



Rensselaer

Lecture 22: Gaussian Mixture Model and Expectation Maximization Algorithm

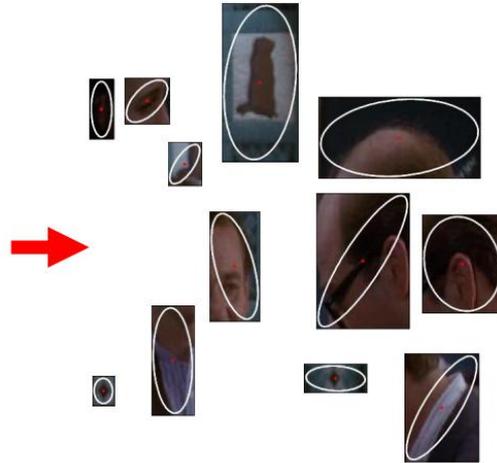
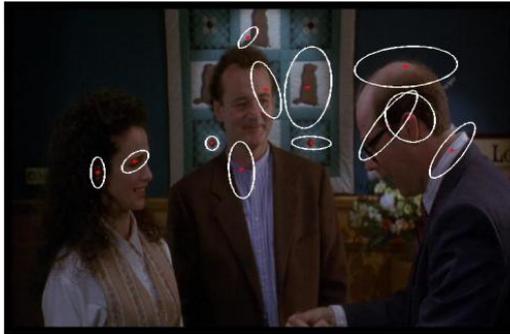
Dr. Chengjiang Long

Computer Vision Researcher at Kitware Inc.

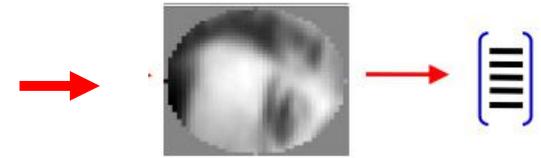
Adjunct Professor at RPI.

Email: longc3@rpi.edu

Recap Previous Lecture



Collection of visual words



Normalize patch

Compute SIFT descriptor

[Lowe'99]

Vector quantize descriptors from a set of training images using k-means

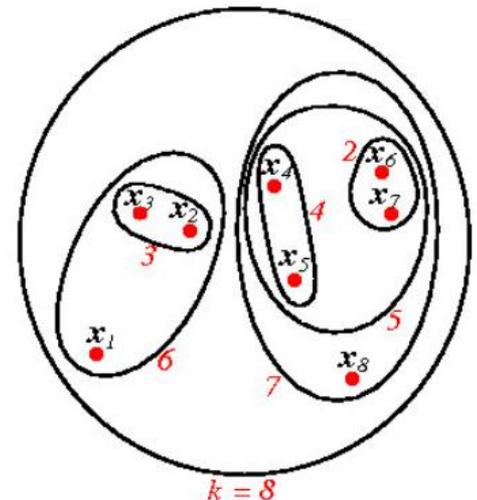
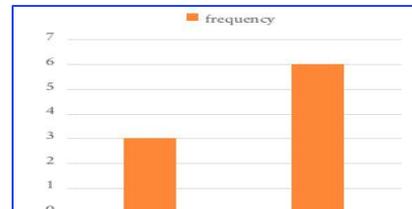
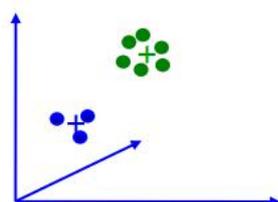
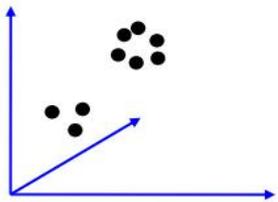


Image representation: a normalized histogram of visual words.

Outline

- Parametric Unsupervised Learning
- Mixture Density Model
- Gaussian Mixture Model
- Expectation Maximization (EM) Algorithm
- Applications

Outline

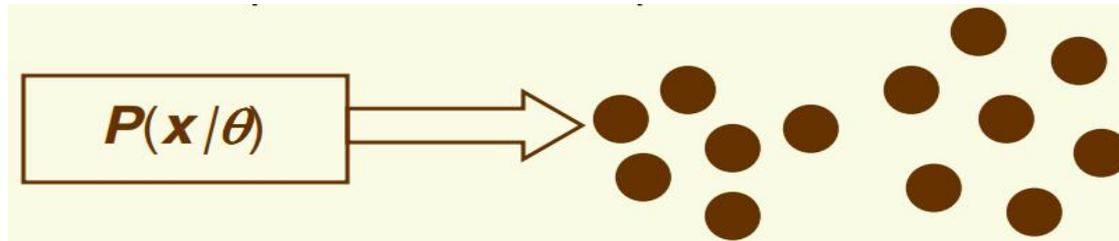
- **Parametric Unsupervised Learning**
- Mixture Density Model
- Gaussian Mixture Model
- Expectation Maximization (EM) Algorithm
- Applications

Unsupervised Learning

- In unsupervised learning, where we are only given samples x_1, \dots, x_n without class labels
- Last lecture: nonparametric approach (clustering)
- Today, parametric approach
 - assume parametric distribution of data
 - estimate parameters of this distribution
 - much “harder” than the supervised learning case

Parametric Unsupervised Learning

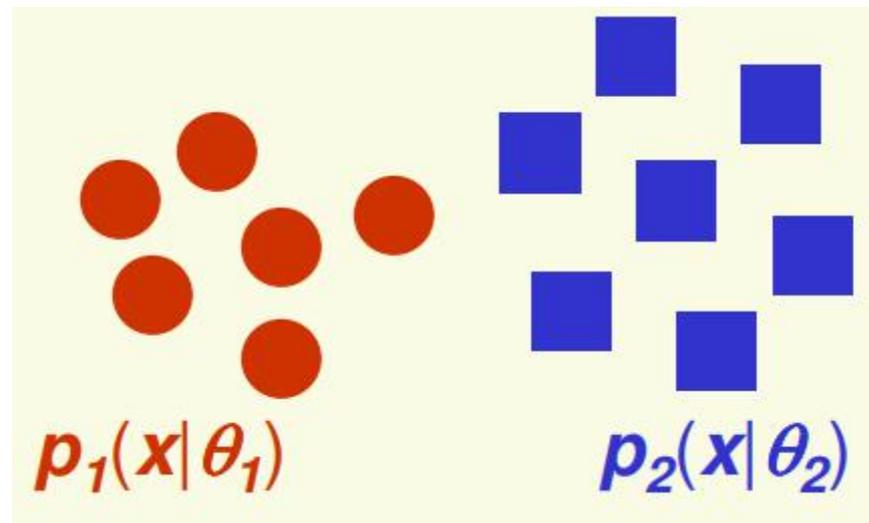
- Assume the data was generated by a model with known shape but unknown parameters



- Advantages of having a model
 - Gives a meaningful way to cluster data
 - adjust the parameters of the model to maximize the probability that the model produced the observed data
 - Can sensibly measure if a clustering is good
 - compute the likelihood of data induced by clustering
 - Can compare 2 clustering algorithms
 - which one gives the higher likelihood of the observed data?

Parametric Supervised Learning

- Let us recall supervised parametric learning
 - have m classes
 - have samples x_1, \dots, x_n each of class $1, 2, \dots, m$
 - suppose D_i holds samples from class i
 - probability distribution for class i is $p_i(x | \theta_i)$

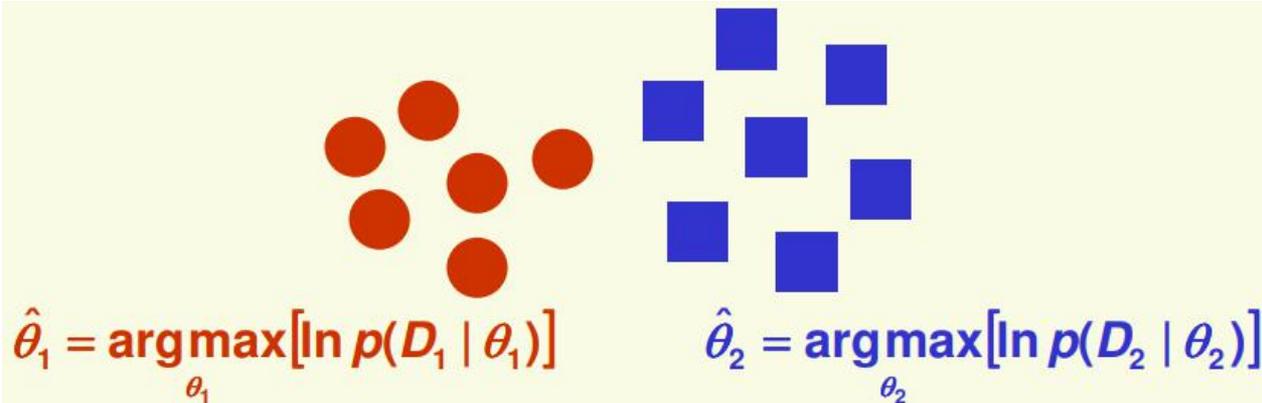


Parametric Supervised Learning

- Use the ML method to estimate parameters θ_i
 - find θ_i which maximizes the likelihood function $F(\theta_i)$
 - or, equivalently, find θ_i which maximizes the log likelihood $l(\theta_i)$

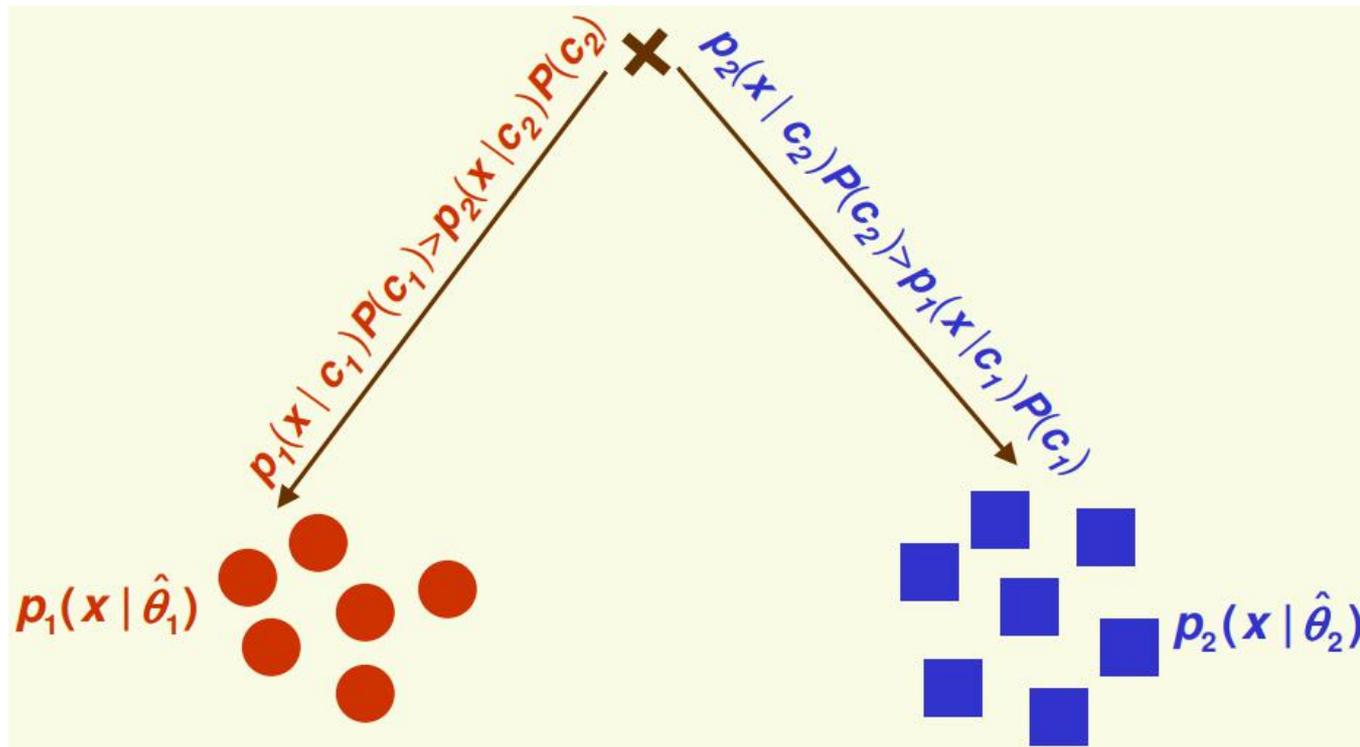
$$p(D_i | \theta_i) = \prod_{x \in D_i} p(x | \theta_i) = F(\theta_i)$$

$$l(\theta_i) = \ln p(D_i | \theta_i) = \sum_{x \in D_i} \ln p(x | \theta_i)$$



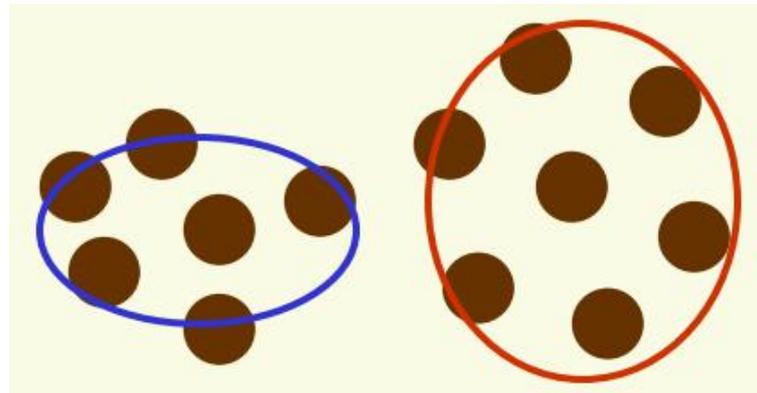
Parametric Supervised Learning

- now the distributions are fully specified
- can classify unknown sample using MAP classifier



Parametric Unsupervised Learning

- In unsupervised learning, no one tells us the true classes for samples. We still know
 - have m classes
 - have samples x_1, \dots, x_n each of unknown class
 - probability distribution for class i is $p_i(x | \theta_i)$
- Can we determine the classes and parameters simultaneously?



