# CpE 646 Pattern Recognition and Classification

## Prof. Hong Man

**Department of Electrical and
Computer Engineering
Stevens Institute of Technology**

STEVENS
Institute of Technology

Visual Information Environment Laboratory

# Bayesian Decision Theory

Chapter 3 (Section 3.1 – 3.5) Outline:

- Introduction
- Maximum-likelihood Estimation
- Bayesian Estimation
- Bayesian Parameter Estimation: Gaussian Case
- Bayesian Parameter Estimation: General Theory

Visual Information Environment Laboratory

# Introduction

- Data availability in a Bayesian framework
  - We could design an optimal classifier if we knew:
    - $P(\omega_i)$ (priors)
    - $p(x|\omega_i)$ (class-conditional densities)
  - Unfortunately, we rarely have this complete information. What we usually have are
    - some knowledge of the probabilities densities
    - some training data sample

# Introduction

- Supervised learning and unsupervised learning
    - In both cases, the observed sample $x$ reflects the prior probabilities $P(\omega_i)$ and the class- conditional densities $p(x|\omega_i)$ independently.
    - Supervised learning assumes that the training sample $x$ is individually labeled with true state of nature, whereas the unsupervised learning only has unlabeled sample.

# Introduction

- Design a classifier from a training sample
  - Estimation of prior probabilities is usually not difficult
  - Estimation of class-conditional densities is difficult
    - Samples are often too small and feature dimension maybe too high
    - If we have some knowledge about the density, e.g. certain model with set of parameters, this estimation can be much simplified.
    - For example Normal density $p(x|\omega_i) \sim N(\mu_i, \Sigma_i)$ can be characterized by 2 parameters
    - The density estimation becomes parameter estimation – parametric methods

# Introduction

- Parameter estimation techniques
  - Maximum-likelihood (ML) estimation and Bayesian estimation
  - Results are nearly identical, but the approaches are conceptually different
    - In ML methods, parameters are quantities with fixed unknown values. The best estimated value will maximize the probability of obtaining the samples observed.
    - In Bayesian methods, parameters are random variables with known prior distribution, and observations will convert this to a posterior density.

# Introduction

- In Bayesian estimation, additional observation sample will sharpen the a posteriori density function, causing it to peak near the true values of the parameters. This phenomenon is called Bayesian learning.

# Maximum-Likelihood Estimation

- 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
  - We have $c$ data sets, $D_1, D_2, \ldots, D_c$
  - $D_j$ contains samples independently drawn from $p(x|\omega_j)$, such samples are called independent and identically distributed random variables, i.i.d.

STEVENS
Institute of Technology

# Maximum-Likelihood Estimation

- We assume $p(x/\omega_j)$ has known parametric form and can be uniquely determined by parameters $\theta_j$. For example $p(x/\omega_j) \sim N(\mu_j, \Sigma_j)$

- Let $p(x/\omega_j) \equiv p(x/\omega_j, \theta_j)$. In Normal density case

$$\boldsymbol{\theta}_j = (\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) = (\boldsymbol{\mu}_j^1, \boldsymbol{\mu}_j^2, ..., \boldsymbol{\sigma}_j^{11}, \boldsymbol{\sigma}_j^{22}, \mathbf{cov}(x_j^m, x_j^n), ...)$$

- The objective is to use training samples to estimate such parameter vector for each class, $\theta_1,\ \theta_2,\ …,\ \theta_c$

- We assume parameters for different classes are independent, so each $\theta_j$ can be estimated independently.

# Maximum-Likelihood Estimation

- Suppose that $D$ contains $n$ samples, $\{x_1, x_2, \ldots, x_n\}$, because the samples are drawn independently

$$p(D \,|\, \theta) = \prod_{k=1}^{n} p(x_k \,|\, \theta)$$

- $p(D|\theta)$ is called the likelihood of $\theta$ w.r.t. the set samples
- The ML estimate of $\theta$ is, by definition the value $\hat{\theta}$ that maximizes $p(D|\theta)$
  - It is the value of $\theta$ that best agrees with the actually observed training sample
- It is usually easier to work with the logarithm of the likelihood – log-likelihood.

