



# 10-708 Probabilistic Graphical Models

Machine Learning Department  
School of Computer Science  
Carnegie Mellon University



## Variational EM

Matt Gormley  
Lecture 18  
Apr. 2, 2021

# Reminders

- **Homework 4: MCMC**
  - Out: Wed, Mar. 24
  - Due: Wed, Apr. 7 at 11:59pm
- **Project Midway Milestones:**
  - **Midway Poster Session:**  
Wed, Apr. 14 at 6:30pm – 8:30pm
  - **Midway Executive Summary**  
**Due: Wed, Apr. 14 at 11:59pm**
- **Homework 5: Variational Inference**
  - Out: Thu, Apr. 8
  - Due: Wed, Apr. 21 at 11:59pm

# **HIDDEN STATE CRFS**

# Case Study: Object Recognition

Data consists of images  $x$  and labels  $y$ .



pigeon



rhinoceros



leopard



llama

# Case Study: Object Recognition

Data consists of images  $x$  and labels  $y$ .

- Preprocess data into “patches”
- Posit a latent labeling  $z$  describing the object’s parts (e.g. head, leg, tail, torso, grass)
- Define graphical model with these latent variables in mind
- $z$  is not observed at train or test time

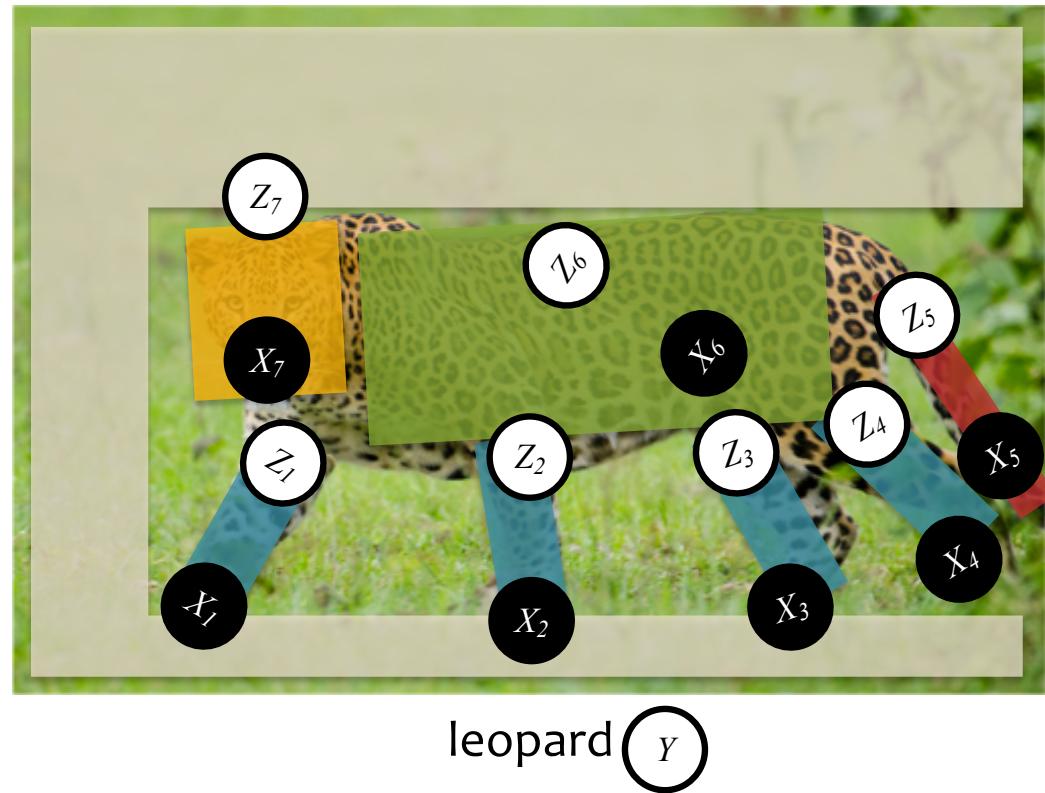


leopard

# Case Study: Object Recognition

Data consists of images  $x$  and labels  $y$ .

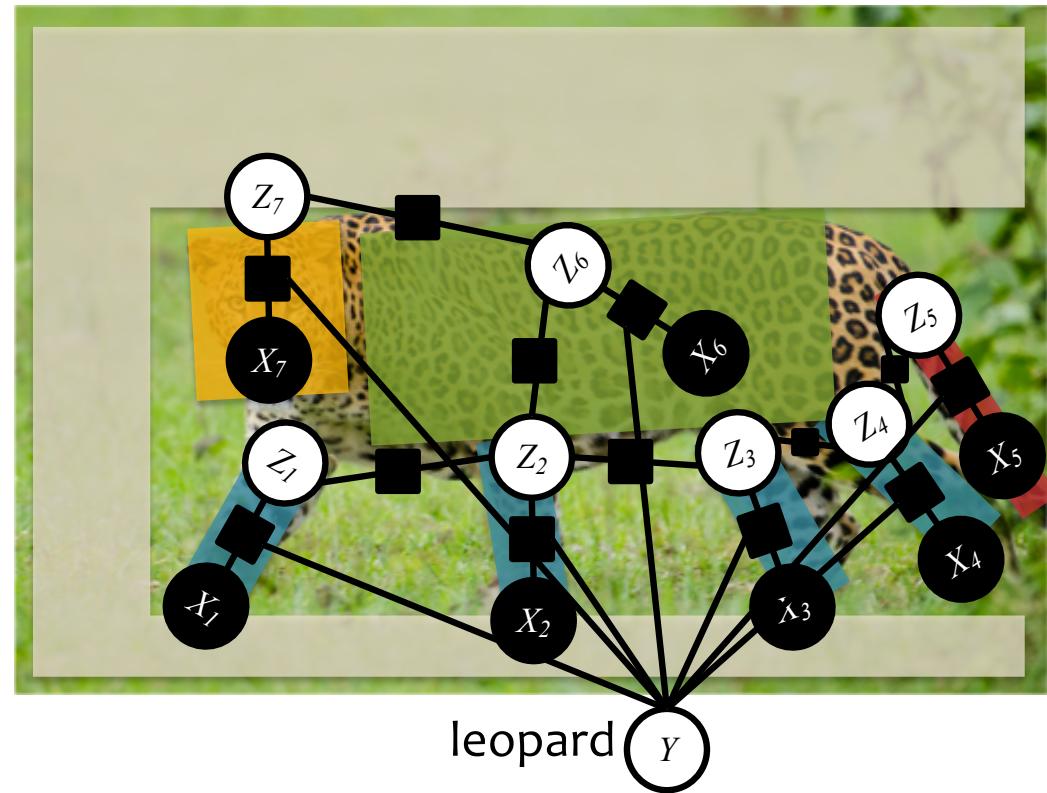
- Preprocess data into “patches”
- Posit a latent labeling  $z$  describing the object’s parts (e.g. head, leg, tail, torso, grass)
- Define graphical model with these latent variables in mind
- $z$  is not observed at train or test time



# Case Study: Object Recognition

Data consists of images  $x$  and labels  $y$ .

- Preprocess data into “patches”
- Posit a latent labeling  $z$  describing the object’s parts (e.g. head, leg, tail, torso, grass)
- Define graphical model with these latent variables in mind
- $z$  is not observed at train or test time

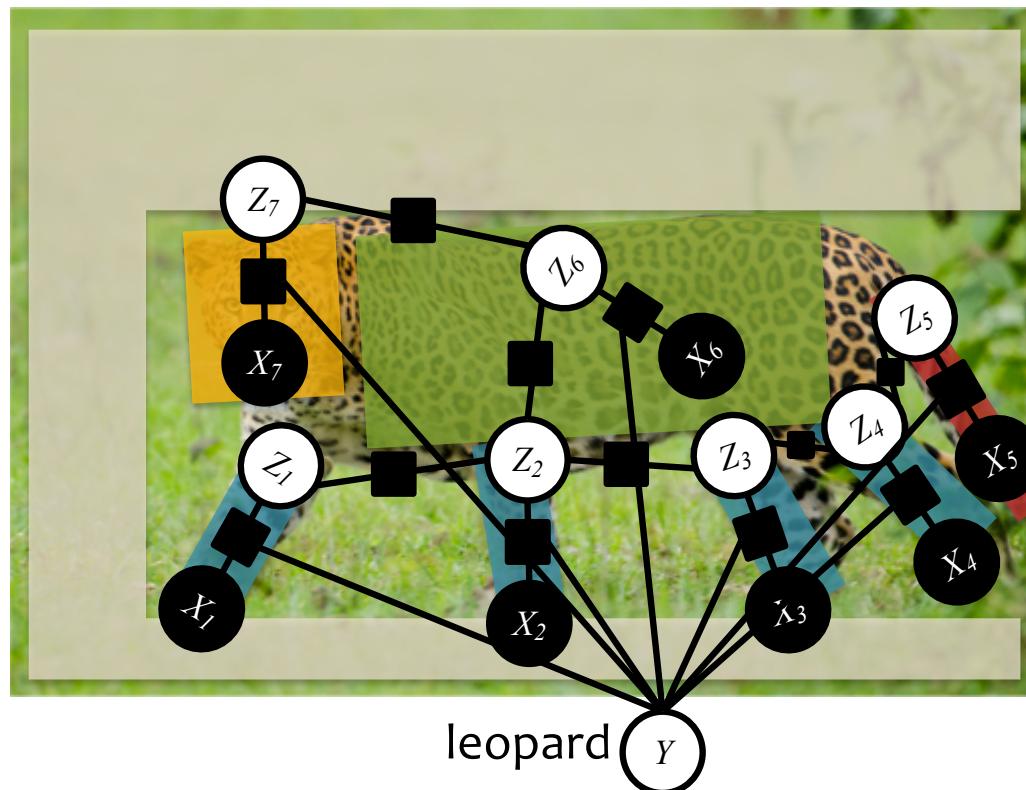


# Hidden-state CRFs

Data:  $\mathcal{D} = \{\mathbf{x}^{(n)}, \mathbf{y}^{(n)}\}_{n=1}^N$

Joint model:  $p_{\theta}(\mathbf{y}, \mathbf{z} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x}, \theta)} \prod_{\alpha} \psi_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{z}_{\alpha}, \mathbf{x})$

Marginalized model:  $p_{\theta}(\mathbf{y} \mid \mathbf{x}) = \sum_{\mathbf{z}} p_{\theta}(\mathbf{y}, \mathbf{z} \mid \mathbf{x})$



# Hidden-state CRFs

Data:  $\mathcal{D} = \{\mathbf{x}^{(n)}, \mathbf{y}^{(n)}\}_{n=1}^N$

Joint model:  $p_{\theta}(\mathbf{y}, \mathbf{z} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x}, \theta)} \prod_{\alpha} \psi_{\alpha}(\mathbf{y}_{\alpha}, \mathbf{z}_{\alpha}, \mathbf{x})$

Marginalized model:  $p_{\theta}(\mathbf{y} \mid \mathbf{x}) = \sum_{\mathbf{z}} p_{\theta}(\mathbf{y}, \mathbf{z} \mid \mathbf{x})$

We can train using gradient based methods:

(the values  $\mathbf{x}$  are omitted below for clarity)

$$\begin{aligned} \frac{d\ell(\theta|\mathcal{D})}{d\theta} &= \sum_{n=1}^N \left( \mathbb{E}_{\mathbf{z} \sim p_{\theta}(\cdot|\mathbf{y}^{(n)})} [f_j(\mathbf{y}^{(n)}, \mathbf{z})] - \mathbb{E}_{\mathbf{y}, \mathbf{z} \sim p_{\theta}(\cdot, \cdot)} [f_j(\mathbf{y}, \mathbf{z})] \right) \\ &= \sum_{n=1}^N \sum_{\alpha} \left( \sum_{\mathbf{z}_{\alpha}} \underbrace{p_{\theta}(\mathbf{z}_{\alpha} \mid \mathbf{y}^{(n)})}_{\text{Inference on clamped factor graph}} f_{\alpha,j}(\mathbf{y}_{\alpha}^{(n)}, \mathbf{z}_{\alpha}) - \sum_{\mathbf{y}_{\alpha}, \mathbf{z}_{\alpha}} \underbrace{p_{\theta}(\mathbf{y}_{\alpha}, \mathbf{z}_{\alpha})}_{\text{Inference on full factor graph}} f_{\alpha,j}(\mathbf{y}_{\alpha}, \mathbf{z}_{\alpha}) \right) \end{aligned}$$

Inference on  
**clamped**  
factor graph

Inference on  
**full**  
factor graph

# **GAUSSIAN MIXTURE MODEL (GMM)**

# Gaussian Mixture-Model

**Data:**  $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$  where  $\mathbf{x}^{(i)} \in \mathbb{R}^M$

**Generative Story:**  $z \sim \text{Categorical}(\boldsymbol{\phi})$   
 $\mathbf{x} \sim \text{Gaussian}(\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$

**Model:** Joint:  $p(\mathbf{x}, z; \boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = p(\mathbf{x}|z; \boldsymbol{\mu}, \boldsymbol{\Sigma})p(z; \boldsymbol{\phi})$   
Marginal:  $p(\mathbf{x}; \boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{z=1}^K p(\mathbf{x}|z; \boldsymbol{\mu}, \boldsymbol{\Sigma})p(z; \boldsymbol{\phi})$

**(Marginal) Log-likelihood:**

$$\begin{aligned}\ell(\boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \log \prod_{i=1}^N p(\mathbf{x}^{(i)}; \boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ &= \sum_{i=1}^N \log \sum_{z=1}^K p(\mathbf{x}^{(i)}|z; \boldsymbol{\mu}, \boldsymbol{\Sigma})p(z; \boldsymbol{\phi})\end{aligned}$$

# Mixture-Model

**Data:**  $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$  where  $\mathbf{x}^{(i)} \in \mathbb{R}^M$

**Generative Story:**  $z \sim \text{Categorical}(\boldsymbol{\phi})$

$$\mathbf{x} \sim p_{\boldsymbol{\theta}}(\cdot|z)$$

**Model:** Joint:  $p_{\boldsymbol{\theta}, \boldsymbol{\phi}}(\mathbf{x}, z) = p_{\boldsymbol{\theta}}(\mathbf{x}|z)p_{\boldsymbol{\phi}}(z)$

Marginal:  $p_{\boldsymbol{\theta}, \boldsymbol{\phi}}(\mathbf{x}) = \sum_{z=1}^K p_{\boldsymbol{\theta}}(\mathbf{x}|z)p_{\boldsymbol{\phi}}(z)$

**(Marginal) Log-likelihood:**

$$\ell(\boldsymbol{\theta}) = \log \prod_{i=1}^N p_{\boldsymbol{\theta}, \boldsymbol{\phi}}(\mathbf{x}^{(i)})$$

$$= \sum_{i=1}^N \log \sum_{z=1}^K p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|z)p_{\boldsymbol{\phi}}(z)$$

# Mixture-Model

**Data:**  $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$  where  $\mathbf{x}^{(i)} \in \mathbb{R}^M$

**Generative Story:**  $z \sim \text{Categorical}(\phi)$

$$\mathbf{x} \sim p_{\theta}(\cdot|z) \quad \leftarrow$$

**Model:** Joint:  $p_{\theta, \phi}(\mathbf{x}, z) = p_{\theta}(z)p_{\theta}(\mathbf{x}|z)$

Marginal:  $p_{\theta, \phi}(\mathbf{x}) = \sum_{z=1}^K p_{\theta}(z)p_{\theta}(\mathbf{x}|z)$

This could be any arbitrary distribution parameterized by  $\theta$ .

Today we're thinking about the case where it is a Multivariate Gaussian.

**(Marginal) Log-likelihood:**

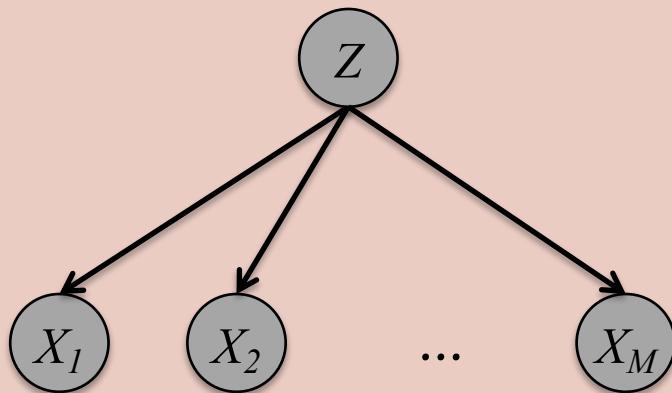
$$\ell(\theta) = \log \prod_{i=1}^N p_{\theta, \phi}(\mathbf{x}^{(i)})$$

$$= \sum_{i=1}^N \log \sum_{z=1}^K p_{\theta}(\mathbf{x}^{(i)}|z)p_{\phi}(z)$$

# Learning a Mixture Model

**Supervised Learning:** The parameters decouple!

$$\mathcal{D} = \{(\mathbf{x}^{(i)}, \mathbf{z}^{(i)})\}_{i=1}^N$$



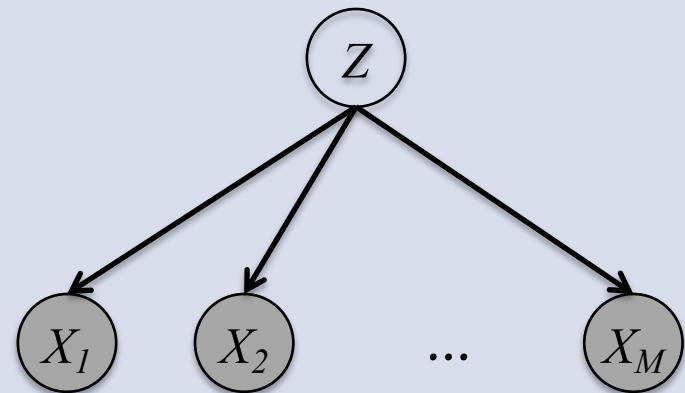
$$\theta^*, \phi^* = \operatorname{argmax}_{\theta, \phi} \sum_{i=1}^N \log p_\theta(\mathbf{x}^{(i)} | z^{(i)}) p_\phi(z^{(i)})$$

$$\theta^* = \operatorname{argmax}_\theta \sum_{i=1}^N \log p_\theta(\mathbf{x}^{(i)} | z^{(i)})$$

$$\phi^* = \operatorname{argmax}_\theta \sum_{i=1}^N \log p_\phi(z^{(i)})$$

**Unsupervised Learning:** Parameters are coupled by marginalization.

$$\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$$



$$\theta^*, \phi^* = \operatorname{argmax}_{\theta, \phi} \sum_{i=1}^N \log \sum_{z=1}^K p_\theta(\mathbf{x}^{(i)} | z) p_\phi(z)$$

# Learning a Mixture Model

**Supervised Learning:** The parameters decouple!

$$\mathcal{D} = \{(\mathbf{x}^{(i)}, \mathbf{z}^{(i)})\}_{i=1}^N$$

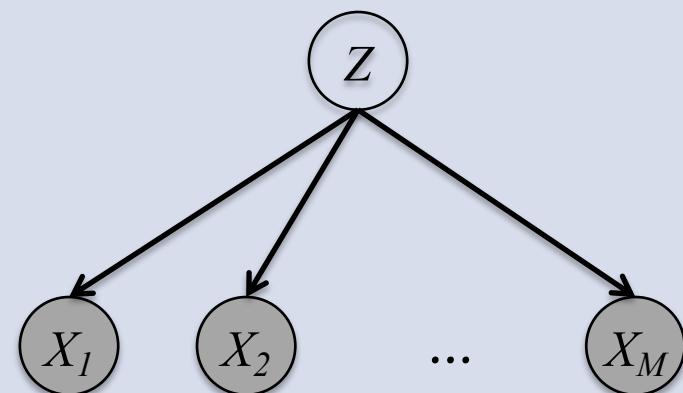
Training certainly isn't as simple as the supervised case.

In many cases, we could still use some black-box optimization method (e.g. Newton-Raphson) to solve this coupled optimization problem.

This lecture is about a more problem-specific method: EM.

