

Unsupervised learning

Lecture 13

David Sontag
New York University

Slides adapted from Carlos Guestrin, Dan Klein, Luke Zettlemoyer,
Dan Weld, Vibhav Gogate, and Andrew Moore

Gaussian Mixture Models

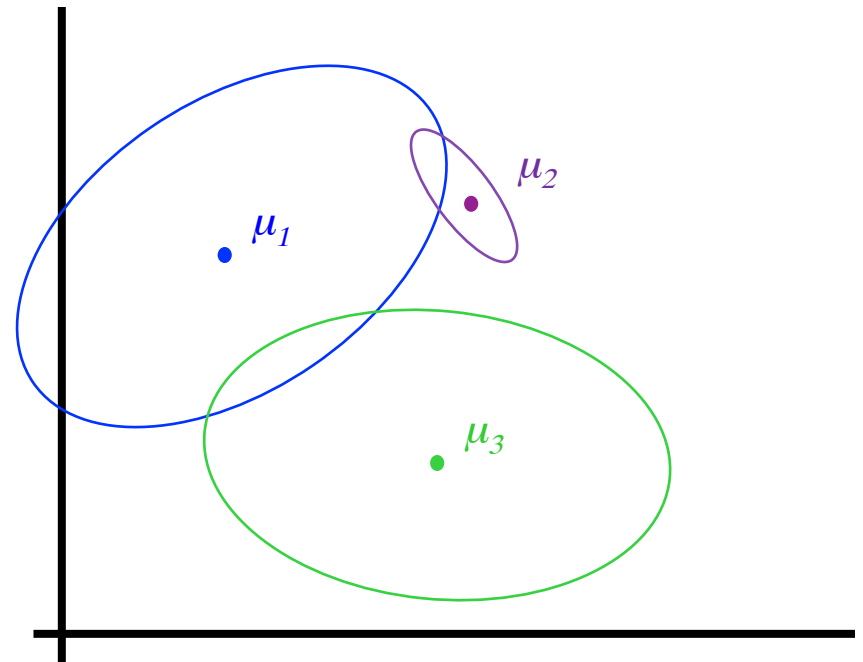
- $P(Y)$: There are k components
- $P(X|Y)$: Each component generates data from a **multivariate Gaussian** with mean μ_i and covariance matrix Σ_i

Each data point is sampled from a **generative process**:

1. Choose component i with probability $P(y=i)$ [Multinomial]
2. Generate datapoint $\sim N(\mu_i, \Sigma_i)$

$$P(X = \mathbf{x}_j | Y = i) =$$

$$\frac{1}{(2\pi)^{m/2} \|\Sigma_i\|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x}_j - \mu_i)^T \Sigma_i^{-1}(\mathbf{x}_j - \mu_i)\right]$$



ML estimation in supervised setting

- Univariate Gaussian

$$\mu_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i \quad \sigma_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2$$

- Mixture of Multivariate Gaussians**

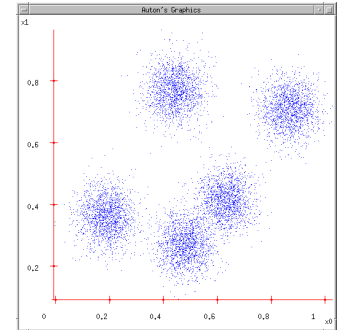
ML estimate for each of the Multivariate Gaussians is given by:

