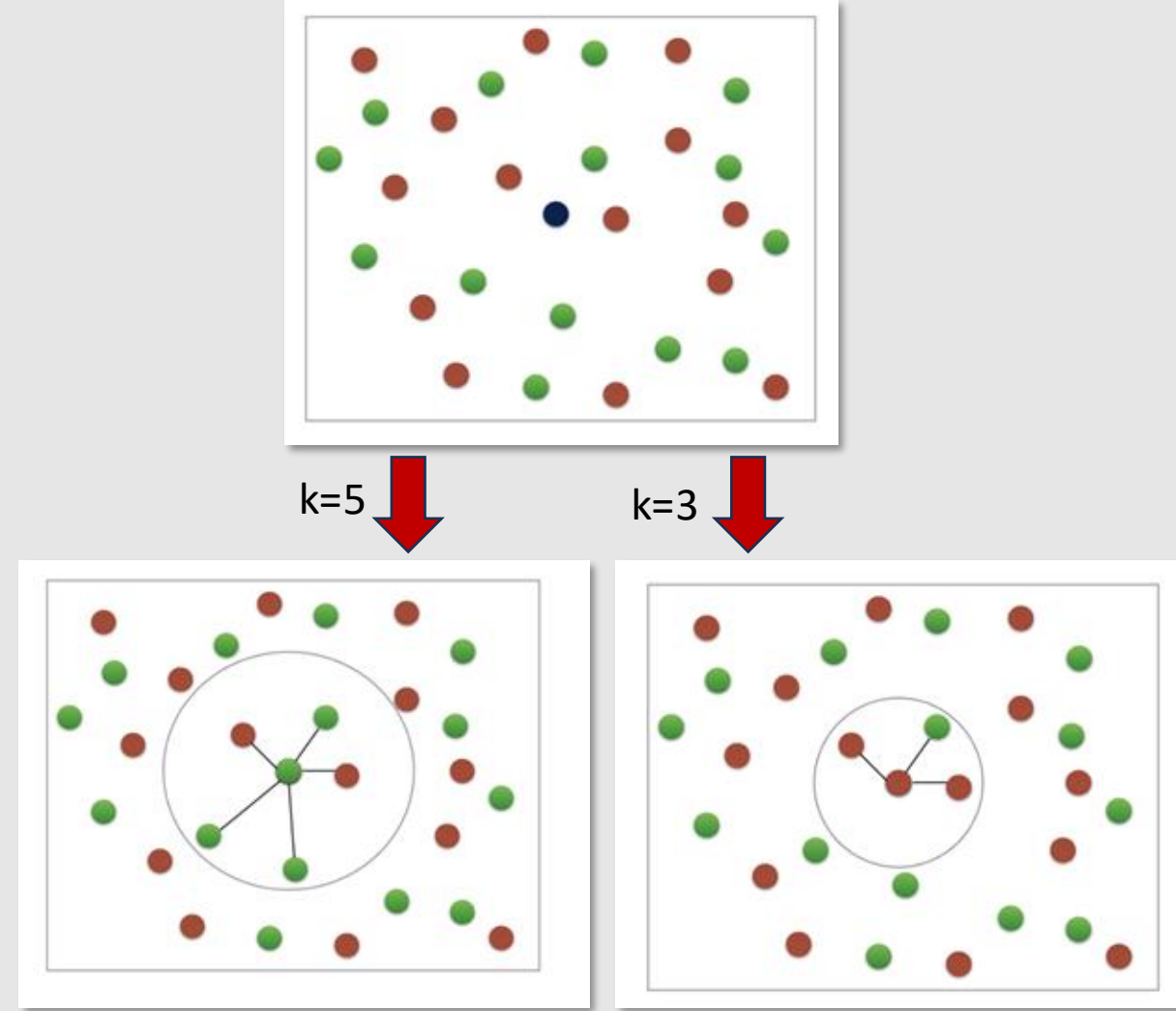
- A. Etiketlenmiş finansal verileri kullanarak, bir hissenin değerinin gelecek hafta **yukarı mı yoksa aşağı mı** gideceğini tahmin etmek.
- B. Etiketlenmiş konut fiyatı verilerini kullanarak, çeşitli özelliklere göre yeni bir evin **fiyatını** tahmin etmek.
- C. Etiketsiz verileri kullanarak, çevrimiçi bir eğitim şirketinin öğrencilerini **öğrenme stillerine göre** farklı kategorilere/kümelere ayırmak.
- D. Etiketlenmiş finansal verileri kullanarak, bir hissenin **gelecek hafta tam olarak kaç değerinde** olacağını tahmin etmek.

1. A. Sınıflandırma

Sınıflandırma

k-En Yakın Komşu (k-Nearest Neighbors)

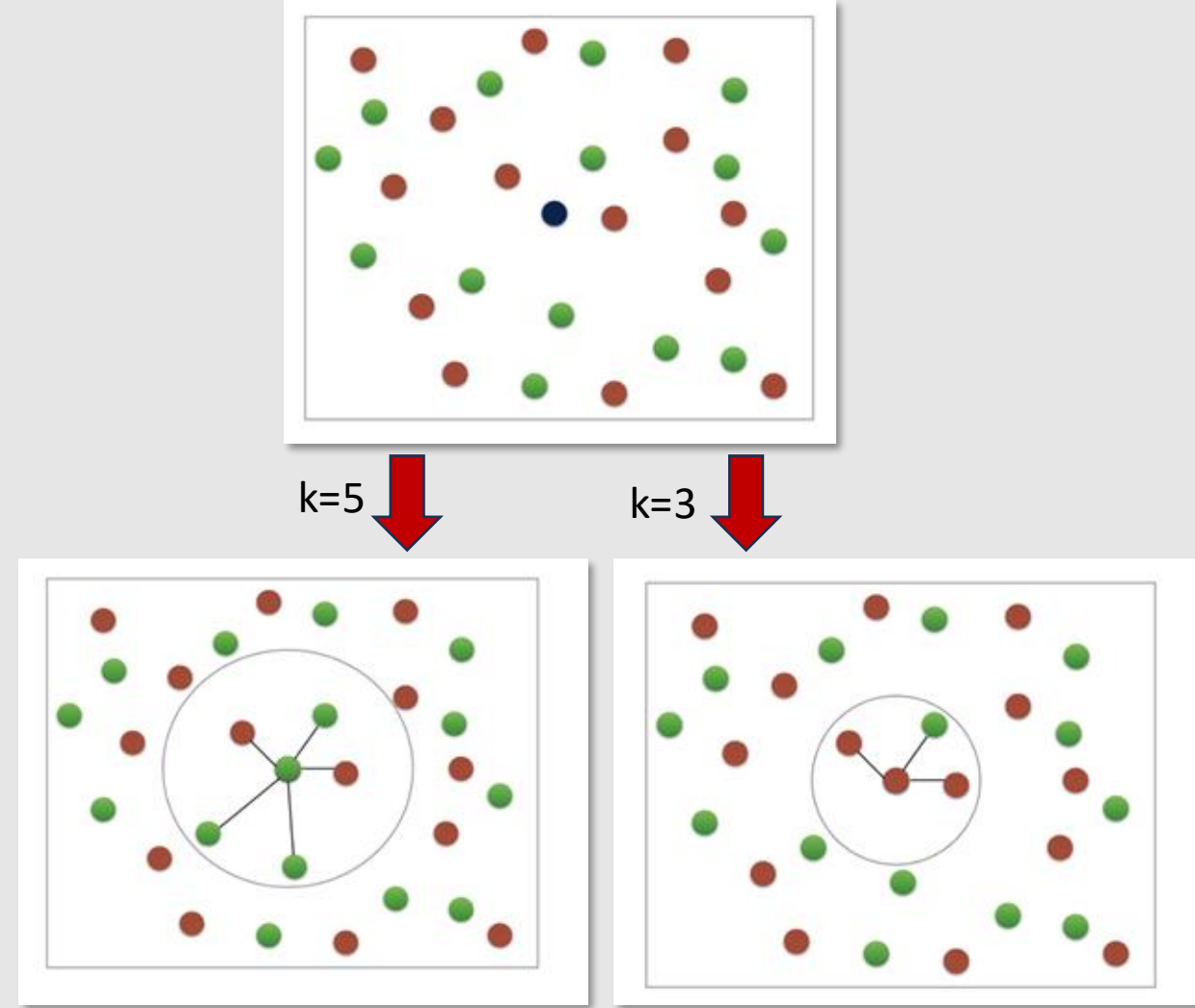
- En yakın **k** tane etiketli veri noktasına bakar ve **çoğunluk oyunu** alır.
- Sınıflar arasında **karar sınırları (decision boundaries)** çizer.
- *5 komşu vs 3 komşu*
- Makine öğrenmesinde iki adım vardır: **öğrenme** ve **tahmin etme**
- Veriler üzerinde bir modeli eğitmek = **.fit()**
- Yeni veri için etiketleri tahmin etmek = **.predict()**