**Unsupervised Learning:** Parameters are coupled by marginalization.

$$\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$$



$$\theta^*, \phi^* = \underset{\theta, \phi}{\operatorname{argmax}} \sum_{i=1}^N \log \sum_{z=1}^K p_\theta(\mathbf{x}^{(i)}|z)p_\phi(z)$$



# **EXPECTATION MAXIMIZATION**

# Hard Expectation-Maximization

- Initialize **parameters** randomly
- **while** not converged

## 1. E-Step:

Set the **latent variables** to the values that maximizes likelihood, treating parameters as observed

Estimate unobserved variables

## 2. M-Step:

Set the **parameters** to the values that maximizes likelihood, treating latent variables as observed

MLE given the estimated values of unobserved variables

# (Soft) Expectation-Maximization

- Initialize **parameters** randomly
- **while** not converged

## 1. E-Step:

Create one training example for each possible value of the **latent variables**

Weight each example according to model's confidence

Treat parameters as observed

## 2. M-Step:

Set the **parameters** to the values that maximizes likelihood

Treat pseudo-counts from above as observed

Estimate unobserved variables

MLE given the estimated values of unobserved variables

# Hard EM vs. Soft EM

## Algorithm 1 Hard EM for GMMs

```

1: procedure HARDEM( $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ )
2:   Randomly initialize parameters,  $\phi, \mu, \Sigma$ 
3:   while not converged do
4:     E-Step:

```

$$z^{(i)} \leftarrow \underset{z}{\operatorname{argmax}} \log p(\mathbf{x}^{(i)}|z; \mu, \Sigma) + \log p(z; \phi)$$

5: M-Step:

$$\phi_k \leftarrow \frac{1}{N} \sum_{i=1}^N \mathbb{I}(z^{(i)} = k), \forall k$$

$$\boldsymbol{\mu}_k \leftarrow \frac{\sum_{i=1}^N \mathbb{I}(z^{(i)} = k) \mathbf{x}^{(i)}}{\sum_{i=1}^N \mathbb{I}(z^{(i)} = k)}, \forall k$$

$$\boldsymbol{\Sigma}_k \leftarrow \frac{\sum_{i=1}^N \mathbb{I}(z^{(i)} = k) (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)^T}{\sum_{i=1}^N \mathbb{I}(z^{(i)} = k)}, \forall k$$

6: **return**  $(\phi, \mu, \Sigma)$

## Algorithm 1 Soft EM for GMMs

```

1: procedure SOFTEM( $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ )
2:   Randomly initialize parameters,  $\phi, \mu, \Sigma$ 
3:   while not converged do
4:     E-Step:

```

$$c_k^{(i)} \leftarrow p(z^{(i)} = k | \mathbf{x}^{(i)}; \phi, \mu, \Sigma)$$

5: M-Step:

$$\phi_k \leftarrow \frac{1}{N} \sum_{i=1}^N c_k^{(i)}, \forall k$$

$$\boldsymbol{\mu}_k \leftarrow \frac{\sum_{i=1}^N c_k^{(i)} \mathbf{x}^{(i)}}{\sum_{i=1}^N c_k^{(i)}}, \forall k$$

$$\boldsymbol{\Sigma}_k \leftarrow \frac{\sum_{i=1}^N c_k^{(i)} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)^T}{\sum_{i=1}^N c_k^{(i)}}, \forall k$$

6: **return**  $(\phi, \mu, \Sigma)$

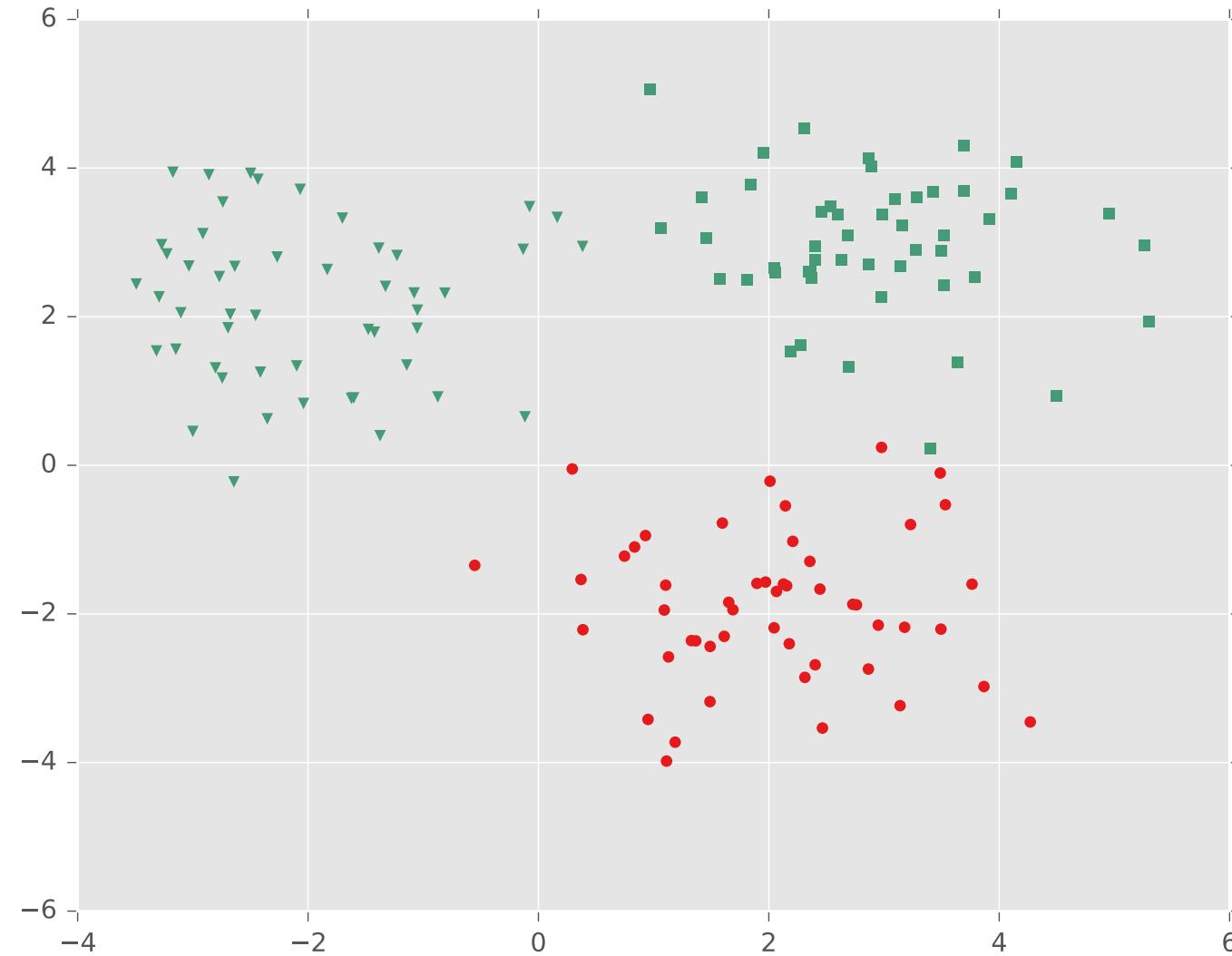
# Posterior Inference for Mixture Model

We obtain the posterior  $p(z^{(i)} = k | x^{(i)}; \phi, \mu, \Sigma)$  as follows:

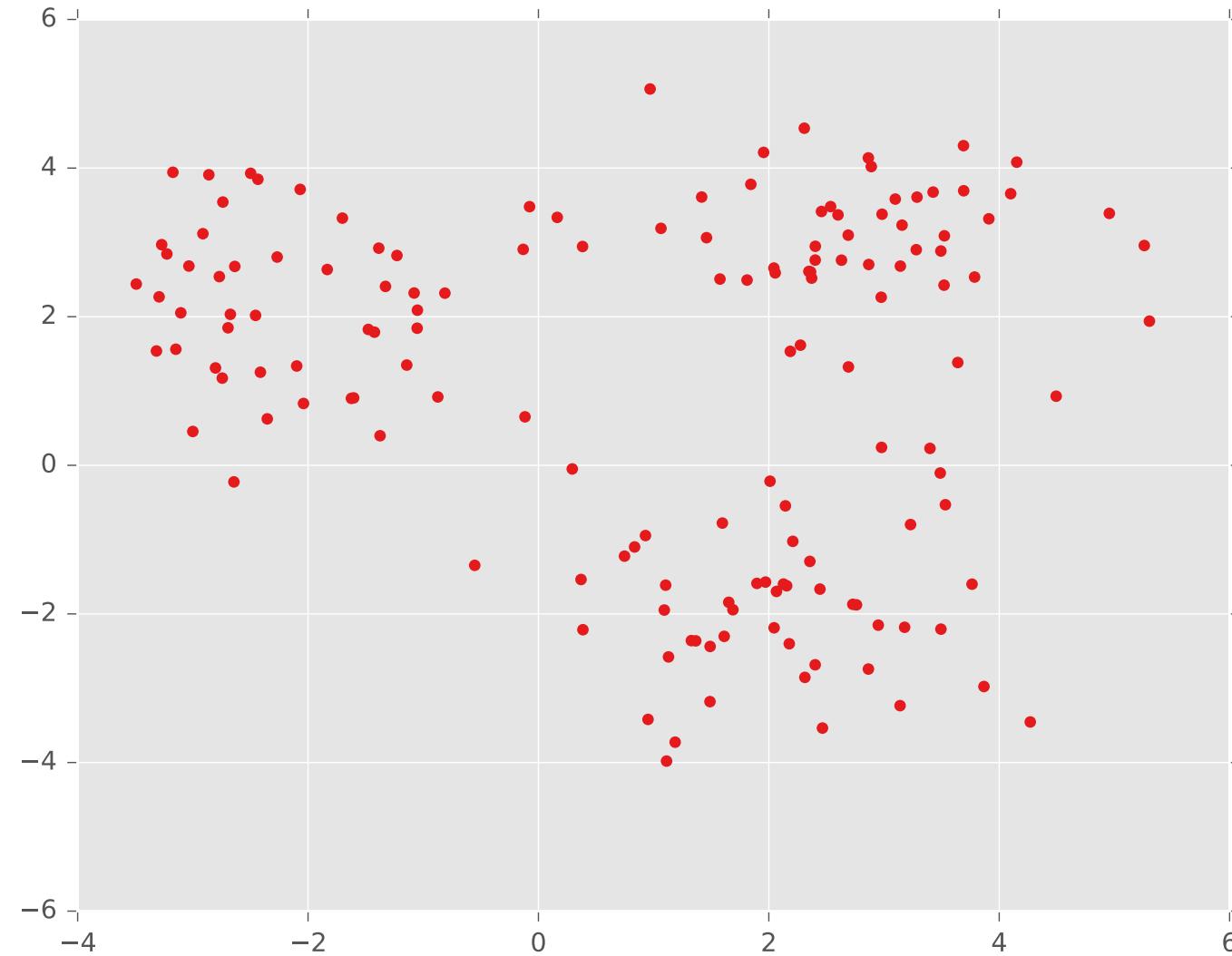
$$p(z^{(i)} = k | \mathbf{x}^{(i)}; \phi, \mu, \Sigma) = \frac{p(\mathbf{x}^{(i)} | z^{(i)} = k; \mu, \Sigma)p(z^{(i)} = k; \phi)}{\sum_{j=1}^K p(\mathbf{x}^{(i)} | z^{(i)} = j; \mu, \Sigma)p(z^{(i)} = j; \phi)} \quad (1)$$

