



10-708 Probabilistic Graphical Models

Machine Learning Department
School of Computer Science
Carnegie Mellon University



Variational Inference

+

Learning Partially Observed Graphical Models

+

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

MEAN FIELD VARIATIONAL INFERENCE

Variational Inference

Whiteboard

- Coordinate Ascent Variational Inference (CAVI) Algorithm
 - Connecting CAVI to BP and Gibbs sampling
 - Computing marginals from a trained mean field approximation
- CAVI algorithm derivation
 - Chain rule decomposition of $\log p(x, z)$
 - Decomposing the entropy
 - Decomposing the ELBO
 - Derivatives and closed form solution

CAVI Algorithm

Coordinate Ascent Variational Inference (CAVI)

- here we assume a **mean field** approximation
- application of **coordinate ascent** to maximization of ELBO
- converges to a **local optimum** of the **nonconvex** ELBO objective

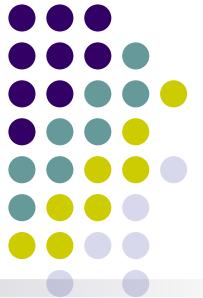
```
1: procedure CAVI( $p_\alpha$ )
2:   Let  $q_\theta(\mathbf{z}) = \prod_{t=1}^T q_t(z_t)$                                  $\triangleright$  Mean field approx.
3:   while ELBO( $q_\theta$ ) has not converged do
4:     for  $t \in \{1, \dots, T\}$  do                                      $\triangleright$  For each variable
5:       Set  $q_t(z_t) \propto \exp(E_{q_{\neg t}}[\log p_\alpha(z_t \mid z_{\neg t}, \mathbf{x})])$ 
6:       while keeping all  $\{q_s(\cdot)\}_{s \neq t}$  fixed
7:     Compute  $\text{ELBO}(q_\theta) = E_{q_\theta(\mathbf{z})} [\log p_\alpha(\mathbf{x}, \mathbf{z})] - E_{q_\theta(\mathbf{z})} [\log q_\theta(\mathbf{z})]$ 
8:   return  $q_\theta$ 
```

Variational Inference

Whiteboard

- Computing the CAVI update
 - Multinomial full conditionals
- Example: two variable factor graph
 - Joint distribution
 - Mean Field Variational Inference
 - Gibbs Sampling

EXPONENTIAL FAMILY DISTRIBUTION



Exponential family, a basic building block

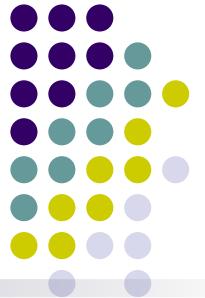
- For a numeric random variable X

$$\begin{aligned} p(x | \eta) &= h(x) \exp\{\eta^T T(x) - A(\eta)\} \\ &= \frac{1}{Z(\eta)} h(x) \exp\{\eta^T T(x)\} \end{aligned}$$

is an **exponential family distribution** with natural (canonical) parameter η

- Function $T(x)$ is a *sufficient statistic*.
- Function $A(\eta) = \log Z(\eta)$ is the log normalizer.
- Examples: Bernoulli, multinomial, Gaussian, Poisson, gamma,...

Example: Multivariate Gaussian Distribution



- For a continuous vector random variable $X \in \mathbb{R}^k$:

$$\begin{aligned}
 p(x|\mu, \Sigma) &= \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right\} \\
 &= \frac{1}{(2\pi)^{k/2}} \exp\left\{-\frac{1}{2} \text{tr}(\Sigma^{-1} xx^T) + \mu^T \Sigma^{-1} x - \frac{1}{2} \mu^T \Sigma^{-1} \mu - \log|\Sigma|\right\}
 \end{aligned}$$

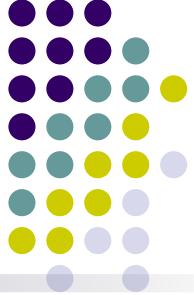
Moment parameter

- Exponential family representation

$$\begin{aligned}
 \eta &= [\Sigma^{-1} \mu; -\frac{1}{2} \text{vec}(\Sigma^{-1})] = [\eta_1, \text{vec}(\eta_2)], \quad \eta_1 = \Sigma^{-1} \mu \text{ and } \eta_2 = -\frac{1}{2} \Sigma^{-1} \\
 T(x) &= [x; \text{vec}(xx^T)] \\
 A(\eta) &= \frac{1}{2} \mu^T \Sigma^{-1} \mu + \log|\Sigma| = -\frac{1}{2} \text{tr}(\eta_2 \eta_1 \eta_1^T) - \frac{1}{2} \log(-2\eta_2) \\
 h(x) &= (2\pi)^{-k/2}
 \end{aligned}$$

Natural parameter

- Note: a k -dimensional Gaussian is a $(d+d)$ -parameter distribution with a $(d+d)$ -element vector of sufficient statistics (but because of symmetry and positivity, parameters are constrained and have lower degree of freedom)



Example: Multinomial distribution

- For a binary vector random variable $X \sim \text{multi}(x | \pi)$,

$$\begin{aligned} p(x|\pi) &= \pi_1^{x_1} \pi_2^{x_2} \cdots \pi_K^{x_K} = \exp \left\{ \sum_k x_k \ln \pi_k \right\} \\ &= \exp \left\{ \sum_{k=1}^{K-1} x_k \ln \pi_k + \left(1 - \sum_{k=1}^{K-1} x_k \right) \ln \left(1 - \sum_{k=1}^{K-1} \pi_k \right) \right\} \\ &= \exp \left\{ \sum_{k=1}^{K-1} x_k \ln \left(\frac{\pi_k}{1 - \sum_{k=1}^{K-1} \pi_k} \right) + \ln \left(1 - \sum_{k=1}^{K-1} \pi_k \right) \right\} \end{aligned}$$

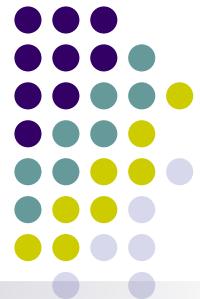
- Exponential family representation

$$\eta = \left[\ln \left(\frac{\pi_k}{\pi_K} \right); 0 \right]$$

$$T(x) = [x]$$

$$A(\eta) = -\ln \left(1 - \sum_{k=1}^{K-1} \pi_k \right) = \ln \left(\sum_{k=1}^K e^{\eta_k} \right)$$

$$h(x) = 1$$



Examples

- Gaussian:

$$\begin{aligned}\eta &= \left[\Sigma^{-1} \mu; -\frac{1}{2} \text{vec}(\Sigma^{-1}) \right] \\ T(x) &= \left[x; \text{vec}(xx^T) \right] \\ A(\eta) &= \frac{1}{2} \mu^T \Sigma^{-1} \mu + \frac{1}{2} \log |\Sigma| \\ h(x) &= (2\pi)^{-k/2}\end{aligned}$$

$$\Rightarrow \mu_{MLE} = \frac{1}{N} \sum_n T_1(x_n) = \frac{1}{N} \sum_n x_n$$

- Multinomial:

$$\begin{aligned}\eta &= \left[\ln\left(\frac{\pi_k}{\pi_K}\right); 0 \right] \\ T(x) &= [x] \\ A(\eta) &= -\ln\left(1 - \sum_{k=1}^{K-1} \pi_k\right) = \ln\left(\sum_{k=1}^K e^{\eta_k}\right) \\ h(x) &= 1\end{aligned}$$

$$\Rightarrow \mu_{MLE} = \frac{1}{N} \sum_n x_n$$

- Poisson:

$$\begin{aligned}\eta &= \log \lambda \\ T(x) &= x \\ A(\eta) &= \lambda = e^\eta \\ h(x) &= \frac{1}{x!}\end{aligned}$$

$$\Rightarrow \mu_{MLE} = \frac{1}{N} \sum_n x_n$$

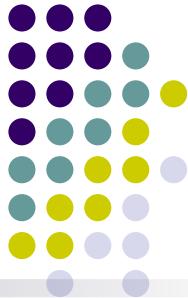


Why exponential family?

- Moment generating property

$$\begin{aligned}\frac{dA}{d\eta} &= \frac{d}{d\eta} \log Z(\eta) = \frac{1}{Z(\eta)} \frac{d}{d\eta} Z(\eta) \\ &= \frac{1}{Z(\eta)} \frac{d}{d\eta} \int h(x) \exp\{\eta^T T(x)\} dx \\ &= \int T(x) \frac{h(x) \exp\{\eta^T T(x)\}}{Z(\eta)} dx \\ &= E[T(x)]\end{aligned}$$

$$\begin{aligned}\frac{d^2 A}{d\eta^2} &= \int T^2(x) \frac{h(x) \exp\{\eta^T T(x)\}}{Z(\eta)} dx - \int T(x) \frac{h(x) \exp\{\eta^T T(x)\}}{Z(\eta)} dx \frac{1}{Z(\eta)} \frac{d}{d\eta} Z(\eta) \\ &= E[T^2(x)] - E^2[T(x)] \\ &= Var[T(x)]\end{aligned}$$



Moment estimation

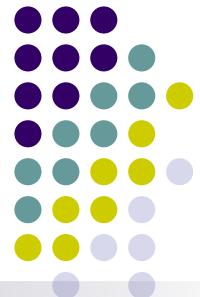
- We can easily compute moments of any exponential family distribution by taking the derivatives of the log normalizer $A(\eta)$.
- The q^{th} derivative gives the q^{th} centered moment.

$$\frac{dA(\eta)}{d\eta} = \text{mean}$$

$$\frac{d^2 A(\eta)}{d\eta^2} = \text{variance}$$

...

- When the sufficient statistic is a stacked vector, partial derivatives need to be considered.



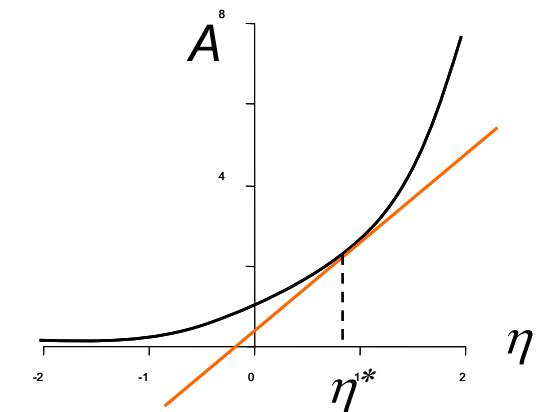
Moment vs canonical parameters

- The moment parameter μ can be derived from the natural (canonical) parameter

$$\frac{dA(\eta)}{d\eta} = E[T(x)] \stackrel{\text{def}}{=} \mu$$

- $A(\eta)$ is convex since

$$\frac{d^2 A(\eta)}{d\eta^2} = \text{Var}[T(x)] > 0$$



- Hence we can invert the relationship and infer the canonical parameter from the moment parameter (1-to-1):

$$\eta \stackrel{\text{def}}{=} \psi(\mu)$$

- A distribution in the exponential family can be parameterized not only by η – the canonical parameterization, but also by μ – the moment parameterization.



MLE for Exponential Family

- For *iid* data, the log-likelihood is

$$\begin{aligned}\ell(\eta; D) &= \log \prod_n h(x_n) \exp\{\eta^T T(x_n) - A(\eta)\} \\ &= \sum_n \log h(x_n) + \left(\eta^T \sum_n T(x_n) \right) - NA(\eta)\end{aligned}$$

- Take derivatives and set to zero:

$$\frac{\partial \ell}{\partial \eta} = \sum_n T(x_n) - N \frac{\partial A(\eta)}{\partial \eta} = 0$$

$$\begin{aligned}\Rightarrow \frac{\partial A(\eta)}{\partial \eta} &= \frac{1}{N} \sum_n T(x_n) \\ \hat{\mu}_{MLE} &= \frac{1}{N} \sum_n T(x_n)\end{aligned}$$

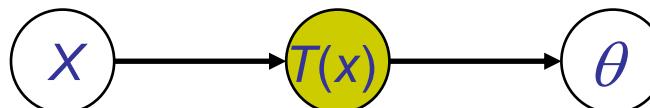
- This amounts to **moment matching**.
- We can infer the canonical parameters using $\hat{\eta}_{MLE} = \psi(\hat{\mu}_{MLE})$



Sufficiency

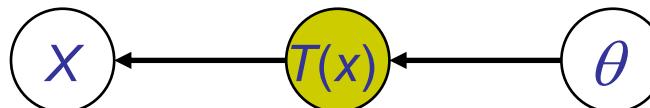
- For $p(x|\theta)$, $T(x)$ is *sufficient* for θ if there is no information in X regarding θ beyond that in $T(x)$.
 - We can throw away X for the purpose of inference w.r.t. θ .

- Bayesian view



$$p(\theta | T(x), x) = p(\theta | T(x))$$

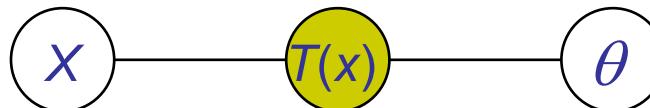
- Frequentist view



$$p(x | T(x), \theta) = p(x | T(x))$$

- The Neyman factorization theorem

- $T(x)$ is *sufficient* for θ if



$$p(x, T(x), \theta) = \psi_1(T(x), \theta) \psi_2(x, T(x))$$

$$\Rightarrow p(x | \theta) = g(T(x), \theta) h(x, T(x))$$

EXPONENTIAL FAMILY AND MEAN FIELD VARIATIONAL INF.

CAVI Algorithm

Coordinate Ascent Variational Inference (CAVI)

- here we assume a **mean field** approximation
- application of **coordinate ascent** to maximization of ELBO
- converges to a **local optimum** of the **nonconvex** ELBO objective

```
1: procedure CAVI( $p_\alpha$ )
2:   Let  $q_\theta(\mathbf{z}) = \prod_{t=1}^T q_t(z_t)$                                  $\triangleright$  Mean field approx.
3:   while ELBO( $q_\theta$ ) has not converged do
4:     for  $t \in \{1, \dots, T\}$  do                                      $\triangleright$  For each variable
5:       Set  $q_t(z_t) \propto \exp(E_{q_{\neg t}}[\log p_\alpha(z_t \mid z_{\neg t}, \mathbf{x})])$ 
6:       while keeping all  $\{q_s(\cdot)\}_{s \neq t}$  fixed
7:     Compute  $\text{ELBO}(q_\theta) = E_{q_\theta(\mathbf{z})} [\log p_\alpha(\mathbf{x}, \mathbf{z})] - E_{q_\theta(\mathbf{z})} [\log q_\theta(\mathbf{z})]$ 
8:   return  $q_\theta$ 
```

Optimizing the ELBO in Mean Field Variational Inference



Notes:

- This coordinate ascent procedure converges to a ***local maximum***.
- The coordinate ascent update for $q(z_j)$ only depends on the other, fixed approximations $q(z_k)$, $k \neq j$.
- While this determines the optimal $q(z_j)$, we haven't yet specified the form (i.e. what specific distribution family) of q we aim to use, only the factorization.
- Depending on what form we use, the coordinate update $q^*(z_j)$ might not be easy to work with (and might not be in the same form as $q(z_j)$...).
 - But in many cases it is!
 - And we will specify what forms yield good coordinate updates.



Optimizing the ELBO in Mean Field Variational Inference

Simple Example: multinomial conditionals

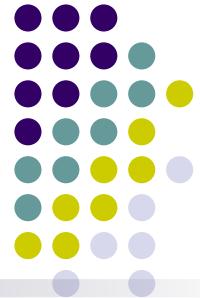
- Suppose we have chosen a model whose conditional distribution is a multinomial, i.e.

$$p(z_j|z_{-j}, x) = \pi(z_j, x)$$

- Then the optimal (coordinate update for) $q(z_j)$ is:

$$q^*(z_j) \propto \exp \{ \mathbb{E} [\log \pi(z_j, x)] \}$$

- Which is also a multinomial, and is easy to compute. So choosing a multinomial family of approximations for each latent variable gives closed form coordinate ascent updates.



Quick Recap

Quick recap on what we've covered:

- We defined a family of approximations called “mean field” approximations, in which there are no dependencies between latent variables (and also a generalized version of this).
- We decomposed the ELBO into a nice form under mean field assumptions.
- We derived coordinate ascent updates to iteratively optimize each local variational approximation under mean field assumptions.
- Next, we will discuss specific forms for the local variational approximations in which we can easily compute (closed-form) coordinate ascent updates.



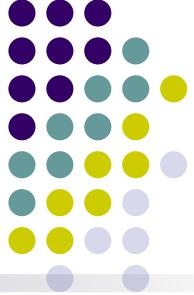
Exponential Family Conditionals

- Is there a general form for models in which the coordinate updates in mean field variational inference are easy to compute and lead to closed-form updates?
- Yes: the answer is exponential family conditionals.
- I.e. models with conditional densities that are in an exponential family, i.e. of the form:

$$p(z_j|z_{-j}, x) = h(z_j) \exp \left\{ \eta(z_{-j}, x)^T t(z_j) - a(\eta(z_{-j}, x)) \right\}$$

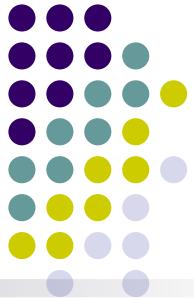
where h , η , t , and a are functions that parameterize the exponential family.

- Different choices of these parameters lead to many popular densities (normal, gamma, exponential, Bernouilli, Dirichlet, categorical, beta, Poisson, geometric, etc.).



Exponential Family Conditionals

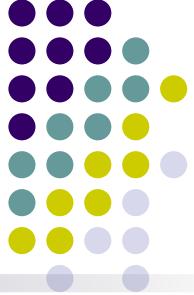
- We call these “exponential-family-conditional” models.
 - Also known as “conditionally conjugate models”.
- Many popular models fall into this category, including:
 - Bayesian mixtures of exponential family models with conjugate priors.
 - Hierarchical hidden Markov models.
 - Kalman filter models and switching Kalman filters.
 - Mixed-membership models of exponential families.
 - Factorial mixtures / hidden Markov models of exponential families.
 - Bayesian linear regression.
 - Any model containing only conjugate pairs and multinomials.
- Some popular models do not fall into this category, including:
 - Bayesian logistic regression and other nonconjugate Bayesian generalized linear models.
 - Correlated topic model, dynamic topic model.
 - Discrete choice models.
 - Nonlinear matrix factorization models.



Exponential Family Conditionals

- We can derive a general formula for the coordinate ascent update for all exponential-family-conditional models.
- First, we will choose the form of our local variational approximation $q(z_j)$ to be the same as the conditional distribution (i.e. in an exponential family).
- When we perform our coordinate ascent update, we will see that the update yields an optimal $q(z_j)$ in the same family.
- Recall from above that we derived the coordinate ascent updates for optimizing the ELBO (under the mean field assumption) as:

$$q^*(z_j) \propto \exp \left\{ \mathbb{E}_{q_{-j}} [\log p(z_j | z_{-j}, x)] \right\}$$



Exponential Family Conditionals

Coordinate ascent updates for exponential-family-conditional models (under the mean field approximation):

- The log of the conditional:

$$\log p(z_j|z_{-j}, x) = \log h(z_j) + \eta(z_{-j}, x)^\top t(z_j) - a(\eta(z_{-j}, x))$$

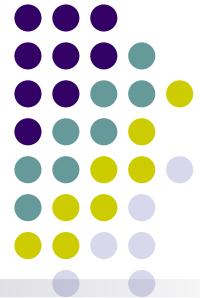
- The expectation of this with respect to $q(z_{-j})$ is:

$$\mathbb{E}_{q_{-j}} [\log p(z_j|z_{-j}, x)] = \log h(z_j) + \mathbb{E}_{q_{-j}} [\eta(z_{-j}, x)]^\top t(z_j) - \mathbb{E}_{q_{-j}} [a(\eta(z_{-j}, x))]$$

- The last term does not depend on $q(z_j)$, so we have the update:

$$q^*(z_j) \propto h(z_j) \exp \left\{ \mathbb{E}_{q_{-j}} [\eta(z_{-j}, x)]^\top t(z_j) \right\}$$

- So the optimal $q(z_j)$ is in the same exponential family as the conditional.



Exponential Family Conditionals

Writing this update in terms of variational parameters ν .

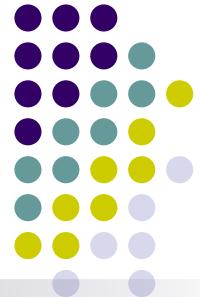
- Give each latent variable a variational parameter ν_j . Under the mean field assumption, we can write the full approximation as :

$$q(z_{1:m}|\nu) = \prod_{j=1}^m q(z_j|\nu_j)$$

where each local variational approximation has an exponential family form.

- Then the coordinate ascent algorithm updates each variational parameter, in turn, as:

$$\nu_j^* = \mathbb{E}_{q_{-j}} [\eta(z_{-j}, x)]$$



Quick Recap

Quick recap on what we've covered:

- We found a family of models (exponential-family-conditional models) in which we have closed form coordinate ascent updates to optimize the ELBO.
 - And we gave a number of examples (and non-examples) of these models.
- We gave an explicit form for the coordinate ascent update for these exponential-family-conditional models.
 - And also looked at the update in terms of the local variational parameters.

MAP INFERENCE AND VARIATIONAL INFERENCE

MAP Inference as Variational Inference

Suppose: We want a family \mathcal{Q} such that the variational inference solution:

$$\hat{q}(\mathbf{z}) = \operatorname{argmin}_{q \in \mathcal{Q}} \text{KL}(q \parallel p)$$

gives back a distribution that is a point mass on the MAP inference solution:

$$\hat{q}(\mathbf{z}) = \begin{cases} \hat{\mathbf{z}} = \operatorname{argmax}_{\mathbf{z} \in \mathcal{Z}(\mathbf{x})} p(\mathbf{z} \mid \mathbf{x}) & \text{w/prob. 1.0} \\ \text{any other } \mathbf{z} \in \mathcal{Z}(\mathbf{x}) & \text{w/prob. 0.0} \end{cases}$$

Question: What is \mathcal{Q} ?