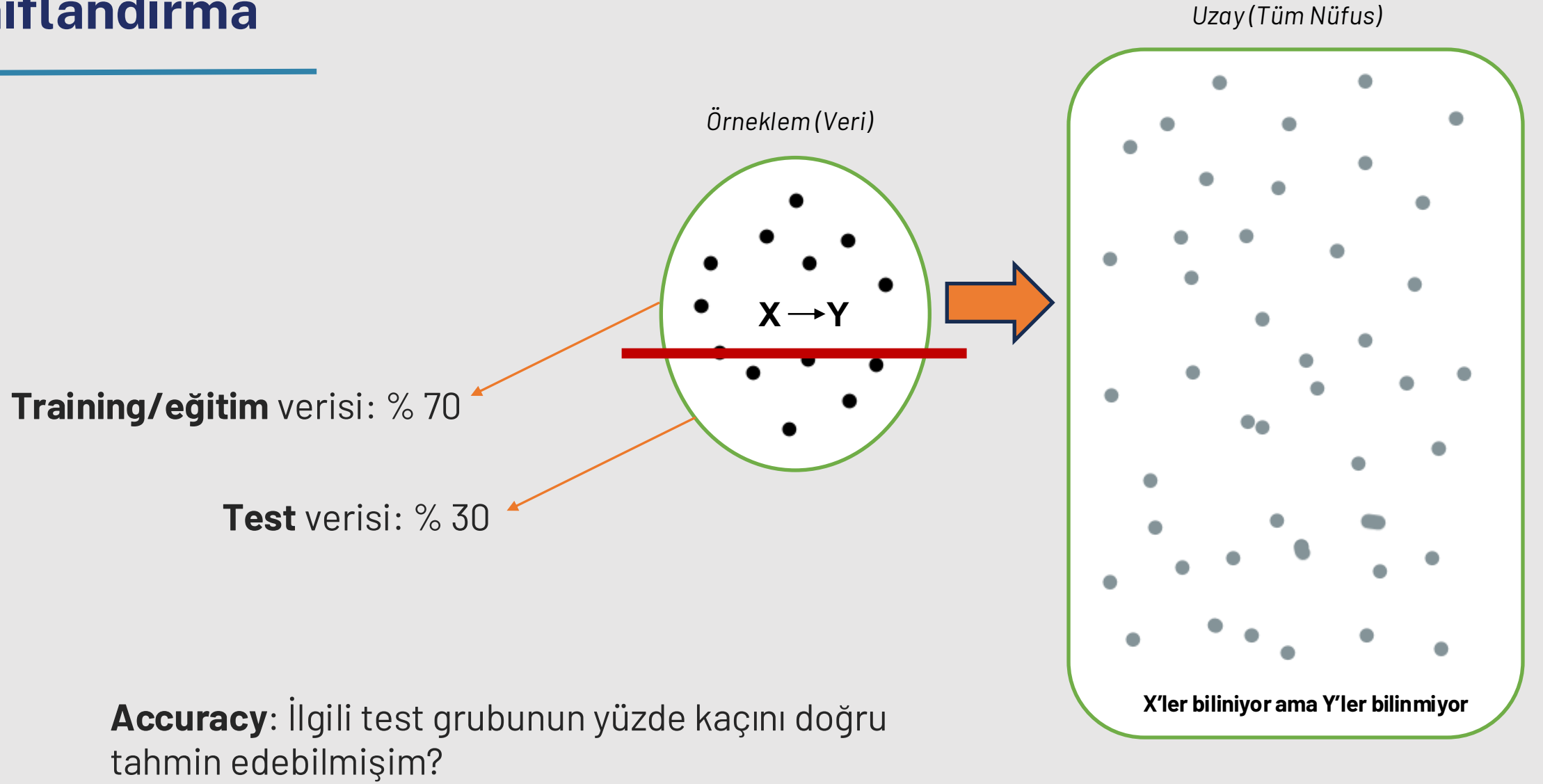
Sınıflandırma

Modelin performansını nasıl ölçeriz?

- **Doğruluk (accuracy)** = Doğru tahminlerin tüm tahminlere oranı
- Peki doğruluğu hangi veri üzerinde kontrol etmeliyim: **modeli bina ettiğim (eğitim) veri mi, yoksa yeni (tahmin) veri mi?**
- Çözüm: **Train/test split (eğitim/test ayrımı)**
 - Sınıflandırıcıyı **eğitim seti** üzerinde fit et
 - **Test seti** üzerinde tahmin yap ve doğruluğu orada ölç