# **EXAMPLE: K-MEANS VS GMM**

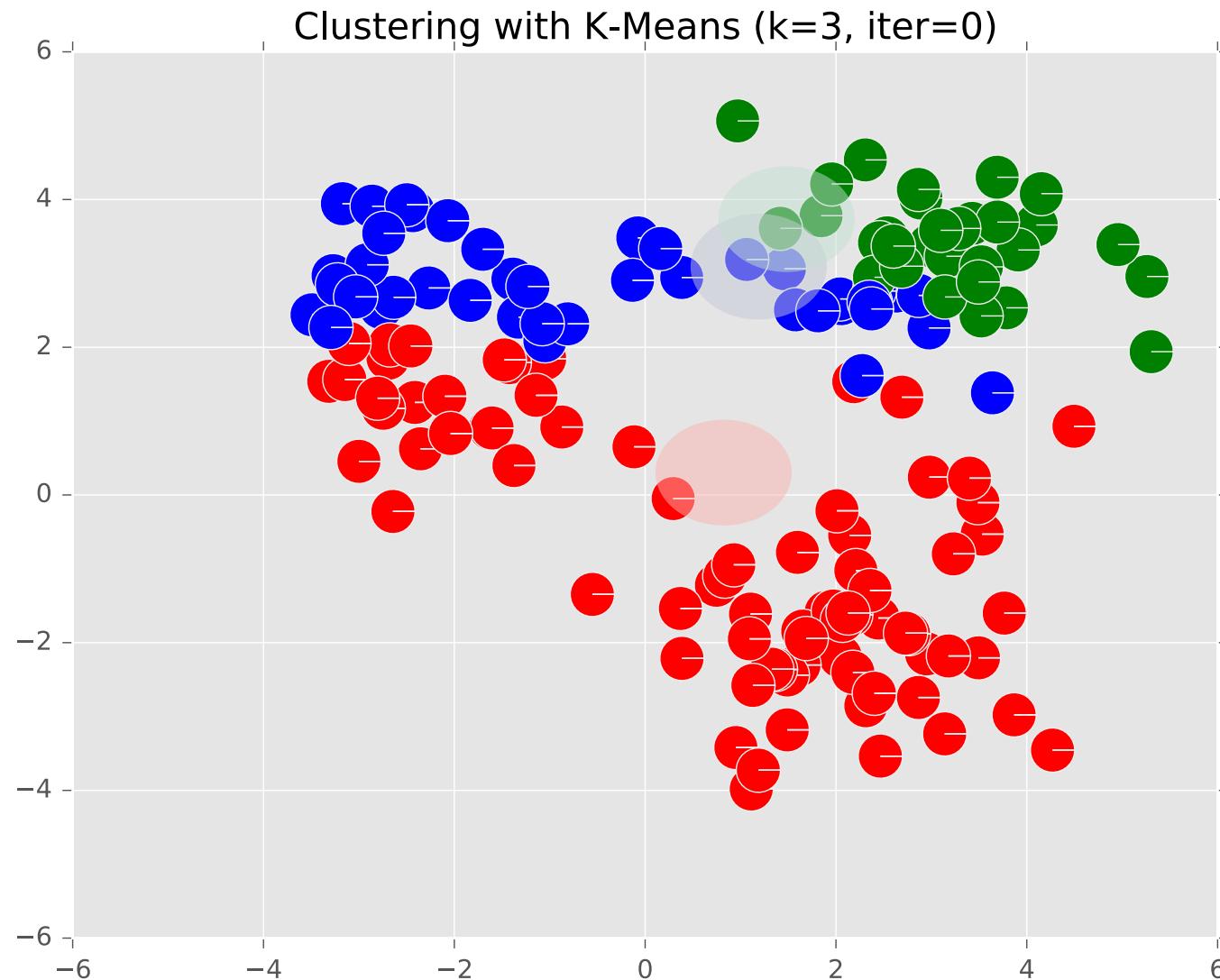
# Example: K-Means



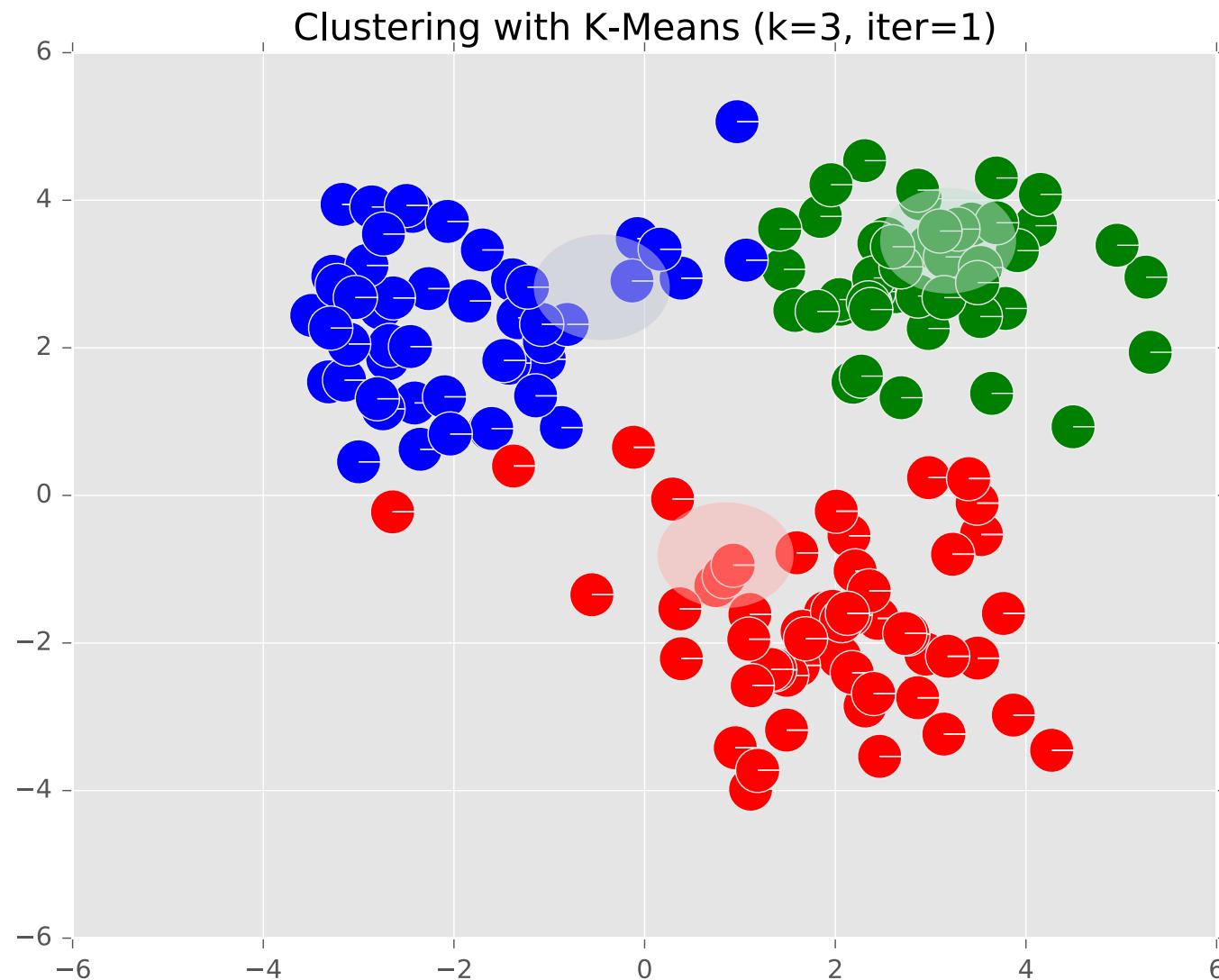
# Example: K-Means



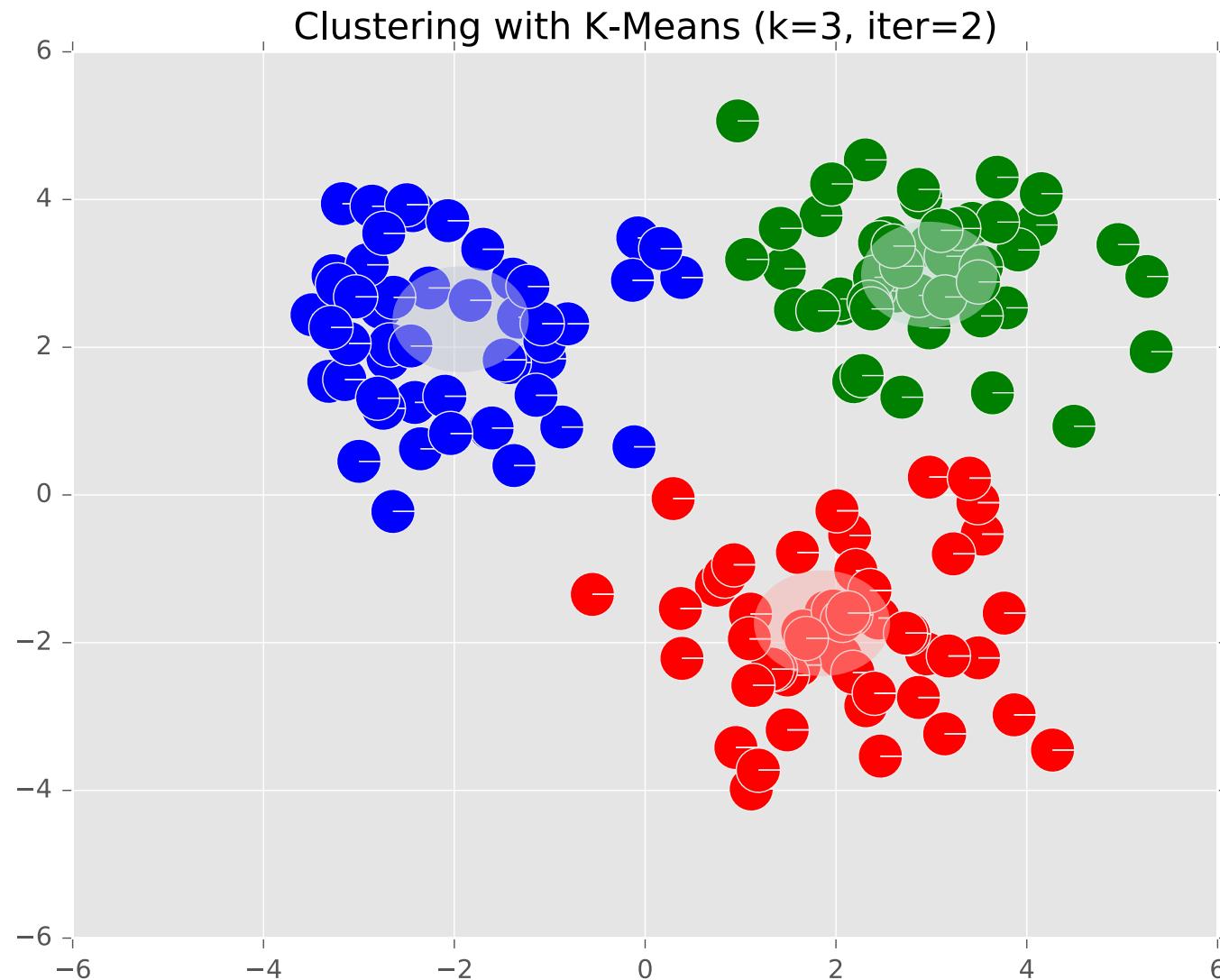
# Example: K-Means



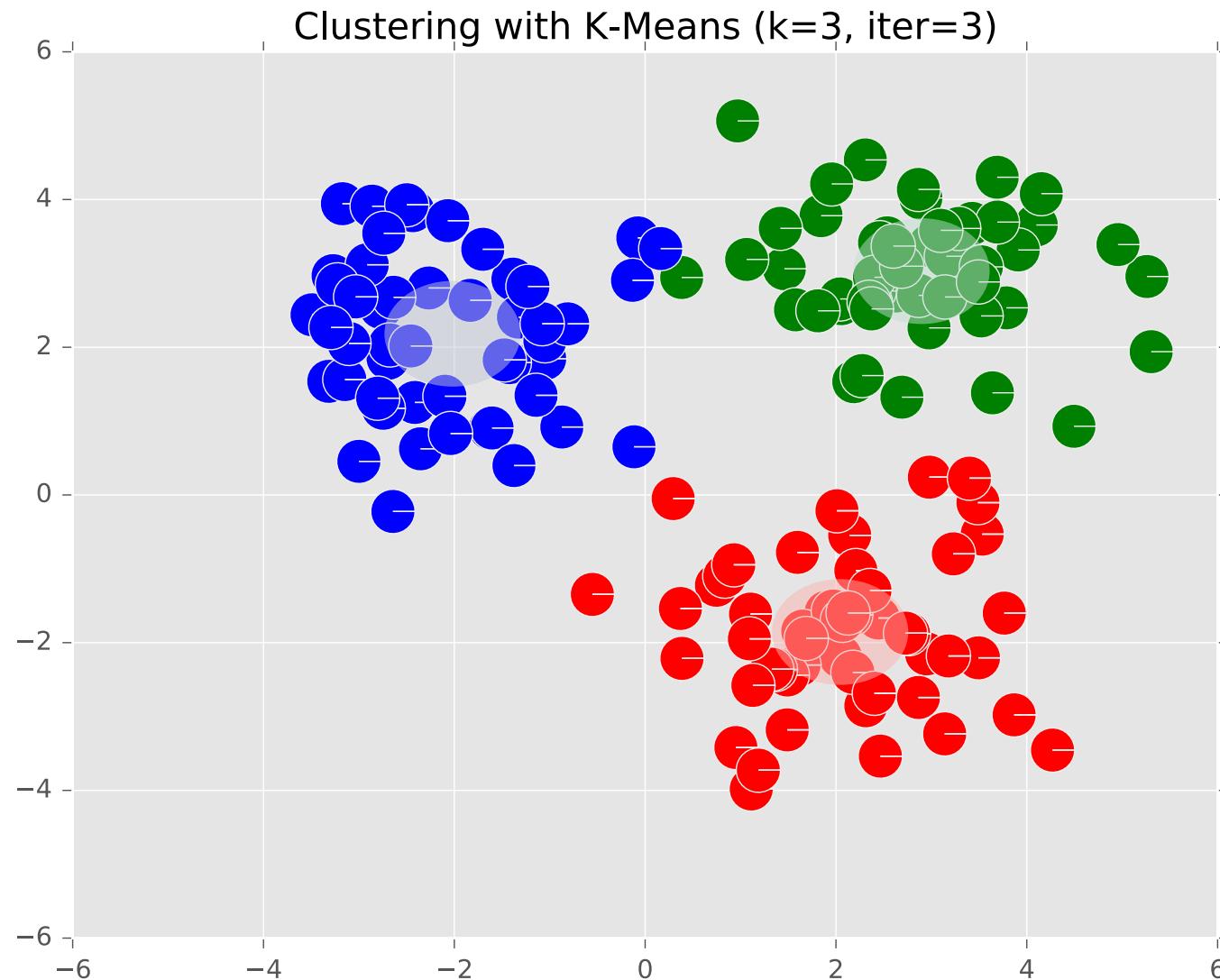
# Example: K-Means



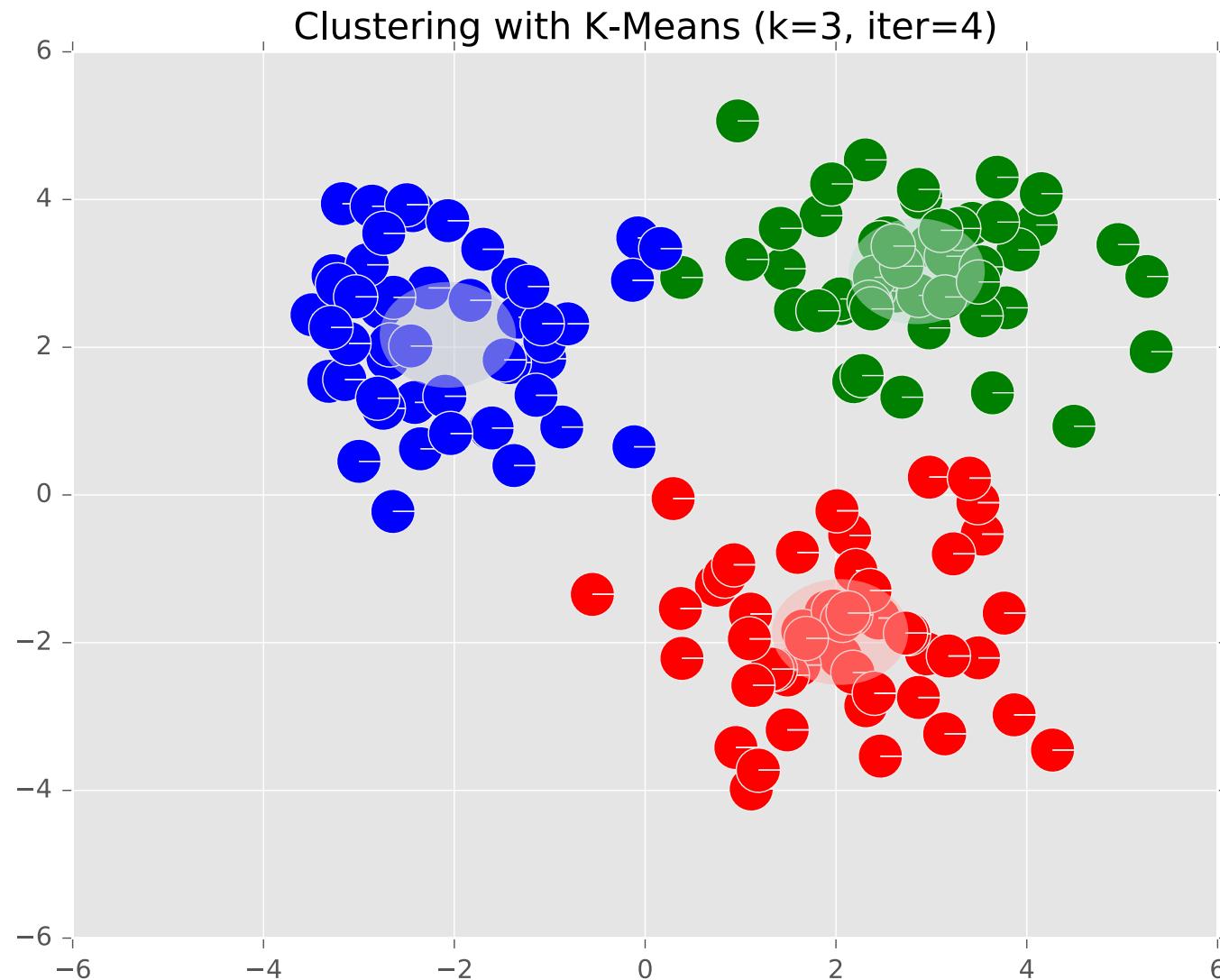
# Example: K-Means



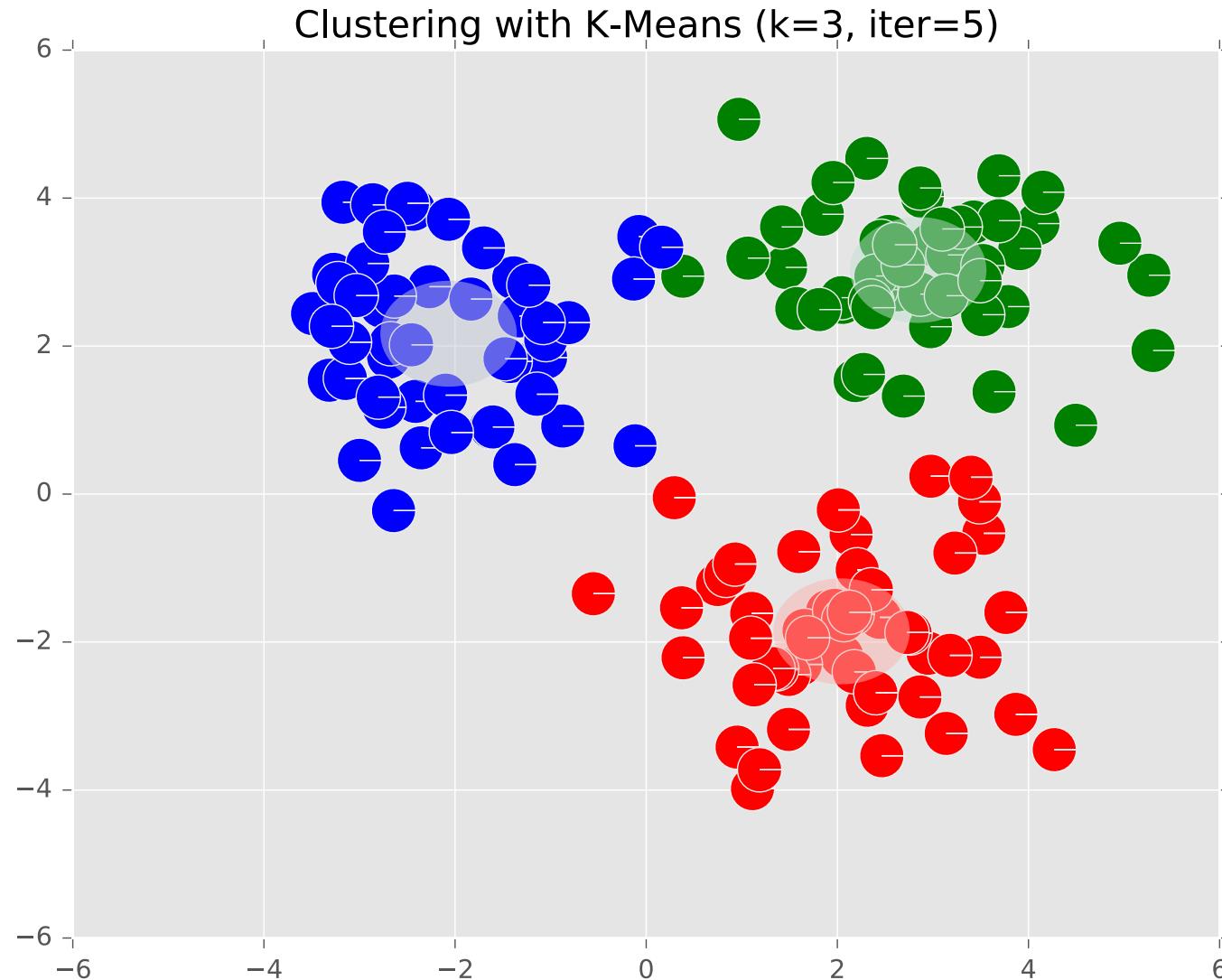
