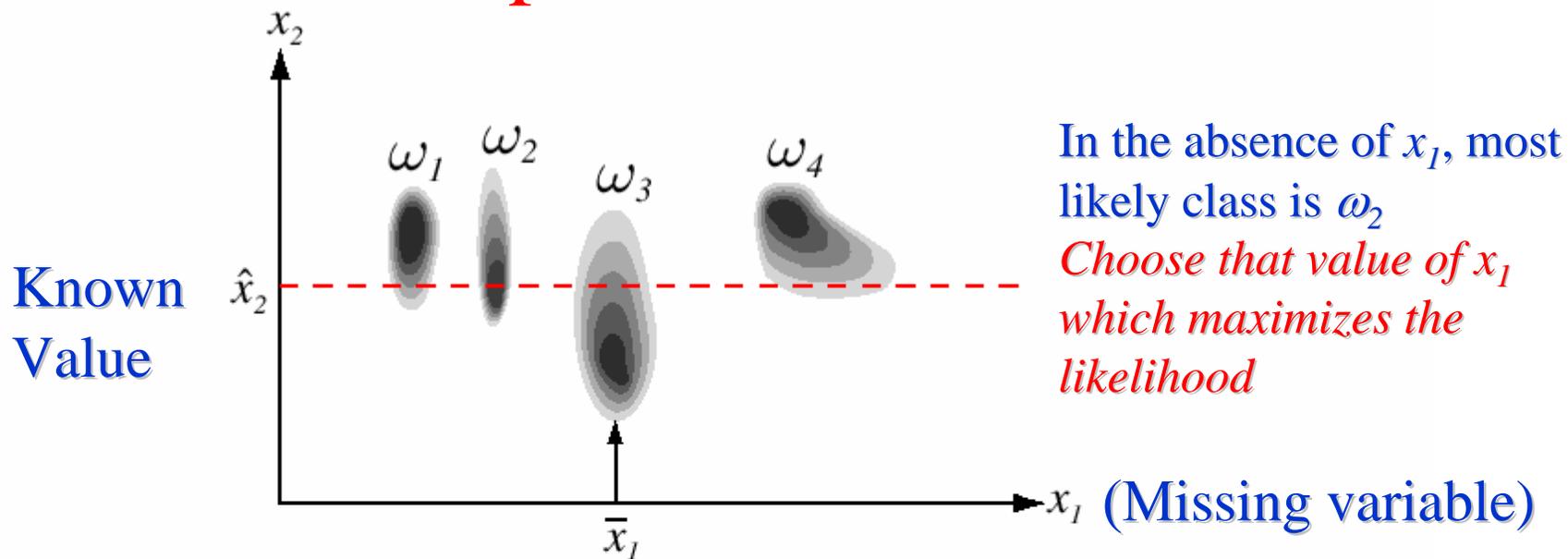


Maximum-Likelihood and Bayesian Parameter Estimation

Expectation Maximization (EM)

Estimating Missing Feature Value

Estimating missing variable with known parameters



Choosing mean of missing feature (over all classes) will result in worse performance!

This is a case of estimating the hidden variable given the parameters.

In EM unknowns are both: parameters and hidden variables

EM Task

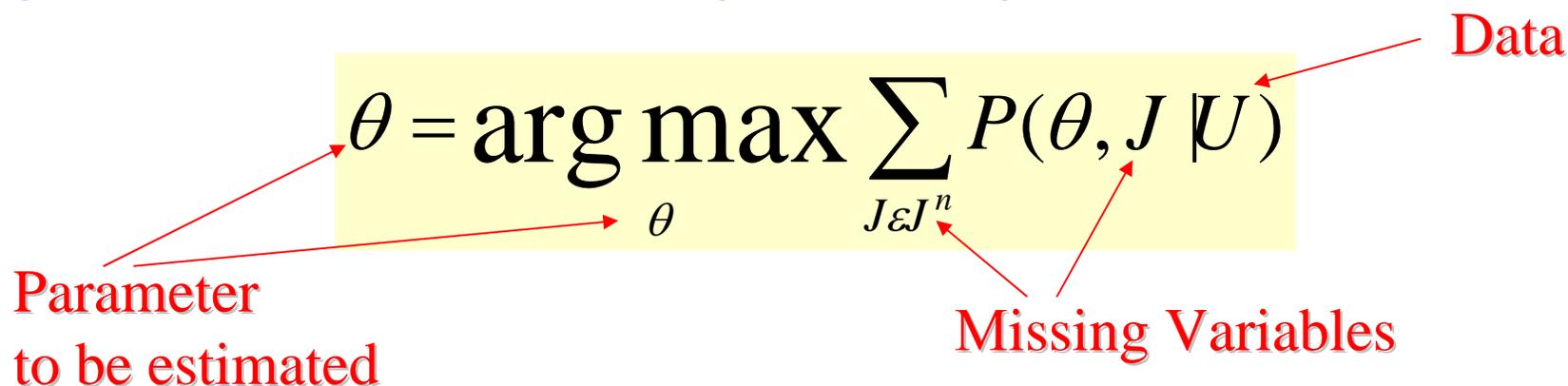
- Estimate unknown parameters θ given measurement data U
- However some variables J are missing which need to be integrated out
- We want to maximize the posterior probability of θ , given the data U , marginalizing over J

$$\theta = \arg \max_{\theta} \sum_{J \in J^n} P(\theta, J | U)$$

Parameter to be estimated

Missing Variables

Data



EM Principle

- Estimate unknown parameters θ given measurement data U and not given nuisance variables J which need to be integrated out

$$\theta = \arg \max_{\theta} \sum_{J \in \mathcal{J}^n} P(\theta, J | U)$$

- Alternate between estimating unknowns θ and the hidden variables J
- At each iteration, instead of finding the best $J \in \mathcal{J}$ given an estimate θ at each iteration, EM computes a distribution over the space J

k-means Algorithm as EM

- Estimate means of k classes when class labels are unknown
 - Parameters: means to be estimated
 - Hidden variables: class labels

begin initialize m_1, m_2, \dots, m_k (E-step)

do classify n samples according to nearest m_i (M-step)

 recompute m_i

until no change in m_i

return m_1, m_2, \dots, m_k

end

An iterative algorithm derivable from EM

EM Importance

- EM algorithm
 - widely used for learning in the presence of unobserved variables, e.g., missing features, class labels
 - used even for variables whose value is never directly observed, provided the general form of the pdf governing these variables is known
 - has been used to train *Radial Basis Function Networks* and *Bayesian Belief Networks*
 - basis for many unsupervised clustering algorithms
 - basis for widely used *Baum-Welch* forward-backward algorithm for *HMMs*

EM Algorithm

- Learning in the presence of unobserved variables
- Only a subset of relevant instance features might be observable
- Case of unsupervised learning or clustering
- How many classes are there?



