

Gözetimli Öğrenme El Kitabı VIP

Afshine AMIDI ve Shervine AMIDI

April 30, 2019

Başak Buluz ve Ayyüce Kızrak tarafından çevrilmiştir

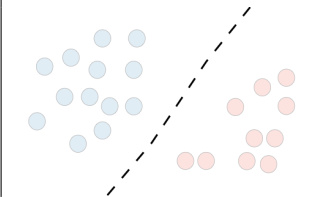
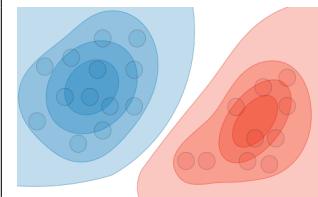
Gözetimli Öğrenmeye Giriş

$\{y^{(1)}, \dots, y^{(m)}\}$ çıktı kümesi ile ilişkili olan $\{x^{(1)}, \dots, x^{(m)}\}$ veri noktalarının kümesi göz önüne alındığında, y 'den x 'i nasıl tahmin edebileceğimizi öğrenen bir sınıflandırıcı tasarlamak istiyoruz.

□ **Tahmin türü** – Farklı tahmin modelleri aşağıdaki tabloda özetlenmiştir:

	Regresyon	Sınıflandırıcı
Çıktı	Sürekli	Sınıf
Örnekler	Lineer regresyon (bağlanım)	Lojistik regresyon (bağlanım), Destek Vektör Makineleri (DVM), Naive Bayes

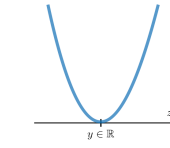
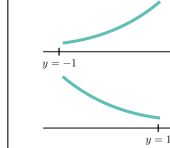
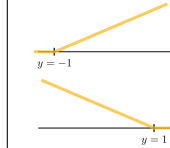
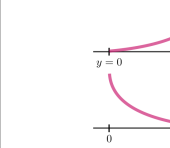
□ **Model türleri** – Farklı modeller aşağıdaki tabloda özetlenmiştir:

	Ayırt edici model	Üretici model
Amaç	Doğrudan tahmin $P(y x)$	$P(y x)$ 'i tahmin etmek için $P(x y)$ 'i tahmin etme
Öğrenilenler	Karar Sınırı	Verilerin olasılık dağılımı
Örnekleme		
Örnekler	Regresyon, DVM	GDA, Naive Bayes

Gösterimler ve genel konsept

□ **Hipotez** – Hipotez h_θ olarak belirtilmiştir ve bu bizim seçtiğimiz modeldir. Verilen $x^{(i)}$ verisi için modelin tahminlediği çıktı $h_\theta(x^{(i)})$ 'dir.

□ **Kayıp fonksiyonu** – $L : (z, y) \in \mathbb{R} \times Y \mapsto L(z, y) \in \mathbb{R}$ şeklinde tanımlanan bir kayıp fonksiyonu y gerçek değerine karşılık geleceği öngörülen z değerini girdi olarak alan ve ne kadar farklı olduklarını gösteren bir fonksiyondur. Yaygın kayıp fonksiyonları aşağıdaki tabloda özetlenmiştir:

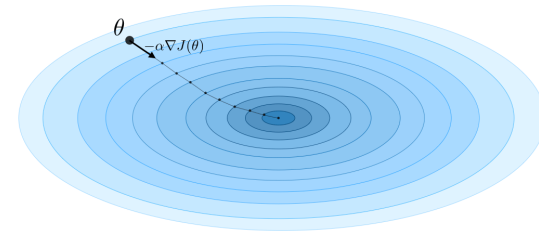
En küçük kareler hatası	Lojistik yitimi (kayıbı)	Menteşe yitimi (kayıbı)	Çapraz entropi
$\frac{1}{2}(y - z)^2$	$\log(1 + \exp(-yz))$	$\max(0, 1 - yz)$	$-[y \log(z) + (1 - y) \log(1 - z)]$
			
