

# Gaussian Mixtures and the EM Algorithm

Reading: Chapter 7.4, Prince book

# Review: The Gaussian Distribution

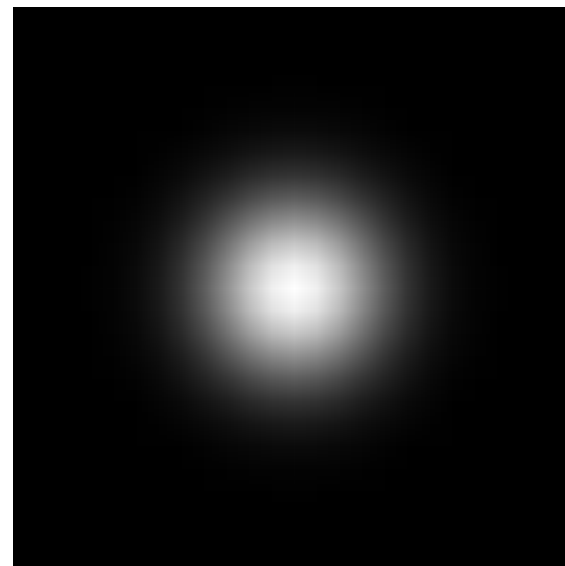
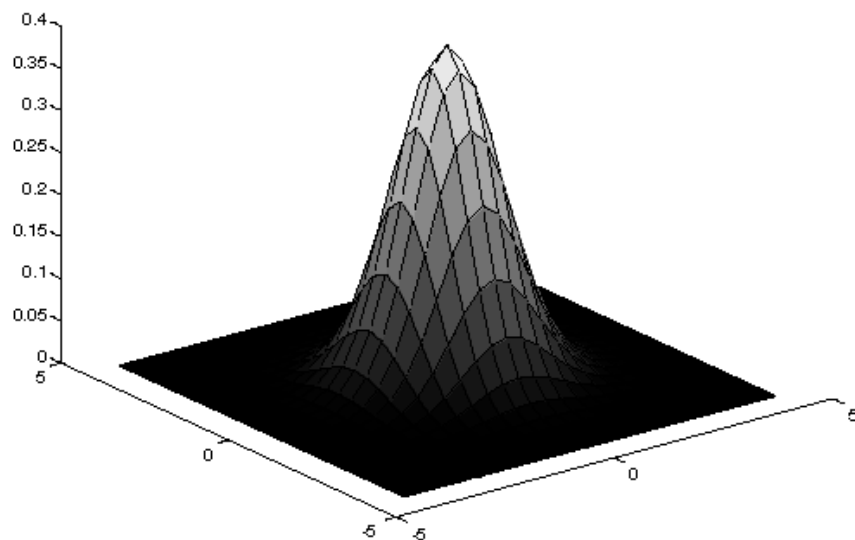
- Multivariate Gaussian

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{((2\pi)^d |\boldsymbol{\Sigma}|)^{1/2}} \exp \left\{ \underbrace{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})}_{\text{Evaluates to a number}} \right\}$$

mean  
dx1 vector

covariance  
dxd matrix

Isotropic (spherical) if covariance is  $\text{diag}(\sigma^2, \sigma^2, \dots, \sigma^2)$



# Likelihood Function

- Data set

$$D = \{\mathbf{x}_n\} \quad n = 1, \dots, N$$

- Assume observed data points generated independently

$$p(D|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- Viewed as a function of the parameters, this is known as the *likelihood function*

# Maximum Likelihood

- Set the parameters by maximizing the likelihood function
- Equivalently maximize the log likelihood

$$\ln p(D|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{N}{2} \ln(2\pi) - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

# Maximum Likelihood Solution

- Maximizing w.r.t. the mean gives the *sample mean*

$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

- Maximizing w.r.t covariance gives the *sample covariance*

$$\Sigma_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mu_{\text{ML}})(\mathbf{x}_n - \mu_{\text{ML}})^{\top}$$

**Note:** if N is small you want to divide by N-1 when computing sample covariance to get an unbiased estimate.

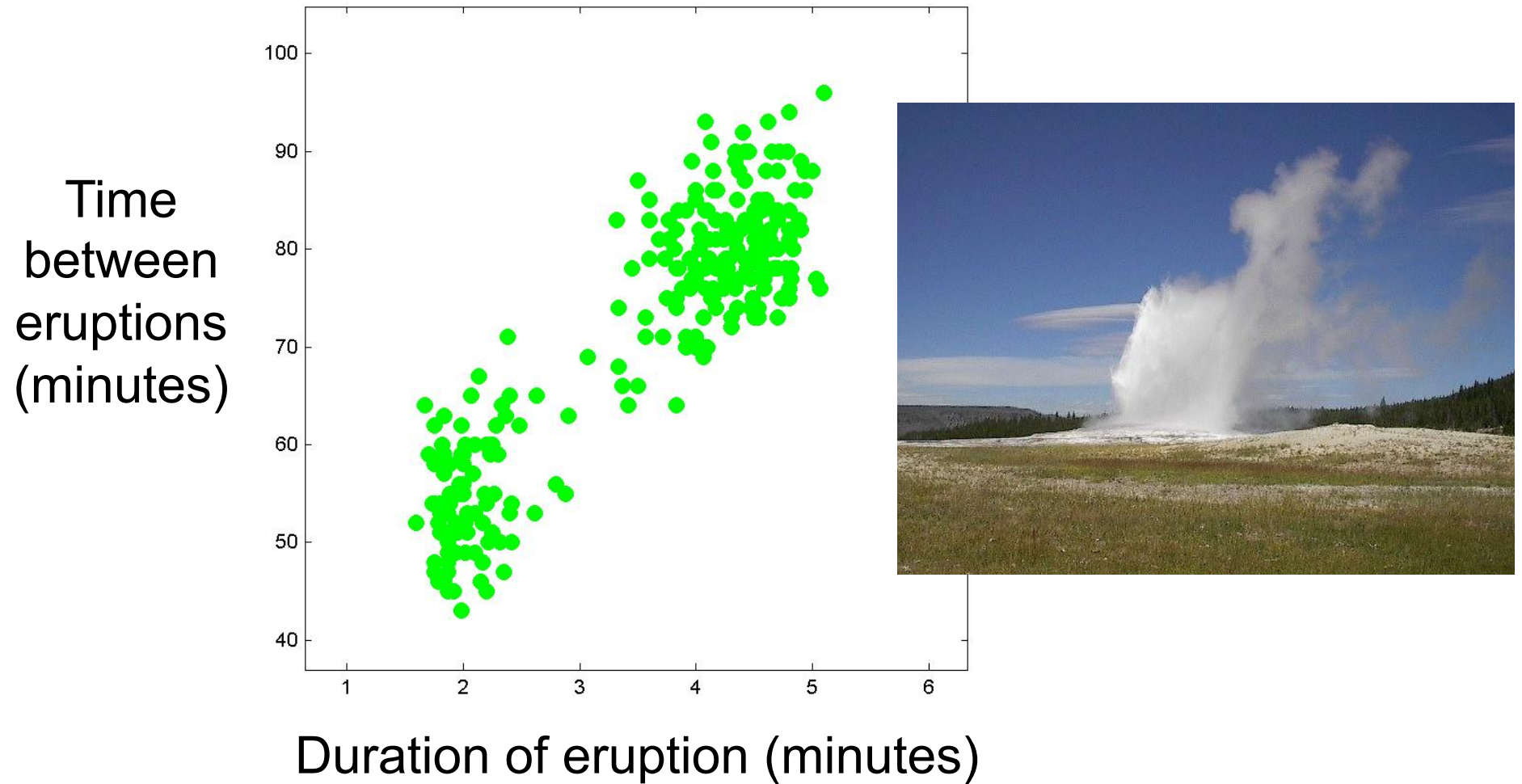
# Comments

Gaussians are well understood and easy to estimate

However, they are unimodal, thus cannot be used to represent inherently multimodal datasets

Fitting a single Gaussian to a multimodal dataset is likely to give a mean value in an area with low probability, and to overestimate the covariance.

