

Refleks Temelli Modeller El Kitabı VIP

Afshine AMIDI ve Shervine AMIDI

September 14, 2019

Yavuz Kömeçoğlu ve Ayyüce Kızrak tarafından çevrilmiştir

Doğrusal öngörücüler

Bu bölümde, girdi-çıktı çiftleri olan örneklerden geçerek, deneyim ile gelişebilecek refleks-temelli modelleri göreceğiz.

□ **Öznitelik vektörü** – x girişinin öznitelik vektörü $\phi(x)$ olarak not edilir ve şöyledir:

$$\phi(x) = \begin{bmatrix} \phi_1(x) \\ \vdots \\ \phi_d(x) \end{bmatrix} \in \mathbb{R}^d$$

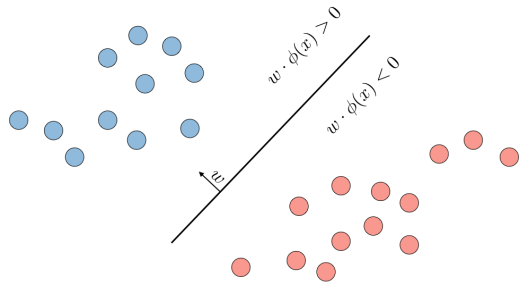
□ **Puan** – Bir örneğin $s(x, w)$ si ni $(\phi(x), y) \in \mathbb{R}^d \times \mathbb{R}$, $w \in \mathbb{R}^d$ doğrusal ağırlık modeline bağlı olarak:

$$s(x, w) = w \cdot \phi(x)$$

Sınıflandırma

□ **Doğrusal sınıflandırıcı** – Bir ağırlık vektörü $w \in \mathbb{R}^d$ ve bir öznitelik vektörü $\phi(x) \in \mathbb{R}^d$ verildiğinde, ikili doğrusal sınıflandırıcı f_w şöyle verilir:

$$f_w(x) = \text{sign}(s(x, w)) = \begin{cases} +1 & \text{ise } w \cdot \phi(x) > 0 \\ -1 & \text{ise } w \cdot \phi(x) < 0 \\ ? & \text{ise } w \cdot \phi(x) = 0 \end{cases}$$



□ **Marj** – $(\phi(x), y) \in \mathbb{R}^d \times \{-1, +1\}$ örneğinin $m(x, y, w) \in \mathbb{R}$ marjları $w \in \mathbb{R}^d$ doğrusal ağırlık modeliyle ilişkili olarak, tahminin güvenilirliği ölçülür: daha büyük değerler daha iyidir. Şöyle ifade edilir:

$$m(x, y, w) = s(x, w) \times y$$

Bağlanım

□ **Doğrusal bağlanım** – $w \in \mathbb{R}^d$ bir ağırlık vektörü ve bir öznitelik vektörü $\phi(x) \in \mathbb{R}^d$ verildiğinde, f_w olarak belirtilen ağırlıkların doğrusal bir bağlanım (linear regression) çıktısı şöyle verilir:

$$f_w(x) = s(x, w)$$

□ **Artık** – Artık (residual) $\text{res}(x, y, w) \in \mathbb{R}$, $f_w(x)$ tahmininin y hedefini aştığı miktar olarak tanımlanır:

$$\text{res}(x, y, w) = f_w(x) - y$$

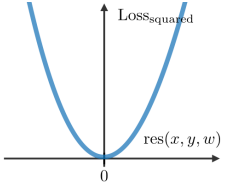
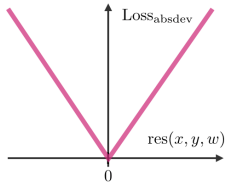
Kayıp minimizasyonu

□ **Kayıp fonksiyonu** – Kayıp fonksiyonu $\text{Loss}(x, y, w)$, x girişinden y çıktısının öngörme görevindeki model ağırlıkları ile ne kadar mutsuz olduğumuzu belirler. Bu değer eğitim sürecinde en aza indirmek istediğimiz bir miktar.

