



Unsupervised Learning: K-Means & Gaussian Mixture Models

Unsupervised Learning

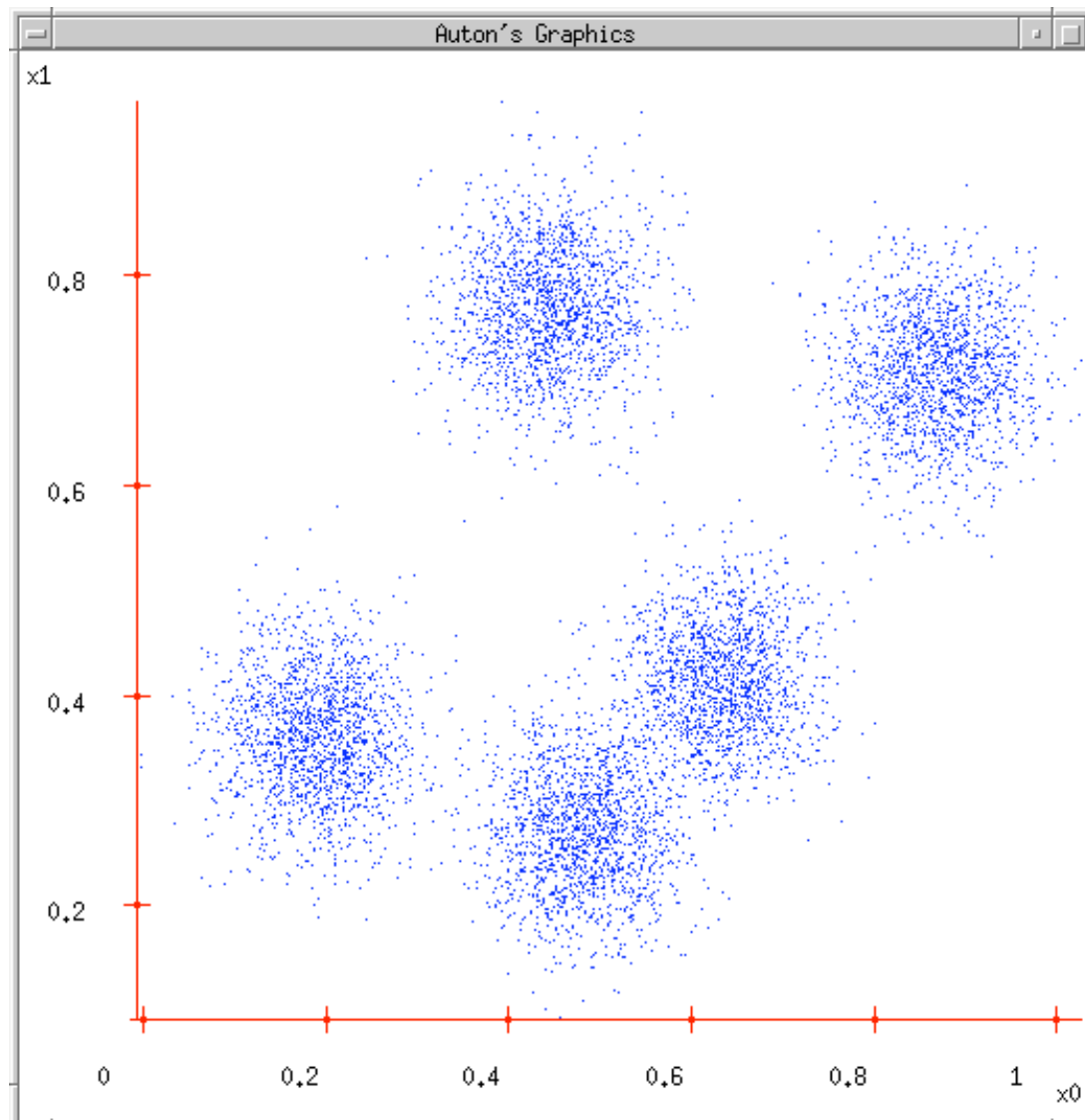
- Supervised learning used labeled data pairs (\mathbf{x}, y) to learn a function $f: X \rightarrow Y$
 - But, what if we don't have labels?
- No labels = **unsupervised learning**
- Only some points are labeled = **semi-supervised learning**
 - Labels may be expensive to obtain, so we only get a few
- **Clustering** is the unsupervised grouping of data points. It can be used for **knowledge discovery**.

K-Means Clustering

Some material adapted from slides by Andrew Moore, CMU.

Visit <http://www.autonlab.org/tutorials/> for
Andrew's repository of Data Mining tutorials.

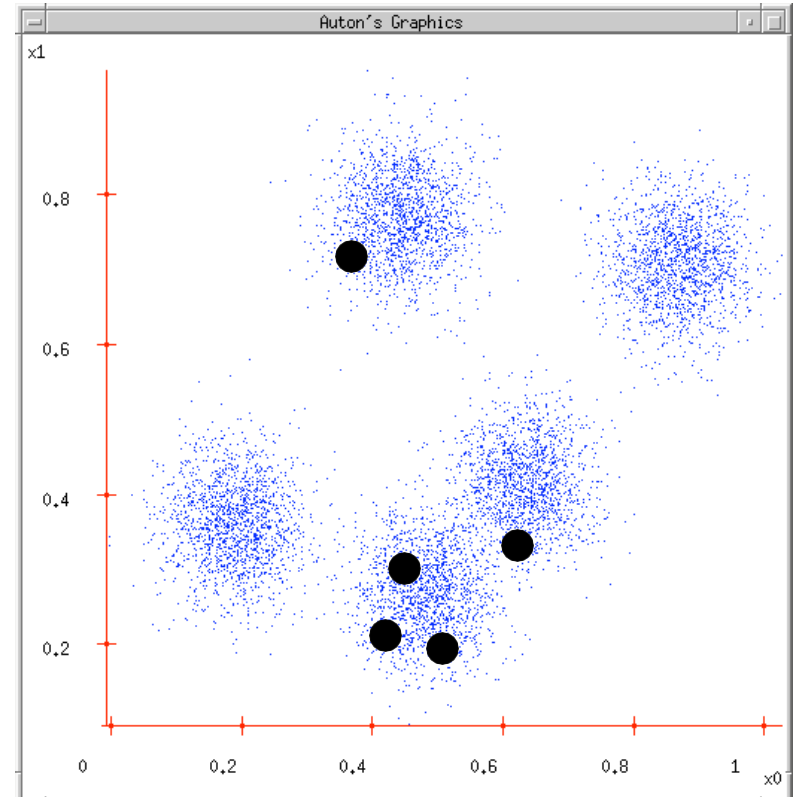
Clustering Data



K-Means Clustering

~~K-Means~~ (k, X)

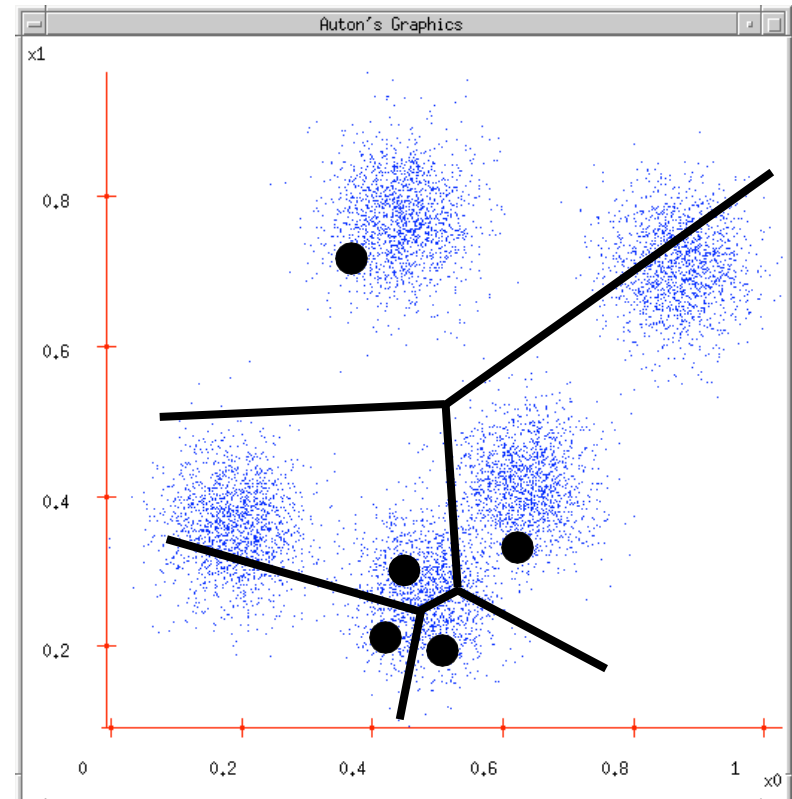
- Randomly choose k cluster center locations (centroids)
- Loop until convergence
 - Assign each point to the cluster of the closest centroid
 - Re-estimate the cluster centroids based on the data assigned to each cluster