# Example: K-Means



# Example: K-Means



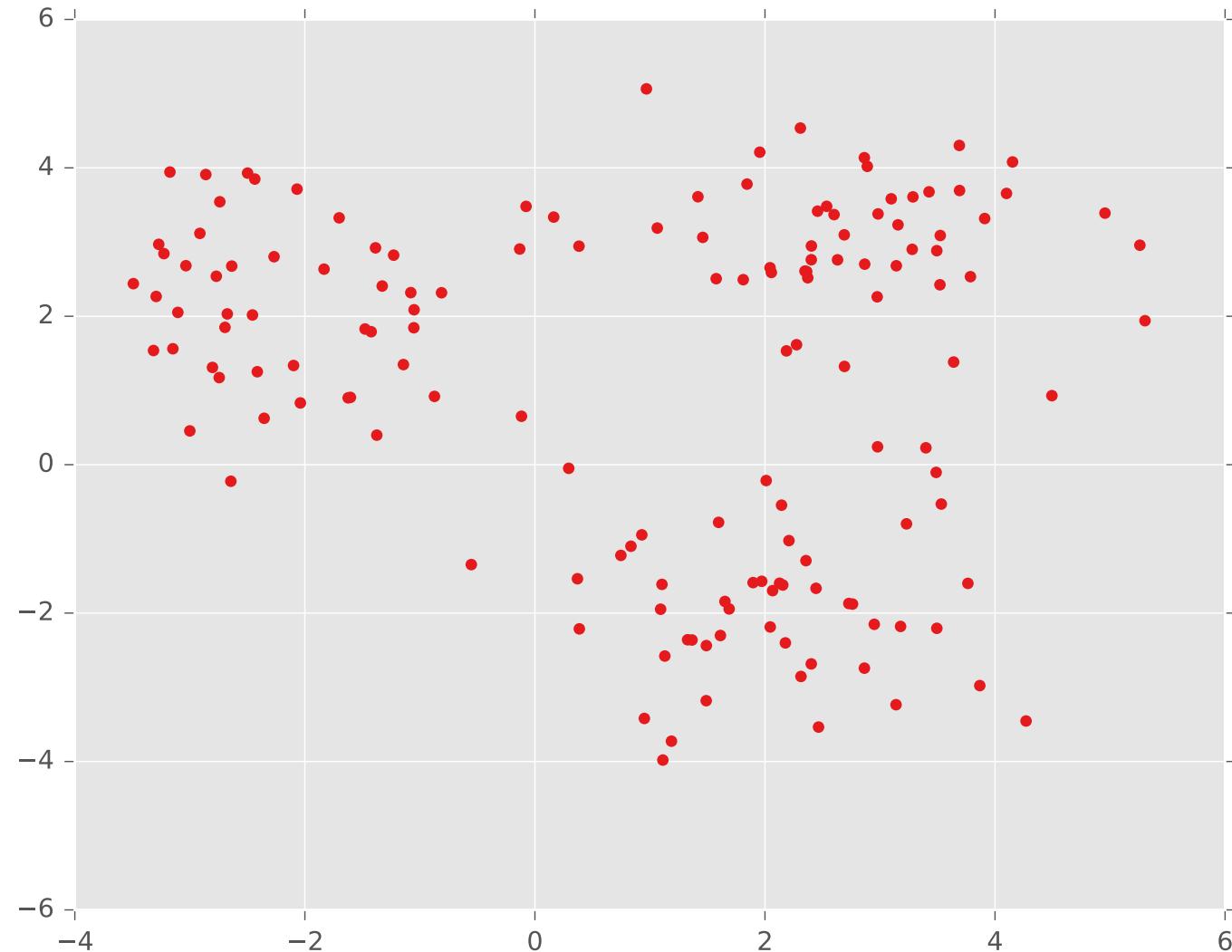
# Example: K-Means



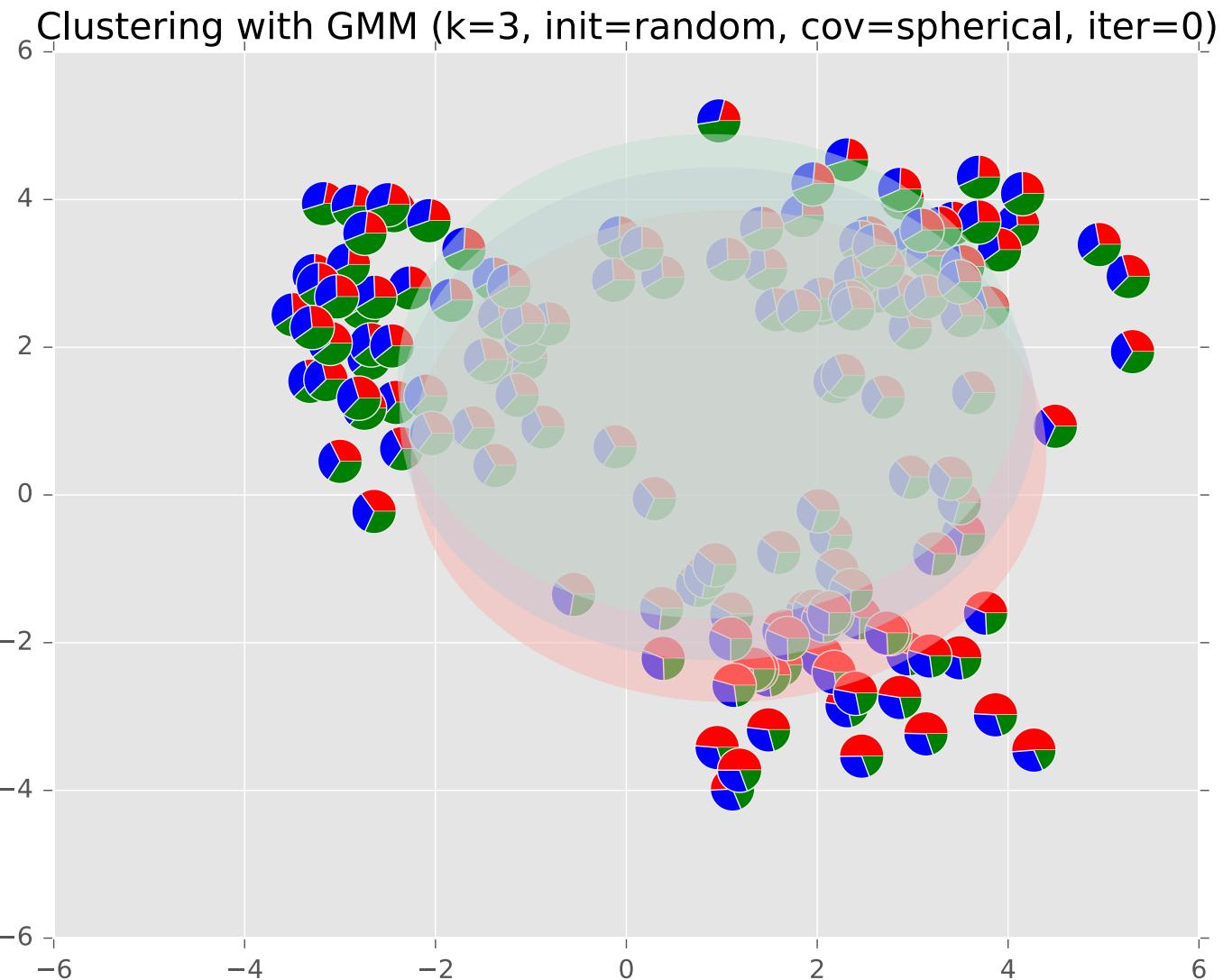
# Example: GMM



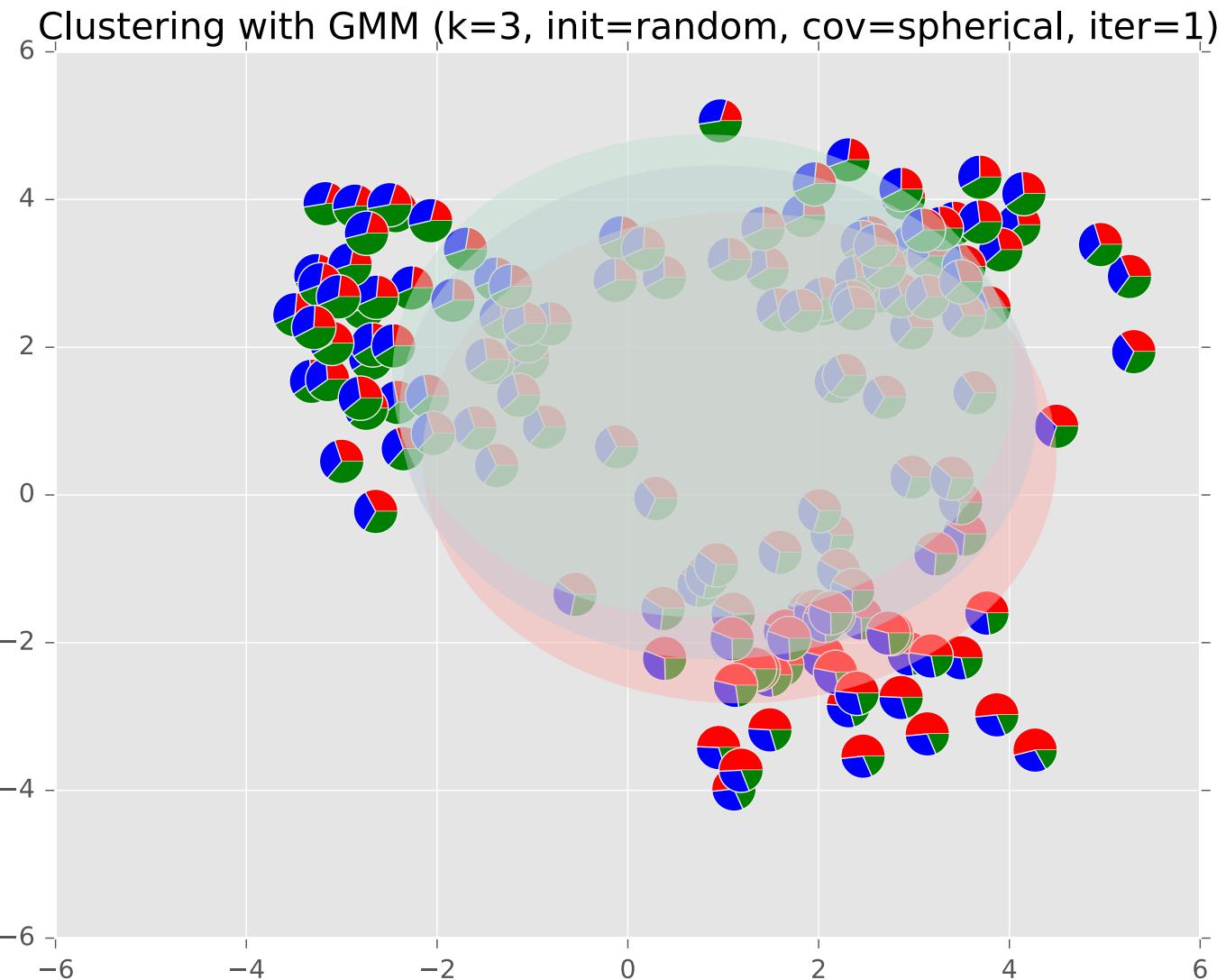
# Example: GMM



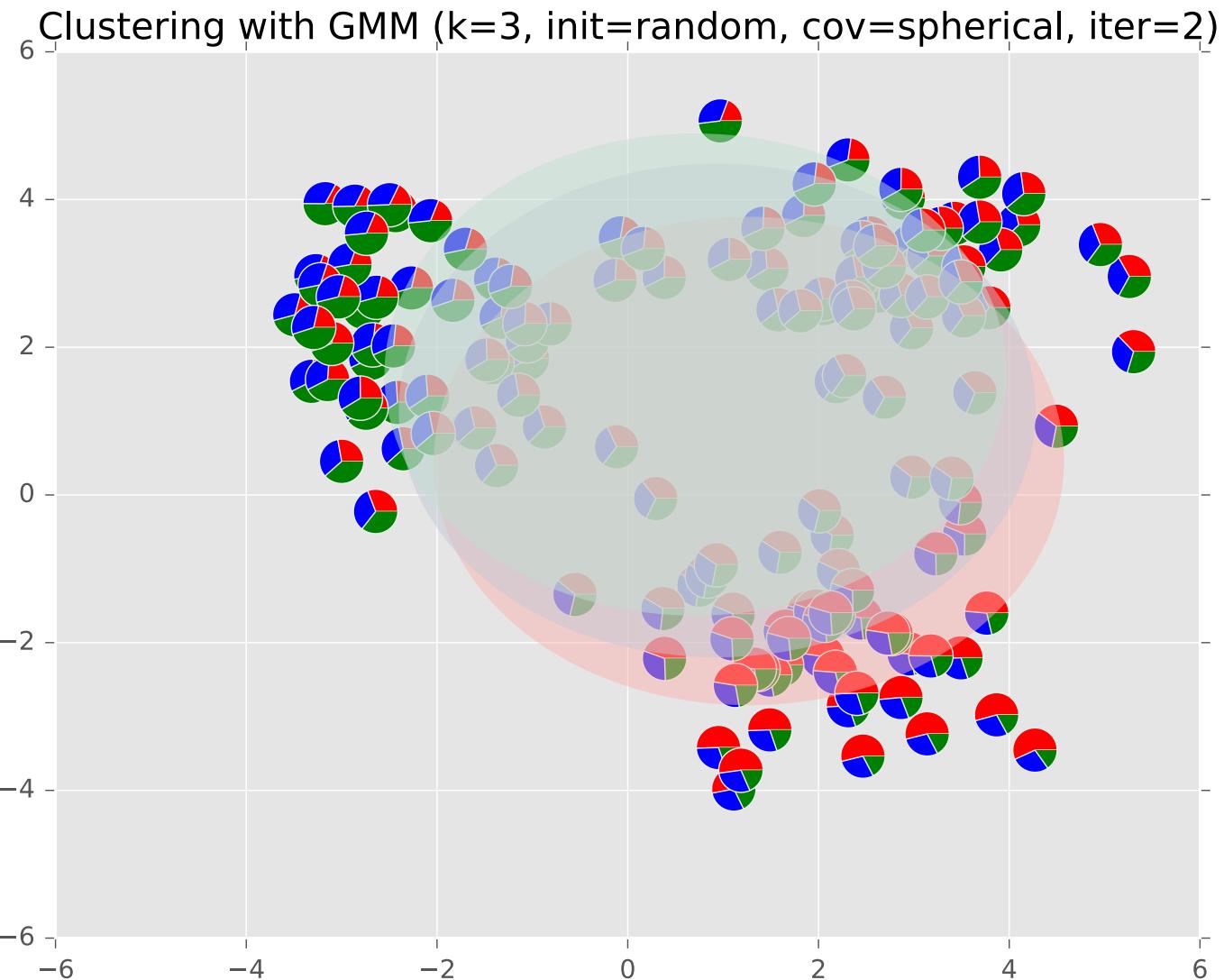
# Example: GMM



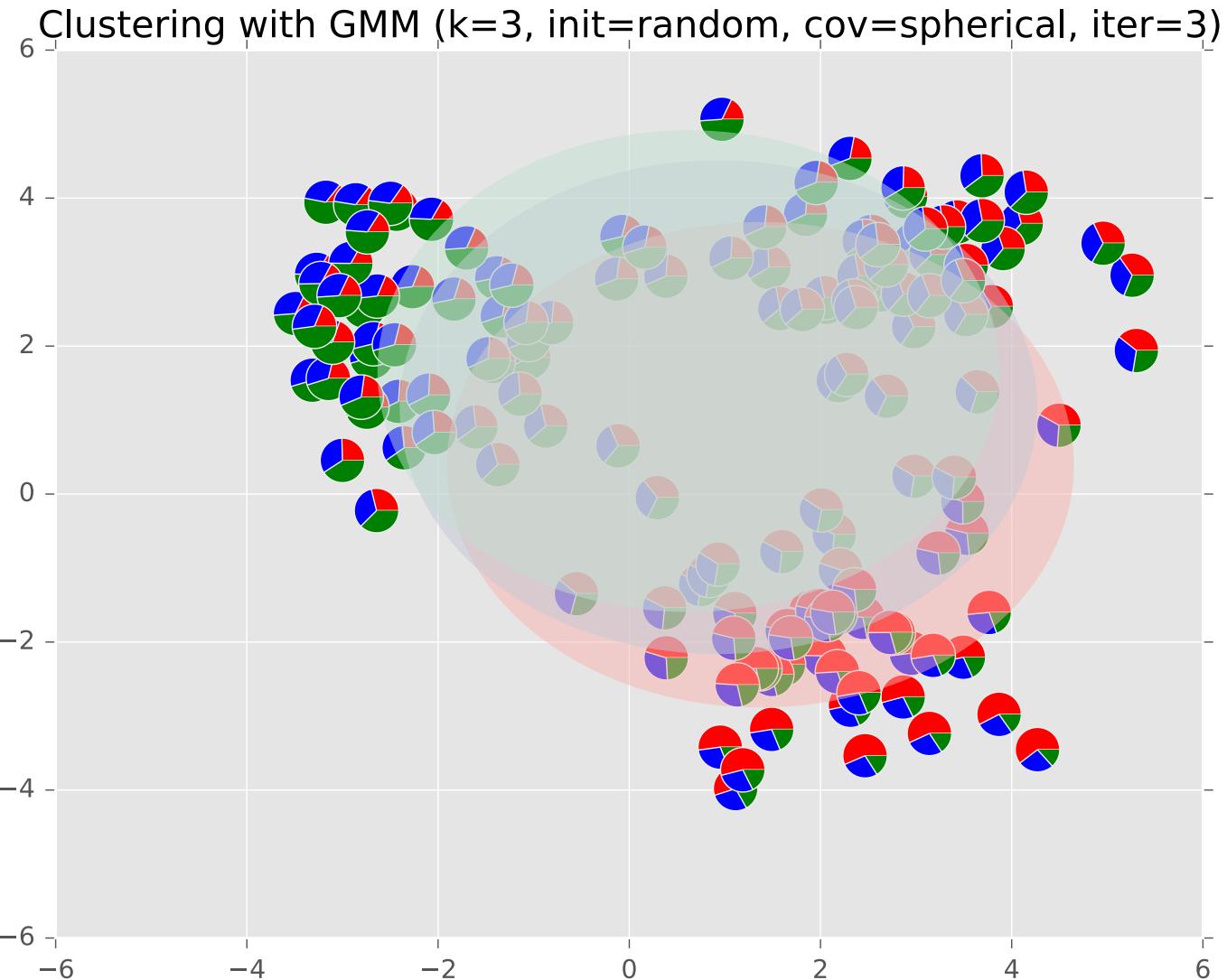
# Example: GMM



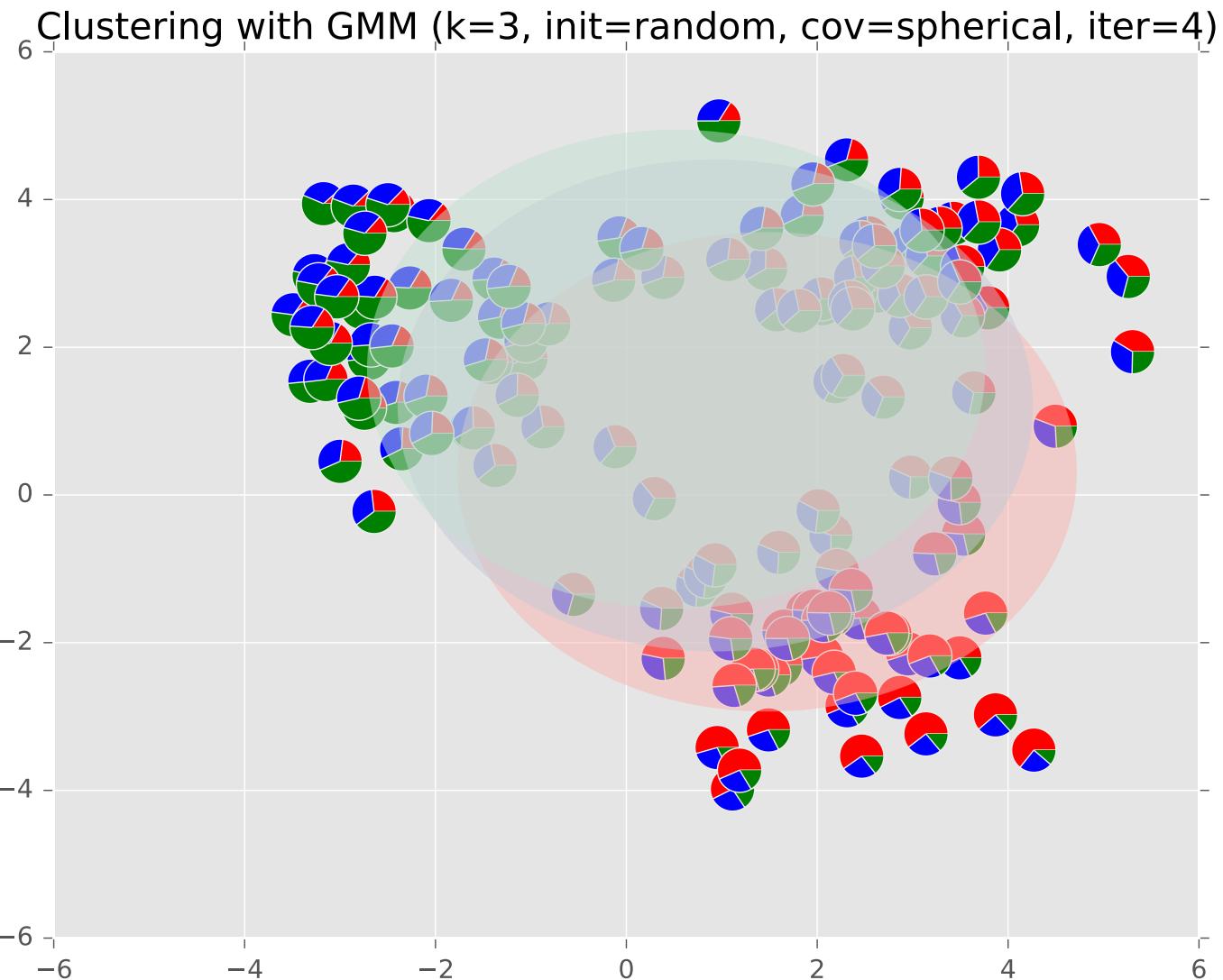
# Example: GMM



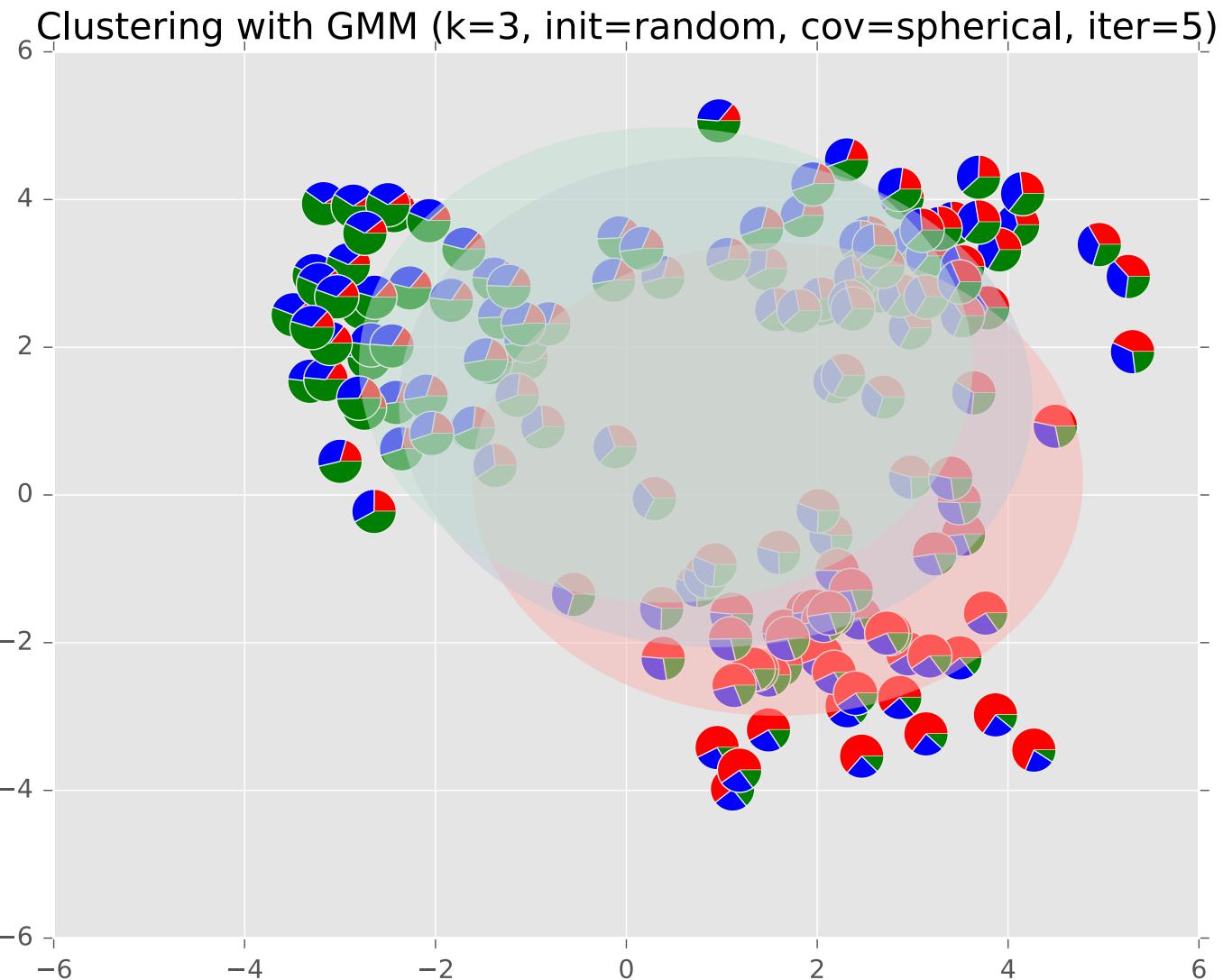
