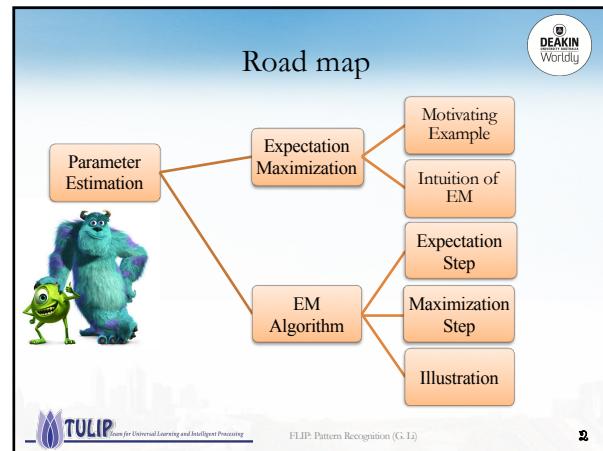


Lecture Notes on
Pattern Recognition

Session 03(B): Parameter Estimation (II)

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

DEAKIN Worldwide



Expectation Maximization

- Motivation Example
- Intuition of EM

DEAKIN Worldwide

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li) 3

Expectation Maximization

- The Expectation Maximization Algorithm is an approach to handle the problem of learning in the presence of unobserved variables.
 - When data is only partially observable
 - Unsupervised Clustering
 - Target Value is unobservable
 - Supervised Learning
 - Some instance attributes unobservable
 - Old idea (late 50's) but formalized by Dempster, Laird and Rubin in 1977

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li) 4

A Silly Example

- Let events be “grades in a class”
 - w_1 = Gets an A
 - w_2 = Gets a B
 - w_3 = Gets a C
 - w_4 = Gets a D

(Note $0 \leq \mu \leq 1/6$)

- Assume we want to estimate μ from data.
- In a given class there were

A	B	C	D
a	b	c	d

- What’s the maximum likelihood estimate of μ given a,b,c,d ?

DEAKIN Worldwide

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li) 5

A Silly Example

- Let events be “grades in a class”
 - w_1 = Gets an A
 - w_2 = Gets a B
 - w_3 = Gets a C
 - w_4 = Gets a D

(Note $0 \leq \mu \leq 1/6$)

- Assume we want to estimate μ from data.
- In a given class there were

A	B	C	D
14	6	9	10

- What’s the maximum likelihood estimate of μ given a,b,c,d ?

DEAKIN Worldwide

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li) 6

Trivial Statistics

$P(A) = \frac{1}{2}$ $P(B) = \mu$ $P(C) = 2\mu$ $P(D) = \frac{1}{2} - 3\mu$

$$P(a, b, c, d | \mu) = K(\frac{1}{2})^a (\mu)^b (2\mu)^c (\frac{1}{2} - 3\mu)^d$$

$$\log P(a, b, c, d | \mu) = \log K + a \log \frac{1}{2} + b \log \mu + c \log 2\mu + d \log (\frac{1}{2} - 3\mu)$$

FOR MAX LIKELIHOOD, SET $\frac{\partial \text{LogP}}{\partial \mu} = 0$

$$\frac{\partial \text{LogP}}{\partial \mu} = \frac{b}{\mu} + \frac{2c}{2\mu} - \frac{3d}{\frac{1}{2} - 3\mu} = 0$$

Gives max likelihood $\mu = \frac{b+c}{6(b+c+d)}$

So if class got

A	B	C	D
14	1	10	0

Max likelihood $\mu = \frac{1}{10}$

Boring, but true!

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

A Not-So-Silly Example

- Someone tells us that
- Number of High grades (A's + B's) = b
- Number of C's = c
- Number of D's = d
- What is the max. likelihood estimate of μ now?

Assume

$P(A) = \frac{1}{2}$
$P(B) = \mu$
$P(C) = 2\mu$
$P(D) = \frac{1}{2} - 3\mu$

EXPECTATION If we know the value of μ we could compute the expected value of a and b

Since the ratio ab should be the same as the ratio $\frac{1}{2} : \mu$

$$a = \frac{1}{2}h \quad b = \frac{\mu}{\frac{1}{2} + \mu}h$$

MAXIMIZATION If we know the expected values of a and b we could compute the maximum likelihood value of μ

$$\mu = \frac{b+c}{6(b+c+d)}$$

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

A Taste of EM Algorithm

- We begin with a guess for μ
- We iterate between E Step and M Step to improve our estimates of μ and a and b .