Answer:

VARIATIONAL INFERENCE RESULTS

Variational Inference & Nonconvexity

- ELBO is a non-convex objective function
- Below shows 10 random initializations of CAVI for Gaussian Mixture Model
- Parameters with higher ELBO are closer to true posterior

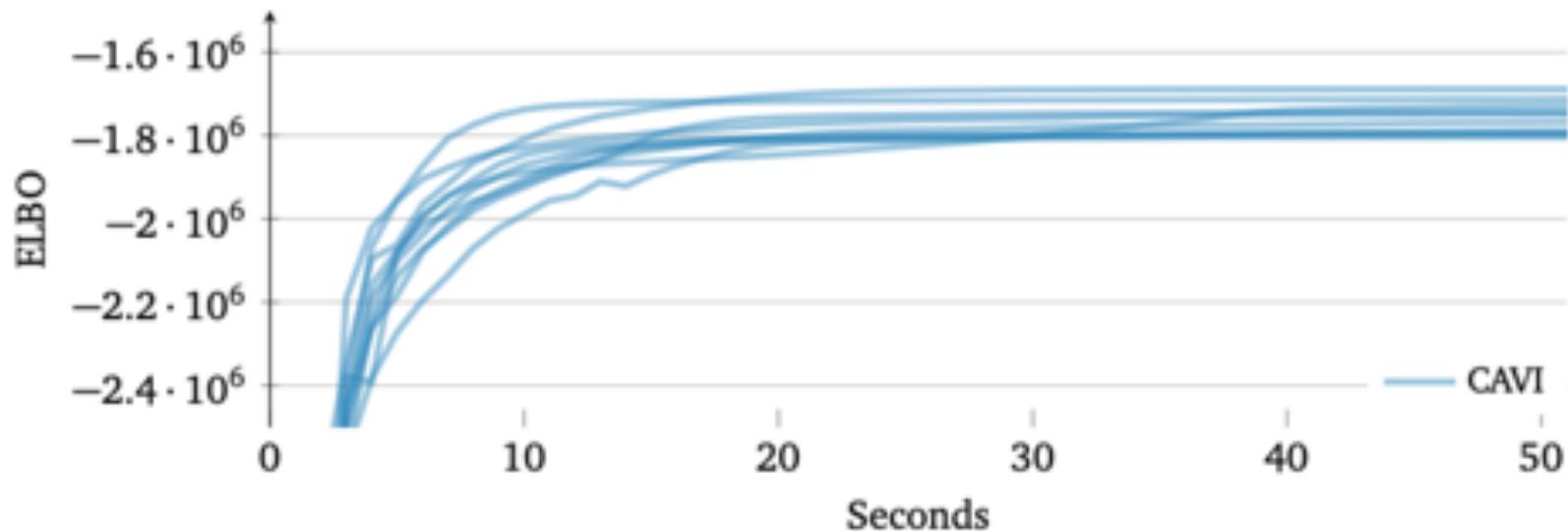
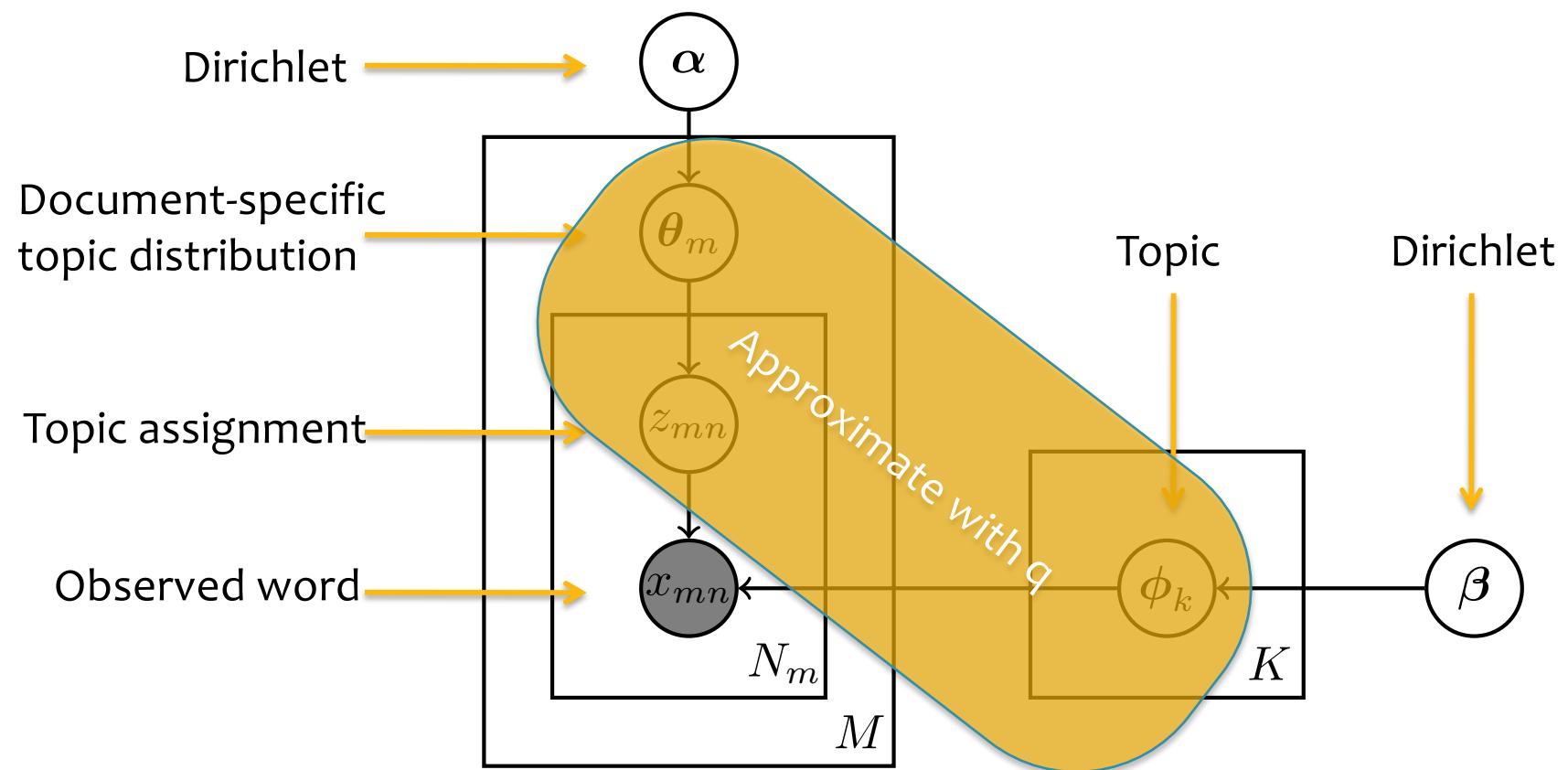


Figure 2: Different initializations may lead CAVI to find different local optima of the ELBO.

Variational Bayesian LDA

- Explicit Variational Inference



Variational Bayesian LDA

- Explicit Variational Inference

Standard VB inference upper bounds the negative log marginal likelihood $-\log p(\mathbf{x}|\alpha, \beta)$ using the variational free energy:

$$-\log p(\mathbf{x}|\alpha, \beta) \leq \tilde{\mathcal{F}}(\bar{q}(\mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\phi})) = E_{\bar{q}}[-\log p(\mathbf{x}, \mathbf{z}, \boldsymbol{\phi}, \boldsymbol{\theta}|\alpha, \beta)] - \mathcal{H}(\bar{q}(\mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\phi})) \quad (2)$$

with $\bar{q}(\mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\phi})$ an approximate posterior, $\mathcal{H}(\bar{q}(\mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\phi})) = E_{\bar{q}}[-\log \bar{q}(\mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\phi})]$ the variational entropy, and $\bar{q}(\mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\phi})$ assumed to be fully factorized:

$$\bar{q}(\mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\phi}) = \prod_{ij} \bar{q}(z_{ij}|\tilde{\gamma}_{ij}) \prod_j \bar{q}(\theta_j|\tilde{\alpha}_j) \prod_k \bar{q}(\phi_k|\tilde{\beta}_k) \quad (3)$$

$\bar{q}(z_{ij}|\tilde{\gamma}_{ij})$ is multinomial with parameters $\tilde{\gamma}_{ij}$ and $\bar{q}(\theta_j|\tilde{\alpha}_j)$, $\bar{q}(\phi_k|\tilde{\beta}_k)$ are Dirichlet with parameters $\tilde{\alpha}_j$ and $\tilde{\beta}_k$ respectively. Optimizing $\tilde{\mathcal{F}}(\bar{q})$ with respect to the variational parameters gives us a set of updates guaranteed to improve $\tilde{\mathcal{F}}(\bar{q})$ at each iteration and converges to a local minimum:

$$\tilde{\alpha}_{jk} = \alpha + \sum_i \tilde{\gamma}_{ijk} \quad (4)$$

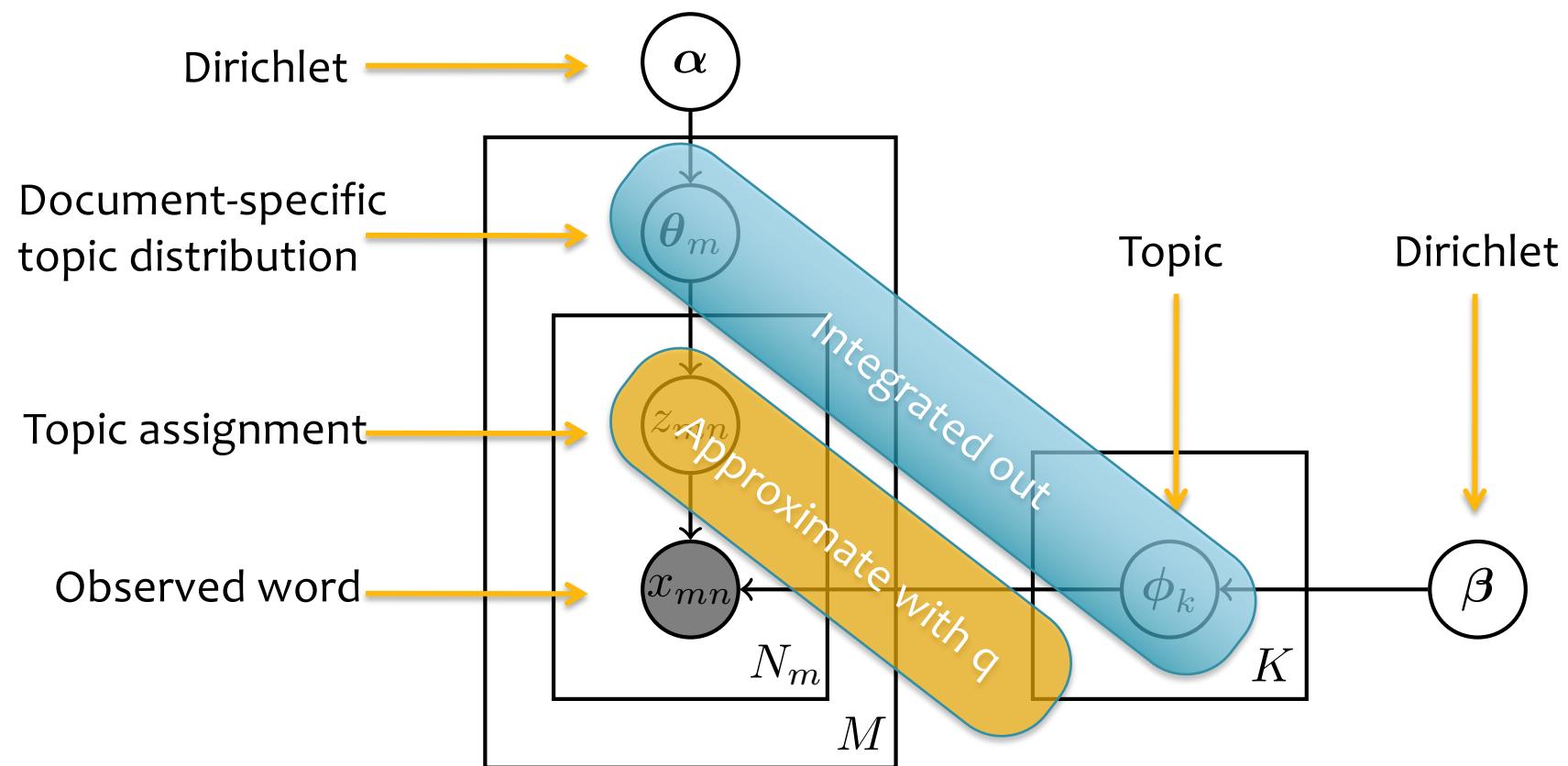
$$\tilde{\beta}_{kw} = \beta + \sum_{ij} \mathbf{1}(x_{ij}=w) \tilde{\gamma}_{ijk} \quad (5)$$

$$\tilde{\gamma}_{ijk} \propto \exp \left(\Psi(\tilde{\alpha}_{jk}) + \Psi(\tilde{\beta}_{kx_{ij}}) - \Psi(\sum_w \tilde{\beta}_{kw}) \right) \quad (6)$$

where $\Psi(y) = \frac{\partial \log \Gamma(y)}{\partial y}$ is the digamma function and $\mathbf{1}$ is the indicator function.

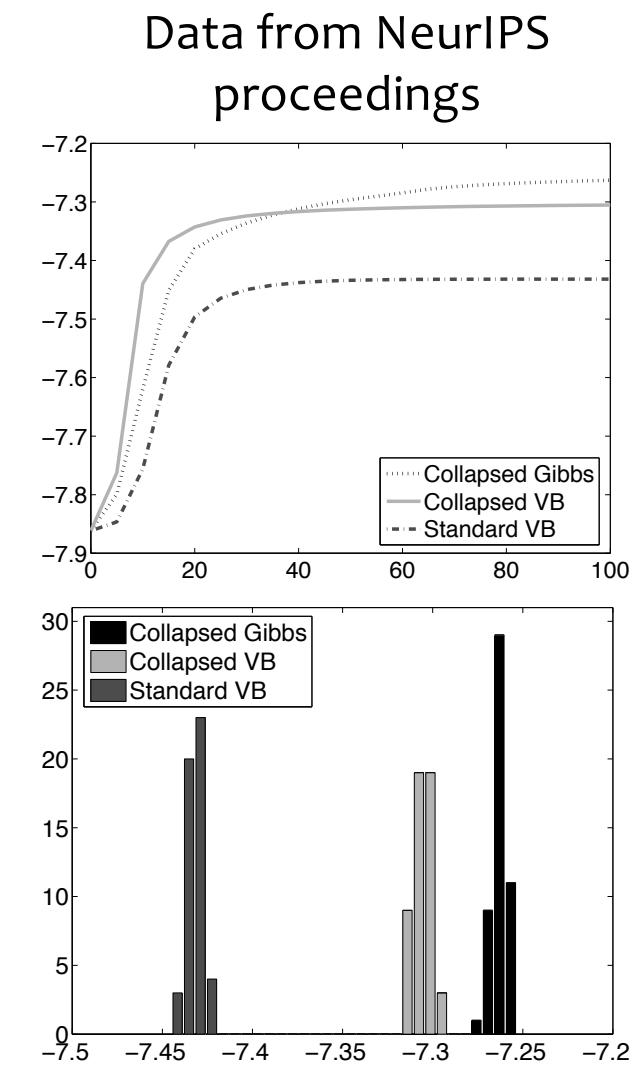
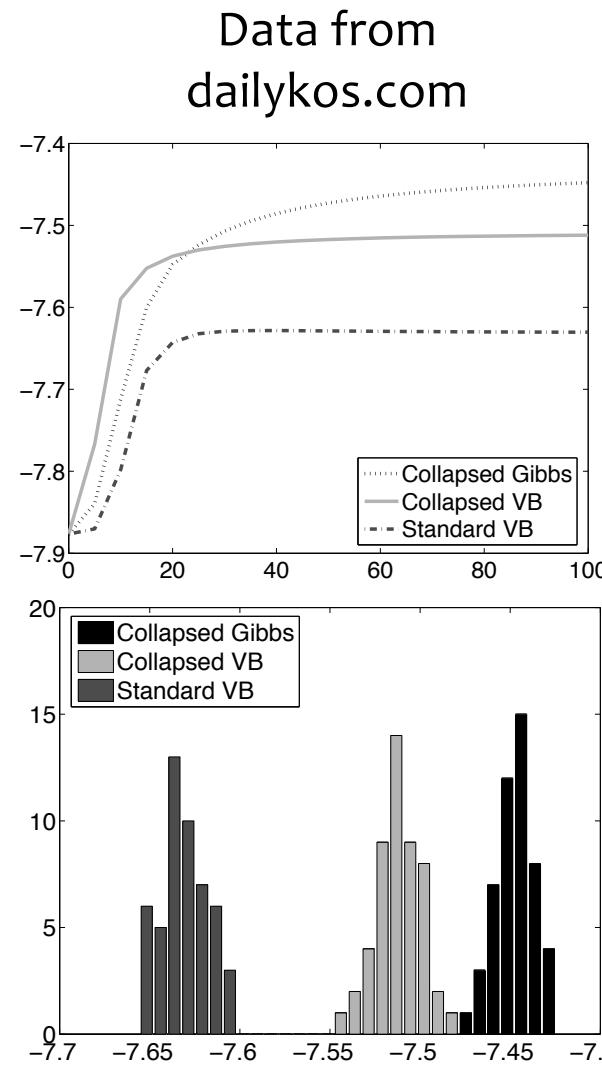
Collapsed Variational Bayesian LDA

- Collapsed Variational Inference

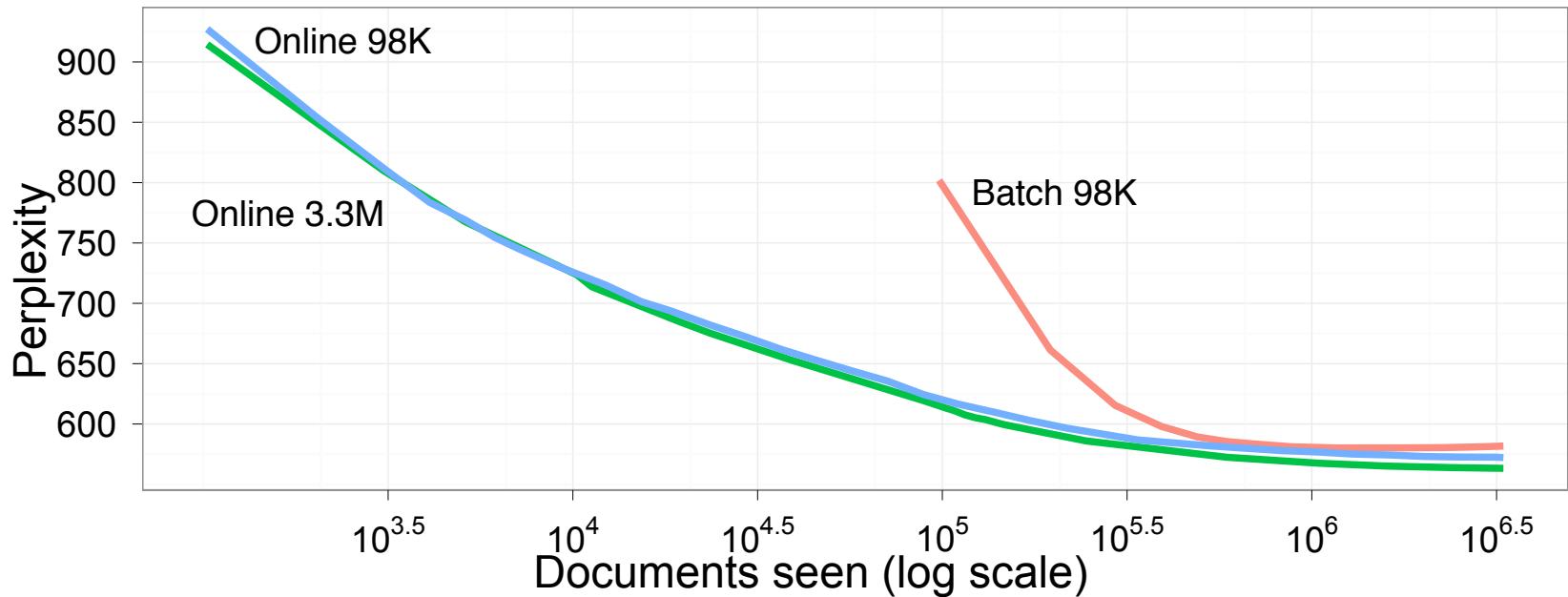


Collapsed Variational Bayesian LDA

- **First row:** test set per word log probabilities as functions of numbers of iterations for VB, CVB and Gibbs.
- **Second row:** histograms of final test set per word log probabilities across 50 random initializations.



Online Variational Bayes for LDA



Documents analyzed	2048	4096	8192	12288	16384	32768	49152	65536
Top eight words	systems road made service announced national west language	systems health communication service billion language care road	service systems health companies market communication company billion	service systems companies business company billion health industry	service companies systems business company industry market billion	business service companies industry company management systems services	business service companies industry services company management public	business industry service companies services company management public

Online Variational Bayes for LDA

Algorithm 1 Batch variational Bayes for LDA

```

Initialize  $\lambda$  randomly.
while relative improvement in  $\mathcal{L}(w, \phi, \gamma, \lambda) > 0.00001$  do
    E step:
    for  $d = 1$  to  $D$  do
        Initialize  $\gamma_{dk} = 1$ . (The constant 1 is arbitrary.)
        repeat
            Set  $\phi_{dwk} \propto \exp\{\mathbb{E}_q[\log \theta_{dk}] + \mathbb{E}_q[\log \beta_{kw}]\}$ 
            Set  $\gamma_{dk} = \alpha + \sum_w \phi_{dwk} n_{dw}$ 
        until  $\frac{1}{K} \sum_k |\text{change in } \gamma_{dk}| < 0.00001$ 
    end for
    M step:
    Set  $\lambda_{kw} = \eta + \sum_d n_{dw} \phi_{dwk}$ 
end while

```

Algorithm 2 Online variational Bayes for LDA

```

Define  $\rho_t \triangleq (\tau_0 + t)^{-\kappa}$ 
Initialize  $\lambda$  randomly.
for  $t = 0$  to  $\infty$  do
    E step:
    Initialize  $\gamma_{tk} = 1$ . (The constant 1 is arbitrary.)
    repeat
        Set  $\phi_{twk} \propto \exp\{\mathbb{E}_q[\log \theta_{tk}] + \mathbb{E}_q[\log \beta_{kw}]\}$ 
        Set  $\gamma_{tk} = \alpha + \sum_w \phi_{twk} n_{tw}$ 
    until  $\frac{1}{K} \sum_k |\text{change in } \gamma_{tk}| < 0.00001$ 
    M step:
    Compute  $\tilde{\lambda}_{kw} = \eta + D n_{tw} \phi_{twk}$ 
    Set  $\lambda = (1 - \rho_t) \lambda + \rho_t \tilde{\lambda}$ .
end for

```

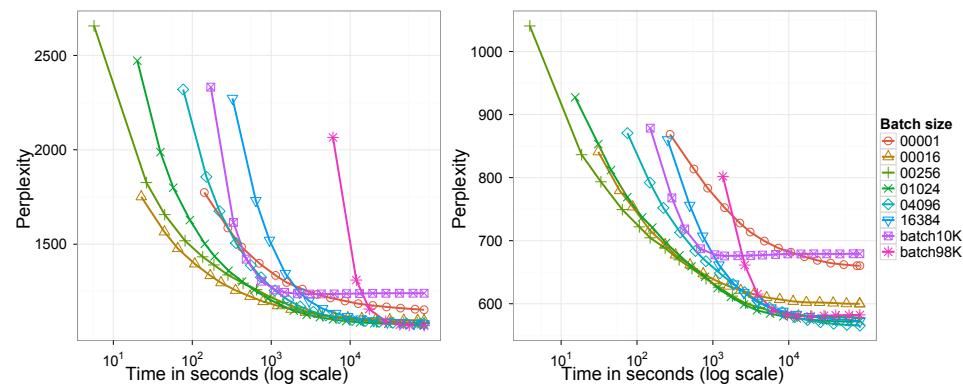


Figure 2: Held-out perplexity obtained on the *Nature* (left) and *Wikipedia* (right) corpora as a function of CPU time. For moderately large mini-batch sizes, online LDA finds solutions as good as those that the batch LDA finds, but with much less computation. When fit to a 10,000-document subset of the training corpus batch LDA's speed improves, but its performance suffers.

Fully-Connected CRF

Model

$$p(\mathbf{x}|\mathbf{i}) = \frac{1}{Z(\mathbf{i})} \exp(-E(\mathbf{x}))$$

$$E(\mathbf{x}) = \sum_i \psi_u(x_i) + \sum_{i < j} \psi_p(x_i, x_j),$$



This is a fully connected graph!

Inference

- Can do MCMC, but slow
- Instead use Variational Inference
- Then filter some variables for speed up

Follow-up Work (combine with CNN)

Published as a conference paper at ICLR 2015

SEMANTIC IMAGE SEGMENTATION WITH DEEP CONVOLUTIONAL NETS AND FULLY CONNECTED CRFs

Liang-Chieh Chen

Univ. of California, Los Angeles
lcchen@cs.ucla.edu

George Papandreou *

Google Inc.
gpapan@google.com

Iasonas Kokkinos

CentraleSupélec and INRIA
iasonas.kokkinos@ecp.fr

Kevin Murphy

Google Inc.
kpmurphy@google.com

Alan L. Yuille

Univ. of California, Los Angeles
yuille@stat.ucla.edu

ABSTRACT

Deep Convolutional Neural Networks (DCNNs) have recently shown state of the art performance in high level vision tasks, such as image classification and object detection. This work brings together methods from DCNNs and probabilistic graphical models for addressing the task of pixel-level classification (also called "semantic image segmentation"). We show that responses at the final layer of DCNNs are not sufficiently localized for accurate object segmentation. This is due to the very invariance properties that make DCNNs good for high level tasks. We overcome this poor localization property of deep networks by combining the responses at the final DCNN layer with a fully connected Conditional Random Field (CRF). Qualitatively, our "DeepLab" system is able to localize segment boundaries at a level of accuracy which is beyond previous methods. Quantitatively, our method sets the new state-of-art at the PASCAL VOC-2012 semantic image segmentation task, reaching 71.6% IOU accuracy in the test set. We show how these results can be obtained efficiently: Careful network re-purposing and a novel application of the 'hole' algorithm from the wavelet community allow dense computation of neural net responses at 8 frames per second on a modern GPU.

Fully-Connected CRF

Model

$$p(\mathbf{x}|\mathbf{i}) = \frac{1}{Z(\mathbf{i})} \exp(-E(\mathbf{x}))$$

$$E(\mathbf{x}) = \sum_i \psi_u(x_i) + \sum_{i < j} \psi_p(x_i, x_j),$$

This is a fully connected graph!