# Example: GMM



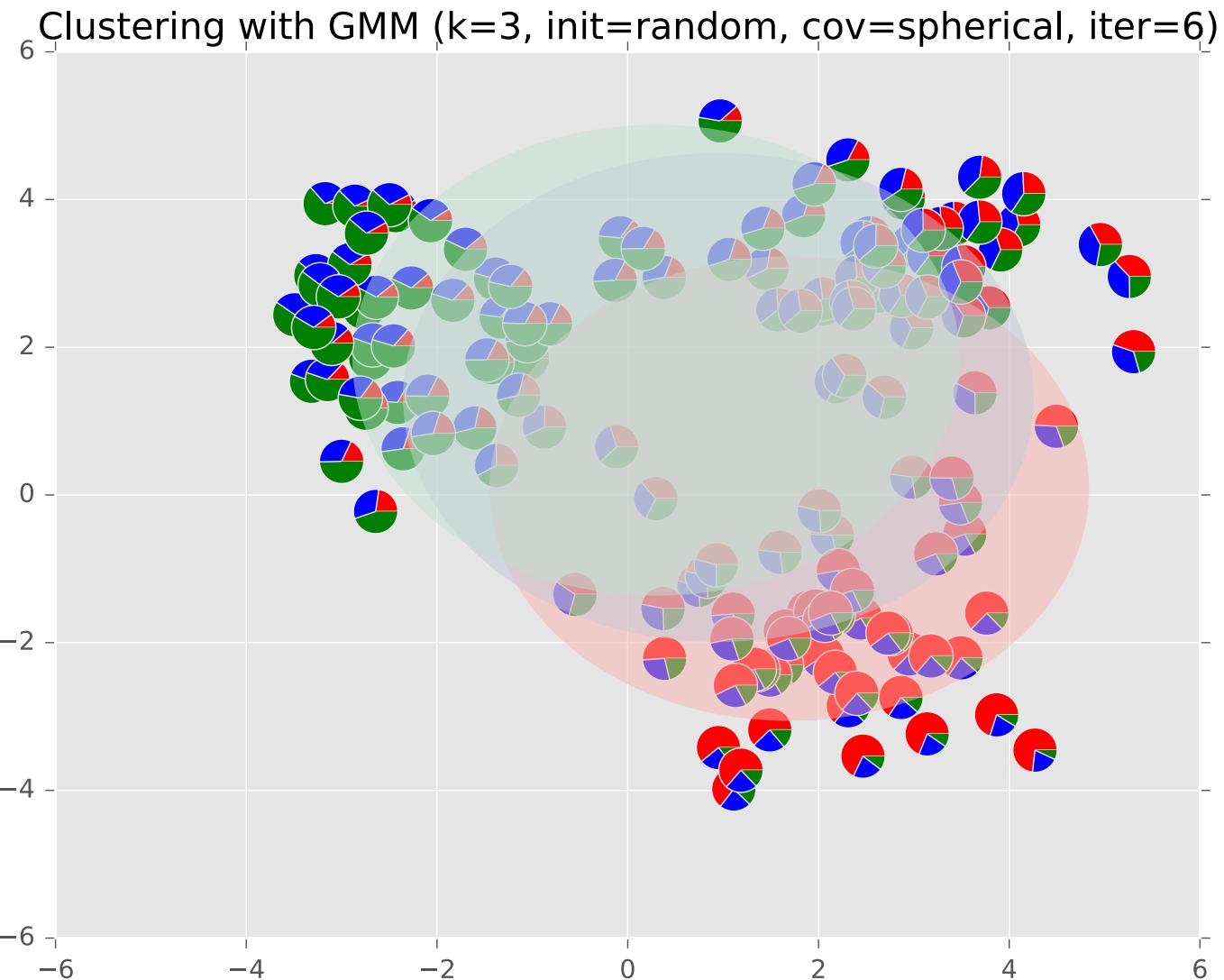
# Example: GMM



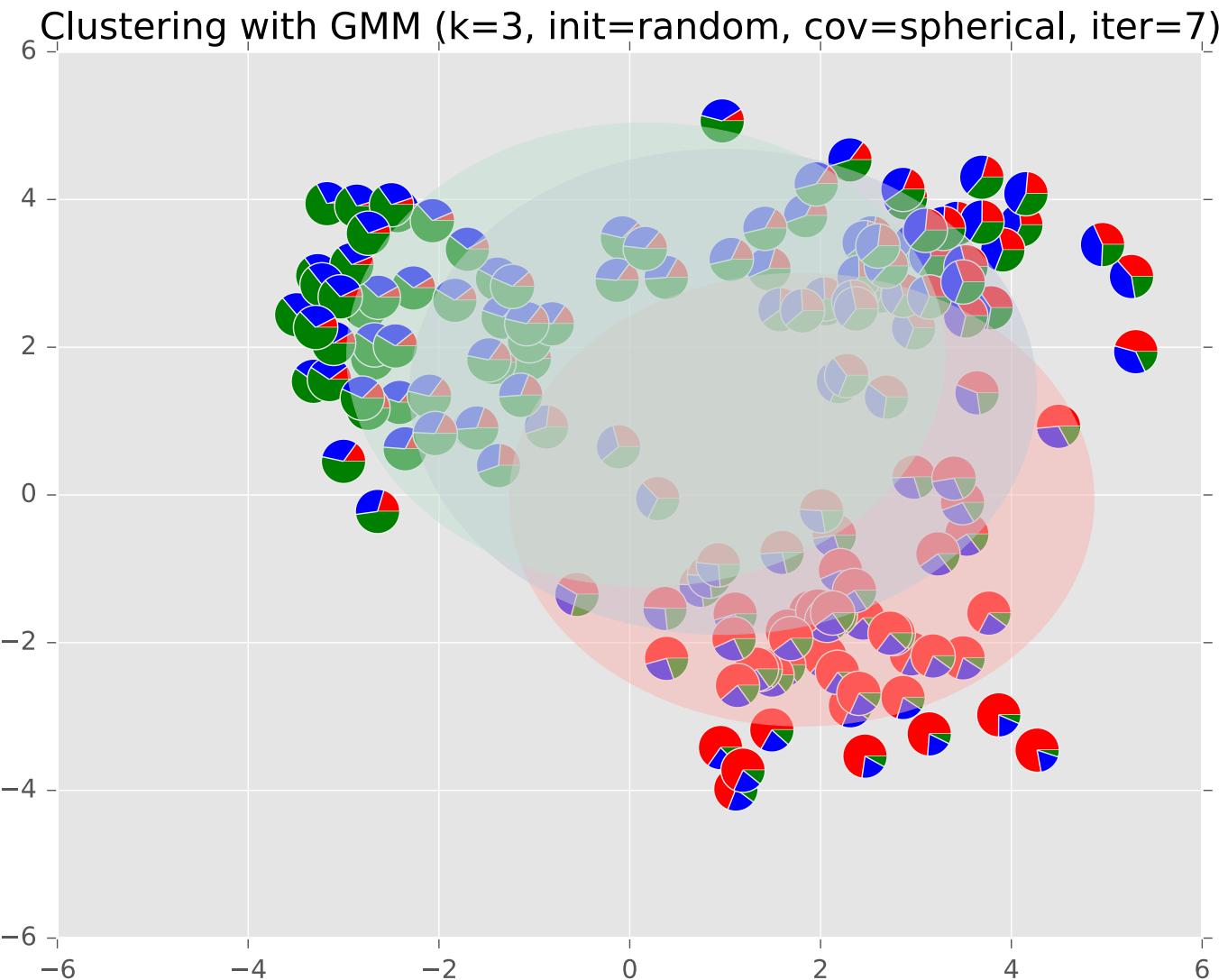
# Example: GMM



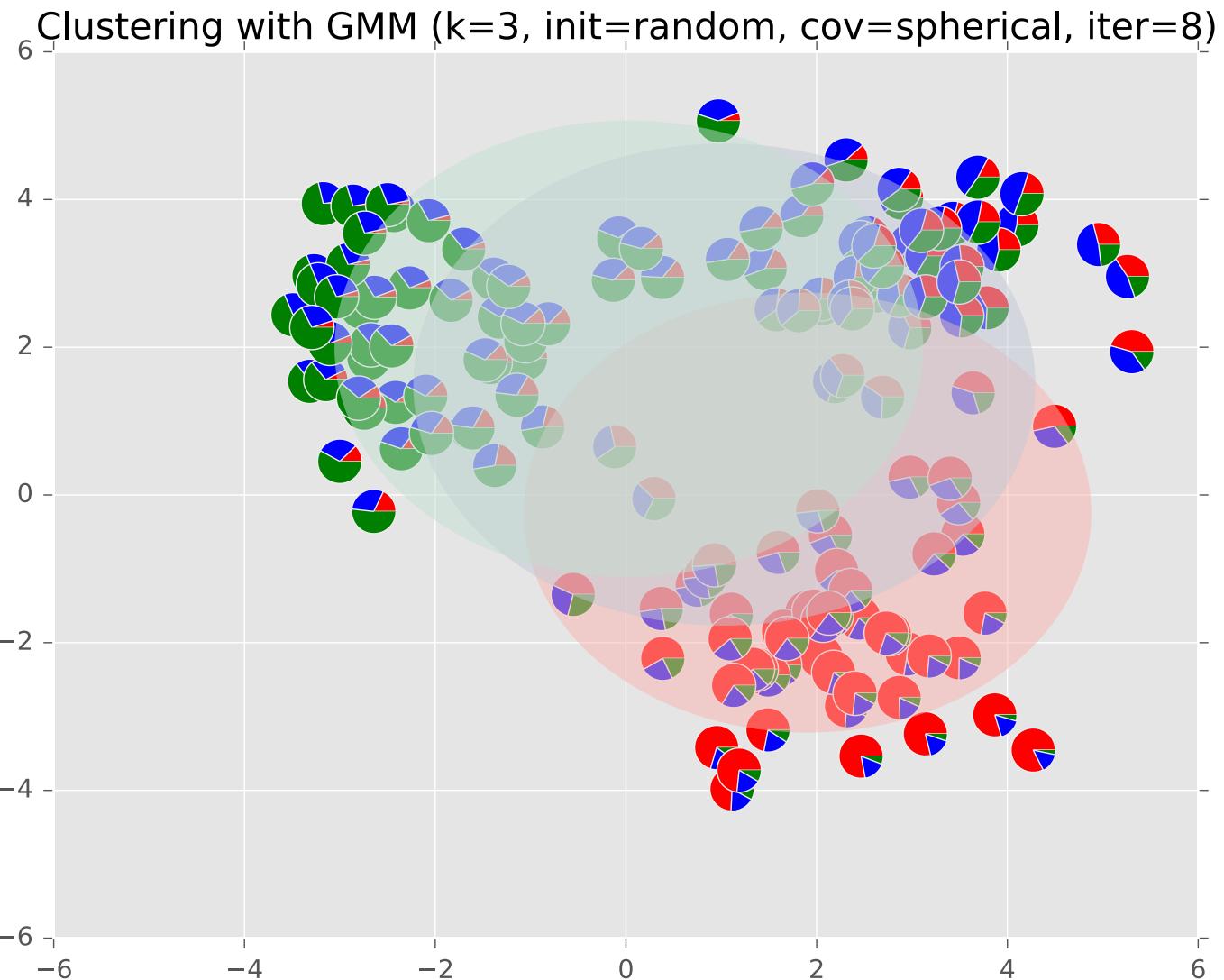
# Example: GMM



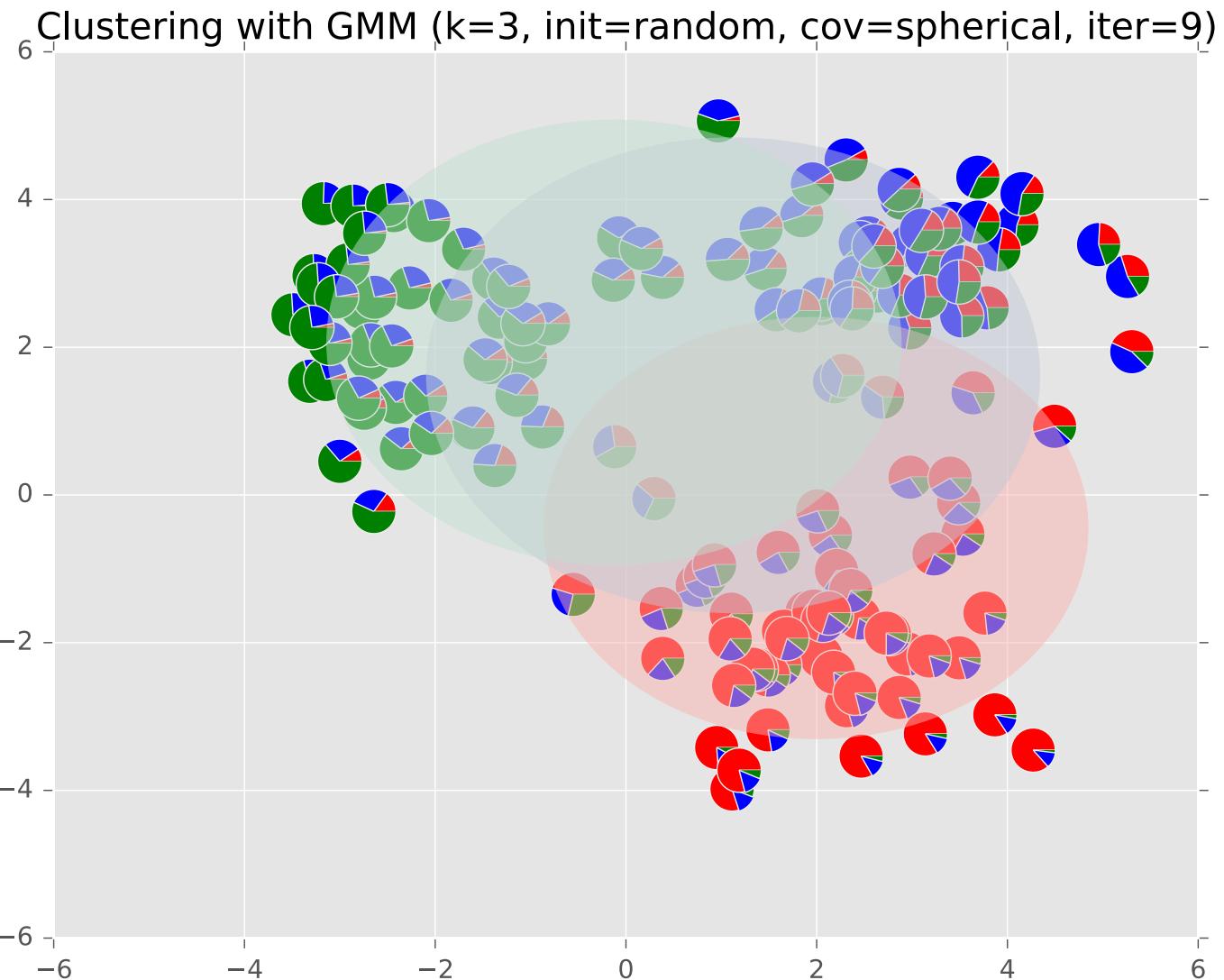
# Example: GMM



# Example: GMM

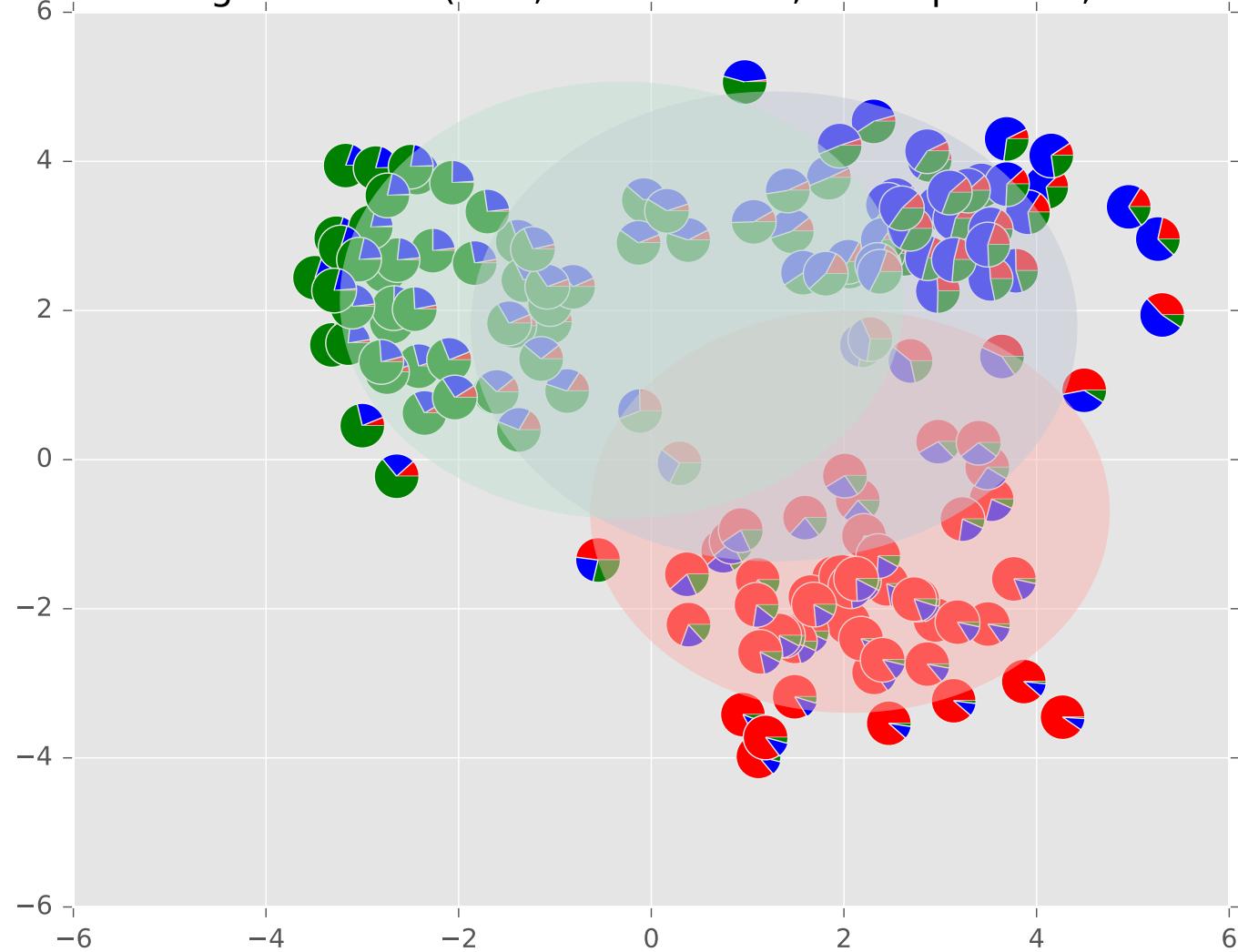


# Example: GMM



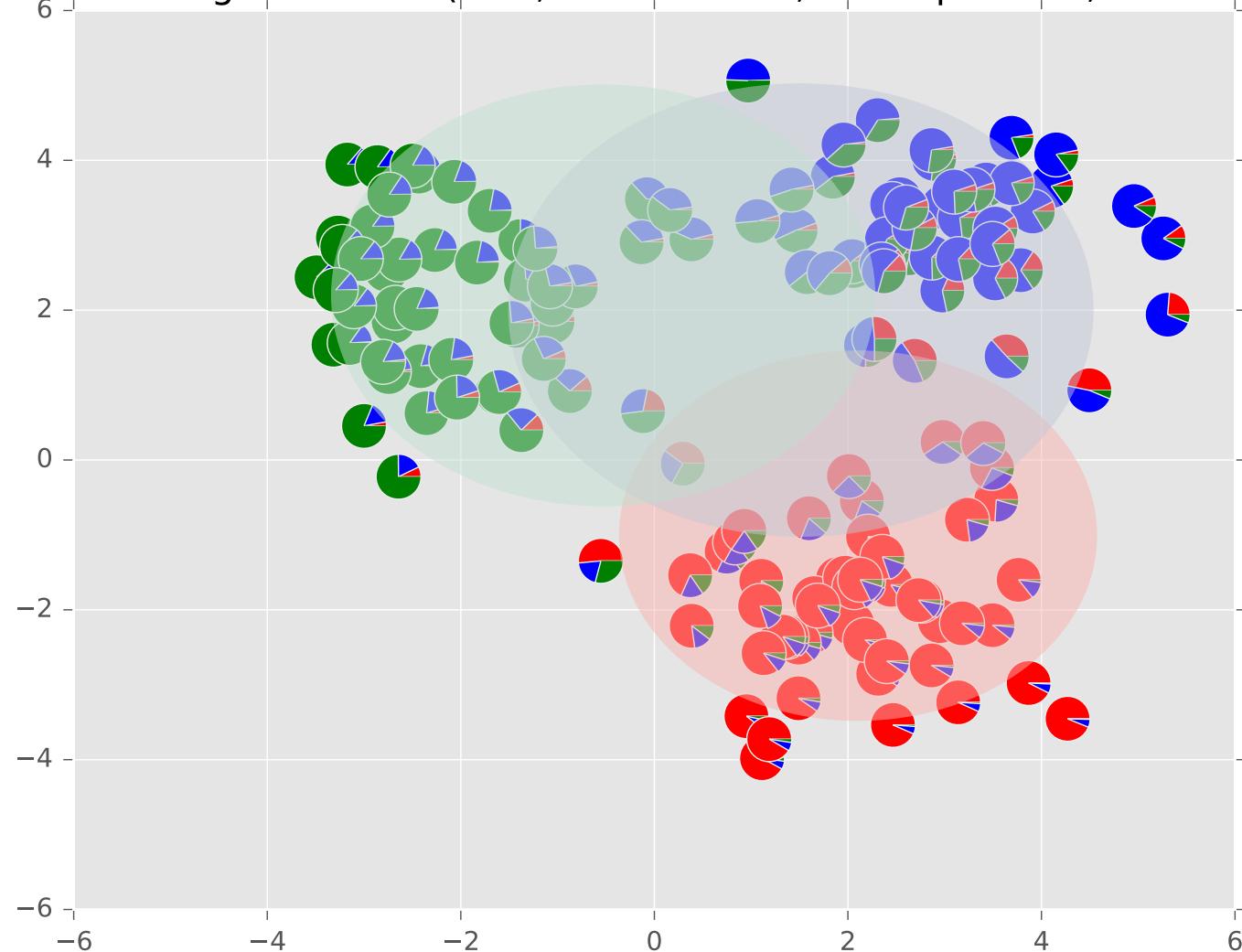
# Example: GMM

Clustering with GMM ( $k=3$ , init=random, cov=spherical, iter=10)



