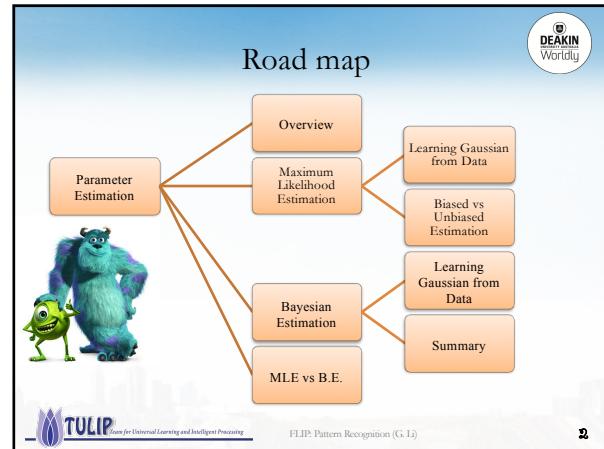


Lecture Notes on
Pattern Recognition

Session 03(A): Parameter Estimation (I)

Gang Li
School of Information Technology
Deakin University, VIC 3125, Australia

DEAKIN
Worldly



Parameter Estimation

- Maximum Likelihood
- Bayesian Estimation

DEAKIN
Worldly

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li) 3

Parameter Estimation

- **Bayesian Decision Theory** allows us to design an optimal classifier if we know the prior probabilities $P(\omega_i)$ and the class-conditional densities $p(x|\omega_i)$.

$$P(\omega_j | x) = \frac{p(x | \omega_j)P(\omega_j)}{p(x)}$$

- Estimate $P(\omega_i)$ and $p(x|\omega_i)$ from training examples.
 - Estimating $P(\omega_i)$ is usually not difficult.
 - Estimating $p(x|\omega_i)$ is more difficult!
 - Number of samples is often too small
 - Dimensionality of feature space is large.

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li) 4

Parameter Estimation

- Assumptions
 - We are given a set of training samples $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, where the samples were drawn according to $p(\mathbf{x} | \omega_j)$
 - $p(\mathbf{x} | \omega_j)$ has known parametric form, e.g.,

$$p(\mathbf{x} | \omega_j) \sim N(\boldsymbol{\mu}_j, \Sigma_j), \quad \boldsymbol{\theta} = (\boldsymbol{\mu}_j, \Sigma_j)$$

- Parameter estimation problem:
Given D , find the best possible $\boldsymbol{\theta}$

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li) 5

Main Methods in Parameter Estimation

- **Maximum Likelihood (ML)**
 - Assumes that the values of the parameters are **fixed** but unknown.
 - Best estimate $\hat{\boldsymbol{\theta}}$ is obtained by **maximizing** the probability of obtaining the samples actually observed (i.e., training data)

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

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li) 6

Main Methods in Parameter Estimation

• Bayesian Estimation (BE)

- Assumes that the parameters θ are **random** variables that have some known a-priori distribution $p(\theta)$.
- Estimates a **distribution** rather than making point estimates like ML.
- BE solution **might not** be of the parametric form assumed.

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

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 7

ML Estimation



- A General Strategy
- Learning Gaussian from Data
- Biased and Unbiased Estimation

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 8

Why we should care

- Maximum Likelihood Estimation** is a very very very very fundamental part of data analysis.
- “**MLE for Gaussians**” is training wheels for our future techniques
 - Learning Gaussians is more useful than you might guess...

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 9

ML Estimation (Assumptions)

- Training data is divided in **c** sets, based on **c** classes.
 D_1, D_2, \dots, D_c
 - Samples in D_i have been drawn independently according to $p(x/\omega_i)$.
 - $p(x/\omega_i)$ has known parametric form with parameters θ_i ;
 - e.g., $\theta_i = (\mu_i, \Sigma_i)$ for Gaussian distributions

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 10

ML Estimation (Problem Definition)