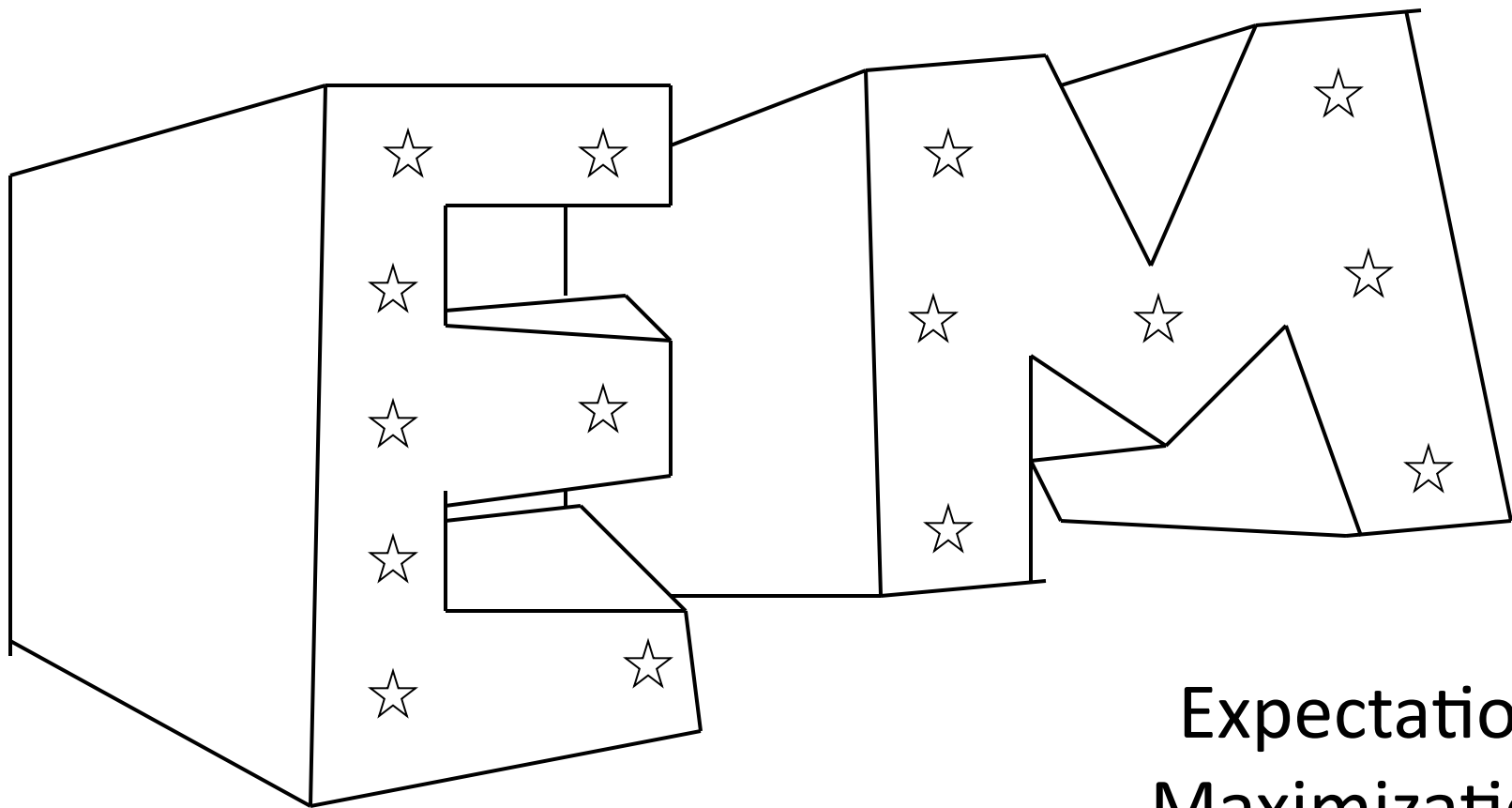
$$\mu_{ML}^k = \frac{1}{n} \sum_{j=1}^n x_j \quad \Sigma_{ML}^k = \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j - \mu_{ML}^k)(\mathbf{x}_j - \mu_{ML}^k)^T$$

Just sums over \mathbf{x} generated from the k 'th Gaussian

What about with unobserved data?

- Maximize ***marginal likelihood***:
 - $\operatorname{argmax}_{\theta} \prod_j P(x_j) = \operatorname{argmax} \prod_j \sum_{k=1}^K P(Y_j=k, x_j)$
- **Almost always a hard problem!**
 - Usually no closed form solution
 - Even when $\lg P(X,Y)$ is convex, $\lg P(X)$ generally isn't...
 - For all but the simplest $P(X)$, we will have to do gradient ascent, in a big messy space with lots of local optimum...





Expectation
Maximization

1977: Dempster, Laird, & Rubin

The EM Algorithm

- A clever method for maximizing marginal likelihood:
 - $\operatorname{argmax}_{\theta} \prod_j P(x_j) = \operatorname{argmax}_{\theta} \prod_j \sum_{k=1}^K P(Y_j=k, x_j)$
 - Based on coordinate descent. Easy to implement (eg, no line search, learning rates, etc.)
- Alternate between two steps:
 - Compute an expectation
 - Compute a maximization
- Not magic: ***still optimizing a non-convex function with lots of local optima***
 - The computations are just easier (often, significantly so!)

