



بازشناسی الگو

درس ۵

طبقه‌بندی مبتنی بر نظریه‌ی تصمیم بیز روش‌های تخمین پارامتری تابع چگالی احتمال مجهول

Classification Based on Bayes Decision Theory
Parametric Estimation Methods for Unknown Probability Density Functions

کاظم فولادی قلعه

دانشکده مهندسی، پردیس فارابی
دانشگاه تهران

<http://courses.fouladi.ir/pr>

تخمین توابع چگالی احتمال مجهول

ESTIMATION OF UNKNOWN PROBABILITY DENSITY FUNCTIONS

روش طبقه‌بندی بیزی نیازمند مدل است.

مدل: توابع چگالی احتمال پیشین و پسین طبقه‌ها

تاکنون فرض شده بود که توابع pdf معلوم هستند، اما معمولاً این‌گونه نیست: در بسیاری از مسائل باید pdf را از روی داده‌های آموزشی موجود، تخمین زد.

رویکردهای تخمین توابع چگالی احتمال مجهول	
نایپارامتری <i>Nonparametric</i>	پارامتری <i>Parametric</i>
pdf را نمی‌دانیم، اما برخی آماره‌های آن (مثلًا μ یا σ^2) معلوم است. روش‌ها: <ul style="list-style-type: none"> ○ Parzen Windows ○ k Nearest Neighbor 	pdf را می‌دانیم، اما پارامترهای آن مجهول است. روش‌ها: <ul style="list-style-type: none"> ○ Maximum Likelihood (ML) ○ Maximum a Posteriori Probability (MAP) ○ Maximum Entropy (ME) ○ Bayesian Inference ○ Mixture Models

تخمین توابع چگالی احتمال مجهول

رویکرد پارامتری

ESTIMATION OF UNKNOWN PROBABILITY DENSITY FUNCTIONS

روش طبقه‌بندی بیزی نیازمند مدل است.

مدل: توابع چگالی احتمال پیشین و پسین طبقه‌ها

تاکنون فرض شده بود که توابع pdf معلوم هستند، اما معمولاً این‌گونه نیست: در بسیاری از مسائل باید pdf را از روی داده‌های آموزشی موجود، تخمین زد.

رویکردهای تخمین توابع چگالی احتمال مجهول

نپارامتری

Nonparametric

pdf را نمی‌دانیم،
اما برخی آمارهای آن (مثلًا μ یا σ^2) معلوم است.

روش‌ها:

- Parzen Windows
- k Nearest Neighbor

پارامتری

Parametric

pdf را می‌دانیم،
اما پارامترهای آن مجهول است.

روش‌ها:

- Maximum Likelihood (ML)
- Maximum a Posteriori Probability (MAP)
- Maximum Entropy (ME)
- Bayesian Inference
- Mixture Models

طبقه‌بندی مبتنی بر نظریه‌ی تصمیم بیز:
روش‌های تخمین پارامتری
تابع چکالی احتمال مجہول

۱

روش تخمین ماکریم درست‌نمایی

تخمین پارامتری توابع چگالی احتمال مجهول

روش ماکزیمم درست‌نمایی

MAXIMUM LIKELIHOOD (ML)

روش ML، پارامترها را به عنوان کمیت‌های ثابت اما مجهول می‌بیند.

مجموعه‌ی داده‌های آموزشی:
 نمونه‌های مستقل با توزیع یکسان (iid)
 $p(\mathbf{x}|\boldsymbol{\theta})$ بیرون کشیده شده از تابع چگالی احتمال $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$

هدف: تخمین پارامتر مجهول $\boldsymbol{\theta}$ با استفاده از نمونه‌های آموزشی \mathcal{D}

تابع درست‌نمایی $\boldsymbol{\theta}$ نسبت به داده‌های \mathcal{D}

$$L(\boldsymbol{\theta}|\mathcal{D})$$

تابع درست‌نمایی
Likelihood Function

$$L(\boldsymbol{\theta}|\mathcal{D}) = p(\mathcal{D}|\boldsymbol{\theta}) = p(\mathbf{x}_1, \dots, \mathbf{x}_n|\boldsymbol{\theta}) = \prod_{i=1}^n p(\mathbf{x}_i|\boldsymbol{\theta})$$

با فرض استقلال آماری

(بررسی وجود استقلال آماری دشوار است: اگر برای وابستگی دو متغیر دلیلی نداشتیم، آن دو را مستقل فرض می‌کنیم.)

تخمین پارامتری توابع چگالی احتمال مجهول

روش ماکزیمم درستنمایی

MAXIMUM LIKELIHOOD (ML)

تخمین درستنمایی MLE برای θ مقداری از $\hat{\theta}$ است تابع درستنمایی را ماکزیمم می‌کند:

$$\hat{\theta} = \arg \max_{\theta} L(\theta | \mathcal{D})$$

معمولًاً کار کردن با لگاریتم تابع درستنمایی ساده‌تر است، و به همان تخمین برای پارامتر می‌رسد.

$$\hat{\theta} = \arg \max_{\theta} \log L(\theta | \mathcal{D}) = \arg \max_{\theta} \sum_{i=1}^n \log p(\mathbf{x}_i | \theta)$$

لگاریتم تابع درستنمایی θ نسبت به داده‌های \mathcal{D}

$$\log L(\theta | \mathcal{D})$$

تابع لگاریتم درستنمایی

Log-Likelihood Function

تخمین پارامتری توابع چگالی احتمال مجهول

روش ماکزیمم درستنامایی: محاسبه

MAXIMUM LIKELIHOOD (ML)

اگر تعداد پارامترها p باشد، آنگاه

$$\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)^T$$

و عملگر گرادیان به صورت زیر تعریف می‌شود:

$$\nabla_{\boldsymbol{\theta}} \equiv \begin{bmatrix} \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_p} \end{bmatrix}$$

در این صورت، تخمین ماکزیمم درستنامایی $\boldsymbol{\theta}$ باید شرط لازم زیر را برآورده کند:

$$\nabla_{\boldsymbol{\theta}} \log L(\boldsymbol{\theta} | \mathcal{D}) = \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}_i | \boldsymbol{\theta}) = 0$$

(مشتق مساوی صفر)

انتخاب ماکزیمم سراسری از میان اکسٹریمم‌های محلی حاصل از مشتق لازم است. همچنین، شرایط مرزی باید بررسی شود.

تخمین پارامتری توابع چگالی احتمال مجهول

روش ماکزیمم درستنمایی: خصوصیات (۱) : نازاریبی مجانبی

MAXIMUM LIKELIHOOD (ML)

تخمین ML به طور مجانبی نازاریب (unbiased) است.

اگر θ_0 مقدار واقعی پارامتر مجهول θ باشد، آن‌گاه

$$\lim_{N \rightarrow \infty} E\{\hat{\theta}_{\text{ML}}\} = \theta_0$$

یعنی: برآورد ML در حالت میانگین به مقدار واقعی پارامتر همگرا می‌شود.

که برای مجموعه داده‌های مختلف \mathcal{D} برآوردهای مختلفی دارد.
 $\hat{\theta}_{\text{ML}}$ به خودی خود یک متغیر تصادفی است

تخمین پارامتری توابع چگالی احتمال مجهول

روش ماکزیمم درست‌نمایی: خصوصیات (۲): سازگاری مجانبی

MAXIMUM LIKELIHOOD (ML)

تخمین ML به‌طور مجانبی سازگار (consistent) است.

اگر θ_0 مقدار واقعی پارامتر مجهول θ باشد، آن‌گاه

$$\lim_{N \rightarrow \infty} E\{||\hat{\theta}_{\text{ML}} - \theta_0||^2\} = 0$$

یعنی: برآوردهای ML در حالت میانگین مربعات همگرا می‌شود \Leftarrow
برای N ‌های بزرگ، واریانس برآوردهای ML به سمت صفر می‌کند.