Outline

- Parametric Unsupervised Learning
- **Mixture Density Model**
- Gaussian Mixture Model
- Expectation Maximization (EM) Algorithm
- Applications

Mixture Density Model

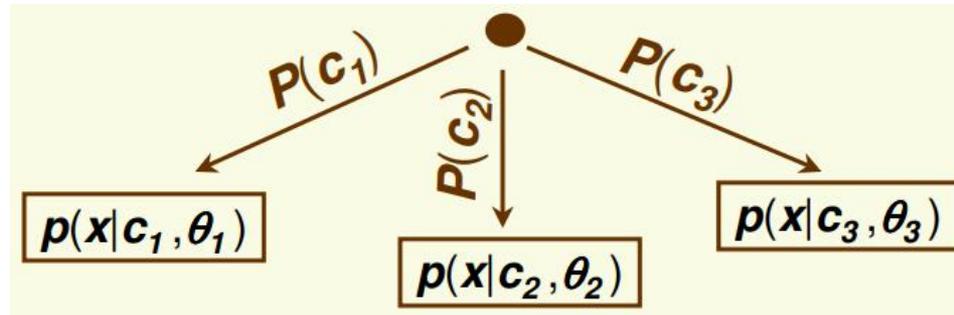
- Model data with **density model**.

component densities

$$p(x|\theta) = \sum_{j=1}^m \underbrace{p(x|c_j, \theta_j)}_{\text{mixing parameters}} P(c_j)$$

where $\theta = \{\theta_1, \dots, \theta_m\}$
 $P(c_1) + P(c_2) + \dots + P(c_m) = 1$

- To generate a sample from distribution $p(x|\theta)$
 - first select j with probability $p(c_j)$
 - then generate x according to probability law $p(x|c_j, \theta_j)$



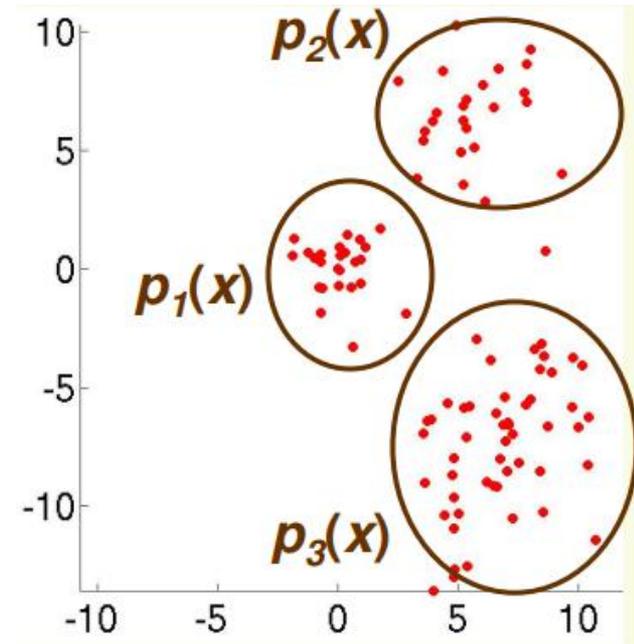
Example: Gaussian Mixture Density

- Mixture of 3 Gaussians

$$p_1(\mathbf{x}) \equiv N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$$

$$p_2(\mathbf{x}) \equiv N\left(\begin{bmatrix} 6 \\ 6 \end{bmatrix}, \begin{bmatrix} 4 & 0 \\ 0 & 4 \end{bmatrix}\right)$$

$$p_3(\mathbf{x}) \equiv N\left(\begin{bmatrix} 7 \\ -7 \end{bmatrix}, \begin{bmatrix} 6 & 0 \\ 0 & 6 \end{bmatrix}\right)$$



$$p(\mathbf{x}) = 0.2p_1(\mathbf{x}) + 0.3p_2(\mathbf{x}) + 0.5p_3(\mathbf{x})$$

Mixture Density

$$p(\mathbf{x} | \theta) = \sum_{j=1}^m p(\mathbf{x} | \mathbf{c}_j, \theta_j) P(\mathbf{c}_j)$$

- $P(\mathbf{c}_1), \dots, P(\mathbf{c}_m)$ can be known or unknown
- Suppose we know how to estimate $\theta_1, \dots, \theta_m$ and $P(\mathbf{c}_1), \dots, P(\mathbf{c}_m)$
- Can “break apart” mixture $p(\mathbf{x} | \theta)$ for classification
- To classify sample \mathbf{x} , use MAP estimation, that is choose class i which maximizes

$$P(\mathbf{c}_i | \mathbf{x}, \theta_i) \propto \underbrace{p(\mathbf{x} | \mathbf{c}_i, \theta_i)}_{\text{probability of component } i \text{ to generate } \mathbf{x}} \underbrace{P(\mathbf{c}_i)}_{\text{probability of component } i}$$

ML Estimation for Mixture Density

$$p(\mathbf{x} | \theta, \rho) = \sum_{j=1}^m p(\mathbf{x} | \mathbf{c}_j, \theta_j) P(\mathbf{c}_j) = \sum_{j=1}^m p(\mathbf{x} | \mathbf{c}_j, \theta_j) \rho_j$$

- Can use Maximum Likelihood estimation for a mixture density; need to estimate

$$\theta_1, \dots, \theta_m$$

$$\rho_1 = P(\mathbf{c}_1), \dots, \rho_m = P(\mathbf{c}_m), \text{ and } \rho = \{\rho_1, \dots, \rho_m\}$$

- As in the supervised case, form the logarithm likelihood function

$$l(\theta, \rho) = \ln p(D | \theta, \rho) = \sum_{k=1}^n \ln p(\mathbf{x}_k | \theta, \rho) = \sum_{k=1}^n \ln \left[\sum_{j=1}^m p(\mathbf{x} | \mathbf{c}_j, \theta_j) \rho_j \right]$$

ML Estimation for Mixture Density

$$l(\theta, \rho) = \sum_{k=1}^n \ln \left[\sum_{j=1}^m p(x | c_j, \theta_j) \rho_j \right]$$

- need to maximize $l(\theta, \rho)$ with respect to θ and ρ
- As you may have guessed, $l(\theta, \rho)$ is not the easiest function to maximize
 - If we take partial derivatives with respect to θ , ρ and set them to 0, typically we have a “coupled” nonlinear system of equation
 - usually closed form solution cannot be found
- We could use the gradient ascent method in general, it is not the greatest method to use, should only be used as last resort
- There is a better algorithm, called Expectation Maximization (EM).

- Before EM, let's look at the mixture density again

$$p(\mathbf{x} | \theta, \rho) = \sum_{j=1}^m p(\mathbf{x} | \mathbf{c}_j, \theta_j) \rho_j$$

- Suppose we know how to estimate θ and ρ
 - Estimating the class of \mathbf{x} is easy with MAP, maximize

$$p(\mathbf{x} | \mathbf{c}_i, \theta_i) P(\mathbf{c}_i) = p(\mathbf{x} | \mathbf{c}_i, \theta_i) \rho_i$$

- Suppose we know the class of samples $\mathbf{x}_1, \dots, \mathbf{x}_n$
 - This is just the supervised learning case, so estimating θ and ρ is easy

$$\hat{\theta}_i = \underset{\theta_i}{\operatorname{argmax}} [\ln p(\mathbf{D}_i | \theta_i)] \quad \hat{\rho}_i = \frac{|\mathbf{D}_i|}{n}$$

- This is an example of chicken-and-egg problem
 - EM algorithm approaches this problem by adding “hidden” variables

Outline

- Parametric Unsupervised Learning
- Mixture Density Model
- **Gaussian Mixture Model**
- Expectation Maximization (EM) Algorithm
- Applications

What Model Should We Use?

- Depends on X !
- Here, maybe Gaussian Naïve Bayes?
 - – Multinomial over clusters Y
 - – (Independent) Gaussian for each X_i given Y

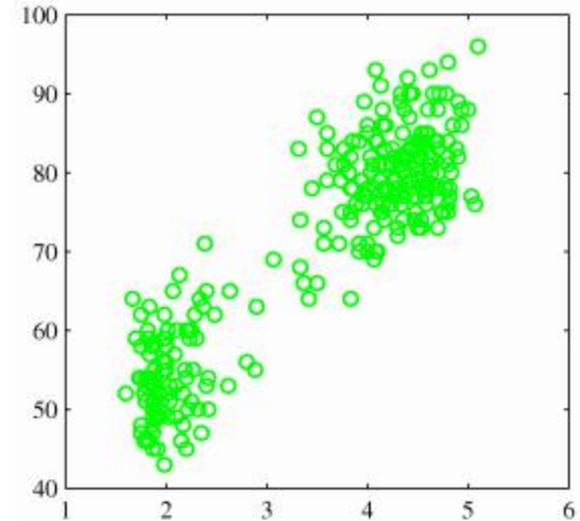
$$p(Y_i = y_k) = \theta_k$$

$$P(X_i = x | Y = y_k) = \frac{1}{\sigma_{ik} \sqrt{2\pi}} e^{-\frac{(x - \mu_{ik})^2}{2\sigma_{ik}^2}}$$

Y	X ₁	X ₂
??	0.1	2.1
??	0.5	-1.1
??	0.0	3.0
??	-0.1	-2.0
??	0.2	1.5
...

Could we make fewer assumptions?

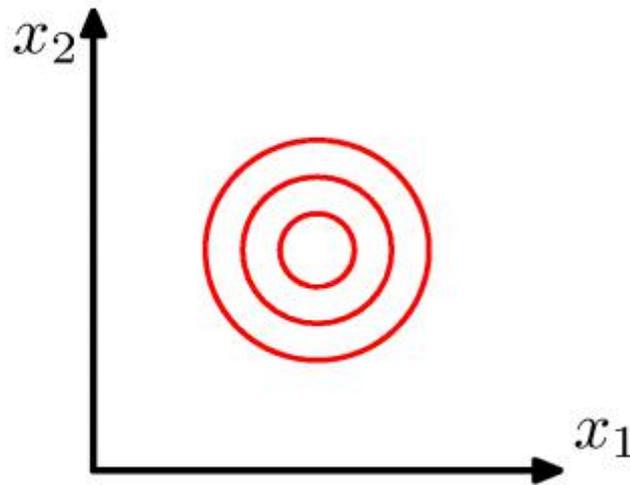
- What if the X_i co-vary?
- What if there are multiple peaks?
- **Gaussian Mixture Models!**
 - – $P(Y)$ still multinomial
 - – $P(X|Y)$ is a multivariate Gaussian distribution:



$$P(X = \mathbf{x}_j | Y = i) = \frac{1}{(2\pi)^{m/2} \|\Sigma_i\|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \mu_i)\right]$$

Multivariate Gaussians

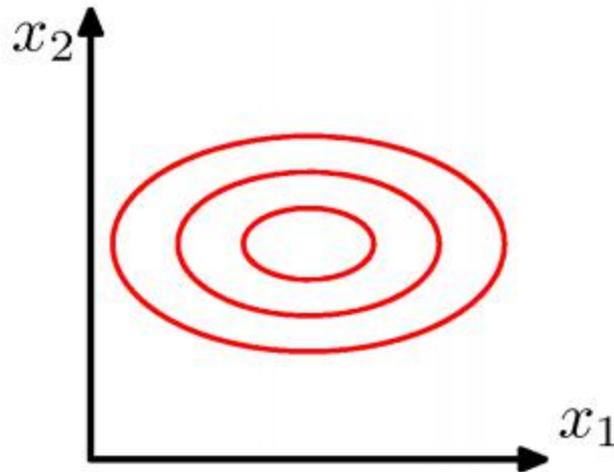
$$P(X=\mathbf{x}_j) = \frac{1}{(2\pi)^{m/2} \|\Sigma\|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}_j - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}_j - \boldsymbol{\mu})\right]$$



$\Sigma \propto$ identity matrix

Multivariate Gaussians

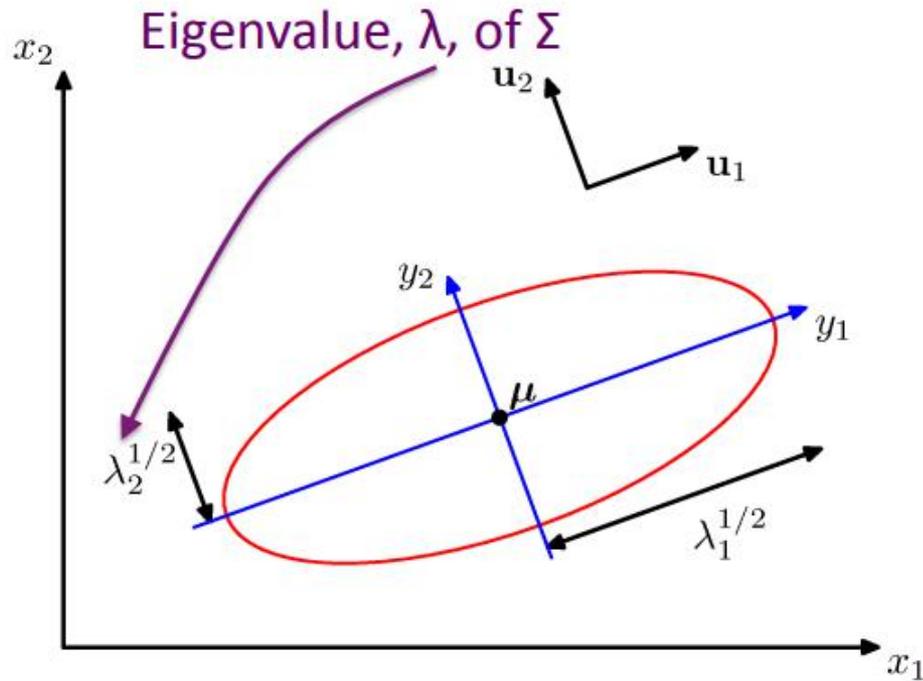
$$P(X=\mathbf{x}_j) = \frac{1}{(2\pi)^{m/2} \|\Sigma\|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}_j - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}_j - \boldsymbol{\mu})\right]$$



Σ = diagonal matrix

X_i are independent *ala* Gaussian NB

Multivariate Gaussians

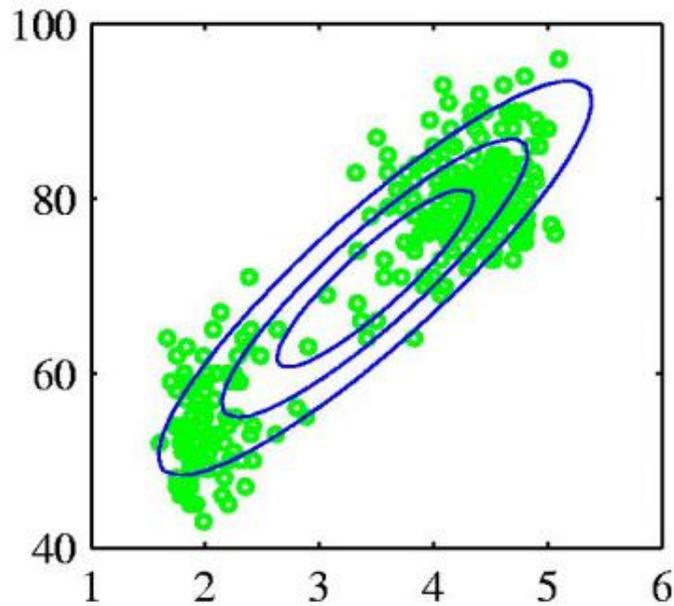


Covariance matrix, Σ , =
degree to which x_i vary
together

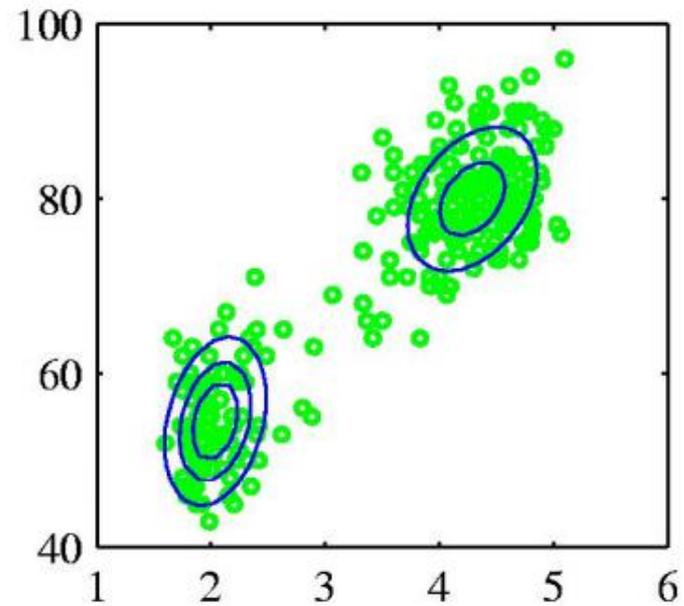
$$P(X=\mathbf{x}_j) = \frac{1}{(2\pi)^{m/2} \|\Sigma\|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}_j - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}_j - \boldsymbol{\mu})\right]$$