EM: Two Easy Steps

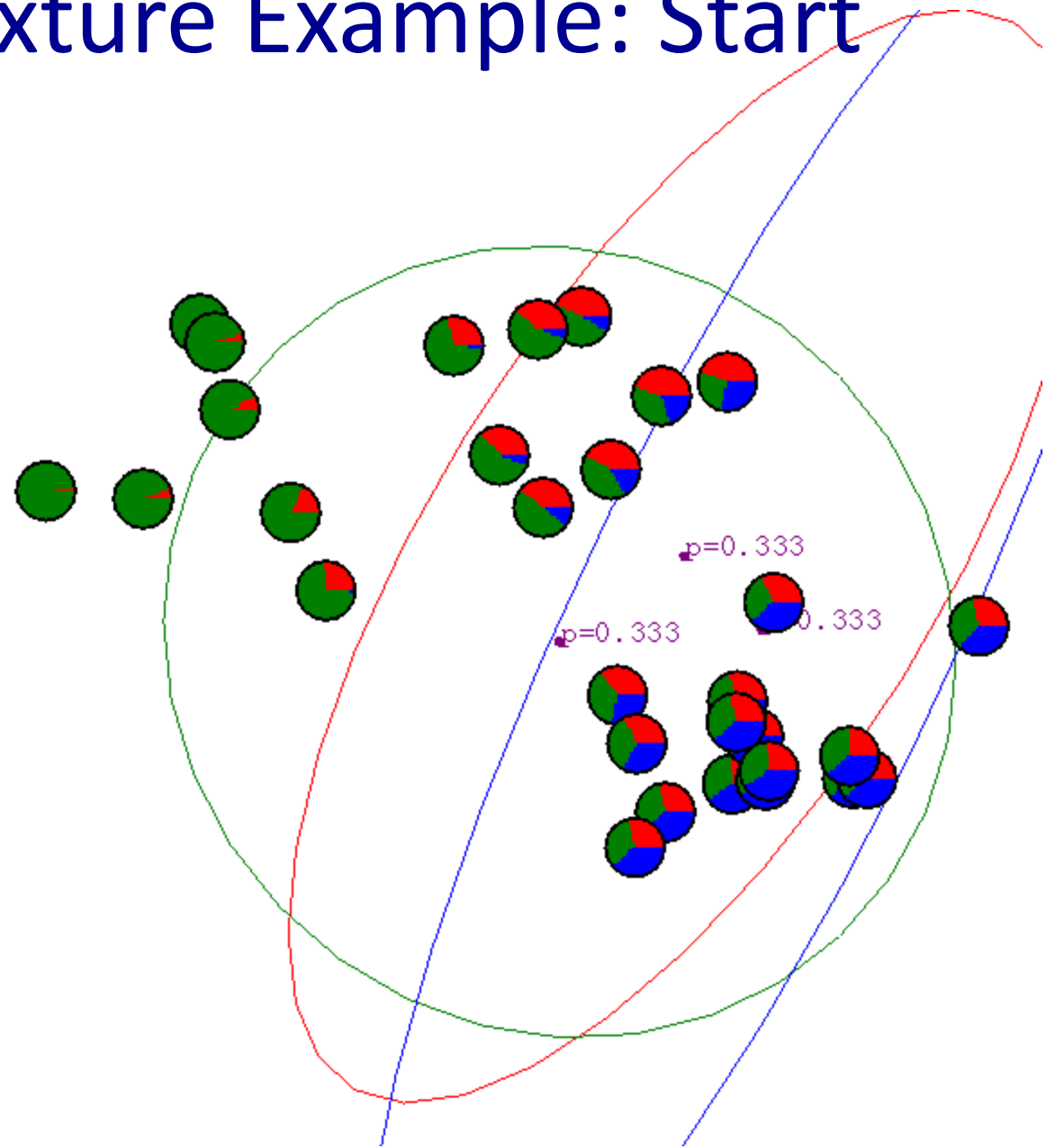
Objective: $\operatorname{argmax}_{\theta} \lg \prod_j \sum_{k=1}^K P(Y_j=k, x_j; \theta) = \sum_j \lg \sum_{k=1}^K P(Y_j=k, x_j; \theta)$