# Old Faithful Data Set



# Idea: Use a Mixture of Gaussians

- Convex Combination of Distributions

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- Normalization and positivity require

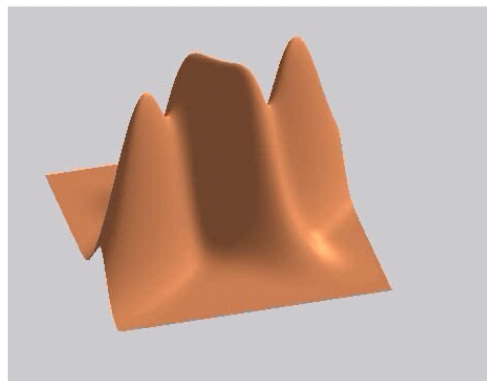
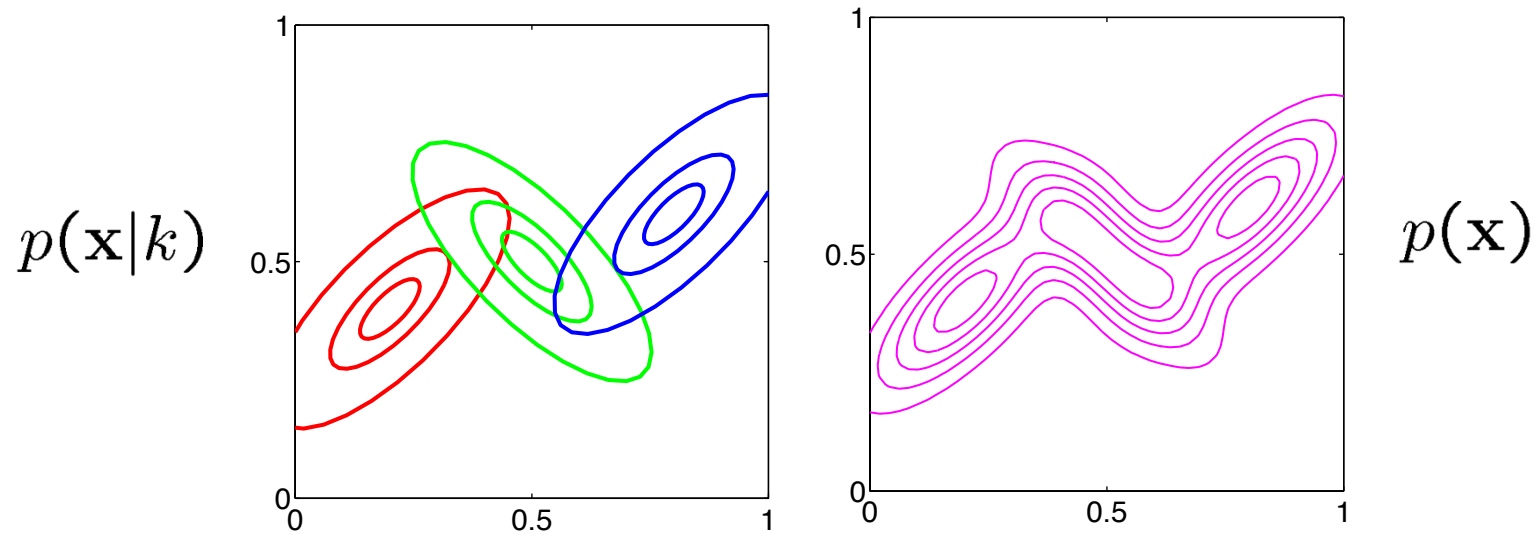
$$\sum_{k=1}^K \pi_k = 1 \quad 0 \leq \pi_k \leq 1$$

- Can interpret the mixing coefficients as prior probabilities

$$p(\mathbf{x}) = \sum_{k=1}^K p(k) p(\mathbf{x} | k)$$



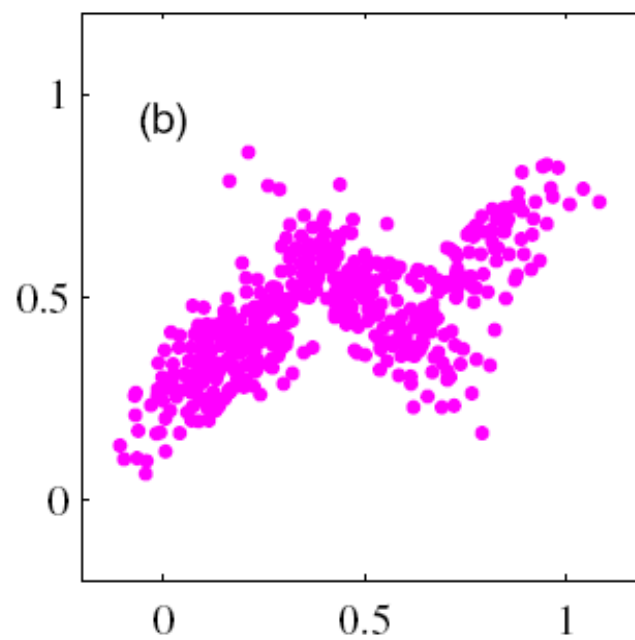
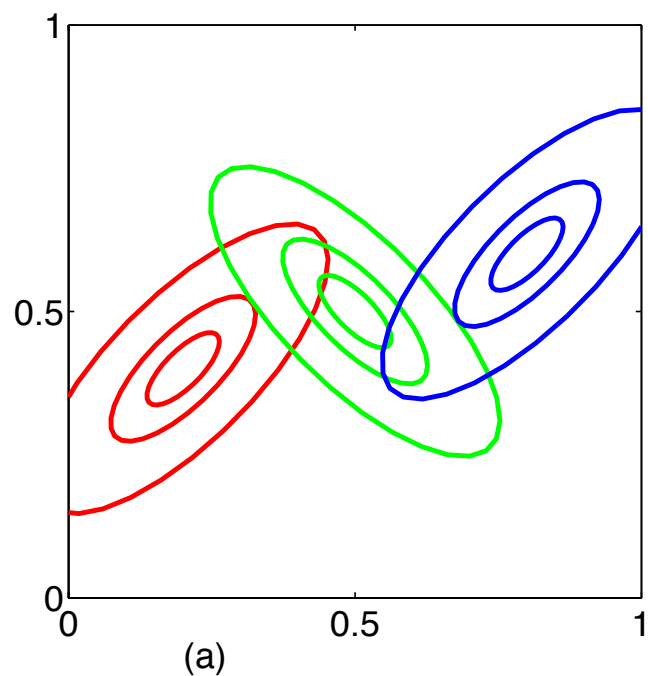
# Example: Mixture of 3 Gaussians



$p(\mathbf{x})$

## Aside: Sampling from a Mixture Model

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$



## Aside: Sampling from a Mixture Model

Generate  $u$  = uniform random number between 0 and 1

If  $u < \pi_1$

    generate  $x \sim N(x \mid \mu_1, \Sigma_1)$

elseif  $u < \pi_1 + \pi_2$

    generate  $x \sim N(x \mid \mu_2, \Sigma_2)$

$\vdots$

elseif  $u < \pi_1 + \pi_2 + \dots + \pi_{K-1}$

    generate  $x \sim N(x \mid \mu_{K-1}, \Sigma_{K-1})$

else

    generate  $x \sim N(x \mid \mu_K, \Sigma_K)$

# MLE of Mixture Parameters

- However, MLE of mixture parameters is HARD!
- Joint distribution:

$$p(\mathbf{X}|\pi, \mu, \Sigma) = \prod_{n=1}^N \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)$$