□ **Sınıflandırma durumu** – Doğru etiket $y \in \{-1, +1\}$ değerinin x örneğinin doğrusal ağırlık w modeliyle sınıflandırılması $f_w(x) \triangleq \text{sign}(s(x, w))$ belirleyicisi ile yapılabilir. Bu durumda, sınıflandırma kalitesini ölçen bir fayda ölçütü $m(x, y, w)$ marjı ile verilir ve aşağıdaki kayıp fonksiyonlarıyla birlikte kullanılabilir:

| Ad | Sıfır-bir kayıp | Menteşe kaybı | Lojistik kaybı |
|------------------------|-----------------------------|---------------------------|-----------------------------|
| $\text{Loss}(x, y, w)$ | $1_{\{m(x, y, w) \leq 0\}}$ | $\max(1 - m(x, y, w), 0)$ | $\log(1 + e^{-m(x, y, w)})$ |
| Örnekleme | | | |

□ **Regresyon durumu** – Doğru etiket $y \in \mathbb{R}$ değerinin x örneğinin bir doğrusal ağırlık modeli w ile öngörülmesi $f_w(x) \triangleq s(x, w)$ öngörüsü ile yapılabilir. Bu durumda, regresyonun kalitesini ölçen bir fayda ölçütü $\text{res}(x, y, w)$ marjı ile verilir ve aşağıdaki kayıp fonksiyonlarıyla birlikte kullanılabilir:

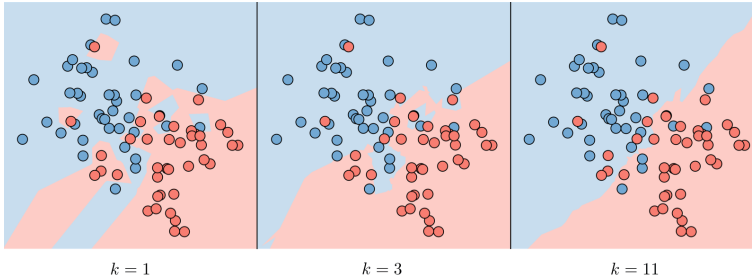
| Ad | Kareler kaybı | Mutlak sapma kaybı |
|------------------------|---|---|
| $\text{Loss}(x, y, w)$ | $(\text{res}(x, y, w))^2$ | $ \text{res}(x, y, w) $ |
| Görselleştirme |  |  |

□ **Kayıp minimize etme çerçevesi** – Bir modeli eğitmek için, eğitim kaybını en aza indirmek istiyoruz:

$$\text{TrainLoss}(w) = \frac{1}{|\mathcal{D}_{\text{train}}|} \sum_{(x, y) \in \mathcal{D}_{\text{train}}} \text{Loss}(x, y, w)$$

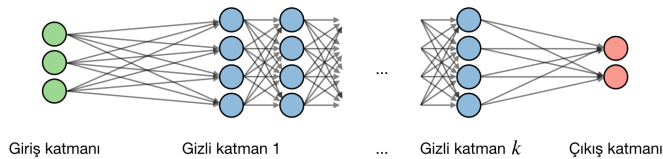
Doğrusal olmayan öngörücüler

□ **k -en yakın komşu** – Yaygın olarak k -NN (k -nearest neighbors) olarak bilinen k -en yakın komşu algoritması, bir veri noktasının tepkisinin eğitim kümesinden k komşularının yapısı tarafından belirlendiği parametrik olmayan bir yaklaşımdır. Hem sınıflandırma hem de regresyon ayarlarında kullanılabilir.



Not: k parametresi ne kadar yüksekse, önyargı (bias) o kadar yüksek ve k parametresi ne kadar düşüğe, varyans o kadar yüksek olur.