- Problem: given D_1, D_2, \dots, D_c and a model for each class, estimate $\theta_1, \theta_2, \dots, \theta_c$
 - If the samples in D_i give no information about θ_i , where $i \neq j$. We need to solve c independent problems
 - i.e., one for each class
 - ML estimate for $D = \{x_1, x_2, \dots, x_n\}$ is the $\hat{\theta}$ value that maximizes $p(D/\theta)$
 - i.e., best supports the training data

$$p(D/\theta) = p(x_1, x_2, \dots, x_n / \theta) = \prod_{k=1}^n p(x_k / \theta)$$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 11

MLE Gaussians from Data (Case of Unknown $\theta = \mu$)

- Suppose you have $x_1, x_2, \dots, x_n \sim (\text{i.i.d.}) N(\mu, \sigma^2)$
- You don't know μ , but you do know σ^2
- MLE:
 - For which μ is x_1, x_2, \dots, x_n most likely?

$$\mu^{mle} = \arg \max_{\mu} p(x_1, x_2, \dots, x_n | \mu, \sigma^2)$$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 12

MLE Gaussians from Data (Case of Unknown $\theta = \mu$)

$$\begin{aligned}\mu^{mle} &= \arg \max_{\mu} p(x_1, x_2, \dots, x_n | \mu, \sigma^2) \\ &= \quad \text{(by i.i.d)} \\ &= \quad \text{(monotonicity of log)} \\ &= \quad \text{(plug in formula for Gaussian)} \\ &= \quad \text{(after simplification)}\end{aligned}$$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 18

MLE Gaussians from Data (Case of Unknown $\theta = \mu$)

$$\begin{aligned}\mu^{mle} &= \arg \max_{\mu} p(x_1, x_2, \dots, x_n | \mu, \sigma^2) \\ &= \arg \max_{\mu} \prod_{i=1}^R p(x_i | \mu, \sigma^2) \quad \text{(by i.i.d)} \\ &= \arg \max_{\mu} \sum_{i=1}^R \log p(x_i | \mu, \sigma^2) \quad \text{(monotonicity of log)} \\ &= \arg \max_{\mu} \frac{1}{\sqrt{2\pi} \sigma} \sum_{i=1}^R \frac{(x_i - \mu)^2}{2\sigma^2} \quad \text{(plug in formula for Gaussian)} \\ &= \arg \min_{\mu} \sum_{i=1}^R (x_i - \mu)^2 \quad \text{(after simplification)}\end{aligned}$$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 19

MLE Gaussians from Data (Case of Unknown $\theta = \mu$)

$$\begin{aligned}\mu^{mle} &= \arg \max_{\mu} p(x_1, x_2, \dots, x_n | \mu, \sigma^2) \\ &= \arg \min_{\mu} \sum_{i=1}^n (x_i - \mu)^2 \\ \text{subject to } 0 &= \frac{\partial LL}{\partial \mu} = \frac{\partial}{\partial \mu} \sum_{i=1}^n (x_i - \mu)^2 = -\sum_{i=1}^n 2(x_i - \mu) \\ \text{Thus } \mu &= \frac{1}{n} \sum_{i=1}^n x_i\end{aligned}$$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 15

MLE Gaussians from Data (Case of Unknown $\theta = \mu$)

- The best estimate of the mean of a distribution is the mean of the sample!

$$\mu^{mle} = \frac{1}{n} \sum_{i=1}^n x_i$$

At first sight:
This kind of pedantic, algebra-filled and ultimately unsurprising fact is exactly the reason people throw down their "Statistics" book and pick up their "Agent-Based Evolutionary Data Mining Using The Neuro-Fuzz Transform" book.

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 16

A General Scalar MLE strategy

Task: Find MLE θ assuming known form for $p(\text{Data} | \theta, \text{stuff})$

- Write $LL = \log P(\text{Data} | \theta, \text{stuff})$
- Work out $\partial LL / \partial \theta$ using high-school calculus
- Set $\partial LL / \partial \theta = 0$ for a maximum, creating an equation in terms of θ
- Solve it*
- Check that you have found a maximum rather than a minimum or saddle-point, and be careful if θ is constrained

*This is a perfect example of something that works perfectly in all textbook examples and usually involves surprising pain if you need it for something new.

