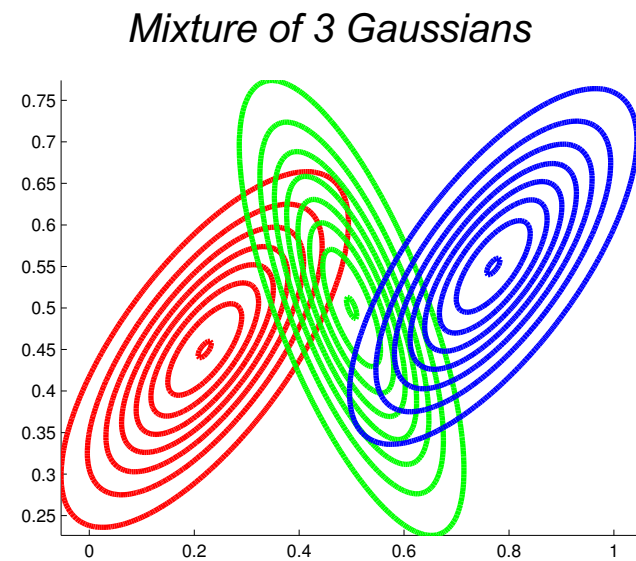
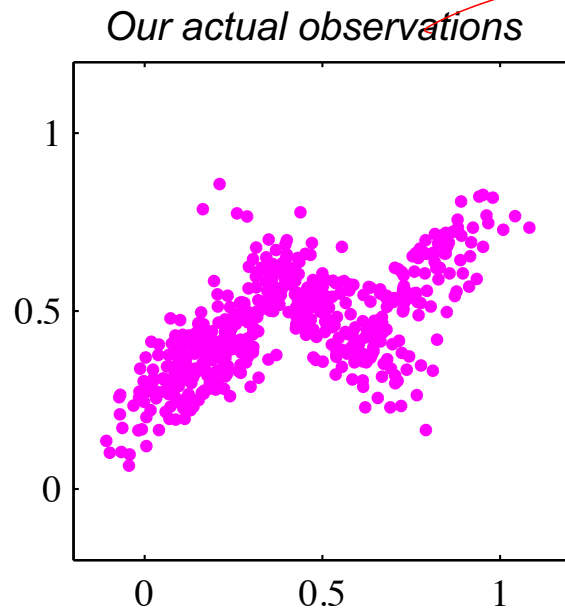


Expectation Maximization for Mixtures of Gaussians

CS229: Machine Learning
Carlos Guestrin
Stanford University

Learning a Mixture of Gaussians



Summary of GMM Components

- Observations $x^i \in \mathbb{R}^d, \quad i = 1, 2, \dots, N$
- Hidden cluster labels $z_i \in \{1, 2, \dots, K\}, \quad i = 1, 2, \dots, N$
- Hidden mixture means $\mu_k \in \mathbb{R}^d, \quad k = 1, 2, \dots, K$
- Hidden mixture covariances $\Sigma_k \in \mathbb{R}^{d \times d}, \quad k = 1, 2, \dots, K$
- Hidden mixture probabilities $\pi_k, \quad \sum_{k=1}^K \pi_k = 1$

Gaussian mixture marginal and conditional likelihood :

$$p(x^i | \pi, \mu, \Sigma) = \sum_{z^i=1}^K \pi_{z^i} p(x^i | z^i, \mu, \Sigma)$$

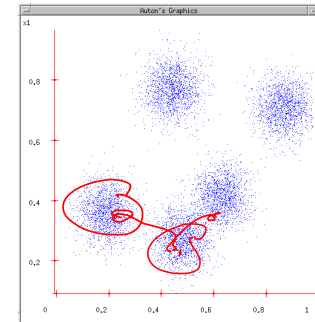
learn from unlabeled data

$$p(x^i | z^i, \mu, \Sigma) = \mathcal{N}(x^i | \mu_{z^i}, \Sigma_{z^i})$$

$P(z=i)$

But we don't see class labels!!!

