

# CSE 312

## Autumn 2012

The Expectation-Maximization  
Algorithm

# 11/19 Puzzler

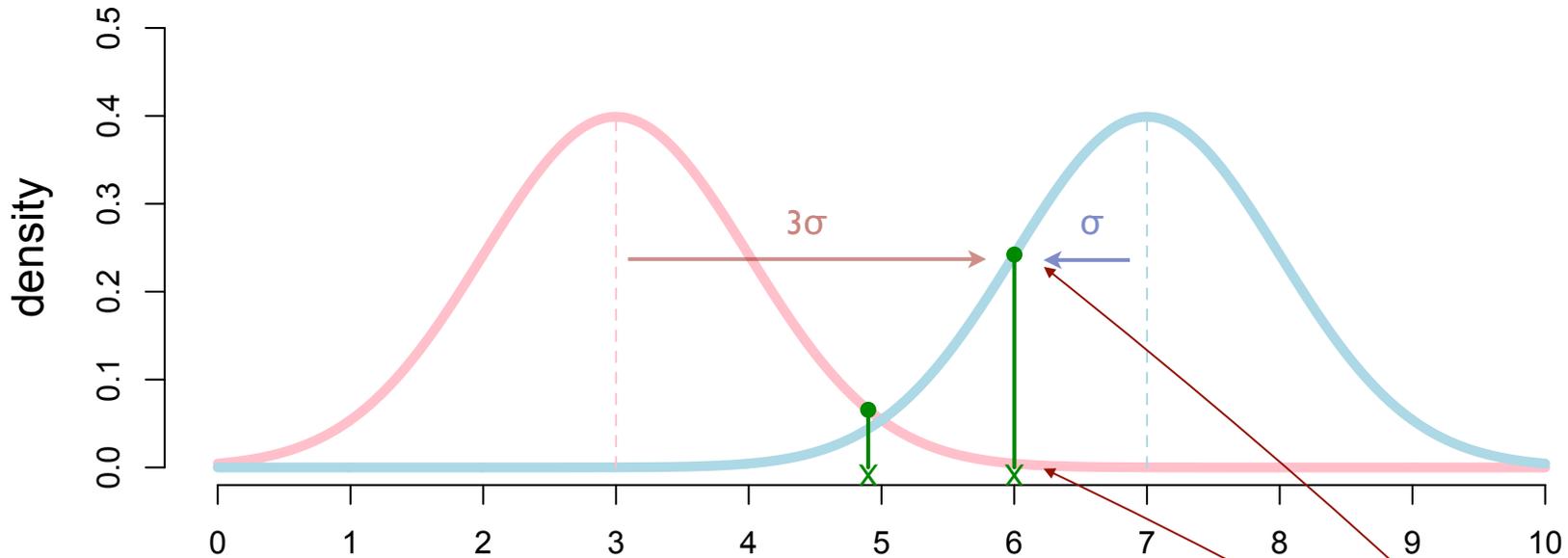
Two slips of paper in a hat:  $\mu=3$  and  $\mu=7$ . You draw one, then (without revealing  $\mu$ ) reveal a single sample  $X \sim \text{Normal}(\text{mean } \mu, \sigma^2 = 1)$ .

You happen to draw  $X = 6.001$ .

Dr. D. says “your slip = 7.” What is  $P(\text{correct})$ ?

What if  $X$  had been 4.9?

# A Hat Trick



Let “ $X \approx 6$ ” be a shorthand for  $6.001 - \delta/2 < X < 6.001 + \delta/2$

$$P(\mu = 7|X = 6) = \lim_{\delta \rightarrow 0} P(\mu = 7|X \approx 6)$$

$$P(\mu = 7|X \approx 6) = \frac{P(X \approx 6|\mu = 7)P(\mu = 7)}{P(X \approx 6)}$$

$$= \frac{0.5P(X \approx 6|\mu = 7)}{0.5P(X \approx 6|\mu = 3) + 0.5P(X \approx 6|\mu = 7)}$$

$$\approx \frac{f(X = 6|\mu = 7)\delta}{f(X = 6|\mu = 3)\delta + f(X = 6|\mu = 7)\delta}, \text{ so}$$

$$P(\mu = 7|X = 6) = \frac{f(X = 6|\mu = 7)}{f(X = 6|\mu = 3) + f(X = 6|\mu = 7)} \approx 0.982$$

$f$  = normal density

# Another Hat Trick

Two secret numbers,  $\mu_{\text{pink}}$  and  $\mu_{\text{blue}}$

On pink slips, many samples of  $\text{Normal}(\mu_{\text{pink}}, \sigma^2 = 1)$ ,

Ditto on blue slips, from  $\text{Normal}(\mu_{\text{blue}}, \sigma^2 = 1)$ .

Based on 16 of each, how would you “guess” the secrets (where “success” means your guess is within  $\pm 0.5$  of each secret)?

Roughly how likely is it that you will succeed?