STEVENS
Institute of Technology
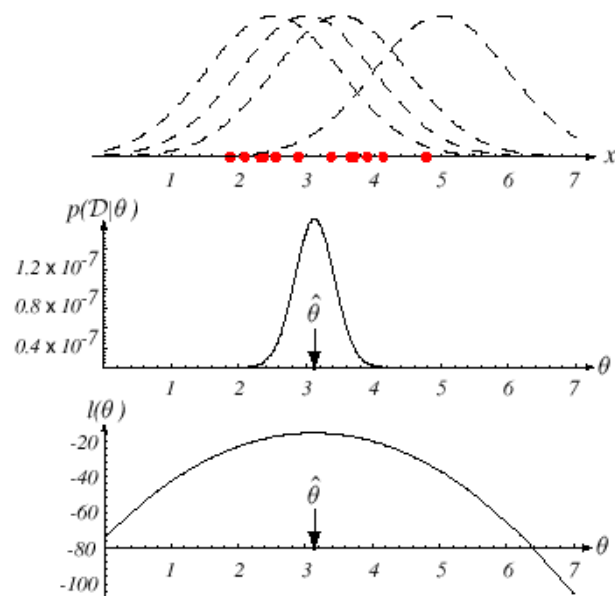
# Maximum-Likelihood Estimation



**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 $l(\theta)$, shown at the bottom. Note that even though they look similar, the likelihood $p(\mathcal{D}|\theta)$ is shown as a function of $\theta$ whereas the conditional density $p(x|\theta)$ is shown as a function of $x$. Furthermore, as a function of $\theta$, 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.

# Maximum-Likelihood Estimation

- Optimal estimation
  - Let $\theta = (\theta_1, \theta_2, \ldots, \theta_p)^t$, i.e. $\theta$ has $p$ parameters, and let $\nabla_\theta$ be the gradient operator

  $$\nabla_\theta \equiv \left[ \frac{\partial}{\partial \boldsymbol{\theta}_1}, \frac{\partial}{\partial \boldsymbol{\theta}_2}, \ldots, \frac{\partial}{\partial \boldsymbol{\theta}_p} \right]^t$$

  - We define $l(\theta)$ as the log-likelihood function

  $$l(\theta) = \ln p(D|\theta)$$

  - New problem statement:
    determine $\theta$ that maximizes the log-likelihood

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

# Maximum-Likelihood Estimation

**Thus**

$$l(\theta) = \ln p(D \mid \theta) = \ln\left(\prod_{k=1}^{n} p(x_k \mid \theta)\right) = \sum_{k=1}^{n} \ln p(x_k \mid \theta)$$

$$\nabla_\theta l = \sum_{k=1}^{n} \nabla_\theta \ln p(x_k \mid \theta)$$

**the necessary condition for an optimum is**

$$\nabla_\theta l = 0$$

- The solution to this equation $\hat{\theta}$ may represent a true global maximum, a local maximum or minimum, or an inflection point of $l(\theta)$

# Maximum *A Posteriori*

- A related estimator is maximum *a posteriori* (MAP) estimator
  - It finds the value of $\theta$ that maximizes $l(\theta)p(\theta)$, where $p(\theta)$ is the prior probability of different parameter values.
  - The MAP estimator becomes the ML estimator for the uniform prior.
  - The drawback of MAP estimator is that if the parameter space is transformed, the densities will change and the MAP solution will no longer hold. (For example the signal is normalized so the mean values is shifted.)

# A Gaussian Case

- A Gaussian case with unknown $\mu$
  - $p(x_i|\mu) \sim N(\mu, \Sigma)$: samples are drawn from a multivariate normal population
  - Let $\theta = \mu$,

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

  and $\quad \nabla_\mu \ln p(x_k \mid \mu) = \Sigma^{-1}(x_k - \mu)$

  - The ML estimate for $\mu$ must satisfy:

  $$\sum_{k=1}^{n} \Sigma^{-1}(x_k - \hat{\mu}) = 0$$

# A Gaussian Case

– Multiplying by $\Sigma$ (because $|\Sigma| \neq 0$) and rearranging, we have:

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^{n} x_k$$

- This is just the arithmetic average of the samples of the training samples – sample mean

• Conclusion:

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

STEVENS
Institute of Technology

# Another Gaussian Case

- Another Gaussian case with unknown $\mu$ and $\Sigma$
    - In univariate case $\Sigma \Rightarrow \sigma^2$ and $\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$

$$l = \ln p(x_k \mid \boldsymbol{\theta}) = -\frac{1}{2}\ln 2\pi\boldsymbol{\theta}_2 - \frac{1}{2\boldsymbol{\theta}_2}(x_k - \boldsymbol{\theta}_1)^2$$

$$\nabla_{\theta} l = \begin{bmatrix} \dfrac{\partial}{\partial \boldsymbol{\theta}_1}(\ln P(x_k \mid \boldsymbol{\theta})) \\[2em] \dfrac{\partial}{\partial \boldsymbol{\theta}_2}(\ln P(x_k \mid \boldsymbol{\theta})) \end{bmatrix} = \begin{bmatrix} \dfrac{1}{\boldsymbol{\theta}_2}(x_k - \boldsymbol{\theta}_1) \\[2em] -\dfrac{1}{2\boldsymbol{\theta}_2} + \dfrac{(x_k - \boldsymbol{\theta}_1)^2}{2\boldsymbol{\theta}_2^2} \end{bmatrix} = 0$$

# Another Gaussian Case

– Conditions:

$$\begin{cases} \sum_{k=1}^{n} \dfrac{1}{\hat{\theta}_2}(x_k - \hat{\theta}_1) = 0 & \textbf{(1)} \\[2em] -\sum_{k=1}^{n} \dfrac{1}{\hat{\theta}_2} + \sum_{k=1}^{n} \dfrac{(x_k - \hat{\theta}_1)^2}{\hat{\theta}_2^2} = 0 & \textbf{(2)} \end{cases}$$

– Combining (1) and (2):

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

# Another Gaussian Case

– In multivariate case

$$\hat{\mu} = \frac{1}{n}\sum_{k=1}^{n} x_k$$

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

- The ML estimate for the mean is the sample mean, and the ML estimate for the covariance matrix is the arithmetic average of the $n$ matrices

$$(x_k - \hat{\mu})(x_k - \hat{\mu})^t$$

# Bias

- The ML estimate for $\sigma^2$ is biased, that is the expected value over all data set of size $n$ of the sample variance is not equal to the true variance (imaging $n=1$)

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

$$=E\left[\frac{1}{n}\sum_{i=1}^{n}(x_i-\mu)^2-2(\overline{x}-\mu)\frac{1}{n}\sum_{i=1}^{n}(x_i-\mu)+(\overline{x}-\mu)^2\right]$$

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

$$=\sigma^2-\frac{1}{n}\sigma^2=\frac{n-1}{n}\sigma^2<\sigma^2$$

# Bias

- An elementary unbiased estimator for $\Sigma$ is:

$$\underbrace{C = \frac{1}{n\text{-}1}\sum_{k=1}^{n}(x_k - \hat{\mu})(x_k - \hat{\mu})^t}_{\text{Sample covariance matrix}}$$

- If an estimator is unbiased for all distributions, it is called absolutely unbiased.

- If an estimator becomes unbiased only when the number of sample is very large, e.g. $(n\text{-}1)/n \rightarrow 1$, it is called asymptotically unbiased.

- Biased estimator $v.s.$ unbiased estimator
  - We just choose one that gives the better classification performance.

# Bayesian Estimation

- Bayesian Estimation (Bayesian learning)
  - In *ML* estimation method $\theta$ was supposed to be fixed
  - In Bayesian estimation, $\theta$ is a random variable with a certain distribution
  - The training data will allow us to convert the distribution of $\theta$ into a posterior probability density $P(\omega_i/x)$
  - Goal: compute $P(\omega_i/x,D)$

# Bayesian Estimation

- The original Bayes formula is

$$P(\omega_i \mid x) = p(x \mid \omega_i) \cdot P(\omega_i) / p(x)$$

- Given the sample $D$, Bayes formula can be written

$$P(\omega_i \mid x, D) = \frac{p(x \mid \omega_i, D)P(\omega_i \mid D)}{\sum_{j=1}^{c} p(x \mid \omega_j, D)P(\omega_j \mid D)}$$

**where**

$$p(x) \approx p(x \mid D) = \sum_{j} p(x, \omega_j \mid D) = \sum_{j=1}^{c} p(x \mid \omega_j, D)P(\omega_j \mid D)$$