- MLE:
 - $\text{argmax} \prod_i P(z^i, x^i)$



latent variable

- But we don't know z^i
- Maximize marginal likelihood:
 - $\text{argmax} \prod_i P(x^i) = \text{argmax} \prod_i \sum_k P(z^i=k, x^i)$

μ, Σ, Π

part of f_{θ}
observed data

$\log \Pi \equiv \sum \log$

$\max f(x) = \max \log f(x) \mid \max -f \equiv \min f$

Special case: spherical Gaussians and hard assignments

$\Sigma = \begin{pmatrix} \sigma^2 & 0 \\ 0 & \dots & \sigma^2 \end{pmatrix}$

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

$N(\mu_k, \Sigma_k)$

cluster probability

same cluster probabilities

- If $P(X|z=k)$ is spherical, with same σ for all classes:

$P(z^i = k) = P(z^i = j)$

$$P(\mathbf{x}^i | z^i = k) \propto \exp\left[-\frac{1}{2\sigma^2} \|\mathbf{x}^i - \mu_k\|^2\right]$$

K-means objective

- If each x^i belongs to one class $C(i)$ (hard assignment), marginal likelihood:

$$\min_{\mu} \sum_{i=1}^N \|\mathbf{x}^i - \mu_{C(i)}\|^2$$

want to max

$$\prod_{i=1}^N \sum_{k=1}^K P(\mathbf{x}^i, z^i = k) \propto \prod_{i=1}^N \exp\left[-\frac{1}{2\sigma^2} \|\mathbf{x}^i - \mu_{C(i)}\|^2\right]$$

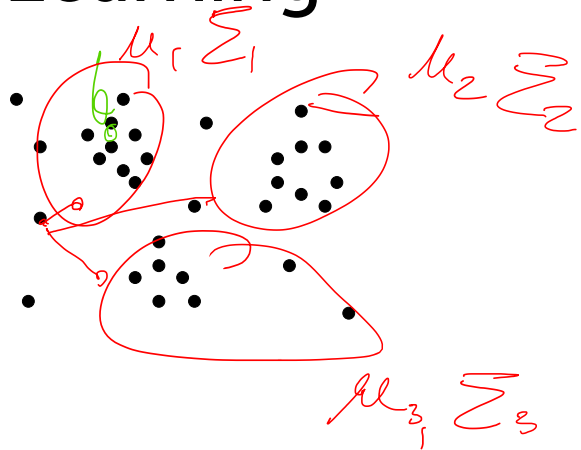
- Same as K-means!!!

$\equiv \max_i \log \Pi e^{\dots}$

$$\max_i \sum_{i=1}^N -\frac{1}{2\sigma^2} \|\mathbf{x}^i - \mu_{C(i)}\|^2$$

EM: "Reducing" Unsupervised Learning to Supervised Learning

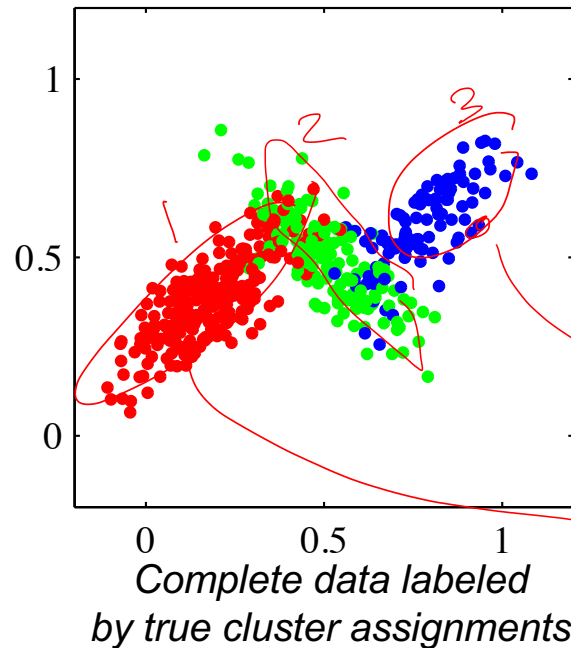
- If we knew assignment of points to classes \rightarrow Supervised Learning!



- Expectation-Maximization (EM)
 - Expectation: Guess assignment of points to classes
 - In standard ("soft") EM: each point associated with prob. of being in each class
 - ~~Maximization: Recompute model parameters~~
 - Iterate