K-Means Clustering

~~K-Means~~ $K\text{-Means}(k, X)$

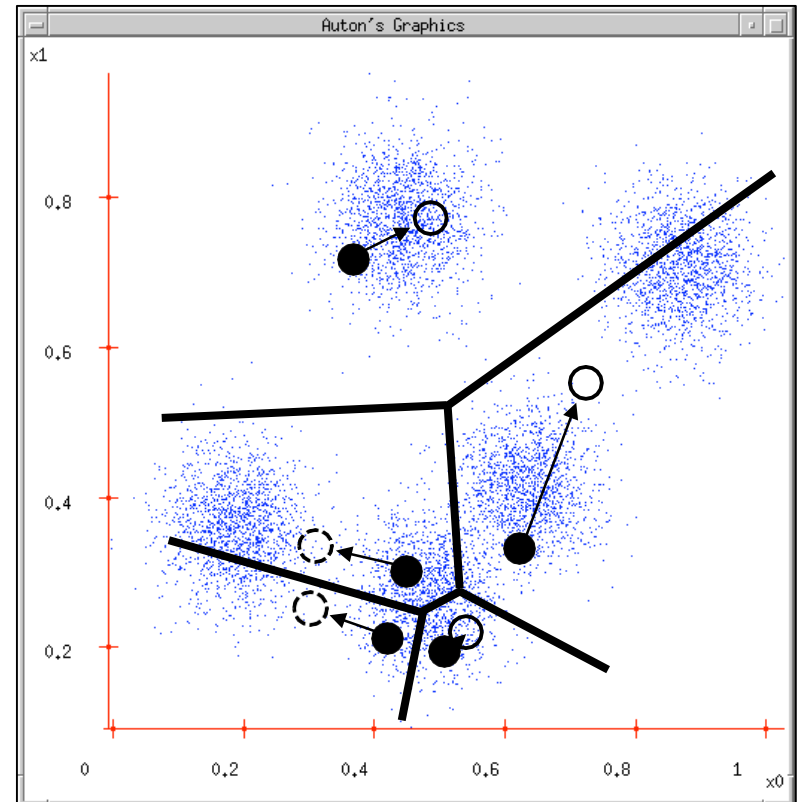
- Randomly choose k cluster center locations (centroids)
- Loop until convergence
 - Assign each point to the cluster of the closest centroid
 - ~~Re-estimate~~ Re-estimate the cluster centroids based on the data assigned to each cluster



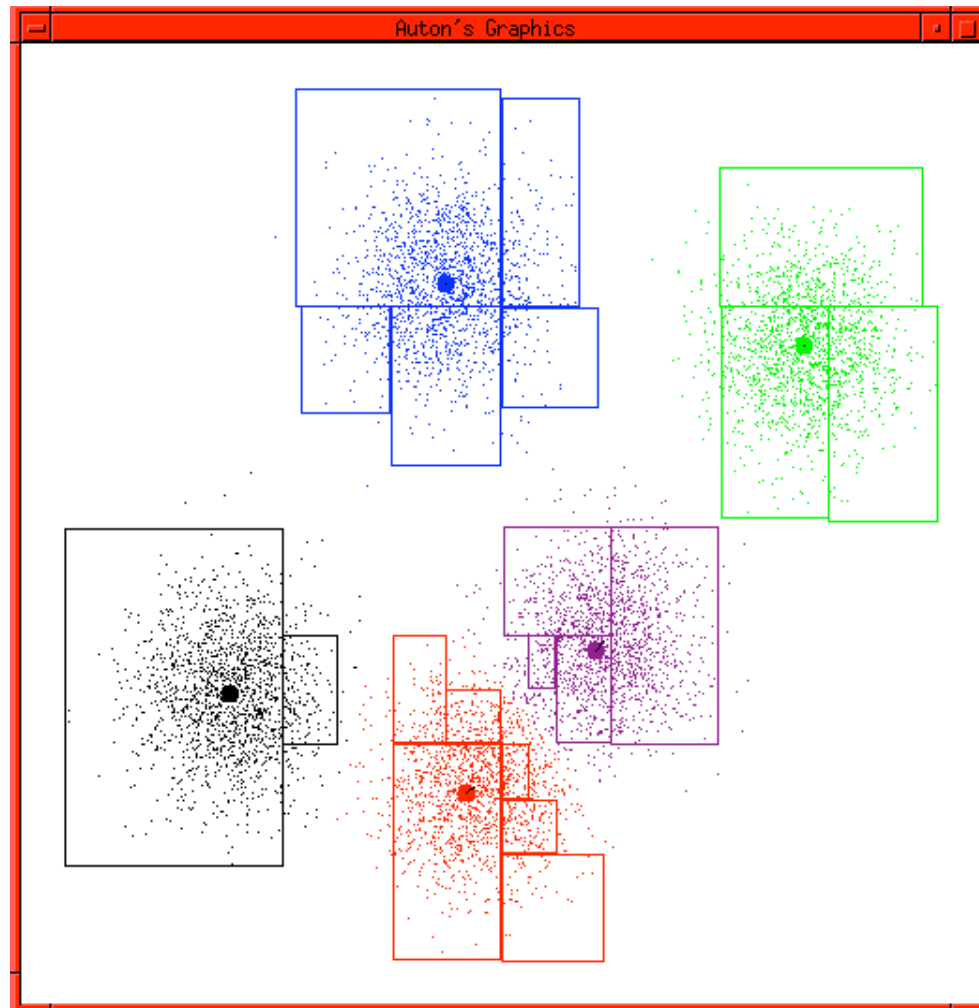
K-Means Clustering

~~K-Means~~ (k, X)

- Randomly choose k cluster center locations (centroids)
- Loop until convergence
 - Assign each point to the cluster of the closest centroid
 - Re-estimate the cluster centroids based on the data assigned to each cluster



K-Means Animation



Example generated by Andrew Moore using Dan Pelleg's super-duper fast K-means system:

Dan Pelleg and Andrew Moore. Accelerating Exact k-means Algorithms with Geometric Reasoning. Proc. Conference on Knowledge Discovery in Databases 1999.

K-Means Objective Function

- K-means finds a local optimum of the following objective function:

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|_2^2$$

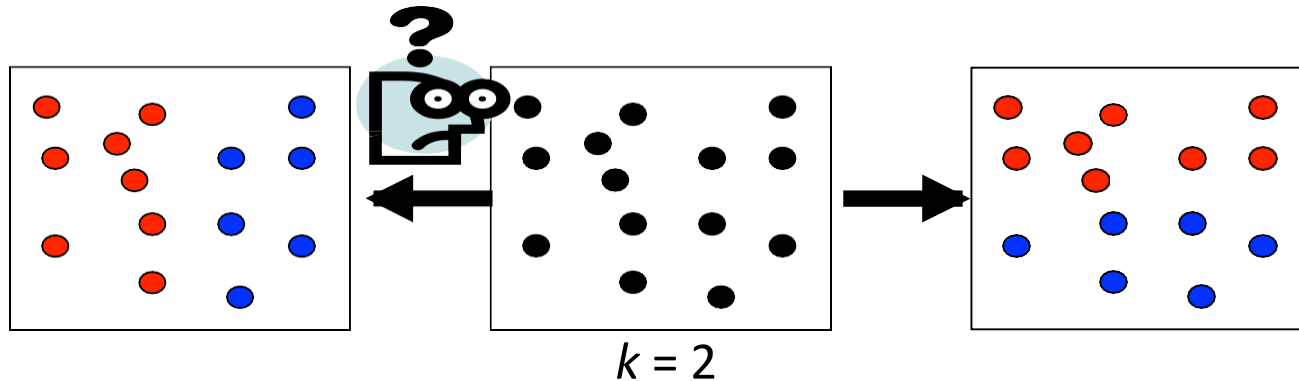
where $\mathbf{S} = \{S_1, \dots, S_k\}$ is a partitioning
 over $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ s.t. $\mathbf{X} = \bigcup_{i=1}^k S_i$
 and $\boldsymbol{\mu}_i = \text{mean}(S_i)$