## Hat Trick (cont.)

Pink/blue = red herrings; separate & independent

Given  $X_1, \dots, X_{16} \sim N(\mu, \sigma^2)$ ,  $\sigma^2 = 1$

Calculate  $Y = (X_1 + \dots + X_{16})/16 \sim N(?, ?)$

$$E[Y] = \mu$$

$$\text{Var}(Y) = 16\sigma^2/16^2 = \sigma^2/16 = 1/16$$

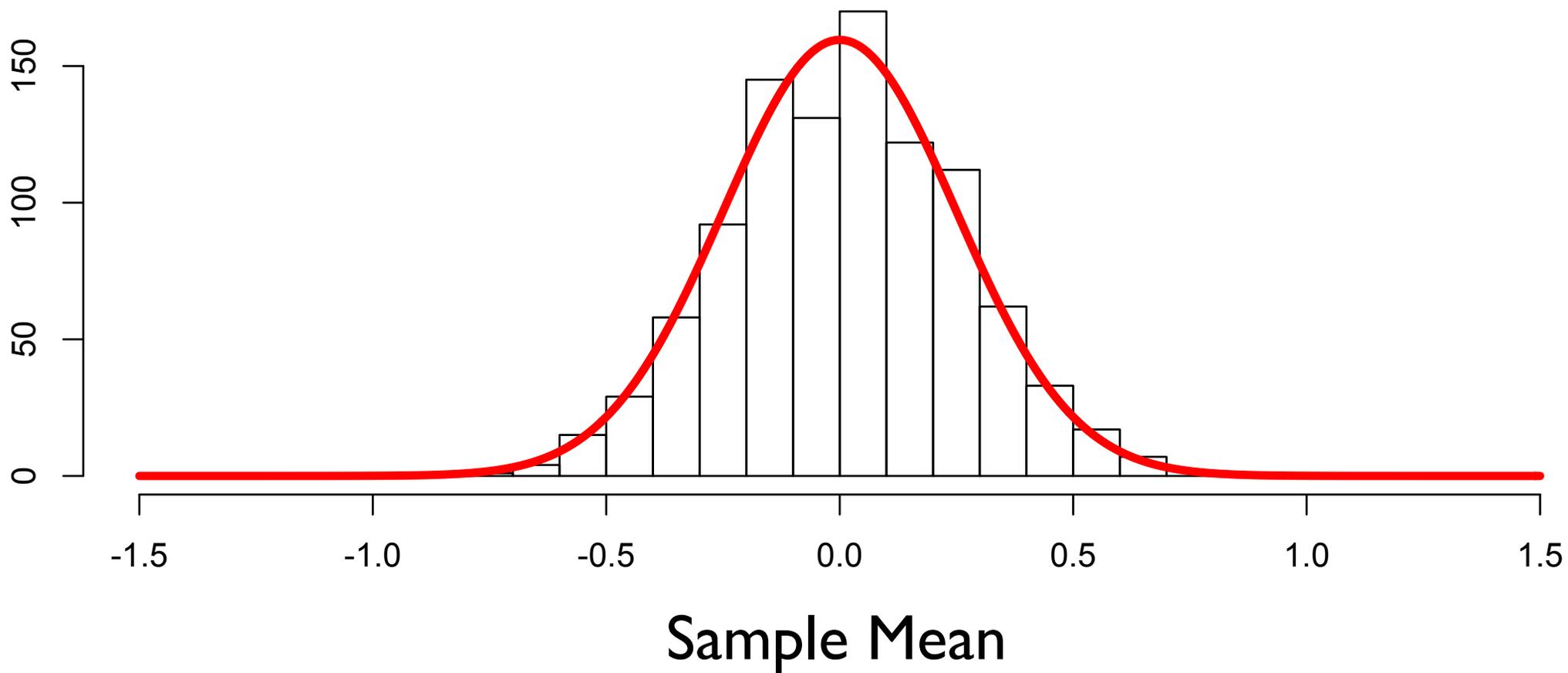
“Y within  $\pm 0.5$  of  $\mu$ ” = “Y within  $\pm 2 \sigma$  of  $\mu$ ”  $\approx 95\%$  prob

Note 1: Y is a *point estimate* for  $\mu$ ;

$Y \pm 2 \sigma$  is a *95% confidence interval* for  $\mu$

(More on this topic later)

**Histogram of 1000 samples of the average of 16  $N(0,1)$  RVs**

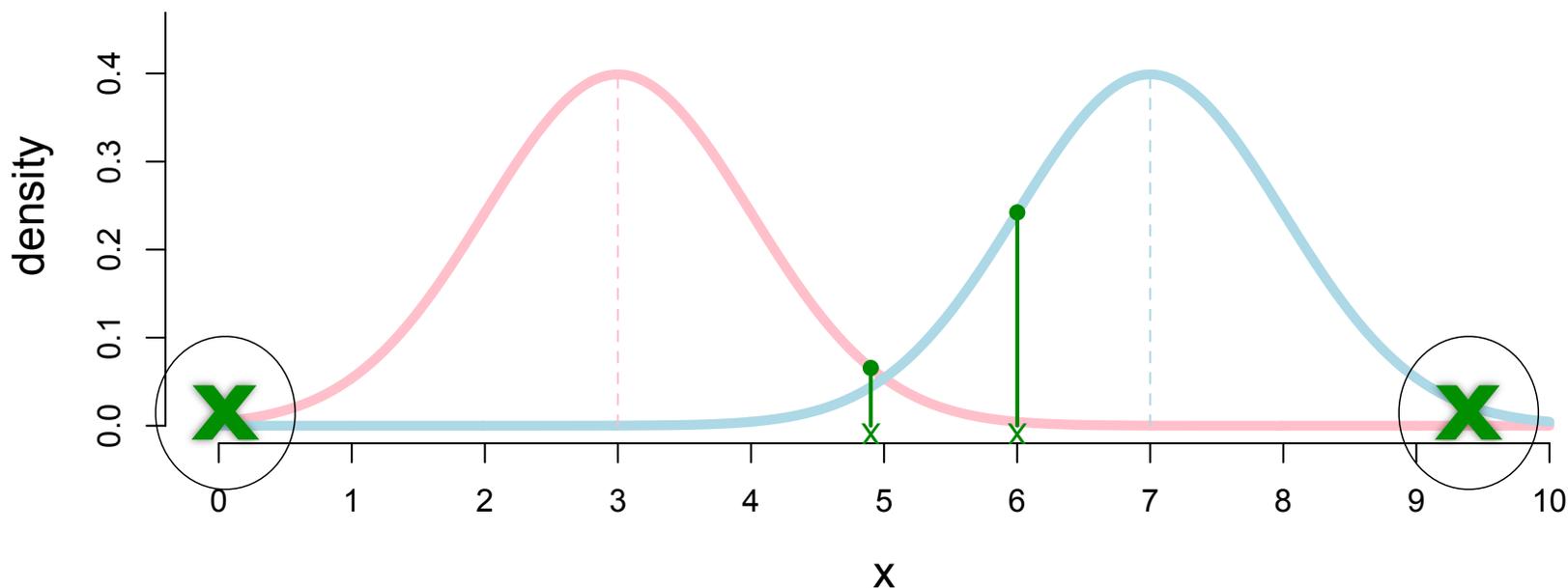


## Hat Trick (cont.)

Note 2: red/blue separation is just like the M-step of EM if values of the hidden variables ( $z_{ij}$ ) were known.

What if they're not? E.g., what would you do if some of the slips you pulled had coffee spilled on them, obscuring color?

