

**PARAMETER ESTIMATION-
PROBABILITY DISTRIBUTION ESTIMATION-
BAYESIAN INFERENCE**

ESTIMATION OF UNKNOWN PROBABILITY DENSITY FUNCTIONS

❖ Maximum Likelihood

- Let $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N$ known and independent
- Let $p(\underline{x})$ known within an unknown vector

parameter $\underline{\theta}$: $p(\underline{x}) \equiv p(\underline{x}; \underline{\theta})$

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

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

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

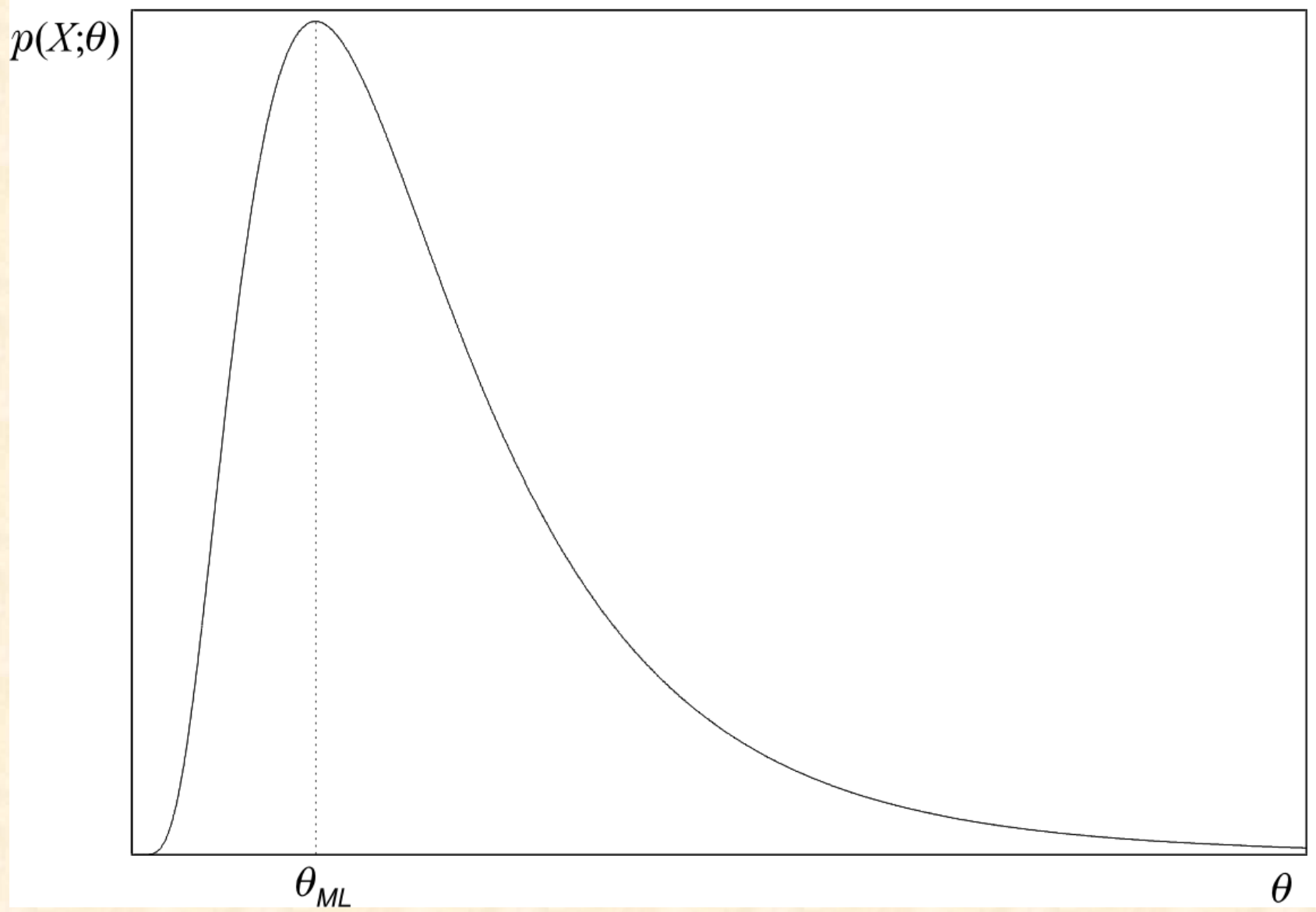
which is known as the Likelihood of $\underline{\theta}$ w.r. to X