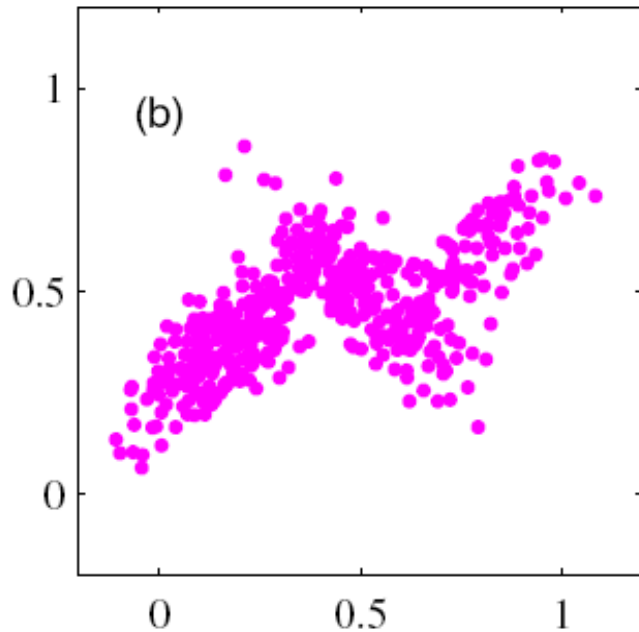
- Log likelihood

$$\ln p(\mathbf{X}|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$



Uh-oh, log of a sum

# EM Algorithm



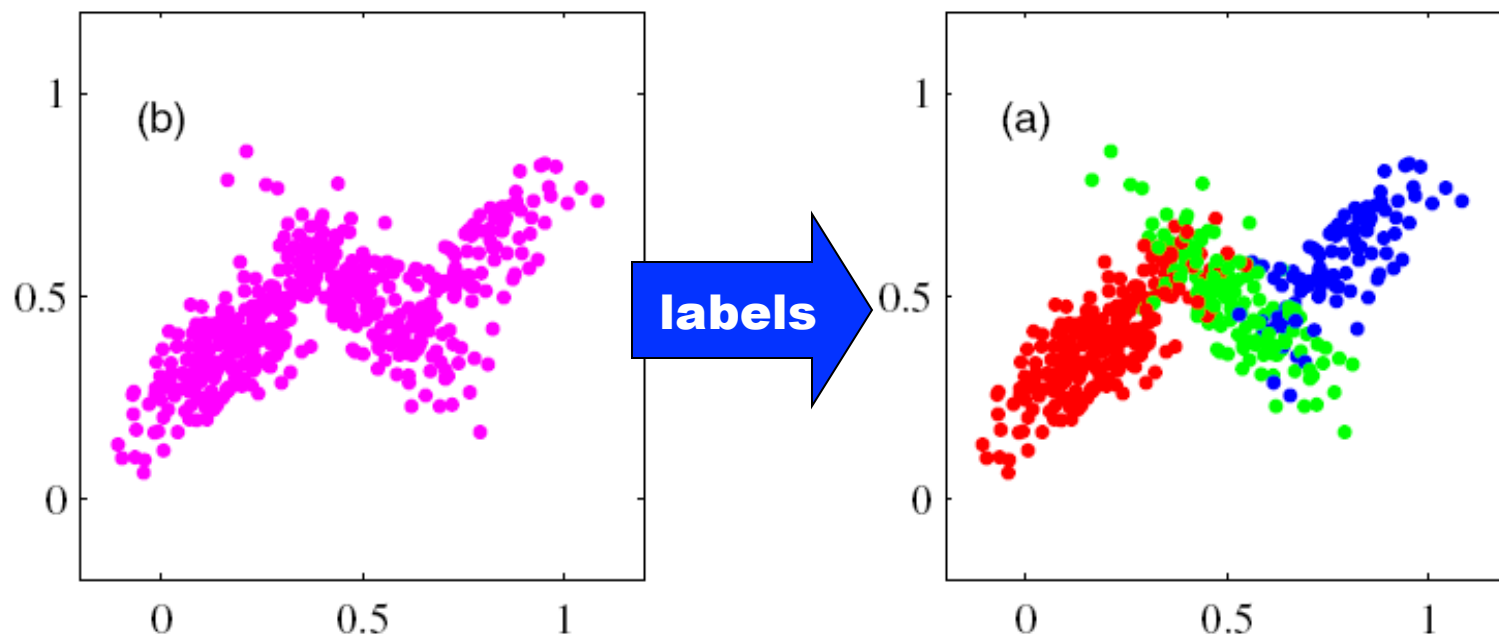
What makes this estimation problem hard?

1) It is a mixture, so log-likelihood is messy

$$\ln p(\mathbf{X}|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

2) We don't directly see what the underlying process is

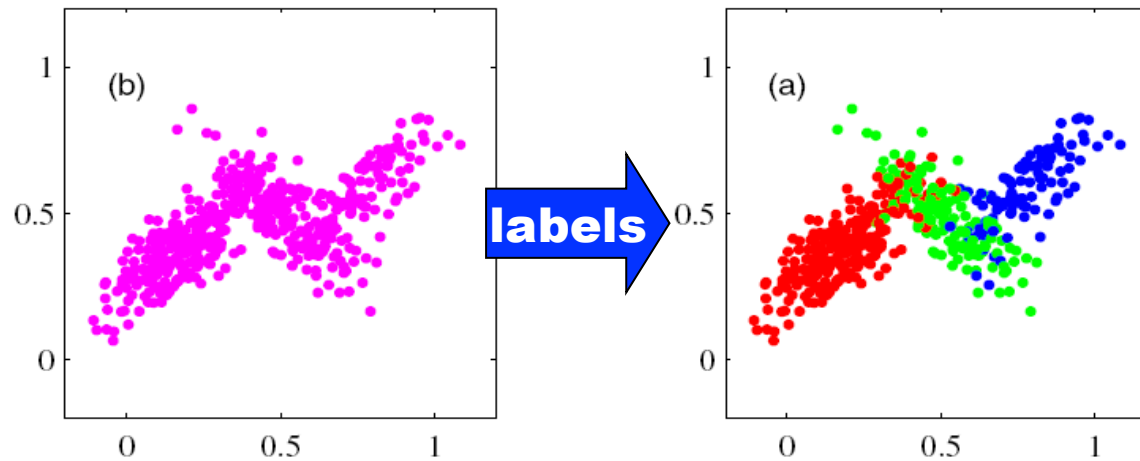
# EM Algorithm



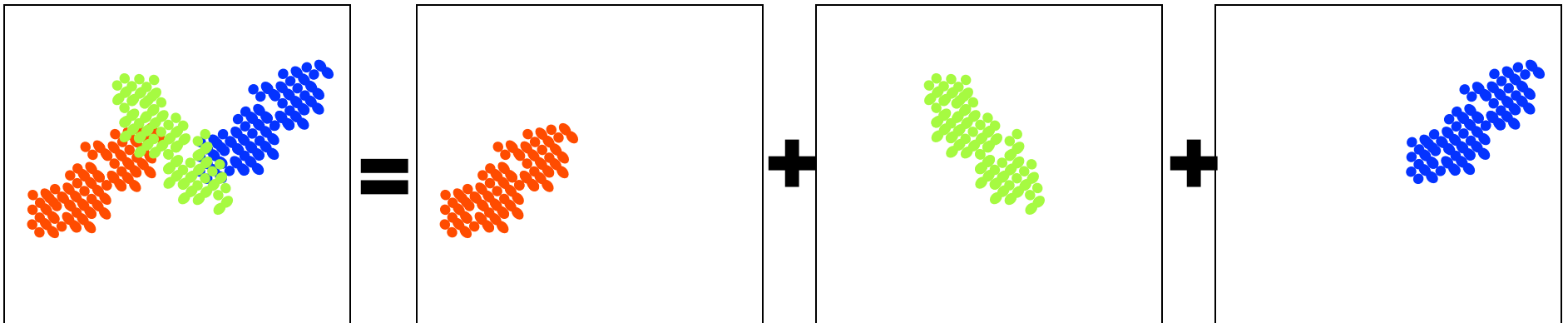
Suppose some oracle told us which point comes from which Gaussian.

How? By providing a “latent” variable  $z_{nk}$  which is 1 if point  $n$  comes from the  $k$ th component Gaussian, and 0 otherwise (a 1 of  $K$  representation)