If they were half way between means of the others?  
If they were on opposite sides of the means of the others

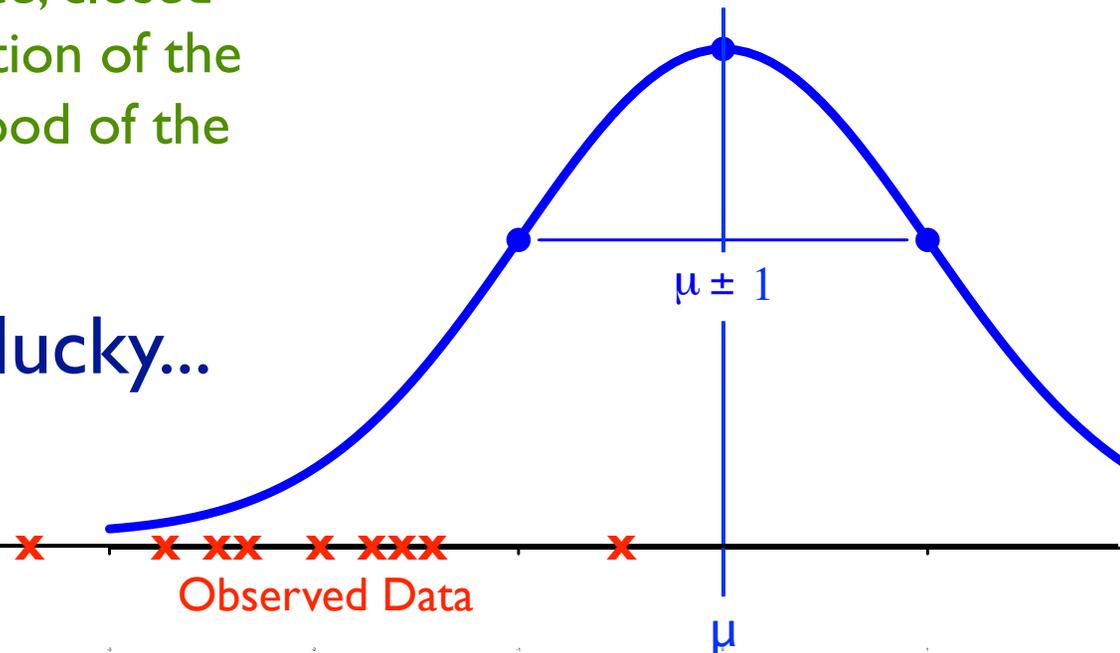


# Previously:

## How to estimate $\mu$ given data

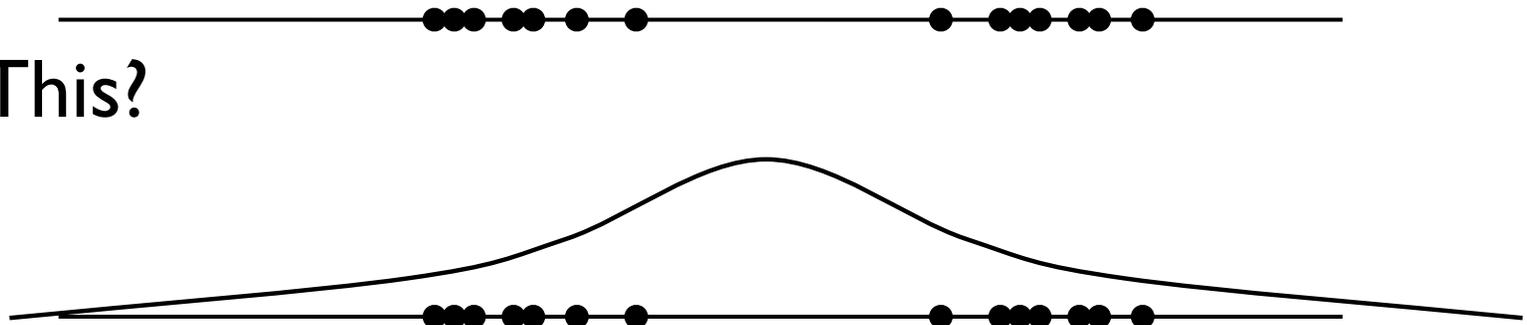
For this problem, we got a nice, closed form, solution, allowing calculation of the  $\mu, \sigma$  that maximize the likelihood of the observed data.

We're not always so lucky...

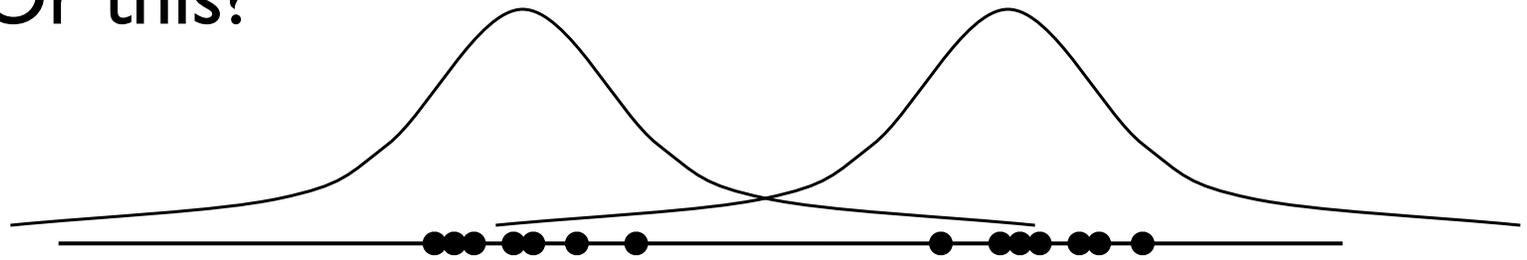


# More Complex Example

This?

