

ICS663: Pattern Recognition

Department of Information and Computer Sciences
University of Hawai'i at Manoa

Kyungim Baek

ICS663 (Fall 2015)

1

Announcement

- Homework assignment #1 **due today by 5:00 PM**
- Project proposal
 - Due: **Wednesday September 23, by 5:00 PM**
 - Presentation (~ 10 minutes)
 - Monday (9/28):
 - BJ, Thomas, Tyson, Sharif, Danny, Jeremy
 - Wednesday (9/30):
 - Tetsuya, Kelly, Nurit
 - See the evaluation criteria for presentations posted in the course website
- **Exam I: Wednesday, October 7**

ICS663 (Fall 2015)

2

Proposal Presentation

- **Title:** A short descriptive title for the project
- **Project overview:**
 - Description of the problem and main idea of the project
 - Motivation and relevant previous work
 - What will come out of this project? (project goals)
- **Technical and experimental plan:**
 - Data: Going to use existing dataset? Will collect a new dataset? Where is the source of the dataset or how will you collect a new dataset?
 - Algorithms to be explored
 - How are you going to incorporate existing techniques? If you are proposing a new approach, what is the basic idea?
 - How is the project going to be evaluated?
- **A schedule of work:** Rough milestones
- **Bibliography:** Key references

ICS663 (Fall 2015)

3

Lecture 7

- Parametric Density Estimation
 - Bayesian Parameter Estimation
 - Formal discussion in general terms
 - Univariate Gaussian Case
 - Multivariate Gaussian Case
 - Recursive Bayes

ICS663 (Fall 2015)

4

Bayesian Estimation (I)

- Bayesian learning to pattern classification problems
 - In MLE θ was supposed to be fixed
 - In Bayesian estimation θ is a random variable
 - The computation of posterior probabilities $P(\omega_i | \mathbf{x})$ lies at the heart of Bayesian classification
 - Goal:** compute $P(\omega_i | \mathbf{x}, D)$
 - Given the sample D , Bayes formula can be written

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

ICS663 (Fall 2015)

5

Bayesian Estimation (II)

- This is a supervised problem so far:
 - $P(\omega_i) = P(\omega_i | D)$: assume that the prior probability is known or obtainable from the training samples
 - $D = \{D_1, D_2, \dots, D_c\}$

$$p(\mathbf{x} | \omega_i, D) = p(\mathbf{x} | \omega_i, \{D_j\}_{j=1, \dots, c}) = p(\mathbf{x} | \omega_i, D_i, \{D_j\}_{j \neq i})$$

$$= p(\mathbf{x} | \omega_i, D_i)$$

- Then,

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

ICS663 (Fall 2015)

6

Bayesian Estimation (III)

- Notationally:

$$p(\mathbf{x} | \omega_i, D_i) \rightarrow p(\mathbf{x} | D)$$

- Now our problem is to estimate density for \mathbf{x} given the data D
- We assume the form of $p(\mathbf{x})$ – the source density for D :

$$p(\mathbf{x}) \rightarrow p(\mathbf{x} | \theta)$$

and treat θ as a *random variable*

ICS663 (Fall 2015)

7

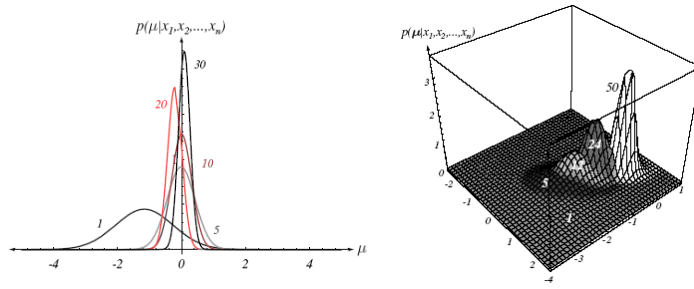
Bayesian Estimation (IV)

- Bayesian estimation attempts to find the most probable values of the parameters, conditioned on the available data. Thus, parameters are regarded as random variables and our uncertainty in the values of the parameters is represented by a probability density function.
 - Before observation: a prior probability $p(\theta)$ describes the parameters
 - As we observe data:
 - prior probability, $p(\theta) \rightarrow$ posterior probability, $p(\theta | D)$
- Bayesian Learning:** since some values of the parameters are more consistent with the data than others, we find that the posterior distribution is narrower than the prior distribution

ICS663 (Fall 2015)

8

Bayesian Learning



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.