The method :

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

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

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



❖ Example:

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

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

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

$$\frac{\partial L(\underline{\theta})}{\partial(\underline{\theta})} \equiv \begin{bmatrix} \frac{\partial L}{\partial \theta_1} \\ \cdot \\ \cdot \\ \cdot \\ \frac{\partial L}{\partial \theta_l} \end{bmatrix} = \sum_{k=1}^N \Sigma^{-1} (\underline{x}_k - \underline{\theta}) = \underline{0} \Rightarrow \underline{\theta}_{ML} = \frac{1}{N} \sum_{k=1}^N \underline{x}_k$$

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

❖ Maximum A Posteriori Probability Estimation

- 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

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

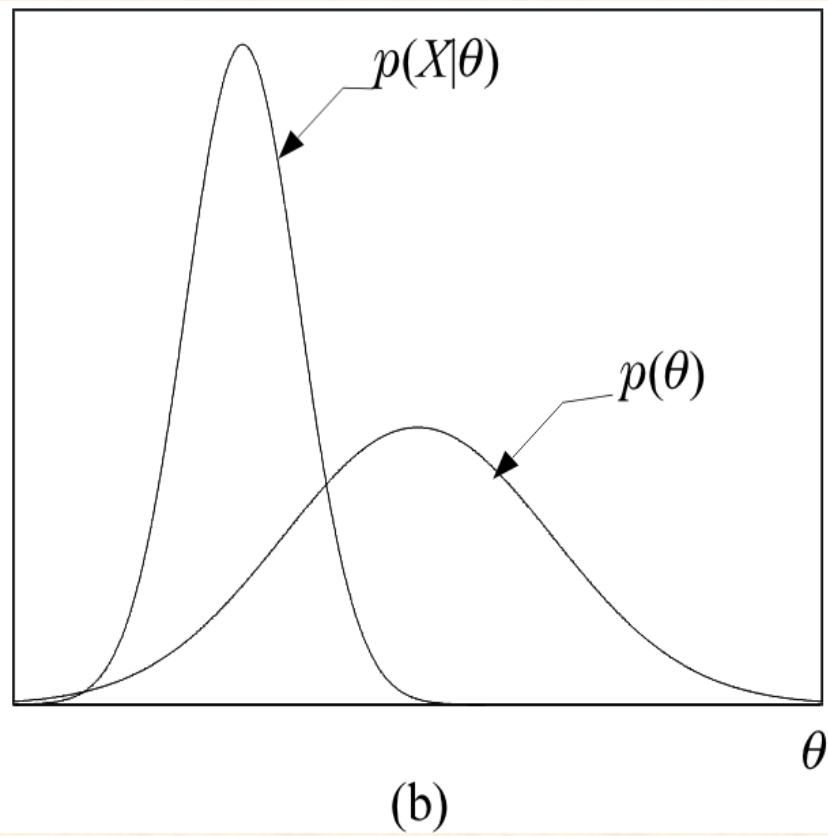
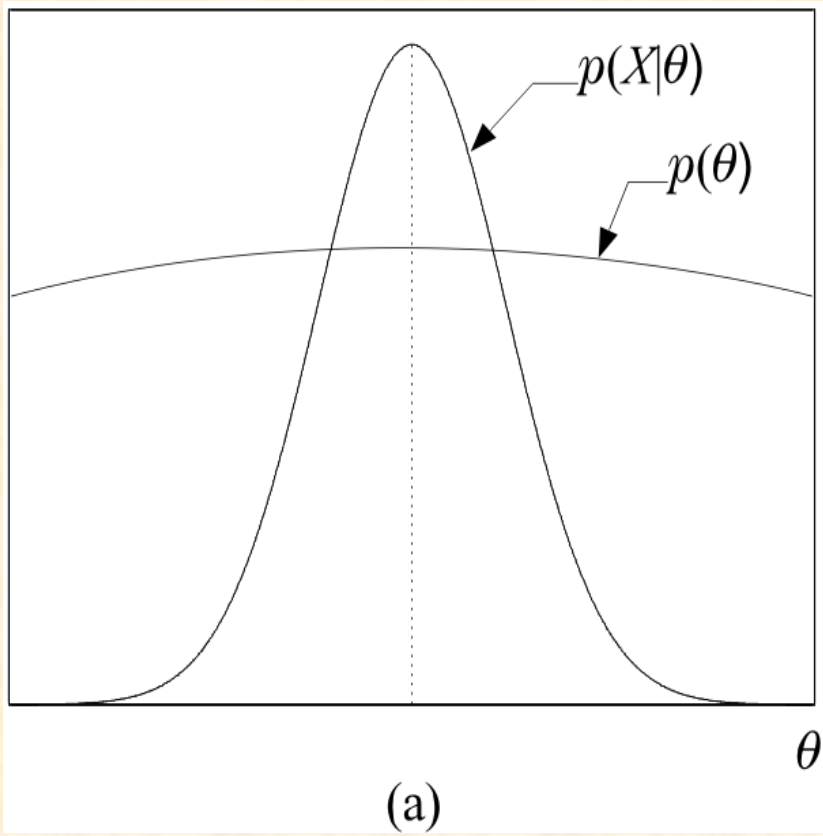
$$p(\underline{\theta} | X) = \frac{p(\underline{\theta}) p(X | \underline{\theta})}{p(X)}$$

