

DENSITY ESTIMATION

Estimating an Unknown Probability Density Function

Th&Ko §2.5 / DHS §3.1-3.4

Density Estimation

- Parametric techniques
 - Maximum Likelihood
 - Maximum A Posteriori
 - Bayesian Inference
 - **Gaussian Mixture Models (GMM)**
 - EM-Algorithm
- Non-parametric techniques
 - Histogram
 - Parzen Windows
 - k-nearest-neighbor rule

Estimation of an unknown PDF

2

Task:

Estimate the parameters $\underline{\theta}$ of a pdf with known structure from a set of data X .

(e.g. $p(\underline{x}) \equiv p(\underline{x}; \underline{\theta})$ is known to be Gaussian $N(\underline{\mu}, \underline{\Sigma})$, with unknown $\underline{\theta} = [\mu_1, \dots, \mu_I, \sigma_{1,1}, \dots, \sigma_{I,I}]^T$)

Formal:

Let $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N$ known and independent (i.i.d.) $X = \{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$

Let $p(\underline{x})$ be known within a vector parameter $\underline{\theta}$: $p(\underline{x}) \equiv p(\underline{x}; \underline{\theta})$

$p(X; \underline{\theta}) \equiv p(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N; \underline{\theta})$

$= \prod_{k=1}^N p(\underline{x}_k; \underline{\theta})$ which is known as the likelihood of $\underline{\theta}$ w.r. to X

Estimation of an unknown PDF

3

❖ Maximum Likelihood Method (ML)

➤ Search for the parameter Θ_{ML} that maximizes

$$\hat{\underline{\theta}}_{ML} := \arg \max_{\underline{\theta}} \prod_{k=1}^N p(\underline{x}_k; \underline{\theta})$$

– a necessary condition $\frac{\partial \prod_{k=1}^N p(\underline{x}_k; \underline{\theta})}{\partial (\underline{\theta})} = \underline{0}$

– since \ln is monotonic we can write the log-likelihood

$$L(\underline{\theta}) \equiv \ln p(X; \underline{\theta}) = \sum_{k=1}^N \ln p(\underline{x}_k; \underline{\theta})$$

$$\hat{\underline{\theta}}_{ML} : \frac{\partial L(\underline{\theta})}{\partial (\underline{\theta})} = \sum_{k=1}^N \frac{1}{p(\underline{x}_k; \underline{\theta})} \frac{\partial p(\underline{x}_k; \underline{\theta})}{\partial (\underline{\theta})} = \underline{0}$$

Properties of Maximum Likelihood Method

If there is a true value Θ_0 , the ML estimate Θ_{ML} has the following properties (no proofs):

a) Θ_{ML} is *asymptotically unbiased and converges in the mean*.

$$p(\underline{x}) = p(\underline{x}; \underline{\theta}_0), \text{ then } \lim_{N \rightarrow \infty} E[\hat{\underline{\theta}}_{ML}] = \underline{\theta}_0$$

b) Θ_{ML} is *asymptotically consistent and converges in probability*.

$$\lim_{N \rightarrow \infty} \text{prob} \left\{ \left\| \hat{\underline{\theta}}_{ML} - \underline{\theta}_0 \right\| \leq \varepsilon \right\} = 1$$

c) Θ_{ML} is *asymptotically consistent and converges in mean square*.

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

ML Example 1:

$p(\underline{x}) : N(\underline{\mu}, \Sigma) : \underline{\mu}$ unknown

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

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

$$\frac{\partial L(\underline{\mu})}{\partial(\underline{\mu})} \equiv \begin{bmatrix} \frac{\partial L}{\partial \mu_1} \\ \frac{\partial L}{\partial \mu_2} \\ \vdots \\ \frac{\partial L}{\partial \mu_i} \end{bmatrix} = \sum_{k=1}^N \Sigma^{-1}(\underline{x}_k - \underline{\mu}) = \underline{0}$$

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

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

ML Example 2 :

$p(\underline{x})$: $N(\underline{\mu}, \Sigma)$: $\underline{\mu}, \Sigma = \sigma^2 I$ unknown