ICS663 (Fall 2015)

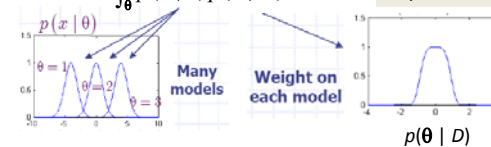
9

Bayesian Estimation (V)

- Assume we have a prior $p(\theta)$ over models
- We have data (can assume i.i.d.) D
- Want to get a model to compute $p(\mathbf{x} | D)$... How to proceed?

$$\begin{aligned} p(\mathbf{x} | D) &= \int_{\theta} p(\mathbf{x}, \theta | D) d\theta \\ &= \int_{\theta} p(\mathbf{x} | \theta, D) p(\theta | D) d\theta \\ &= \int_{\theta} p(\mathbf{x} | \theta) p(\theta | D) d\theta \end{aligned}$$

average densities $p(\mathbf{x} | \theta)$ for ALL possible values of θ weighted by its posterior probability



ICS663 (Fall 2015)

10

Bayesian Estimation (VI)

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

- If $p(\theta | D)$ is sharply peaked at $\theta_p \rightarrow p(\mathbf{x} | D) \approx p(\mathbf{x} | \theta_p)$
- Posterior distribution for θ

$$p(\theta | D) = \frac{p(D | \theta) p(\theta)}{p(D)} = \frac{p(\theta)}{p(D)} \prod_{k=1}^n p(\mathbf{x}_k | \theta)$$

- The equation shows how the observation of a set of training samples is incorporated to convert prior probability to posterior probability
- Bayesian method does not commit to a particular value of θ , but uses the entire distribution

ICS663 (Fall 2015)

11

Bayesian Estimation Summary

$$P(\omega_i | \mathbf{x}) = P(\omega_i | \mathbf{x}, D) = \frac{p(\mathbf{x} | \omega_i, D) P(\omega_i | D)}{p(\mathbf{x} | D)}$$

$$\int_{\theta} p(\mathbf{x} | \theta) p(\theta | D) d\theta$$

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

$$\prod_{k=1}^n p(\mathbf{x}_k | \theta)$$

ICS663 (Fall 2015)

12

BE: Univariate Gaussian (I)

- Univariate Gaussian case: $p(\mu|D)$

$$p(x|\mu) \sim N(\mu, \sigma^2) \quad \text{likelihood}$$

$$p(\mu) \sim N(\mu_0, \sigma_0^2) \quad \text{prior density of the parameter}$$

- Assume that
 - μ is the only unknown parameter
 - σ^2 , μ_0 and σ_0^2 are known

- Need to find $p(\mu|D)$

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

ICS663 (Fall 2015)

13

BE: Univariate Gaussian (II)

$$\begin{aligned} p(\mu|D) &= \alpha p(\mu) \prod_{k=1}^n p(x_k|\mu) \\ &= \alpha \frac{1}{\sqrt{2\pi}\sigma_0} \exp\left[-\frac{1}{2}\left(\frac{\mu-\mu_0}{\sigma_0}\right)^2\right] \prod_{k=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x_k-\mu}{\sigma}\right)^2\right] \\ &= \alpha' \exp\left[-\frac{1}{2}\left(\left(\frac{\mu-\mu_0}{\sigma_0}\right)^2 + \sum_{k=1}^n \left(\frac{x_k-\mu}{\sigma}\right)^2\right)\right] \\ &= \alpha' \exp\left[-\frac{1}{2}\left[\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}\right)\mu^2 - 2\left(\frac{1}{\sigma^2}\sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2}\right)\mu\right]\right] \end{aligned} \quad (1)$$