که برای مجموعه داده‌های مختلف \mathcal{D} برآوردهای مختلفی دارد.
 $\hat{\theta}_{\text{ML}}$ به خودی خود یک متغیر تصادفی است

تخمین پارامتری توابع چگالی احتمال مجهول

روش ماکزیمم درست نمایی: خصوصیات (۲): همگرایی به توزیع گاووسی

MAXIMUM LIKELIHOOD (ML)

تابع چگالی تخمین ML برای $\infty \rightarrow N$ به سمت توزیع گاووسی با میانگین θ_0 میل می‌کند.

زیرا:

۱) قضیه‌ی حد مرکزی

۲) ML وابسته به مجموع N متغیر تصادفی است.

تخمین پارامتری توابع چگالی احتمال مجهول

روش ماکریزم درست‌نمایی: خصوصیات (۴): تغییرناپذیری

MAXIMUM LIKELIHOOD (ML)

اگر $\hat{\theta}$ تخمین ML پارامتر θ باشد،
آن‌گاه برای هر تابع f دلخواه،
تخمین ML پارامتر $f(\theta)$ برابر با $f(\hat{\theta})$ است.

تخمین پارامتری توابع چگالی احتمال مجهول

روش ماکزیمم درستنمایی: خصوصیات

MAXIMUM LIKELIHOOD (ML)

تخمین ML فقط برای زمانی که $N \rightarrow \infty$ خصوصیات زیر را دارد:

ناریبی

Unbiased

می‌نیم واریانس ممکن

Minimum Variance

همگرایی به گاوی

Convergence to Gaussian pdf

تخمین پارامتری توابع چگالی احتمال مجهول

روش ماکزیمم درست‌نمایی: مثال

MAXIMUM LIKELIHOOD (ML)

داده‌ی تولید شده با یک pdf گاوی یک بعدی داریم

$$x_1, x_2, \dots, x_N$$

که میانگین آن μ و واریانس آن مجهول است.

تخمین ماکزیمم درست‌نمایی واریانس مجهول را به دست آورید.

تخمین پارامتری توابع چگالی احتمال مجهول

روش ماکزیمم درست‌نمایی: مثال

MAXIMUM LIKELIHOOD (ML)

N داده‌ی تولید شده با یک pdf گاوی \sim -بعدی داریم
 $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N$

که بردار میانگین آن μ مجهول و ماتریس کوواریانس آن Σ معلوم است.
تخمین ماکزیمم درست‌نمایی بردار میانگین مجهول را به دست آورید.

تخمین پارامتری توابع چگالی احتمال مجهول

روش ماکزیمم درستنمایی: توزیع گاووسی

MAXIMUM LIKELIHOOD (ML)

Suppose that $p(\mathbf{x}|\boldsymbol{\theta}) = N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

- ▶ When $\boldsymbol{\Sigma}$ is known but $\boldsymbol{\mu}$ is unknown:

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

- ▶ When both $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ are unknown:

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \text{and} \quad \hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$$

تخمین پارامتری توابع چگالی احتمال مجهول

روش ماکزیمم درستنمایی: توزیع برنولی

MAXIMUM LIKELIHOOD (ML)

- Suppose that $P(x|\theta) = \text{Bernoulli}(\theta) = \theta^x(1-\theta)^{1-x}$ where $x = 0, 1$ and $0 \leq \theta \leq 1$.
- The MLE of θ can be computed as

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n x_i$$

تخمین پارامتری توابع چگالی احتمال مجهول

بایاس تخمین‌گر

BIAS OF ESTIMATOR

- *Bias* of an estimator $\hat{\theta}$ is the difference between the expected value of $\hat{\theta}$ and θ .
- The MLE of μ is an unbiased estimator for μ because $E[\hat{\mu}] = \mu$.
- The MLE of Σ is not an unbiased estimator for Σ because $E[\hat{\Sigma}] = \frac{n-1}{n}\Sigma \neq \Sigma$.
- The *sample covariance*

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$$

is an unbiased estimator for Σ .



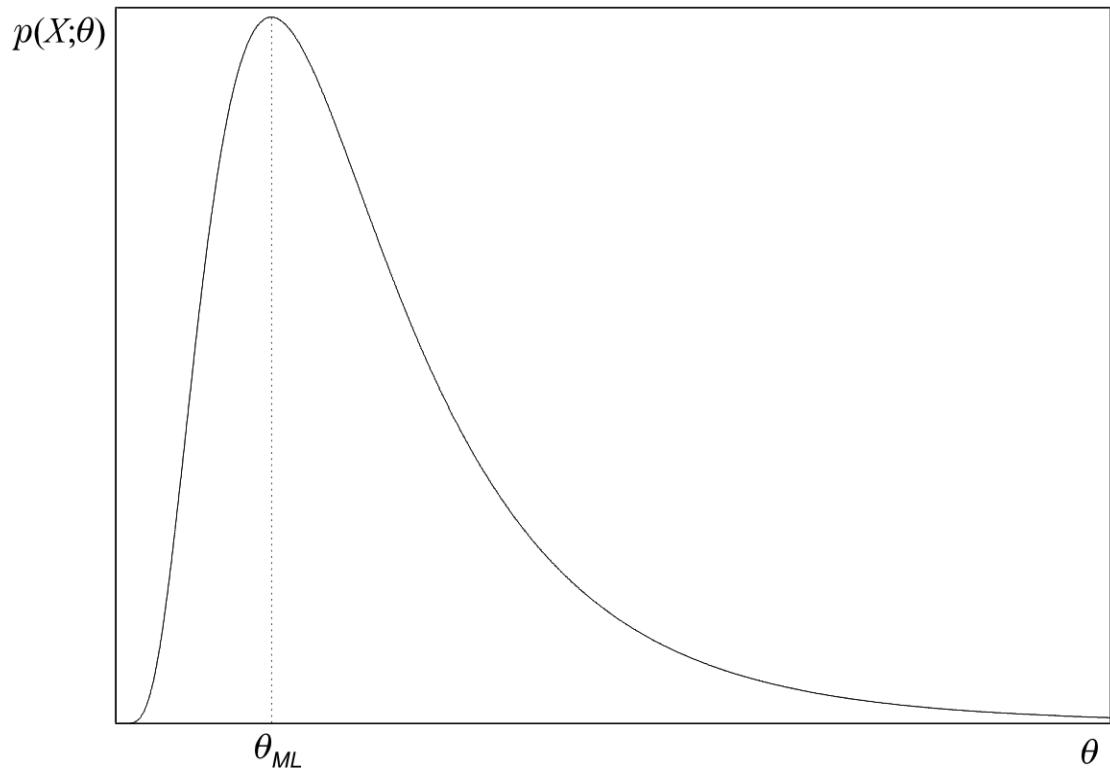
❖ Maximum Likelihood (ML)

- Let $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N$ known and independent
- Let $p(\underline{x})$ known within an unknown vector parameter $\underline{\theta}$: $p(\underline{x}) \equiv p(\underline{x}; \underline{\theta})$
- $X = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$
- $p(X; \underline{\theta}) \equiv p(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N; \underline{\theta})$
 $= \prod_{k=1}^N p(\underline{x}_k; \underline{\theta})$

which is known as the Likelihood of $\underline{\theta}$ w.r. to X

The method :