➤ 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 $\hat{\underline{\theta}}_{MAP} \cong \underline{\theta}_{ML}$



❖ Example:

θ unknown, let $p(x|\theta) \rightarrow N(\theta, \sigma_\eta^2)$, $X = \{x_1, \dots, x_N\}$

$$p(\theta) = \frac{1}{(2\pi)^{\frac{1}{2}} \sigma_0} \exp\left(-\frac{\|\theta - \theta_0\|^2}{2\sigma_0^2}\right)$$

$$\hat{\theta}_{MAP} : \frac{\partial}{\partial \theta} \ln\left(\prod_{k=1}^N p(x_k | \theta) p(\theta)\right) = 0 \quad \text{or} \quad \sum_{k=1}^N \frac{1}{\sigma_\eta^2} (x_k - \theta) - \frac{1}{\sigma_0^2} (\theta - \theta_0) = 0 \Rightarrow$$

$$\hat{\theta}_{MAP} = \frac{N\bar{x} + \theta_0 \frac{\sigma_\eta^2}{\sigma_0^2}}{N + \frac{\sigma_\eta^2}{\sigma_0^2}} \quad \text{For} \quad \frac{\sigma_\eta^2}{\sigma_0^2} \ll 1, \quad \text{or for } N \rightarrow \infty, \quad \hat{\theta}_{MAP} \cong \hat{\theta}_{ML} = \bar{x} = \frac{1}{N} \sum_{k=1}^N x_k$$

$$\text{For } \frac{\sigma_\eta^2}{\sigma_0^2} \gg 1, \quad \hat{\theta}_{MAP} \cong \theta_0$$

❖ Bayesian Inference

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

Here a different root is followed.

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

The goal : estimate $p(\underline{x}|X)$

How??

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

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

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

A bit more insight via an example

- Let $p(x|\theta) \rightarrow N(\theta, \sigma_\eta^2)$
- $p(\theta) \rightarrow N(\theta_0, \sigma_0^2)$
- It turns out that: $p(\theta|X) \rightarrow N(\theta_N, \sigma_N^2)$