Data: $\{x_j \mid j=1 \dots n\}$

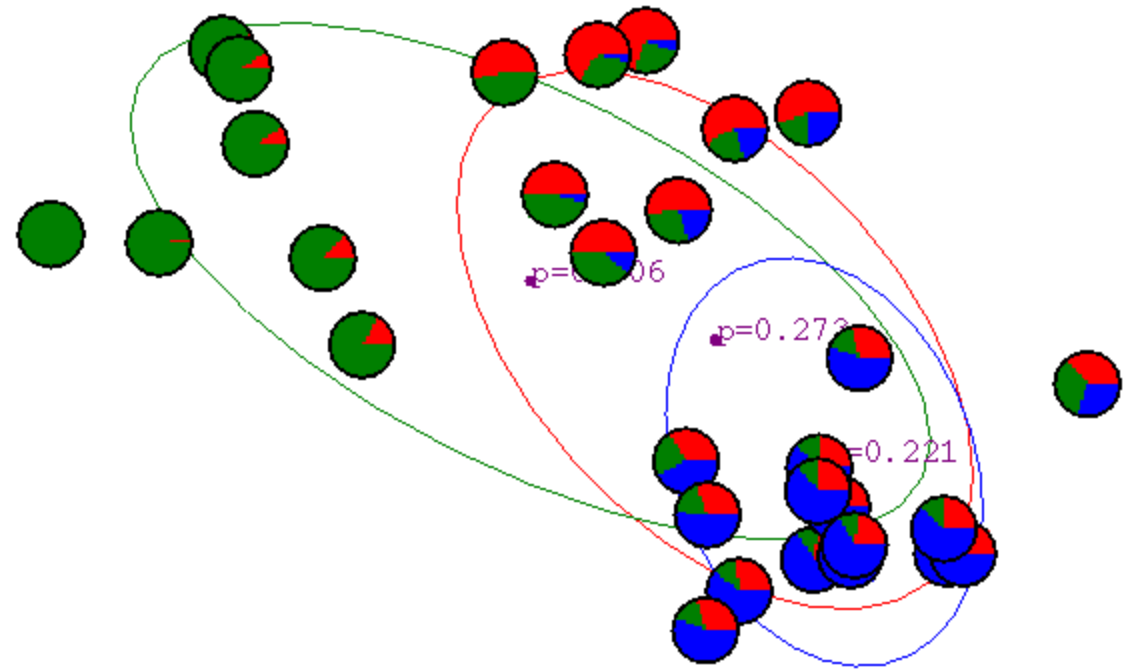
- **E-step:** Compute expectations to “fill in” missing y values according to current parameters, θ
 - For all examples j and values k for Y_j , compute: $P(Y_j=k \mid x_j; \theta)$
- **M-step:** Re-estimate the parameters with “weighted” MLE estimates
 - Set $\theta^{\text{new}} = \operatorname{argmax}_{\theta} \sum_j \sum_k P(Y_j=k \mid x_j; \theta^{\text{old}}) \log P(Y_j=k, x_j; \theta)$