EM Principle

- EM algorithm iteratively estimates the likelihood given the data that is present

Likelihood Formulation

- Sample points from a single distribution

$$\mathcal{D} = \{x_1, \dots, x_n\}$$

- Any sample has good and missing (bad) features

$$\mathbf{X}_k = \{\mathbf{X}_{kg}, \dots, \mathbf{X}_{kb}\}$$

- Features are divided into two sets

$$\mathcal{D} = \mathcal{D}_g \cup \mathcal{D}_b$$

Likelihood Formulation

- Central equation in EM

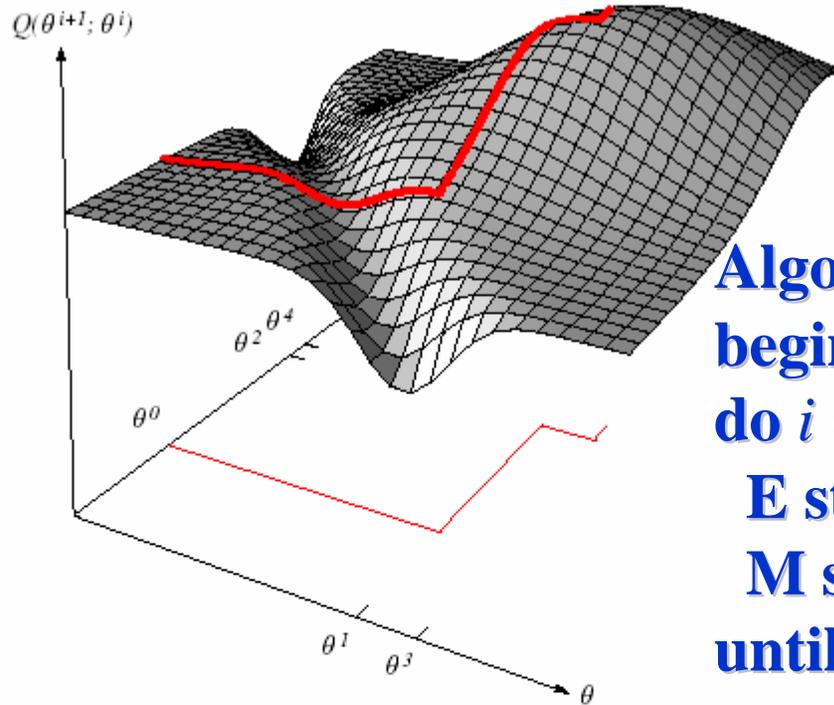
$$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^i) = \mathcal{E}_{\mathcal{D}_b} [\ln p(\mathcal{D}_g, \mathcal{D}_b; \boldsymbol{\theta}) | \mathcal{D}_g; \boldsymbol{\theta}^i],$$

Expected Value is over missing features

Current best estimate for the full distribution

Candidate Vector for an improved estimate

- Algorithm will select the best candidate $\boldsymbol{\theta}$ and call it $\boldsymbol{\theta}^{i+1}$



Algorithm EM

begin initialize $\theta^0, T, i = 0$

do $i = i + 1$

E step: compute $Q(\theta; \theta^i)$

M step: $\theta^{i+1} = \arg \max_{\theta} Q(\theta; \theta^i)$

until $Q(\theta^{i+1}; \theta^i) - Q(\theta^i; \theta^{i-1}) \leq T$

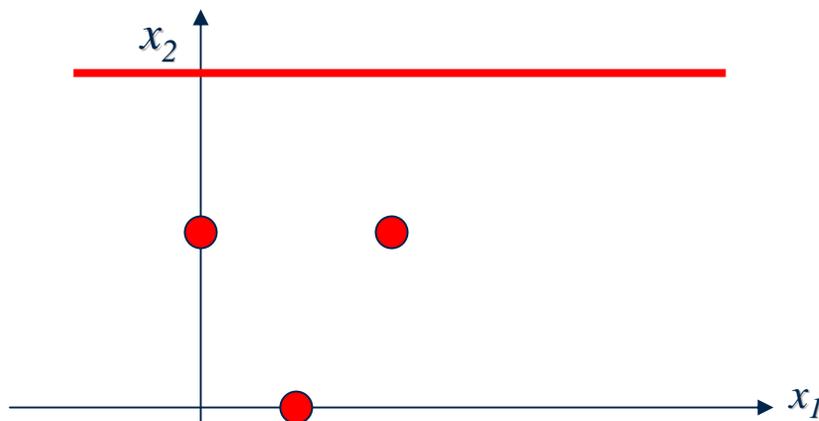
The search for the best model via the EM algorithm starts with some initial value of the model parameters, θ^0 . Then, via the **M step**, the optimal θ^1 is found. Next, θ^1 is held constant and the value θ^2 is found that optimizes $Q(\cdot; \cdot)$. This process iterates until no value of θ can be found that will increase $Q(\cdot; \cdot)$. Note in particular that this is different from a gradient search. For example here θ^1 is the global optimum (given fixed θ^0), and would not necessarily have been found via gradient search. (In this illustration, $Q(\cdot; \cdot)$ is shown symmetric in its arguments; this need not be the case in general, however.)

