

Model - veri ilişkisi

Sinan Yıldırım

MDBF, Sabancı Üniversitesi

4 Şubat 2021

Bu sunumdaki şekiller, şu kitaplardan alınmıştır.

Pattern Recognition and Machine Learning, Christopher M. Bishop;

Mathematics for Machine Learning, Deisenroth v.d.

Örnek-etiket tipi veriler

Veri:

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$$

- ▶ \mathbf{x}_n : örnek
- ▶ y_n : etiket

Amaçlar

- ▶ Veriyi açıklama
- ▶ Yeni bir \mathbf{x} verildiğinde y 'yi tahminleme.
- ▶ vb

Örnek veri

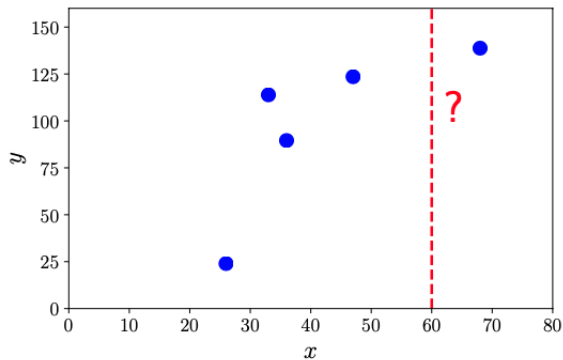
Name	Gender	Degree	Postcode	Age	Annual salary
Aditya	M	MSc	W21BG	36	89563
Bob	M	PhD	EC1A1BA	47	123543
Chloé	F	BEcon	SW1A1BH	26	23989
Daisuke	M	BSc	SE207AT	68	138769
Elisabeth	F	MBA	SE10AA	33	113888

Örnek veri - sayısallaştırılmış

Gender ID	Degree	Latitude (in degrees)	Longitude (in degrees)	Age	Annual Salary (in thousands)
-1	2	51.5073	0.1290	36	89.563
-1	3	51.5074	0.1275	47	123.543
+1	1	51.5071	0.1278	26	23.989
-1	1	51.5075	0.1281	68	138.769
+1	2	51.5074	0.1278	33	113.888

Tahminleme

x : yaş,
 y : maaş



Tahminleyici

Aday fonksiyonlar:

$$f(\cdot, \boldsymbol{\theta}) : \mathbb{R}^d \mapsto \mathbb{R}.$$

$\boldsymbol{\theta}$: Tahminleyicinin parametresi

Amaç:

$$f(\mathbf{x}_n, \boldsymbol{\theta}^*) \approx y_n, \quad \forall n = 1, \dots, N.$$

olacak şekilde bir $\boldsymbol{\theta}^*$ belirlemek.

Tahmin: $\hat{y}_n = f(\mathbf{x}_n, \boldsymbol{\theta}^*)$.

Örnek: Doğrusal regresyon

$$\mathbf{x}_n = \begin{bmatrix} 1 & x_n^{(1)} & x_n^{(2)} & \dots & x_n^{(D)} \end{bmatrix}^T, \quad y \in \mathbb{R}$$

$$\boldsymbol{\theta} = [\theta_0 \quad \theta_1 \quad \dots \quad \theta_D]^T$$

$$f(\cdot, \boldsymbol{\theta}) : \mathbb{R}^{D+1} \mapsto \mathbb{R},$$

$$f(\mathbf{x}_n, \boldsymbol{\theta}) = \boldsymbol{\theta}^T \mathbf{x}_n$$

Ampirik risk enküçültme

Ampirik risk

$$\mathbf{X} = [\mathbf{x}_1 \quad \dots \quad \mathbf{x}_N]^T \in \mathbb{R}^{N \times D}, \quad \mathbf{y} = [y_1 \quad \dots \quad y_N]^T$$

Tahmin

$$\hat{y}_n = f(\mathbf{x}_n, \boldsymbol{\theta}).$$

Kayıp

$$\ell(y_n, \hat{y}_n).$$

Ampirik risk:

$$\mathbf{R}_{\text{amp}}(f, \mathbf{X}, \mathbf{y}) = \frac{1}{N} \sum_{n=1}^N \ell(y_n, \hat{y}_n).$$

Ampirik risk enküçültmesi:

$$f^* = \arg \min_f \mathbf{R}_{\text{emp}}(f, \mathbf{X}, \mathbf{y})$$

Örnek: En küçük kareler

Karesel kayıp:

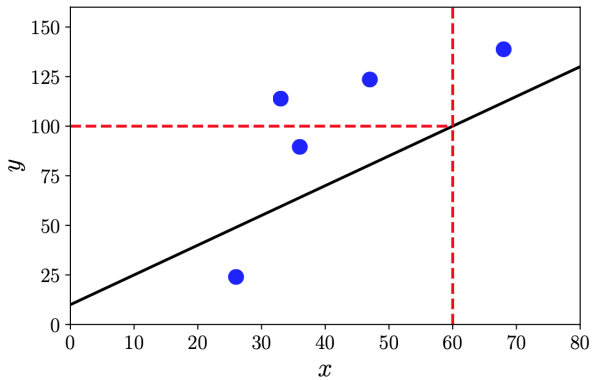
$$\ell(y_n, \hat{y}_n) = (y_n - \hat{y}_n)^2$$

Doğrusal tahminleyici:

$$f(\mathbf{x}_n, \boldsymbol{\theta}) = \boldsymbol{\theta}^T \mathbf{x}_n.$$

Ampirik risk enküçültmesi:

$$\begin{aligned}\boldsymbol{\theta}^* &= \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^D} \frac{1}{N} \sum_{n=1}^N (y_n - \boldsymbol{\theta}^T \mathbf{x}_n)^2 \\ &= \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^D} \frac{1}{N} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2.\end{aligned}$$



Gerçek ve ampirik risk

Gerçek risk

$$R_{\text{ideal}}(f) = \mathbb{E}_{\mathbf{x}, y}[\ell(y, f(\mathbf{x}))]$$

Ampirik risk, gerçek riski kestirir.

$$R_{\text{amp}}(f, \mathbf{X}, \mathbf{y}) \approx R_{\text{ideal}}(f).$$

Aşırı uyma problemi

Eğitim sonucu *görünmeyen* veriye ne kadar uyumlu?

Bunu kestirmek için genelde veri ikiye bölünür:

$$(\mathbf{X}_{\text{train}}, y_{\text{train}}), \quad (\mathbf{X}_{\text{test}}, y_{\text{test}})$$

İstenen:

$$R_{\text{amp}}(f, \mathbf{X}_{\text{train}}, y_{\text{train}}) \approx R_{\text{amp}}(f, \mathbf{X}_{\text{test}}, y_{\text{test}})$$

Aşırı uyma:

$$R_{\text{amp}}(f, \mathbf{X}_{\text{train}}, y_{\text{train}}) < R_{\text{amp}}(f, \mathbf{X}_{\text{test}}, y_{\text{test}})$$

Aşırı uyumun belirtilerinden biri θ bileşenlerinin yüksek değerler almasıdır.

Düzenlenmiş problem

$$\min_{\theta \in \mathbb{R}^D} \frac{1}{N} \|y - \mathbf{X}\theta\|^2 + \lambda \|\theta\|^2.$$

- ▶ $\|\theta\|^2$: düzenleyici
- ▶ λ : düzenleme katsayısı
- ▶ $\lambda \|\theta\|^2$: ceza

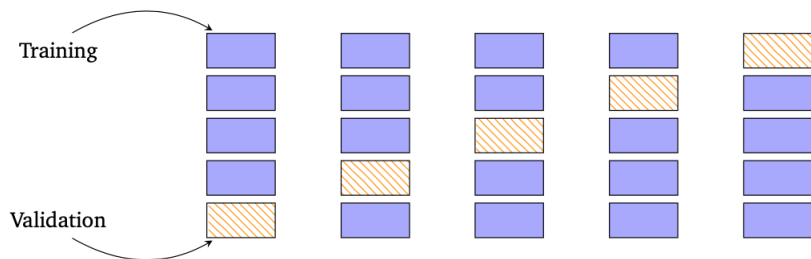
Belli bir veri üzerinde eğitilen bir model aynı kaynaktan başka verileri iyi tahminler mi?

Seçilmiş olan aday fonksiyon ailesine, yani modele dair bir ölçüt.

Belli bir θ 'yı değil, $\{f(\cdot, \theta) : \theta \in \Theta\}$ ailesini ilgilendiren bir ölçüt.

Genelleme başarımı nasıl ölçülür?

K katlı çapraz doğrulama



$K - 1$ parçada eğit, kalan parçada sına. K kez tekrarla.

$$\mathbb{E}_{\mathcal{V}}[R(f, \mathcal{V})] \approx \frac{1}{K} \sum_{k=1}^K R(f^{(k)}, \mathcal{V}^{(k)})$$

Parametre kestirimi

Parametre kestirimi

Kayıp fonksiyonları yerine olasılık dağılımları kullanılır.