Problems with K-Means

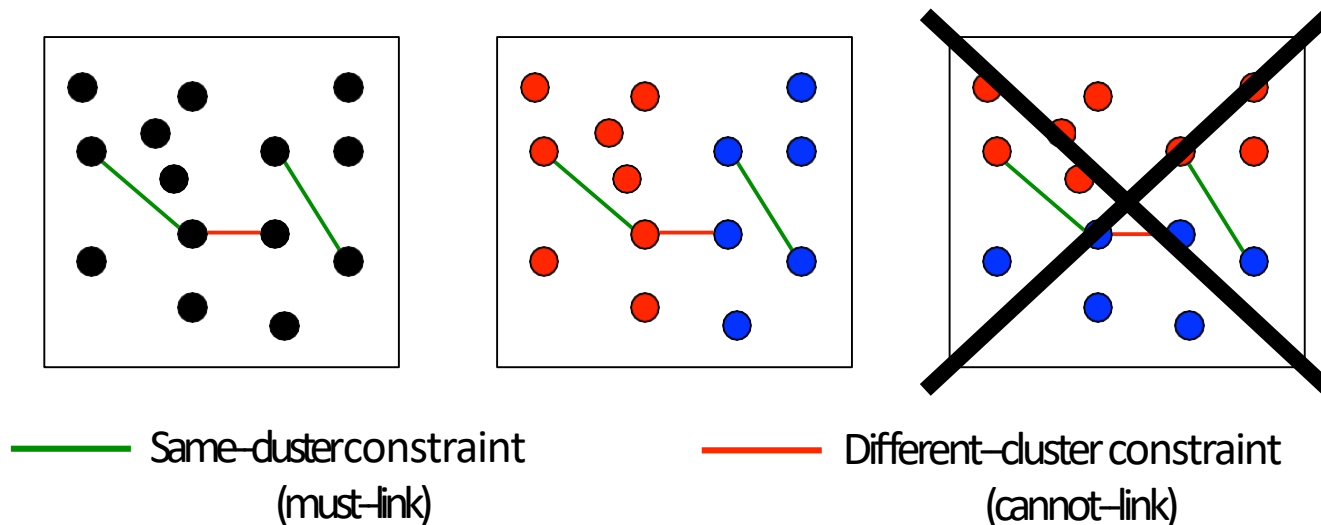
- **Very** sensitive to the initial points
 - Do many runs of K-Means, each with different initial centroids
 - Seed the centroids using a better method than randomly choosing the centroids
 - e.g., Farthest-first sampling
- Must manually choose k
 - Learn the optimal k for the clustering
 - Note that this requires a performance measure

Problems with K-Means

- How do you tell it which clustering you want?



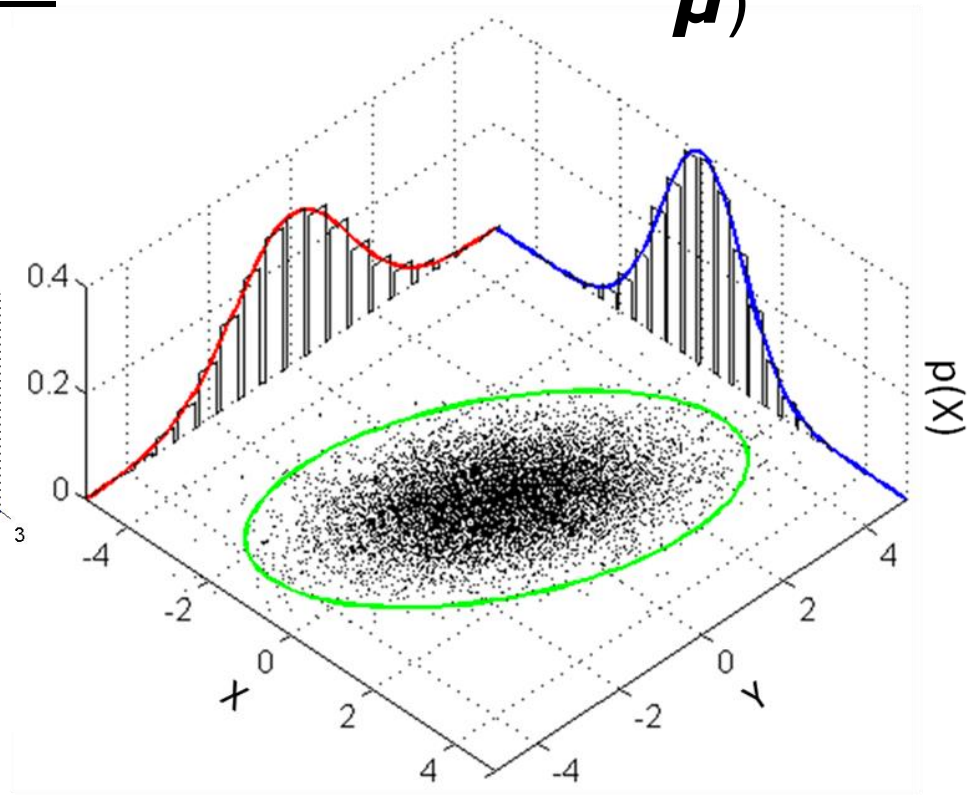
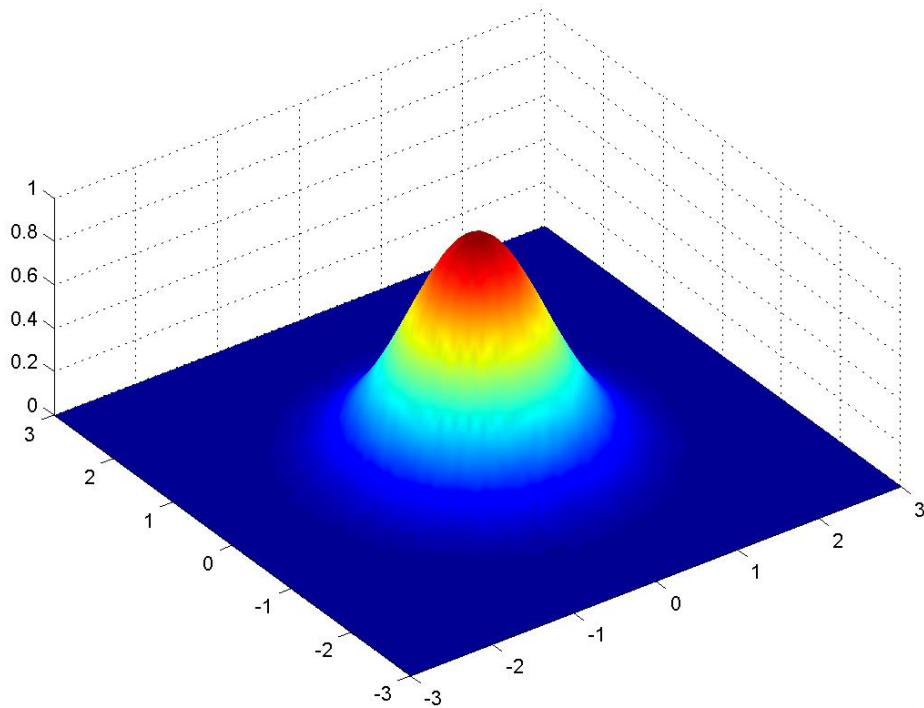
Constrained clustering techniques (semi-supervised)