as if ~~we~~ we had fully observed data

Imagine we have an assignment of each x^i to a Gaussian



- Introduce latent cluster indicator variable z^i

$$z^i \in \{1, \dots, k\}$$

$$P(z^i = k) = \pi_k$$

- Then we have

$$p(x^i | z^i, \pi, \mu, \Sigma) = N(x^i | \mu_{z^i}, \Sigma_{z^i})$$

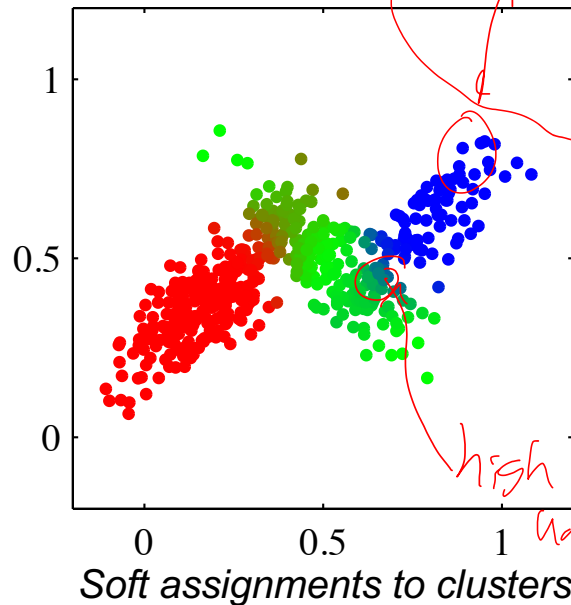
$N(\mu_1, \Sigma_1)$

$N(\mu_3, \Sigma_3)$

Expectation: infer cluster assignments from observations

Normalization:

$$\sum_{k=1}^K r_{ik} = 1$$



- Posterior probabilities of assignments to each cluster *given* model parameters: π, μ, Σ

$$r_{ik} = p(z^i = k | x^i, \pi, \mu, \Sigma) = \frac{P(z^i = k) P(x^i | z^i = k)}{P(x^i)} \leftarrow \text{normalize}$$

$$= \frac{\pi_k N(x^i | \mu_k, \Sigma_k)}{P(x^i)}$$

ML Estimate of Mixture Model Params

- Log likelihood

$$L_x(\theta) \triangleq \log p(\{x^i\} | \theta) = \sum_{i=1}^N \log \sum_{j=1}^K p(x^i, z = j | \theta)$$

$\{\pi, \mu, \Sigma\}$

- Want ML estimate

$$\hat{\theta}^{ML} = \underset{\theta}{\operatorname{argmax}} L_x(\theta)$$

- ~~Neither convex nor concave and local optima~~

$$\theta = \{\pi, \mu, \Sigma\}$$

Maximization: If "complete" data were observed...

Solve as two learning problems

- Assume class labels z^i were observed in addition to x^i

MLE:

$$\hat{\pi}_j = \frac{\text{count}(z^i = j)}{N}$$

$$L_{x,z}(\theta) = \sum_{i=1}^N \log p(x^i, z^i | \theta) = \sum_{i=1}^N \log P(z^i | \theta) P(x^i | z^i, \theta)$$

$$= \sum_{i=1}^N \log P(z^i | \pi) + \sum_{i=1}^N \log P(x^i | z^i, \mu, \Sigma)$$

- Compute ML estimates
 - Separates over clusters $k!$

K indep. learning problems

$$\sum_{j=1}^K \sum_{i: z^i = j} \log P(x^i | \mu_j, \Sigma_j)$$

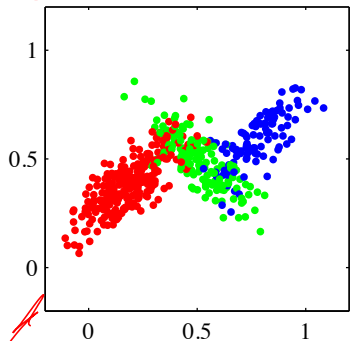
- Example: mixture of Gaussians (MoG)

$$\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$$

MLE $\hat{\mu}_j = \text{mean of } x^i \text{ assigned to cluster } j$

Maximization: if inferred cluster assignments from observations

fully observed data



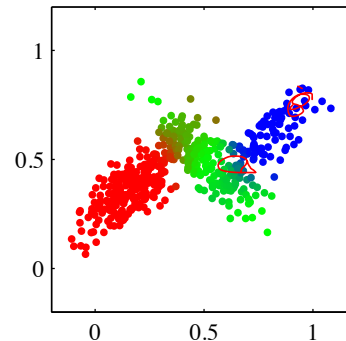
Complete data labeled by true cluster assignments

estimate π, μ, Σ using simple counts & averages

e.g.

$$\hat{\pi}_j = \text{weighted count of fraction of points in cluster } j = \frac{\sum_{i=1}^N r_{ij}}{N}$$

$$\hat{\mu}_j = \text{weighted avg. of } x^i = \frac{\sum_{i=1}^N r_{ij} x^i}{N}$$



Soft assignments to clusters

- Posterior probabilities of assignments to each cluster *given* model parameters:

$$r_{ik} = p(z^i = k | x^i, \pi, \mu, \Sigma)$$

weighted counts!!