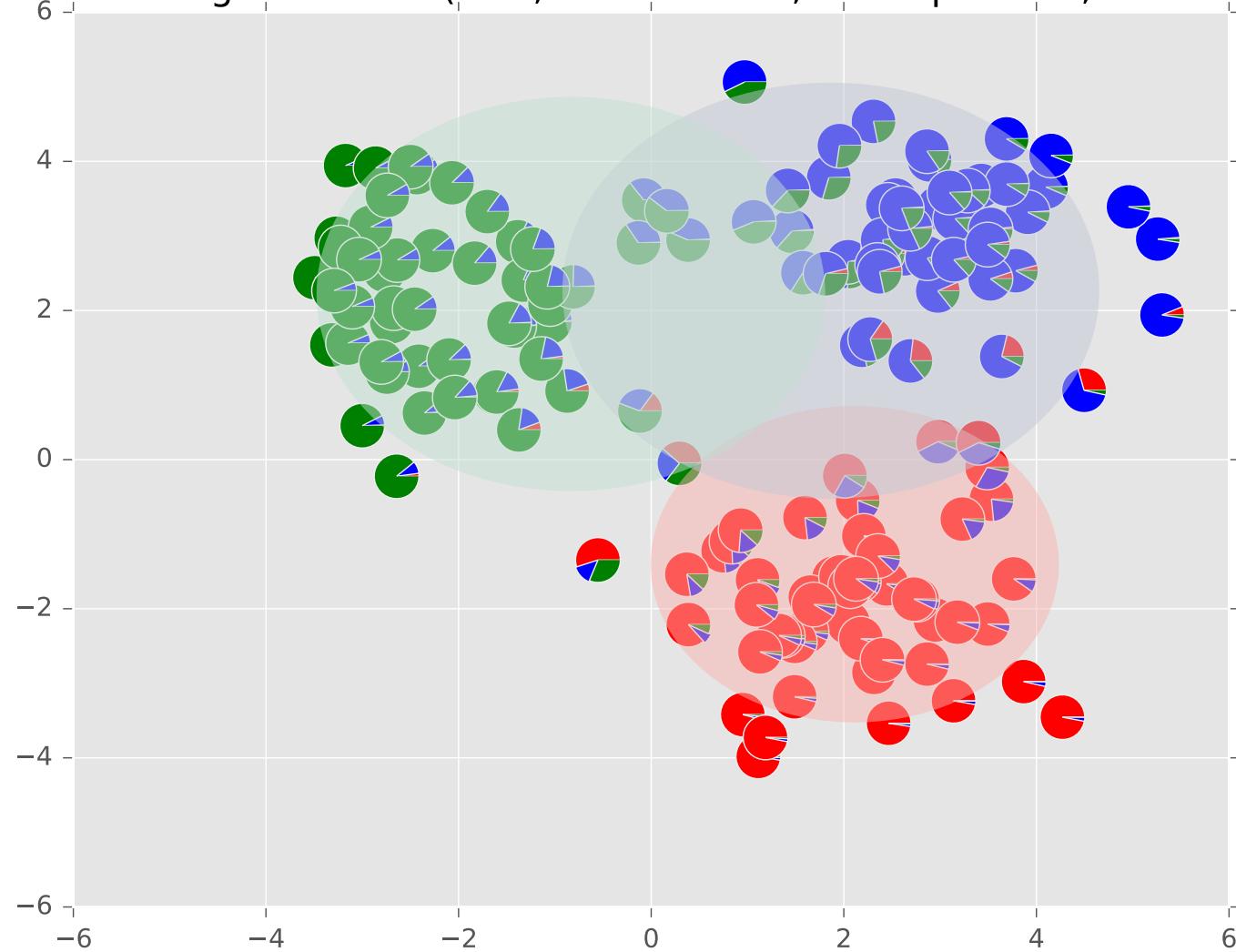
# Example: GMM

Clustering with GMM (k=3, init=random, cov=spherical, iter=11)

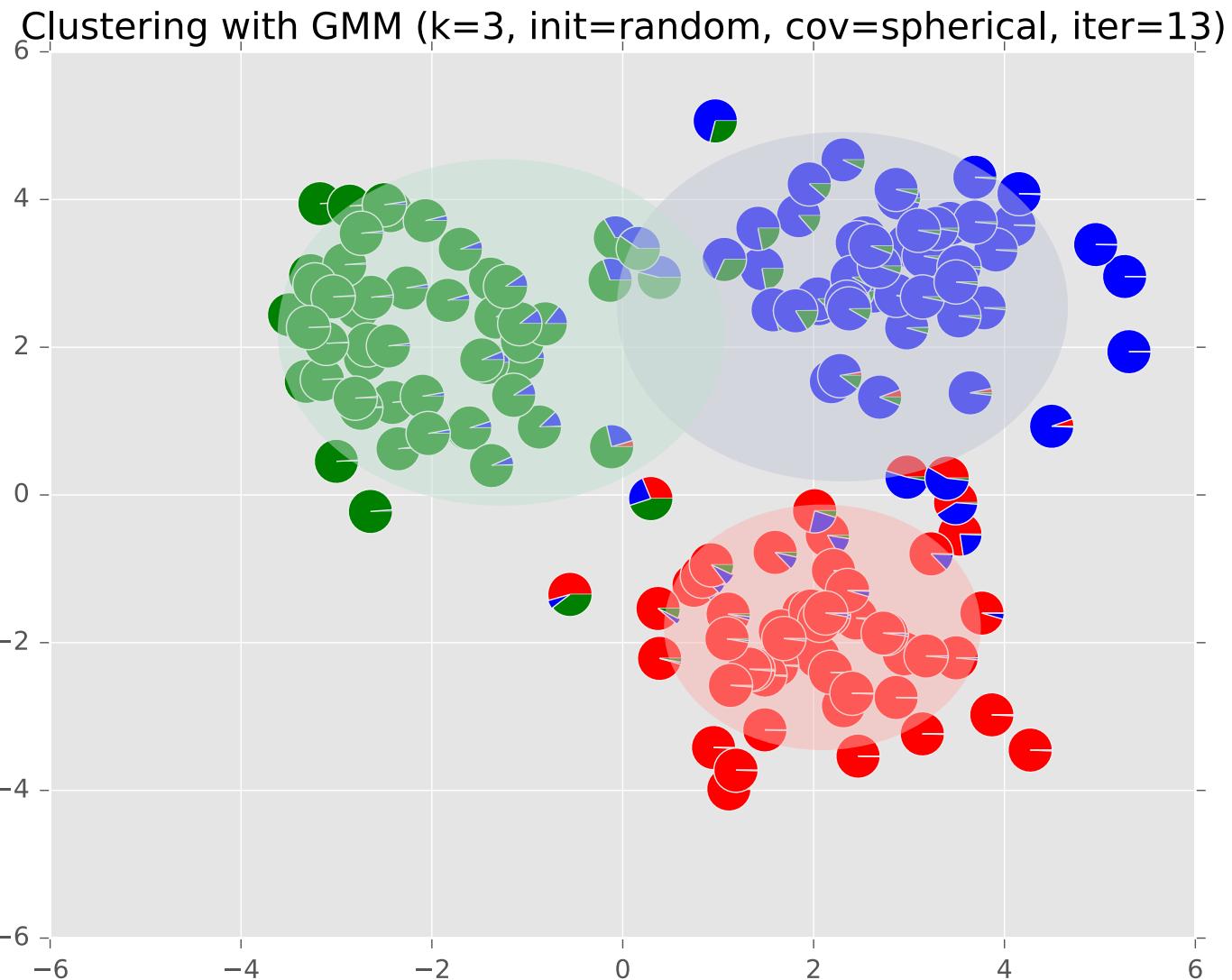


# Example: GMM

Clustering with GMM (k=3, init=random, cov=spherical, iter=12)

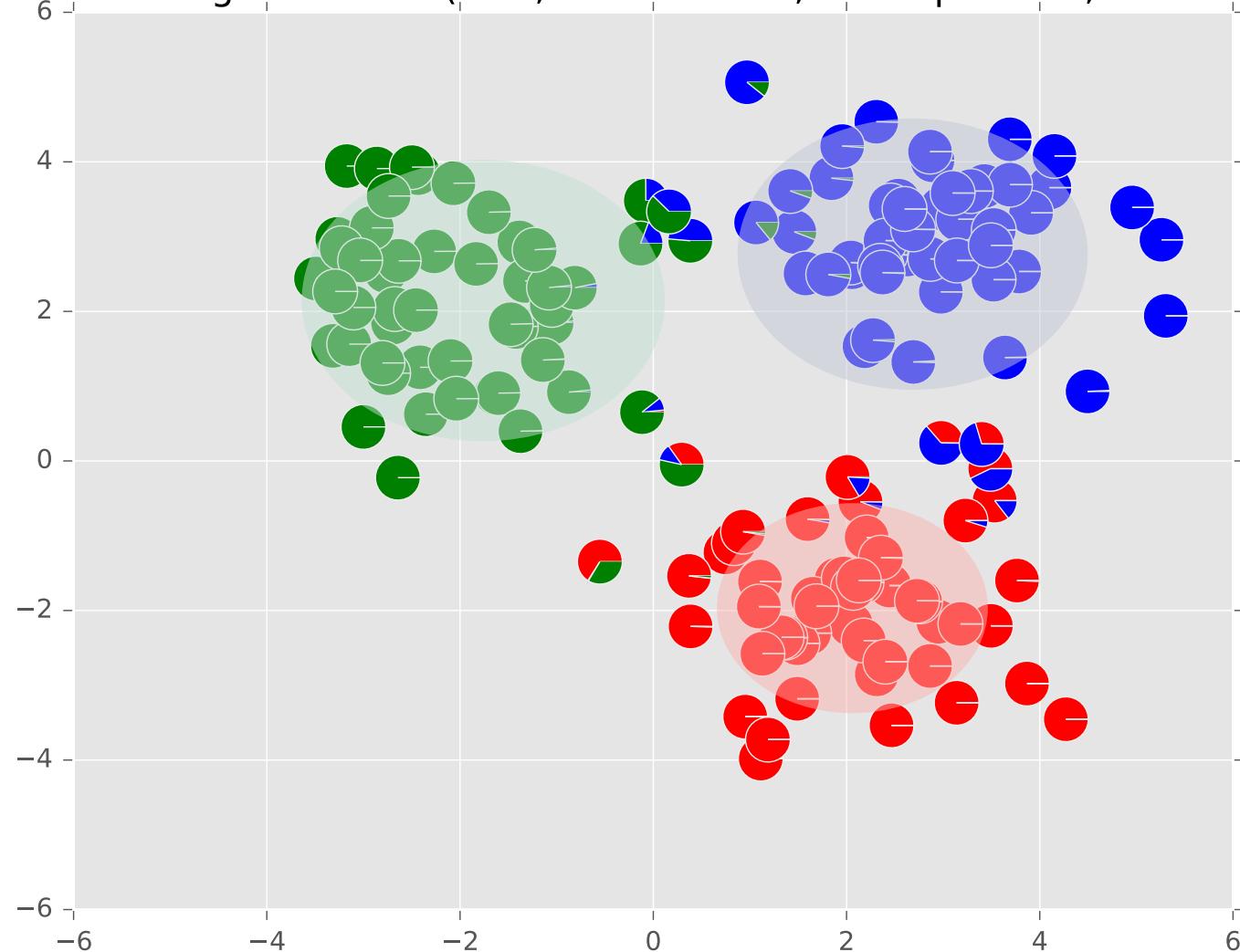


# Example: GMM

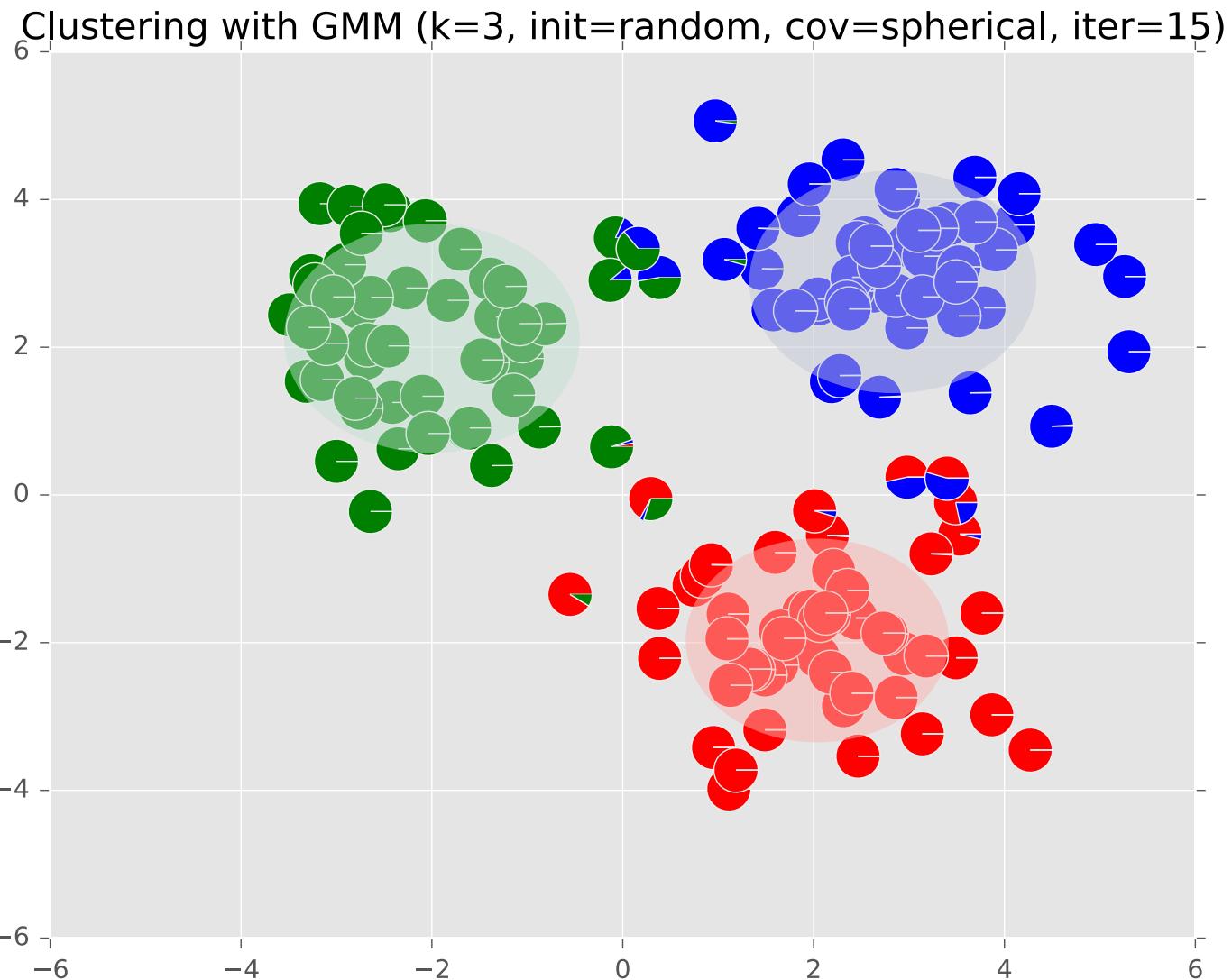


# Example: GMM

Clustering with GMM (k=3, init=random, cov=spherical, iter=14)