# Bayesian Estimation

- We can assume that
    - $P(\omega_i) = P(\omega_i|D)$, i.e. the true prior probabilities are known from the training data,
    - The samples in $D$ are labeled, so we have $\{D_1, D_2, \ldots, D_c\}$ where $D_i$ belongs to $\omega_i$,
    - $D_i$ will have no influence to $p(x|\omega_j,D)$ for $i \neq j$
    - then

$$P(\omega_i \mid x, D) = \frac{p(x \mid \omega_i, D_i)P(\omega_i)}{\sum_{j=1}^{c} p(x \mid \omega_j, D_j)P(\omega_j)}$$

# Bayesian Estimation

- General theory:
  - Ultimately we want to estimate $p(x)$, and the best we can achieve is $p(x/D)$
  - We assume $p(x)$ takes known parametric form, i.e. $p(x|\theta)$ function is known, but parameters are unknown
  - We assume some knowledge of $p(\theta)$
  - Observation samples $D=\{x_1, x_2, \ldots, x_n\}$ will convert prior density $p(\theta)$ to posterior density $p(\theta|D)$ which we hope will peak at true value of $\theta$.

# Bayesian Estimation

- Given training sample samples $D=\{x_1, x_2, \ldots, x_n\}$, because sample are drawn independently, we have

$$p(D\,|\,\theta) = \prod_{k=1}^{n} p(x_k\,|\,\theta)$$

- With $p(\theta)$ and $p(D|\theta)$, the posterior density $p(\theta|D)$ is

$$p(\theta\,|\,D) = \frac{p(D\,|\,\theta)\,p(\theta)}{\int p(D\,|\,\theta)\,p(\theta)d\theta} = \frac{\prod_{k=1}^{n} p(x_k\,|\,\theta)\,p(\theta)}{\int p(D\,|\,\theta)\,p(\theta)d\theta}$$

**note that** $\int p(D\,|\,\theta)\,p(\theta)d\theta$ **is a normalization factor independent of** $\theta$

# Bayesian Estimation

– Finally to compute $p(x/D)$, where $D$ are training samples and $x$ are testing samples

$$p(x \mid D) = \int p(x, \theta \mid D) d\theta$$

$$= \int p(x \mid \theta, D) p(\theta \mid D) d\theta$$

$$= \int p(x \mid \theta) p(\theta \mid D) d\theta$$

note that $p(x \mid \theta, D) = p(x \mid \theta)$ because $x$ and $D$ are drawn independently.

- This is formal Bayesian solution to this problem.

# Bayesian Estimation

– Comments:

- If $p(\theta|D)$ will peak around some value of $\hat{\boldsymbol{\theta}}$, we will have $p(x|D) \simeq p(x|\hat{\boldsymbol{\theta}})$, i.e. the result is equivalent to substituting the estimate $\hat{\boldsymbol{\theta}}$ for the true parameter vector $\theta$.

- If there is no clear best estimate $\hat{\boldsymbol{\theta}}$, the $p(x/D)$ expression represents an average over all possible values of $\theta$. This expression is then usually evaluated through Monte-Carlo simulations.

# Gaussian Case

- Gaussian case: to calculate the a posteriori density $p(\theta|D)$ and the desired density $p(x/D)$ for the case where $p(x/\mu) \sim N(\mu, \Sigma)$

- The univariate case where $\mu$ is the only unknown parameter: $p(x/\mu) \sim N(\mu, \sigma^2)$

  - Compute $p(\mu|D)$

    - We assume prior density for $\mu$ as $p(\mu) \sim N(\mu_0, \sigma_0^2)$ where $\mu_0$, $\sigma_0^2$ are known (they represent our knowledge on $\mu$)

# Gaussian Case

- Let $D = \{x_1, x_2, \ldots, x_n\}$, we have

$$p(\mu \mid D) = \frac{p(D \mid \mu) p(\mu)}{\int p(D \mid \mu) p(\mu) d\mu}$$

$$= \alpha \prod_{k=1}^{n} p(x_k \mid \mu) p(\mu) \qquad (1)$$

where $\alpha$ represents the normalization factor that is dependent of $D$ but is independent of $\mu$

# Gaussian Case

- Also we have

$$p(x_k \mid \mu) = \frac{1}{\sqrt{2\pi}\sigma}\exp\left[-\frac{1}{2}\left(\frac{x_k - \mu}{\sigma}\right)^2\right]$$

$$p(\mu) = \frac{1}{\sqrt{2\pi}\sigma_0}\exp\left[-\frac{1}{2}\left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right]$$