Lineer regresyon (bağlanım)	Lojistik regresyon (bağlanım)	DVM	Sinir Ağı

□ **Maliyet fonksiyonu** – J maliyet fonksiyonu genellikle bir modelin performansını değerlendirmek için kullanılır ve L kayıp fonksiyonu aşağıdaki gibi tanımlanır:

$$J(\theta) = \sum_{i=1}^m L(h_\theta(x^{(i)}), y^{(i)})$$

□ **Bayır inişi** – $\alpha \in \mathbb{R}$ öğrenme oranı olmak üzere, bayır inişi için güncelleme kuralı olarak ifade edilen öğrenme oranı ve J maliyet fonksiyonu aşağıdaki gibi ifade edilir:

$$\theta \leftarrow \theta - \alpha \nabla J(\theta)$$



Not: Stokastik bayır inişi her eğitim örneğine bağlı olarak parametreyi günceller, ve yığın bayır inişi bir dizi eğitim örneği üzerindedir.

□ **Olabilirlik** – θ parametreleri verilen bir $L(\theta)$ modelinin olabilirliğini, olabilirliği maksimize ederek en uygun θ parametrelerini bulmak için kullanılır. bulmak için kullanılır. Uygulamada, optimize edilmesi daha kolay olan log-olabilirlik $\ell(\theta) = \log(L(\theta))$ 'i kullanıyoruz. Sahip olduğumuz:

$$\theta^{\text{opt}} = \arg \max_{\theta} L(\theta)$$

□ **Newton'un algoritması** – $\ell'(\theta) = 0$ olacak şekilde bir θ bulan nümerik bir yöntemdir. Güncelleme kuralı aşağıdaki gibidir:

$$\theta \leftarrow \theta - \frac{\ell'(\theta)}{\ell''(\theta)}$$

Not: Newton-Raphson yöntemi olarak da bilinen çok boyutlu genelleme aşağıdaki güncelleme kuralına sahiptir:

$$\theta \leftarrow \theta - (\nabla_{\theta}^2 \ell(\theta))^{-1} \nabla_{\theta} \ell(\theta)$$

Lineer regresyon

$y|x; \theta \sim \mathcal{N}(\mu, \sigma^2)$ olduğunu varsayıyoruz

□ **Normal denklemler** – X matris tasarımı olmak üzere, maliyet fonksiyonunu en aza indiren θ değeri X 'in matris tasarımını not ederek, maliyet fonksiyonunu en aza indiren θ değeri kapalı formulu çözümdür:

$$\theta = (X^T X)^{-1} X^T y$$

□ **En Küçük Ortalama Kareler algoritması** – α öğrenme oranı olmak üzere, m veri noktasını içeren eğitim kümesi için Widrow-Hoff öğrenme oranı olarak bilinen En Küçük Ortalama Kareler Algoritmasının güncelleme kuralı aşağıdaki gibidir:

$$\forall j, \theta_j \leftarrow \theta_j + \alpha \sum_{i=1}^m [y^{(i)} - h_{\theta}(x^{(i)})] x_j^{(i)}$$

Not: güncelleme kuralı, bayır yükselişinin özel bir halidir.

□ **Yerel Ağırlıklı Regresyon** – LWR olarak da bilinen Yerel Ağırlıklı Regresyon ağırlıkları her eğitim örneğini maliyet fonksiyonunda $w^{(i)}(x)$ ile ölçen doğrusal regresyonun bir çeşididir.

$$w^{(i)}(x) = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right)$$

Sınıflandırma ve lojistik regresyon

□ **Sigmoid fonksiyonu** – Lojistik fonksiyonu olarak da bilinen sigmoid fonksiyonu g , aşağıdaki gibi tanımlanır:

$$\forall z \in \mathbb{R}, g(z) = \frac{1}{1 + e^{-z}} \in]0,1[$$

□ **Lojistik regresyon** – $y|x; \theta \sim \text{Bernoulli}(\phi)$ olduğunu varsayıyoruz. Aşağıdaki forma sahibiz:

$$\phi = p(y = 1|x; \theta) = \frac{1}{1 + \exp(-\theta^T x)} = g(\theta^T x)$$

Not: Lojistik regresyon durumunda kapalı form çözümü yoktur.