Particularly useful when the E and M steps have closed form solutions

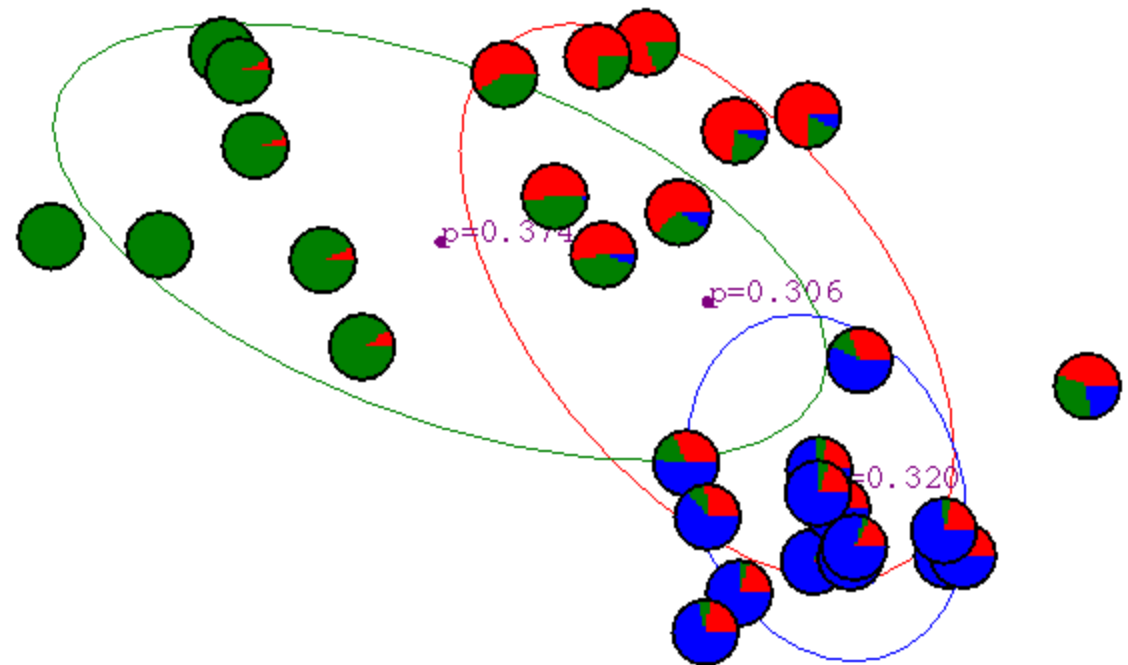
Gaussian Mixture Example: Start



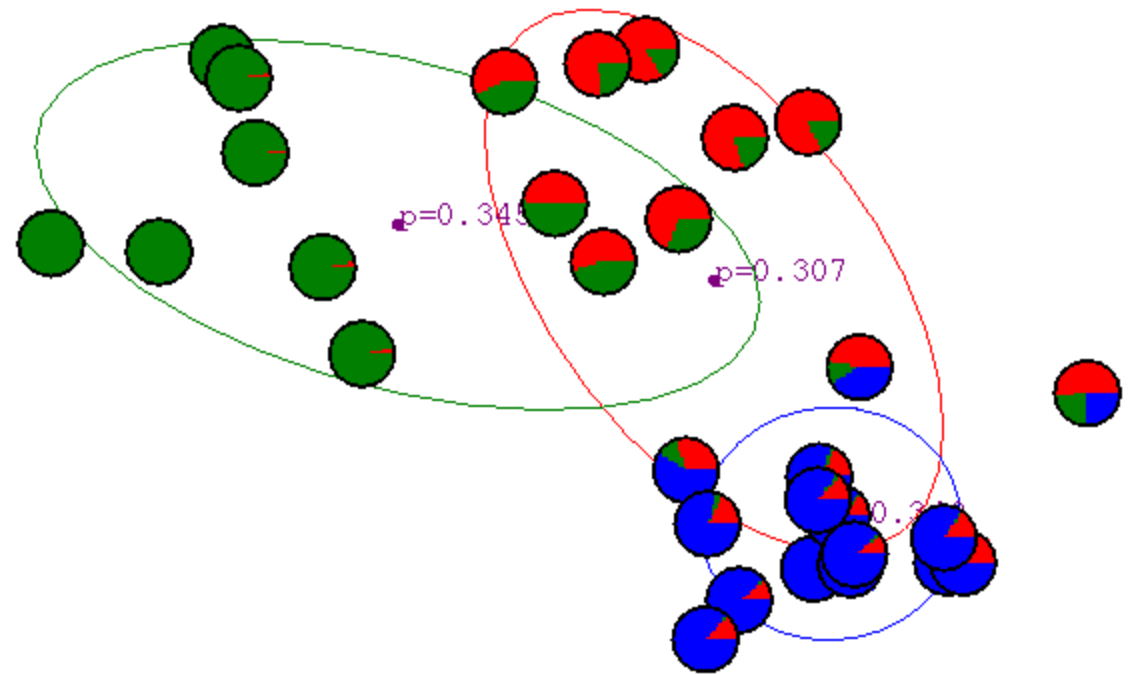