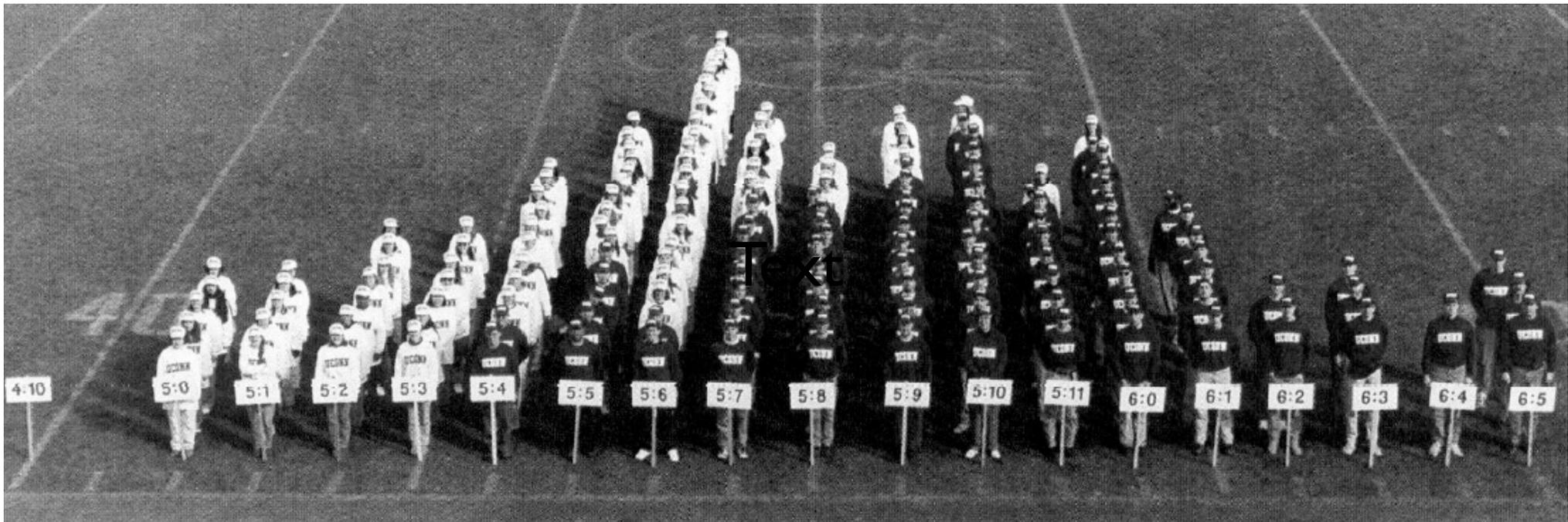


Or this?



(A modeling decision, not a math problem...,  
but if the later, what math?)

# A Living Histogram

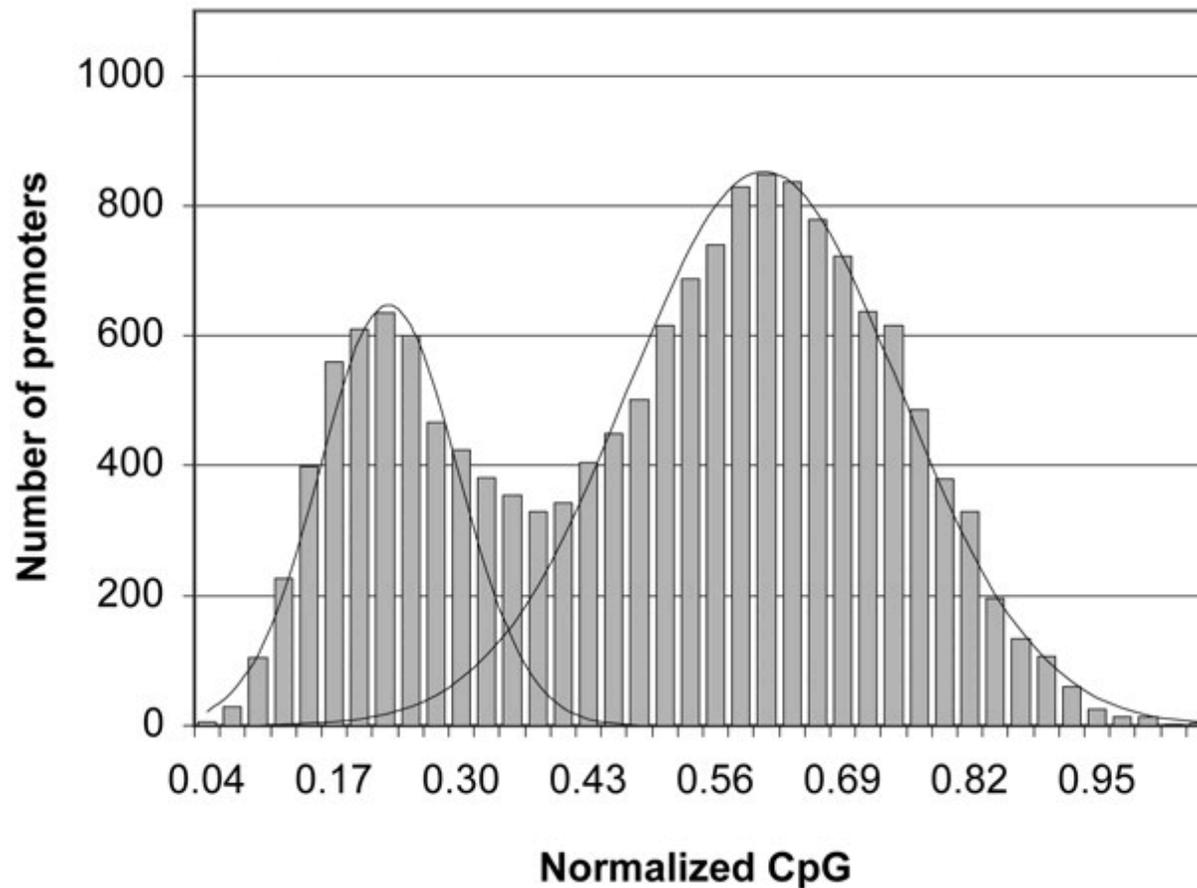


male and female genetics students, University of Connecticut in 1996

<http://mindprod.com/jgloss/histogram.html>

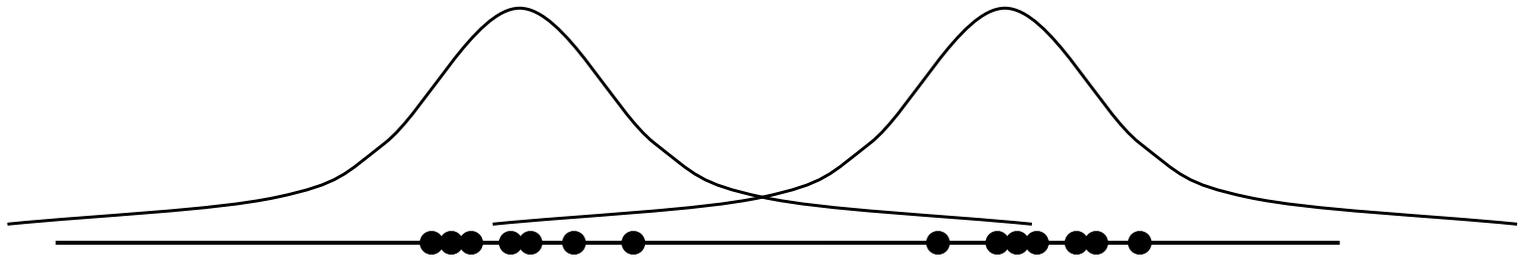
# A Real Example:

## CpG content of human gene promoters



“A genome-wide analysis of CpG dinucleotides in the human genome distinguishes two distinct classes of promoters” Saxonov, Berg, and Brutlag, PNAS 2006;103:1412-1417

# Gaussian Mixture Models / Model-based Clustering



Parameters  $\theta$

means	$\mu_1$	$\mu_2$
variances	$\sigma_1^2$	$\sigma_2^2$
mixing parameters	$\tau_1$	$\tau_2 = 1 - \tau_1$

P.D.F.  $\xrightarrow{\text{separately}}$   $f(x|\mu_1, \sigma_1^2)$     $f(x|\mu_2, \sigma_2^2)$

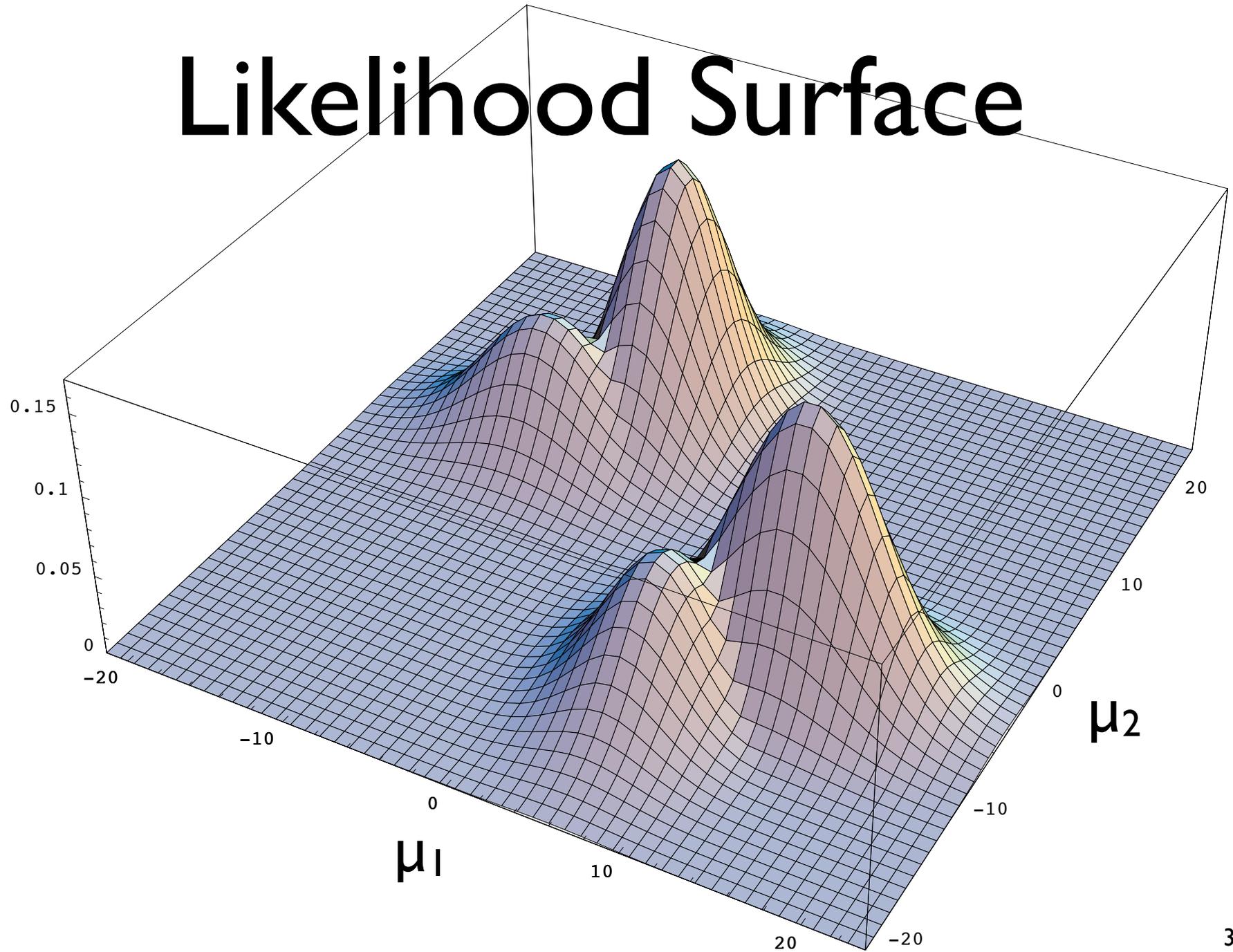
Likelihood  $\xrightarrow{\text{together}}$   $\tau_1 f(x|\mu_1, \sigma_1^2) + \tau_2 f(x|\mu_2, \sigma_2^2)$