- $$\hat{\underline{\theta}}_{ML} : \arg \max_{\underline{\theta}} \prod_{k=1}^N p(\underline{x}_k; \underline{\theta})$$
- $$L(\underline{\theta}) \equiv \ln p(X; \underline{\theta}) = \sum_{k=1}^N \ln p(\underline{x}_k; \underline{\theta})$$
- $$\hat{\underline{\theta}}_{ML} : \frac{\partial L(\underline{\theta})}{\partial (\underline{\theta})} = \sum_{k=1}^N \frac{1}{p(\underline{x}_k; \underline{\theta})} \frac{\partial p(\underline{x}_k; \underline{\theta})}{\partial (\underline{\theta})} = \underline{0}$$



If, indeed, there is a $\underline{\theta}_0$ such that

$$p(\underline{x}) = p(\underline{x}; \underline{\theta}_0), \text{ then}$$

$$\lim_{N \rightarrow \infty} E[\underline{\theta}_{ML}] = \underline{\theta}_0$$

$$\lim_{N \rightarrow \infty} E\left\| \hat{\underline{\theta}}_{ML} - \underline{\theta}_0 \right\|^2 = 0$$

Asymptotically **unbiased** and **consistent**

❖ Example: $p(\underline{x}) : N(\underline{\mu}, \Sigma) : \underline{\mu}$ unknown, $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N$ $p(\underline{x}_k) \equiv p(\underline{x}_k; \underline{\mu})$

$$L(\underline{\mu}) = \ln \prod_{k=1}^N p(\underline{x}_k; \underline{\mu}) = C - \frac{1}{2} \sum_{k=1}^N (\underline{x}_k - \underline{\mu})^T \Sigma^{-1} (\underline{x}_k - \underline{\mu})$$

$$p(\underline{x}_k; \underline{\mu}) = \frac{1}{(2\pi)^{\frac{l}{2}} |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2} (\underline{x}_k - \underline{\mu})^T \Sigma^{-1} (\underline{x}_k - \underline{\mu}))$$

$$\frac{\partial L(\underline{\mu})}{\partial (\underline{\mu})} \equiv \begin{bmatrix} \frac{\partial L}{\partial \mu_1} \\ \vdots \\ \frac{\partial L}{\partial \mu_l} \end{bmatrix} = \sum_{k=1}^N \Sigma^{-1} (\underline{x}_k - \underline{\mu}) = \underline{0} \Rightarrow \underline{\mu}_{ML} = \frac{1}{N} \sum_{k=1}^N \underline{x}_k$$

Remember: if $A = A^T \Rightarrow \frac{\partial (\underline{\alpha}^T A \underline{\alpha})}{\partial \underline{\alpha}} = 2A\underline{\alpha}$

طبقه‌بندی مبتنی بر نظریه‌ی تصمیم بیز:
روش‌های تخمین پارامتری
تابع چکالی احتمال مجہول

۲

روش تخمین ماکریم احتمال پسین

تخمین پارامتری توابع چگالی احتمال مجهول

روش ماکزیمم احتمال پسین

MAXIMUM A POSTERIORI PROBABILITY (MAP)

در روش MAP، به پارامتر به عنوان یک متغیر تصادفی نگاه می‌کنیم که دارای یک توزیع پیشین معلوم است. مشاهده‌ی نمونه‌های جدید، این چگالی پیشین را به یک چگالی پسین تبدیل می‌کند.

❖ Maximum Aposteriori Probability Estimation (MAP)

- In ML method, $\underline{\theta}$ was considered as a parameter
- Here we shall look at $\underline{\theta}$ as a random vector described by a pdf $p(\underline{\theta})$, assumed to be known
- Given

$$X = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$$

Compute the maximum of

$$p(\underline{\theta}|X)$$

- From Bayes theorem $p(\underline{\theta})p(X|\underline{\theta}) = p(X)p(\underline{\theta}|X)$ or

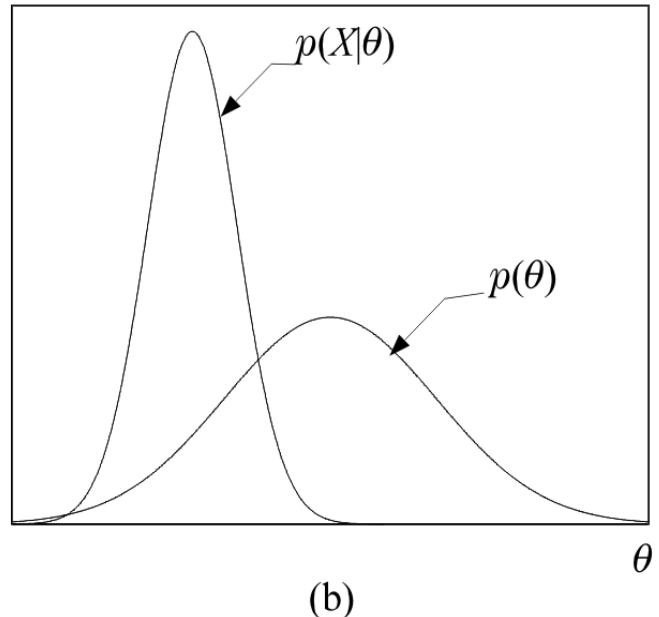
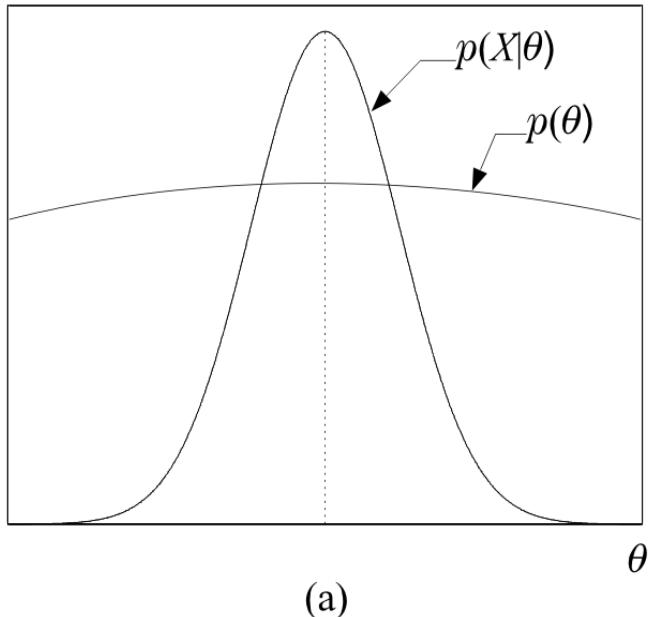
$$p(\underline{\theta}|X) = \frac{p(\underline{\theta})p(X|\underline{\theta})}{p(X)}$$

➤ The method:

$$\hat{\underline{\theta}}_{MAP} = \arg \max_{\underline{\theta}} p(\underline{\theta}|X) \text{ or}$$

$$\hat{\underline{\theta}}_{MAP} : \frac{\partial}{\partial \underline{\theta}} (P(\underline{\theta}) p(X|\underline{\theta}))$$

If $p(\underline{\theta})$ is uniform or broad enough $\hat{\underline{\theta}}_{MAP} \cong \underline{\theta}_{ML}$



ML and MAP estimates of θ will be approximately the same in (a) and different in (b).

تخمین پارامتری توابع چگالی احتمال مجهول

روش ماکزیمم احتمال پسین: مثال

MAXIMUM A POSTERIORI PROBABILITY (MAP)

N داده‌ی تولید شده با یک pdf گاوی l -بعدی داریم
 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$

که بردار میانگین آن $\boldsymbol{\mu}$ مجهول و ماتریس کوواریانس آن Σ معلوم است.
اگر توزیع پیشین بردار میانگین مجهول به صورت زیر باشد،

$$P(\boldsymbol{\mu}) = \frac{1}{(2\pi)^{l/2} \sigma_{\mu}^l} \exp\left(-\frac{1}{2} \frac{\|\boldsymbol{\mu} - \boldsymbol{\mu}_0\|^2}{\sigma_{\mu}^2}\right)$$

تخمین ماکزیمم احتمال پسین بردار میانگین مجهول را به دست آورید.

