

Pattern Classification

All materials in these slides were taken from Pattern Classification (2nd ed) by R. O. Duda, P. E. Hart and D. G. Stork, John Wiley & Sons, 2000 with the permission of the authors and the publisher

Chapter 3: Maximum-Likelihood & Bayesian Parameter Estimation (part 1)

Introduction

Maximum-Likelihood Estimation

- Example of a Specific Case
- \bullet The Gaussian Case: unknown μ and σ
- Bias
- Appendix: ML Problem Statement

Introduction

Data availability in a Bayesian framework

- We could design an optimal classifier if we knew:
 - P(ω_i) (priors)
 - P(x | ω_i) (class-conditional densities)
 Unfortunately, we rarely have this complete information!

• Design a classifier from a training sample

- No problem with prior estimation
- Samples are often too small for class-conditional estimation (large dimension of feature space!)

• A priori information about the problem

- Do we know something about the distribution?
- \rightarrow find parameters to characterize the distribution
- Example: Normality of $P(x \mid \omega_i)$

 $\mathsf{P}(\mathsf{x} \mid \omega_{\mathsf{i}}) \sim \mathsf{N}(\mu_{\mathsf{i}}, \Sigma_{\mathsf{i}})$

- Characterized by 2 parameters
- Estimation techniques
 - Maximum-Likelihood (ML) and the Bayesian estimations
 - Results are nearly identical, but the approaches are different

- Parameters in ML estimation are fixed but unknown!
 - Best parameters are obtained by maximizing the probability of obtaining the samples observed
- Bayesian methods view the parameters as random variables having some known distribution
- In either approach, we use P(ω_i | x) for our classification rule!

Maximum-Likelihood Estimation

- Has good convergence properties as the sample size increases
- Simpler than any other alternative techniques

General principle

• Assume we have *c* classes and $P(x \mid \omega_j) \sim N(\mu_j, \Sigma_j)$ $P(x \mid \omega_j) \equiv P(x \mid \omega_j, \theta_j)$ where:

$$\theta = (\mu_j, \Sigma_j) = (\mu_j^1, \mu_j^2, ..., \sigma_j^{11}, \sigma_j^{22}, cov(x_j^m, x_j^n)...)$$

Use the information

provided by the training samples to estimate $\theta = (\theta_1, \theta_2, ..., \theta_c)$, each θ_i (i = 1, 2, ..., c) is associated with each category

Suppose that D contains n samples, x₁, x₂,..., x_n

$$P(D \mid \theta) = \prod_{k=1}^{k=n} P(x_k \mid \theta) = F(\theta)$$

P(D \mid \theta) is called the likelihood of θ w.r.t. the set of samples)

ML estimate of θ is, by definition the value that θ maximizes P(D | θ)
 "It is the value of θ that best agrees with the actually observed training sample"



FIGURE 3.1. The top graph shows several training points in one dimension, known or assumed to be drawn from a Gaussian of a particular variance, but unknown mean. Four of the infinite number of candidate source distributions are shown in dashed lines. The middle figure shows the likelihood $p(\mathcal{D}|\theta)$ as a function of the mean. If we had a very large number of training points, this likelihood would be very narrow. The value that maximizes the likelihood is marked $\hat{\theta}$; it also maximizes the logarithm of the likelihood—that is, the log-likelihood $I(\theta)$, shown at the bottom. Note that even though they look similar, the likelihood $p(\mathcal{D}|\theta)$ is shown as a function of θ whereas the conditional density $p(x|\theta)$ is shown as a function of x. Furthermore, as a function of θ , the likelihood $p(\mathcal{D}|\theta)$ is not a probability density function and its area has no significance. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Optimal estimation

• Let $\theta = (\theta_1, \theta_2, ..., \theta_p)^t$ and let ∇_{θ} be the gradient operator

$$\nabla_{\boldsymbol{\theta}} = \left[\frac{\partial}{\partial \theta_1}, \frac{\partial}{\partial \theta_2}, \dots, \frac{\partial}{\partial \theta_p}\right]^t$$

- We define *l*(θ) as the log-likelihood function
 l(θ) = In P(D | θ)
 (recall D is the training data)
- New problem statement:
 determine θ that maximizes the log-likelihood

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} l(\boldsymbol{\theta})$$

The definition of *I*() is:

$$l(\mathbf{\theta}) = \sum_{k=1}^{n} \ln p(\mathbf{x}_k \mid \mathbf{\theta})$$

and

$$(\nabla_{\boldsymbol{\theta}} l = \sum_{k=1}^{k=n} \nabla_{\boldsymbol{\theta}} \ln P(\mathbf{x}_k \mid \boldsymbol{\theta})) \quad (\text{eq } 6)$$