EM for 2D Normal Model

Suppose data consists of 4 points in 2 dimensions, one point of which is missing a feature

$$D = \{x_1, x_2, x_3, x_4\} = \left\{ \begin{pmatrix} 0 \\ 2 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \begin{pmatrix} * \\ 4 \end{pmatrix} \right\}, \text{ where } * \text{ represents}$$

the unknown value of the first feature of point x_4 . Thus our bad data D_b consists of the single feature x_{41} , and the good data D_g consists of the rest.



EM for 2D Normal Model

Assuming that the model is a Gaussian with diagonal covariance and arbitrary mean, it can be described by the parameter vector

$$\theta = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \sigma_1^2 \\ \sigma_2^2 \end{pmatrix}$$

EM for 2D Normal model

We take our initial initial guess to be a Gaussian centered on the origin having that is, $\Sigma = 1$

$$\theta^0 = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \end{pmatrix}$$

EM for a 2D Normal Model

$$\boldsymbol{\theta} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \sigma_1^2 \\ \sigma_2^2 \end{pmatrix}.$$

$$\boldsymbol{\theta}^0 = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \end{pmatrix}.$$

To find improved estimate, must calculate

$$\begin{aligned} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^0) &= \mathcal{E}_{x_{41}} [\ln p(\mathbf{x}_g, \mathbf{x}_b; \boldsymbol{\theta}) | \boldsymbol{\theta}^0; \mathcal{D}_g] \\ &= \int_{-\infty}^{\infty} \left[\sum_{k=1}^3 \ln p(\mathbf{x}_k | \boldsymbol{\theta}) + \ln p(\mathbf{x}_4 | \boldsymbol{\theta}) \right] p(x_{41} | \boldsymbol{\theta}^0; x_{42} = 4) dx_{41} \\ &= \sum_{k=1}^3 [\ln p(\mathbf{x}_k | \boldsymbol{\theta})] + \int_{-\infty}^{\infty} \ln p \left(\begin{pmatrix} x_{41} \\ 4 \end{pmatrix} \middle| \boldsymbol{\theta} \right) \underbrace{\frac{p \left(\begin{pmatrix} x_{41} \\ 4 \end{pmatrix} | \boldsymbol{\theta}^0 \right)}{\left(\int_{-\infty}^{\infty} p \left(\begin{pmatrix} x'_{41} \\ 4 \end{pmatrix} | \boldsymbol{\theta}^0 \right) dx'_{41} \right)}}_{\equiv K} dx_{41}, \end{aligned}$$

EM for a 2D Normal Model

Simplifying

$$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^0) = \sum_{k=1}^3 [\ln p(\mathbf{x}_k | \boldsymbol{\theta})] + \frac{1}{K} \int_{-\infty}^{\infty} \ln p \left(\begin{pmatrix} x_{41} \\ 4 \end{pmatrix} \middle| \boldsymbol{\theta} \right) \frac{1}{2\pi \begin{vmatrix} 1 & 0 \\ 0 & 1 \end{vmatrix}} \exp \left[-\frac{1}{2} (x_{41}^2 + 4^2) \right] dx_{41}$$
$$= \sum_{k=1}^3 [\ln p(\mathbf{x}_k | \boldsymbol{\theta})] - \frac{1 + \mu_1^2}{2\sigma_1^2} - \frac{(4 - \mu_2)^2}{2\sigma_2^2} - \ln(2\pi\sigma_1\sigma_2).$$

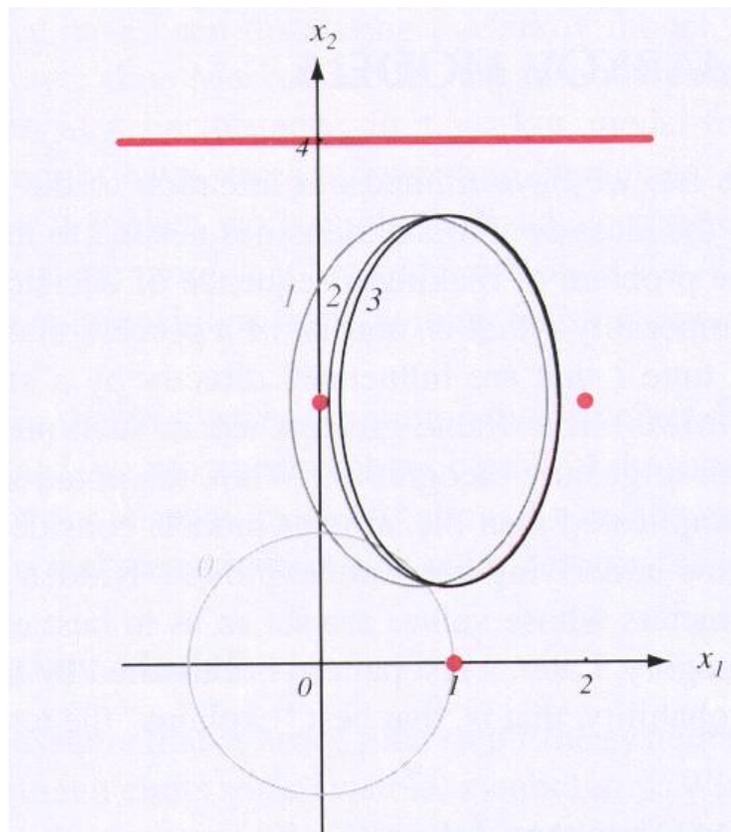
Completes E step. Gives the next estimate

$$\boldsymbol{\theta}^1 = \begin{pmatrix} 0.75 \\ 2.0 \\ 0.938 \\ 2.0 \end{pmatrix}.$$

Final Solution

$$\boldsymbol{\theta} = \begin{pmatrix} 1.0 \\ 2.0 \\ 0.667 \\ 2.0 \end{pmatrix}$$

EM for 2D normal model



Four data points, one missing the value of x_1 , are in red.

Initial estimate is a circularly symmetric Gaussian, centered on the origin (gray).