Gaussian Mixture Models

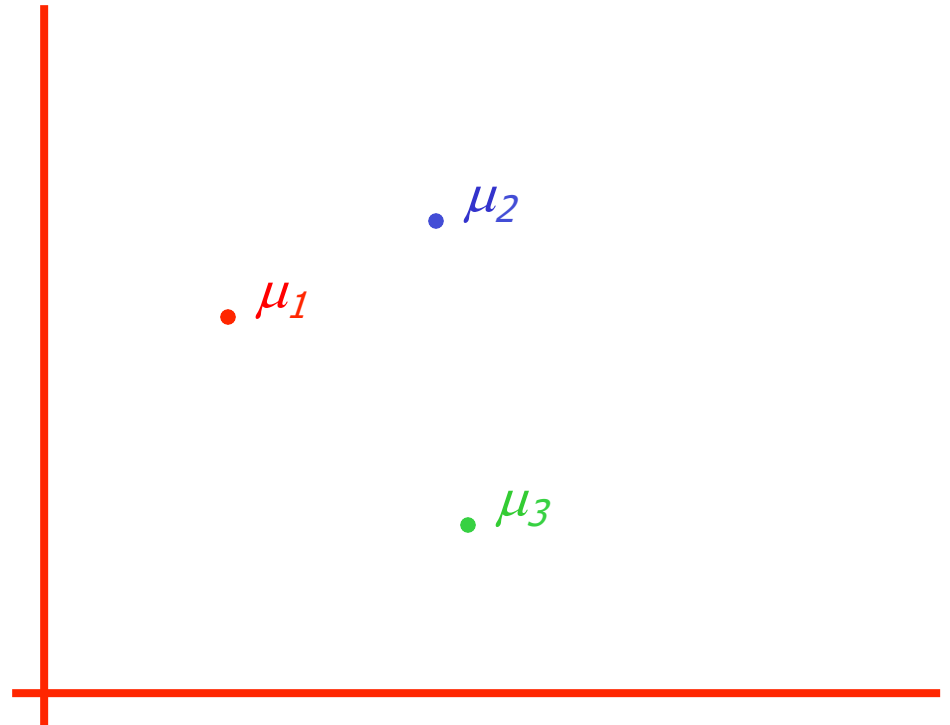
- Recall the Gaussian distribution:

$$P(\mathbf{x} | \boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$



The GMM assumption

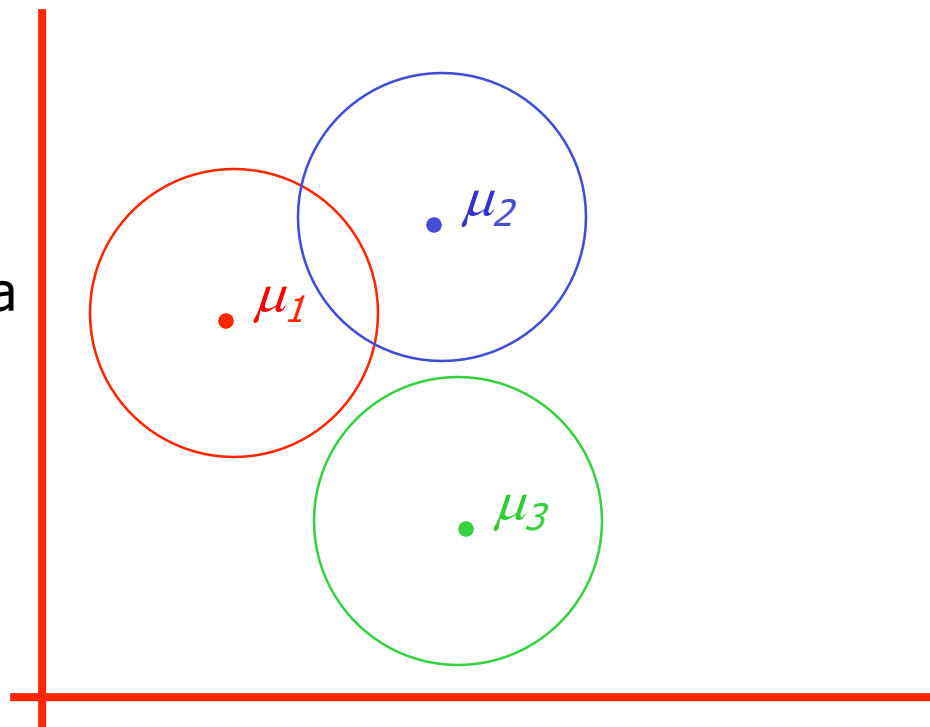
- There are k components. The i th component is called ω_i
- Component ω_i has an associated mean vector μ_i



The GMM assumption

- There are k components. The i th component is called ω_i
- Component ω_i has an associated mean vector μ_i
- Each component generates data from a Gaussian with mean μ_i and covariance matrix $\sigma^2 \mathbf{I}$

Assume that each datapoint is generated according to the following recipe:

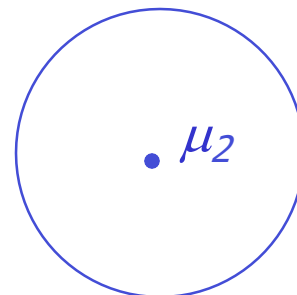


The GMM assumption

- There are k components. The i th component is called ω_i
- Component ω_i has an associated mean vector μ_i
- Each component generates data from a Gaussian with mean μ_i and covariance matrix $\sigma^2 \mathbf{I}$

Assume that each datapoint is generated according to the following recipe:

1. Pick a component at random. Choose component i with probability $P(\omega_i)$.

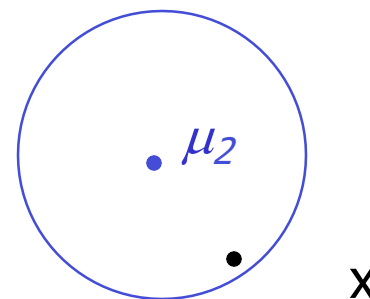


The GMM assumption

- There are k components. The i th component is called ω_i
- Component ω_i has an associated mean vector μ_i
- Each component generates data from a Gaussian with mean μ_i and covariance matrix $\sigma^2 \mathbf{I}$

Assume that each datapoint is generated according to the following recipe:

1. Pick a component at random. Choose component i with probability $P(\omega_i)$.
2. Datapoint $\sim N(\mu_i, \sigma^2 \mathbf{I})$

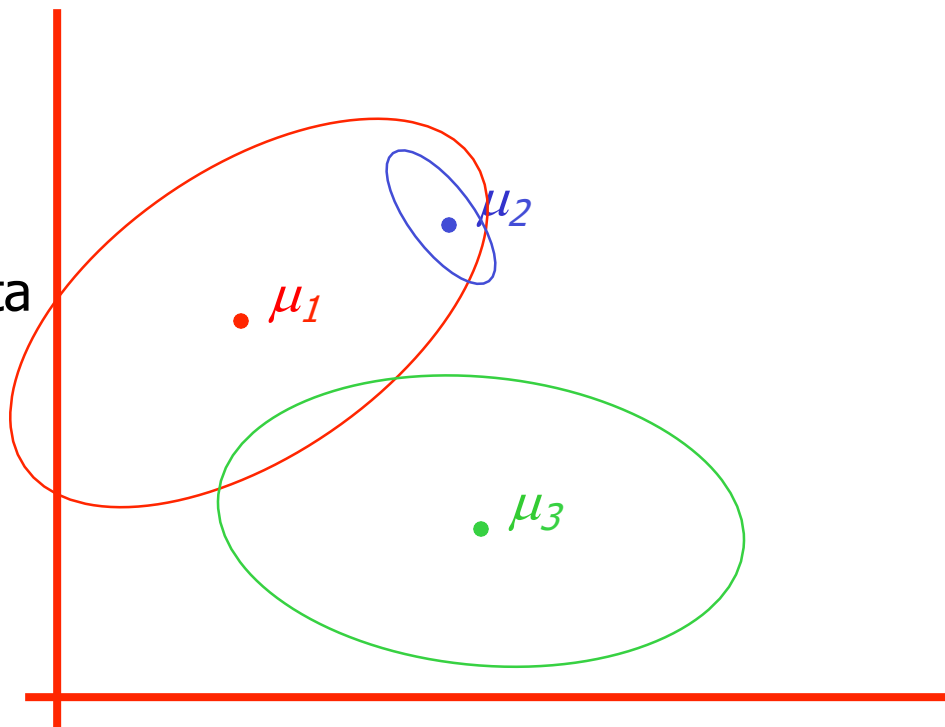


The **General** GMM assumption

- There are k components. The i th component is called ω_i
- Component ω_i has an associated mean vector μ_i
- Each component generates data from a Gaussian with mean μ_i and covariance matrix Σ_i