θ : veriyi açıklamak için kullanılan olasılık dağılımının parametresi.

Parametre kestirimi: En iyi θ 'yı bulmak.

Örnek - etiket veri modeli

Bağımsız özdeş ikililer:

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$$

$$\mathcal{Y} = \{y_1, \dots, y_N\}, \quad \mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}.$$

Genelde \mathbf{x}_n rassal kabul edilmez:

$$p(\mathcal{Y}|\mathcal{X}, \boldsymbol{\theta}) = \prod_{n=1}^N p(y_n|\mathbf{x}_n, \boldsymbol{\theta}).$$

Olabilirlik fonksiyonu

Negatif log-olabilirlik

$$\mathcal{L}(\boldsymbol{\theta}) = -\log p(\mathcal{Y}|\mathcal{X}, \boldsymbol{\theta}) = -\sum_{n=1}^N \log p(y_n|\mathbf{x}_n, \boldsymbol{\theta}).$$

Enbüyük olabilirlik kestirimi:

$$\boldsymbol{\theta}_{\text{ML}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})$$

Gauss dağılımı ve en küçük kareler

$$p(y_n|\mathbf{x}_n, \boldsymbol{\theta}) = \mathcal{N}(y_n|\mathbf{x}_n^T \boldsymbol{\theta}, \sigma^2).$$

Negatif log-olabilirlik

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{x}_n^T \boldsymbol{\theta})^2 - \sum_{n=1}^N \log \frac{1}{\sqrt{2\pi\sigma^2}}.$$

Enbüyük sonsal dağılım kestirimi

Veri \mathcal{D} , parametre: θ

$$p(\theta|\mathcal{D}) = \frac{p(\theta)p(\mathcal{D}|\theta)}{p(\mathcal{D})} \propto p(\theta)p(\mathcal{D}|\theta).$$

$$\theta^* = \arg \max_{\theta} p(\theta|\mathcal{D})$$

Önceki örneklerde, $\mathcal{D} = (\mathcal{X}, \mathcal{Y})$.

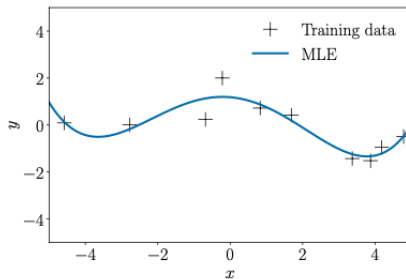
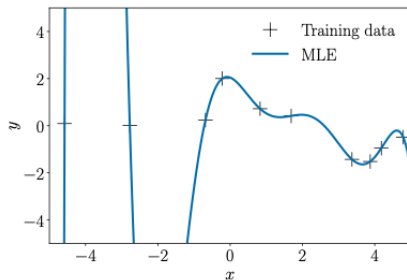
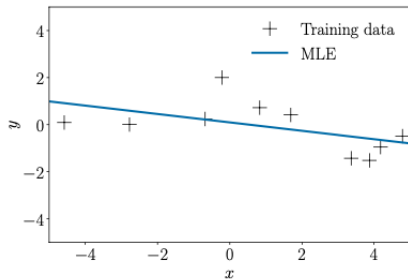
Hedef fonksiyon:

$$p(\mathcal{Y}|\mathcal{X}, \theta)p(\theta)$$

(\mathcal{X} 'in olasılık dağılımı olsa da aynı hedef kullanılabilir.)

Aşırı uyum, yetersiz uyum, kararında uyum

Olasılık modeli (görünmeyen) veri için ne kadar uygun?



Olasılıksal modelleme ve çıkarım

Olasılıksal modelleme

Hem parametre hem de veri rassal değişken:

- ▶ \mathbf{x} : gözlemlenen değişken (veri)
- ▶ θ bilinmeyen değişken, parametre

Çıkış noktası: Ortak dağılım

$$p(\mathbf{x}, \theta)$$

Ortak dağılımın içerdiği bilgiler:

- ▶ Öncül dağılım $p(\theta)$ ve olabilirlik $p(\mathbf{x}|\theta)$.
- ▶ Marjinal olabilirlik $p(\mathbf{x})$ (model seçimi için gerekli)
- ▶ Sonsal dağılım $p(\theta|\mathbf{x})$

Bayesci çıkarım

Hedef: Sonsal dağılımın belirlenmesi.