Expectation-Maximization Algorithm

- Motivates a coordinate ascent-like algorithm:

- E 1. Infer missing values z^i given estimate of parameters $\hat{\theta}$
- M 2. Optimize parameters to produce new $\hat{\theta}$ given “filled in” data z^i
- 3. Repeat

- Example: MoG

- 1. Infer “responsibilities”

$$\underline{r_{ik} = p(z^i = k \mid x^i, \hat{\theta}^{(t-1)})} \leftarrow \text{Bayes rule}$$

- 2. Optimize parameters

max w.r.t. π_k :

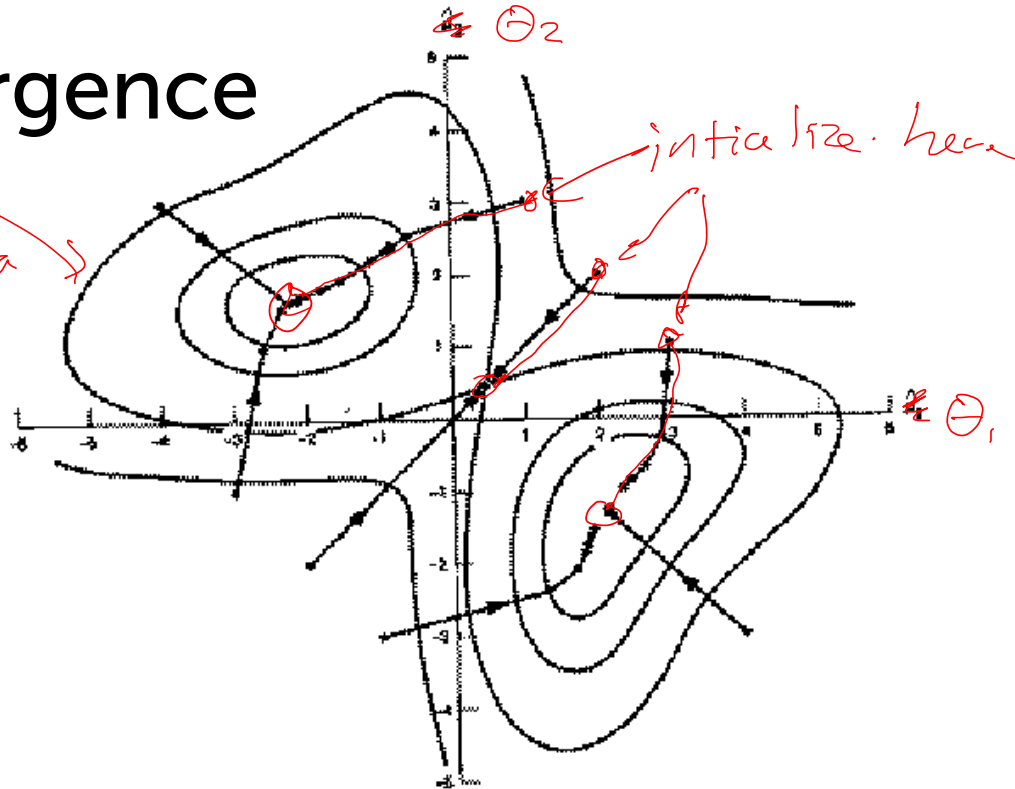
max w.r.t. μ_k, Σ_k :