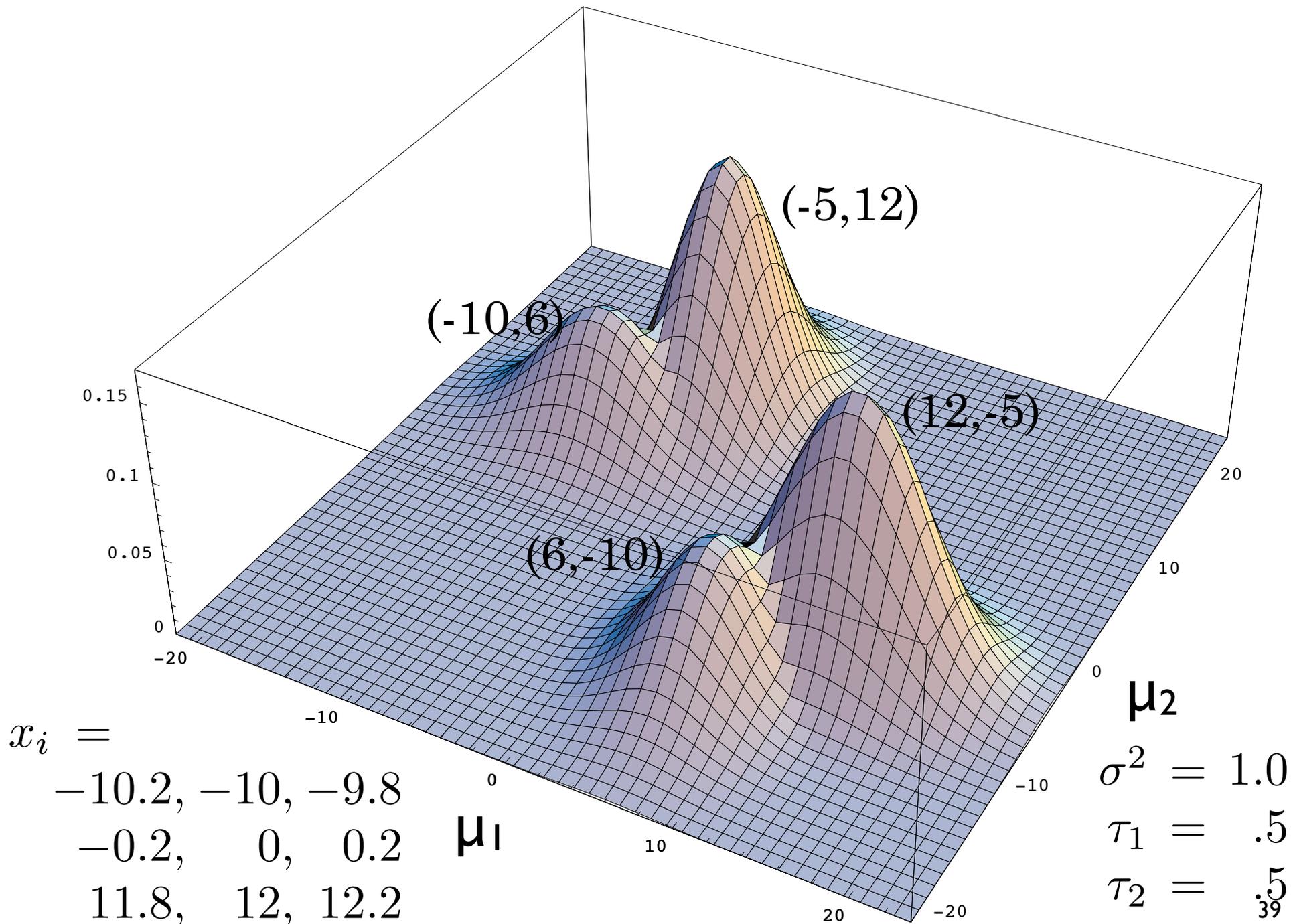
$$L(x_1, x_2, \dots, x_n | \mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \tau_1, \tau_2)$$

$$= \prod_{i=1}^n \sum_{j=1}^2 \tau_j f(x_i | \mu_j, \sigma_j^2)$$

No  
closed-  
form  
max

# Likelihood Surface





# A What-If Puzzle

Likelihood

$$L(x_1, x_2, \dots, x_n \mid \overbrace{\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \tau_1, \tau_2}^{\theta})$$
$$= \prod_{i=1}^n \sum_{j=1}^2 \tau_j f(x_i \mid \mu_j, \sigma_j^2)$$

Messy: no closed form solution known for finding  $\theta$  maximizing L

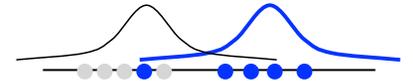
But *what if* we knew the *hidden data*?

$$z_{ij} = \begin{cases} 1 & \text{if } x_i \text{ drawn from } f_j \\ 0 & \text{otherwise} \end{cases}$$

# EM as Egg vs Chicken

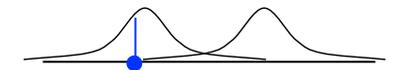
*IF*  $z_{ij}$  known, could estimate parameters  $\theta$

E.g., only points in cluster 2 influence  $\mu_2, \sigma_2$



*IF* parameters  $\theta$  known, could estimate  $z_{ij}$

E.g.,  $|x_i - \mu_1|/\sigma_1 \ll |x_i - \mu_2|/\sigma_2 \Rightarrow P[z_{i1}=1] \gg P[z_{i2}=1]$



**But we know neither;** (optimistically) iterate:

E: calculate expected  $z_{ij}$ , given parameters

M: calc “MLE” of parameters, given  $E(z_{ij})$

Overall, a clever “hill-climbing” strategy

Not what's needed for  
homework, but may  
help clarify concepts

# Simple Version: “Classification EM”

If  $E[z_{ij}] < .5$ , pretend  $z_{ij} = 0$ ;  $E[z_{ij}] > .5$ , pretend it's 1

I.e., *classify* points as component 0 or 1

Now recalc  $\theta$ , assuming that partition (standard MLE)

Then recalc  $E[z_{ij}]$ , assuming that  $\theta$

Then re-recalc  $\theta$ , assuming new  $E[z_{ij}]$ , etc., etc.

“Full EM” is a bit more involved, (to account for uncertainty in classification) but this is the crux.

# Full EM

$x_i$ 's are known;  $\theta$  unknown. Goal is to find MLE  $\theta$  of:

$$L(x_1, \dots, x_n \mid \theta) \quad \text{(hidden data likelihood)}$$

Would be easy *if*  $z_{ij}$ 's were known, i.e., consider:

$$L(x_1, \dots, x_n, z_{11}, z_{12}, \dots, z_{n2} \mid \theta) \quad \text{(complete data likelihood)}$$

But  $z_{ij}$ 's aren't known.

Instead, maximize *expected* likelihood of visible data

$$E(L(x_1, \dots, x_n, z_{11}, z_{12}, \dots, z_{n2} \mid \theta)),$$

where expectation is over distribution of hidden data ( $z_{ij}$ 's)

# The E-step:

Find  $E(z_{ij})$ , i.e.,  $P(z_{ij}=1)$

Assume  $\theta$  known & fixed

A (B): the event that  $x_i$  was drawn from  $f_1$  ( $f_2$ )