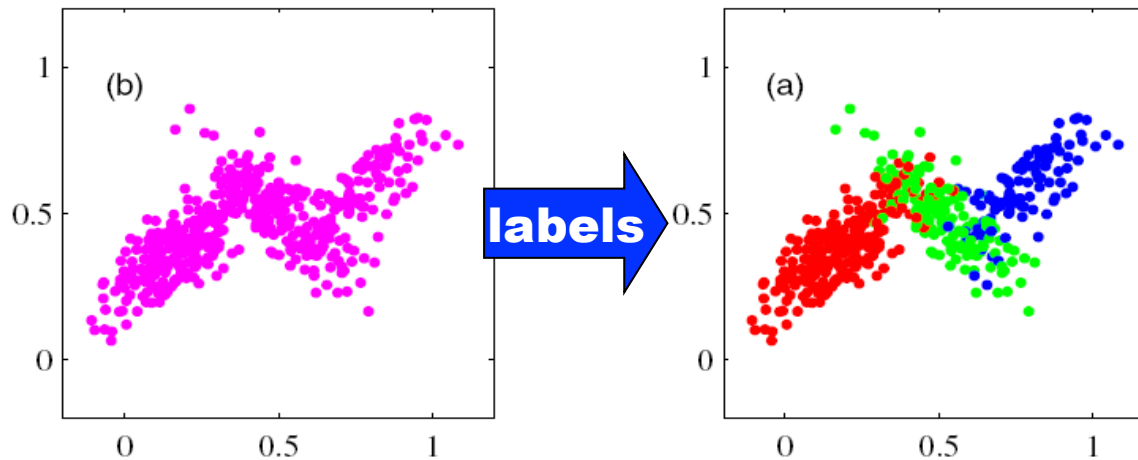
# EM Algorithm



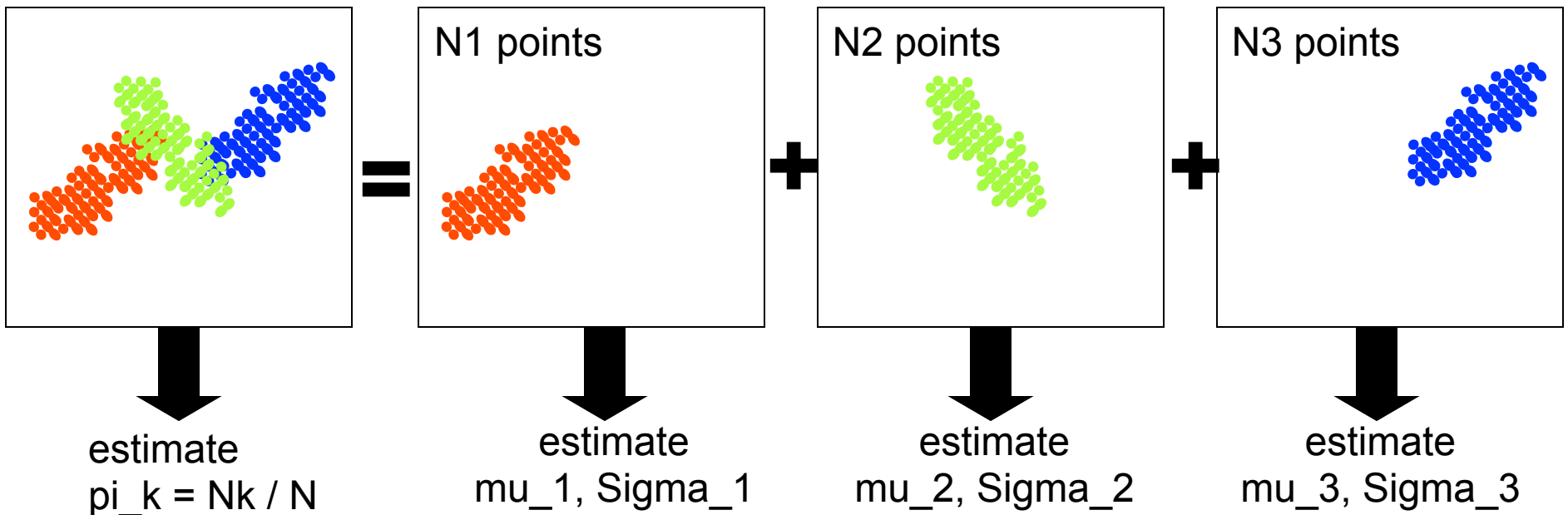
This lets us recover the underlying generating process decomposition:



# EM Algorithm



And we can easily estimate each Gaussian, along with the mixture weights!





# EM Algorithm

Remember that this was a problem...

$$\ln p(\mathbf{X}|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

**how can I make that  
inner sum be a product  
instead???**



# EM Algorithm

Remember that this was a problem...

$$\ln p(\mathbf{X}|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

Again, if an oracle gave us the values of the latent variables (component that generated each point) we could work with the complete log likelihood

$$p(\mathbf{X}, \mathbf{Z} | \mu, \Sigma, \pi) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)^{z_{nk}}$$

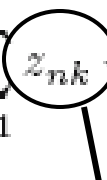
and the log of that looks much better!

$$\ln p(\mathbf{X}, \mathbf{Z} | \mu, \Sigma, \pi) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \{ \ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \}.$$

how can I make that inner sum be a product instead???

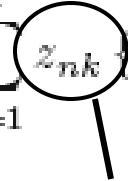


# Latent Variable View

$$\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \{ \ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \}.$$


note: for a given  $n$ , there are  $K$  of these latent variables,  
and only ONE of them is 1 (all the rest are 0)

# Latent Variable View

$$\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \{ \ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \}.$$


note: for a given  $n$ , there are  $K$  of these latent variables, and only ONE of them is 1 (all the rest are 0)

This is thus equivalent to

$$\begin{aligned} & \sum_{\substack{\text{all } n \text{ for which} \\ z_{n,1}=1}} \ln \pi_1 + \ln \mathcal{N}(x_n | \mu_1, \Sigma_1) \\ & + \sum_{\substack{\text{all } n \text{ for which} \\ z_{n,2}=1}} \ln \pi_2 + \ln \mathcal{N}(x_n | \mu_2, \Sigma_2) \quad + \quad \dots \quad + \\ & + \sum_{\substack{\text{all } n \text{ for which} \\ z_{n,K}=1}} \ln \pi_K + \ln \mathcal{N}(x_n | \mu_K, \Sigma_K) \end{aligned}$$

# Latent Variable View

$$\begin{aligned} & \sum_{\substack{\text{all } n \text{ for which} \\ z_{n,1}=1}} \ln \pi_1 + \sum_{\substack{\text{all } n \text{ for which} \\ z_{n,1}=1}} \ln \mathcal{N}(x_n | \mu_1, \Sigma_1) \\ & + \sum_{\substack{\text{all } n \text{ for which} \\ z_{n,2}=1}} \ln \pi_2 + \sum_{\substack{\text{all } n \text{ for which} \\ z_{n,2}=1}} \ln \mathcal{N}(x_n | \mu_2, \Sigma_2) \\ & + \dots + \\ & + \sum_{\substack{\text{all } n \text{ for which} \\ z_{n,K}=1}} \ln \pi_K + \sum_{\substack{\text{all } n \text{ for which} \\ z_{n,K}=1}} \ln \mathcal{N}(x_n | \mu_K, \Sigma_K) \end{aligned}$$

# Latent Variable View

$$\begin{aligned} & \sum_{\substack{\text{all } n \text{ for which} \\ z_{n,1}=1}} \ln \pi_1 + \boxed{\sum_{\substack{\text{all } n \text{ for which} \\ z_{n,1}=1}} \ln \mathcal{N}(x_n | \mu_1, \Sigma_1)} \quad \begin{array}{l} \text{can be} \\ \text{estimated} \\ \text{separately} \end{array} \\ & + \sum_{\substack{\text{all } n \text{ for which} \\ z_{n,2}=1}} \ln \pi_2 + \boxed{\sum_{\substack{\text{all } n \text{ for which} \\ z_{n,2}=1}} \ln \mathcal{N}(x_n | \mu_2, \Sigma_2)} \quad \begin{array}{l} \text{can be} \\ \text{estimated} \\ \text{separately} \end{array} \\ & + \dots + \\ & + \sum_{\substack{\text{all } n \text{ for which} \\ z_{n,K}=1}} \ln \pi_K + \boxed{\sum_{\substack{\text{all } n \text{ for which} \\ z_{n,K}=1}} \ln \mathcal{N}(x_n | \mu_K, \Sigma_K)} \quad \begin{array}{l} \text{can be} \\ \text{estimated} \\ \text{separately} \end{array} \end{aligned}$$

# Latent Variable View

$$\begin{aligned} & \sum_{\text{all } n \text{ for which } z_{n,1}=1} \ln \pi_1 + \sum_{\text{all } n \text{ for which } z_{n,1}=1} \ln \mathcal{N}(x_n | \mu_1, \Sigma_1) && \text{can be estimated separately} \\ + & \sum_{\text{all } n \text{ for which } z_{n,2}=1} \ln \pi_2 + \sum_{\text{all } n \text{ for which } z_{n,2}=1} \ln \mathcal{N}(x_n | \mu_2, \Sigma_2) && \text{can be estimated separately} \\ + & \dots + \\ + & \sum_{\text{all } n \text{ for which } z_{n,K}=1} \ln \pi_K + \sum_{\text{all } n \text{ for which } z_{n,K}=1} \ln \mathcal{N}(x_n | \mu_K, \Sigma_K) && \text{can be estimated separately} \end{aligned}$$

these are coupled because the mixing weights all sum to 1, but it is no big deal to solve

# EM Algorithm

Unfortunately, oracles don't exist (or if they do, they won't talk to us)

So we don't know values of the  $z_{nk}$  variables

What EM proposes to do:

- 1) compute  $p(Z|X, \theta)$ , the posterior distribution over  $z_{nk}$ , given our current best guess at the values of  $\theta$
- 2) compute the expected value of the log likelihood  $\ln(p(X, Z|\theta))$  with respect to the distribution  $p(Z|X, \theta)$
- 3) find  $\theta_{\text{new}}$  that maximizes that function.  
This is our new best guess at the values of  $\theta$ .
- 4) iterate...