After first iteration



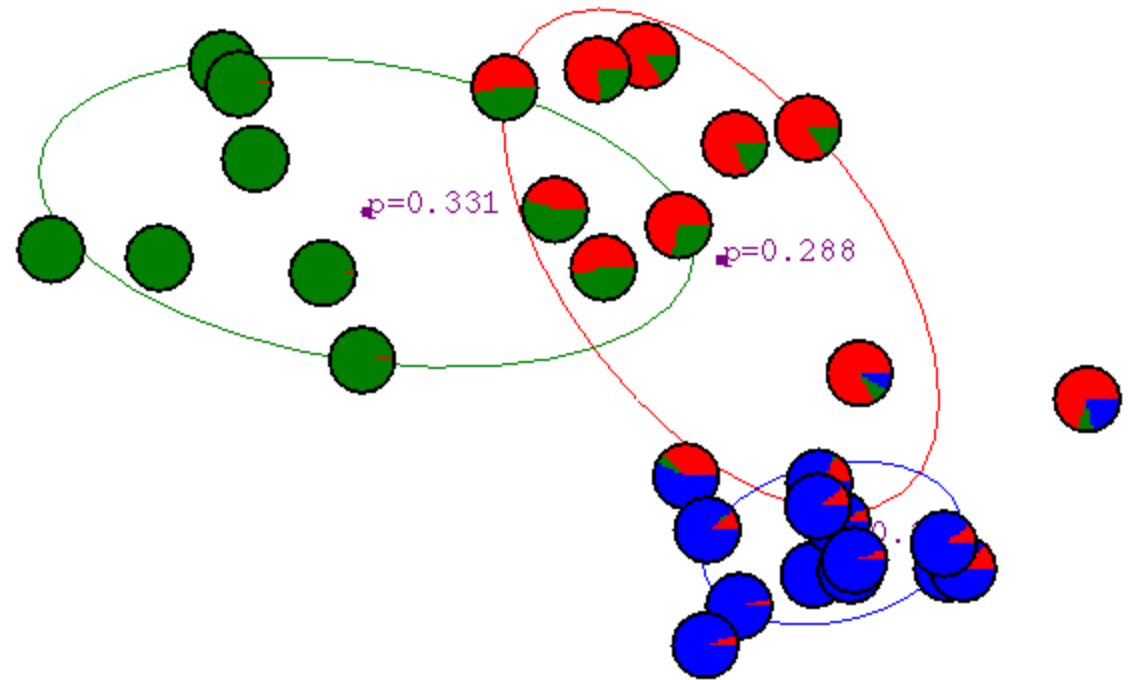
After 2nd iteration



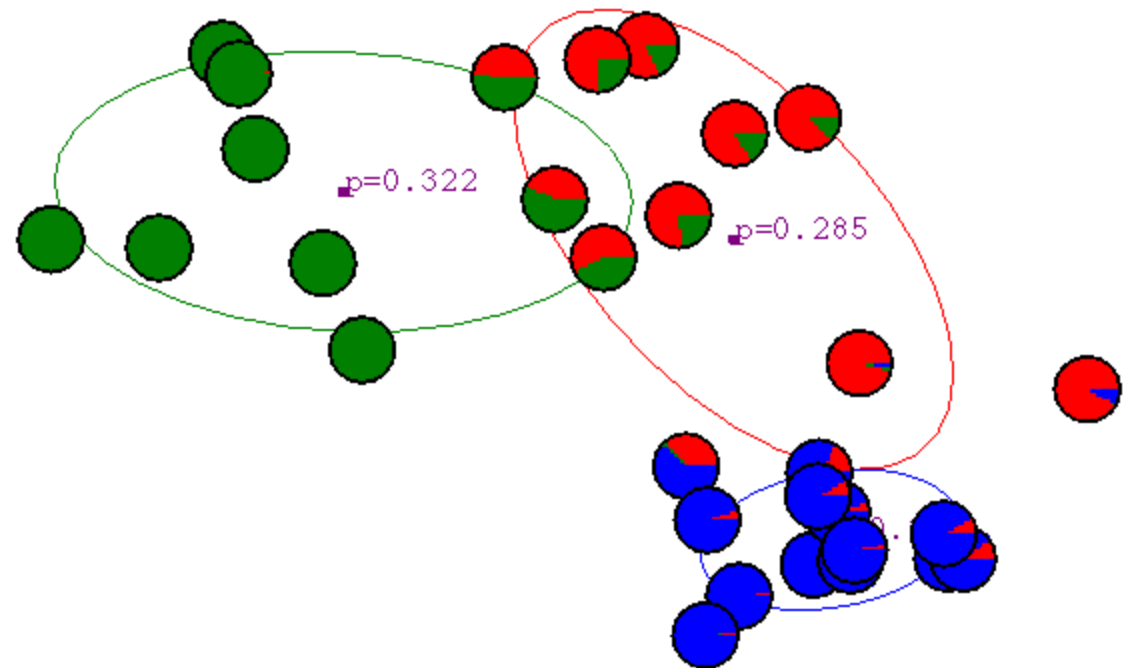
After 3rd iteration



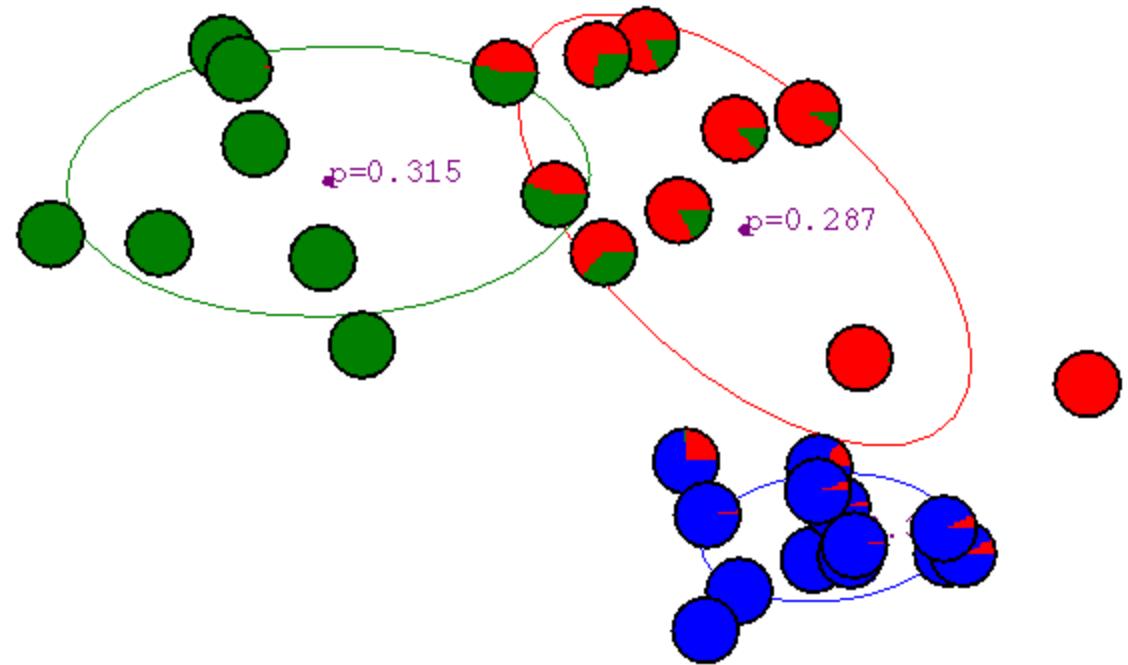
After 4th iteration