Mixtures of Gaussians (1)

Old Faithful Data Set

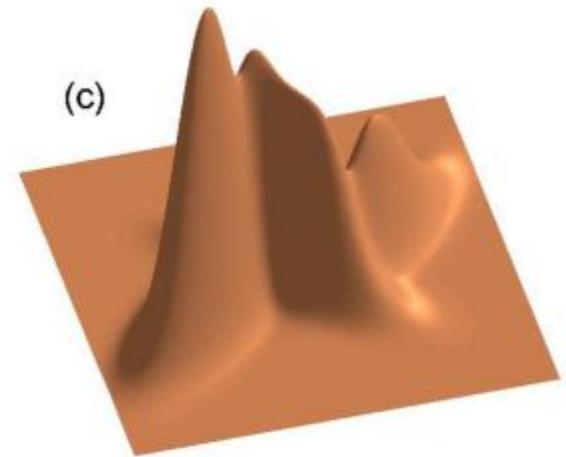
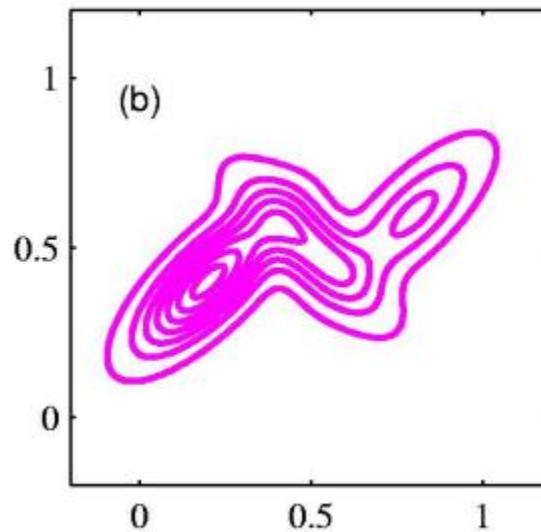
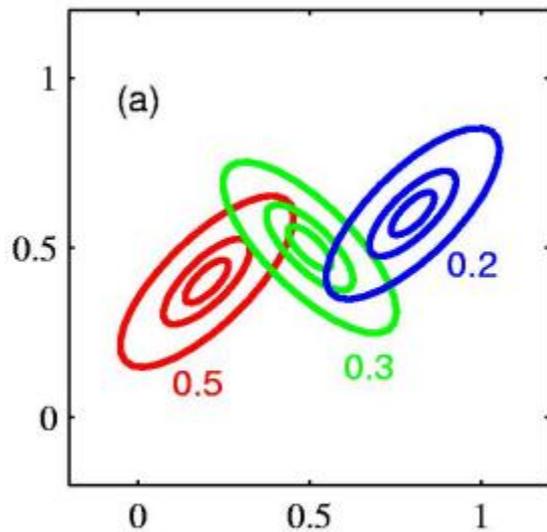


Single Gaussian



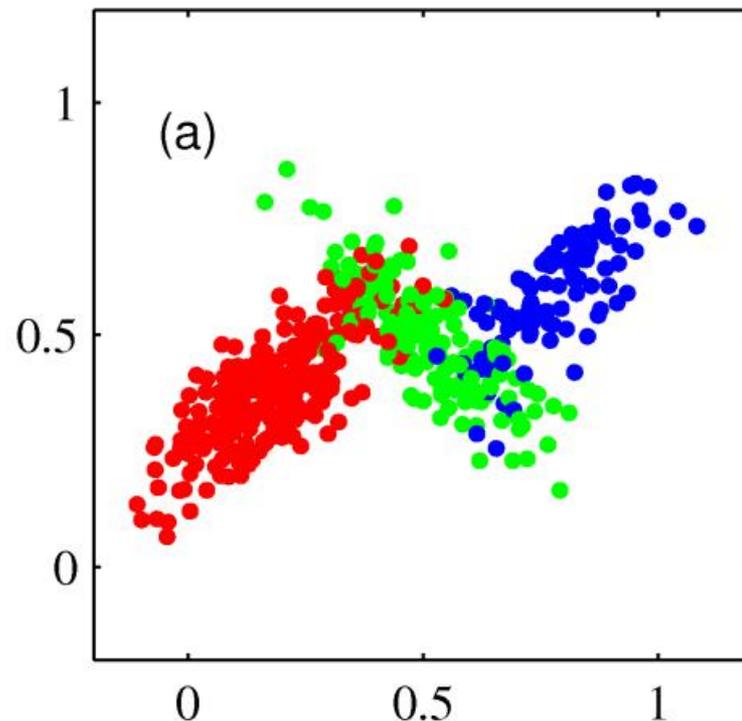
Mixture of two Gaussians

Mixtures of Gaussians (3)



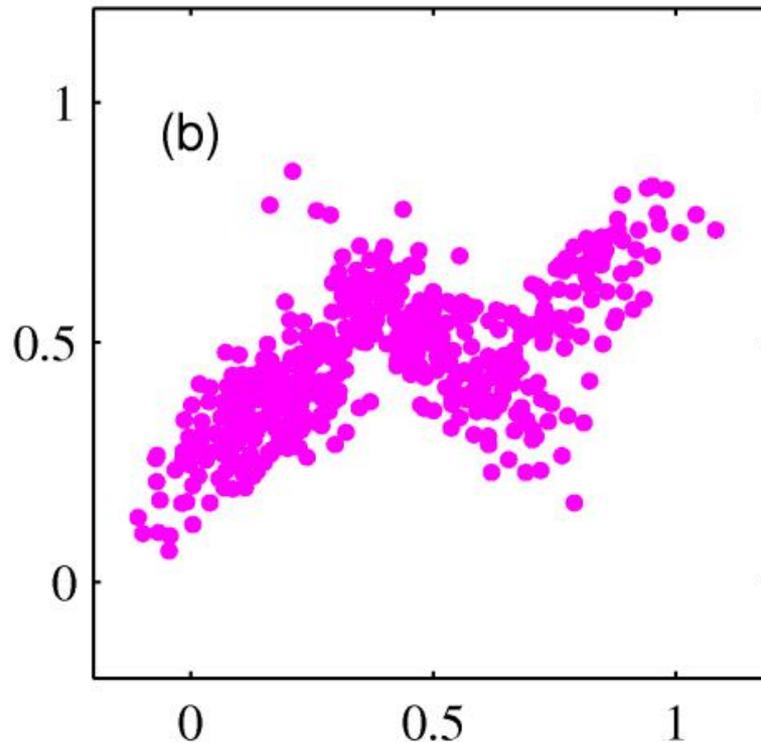
Eliminating Hard Assignments to Clusters

- Model data as mixture of multivariate Gaussians



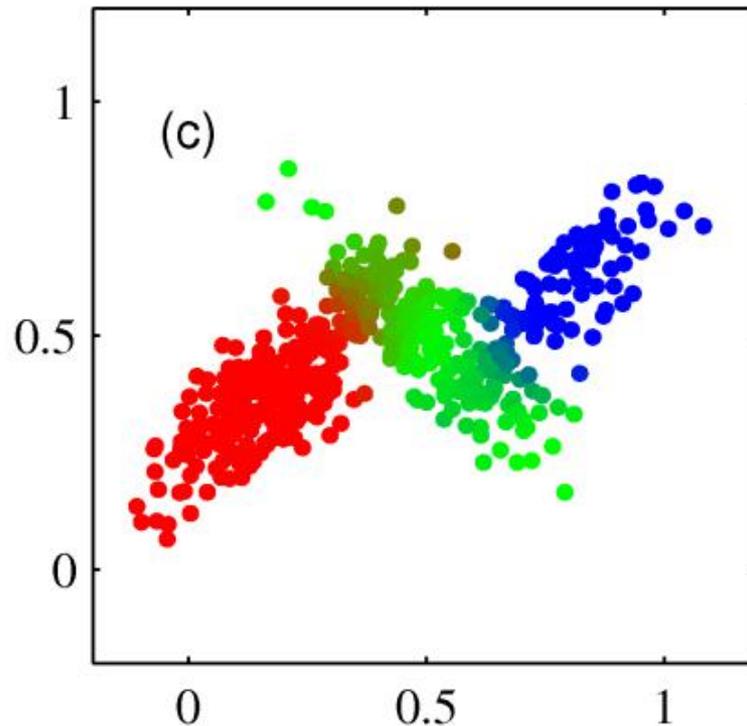
Eliminating Hard Assignments to Clusters

- Model data as mixture of multivariate Gaussians



Eliminating Hard Assignments to Clusters

- Model data as mixture of multivariate Gaussians



Shown is the **posterior probability** that a point was generated from i -th Gaussian: $\Pr(Y = i | x)$

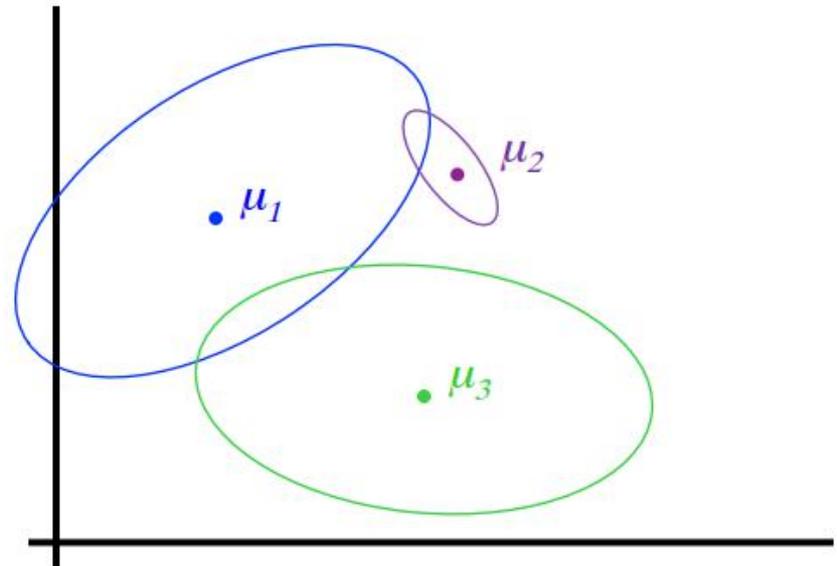
The General GMM Assumption

- $P(Y)$: There are k components
- $P(X|Y)$: Each component generates data from a **multivariate Gaussian** with mean μ_i and covariance matrix Σ_i

Each data point is sampled from a **generative process**:

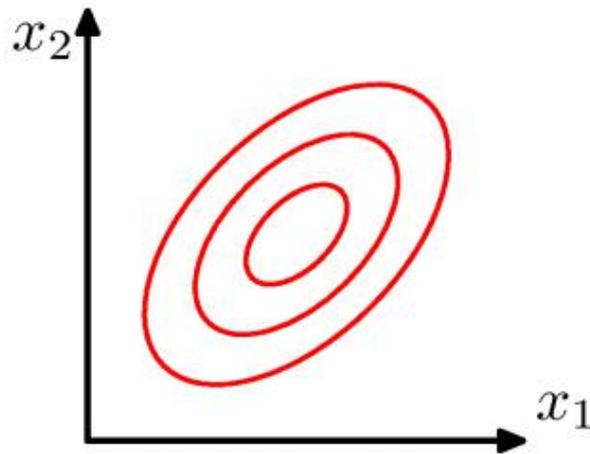
1. Choose component i with probability $P(y=i)$
2. Generate datapoint $\sim N(m_i, \Sigma_i)$

Gaussian mixture model
(GMM)



Multivariate Gaussians

$$P(X=\mathbf{x}_j) = \frac{1}{(2\pi)^{m/2} \|\Sigma\|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}_j - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x}_j - \boldsymbol{\mu})\right]$$



Σ = arbitrary (semidefinite) matrix:

- specifies rotation (change of basis)
- eigenvalues specify relative elongation

ML Estimation in Supervised Setting

- Univariate Gaussian

$$\mu_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i \quad \sigma_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2$$

- Mixture of Multivariate Gaussians
 - ML estimate for each of the Multivariate Gaussians is given by:

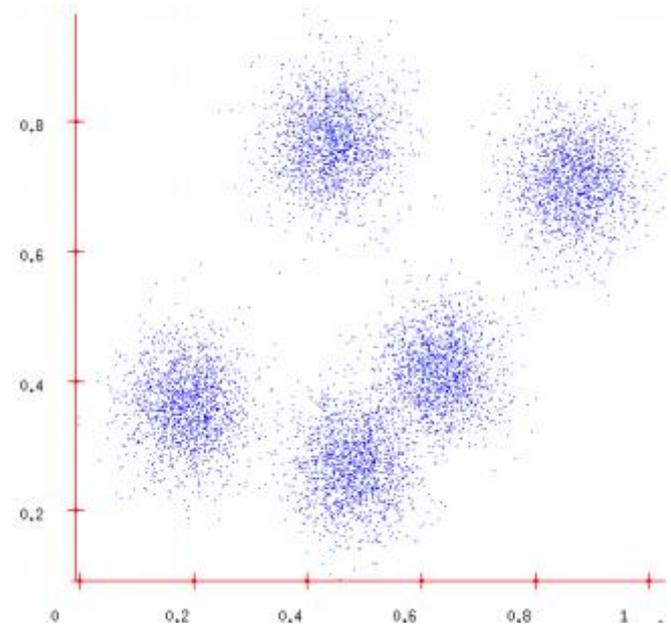
$$\mu_{ML}^k = \frac{1}{n} \sum_{j=1}^n x_j \quad \Sigma_{ML}^k = \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j - \mu_{ML}^k)(\mathbf{x}_j - \mu_{ML}^k)^T$$

Just sums over \mathbf{x} generated from the k 'th Gaussian

That was easy! But what if unobserved data?