□ **Yapay sinir ağları** – Neural networks are a class of models that are built with layers. Commonly used types of neural networks include convolutional and recurrent neural networks. The vocabulary around neural networks architectures is described in the figure below:



i ağın i . katmanı ve j . katmanın j . gizli birimi olacak şekilde aşağıdaki gibi ifade edilir:

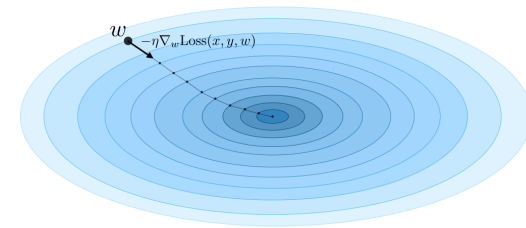
$$z_j^{[i]} = w_j^{[i]T} x + b_j^{[i]}$$

w , b , x , z değerlerinin sırasıyla nöronun ağırlık, önyargı (bias), girdi ve aktive edilmemiş çıkışı olarak ifade eder.

Stokastik gradyan inişi

□ **Gradyan inişi** – $\eta \in \mathbb{R}$ öğrenme oranını (aynı zamanda adım boyutu olarak da bilinir) dikkate alınarak, gradyan (bayır) inişine ilişkin güncelleme kuralı, öğrenme oranı ve $\text{Loss}(x, y, w)$ kayıp fonksiyonu ile aşağıdaki şekilde ifade edilir:

$$w \leftarrow w - \eta \nabla_w \text{Loss}(x, y, w)$$



□ **Stokastik güncellemeler** – Stokastik gradyan inişi (SGİ/SGD, stochastic gradient descent), bir seferde bir eğitim örneğinin $(\phi(x), y) \in \mathcal{D}_{\text{train}}$ parametrelerini günceller. Bu yöntem bazen gürültülü, ancak hızlı güncellemeler yol açar.

□ **Yığın güncellemeler** – Yığın gradyan inişi (YĞİ/BGD, batch gradient descent), bir seferde bir grup örnek (örneğin, tüm eğitim kümesi) parametrelerini günceller. Bu yöntem daha yüksek bir hesaplama maliyetiyle kararlı güncelleme talimatlarını hesaplar.

İnce ayar modelleri

□ **Hipotez sınıfı** – Bir hipotez sınıfı \mathcal{F} , sabit bir $\phi(x)$ ve değişken w ile olası öngörücü kümesidir:

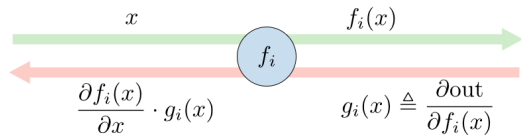
$$\mathcal{F} = \{f_w : w \in \mathbb{R}^d\}$$

□ **Lojistik fonksiyon** – Ayrıca sigmoid fonksiyon olarak da adlandırılan lojistik fonksiyon (sigmoid function) σ , şöyle tanımlanır:

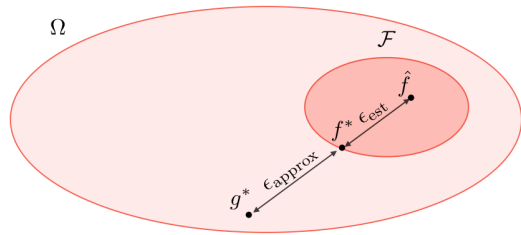
$$\forall z \in]-\infty, +\infty[, \quad \sigma(z) = \frac{1}{1 + e^{-z}}$$

Not: $\sigma'(z) = \sigma(z)(1 - \sigma(z))$ şeklinde ifade edilir.

□ **Geri yayılım** – İleriye geçiş, i 'de yer alan alt ifadenin değeri olan f_i ile yapılırken, geriye doğru geçiş $g_i = \frac{\partial \text{out}}{\partial f_i}$ aracılığıyla yapılır ve f_i 'nin çıkışı nasıl etkilendiğini gösterir.