# Insight

Since we don't know the latent variables, we instead take the expected value of the log likelihood with respect to their posterior distribution  $P(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})$ . In the GMM case, this is equivalent to “softening” the binary latent variables to continuous ones (the expected values of the latent variables)

$$\ln p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta}) = \sum_{n=1}^N \sum_{k=1}^K \underbrace{z_{nk}}_{\text{unknown discrete value 0 or 1}} \{ \ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \}$$

unknown discrete value 0 or 1

$$E_{\mathbf{z}}[\ln p(\mathbf{x}, \mathbf{z}|\boldsymbol{\theta})] = \sum_{n=1}^N \sum_{k=1}^K \underbrace{\gamma_k(\mathbf{x}_n)}_{\text{known continuous value between 0 and 1}} \{ \ln \pi_k + \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \}$$

known continuous value between 0 and 1

Where  $\gamma_j(\mathbf{x}_n)$  is  $P(z_{nk} = 1)$

# Insight

So now, after replacing the binary latent variables with their continuous expected values:

all points contribute to the estimation of all components

each point has unit mass to contribute, but splits it across the  $K$  components

the amount of weight a point contributes to a component is proportional to the relative likelihood that the point was generated by that component

# Latent Variable View (with an oracle)

$$\begin{array}{l} \boxed{\sum_{\substack{\text{all } n \text{ for which} \\ z_{n,1}=1}} \ln \pi_1} + \boxed{\sum_{\substack{\text{all } n \text{ for which} \\ z_{n,1}=1}} \ln \mathcal{N}(x_n | \mu_1, \Sigma_1)} \quad \begin{array}{l} \text{can be} \\ \text{estimated} \\ \text{separately} \end{array} \\ + \sum_{\substack{\text{all } n \text{ for which} \\ z_{n,2}=1}} \ln \pi_2 + \boxed{\sum_{\substack{\text{all } n \text{ for which} \\ z_{n,2}=1}} \ln \mathcal{N}(x_n | \mu_2, \Sigma_2)} \quad \begin{array}{l} \text{can be} \\ \text{estimated} \\ \text{separately} \end{array} \\ + \dots + \\ + \sum_{\substack{\text{all } n \text{ for which} \\ z_{n,K}=1}} \ln \pi_K + \boxed{\sum_{\substack{\text{all } n \text{ for which} \\ z_{n,K}=1}} \ln \mathcal{N}(x_n | \mu_K, \Sigma_K)} \quad \begin{array}{l} \text{can be} \\ \text{estimated} \\ \text{separately} \end{array} \end{array}$$

these are coupled because the mixing weights  
all sum to 1, but it is no big deal to solve

# Latent Variable View (with EM, $\gamma_{n,k}^i$ a constant at iteration $i$ )

$$\begin{aligned}
 & \sum_N \sum_K \gamma_{n,k}^i \ln \pi_1 + \sum_N \sum_K \gamma_{n,k}^i \ln \mathcal{N}(x_n | \mu_1, \Sigma_1) && \text{can be estimated separately} \\
 + & \sum_N \sum_K \gamma_{n,k}^i \ln \pi_2 + \sum_N \sum_K \gamma_{n,k}^i \ln \mathcal{N}(x_n | \mu_2, \Sigma_2) && \text{can be estimated separately} \\
 + & \dots + \\
 + & \sum_N \sum_K \gamma_{n,k}^i \ln \pi_K + \sum_N \sum_K \gamma_{n,k}^i \ln \mathcal{N}(x_n | \mu_K, \Sigma_K) && \text{can be estimated separately}
 \end{aligned}$$

these are coupled because the mixing weights all sum to 1, but it is no big deal to solve

# EM Algorithm for GMM

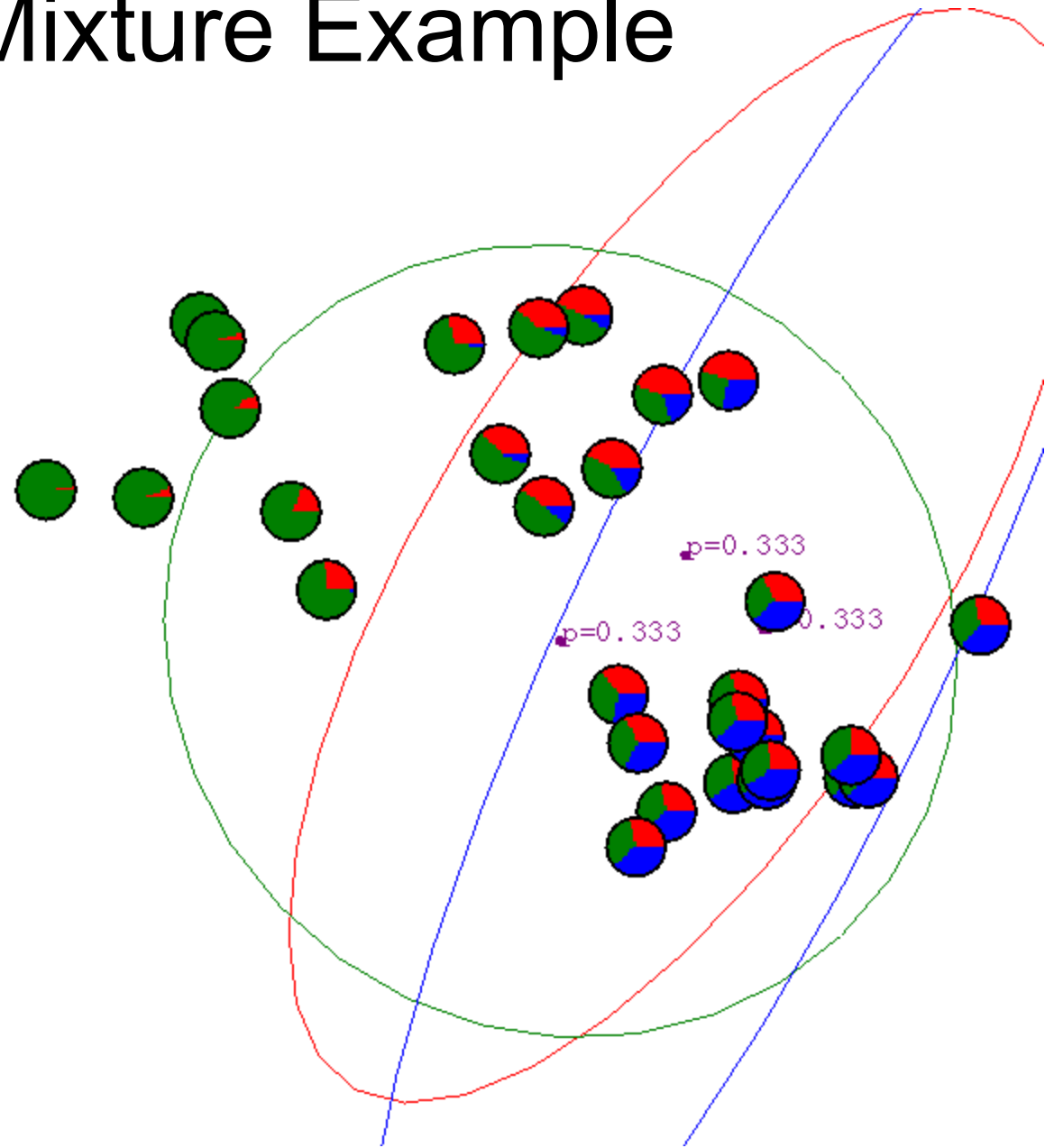
$$\mathbf{E} \quad \gamma_j(\mathbf{x}_n) = \frac{\pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \quad \text{ownership weights}$$

---

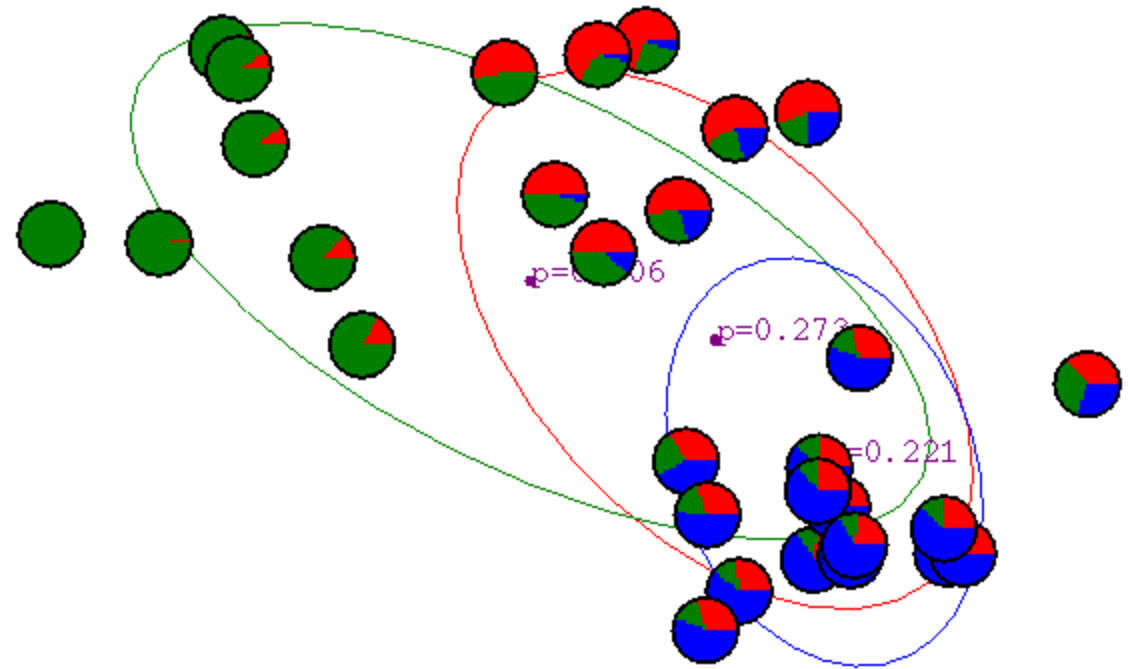
$$\mathbf{M} \quad \boldsymbol{\mu}_j = \frac{\sum_{n=1}^N \gamma_j(\mathbf{x}_n) \mathbf{x}_n}{\sum_{n=1}^N \gamma_j(\mathbf{x}_n)} \quad \text{means} \quad \boldsymbol{\Sigma}_j = \frac{\sum_{n=1}^N \gamma_j(\mathbf{x}_n) (\mathbf{x}_n - \boldsymbol{\mu}_j)(\mathbf{x}_n - \boldsymbol{\mu}_j)^\top}{\sum_{n=1}^N \gamma_j(\mathbf{x}_n)} \quad \text{covariances}$$