MLE on weighted data

E.M. Convergence

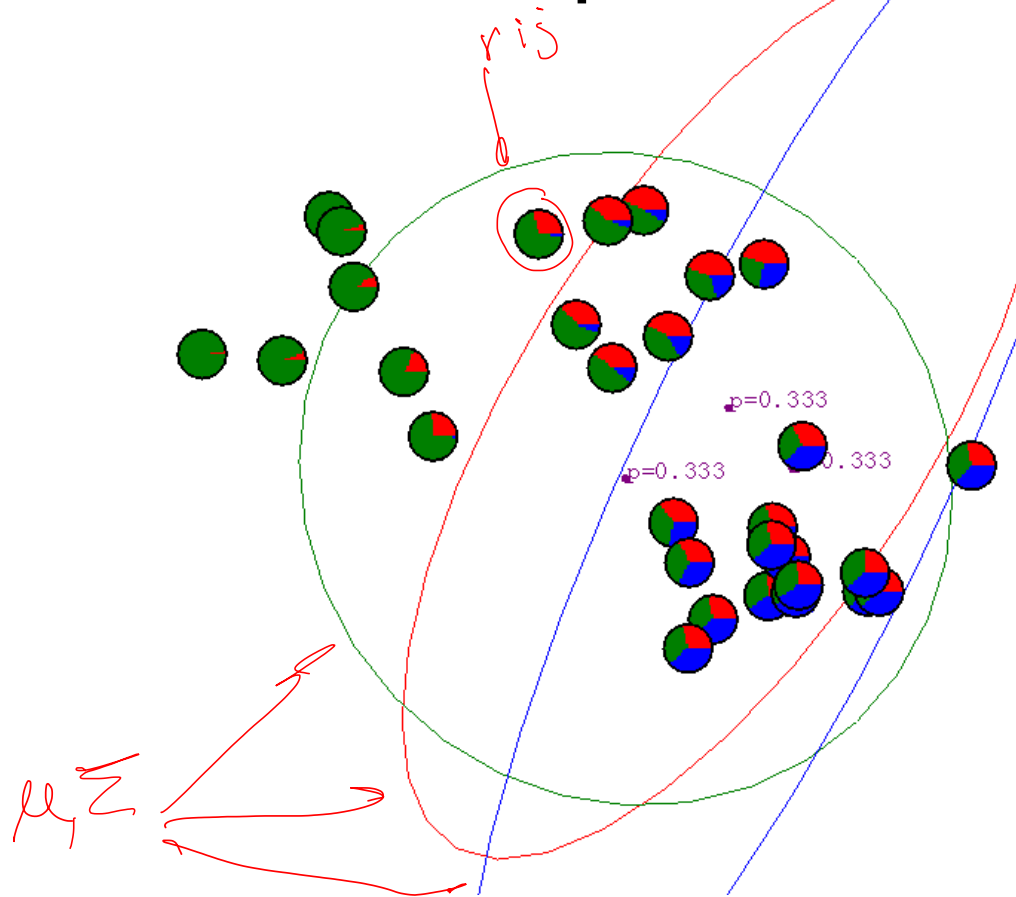
level sets
of loss function
on unlabeled data

- EM is coordinate ascent on an interesting potential function
- Coord. ascent for bounded pot. func. → convergence to a local optimum guaranteed

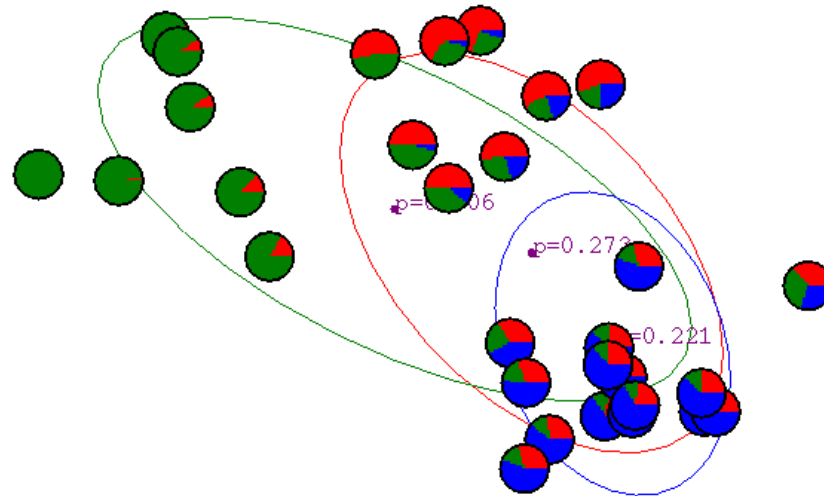


- This algorithm is REALLY USED. And in high dimensional state spaces, too.