Set of necessary conditions for an optimum is:

$$\nabla_{\theta} I = 0 \quad (eq. 7)$$

- Example, the Gaussian case: unknown μ
 - We assume we know the covariance
 - *p*(x_i | μ) ~ N(μ, Σ) (Samples are drawn from a multivariate normal population)

$$\ln p(\mathbf{x}_{k} | \mathbf{\mu}) = -\frac{1}{2} \ln \left[(2\pi)^{d} | \Sigma | \right] - \frac{1}{2} (\mathbf{x}_{k} - \mathbf{\mu})^{t} \Sigma^{-1} (\mathbf{x}_{k} - \mathbf{\mu})$$

and $\nabla_{\mathbf{\mu}} \ln p(\mathbf{x}_{k} | \mathbf{\mu}) = \Sigma^{-1} (\mathbf{x}_{k} - \mathbf{\mu})$ (eq. 9)

 $\theta = \mu$ therefore: The ML estimate for μ must satisfy: $\sum_{k=1}^{n} \Sigma^{-1}(\mathbf{x}_{k} - \hat{\boldsymbol{\mu}}) = \mathbf{0} \text{ from eqs 6,7 \& 9}$

Multiplying by Σ and rearranging, we obtain:



Just the arithmetic average of the samples of the training samples!

Conclusion:

If $P(x_k \mid \omega_j)$ (j = 1, 2, ..., c) is supposed to be Gaussian in a *d*dimensional feature space; then we can estimate the vector $\theta = (\theta_1, \theta_2, ..., \theta_c)^t$ and perform an optimal classification!

• Example, Gaussian Case: unknown μ and Σ • First consider univariate case: unknown μ and σ $\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$ $l = \ln p(x_k | \mathbf{\theta}) = -\frac{1}{2} \ln 2\pi\theta_2 - \frac{1}{2\theta_2} (x_k - \theta_1)^2$ $\nabla_{\theta} l = \begin{pmatrix} \frac{\sigma}{\sigma \theta_1} (\ln P(x_k \mid \boldsymbol{\theta})) \\ \frac{\sigma}{\sigma \theta_2} (\ln P(x_k \mid \boldsymbol{\theta})) \end{pmatrix}$ $\begin{cases} \frac{1}{\theta_2} (x_k - \theta_1) \\ -\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2} \end{cases}$

12

Summation (over the training set):

$$\begin{cases} \sum_{k=1}^{n} \frac{1}{\hat{\theta}_{2}} (x_{k} - \hat{\theta}_{1}) = 0 \\ -\sum_{k=1}^{n} \frac{1}{\hat{\theta}_{2}} + \sum_{k=1}^{n} \frac{(x_{k} - \hat{\theta}_{1})^{2}}{\hat{\theta}_{2}^{2}} = 0 \end{cases}$$
(1)

Combining (1) and (2), one obtains:

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^{n} x_k$$
; $\hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^{n} (x_k - \hat{\mu})^2$

Pattern Classification, Chapte23

- The ML estimates for the multivariate case is similar
 - The scalars χ and μ are replaced with vectors
 - The variance σ^2 is replaced by the covariance matrix

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{k=1}^{n} \mathbf{x}_{k}$$
$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{k=1}^{n} (\mathbf{x}_{k} - \hat{\boldsymbol{\mu}}) (\mathbf{x}_{k} - \hat{\boldsymbol{\mu}})^{t}$$

Bias

• ML estimate for σ^2 is biased

$$E\left[\frac{1}{n}\sum_{i=1}^{n}(x_{i}-\overline{x})^{2}\right] = \frac{n-1}{n}\sigma^{2} \neq \sigma^{2}$$

• Extreme case: n=1, E[] = $0 \neq \sigma^2$

As *n* increases the bias is reduced
 → this type of estimator is called *asymptotically* unbiased

• An elementary unbiased estimator for Σ is:

$$\mathbf{C} = \frac{1}{n-1} \sum_{k=1}^{n} (\mathbf{x}_{k} - \hat{\boldsymbol{\mu}}) (\mathbf{x}_{k} - \hat{\boldsymbol{\mu}})^{t}$$

Sample covariance matrix

This estimator is unbiased for all distributions
→ Such estimators are called absolutely unbiased

• Our earlier estimator for Σ is biased:

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{k=1}^{n} (\mathbf{x}_k - \hat{\boldsymbol{\mu}}) (\mathbf{x}_k - \hat{\boldsymbol{\mu}})^t$$

In fact it is asymptotically unbiased:

Observe that

$$\hat{\Sigma} = \frac{n-1}{n}C$$

• Appendix: ML Problem Statement

• Let $D = \{x_1, x_2, ..., x_n\}$

 $\mathsf{P}(\mathsf{x}_1, \dots, \mathsf{x}_n \mid \theta) = \Pi^{1,n} \mathsf{P}(\mathsf{x}_k \mid \theta); \mid D \mid = n$

Our goal is to determine $\hat{\theta}$ (value of θ that maximizes the likelihood of this sample set!)



Pattern Classification, Chapte23

19

$\theta = (\theta_1, \theta_2, ..., \theta_c)$ Problem: find $\hat{\theta}$ such that:

$$\begin{aligned} \underset{\theta}{\text{Max}P(D \mid \theta) &= \text{Max}P(x_1, ..., x_n \mid \theta) \\ &= \text{Max} \prod_{k=1}^{n} P(x_k \mid \theta) \end{aligned}$$

20

Pattern Classification, Chapte23