Results

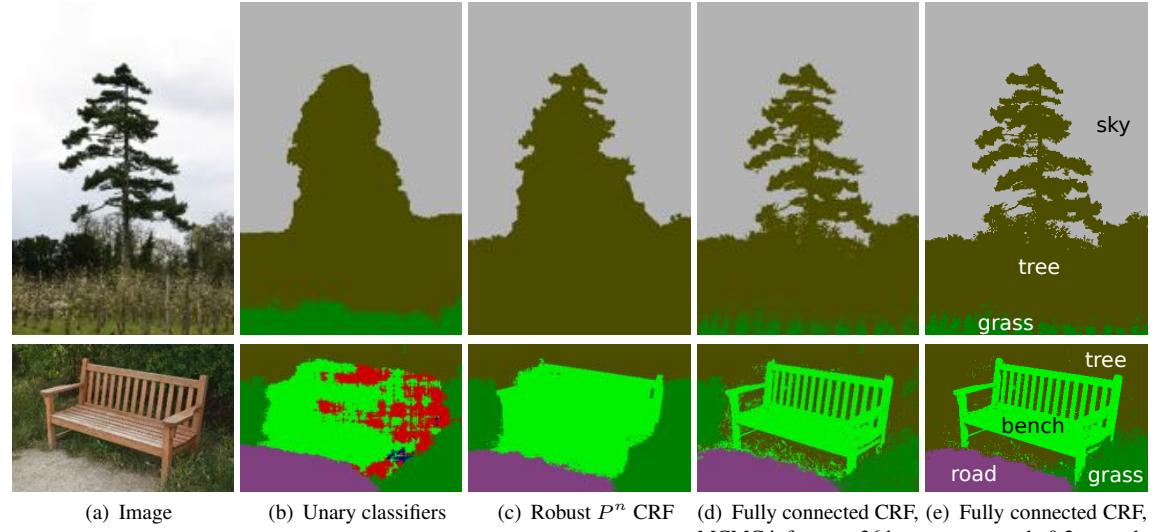


Figure 1: Pixel-level classification with a fully connected CRF. (a) Input image from the MSRC-21 dataset. (b) The response of unary classifiers used by our models. (c) Classification produced by the Robust P^n CRF [9]. (d) Classification produced by MCMC inference [17] in a fully connected pixel-level CRF model; the algorithm was run for 36 hours and only partially converged for the bottom image. (e) Classification produced by our inference algorithm in the fully connected model in 0.2 seconds.

Inference

- Can do MCMC, but slow
- Instead use Variational Inference
- Then filter some variables for speed up

Figures from Krähenbühl & Koltun (2011)

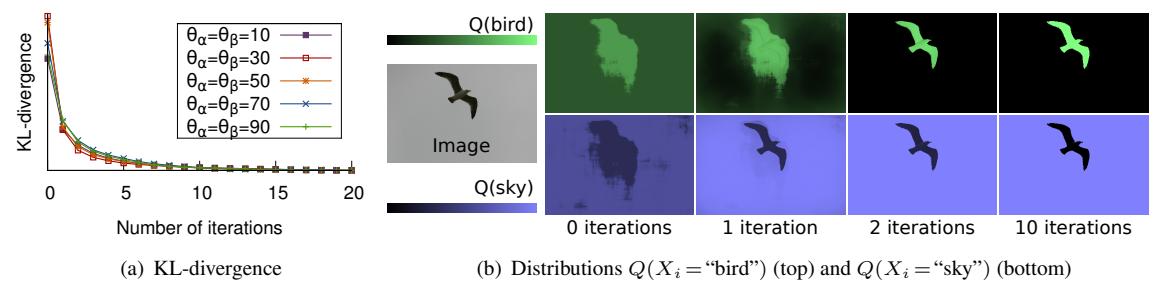


Figure 2: Convergence analysis. (a) KL-divergence of the mean field approximation during successive iterations of the inference algorithm, averaged across 94 images from the MSRC-21 dataset. (b) Visualization of convergence on distributions for two class labels over an image from the dataset.

Joint Parsing and Alignment with Weakly Synchronized Grammars

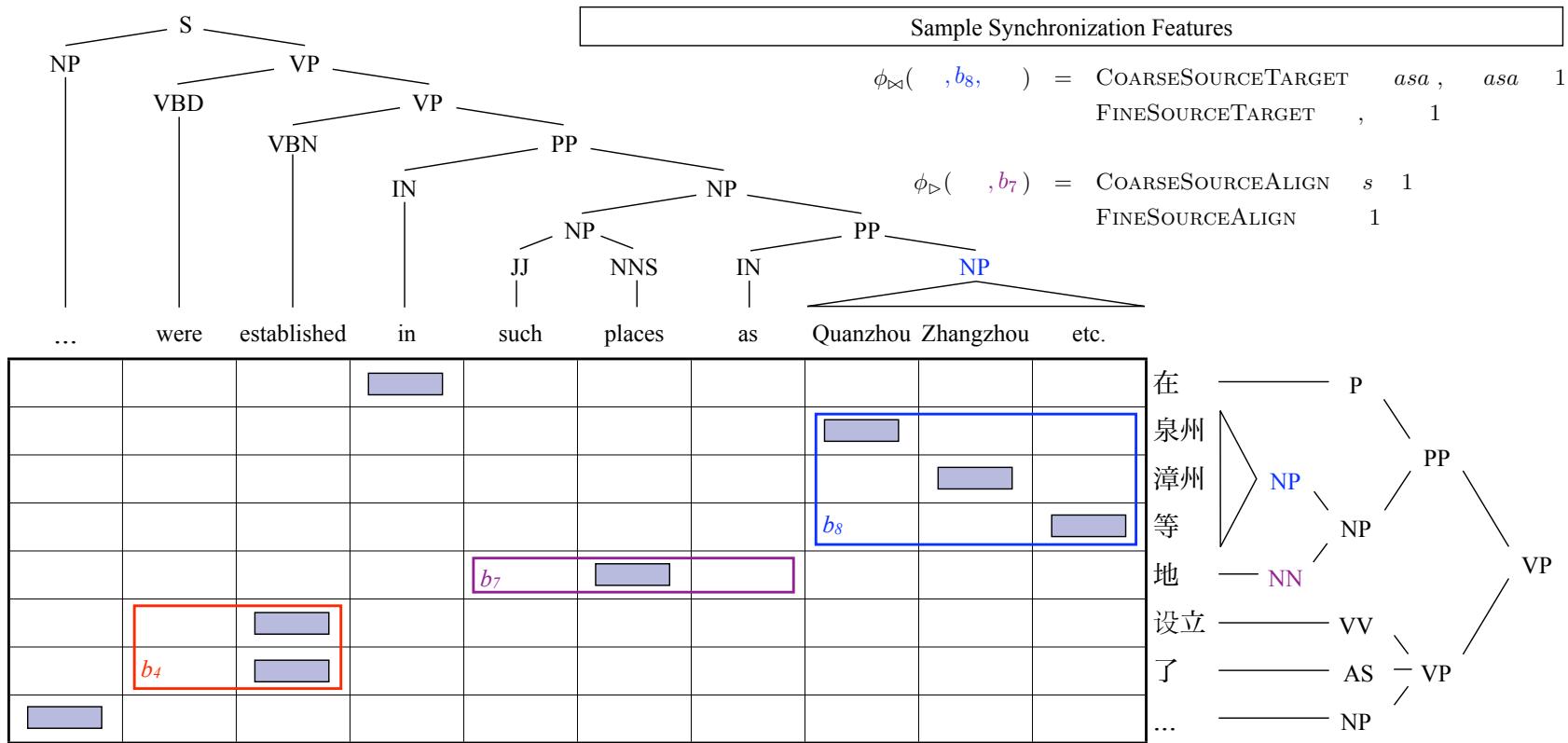
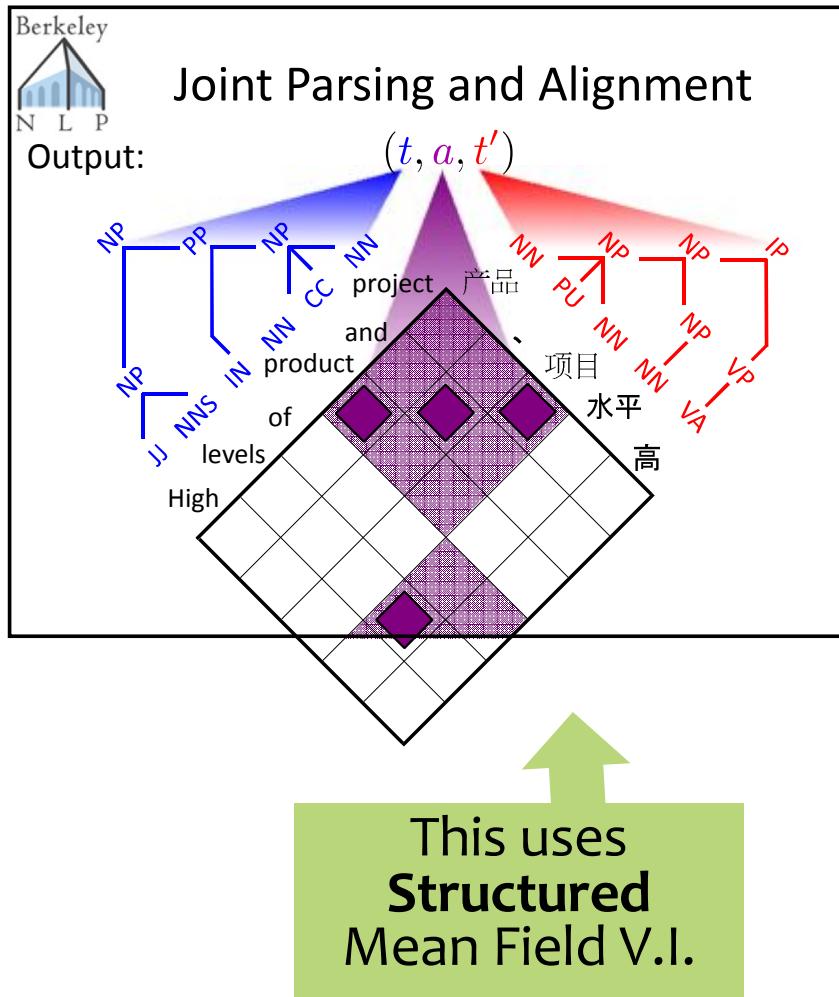


Figure 2: An example of a Chinese-English sentence pair with parses, word alignments, and a subset of the full optimal ITG derivation, including one totally unsynchronized bispans (b_4), one partially synchronized bispans (b_7), and one fully synchronized bispans (b_8). The inset provides some examples of active synchronization features (see Section 4.3) on these bispans. On this example, the monolingual English parser erroneously attached the lower PP to the VP headed by *established*, and the non-syntactic ITG word aligner misaligned 等 to *such* instead of to *etc*. Our joint model corrected both of these mistakes because it was rewarded for the synchronization of the two NPs joined by b_8 .

Joint Parsing and Alignment with Weakly Synchronized Grammars

Figures from Burkett & Klein (ACL 2013 tutorial)



	Test Results		
	Ch F ₁	Eng F ₁	Tot F ₁
Monolingual	83.6	81.2	82.5
Reranker	86.0	83.8	84.9
Joint	85.7	84.5	85.1

Table 1: Parsing results. Our joint model has the highest reported F₁ for English-Chinese bilingual parsing.

	Test Results			
	Precision	Recall	AER	F ₁
HMM	86.0	58.4	30.0	69.5
ITG	86.8	73.4	20.2	79.5
Joint	85.5	84.6	14.9	85.0

Table 2: Word alignment results. Our joint model has the highest reported F₁ for English-Chinese word alignment.

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

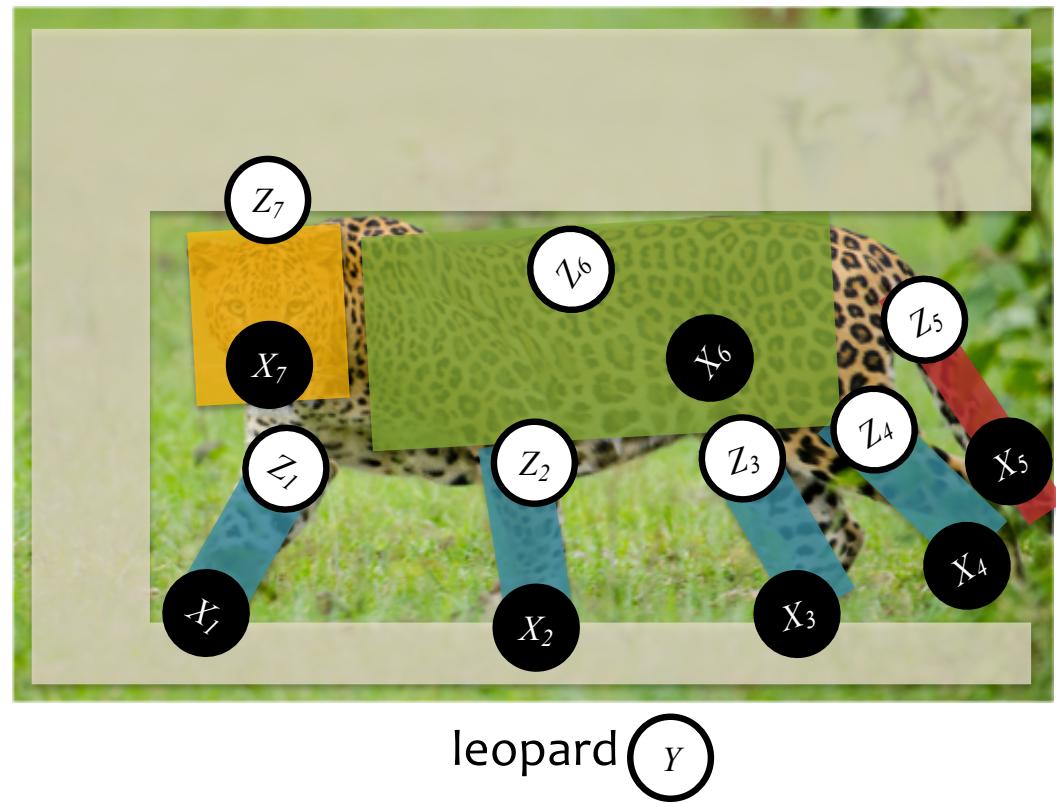


leopard

Case Study: Object Recognition

Data consists of images x and labels y .

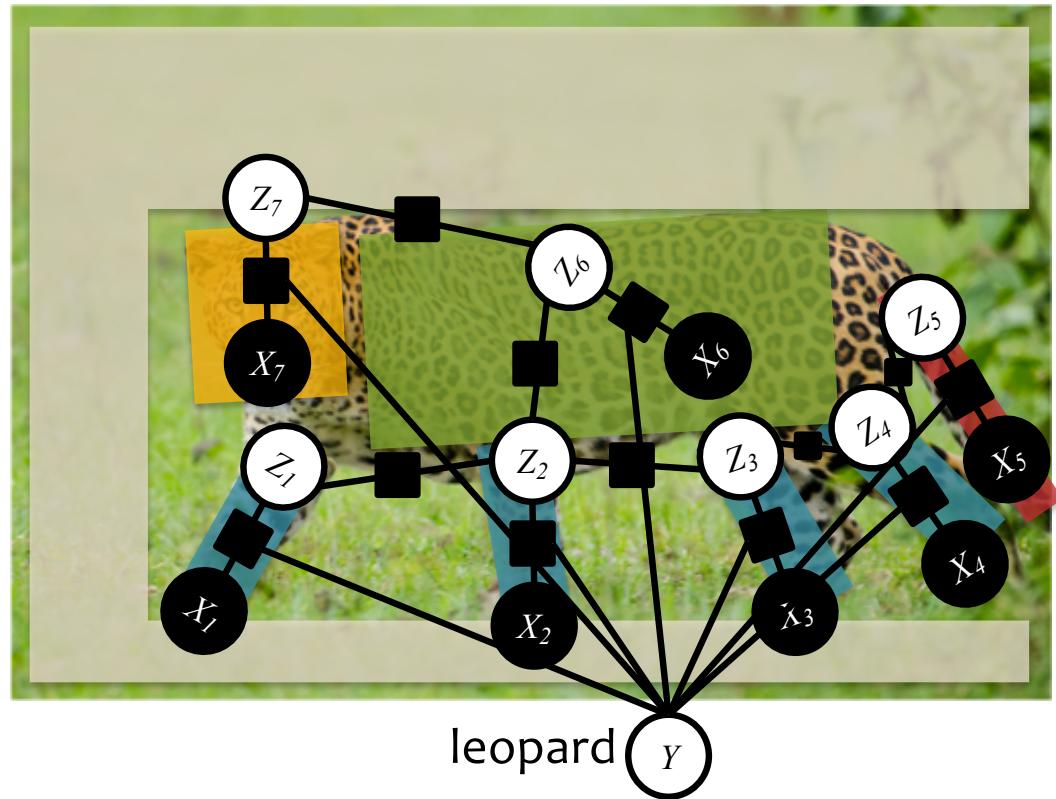
- Preprocess data into “patches”
- Posit a latent labeling z describing the object’s parts (e.g. head, leg, tail, torso, grass)
- Define graphical model with these latent variables in mind
- z is not observed at train or test time



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

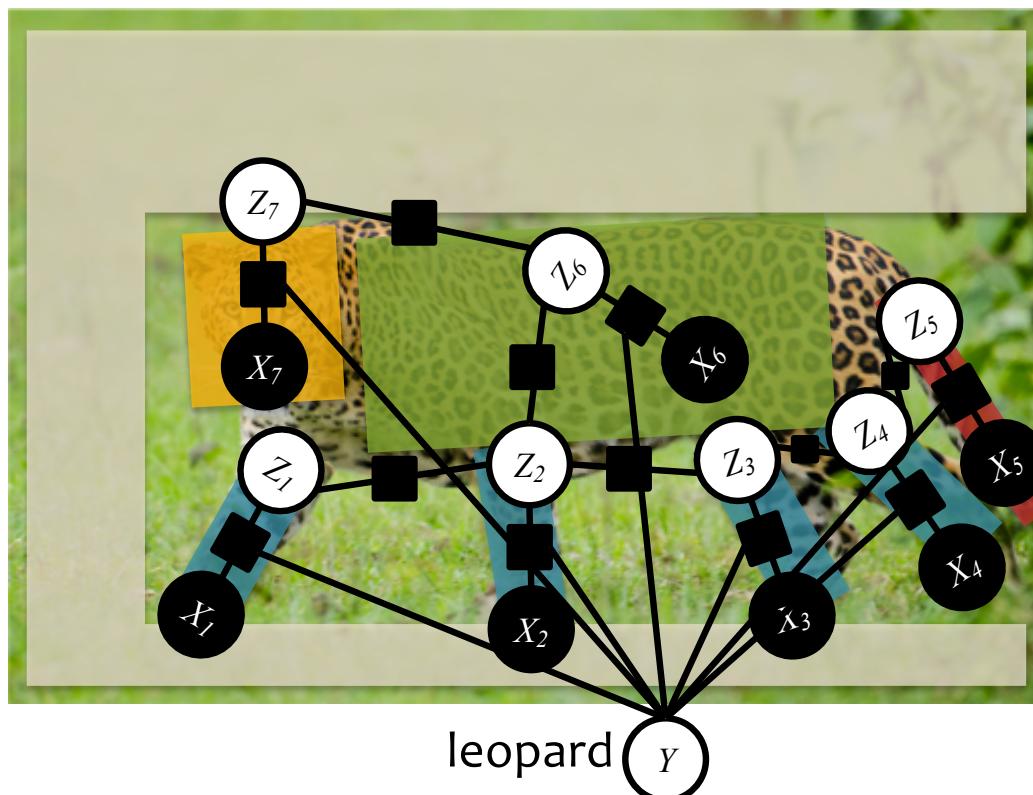


Hidden-state CRFs

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)

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

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:

$$\ell(\boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log \prod_{i=1}^N p(\mathbf{x}^{(i)}; \boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

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

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

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

This could be any arbitrary distribution parameterized by θ .

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

(Marginal) Log-likelihood:

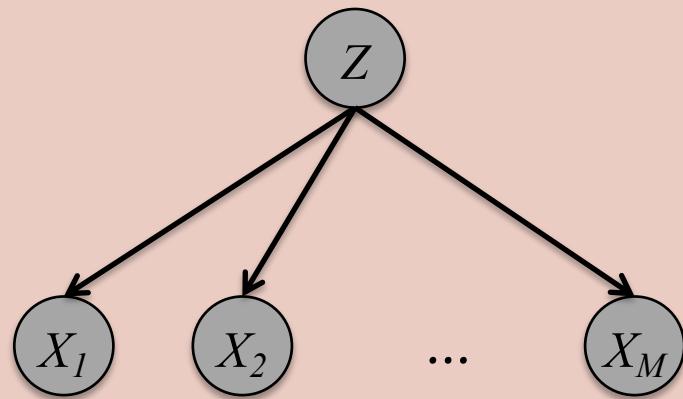
$$\ell(\theta) = \log \prod_{i=1}^N p_{\theta, \phi}(\mathbf{x}^{(i)})$$

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

Learning a Mixture Model

Supervised Learning: The parameters decouple!

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



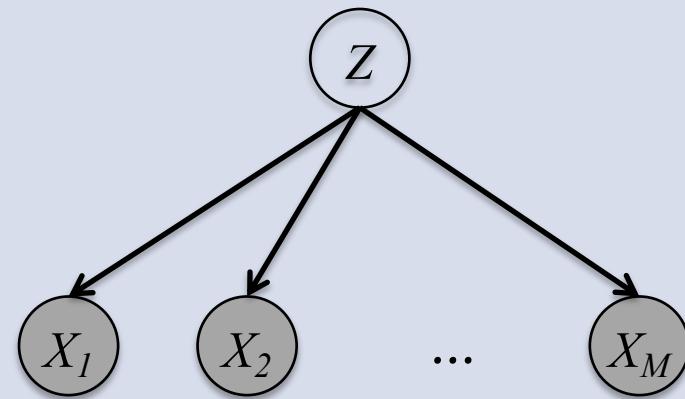
$$\theta^*, \phi^* = \operatorname{argmax}_{\theta, \phi} \sum_{i=1}^N \log p_\theta(\mathbf{x}^{(i)} | z^{(i)}) p_\phi(z^{(i)})$$

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

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

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!

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

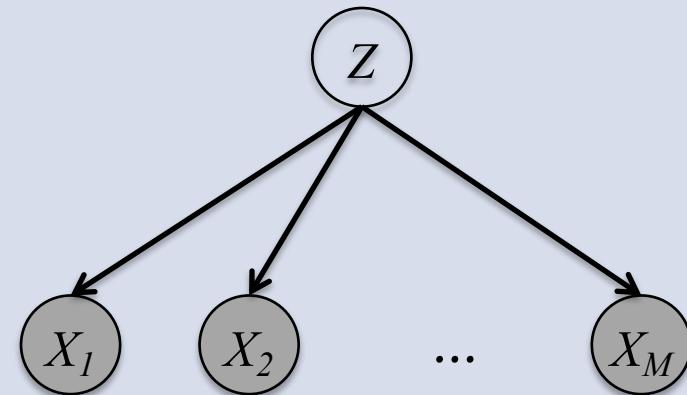
Training certainly isn't as simple as the supervised case.