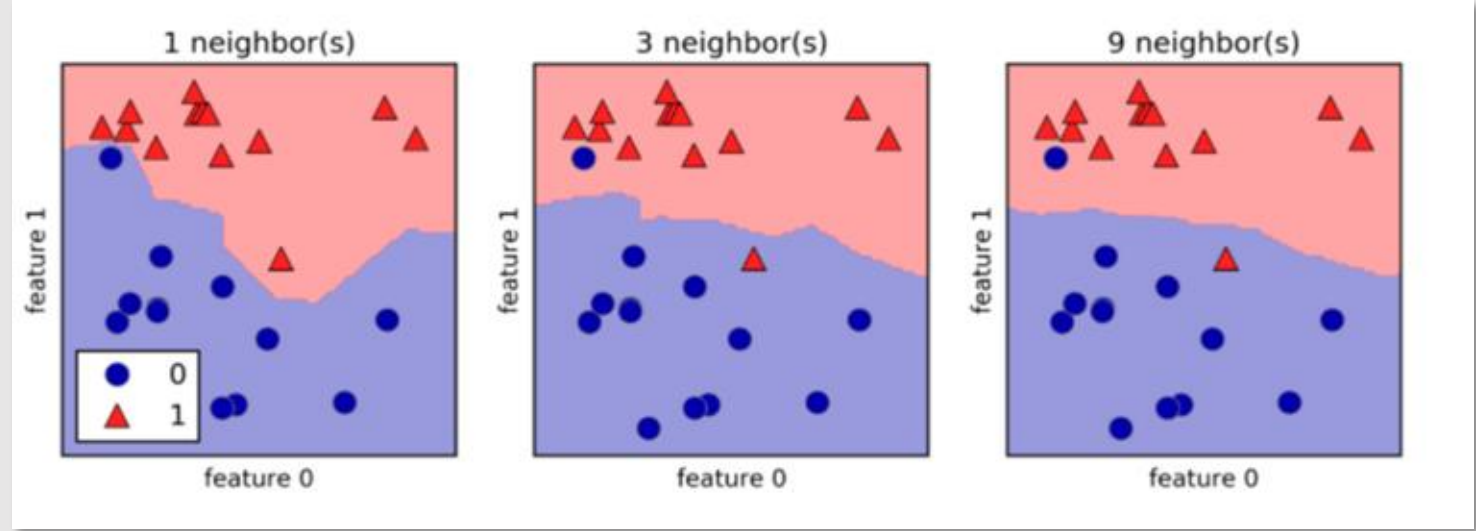
Sınıflandırma



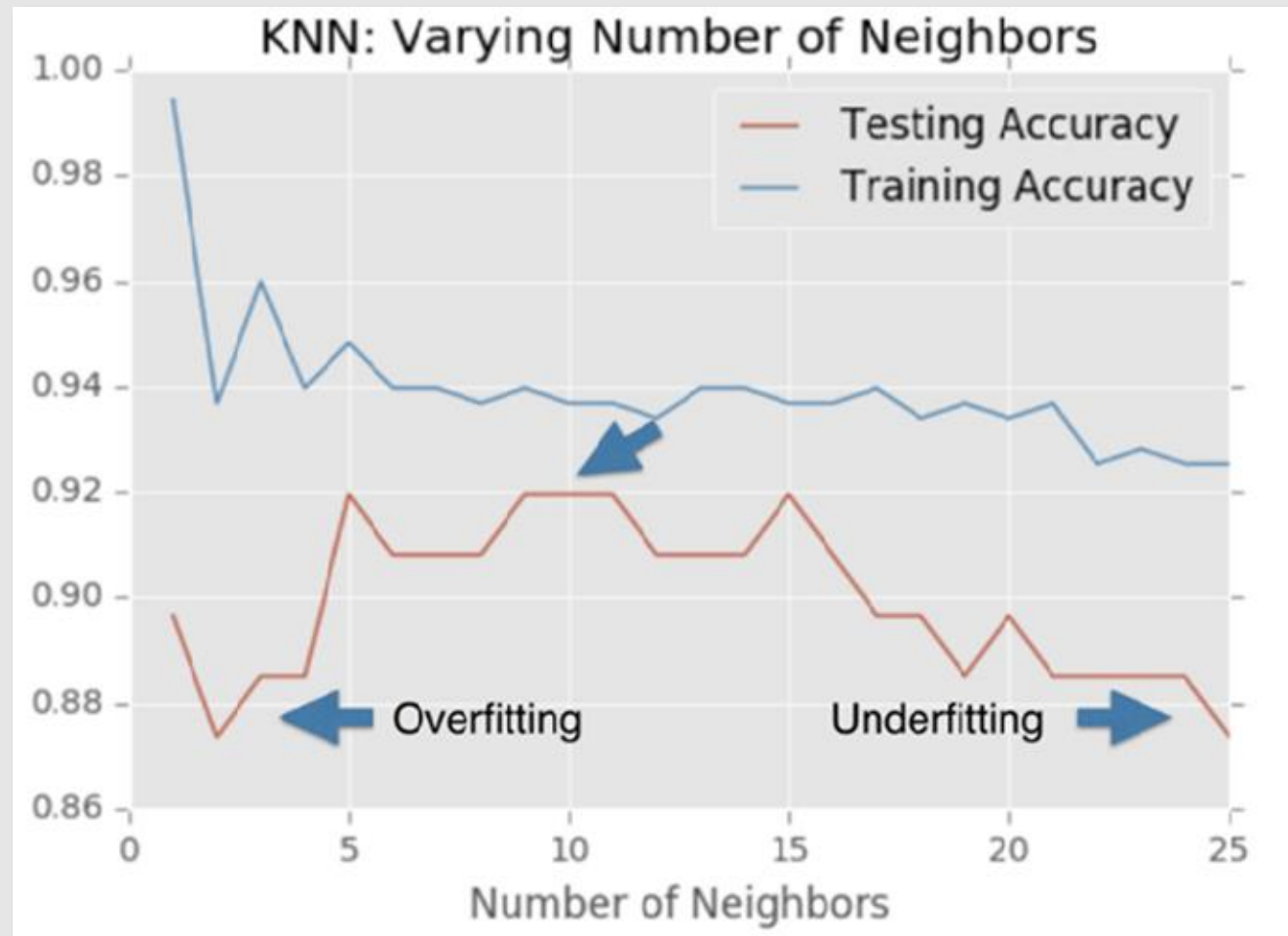
Accuracy: İlgili test grubunun yüzde kaçını doğru tahmin edebilmişim?

Sınıflandırma

Model Kompleksitesi



- Daha büyük **k** \Rightarrow daha pürüzsüz karar sınırları \Rightarrow **daha az karmaşık model**
- Daha küçük **k** \Rightarrow **daha karmaşık model** \Rightarrow aşırı uyuma (*overfitting*) yol açabilir



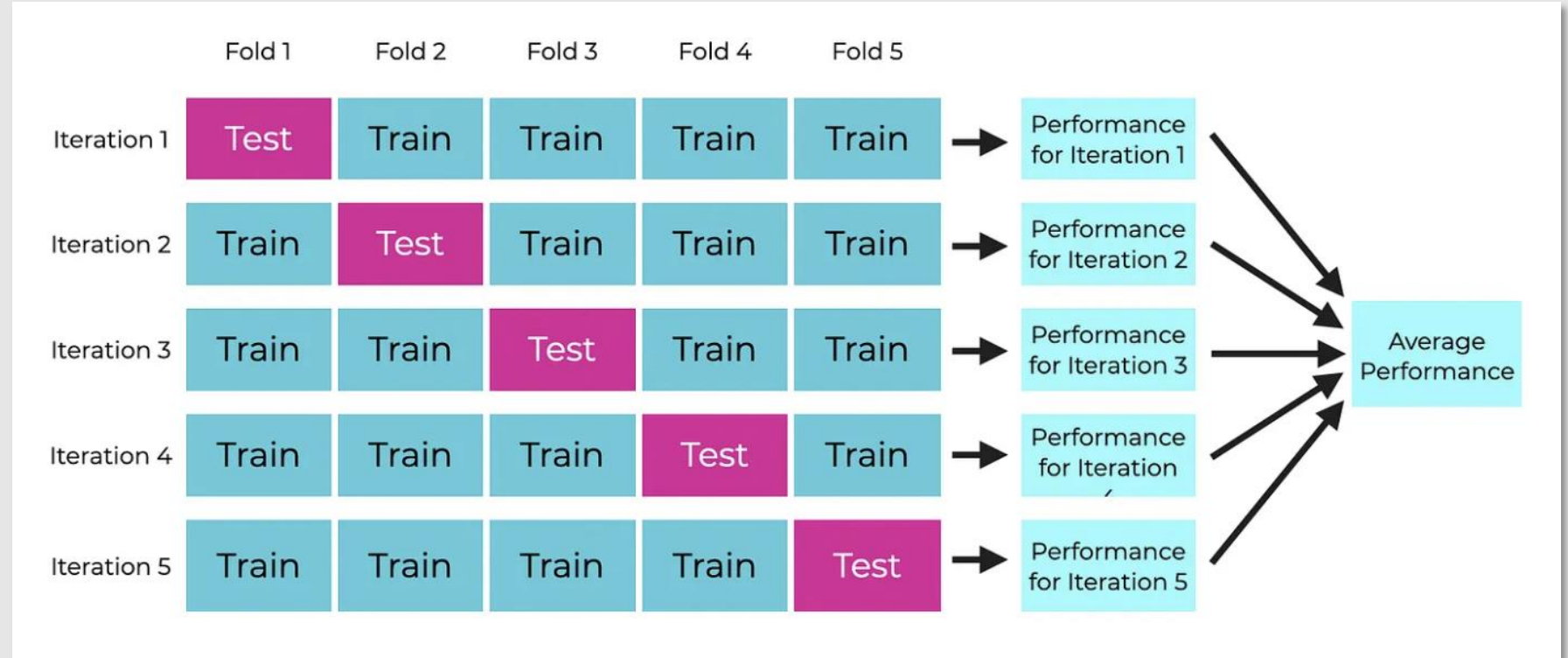
1. B. Regresyon

Makine öğrenmesi modellerinde gerçekten “yeni” olan ne?

- Çapraz doğrulama (cross-validation)
- Düzenleştiriciler (regularizers)
 - Lasso
 - Ridge
- Modeli ince ayarlamak (fine-tuning)
 - ROC eğrisi
 - AUC (eğri altındaki alan) hesaplama
- Ön-işleme adımları ve bunların **pipeline**’lar içinde otomatikleştirilmesi

Doğrusal Regresyon

Çapraz doğrulama



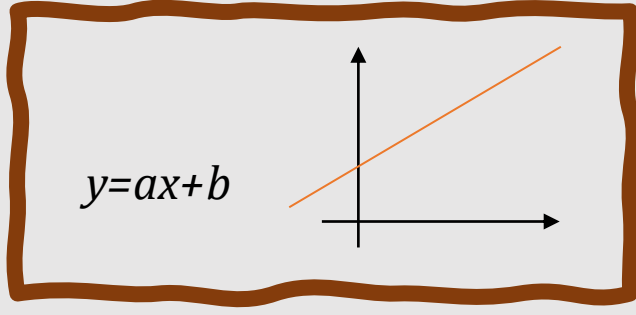
Doğrusal Regresyon

Aşağıdakilerden hangisi bir doğrusal regresyon problemidir?

1. Bir e-ticaret şirketinin, etiketlenmiş müşteri verilerini kullanarak bir müşterinin belirli bir ürünü satın alıp almayacağını tahmin etmesi.
2. Bir sağlık şirketinin, kanser tümörlerine dair verileri (örneğin geometrik ölçümlerini) kullanarak yeni bir tümörün iyi huylu mu kötü huylu mu olduğunu tahmin etmesi.
3. Bir restoranın, yorum verilerini kullanarak bir yoruma pozitif ya da negatif duygu ataması.
4. Bir bisiklet paylaşım şirketinin, zaman ve hava durumu verilerini kullanarak herhangi bir saatte kiralanan bisiklet sayısını tahmin etmesi.