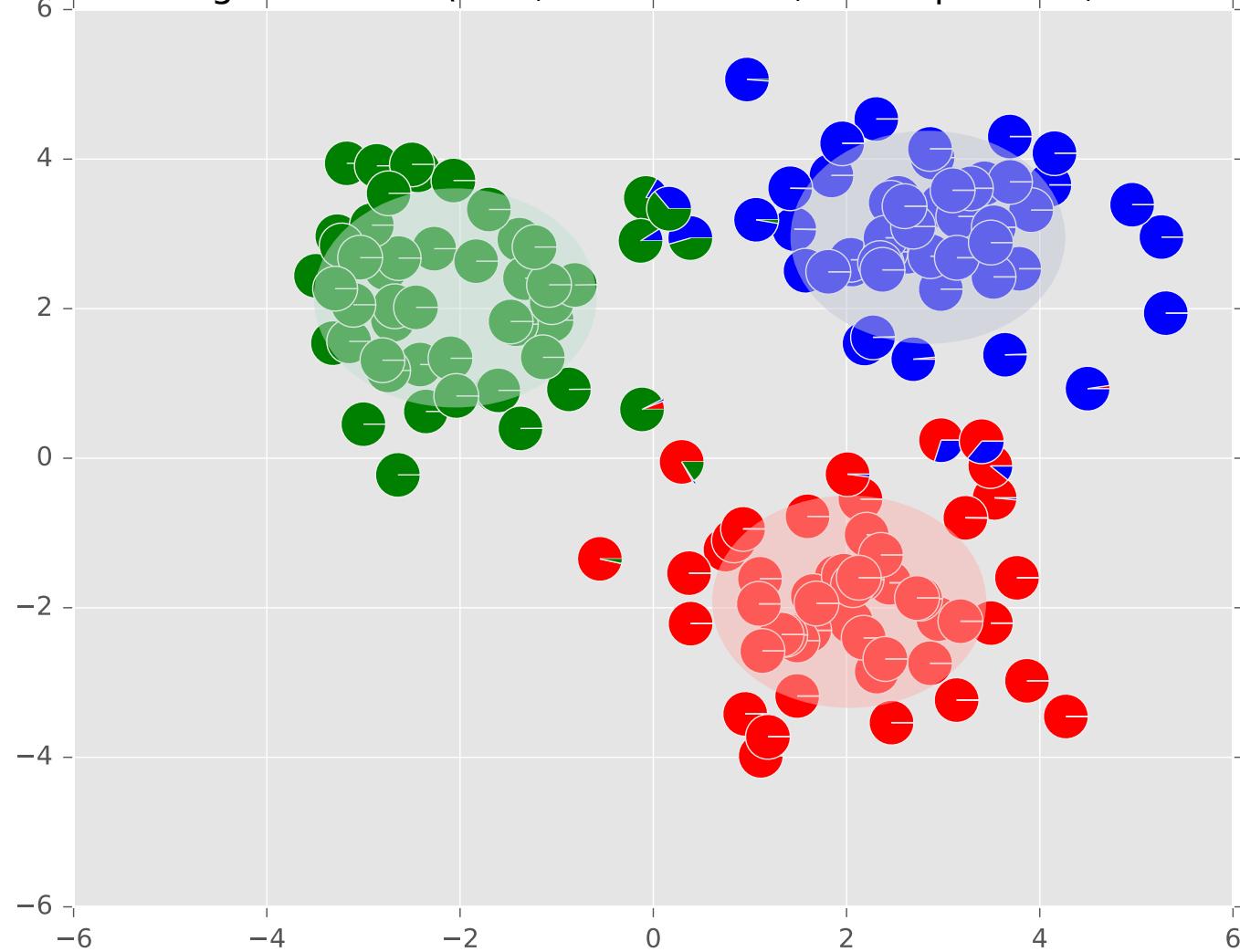


# Example: GMM



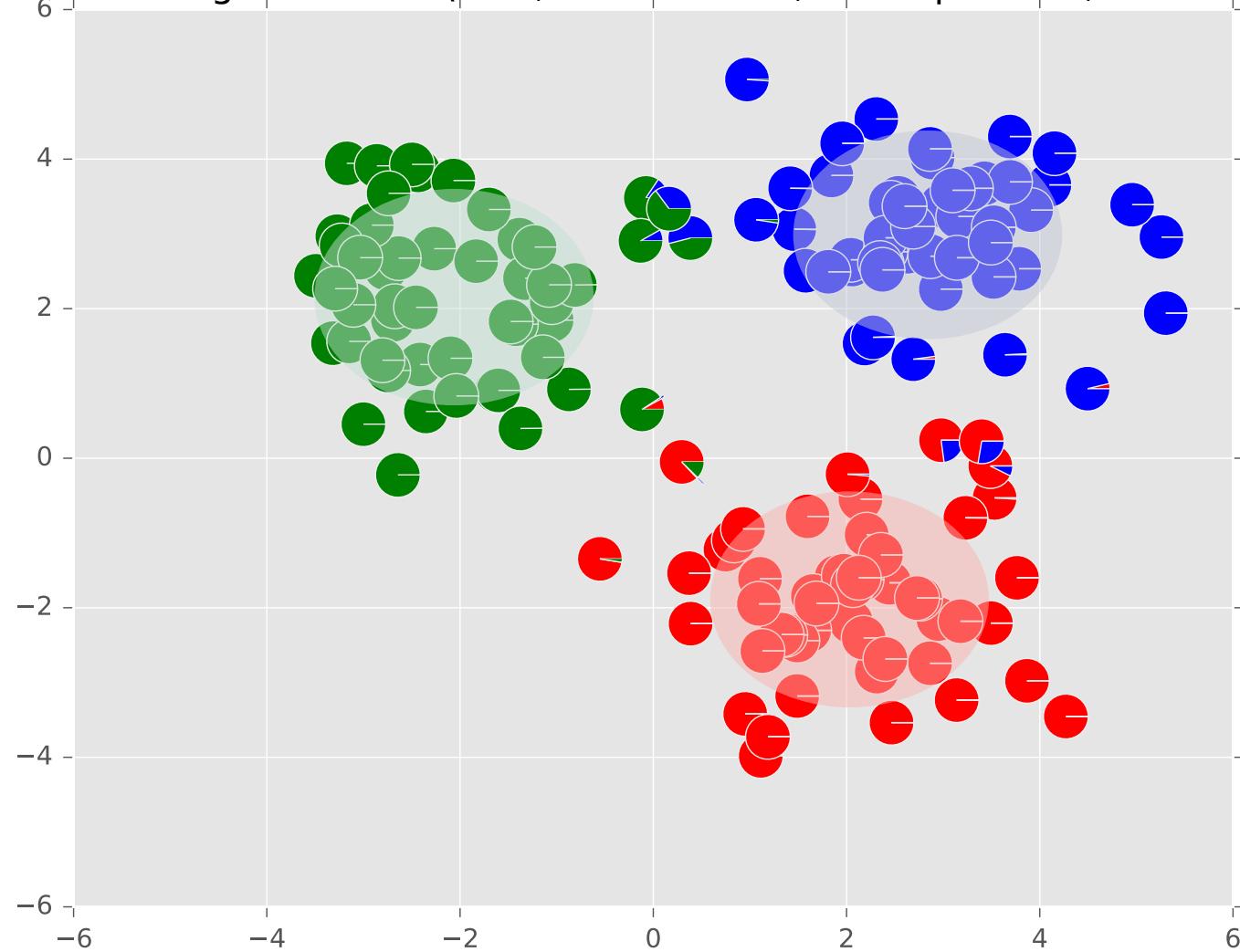
# Example: GMM

Clustering with GMM (k=3, init=random, cov=spherical, iter=16)

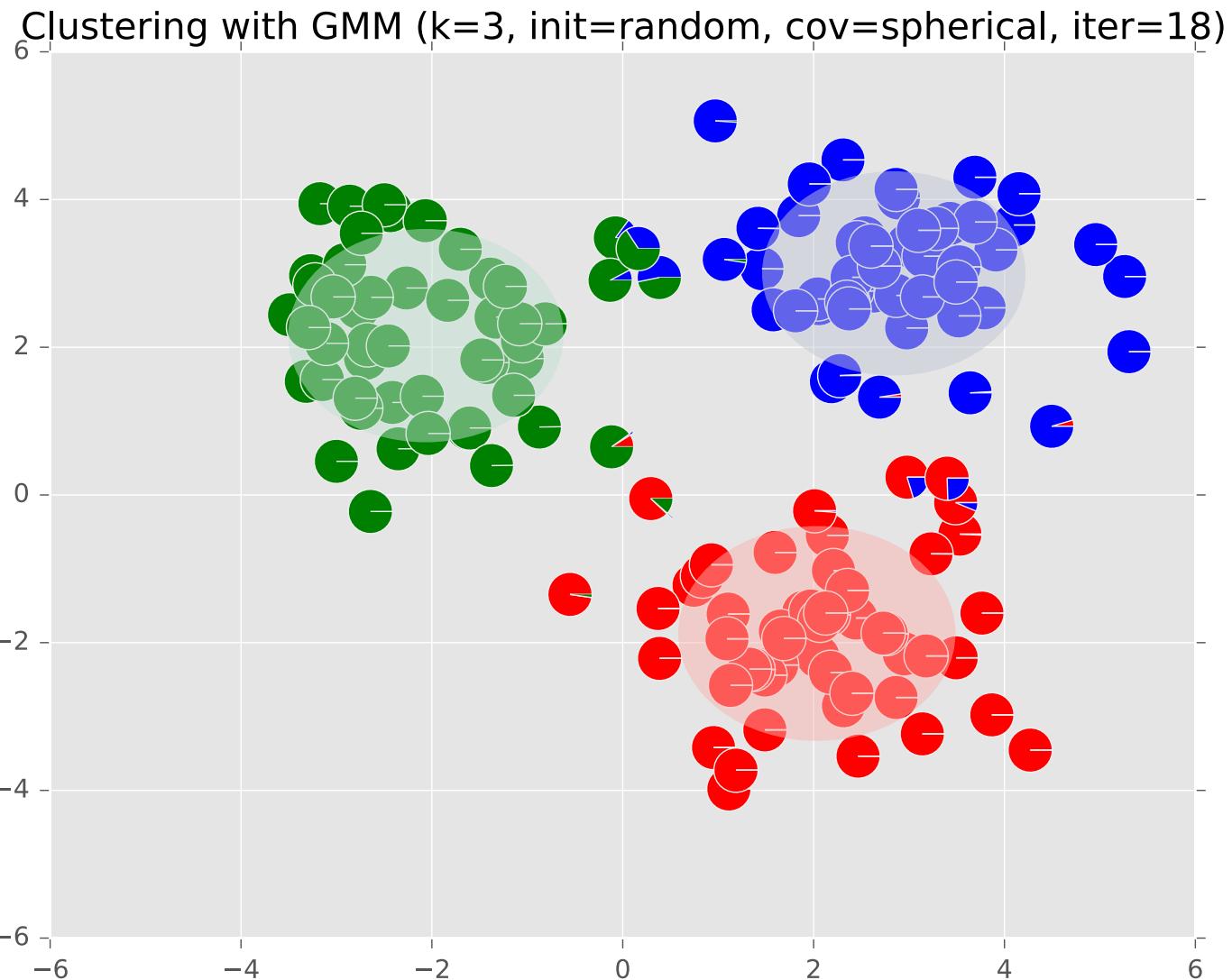


# Example: GMM

Clustering with GMM (k=3, init=random, cov=spherical, iter=17)

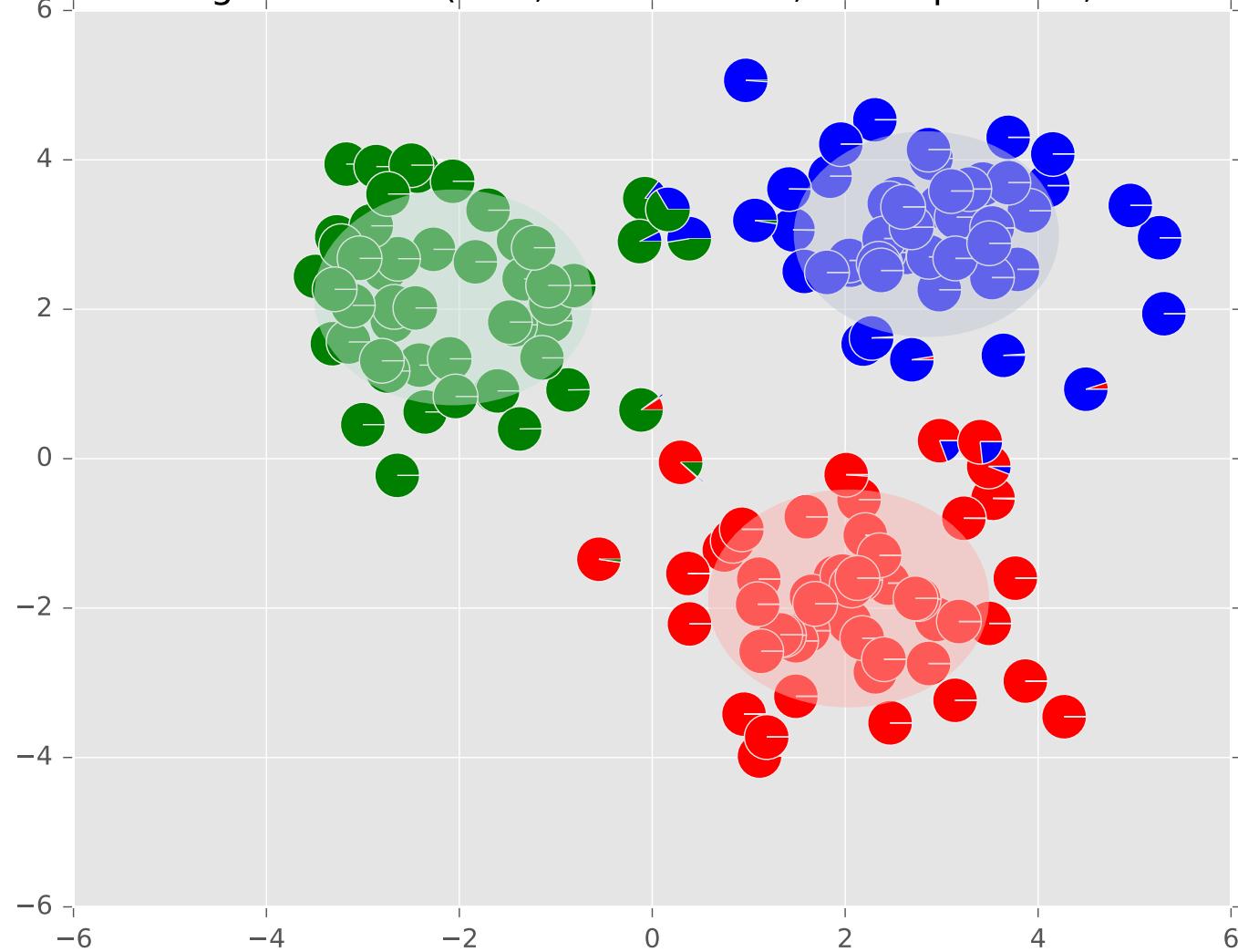


# Example: GMM



# Example: GMM

Clustering with GMM (k=3, init=random, cov=spherical, iter=19)



# K-Means vs. GMM

Convergence:

**K-Means** tends to **converge** much faster than a **GMM**

Speed:

Each iteration of **K-Means** is **computationally less intensive** than each iteration of a **GMM**

Initialization:

To **initialize** a **GMM**, we typically first run **K-Means** and use the resulting cluster centers as the means of the Gaussian components

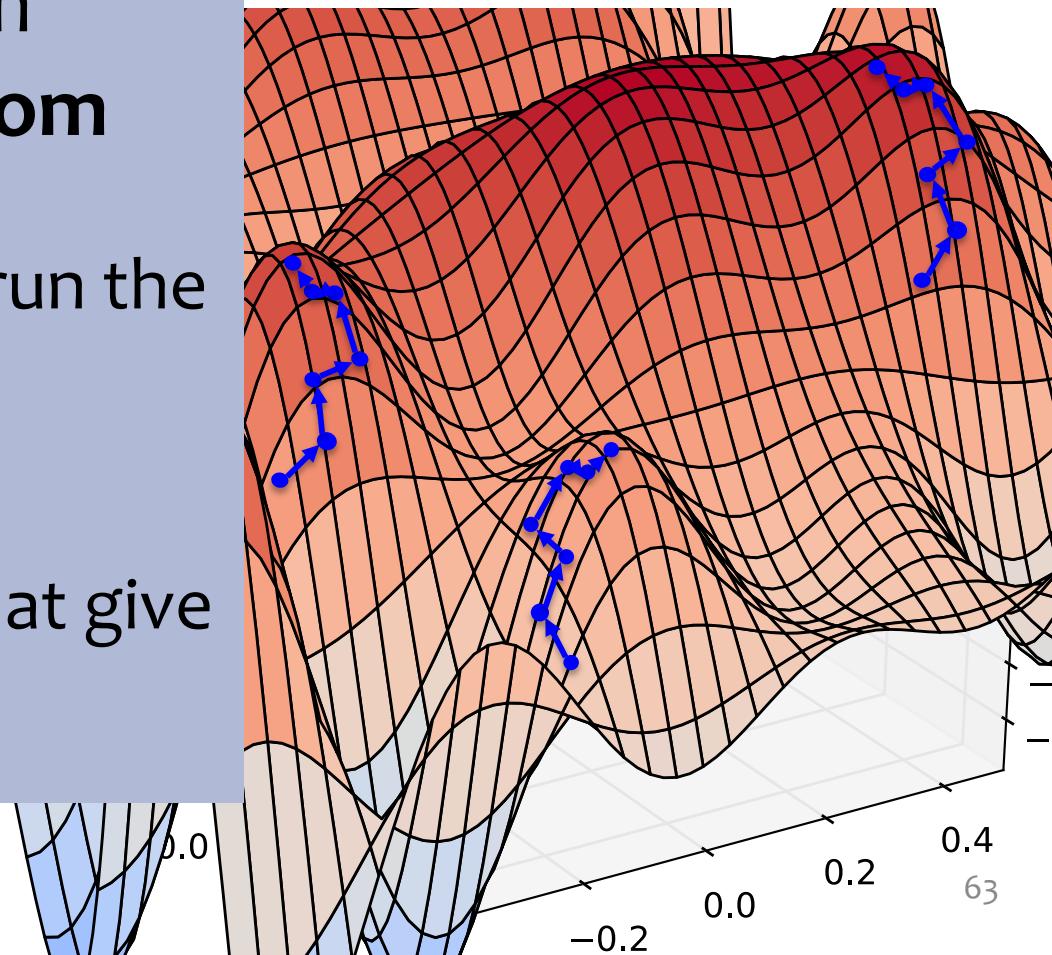
Output:

A **GMM** yields a **probability distribution** over the cluster assignment for each point; whereas **K-Means** gives a single **hard assignment**

# **PROPERTIES OF EM**

# Properties of (Variational) EM

- EM is trying to optimize a **nonconvex** function
- But EM is a **local** optimization algorithm
- Typical solution: **Random Restarts**
  - Just like K-Means, we run the algorithm many times
  - Each time initialize parameters randomly
  - Pick the parameters that give highest likelihood



# Variants of EM

- **Generalized EM:** Replace the M-Step by a single gradient-step that improves the likelihood
- **Monte Carlo EM:** Approximate the E-Step by sampling
- **Sparse EM:** Keep an “active list” of points (updated occasionally) from which we estimate the expected counts in the E-Step
- **Incremental EM / Stepwise EM:** If standard EM is described as a *batch* algorithm, these are the *online* equivalent
- **etc.**

# A Report Card for EM

- Some good things about EM:
  - no learning rate (step-size) parameter
  - automatically enforces parameter constraints
  - very fast for low dimensions
  - each iteration guaranteed to improve likelihood
- Some bad things about EM:
  - can get stuck in local minima
  - can be slower than conjugate gradient (especially near convergence)
  - requires expensive inference step
  - is a maximum likelihood/MAP method

# **VARIATIONAL EM**

# Variational EM

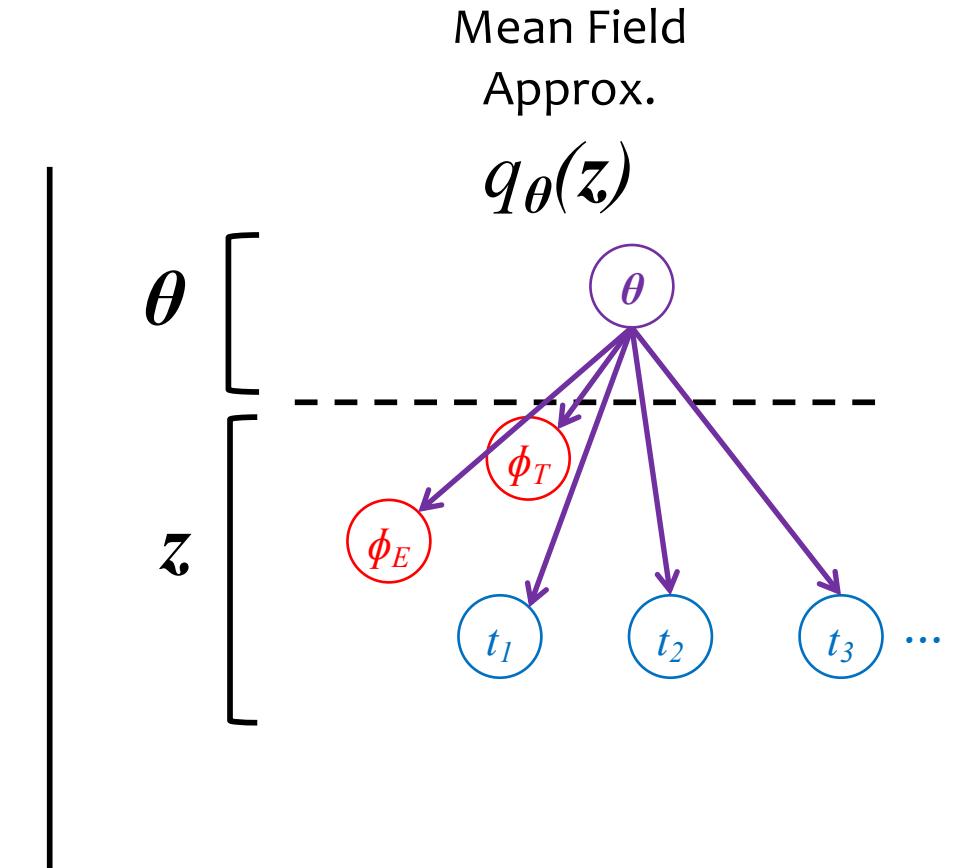
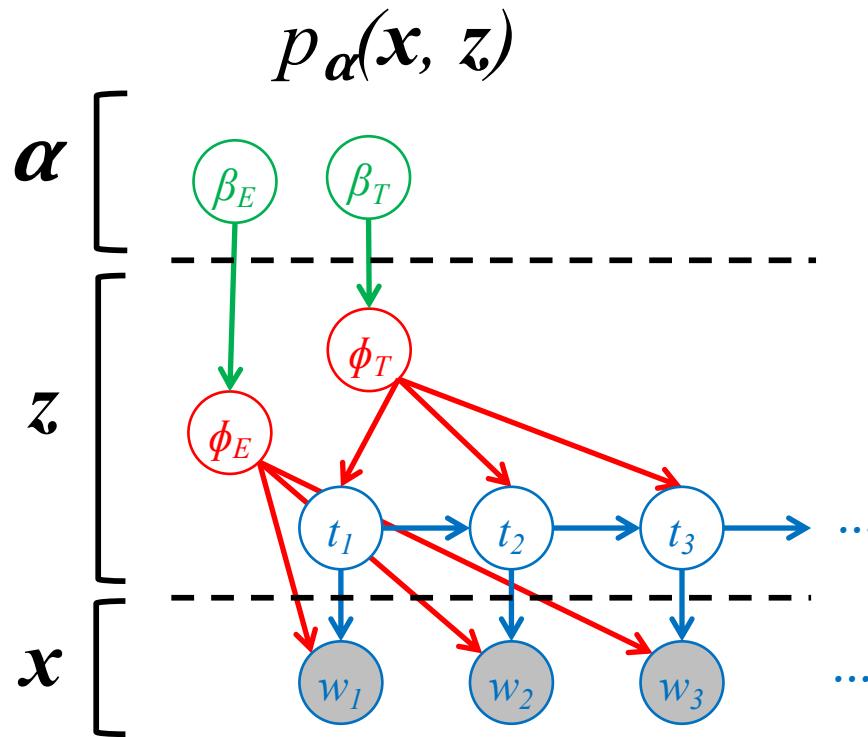
## **Whiteboard**

- Example: Unsupervised POS Tagging
- Variational Bayes
- Variational EM

# **UNDERSTANDING THE VARIATIONAL FRAMEWORK**

# Bayesian Unsupervised POS Tagging

- Given: sentences only (concatenated together into one long string)
- Goal: infer the POS tags for unlabeled sentences
- Model: Bayesian HMM



# Variational Inference Framework

	$\beta$	$\phi_E, \phi_T$	$t_{1:T}$	optimization problem	
1)	given	given	max		
2)	given	given	sum		
3)	given	max	sum		
4)	given	sum	sum		
5)	max	max	sum		

**Question:**

What is the corresponding optimization problem?