- MLE:
 - $\operatorname{argmax}_{\theta} \prod_j P(y_j, x_j)$
 - θ : all model parameters
 - eg, class probs, means, and variances
- But we don't know !!!
- Maximize **marginal likelihood**

$$- \operatorname{argmax}_{\theta} \prod_j P(x_j) = \operatorname{argmax}_{\theta} \prod_j \sum_{k=1}^K P(Y_j=k, x_j)$$



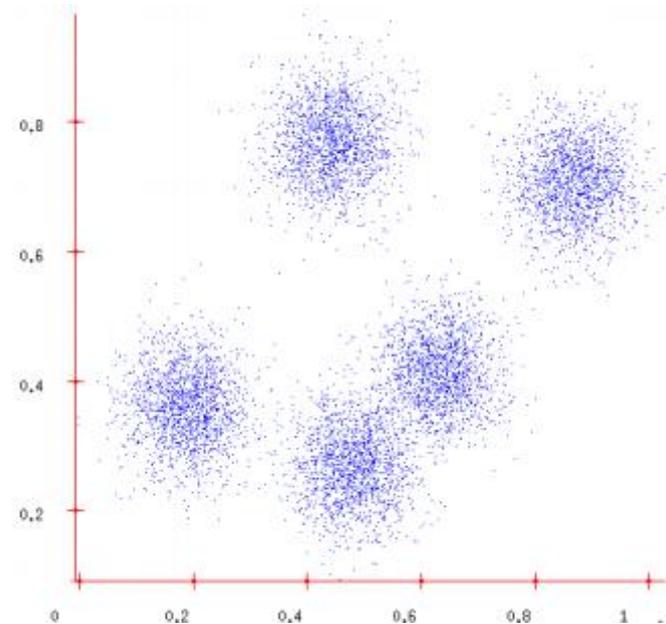
How Do We Optimize? Closed Form?

- Maximize **marginal likelihood**:

$$- \operatorname{argmax}_{\theta} \prod_j P(x_j) = \operatorname{argmax} \prod_j \sum_{k=1}^K P(Y_j=k, x_j)$$

- **Almost always a hard problem!**

- – Usually no closed-form solution
- – Even when $\lg P(X, Y)$ is convex, $\lg P(X)$ generally isn't...
- – For all but the simplest $P(X)$, we will have to do gradient ascent, in a big messy space with lots of local optimum...



Learning General Mixtures of Gaussian

$$P(y = k | \mathbf{x}_j) \propto \frac{1}{(2\pi)^{m/2} \|\Sigma_k\|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}_j - \mu_k)^T \Sigma_k^{-1}(\mathbf{x}_j - \mu_k)\right] P(y = k)$$

- Marginal likelihood:

$$\begin{aligned} \prod_{j=1}^m P(\mathbf{x}_j) &= \prod_{j=1}^m \sum_{k=1}^K P(\mathbf{x}_j, y = k) \\ &= \prod_{j=1}^m \sum_{k=1}^K \frac{1}{(2\pi)^{m/2} \|\Sigma_k\|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}_j - \mu_k)^T \Sigma_k^{-1}(\mathbf{x}_j - \mu_k)\right] P(y = k) \end{aligned}$$

- Need to differentiate and solve for $\mu_{k'}$, $\Sigma_{k'}$, and $P(Y=k)$ for $k=1..K$
- There will be no closed form solution, gradient is complex, lots of local optimum
- ***Wouldn't it be nice if there was a better way!?!***

Outline

- Parametric Unsupervised Learning
- Mixture Density Model
- Gaussian Mixture Model
- **Expectation Maximization (EM) Algorithm**
- Applications

Expectation Maximization Algorithm

- EM is an algorithm for ML parameter estimation when the data has missing values. It is used when
 - 1. data is incomplete (has missing values)
 - some features are missing for some samples due to data corruption, partial survey responses, etc.
 - This scenario is very useful
 - 2. Suppose data X is complete, but $p(X|\theta)$ is hard to optimize. Suppose further that introducing certain hidden variables U whose values are missing, and suppose it is easier to optimize the “complete” likelihood function $p(X,U|\theta)$. Then EM is useful.
 - This scenario is useful for the mixture density estimation, and is subject of our lecture today
- Notice that after we introduce artificial (hidden) variables U with missing values, case 2 is completely equivalent to case 1

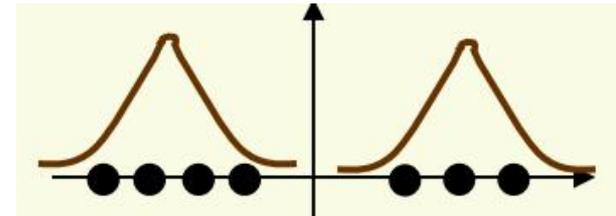
EM: Hidden Variables for Mixture Density

$$p(x | \theta) = \sum_{j=1}^m p(x | c_j, \theta_j) \rho_j$$

- For simplicity, assume component densities are

$$p(x | c_j, \theta_j) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x - \mu_j)^2}{2\sigma^2}\right)$$

- assume for now that the variance is known
- need to estimate $\theta = \{\mu_1, \dots, \mu_m\}$



- If we knew which sample came from which component (that is the class label), the ML parameter estimation is easy
- Thus to get an easier problem, introduce hidden variables which indicate which component each sample belongs to

EM: Hidden Variables for Mixture Density

- For $1 \leq i \leq n$, $1 \leq k \leq m$, define hidden variables

$$z_i^{(k)} = \begin{cases} 1 & \text{if sample } i \text{ was generated by component } k \\ 0 & \text{otherwise} \end{cases}$$

$$x_i \rightarrow \{x_i, z_i^{(1)}, \dots, z_i^{(m)}\}$$

- $z_i^{(k)}$ are indicator random variables, they indicate which Gaussian component generated sample x_i
- Let $z_i = \{z_i^{(1)}, \dots, z_i^{(m)}\}$ indicator r.v. corresponding to sample x_i
- Conditioned on z_i , distribution of x_i is Gaussian

$$p(x_i | z_i, \theta) \sim N(\mu_k, \sigma^2)$$

$$\text{where } k \text{ is s.t. } z_i^{(k)} = 1$$

EM: Joint Likelihood

- Let $\mathbf{z}_i = \{\mathbf{z}_i^{(1)}, \dots, \mathbf{z}_i^{(m)}\}$, and $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_n\}$
- The complete likelihood is

$$\begin{aligned} p(\mathbf{X}, \mathbf{Z} | \theta) &= p(x_1, \dots, x_n, z_1, \dots, z_n | \theta) = \prod_{i=1}^n p(x_i, z_i | \theta) \\ &= \prod_{i=1}^n \underbrace{p(x_i | z_i, \theta)}_{\text{gaussian}} \underbrace{p(z_i)}_{\text{part of } \rho_c} \end{aligned}$$

- If we actually observed \mathbf{Z} , the log likelihood $\ln[p(\mathbf{X}, \mathbf{Z} | \theta)]$ would be trivial to maximize with respect to θ and ρ_i
- The problem, of course, is that the values of \mathbf{Z} are missing, since we made it up (that is \mathbf{Z} is hidden)

EM Derivation

- Instead of maximizing $\ln[p(X,Z|\theta)]$ the idea behind EM is to maximize some function of $\ln[p(X,Z|\theta)]$, usually its expected value

$$E_Z[\ln p(X,Z|\theta)]$$

- If θ makes $\ln[p(X,Z|\theta)]$ large, then θ tends to make $E[\ln p(X,Z|\theta)]$ large
- the expectation is with respect to the missing data Z
- that is with respect to density $p(Z|X,\theta)$
- however θ is our ultimate goal, we don't know θ !

EM Algorithm

- EM solution is to iterate

1. start with initial parameters $\theta^{(0)}$

iterate the following 2 step until convergence

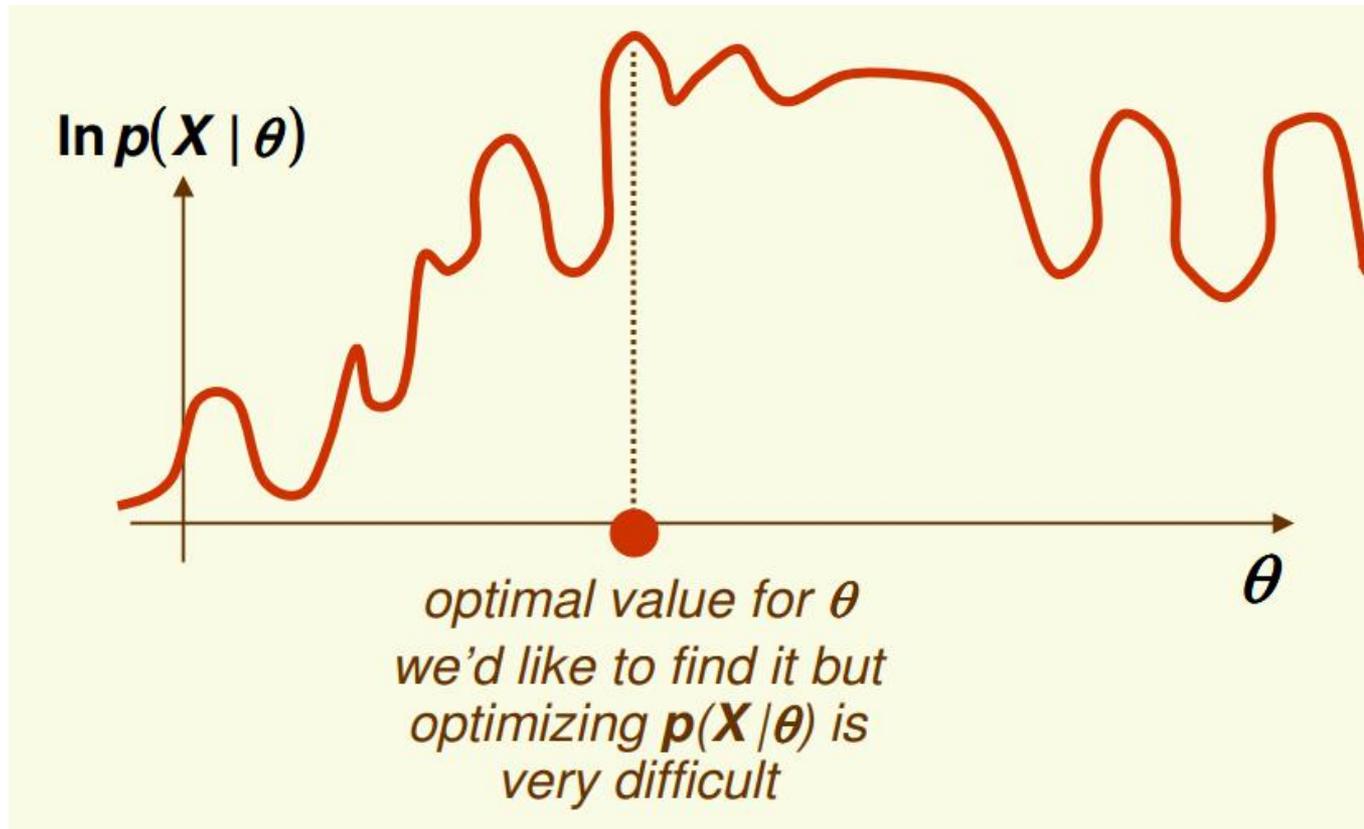
E. compute the expectation of log likelihood with respect to current estimate $\theta^{(t)}$ and \mathbf{X} . Let's call it $Q(\theta | \theta^{(t)})$

$$Q(\theta | \theta^{(t)}) = E_Z [\ln p(\mathbf{X}, \mathbf{Z} | \theta) | \mathbf{X}, \theta^{(t)}]$$

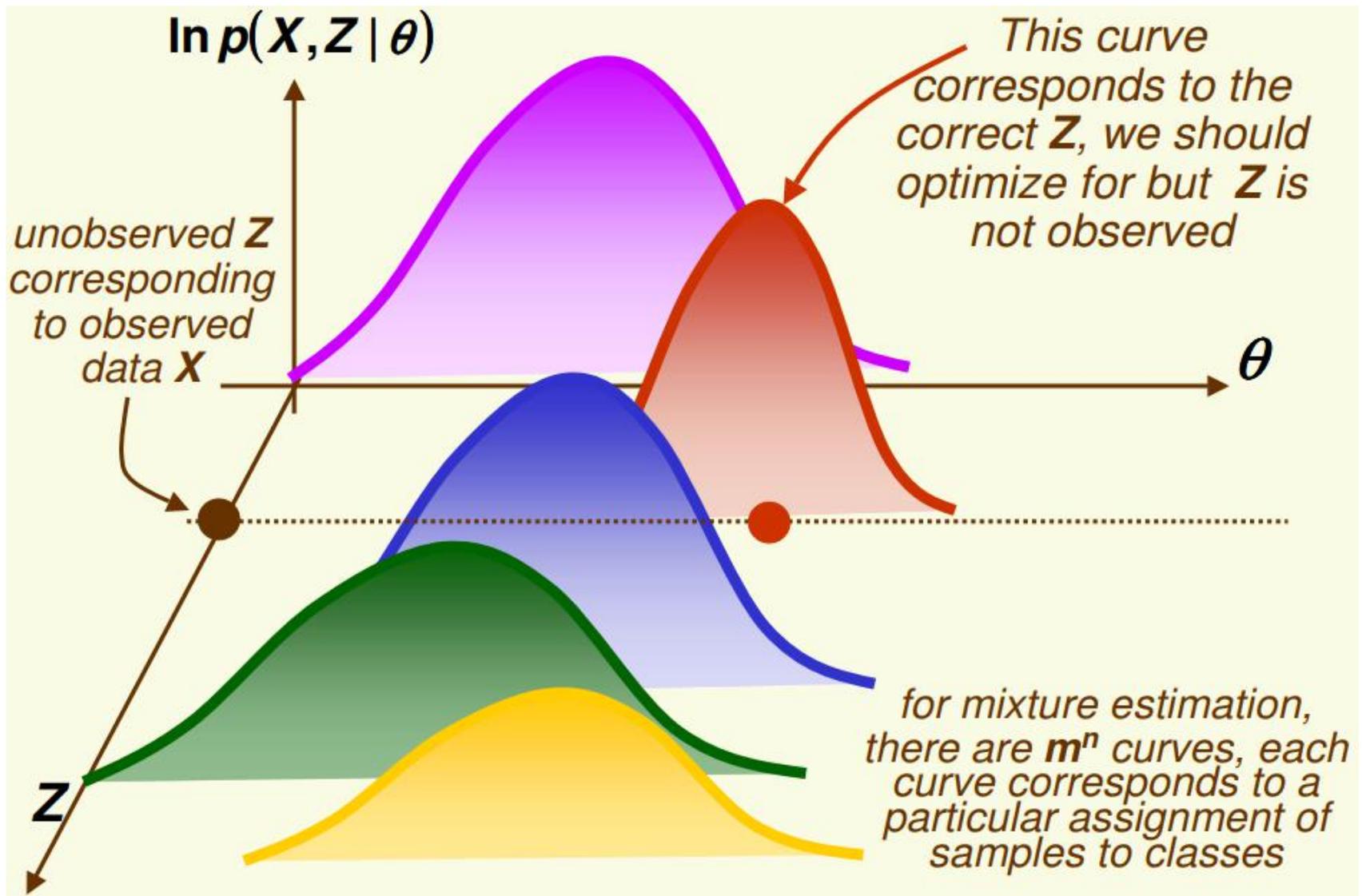
M. maximize $Q(\theta | \theta^{(t)})$

$$\theta^{(t+1)} = \underset{\theta}{\operatorname{argmax}} Q(\theta | \theta^{(t)})$$