$$\pi_j = \frac{1}{N} \sum_{n=1}^N \gamma_j(\mathbf{x}_n) \quad \text{mixing probabilities}$$

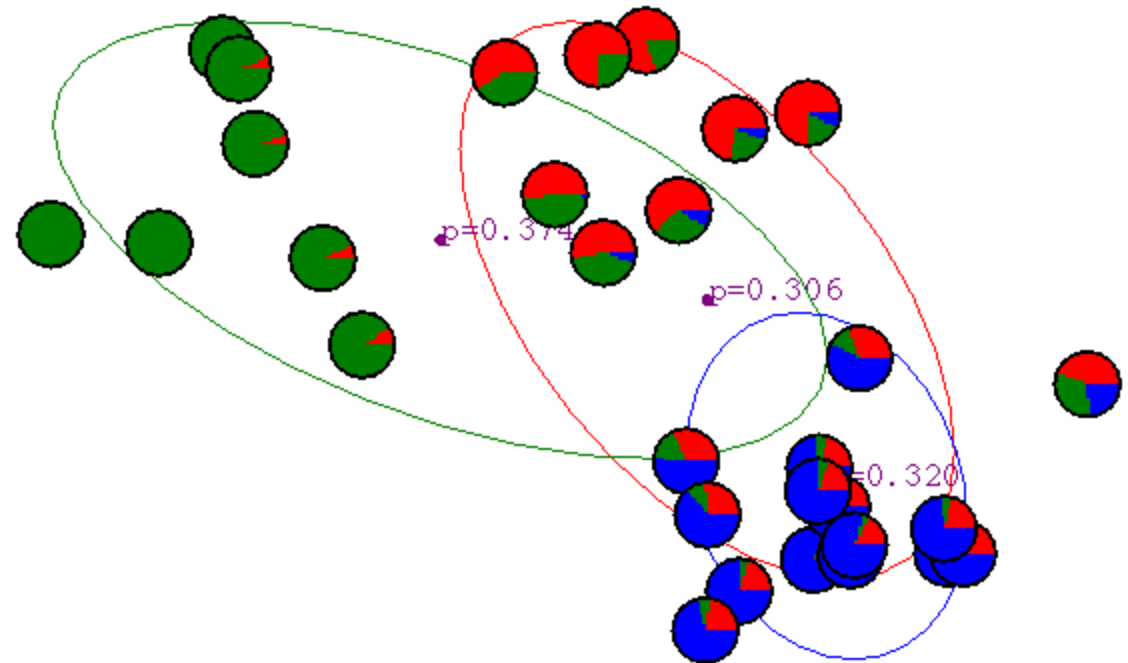
# Gaussian Mixture Example



# After first iteration

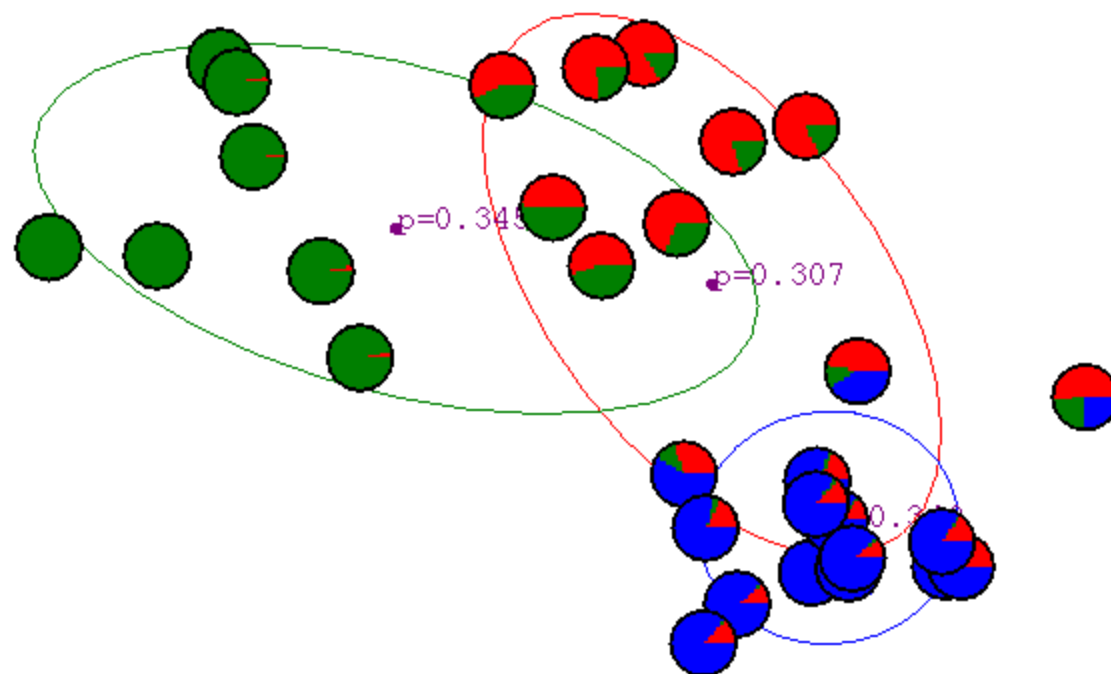


# After 2nd iteration

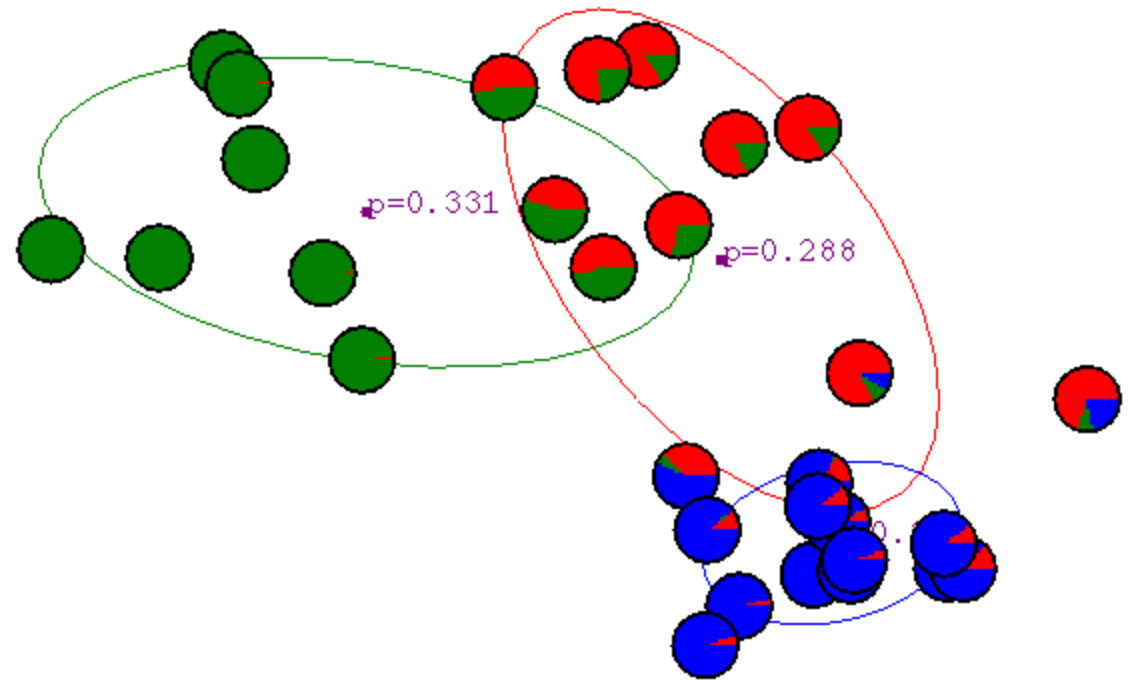




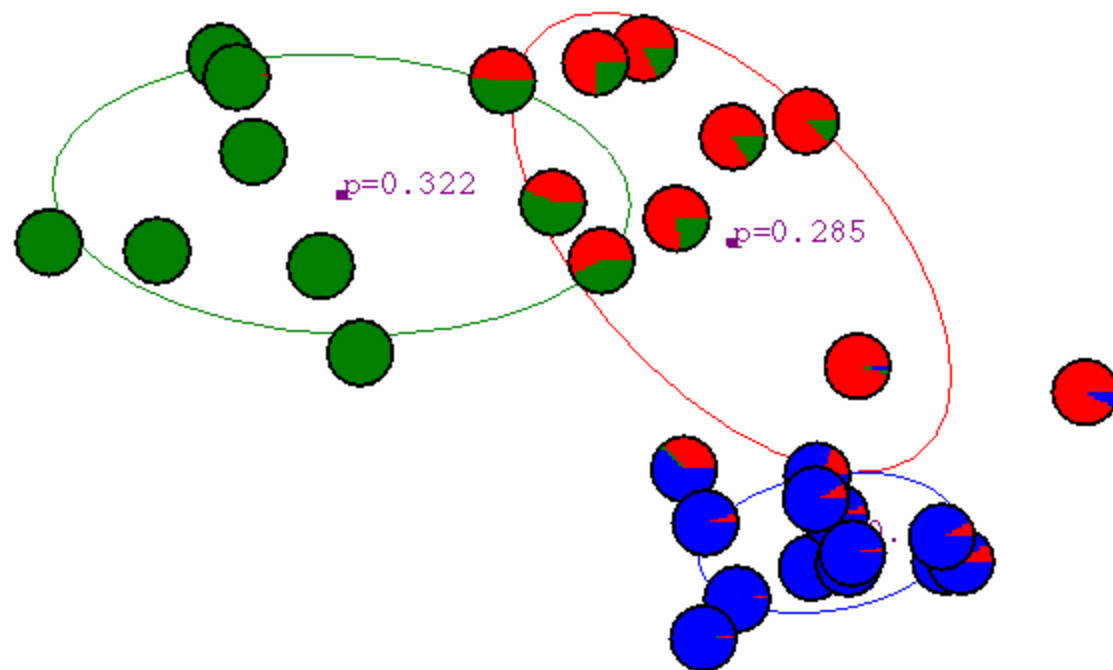
# After 3rd iteration



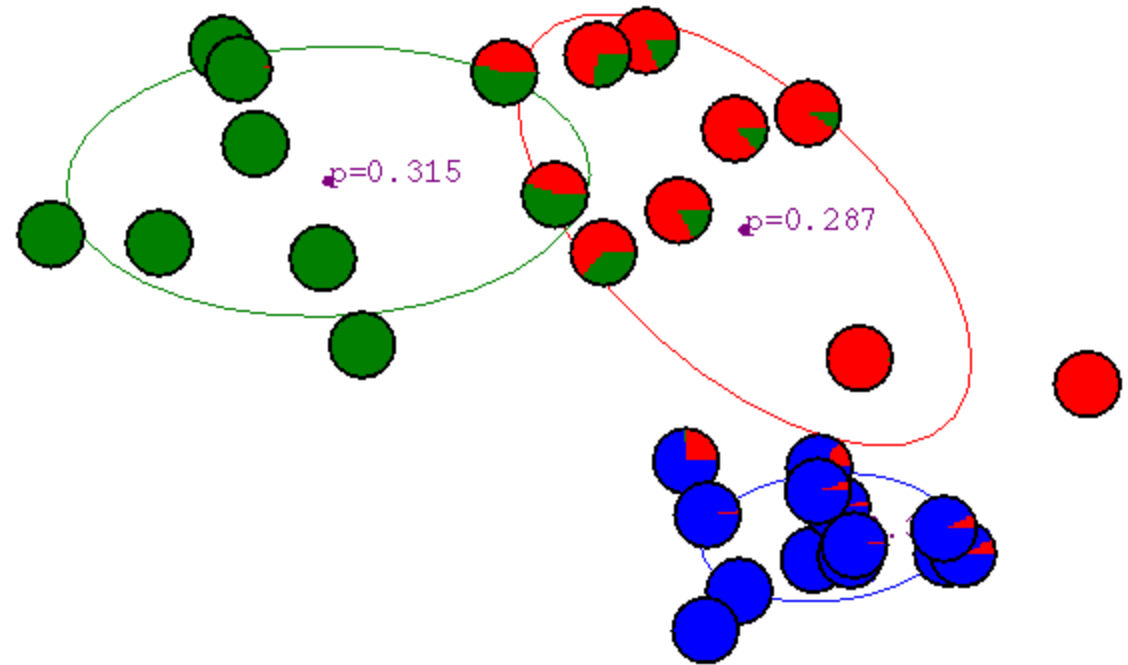
# After 4th iteration



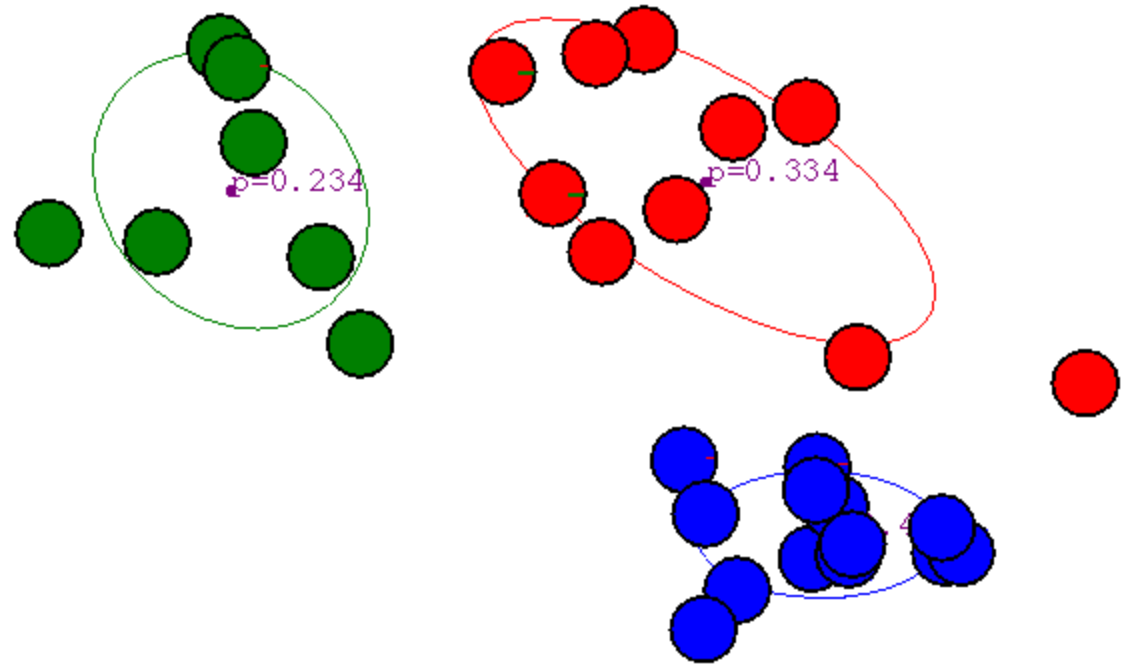
# After 5th iteration



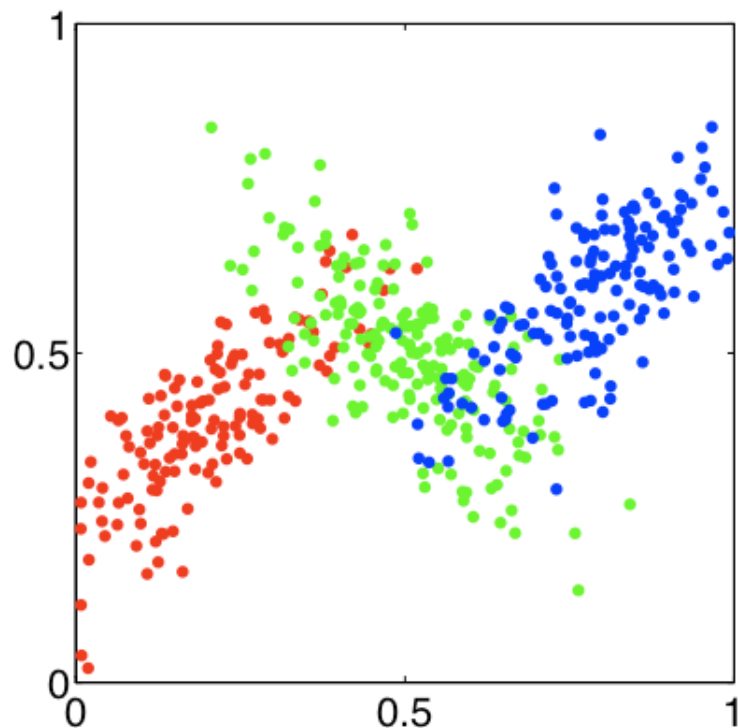
# After 6th iteration



# After 20th iteration

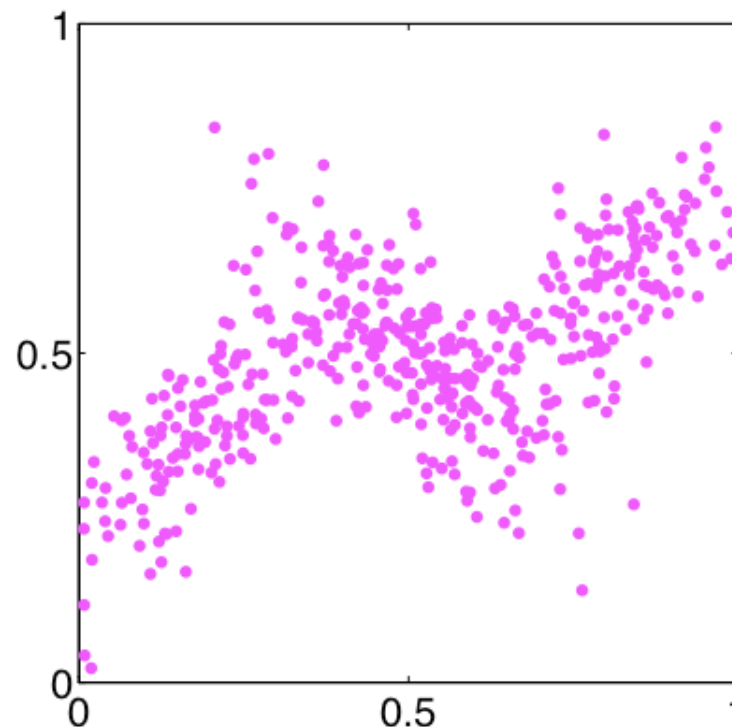


# Recall: Labeled vs Unlabeled Data



labeled

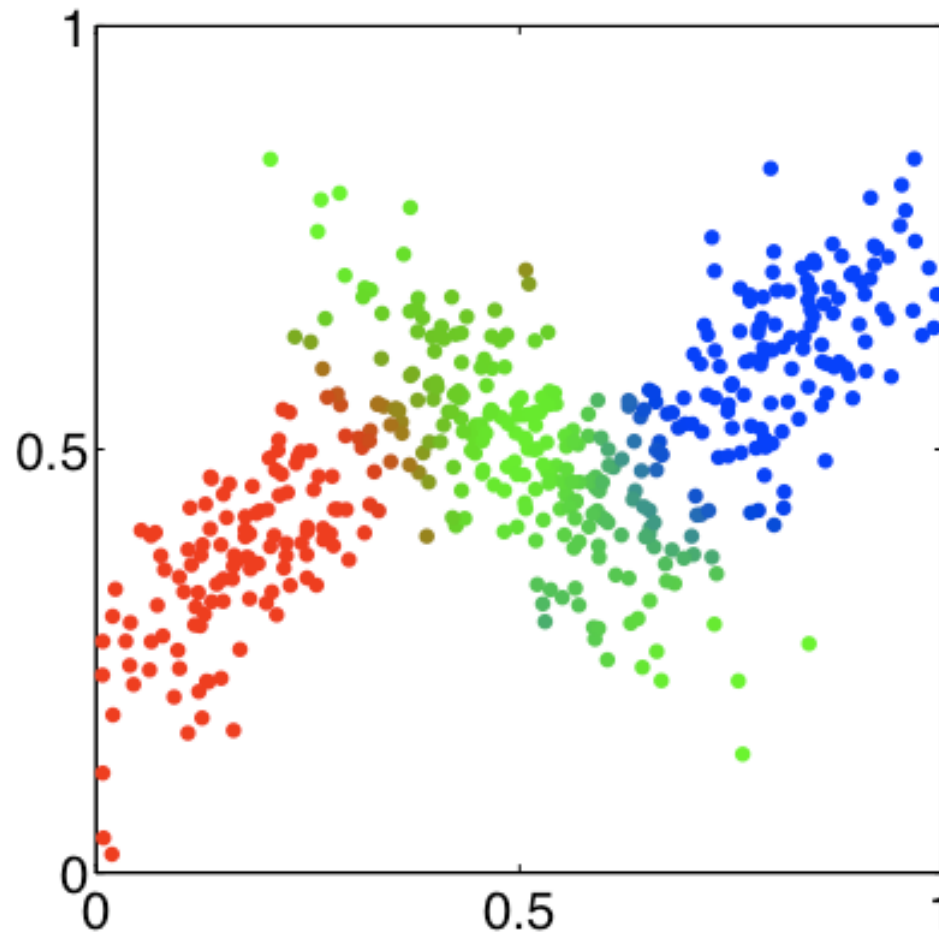
Easy to estimate params  
(do each color separately)



unlabeled

Hard to estimate params  
(we need to assign colors)

# EM produces a “Soft” labeling



each point makes a weighted contribution  
to the estimation of ALL components

# Review of EM for GMMs

**E**  $\gamma(z_{nj}) = \frac{\pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}$  ownership weights (soft labels)

---