□ **Softmax regresyonu** – Çok sınıflı lojistik regresyon olarak da adlandırılan Softmax regresyonu 2'den fazla sınıf olduğunda lojistik regresyonu genelleştirmek için kullanılır. Genel kabul olarak, her i sınıfı için Bernoulli parametresi ϕ_i 'nin eşit olmasını sağlaması için $\theta_K = 0$ olarak ayarlanır.

$$\phi_i = \frac{\exp(\theta_i^T x)}{\sum_{j=1}^K \exp(\theta_j^T x)}$$

Genelleştirilmiş Lineer Modeller

□ **Üstel aile** – Eğer kanonik parametre veya bağlantı fonksiyonu olarak adlandırılan doğal bir parametre η , yeterli bir istatistik $T(y)$ ve aşağıdaki gibi bir log-partition fonksiyonu $a(\eta)$ şeklinde yazılabilirse, dağılım sınıfının üstel ailede olduğu söylenir:

$$p(y; \eta) = b(y) \exp(\eta T(y) - a(\eta))$$

Not: Sık sık $T(y) = y$ olur. Ayrıca, $\exp(-a(\eta))$, olasılıkların birleştiğinden emin olan normalleştirme parametresi olarak görülebilir.

Aşağıdaki tabloda özetlenen en yaygın üstel dağılımlar:

Dağılım	η	$T(y)$	$a(\eta)$	$b(y)$
Bernoulli	$\log\left(\frac{\phi}{1-\phi}\right)$	y	$\log(1 + \exp(\eta))$	1
Gauss	μ	y	$\frac{\eta^2}{2}$	$\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right)$
Poisson	$\log(\lambda)$	y	e^{η}	$\frac{1}{y!}$
Geometrik	$\log(1 - \phi)$	y	$\log\left(\frac{e^{\eta}}{1 - e^{\eta}}\right)$	1

□ **Genelleştirilmiş Lineer Modellerin Yaklaşımları** – Genelleştirilmiş Lineer Modeller $x \in \mathbb{R}^{n+1}$ için rastgele bir y değişkenini tahminlemeyi hedeflen ve aşağıdaki 3 varsayıma dayanan bir fonksiyondur:

$$(1) \quad y|x; \theta \sim \text{ExpFamily}(\eta) \quad (2) \quad h_{\theta}(x) = E[y|x; \theta] \quad (3) \quad \eta = \theta^T x$$

Not: sıradan en küçük kareler ve lojistik regresyon, genelleştirilmiş doğrusal modellerin özel durumlarıdır.

Destek Vektör Makineleri

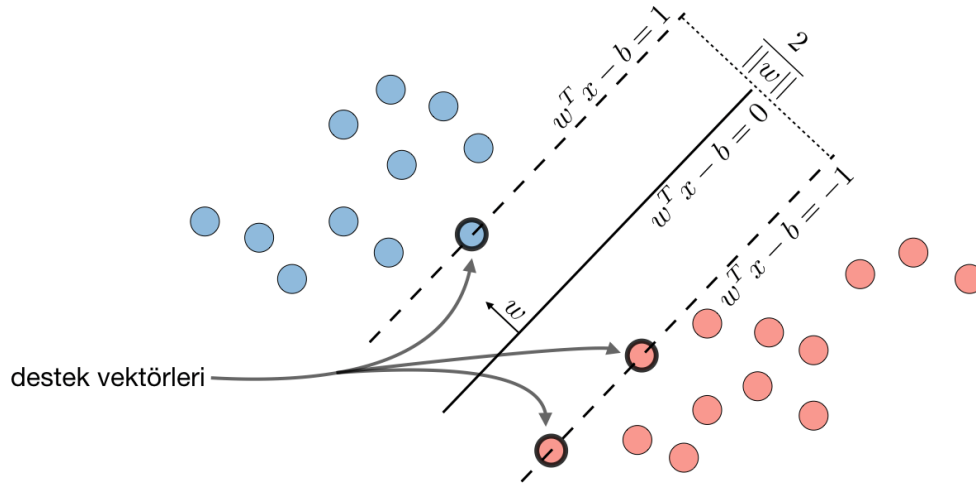
Destek Vektör Makinelerinin amacı minimum mesafeyi maksimuma çıkaran doğruyu bulmaktır.