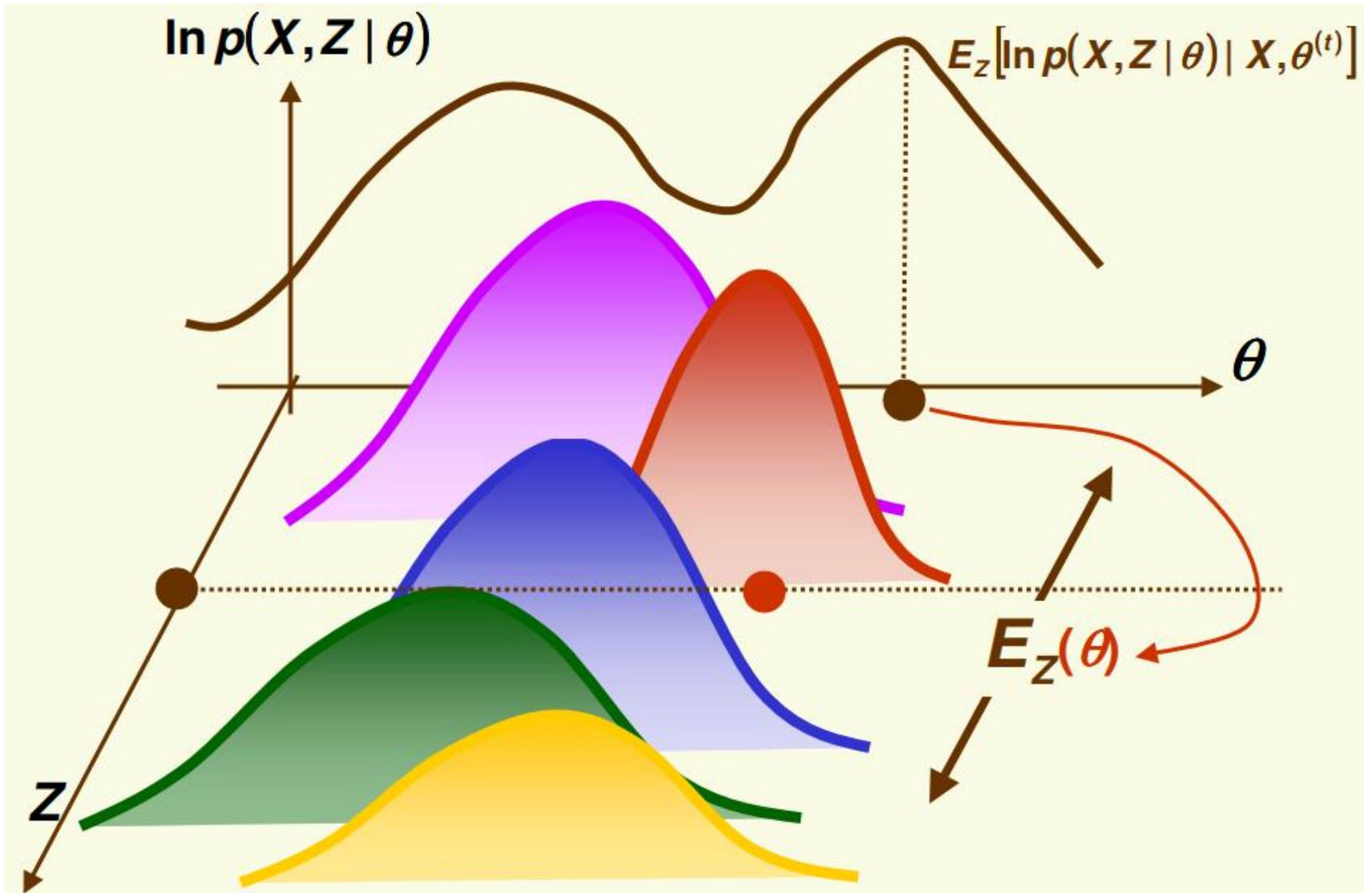
EM Algorithm: Picture



EM Algorithm: Picture



EM Algorithm: Picture



EM Algorithm

- It can be proven that EM algorithm converges to the local maximum of the log-likelihood

$$\ln p(X | \theta)$$

- Why is it better than gradient ascent?
 - Convergence of EM is usually significantly faster, in the beginning, very large steps are made (that is likelihood function increases rapidly), as opposed to gradient ascent which usually takes tiny steps
 - gradient descent is not guaranteed to converge
 - recall all the difficulties of choosing the appropriate
 - learning rate

EM for Mixture of Gaussians: E step

- Let's come back to our example $p(\mathbf{x} | \theta) = \sum_{j=1}^m p(\mathbf{x} | \mathbf{c}_j, \theta_j) \rho_j$

$$p(\mathbf{x} | \mathbf{c}_j, \theta_j) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(\mathbf{x} - \mu_j)^2}{2\sigma^2}\right)$$

need to estimate $\theta = \{\mu_1, \dots, \mu_m\}$ and ρ_1, \dots, ρ_m

- For $1 \leq i \leq n$, $1 \leq k \leq m$, we define

$$\mathbf{z}_i^{(k)} = \begin{cases} 1 & \text{if sample } i \text{ was generated by component } k \\ 0 & \text{otherwise} \end{cases}$$

- As before, $\mathbf{z}_i = \{\mathbf{z}_i^{(1)}, \dots, \mathbf{z}_i^{(m)}\}$ and $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_n\}$
- We need log-likelihood of observed \mathbf{X} and hidden \mathbf{Z}

$$\ln p(\mathbf{X}, \mathbf{Z} | \theta) = \ln \prod_{i=1}^n p(\mathbf{x}_i, \mathbf{z}_i | \theta) = \sum_{i=1}^n \ln p(\mathbf{x}_i | \mathbf{z}_i, \theta) P(\mathbf{z}_i)$$

EM for Mixture of Gaussians: E step

- We need log-likelihood of observed X and hidden Z

$$\ln p(X, Z | \theta) = \ln \prod_{i=1}^n p(x_i, z_i | \theta) = \sum_{i=1}^n \ln p(x_i | z_i, \theta) P(z_i)$$

- First let's rewrite $p(x_i | z_i, \theta) P(z_i)$

$$p(x_i | z_i, \theta) P(z_i) = \begin{cases} p(x_i | z_i^{(1)} = 1, \theta) P(z_i^{(1)} = 1) & \text{if } z_i^{(1)} = 1 \\ \vdots \\ p(x_i | z_i^{(m)} = 1, \theta) P(z_i^{(m)} = 1) & \text{if } z_i^{(m)} = 1 \end{cases}$$

$$= \prod_{k=1}^m [p(x_i | z_i^{(k)} = 1, \theta) P(z_i^{(k)} = 1)]^{z_i^{(k)}}$$

$$= \prod_{k=1}^m \left[\frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x_i - \mu_k)^2}{2\sigma^2}\right) P(z_i^{(k)} = 1) \right]^{z_i^{(k)}}$$

EM for Mixture of Gaussians: E step

- log-likelihood of observed X and hidden Z is

$$\begin{aligned}\ln p(X, Z | \theta) &= \sum_{i=1}^n \ln p(x_i | z_i, \theta) P(z_i) \\ &= \sum_{i=1}^n \ln \prod_{k=1}^m \left[\frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x_i - \mu_k)^2}{2\sigma^2}\right) P(z_i^{(k)} = 1) \right]^{z_i^{(k)}} \\ &= \sum_{i=1}^n \sum_{k=1}^m \ln \left[\frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x_i - \mu_k)^2}{2\sigma^2}\right) P(z_i^{(k)} = 1) \right]^{z_i^{(k)}} \\ &= \sum_{i=1}^n \sum_{k=1}^m z_i^{(k)} \left[\ln \frac{1}{\sigma \sqrt{2\pi}} - \frac{(x_i - \mu_k)^2}{2\sigma^2} + \underbrace{\ln P(z_i^{(k)} = 1)}_{P(\text{sample } x_i \text{ from class } k) = P(c_k) = \rho_k} \right] \\ &= \sum_{i=1}^n \sum_{k=1}^m z_i^{(k)} \left[\ln \frac{1}{\sigma \sqrt{2\pi}} - \frac{(x_i - \mu_k)^2}{2\sigma^2} + \ln \rho_k \right]\end{aligned}$$

EM for Mixture of Gaussians: E step

- log-likelihood of observed X and hidden Z is

$$\ln p(X, Z | \theta) = \sum_{i=1}^n \sum_{k=1}^m z_i^{(k)} \left[\ln \frac{1}{\sigma \sqrt{2\pi}} - \frac{(x_i - \mu_k)^2}{2\sigma^2} + \ln \rho_k \right]$$

- For the E step, we must compute

$$\begin{aligned} Q(\theta | \theta^{(t)}) &= Q(\theta | \mu_1^{(t)}, \dots, \mu_m^{(t)}, \rho_1^{(t)}, \dots, \rho_m^{(t)}) = E_Z [\ln p(X, Z | \theta) | X, \theta^{(t)}] \\ &= E_Z \left(\sum_{i=1}^n \sum_{k=1}^m z_i^{(k)} \left[\ln \frac{1}{\sigma \sqrt{2\pi}} - \frac{(x_i - \mu_k)^2}{2\sigma^2} + \ln \rho_k^{(t)} \right] \right) \\ &\quad \Downarrow E_X \left[\sum_i a_i x_i + b \right] = \sum_i a_i E_X [x_i] + b \\ &= \sum_{i=1}^n \sum_{k=1}^m E_Z [z_i^{(k)}] \left(\ln \frac{1}{\sigma \sqrt{2\pi}} - \frac{(x_i - \mu_k)^2}{2\sigma^2} + \ln \rho_k \right) \end{aligned}$$

EM for Mixture of Gaussians: E step

$$Q(\theta | \theta^{(t)}) = \sum_{i=1}^n \sum_{k=1}^m E_Z[z_i^{(k)}] \left(\ln \frac{1}{\sigma \sqrt{2\pi}} - \frac{(x_i - \mu_k)^2}{2\sigma^2} + \ln \rho_k \right)$$

- need to compute $E_Z[z_i^{(k)}]$ in the above expression

$$\begin{aligned} E_Z[z_i^{(k)}] &= 0 * P(z_i^{(k)} = 0 | \theta^{(t)}, x_i) + 1 * P(z_i^{(k)} = 1 | \theta^{(t)}, x_i) \\ &= P(z_i^{(k)} = 1 | \theta^{(t)}, x_i) = \frac{p(x_i | \theta^{(t)}, z_i^{(k)} = 1) P(z_i^{(k)} = 1 | \theta^{(t)})}{p(x_i | \theta^{(t)})} \end{aligned}$$

$$= \frac{\rho_k^{(t)} \exp\left(-\frac{1}{2}(x_i - \mu_k^{(t)})^2\right)}{\sum_{j=1}^m P(x_i | \theta^{(t)}, z_i^{(j)} = 1) P(z_i^{(j)} = 1 | \theta^{(t)})} = \frac{\rho_k^{(t)} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu_k^{(t)})^2\right)}{\sum_{j=1}^m \rho_j^{(t)} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu_j^{(t)})^2\right)}$$

- We are finally done with the E step
 - for implementation, just need to compute $E_Z[z_i^{(k)}]$ s don't need to compute Q

EM for Mixture of Gaussians: M step

$$Q(\theta | \theta^{(t)}) = \sum_{i=1}^n \sum_{k=1}^m \mathbf{E}_Z[z_i^{(k)}] \left(\ln \frac{1}{\sigma \sqrt{2\pi}} - \frac{(x_i - \mu_k)^2}{2\sigma^2} + \ln \rho_k \right)$$

- Need to maximize Q with respect to all parameters
- First differentiate with respect to μ_k

$$\frac{\partial}{\partial \mu_k} Q(\theta | \theta^{(t)}) = \sum_{i=1}^n \mathbf{E}_Z[z_i^{(k)}] \frac{(x_i - \mu_k)}{\sigma^2} = 0$$

$$\Rightarrow \text{new } \mu_k = \mu_k^{(t+1)} = \frac{1}{n} \sum_{i=1}^n \mathbf{E}_Z[z_i^{(k)}] x_i$$



the mean for class k is weighted average of all samples, and this weight is proportional to the current estimate of probability that the sample belongs to class k

EM for Mixture of Gaussians: M step

$$Q(\theta | \theta^{(t)}) = \sum_{i=1}^n \sum_{k=1}^m E_Z[z_i^{(k)}] \left(\ln \frac{1}{\sigma \sqrt{2\pi}} - \frac{(x_i - \mu_k)^2}{2\sigma^2} + \ln \rho_k \right)$$

- For ρ_k we have to use Lagrange multipliers to preserve constraint

$$\sum_{j=1}^m \rho_j = 1$$

- Thus we need to differentiate

$$F(\lambda, \rho) = Q(\theta | \theta^{(t)}) - \lambda \left(\sum_{j=1}^m \rho_j - 1 \right)$$

$$\frac{\partial}{\partial \rho_k} F(\lambda, \rho) = \sum_{i=1}^n \frac{1}{\rho_k} E_Z[z_i^{(k)}] - \lambda = 0 \Rightarrow \sum_{i=1}^n E_Z[z_i^{(k)}] - \lambda \rho_k = 0$$

- Summing up over all components $\sum_{k=1}^m \sum_{i=1}^n E_Z[z_i^{(k)}] = \sum_{k=1}^m \lambda \rho_k$
- Since $\sum_{k=1}^m \sum_{i=1}^n E_Z[z_i^{(k)}] = n$ and $\sum_{k=1}^m \rho_k = 1$, we get $\lambda = n$

$$\rho_k^{(t+1)} = \frac{1}{n} \sum_{i=1}^n E_Z[z_i^{(k)}]$$

EM Algorithm

- The algorithm on this slide applies ONLY to univariate gaussian case with known variances

1. Randomly initialize $\mu_1, \dots, \mu_m, \rho_1, \dots, \rho_m$ (with constraint $\sum \rho_i = 1$)

iterate until no change in $\mu_1, \dots, \mu_m, \rho_1, \dots, \rho_m$

E. for all i, k , compute

$$E_z[z_i^{(k)}] = \frac{\rho_k \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu_k)^2\right)}{\sum_{j=1}^m \rho_j \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu_j)^2\right)}$$