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,  $\phi, \boldsymbol{\mu}, \Sigma$ 
3:   while not converged do
4:     E-Step:

```

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

```
5:   M-Step:
```

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

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

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

```
6:   return  $(\phi, \boldsymbol{\mu}, \Sigma)$ 
```

Algorithm 1 Soft EM for GMMs

```

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

```

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

```
5:   M-Step:
```

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

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

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

```
6:   return  $(\phi, \boldsymbol{\mu}, \Sigma)$ 
```

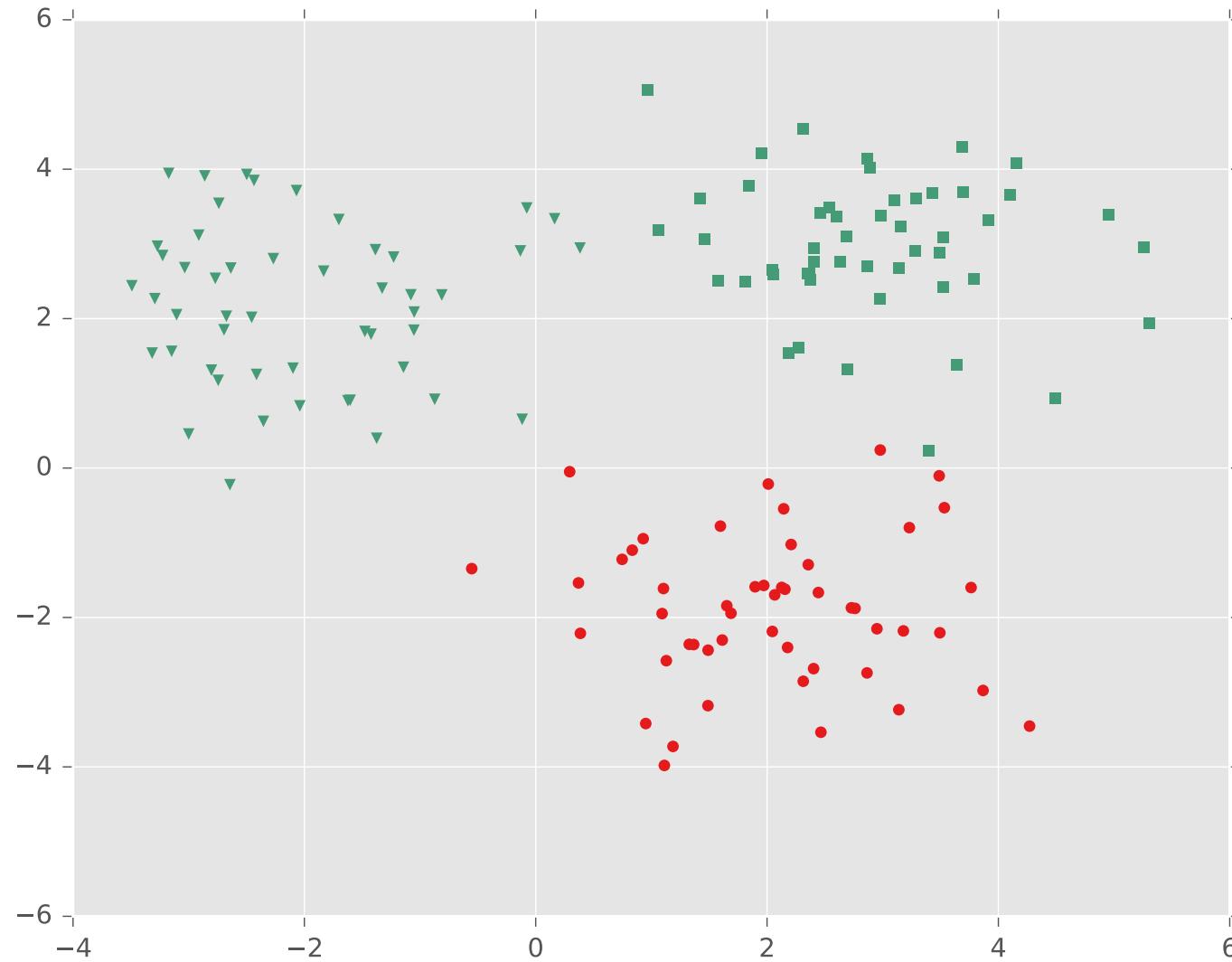
Posterior Inference for Mixture Model

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

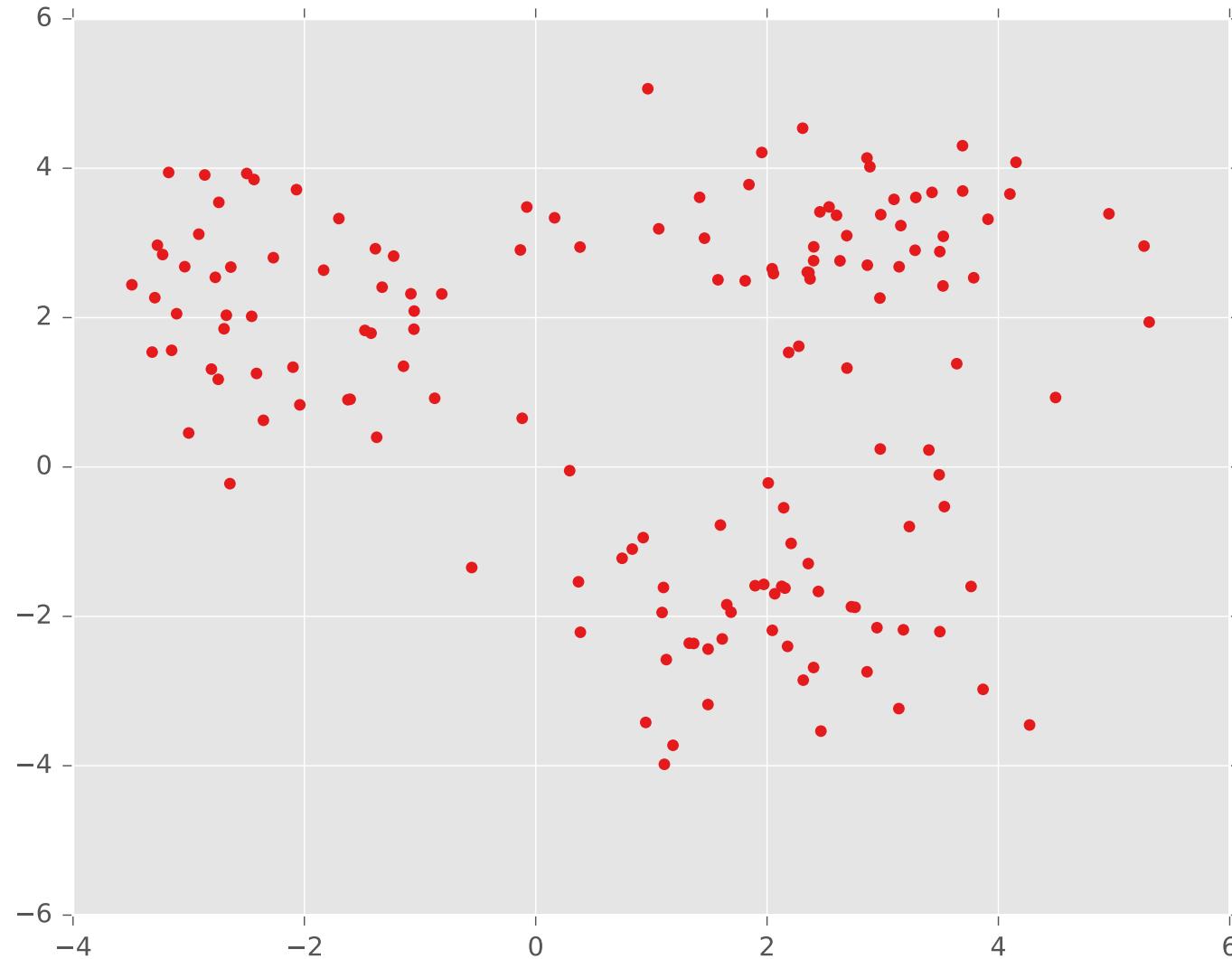