Let $p(\mu|D) \sim N(\mu_n, \sigma_n^2)$ i.e. This is a Gaussian

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

ICS663 (Fall 2015)

14

BE: Univariate Gaussian (III)

- Identifying coefficients in (1) and (2) yields:

$$\begin{aligned} \frac{1}{\sigma_n^2} &= \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \quad \text{and} \quad \frac{\mu_n}{\sigma_n^2} = \frac{1}{\sigma^2} \sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2} \\ &= \frac{n\sigma_0^2 + \sigma^2}{\sigma^2\sigma_0^2} \quad \quad \quad = \frac{n}{\sigma^2} \hat{\mu}_n + \frac{\mu_0}{\sigma_0^2} \end{aligned}$$

where $\hat{\mu}_n$ is the sample mean

- Solve the equations for μ_n and σ_n^2 . Then

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

$$\sigma_n^2 = \frac{\sigma_0^2\sigma^2}{n\sigma_0^2 + \sigma^2}$$

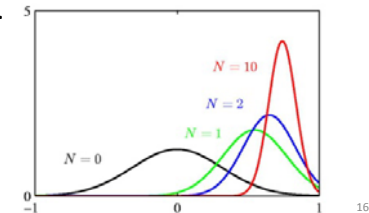
ICS663 (Fall 2015)

15

BE: Univariate Gaussian (IV)

$$\mu_n = \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0 \quad \text{and} \quad \sigma_n^2 = \frac{\sigma_0^2\sigma^2}{n\sigma_0^2 + \sigma^2}$$

- If $n \rightarrow \infty$, then $\sigma_n^2 \rightarrow 0$. That is, additional observation decreases the uncertainty about μ_n , our best guess for μ , and eventually $p(\mu|D)$ becomes sharply peaked at μ_n (Bayesian learning).



ICS663 (Fall 2015)

16

BE: Univariate Gaussian (V)

$$\mu_n = \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0 \quad \text{and} \quad \sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}$$

- μ_n : linear combination of $\hat{\mu}_n$ and μ_0
 - If $\sigma_0^2 \neq 0$ then $\mu_n \rightarrow \hat{\mu}_n$ as $n \rightarrow \infty$
 - If $\sigma_0^2 = 0$ then $\mu_n = \mu_0$: very strong prior for $\mu = \mu_0$, observation cannot change it
 - If $\sigma_0^2 \gg \sigma^2$ then $\mu_n \rightarrow \hat{\mu}_n$: uncertainty about the prior is large so does not have much influence (our guess for the mean becomes the sample mean)
- If $\sigma^2/\sigma_0^2 \neq \infty$, then μ_n will converge to the sample mean after enough samples are taken

ICS663 (Fall 2015)

17

BE: Univariate Gaussian (VI)

- Class-conditional density for $p(x|D)$

- $p(\mu|D)$ computed
- $p(x|D)$ remains to be computed!

$$p(x|D) = \int p(x|\mu) p(\mu|D) d\mu = \int N(\mu, \sigma^2) N(\mu_n, \sigma_n^2) d\mu$$

$$= \frac{1}{2\pi\sigma\sigma_n} \exp\left[-\frac{1}{2} \frac{(x - \mu_n)^2}{\sigma^2 + \sigma_n^2}\right] f(\sigma, \sigma_n)$$

i.e. $p(x|D)$ is normally distributed and $p(x|D) \sim N(\mu_n, \sigma^2 + \sigma_n^2)$

(Desired class-conditional density $p(x|D_i, \omega_i)$)

- $p(x|D_i, \omega_i)$ together with $P(\omega_i)$, and using Bayes formula, we obtain the posterior probabilities $P(\omega_i|x, D)$