Assume that each datapoint is generated according to the following recipe:

1. Pick a component at random. Choose component i with probability $P(\omega_i)$.
2. Datapoint $\sim N(\mu_i, \Sigma_i)$



Fitting a Gaussian Mixture Model

(Optional)

Expectation-Maximization for GMMs

Iterate until convergence:

On the t th iteration let our estimates be

$$\lambda_t = \{ \mu_1(t), \mu_2(t) \dots \mu_c(t) \}$$

*Just evaluate a
Gaussian at x_k*

E-step: Compute “expected” classes of all datapoints for each class

$$P(w_i | x_k, \lambda_t) = \frac{p(x_k | w_i, \lambda_t) P(w_i | \lambda_t)}{p(x_k | \lambda_t)} = \frac{p(x_k | w_i, \mu_i(t), \sigma^2 \mathbf{I}) p_i(t)}{\sum_{j=1}^c p(x_k | w_j, \mu_j(t), \sigma^2 \mathbf{I}) p_j(t)}$$

M-step: Estimate μ given our data's class membership distributions

$$\mu_i(t+1) = \frac{\sum_k P(w_i | x_k, \lambda_t) x_k}{\sum_k P(w_i | x_k, \lambda_t)}$$

E.M. for General GMMs

$p_i(t)$ is shorthand
for estimate of
 $P(w_i)$ on t th
iteration

Iterate. On the t th iteration let our estimates be

$$\lambda_t = \{ \mu_1(t), \mu_2(t) \dots \mu_c(t), \Sigma_1(t), \Sigma_2(t) \dots \Sigma_c(t), p_1(t), p_2(t) \dots p_c(t) \}$$

E-step: Compute “expected” clusters of all datapoints

*Just evaluate a
Gaussian at x_k*

$$P(w_i | x_k, \lambda_t) = \frac{p(x_k | w_i, \lambda_t) P(w_i | \lambda_t)}{p(x_k | \lambda_t)} = \frac{p(x_k | w_i, \mu_i(t), \Sigma_i(t)) p_i(t)}{\sum_{j=1}^c p(x_k | w_j, \mu_j(t), \Sigma_j(t)) p_j(t)}$$

M-step: Estimate μ, Σ given our data's class membership distributions

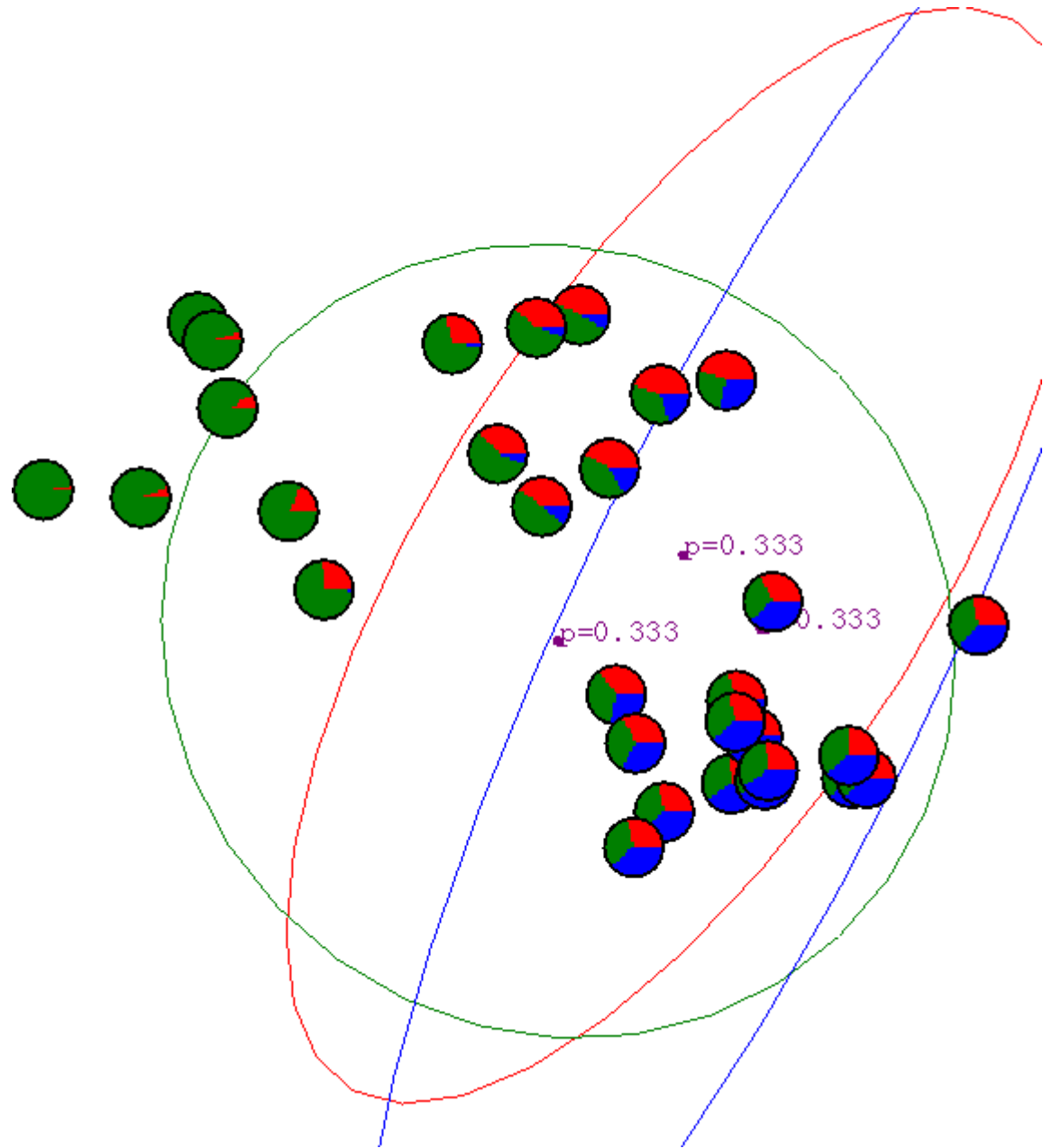
$$\mu_i(t+1) = \frac{\sum_k P(w_i | x_k, \lambda_t) x_k}{\sum_k P(w_i | x_k, \lambda_t)} \quad \Sigma_i(t+1) = \frac{\sum_k P(w_i | x_k, \lambda_t) [x_k - \mu_i(t+1)][x_k - \mu_i(t+1)]^T}{\sum_k P(w_i | x_k, \lambda_t)}$$