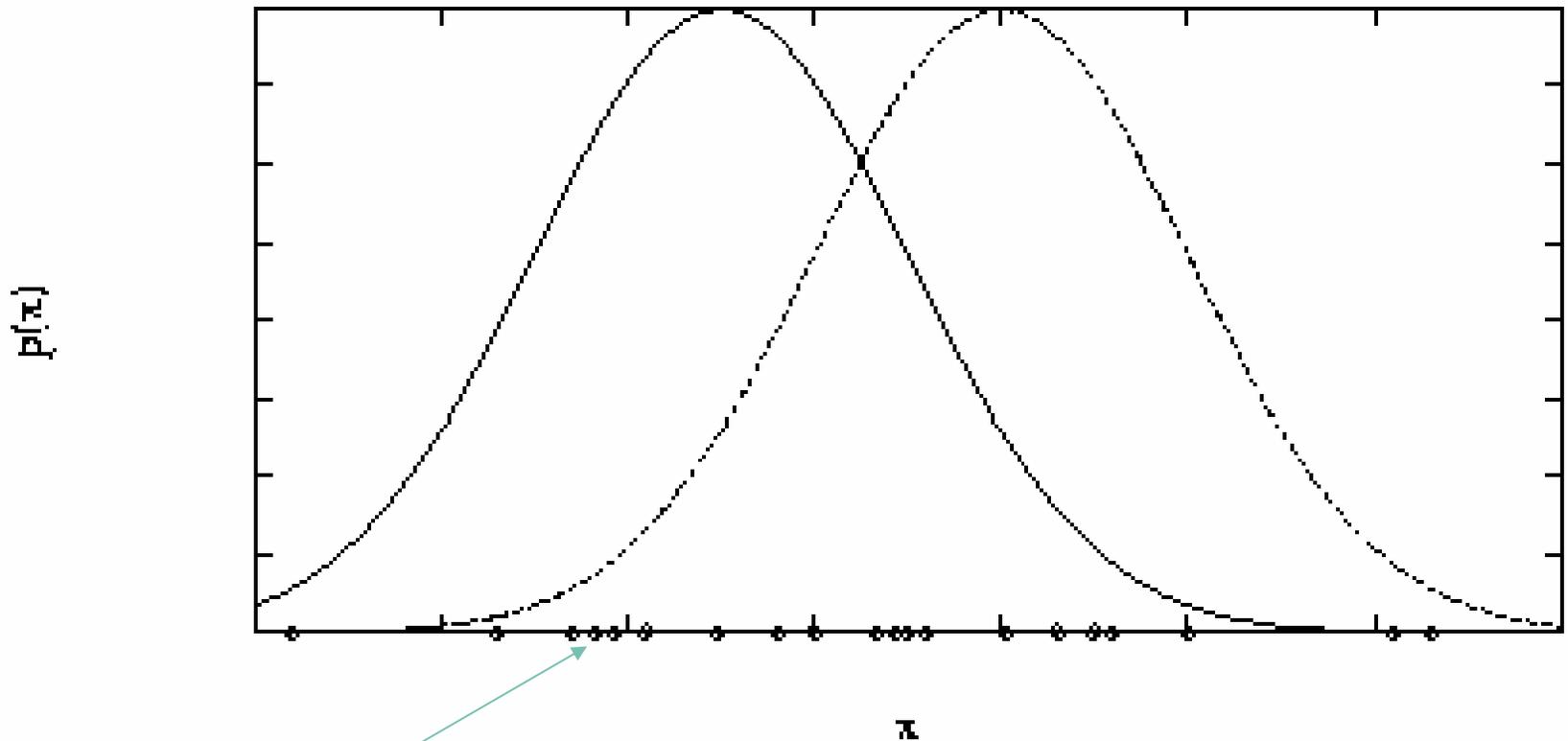
(A better initial estimate could have been derived from the 3 known points.)

Each iteration leads to an improved estimate, labeled by iteration number i ; after 3 iterations, the algorithm converged.

EM to Estimate Means of k Gaussians

- Data D drawn from mixture of k distinct normal distributions
- Two-step process generates samples
 - One of the k distributions is selected at random
 - A single random instance x_i is generated according to selected distribution

Instances Generated by a Mixture of Two Normal Distributions



instances

Example of EM to Estimate Means of k Gaussians

- Each instance is generated by
 - 1 Choosing one of the k Gaussians with uniform probability
 - 2 Generating an instance at random according to that Gaussian
- The single normal distribution is chosen with uniform probability
- Each of the k Normal distributions has the same known variance
- Learning task: output a hypothesis that describes the means of each of $h = \langle \mu_1, \dots, \mu_k \rangle$

Estimating Means of k Gaussians

- We would like to find a maximum likelihood hypothesis for these means: a hypothesis h that maximizes $p(D/h)$

Maximum Likelihood Estimate of Mean of *a* Single Gaussian

- Given observed data instances

$$x_1, x_2, \dots, x_m$$

- Drawn from a single distribution that is normally distributed
- Problem is to find the mean of that distribution

Maximum Likelihood Estimate of Mean of a *Single Gaussian*

- Maximum likelihood estimate of the mean of a normal distribution can be shown to be one that minimizes the sum of squared errors