Define $\mu(t)$ the estimate of μ on the t-th iteration
 $b(t)$ the estimate of b on t-th iteration

$\mu(0) = \text{initial guess}$
 $b(t) = \frac{\mu(t)h}{\frac{1}{2} + \mu(t)} = E[b | \mu(t)]$
 $\mu(t+1) = \frac{b(t) + c}{6(b(t) + c + d)}$
= max like est of μ given $b(t)$

E-step
M-step

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

A Taste of EM Algorithm

- We begin with a guess for μ
- We iterate between E Step and M Step to improve our estimates of μ and a and b .

Define $\mu(t)$ the estimate of μ on the t-th iteration
 $b(t)$ the estimate of b on t-th iteration

suppose we had

t	$\mu(t)$	$b(t)$
0	0	0
1	0.0833	2.857
2	0.0937	3.158
3	0.0947	3.185
4	0.0948	3.187
5	0.0948	3.187
6	0.0948	3.187

Convergence proof based on fact that $\text{Prob}(\text{data} | \mu)$ must increase or remain same between each iteration, and \leq

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

Probabilistic model

Imagine model generating data

- Need to introduce label, z , for each data point
- Label is called a *latent variable*, also called *hidden, unobserved, missing*

Simplification:
if we knew the labels, we can decouple the components as estimate parameters separately for each one

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

Motivating example

- Task: Fit mixture of Gaussian model with $C=2$ components $\theta^* = \arg\max_{\theta} \prod_{i=1}^N p(x_i | \theta)$

Data: $x = (x_1, x_2, \dots, x_N)$

$$p(x_i | \theta) = \sum_c p(c | \theta) p(x_i | c, \theta) = \sum_c w(c) N(x_i | \mu_c, \sigma_c)$$

where $\sum_{c=1}^C w(c) = 1$

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

Motivating example

- Task: Fit mixture of Gaussian model with C=2 components $\theta^* = \operatorname{argmax}_{\theta} \prod_{n=1}^N p(x_i|\theta)$

Data: $x = (x_1, x_2, \dots, x_N)$

Parameters: $\theta = \{w, \mu, \sigma\}$
keep w, σ fixed
i.e. only estimate μ

TULIP Team for Universal Learning and Intelligent Processing

Intuition of EM

E-step: Compute a distribution on the labels of the points, using current parameters

M-step: Update parameters using current guess of label distribution.

TULIP Team for Universal Learning and Intelligent Processing

Expectation Maximization

- Consider a model with parameters θ and a number of observations (and hidden variables) D from which we want to estimate θ .
- MLE uses the probability distributions $P(D|\theta)$ in an unusual way.
 - D is treated as the constant and θ as the variable.
 - Using this perspective, the function value of p is viewed as a likelihood rather than a probability distribution value.

TULIP Team for Universal Learning and Intelligent Processing

Likelihood function

- Likelihood** is a function of parameters θ
Probability is a function of random variable x or D

TULIP Team for Universal Learning and Intelligent Processing

Expectation Maximization

- Task: find θ that maximizes the likelihood of D.
 $\theta^* = \operatorname{argmax} P(D|\theta) = \operatorname{argmax} \ln P(D|\theta)$
 - $\ln P(D|\theta)$ is used for simplicities sake.
 - Depending on the form of $P(D|\theta)$ this problem can be easy or hard.
 - The more complex the distribution, the more complex the technique of using the MLE.
 - While MLE is popular, other estimators may be used