□ **Optimal marj sınıflandırıcısı** – h optimal marj sınıflandırıcısı şöyledir:

$$h(x) = \text{sign}(w^T x - b)$$

burada $(w, b) \in \mathbb{R}^n \times \mathbb{R}$, aşağıdaki optimizasyon probleminin çözümüdür:

$$\min \frac{1}{2} \|w\|^2 \quad \text{öyle ki} \quad y^{(i)}(w^T x^{(i)} - b) \geq 1$$



Not: doğru $w^T x - b = 0$ şeklinde tanımlanır.

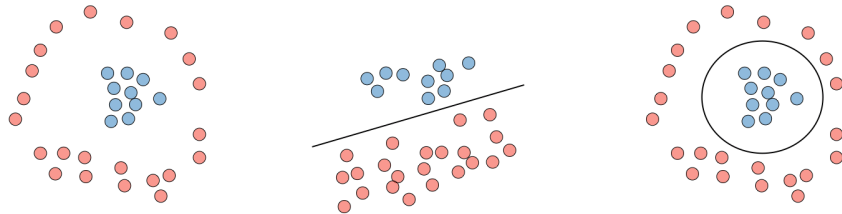
□ **Menteşe yitimi (kayıbı)** – Mentеше yitimi Destek Vektör Makinelerinin ayarlarında kullanılır ve aşağıdaki gibi tanımlanır:

$$L(z, y) = [1 - yz]_+ = \max(0, 1 - yz)$$

□ **Çekirdek** – ϕ gibi bir özellik haritası verildiğinde, K olarak tanımlanacak çekirdeği tanımlarız:

$$K(x, z) = \phi(x)^T \phi(z)$$

Uygulamada, $K(x, z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$ tarafından tanımlanan çekirdek K , Gauss çekirdeği olarak adlandırılır ve yaygın olarak kullanılır.



Lineer olmayan ayrılabilirlik \rightarrow Çekirdek Haritalarının Kullanımı ϕ \rightarrow Orjinal uzayda karar sınırı

Not: Çekirdeği kullanarak maliyet fonksiyonunu hesaplamak için "çekirdek numarası" nı kullandığımızı söylüyoruz çünkü genellikle çok karmaşık olan ϕ açık haritalamasını bilmeye gerek yok. Bunun yerine, yalnızca $K(x, z)$ değerlerine ihtiyacımız vardır.

□ **Lagranj** – Lagranj $\mathcal{L}(w, b)$ şeklinde şöyle tanımlanır:

$$\mathcal{L}(w, b) = f(w) + \sum_{i=1}^l \beta_i h_i(w)$$

Not: β_i katsayılarına Lagranj çarpanları denir.

Üretici Öğrenme

Üretken bir model, önce Bayes kuralını kullanarak $P(y|x)$ değerini tahmin etmek için kullanabileceğimiz $P(x|y)$ değerini tahmin ederek verilerin nasıl üretildiğini öğrenmeye çalışır.

Gauss Diskriminant (Ayrıtç) Analizi

□ **Yöntem** – Gauss Diskriminant Analizi y ve $x|y = 0$ ve $x|y = 1$ 'in şu şekilde olduğunu varsayar:

$$y \sim \text{Bernoulli}(\phi)$$

$$x|y = 0 \sim \mathcal{N}(\mu_0, \Sigma) \quad \text{ve} \quad x|y = 1 \sim \mathcal{N}(\mu_1, \Sigma)$$

□ **Tahmin** – Aşağıdaki tablo, olasılığı en üst düzeye çıkarırken bulduğumuz tahminleri özetlemektedir:

$\hat{\phi}$	$\hat{\mu}_j \quad (j = 0, 1)$	$\hat{\Sigma}$
$\frac{1}{m} \sum_{i=1}^m 1_{\{y^{(i)}=1\}}$	$\frac{\sum_{i=1}^m 1_{\{y^{(i)}=j\}} x^{(i)}}{\sum_{i=1}^m 1_{\{y^{(i)}=j\}}}$	$\frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T$

Naive Bayes

□ **Varsayım** – Naive Bayes modeli, her veri noktasının özelliklerinin tamamen bağımsız olduğunu varsayar:

$$P(x|y) = P(x_1, x_2, \dots | y) = P(x_1|y)P(x_2|y)\dots = \prod_{i=1}^n P(x_i|y)$$

□ **Çözümler** – Log-olabilirliğinin $k \in \{0, 1\}, l \in [1, L]$ ile birlikte aşağıdaki çözümlerle maksimize edilmesi:

$$P(y = k) = \frac{1}{m} \times \#\{j|y^{(j)} = k\}$$

$$P(x_i = l | y = k) = \frac{\#\{j|y^{(j)} = k \text{ ve } x_i^{(j)} = l\}}{\#\{j|y^{(j)} = k\}}$$

Not: Naive Bayes, metin sınıflandırması ve spam tespitinde yaygın olarak kullanılır.

Ağaç temelli ve topluluk yöntemleri

Bu yöntemler hem regresyon hem de sınıflandırma problemleri için kullanılabilir.

□ **CART** – Sınıflandırma ve Regresyon Ağaçları (Classification and Regression Trees (CART)), genellikle karar ağaçları olarak bilinir, ikili ağaçlar olarak temsil edilirler.

□ **Rastgele orman** – Rastgele seçilen özelliklerden oluşan çok sayıda karar ağacı kullanan ağaç tabanlı bir tekniktir. Basit karar ağacının tersine, oldukça yorumlanamaz bir yapıdadır ancak genel olarak iyi performansı onu popüler bir algoritma yapar.

Not: Rastgele ormanlar topluluk yöntemlerindedir.

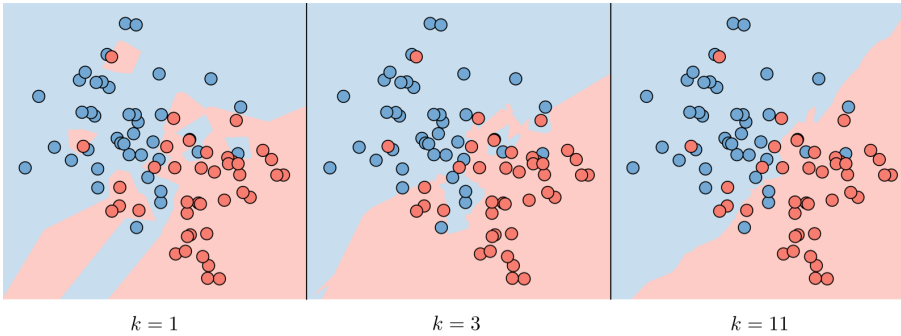
□ **Artırım** – Artırım yöntemlerinin temel fikri bazı zayıf öğrencileri biraraya getirerek güçlü bir öğrenci oluşturmaktır. Temel yöntemler aşağıdaki tabloda özetlenmiştir:

Adaptif artırım	Gradyan artırım
Yüksek ağırlıklar bir sonraki artırım adımında iyileşmesi için hatalara maruz kalır.	Zayıf öğrenciler kalan hatalar üzerinde eğitildi

Diğer parametrik olmayan yaklaşımlar

□ **k -en yakın komşular** – Genellikle k -NN olarak adlandırılan k -en yakın komşular algoritması, bir veri noktasının tepkisi eğitim kümesindeki kendi k komşularının doğası ile belirlenen parametrik olmayan bir yaklaşımdır. Hem sınıflandırma hem de regresyon yöntemleri için kullanılabilir.