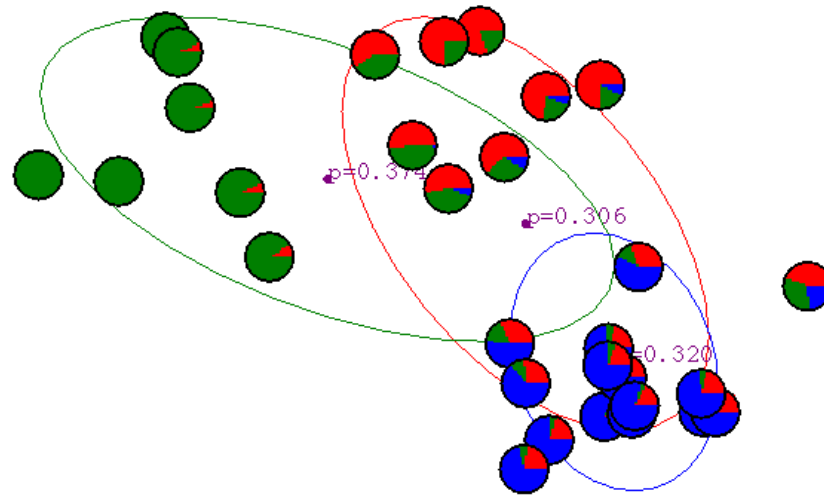
Gaussian Mixture Example: Start



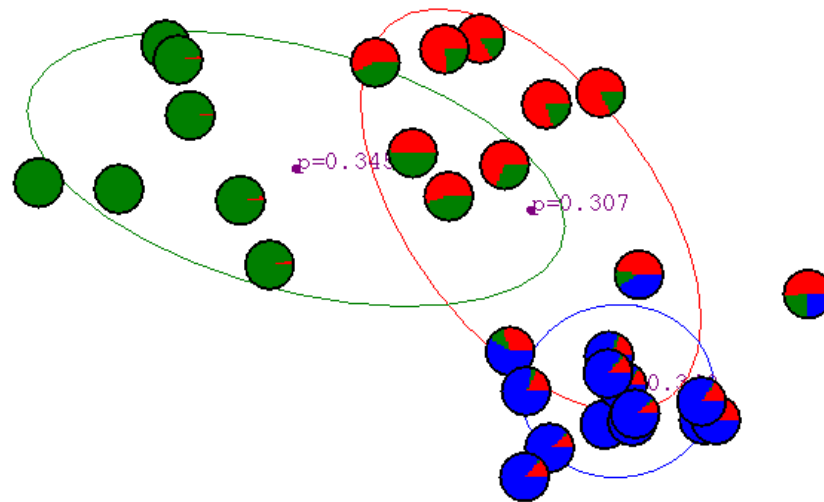
After first iteration



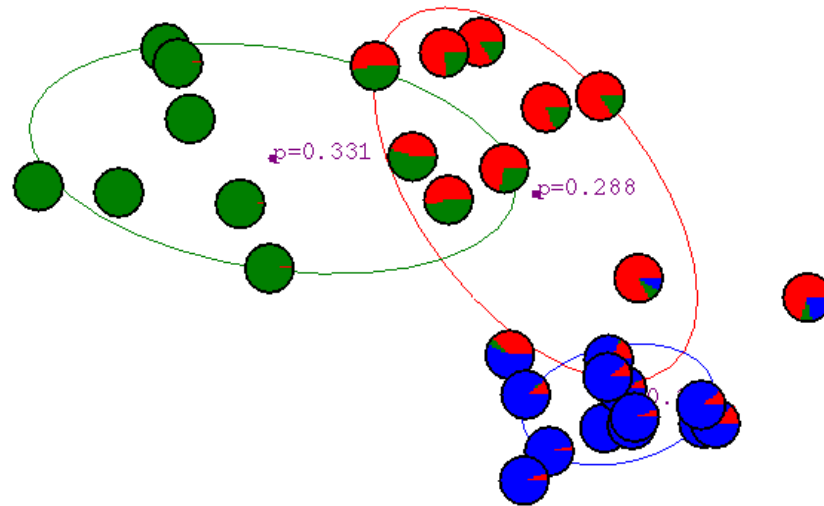
After 2nd iteration