$$\mu_{\text{ML}} = \arg \min_{\mu} \frac{1}{2} \sum_{i=1}^m (x_i - \mu)^2$$

- Right hand side has a maximum value at

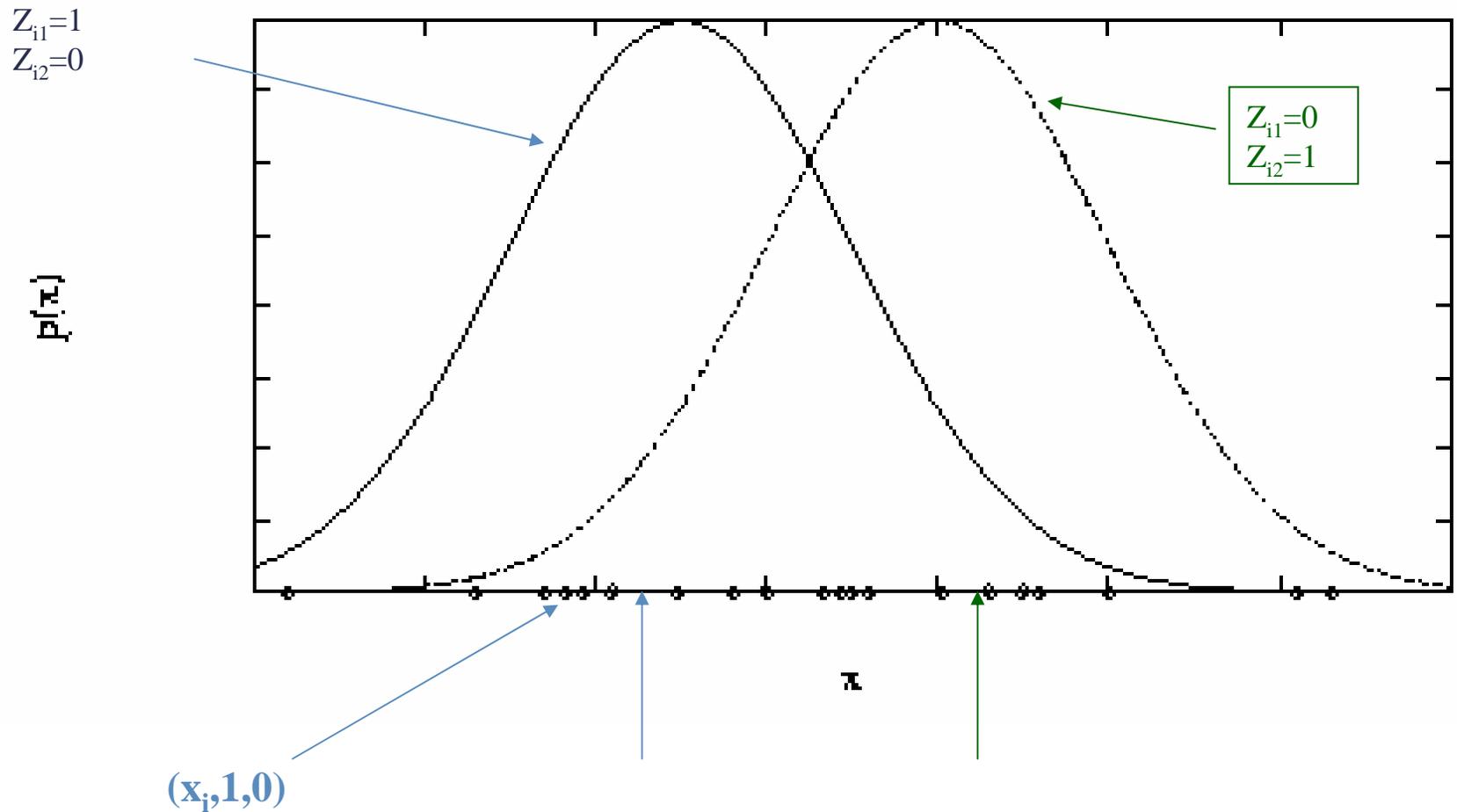
$$\mu_{\text{ML}} = \frac{1}{m} \sum_{i=1}^m x_i$$

- which is the sample mean

Mixture of Two Normal Distributions

- We cannot observe as to which instances were generated by which distribution
- Full description of instance $\langle x_i, z_{i1}, z_{i2} \rangle$
- x_i = observed value of i^{th} instance
- z_{i1} and z_{i2} indicate which of two normal distributions was used to generate x_i
- $z_{ij} = 1$ if z_{ij} was used to generate x_i , 0 otherwise
- z_{i1} and z_{i2} are hidden variables, which have probability distributions associated with them

Hidden variables specify distribution



2-Means Problem

- Full description of instance

$$\langle x_i, z_{i1}, z_{i2} \rangle$$

- x_i = observed variable
- z_{i1} and z_{i2} are hidden variables
- If z_{i1} and z_{i2} were observed, we could use maximum likelihood estimates for the means:

$$\mu_1 = \frac{1}{m} \sum_{\substack{i=1 \\ \exists z_{i1}=1}}^m x_i$$

$$\mu_2 = \frac{1}{m} \sum_{\substack{i=1 \\ \exists z_{i2}=1}}^m x_i$$

- Since we do not know z_{i1} and z_{i2} , we will use EM instead

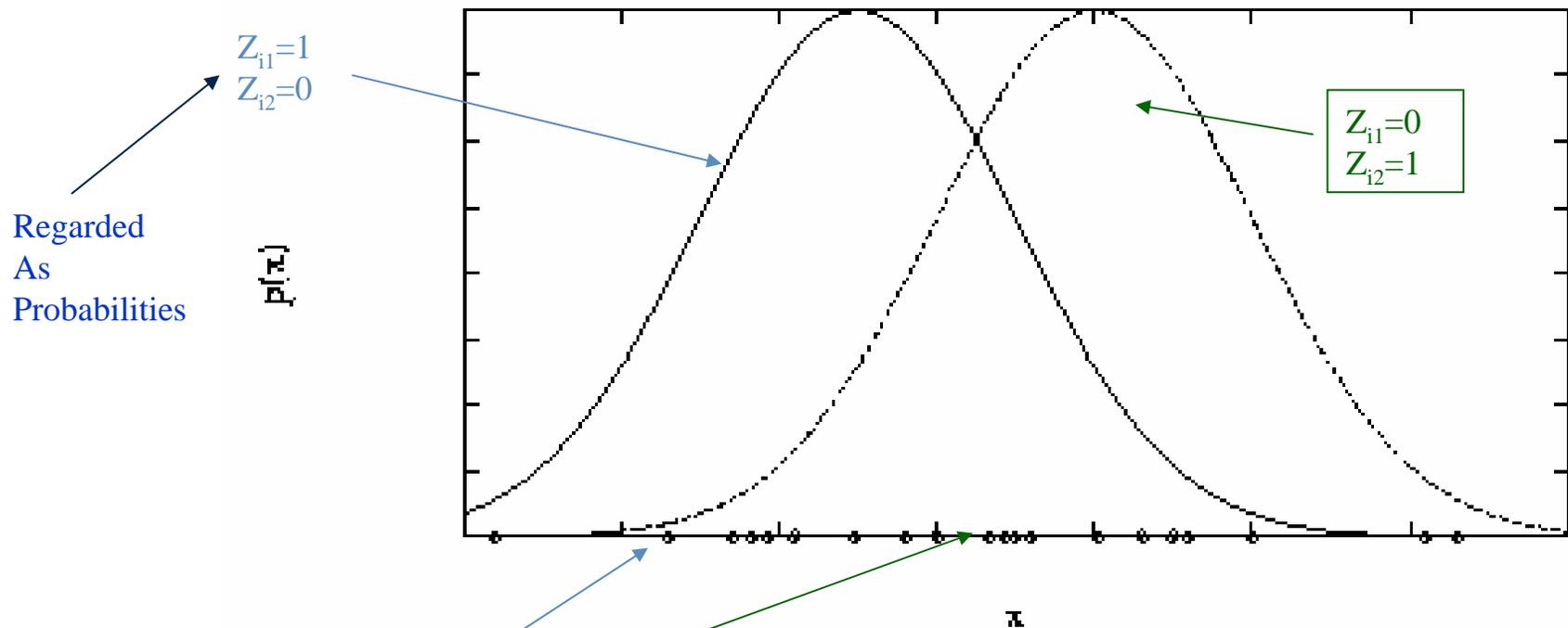
EM Algorithm Applied to k -Means Problem