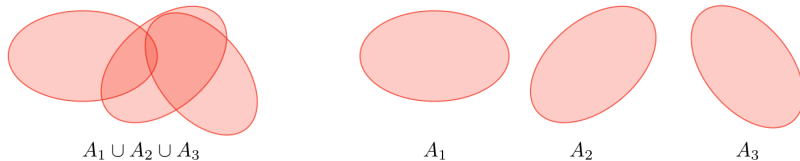
Not: k parametresi ne kadar yüksekse, yanlılık okadar yüksek ve k parametresi ne kadar düşükse, varyans o kadar yüksek olur.



Öğrenme Teorisi

□ **Birleşim sınırı** – A_1, \dots, A_k k olayları olsun. Sahip olduklarımız:

$$P(A_1 \cup \dots \cup A_k) \leq P(A_1) + \dots + P(A_k)$$



□ **Hoeffding eşitsizliği** – Z_1, \dots, Z_m , ϕ parametresinin Bernoulli dağılımından çizilen değişkenler olsun. Örnek ortalamaları mean ve $\gamma > 0$ sabit olsun. Sahip olduklarımız:

$$P(|\phi - \hat{\phi}| > \gamma) \leq 2 \exp(-2\gamma^2 m)$$

Not: Bu eşitsizlik, Chernoff sınırı olarak da bilinir.

□ **Eğitim hatası** – Belirli bir h sınıflandırıcısı için, ampirik risk veya ampirik hata olarak da bilinen eğitim hatasını $\hat{\epsilon}(h)$ şöyle tanımlarız:

$$\hat{\epsilon}(h) = \frac{1}{m} \sum_{i=1}^m 1_{\{h(x^{(i)}) \neq y^{(i)}\}}$$

□ **Olası Yaklaşık Doğru** – PAC, öğrenme teorisi üzerine sayısız sonuçların kanıtlandığı ve aşağıdaki varsayımlara sahip olan bir çerçevedir:

- eğitim ve test kümeleri aynı dağılımı takip ediyor
- eğitim örnekleri bağımsız olarak çizilir

□ **Parçalanma** – $S = \{x^{(1)}, \dots, x^{(d)}\}$ kümesi ve \mathcal{H} sınıflandırıcıların kümesi verildiğinde, \mathcal{H} herhangi bir etiketler kümesi S' 'e parçalar.

$$\exists h \in \mathcal{H}, \quad \forall i \in [1, d], \quad h(x^{(i)}) = y^{(i)}$$

□ **Üst sınır teoremi** – $|\mathcal{H}| = k$, δ ve örneklem sayısı m 'nin sabit olduğu sonlu bir hipotez sınıfı \mathcal{H} olsun. Ardından, en az $1 - \delta$ olasılığı ile elimizde:

$$\epsilon(\hat{h}) \leq \left(\min_{h \in \mathcal{H}} \epsilon(h) \right) + 2 \sqrt{\frac{1}{2m} \log \left(\frac{2k}{\delta} \right)}$$

□ **VC boyutu** – $VC(\mathcal{H})$ olarak ifade edilen belirli bir sonsuz \mathcal{H} hipotez sınıfının Vapnik-Chervonenkis (VC) boyutu, \mathcal{H} tarafından parçalanabilen en büyük kümenin boyutudur.

Not: $\mathcal{H} = \{2 \text{ boyutta doğrusal sınıflandırıcılar kümesi}\}$ 'nin VC boyutu 3'tür.



□ **Teorem (Vapnik)** – \mathcal{H} , $VC(\mathcal{H}) = d$ ve eğitim örneği sayısı m verilmiş olsun. En az $1 - \delta$ olasılığı ile, sahip olduklarımız:

$$\epsilon(\hat{h}) \leq \left(\min_{h \in \mathcal{H}} \epsilon(h) \right) + O \left(\sqrt{\frac{d}{m} \log \left(\frac{m}{d} \right)} + \frac{1}{m} \log \left(\frac{1}{\delta} \right) \right)$$