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 17

A General MLE strategy

- Suppose $\theta = (\theta_1, \theta_2, \dots, \theta_n)^T$ is a vector of parameters.
Find MLE θ assuming known form for $p(\text{Data} | \theta, \text{stuff})$
 - Write $LL = \log P(\text{Data} | \theta, \text{stuff})$
 - Work out $\partial LL / \partial \theta$ using high-school calculus

$$\frac{\partial LL}{\partial \theta} = \begin{pmatrix} \frac{\partial LL}{\partial \theta_1} \\ \frac{\partial LL}{\partial \theta_2} \\ \vdots \\ \frac{\partial LL}{\partial \theta_n} \end{pmatrix}$$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 18

A General MLE strategy

- Suppose $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_n)^T$ is a vector of parameters.
Find MLE $\hat{\boldsymbol{\theta}}$ assuming known form for $p(\text{Data} | \boldsymbol{\theta}, \text{stuff})$

 - Write $LL = \log P(\text{Data} | \boldsymbol{\theta}, \text{stuff})$
 - Work out $\partial LL / \partial \boldsymbol{\theta}$ using high-school calculus
 - Solve the set of simultaneous equations

$$\frac{\partial LL}{\partial \boldsymbol{\theta}} = \begin{pmatrix} \frac{\partial LL}{\partial \theta_1} & \frac{\partial LL}{\partial \theta_1} \\ \frac{\partial LL}{\partial \theta_2} & \frac{\partial LL}{\partial \theta_2} \\ \vdots & \vdots \\ \frac{\partial LL}{\partial \theta_n} & \frac{\partial LL}{\partial \theta_n} \end{pmatrix} = 0$$



FLIP: Pattern Recognition (G. Li)

19



A General MLE strategy

- Suppose $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_n)^T$ is a vector of parameters.
Find MLE $\hat{\boldsymbol{\theta}}$ assuming known form for $p(\text{Data} | \boldsymbol{\theta}, \text{stuff})$

 - Write $LL = \log P(\text{Data} | \boldsymbol{\theta}, \text{stuff})$
 - Work out $\partial LL / \partial \boldsymbol{\theta}$ using high-school calculus
 - Solve the set of simultaneous equations
 - Check that you are at a maximum



FLIP: Pattern Recognition (G. Li)

20



MLE Gaussians from Data (Case of Unknown $\boldsymbol{\theta} = (\theta_1, \theta_2) = (\mu, \sigma^2)$)

- Suppose you have $x_1, x_2, \dots, x_n \sim (\text{i.i.d.}) N(\mu, \sigma^2)$
 - You don't know μ or σ^2
 - MLE:
- For which $\boldsymbol{\theta} = (\mu, \sigma^2)$ is x_1, x_2, \dots, x_n most likely?

$$\log p(x_1, x_2, \dots, x_n | \mu, \sigma^2) = -n(\log \pi + \frac{1}{2} \log \sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\frac{\partial LL}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) \quad \frac{\partial LL}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2$$



FLIP: Pattern Recognition (G. Li)

21



MLE Gaussians from Data (Case of Unknown $\boldsymbol{\theta} = (\theta_1, \theta_2) = (\mu, \sigma^2)$)

- Suppose you have $x_1, x_2, \dots, x_n \sim (\text{i.i.d.}) N(\mu, \sigma^2)$
 - You don't know μ or σ^2
 - MLE:
- For which $\boldsymbol{\theta} = (\mu, \sigma^2)$ is x_1, x_2, \dots, x_n most likely?

$$\log p(x_1, x_2, \dots, x_n | \mu, \sigma^2) = -n(\log \pi + \frac{1}{2} \log \sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$0 = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) \quad 0 = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2$$



FLIP: Pattern Recognition (G. Li)

22

MLE Gaussians from Data (Case of Unknown $\boldsymbol{\theta} = (\theta_1, \theta_2) = (\mu, \sigma^2)$)

- Suppose you have $x_1, x_2, \dots, x_n \sim (\text{i.i.d.}) N(\mu, \sigma^2)$
 - You don't know μ or σ^2
 - MLE:
- For which $\boldsymbol{\theta} = (\mu, \sigma^2)$ is x_1, x_2, \dots, x_n most likely?