- Search for a Maximum Likelihood Hypothesis by repeatedly re-estimating expected values of hidden binary variables z_{ij} given its current hypothesis

$$\langle \mu_1, \dots, \mu_k \rangle$$

- Then recalculate the maximum likelihood hypothesis using these expected values for the hidden variables

EM algorithm for two means



1. Hypothesize means, then determine expected value of hidden variables for all samples
2. Use these hidden variable values to recalculate the means

EM Applied to Two-Means Problem

- Initialize the hypothesis to

$$h = \langle \mu_1, \mu_2 \rangle$$

- Estimate the expected values of hidden variables z_{ij} given its current hypothesis
- Recalculate the maximum likelihood hypothesis using these expected values for the hidden variables
- Re-estimate h repeatedly until the procedure converges to a stationary value of h

EM Algorithm for 2-Means

Step 1

Calculate the expected value $E[z_{ij}]$ of each hidden variable z_{ij} , assuming the current hypothesis

$$h = \langle \mu_1, \mu_2 \rangle$$

Step 2

Calculate new maximum likelihood hypothesis assuming the value taken on by each hidden variable expected value $E[z_{ij}]$ calculated in Step 1.

$$h' = \langle \mu'_1, \mu'_2 \rangle$$

Then replace hypothesis by the new hypothesis and iterate.

$$h = \langle \mu_1, \mu_2 \rangle$$

$$h' = \langle \mu'_1, \mu'_2 \rangle$$

EM First Step

Calculate the expected value $E[z_{ij}]$ of each hidden variable z_{ij} , assuming the current hypothesis $h = \langle \mu_1, \mu_2 \rangle$

$$E[z_{ij}] = \frac{p(x = x_i | \mu = \mu_j)}{\sum_{n=1}^2 p(x = x_i | \mu = \mu_n)}$$

= Probability that instance x_i was generated by the j^{th} Gaussian

$$= \frac{e^{-\frac{1}{2\sigma^2}(x_i - \mu_j)^2}}{\sum_{n=1}^2 e^{-\frac{1}{2\sigma^2}(x_i - \mu_n)^2}}$$

EM Second Step

Calculate a new maximum likelihood hypothesis

$$h' = \langle \mu'_1, \mu'_2 \rangle$$

$$\mu_j \leftarrow \frac{\sum_{i=1}^m E[z_{ij}] x_i}{\sum_{i=1}^m E[z_{ij}]}$$

Observations:

1. Similar to earlier sample mean calculation for a single Gaussian.

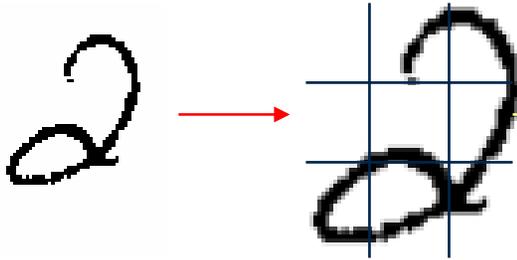
$$\mu_{\text{ML}} = \frac{1}{m} \sum_{i=1}^m x_i$$

Clustering using EM

A handwritten digit '2' in a cursive style, written in black ink on a white background.A handwritten digit '2' in a cursive style, written in black ink on a white background.A handwritten digit '2' in a simple, blocky style, written in black ink on a white background.A handwritten digit '2' in a cursive style, written in black ink on a white background.A handwritten digit '2' in a cursive style, written in black ink on a white background.

Feature Extraction

Image I_1



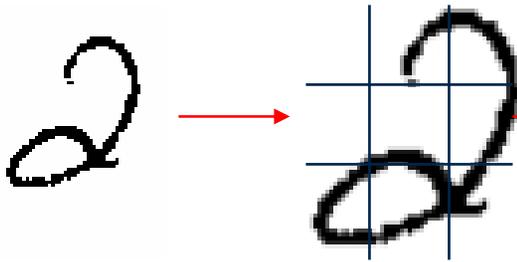
Feature
Extraction

Feature Vector (74 values)

0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000	0.367379 0.715380	0.175000 0.962500 0.688889 0.129167 0.419643 1.000000 0.387500 0.000000
1.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000	0.727273 1.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000	0.000000 0.000000 1.000000 1.000000 0.108295 0.000000 1.000000 0.599078
0.207547 0.207547 0.194969 0.779874 0.760108 0.215633 0.350943 0.630728	0.127660 0.351064 0.073286 0.879433 1.000000 0.607903 0.263830 0.474164	0.111111 0.407407 0.637860 0.765432 0.870370 0.000000 0.114815 0.687831

Feature Extraction

Image I_1



Feature
Extraction

2 Global Features: Aspect ratio
Stroke ratio

72 Local Features

$$F(i,j) = \frac{s(i,j)}{N(i)*S(j)} \quad \begin{matrix} i= 1,9 \\ j = 0,7 \end{matrix}$$