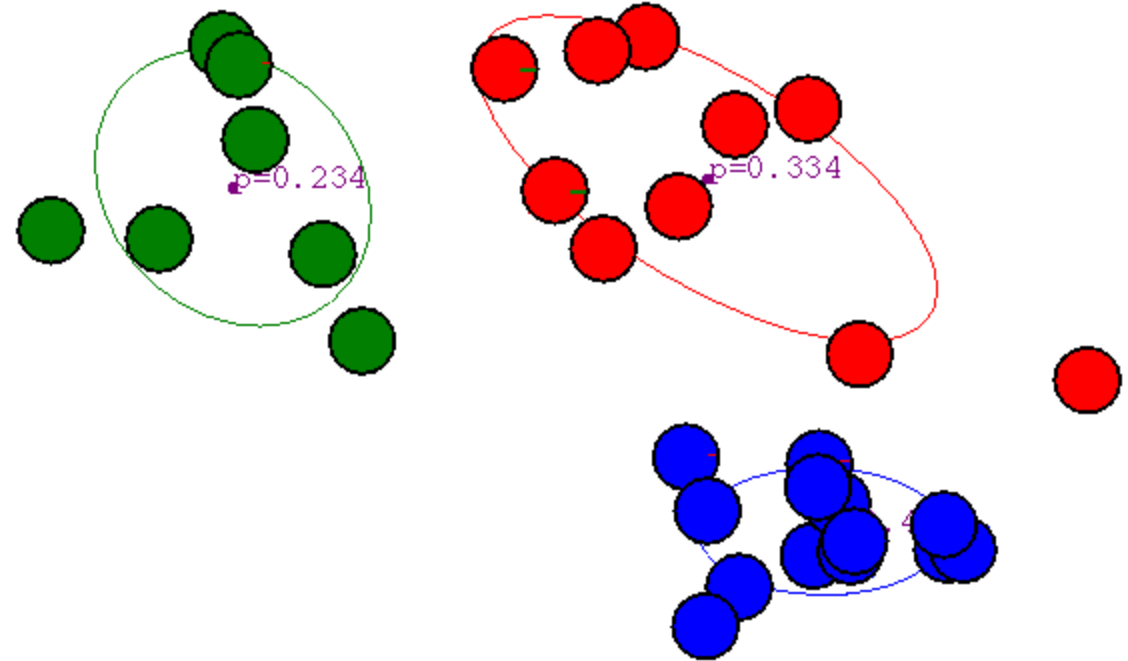
After 5th iteration



After 6th iteration



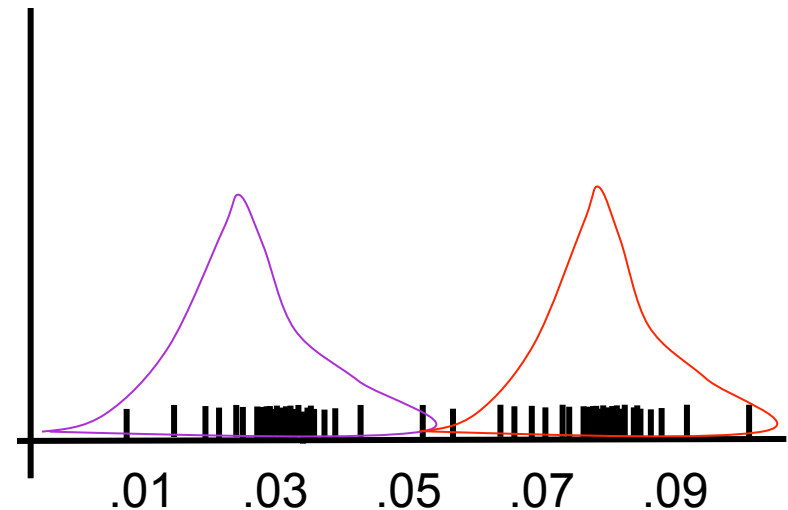
After 20th iteration



Simple example: learn means only!