طبقه‌بندی مبتنی بر نظریه‌ی تصمیم بیز:
روش‌های تخمین پارامتری
تابع چکالی احتمال مجهول

۳

روش تخمین استنتاج بیزی

تخمین پارامتری توابع چگالی احتمال مجهول

روش استنتاج بیزی

BAYESIAN INFERENCE

روش ML و MAP، یک تخمین خاص از بردار پارامتر مجهول θ ارائه می‌کردند،

اما در استنتاج بیزی، مسیر دیگری دنبال می‌شود:

با داشتن مجموعه \mathcal{D} از بردارهای آموختشی و اطلاعات پیشین $p(\theta)$ می‌خواهیم $p(\mathbf{x}|\mathcal{D})$ را محاسبه کنیم.

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

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

$$p(\mathcal{D}|\theta) = \prod_{i=1}^n p(\mathbf{x}_i|\theta)$$

$$p(\mathbf{x}|\mathcal{D}) = \int p(\mathbf{x}|\theta) p(\theta|\mathcal{D}) d\theta$$

توزیع شرطی $p(\theta|\mathcal{D})$ نیز به عنوان تخمین pdf پسین معلوم است:
 (زیرا دانایی به روز شده در مورد خواص آماری θ پس از مشاهده مجموعه داده \mathcal{D} است.)

❖ Bayesian Inference

➤ ML, MAP \Rightarrow a single estimate for $\underline{\theta}$.

Here a different root is followed.

Given :

$$X = \{\underline{x}_1, \dots, \underline{x}_N\}, p(\underline{x}|\underline{\theta}) \text{ and } p(\underline{\theta})$$

The goal :

$$\text{estimate } p(\underline{x}|X)$$

How??

$$p(\underline{x}|X) = \int p(\underline{x}|\underline{\theta})p(\underline{\theta}|X)d\underline{\theta}$$

$$p(\underline{\theta}|X) = \frac{p(X|\underline{\theta})p(\underline{\theta})}{p(X)} = \frac{p(X|\underline{\theta})p(\underline{\theta})}{\int p(X|\underline{\theta})p(\underline{\theta})d\underline{\theta}}$$

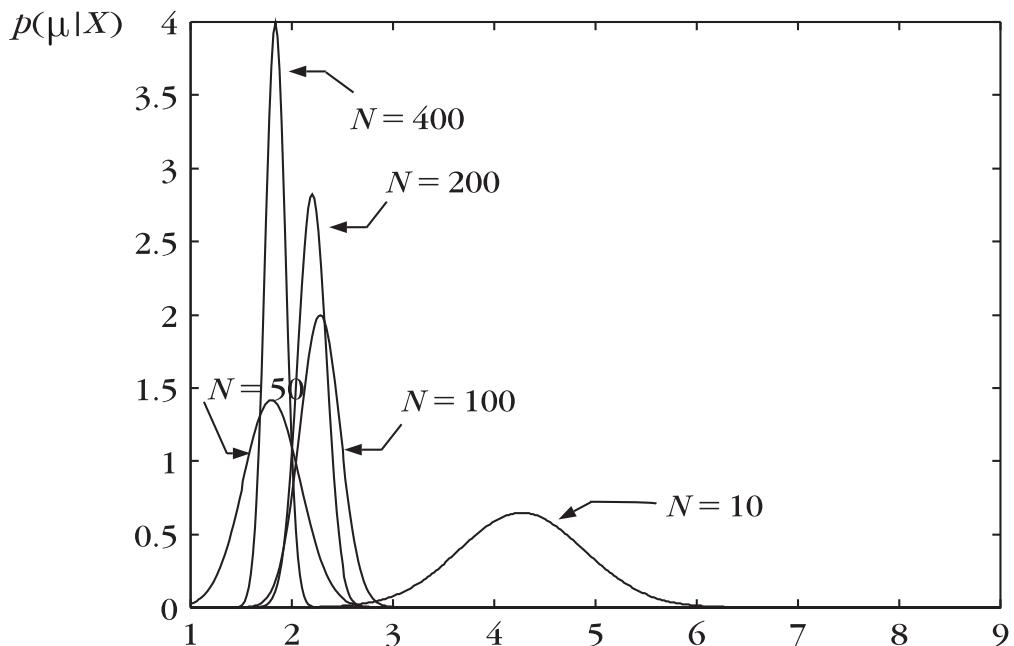
$$p(X|\underline{\theta}) = \prod_{k=1}^N p(\underline{x}_k|\underline{\theta})$$

A bit more insight via an example

- Let $p(x|\mu) \rightarrow N(\mu, \sigma^2)$
- $p(\mu) \rightarrow N(\mu_0, \sigma_0^2)$
- It turns out that: $p(\mu|X) \rightarrow N(\mu_N, \sigma_N^2)$

$$\mu_N = \frac{N\bar{x} + \sigma^2 \mu_0}{N\sigma_0^2 + \sigma^2}, \quad \sigma_N^2 = \frac{\sigma^2 \sigma_0^2}{N\sigma_0^2 + \sigma^2}, \quad \bar{x} = \frac{1}{N} \sum_{k=1}^N x_k$$

➤ The above is a sequence of Gaussians as $N \rightarrow \infty$



❖ Example:

$p(\underline{x}) : N(\underline{\mu}, \Sigma)$, $\underline{\mu}$ unknown, $X = \{\underline{x}_1, \dots, \underline{x}_N\}$

$$p(\underline{\mu}) = \frac{1}{(2\pi)^{\frac{l}{2}} \sigma_{\mu}^l} \exp\left(-\frac{\|\underline{\mu} - \underline{\mu}_0\|^2}{2\sigma_{\mu}^2}\right)$$

$$\hat{\underline{\theta}}_{MAP} : \frac{\partial}{\partial \underline{\mu}} \ln(\prod_{k=1}^N p(\underline{x}_k | \underline{\mu}) p(\underline{\mu})) = \underline{0} \text{ or } \sum_{k=1}^N \frac{1}{\sigma^2} (\underline{x}_k - \underline{\mu}) - \frac{1}{\sigma_{\mu}^2} (\hat{\underline{\mu}} - \underline{\mu}_0) = \underline{0} \Rightarrow$$

$$\hat{\underline{\mu}}_{MAP} = \frac{\underline{\mu}_0 + \frac{\sigma_{\mu}^2}{\sigma^2} \sum_{k=1}^N \underline{x}_k}{1 + \frac{\sigma_{\mu}^2}{\sigma^2} N} \text{ for } \frac{\sigma_{\mu}^2}{\sigma^2} \gg 1, \text{ or for } N \rightarrow \infty$$

$$\hat{\underline{\mu}}_{MAP} \cong \hat{\underline{\mu}}_{ML} = \frac{1}{N} \sum_{k=1}^N \underline{x}_k$$

طبقه‌بندی مبتنی بر نظریه‌ی تصمیم بیز:
روش‌های تخمین پارامتری
تابع چکالی احتمال مجهول

۴

روش تخمین ماکزیمم آنتروپی

تخمین پارامتری توابع چگالی احتمال مجهول

روش ماکزیمم آنتروپی

MAXIMUM ENTROPY (ME)

روش ME تابع چگالی احتمالی را می‌یابد که آنتروپی را ماکزیمم می‌کند.
(منتظر با توزیعی که بالاترین میزان تصادفی بودن ممکن را نشان می‌دهد)