**M**  $\boldsymbol{\mu}_j = \frac{\sum_{n=1}^N \gamma(z_{nj}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nj})}$  means

$\boldsymbol{\Sigma}_j = \frac{\sum_{n=1}^N \gamma(z_{nj}) (\mathbf{x}_n - \boldsymbol{\mu}_j)(\mathbf{x}_n - \boldsymbol{\mu}_j)^\top}{\sum_{n=1}^N \gamma(z_{nj})}$  covariances

$\pi_j = \frac{1}{N} \sum_{n=1}^N \gamma(z_{nj})$  mixing weights

Alternate E and M steps to convergence.



# From EM to K-Means

Alternative explanation of K-means!

- Fix all mixing weights to  $1/K$   
*[drop out of the estimation]*
- Fix all covariances to  $\sigma^2 \mathbf{I}$   
*[drop out of the estimation so we only have to estimate the means; each Gaussian likelihood becomes inversely proportional to distance from a mean]*
- Take limit as  $\sigma^2$  goes to 0  
*[this forces soft weights to become binary]*

# From EM to K-Means

$$\mathbf{E} \quad \gamma(z_{nj}) = \frac{\pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_k \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}$$

ownership weights  
(soft labels)

$$= \frac{\exp(-\frac{1}{2\sigma^2} \|x_n - \mu_j\|^2)}{\sum_{k=1}^K \exp(-\frac{1}{2\sigma^2} \|x_n - \mu_k\|^2)}$$