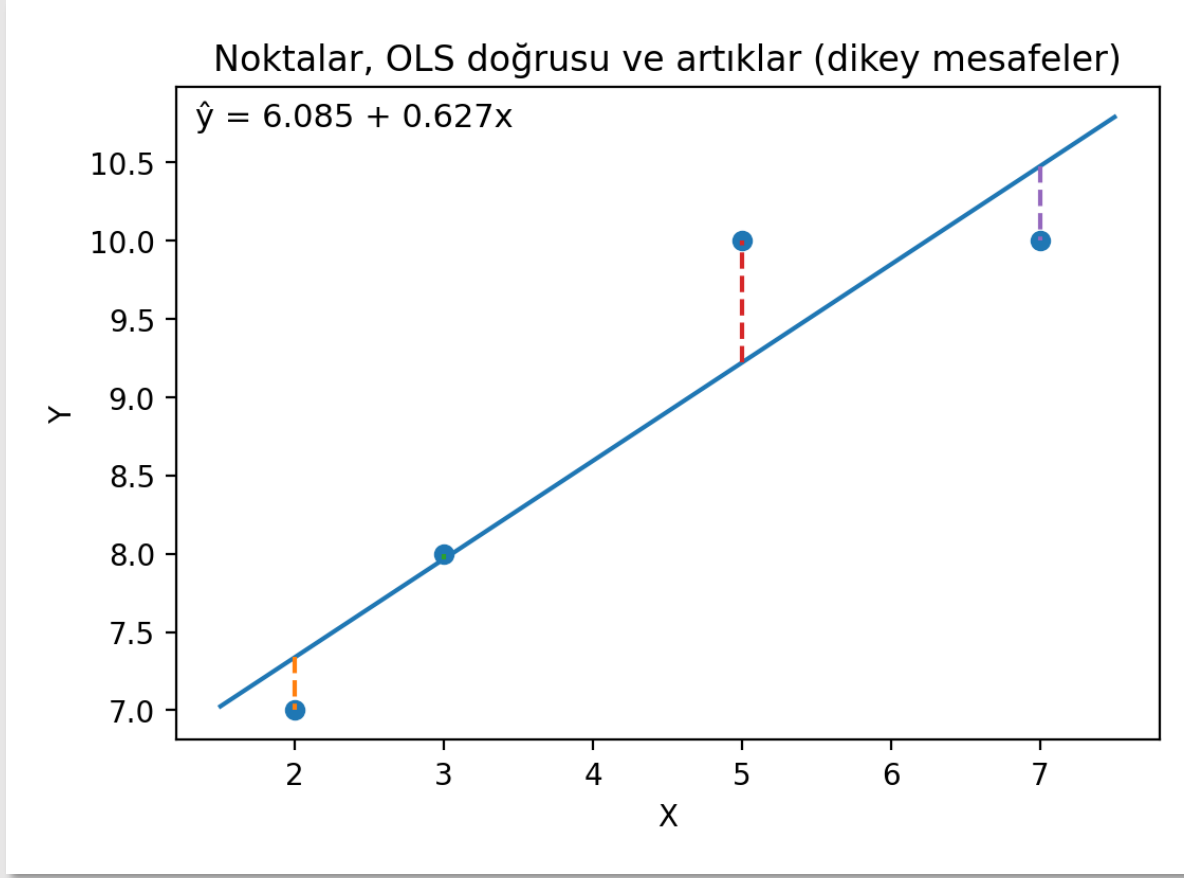
Doğrusal Regresyon



- Amacımız, **hata / kayıp fonksiyonunu en aza indiren** bir doğru bulmak.
- NEDEN doğrusal? Çünkü şunu varsayıyoruz: X'teki bir birim artış Y'yi sabit miktarda arttırıyor.
- OLS (En Küçük Kareler Yöntemi), **artıkların (hataların) kareleri toplamını** mümkün olduğunca küçültmeye çalışır.
- Klasik bir OLS modelini sklearn kullanarak çalıştırabiliriz.

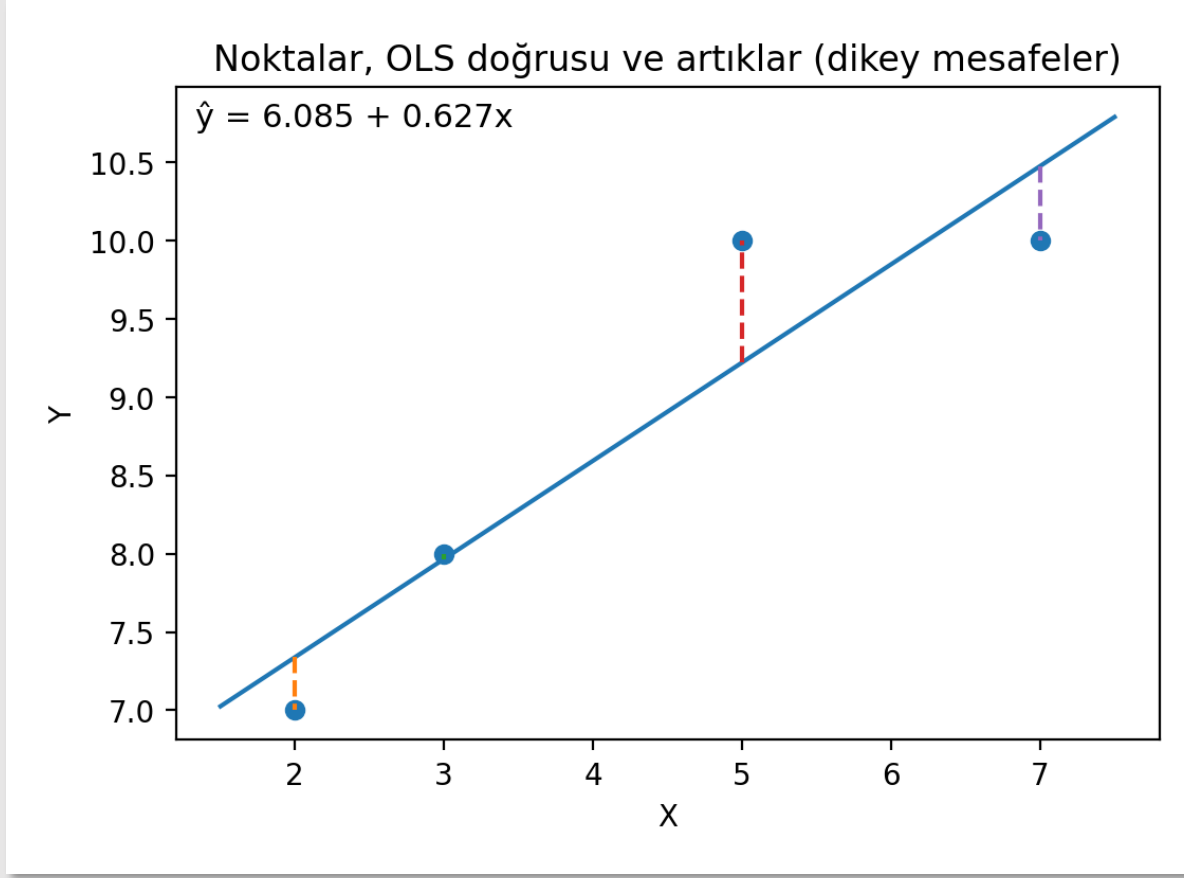
Doğrusal Regresyon

- Yaşları sırasıyla 2, 3, 5 ve 7 olan 4 kardeşin baş parmak boyları sırasıyla 7, 8, 10 ve 10 cm'dir.



Doğrusal Regresyon

- Yaşları sırasıyla 2, 3, 5 ve 7 olan 4 kardeşin baş parmak boyları sırasıyla 7, 8, 10 ve 10 cm'dir.



1 yaş artışı, yaklaşık 0.627 cm ekstra parmak boyu getirir.