$$p(\theta|\mathbf{x}) = \frac{p(\theta)p(\mathbf{x}|\theta)}{p(\mathbf{x})}$$

Tahminleme:

$$p(\mathbf{x}_{yeni}|\mathbf{x}) = \int_{\theta} p(\theta|\mathbf{x})p(\mathbf{x}_{new}|\theta, \mathbf{x})d\theta$$

Noktasal kestirimler ile tahminleme:

$$p(\mathbf{x}_{yeni}|\mathbf{x}) \approx p(\mathbf{x}_{new}|\theta^*, \mathbf{x})$$

Saklı değişkenli modeller

z : saklı değişkeni

$$p(\mathbf{x}|\boldsymbol{\theta}) = \int p(\mathbf{x}|z, \boldsymbol{\theta})p(z)dz$$

Ortak dağılım:

$$p(\boldsymbol{\theta}, \mathbf{x}, z) = p(\boldsymbol{\theta})p(z)p(\mathbf{x}|\boldsymbol{\theta}, z)$$

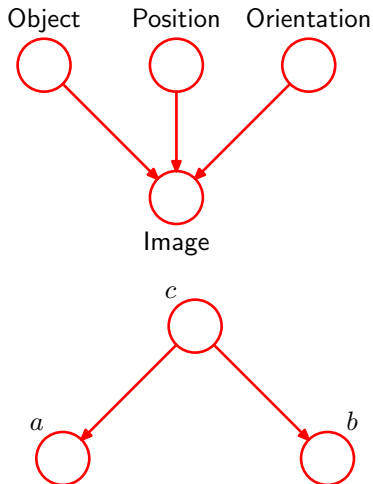
Kullanımı?

$$p(\boldsymbol{\theta}|z, \mathbf{x}), \quad p(z|\boldsymbol{\theta}, \mathbf{x})$$

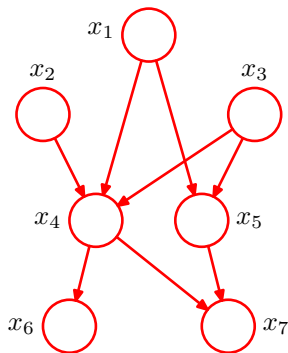
Çizge modelleri

Yönlü düz çizgeler (Bayes ağları)

Rassal değişkenlerin arasındaki koşullu bağımsızlık ilişkilerini sebep sonuç ilişkisini temel alarak verir.



Yönlü düz çizgeler - Olasılık dağılımı



$$p(x_{1:7}) = p(x_1)p(x_2)p(x_3)p(x_4|x_1, x_2, x_3)p(x_5|x_1, x_3)p(x_6|x_4)p(x_7|x_4, x_5)$$

Genel:

$$p(\mathbf{x}) = \prod_{k=1}^K p(x_k | \text{Pa}_k)$$

Örnek

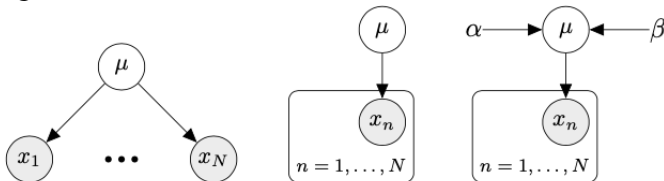
Bir madeni para N kere atılıyor. Yazı: $x = 1$, Tura: $x = 0$

$$p(x|\mu) = \text{Ber}(x|\mu)$$

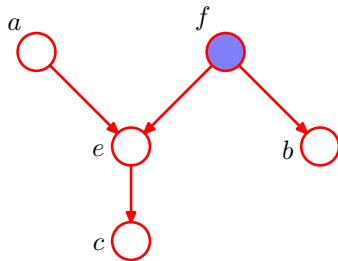
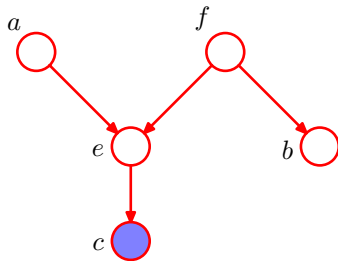
$$p(x_1, \dots, x_N|\mu) = \prod_{n=1}^N p(x_n|\mu)$$

Öncül dağılım $p(\mu) = \text{Beta}(\mu|\alpha, \beta)$.

Üç farklı gösterim



Koşullu bağımsızlık



► $a \perp\!\!\!\perp b|c$ midir?

► $a \perp\!\!\!\perp b|f$ midir?

Koşullu bağımsızlık - d -ayrılık

A , B ve C düğüm kümeleri olsun.

A ve B C 'ye koşullu bağımsız mıdır?

$$A \perp\!\!\!\perp B | C$$

d -ayrılık

A 'daki bir düğümden B 'deki bir düğüme giden bir yol ele alalım.