$$p_i(t+1) = \frac{\sum_k P(w_i | x_k, \lambda_t)}{R}$$

$R = \# \text{records}$

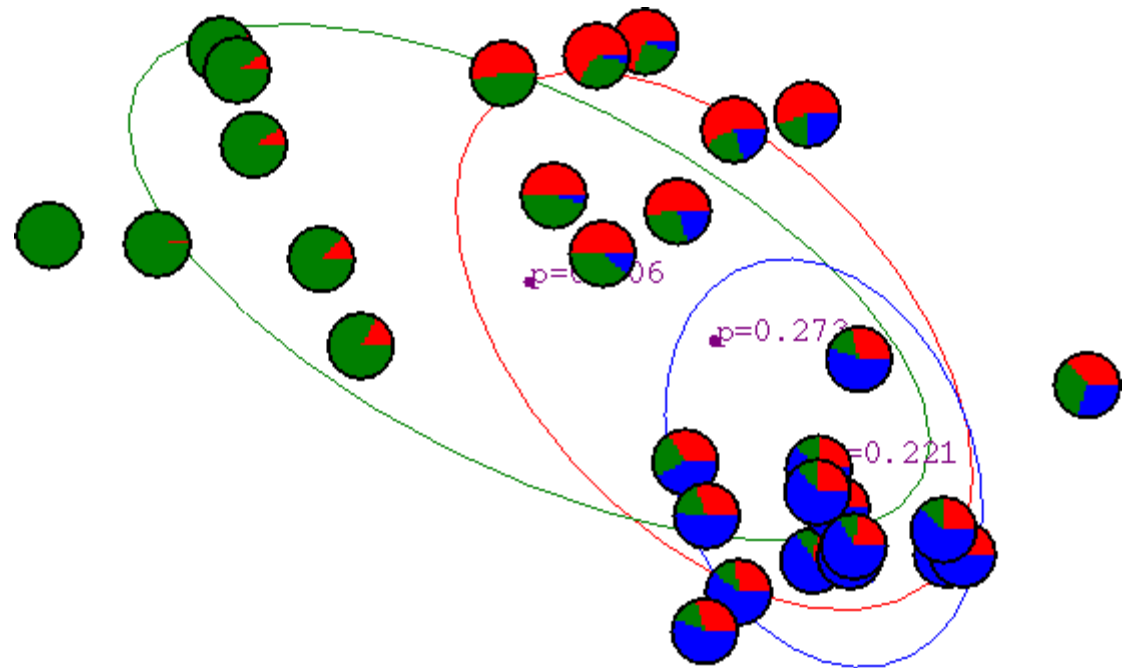
(End optional section)