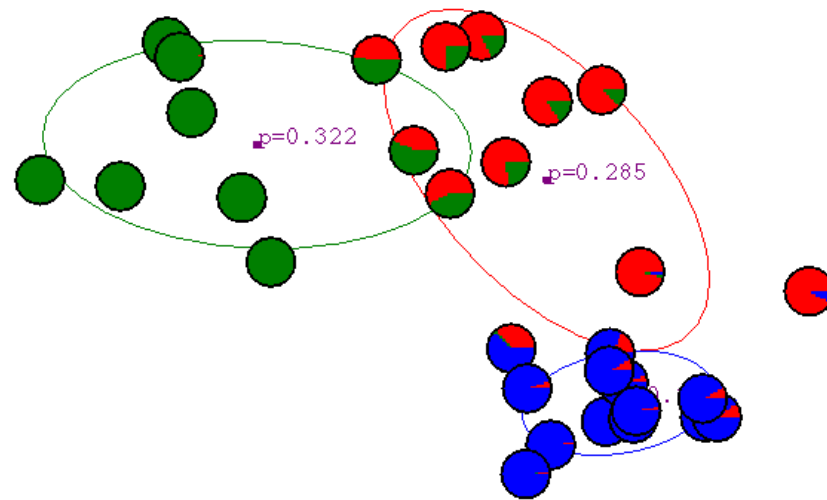
After 3rd iteration



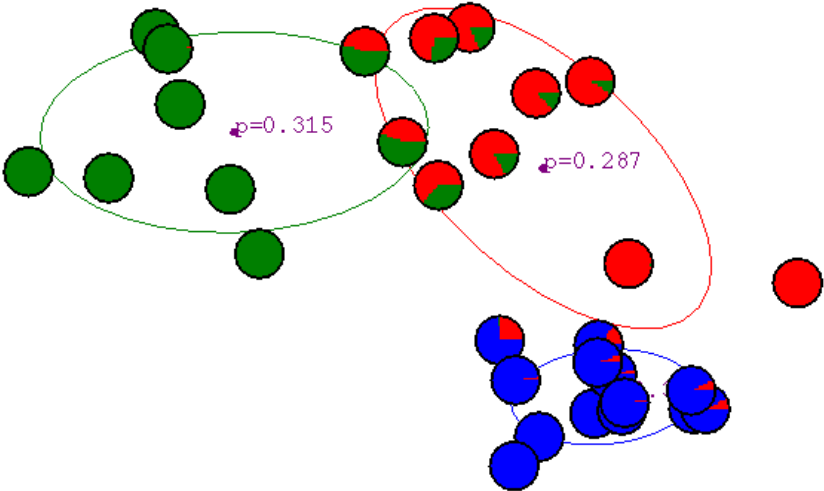
After 4th iteration



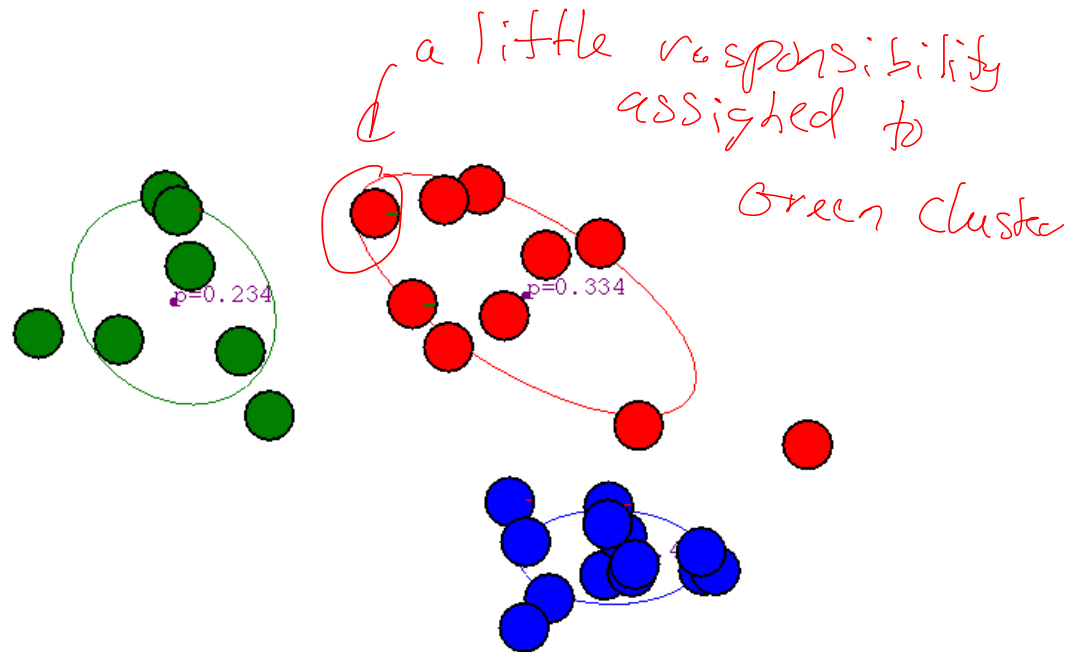
After 5th iteration