Consider:

- 1D data
- Mixture of $k=2$ Gaussians
- Variances fixed to $\sigma=1$
- Distribution over classes is uniform
- Just need to estimate μ_1 and μ_2



$$\prod_{j=1}^m \sum_{k=1}^K P(x, Y_j = k) \propto \prod_{j=1}^m \sum_{k=1}^{K=2} \exp\left[-\frac{1}{2\sigma^2} \|x - \mu_k\|^2\right] P(Y_j = k)$$

EM for GMMs: only learning means

Iterate: On the t 'th iteration let our estimates be

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

E-step

Compute “expected” classes of all datapoints

$$P(Y_j = k | x_j, \mu_1 \dots \mu_K) \propto \exp\left(-\frac{1}{2\sigma^2} \|x_j - \mu_k\|^2\right) P(Y_j = k)$$

M-step

Compute most likely new μ s given class expectations

$$\mu_k = \frac{\sum_{j=1}^m P(Y_j = k | x_j) x_j}{\sum_{j=1}^m P(Y_j = k | x_j)}$$

E.M. for General GMMs

$p_k^{(t)}$ is shorthand for estimate of $P(y=k)$ on t 'th iteration

Iterate: On the t 'th iteration let our estimates be

$$\lambda_t = \{ \mu_1^{(t)}, \mu_2^{(t)} \dots \mu_K^{(t)}, \Sigma_1^{(t)}, \Sigma_2^{(t)} \dots \Sigma_K^{(t)}, p_1^{(t)}, p_2^{(t)} \dots p_K^{(t)} \}$$

E-step

Compute “expected” classes of all datapoints for each class

$$P(Y_j = k | x_j; \lambda_t) \propto p_k^{(t)} p(x_j; \mu_k^{(t)}, \Sigma_k^{(t)})$$

Just evaluate a Gaussian at x_j

M-step

Compute weighted MLE for μ given expected classes above

$$\mu_k^{(t+1)} = \frac{\sum_j P(Y_j = k | x_j; \lambda_t) x_j}{\sum_j P(Y_j = k | x_j; \lambda_t)} \quad \Sigma_k^{(t+1)} = \frac{\sum_j P(Y_j = k | x_j; \lambda_t) [x_j - \mu_k^{(t+1)}][x_j - \mu_k^{(t+1)}]^T}{\sum_j P(Y_j = k | x_j; \lambda_t)}$$

$$p_k^{(t+1)} = \frac{\sum_j P(Y_j = k | x_j; \lambda_t)}{m}$$

m = #training examples

What if we do hard assignments?

Iterate: On the t 'th iteration let our estimates be

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

E-step

Compute “expected” classes of all datapoints

$$P(Y_j = k | x_j; \mu_1 \dots \mu_K) \propto \exp\left(-\frac{1}{2\sigma^2} \|x_j - \mu_k\|^2\right) \cancel{P(Y_j = k)}$$

M-step

Compute most likely new μ s given class expectations

$$\mu_k = \frac{\sum_{j=1}^m P(Y_j = k | x_j) x_j}{\sum_{j=1}^m P(Y_j = k | x_j)} \quad \mu_k = \frac{\sum_{j=1}^m \delta(Y_j = k, x_j) x_j}{\sum_{j=1}^m \delta(Y_j = k, x_j)}$$

δ represents hard assignment to “most likely” or nearest cluster

Equivalent to k-means clustering algorithm!!!

Properties of EM

- We will prove that
 - EM converges to a local maxima
 - Each iteration improves the log-likelihood
- How? (Same as k-means)
 - E-step can never decrease likelihood
 - M-step can never decrease likelihood

The general learning problem with missing data

- **Marginal likelihood:** \mathbf{X} is observed,

\mathbf{Z} (e.g. the class labels \mathbf{Y}) is missing:

$$\begin{aligned}\ell(\theta : \mathcal{D}) &= \log \prod_{j=1}^m P(\mathbf{x}_j \mid \theta) \\ &= \sum_{j=1}^m \log P(\mathbf{x}_j \mid \theta) \\ &= \sum_{j=1}^m \log \sum_{\mathbf{z}} P(\mathbf{x}_j, \mathbf{z} \mid \theta)\end{aligned}$$

- **Objective:** Find $\text{argmax}_{\theta} \ell(\theta : \text{Data})$
- Assuming hidden variables are *missing completely at random* (otherwise, we should explicitly model *why* the values are missing)