# Gaussian Case

- Substituting $p(x/\mu)$ and $p(\mu)$ into (1), we find the *a posteriori* density $p(\mu|D)$ remains a normal density

$$p(\mu \mid D) = \frac{1}{\sqrt{2\pi}\sigma_n} \exp\left[-\frac{1}{2}\left(\frac{\mu - \mu_n}{\sigma_n}\right)^2\right]$$

**where**

$$\begin{cases} \mu_n = \left(\frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2}\right)\hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2}\mu_0 \\ \\ \sigma_n^2 = \frac{\sigma_0^2\sigma^2}{n\sigma_0^2 + \sigma^2} \end{cases}$$

**where** $\hat{\mu}_n = \frac{1}{n}\sum_{k=1}^{n} x_k$

# Gaussian Case

- Because $p(\mu|D)$ remains a normal density as number $n$ of training samples increases, $p(\mu|D)$ is said to be a reproducing density and the prior $p(\mu)$ is called a conjugate prior.

- $\mu_n$ roughly represents the best guess of $\mu$ after observing $D$ with $n$ samples, and $\sigma_n^2$ represents our uncertainty about this guess.

- As $n \to \infty$, $\sigma_n^2 \to \sigma^2/n$, i.e. $p(\mu|D)$ will peak around $\mu_n$ as a Dirac delta function. This behavior is called Bayesian learning.
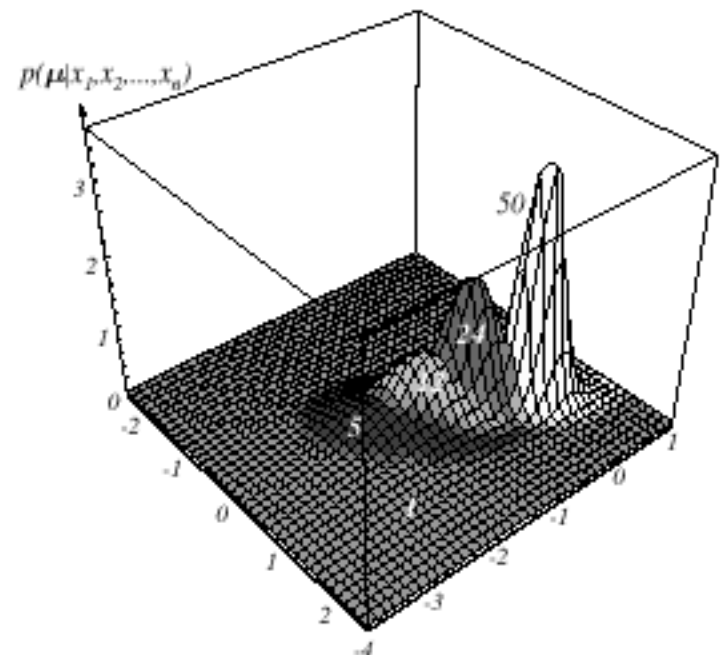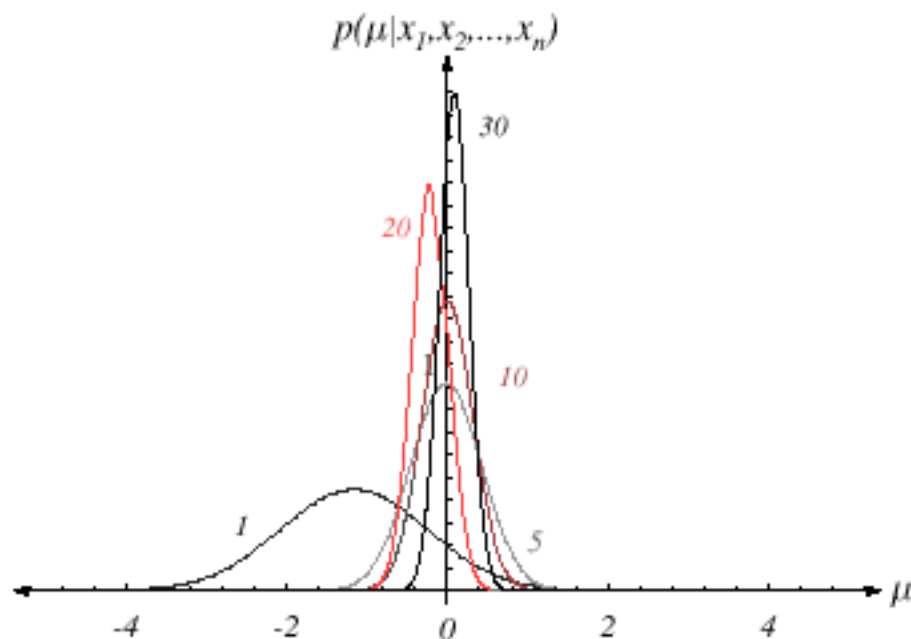
# Gaussian Case