Doğrunun eğimi: 0.627

Doğrunun Y eksenini kestiği değer: 6.085

Doğrusal Regresyon

X_i	Y_i
2	7
3	8
5	10
7	10

$$y_i = \hat{\beta}_0 + \hat{\beta}_1 \cdot X_i + \varepsilon_i$$

hata fonksiyonu

$$\sum_i \underbrace{(y_i - \hat{\beta}_0 - \hat{\beta}_1 \cdot X_i)}_{\varepsilon_i}^2 = \min$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \cdot \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\bar{y} = 8.75$$

$$\bar{x} = 4.25$$

$$\beta_0 = 6.085$$

$$\beta_1 = 0.627$$

X_i	Y_i	\hat{Y}_i	ε_i
2	7	7.3	-0.3
3	8	7.96	0.004
5	10	9.2	0.8
7	10	10.47	0.47

Doğrusal Regresyon

X_i	y_i
2	7
3	8
5	10
7	10

$$y_i = \hat{\beta}_0 + \hat{\beta}_1 \cdot X_i + \varepsilon_i$$

hata fonksiyonu

$$\sum_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 \cdot X_i)^2 = \min$$

ε_i

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \cdot \bar{X}$$

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{X})(y_i - \bar{y})}{\sum (x_i - \bar{X})^2}$$

$$\bar{y} = 8.75$$

$$\bar{X} = 4.25$$

$$\beta_0 = 6.085$$

$$\beta_1 = 0.627$$

X_i	y_i	\hat{y}_i	ε_i
2	7	7.3	-0.3
3	8	7.96	0.004
5	10	9.2	0.8
7	10	10.47	0.47

- Residual sum of square = RSS

$$\min_{\beta} \text{RSS}(\beta) = \sum_{i=1}^n (y_i - x_i^T \beta)^2$$

Düzenlileştiriciler (regularizers)

- OLS mükemmel bir "overfitting" örneğidir.
- Düzenlileştirmek, temelde katsayıları cezalandırmak anlamına gelir.

Ridge:

- Kayıp fonksiyonu = OLS kayıp fonksiyonu + Ridge ceza terimi
- **Alpha (α)** modelin karmaşıklığını kontrol eder
- $\alpha = 0$ olduğunda model klasik OLS'tir (overfitting riski yüksek)
- Çok yüksek α ise modeli fazla sadeleştirip **underfitting'e** yol açabilir

Lasso:

- Kayıp fonksiyonu = OLS kayıp fonksiyonu + Lasso ceza terimi
- **Alpha (α)** modelin karmaşıklığını kontrol eder
- $\alpha = 0$ olduğunda model klasik OLS'tir (overfitting riski yüksek)
- Çok yüksek α ise modeli aşırı sadeleştirerek **underfitting'e** yol açabilir
- Lasso, **en önemli özellikleri** ortaya çıkarmaya ve **görel olarak daha önemsiz** özellikleri bastırmaya / sıfırlamaya yardımcı olur.



Düzenlileştiriciler (regularizers)

- OLS mükemmel bir "overfitting" örneğidir.
- Düzenlileştirmek, temelde katsayıları cezalandırmak anlamına gelir.

Ridge:

- Kayıp fonksiyonu = OLS kayıp fonksiyonu + Ridge ceza terimi
- **Alpha (α)** modelin karmaşıklığını kontrol eder
- $\alpha = 0$ olduğunda model klasik OLS'tir (overfitting riski yüksek)
- Çok yüksek α ise modeli fazla sadeleştirip **underfitting'e** yol açabilir

Lasso:

- Kayıp fonksiyonu = OLS kayıp fonksiyonu + Lasso ceza terimi
- **Alpha (α)** modelin karmaşıklığını kontrol eder
- $\alpha = 0$ olduğunda model klasik OLS'tir (overfitting riski yüksek)
- Çok yüksek α ise modeli aşırı sadeleştirerek **underfitting'e** yol açabilir
- Lasso, **en önemli özellikleri** ortaya çıkarmaya ve **görel olarak daha önemsiz** özellikleri bastırmaya / sıfırlamaya yardımcı olur.

$$\min_{\beta} \text{RSS}(\beta) = \sum_{i=1}^n (y_i - x_i^T \beta)^2$$

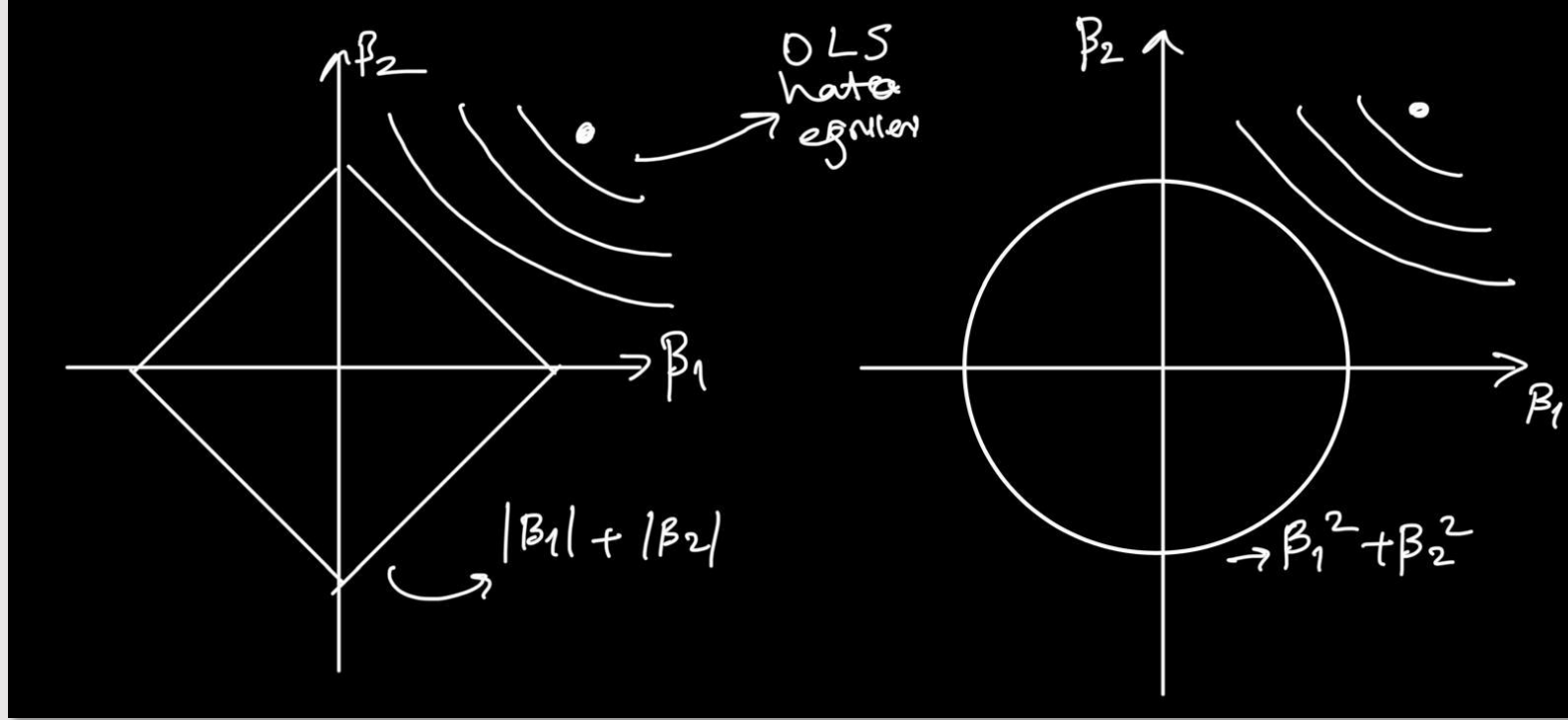
- Lasso (L1 ceza):

$$\min_{\beta} \text{RSS}(\beta) + \lambda \sum_{j=1}^p |\beta_j|$$

- Ridge (L2 ceza):

$$\min_{\beta} \text{RSS}(\beta) + \lambda \sum_{j=1}^p \beta_j^2$$

Düzenlileştiriciler (regularizers)