M. for all k , do parameter update

$$\mu_k = \frac{1}{n} \sum_{i=1}^n E_z[z_i^{(k)}] x_i \quad \rho_k = \frac{1}{n} \sum_{i=1}^n E_z[z_i^{(k)}]$$

EM Algorithm

- For the more general case of multivariate
- Gaussians with unknown means and variances

$$\mathbf{E} \text{ step: } \mathbf{E}_Z[\mathbf{z}_i^{(k)}] = \frac{\rho_k p(\mathbf{x} | \mu_k, \Sigma_k)}{\sum_{j=1}^m \rho_j p(\mathbf{x} | \mu_j, \Sigma_j)}$$

$$\text{where } p(\mathbf{x} | \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{d/2} |\Sigma_k^{-1}|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mu_k)^t \Sigma_k^{-1}(\mathbf{x} - \mu_k)\right]$$

M step:

$$\rho_k = \frac{1}{n} \sum_{i=1}^n \mathbf{E}_Z[\mathbf{z}_i^{(k)}]$$

$$\mu_k = \frac{\sum_{i=1}^n \mathbf{E}_Z[\mathbf{z}_i^{(k)}] \mathbf{x}_i}{\sum_{i=1}^n \mathbf{E}_Z[\mathbf{z}_i^{(k)}]}$$

$$\Sigma_k = \frac{\sum_{i=1}^n \mathbf{E}_Z[\mathbf{z}_i^{(k)}] (\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^T}{\sum_{i=1}^n \mathbf{E}_Z[\mathbf{z}_i^{(k)}]}$$

EM Algorithm and K-means

- k-means can be derived from EM algorithm
- Setting mixing parameters equal for all classes,

$$E_Z[z_i^{(k)}] = \frac{\rho_k \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu_k)^2\right)}{\sum_{j=1}^m \rho_j \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu_j)^2\right)} = \frac{\exp\left(-\frac{1}{2\sigma^2}(x_i - \mu_k)^2\right)}{\sum_{j=1}^m \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu_j)^2\right)}$$

- If we let $\sigma \rightarrow \infty$, then

$$E_Z[z_i^{(k)}] = \begin{cases} 1 & \text{if } \forall j, \|x_i - \mu_k\| > \|x_i - \mu_j\| \\ 0 & \text{otherwise} \end{cases}$$

- so at the E step, for each current mean, we find all points closest to it and form new clusters
- at the M step, we compute the new means inside current clusters

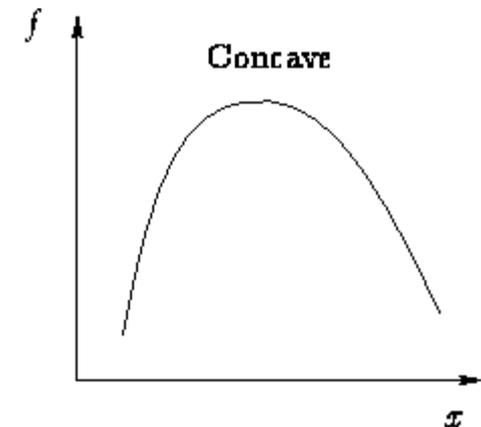
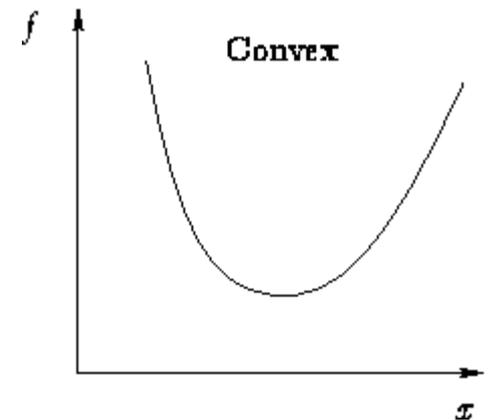
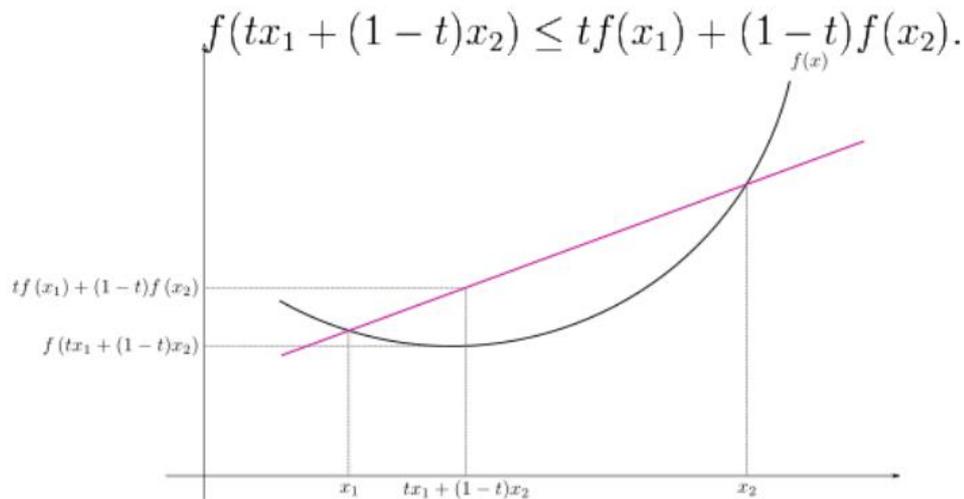
$$\mu_k = \frac{1}{n} \sum_{i=1}^n E_Z[z_i^{(k)}] x_i$$

Properties of EM

- We will prove that
 - – EM converges to a local minima
 - – Each iteration improves the $\log Z$ likelihood
- How? (Same as k-means)
 - – E-step can never decrease likelihood
 - – M-step can never decrease likelihood

Jensen's Inequality

- **Theorem:** $\log \sum_z P(z) f(z) \geq \sum_z P(z) \log f(z)$
 - e.g., Binary case for convex function f :



Applying Jensen's Inequality

- Use: $\log \sum_{\mathbf{z}} P(\mathbf{z}) f(\mathbf{z}) \geq \sum_{\mathbf{z}} P(\mathbf{z}) \log f(\mathbf{z})$

$$\begin{aligned} \ell(\theta^{(t)} : \mathcal{D}) &= \sum_{j=1}^m \log \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \frac{P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)})}{Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j)} \\ &\geq \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log \left(\frac{p(\mathbf{z}, \mathbf{x}_j | \theta^{(t)})}{Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j)} \right) \\ &= \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log \left(p(\mathbf{z}, \mathbf{x}_j | \theta^{(t)}) \right) - \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log \left(Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \right) \\ \ell(\theta^{(t)} : \mathcal{D}) &\geq \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)}) + \sum_{j=1}^m H(Q^{(t+1),j}) \end{aligned}$$

The M-step

- Lower bound:

$$\ell(\theta^{(t)} : \mathcal{D}) \geq \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta^{(t)}) + \sum_{j=1}^m H(Q^{(t+1),j})$$



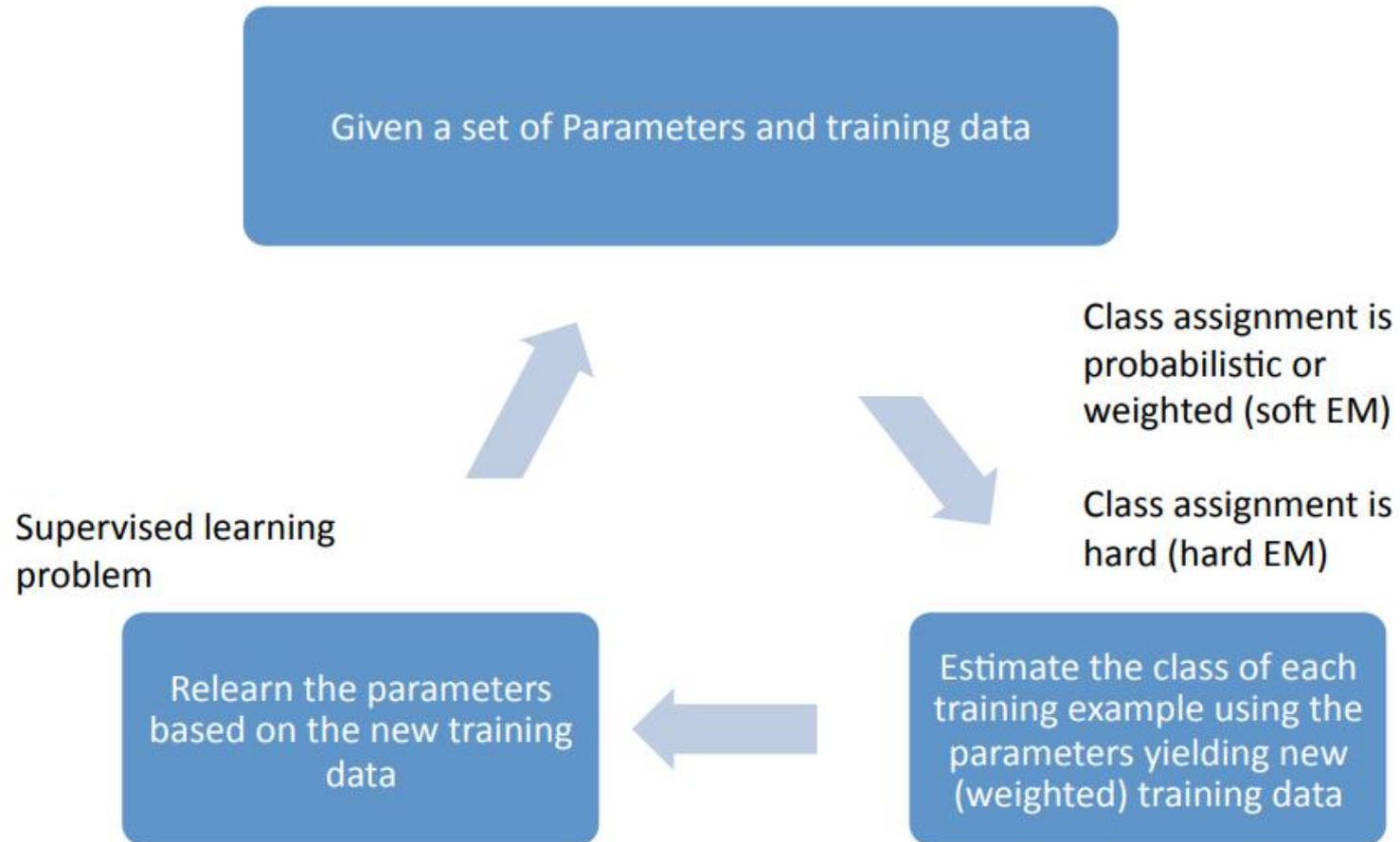
This term is a constant with respect to θ

- Maximization step

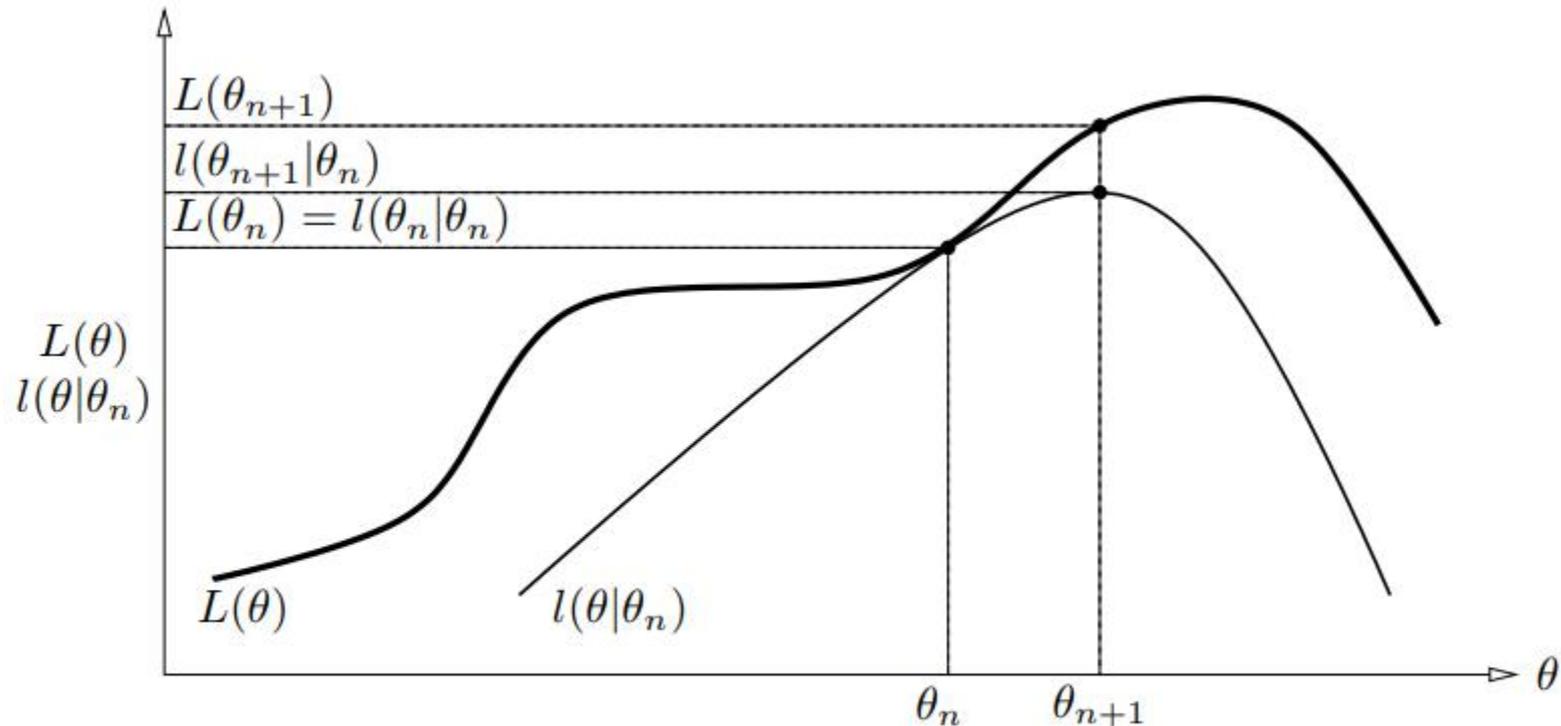
$$\theta^{(t+1)} \leftarrow \arg \max_{\theta} \sum_{j=1}^m \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} | \mathbf{x}_j) \log P(\mathbf{z}, \mathbf{x}_j | \theta)$$

- We are optimizing a lower bound!