**FIGURE 3.2.** Bayesian learning of the mean of normal distributions in one and two dimensions. The posterior distribution estimates are labeled by the number of training samples used in the estimation. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Visual Information Environment Laboratory

# Gaussian Case

– Compute $p(x|D)$ we have

$$p(x\,|\,D) = \int p(x\,|\,\mu)p(\mu\,|\,D)d\mu$$

$$= \int \frac{1}{\sqrt{2\pi}\sigma}\exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right]\frac{1}{\sqrt{2\pi}\sigma_n}\exp\left[-\frac{1}{2}\left(\frac{\mu-\mu_n}{\sigma_n}\right)^2\right]d\mu$$

$$\sim N(\mu_n,\sigma^2+\sigma_n^2)$$

- Comment: to compute $p(x|D)$ whose parametric form is known to be $p(x/\mu)\sim N(\mu,\sigma^2)$, we can just replace $\mu$ with $\mu_n$ and $\sigma^2$ as $(\sigma^2+\sigma_n^2)$

# Recursive Bayes Learning

- Let $D^n = \{x_1, x_2, \ldots, x_n\}$ where $n$ is number of training samples,

  **becaues** $p(D \mid \theta) = \displaystyle\prod_{k=1}^{n} p(x_k \mid \theta)$

  **when $n > 1$ for every new sample $x_n$ we have**

  $$p(D^n \mid \theta) = p(x_n \mid \theta) p(D^{n-1} \mid \theta)$$

  **and** $p(D^0 \mid \theta) = p(\theta)$,

  **then we have the recursive relation**

  $$p(\theta \mid D^n) = \frac{p(x_n \mid \theta) p(\theta \mid D^{n-1})}{\int p(x_n \mid \theta) p(\theta \mid D^{n-1}) d\theta}$$

# Recursive Bayes Learning

- – Use this recursive relation, we can calculate the sequence of densities $p(\theta)$, $p(\theta|x_1)$, $p(\theta|x_1, x_2)$ ,…, $p(\theta|D^n)$. This is called recursive Bayes approach.
- – When we reach a Dirac delta function around the true parameter, we have Bayesian learning

# Example 1 Recursive Bayes Learning

- Problem statement:
  - We assume our 1-D samples coming from a uniform distribution with unknown parameter $\theta$

$$p(x \mid \theta) \sim U(0, \theta) = \begin{cases} 1/\theta & 0 \le x \le \theta \\ \\ 0 & \textbf{otherwise} \end{cases}$$

  - We also assume the prior knowledge about $\theta$

$$p(\theta) \sim U(0, 10)$$

  - Now we observed four samples $D=\{4, 7, 2, 8\}$
  - We like to estimate $\theta$

# Example 1 Recursive Bayes Learning

- Recursive Bayes learning:
    - Before the first sample arrives, we have
      $$p(\theta|D^0) = p(\theta) = U(0,10)$$
    - When the first sample $x_1 = 4$ arrives, we have
      $$p(\theta|D^1) \propto p(x_1|\theta)p(\theta|D^0) = \begin{cases} 1/\theta & 4 \leq \theta \leq 10 \\ 0 & \textbf{otherwise} \end{cases}$$
    - When the second sample $x_2 = 7$ arrives, we have
      $$p(\theta|D^2) \propto p(x_2|\theta)p(\theta|D^1) = \begin{cases} 1/\theta^2 & 7 \leq \theta \leq 10 \\ 0 & \textbf{otherwise} \end{cases}$$