□ **Yaklaşım ve kestirim hatası** – Yaklaşım hatası ϵ_{approx} , \mathcal{F} tüm hipotez sınıfının hedef öngörücü g^* ne kadar uzak olduğunu gösterirken, kestirim hatası ϵ_{est} öngörücüsü \hat{f} , \mathcal{F} hipotez sınıfının en iyi yordayıcısı f^* 'ya göre ne kadar iyi olduğunu gösterir.



□ **Düzenleştirme** – Düzenleştirme (regularization) prosedürü, modelin verilerin aşırı öğrenmesinden kaçınmayı amaçlar ve böylece yüksek değişkenlik sorunlarıyla ilgilenir. Aşağıdaki tablo, yaygın olarak kullanılan düzenleştirme tekniklerinin farklı türlerini özetlemektedir:

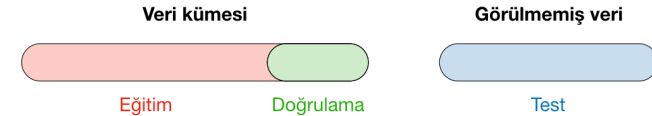
| LASSO | Ridge | Elastic Net |
|---|--|--|
| <ul style="list-style-type: none"> - Katsayıları 0'a düşürür - Değişken seçimi için iyi | Katsayıları daha küçük yapar | Değişken seçimi ile küçük katsayılar arasında ödünleşim |
| <p>$\theta _1 \leq 1$</p> | <p>$\theta _2 \leq 1$</p> | <p>$(1-\alpha) \theta _1 + \alpha \theta _2^2 \leq 1$</p> |
| $\dots + \lambda \theta _1$ $\lambda \in \mathbb{R}$ | $\dots + \lambda \theta _2^2$ $\lambda \in \mathbb{R}$ | $\dots + \lambda \left[(1-\alpha) \theta _1 + \alpha \theta _2^2 \right]$ $\lambda \in \mathbb{R}, \alpha \in [0,1]$ |

□ **Hiperparametreler** – Hiperparametreler öğrenme algoritmasının özellikleridir ve öznitelikler dahildir, λ normalizasyon parametresi, yinleme sayısı T , adım büyüklüğü η , vb.

□ **Kümeler** – Bir model seçerken, veriyi aşağıdaki gibi 3 farklı parçaya ayırırız:

| Eğitim kümesi | Doğrulama kümesi | Test kümesi |
|---|---|--|
| <ul style="list-style-type: none"> - Model eğitilir - Veri kümesinin genellikle %80'i | <ul style="list-style-type: none"> - Model değerlendirilir - Veri kümesinin genellikle %20'si - Ayrıca tutma veya geliştirme kümesi olarak da adlandırılır | <ul style="list-style-type: none"> - Model tahminlerini verir - Görünmeyen veriler |

Model seçildikten sonra, tüm veri kümesi üzerinde eğitilir ve görünmeyen test kümesinde test edilir. Bunlar aşağıdaki şekilde gösterilmektedir:



Gözetimsiz Öğrenme

Gözetimsiz öğrenme yöntemlerinin sınıfı, zengin gizli yapılar sahip olabilecek verilerin yapısını keşfetmeyi amaçlamaktadır.