$$\log p(x_1, x_2, \dots, x_n | \mu, \sigma^2) = -n(\log \pi + \frac{1}{2} \log \sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\mu^{mle} = \frac{1}{n} \sum_{i=1}^n x_i \quad \sigma^2_{mle} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu^{mle})^2$$



FLIP: Pattern Recognition (G. Li)

23



Unbiased Estimators

- An estimator of a parameter is **unbiased** if the expected value of the estimate is the **same** as the true value of the parameters.
- If $x_1, x_2, \dots, x_n \sim (\text{i.i.d.}) N(\mu, \sigma^2)$ then

$$E[\mu^{mle}] = E\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \mu$$

- μ^{mle} is unbiased



FLIP: Pattern Recognition (G. Li)

24

Biased Estimators

- An estimator of a parameter is **biased** if the expected value of the estimate is the **different from** the true value of the parameters.
- If $x_1, x_2, \dots, x_R \sim (\text{i.i.d.}) N(\mu, \sigma^2)$ then

$$E[\sigma_{\text{mle}}^2] = E\left[\frac{1}{n} \sum_{i=1}^n (x_i - \mu^{\text{mle}})^2\right] = E\left[\frac{1}{n} \left(\sum_{i=1}^n x_i - \frac{1}{n} \sum_{j=1}^n x_j \right)^2\right] \neq \sigma^2$$

• σ^2_{mle} is biased

 TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 25

MLE Variance Bias

- If $x_1, x_2, \dots, x_R \sim (\text{i.i.d.}) N(\mu, \sigma^2)$ then

$$E[\sigma_{\text{mle}}^2] = E\left[\frac{1}{n} \left(\sum_{i=1}^n x_i - \frac{1}{n} \sum_{j=1}^n x_j \right)^2\right] = \left(1 - \frac{1}{n}\right) \sigma^2 \neq \sigma^2$$

Intuition check: consider the case of n=1
Why should our guts expect that σ^2_{mle} would be an underestimate of true σ^2 ?
How could you prove that?

• σ^2_{mle} is biased

 TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 26

MLE Variance Bias

- If $x_1, x_2, \dots, x_R \sim (\text{i.i.d.}) N(\mu, \sigma^2)$ then

$$E[\sigma_{\text{mle}}^2] = E\left[\frac{1}{n} \left(\sum_{i=1}^n x_i - \frac{1}{n} \sum_{j=1}^n x_j \right)^2\right] = \left(1 - \frac{1}{n}\right) \sigma^2 \neq \sigma^2$$

$$\sigma_{\text{unbiased}}^2 = \frac{\sigma^2}{\left(1 - \frac{1}{n}\right)} \quad \text{So } E[\sigma_{\text{unbiased}}^2] = \sigma^2$$

Which one is better?
• It depends on the task
• no much difference when R \Rightarrow large

 TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 27

MLE Gaussians from Data

(m-Dimensional $\theta = (\theta_1, \theta_2) = (\mu, \Sigma)$)

- Suppose you have $x_1, x_2, \dots, x_n \sim (\text{i.i.d.}) N(\mu, \Sigma)$
- You don't know μ or Σ
- MLE:
 - For which $\theta = (\mu, \Sigma)$ is x_1, x_2, \dots, x_n most likely?

$$\mu^{\text{mle}} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$$

 TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 28

MLE Gaussians from Data

(m-Dimensional $\theta = (\theta_1, \theta_2) = (\mu, \Sigma)$)

- Suppose you have $x_1, x_2, \dots, x_n \sim (\text{i.i.d.}) N(\mu, \Sigma)$
- You don't know μ or Σ
- MLE:
 - For which $\theta = (\mu, \Sigma)$ is x_1, x_2, \dots, x_n most likely?

$$\mu^{\text{mle}} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k \quad \mu_i^{\text{mle}} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_{ki}$$

Where $1 \leq i \leq m$
And x_{ki} is value of the i^{th} component of \mathbf{x}_k .
(the i^{th} attribute of the k^{th} record)
And μ_i^{mle} is the i^{th} component of μ^{mle}

 TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 29

MLE Gaussians from Data

(m-Dimensional $\theta = (\theta_1, \theta_2) = (\mu, \Sigma)$)

- Suppose you have $x_1, x_2, \dots, x_n \sim (\text{i.i.d.}) N(\mu, \Sigma)$
- You don't know μ or Σ
- MLE:
 - For which $\theta = (\mu, \Sigma)$ is x_1, x_2, \dots, x_n most likely?

$$\sigma_{ij}^{\text{mle}} = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_{ki} - \mu_i^{\text{mle}})(\mathbf{x}_{kj} - \mu_j^{\text{mle}})$$

Where $1 \leq i \leq m, 1 \leq j \leq m$
And x_{ki} is value of the i^{th} component of \mathbf{x}_k .
(the i^{th} attribute of the k^{th} record)
And σ_{ij}^{mle} is the $(i,j)^{\text{th}}$ component of Σ^{mle}

 TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 30

MLE Gaussians from Data

(m-Dimensional $\theta = (\theta_1, \theta_2) = (\mu, \Sigma)$)

- Suppose you have $x_1, x_2, \dots, x_n \sim (\text{i.i.d.}) N(\mu, \Sigma)$
- You don't know μ or Σ
- MLE:
 - For which $\theta = (\mu, \Sigma)$ is x_1, x_2, \dots, x_n most likely?

$$\sigma_{ij}^{mle} = \frac{1}{n} \sum_{k=1}^n (x_{ki} - \mu_i^{mle})(x_{kj} - \mu_j^{mle})$$

Where $1 \leq i \leq m, 1 \leq j \leq m$
And x_{ki} is value of the i^{th} component of x_k (the i^{th} attribute of the k^{th} record)
And σ_{ij}^{mle} is the $(i,j)^{\text{th}}$ component of Σ^{mle}

FLIP: Pattern Recognition (G. Li)
31

MLE Gaussians from Data

(m-Dimensional $\theta = (\theta_1, \theta_2) = (\mu, \Sigma)$)

- Suppose you have $x_1, x_2, \dots, x_n \sim (\text{i.i.d.}) N(\mu, \Sigma)$
- You don't know μ or Σ
- MLE:
 - For which $\theta = (\mu, \Sigma)$ is x_1, x_2, \dots, x_n most likely?

$$\Sigma^{\text{unbiased}} = \frac{\Sigma^{mle}}{1 - \frac{1}{n}} = \frac{1}{n-1} \sum_{k=1}^n (x_k - \mu^{mle})(x_k - \mu^{mle})^T$$

Note how Σ^{mle} is forced to be symmetric non-negative definite

FLIP: Pattern Recognition (G. Li)
32

Summary on MLE

- ML estimation is usually simpler than alternative methods.
 - More accurate estimates as the number of training samples increases.
 - If the model chosen for $p(x|\theta)$ is correct, and independence assumptions among variables are true, ML will give very good results.
 - Otherwise, ML will give poor results.

FLIP: Pattern Recognition (G. Li)
33

Bayesian Estimation

- A General Strategy
- Learning Gaussian from Data
- Summary

FLIP: Pattern Recognition (G. Li)
34

Bayesian Estimation

- Assumes that the parameters θ are random variables that have some known a-priori distribution $p(\theta)$.
- Using the training examples D , BE converts $p(\theta)$ to $p(\theta|D)$.
 - BE estimates a distribution rather than making point estimates like ML.

$$p(\mathbf{x}/D) = \int p(\mathbf{x}/\theta)p(\theta/D)d\theta$$

- Note:** BE solution might not be of the parametric form assumed.