# Example 1 Recursive Bayes Learning

– When the third sample $x_3=2$ arrives, we have

$$p(\theta \mid D^3) \propto p(x_3 \mid \theta) p(\theta \mid D^2) = \begin{cases} 1/\theta^3 & 7 \le \theta \le 10 \\ 0 & \textbf{otherwise} \end{cases}$$
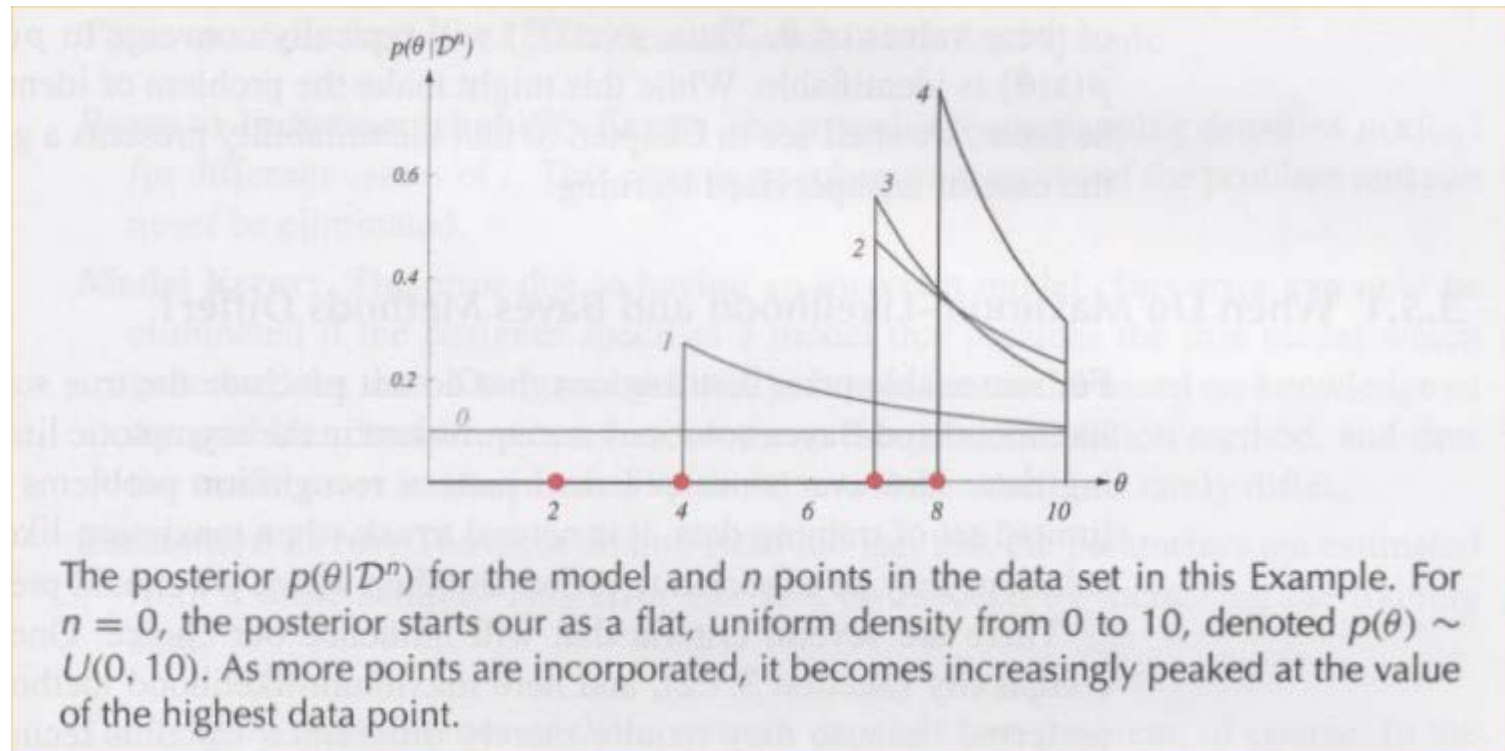
– When the forth sample $x_4=8$ arrives, we have