k-ortalama

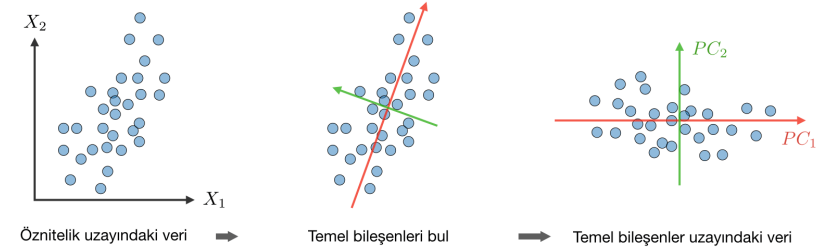
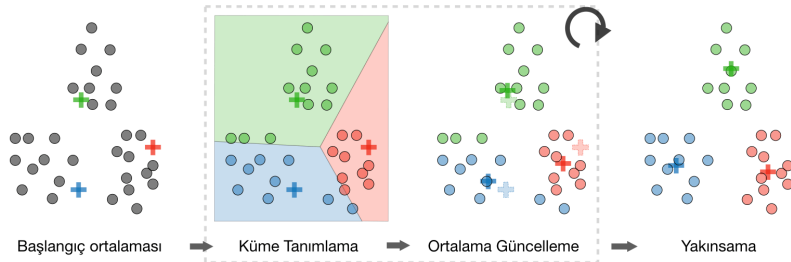
□ **Kümeleme** – $\mathcal{D}_{\text{train}}$ giriş noktalarından oluşan bir eğitim kümesi göz önüne alındığında, kümeleme (clustering) algoritmasının amacı, her bir $\phi(x_i)$ noktasını $z_i \in \{1, \dots, k\}$ kümesine atamaktır.

□ **Amaç fonksiyonu** – Ana kümeleme algoritmalarından biri olan k -ortalama için kayıp fonksiyonu şöyle ifade edilir:

$$\text{Loss}_{k\text{-means}}(x, \mu) = \sum_{i=1}^n ||\phi(x_i) - \mu_{z_i}||^2$$

□ **Algoritma** – Küme merkezlerini $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^n$ kümesini rasgele başlattıktan sonra, k -ortalama algoritması yakınsayana kadar aşağıdaki adımı tekrarlar:

$$z_i = \arg \min_j ||\phi(x_i) - \mu_j||^2 \quad \text{and} \quad \mu_j = \frac{\sum_{i=1}^m 1_{\{z_i=j\}} \phi(x_i)}{\sum_{i=1}^m 1_{\{z_i=j\}}}$$



Temel Bileşenler Analizi

□ **Özdeğer, özvektör** – Bir $A \in \mathbb{R}^{n \times n}$ matrisi verildiğinde, $z \in \mathbb{R}^n \setminus \{0\}$ olacak şekilde bir vektör varsa λ , A 'nın bir öz değeri olduğu söylenir, aşağıdaki gibi ifade edilir:

$$Az = \lambda z$$

□ **Spektral teoremi** – $A \in \mathbb{R}^{n \times n}$ olsun. A simetrik ise, o zaman A gerçek ortogonal matris $U \in \mathbb{R}^{n \times n}$ olacak şekilde köşegenleştirilebilir. $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ formülü dikkate alınarak aşağıdaki gibi ifade edilir:

$$\exists \Lambda \text{ diagonal, } A = U\Lambda U^T$$

Not: en büyük özdeğerle ilişkilendirilen özvektör, A matrisinin temel özvektörüdür.

□ **Algoritma** – Temel Bileşenler Analizi (PCA, principal component analysis) prosedürü, verilerin varyansını en üst düzeye çıkararak k boyutlarına indirgeyen bir boyut küçültme tekniğidir:

- **Adım 1:** Verileri ortalama 0 ve 1 standart sapma olacak şekilde normalize edin.

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{\sigma_j} \quad \text{where} \quad \mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)} \quad \text{and} \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

- **Adım 2:** Hesaplama $\Sigma = \frac{1}{m} \sum_{i=1}^m \phi(x_i) \phi(x_i)^T \in \mathbb{R}^{n \times n}$, ki bu, gerçek özdeğerlerle simetrik tir.
- **Adım 3:** Hesaplama $u_1, \dots, u_k \in \mathbb{R}^n$ k 'nin ortogonal ana özvektörleri, yani k en büyük özdeğerlerin ortogonal özvektörleri.
- **Adım 4:** $\text{span}_{\mathbb{R}}(u_1, \dots, u_k)$ 'daki verilerin izdüşümünü al. Bu prosedür, tüm k boyutlu uzaylar arasındaki farkı en üst düzeye çıkarır.