Gaussian Mixture Example: Start

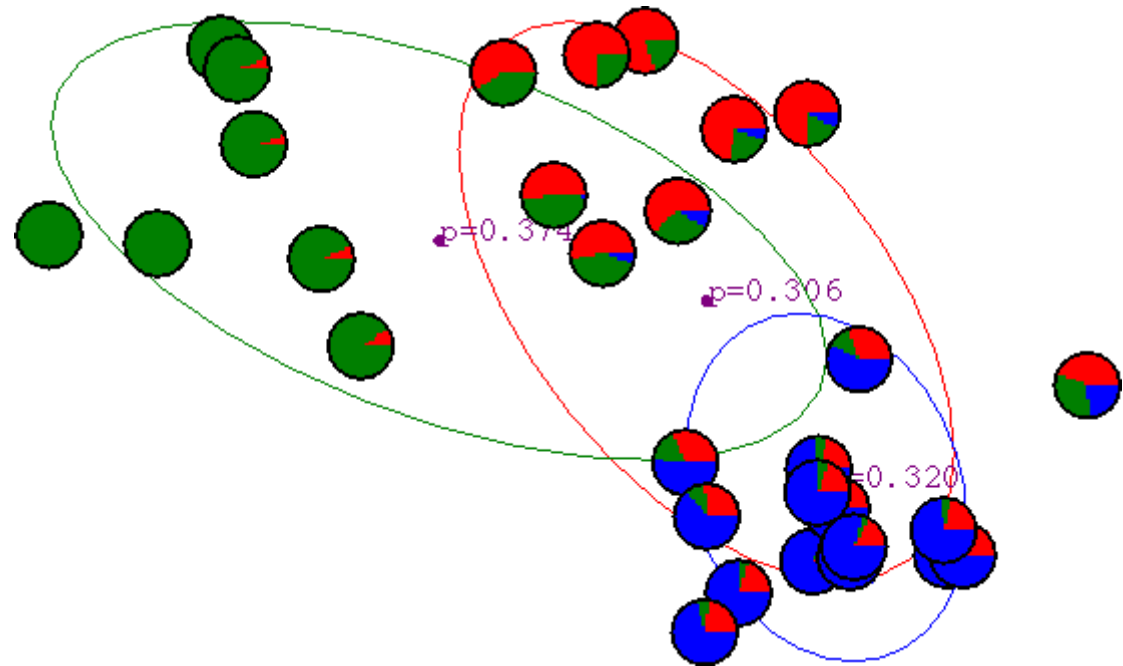


*Advance apologies: in Black
and White this example will be
incomprehensible*

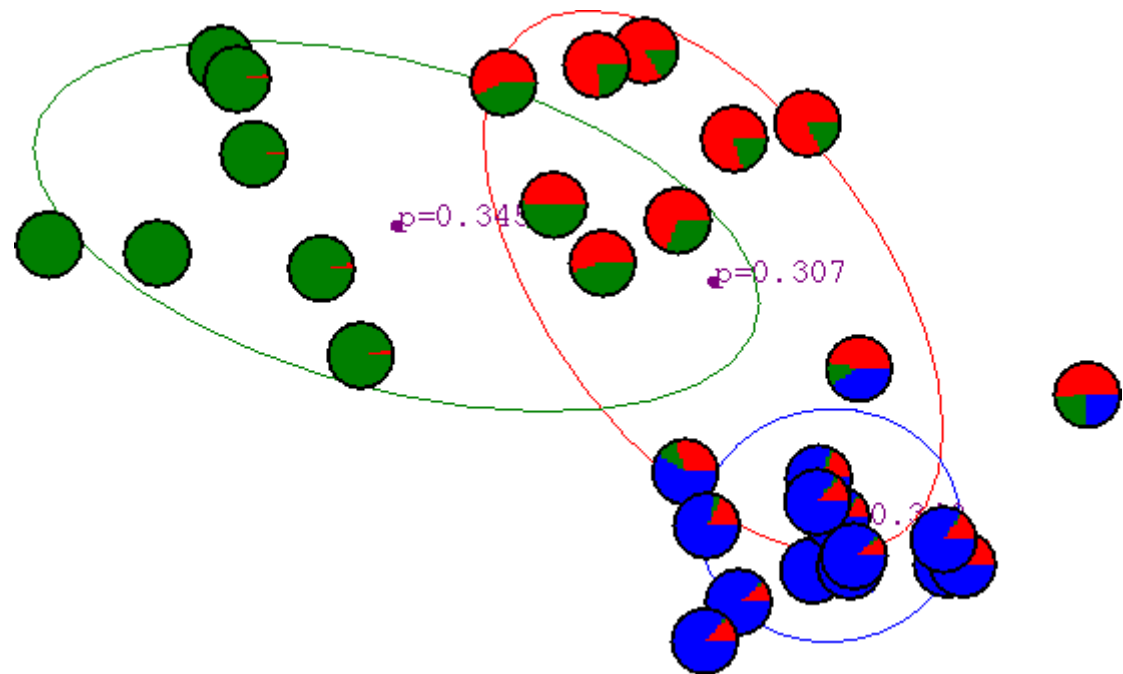
After first iteration



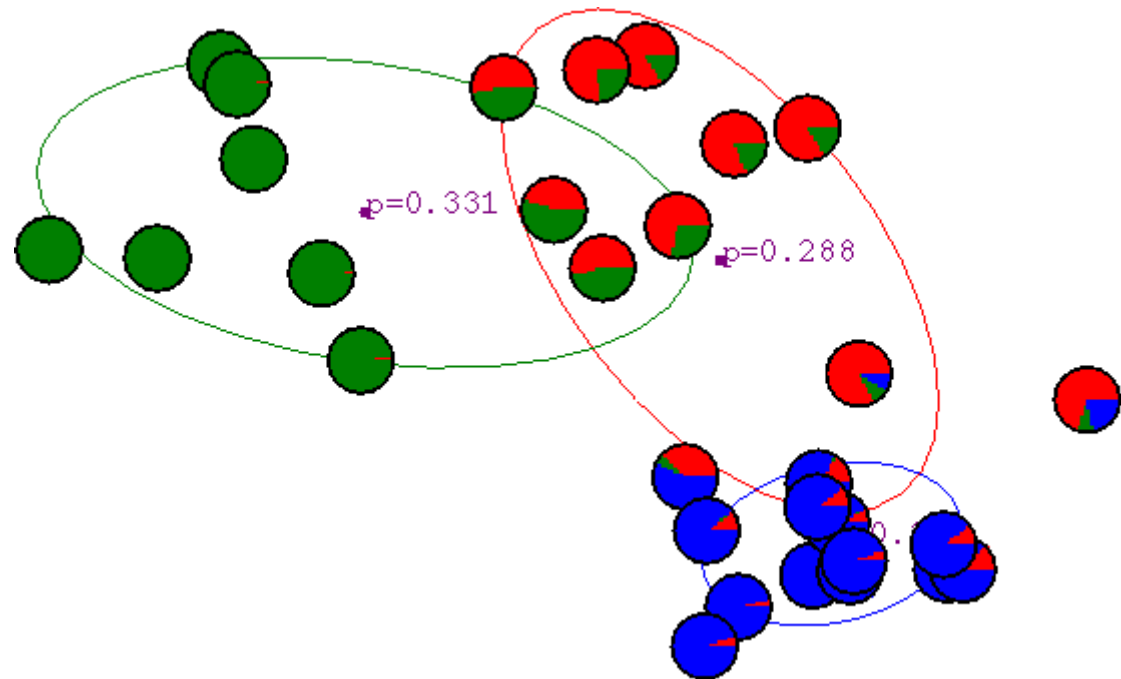
After 2nd iteration



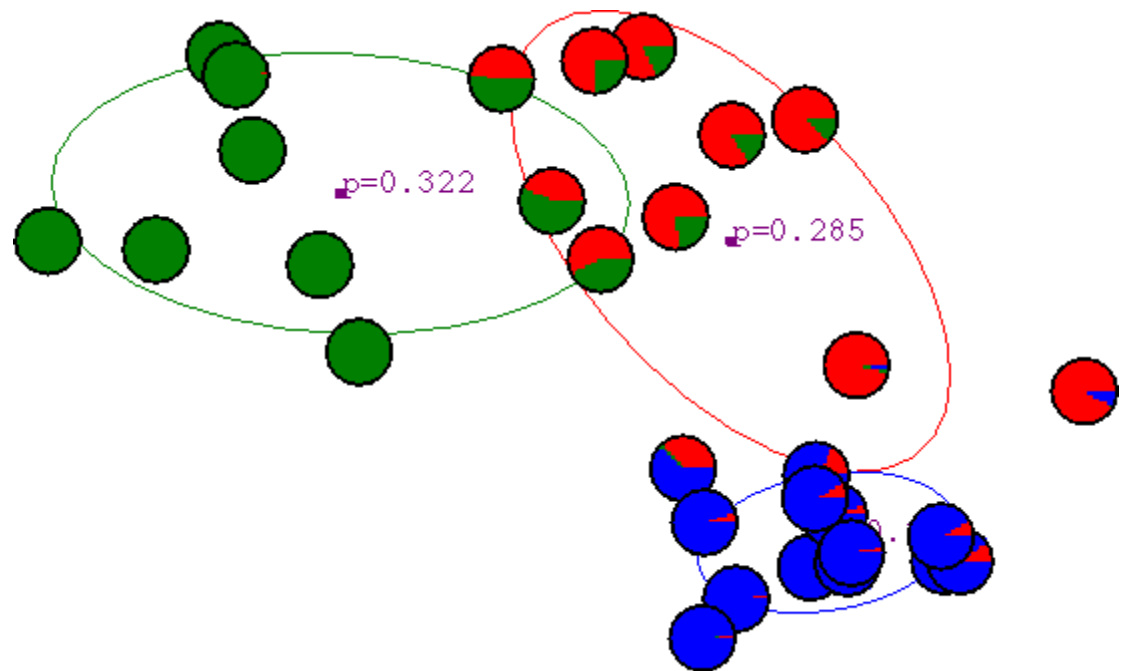
After 3rd iteration