D: the observed datum  $x_i$

Expected value of  $z_{i1}$  is  $P(A|D)$  —  $E = 0 \cdot P(0) + 1 \cdot P(1)$

$$P(A|D) = \frac{P(D|A)P(A)}{P(D)}$$

$$P(D) = P(D|A)P(A) + P(D|B)P(B)$$

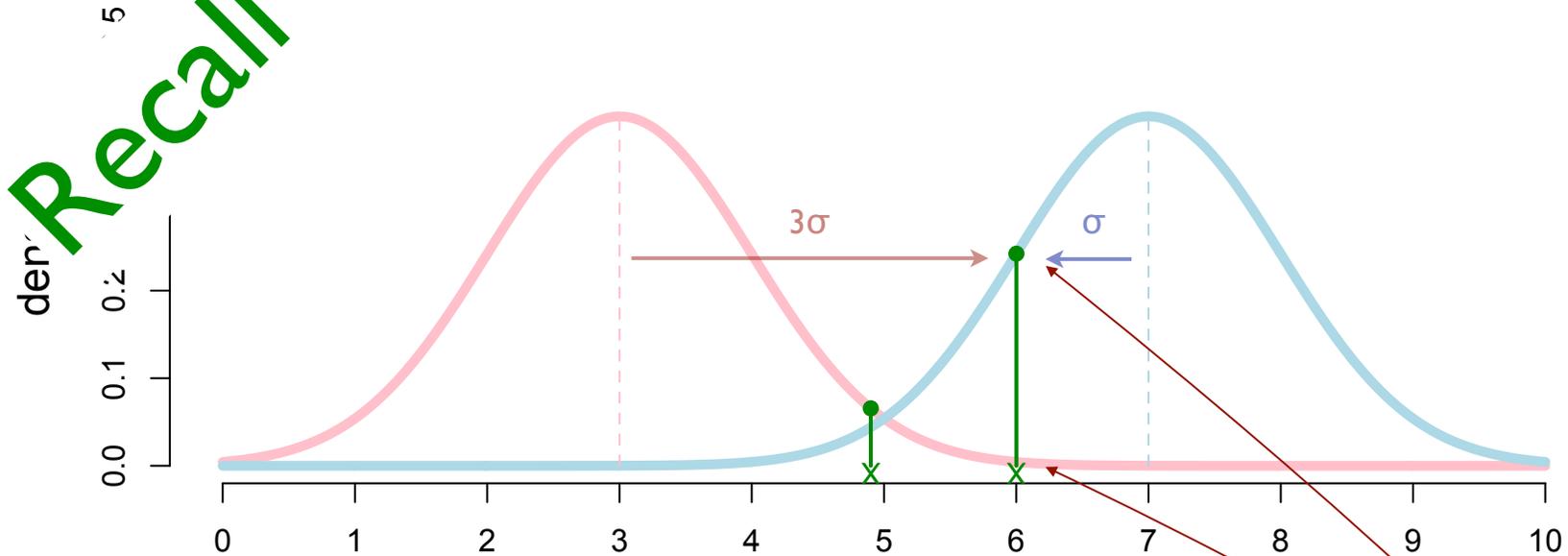
$$= f_1(x_i|\theta_1) \tau_1 + f_2(x_i|\theta_2) \tau_2$$

Repeat  
for  
each  
 $x_i$

Note: denominator = sum of numerators - i.e. that which normalizes sum to 1 (typical Bayes)

# A Hat Trick

der  
Recall



Let “ $X \approx 6$ ” be a shorthand for  $6.001 - \delta/2 < X < 6.001 + \delta/2$

$$P(\mu = 7|X = 6) = \lim_{\delta \rightarrow 0} P(\mu = 7|X \approx 6)$$

$$P(\mu = 7|X \approx 6) = \frac{P(X \approx 6|\mu = 7)P(\mu = 7)}{P(X \approx 6)}$$

$$= \frac{0.5P(X \approx 6|\mu = 7)}{0.5P(X \approx 6|\mu = 3) + 0.5P(X \approx 6|\mu = 7)}$$

$$\approx \frac{f(X = 6|\mu = 7)\delta}{f(X = 6|\mu = 3)\delta + f(X = 6|\mu = 7)\delta}, \text{ so}$$

$$P(\mu = 7|X = 6) = \frac{f(X = 6|\mu = 7)}{f(X = 6|\mu = 3) + f(X = 6|\mu = 7)} \approx 0.982$$

$f$  = normal density

# Complete Data Likelihood

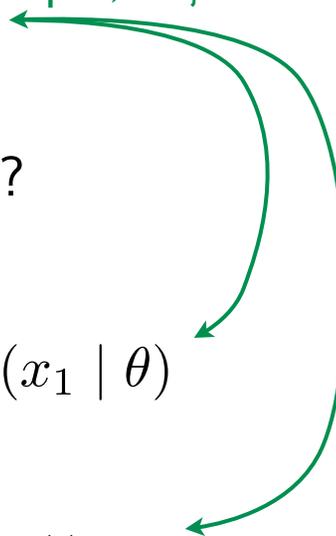
Recall:

$$z_{1j} = \begin{cases} 1 & \text{if } x_1 \text{ drawn from } f_j \\ 0 & \text{otherwise} \end{cases}$$

so, correspondingly,

$$L(x_1, z_{1j} | \theta) = \begin{cases} \tau_1 f_1(x_1 | \theta) & \text{if } z_{11} = 1 \\ \tau_2 f_2(x_1 | \theta) & \text{otherwise} \end{cases}$$

equal, if  $z_{ij}$  are 0/1



Formulas with “if’s” are messy; can we blend more smoothly?

Yes, many possibilities. Idea 1:

$$L(x_1, z_{1j} | \theta) = z_{11} \cdot \tau_1 f_1(x_1 | \theta) + z_{12} \cdot \tau_2 f_2(x_1 | \theta)$$

Idea 2 (Better):

$$L(x_1, z_{1j} | \theta) = (\tau_1 f_1(x_1 | \theta))^{z_{11}} \cdot (\tau_2 f_2(x_1 | \theta))^{z_{12}}$$

# M-step:

Find  $\theta$  maximizing  $E(\log(\text{Likelihood}))$

(For simplicity, assume  $\sigma_1 = \sigma_2 = \sigma; \tau_1 = \tau_2 = .5 = \tau$ )

$$L(\vec{x}, \vec{z} | \theta) = \prod_{1 \leq i \leq n} \left( \frac{\tau}{\sqrt{2\pi\sigma^2}} \exp \left( - \sum_{1 \leq j \leq 2} z_{ij} \frac{(x_i - \mu_j)^2}{2\sigma^2} \right) \right)$$

$$E[\log L(\vec{x}, \vec{z} | \theta)] = E \left[ \sum_{1 \leq i \leq n} \left( \log \tau - \frac{1}{2} \log 2\pi\sigma^2 - \sum_{1 \leq j \leq 2} z_{ij} \frac{(x_i - \mu_j)^2}{2\sigma^2} \right) \right]$$

wrt dist of  $z_{ij}$

$$= \sum_{1 \leq i \leq n} \left( \log \tau - \frac{1}{2} \log 2\pi\sigma^2 - \sum_{1 \leq j \leq 2} E[z_{ij}] \frac{(x_i - \mu_j)^2}{2\sigma^2} \right)$$