$$p(\theta \mid D^4) \propto p(x_4 \mid \theta) p(\theta \mid D^3) = \begin{cases} 1/\theta^4 & 8 \le \theta \le 10 \\ 0 & \textbf{otherwise} \end{cases}$$
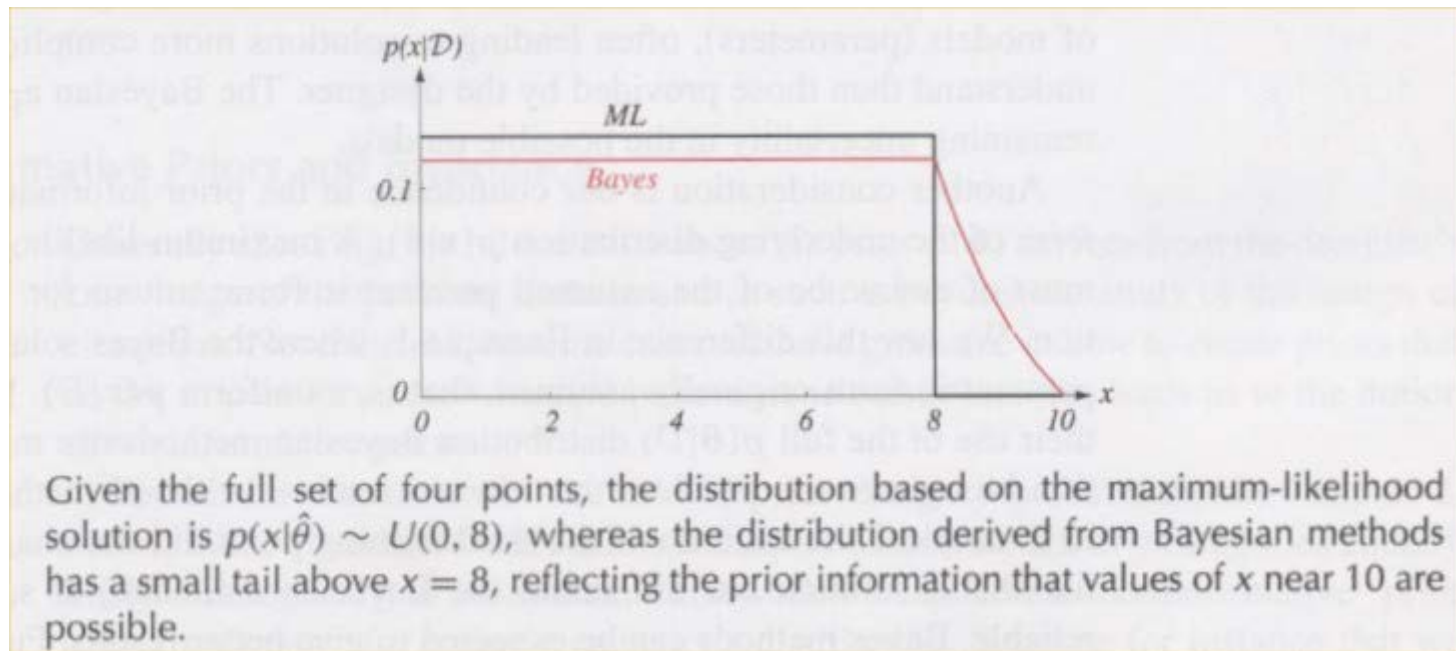
– The general solution is

$$p(\theta \mid D^n) \propto 1/\theta^n \ \ \textbf{for} \ \max_x[D^n] \le \theta \le 10$$

# Example 1 Recursive Bayes Learning



The posterior $p(\theta|\mathcal{D}^n)$ for the model and $n$ points in the data set in this Example. For $n = 0$, the posterior starts our as a flat, uniform density from 0 to 10, denoted $p(\theta) \sim U(0, 10)$. As more points are incorporated, it becomes increasingly peaked at the value of the highest data point.

# Example 1 Recursive Bayes Learning



Given the full set of four points, the distribution based on the maximum-likelihood solution is $p(x|\hat{\theta}) \sim U(0, 8)$, whereas the distribution derived from Bayesian methods has a small tail above $x = 8$, reflecting the prior information that values of $x$ near 10 are possible.

# ML Method v.s. Bayes Method

- If the prior distribution well includes the true solution, ML solution and Bayes solution are equivalent in the asymptotic limit of infinite training data.

- ML method has less computational complexity and better interpretability

- ML solution is always consistent with the model structure assumption, Bayes solution may be a little different (example 1).

- Bayes method uses more information from data than ML method.

# Sources of Classification Errors

- Bayes error: due to overlapping densities $p(x|\omega_i)$ for different $i$. This error can never be eliminated.

- Model error: due to incorrect model assumption. This error can be reduced with better prior knowledge.

- Estimation error: due to the limited number of training sample in parameter estimation. This error can be reduced by increasing sample size.