FLIP: Pattern Recognition (G. Li)
35

The Role of Training Examples

- If $p(\mathbf{x}/\omega_i)$ and $P(\omega_i)$ are known, Bayes' rule allows us to compute the posterior probabilities $P(\omega_i | \mathbf{x})$:

$$P(\omega_i | \mathbf{x}) = \frac{p(\mathbf{x}/\omega_i)p(\omega_i)}{\sum_j p(\mathbf{x}/\omega_j)p(\omega_j)}$$

- The role of the training examples D_i can be emphasized by using them in the computation of the posterior probabilities:

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

FLIP: Pattern Recognition (G. Li)
36

The Role of Training Examples



$$\begin{aligned} P(\omega_i / \mathbf{x}, D_i) &= \frac{p(\mathbf{x}, D_i / \omega_i)P(\omega_i)}{p(\mathbf{x}, D_i)} = \frac{p(\mathbf{x} / D_i, \omega_i)p(D_i / \omega_i)P(\omega_i)}{p(\mathbf{x} / D_i)p(D_i)} \\ &= \frac{p(\mathbf{x} / \omega_i, D_i)P(\omega_i / D_i)}{p(\mathbf{x} / D_i)} = \frac{p(\mathbf{x} / \omega_i, D_i)P(\omega_i / D_i)}{\sum_j p(\mathbf{x}, \omega_j / D_j)} \\ &= \frac{p(\mathbf{x} / \omega_i, D_i)P(\omega_i / D_i)}{\sum_j p(\mathbf{x} / \omega_j, D_j)P(\omega_j / D_j)} \end{aligned}$$

 FLIP: Pattern Recognition (G. Li) 37

The Role of Training Examples

- This implies that the training examples D_i can help us to determine both the class-conditional densities and the prior probabilities:

$$P(\omega_i / \mathbf{x}, D_i) = \frac{p(\mathbf{x} / \omega_i, D_i)P(\omega_i / D_i)}{\sum_j p(\mathbf{x} / \omega_j, D_j)P(\omega_j / D_j)}$$

- To simplify things, let's assume that $P(\omega_i / D) = P(\omega_i)$:

$$P(\omega_i / \mathbf{x}, D_i) = \frac{p(\mathbf{x} / \omega_i, D_i)P(\omega_i)}{\sum_j p(\mathbf{x} / \omega_j, D_j)P(\omega_j)}$$

 FLIP: Pattern Recognition (G. Li) 38

Bayesian Estimation (Problem Definition)



- Problem: need to estimate $p(\mathbf{x} | \omega_i, D)$ for every class ω_i
 - If samples in D_i give no information about θ_i ($i \neq j$), we need to solve c independent problems (i.e., one for each class)

“given D , estimate $p(\mathbf{x} | D)$

 FLIP: Pattern Recognition (G. Li) 39

A General BE strategy



- Estimate $p(\mathbf{x} | D)$ as follows:

$$p(\mathbf{x} / D) = \int p(\mathbf{x}, \boldsymbol{\theta} / D) d\boldsymbol{\theta} = \int p(\mathbf{x} / \boldsymbol{\theta}, D) p(\boldsymbol{\theta} / D) d\boldsymbol{\theta}$$

- Since the distribution is known completely given $\boldsymbol{\theta}$, we have:

$$p(\mathbf{x} / D) = \int p(\mathbf{x} / \boldsymbol{\theta}) p(\boldsymbol{\theta} / D) d\boldsymbol{\theta}$$

- Important equation to link $p(\mathbf{x} | D)$ with $p(\boldsymbol{\theta} | D)$

 FLIP: Pattern Recognition (G. Li) 40

A General BE strategy



- Steps
 - Compute $p(\boldsymbol{\theta} / D)$

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

- Compute $p(\mathbf{x} / D)$

$$p(\mathbf{x} / D) = \int p(\mathbf{x} / \boldsymbol{\theta})p(\boldsymbol{\theta} / D) d\boldsymbol{\theta}$$

- $p(\mathbf{x} / D)$ will replace $p(\mathbf{x} / \omega_i, D_i)$ in Bayes' rule:

$$P(\omega_i / \mathbf{x}, D_i) = \frac{p(\mathbf{x} / \omega_i, D_i)P(\omega_i)}{\sum_j p(\mathbf{x} / \omega_j, D_j)P(\omega_j)}$$

 FLIP: Pattern Recognition (G. Li) 41

BE Gaussians from Data (Case of Unknown $\theta = \mu$)