❖ Maximum Entropy

➤ Entropy

$$H = - \int p(\underline{x}) \ln p(\underline{x}) d\underline{x}$$

➤ $\hat{p}(x)$: maximum H

subject to the available constraints

➤ Example: x is nonzero in the interval $[x_1, x_2]$ and zero otherwise. Compute the ME pdf

- The constraint: $\int_{x_1}^{x_2} p(x)dx = 1$
- Lagrange Multipliers $H_L = H + \lambda(\int_{x_1}^{x_2} p(x)dx - 1)$

$$\hat{p}(x) = \exp(\lambda - 1)$$

- $\hat{p}(x) = \begin{cases} \frac{1}{x_2 - x_1} & x_1 \leq x \leq x_2 \\ 0 & \text{otherwise} \end{cases}$

طبقه‌بندی مبتنی بر نظریه‌ی تصمیم بیز:
روش‌های تخمین پارامتری
تابع چکالی احتمال مجہول

۵

روش مدل‌های مخلط

تخمین پارامتری توابع چگالی احتمال مجهول

روش مدل‌های چگالی مخلوط

MIXTURE DENSITY MODELS

- A mixture model is a linear combination of m densities

$$p(\mathbf{x}|\Theta) = \sum_{j=1}^m \alpha_j p_j(\mathbf{x}|\theta_j)$$

where $\Theta = (\alpha_1, \dots, \alpha_m, \theta_1, \dots, \theta_m)$ such that $\alpha_j \geq 0$
and $\sum_{j=1}^m \alpha_j = 1$.

- $\alpha_1, \dots, \alpha_m$ are called the mixing parameters.
- $p_j(\mathbf{x}|\theta_j)$, $j = 1, \dots, m$ are called the component densities.

فرض می‌شود یک توزیع pdf مجهول، از طریق ترکیب خطی m تابع چگالی ساخته می‌شود.

* این مدل به طور ضمنی فرض می‌کند که

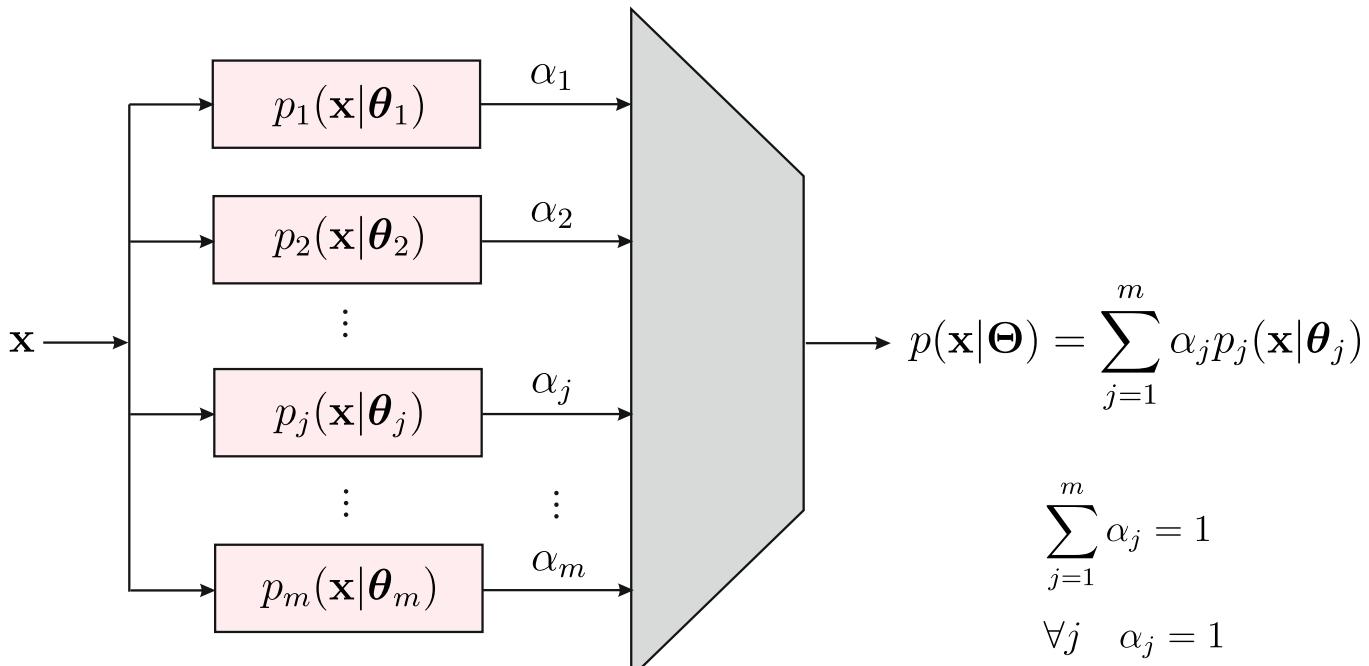
هر نقطه‌ی x می‌تواند از هر یک از m توزیع با احتمال α_j بیرون کشیده شود.

* این مدل می‌تواند هر تابع چگالی پیوسته را با تعداد کافی مؤلفه و با پارامترهای مناسب تقریب بزند.

تخمین پارامتری توابع چگالی احتمال مجهول

روش مدل‌های چگالی مخلوط

MIXTURE DENSITY MODELS

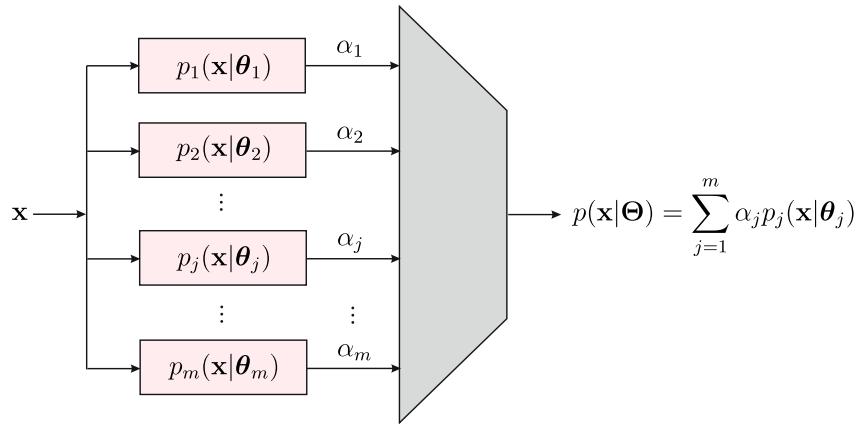


$$\boldsymbol{\Theta} = [\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_j, \dots, \boldsymbol{\theta}_m, \alpha_1, \alpha_2, \dots, \alpha_j, \dots, \alpha_m]$$

تخمین پارامتری توابع چگالی احتمال مجهول

روش مدل‌های چگالی مخلوط: تخمین پارامترها

MIXTURE DENSITY MODELS



چون مشخص نیست هر داده‌ی آموزشی \mathbf{x} از کدام مؤلفه‌ی مخلوط می‌آید، استفاده از ML برای تخمین پارامترها به یک مسئله‌ی بهینه‌سازی غیرخطی پیچیده منجر می‌شود.

(اگر مشخص بود، m مسئله‌ی ML جداگانه حل می‌شد.)



با یک مسئله با مجموعه داده‌ی آموزشی ناکامل سروکار داریم (اطلاعات برچسب ناقص است)



از روش **ماکزیمم‌سازی امید (EM)** استفاده می‌کنیم.