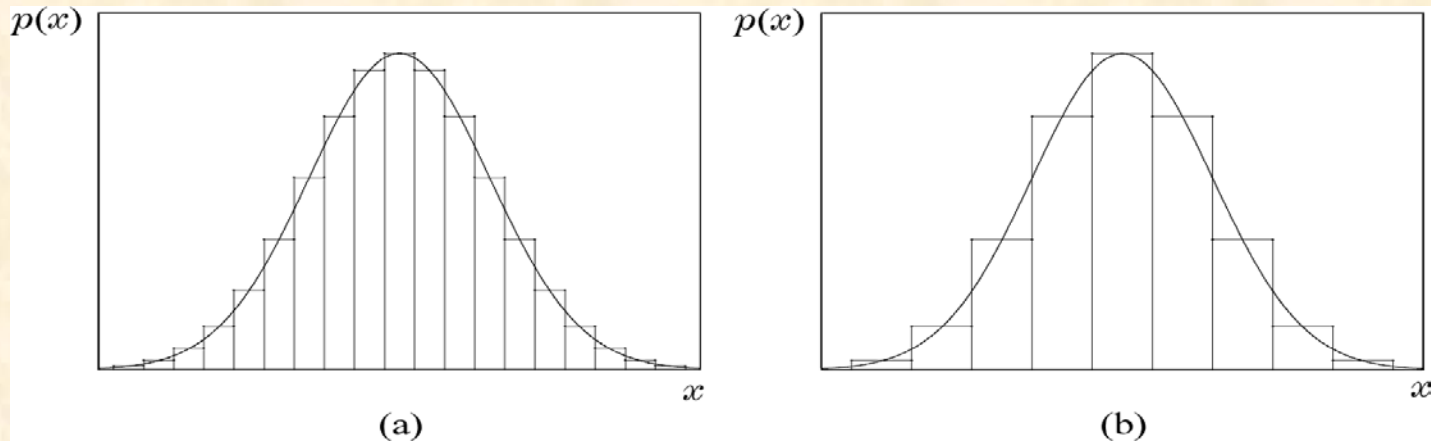
$$\theta_N = \frac{N\sigma_0^2\bar{x} + \sigma_\eta^2\theta_0}{N\sigma_0^2 + \sigma_\eta^2}, \quad \sigma_N^2 = \frac{\sigma_\eta^2\sigma_0^2}{N\sigma_0^2 + \sigma_\eta^2}, \quad \bar{x} = \frac{1}{N} \sum_{k=1}^N x_k$$

- Also: $p(x|X) \rightarrow N(\theta_N, \sigma_x^2)$, $\sigma_x^2 = \sigma_\eta^2 + \frac{\sigma_\eta^2\sigma_N^2}{\sigma_N^2 + \sigma_\eta^2}$
- Bayesian inference has given us a result similar to MAP
- Same as MAP concerning the estimation of the parameter
- But: Additional information as regards the uncertainty of the estimate

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



❖ Parzen Windows



$$P \approx \frac{k_N}{N} \begin{array}{l} \nearrow k_N \text{ in } h \\ \searrow N \text{ total} \end{array}$$

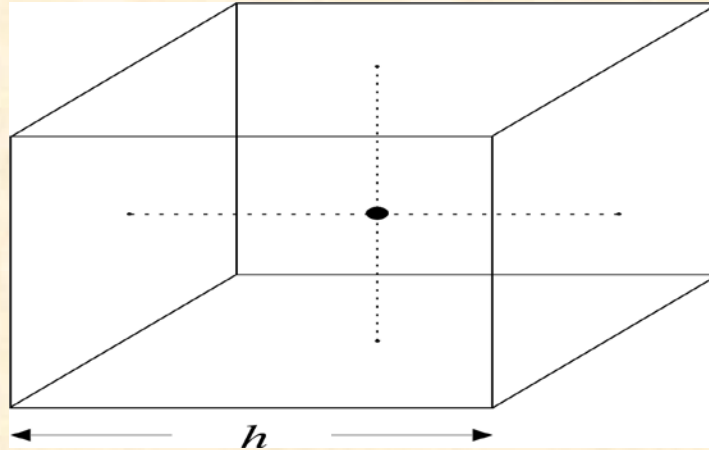
$$\hat{p}(x) \equiv \hat{p}(\hat{x}) = \frac{1}{h} \frac{k_N}{N}, \quad |x - \hat{x}| \leq \frac{h}{2}$$

$$\hat{x} - \frac{h}{2} \quad \hat{x} \quad \hat{x} + \frac{h}{2}$$

If $p(x)$ continuous, $\hat{p}(x) \rightarrow p(x)$ as $N \rightarrow \infty$, if

$$h_N \rightarrow 0, \quad k_N \rightarrow \infty, \quad \frac{k_N}{N} \rightarrow 0$$

- Divide the multidimensional space in hypercubes



➤ Define

$$\phi(\underline{x}) = \left\{ \begin{array}{l} 1, \quad |x_j| \leq \frac{1}{2} \quad \forall j \\ 0, \quad \text{otherwise} \end{array} \right\}$$