Lasso

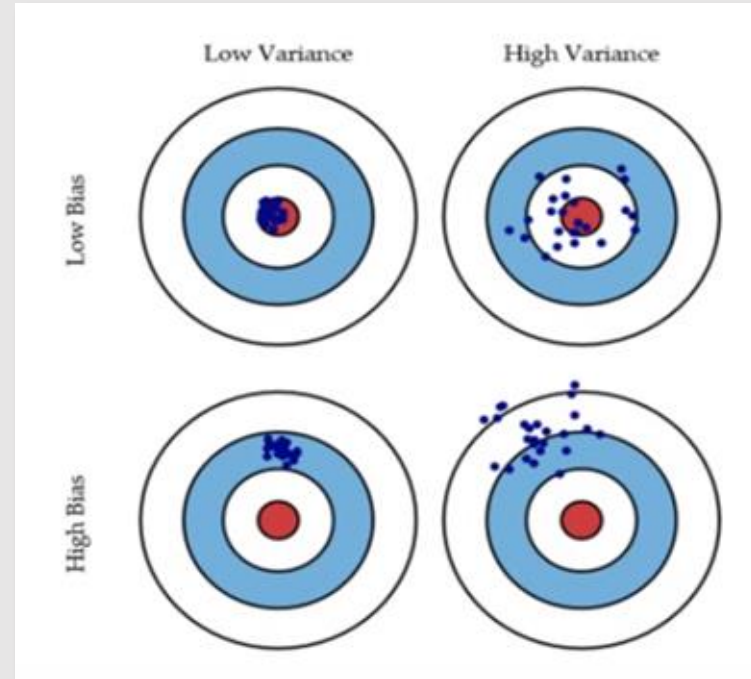
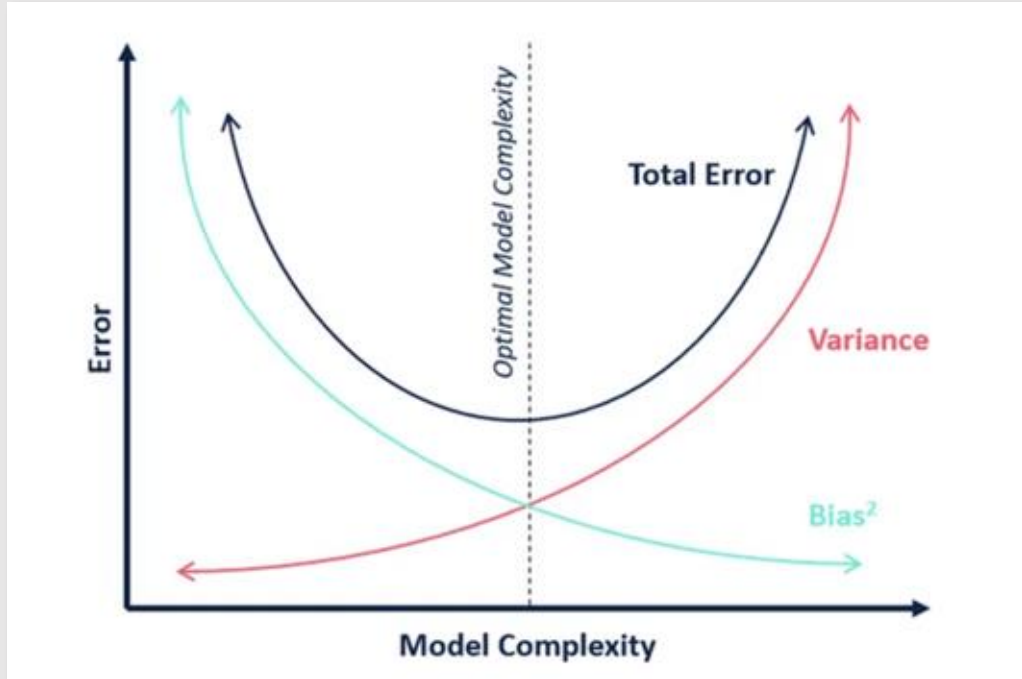
Ridge

Lasso'nun mutlak değerli toplamı gösteren düz çizgileriyle OLS'in eliptik hata eğrilerinin, Beta'lardan bazılarının sıfır olduğu yerlerde kesişmesi matematiksel olarak mümkündür. Ancak bu durum Ridge için geçerli değildir. Pratikteyse bu ikisi genelde birlikte kullanılır (bkz: elastik net modelleri).

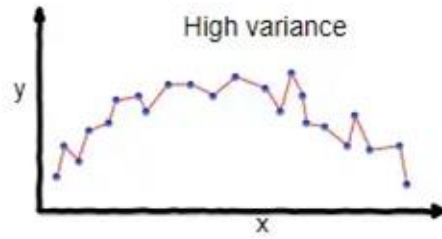
$$\min_{\beta} \text{RSS}(\beta) = \sum_{i=1}^n (y_i - x_i^T \beta)^2$$

- Lasso (L1 ceza):
$$\min_{\beta} \text{RSS}(\beta) + \lambda \sum_{j=1}^p |\beta_j|$$
- Ridge (L2 ceza):
$$\min_{\beta} \text{RSS}(\beta) + \lambda \sum_{j=1}^p \beta_j^2$$

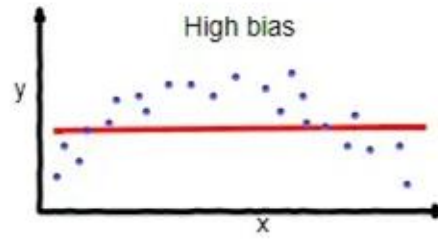
Düzenleştiriciler (regularizers)



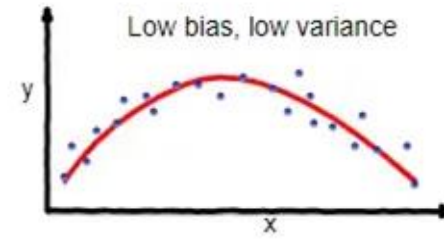
Düzenlileştiriciler (regularizers)



overfitting



underfitting



Good balance

Düzenlileştiriciler (regularizers)

Bir şehir-ay panelinde **protesto sayısını** (Y) tahmin etmek için 20 değişkenli bir regresyon kuruluyor (işsizlik, enflasyon, seçim ayı, önceki ay protesto, sendikalaşma vb.). **OLS** ile eğitim datasındaki açıklama gücü (örn: R^2) oldukça yüksek ama test aşamasındaki performans çok düşük. Bu tabloyu “genelleme” ve yanlılık-varyans dengesi (**bias-variance trade-off**) açısından en tutarlı şekilde yorumlayıp en uygun müdahaleyi seçmek istiyoruz. Aşağıdakilerden hangisi en doğru açıklama + çözüm paketidir?

- A) Bu **yüksek bias** (underfitting) göstergesidir; çözüm olarak **λ 'yı artırıp** Ridge veya Lasso ile daha esnek bir model kurmak gerekir.
- B) Bu **yüksek variance** (overfitting) göstergesidir; çözüm olarak **Ridge veya Lasso** ile ceza ekleyip **λ 'yı cross-validation ile seçmek** gerekir: variance düşer, bias biraz artabilir, ama test performansı iyileşir.
- C) OLS zaten “en iyi” tahmini verdiği için problem veri toplamadadır. Çözüm olarak **test setini train'e katıp** tek seferde fit etmek gerekir.
- D) Lasso ve Ridge sadece sınıflandırmada işe yarar. Regresyonda bias-variance ile ilişkileri yoktur.



1. C. Performans Testleri ve Model Geliştirme

Modelin Performansı

Modelin ne kadar iyi?

Doğruluğu (**accuracy**) ne kadar?

- Spam sınıflandırması örneği: E-postaların %1'i spam, %99'u gerçek (spam değil).
 - Eğer model tüm e-postaları "**gerçek**" diye sınıflandırırsa, model %99 doğru görünür!
- Buna sınıf dengesizliği (**class imbalance**) denir.
Bu yüzden daha nüanslı/ayırt edici metriklere ihtiyacımız var.

Karmaşıklık matrisi (Confusion matrix):

- Accuracy** (Doğruluk) = $(TP + TN) / (TP + TN + FP + FN)$
- Precision** (Kesinlik / Pozitif öngörü değeri) = $TP / (TP + FP)$
- Recall** (Duyarlılık / Yakalama oranı) = $TP / (TP + FN)$
- F1-score** = $2 \times (Precision \times Recall) / (Precision + Recall)$

Kısaltmalar: **TP**=Doğru Pozitif, **TN**=Doğru Negatif, **FP**=Yanlış Pozitif, **FN**=Yanlış Negatif.

	Predicted: Spam Email	Predicted: Real Email
Actual: Spam Email	True Positive	False Negative
Actual: Real Email	False Positive	True Negative

Modelin Performansı

Toplumsal hareketler üzerine çalışan bir araştırmacı, sosyal medyadan “**eylem var!**” diye otomatik alarm üreten bir sistem kullanıyor. Sistem protestoların büyük kısmını doğru yakalıyor gibi görünüyor ama ekip şunu fark ediyor: “Alarm verdiğimiz her 10 durumun sadece 4’ü gerçekten protesto; öte yandan, gerçekte yaşanan protestoları ise neredeyse tamamen yakalayabiliyoruz.”

Bu durumda aşağıdaki yorumlardan hangisi en doğrudur?