TULIP Team for Universal Learning and Intelligent Processing

Maximum Likelihood Estimation

TULIP Team for Universal Learning and Intelligent Processing

EM Algorithm

- Expectation Step
- Maximization Step
- Illustration

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

19

Jensen's Inequality

- Jensen's Inequality:**
- For a real continuous concave function $f()$ and $\lambda_j \geq 0$, $\sum_j \lambda_j = 1$, we have

$$f(\sum_j \lambda_j x_j) \geq \sum_j \lambda_j f(x_j)$$

Equality holds when all x are the same

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

20

Terminology

- Observed data:**
$$x = (x_1, x_2, \dots, x_N)$$

Continuous I.I.D

- Latent variables:**
$$z = (z_1, z_2, \dots, z_N)$$

Discrete 1 ... C

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

21

Terminology

- Incomplete Log-Likelihood (ILL):**
$$l(\theta; x) = \log p(x|\theta) = \log \prod_x p(x|\theta) = \sum_x \log \sum_z p(x, z|\theta)$$
- Complete Log-Likelihood (CLL):**
$$l_c(\theta; x, z) \triangleq \sum_x \log p(x, z | \theta)$$
- Expected Complete Log-Likelihood (ECLL):**
$$\langle l_c(\theta; x, z) \rangle_q \triangleq \sum_x \sum_z q(z | x, \theta) \log p(x, z | \theta)$$

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

22

Lower bound on log-likelihood

$$\begin{aligned}
 l(\theta; x) &= \log p(x | \theta) \\
 &= \log \sum_z p(x, z | \theta) \\
 &= \log \sum_z q(z | x) \frac{p(x, z | \theta)}{q(z | x)} \\
 &\stackrel{\text{Use Jensen's inequality}}{\geq} \sum_z q(z | x) \log \frac{p(x, z | \theta)}{q(z | x)} \\
 &\triangleq \mathcal{L}(q, \theta), \quad \text{AUXILIARY FUNCTION}
 \end{aligned}$$

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

23

EM is alternating ascent

- Key: ILL is always larger than the Auxiliary function
$$\mathcal{L}(q, \theta) \leq l(\theta; x)$$
- EM Algorithm:
 - Alternately improving q then θ , is guaranteed to improve likelihood itself....

(E step)	$q^{(t+1)} = \arg \max_q \mathcal{L}(q, \theta^{(t)})$
(M step)	$\theta^{(t+1)} = \arg \max_{\theta} \mathcal{L}(q^{(t+1)}, \theta)$.

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

24

Expectation-step:

Choosing the optimal $q(z|x, \theta)$

- Turns out that $q(z|x, \theta) = p(z|x, \theta^t)$ is the best

$$\begin{aligned}\mathcal{L}(p(z|x, \theta^{(t)}), \theta^{(t)}) &= \sum_z p(z|x, \theta^{(t)}) \log \frac{p(x, z|\theta)}{p(z|x, \theta^{(t)})} \\ &= \sum_z p(z|x, \theta^{(t)}) \log p(x|z, \theta^{(t)}) \\ &= \log p(x|\theta^{(t)}) \\ &= l(\theta^{(t)}; x).\end{aligned}$$



Expectation-step:

Choosing the optimal $q(z|x, \theta)$

- Turns out that $q(z|x, \theta) = p(z|x, \theta^t)$ is the best

$$\begin{aligned}l(\theta; x) - \mathcal{L}(q, \theta) &= l(\theta; x) - \sum_z q(z|x) \log \frac{p(x, z|\theta)}{q(z|x)} \\ &= \sum_z q(z|x) \log p(x|\theta) - \sum_z q(z|x) \log \frac{p(x, z|\theta)}{q(z|x)} \\ &= \sum_z q(z|x) \log p(x|\theta) - \log \frac{q(z|x)}{p(z|x, \theta)} \\ &= D(q(z|x) \| p(z|x, \theta)).\end{aligned}$$



ALTERNATE DERIVATION

Expectation-step:
Choosing the optimal $q(z|x, \theta)$

- Turns out that $q(z|x, \theta) = p(z|x, \theta^t)$ is the best

$$\begin{aligned}p(z_i = c|x_i, \theta^t) &= \frac{p(x_i, z_i = c|\theta^t)}{\sum_c p(x_i, z_i = c|\theta^t)} \\ &= \frac{\pi(z_i = c) N(x_i|\mu_c, \sigma_c)}{\sum_d \pi(z_i = d) N(x_i|\mu_d, \sigma_d)} \\ &\quad \begin{array}{c} \bar{x} \quad \bar{z} \\ \text{Component 1} \\ \text{Component 2} \end{array} \quad p(z|x, \theta^t)\end{aligned}$$



Maximization-step:
Choosing the optimal θ

- Auxiliary function separates into **ECLL** and **entropy** term:

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_z q(z|x) \log \frac{p(x, z|\theta)}{q(z|x)} \\ &= \sum_z q(z|x) \log p(x, z|\theta) - \sum_z q(z|x) \log q(z|x) \\ &= \langle l_c(\theta; x, z) \rangle_q - \sum_z q(z|x) \log q(z|x),\end{aligned}$$



Maximization-step:
Choosing the optimal θ

- Take derivatives to maximize the **ECLL**:

$$\begin{aligned}\langle l_c(\theta; x, z) \rangle_q &\triangleq \sum_x \sum_z q(z|x, \theta) \log p(x, z|\theta) \\ \frac{\partial \mathcal{L}(q^{t+1}, \theta)}{\partial \theta} &= \frac{\partial \langle l_c(\theta; x, z) \rangle_{q^{t+1}}}{\partial \theta} \\ &= \frac{\partial \sum_x \sum_z p^{t+1}(z|x, \theta^t) \log p(x, z|\theta)}{\partial \theta} = 0\end{aligned}$$

From E-step



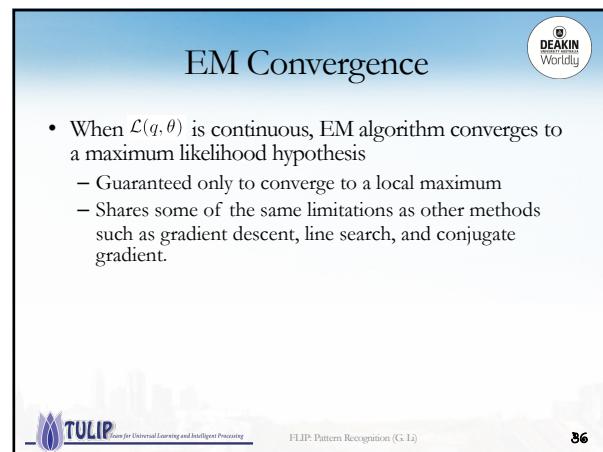
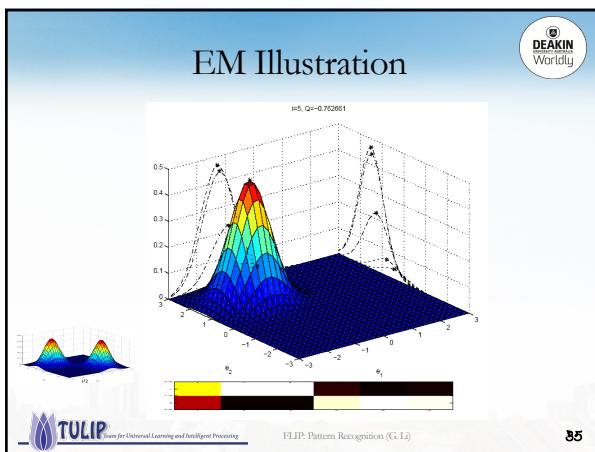
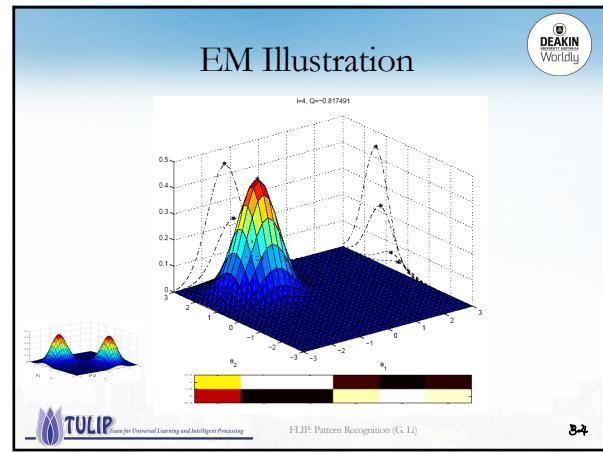
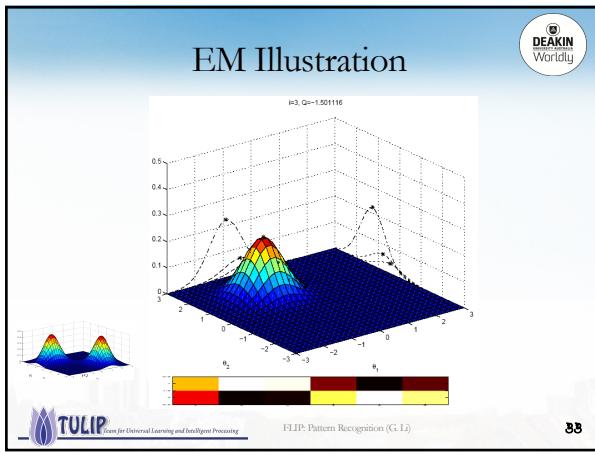
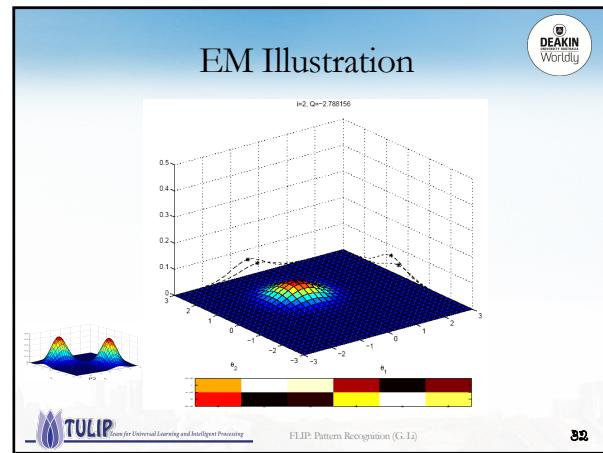
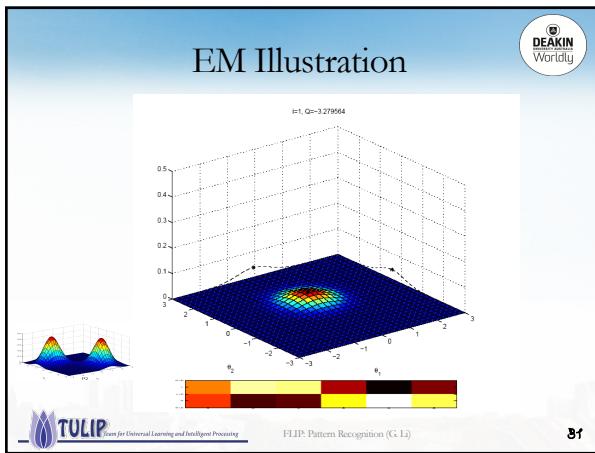
Maximization-step:
Choosing the optimal θ