- That is, it is 1 inside a unit side hypercube centered at 0

$$\hat{p}(\underline{x}) = \frac{1}{h^l} \left(\frac{1}{N} \sum_{i=1}^N \phi\left(\frac{\underline{x}_i - \underline{x}}{h}\right) \right)$$

- $\frac{1}{\text{volume}} * \frac{1}{N} * \text{number of points inside an } h\text{-side hypercube centered at } \underline{x}$

- The problem: $p(\underline{x})$ continuous
 $\phi(\cdot)$ discontinuous

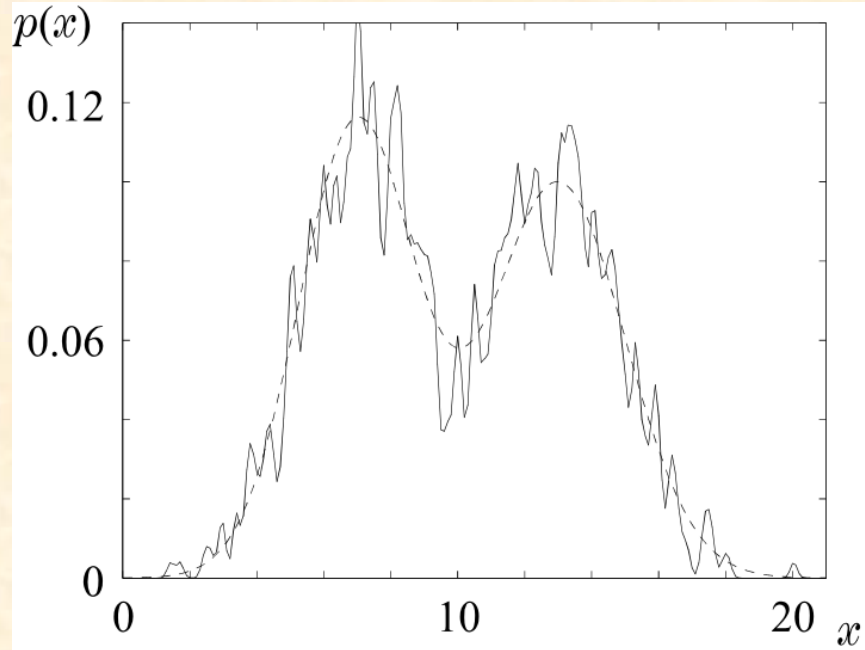
- Parzen windows-kernels-potential functions
 $\varphi(\underline{x})$ is smooth

$$\varphi(\underline{x}) \geq 0, \quad \int_{\underline{x}} \varphi(\underline{x}) d\underline{x} = 1$$

➤ Variance

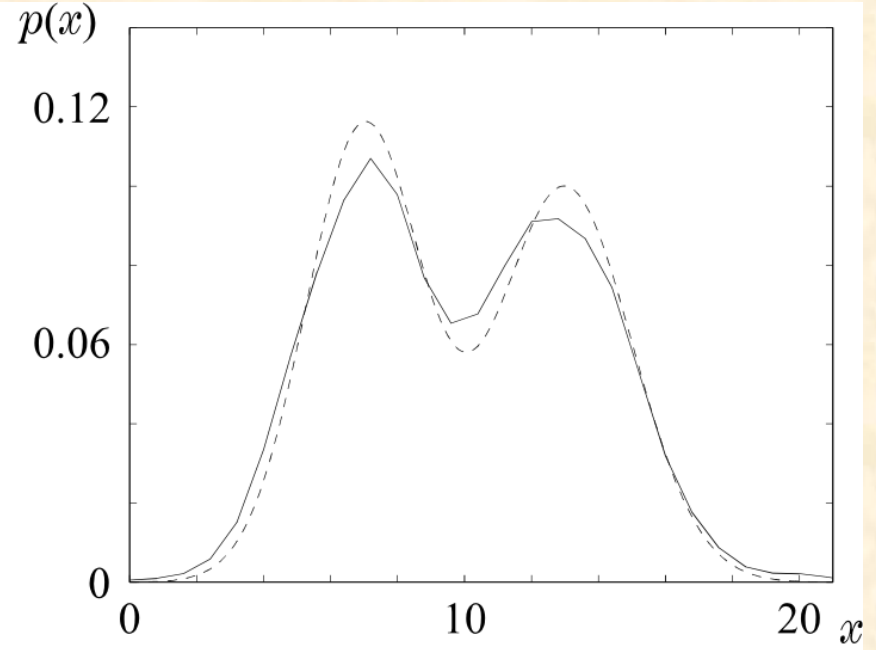
- The **smaller** the h the **higher** the variance