Find  $\theta$  maximizing this as before, using  $E[z_{ij}]$  found in E-step. Result:

$$\boxed{\mu_j = \sum_{i=1}^n E[z_{ij}] x_i / \sum_{i=1}^n E[z_{ij}]} \quad (\text{intuit: avg, weighted by subpop prob})$$

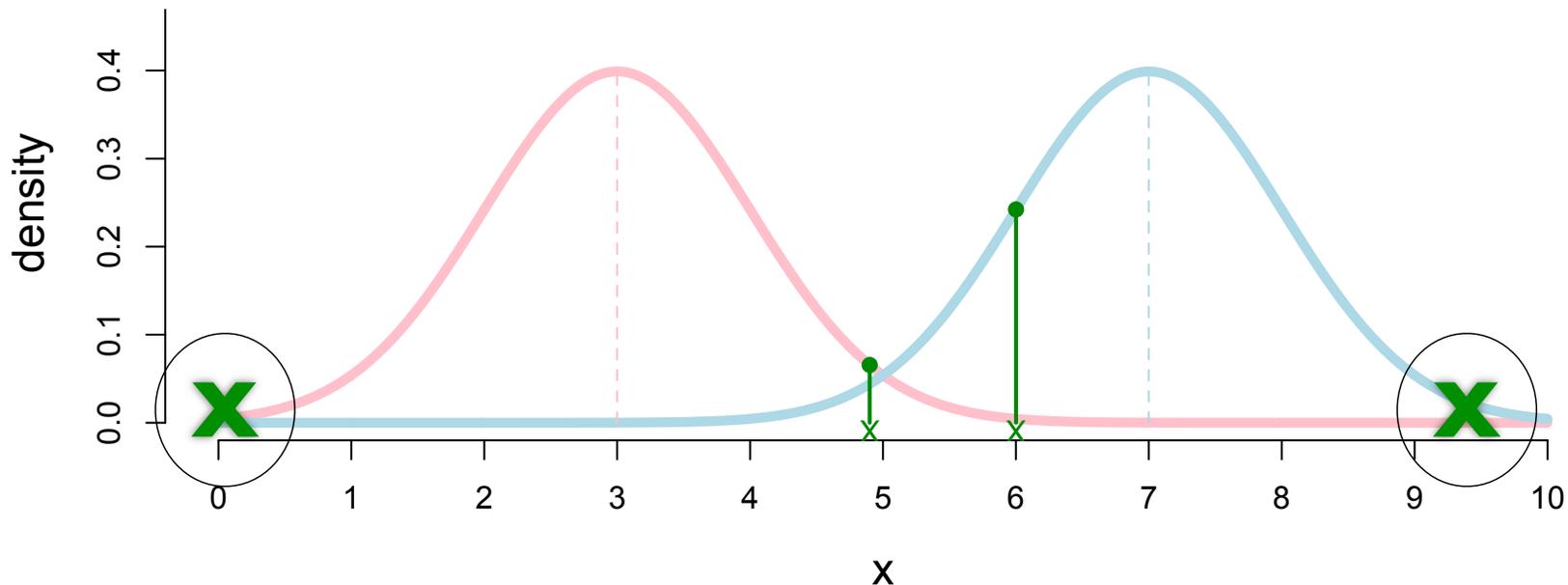
Recall

## Hat Trick (cont.)

Note 2: red/blue separation is just like the M-step of EM if values of the hidden variables ( $z_{ij}$ ) were known.

What if they're not? E.g., what would you do if some of the slips you pulled had coffee spilled on them, obscuring color?

If they were half way between means of the others?  
If they were on opposite sides of the means of the others



# 2 Component Mixture

$$\sigma_1 = \sigma_2 = 1; \tau = 0.5$$

		<b>mu1</b>	-20.00		-6.00		-5.00		-4.99
		<b>mu2</b>	6.00		0.00		3.75		3.75
<b>x1</b>	<b>-6</b>	<b>z11</b>		5.11E-12		1.00E+00		1.00E+00	
<b>x2</b>	<b>-5</b>	<b>z21</b>		2.61E-23		1.00E+00		1.00E+00	
<b>x3</b>	<b>-4</b>	<b>z31</b>		1.33E-34		9.98E-01		1.00E+00	
<b>x4</b>	<b>0</b>	<b>z41</b>		9.09E-80		1.52E-08		4.11E-03	
<b>x5</b>	<b>4</b>	<b>z51</b>		6.19E-125		5.75E-19		2.64E-18	
<b>x6</b>	<b>5</b>	<b>z61</b>		3.16E-136		1.43E-21		4.20E-22	
<b>x7</b>	<b>6</b>	<b>z71</b>		1.62E-147		3.53E-24		6.69E-26	

Essentially converged in 2 iterations

(Excel spreadsheet on course web)

# Applications

Clustering is a remarkably successful exploratory data analysis tool

Web-search, information retrieval, gene-expression, ...

Model-based approach above is one of the leading ways to do it

Gaussian mixture models widely used

With many components, empirically match arbitrary distribution

Often well-justified, due to “hidden parameters” driving the visible data

EM is extremely widely used for “hidden-data” problems

Hidden Markov Models – speech recognition, DNA analysis, ...

# EM Summary

Fundamentally a maximum likelihood parameter estimation problem

Useful if 0/1 hidden data, and if analysis would be more tractable if hidden data  $z$  were known

Iterate:

E-step: estimate  $E(z)$  for each  $z$ , given  $\theta$

M-step: estimate  $\theta$  maximizing  $E[\log \text{likelihood}]$

given  $E[z]$  [where “ $E[\log L]$ ” is wrt random  $z \sim E[z] = p(z=1)$ ]

# EM Issues

Under mild assumptions, EM is guaranteed to increase likelihood with every E-M iteration, hence will *converge*.

*But* it may converge to a *local*, not global, max.  
(Recall the 4-bump surface...)

Issue is intrinsic (probably), since EM is often applied to problems (including clustering, above) that are *NP-hard* (so fast alg is unlikely)

Nevertheless, widely used, often effective