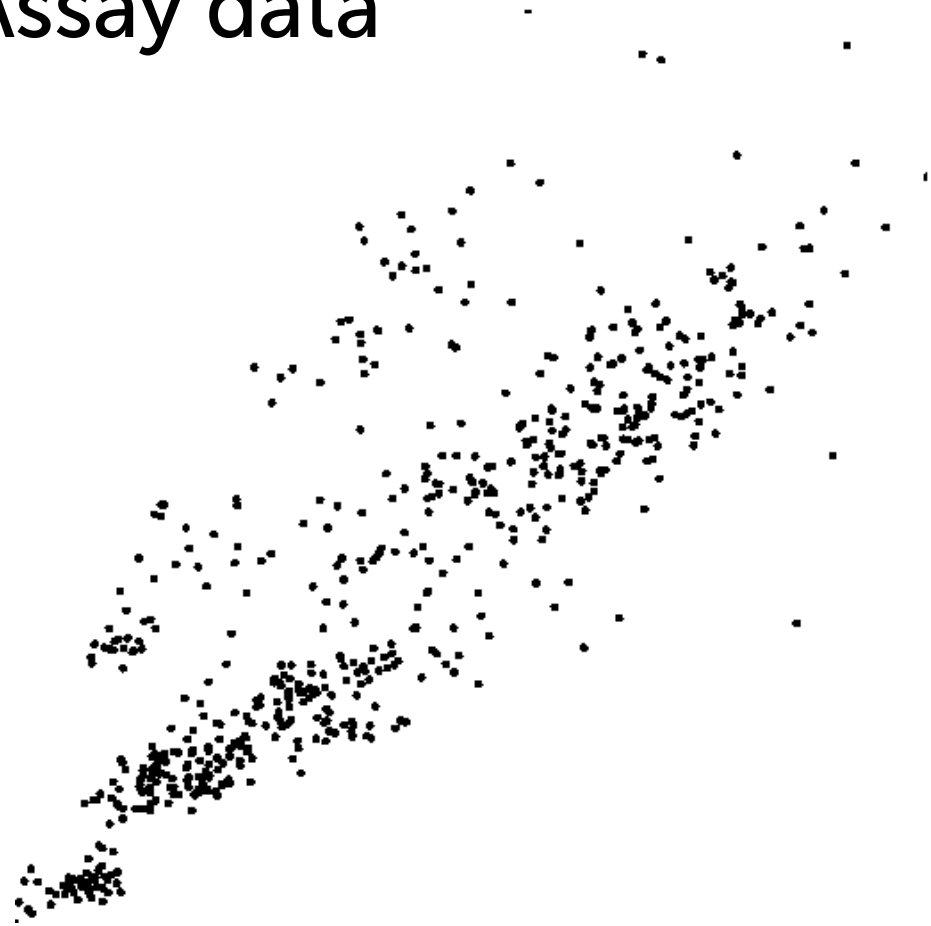
After 6th iteration



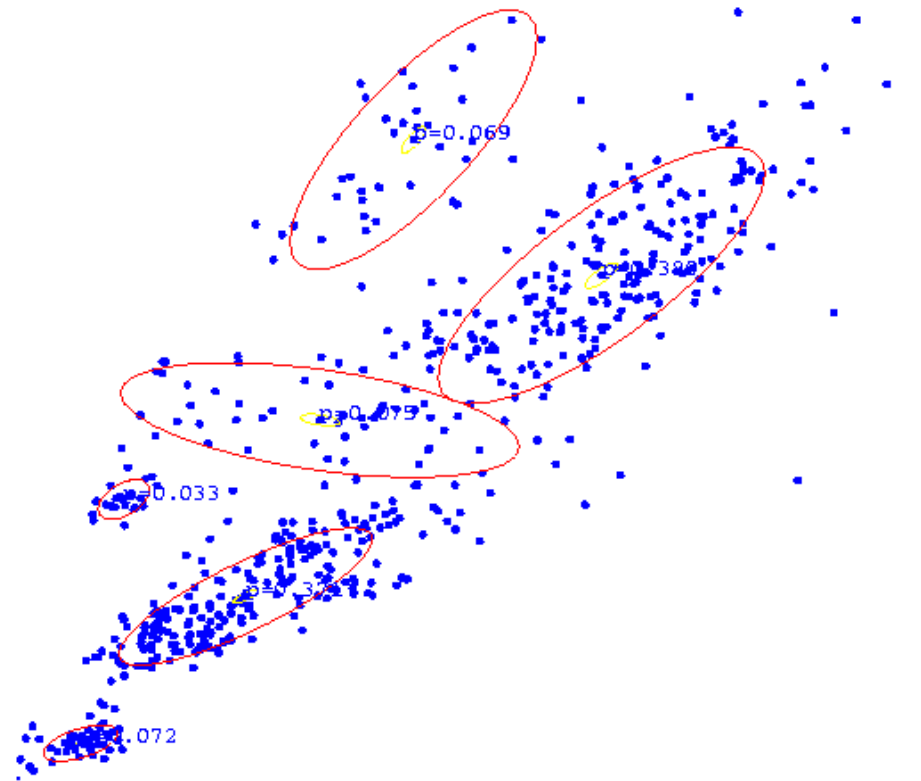
After 20th iteration



Some Bio Assay data



GMM clustering of the assay data



Resulting Density Estimator