- A. Bayesian inference
- B. empirical Bayes
- C. MAP inference
- D. marginal inference
- E. MAP estimation

# **VARIATIONAL EM RESULTS**

# Unsupervised POS Tagging

## Bayesian Inference for HMMs

- **Task:** unsupervised POS tagging
- **Data:** 1 million words (i.e. unlabeled sentences) of WSJ text
- **Dictionary:** defines legal part-of-speech (POS) tags for each word type
- **Models:**
  - EM: standard HMM
  - VB: uncollapsed variational Bayesian HMM
  - Algo 1 (CVB): collapsed variational Bayesian HMM (strong indep. assumption)
  - Algo 2 (CVB): collapsed variational Bayesian HMM (weaker indep. assumption)
  - CGS: collapsed Gibbs Sampler for Bayesian HMM

Algo 1 mean field update:

$$q(z_t = k) \propto \frac{\mathbb{E}_{q(\mathbf{z}^{-t})}[C_{k,w}^{-t}] + \beta}{\mathbb{E}_{q(\mathbf{z}^{-t})}[C_{k,\cdot}^{-t}] + W\beta} \cdot \frac{\mathbb{E}_{q(\mathbf{z}^{-t})}[C_{z_{t-1},k}^{-t}] + \alpha}{\mathbb{E}_{q(\mathbf{z}^{-t})}[C_{z_{t-1},\cdot}^{-t}] + K\alpha} \cdot \frac{\mathbb{E}_{q(\mathbf{z}^{-t})}[C_{k,z_{t+1}}^{-t}] + \alpha + \mathbb{E}_{q(\mathbf{z}^{-t})}[\delta(z_{t-1} = k = z_{t+1})]}{\mathbb{E}_{q(\mathbf{z}^{-t})}[C_{k,\cdot}^{-t}] + K\alpha + \mathbb{E}_{q(\mathbf{z}^{-t})}[\delta(z_{t-1} = k)]}$$

CGS full conditional:

$$p(z_t = k | \mathbf{x}, \mathbf{z}^{-t}, \alpha, \beta) \propto \frac{C_{k,w}^{-t} + \beta}{C_{k,\cdot}^{-t} + W\beta} \cdot \frac{C_{z_{t-1},k}^{-t} + \alpha}{C_{z_{t-1},\cdot}^{-t} + K\alpha} \cdot \frac{C_{k,z_{t+1}}^{-t} + \alpha + \delta(z_{t-1} = k = z_{t+1})}{C_{k,\cdot}^{-t} + K\alpha + \delta(z_{t-1} = k)}$$

# Unsupervised POS Tagging

## Bayesian Inference for HMMs

- **Task:** unsupervised POS tagging
- **Data:** 1 million words (i.e. unlabeled sentences) of WSJ text
- **Dictionary:** defines legal part-of-speech (POS) tags for each word type
- **Models:**
  - EM: standard HMM
  - VB: uncollapsed variational Bayesian HMM
  - Algo 1 (CVB): collapsed variational Bayesian HMM (strong indep. assumption)
  - Algo 2 (CVB): collapsed variational Bayesian HMM (weaker indep. assumption)
  - CGS: collapsed Gibbs Sampler for Bayesian HMM

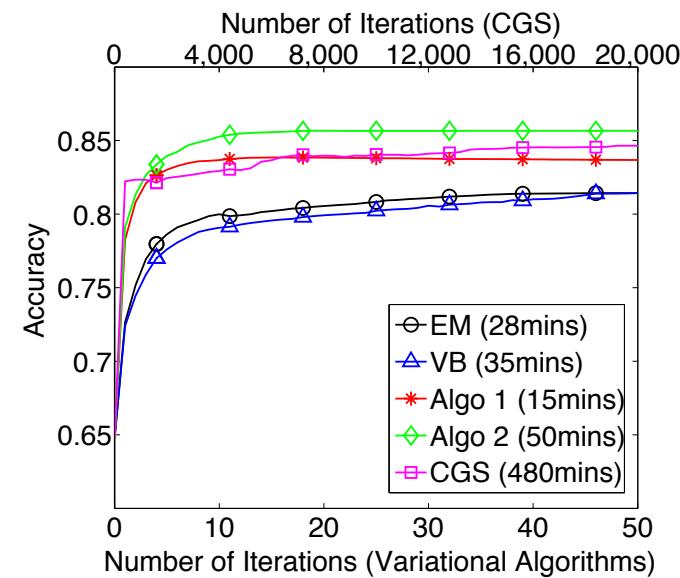
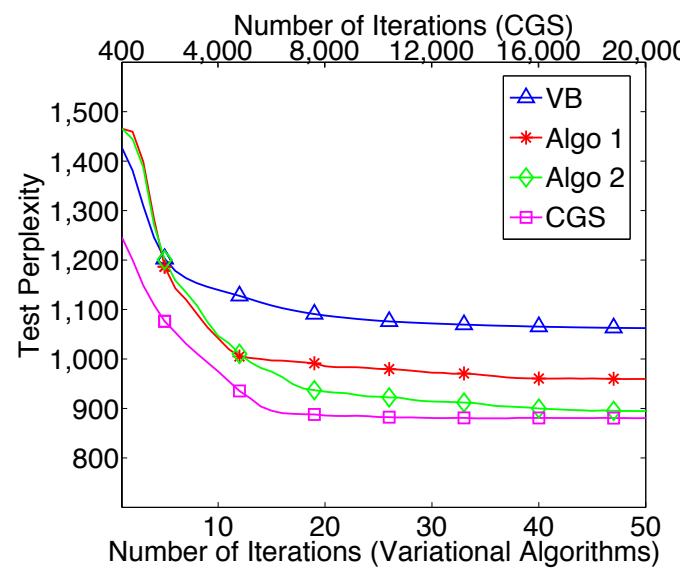


Figure from Wang & Blunsom (2013)

# Unsupervised POS Tagging

## Bayesian Inference for HMMs

- **Task:** unsupervised POS tagging
- **Data:** 1 million words (i.e. unlabeled sentences) of WSJ text
- **Dictionary:** defines legal part-of-speech (POS) tags for each word type
- **Models:**
  - EM: standard HMM
  - VB: uncollapsed variational Bayesian HMM
  - Algo 1 (CVB): collapsed variational Bayesian HMM (strong indep. assumption)
  - Algo 2 (CVB): collapsed variational Bayesian HMM (weaker indep. assumption)
  - CGS: collapsed Gibbs Sampler for Bayesian HMM

## Speed:

- ⊖ EM (28mins)
- △ VB (35mins)
- \* Algo 1 (15mins)
- ◇ Algo 2 (50mins)
- CGS (480mins)

- EM is slow b/c of log-space computations
- VB is slow b/c of digamma computations
- Algo 1 (CVB) is the fastest!
- Algo 2 (CVB) is slow b/c it computes dynamic parameters
- CGS: an order of magnitude slower than any deterministic algorithm

# Stochastic Variational Bayesian HMM

- **Task:** Human Chromatin Segmentation
- **Goal:** unsupervised segmentation of the genome
- **Data:** from ENCODE, “250 million observations consisting of twelve assays carried out in the chronic myeloid leukemia cell line K562”
- **Metric:** “the false discovery rate (FDR) of predicting active promoter elements in the sequence”
- **Models:**
  - DBN HMM: dynamic Bayesian HMM trained with standard EM
  - SVIHMM: stochastic variational inference for a Bayesian HMM
- **Main Takeaway:**
  - the two models perform at similar levels of FDR
  - SVIHMM takes **one hour**
  - DBNHMM takes **days**

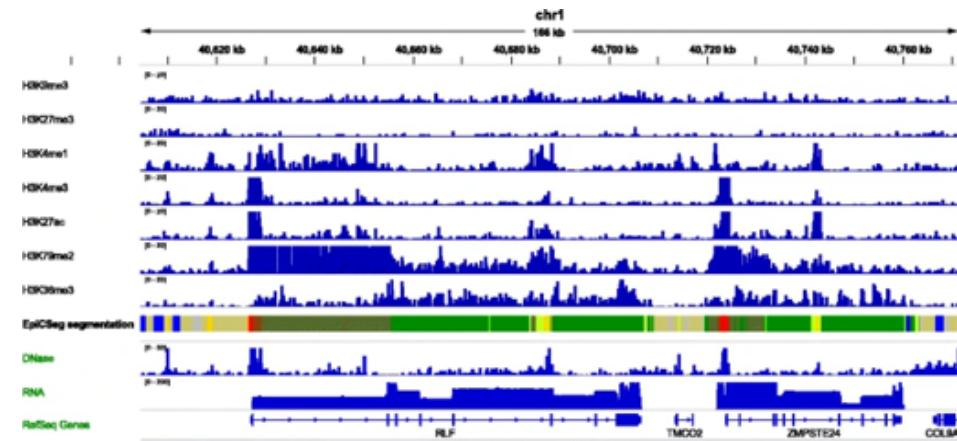


Figure from Foti et al. (2014)

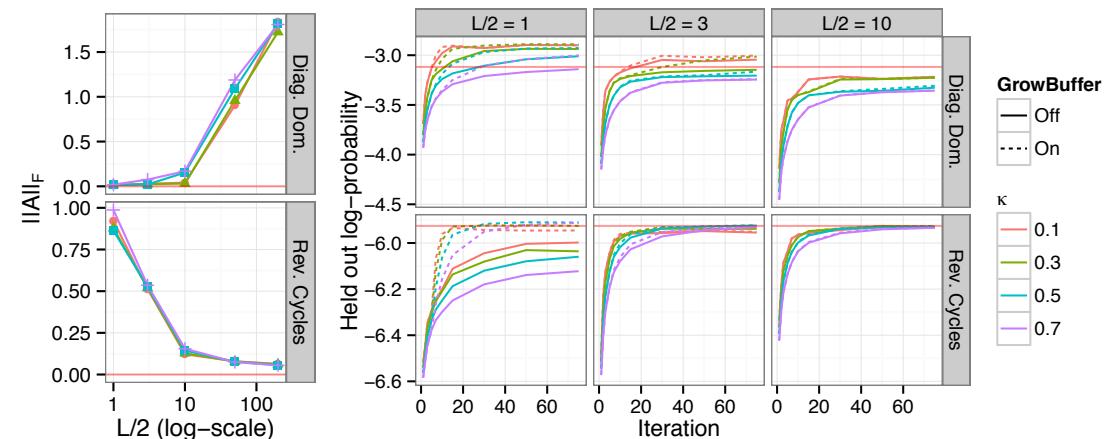


Figure from Mammana & Chung (2015)

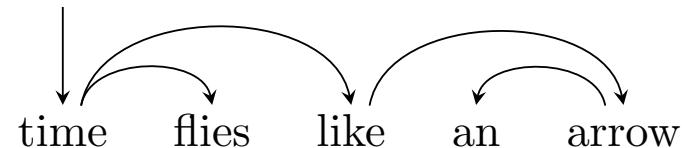
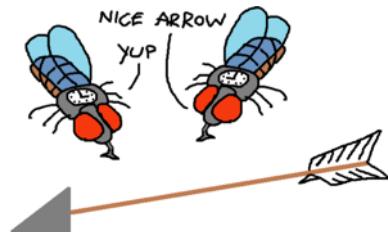
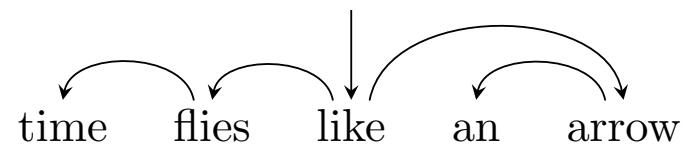
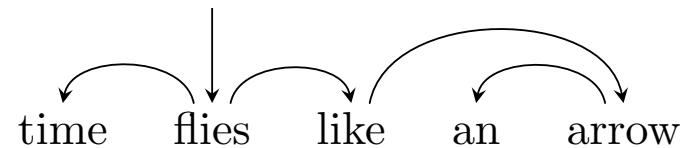
# Grammar Induction

**Question:** Can maximizing (unsupervised) marginal likelihood produce useful results?

**Answer:** Let's look at an example...

- **Babies** learn the syntax of their **native language** (e.g. English) just by **hearing** many sentences
- Can a **computer** similarly learn syntax of a **human language** just by looking at lots of example sentences?
  - This is the problem of Grammar Induction!
  - It's an unsupervised learning problem
  - We try to recover the **syntactic structure** for each sentence without any supervision

# Grammar Induction



...



No semantic interpretation

# Grammar Induction

**Training Data:** Sentences only, without parses

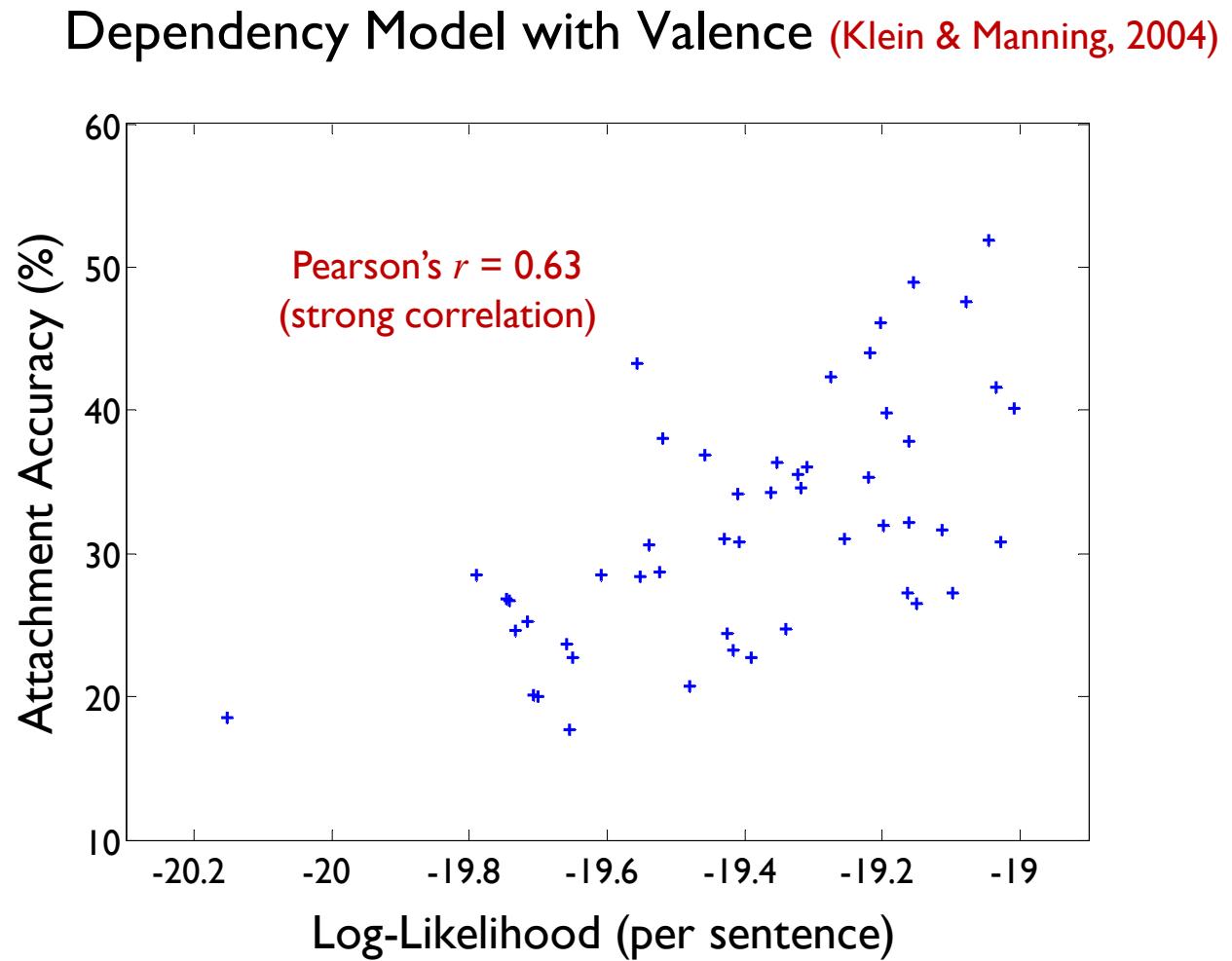
Sample 1:	time	flies	like	an	arrow	$x^{(1)}$
Sample 2:	real	flies	like	soup		$x^{(2)}$
Sample 3:	flies	fly	with	their	wings	$x^{(3)}$
Sample 4:	with	time	you	will	see	$x^{(4)}$

**Test Data:** Sentences with parses, so we can evaluate accuracy

# Grammar Induction

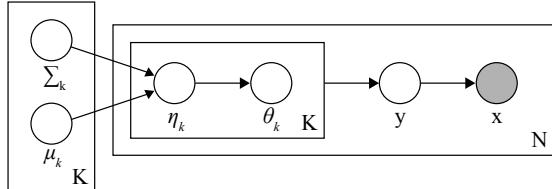
**Q:** Does likelihood correlate with accuracy on a task we care about?

**A:** Yes, but there is still a wide range of accuracies for a particular likelihood value



# Grammar Induction

## Graphical Model for Logistic Normal Probabilistic Grammar



$y = \text{syntactic parse}$   
 $x = \text{observed sentence}$

## Settings:

**EM** Maximum likelihood estimate of  $\theta$  using the EM algorithm to optimize  $p(\mathbf{x} | \theta)$  [14].

**EM-MAP** Maximum *a posteriori* estimate of  $\theta$  using the EM algorithm and a fixed symmetric Dirichlet prior with  $\alpha > 1$  to optimize  $p(\mathbf{x}, \theta | \alpha)$ . Tune  $\alpha$  to maximize the likelihood of an unannotated development dataset, using grid search over [1.1, 30].

**VB-Dirichlet** Use variational Bayes inference to estimate the posterior distribution  $p(\theta | \mathbf{x}, \alpha)$ , which is a Dirichlet. Tune the symmetric Dirichlet prior's parameter  $\alpha$  to maximize the likelihood of an unannotated development dataset, using grid search over [0.0001, 30]. Use the mean of the posterior Dirichlet as a point estimate for  $\theta$ .

**VB-EM-Dirichlet** Use variational Bayes EM to optimize  $p(\mathbf{x} | \alpha)$  with respect to  $\alpha$ . Use the mean of the learned Dirichlet as a point estimate for  $\theta$  (similar to [5]).

**VB-EM-Log-Normal** Use variational Bayes EM to optimize  $p(\mathbf{x} | \mu, \Sigma)$  with respect to  $\mu$  and  $\Sigma$ . Use the (exponentiated) mean of this Gaussian as a point estimate for  $\theta$ .

## Results:

	attachment accuracy (%)					
	Viterbi decoding			MBR decoding		
	$ \mathbf{x}  \leq 10$	$ \mathbf{x}  \leq 20$	all	$ \mathbf{x}  \leq 10$	$ \mathbf{x}  \leq 20$	all
Attach-Right	38.4	33.4	31.7	38.4	33.4	31.7
EM	45.8	39.1	34.2	46.1	39.9	35.9
EM-MAP, $\alpha = 1.1$	45.9	39.5	34.9	46.2	40.6	36.7
VB-Dirichlet, $\alpha = 0.25$	46.9	40.0	35.7	47.1	41.1	37.6
VB-EM-Dirichlet	45.9	39.4	34.9	46.1	40.6	36.9
VB-EM-Log-Normal, $\Sigma_k^{(0)} = \mathbf{I}$	56.6	43.3	37.4	59.1	<b>45.9</b>	39.9
VB-EM-Log-Normal, families	<b>59.3</b>	<b>45.1</b>	<b>39.0</b>	<b>59.4</b>	<b>45.9</b>	<b>40.5</b>

Table 1: Attachment accuracy of different learning methods on unseen test data from the Penn Treebank of varying levels of difficulty imposed through a length filter. Attach-Right attaches each word to the word on its right and the last word to \$. EM and EM-MAP with a Dirichlet prior ( $\alpha > 1$ ) are reproductions of earlier results [14, 18].