❖ Mixture Models

➤
$$p(\underline{x}) = \sum_{j=1}^J p(\underline{x}|j)P_j \quad \sum_{j=1}^M P_j = 1, \int_{\underline{x}} p(\underline{x}|j)d\underline{x} = 1$$

- Assume parametric modeling, i.e., $p(\underline{x}|j; \underline{\theta})$
- The goal is to estimate $\underline{\theta}$ and P_1, P_2, \dots, P_j
given a set $X = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$
- Why not ML? As before? $\max_{\underline{\theta}, P_i, \dots, P_j} \prod_{k=1}^N P(\underline{x}_k; \underline{\theta}, P_i, \dots, P_j)$

- This is a **nonlinear problem** due to the missing label information.
This is a typical problem with **an incomplete data set**.

Solution??

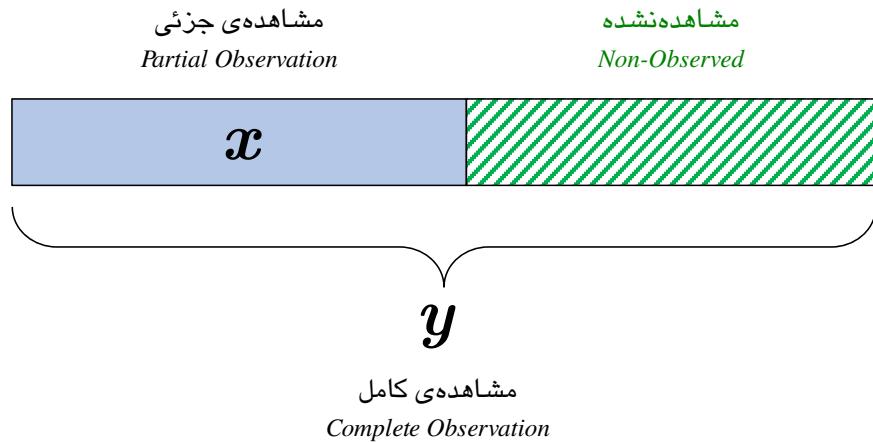
- The **Expectation-Maximization (EM)** algorithm.

تخمین پارامتری توابع چگالی احتمال مجهول

ماکریزم‌سازی امید

EXPECTATION MAXIMIZATION (EM)

روشی کارآمد برای حل عددی مسئله‌ی ML است، وقتی که مشاهدات جزئی در دست است.



الگوریتم EM مقدار امید تابع درست‌نمایی را به شرط نمونه‌های مشاهده شده و تخمین فعلی θ ماکریزم می‌کند:

$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta} E \left\{ \sum_k \log p_y(\mathbf{y}_k | \boldsymbol{\theta}) \right\}$$

که در آن امید ریاضی بر روی بخش تصادفی y (قسمت مشاهده نشده) محاسبه می‌شود.

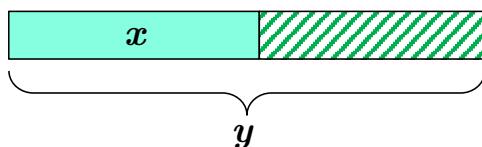
➤ The Expectation-Maximization (EM) algorithm

- General formulation
 - \underline{y} the complete data set $\underline{y} \in Y \subseteq \Re^m$, with $p_{\underline{y}}(\underline{y}; \underline{\theta})$, which are **not observed directly**.

We observe

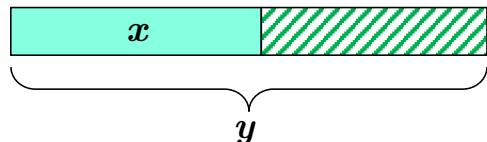
$$\underline{x} = g(\underline{y}) \in X_{ob} \subseteq \Re^l, l < m \text{ with } P_x(\underline{x}; \underline{\theta}),$$

a many to one transformation



- Let $Y(\underline{x}) \subseteq Y$ all \underline{y} 's \rightarrow to a specific \underline{x}

$$p_{\underline{x}}(\underline{x}; \underline{\theta}) = \int_{Y(\underline{x})} p_{\underline{y}}(\underline{y}; \underline{\theta}) d\underline{y}$$



- What we need is to compute

$$\hat{\theta}_{ML} : \sum_k \frac{\partial \ln(p_{\underline{y}}(\underline{y}_k; \underline{\theta}))}{\partial \underline{\theta}} = 0$$

- But \underline{y}_k 's are not observed.
Here comes the EM.
Maximize the **expectation** of the log-likelihood
conditioned on the observed samples
and the current iteration estimate of $\underline{\theta}$.

➤ The algorithm:

- E-step:

$$Q(\underline{\theta}; \underline{\theta}(t)) = E\left[\sum_k \ln(p_{\underline{y}}(y_k; \underline{\theta}|X; \underline{\theta}(t)))\right]$$

- M-step:

$$\underline{\theta}(t+1) = \arg \max_{\underline{\theta}} Q(\underline{\theta}; \underline{\theta}(t))$$

$$\underline{\theta}(t+1) : \frac{\partial Q(\underline{\theta}; \underline{\theta}(t))}{\partial \underline{\theta}} = 0$$

تخمین پارامتری توابع چگالی احتمال مجهول

ماکزیمم‌سازی امید: الگوریتم

EXPECTATION MAXIMIZATION (EM)

Algorithm (Expectation-Maximization)

```

1 begin initialize  $\theta^0, T, i = 0$ 
2           do  $i \leftarrow i + 1$ 
3             E step : compute  $Q(\theta; \theta^i)$ 
5             M step :  $\theta^{i+1} \leftarrow \arg \max_{\theta} Q(\theta; \theta^i)$ 
6             until  $Q(\theta^{i+1}; \theta^i) - Q(\theta^i; \theta^{i-1}) \leq T$ 
7           return  $\hat{\theta} \leftarrow \theta^{i+1}$ 
8 end

```

تخمین پارامتری توابع چگالی احتمال مجهول

ماکریزم‌سازی امید: نکته‌ها

EXPECTATION MAXIMIZATION (EM)



در الگوریتم EM تخمین‌های پی‌درپی ($\theta(t)$) هرگز مقدار تابع درستنمایی را کاهش نمی‌دهد.

تابع درستنمایی صعودی می‌ماند تا به یک ماکریزم (محلی/سراسری) برسد و EM همگرا شود.



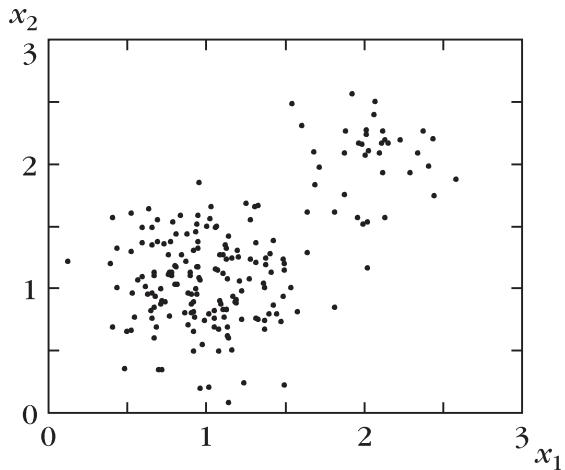
همگرایی الگوریتم EM کندر از همگرایی الگوریتم‌های جستجو مانند نیوتن است.

اما همگرایی آن ملايم است و نسبت به ناپایداری آسیب‌پذیری ندارد (قوام بالاتر).