EM Algorithm: Pictorial View



EM Pictorially

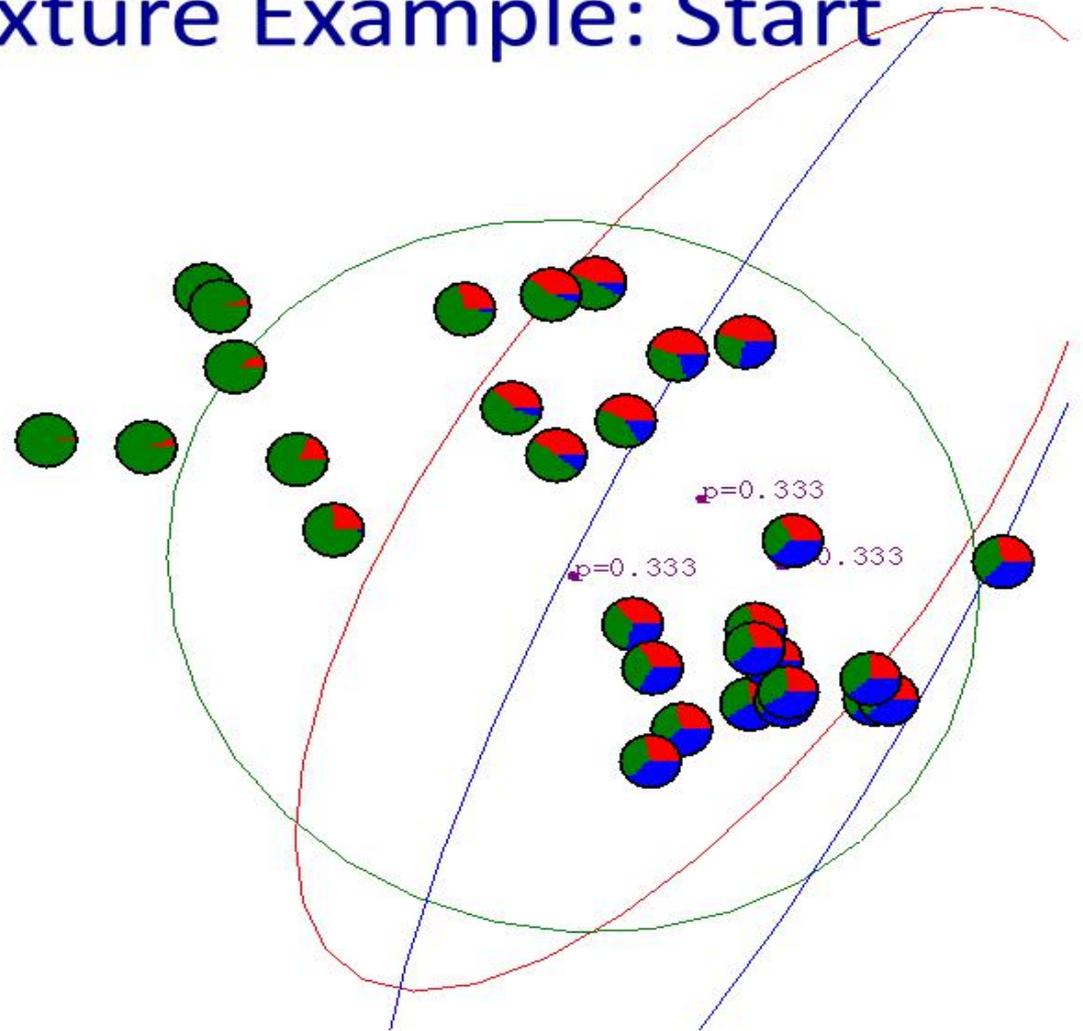


(Figure from tutorial by Sean Borman)

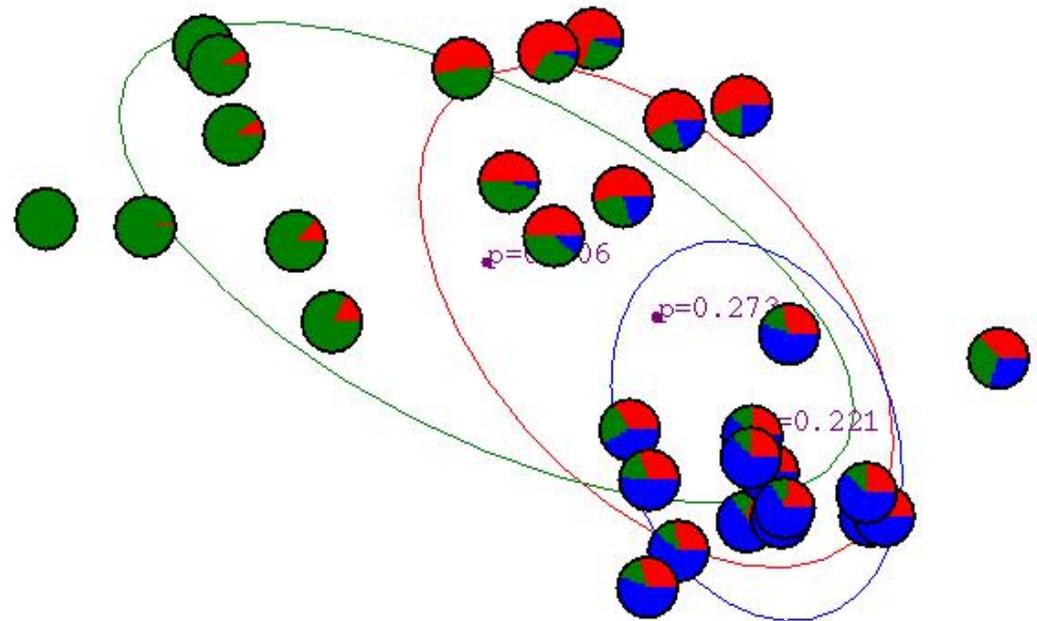
Outline

- Parametric Unsupervised Learning
- Mixture Density Model
- Gaussian Mixture Model
- Expectation Maximization (EM) Algorithm
- **Applications**