$$p(z^{(i)} = k | \mathbf{x}^{(i)}; \phi, \mu, \Sigma) = \frac{p(\mathbf{x}^{(i)} | z^{(i)} = k; \mu, \Sigma)p(z^{(i)} = k; \phi)}{\sum_{j=1}^K p(\mathbf{x}^{(i)} | z^{(i)} = j; \mu, \Sigma)p(z^{(i)} = j; \phi)} \quad (1)$$

EXAMPLE: K-MEANS VS GMM

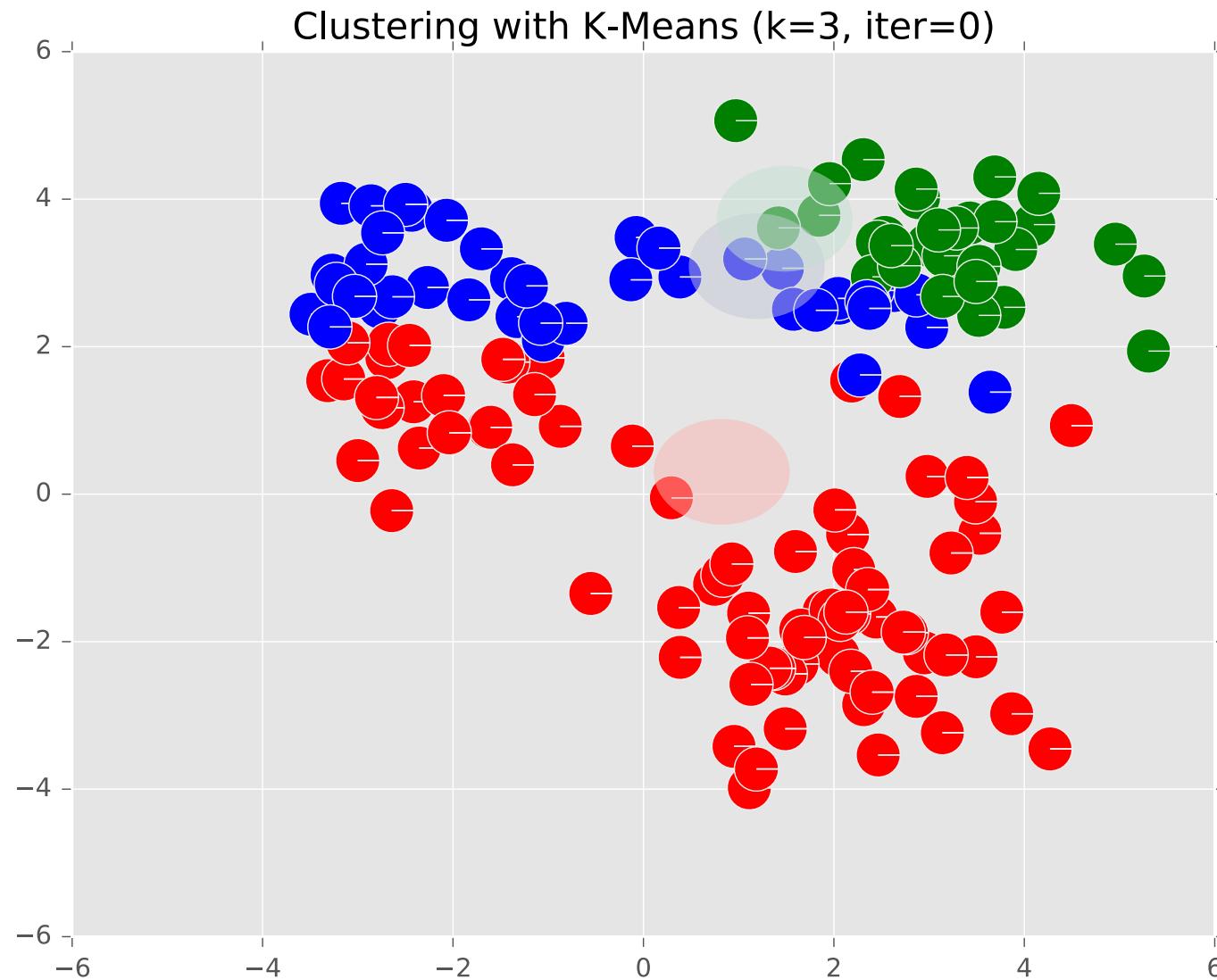
Example: K-Means



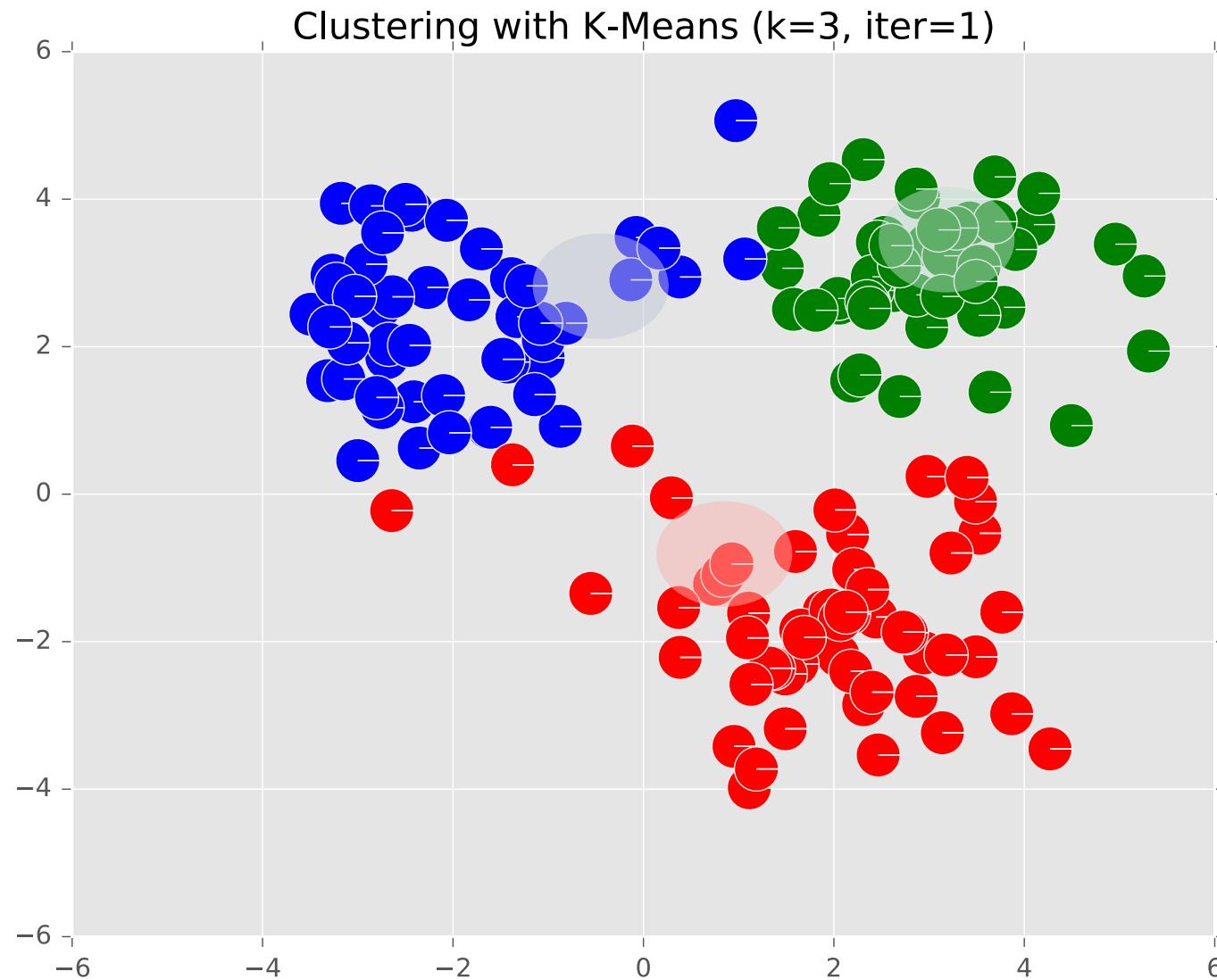
Example: K-Means



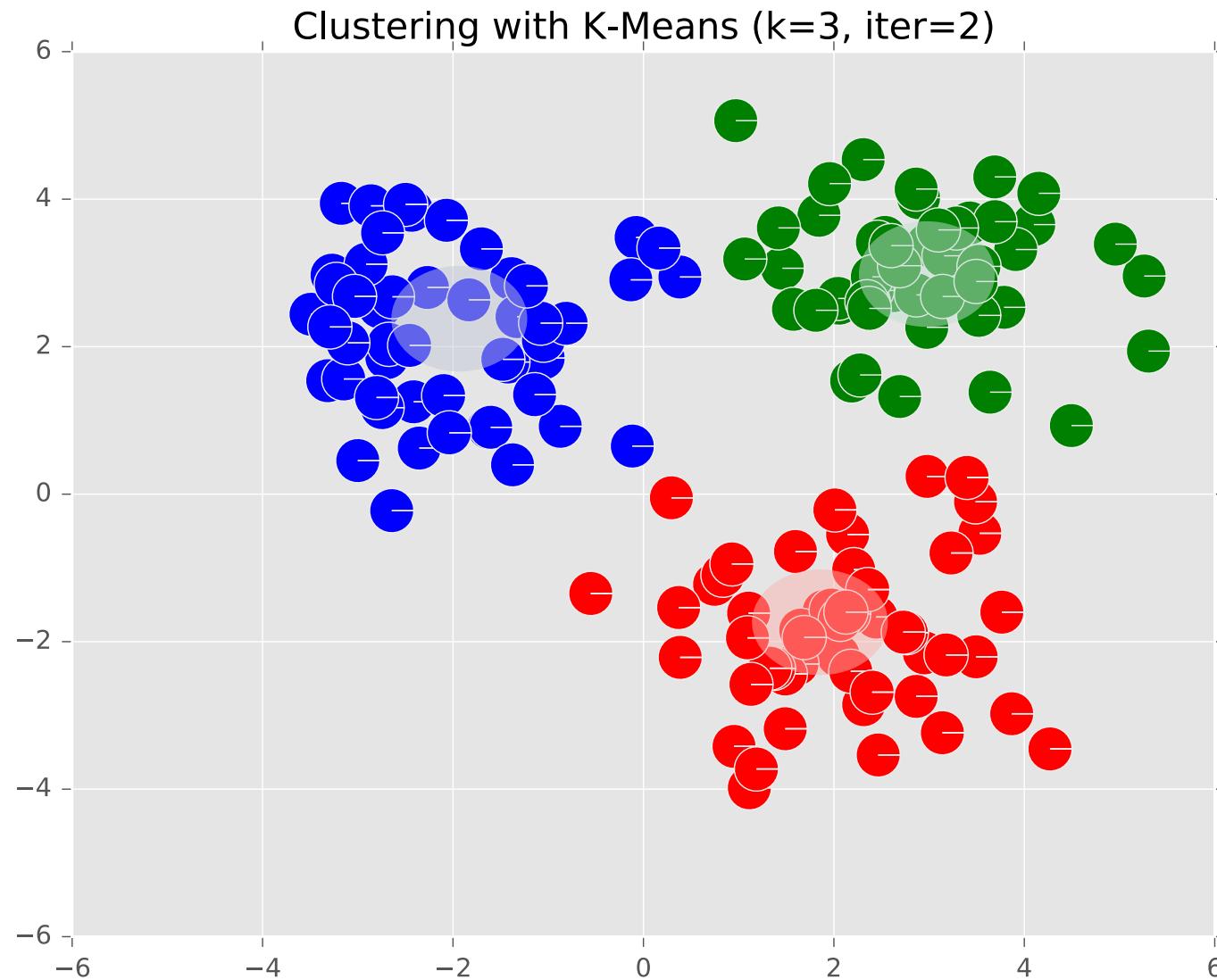
Example: K-Means



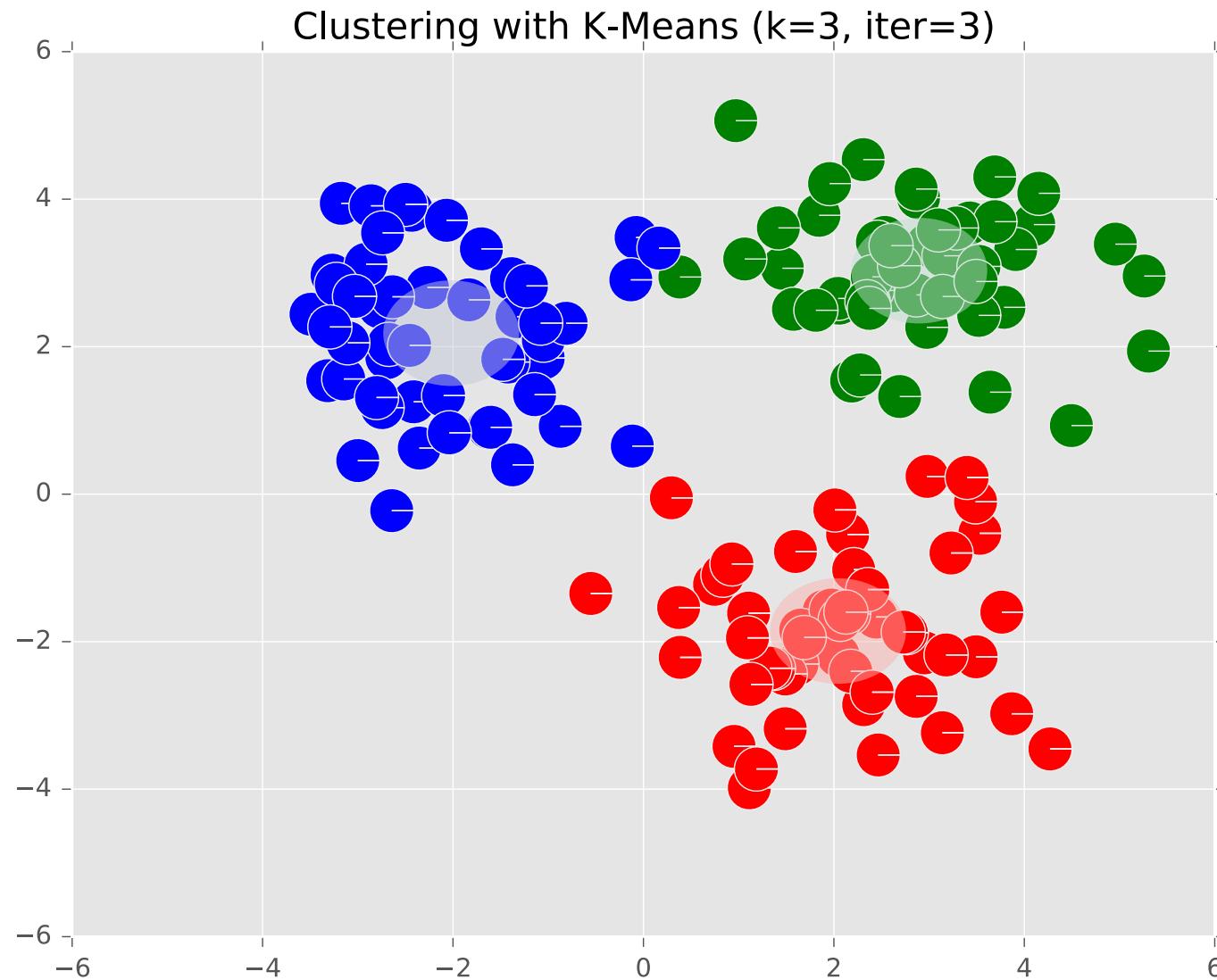
Example: K-Means



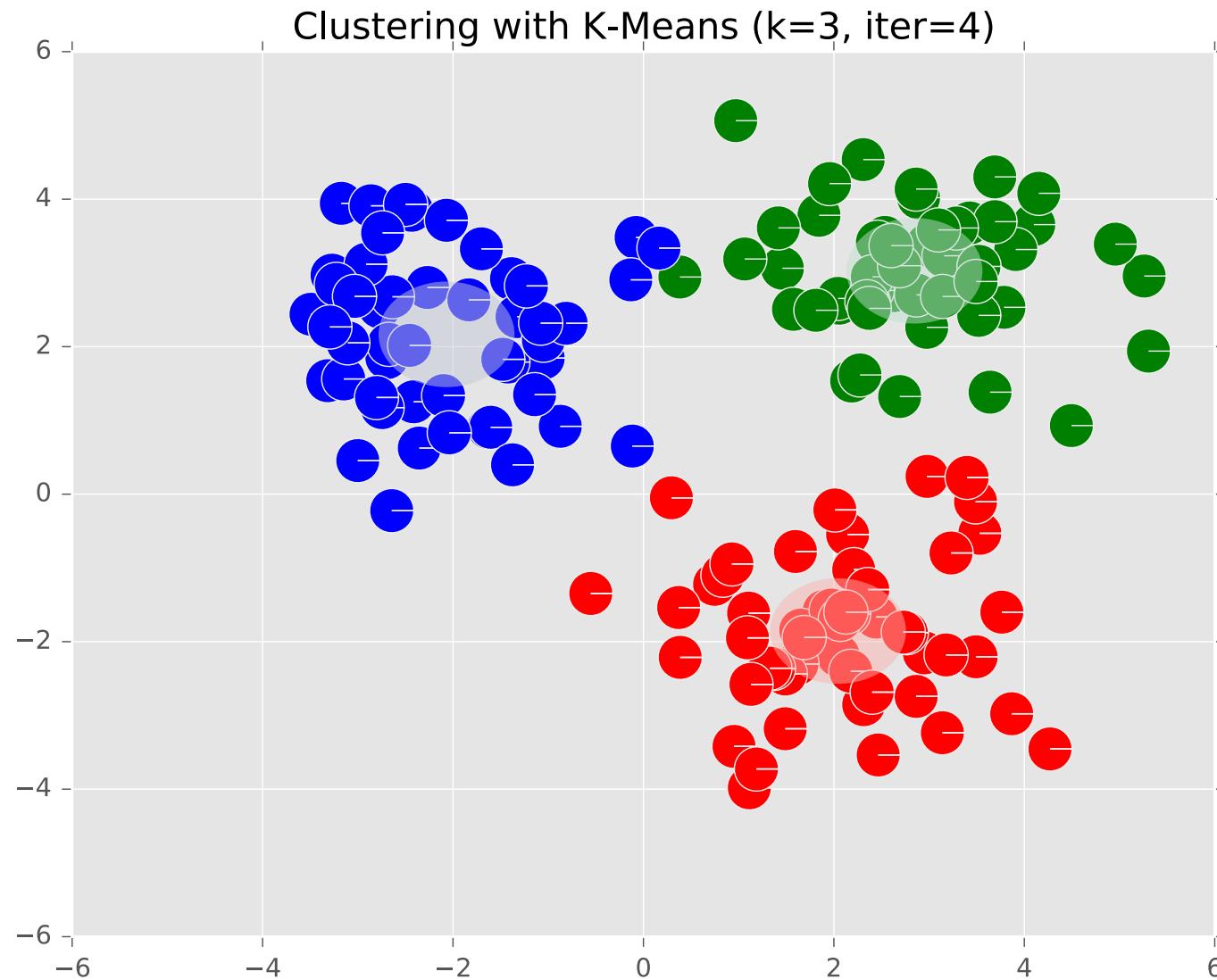
Example: K-Means



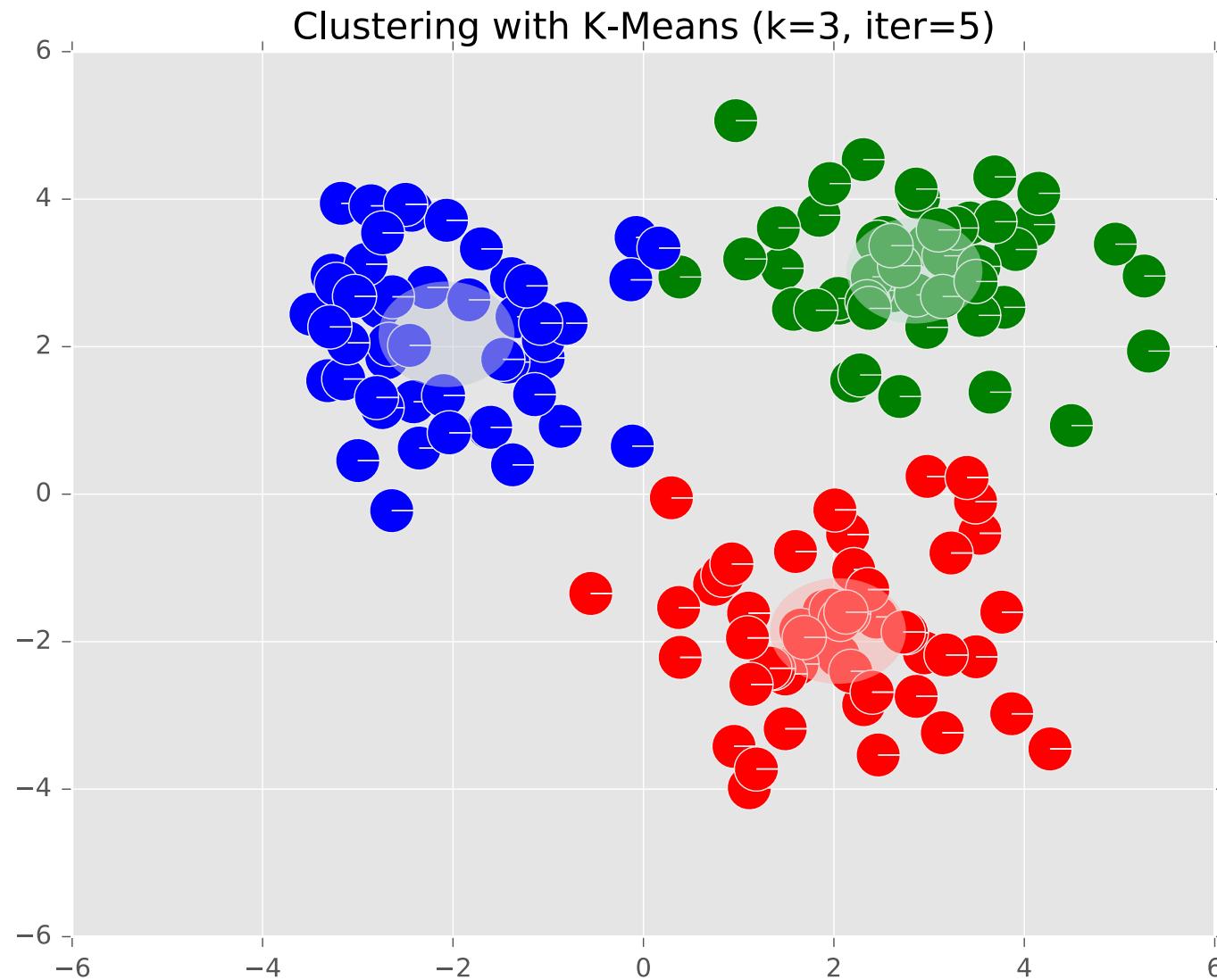
Example: K-Means



Example: K-Means



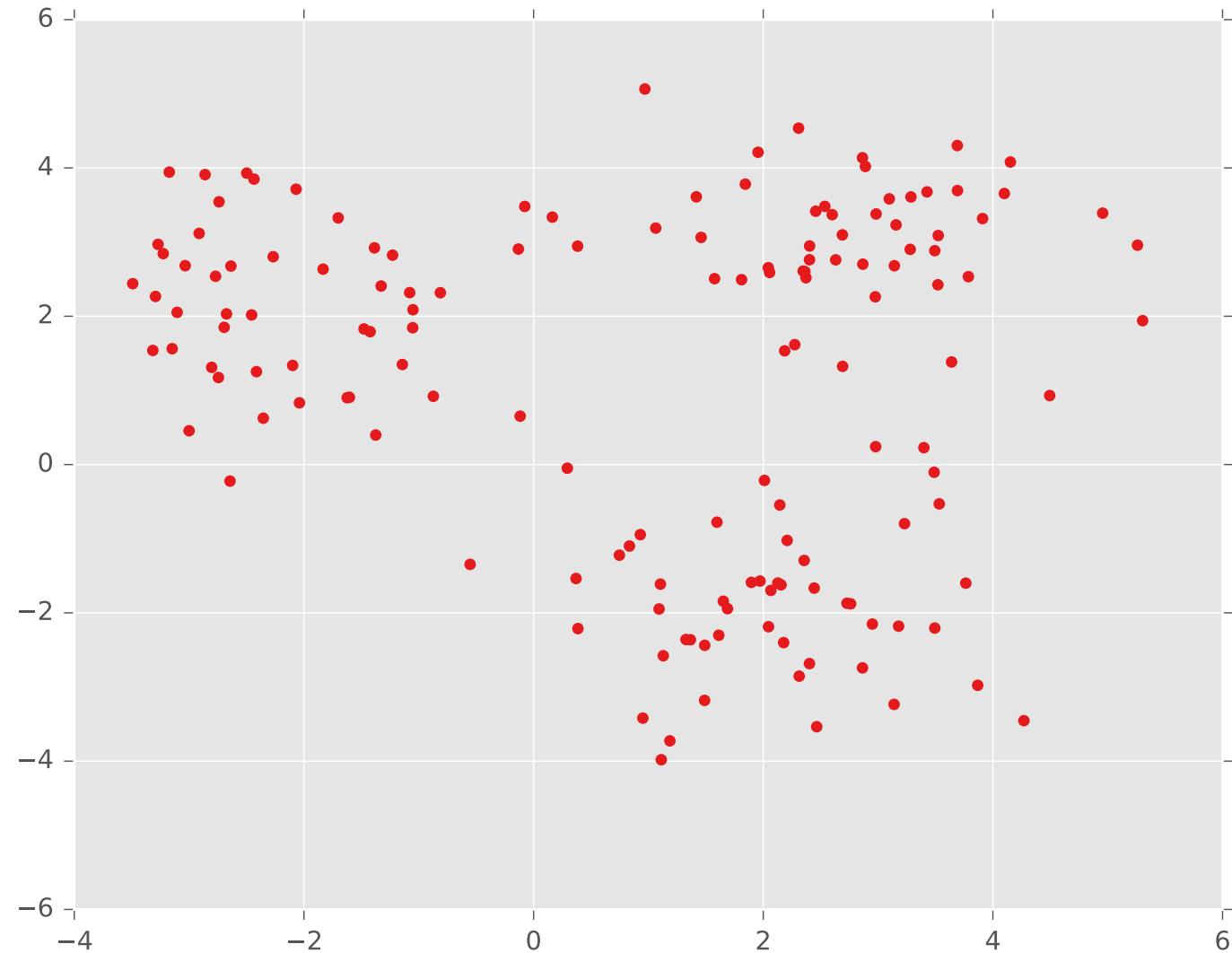
Example: K-Means



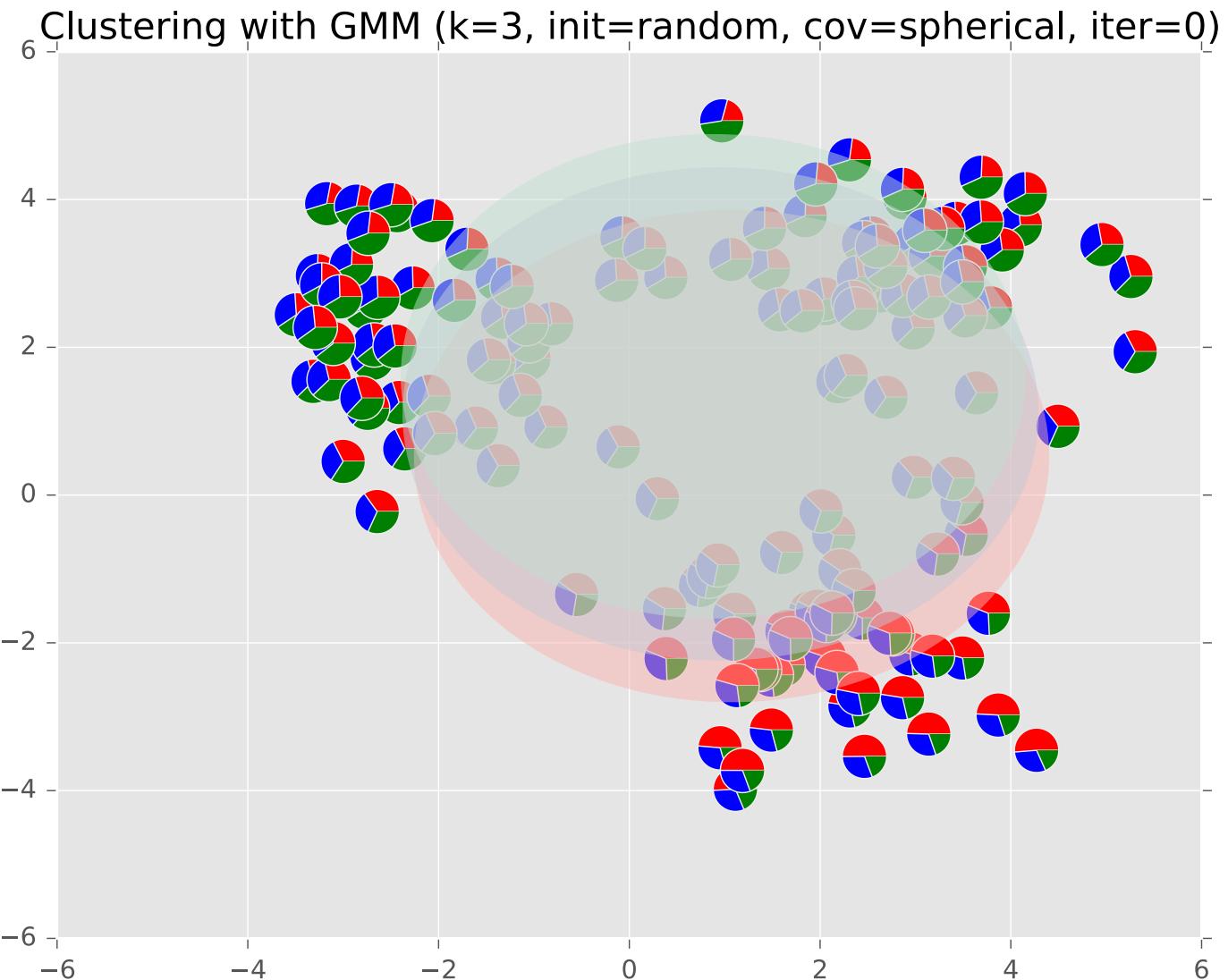
Example: GMM



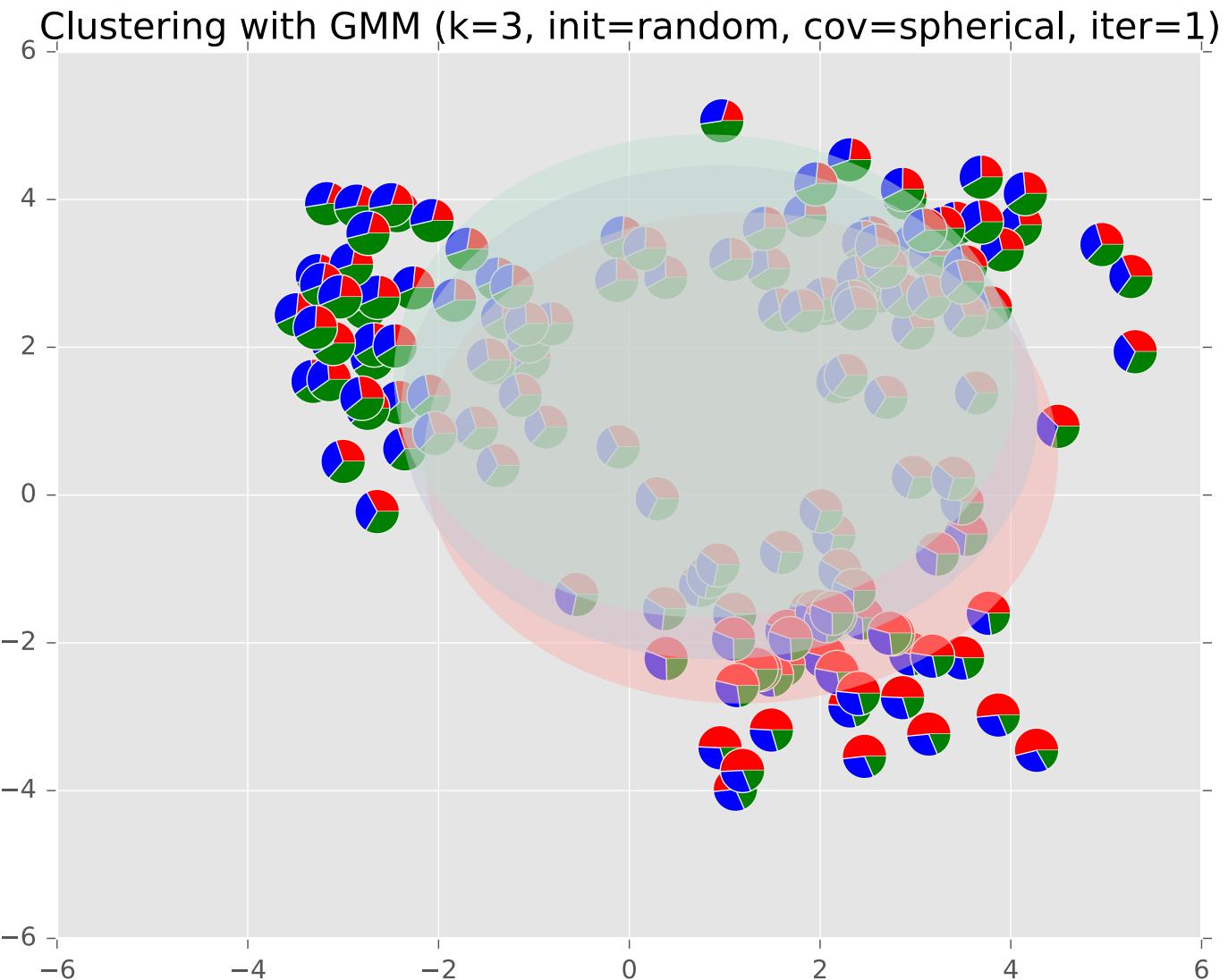
Example: GMM



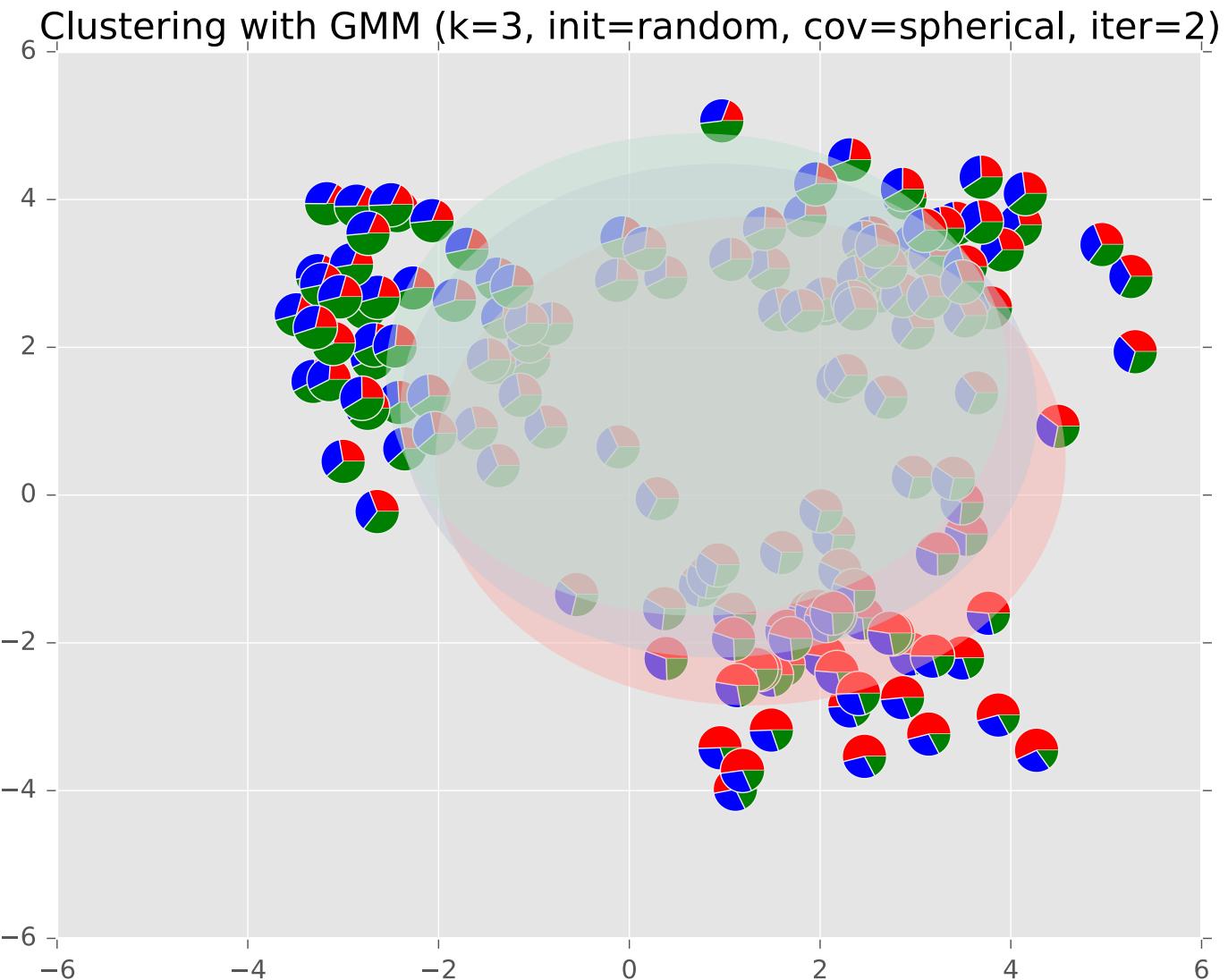
Example: GMM



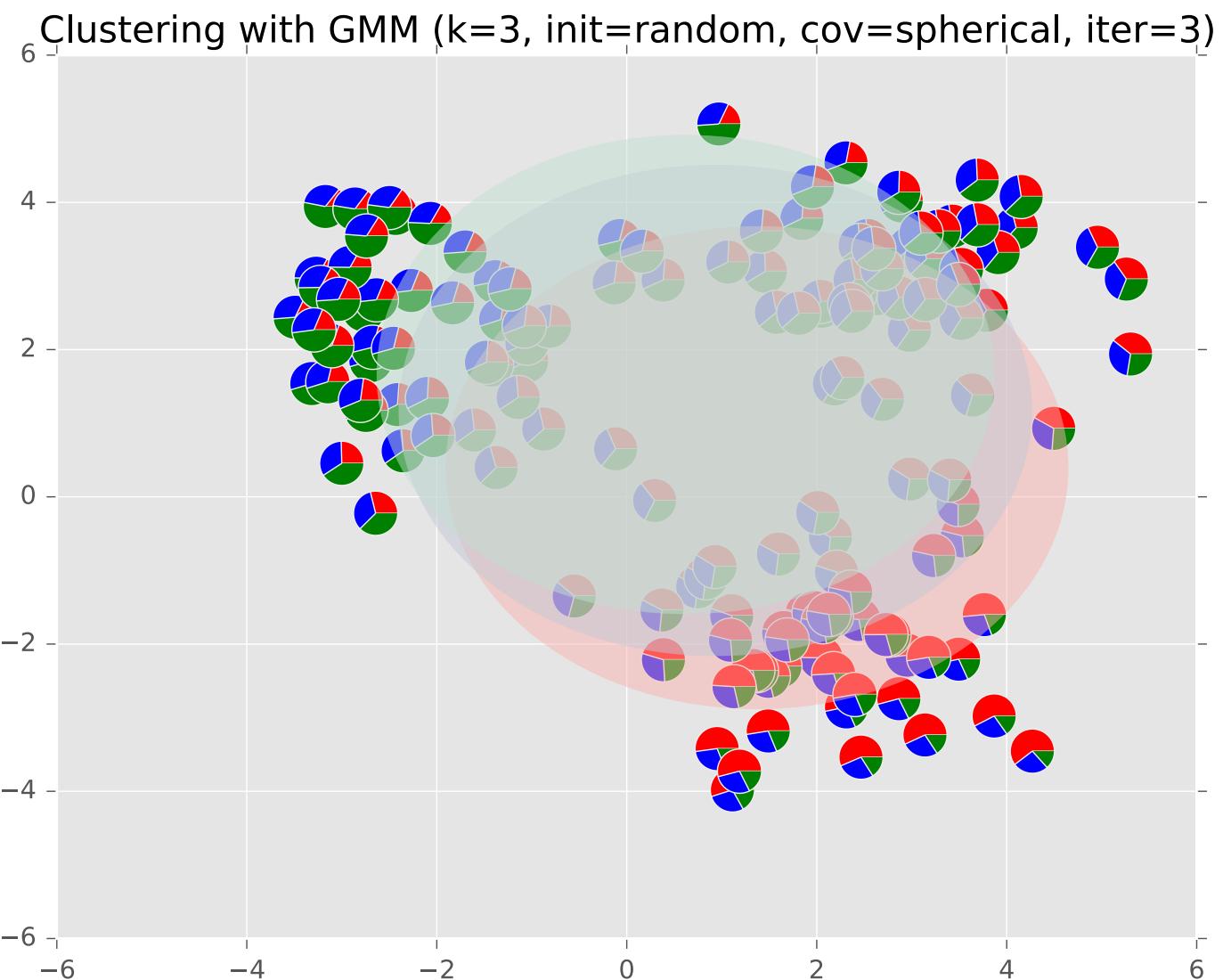
Example: GMM



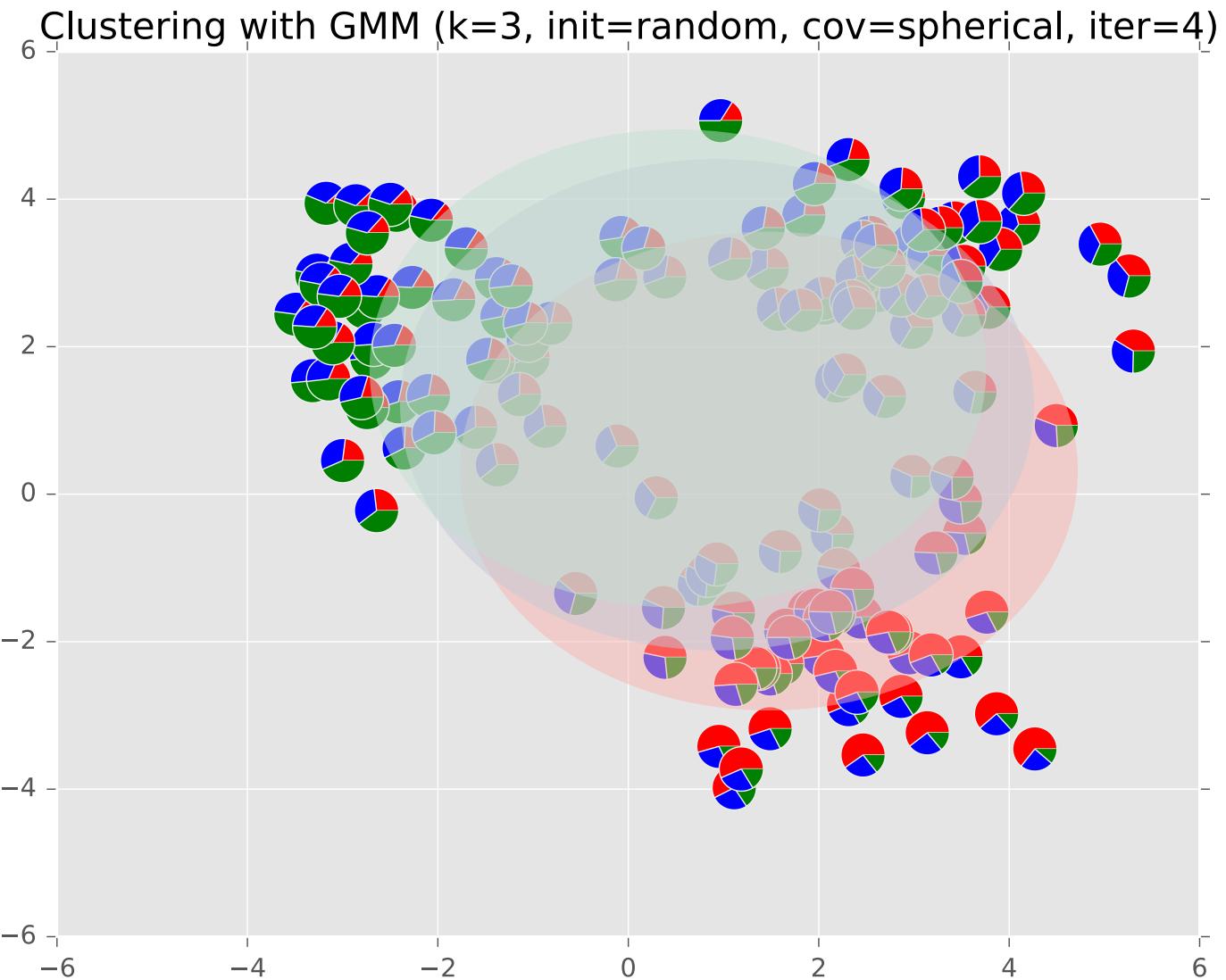
Example: GMM