- Solution to the motivation example:

- Calculate the new $\theta^{(t+1)} = (\mathbf{u}_1^{(t+1)}, \mathbf{u}_2^{(t+1)})$, assuming the value taken by setting the hidden variable z as $p^{(t+1)}(z|x, \theta^t)$
- Replace $\theta^t = (\mathbf{u}_1^t, \mathbf{u}_2^t)$ by $\theta^{(t+1)} = (\mathbf{u}_1^{(t+1)}, \mathbf{u}_2^{(t+1)})$
- Iterate...

$$\mathbf{u}_c^{(t+1)} = \frac{\sum_x p(z_i = c|x_i, \theta^t) x_i}{\sum_x p(z_i = c|x_i, \theta^t)}$$





Seminar S04



DEAKIN
Worldly

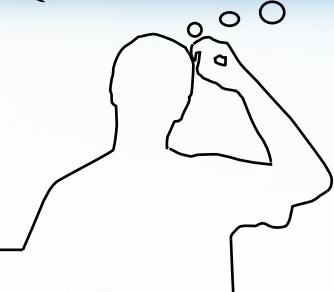
- **Topics**
 - Present one specific example of the EM algorithm used to solve a real world problem.
 - Connection between EM and K-Means
- **Requirements**
 - Prepare a **15 minutes** talk on your chosen topic
 - Make **ppt** to assist your talk
 - Prepare **at least 3 questions** to ask the audience after your talk
 - Get ready to **take questions** from the audience
- **Hints**
 - You can search for research articles from Google Scholar
 - Or illustrate a solution step by step using EM

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

37

Questions?



DEAKIN
Worldly

TULIP Team for Universal Learning and Intelligent Processing

FLIP: Pattern Recognition (G. Li)

38