$h=0.1, N=1000$



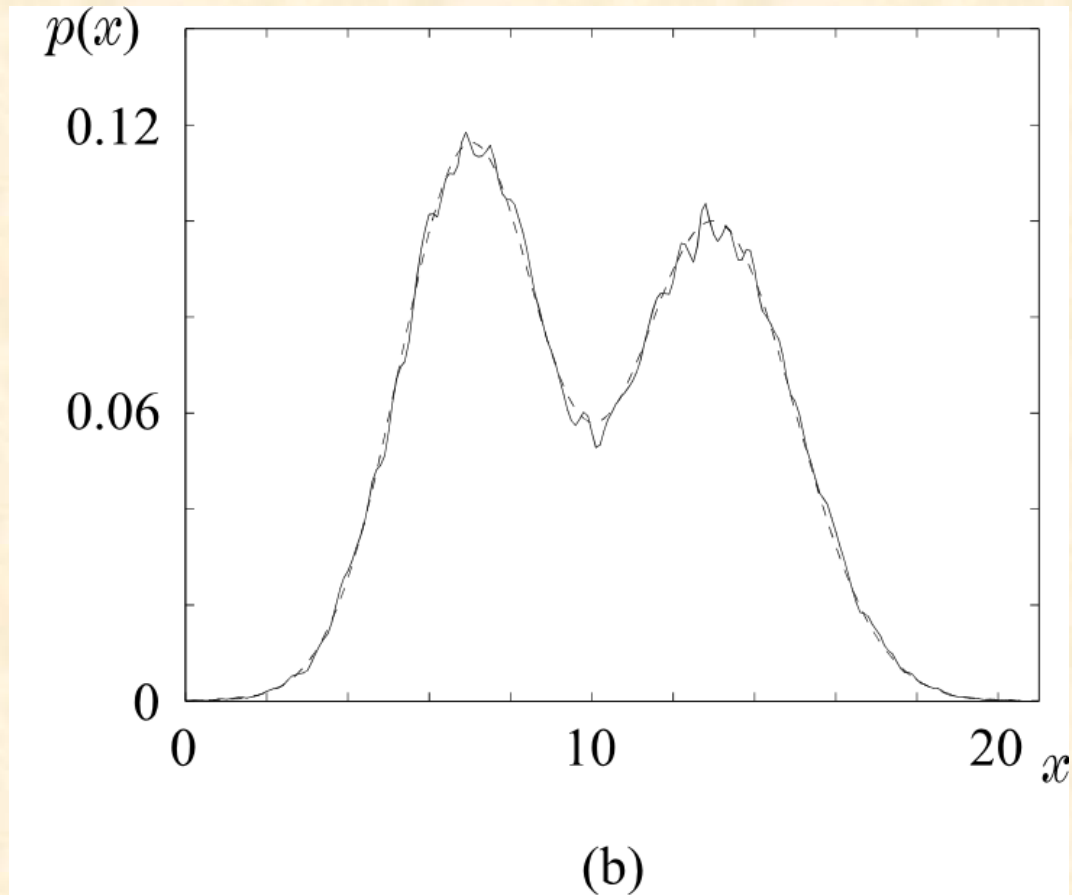
(a)

$h=0.8, N=1000$



(b)

$h=0.1, N=10000$



➤ The **higher** the N the **better** the accuracy

❖ 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 $x_1 \leq x \leq x_2$ and zero otherwise. Compute the ME pdf

- The constraint:

$$\int_{x_1}^{x_2} p(x) dx = 1$$

- Lagrange Multipliers

$$H_L = H + \lambda \left(\int_{x_1}^{x_2} p(x) dx - 1 \right)$$

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

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