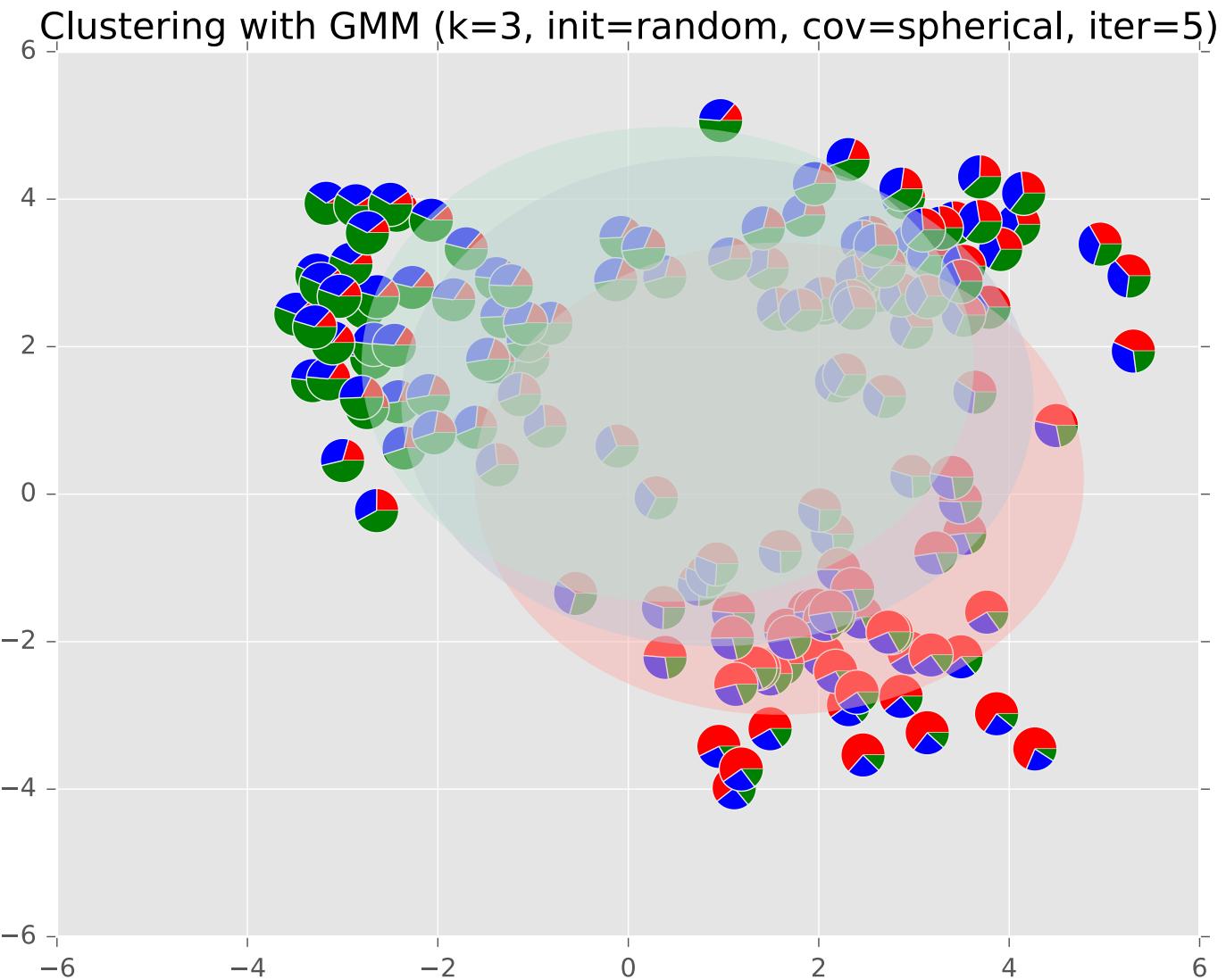
Example: GMM



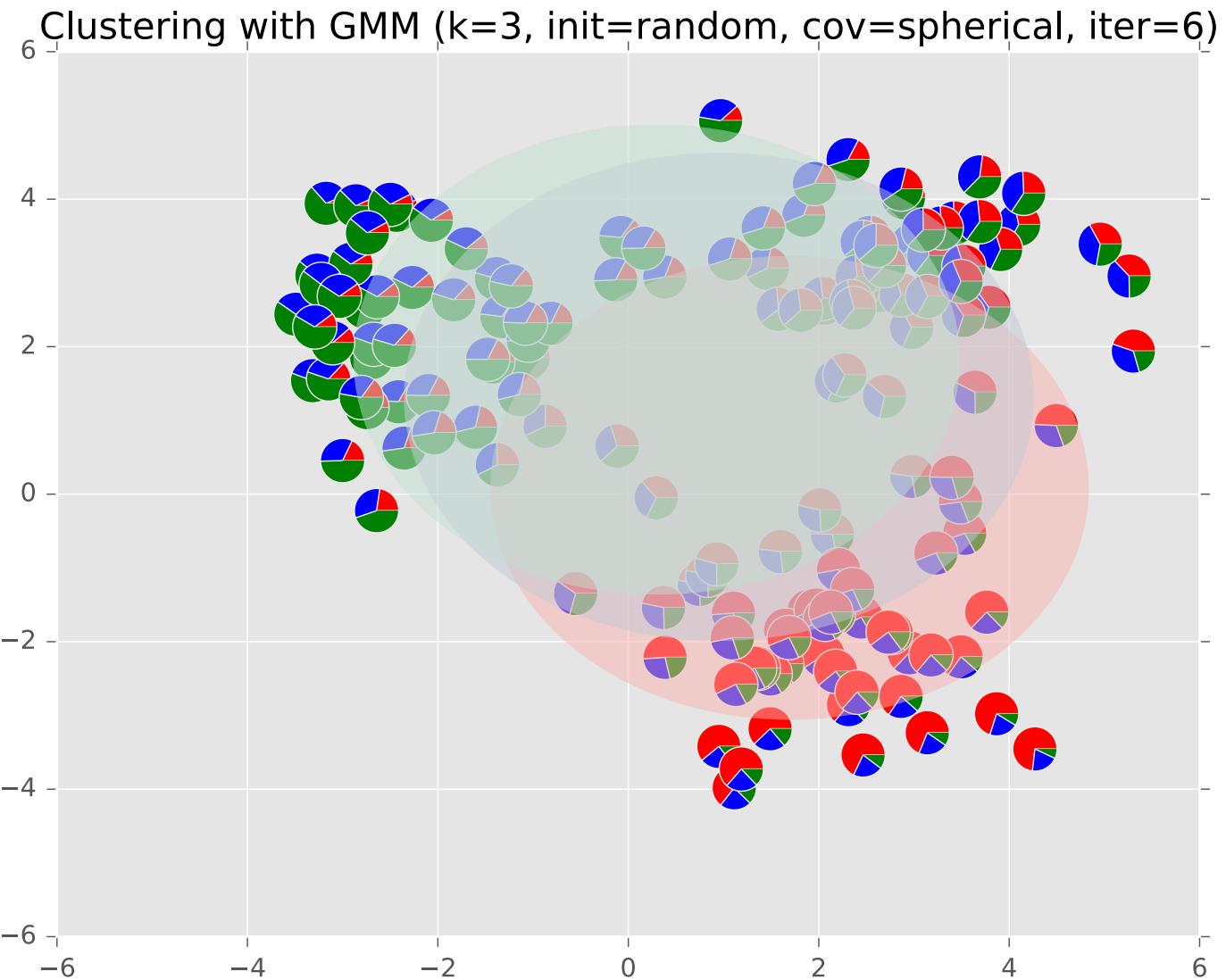
Example: GMM



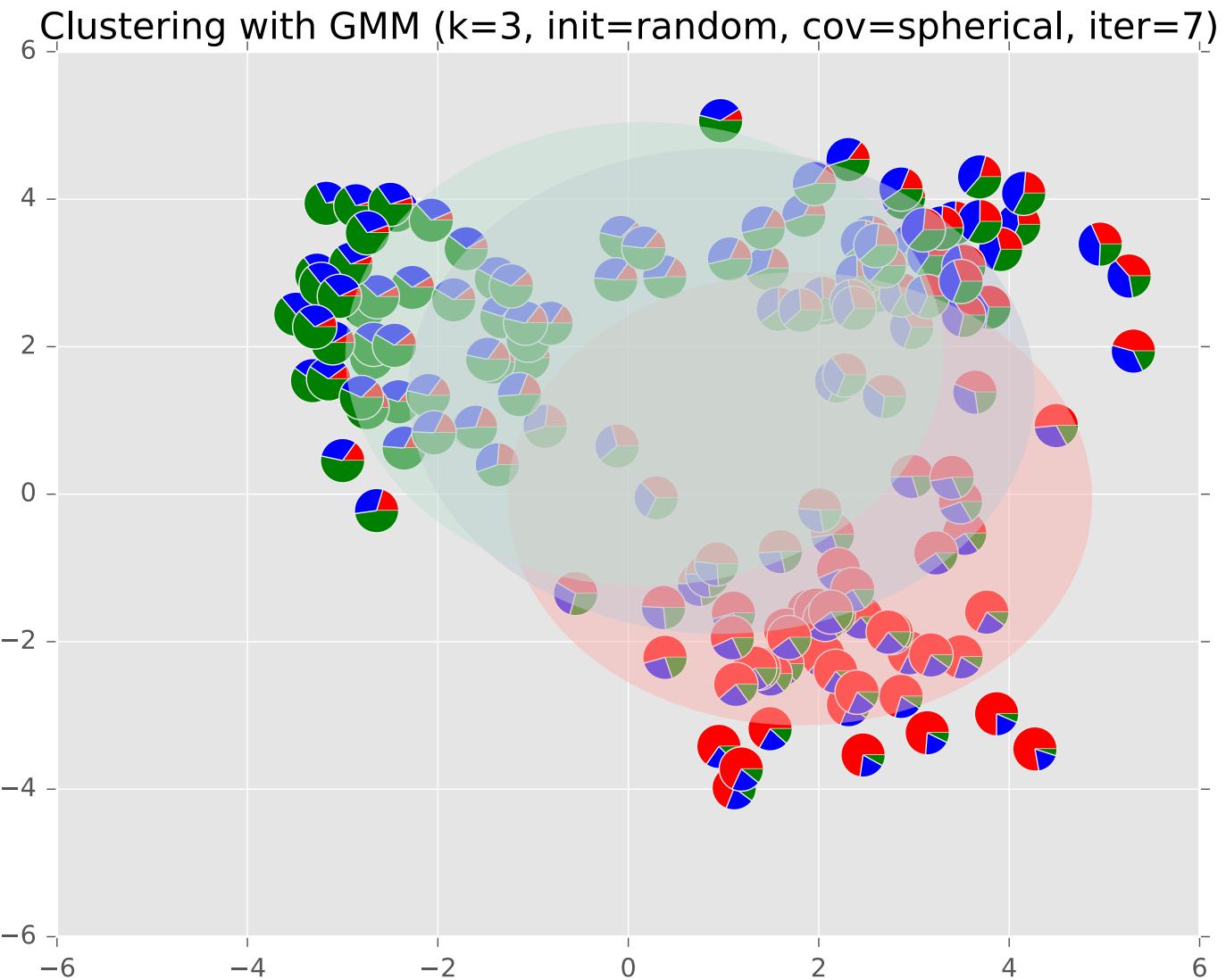
Example: GMM



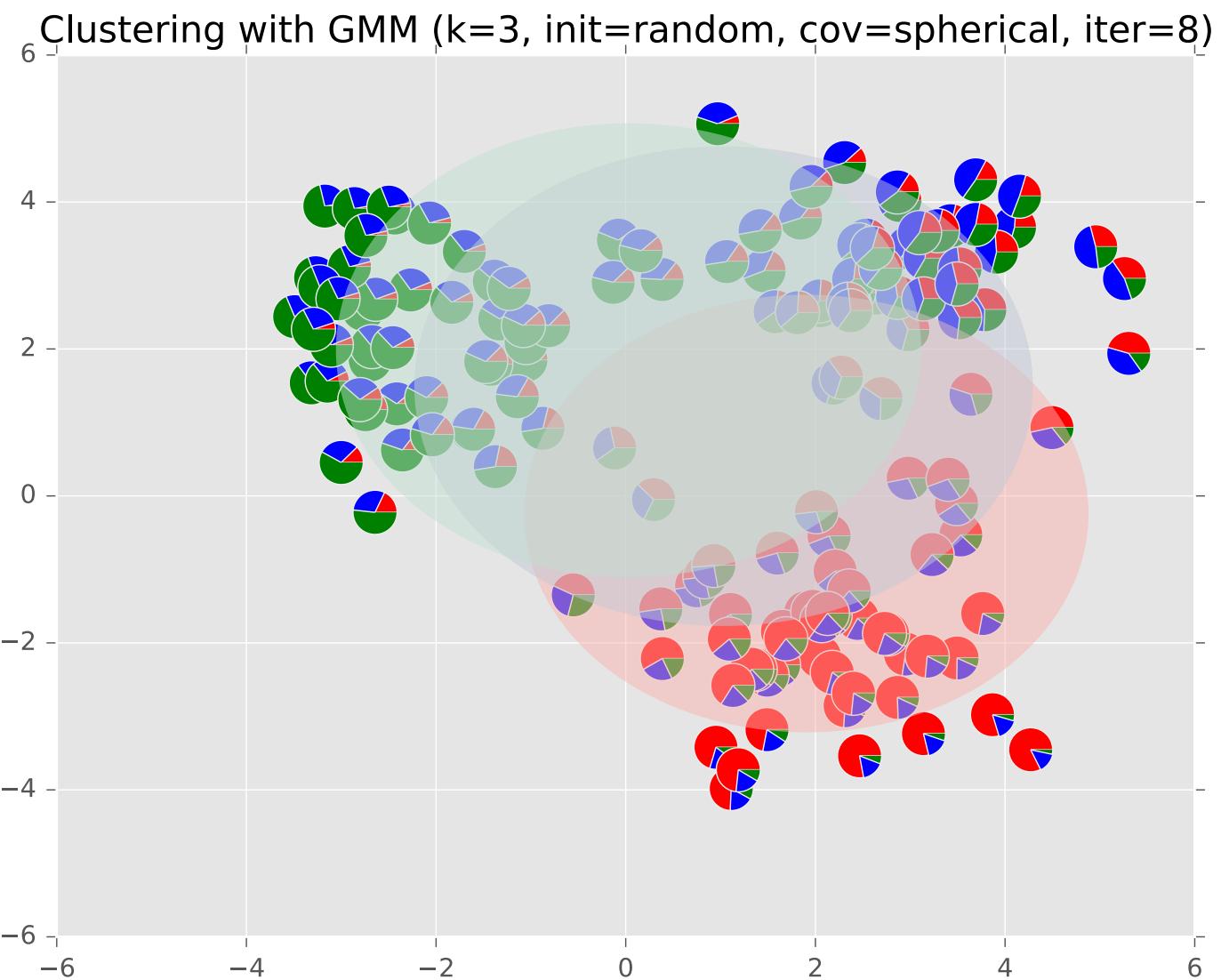
Example: GMM



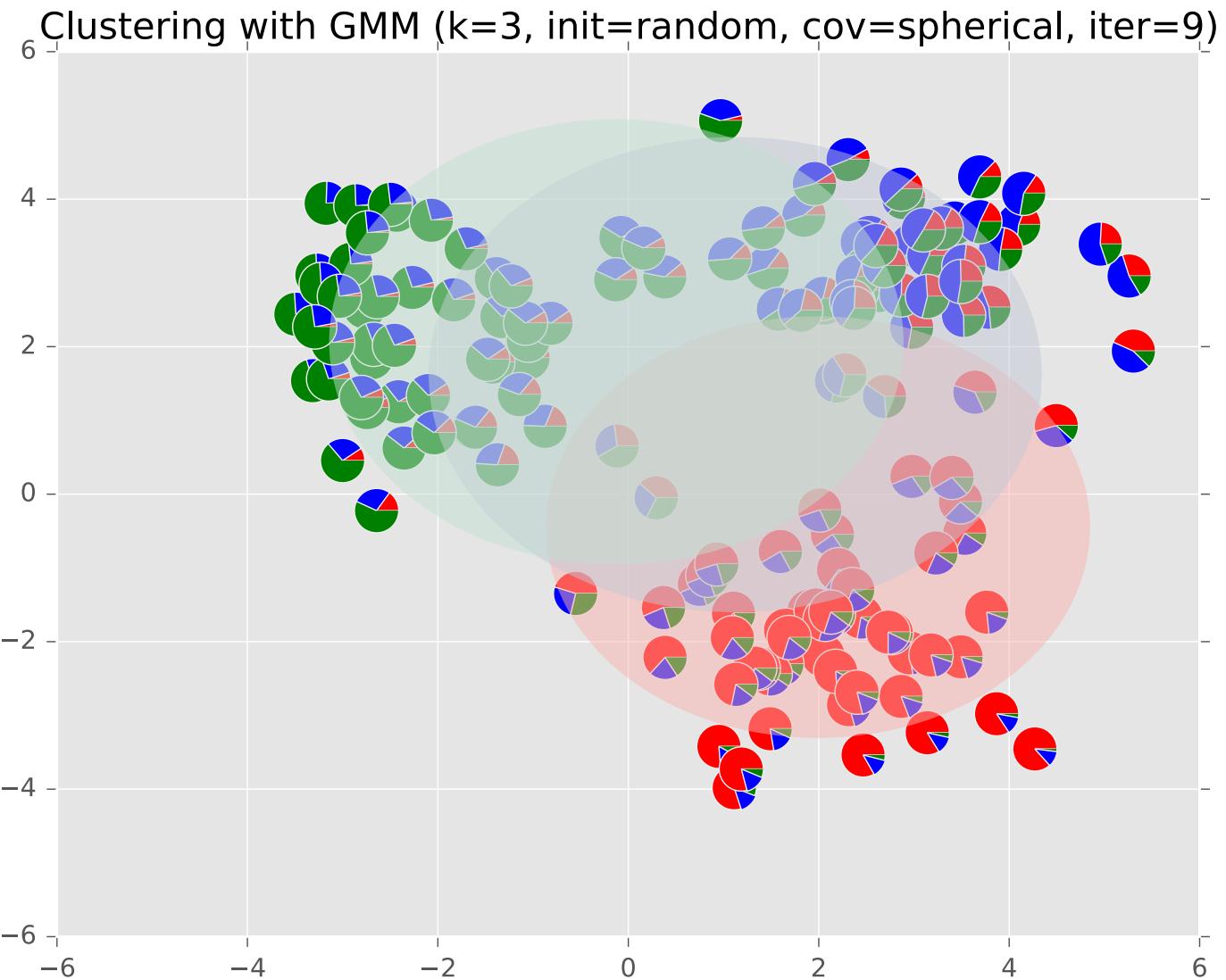
Example: GMM



Example: GMM

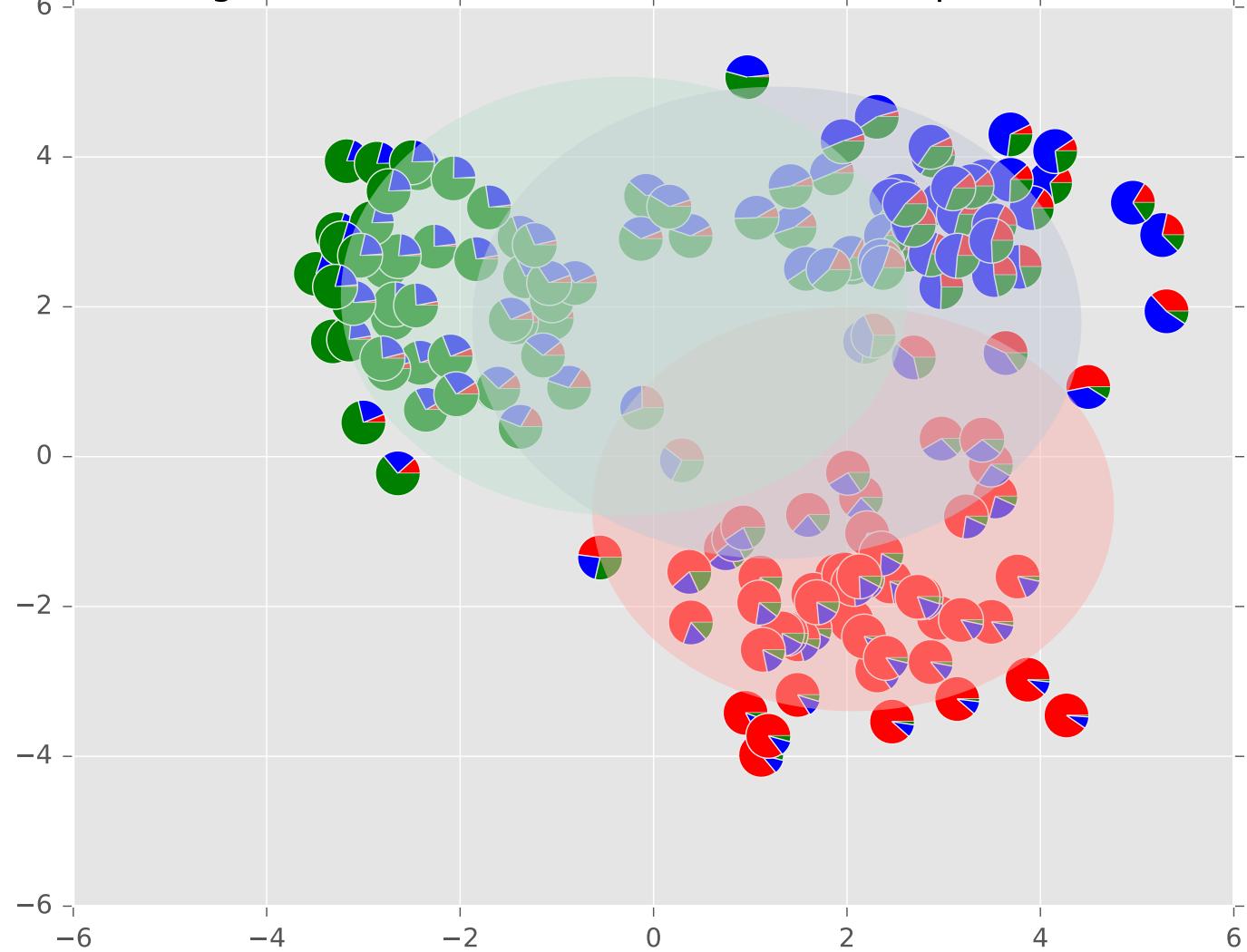


Example: GMM

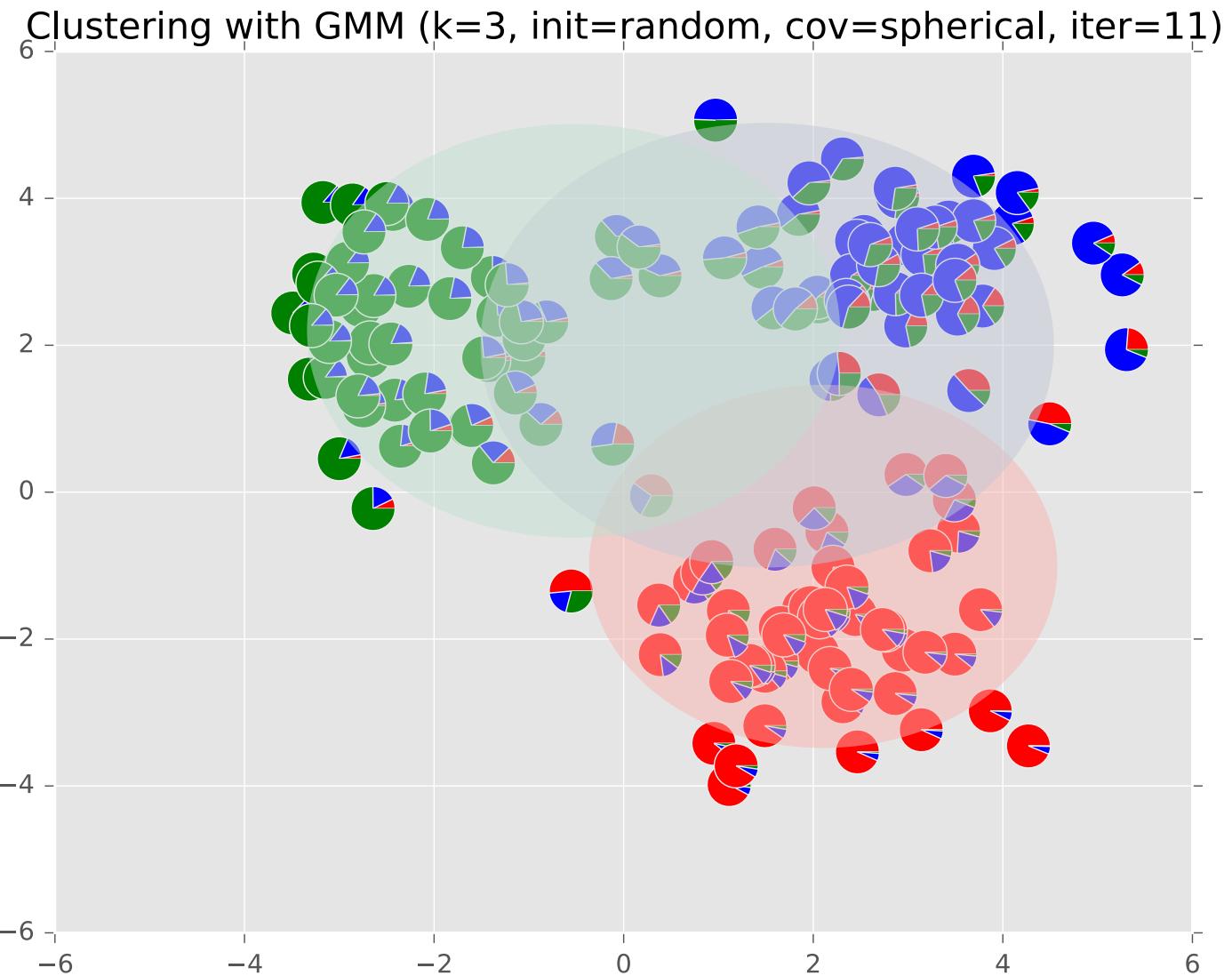


Example: GMM

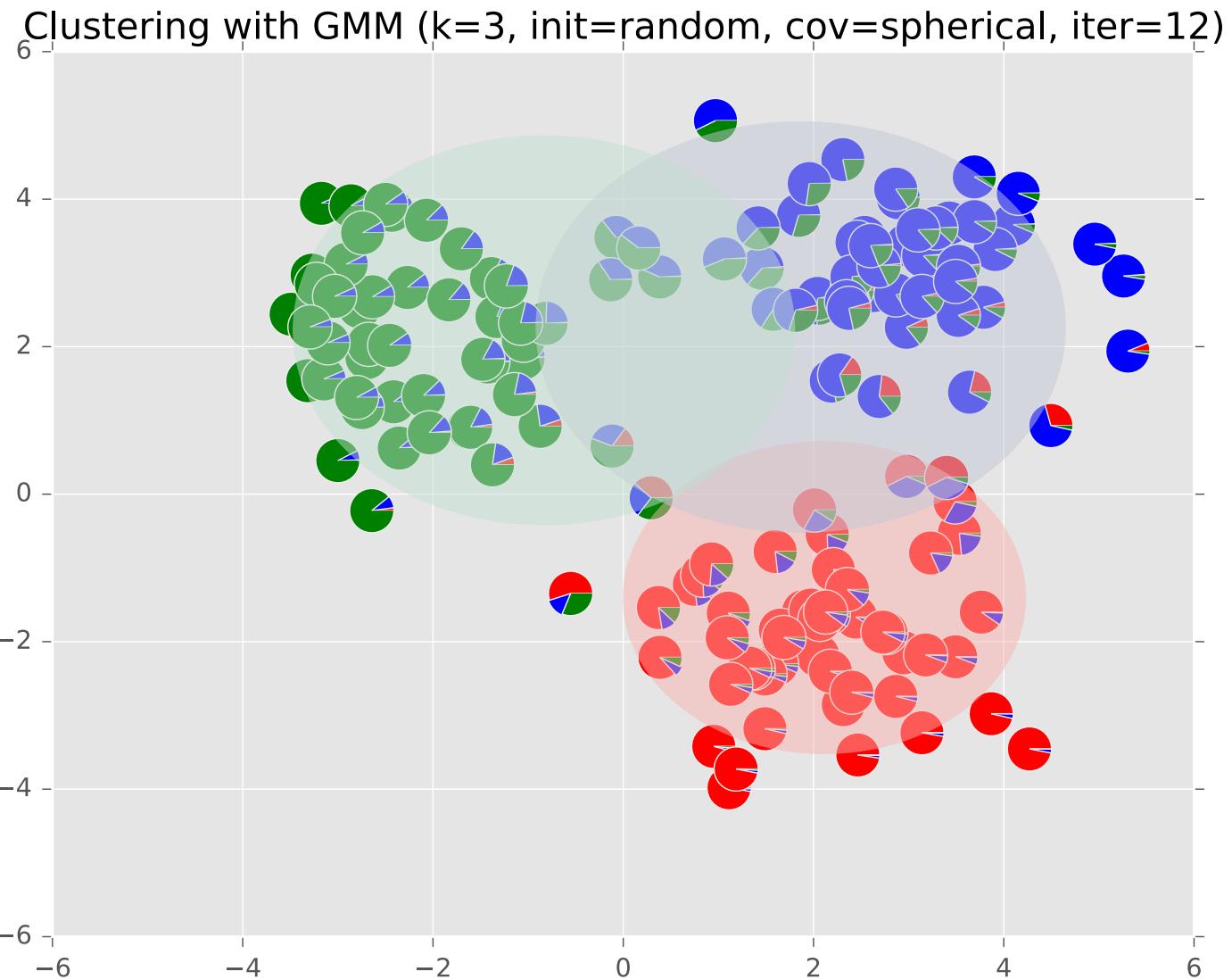
Clustering with GMM ($k=3$, init=random, cov=spherical, iter=10)



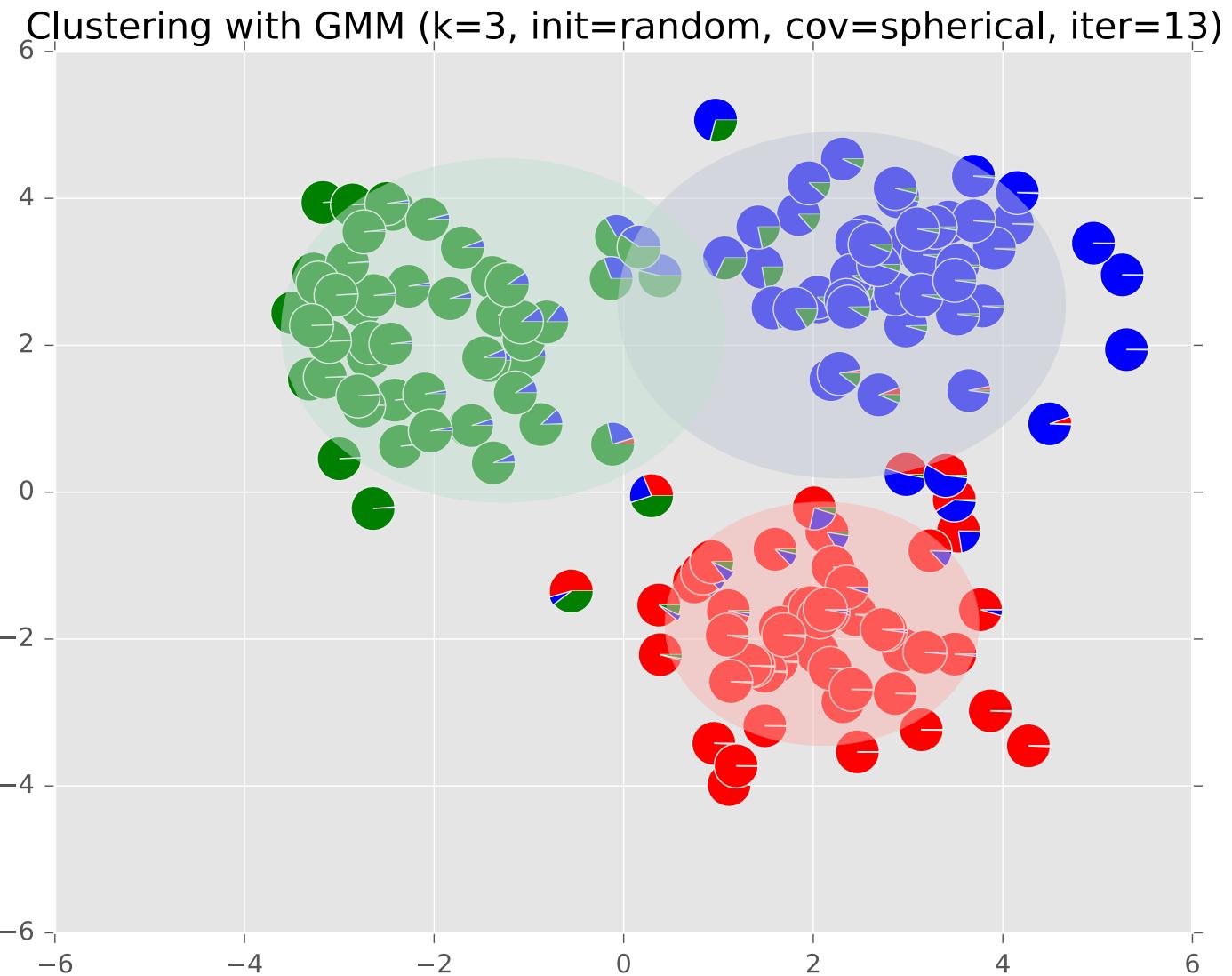
Example: GMM



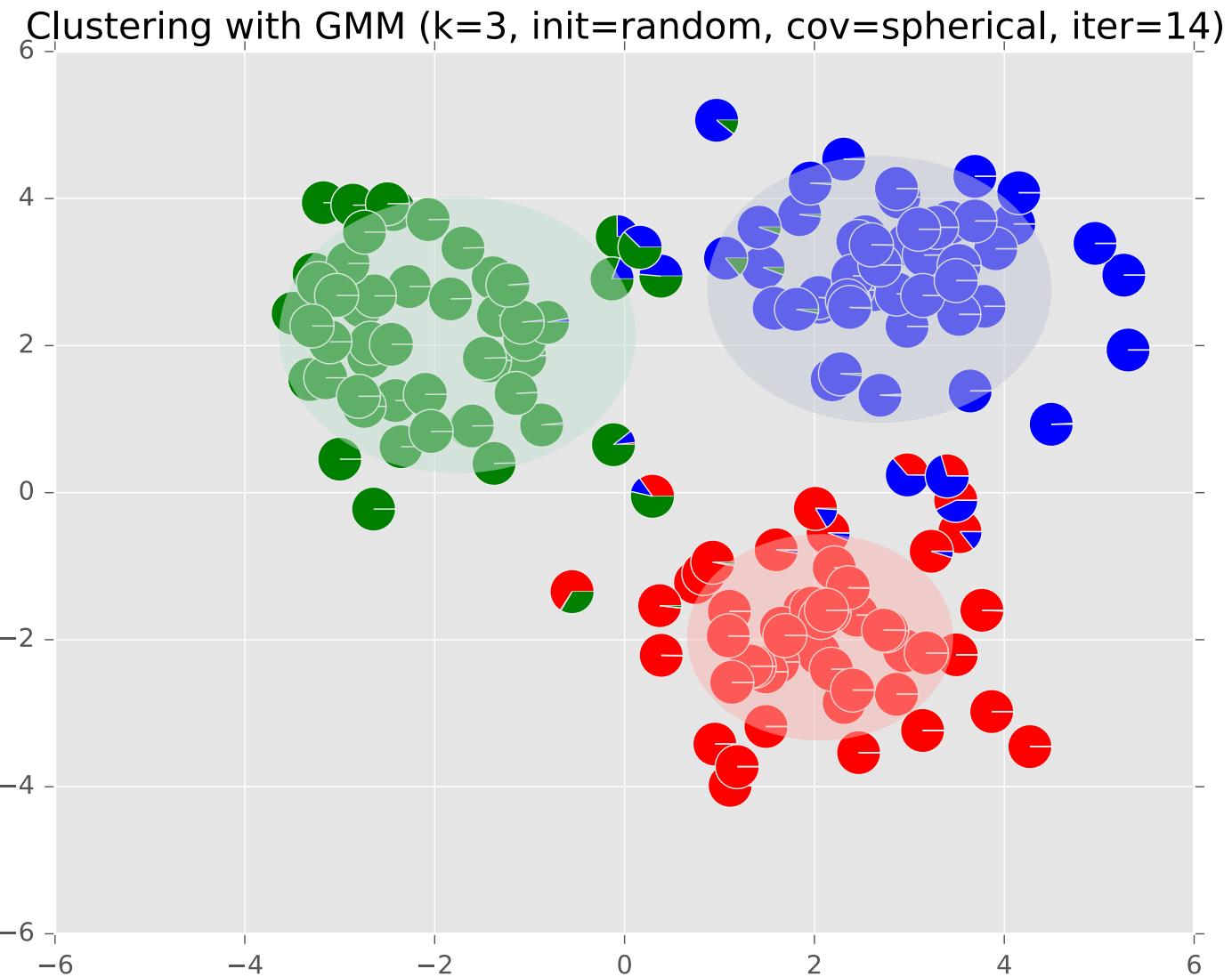
Example: GMM



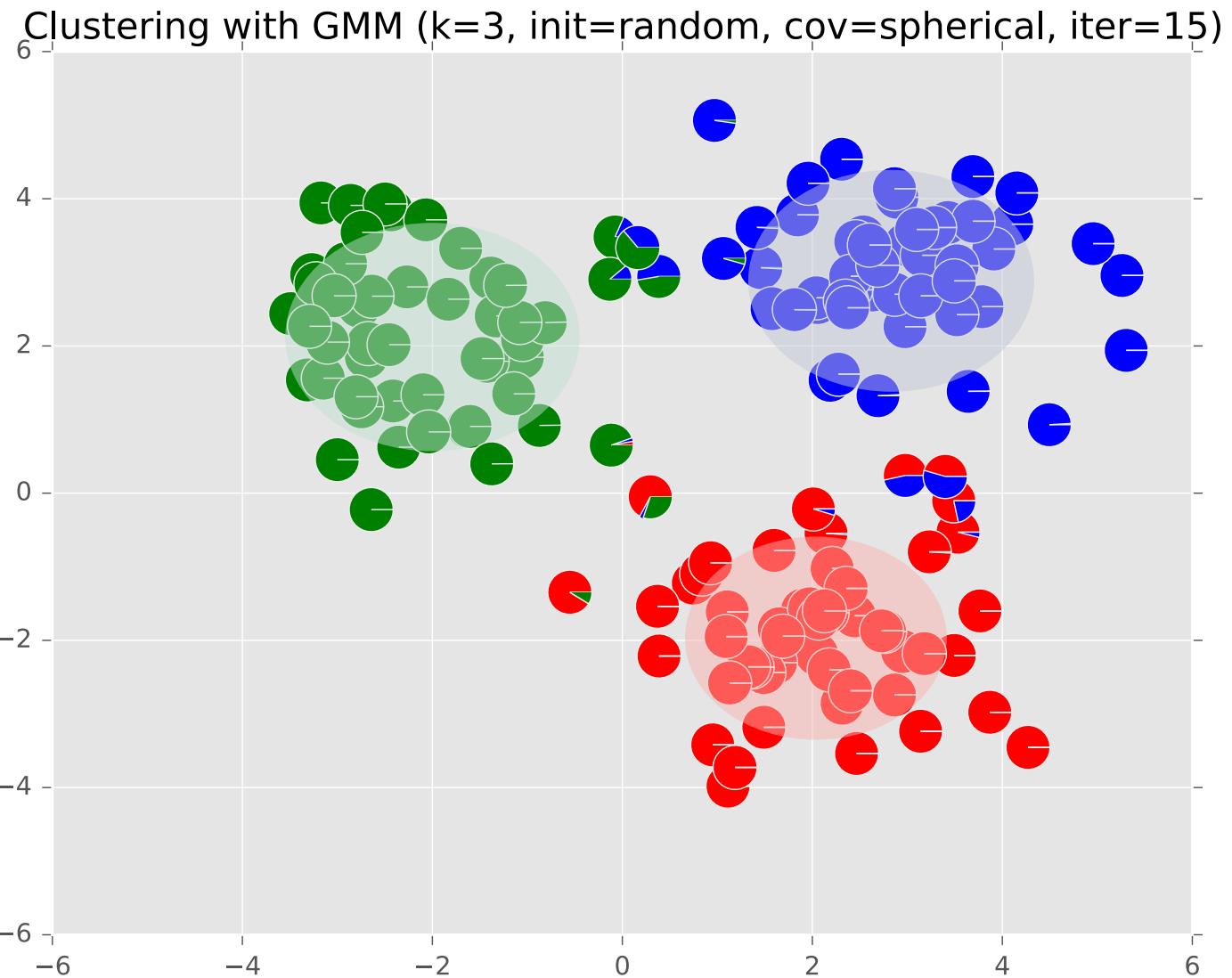
Example: GMM



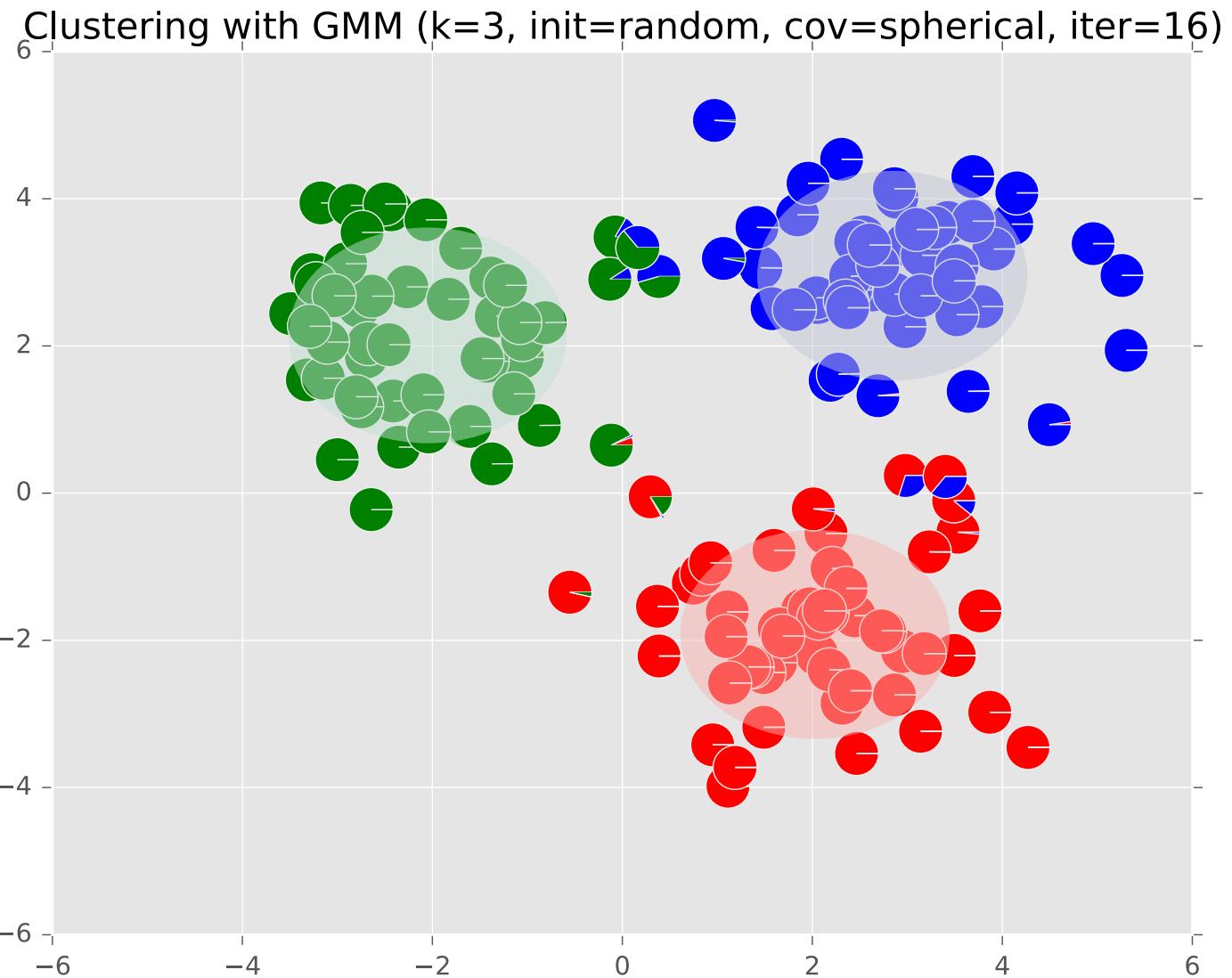
Example: GMM



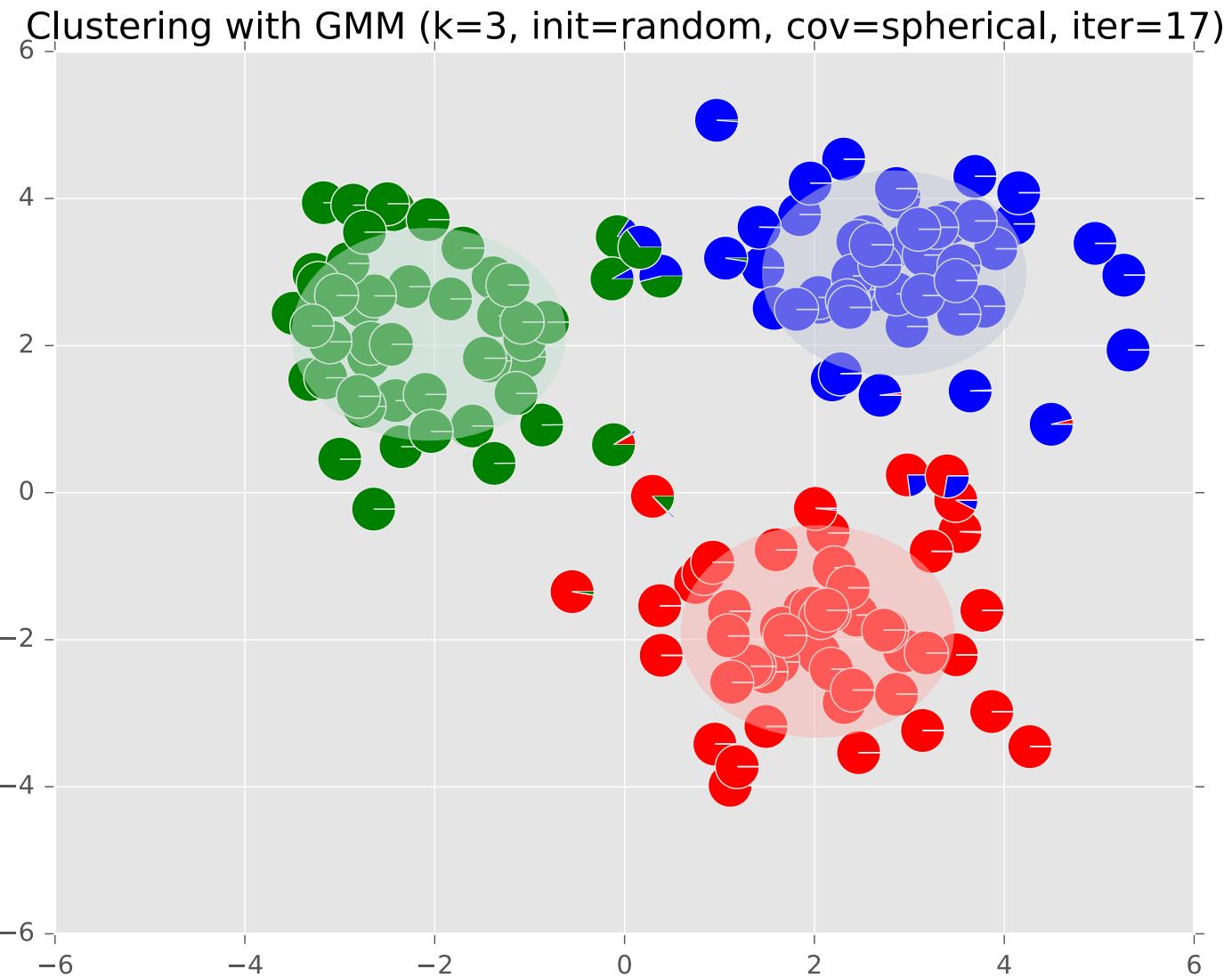
Example: GMM



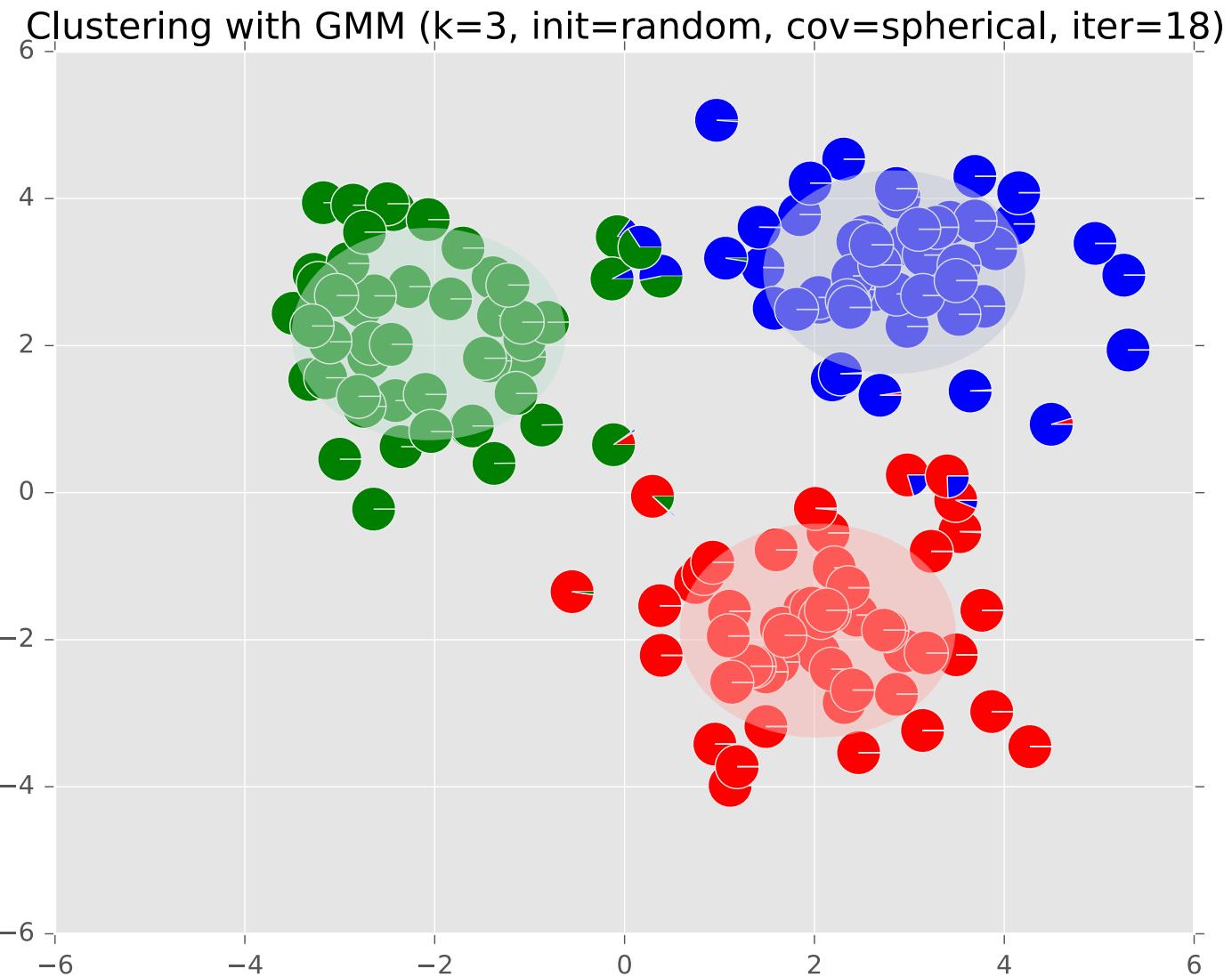
Example: GMM



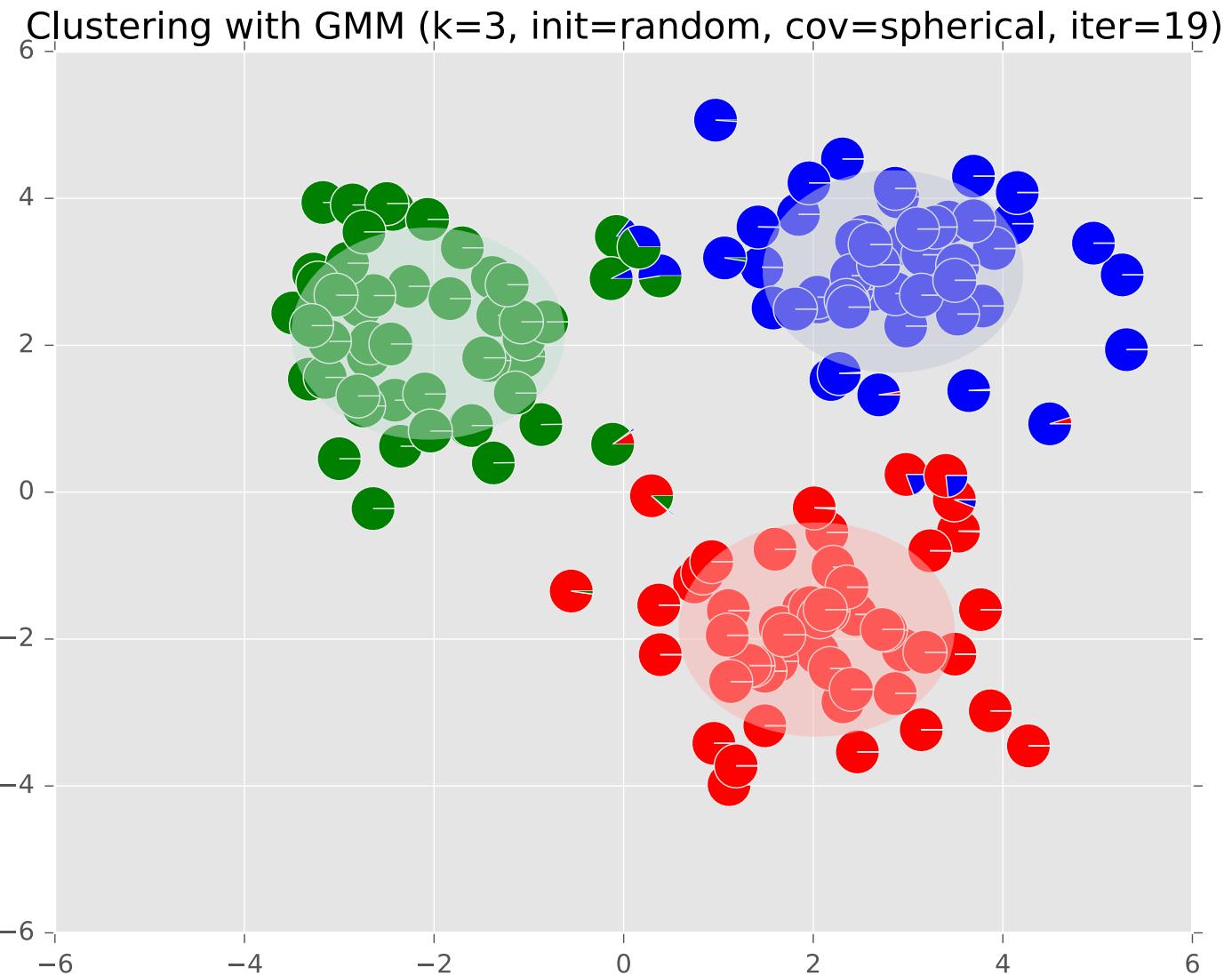
Example: GMM



Example: GMM



Example: GMM



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