Bu yolun üstündeki herhangi bir düğümü ele alalım:

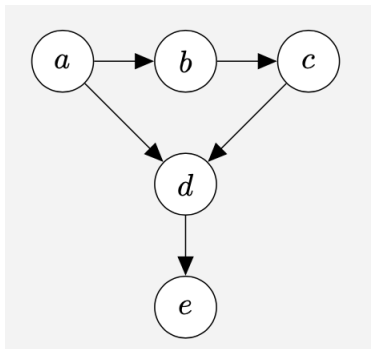
- Bu düğüme ulaşan oklar **kuyruk-kuyruğa** ve **kafa-kuyruğa** ise ve bu düğüm C 'nin içindeyse; veya
- Bu düğüme gelen oklar **kafa-kafaya** ise ve ne bu düğüm ne de onun alt-düğümü C 'nin içinde ise (neither-nor);

bu yol C tarafından engellenmiş sayılır.

Eğer A 'daki her bir düğümden B 'deki her bir düğüme giden bütün yollar C tarafından engellenmişse, A ve B , C tarafından d -ayrılmıştır, ve

$$A \perp\!\!\!\perp B | C$$

sağlanır.



$$b \perp\!\!\!\perp d \mid a, c \quad ?$$

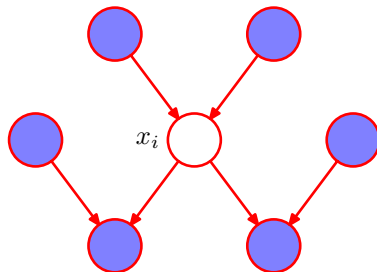
$$a \perp\!\!\!\perp c \mid a \quad ?$$

$$b \perp\!\!\!\perp d \mid c \quad ?$$

$$a \perp\!\!\!\perp c \mid b, e \quad ?$$

Markov battaniyesi

- ▶ Bir üst-düğüm (ebeveynler)
- ▶ bir alt-düğüm (çocuklar),
- ▶ çocukların beraber yapıldığı partnerler



Bir düğüm, Markov battaniyesi verildiğinde diğer düğümlerden bağımsızdır.

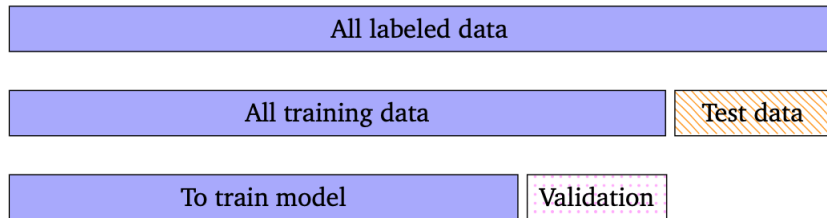
Model seçimi

Model seçimi problemleri

Örnekler:

- ▶ Bir örneklem için popülasyon dağılımı seçimi
- ▶ Regresyon için kullanılan polinomun derecesi seçimi
- ▶ Regresyonda açıklayıcı değişkenlerin seçimi
- ▶ Karışım dağılımlarında bileşen sayısı seçimi
- ▶ Temel bileşen analizinde boyut seçimi
- ▶ Destek vektör makinelerinde çekirdek seçimi
- ▶ Derin öğrenmede ağ yapısının seçimi

İç içe geçmiş çapraz doğrulama



$$\mathbb{E}_{\mathcal{V}}[R(\mathcal{V}|M)] \approx \frac{1}{K} \sum_{k=1}^K R(\mathcal{V}^{(k)}|M)$$

Tüm M modelleri için hesapla ve en iyisini seç.

Bayesci model seçimi

Modelin kendisi de bir değişken:

$$M \sim p(M)$$

$$\theta|M \sim p(\theta|M)$$

$$\mathcal{D}|M, \theta \sim p(\mathcal{D}|\theta, M)$$

Marjinal olabilirlik:

$$p(\mathcal{D}|M) = \int p(\mathcal{D}|\theta, M)p(\theta|M)d\theta$$

Modelin sonsal dağılımı:

$$p(M|\mathcal{D}) \propto p(M)p(\mathcal{D}|M)$$

Hedef:

$$M^* = \arg \max_M p(M|\mathcal{D}).$$

Marjinal olabilirliğin düzenleyici özelliği

Occam'ın usturası: Veriyi açıklayabilen iki modelden basit olanının marjinal olabilirliği genelde daha yüksektir.