ICS663 (Fall 2015)

18

BE: Multivariate Gaussian

- Direct generalization of the univariate case
- Read Section 3.4.3 of the textbook

ICS663 (Fall 2015)

19

BE: General Theory (I)

- $p(\mathbf{x}|D)$ computation can be applied to any situation in which the unknown density can be parameterized. The basic assumptions are:
 - Known form of $p(\mathbf{x}|\boldsymbol{\theta})$, but unknown value of $\boldsymbol{\theta}$
 - Known prior density $p(\boldsymbol{\theta})$:
 - contains our knowledge about $\boldsymbol{\theta}$
 - The rest of our knowledge $\boldsymbol{\theta}$ is contained in a set D of n random variables $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ that follows unknown $p(\mathbf{x})$

ICS663 (Fall 2015)

20

BE: General Theory (II)

- Basic problem:
"Compute the posterior density $p(\theta|D)$ " and then
"Derive $p(x|D)$ "

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

- Using Bayes formula, we have:

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

and by independence assumption: $p(D|\theta) = \prod_{k=1}^n p(x_k|\theta)$

ICS663 (Fall 2015)

21

Maximum-a-Posteriori vs. Maximum-likelihood

- Integration is difficult and rarely solvable, so approximate it...

$$p(x|D) = \int_{\theta} p(x|\theta)p(\theta|D)d\theta = \int_{\theta} p(x|\theta) \frac{p(D|\theta)p(\theta)}{p(D)}d\theta$$

$$\approx \int_{\theta} p(x|\theta)\delta(\theta - \hat{\theta})d\theta = p(x|\hat{\theta})$$

where

$$\hat{\theta} = \begin{cases} \arg \max_{\theta} p(D|\theta)p(\theta) & \text{MAP} \\ \arg \max_{\theta} p(D|\theta) & \text{ML (uniform } p(\theta)) \end{cases}$$

ICS663 (Fall 2015)

22

MAP vs. ML (cont'd)

- Maximum a posteriori (MAP)** estimates: maximize the log of the posterior estimate with respect to parameters θ

$$\theta_{MAP} = \arg \max_{\theta} [\ln p(D|\theta) + \ln p(\theta)]$$

- MAP is like ML, but also has penalty function on different models... $p(\theta)$ favors some over others. (If the prior distribution is flat (i.e. independent of θ), then the MAP estimate matches the ML estimate)

ICS663 (Fall 2015)

23

Recursive Bayes

- Given a training set of n samples $D^n = \{x_1, x_2, \dots, x_n\}$

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

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

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

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

Posterior is used as a prior in the next step

- Recursive Bayes** approach to parameter estimation

- Start with $p(\theta|D^0) = p(\theta)$ See Example 1 in page 98.
- Apply eq.(1): produces the sequence of densities $p(\theta)$, $p(\theta|x_1)$, $p(\theta|x_1, x_2)$, ... (incremental or on-line learning)

ICS663 (Fall 2015)

24

ML vs. Bayes Parameter Estimation

	ML	Bayes
parameters	Fixed value	Random variables
assumptions	It avoids assumptions on prior information and it is analytically easier to solve although some estimates can be biased	It permits including a priori information about the unknown, but the analytical derivation are cumbersome
computational complexity	Low (differential calculus)	possibly complex multidimensional integration
accuracy		More accurate statistically because it considers the uncertainty in estimating the parameters.
size of training set	Estimates do not change significantly as the training set size increases.	Estimates continuously improves as the training set size increases.
	For ordinary cases, both approaches give similar results with sufficient sample data	
interpretability	Easier (single best solution)	Complicated and harder (weighted average of models (parameters))
Confidence in the prior information	Solution must be of the assumed parametric form (e.g. Example 1)	Solution may not be of the parametric form originally assumed (e.g. Example 1).

ICS663 (Fall 2015)

partly from S. Iliescu