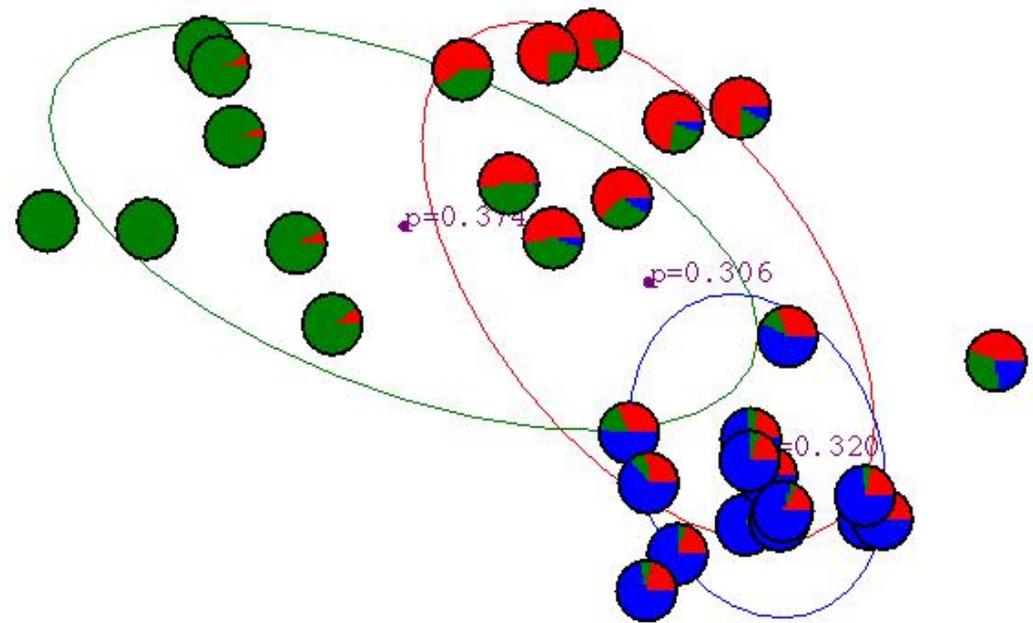
Gaussian Mixture Example: Start



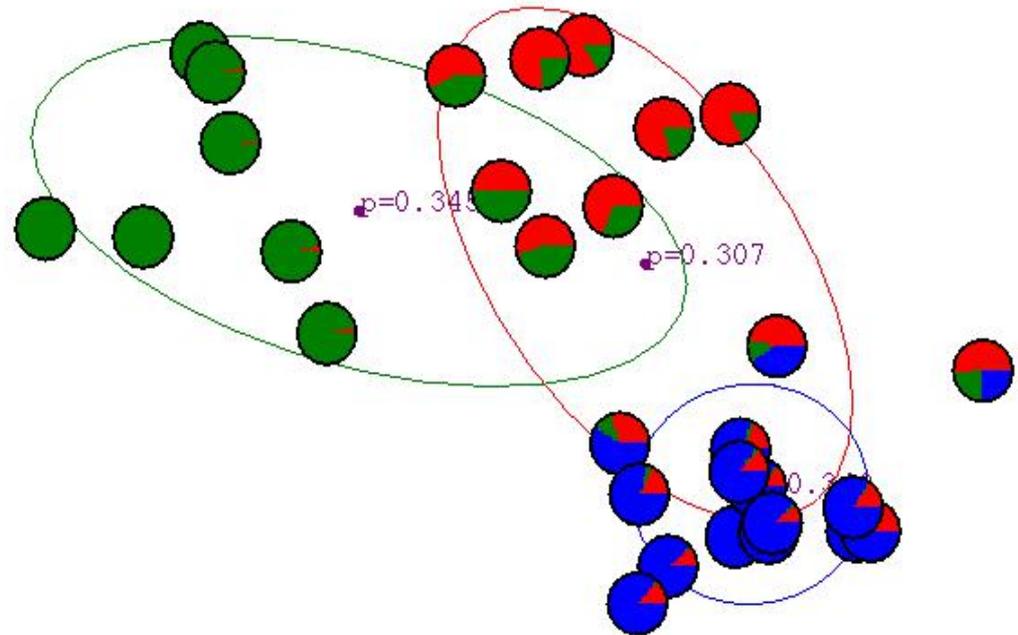
After first iteration



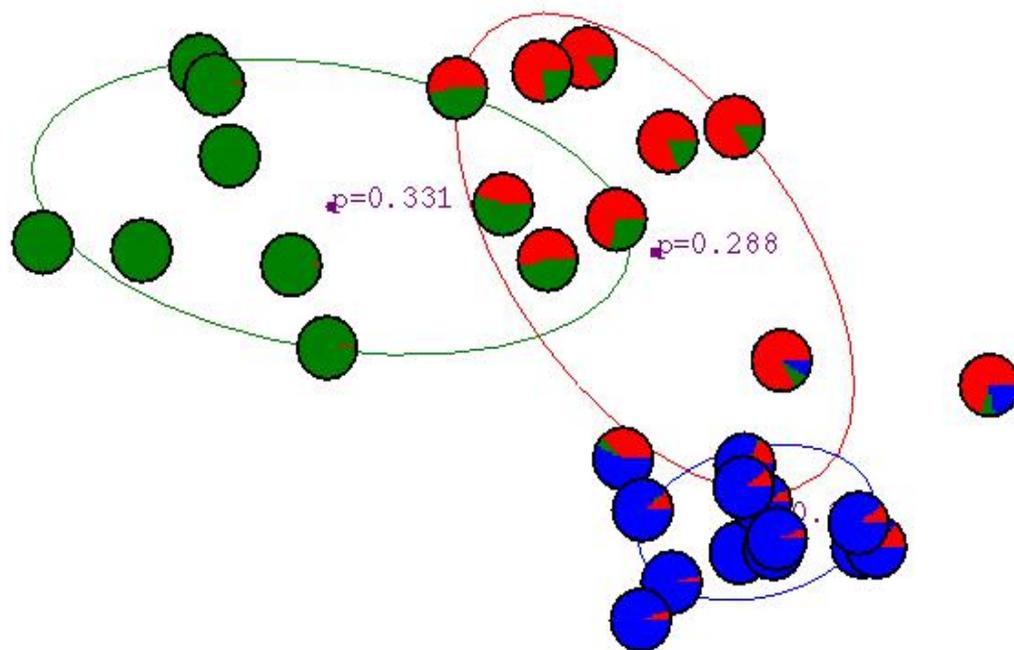
After 2nd iteration



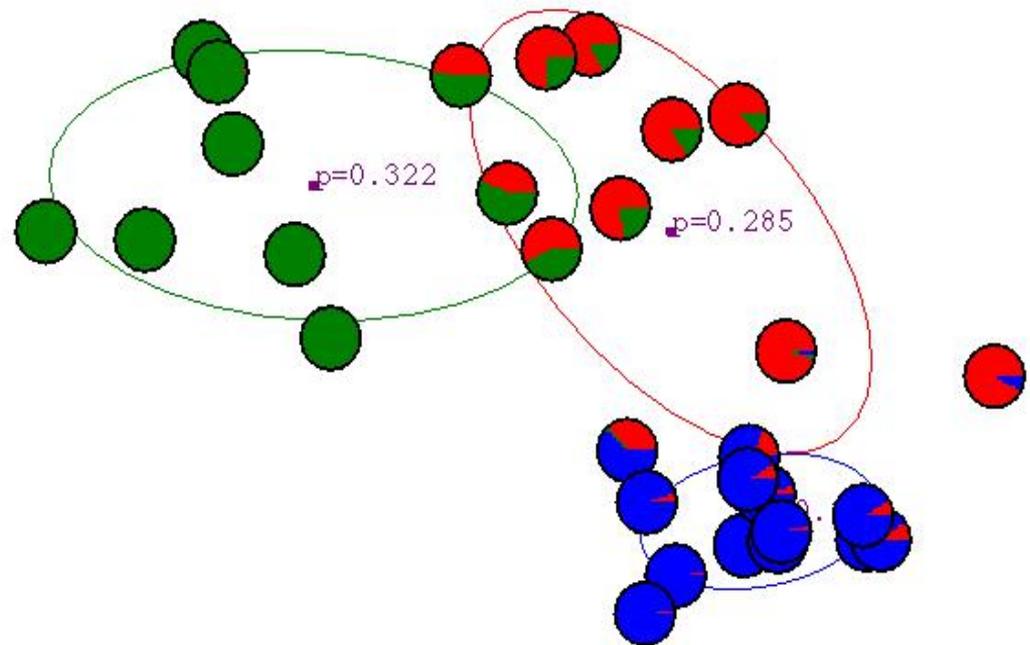
After 3rd iteration



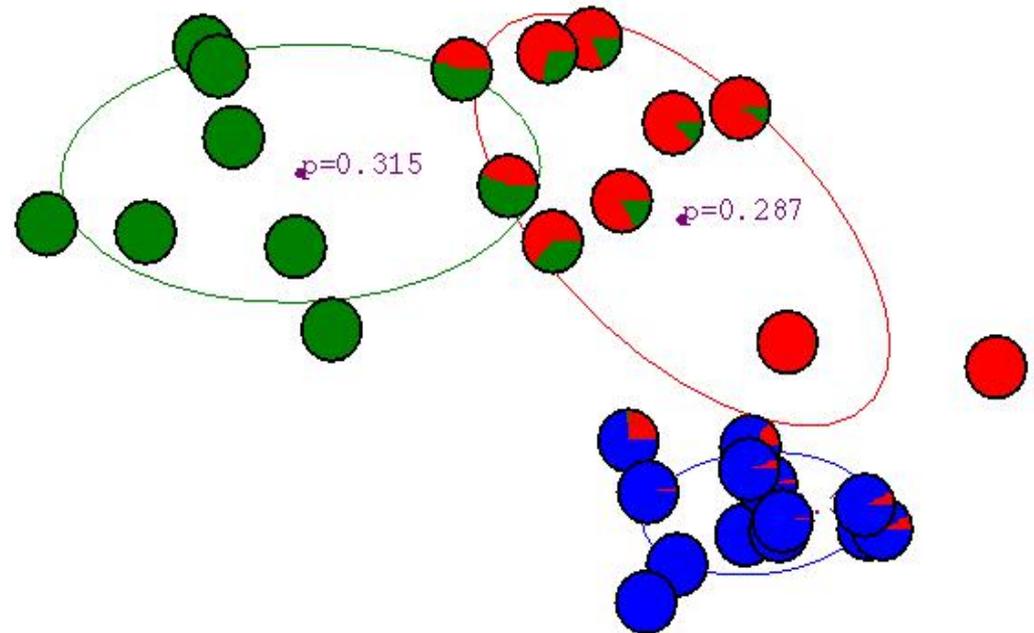
After 4th iteration



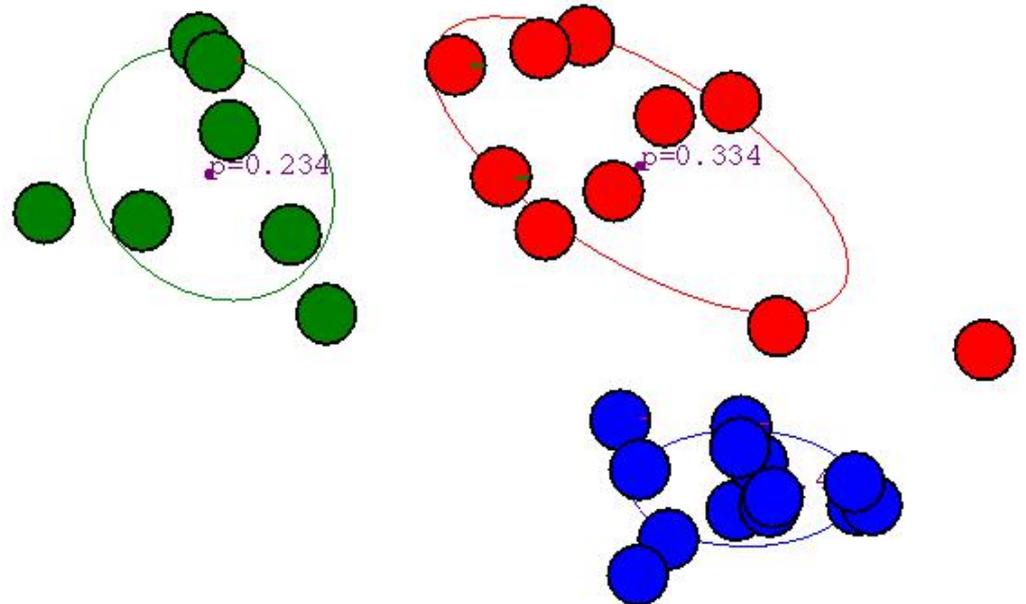
After 5th iteration



After 6th iteration



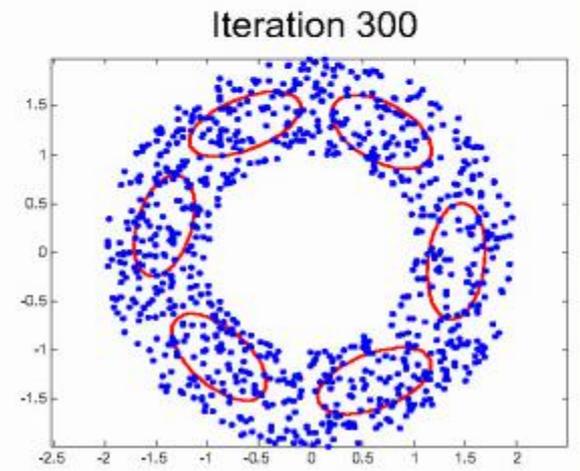
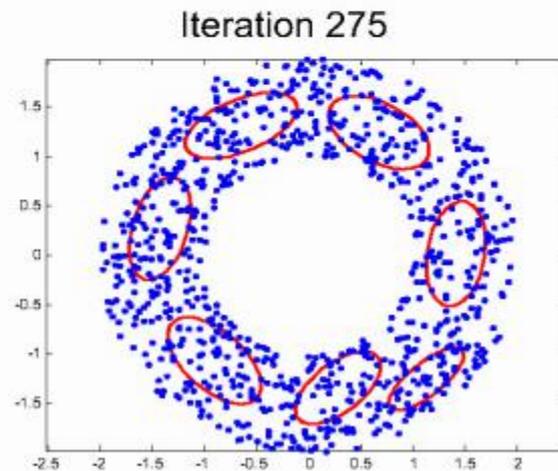
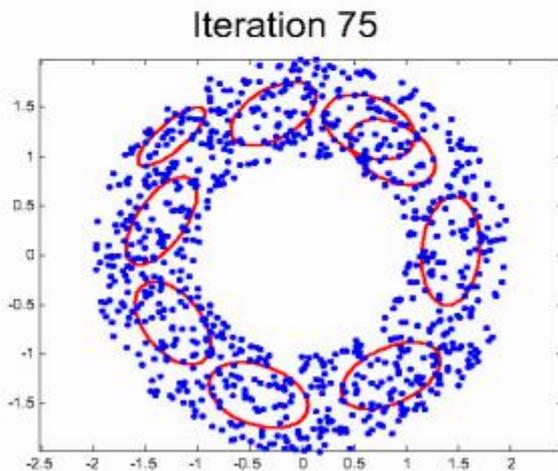
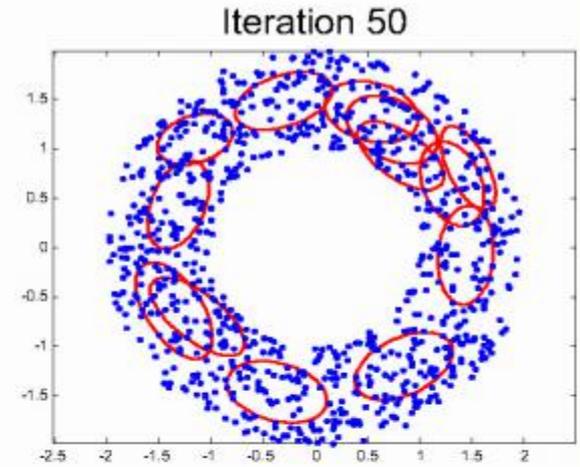
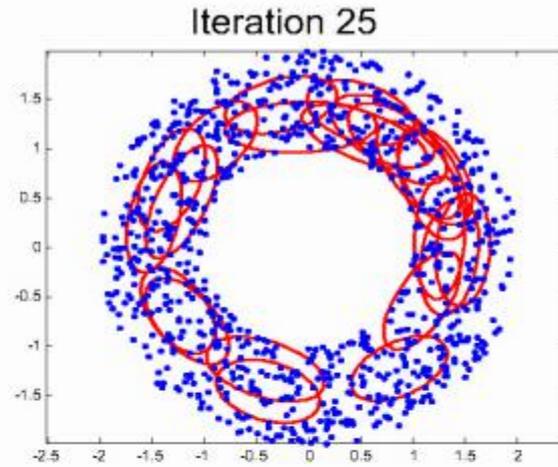
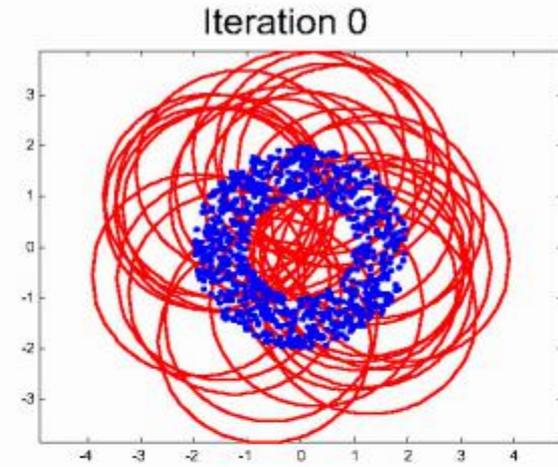
After 20th iteration



EM Example

- Example from R. Gutierrez-Osuna
- Training set of 900 examples forming an annulus
- Mixture model with $m = 30$ Gaussian components of unknown mean and variance is used
- Training:
 - Initialization:
 - means to 30 random examples
 - covariance matrices initialized to be diagonal, with large variances on the diagonal (compared to the training data variance)
 - During EM training, components with small mixing coefficients were trimmed
 - This is a trick to get in a more compact model, with fewer than 30 Gaussian components

EM Example



from R. Gutierrez-Osuna

EM Texture Segmentation Example

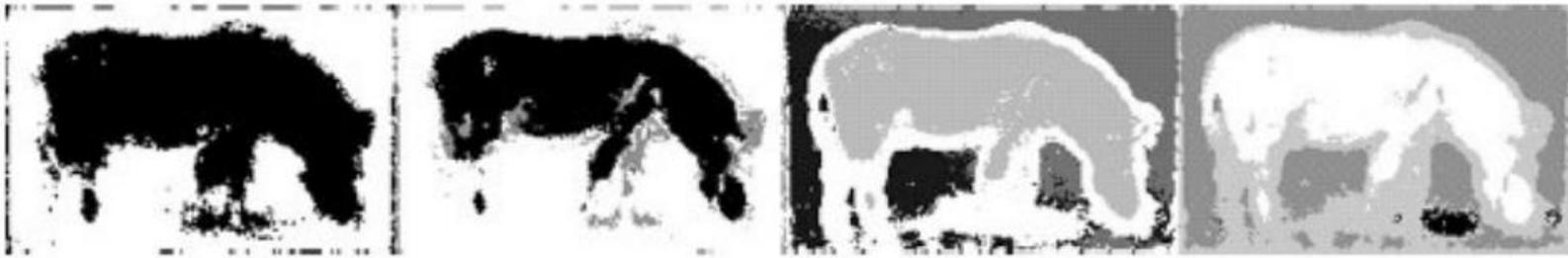
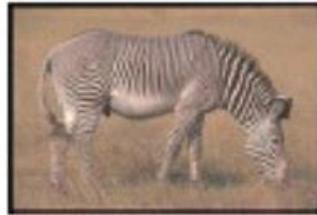


Figure from “Color and Texture Based Image Segmentation Using EM and Its Application to Content Based Image Retrieval”, S.J. Belongie et al., ICCV 1998

EM Motion Segmentation Example

Three frames from the MPEG “flower garden” sequence

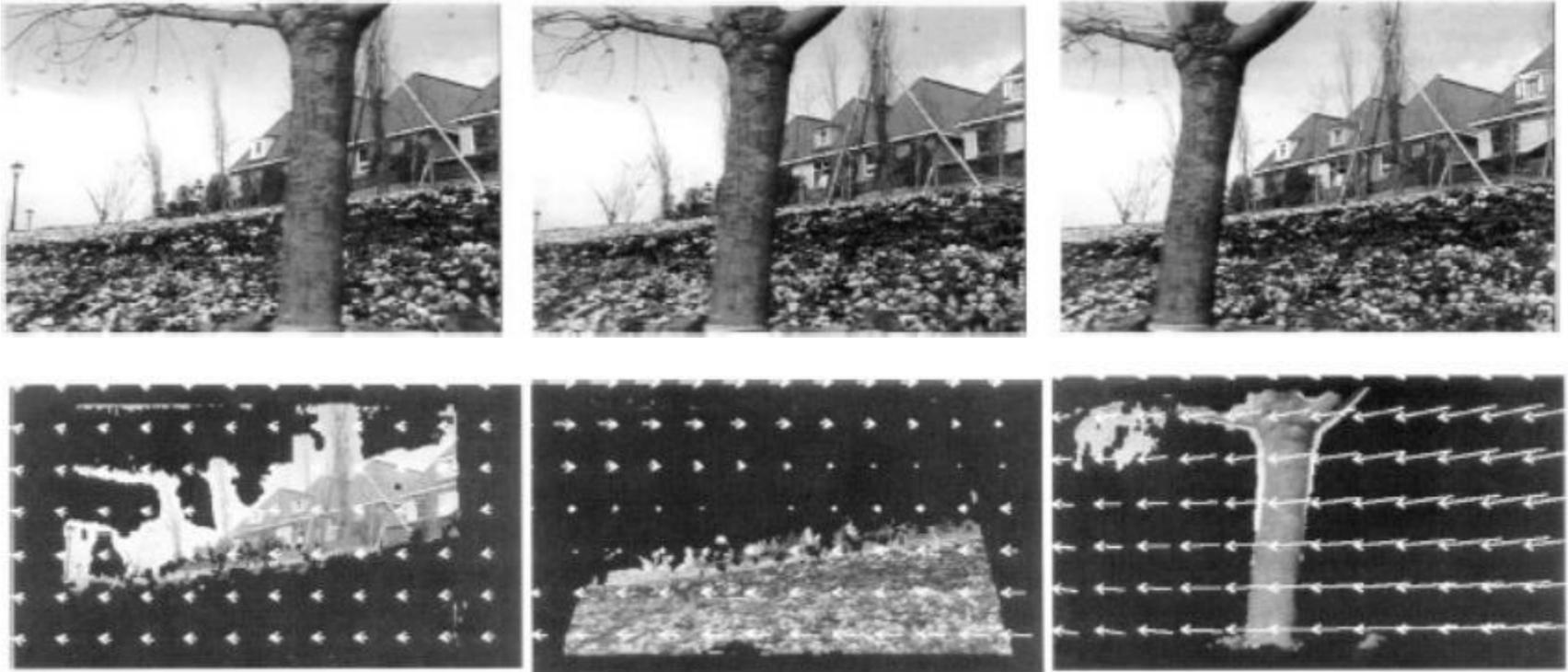


Figure from “Representing Images with layers,” by J. Wang and E.H. Adelson, IEEE Transactions on Image Processing, 1994.

Summary

- Advantages
 - If the assumed data distribution is correct, the algorithm works well
- Disadvantages
 - If assumed data distribution is wrong, results can be quite bad.
 - In particular, bad results if use incorrect number of classes (i.e. the number of mixture components)

What You Should Know

- Mixture of Gaussians
- EM for mixture of Gaussians:
 - – Coordinate ascent, just like k-means
 - – How to “learn” maximum likelihood parameters (locally max. like.) in the case of unlabeled data
 - – Relation to k-means
- Hard / soft clustering
- Probabilistic model
- Remember, EM can get stuck in local minima,
 - – And empirically it DOES

Q & A