- Estimate $\boldsymbol{\theta}$ using the a-posteriori density $P(\boldsymbol{\theta} | D)$
 - The univariate Gaussian case: $P(\mu | D)$
 - μ is the only unknown parameter
 - Prior distribution $P(\mu)$ is provided

$$P(\mathbf{x} | \mu) \sim N(\mu, \sigma^2)$$

$$P(\mu) \sim N(\mu_0, \sigma_0^2)$$

 FLIP: Pattern Recognition (G. Li) 42

BE Gaussians from Data (Case of Unknown $\theta = \mu$)

$$P(\mu | D) = \frac{P(D | \mu).P(\mu)}{\int P(D | \mu).P(\mu) d\mu} \quad (1)$$

$$= \alpha \prod_{k=1}^{k=n} P(x_k | \mu).P(\mu)$$

$$P(\mu | D) \sim N(\mu_n, \sigma_n^2) \quad (2)$$

- The updated parameters of the prior:

$$\mu_n = \left(\frac{n\sigma_0^2}{n_0\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \cdot \mu_0, \text{ and } \sigma_n^2 = \frac{\sigma_0^2\sigma^2}{n\sigma_0^2 + \sigma^2}$$

- μ_n is our best guess for μ ; σ_n^2 measures the uncertainty.

BE Gaussians from Data (Case of Unknown $\theta = \mu$)

- The updated parameters of the prior:

$$\mu_n = \left(\frac{n\sigma_0^2}{n_0\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \cdot \mu_0, \text{ and } \sigma_n^2 = \frac{\sigma_0^2\sigma^2}{n\sigma_0^2 + \sigma^2}$$

- σ_n^2 approaches σ^2/n as n increases (more observations will decrease our uncertainty about μ)

- $p(\mu | D)$ becomes sharply peaked (*Bayesian Learning*) as n increases.

BE Gaussians from Data (Case of Unknown $\theta = \mu$)

- The univariate case $P(x | D)$
 - $P(\mu | D)$ has been computed
 - $P(x | D)$ remains to be computed!

$$P(x | D) = \int P(x | \mu).P(\mu | D) d\mu \text{ is Gaussian}$$

$$P(x | D) \sim N(\mu_n, \sigma_n^2 + \sigma_x^2)$$

BE Gaussians from Data (Case of Unknown $\theta = \mu$)

- The updated parameters of the prior:

$$\mu_n = \left(\frac{n\sigma_0^2}{n_0\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \cdot \mu_0, \text{ and } \sigma_n^2 = \frac{\sigma_0^2\sigma^2}{n\sigma_0^2 + \sigma^2}$$

- μ_n is a linear combination of μ_n and μ_0

- If $\sigma \neq 0$, then μ_n approaches μ_n (maximum likelihood estimate !!)

- If $\sigma_0=0$, then $\mu = \mu_0$; If $\sigma \gg \sigma$ then $\mu_n = \mu_0$

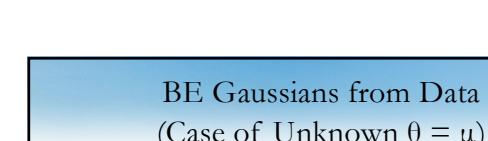


FIGURE 3.2. Bayesian learning of the mean of normal distributions in one and two dimensions. The posterior distribution estimates are labeled by the number of training samples used in the estimation. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001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

BE Gaussians from Data (Case of Unknown $\theta = \mu$)

- Bayesian Classification Rule
 - Desired class-conditional density $P(x | D_j, \omega_j)$
 - $P(x | D_j, \omega_j)$ together with $P(\omega_j)$ and using Bayes formula, we obtain the Bayesian classification rule:

$$\underset{\omega_j}{\operatorname{Max}} [P(\omega_j | x, D)] = \underset{\omega_j}{\operatorname{Max}} [P(x | D_j, \omega_j) \cdot P(\omega_j)]$$

Summary on BE

- $P(x|D)$ computation can be applied to any situation in which the unknown density can be parameterized. The basic assumptions are:
 - The form of $P(x|\theta)$ is assumed known, but the value of θ is not known exactly
 - Our knowledge about θ is assumed to be contained in a known prior density $P(\theta)$
 - The rest of our knowledge about θ is contained in a set D of n random variables x_1, x_2, \dots, x_n that follows $P(x)$