1. Precision yüksek, Recall düşük, Precision odaklı bir model kurulmuş.
2. Precision da düşük, Recall da düşüktür. F1 da bu sebeple düşük olur.
3. Accuracy yüksekse precision/recall’a bakmaya gerek yoktur.
4. Precision düşük, Recall yüksek, Recall odaklı bir model kurulmuş.



Lojistik Regresyon ve ROC eğrisi

- Lojistik regresyon, sınıflandırmaya benzer bir mantıkla çalışır: her gözlem için bir olasılık tahmin eder (ör. $p > 0.5$ ise etiketi 1 yap).
- ROC eğrisi (eşik değeri değiştirerek)

- **FPR (False Positive Rate / Yanlış Pozitif Oranı):**

0'a gitmesi gerektiği halde 1'e gidenlerin oranı

(Yani "negatifleri ne kadar yanlış alarm vererek pozitif ilan ediyorum?")

- **TPR (True Positive Rate / Doğru Pozitif Oranı, Recall / Duyarlılık):**

1'e gitmesi gerekenler içerisinde 1'e gidebilmiş olanların oranı

(Yani "pozitifleri ne kadarını yakalayabiliyorum?")



Lojistik Regresyon ve ROC eğrisi

- Lojistik regresyon, sınıflandırmaya benzer bir mantıkla çalışır: her gözlem için bir olasılık tahmin eder (ör. $p > 0.5$ ise etiketi 1 yap).
- ROC eğrisi (eşik değeri değiştirerek)

- FPR (False Positive Rate / Yanlış Pozitif Oranı):**

0'a gitmesi gerektiği halde 1'e gidenlerin oranı

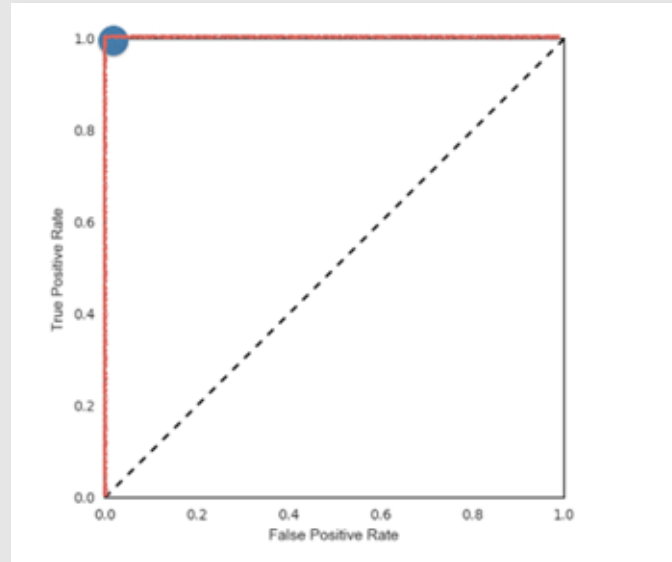
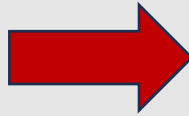
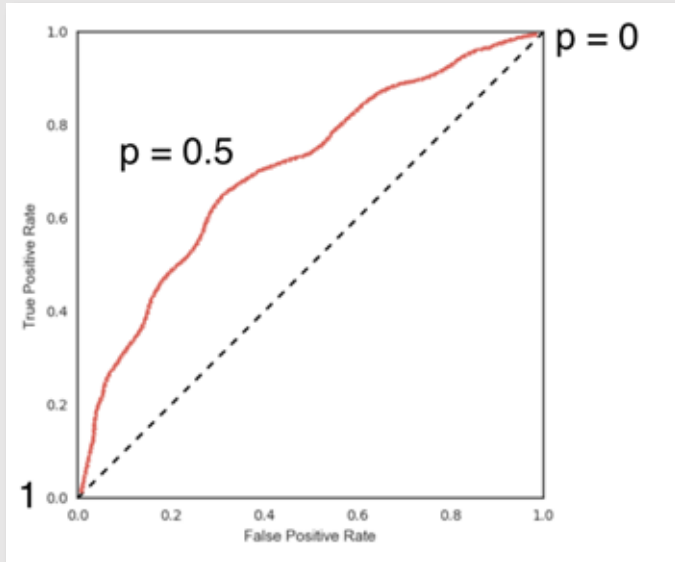
(Yani "negatifleri ne kadar yanlış alarm vererek pozitif ilan ediyorum?")

- TPR (True Positive Rate / Doğru Pozitif Oranı, Recall / Duyarlılık):**

1'e gitmesi gerekenler içerisinde 1'e gidebilmiş olanların oranı

(Yani "pozitifleri ne kadarını yakalayabiliyorum?")

ROC eğrisinin altındaki alan (AUC) ne kadar büyükse, model o kadar iyidir.



Lojistik Regresyon ve ROC eğrisi

$A = 0$ (negative)

$B = 1$ (positive)

Gerçek pozitif sayısı $P = 5$

Gerçek negatif sayısı $N = 5$

İndeks:	1	2	3	4	5	6	7	8	9	10
Gerçek sınıf:	A	A	B	A	A	B	A	B	B	B
Modelin tahmin ettiği olasılık:	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1

Lojistik Regresyon ve ROC eğrisi

$A = 0$ (negative)

$B = 1$ (positive)

Gerçek pozitif sayısı $P = 5$

Gerçek negatif sayısı $N = 5$

İndeks:	1	2	3	4	5	6	7	8	9	10
Gerçek sınıf:	A	A	B	A	A	B	A	B	B	B
Modelin tahmin ettiği olasılık:	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1

Threshold (t)	Pozitif tahmin edilenler (score $\geq t$)	TP	FP	TPR	FPR
>1.0	—	0	0	0.00	0.00
≥ 1.0	10	1	0	0.20	0.00
≥ 0.9	9–10	2	0	0.40	0.00
≥ 0.8	8–10	3	0	0.60	0.00
≥ 0.7	7–10	3	1	0.60	0.20
≥ 0.6	6–10	4	1	0.80	0.20
≥ 0.5	5–10	4	2	0.80	0.40
≥ 0.4	4–10	4	3	0.80	0.60
≥ 0.3	3–10	5	3	1.00	0.60
≥ 0.2	2–10	5	4	1.00	0.80
≥ 0.1	1–10	5	5	1.00	1.00



Lojistik Regresyon ve ROC eğrisi

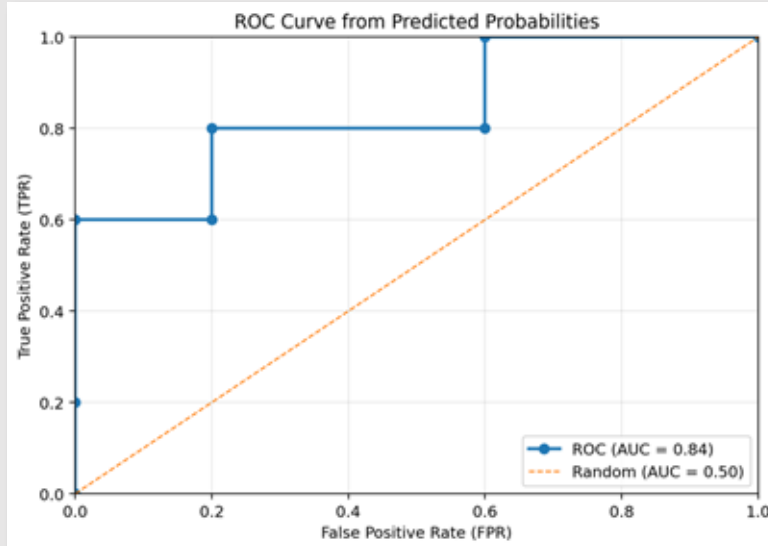
$A = 0$ (negative)

$B = 1$ (positive)