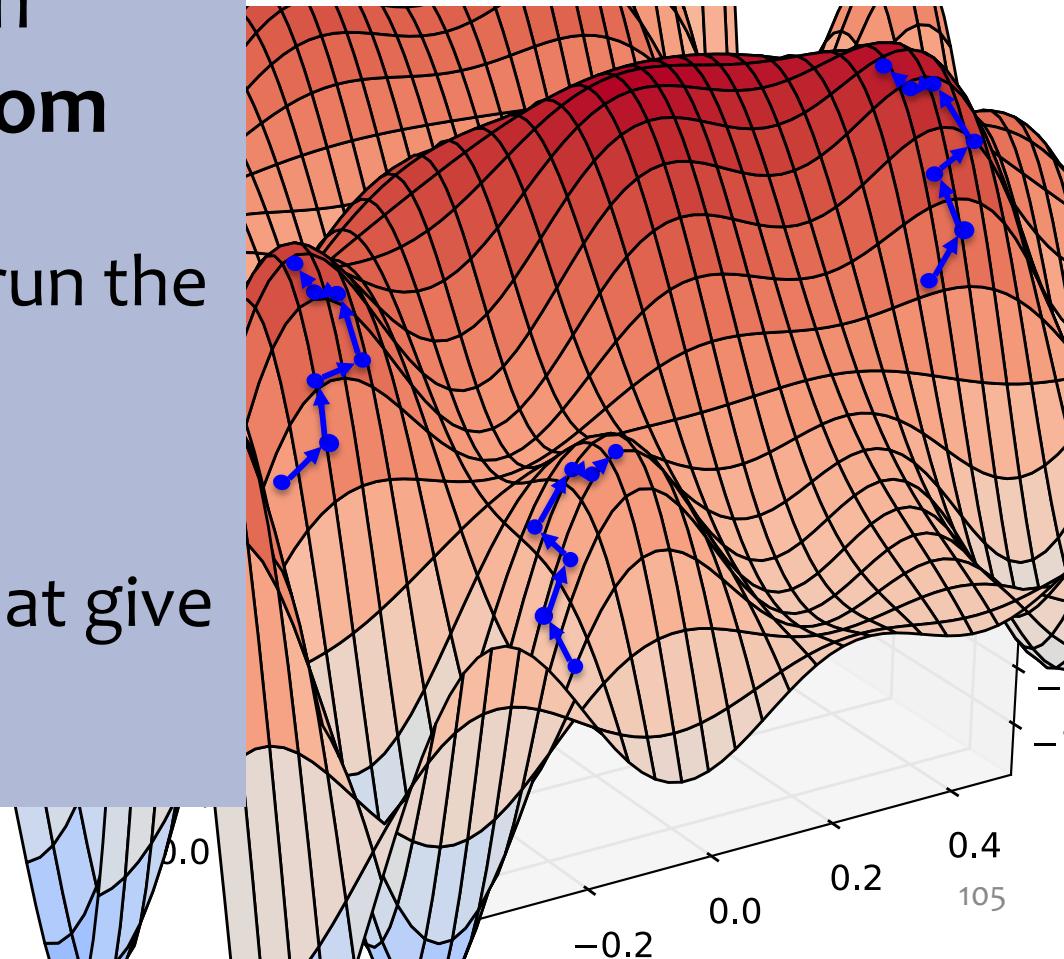
Output:

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

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:

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

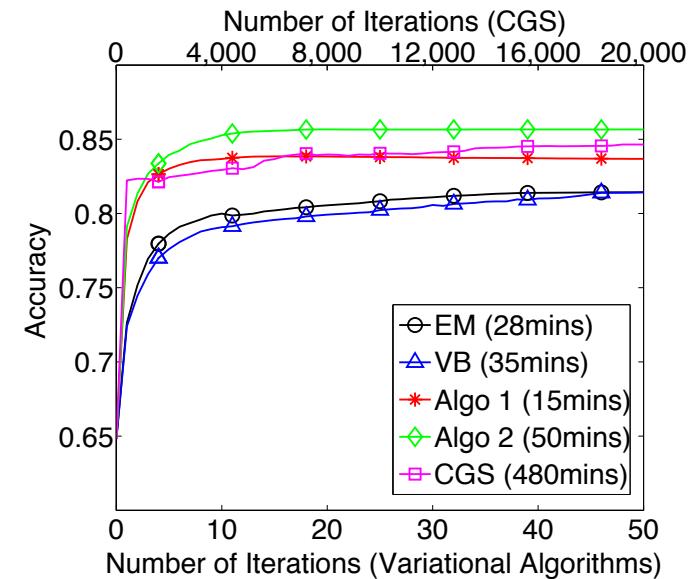
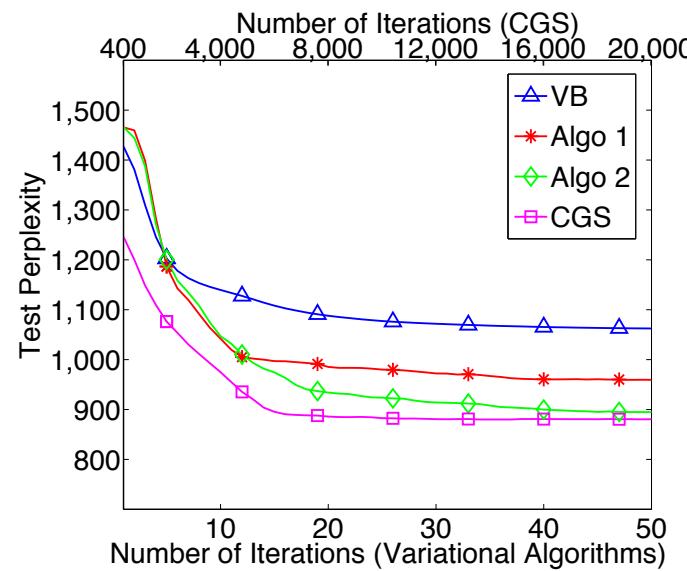
CGS full conditional:

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

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



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**

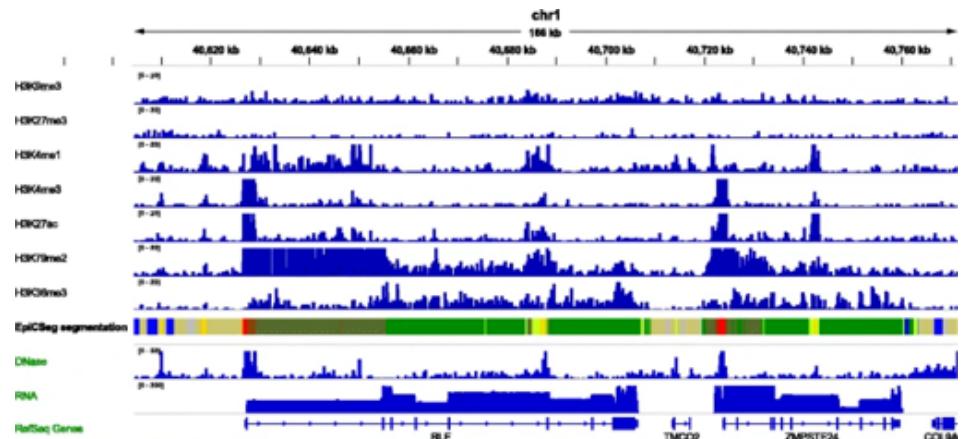


Figure from Foti et al. (2014)

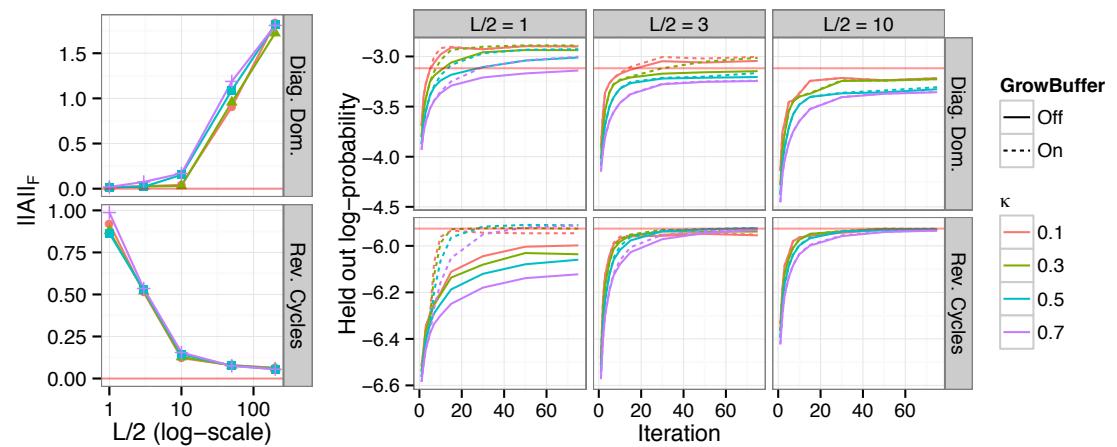


Figure from Mammana & Chung (2015)

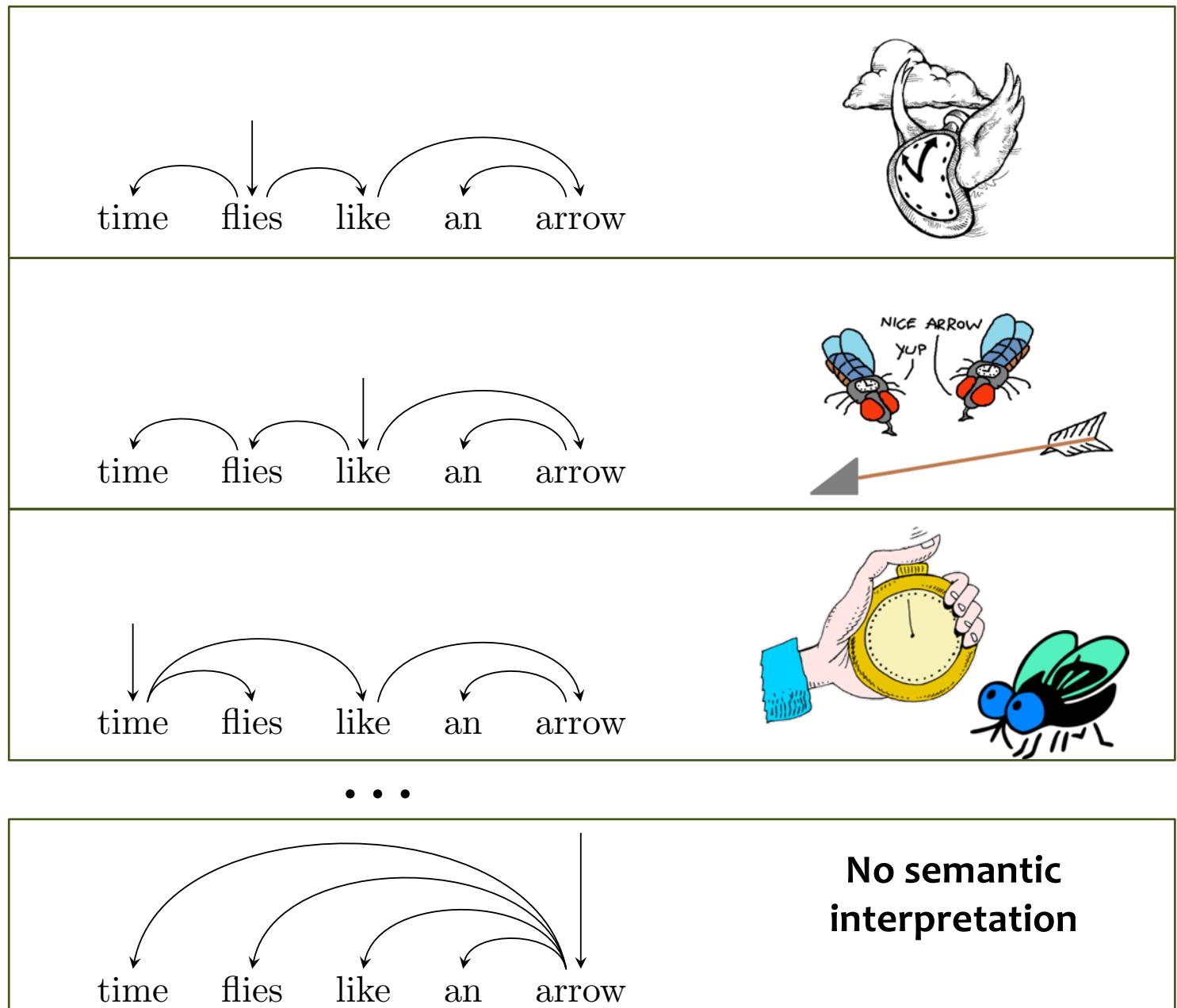
Grammar Induction

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



Grammar Induction

Training Data: Sentences only, without parses

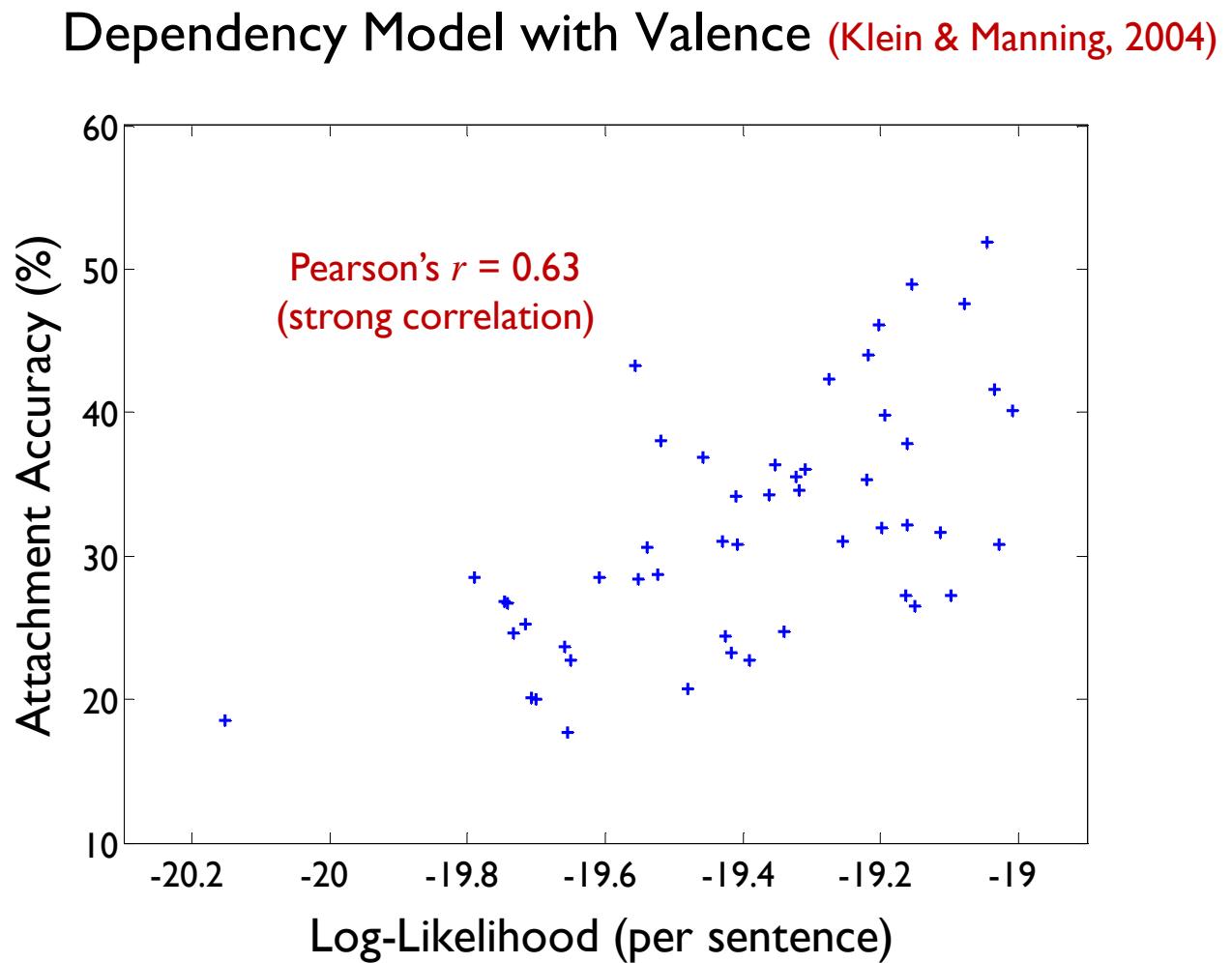
Sample 1:	time	flies	like	an	arrow	$x^{(1)}$
Sample 2:	real	flies	like	soup		$x^{(2)}$
Sample 3:	flies	fly	with	their	wings	$x^{(3)}$
Sample 4:	with	time	you	will	see	$x^{(4)}$

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

Grammar Induction

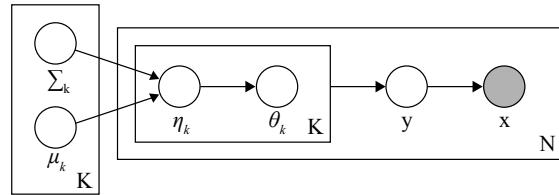
Q: Does likelihood correlate with accuracy on a task we care about?

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 = syntactic parse
 x = observed sentence

Settings:

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

EM-MAP Maximum *a posteriori* estimate of θ using the EM algorithm and a fixed symmetric Dirichlet prior with $\alpha > 1$ to optimize $p(\mathbf{x}, \theta | \alpha)$. Tune α 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 α 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 θ .

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

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

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	45.9	39.9
VB-EM-Log-Normal, families	59.3	45.1	39.0	59.4	45.9	40.5

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].