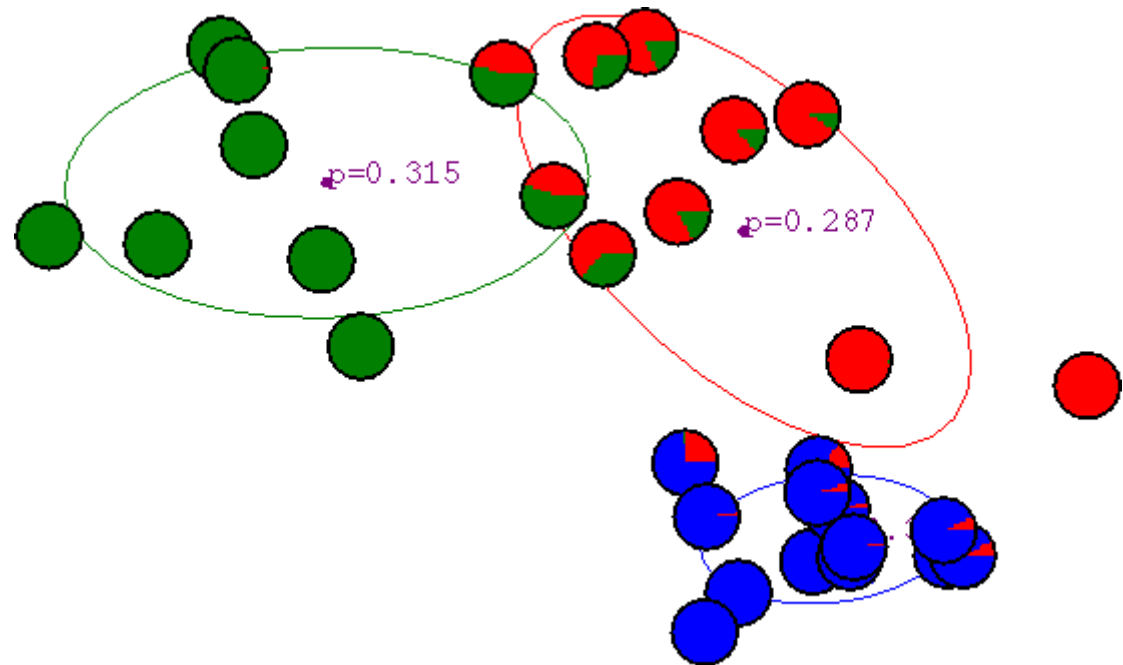
After 4th iteration



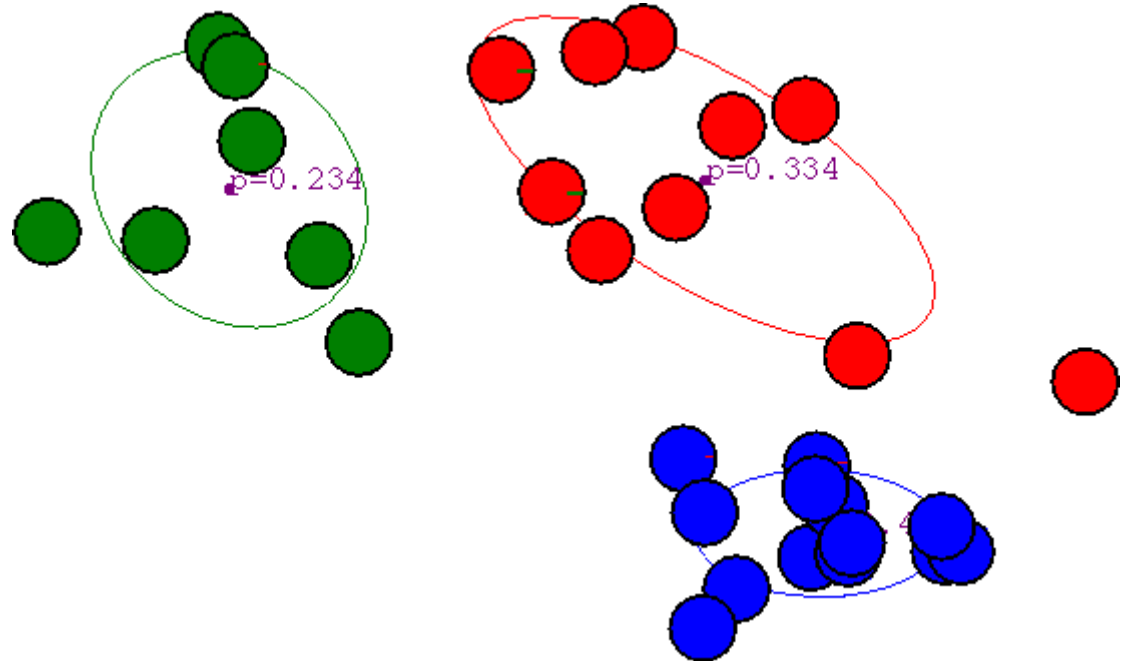
After 5th iteration



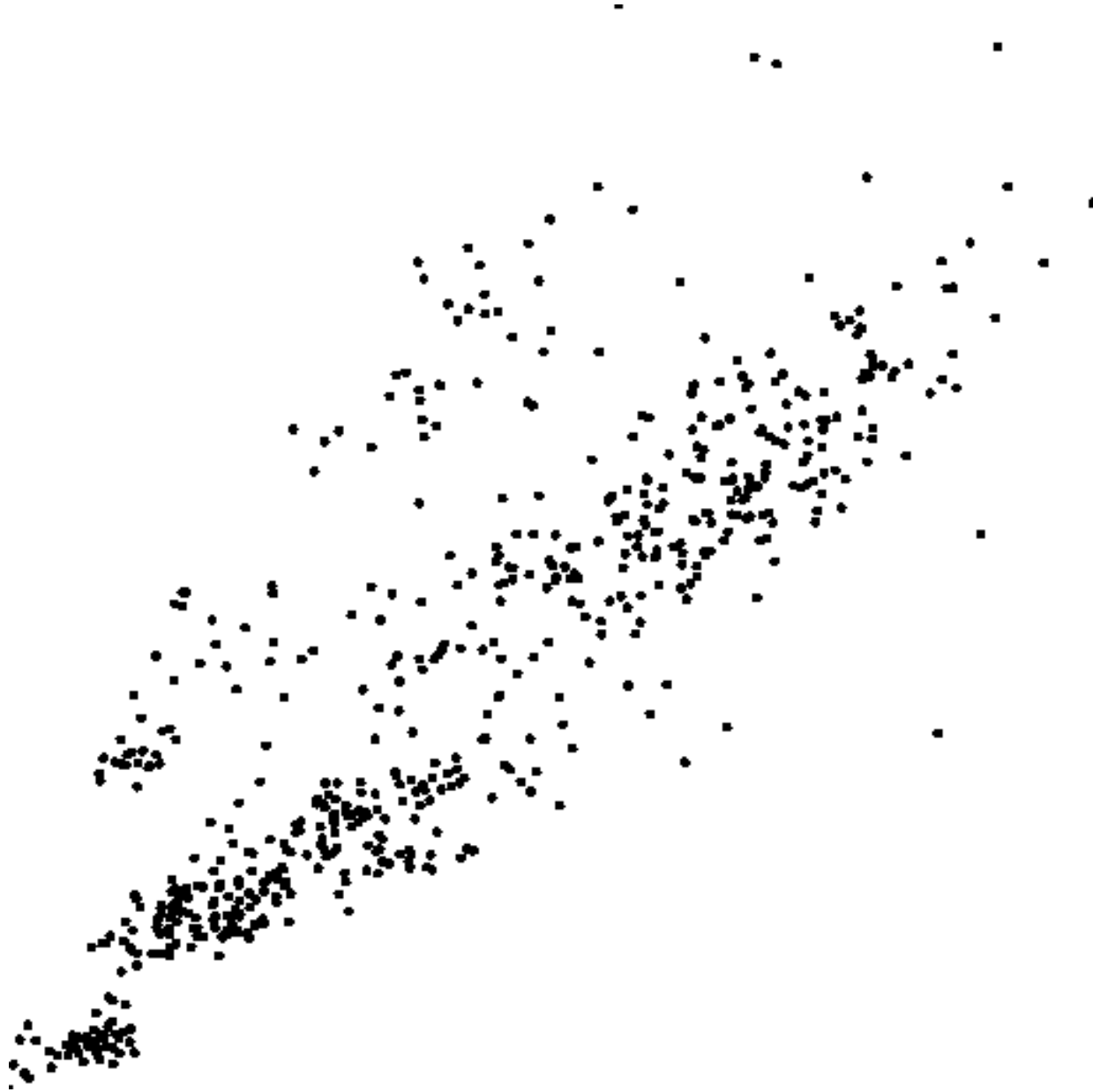
After 6th iteration



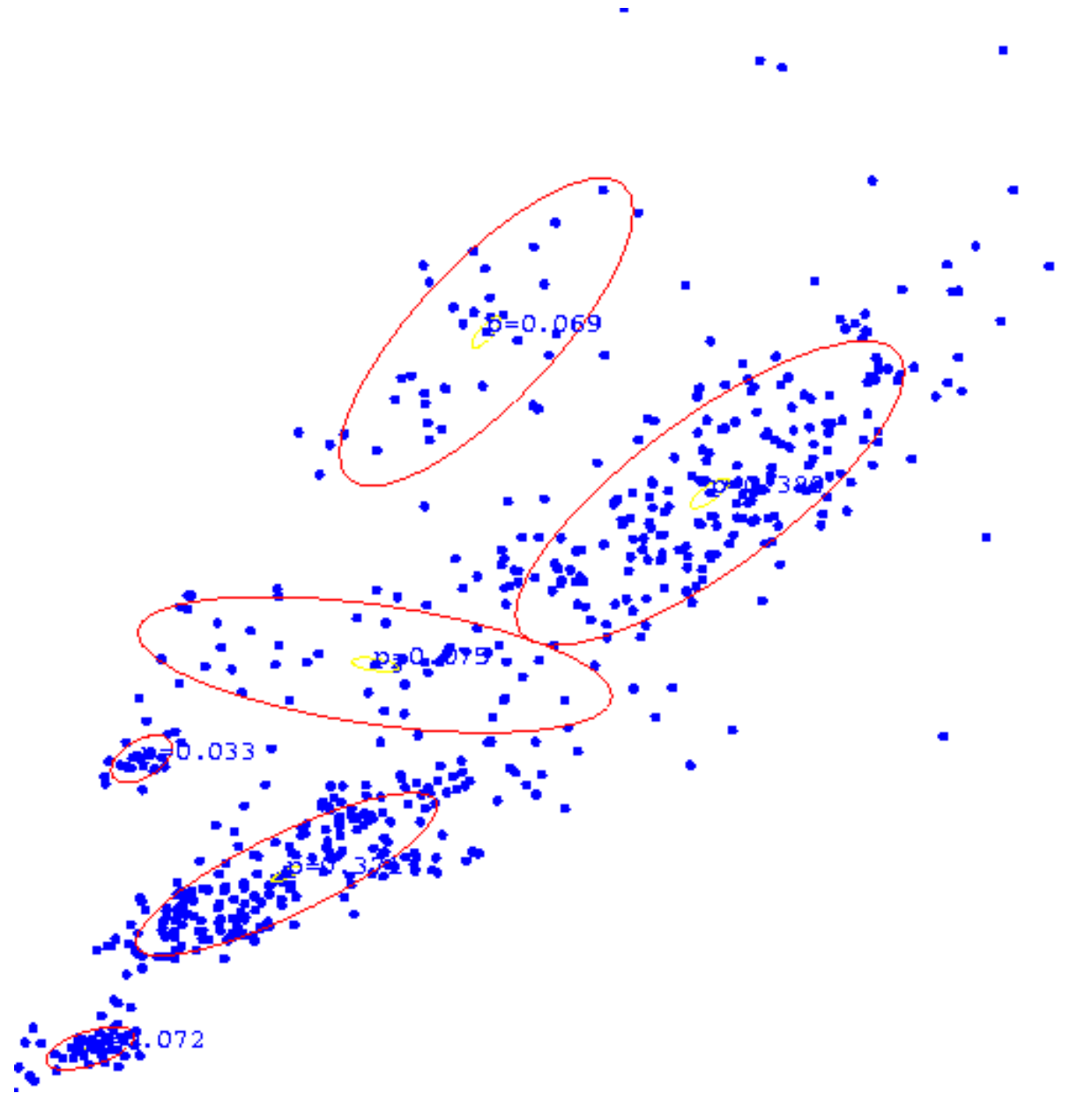
After 20th iteration



Some Bio Assay data



GMM clustering of the assay data



Resulting Density Estimator