$s(i,j)$ - no. of components with slope j in subimage i

$N(i)$ - no. of components in subimage i

$$S(j) = \max_i \frac{s(i,j)}{N(i)}$$

0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000	0.367379 0.715380	0.175000 0.962500 0.688889 0.129167 0.419643 1.000000 0.387500 0.000000
1.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000	0.727273 1.000000 0.000000 0.000000 0.000000 0.000000 0.000000 0.000000	0.000000 0.000000 1.000000 1.000000 0.108295 0.000000 1.000000 0.599078
0.207547 0.207547 0.194969 0.779874 0.760108 0.215633 0.350943 0.630728	0.127660 0.351064 0.073286 0.879433 1.000000 0.607903 0.263830 0.474164	0.111111 0.407407 0.637860 0.765432 0.870370 0.000000 0.114815 0.687831

Feature Vector (74 values)

For One Feature

	$P(I_1 C_1)$	$P(I_1 C_2)$	$P(I_1 C_3)$	$P(I_1 C_4)$	$P(I_1 C_5)$
Cycle 1	1.000000	0.000000	0.000000	0.000000	0.000000
Cycle 2	1.000000	0.000000	0.000000	0.000000	0.000000
...					
Cycle N	0.990915	0.009083	0.000001	0.000000	0.000000

Final Clusters Centers (Images closer to the means)



$P(C_1) = 0.206088$



$P(C_2) = 0.228980$



$P(C_3) = 0.203198$

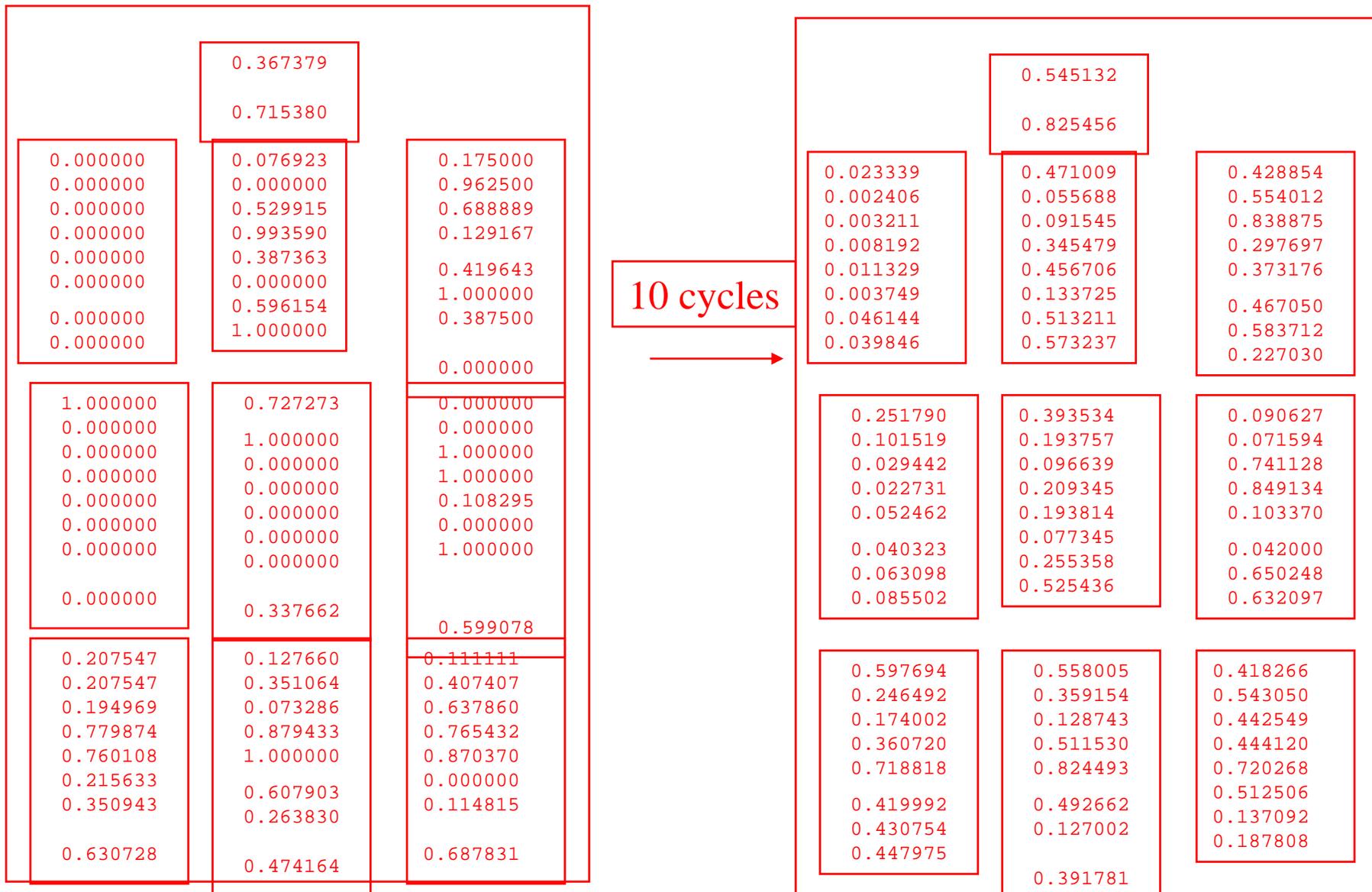


$P(C_3) = 0.218480$



$P(C_3) = 0.143254$

Clustering



Initial Mean for Cluster 1
(Initialized with FV for Image 1)

CSE 555: Srihari

Final Mean for Cluster 1

Essence of EM Algorithm

- The current hypothesis is used to estimate the unobserved variables
- Expected values of these variables are used to calculate an improved hypothesis
- It can be shown that on each iteration through the loop EM increases the likelihood $P(D/h)$ unless it is at a local maximum
- Algorithm converges on a local maximum likelihood hypothesis for

$$\langle \mu_1, \mu_2 \rangle$$

General Statement of EM Algorithm

Parameters of interest were $\theta = \langle \mu_1, \mu_2 \rangle$

Full data were the triples $\langle x_i, z_{i1}, z_{i2} \rangle$, of which only x_i is observed

General Statement of EM Algorithm (cont'd.)