تخمین پارامتری توابع چگالی احتمال مجهول

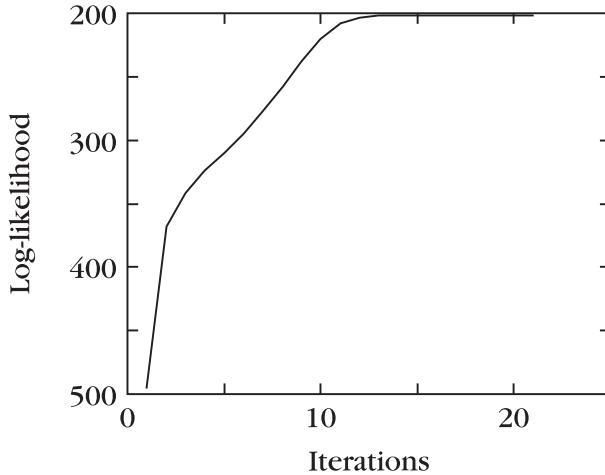
ماکریزم‌سازی امید: مثال

EXPECTATION MAXIMIZATION (EM)



(a)

مجموعه‌ی داده‌های آموزشی (دوبعدی)



(b)

تابع log-likelihood محاسبه شده در تکرارهای پی‌درپی الگوریتم EM (همیشه غیرنزولی)

تخمین پارامتری توابع چگالی احتمال مجهول

مدل‌های مخلوط گاوی

GAUSSIAN MIXTURE MODELS (GMM)

یک مدل چگالی مخلوط که مؤلفه‌های آن گاوی باشد.

مدل مخلوط گاوی

Gaussian Mixture Model

یک راه برآورد تعداد مؤلفه‌ها، تعداد قله‌های موجود در pdf است.

➤ Application to the mixture modeling problem

- Complete data

$$(\underline{x}_k, j_k), k = 1, 2, \dots, N$$

- Observed data

$$\underline{x}_k, k = 1, 2, \dots, N$$

$$p(\underline{x}_k, j_k; \underline{\theta}) = p(\underline{x}_k | j_k; \underline{\theta}) P_{jk}$$

- Assuming mutual independence, the log-likelihood is

$$L(\underline{\theta}) = \sum_{k=1}^N \ln(p(\underline{x}_k | j_k; \underline{\theta}) P_{jk})$$

➤ Unknown parameters

$$\underline{\Theta}^T = [\underline{\theta}^T, \underline{P}^T]^T, \quad \underline{P} = [P_1, P_2, \dots, P_j]^T$$

➤ E-step

$$Q(\underline{\Theta}; \underline{\Theta}(t)) = E\left[\sum_{k=1}^N \ln(p(\underline{x}_k | j_k; \underline{\theta}) P_{jk})\right] = \sum_{k=1}^N E\left[\sum_{j_k=1}^J P(j_k | \underline{x}_k; \underline{\Theta}(t)) \ln(p(\underline{x}_k | j_k; \underline{\theta}) P_{jk})\right]$$

➤ M-step

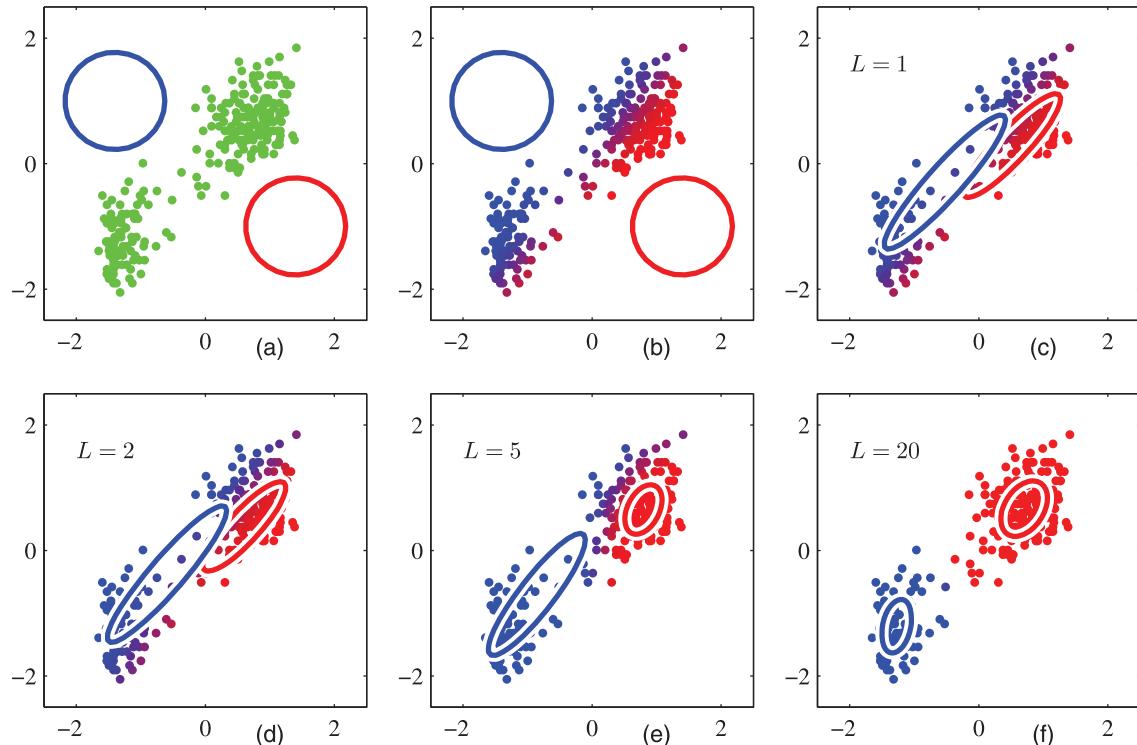
$$\frac{\partial Q}{\partial \underline{\theta}} = 0 \quad \frac{\partial Q}{\partial P_{jk}} = 0, \quad j_k = 1, 2, \dots, J$$

$$P(j | \underline{x}_k; \Theta(t)) = \frac{p(\underline{x}_k | j; \underline{\Theta}(t)) P_j}{P(\underline{x}_k; \underline{\Theta}(t))} \quad p(\underline{x}_k; \underline{\Theta}(t)) = \sum_{j=1}^J p(\underline{x}_k | j; \underline{\Theta}(t)) P_j$$

تخمین پارامتری توابع چگالی احتمال مجهول

ماکزیمم‌سازی امید: مثال

EXPECTATION MAXIMIZATION (EM)



مدل‌های چگالی مخلوط

کاربرد در خوشبندی

MIXTURE DENSITY MODELS: APPLICATION IN CLUSTERING

می‌توان از مدل‌های مخلوط در کاربرد خوشبندی استفاده کرد.

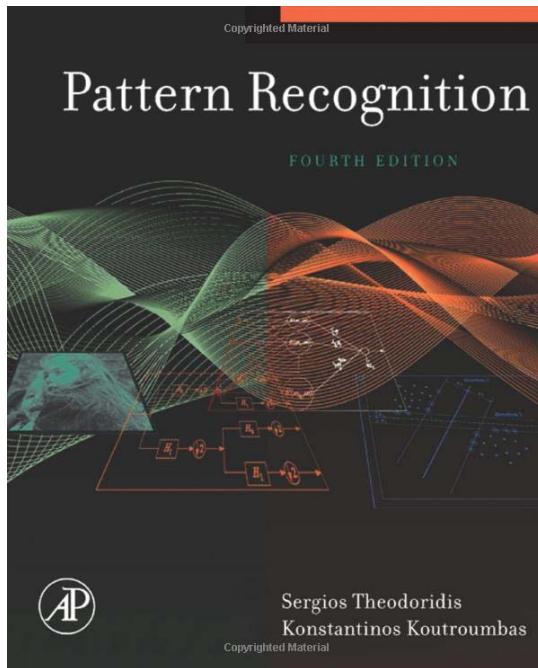
هر مؤلفه را می‌توان یک خوش در نظر گرفت:
وقتی مشخص شود که احتمال تولید هر \mathbf{x} توسط کدام مؤلفه (j) بیشتر است،
عملاً خوشی آن \mathbf{x} مشخص می‌شود.