$\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N \in \mathbb{R}^l$

$$p(\underline{x}_k) \equiv p(\underline{x}_k; \underline{\mu}, \sigma^2) = \frac{1}{(2\pi^l \sigma^2)^{l/2}} \exp\left(-\frac{1}{2\sigma^2}(\underline{x}_k - \underline{\mu})^2\right)$$

$$L(\underline{\mu}, \sigma^2) = \ln \prod_{k=1}^N p(\underline{x}_k; \underline{\mu}, \sigma^2) = \sum_{k=1}^N -\frac{1}{2} \ln(2\pi^l \sigma^2) - \sum_{k=1}^N \frac{1}{2\sigma^2} (\underline{x}_k - \underline{\mu})^2$$

$$\frac{\partial L(\theta)}{\partial(\theta)} \equiv \begin{bmatrix} \frac{\partial L}{\partial \underline{\mu}} \\ \frac{\partial L}{\partial \sigma^2} \end{bmatrix} = \begin{bmatrix} \sum_{k=1}^N \frac{1}{\sigma^2} (\underline{x}_k - \underline{\mu}) \\ -\sum_{k=1}^N \frac{1}{2\sigma^2} + \sum_{k=1}^N \frac{1}{2\sigma^4} (\underline{x}_k - \underline{\mu})^2 \end{bmatrix} = \underline{0}$$

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

$$\Rightarrow \sigma_{ML}^2 = \frac{1}{N} \sum_{k=1}^N (\underline{x}_k - \underline{\mu})^2$$

However, the true $\underline{\mu}$ is unknown, therefore we have to use $\underline{\mu}_{ML}$

Maximum Likelihood estimates are only asymptotically unbiased, so N should be large enough !

ML Example (3) :

$p(\underline{x})$: $N(\underline{\mu}, \Sigma)$: $\underline{\mu}, \Sigma = \sigma^2 I$ unknown

$\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N \in \mathbb{R}^l$

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

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

$$\frac{\partial L(\theta)}{\partial(\theta)} \equiv \begin{bmatrix} \frac{\partial L}{\partial \underline{\mu}} \\ \frac{\partial L}{\partial \Sigma_{11}} \\ \vdots \\ \frac{\partial L}{\partial \Sigma_{ll}} \end{bmatrix} = \dots = \underline{0}$$

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

$$\Rightarrow \Sigma_{ML} = \frac{1}{N} \sum_{k=1}^N (\underline{x}_k - \underline{\mu}_{ML})(\underline{x}_k - \underline{\mu}_{ML})^T$$

Maximum Likelihood estimates are only asymptotic unbiased, so we need a large N !

WARNING: An unbiased estimator is also no guarantee for a correct result!!

Estimation of an unknown PDF

8

❖ Maximum A Posteriori Probability estimation (MAP)

➤ (In ML $\underline{\theta}$ 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 $\mathbf{X} = \{ \underline{x}_1, \underline{x}_2, \dots, \underline{x}_N \}$

Compute the maximum of $p(\underline{\theta} | \mathbf{X})$

➤ From Bayes theorem

$$p(\underline{\theta}) p(\mathbf{X} | \underline{\theta}) = p(\mathbf{X}) p(\underline{\theta} | \mathbf{X}) \quad \text{or} \quad p(\underline{\theta} | \mathbf{X}) = \frac{p(\underline{\theta}) p(\mathbf{X} | \underline{\theta})}{p(\mathbf{X})}$$

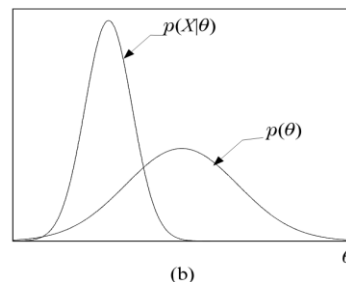
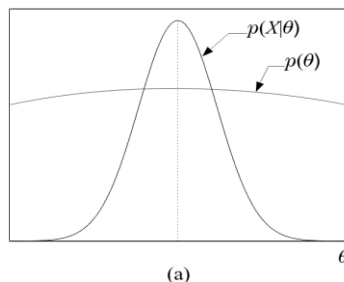
9

➤ The method:

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

$$\hat{\underline{\theta}}_{MAP} : \frac{\partial}{\partial \underline{\theta}} (P(\underline{\theta}) p(\mathbf{X} | \underline{\theta})) = 0$$

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



❖ MAP Example:

$p(\underline{x}) : N(\underline{\mu}, \Sigma = \sigma^2 I), \underline{\mu}$ unknown

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

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

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

$$\Rightarrow \hat{\underline{\mu}}_{MAP} = \frac{\underline{\mu}_0 + \frac{\sigma_{\mu}^2}{\sigma^2} \sum_{k=1}^N \underline{x}_k}{1 + \frac{\sigma_{\mu}^2}{\sigma^2} N}$$

For $\frac{\sigma_{\mu}^2}{\sigma^2} \gg 1$, or for $N \rightarrow \infty$

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