- Let $X = \{x_1, x_2, \dots, x_m\}$
 - denote the observed data in a set of m independently drawn instances

- Let $Z = \{z_1, z_2, \dots, z_m\}$

- denote the unobserved data in these same instances
- Let

$$Y = X \cup Z$$

- denote the full data

General Statement of EM Algorithm (cont'd.)

- The unobserved Z can be treated as a random variable whose p.d.f. depends on the unknown parameters θ and on the observed data X
- Similarly Y is a random variable because it is defined in terms of the random variable Z
- h denotes current hypothesized values of the parameters θ
- h' denotes revised hypothesis estimated on each iteration of EM algorithm

General Statement of EM Algorithm (cont'd.)

- EM algorithm searches for m.l.e. hypothesis by seeking the h' that maximizes h'

$$E[\ln P(Y | h')]$$

General Statement of EM Algorithm (cont'd.)

- Define the function $Q(h' | h)$ that gives

$$E[\ln P(Y | h')]$$

- as a function of h' under the assumption that $\theta = h$

$$Q(h' | h) \leftarrow E[\ln P(Y | h') | h, X]$$

- And given the observed portion X of the full data Y

General Statement of EM Algorithm

- Repeat until convergence
 - Step 1: Estimation (E) Step: Calculate $Q(h' | h)$ using the current hypothesis h and the observed data X to estimate the probability distribution over Y

$$Q(h' | h) \leftarrow E[\ln P(Y | h') | h, X]$$

- Step 2: Maximization (M) Step: Replace hypothesis h by the hypothesis h' that maximizes this Q function

$$h \leftarrow \arg \max_{h'} Q(h' | h)$$

Derivation of k Means Algorithm from General EM algorithm

- Derive previously seen algorithm for estimating the means of a mixture of k Normal distributions.
Estimate the parameters

$$\theta = \langle \mu_1, \dots, \mu_k \rangle$$

- We are given the observed data

$$X = \{\langle x_i \rangle\}$$

- The hidden variables

$$Z = \{\langle z_{i1}, \dots, z_{ik} \rangle\}$$

- Indicate which of the k Normal distributions was used to generate x_i

Derivation of the k Means Algorithm from General EM Algorithm (cont'd.)

- Need to derive an expression for

$$Q(h/h')$$

- First we derive an expression for
 $\ln P(Y | h')$

Derivation of the k Means Algorithm

- The probability $p(y_i | h')$ of a single instance of the full data can be written

$$y_i = \langle x_1, z_{i1}, \dots, z_{ik} \rangle$$

$$p(y_i | h') = p(x_i, z_{i1}, \dots, z_{ik} | h') = \frac{1}{\sqrt{2\pi \sigma^2}} e^{-\frac{1}{2\sigma^2} \sum_{j=1}^k z_{ij} (x_1 - \mu'_j)^2}$$

Derivation of the k Means Algorithm (cont'd.)

- Given this probability for a single instance $p(y_i | h')$ logarithmic probability $\ln P(Y | h')$ for all m instances in the data is

$$\ln P(Y | h') = \ln \prod_{i=1}^m p(y_i | h')$$

$$= \sum_{i=1}^m \ln p(y_i | h')$$

$$= \sum_{i=1}^m \left(\ln \frac{1}{\sqrt{2\pi \sigma^2}} - \frac{1}{2\sigma^2} \sum_{j=1}^k z_{ij} (x_i - \mu'_j)^2 \right)$$

Derivation of the k Means Algorithm (cont'd.)

- In general, for any function $f(z)$ that is a linear function of z , the following equality holds

$$E[f(z)] = f(E[z])$$

$$E[\ln P(Y | h')] = E \left[\sum_{i=1}^m \left(\ln \frac{1}{\sqrt{2\pi} \sigma} - \frac{1}{2\sigma^2} \sum_{j=1}^k z_{ij} (x_i - \mu'_j)^2 \right) \right]$$

$$= \sum_{i=1}^m \left(\ln \frac{1}{\sqrt{2\pi} \sigma} - \frac{1}{2\sigma^2} \sum_{j=1}^k E[z_{ij}] (x_i - \mu'_j)^2 \right)$$

Derivation of the k Means Algorithm

- To summarize, where $h' = \langle \mu'_1, \dots, \mu'_k \rangle$ and $E[z_{ij}]$ calculated on the current hypothesis h and the observed data X , the function for the $Q(h' | h)$ problem is

$$Q(h' | h) = \sum_{i=1}^m \left(\ln \frac{1}{\sqrt{2\pi} \sigma} - \frac{1}{2\sigma^2} \sum_{j=1}^k E[z_{ij}] (x_i - \mu'_j)^2 \right)$$

$$E[z_{ij}] = \frac{e^{-\frac{1}{2\sigma^2}(x_i - \mu_j)^2}}{\sum_{n=1}^k e^{-\frac{1}{2\sigma^2}(x_i - \mu_n)^2}}$$

Derivation of k Means Algorithm:

Second (Maximization Step) To Find the Values

$$h' = \langle \mu'_1, \dots, \mu'_k \rangle$$

$$\arg \max_{h'} Q(h' | h) = \arg \max_{h'} \sum_{i=1}^m \left(\ln \frac{1}{\sqrt{2\pi} \sigma^2} - \frac{1}{2\sigma^2} \sum_{j=1}^k E[z_{ij}] (x_i - \mu'_j)^2 \right)$$

$$= \arg \min_{h'} \sum_{i=1}^m \sum_{j=1}^k E[z_{ij}] (x_i - \mu'_j)^2$$

$$\mu_j \leftarrow \frac{\sum_{i=1}^m E[z_{ij}] x_i}{\sum_{i=1}^m E[z_{ij}]}$$

Summary

- In many parameter estimation tasks, some of the relevant instance variables may be unobservable.
- In this case, the EM algorithm is useful.