طبقه‌بندی مبتنی بر نظریه‌ی تصمیم بیز:
روش‌های تخمین پارامتری
تابع چکالی احتمال مجہول

۶

منابع

منبع اصلی



S. Theodoridis, K. Koutroumbas,
Pattern Recognition,
Fourth Edition, Academic Press, 2009.

Chapter 2

CHAPTER

2

Classifiers Based on Bayes Decision Theory

2.1 INTRODUCTION

This is the first chapter, out of three, dealing with the design of the classifier in a pattern recognition system. The approach to be followed builds upon probabilistic arguments stemming from the statistical nature of the generated features. As has already been pointed out in the introductory chapter, this is due to the statistical variation of the patterns as well as to the noise in the measuring sensors. Adopting this reasoning as our kickoff point, we will design classifiers that classify an unknown pattern in the most probable of the classes. Thus, our task now becomes that of defining what “most probable” means.

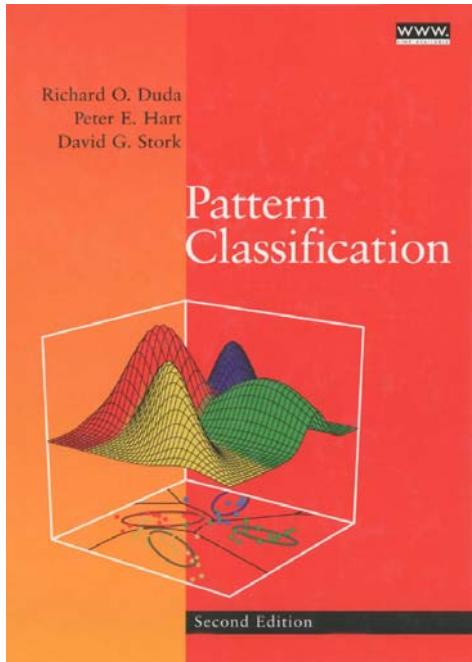
Given a classification task of M classes, $\omega_1, \omega_2, \dots, \omega_M$, and an unknown pattern, which is represented by a feature vector x , we form the M conditional probabilities $P(\omega_i|x), i = 1, 2, \dots, M$. Sometimes, these are also referred to as *a posteriori probabilities*. In words, each of them represents the probability that the unknown pattern belongs to the respective class ω_i , given that the corresponding feature vector takes the value x . Who could then argue that these conditional probabilities are not sensible choices to quantify the term *most probable*? Indeed, the classifiers to be considered in this chapter compute either the maximum of these M values or, equivalently, the maximum of an appropriately defined function of them. The unknown pattern is then assigned to the class corresponding to this maximum.

The first task we are faced with is the computation of the conditional probabilities. The Bayes rule will once more prove its usefulness! A major effort in this chapter will be devoted to techniques for estimating probability density functions (pdf), based on the available experimental evidence, that is, the feature vectors corresponding to the patterns of the training set.

2.2 BAYES DECISION THEORY

We will initially focus on the two-class case. Let ω_1, ω_2 be the two classes in which our patterns belong. In the sequel, we assume that the *a priori probabilities*

13



R.O. Duda, P.E. Hart, and D.G. Stork,
Pattern Classification,
 Second Edition, John Wiley & Sons, Inc., 2001.

Chapter 3

CHAPTER
3

MAXIMUM-LIKELIHOOD AND BAYESIAN PARAMETER ESTIMATION

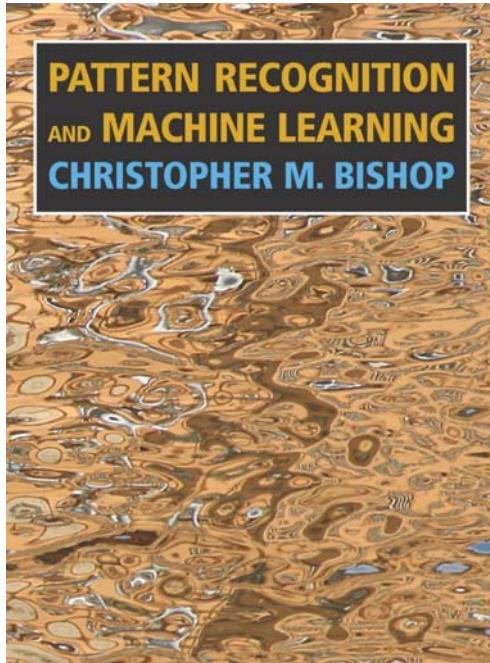
3.1 INTRODUCTION

In Chapter 2 we saw how we could design an optimal classifier if we knew the prior probabilities $P(\omega_i)$ and the class-conditional densities $p(x|\omega_i)$. Unfortunately, in pattern recognition applications we rarely, if ever, have this kind of complete knowledge about the probabilistic structure of the problem. In a typical case we merely have some vague, general knowledge about the situation, together with a number of *design samples* or *training data*—particular representatives of the patterns we want to classify. The problem, then, is to find some way to use this information to design or train the classifier.

One approach to this problem is to use the samples to estimate the unknown probabilities and probability densities, and then use the resulting estimates as if they were the true values. In typical supervised pattern classification problems, the estimation of the prior probabilities presents no serious difficulties (Problem 3). However, estimation of the class-conditional densities is quite another matter. The number of available samples always seems too small, and serious problems arise when the dimensionality of the feature vector x is large. If we know the number of parameters in advance and our general knowledge about the problem permits us to parameterize the conditional densities, then the severity of these problems can be reduced significantly. Suppose, for example, that we can reasonably assume that $p(x|\omega_i)$ is a normal density with mean μ_i and covariance matrix Σ_i , although we do not know the exact values of these quantities. This knowledge simplifies the problem from one of estimating an unknown function $p(x|\omega_i)$ to one of estimating the parameters μ_i and Σ_i .

The problem of parameter estimation is a classical one in statistics, and it can be approached in several ways. We shall consider two common and reasonable procedures, namely, *maximum-likelihood* estimation and *Bayesian* estimation. Although the results obtained with these two procedures are frequently nearly identical,

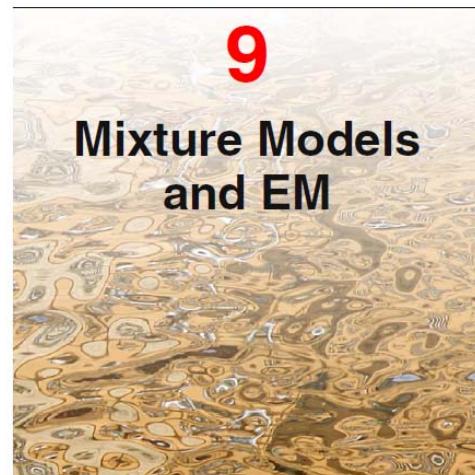
MAXIMUM-
 LIKELIHOOD
 BAYESIAN
 ESTIMATION



C.M. Bishop,
Pattern Recognition and Machine Learning,
Springer, 2006.

Chapter 9

Section 9.1



If we define a joint distribution over observed and latent variables, the corresponding distribution of the observed variables alone is obtained by marginalization. This allows relatively complex marginal distributions over observed variables to be expressed in terms of more tractable joint distributions over the expanded space of observed and latent variables. The introduction of latent variables thereby allows complicated distributions to be formed from simpler components. In this chapter, we shall see that mixture distributions, such as the Gaussian mixture discussed in Section 2.3.9, can be interpreted in terms of discrete latent variables. Continuous latent variables will form the subject of Chapter 12.

As well as providing a framework for building more complex probability distributions, mixture models can also be used to cluster data. We therefore begin our discussion of mixture distributions by considering the problem of finding clusters in a set of data points, which we approach first using a nonprobabilistic technique called the *K*-means algorithm (Lloyd, 1982). Then we introduce the latent variable