FLIP: Pattern Recognition (G. Li)

49



Summary on BE

- The basic problem is:
“Compute the posterior density $P(\theta|D)$ ”
 then “**Derive $P(x|D)$** ”
- Using Bayes formula, we have:

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

– by independence assumption:

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



FLIP: Pattern Recognition (G. Li)

50



Summary on BE

- If we are less certain about the exact value of θ , we should consider a weighted average of $p(x|\theta)$ over the possible values of θ

$$p(x/D) = \int p(x/\theta)p(\theta/D)d\theta$$
 - Bayesian estimation approach estimates a distribution for $p(x|D)$ rather than making point estimates like ML



FLIP: Pattern Recognition (G. Li)

51



Summary on BE

- Suppose $p(\theta|D)$ peaks very sharply at $\theta = \hat{\theta}$ and $p(\hat{\theta}) \neq 0$ then $p(x|D)$ can be approximated as follows:

$$p(x/D) \approx p(x/\hat{\theta})$$
 - i.e., the best estimate is obtained by setting $\theta = \hat{\theta}$
 - This is the ML solution
 - i.e., $p(D|\theta)$ peaks at $\hat{\theta}$ too

$$\text{since } p(\theta/D) = \frac{p(D/\theta)p(\theta)}{p(D)}$$



FLIP: Pattern Recognition (G. Li)

52



Summary on BE

- Given a large number of samples, $p(\theta|D)$ will have a very strong peak at $\hat{\theta}$; in this case:

$$p(x/D) \approx p(x/\hat{\theta})$$
 - When $p(\theta|D)$ contains more than one peaks (i.e., more than one θ explains the data), the solution $p(x|\theta)$ should be obtained by integration.

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



FLIP: Pattern Recognition (G. Li)

53



ML vs BE

- Comparisons



FLIP: Pattern Recognition (G. Li)

54



MLE vs Bayesian Estimation

<u>Number of training data</u>	<u>Computational complexity</u>
<ul style="list-style-type: none"> The two methods are equivalent assuming infinite number of training data <ul style="list-style-type: none"> prior distributions do not exclude the true solution For small training data sets, they give different results in most cases. 	<ul style="list-style-type: none"> ML uses differential calculus or gradient search for maximizing the likelihood. Bayesian estimation requires complex multidimensional integration techniques.

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 55

MLE vs Bayesian Estimation

<u>Solution complexity</u>	<u>Prior distribution</u>
<ul style="list-style-type: none"> Easier to interpret ML solutions (i.e., must be of the assumed parametric form). A Bayesian estimation solution might not be of the parametric form assumed. 	<ul style="list-style-type: none"> If the prior distribution $p(\theta)$ is uniform, Bayesian estimation solutions are equivalent to ML solutions. Otherwise, the two methods will give different solutions.

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 56

MLE vs Bayesian Estimation

- General comments
 - There are strong theoretical and methodological arguments supporting Bayesian estimation.
 - In practice, ML estimation is simpler and can lead to comparable performance.

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 57

Seminar S03

- Topic**
 - Give a presentation on how to choose the prior density in Bayesian Estimation
- Requirements**
 - Prepare a **15 minutes** talk on your chosen topic
 - Make **ppt** to assist your talk
 - Prepare **at least 3 questions** to ask the audience after your talk
 - Get ready to **take questions** from the audience
- Hints**
 - Check concepts such as
 - Conjugate Prior**
 - Beta distribution, Gamma distribution**
 - Dirichlet distribution**
 - Inverse-Wishart distribution**

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 58

Exercise

- Suppose categorical arity-n inputs $x_1, x_2, \dots, x_n \sim$ (i.i.d.) from a multinomial $M(p_1, p_2, \dots, p_n)$ where $P(x_k=j|p)=p_j$
- What is the MLE $\hat{p}=(\hat{p}_1, \hat{p}_2, \dots, \hat{p}_n)$?

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 59

Questions?

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 60