Gerçek pozitif sayısı $P = 5$

Gerçek negatif sayısı $N = 5$

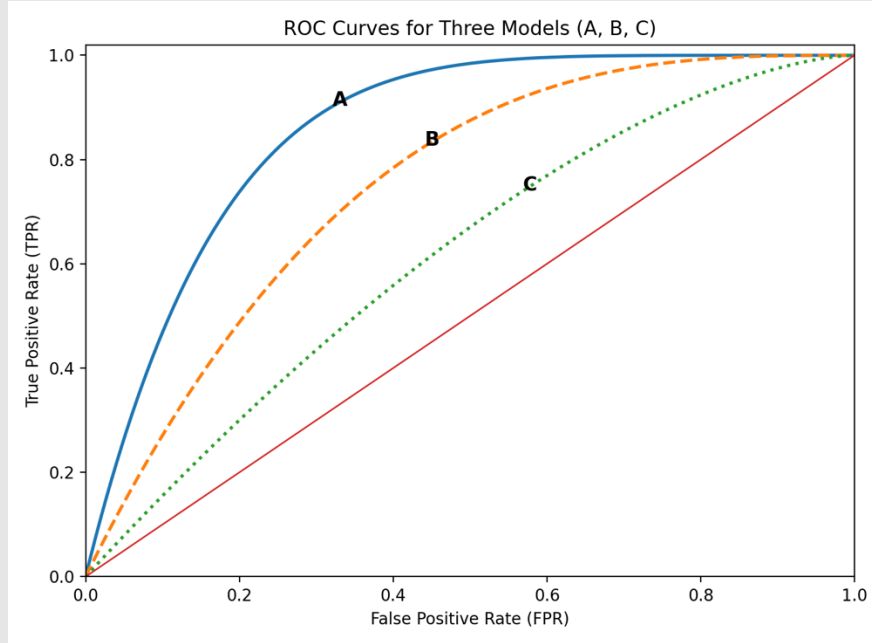
İndeks:	1	2	3	4	5	6	7	8	9	10
Gerçek sınıf:	A	A	B	A	A	B	A	B	B	B
Modelin tahmin ettiği olasılık:	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1



Threshold (t)	Pozitif tahmin edilenler (score $\geq t$)	TP	FP	TPR	FPR
>1.0	—	0	0	0.00	0.00
≥ 1.0	10	1	0	0.20	0.00
≥ 0.9	9–10	2	0	0.40	0.00
≥ 0.8	8–10	3	0	0.60	0.00
≥ 0.7	7–10	3	1	0.60	0.20
≥ 0.6	6–10	4	1	0.80	0.20
≥ 0.5	5–10	4	2	0.80	0.40
≥ 0.4	4–10	4	3	0.80	0.60
≥ 0.3	3–10	5	3	1.00	0.60
≥ 0.2	2–10	5	4	1.00	0.80
≥ 0.1	1–10	5	5	1.00	1.00

Lojistik Regresyon ve ROC eğrisi

Şehir-ay düzeyinde bir veri setinde, gerçekleşen bir işçi grevinin başarıya ulaşp ulaşmayacağını $Y \in \{0, 1\}$ ile kodlanıyor (1=karşılandı, 0=karşılanmadı). X 'lerde ekonomik göstergeler, önceki ay mobilizasyon düzeyi, seçim ayı, baskı göstergeleri vb. var. **Lojistik regresyon** modeli her gözlem için $P(Y = 1)$ üretiyor. Aşağıdaki A, B ve C (**ROC**) eğrilerinden hangisinde gösterilen modelin performansı daha yüksektir?



Hiperparametre Ayarı

- **Ridge** ve **Lasso** regresyonu / α seçimi
- **k-En Yakın Komşu** (k-NN) / $n_neighbors$ seçimi
 - **k** ve **α** gibi parametreler hiperparametredir.

Doğru hiperparametre nasıl seçilir?

- Birçok hiperparametre değeri dene; her birini ayrı ayrı eği; performanslarını karşılaştı; en iyi performans vereni seç.
 - **Grid search** (çapraz doğrulama ile)
 - **Randomized search** (çapraz doğrulama ile)



Ön İşleme ve Boru Hattı (Pipeline)

- **Gerçek dünyadaki veri**, en iyi modeli hemen çalıştırabilmen için **hazır gelmez**.
- **Kategorik değişkenleri** dummy (one-hot) değişkenlere çevir.
- **Eksik veriyi yönet**
 - Satırları/gözlemleri silersen: **veri kaybı**
 - 0 yüzden: **impute (tamamlama) yapabilirsin!**
 - Eksik olmayan değerlerin ortalamasıyla doldurma
 - Zaman boyutu varsa doğrusal (linear) imputasyon
 - Scikit-learn: Imputer

```
from sklearn.preprocessing import Imputer  
imp = Imputer(missing_values='NaN', strategy='mean', axis=0)  
imp.fit(X)  
X = imp.transform(X)
```

Pipeline: Hepsini tek seferde yap!

Impute, normalize et, çalıştır, fit et, tahmin et, test et...



2. Gözetimsiz Öğrenme

Gözetimsiz Öğrenme

- Gözetimli öğrenme (supervised learning) belirli bir tahmin görevi için örüntüler bulurken, gözetimsiz öğrenme (unsupervised learning) belirli bir tahmin hedefi olmadan verinin içindeki örüntüleri keşfeder.
- Diyelim ki elimizde sadece özellikler (target/etiket olmadan) var:
- Gözetimsiz bir model çalıştırırız ve model, veride zaten var olan doğal kümeleri (clusters) kendisi keşfeder.

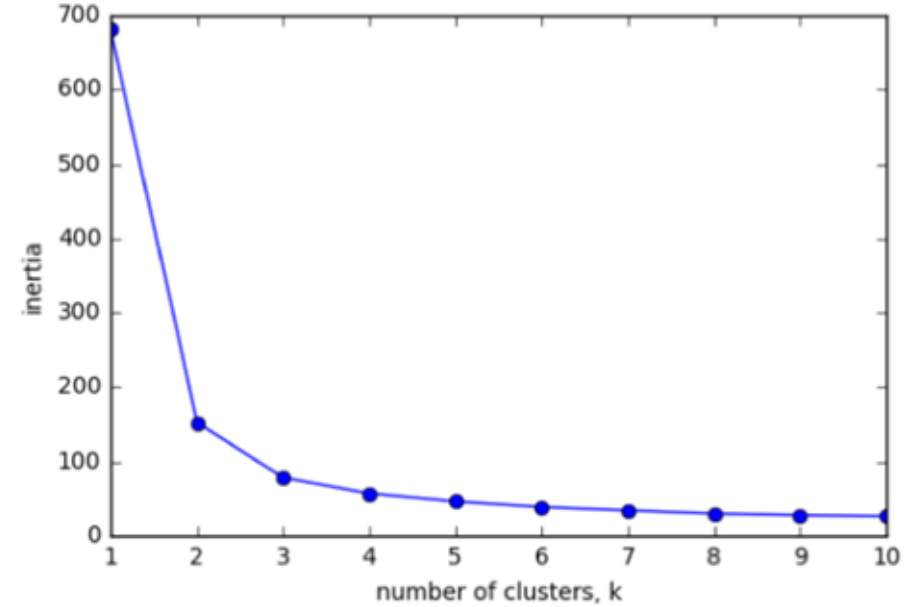
k-means kümeleme (clustering)

- kNN sınıflandırmaya benzer; ama burada etiket yoktur.
- Model etiketleri/küme atamalarını kendi kendine üretir.



Kümeleme

- **Kümelerimizin kalitesini ölçmek:**
Inertia!
- İyi bir kümeleme, **sıkı (tight) kümelere** sahiptir (bu genelde **daha fazla küme** ve **daha düşük inertia** demektir), ama **çok fazla küme** de istemeyiz.
- Bu yüzden (sıklıkla önerildiği gibi) **dirsek (elbow) noktasını** seçeriz.



Kümeleme

Bir araştırma ekibi, tweetlerden çıkarılan metin-özellikleriyle (TF-IDF/embedding) belli bir anahtar kelimeyi içeren (örn: **boykot**) paylaşımlarını otomatik olarak temalara ayırmak istiyor. Etiket (ground truth) yok; amaç “benzer içerikleri” aynı grupta toplamak. k-means’i deniyorlar ve **k=2’den k=10’a** kadar farklı değerlerle modeli çalıştırıp **inertia** değerlerini izliyorlar. Aşağıdaki seçeneklerden hangisi en doğru yorumdur?

- a) Inertia, k arttıkça genelde artar; bu yüzden en büyük inertia veren k seçilmelidir.
- b) Inertia, k arttıkça genelde azalır; “dirsek (elbow)” yöntemi, inertia’daki azalışın belirgin biçimde yavaşladığı noktayı makul k olarak seçmeye yardım eder.
- c) k-means gözetimli bir yöntemdir; doğru k, test setinde accuracy ile bulunur.
- d) Inertia ile hiçbir şey seçilemez; k her zaman gözlem sayısına eşit olmalıdır.



Neler var neler?

- **Naive Bayes Classification**
- **Support Vector Machines**
- **Decision Trees and Random Forests**
- **Principal Component Analysis**
- Manifold Learning
- Gaussian Mixture Models
- Kernel Density Estimation

Derin Öğrenme

- **Neural Networks**