after fixing mixing weights  
and covariances as described  
on last slide

now divide top and bottom by  $\exp(-\frac{1}{2\sigma^2} d_{min}^2)$  where  $d_{min}^2 = \min_k \|x_n - \mu_k\|^2$   
and take limit as  $\sigma^2$  goes to 0

$$\gamma(z_{nj}) = z_{nj} = \begin{cases} 1 & \text{if } \mu_j \text{ is closest mean to } x_n \\ 0 & \text{otherwise} \end{cases}$$

hard labels, as in the  
K-means algorithm

# K-Means Algorithm

- Given  $N$  data points  $x_1, x_2, \dots, x_N$
- Find  $K$  cluster centers  $\mu_1, \mu_2, \dots, \mu_K$  to minimize  $\sum_{n=1}^N \sum_{k=1}^K z_{nk} \|x_n - \mu_j\|^2$   
( $z_{nk}$  is 1 if point  $n$  belongs to cluster  $k$ ; 0 otherwise)
- Algorithm:
  - initialize  $K$  cluster centers  $\mu_1, \mu_2, \dots, \mu_K$
  - repeat
    - set  $z_{nk}$  labels to assign each point to closest cluster center
    - revise each cluster center  $\mu_j$  to be center of mass of points in that cluster  $\mu_j = \frac{\sum_{n=1}^N z_{nj} x_n}{\sum_{n=1}^N z_{nj}}$
  - until convergence (e.g.  $z_{nk}$  labels don't change)

E

M